

**EPA 540-R-01-003
OSWER 9285.7-41
September 2002**

Guidance for Comparing Background and Chemical Concentrations in Soil for CERCLA Sites

**Office of Emergency and Remedial Response
U.S. Environmental Protection Agency
Washington, DC 20460**



Recycled/Recyclable
Printed with Soy/Canola Ink on paper that
contains at least 50% recycled fiber

PREFACE

This document provides guidance to the U.S. Environmental Protection Agency Regions concerning how the Agency intends to exercise its discretion in implementing one aspect of the CERCLA remedy selection process. The guidance is designed to implement national policy on these issues.

Some of the statutory provisions described in this document contain legally binding requirements. However, this document does not substitute for those provisions or regulations, nor is it a regulation itself. Thus, it cannot impose legally binding requirements on EPA, States, or the regulated community, and may not apply to a particular situation based upon the circumstances. Any decisions regarding a particular remedy selection decision will be made based on the statute and regulations, and EPA decision makers retain the discretion to adopt approaches on a case-by-case basis that differ from this guidance where appropriate. EPA may change this guidance in the future.

ACKNOWLEDGMENTS

The EPA working group, chaired by Jayne Michaud (Office of Emergency and Remedial Response), included Tom Bloomfield (Region 9), Clarence Callahan (Region 9), Sherri Clark (Office of Emergency and Remedial Response), Steve Ells (Office of Emergency and Remedial Response), Audrey Galizia (Region 2), Cynthia Hanna (Region 1), Jennifer Hubbard (Region 3), Dawn Ioven (Region 3), Julius Nwosu (Region 10), Sophia Serda (Region 9), Ted Simon (Region 4), and Paul White (Office of Research and Development). David Bennett (Office of Emergency and Remedial Response) was the senior advisor for this working group. Comments and suggestions provided by Agency staff are gratefully acknowledged.

Technical and editorial assistance by Mary Deardorff and N. Jay Bassin of Environmental Management Support, Inc., and Harry Chmelynski of Sanford Cohen & Associates, Inc., are gratefully acknowledged.

<p>This report was prepared for the Office of Emergency and Remedial Response, United States Environmental Protection Agency. It was edited and revised by Environmental Management Support, Inc., of Silver Spring, Maryland, under contract 68-W6-0046, work assignment 007, and contract 68-W-02-033, work assignment 004, managed by Jayne Michaud. Mention of trade names or specific applications does not imply endorsement or acceptance by EPA. For further information, contact Jayne Michaud, U.S. EPA, Office of Emergency and Remedial Response, Mail Code 5202G, 1200 Pennsylvania Avenue, Washington, DC 20460.</p>
--

CONTENTS

	<u>Page</u>
ACRONYMS AND ABBREVIATIONS	vi
GLOSSARY	vii
CHAPTER 1: INTRODUCTION	1-1
1.1 Application of Guidance	1-1
1.2 Goals	1-2
1.3 Scope of Guidance	1-2
1.4 Intended Audience	1-2
1.5 Definition of Background	1-2
CHAPTER 2: SCOPING	2-1
2.1 When Background Samples Are Not Needed	2-1
2.2 When Background Samples Are Needed	2-2
2.3 Selecting a Reference Area	2-2
CHAPTER 3: HYPOTHESIS TESTING AND DATA QUALITY OBJECTIVES	3-1
3.1 Hypothesis Testing	3-1
3.1.1 Background Test Form 1	3-5
3.1.2 Background Test Form 2	3-7
3.1.3 Selecting a Background Test Form	3-8
3.2 Errors Tests and Confidence Levels	3-8
3.3 Test Performance Plots	3-9
3.4 DQO Steps for Characterizing Background	3-11
3.5 Sample Size	3-14
3.6 An Example of the DQO Process	3-15
CHAPTER 4: PRELIMINARY DATA ANALYSIS	4-1
4.1 Tests for Normality	4-2
4.2 Graphical Displays	4-2
4.2.1 Quantile Plot	4-3
4.2.2 Quantile-Quantile Plots	4-4
4.2.3 Quantile Difference Plot	4-5
4.3 Outliers	4-6
4.4 Censored Data (Non-Detects)	4-7
CHAPTER 5: COMPARING SITE AND BACKGROUND DATA	5-1
5.1 Descriptive Summary Statistics	5-2
5.2 Simple Comparison Methods	5-3
5.3 Statistical Methods for Comparisons with Background	5-3

5.3.1	Parametric Tests	5-4
5.3.2	Nonparametric Tests	5-6
5.4	Hypothesis Testing	5-12
5.4.1	Initial Considerations	5-12
5.4.2	Examples	5-12
5.4.3	Conclusions	5-14

APPENDIX A: SUPPLEMENTAL INFORMATION FOR DETERMINING “SUBSTANTIAL DIFFERENCE”		A-1
A.1	Precedents for Selecting a Background Test Form	A-1
A.2	Options for Establishing the Value of a Substantial Difference	A-4
A.2.1	Proportion of Mean Background Concentration	A-4
A.2.2	A Selected Percentile of the Background Distribution	A-5
A.2.3	Proportion of Background Variability	A-5
A.2.4	Proportion of Preliminary Remediation Goal	A-5
A.2.5	Proportion of Soil Screening Level	A-5
A.3	Statistical Tests and Confidence Intervals for Background Comparisons	A-6
A.3.1	Comparisons Based on the t-Test	A-6
A.3.2	Comparisons Based on the Wilcoxon Rank Sum Test	A-8

APPENDIX B: POLICY CONSIDERATIONS FOR THE APPLICATION OF BACKGROUND DATA IN RISK ASSESSMENT AND REMEDY SELECTION		B-1
	Purpose	B-3
	History	B-3
	Definitions of Terms	B-4
	Consideration of Background in Risk Assessment	B-5
	Consideration of Background in Risk Management	B-6
	Consideration of Background in Risk Communication	B-7
	Hypothetical Case Examples	B-7
	Hypothetical Case 1	B-8
	Hypothetical Case 2	B-9
	Hypothetical Case 3	B-9
	References	B-10

EXHIBITS

	<u>Page</u>
Figure 2.1 Determining the need for background sampling	2-1
Figure 3.1 Test performance plot: site is not significantly different from background	3-10
Figure 3.2 Test performance plot: site does not exceed background by more than S	3-11
Figure 4.1 Example of a double quantile plot	4-3
Figure 4.2 Example of an empirical quantile-quantile plot	4-4
Figure 4.3 Example of a quantile difference plot	4-5
Table 3.1 Required sample size for selected values of σ	3-4
Table 3.2 Achievable values of $\alpha = \beta$ for selected values of N	3-5
Table 3.3 Hypothesis Testing: Type I and Type II Errors	3-8
Table 5.1 Site data	5-8
Table 5.2 Background data	5-8
Table 5.3 WRS test for Test Form 1 (H_0 : site < background)	5-8
Table 5.4 Critical Values for the WRS Test	5-9
Table 5.5 WRS test for Test Form 2 (H_0 : site > background + 100)	5-10
Table 5.6 WRS test for Test Form 2 (H_0 : site > background + 50)	5-10
Table 5.7 Summary of hypothesis tests	5-15

ACRONYMS AND ABBREVIATIONS

α	Alpha Decision Error (Type I)
ANOVA	Analysis of Variance Procedure
ARAR	Applicable or Relevant and Appropriate Requirements
β	Beta Decision Error (Type II)
CERCLA	Comprehensive Environmental Response, Compensation, and Liability Act
cm	Centimeter
COC	Chemical of Concern
COPC	Chemical of Potential Concern
CV	Coefficient of Variation
Δ	Delta (Difference)
DCGL	Design Concentration Guideline Level
DQO	Data Quality Objective
EPA	U.S. Environmental Protection Agency
H_A	Alternative Hypothesis
H_0	Null Hypothesis
HRS	Hazard Ranking System
K	Tolerance Coefficient
kg	Kilogram
LBGR	Lower Bound of the Gray Region
m	Meter
M_b	Mean Background Concentration
MDD	Minimum Detectable Difference
mg	Milligram
N (n)	Number of Samples
ND	Non-Detect Measurement
NRC	U.S. Nuclear Regulatory Commission
PA/SI	Preliminary Assessment/Site Investigation
PRG	Preliminary Remediation Goal
RAGS	Risk Assessment Guidance for Superfund Vol. I, Human Health Evaluation Manual (Part A)
RCRA	Resource Conservation and Recovery Act
RPM	Remedial Project Manager
σ	Standard Deviation
S	Substantial Difference
SAP	Sampling and Analysis Plan
SSL	Soil Screening Level
TL	Tolerance Limit
TRW	Technical Review Workgroup for Lead
UTL	Upper Tolerance Limit
WRS	Wilcoxon Rank Sum

GLOSSARY

Background	Substances or locations that are not influenced by the releases from a site and are usually described as naturally occurring or anthropogenic: (1) Naturally occurring substances present in the environment in forms that have not been influenced by human activity. (2) Anthropogenic substances are natural and human-made substances present in the environment as a result of human activities (not specifically related to the CERCLA site in question).
Background reference area	The area where background samples are collected for comparison with samples collected on site. The reference area should have the same physical, chemical, geological, and biological characteristics as the site being investigated, but has not been affected by activities on the site.
Background Test Form 1	Within this guidance, the null hypothesis that the mean concentration in potentially impacted areas <i>is less than or equal to</i> the mean background concentration.
Background Test Form 2	Within this guidance, the null hypothesis that the mean concentration in potentially impacted areas <i>exceeds</i> the mean background concentration.
Coefficient of variation	The ratio of the standard deviation to the mean. A unitless measure that allows the comparison of dispersion across several sets of data. It is often used instead of the standard deviation in environmental applications because the standard deviation is often proportional to the mean.
Δ (delta)	The true difference between the mean concentration of chemical X in potentially impacted areas and the mean background concentration of chemical X. Delta is an unknown parameter which describes the true state of nature. Hypotheses about its value are evaluated using statistical hypothesis tests. In principle, we can select any specific value for Δ and then test if the observed difference is as large as Δ or not with a given confidence and power.
Detection limit	Smallest concentration of a substance that can be distinguished from zero.
Gehan test	The Gehan test is a generalized version of the WRS test. The Gehan test addresses multiple detection limits using a modified ranking procedure rather than relying on the “all ties get the same rank” approach used in the WRS test.
Gray region	A range of possible values of Δ where the consequences of making a decision error are relatively minor—where the statistical test will yield inconclusive results. The width of the gray region is equal to the MDD for the test. The location of the gray region depends on the type of statistical test selected.
Hypothesis (statistical)	A statement that may be supported or rejected by examining relevant data. To determine if we should accept a hypothesis, it is commonly easier to attempt to reject its converse (that is, first assume that the hypothesis is not true). This assumption to be tested is called the null hypothesis (H_0), which is any testable presumption set up to be rejected. An alternative hypothesis (H_A) is the logical opposite of the null hypothesis.

Hypothesis testing	A quantitative method to determine whether a specific statement concerning Δ (called the null hypothesis) can be rejected or not by examining data. The hypothesis testing process provides a formal procedure to quantify the decision maker's acceptable limits for decision errors.
Judgmental (or authoritative) samples	Samples collected in areas suspected to have higher contaminant concentrations due to operational or historical knowledge. Judgmental samples cannot be extrapolated to represent the entire site.
MDD (minimum detectable difference)	The smallest difference in means that the statistical test can resolve. The MDD depends on sample-to-sample variability, the number of samples, and the power of the statistical test. The MDD is a property of the survey design.
Nonparametric data analysis	A distribution-free statistical method that does not depend on knowledge of the population distribution.
Outliers	Measurements (usually larger or smaller than other data values) that are not representative of the sample population from which they were drawn. They distort statistics if used in any calculations.
Parametric data analysis	A statistical method that relies on a known probability distribution for the population from which data are selected. Parametric statistical tests are used to evaluate statements (hypotheses) concerning the parameters of the distribution. They are usually based on the assumption that the raw data are normally or lognormally distributed.
P-value	The smallest value of α at which the null hypothesis would be rejected for the given observations. The p-value of the test is sometimes called the critical level, or the significance level, of the test.
Quantile plot	A graph that displays the entire distribution of data, ranging from the lowest to the highest value. The vertical axis is the measured concentration, and the horizontal axis is the percentile of the distribution.
Quantile test	The quantile test is a nonparametric test specifically designed to compare the upper tails of two distributions. The quantile test may detect differences that are not detected by the Wilcoxon rank sum test.
Robustness	A method of comparing statistical tests. A robust test is one with good performance (that is not unduly affected by outliers) for a wide variety of data distributions.
S (substantial difference)	A difference in mean concentrations that is sufficiently large to warrant additional interest based on health or ecological information. S is the investigation level. If Δ exceeds S, the difference in concentrations is judged to be sufficiently large to be of concern, for the purpose of the analysis. A hypothesis test uses measurements from the site and from background to determine if Δ exceeds S.
Test performance plot	A graph that displays the combined effects of the decision error rates, the gray region for the decision-making process, and the level of substantial difference between site and background. It is used in the data quality objective process during scoping to aid in the selection of reasonable values for the decision error rates (α and β), the MDD, and the required number of samples.

Tolerance limit	A confidence limit on a percentile of the population rather than a confidence limit on the mean. For example, a 95 percent one-sided TL for 95 percent coverage represents the value below which 95 percent of the population are expected to fall (with 95 percent confidence).
Type I error	The probability, referred to as α (<i>alpha</i>), that the null hypothesis will be rejected when in fact it is true (false positive).
Type II error	The probability, referred to as β (<i>beta</i>), that the null hypothesis will be accepted when in fact it is false (false negative).
Walsh's test for outliers	A nonparametric test for determining the presence of outliers in either the background or onsite data sets.
Wilcoxon rank sum (WRS) test	A nonparametric test that examines whether measurements from one population consistently tend to be larger (or smaller) than those from the other population. It is used for determining whether a substantial difference exists between site and background population distributions.

CHAPTER 1

INTRODUCTION

The U.S. Environmental Protection Agency (EPA) developed this document to assist CERCLA remedial project managers (RPMs) and human health and ecological risk assessors during the remedial investigation process to evaluate background concentrations at CERCLA sites. An issue that is often raised at CERCLA sites is whether a reliable representation of background has been established.¹ This document recommends statistical methods for characterizing background concentrations of chemicals in soil.

The general application of background concentrations during the CERCLA remedial investigation process is addressed in EPA policy.² Ecological risk assessment guidance also provides specific recommendations for applying background concentration data.³

This document supplements Agency guidance included in the *Risk Assessment Guidance for Superfund Vol. I, Human Health Evaluation Manual (Part A)* (RAGS).¹ RAGS contains useful guidance on background issues that the reader should also consult:

- ▶ Sampling needs (Sections 4.4 and 4.6);
- ▶ Statistical methods (Section 4.4);
- ▶ Exposure assessment (Section 6.5); and
- ▶ Risk characterization (Section 8.6).

This document draws upon many other publications and statistical references. In general, background may play a role in the CERCLA process when:

- ▶ Determining whether a release falls within the limitation contained in Section 104(a)(3)(A) of

the Comprehensive Environmental Response, Compensation, and Liability Act (CERCLA), which addresses naturally occurring substances in their unaltered form from a location where they are naturally found;⁴

- ▶ Developing remedial goals;⁵
- ▶ Characterizing risks from contaminants that may also be attributed to background sources; and
- ▶ Communicating cumulative risks associated with the CERCLA site.

As stated in RAGS, a statistically significant difference between background samples and site-related contamination should not, by itself, trigger a cleanup action. Risk assessment methods should be applied to ascertain the significance of the chemical concentrations. EPA's national policy clarifying the role of background characterization in the CERCLA risk assessment and remedy selection process is included as Appendix B in this document.

1.1 Application of Guidance

This guidance should be applied on a site-specific basis, with assistance from a statistician who is familiar with the CERCLA remedial investigation process. Not every CERCLA site investigation will need to characterize chemicals in background areas. A background evaluation usually is considered when certain contaminants that pose risks and may drive an action are believed to be attributable to background. The need for background characterization, the timing of sampling efforts, and the required level

of effort should be determined on a site-specific basis. EPA should consider whether collecting background samples is necessary (Chapter 2); when, where, and how to collect background samples (Chapter 3); and how to evaluate the data (Chapters 4 and 5).

To the extent practicable, this guidance may also be applicable to sites addressed under removal actions, especially non-time-critical removal actions, and Resource Conservation and Recovery Act (RCRA) corrective actions.

1.2 Goals

The general goals of this guidance are to:

- ▶ Provide a practical guide for characterizing background concentrations at CERCLA sites; and
- ▶ Present sound options for evaluating background data sets in comparison to site contamination data.

1.3 Scope of Guidance

This guidance pertains to the evaluation of chemical contamination in soil at CERCLA sites. This guidance may be updated in the future to address non-soil media. Non-soil media are dynamic and influenced by upstream or upgradient sources. Such media—air, groundwater, surface water, and sediments—typically require additional analyses of release and transport, involve more complex spatial and temporal sampling strategies, and require different ways of combining and analyzing data.¹

The user should consult the available Agency guidances and policies when dealing with sites with radioactive contaminants. Certain types of CERCLA sites, such as mining or dioxin-contaminated sites, may require consideration of specific Agency policies and regulations. Therefore, this guidance should be applied on a case-by-case basis, with

consideration of Agency statutes, regulations, and policies.

1.4 Intended Audience

The intended audience of this guidance is EPA human health and ecological risk assessors, RPMs, and decision makers.

1.5 Definition of Background

For the purposes of this guidance, *background* refers to substances or locations that are not influenced by the releases from a site, and are usually described as naturally occurring or anthropogenic:^{1,6}

- 1) *Naturally occurring* – substances present in the environment in forms that have not been influenced by human activity; and,
- 2) *Anthropogenic* – natural and human-made substances present in the environment as a result of human activities (not specifically related to the CERCLA site in question).

Some chemicals may be present in background as a result of both natural and man-made conditions (such as naturally occurring arsenic and arsenic from pesticide applications or smelting operations).

CERCLA site activity (such as waste disposal practices) may cause naturally occurring substances to be released into other environmental media or chemically transformed. The concentrations of the released naturally occurring substance may not be considered as representative of natural background according to CERCLA 104(a)(3)(A).

Generally, the type of background substance (natural or anthropogenic) does not influence the statistical or technical method used to characterize background concentrations. For comparison purposes soil samples should have the same basic characteristics as the site sample (i.e., similar soil depths and soil types).⁷ (See Section 2.3).

CHAPTER NOTES

1. U.S. Environmental Protection Agency (EPA). 1989. *Risk Assessment Guidance for Superfund Vol. I, Human Health Evaluation Manual (Part A)*. Office of Emergency and Remedial Response, Washington, DC. EPA 540-1-89-002. Hereafter referred to as “RAGS.” For information on non-soil media, see Sections 4.5 and 6.5.
2. U.S. Environmental Protection Agency (EPA). April 2002. *Role of Background in the CERCLA Cleanup Program*. Office of Emergency and Remedial Response, Washington, DC. OSWER 9285.6-07P (see Appendix B of this guidance).
3. U.S. Environmental Protection Agency (EPA). 2001. *The Role of Screening-Level Risk Assessments and Refining Contaminants of Concern in Baseline Ecological Risk Assessments*. Office of Solid Waste and Emergency Response, Washington, DC.
4. CERCLA 104(a)(3)(A) restricts the authority to take an action in response to the release or threat of release of a “naturally occurring substance in its unaltered form or altered solely through naturally occurring processes or phenomena, from a location where it is naturally found.”
5. The National Oil and Hazardous Substances Pollution Contingency Plan (NCP) (40 CFR Part 300) is the primary regulation that implements CERCLA. The preamble to the NCP discusses the use of background levels for setting cleanup levels for constituents at CERCLA sites.

“...In some cases, background levels are not necessarily protective of human health, such as in urban or industrial areas; in other cases, cleaning up to background levels may not be necessary to achieve protection of human health because the background level for a particular contaminant may be close to zero, as in pristine areas” (55 FR 8717-8718).

The preamble to the NCP also identifies background as a technical factor to consider when determining an appropriate remedial level:

“Preliminary remediation goals...may be revised to a different risk level within the acceptable risk range based on the consideration of appropriate factors including, but not limited to: exposure factors, uncertainty factors, and technical factors...Technical factors may include...background levels of contaminants...”(55 FR 8717).

6. U.S. Environmental Protection Agency (EPA). 1995. *Engineering Forum Issue Paper. Determination of Background Concentrations of Inorganics in Soils and Sediments at Hazardous Waste Sites*. R.P Breckenridge and A.B. Crockett. Office of Research and Development, Washington, DC. EPA/540/S-96/500.
7. U.S. Environmental Protection Agency (EPA). July 2000. Draft *Ecological Soil Screening Level Guidance*. Office of Emergency and Remedial Response, Washington, DC. EPA 540/F-01/014. OSWER 9345.0-14.

CHAPTER 2

SCOPING

A first step in determining the need for background sampling data is gathering and evaluating all of the available data. Some information gathered during the Preliminary Assessment/Site Investigation (PA/SI) may provide data on background levels of chemicals. The SI usually provides the first opportunity to collect some background samples. Data collected and assessed for the hazard ranking system (HRS) process may include both site-related contaminants and off-site (or estimated background) substances. These data are generally limited in quantity and sample location and may have limited value in the remedial investigation. The locations of all data should be identified and reported when these data are considered during the remedial investigation. Sampling locations should be recorded with sufficient precision to permit follow up confirmatory measurements if required at a later date. The general types of information to consider when determining the need for background sampling are highlighted in the box below and in Figure 2.1.

Background Sampling Considerations

- ▶ Natural variability of soil types
- ▶ Operational practices
- ▶ Waste type
- ▶ Contaminant mobility

Information from preliminary site studies or published sources (regional or local data from the state or U.S. Geological Survey) may be useful for identifying local soil, water, and air quality characteristics.¹ Data from these resources may be useful for qualitative analyses of regional conditions. However, usually they are not sufficient to assess site-specific conditions in a quantitative manner.²

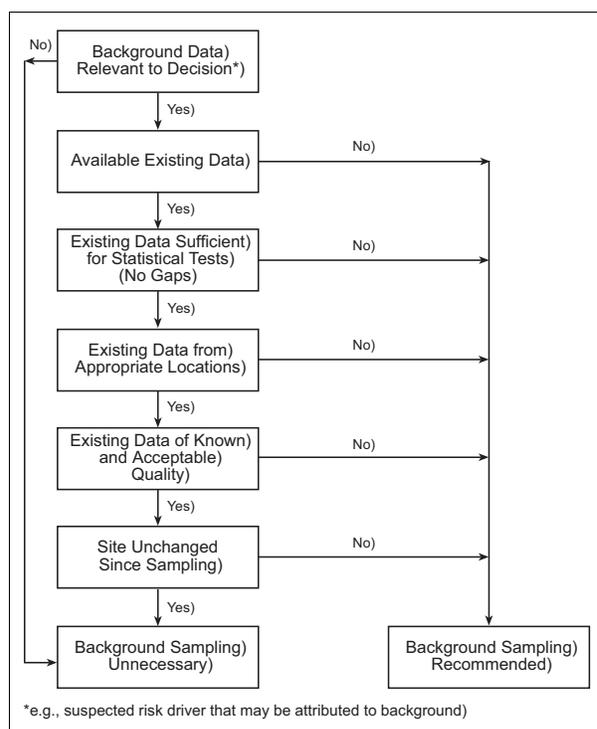


Figure 2.1 Determining the need for background sampling

After compiling and considering the relevant information, EPA should determine if the data are sufficient for the risk assessment and risk management decisions, or if additional site-specific data should be collected to characterize background.

2.1 When Background Samples Are Not Needed

If the sample quantity, location, and quality of existing data can be used to characterize background chemical concentrations and compare them to site data, then additional samples may not be needed. In some cases, background chemical concentration

levels are irrelevant to the decision-making process. For example, for a chemical release whose constituents are known and not expected to have been released to the environment from any source other than the site, background data would not be necessary. In other cases, levels of background constituents may not exceed risk-based cleanup goals, and, therefore, further background analysis would not be relevant.

2.2 When Background Samples Are Needed

In some cases, the existing data may be inadequate to characterize background. The reasons for this include, but are not limited to, the following:

- ▶ Insufficient number of samples to perform the desired statistical analysis or to perform the tests with the desired level of statistical power;
- ▶ Inappropriate background sample locations (such as those affected by another contamination source, or in soil types that do not reflect onsite soil types of interest);^{3,4}
- ▶ Unknown or suspect data quality;
- ▶ Alterations in the land since the samples were collected (such as by filling, excavation, or introduction of new anthropogenic sources); and
- ▶ Gaps in the available data (certain chemicals were excluded from the sample analyses, or certain soil types were not collected).

2.3 Selecting a Reference Area

A background reference area is the area where

background samples will be collected for comparison with the samples collected on the site. A background reference area should have the same physical, chemical, geological, and biological characteristics as the site being investigated, but has not been affected by activities on the site. RAGS states that "...the locations of the background samples must be areas that could not have received contamination from the site, but that do have the same basic characteristics as the medium of concern at the site."²

The ideal background reference area would have the same distribution of concentrations of the chemicals of concern as those which would be expected on the site if the site had never been impacted. In most situations, this ideal reference area does not exist. If necessary, more than one reference area may be selected if the site exhibits a range of physical, chemical, geological, or biological variability. Background reference areas are normally selected from off-site areas, but are not limited to natural areas undisturbed by human activities. It may be difficult to find a suitable background reference area in an industrial complex. In some cases, a non-impacted onsite area may be suitable as a background reference area.⁵

Complete discussion of the role of geochemical properties of soils in the conduct of background investigations is beyond the scope of this document. In most cases, geochemical methods require more detailed site-specific analysis of the local soil types, biology, and geology than is required for the background comparison methods discussed in this document. The methods in this guidance are based on randomly sampled concentrations of the chemicals of concern.³

CHAPTER NOTES

1. U.S. Environmental Protection Agency (EPA). October 1988. *Guidance for Conducting Remedial Investigations and Feasibility Studies Under CERCLA; Interim Final*. (NTIS PB89-184626, EPA 540-G-89-004, OSWER 9355.3-01).
2. U.S. Environmental Protection Agency (EPA). 1989. *Risk Assessment Guidance for Superfund Vol. I, Human Health Evaluation Manual (Part A)*. Office of Emergency and Remedial Response, Washington, DC. EPA 540-1-89-002. Hereafter referred to as “RAGS.”
3. U.S. Environmental Protection Agency (EPA). December 1995. *Determination of Background Concentrations of Inorganics in Solids and Sediments at Hazardous Waste Sites*. R.P. Breckenridge and A.B. Crockett, National Engineering Forum Issue. Office of Research and Development, Washington, DC. EPA/540/S-96/500.
4. U.S. Environmental Protection Agency (EPA). July 2000. Draft *Ecological Soil Screening Level Guidance*. Office of Emergency and Remedial Response, Washington, DC. EPA 540/F-01/014. OSWER 9345.0-14.
5. Statistical methods based only on sample data collected from both impacted and non-impacted areas on the site are addressed by A. Singh, A.K. Singh, and G. Flatman, “Estimation of background levels of contaminants,” *Mathematical Geology*, Vol. 26, No. 3, 1994.

CHAPTER 3

HYPOTHESIS TESTING AND DATA QUALITY OBJECTIVES

3.1 Hypothesis Testing

The first step in developing a hypothesis test is to transform the problem into statistical terminology by developing a *null hypothesis* and an *alternative hypothesis* (see box on next page). These hypotheses form the two alternative decisions that the hypothesis test will evaluate.

In comparisons with background, the parameter of interest is symbolized by the Greek letter *delta* (Δ), the amount by which the mean of the distribution of concentrations in potentially impacted areas exceeds the mean of the background distribution (see definitions below). Delta is an unknown parameter,

and statistical tests may be used to evaluate hypotheses relating to its possible values. The statistical tests are designed to reject or not reject hypotheses about Δ based on test statistics computed from limited sample data.

The action level for background comparisons is the largest value of the difference in means that is acceptable to the decision maker. In this guidance, the action level for the difference in means is defined as a substantial difference (*S*), which may be zero or a positive value based on the risk assessment, an applicable regulation, a screening level, or guidance. In some cases, the largest acceptable value for the difference in means may be $S = 0$. This

Definitions

Δ (delta): The true difference between the mean concentration of chemical X in potentially impacted areas and the mean background concentration of chemical X. Delta is an unknown parameter which describes the true state of nature. Hypotheses about its value are evaluated using statistical hypothesis tests. In principle, we can select any specific value for Δ and then test if this difference is statistically significant or not with a given confidence and power.

S (substantial difference): A difference in mean concentrations that is sufficiently large to warrant additional interest based on health or ecological information. *S* is the investigation level. If Δ exceeds *S*, the difference in concentrations is judged to be sufficiently large to be of concern, for the purpose of the analysis. A hypothesis test uses measurements from the site and from background to determine if Δ exceeds *S*. The *S* value is discussed further in Appendix A.

MDD (minimum detectable difference): The smallest difference in means that the statistical test can resolve. The MDD depends on sample-to-sample variability, the number of samples, and the power of the statistical test. The MDD is a property of the survey design.

Gray Region: A range of values of Δ where the statistical test will yield inconclusive results. The width of the gray region is equal to the MDD for the test. The location of the gray region depends on the type of statistical test selected.

Null and Alternative Hypotheses

A *statistical hypothesis* is a statement that may be supported or rejected by examining relevant data. Conventionally, hypotheses are stated in such a way that we know what to expect if they are true. However, to determine if we should accept a proposed hypothesis, it is commonly easier to reject its converse (that is, first assume that the hypothesis is *not* true). This assumption to be tested is called the *null hypothesis* (H_0)—if the null hypothesis is rejected, then the initial presumption is *accepted*. A null hypothesis, then, is any testable presumption set up to be rejected. If we want to show that site concentration exceeds background, we formulate a null hypothesis that the site concentration is less than or equal to the background concentration. Similarly, if we want to show that the site concentration is less than or equal to the background concentration, we formulate a null hypothesis that the site concentration exceeds the background concentration.

An *alternative hypothesis* (H_A) is the logical opposite of the null hypothesis: if H_0 is true, H_A is false, and vice-versa. Consequently, the alternative hypothesis is usually logically the same as the investigator's research hypothesis. H_A is the conclusion we accept if we find sufficient evidence to reject the null hypothesis H_0 .

A null hypothesis that specifies the unknown parameter (Δ) as an equality (" $H_0: \Delta = 0$ ") has a corresponding alternative hypothesis that can be higher or lower (" $H_0: \Delta < 0$ " or " $H_0: \Delta > 0$ "). Such a null hypothesis is termed "two tailed" or "two sided" because the alternative hypothesis has two possibilities. A hypothesis test that uses a null hypothesis like " $H_0: \Delta < 0$ " is called "one sided" or "one tailed" because the corresponding alternative hypothesis is true only if the values are greater than zero. One-sided tests are most often used in background comparisons.

guidance does not establish a value for "S"; the value for "S" should be considered on a case-by-case basis. The S value is discussed further in Appendix A. The determination of S should be considered during the development of a Quality Assurance Project Plan as part of the planning process for the background evaluation.¹

Estimates of Δ are obtained by measuring contaminant concentrations in potentially impacted areas and in background areas. For example, one estimate of the mean concentration in potentially impacted areas is the simple arithmetic average of the measurements from these areas. An estimate of the mean background concentration is similarly calculated. An estimate of the difference in means (Δ) is obtained by subtracting the mean background concentration from the mean concentration in potentially impacted areas. In most cases of interest, the estimate of Δ will be a positive number. If there is little or no contamination on the site, then the esti-

mate for Δ may be near zero or slightly negative. Note that the estimated value for Δ calculated by this simple procedure (or by any more complicated procedure) is only an approximation of the true value of Δ . Hence, decisions based on any estimated value for Δ may be incorrect due to uncertainty concerning its true value.

Adopting hypothesis tests and a Data Quality Objective (DQO) approach (Section 3.4) can control the probability of making decision errors. However, incorrect use of hypothesis tests can lead to erratic decisions. Each type of hypothesis test is based on a set of assumptions that should be verified to confirm proper use of the test. Procedures for verifying the selection and proper use of parametric tests, such as the t-tests, are provided in EPA QA/G-9, Chapter 4.² Nonparametric tests generally have fewer assumptions to verify.

Hypothesis testing is a quantitative method to

determine whether a specific statement concerning Δ (called the null hypothesis) can be rejected. Decisions concerning the true value of Δ reduce to a choice between “yes” or “no.” When viewed in this way, two types of incorrect decisions, or decision errors, may occur:

- ▶ Incorrectly deciding the answer is “yes” when the true answer is “no;” and
- ▶ Incorrectly deciding the answer is “no” when the true answer is “yes.”

While the possibility of decision errors can never be totally eliminated, it can be controlled. To control decision errors, it is necessary to control the uncertainty in the estimate of Δ . Uncertainty arises from three sources:

- ▶ Sampling error;
- ▶ Measurement error; and
- ▶ Natural variability.

The decision maker has some control of the first two sources of uncertainty. For example, a larger number of samples may lead to fewer decision errors because the probability of a decision error decreases as the number of samples increases. Use of more precise measurement techniques or duplicate measurements can reduce measurement error, thus minimizing the likelihood of a decision error. The third source of uncertainty is more difficult to control.

Natural variability arises from the uneven distribution of chemical concentrations on the site and in background areas. Natural variability is measured by the true standard deviation (σ) of the distribution. A large value of σ indicates that a large number of measurements will be needed to achieve a desired limit on decision errors. Since variability is usually higher in impacted areas of the site than in background locations, data collected on the site is used to estimate σ . An estimate for σ frequently is obtained from historical data, if available. Estimates of variability reported elsewhere at similar sites with similar contamination problems may be used. If an estimate of the mean concentration in contaminated areas is available, the coefficient of variation obser-

ved at other sites may be multiplied by the mean to estimate the standard deviation. If no acceptable historical source for an estimate of σ is available, it may be necessary to conduct a small-scale pilot survey on site using 20 or more random samples to estimate σ . Due to the small sample size of the pilot, it is advisable to use an 80 or 90 percent upper confidence limit for the estimate of σ rather than an unbiased estimate to avoid underestimating the true variability. A very crude approximation for σ may be made by dividing the anticipated range (maximum - minimum) by 6. It is important that overly optimistic estimates for σ be avoided because this may result in a design that fails to generate data with sufficient power for the decision.

The hypothesis testing process provides a formal procedure to quantify the decision maker’s acceptable limits for decision errors. The decision maker’s limits on decision errors are used to establish performance goals for data collection that reduce the chance of making decision errors of both types. The gray region is a range of possible values of Δ where the consequences of making a decision error are relatively minor. Examples of the gray region are shown in Figures 3.1 and 3.2 (Section 3.3).

Any useful statistical test has a low probability of reflecting a substantial difference when the site and background distributions are identical (false positive) but has a high probability of reflecting a substantial difference when the distribution of contamination in potentially impacted areas greatly exceeds the background distribution. In the gray region between these two extremes, the statistical test has relatively poor performance. When the test procedure is applied to a site with a true mean concentration in the gray region, the test may indicate that the site exceeds background, or may indicate that the site does not exceed background, depending on random fluctuations in the sample.

It is necessary to specify a gray region for the test because the decision may be “too close to call” due to uncertainty in the estimate of Δ . This may occur when the difference in means is small compared to the MDD for the test. In the gray region, the uncer-

tainty in the measurement of Δ is larger than the difference between Δ and the action level, so it may not be possible for the test to yield a correct decision with a high probability. One step in the hypothesis test procedure is to assign upper bounds on the decision error rates for values of Δ above and below the gray region. These bounds limit the probability of occurrence of decision errors.

The exact definition of the gray region is determined by the type of hypothesis test that is selected by the decision maker (see Figures 3.1 and 3.2 in Section 3.3). In general, the gray region for Δ is to the right of the origin ($\Delta = 0$) and bounded from above by the substantial difference ($\Delta = S$). Additional guidance on specifying a gray region for the test is available in Chapter 6 of *Guidance for the Data Quality Objectives Process*.³ The size of the gray region may also depend on specific regulatory requirements or policy decisions that may not be addressed in the DQO guidance.

The width of the gray region is called the “minimum detectable difference” for the statistical test, indicating that differences smaller than the MDD cannot be detected reliably by the test. If the test is used to determine if concentrations in the potentially impacted areas exceed background concentrations by more than S , it is necessary to ensure that MDD for the test is less than S . In the planning stage, this requirement is met by designing a sampling plan with sufficient power to detect differences as small as S . If data were collected without the benefit of a sampling plan, retrospective calculation of the power of the test may be necessary before making a decision.

In the planning stage, the absolute size of the MDD is of less importance than the ratio of the MDD to the natural variability of the contaminant concentrations in the potentially impacted area. This ratio is termed the “relative difference” and defined as MDD/σ , where σ is an estimate of the standard deviation of the distribution of concentrations on the site. The relative difference expresses the power of resolution of the statistical test in units of uncertainty. Relative differences much less than one standard deviation ($MDD/\sigma \ll 1$) are more difficult to

resolve unless a larger number of measurements are available. Relative differences of more than three standard deviations ($MDD/\sigma > 3$) are easier to resolve. As a general rule, values of MDD/σ near 1 will result in acceptable sample sizes. The required number of samples may increase dramatically when MDD/σ is much smaller than one. Conversely, designs with MDD/σ larger than three may be inefficient. If MDD/σ is greater than three, additional measurement precision is available at minimal cost by reducing the width of the gray region. The cost of the data collection plan should be examined quantitatively for a range of possible values of the MDD before selecting a final value. A tradeoff exists between cost (number of samples required) and benefit (better power of resolution of the test).

σ (mg/kg)	MDD/ σ	n	N
25	2	3.70	5
50	1	13.55	16
75	0.67	29.97	35
100	0.50	52.97	62
125	0.40	82.53	96
150	0.33	118.66	138
175	0.29	161.36	188
200	0.25	210.63	245

Table 3.1 Required sample size for selected values of σ ($\alpha = \beta = 0.10$ and MDD = 50 mg/kg)

The number of measurements required to achieve the specified decision error rates has a strong inverse relationship with the value of MDD/σ . An example of this inverse relationship is demonstrated in Table 3.1 for hypothetical values of $\alpha = \beta = 0.10$ and MDD = 50 mg/kg. Sample sizes may be obtained using the approximate formula given in EPA QA/G-9² (Section 3.3.3.1, Box 3-22, Step 5 of that document), written here as:

$$n = (0.25) z_{1-\alpha}^2 + 2 (z_{1-\alpha} + z_{1-\beta})^2 \sigma^2 / (MDD)^2,$$

where z_p is the p^{th} percentile of the standard normal distribution. Note the inverse-squared dependence

of n on MDD/σ . Smaller values of α and β (leading to larger values for the z terms) magnify the strength of this inverse relationship. A recommended sample size of $N = (1.16)n$ is tabulated for a variety of σ values in the table. Note the dramatic increase in the sample size as the value of MDD/σ is lowered from 1 to 0.25.

Letting $\alpha = \beta$, we can solve for $z_{1-\alpha} = z_{1-\beta}$:

$$z_{1-\alpha}^2 = n / [0.25 + 8\sigma^2 / (MDD)^2].$$

For any fixed value of MDD/σ , the decision error rate α is a function of n :

$$\alpha = 1 - \Phi[z_{1-\alpha}(n)],$$

where Φ is the cumulative normal distribution function. Achievable values of α (and β) for selected sample sizes with a hypothetical value of $MDD/\sigma = 1/2$ are shown in Table 3.2.

N	n	$Z_{1-\alpha}$	$\alpha = \beta$
10	8.62	0.517	0.303
15	12.93	0.633	0.263
20	17.24	0.731	0.232
25	21.55	0.817	0.207
30	25.86	0.896	0.185
40	34.48	1.034	0.151
50	43.10	1.156	0.124
60	51.72	1.266	0.103
70	60.34	1.368	0.086
100	86.21	1.635	0.051
150	129.31	2.002	0.023
200	172.41	2.312	0.010

Table 3.2 Achievable values of $\alpha = \beta$ for selected values of N with $MDD/\sigma = 1/2$

A tradeoff analysis should begin with analysis of the choice $MDD = S$, where S is a substantial difference. Note that a choice of $MDD > S$ would lead to a sample size that does not have sufficient power to distinguish a difference between the site and background means as small as S . Hence the minimum acceptable number of samples for the decision is obtained when $MDD = S$. If S/σ is less

than one, this indicates that MDD/σ is also less than one, and a relatively large number of samples will be required to make the decision. If S/σ exceeds three, then a reasonably small number of samples are required for this minimally acceptable test design. Additional measurement precision is available at minimal cost by choosing $MDD < S$. A binary search procedure would indicate the choice of $MDD = S/2$ as the next trial in the cost tradeoff comparison. If S/σ is between one and three, then selecting $MDD = S$ is a reasonable alternative. If $S/\sigma < 1$, then selecting $MDD = S$ is the most cost-effective choice consistent with the requirement that $MDD \leq S$.

The MDD , in conjunction with the values selected for the decision error rates, determines the cost of the survey design and the success of the survey in determining which areas present unacceptable risks. From a risk assessment perspective, selection of the proper width of the gray region is one of the most difficult tasks. The goal is to make the MDD as small as possible within the goals and resources of the cleanup effort.

Two forms of the statistical hypothesis test are useful for comparisons with background. The null hypothesis in the first form of the test states that there is *no statistically significant difference* between the means of the concentration distributions measured at the site and in the selected background areas. The null hypothesis in the second form of the test is that the impacted area of the site *exceeds background by a substantial difference*. RAGS⁴ provides guidance for the first form of the background hypothesis test. Both forms are described in the next section.

3.1.1 Background Test Form 1

The null hypothesis for background comparisons, “the concentration in potentially impacted areas does not exceed background concentration,” is formulated for the express purpose of being rejected:

- ▶ *The null hypothesis (H_0).* The mean contaminant

concentration in samples from potentially impacted areas is less than or equal to the mean concentration in background areas ($\Delta \leq 0$).⁵

- ▶ *The alternative hypothesis (H_a)*. The mean contaminant concentration in samples from potentially impacted areas is greater than the mean in background areas ($\Delta > 0$).

When using this form of hypothesis test, the data should provide statistically significant evidence that the null hypothesis is false—the site does exceed background. Otherwise, the null hypothesis cannot be rejected based on the available data, and the concentrations found in the potentially impacted areas are considered equivalent to background.

An easy way to think about the decision errors that may occur using Background Test Form 1 is to think about the criminal justice system in this country and consider what a jury must weigh to determine guilt. The only choices are “guilty” and “not guilty.” A person on trial is presumed “innocent until proven guilty.” When the evidence (data) is clearly not consistent with the presumption of innocence, a jury reaches a “guilty” verdict. Otherwise the verdict of “not guilty” is rendered when the evidence is not sufficient to reject the presumption of innocence. A jury does not have to be convinced that the defendant is innocent to reach a verdict of “not guilty.” Similarly, when using Background Test Form 1, the null hypothesis is presumed true until it is rejected.

Two serious problems arise when using Background Tests Form 1. One type of problem arises when there is a very large amount of data. In this case, the MDD for the test will be very small, and the test may reject the null hypothesis when there is only a very small difference between the site and background mean concentrations. If the site exceeds background by only a small amount, there is a very high probability that the null hypothesis will be rejected if a sufficiently large number of samples is taken. This case can be avoided by selecting Background Test Form 2, which incorporates an acceptable level for the difference between site and background concentrations.

A second type of problem may arise in the use of Background Test Form 1 when insufficient data are available. This may occur, for example, when the onsite or background variability was underestimated in the design phase. An estimated value for σ is used during the preliminary phase of the DQO planning process to determine the required number of samples. When the samples are actually collected, σ can then be re-estimated, and the power of the analysis should be re-evaluated. If the variance estimate used in the planning stage was too low, the statistical test is unlikely to reject the null hypothesis due to the lack of sufficient power. Hence, when using Background Test Form 1, it is always best to conduct a *retrospective power analysis* to ensure that the power of the test was adequate to detect a site with mean contamination that exceeds background by more than the MDD. A simple way to do this is to recompute the required sample size using the sample variance in place of the estimated variance that was used to determine the required sample size in the planning phase. If the actual sample size is greater than this post-calculated size, then it is likely that the test has adequate power. The exact power of the WRS test used for Background Test Form 1 is difficult to calculate. See Section 5.3.2 for more information on the power of the WRS test. If the retrospective analysis indicates that adequate power was not obtained, it may be necessary to collect more samples. Hence, if large uncertainties exist concerning the variability of the contaminant concentration in potentially impacted areas, Background Test Form 1 may lead to inconclusive results. Therefore, the sample size should exceed the minimum number of samples required to give the test sufficient power.

Detailed information on the application and characteristics of Background Test Form 1 is available in the document series *Statistical Methods for Evaluating the Attainment of Cleanup Standards*. Volume 3, subtitled *Reference-Based Standards for Soils and Solid Media*⁶ contains detailed procedures for comparing site measurements with background reference area data using parametric and nonparametric tests based on Background Test Form 1.

Interpretation of the Statistical Measures

Background Test Form 1

Confidence level = 80%: On average, in 80 out of 100 cases, chemical concentrations in potentially contaminated areas will be correctly identified as being no different (statistically) from background concentrations, while in 20 out of 100 cases, concentrations in potentially contaminated areas will be incorrectly identified as being greater than background concentrations.

Power = 90%: On average, in 90 out of 100 cases, concentrations in potentially contaminated areas will be correctly identified as being greater than background concentrations, while in 10 out of 100 cases, concentrations in potentially contaminated areas will be incorrectly identified as being less than or equal to background concentrations.

Background Test Form 2

Confidence level = 90%: On average, in 90 out of 100 cases, concentrations in potentially contaminated areas will be correctly identified as exceeding background concentrations by more than S, while in 10 out of 100 cases, concentrations in potentially contaminated areas will be incorrectly identified as not exceeding background concentrations by more than S.

Power = 80%: On average, in 80 out of 100 cases, concentrations in potentially contaminated areas will be correctly identified as not exceeding background concentrations by more than S, while in 20 out of 100 cases, concentrations in potentially contaminated areas will be incorrectly identified as exceeding background concentrations by more than S.

3.1.2 Background Test Form 2

An alternative form of hypotheses test for comparing two distributions is presented in *Guidance for the Data Quality Objectives Process, EPA QA/G-4*.³ When adapted to the background problem, the null hypothesis, “the concentration in potentially impacted areas exceeds background concentration,” again is formulated for the express purpose of being rejected:

- ▶ *The null hypothesis* (H_0): The mean contaminant concentration in potentially impacted areas exceeds background by more than S. Symbolically, the null hypothesis is written as $H_0: \Delta > S$.⁷
- ▶ *The alternative hypothesis* (H_A): The mean contaminant concentration in potentially impacted areas does not exceed background by more than S ($H_A: \Delta \leq S$).

Here, S is the background investigation level. Although there is no explicit use of the quantity S in the hypothesis statement used in Background Test Form 1, an estimate of S is important for determining an upper limit for the MDD for Background Test Form 1, as discussed below. Issues affecting the determination of site-specific values for S are not the subject of this guidance. The background investigation level is determined on a case-by-case basis by EPA and other stakeholders. Several approaches for determining a background investigation level are discussed in more detail in Appendix A.2.

Detailed information on the application and characteristics of parametric statistical tests based on Background Test Form 2 is available in Volumes 1 and 2 of the EPA document series *Statistical Methods for Evaluating the Attainment of Cleanup Standards*.⁶

3.1.3 Selecting a Background Test Form

When comparing Background Test Forms 1 and 2, it is important to distinguish between the selection of the null hypothesis, which is a burden-of-proof issue, and the selection of the investigation level, which involves determination of an action level.

Background Test Form 1 uses a conservative investigation level of $\Delta = 0$, but relaxes the burden of proof by selecting the null hypothesis that the contaminant concentration in potentially impacted areas is not statistically different from background. Background Test Form 2 requires a stricter burden of proof, but relaxes the investigation level from O to S. Section 5.4 includes further discussion of how to choose between Test Forms 1 and 2, and gives additional guidance for setting up the hypotheses. See the box on the previous page about *Interpretation of the Statistical Measures*.

Regardless of the choice of hypothesis, an incorrect conclusion could be drawn from the data analysis using either form of the test. To account for this inherent uncertainty, one should specify the limits on the Type I and Type II decision errors. This task is addressed in Step 6 of the DQO process and described in Section 3.4.

3.2 Errors Tests and Confidence Levels

A key step in developing a sampling and analysis plan is to establish the level of precision required of the data.³ Whether the null hypothesis (Section 3.1) will be rejected or not depends on the results of the

sampling. Due to the uncertainties that result from sampling variation, decisions made using hypothesis tests will be subject to errors. Decisions should be made about the width of the gray region and degree of decision error that is acceptable. These topics are discussed below and in more detail in Chapter 5. There are two ways to err when analyzing data (Table 3.3):

- ▶ *Type I Error*: Based on the observed data, the test may reject the null hypothesis when in fact the null hypothesis is true (a false positive). This is a *Type I error*. The probability of making a Type I error is α (*alpha*); and
- ▶ *Type II Error*: On the other hand, the test may fail to reject the null hypothesis when the null hypothesis is in fact false (a false negative). This is a *Type II error*. The probability of making a Type II error is β (*beta*).

The *acceptable level of decision error* associated with hypothesis testing is defined by two key parameters—*confidence level* and *power* (see the box at the bottom of the previous page). These parameters are closely related to the two error probabilities, α and β .

- ▶ *Confidence level* $100(1 - \alpha)\%$: As the confidence level is lowered (or alternatively, as α is increased), the likelihood of committing a Type I error increases.
- ▶ *Power* $100(1 - \beta)\%$: As the power is lowered (or alternatively, as β is increased), the likelihood of committing a Type II error increases.

Decision Based on Sample Data	Actual Site Condition	
	H_0 is True	H_0 is not True
H_0 is not rejected	Correct Decision: $(1 - \alpha)$	Type II Error: False Negative (β)
H_0 is rejected	Type I Error: False Positive (α)	Correct Decision: $(1 - \beta)$

Table 3.3 Hypothesis Testing: Type I and Type II Errors

Although a range of values can be selected for these two parameters, as the demand for precision increases, the number of samples and the cost will generally also increase. The cost of sampling is often an important determining factor in selecting the acceptable level of decision errors. However, unwarranted cost reduction at the sampling stage may incur greater costs later. The number of samples, and hence the cost of sampling, can be reduced but at the expense of a higher possibility of making decision errors that may result in the need for additional sampling, unnecessary remediation, or increased risk. The selection of appropriate levels for decision errors and the resulting number of samples is a critical component of the DQO process that should concern all stakeholders.

Because there is an inherent tradeoff between the probability of committing a Type I or Type II error, a simultaneous reduction in both types can only occur by increasing the number of samples. If the probability of committing a false positive is reduced by increasing the level of confidence of the test (in other words, by decreasing α), the probability of committing a false negative is increased because the power of the test is reduced (increasing β).

For the purposes of this guidance, minimum recommended performance measures are:⁸

- ▶ For Background Test Form 1, confidence level at least 80% ($\alpha = 0.20$) and power at least 90% ($\beta = 0.10$).
- ▶ For Background Test Form 2, confidence level at least 90% ($\alpha = 0.10$) and power at least 80% ($\beta = 0.20$).

When using Background Test Form 1, a Type I error (false positive) is less serious than a Type II error (false negative). *This approach favors the protection of human health and the environment.* To ensure that there is a low probability of Type II errors, a Test Form 1 statistical test should have adequate power at the right edge of the gray region.

When Background Test Form 2 is used, a Type II

error is preferable to committing a Type I error. *This approach favors the protection of human health and the environment.* The choice of hypotheses used in Background Test Form 2 is designed to be protective of human health and the environment by requiring that the data contain evidence of *no substantial contamination*. This approach may be contrasted to the “innocent until proven guilty” approach used in Background Test Form 1.

3.3 Test Performance Plots

During the scoping stage for the development of the sampling plan, the interrelationships among the decision parameters can be visualized using a *test performance plot*. The test performance plot is a graph that displays the combined effects of the decision error rates, the gray region for the decision-making process, and the level of a substantial difference between site and background. In short, it displays most of the important parameters developed in the DQO process.

A test performance plot is used in the planning stages of the DQO process to aid in the selection of reasonable values for the decision error rates (α and β), the MDD, and the required number of samples. Selection of these parameters is usually an iterative process. Trial values of the decision error rates, the location of the gray region, and its width (the MDD) are used to generate initial estimates of the required number of samples and the resulting test performance curve. Adjustments to the inputs are made until a design is achieved that offers acceptable test performance at an acceptable cost.

Figure 3.1 illustrates an example of a test performance plot for decision making on a statistical test based on the null hypothesis that the mean concentration in the potentially impacted area does not exceed mean background concentration (Background Test Form 1). At the origin of the plot, the true difference between the means of the site and background distributions is zero ($\Delta=0$). Positive values of the difference between the site and background mean concentrations ($\Delta > 0$) are plotted on

the horizontal axis to the right of the origin, negative values ($\Delta < 0$) to the left. The vertical axis shows the value of the test performance measure, defined as the power of the test. The power of the test is the probability of rejecting the null hypothesis which, for this test form, equals the probability of deciding the mean concentration in potentially impacted areas exceeds the mean background concentration. This probability ranges from 0 to 1.0 (0 to 100 percent).

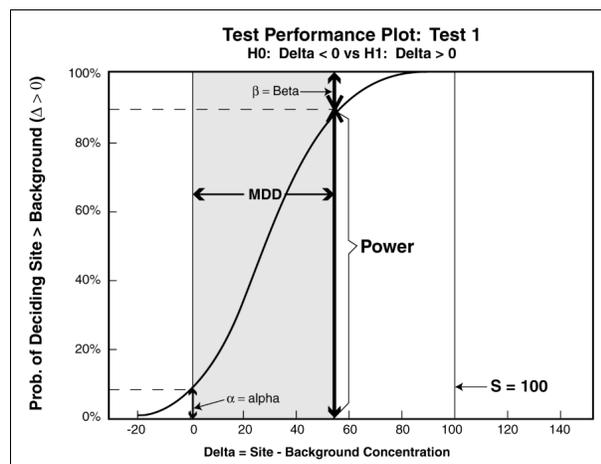


Figure 3.1 Test performance plot: site is not significantly different from background

At the left edge of the gray region, the test performance curve is no greater than α for potentially impacted areas with mean contaminant concentrations less than or equal to background mean concentration ($\Delta \leq 0$) and greater than α for potentially impacted areas with mean concentration exceeding the mean background concentration ($\Delta > 0$). The test performance curve increases as the difference between the potentially impacted area and background means increases. The number of samples and the standard deviation, σ , determine the rate of increase. The right edge of the gray region is located at the MDD ($\Delta = \text{MDD}$). At this value of the difference between the mean potentially impacted area and background concentrations, the probability of deciding that the potentially impacted area exceeds background is equal to $1 - \beta$. When using Background Test Form 1, the test performance curve equals the power of the test. A statistical software package for plotting the power of a statistical test may be used to generate a test performance plot.

EPA has developed two software packages that generate power curves for the two-sample t-test: DEFT⁹ and DataQUEST.¹⁰

Figure 3.1 also shows a hypothetical value of a substantial difference for this chemical of $S = 100$. The value of S was developed by conducting an evaluation of the risks presented by the site. The value of S is used in the DQO process as an upper limit for the width of the gray region (MDD). In some cases, an MDD less than S may be selected for the test. This is determined by site-specific conditions, summarized by the standard deviation, σ . If the ratio S/σ exceeds 3, then a sample design with an MDD less than S may offer a test with better power of resolution at little additional cost of sampling, a strategy often described using the term “ALARA”—“As Low As Reasonably Achievable.” If the MDD is selected to be smaller than S , then the design is conservative in the sense that potentially impacted areas with differences from background smaller than S can be identified by the test. The test will have a higher power to reject the null hypothesis for sites with mean concentrations that are in the range between the MDD and S higher than background. In statistical terms, the power of rejection will be $(1 - \beta)$ at $\Delta = \text{MDD}$, and higher than $(1 - \beta)$ for all $\Delta > \text{MDD}$.

Selecting an MDD less than S is also useful for screening a large number of areas using a low cost sample measurement procedure, with subsequent confirmatory testing using more expensive procedures before making a final decision. Finally, before using previously collected data for decision making, the power of the test should be calculated to determine if the MDD is less than S .

An equivalent plot in Figure 3.2 shows the test performance curve for a statistical test using the null hypothesis that *the potentially impacted area does not exceed background by more than a substantial difference* (Background Test Form 2). For this Test Form, the MDD again measures the width of the gray region, but the gray region now extends from a difference of $\Delta = S - \text{MDD}$ on the left to a difference $\Delta = S$ on the right.

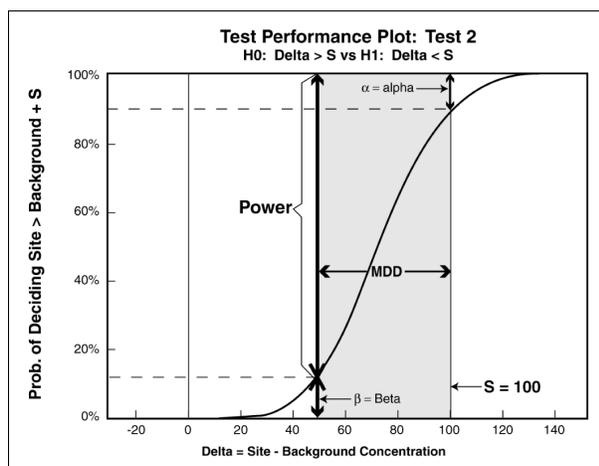


Figure 3.2 Test performance plot: site does not exceed background by more than S

When using Background Test Form 2, the MDD may be selected to be as large as S or smaller. The implications of making the MDD smaller than S for this Test Form differ from those that occur when using Background Test Form 1. As the MDD decreases below S, the test will identify more potentially impacted areas as not having mean concentrations that exceed background by more than S. The sites with mean concentration in the range between $\Delta = 0$ and $\Delta = S - \text{MDD}$ (those with mean concentrations only slightly higher than background) will have a higher probability of being classified correctly. With this Test Form, a tradeoff exists between taking more samples and making more errors. Since the errors tend to occur in sites that are marginally acceptable, it would be beneficial for responsible parties to increase the number of samples and the power of the test.

This second form of background test requires switching the location of α and β . The Type I error (α) for Background Test Form 2 is measured by the difference between 100% and the test performance curve at the right of the gray region, while the Type II error (β) is measured by the value of the test performance curve at a difference equal to $\Delta = S - \text{MDD}$, located at the left of the gray region. When using Background Test Form 2, the test performance curve equals 100% minus the power of the test.

When using Background Test Form 1, a Type I error could lead to unnecessary remediation while a Type II error could lead to unacceptable health risks. If Background Test Form 2 is used, a Type II error could lead to unnecessary remediation while a Type I error could lead to unacceptable health risks. Therefore, one should attempt to reduce the chance of making either of these errors.

Comparison of Figures 3.1 and 3.2 demonstrates that the choice $\alpha_2 = \beta_1$, $\beta_2 = \alpha_1$, and $\text{MDD} = S$ will result in almost identical test performance plots for Background Test Form 1 and Background Test Form 2. If MDD is less than S, then Background Test Form 1 will indicate that more potentially impacted areas require remediation than Background Test Form 2. In general, α will differ from β , and the value selected for the MDD may be smaller than S.

The selection of acceptable decision error rates for hypothesis testing is a decision that should be made on a site-specific basis. The consequences of making a wrong decision (such as failing to reject the null hypothesis when it is false) should be considered when specifying acceptable values for the confidence and power factors ($\alpha = 0.20$ and $\beta = 0.10$ are maximum values for Background Test Form 1).

3.4 DQO Steps for Characterizing Background

DQOs should be used when developing sampling and analysis plans (SAPs) to ensure that reliable data are acquired. The process is outlined here with a case example for purposes of developing background sampling plans. For further details, consult Section 6 of *Guidance for the Data Quality Objectives Process*³ and *Guidance for Data Quality Assessment: Practical Methods for Data Analysis*.²

The DQO process is the starting point for many decisions that shape the sampling plan. It involves a series of steps for making optimal decisions based on limited data. A careful statement of the DQOs for

a study will clarify the study objectives, define the most appropriate type of data to collect, determine the most appropriate conditions for collecting data, and specify limits on decision errors. Use of the DQO process ensures that the type, quantity, and quality of environmental data used in decision making will be appropriate for the intended application. It improves efficiency by eliminating unnecessary, duplicative, or overly precise data. The DQO process provides a systematic process for defining an acceptable level for decision errors. The DQO process and decision parameters establish the quantity and quality of data needed. A sampling design is developed to implement these requirements by defining the specific measurement protocol, sample locations, and number of samples that will be collected. Detailed procedures for developing the sampling design are presented in EPA QA/G-5S¹¹. Many new sampling approaches are discussed in this document, including ranked set sampling and adaptive cluster sampling.

Each of the seven steps of the DQO process, listed below, may be phrased as a question about background issues:

1. State the Problem
2. Identify the Decision
3. Identify Inputs to the Decision
4. Define Boundaries of Study
5. Develop a Decision Rule
6. Specify Limits on Decision Errors
7. Optimize the Design for Obtaining Data

The examples provided in this section should be modified to fit the site of concern. A statistician familiar with the challenges posed by environmental data should be consulted *before* data are collected. The statistician should be involved in discussions about the goals of the background analysis, time and cost constraints, limitations of the measurement techniques, and the availability of preliminary data.

Step 1. State the Problem: *Example: Are there differences between the concentrations of a site contaminant and those concentrations that are found in background samples?*

Tasks include:

- ▶ Identifying the resources available to resolve the problem. The team should include the decision makers, technical staff and data users, and stakeholders. Members of the technical staff may include quality assurance managers, chemists, modelers, soil scientists, engineers, geologists, health physicists, risk assessors, field personnel, and regulators.
- ▶ Developing or refining the comprehensive conceptual site model.

Step 2. Identify the Decision: *Example: Are the chemicals associated with a site-related source or are they associated with background?*

Tasks include:

- ▶ Identifying the chemicals to analyze; and
- ▶ Determining if these chemicals are expected to occur in reference areas selected to reflect background conditions.

Step 3. Identify Inputs into the Decision: *Example: What kinds of data are needed? What kinds of data are available?*

Tasks include identifying:

- ▶ Which chemicals need to be analyzed;
- ▶ Which soil types and depths need to be sampled;
- ▶ Which comparison tests are likely to be used (see Chapter 5 for details about comparison tests);
- ▶ What coefficient of variation is expected for the data (based on previous samples if possible);
- ▶ What preliminary remediation goals (PRGs) or applicable or relevant and appropriate requirements (ARARs) should be considered; and
- ▶ What are the desired power and confidence levels?

Decision outputs for background characterizations are discussed in detail in Chapter 5.

Step 4. Define Boundaries of the Study: *Example: What are the spatial and temporal aspects of the environmental media that the data should represent to support the decision?*

Tasks include:

- ▶ Defining the geographic areas for field investigation;
- ▶ Defining the characteristics of the soil data population of interest;
- ▶ Dividing the soil data population of interest into strata having relatively homogeneous characteristics;
- ▶ Determining the timeframe to which the decision applies; and
- ▶ Identifying practical constraints that may hinder sample collection.

Step 5. Develop a Decision Rule: *Example: If the mean concentration in potentially impacted areas exceeds the mean background concentration, then the chemical will be treated as site-related. Otherwise, if the mean concentration in potentially impacted areas does not exceed the background mean, the chemical will be treated as background-related.*

Tasks include:

- ▶ Choosing the null hypothesis, H_0 ;
- ▶ Specifying the alternative hypothesis, H_A ;
- ▶ Specifying the gray region for the hypothesis test; and
- ▶ Determining the level of a substantial difference above background, S .

Hypothesis testing is an approach that helps the decision maker through the analysis of data. Chapter 5 discusses the application of hypothesis testing at CERCLA sites. General information on hypothesis testing is provided in Section 3.1.

Step 6. Specify the Limits on Decision Errors: *Example: What level of uncertainty is acceptable for this decision? (For definitions, see Section 3.1 on*

Hypothesis Testing, Section 3.2 on Errors, and Confidence Levels, and Figures 3.1 and 3.2.):

- ▶ Test Form 1—The gray region extends from a difference of $\Delta = 0$ on the left to $\Delta = \text{MDD}$ on the right. Acceptable limits on decision errors are α_1 at the left edge of the gray region, and β_1 at the right edge. Here, α_1 measures the Type I error rate for Test Form 1, which is the probability of rejecting the null hypothesis when it is true, i.e., the probability of wrongly concluding that the mean concentration on site exceeds the background mean when it does not. β_1 measures the Type II error rate for Test Form 1, which is the probability of not rejecting the null hypothesis when it is false, i.e., wrongly concluding the mean concentration on the site does not exceed background when it does.
- ▶ Test Form 2—The gray region extends from a difference of $\Delta = (S - \text{MDD})$ on the left to $\Delta = S$ on the right. The acceptable limits on decision errors are α_2 at the right edge of the gray region, and β_2 at the left edge. Here, α_2 measures the Type I error rate for Test Form 2, which is the probability of rejecting the null hypothesis when it is true. For this test, the Type I error rate is the probability of concluding (wrongly) that the mean concentration on the site does not exceed the background mean by more than S when it does. Similarly, β_2 measures the Type II error rate for Test Form 2, which is the probability of not rejecting the null hypothesis when it is false. In this case, the Type II error rate is the probability of concluding (wrongly) that the mean concentration on site exceeds the background mean by more than S when it does not.

Tasks include:

- ▶ Determining the possible range for Δ ;
- ▶ Specifying both types of decision errors (Type I and Type II—see Section 3.2);¹²
- ▶ Identifying the potential consequences of each type of error, specifying a range of possible values for Δ (the gray region—see Figures 3.1 and 3.2) where consequences of decision errors

are relatively minor; and

- ▶ Selecting the limits on decision errors (α and β) to reflect the decision-maker's concern about the relative consequences for each type of decision error (Section 3.2).

Step 7. Optimize the Sampling Design: *Example: What is the most resource-effective sampling and analysis design for generating data that are expected to satisfy the DQOs?*

Tasks include:

- ▶ Reviewing the DQO outputs and existing environmental data;
- ▶ Developing general sampling and analysis design alternatives;
- ▶ Verifying that DQOs are satisfied for each design alternative;
- ▶ Selecting the most resource-effective design that satisfies all of the DQOs; and
- ▶ Documenting the operational details and theoretical assumptions of the selected design in the sampling and analysis plan.

More information may be required to make a decision. If the required sample size is too large, it may be necessary to modify the original DQO parameters. To reduce sampling cost while maximizing utility of the available resources, one or more of the constraints used to develop the sampling design may be relaxed. Gilbert presents useful information on how to factor cost into a sampling design.¹³

3.5 Sample Size

The RPM should consult with a statistician who has experience in designing environmental sampling programs to select the appropriate sampling design. Several sampling design options are available. See EPA QA/G5S¹¹ for guidance on sampling design. A consistent grid to cover the entire site and areas considered as background should provide a reasonable characterization of the concentrations onsite and in background areas. The ideal data sets should

be independent (spatially uncorrelated), unbiased, and representative of the underlying site and background populations. These assumptions favor widespread random samples. However, in many instances, the background analyses should rely on existing site data collected using judgmental sampling. Such data sets are often biased, clustered, and correlated. In certain cases, the existing clustered data set may be declustered for background analyses. A variety of de-clustering alternatives exist. For example, the investigated area can be divided into equally spaced grids. Each grid can then be represented by average concentration of measured values within the grid, or a predefined number of samples can be selected randomly from each grid. Additional options are described in other guidance, including Chapter 4 of RAGS.⁴

In most DQO applications, after electing to use a test with confidence level $100(1 - \alpha)$ percent, the required number of samples is determined by simultaneously selecting:

- ▶ the MDD for the test; and
- ▶ the power $(1 - \beta)$ of the test at the MDD.

Therefore, limits on the probability of committing Type I and Type II errors can be used as constraints on the number and location of samples. The DQO process is meant to be an iterative process. If the number of samples determined with the selected error probabilities is too large for the available resources, the DQO procedure should be repeated with more reasonable error objectives until an acceptable number of samples is determined. To determine realistic limits for the decision errors, the number of samples (and the corresponding cost of sampling) could be estimated for a range of error probability values, which would indicate the likelihood of making either type of error. Reports of the results of the DQO process should specify the number of samples selected and the expected error probabilities that result from this selection.

Several reference documents give formulas or tables for selecting the number of samples, given the specific confidence and power limits.¹⁴ Chapter 5

offers guidance for selecting appropriate statistical techniques for comparing onsite and background contaminant concentrations in soil.

Examples of constraints that may be adjusted to influence the required sample size include:

- ▶ Increasing the decision-error rates, α and β , while considering the increased costs and risks associated with the increased probability of making an incorrect decision;
- ▶ Increasing the width of the gray region (MDD), but not to exceed a substantial difference ($MDD \leq S$); and
- ▶ Changing the boundaries. It may be possible to reduce measurement costs by segregating the site into subunits that require different decision parameters due to different risks.

3.6 An Example of the DQO Process

This section presents a hypothetical application of the DQO process for comparing lead concentrations in a potentially impacted area to background. The conceptual site model and remedial goals for individual sites will determine what sampling and analysis is done at any site. The example will illustrate some outputs of the DQO process and will be extended to the preliminary data analysis stage in Chapter 4 and to the hypothesis testing stage in Chapter 5. RPMs should consult the Technical Review Workgroup for Lead (TRW) for technical assistance with lead-contaminated sites.¹⁵ This example only illustrates the DQO process and does not establish guidance pertaining to subsurface soil sampling or lead cleanup goals.

Step 1. State the Problem

An abandoned storage yard has been identified as the previous location of a battery distributorship. Concerns have focused on this storage area as a possible source of lead contamination. Other sources of background lead are present in the vicinity of the

storage yard due to nearby highways and industrial facilities. Available data are not sufficient to determine that the concentrations in the potentially impacted area are different from background chemical concentrations. The study team has decided to conduct field measurements.

- a. *What resources (including necessary personnel) are available to resolve the problem?*

The members of the study team will include the plant manager, a plant engineer, a chemist with field sampling experience, a quality assurance officer, a statistician, a risk assessor, and the remedial project manager.

- b. *What characteristics or data will determine the comprehensive conceptual site model?*

Historical site assessment was used to develop a comprehensive conceptual site model. Due to nearby highways and industrial sources in the vicinity of the yard, background lead concentrations in soil are expected to be above the national average. Also, because of run-off from paved areas, background concentration near the paved areas are likely to be higher than background concentrations in soils distant from the paved areas. The selection of appropriate background areas for the comparison was restricted to areas at least 1,000 meters from heavily used highways and 30 meters from paved surfaces. These requirements were selected to match the relative location of the site with respect to the surrounding roads and highways.

Step 2. Identify the Decision

Do soils in the storage area have higher lead concentrations than found in soils in the surrounding area?

- a. *What chemical(s) should be analyzed?*

The purpose of the study is to compare total lead concentrations at the storage yard and in surrounding background areas.

- b. *Is the chemical likely to be a background*

constituent?

Because of the nearby highways and other industrial sources in the vicinity of the yard, background lead concentrations are expected to be elevated. This example will include statistical evaluation of only unpaved areas.

Step 3. Identify Inputs into the Decision

a. *Which chemicals will be analyzed?*

EPA decides to focus on total lead concentration.

b. *Which soil types and depths need to be sampled?*

Because there is neither surface evidence, nor historical record, of excavation in the storage area, EPA decides to measure total lead concentration in the first 12 inches of surface soils. Soils in background locations will be sampled in the same way. The TRW has recommended soil sieving at 250 μm to assess exposures to lead on the fine fraction of soil and dust.¹⁶ For background sampling of lead, this fractionation may be appropriate as it relates to human health risks.

c. *Which comparison tests are likely to be used?*

EPA expects that lead concentrations may not be normally or lognormally distributed. The study team decides to use a nonparametric statistical test for differences in the soil lead concentration distribution in the storage yard and in the surrounding areas.

d. *What coefficient of variation is expected?*

Based on previous sampling in other areas, a coefficient of variation ranging from 50% to 200% is expected. Preliminary data collected at the site indicate a standard deviation of approximately 50 mg/kg. Since this estimate is based on very limited data, the team decides to use a more conservative, preliminary estimate of $\sigma = 75$ mg/kg in the first stage of planning the survey design.

e. *What preliminary remediation goals (PRGs) may need to be met?*

A PRG of 400 mg/kg is available for residential sites.¹⁷

f. *What are the desired power and confidence levels?*

The study team decides initially on a Type I decision error limit of $\alpha = 0.10$ and a Type II decision error limit of $\beta = 0.10$ (power = 90%). The team agrees to review this decision, depending on the overall cost estimates produced by these objectives.

Step 4. Define Boundaries of the Study

a. *What geographic areas should be investigated?*

The study team decides that the entire storage yard area, approximately 5 acres, will be included in the study. Four different background areas of approximately 10,000 m^2 were selected at distances of between 1,000 m and 10,000 m from the storage yard boundaries.

b. *What are the characteristics of the soil data or population of interest?*

Soil samples should be collected in dry, unpaved areas. Prepared samples should be free of roots, leaves, and rocks or other consolidated materials. When preparing the samples, these materials should be removed using a 3 cm diameter sieve. Oversized materials should be retained for additional weighing and analysis, if necessary.

c. *How should the soil data be stratified statistically into relatively homogeneous characteristics?*

No stratification is planned for this study.

d. *What is the time frame to which the decision applies?*

Sampling will be conducted during a four-week period in the fall. Lead concentrations in soil are relatively static, and decisions based on the sampling results will remain applicable for many years, barring additional contamination.

e. *What practical constraints may hinder sample collection?*

The plant manager agreed to permit EPA sampling on the storage yard. Permission must be obtained from the owners of the selected background sampling areas for permission to enter and to collect background samples on their property.

Step 5. Develop a Decision Rule

If the selected statistical test indicates that the mean concentration in potentially impacted areas exceeds the mean background concentration by more than a substantial difference, then the chemical will be treated as site-related. Otherwise, if the statistical test indicates that the mean concentration in potentially impacted areas does not exceed the background mean, the chemical will be treated as background-related.

a. *What should the null hypothesis be?*

The study team chooses a null hypothesis that the lead concentrations in the storage yard exceed background concentrations.

- ▶ H_0 : Lead concentrations in the storage yard samples exceed background concentrations by more than $S = 100$ mg/kg (see paragraphs c and d, below, and Appendix A for how 100 mg/kg was chosen).

b. *What is the alternative hypothesis?*

The alternative hypothesis is the opposite of the null hypothesis.

- ▶ H_A : Lead concentrations in the storage yard samples do not exceed the background concentrations by more than $S = 100$ mg/kg.

c. *What level constitutes a substantial difference above background?*

The study team decided to use 100 mg/kg as the value for a substantial difference in lead concentrations between the storage yard and background areas. Issues pertaining to the selection of a value for a substantial difference are discussed in Appendix A.

d. *Specify the gray region for the hypothesis test*

When using Background Test Form 2, the gray region of width MDD starts at a difference of $\Delta = S = 100$ mg/kg and extends on the left down to $\Delta = (S - \text{MDD})$. As a trial value, the study team chose to use an MDD that is one-half of S , 50 mg/kg (refer to Table 3.1). This MDD represents a balance between the cost of extra sampling and the expected cost of remediating the site unnecessarily.

Step 6. Specify the Limits on Decision Errors

a. *What is the possible range of the parameter of interest?*

The possible range of lead concentrations in industrial soil is very wide, ranging from 0 to many grams per kilogram.

b. *What are the acceptable decision errors (Type I and Type II)?*

The team decides that the acceptable limits on decision errors are $\alpha = 0.10$ for Type I errors at a difference of $\Delta = S = 100$ mg/kg, and $\beta = 0.10$ for Type II errors at a difference of $\Delta = S/2 = 50$ mg/kg.

In Figure 3.2, the test performance curve achieves a probability of 90% of detecting significant difference ($\Delta = S$). The study team is comfortable with the choice of a 90% confidence level for the test, because this reduces the chance of a false negative—deciding that the yard does not exceed background by more than S .

The choice of $\beta = 0.10$ and the selected value for the

MDD equal to one-half the width of the gray region means that the power of 90% will be required at $\Delta = S/2$. The plant manager recognizes that a lower value of β (higher power) would result in a lower probability of a Type II error and improve his chances of passing the test, but he has decided that the extra sampling costs required to achieve a higher power are not necessary.

- c. *What are the potential consequences of each type or error, specifying a range of possible parameter values (gray region) where consequences of decision errors are relatively minor?*

The team decides that the decision errors are $\alpha = 0.10$ at $\Delta = S$, and $\beta = 0.10$ at $\Delta = S/2$. The gray region extends from a difference $\Delta = 50$ mg/kg to a difference of $\Delta = 100$ mg/kg (refer to Figure 3.2 and decisions made in Steps 5d and 6b).

Test Form 2 has at least $100(1-\alpha)\%$ confidence of correctly detecting a site that exceeds background by more than S , regardless of the sample size. Greater sample size increases the power of the test and reduces β , which reduces the chance that a site is remediated unnecessarily. When using Test Form 2, extra samples represent the cost of increasing the chance that the site is determined to be acceptable when the true Δ is less than S . The study team agrees to review this decision, depending on the overall cost estimates produced by the decision objectives.

- d. *Do the limits on decision errors ensure that they accurately reflect the study team's concern about the relative consequences for each type of decision error?*

The study team is satisfied with the choice of the 90% confidence level for the statistical test, because this will reduce to 10% the chance of falsely deciding that the yard does not exceed background by more than 100 mg/kg when it truly does. The use of a level- α test will provide 90% confidence for all sample sizes, but may have poor power if the sample size is too low.

The sample size is fixed by the choice of MDD and β . Choosing $\beta = 0.10$ at a difference of $\Delta = 50$ mg/kg means that a power of at least 90% will be obtained if the true lead concentration on the yard is at or below that value. The plant manager recognizes that a lower value of β (higher power) would result in a lower probability that the test will decide the yard exceeds background lead concentrations if the yard is only 50 mg/kg higher than background. However, the manager has decided that this extra power would require more sampling and unwanted additional sampling costs.

The DQO parameters α , β , S , MDD, and σ provide the information needed to calculate the number of samples (N) required from each population. N samples will be collected in contaminated areas, and N samples will be collected from background areas.

The sample size may be calculated using the approximate formulas presented in Chapter 3 of EPA QA/G-9.² The approximate sample size calculated with the values $\alpha = 0.10$, $\beta = 0.10$, MDD = 50 mg/kg, with a conservative estimate for σ of 75 mg/kg, is $N = 35$, as shown in Table 3.1. If the actual σ is measured to be only 50 mg/kg, as indicated by preliminary data, then only 16 samples would be required in each area. In this case, a retrospective power analysis would show that the design had more than adequate power. If the actual σ is measured to be 100 mg/kg, then 62 samples would be required. In this latter case, the retrospective power analysis would indicate that the design did not have adequate power to make the decision and additional samples should be collected. The estimate of σ is one of the most important design parameters, and the success of the survey design will depend strongly on the accuracy of this estimate. More specific sample-size calculation procedures are given in MARSSIM.¹⁸

Step 7. Optimize the Sampling Design

What is the most resource effective sampling and analysis design for generating data that are expected to satisfy the DQOs?

- a. *Review the DQO outputs and existing environmental data*

The statistician, chemist, and plant engineer on the study team have reviewed the outputs developed at each stage of the DQO process.

b. *Develop general sampling and analysis design alternatives*

The study team decides to use a randomly-oriented, rectangular grid sampling strategy for the storage yard and selected background area. Two random numbers (x and y) randomly will determine the starting point selected for the grid. The grid orientation will be determined by a third random number. The size of the grid will be calculated based on the number of samples required for each area.

c. *Verify that DQOs are satisfied for each design alternative*

Only one sample design is used in this study.

d. *Select the most resource-effective design that satisfies all of the DQOs*

Alternative sampling designs may result in lower sampling costs. The study team agrees to consider the alternative sample designs suggested in EPA QA/G-5S before the sampling program begins.¹²

e. *Document the operational details and theoretical assumptions of the selected design in the sampling and analysis plan*

The EPA team has documented the discussions leading to each DQO parameter.

CHAPTER NOTES

1. U.S. Environmental Protection Agency (EPA). 2001. *Requirements for Quality Assurance Project Plans, EPA QA/R-5*. <http://www.epa.gov/quality/qapps.html>.
2. U.S. Environmental Protection Agency (EPA). 2000. *Guidance for Data Quality Assessment: Practical Methods for Data Analysis, EPA QA/G-9, QA00 Version*. Quality Assurance Management Staff, Washington, DC, EPA 600-R-96-084. Available at http://www.epa.gov/quality/qa_docs.html.
 - ▶ Equations for computing retrospective power are provided in the detailed step-by-step instructions for each hypothesis test procedure in Chapter 3.
3. U.S. Environmental Protection Agency (EPA). 1994. *Guidance for the Data Quality Objectives Process, EPA QA/G-4*, EPA 600-R-96-065. Washington DC.
4. U.S. Environmental Protection Agency (EPA). 1989. *Risk Assessment Guidance for Superfund Vol. I, Human Health Evaluation Manual (Part A)*. Office of Emergency and Remedial Response, Washington, DC. EPA 540-1-89-002. Hereafter referred to as “RAGS.”
5. Mathematically, Background Test Form 1 is written:
$$H_0: \Delta \leq 0 \text{ vs } H_A: \Delta > 0$$
with $\Delta = \theta_S - \theta_B$, where θ_S is the selected decision parameter (mean, median, etc.) for the site distribution, and θ_B is the same parameter for the background distribution.
6. U.S. Environmental Protection Agency (EPA). 1989. *Statistical Methods for Evaluating the Attainment of Cleanup Standards*, EPA 230/02-89-042, Washington DC.
7. Mathematically, Background Test Form 2 uses the substantial difference S as a non-zero action level:
$$H_0: \Delta > S \text{ vs } H_A: \Delta \leq S$$
with $\Delta = \theta_S - \theta_B$, where θ_S is the selected decision parameter (mean, median, etc.) for the site distribution, and θ_B is the same parameter for the background distribution.
8. U.S. Environmental Protection Agency (EPA). 1990. *Guidance for Data Usability in Risk Assessment: Interim Final, October 1990*. EPA 540-G-90-008, PB91-921208, Washington, DC.
9. U.S. Environmental Protection Agency (EPA). 1994. *The Data Quality Objectives Decision Error Feasibility Trials (DEFT) Software (EPA QA/G-4D)*, EPA/600/R-96/056, Office of Research and Development, Washington, DC.
10. U.S. Environmental Protection Agency (EPA). 1996. *The Data Quality Evaluation Statistical Toolbox (DataQUEST) Software (EPA QA/G-9D)*, Office of Research and Development, Washington, DC.
11. *Guidance for Choosing a Sampling Design for Environmental Data Collection*, EPA QA/G5S, U.S.EPA, Office of Environmental Information, Peer Review Draft, Aug. 2000.

12. For further guidance on the use of hypothesis tests in environmental decision making, see *EPA QA/G-4, Guidance for the Data Quality Objectives Process*, EPA/600/R-96/055. The theory of hypothesis testing is discussed in many introductory statistics textbooks, including the popular text by Mood et al. (1974) *Introduction to the Theory of Statistics, 3rd Ed.*, McGraw Hill, Chapter IX. Readers with some background in statistics may refer to Chapter 5 of *Mathematical Statistics: Basic Ideas and Selected Topics*, P.J. Bickel and K.A. Doksum, Holden-Day, 1977, for a discussion of error rates and relative importance of the errors (p. 168) that can be committed in hypothesis testing.
13. Gilbert, Richard O. 1987. *Statistical Methods for Environmental Pollution Monitoring*, Van Nostrand Reinhold.
14. Common references for sample selection include:
 - ▶ Cochran, W. 1977. *Sampling Techniques*. New York: John Wiley.
 - ▶ Gilbert, Richard O. 1987. *Statistical Methods for Environmental Pollution Monitoring*. New York: Van Nostrand Reinhold.
 - ▶ U.S. Environmental Protection Agency (EPA). 1989. *Statistical Methods for Evaluating the Attainment of Cleanup Standards*. *Op. cit.*
 - ▶ U.S. Environmental Protection Agency (EPA). 1990. *Guidance for Data Usability in Risk Assessment*. *Op. cit.*
15. EPA's Technical Review Workgroup for Lead provides technical assistance for people working on lead-contaminated sites. For assistance or more information, the reader should refer to their website (<http://epa.gov/superfund/programs/lead>) or call the Lead Hotline (800-680-5323).
16. U.S. Environmental Protection Agency (EPA). *TRW Recommendations for Sampling and Analysis of Soil at Lead (Pb) Sites*. Office of Emergency and Remedial Response, Washington, DC. EPA 540-F-00-010, OSWER 9285.7-38.
17. U.S. Environmental Protection Agency (EPA). 1994. *Revised Interim Soil Lead Guidance for CERCLA Sites and RCRA Corrective Action Facilities*. OSWER Directive 9355.4-12.
18. U.S. Environmental Protection Agency (EPA), U.S. Nuclear Regulatory Commission, et al. 2000. *Multi-Agency Radiation Survey and Site Investigation Manual (MARSSIM)*. Revision 1. EPA 402-R-97-016. Available at <http://www.epa.gov/radiation/marssim/> or from <http://bookstore.gpo.gov/index.html> (GPO Stock Number for Revision 1 is 052-020-00814-1).

CHAPTER 4

PRELIMINARY DATA ANALYSIS

This chapter provides guidance for preliminary data analysis using graphs and distributions of the data. Depending upon the quality of existing site and background data, quantitative analysis used to establish background concentration may involve a combination of comparative statistical analysis and graphical methods. The preliminary data analysis is an integral part of choosing the appropriate methods for making statistically valid comparisons of site and background concentrations.

Preliminary data analysis should include a detailed “hands-on” inspection of the site and background data before proceeding to the statistical tests. Graphs are used to identify patterns and relationships within the onsite and background data sets, and to compare the two data sets. Preliminary data analysis should be focused on verifying assumptions, such as normality, made in the DQO process. The review should identify anomalies in the data, including potential outliers. This step is formally a part of the Data Quality Assessment.¹

The preliminary inspection may include development of a posting plot,¹ which is a map showing the measured concentration and location of each sample. The posting plot may reveal likely sources of contamination, important areas that have not been sampled, spatial correlations or trends in the data, and the location of suspected outliers. Note that one possible outcome of the preliminary data inspection is that the chemical concentrations detected at the site are much higher than background ranges reported for similar soil types. In this case, a formal background analysis may not be necessary if all or most of the detected concentrations are well above

the range likely to represent background. Another possible outcome of the preliminary analysis is that all chemical concentrations are well below risk-based screening levels. In this case, background analysis is not likely to be necessary.

This chapter presents information useful for both parametric and nonparametric data analysis (defined in the box below). Parametric statistical methods are based on the assumption of a known mathematical form for the probability distributions that represent the site and background populations. For many parametric methods, the data user should first determine whether the data are normally distributed, using any of several tests for normality.

Parametric and Nonparametric Methods

Parametric: A statistical method that relies on a known probability distribution for the population from which the data are selected. Parametric statistical tests are used to evaluate statements (hypotheses) concerning the parameters of the distribution.

Nonparametric: A distribution-free statistical method that does not depend on knowledge of the population distribution.

Nonparametric methods do not require that the data distribution be characterized by a known family of distributions. Several graphical methods are presented for nonparametric comparisons.

4.1 Tests for Normality

Tests for normality are an important step in assessing the type of statistical test to use for comparison with background. Parametric tests, such as the t-test for comparing the means of the site and background distributions, are usually based on the assumption of normality for both data sets. Before using a parametric test for a background comparison, tests should be conducted on each data set to show whether it meets the assumption of normality. If the raw data are not normally or lognormally distributed, other types of transformations should be conducted to approximate normality prior to using the data sets in parametric statistical comparisons.

Since it is unusual to encounter environmental data sets that are normally distributed,² these tests are most commonly applied after a transformation. Usually the logarithms of the data have been applied to the raw data. The test for normality is then applied to the transformed data sets. In most cases, direct application of a nonparametric background comparison test using the raw data is preferred to using a parametric test on transformed data. This is particularly true when there are outliers and/or non-detect values in the raw data. The assumption of normality is very important as it is the mathematical basis for the majority of parametric statistical tests. Examples of how to perform each of these tests can be found in Chapter 4 of EPA QA/G-9.¹

The *Shapiro-Wilk* test is a powerful general purpose test for normality or lognormality when the sample size is less than or equal to 50, and is highly recommended. The Shapiro-Wilk test is an effective method for testing whether a data set has been drawn from an underlying normal distribution. It can also evaluate lognormality if the test is conducted on logarithms of the data. If the normal probability plot is approximately linear—the distribution follows a normal curve—the test statistic will be relatively high. If the normal probability plot contains significant curves, the test statistic will be relatively low.

Another test related to the Shapiro-Wilk test is the

Filliben statistic, also called the “probability plot correlation coefficient.” If the normal probability plot is approximately linear, the correlation coefficient is relatively high. If the normal probability plot contains significant curves—the distribution does not follow a normal curve—the correlation coefficient will be relatively low. The Filliben test is recommended for sample sizes less than or equal to 100.

D’Agostino’s test for normality or lognormality is used when sample sizes are greater than 50. This test is based on an estimate of the standard deviation obtained using the ranks of the data. This estimate is compared to the usual mean square estimate of the standard deviation, which is appropriate for the normal distribution.

The *studentized range test* for normality is based on the fact that almost 100 percent of the area of a normal curve lies within ± 5 standard deviations from the mean. The studentized range test compares the range of the sample to the sample standard deviation. For example, if the minimum of 50 data points is 40.2, the maximum is 62.7, and the standard deviation is 4.2, then the studentized range is $(62.7 - 40.2)/4.2 = 5.4$. Tables of critical sizes up to 1,000 are available for determining whether the absolute value of the studentized range is significantly large. Using, for example, Table A-2 in EPA QA/G-9¹ the upper critical values for the studentized range test with $n = 50$ are 5.35 for $\alpha = 0.05$ and 5.77 for $\alpha = 0.01$. In this example, the assumption of normality would be accepted at the 95% confidence level, but rejected at the 99% confidence level. The studentized range test does not perform well if the distribution is asymmetric and if the tails of the distribution are heavier than the normal distribution. In most cases, this test performs as well as the Shapiro-Wilk test and is easier to apply.

4.2 Graphical Displays

Graphical methods provide visual examination of the site and background distributions, and compari-

sons of the two. Graphical methods supplement the statistical tests described in Chapter 5. Graphical methods also may be used to verify that the assumptions of statistical tests are satisfied, to identify outliers, and to estimate parameters of probability distributions that fit to the data. The methods described in this section assume that separate data sets are collected on site and in background. In some situations, an appropriate background area with similar soil types and chemistry cannot be identified. Graphical methods designed for analyzing sample data collected from both impacted and non-impacted areas on the site are addressed by Singh et al. (1994).³

4.2.1 Quantile Plot

A *quantile plot* displays the entire distribution of the data, ranging from the lowest value to the highest value. The vertical axis for the quantile plot is the measured concentration, and the horizontal axis is the percentile of the distribution. Each ranked data value is plotted against the percentage of the data with that value or less.

To construct a quantile plot, the data set is ranked from smallest to largest. The percentage value for each data point of rank_j is computed as

$$\text{Percent}_j = 100 (\text{rank}_j - 0.5) / n$$

where n is the number of values in the data set. There are two quantile plots in Figure 4.1, one for the site data and another for background data. If one or more data values are non-detects, all non-detects are ranked first, below the first numerical value. The plot starts with the first numerical value. For example, if a data set with 10 observations has 2 non-detect values, then the smallest detected value has rank 3 and a percentage of $100(3 - 0.5)/10 = 25$. The highest 8 data points would be shown on the plot, starting at the 25th percentile.

The slope of the curve in the quantile plot is an indication of the amount of data in a given range of values. A small amount of data in a given range will

result in a large slope for the quantile plot. A large amount of data in a range will result in a more horizontal slope. A sharp rise near the bottom or the top of the curve may indicate the presence of outliers.

A graph may contain more than one quantile plot. In a *double-quantile plot*, the site and background data are each plotted in a single graph, providing a direct visual comparison of the two distributions. A curve that is higher in the vertical direction indicates a higher distribution of data values.

An example of the double-quantile plot is shown in Figure 4.1. The lower curve shows the distribution of the background data, and the middle curve (indicated by symbols only) shows the quantile plot for the site data. In this example, the entire site distribution is higher than the background distribution indicating that some degree of contamination is likely. The close proximity of the site and background quantile plots near the 70th percentile and rapid divergence above indicate a larger difference between the two distributions in the upper 30 percent of the distributions. At the left end of the plot, the background data distribution falls off more rapidly to zero concentration below the 20th percentile than the site data distribution, which has a y-intercept substantially above zero. The positive intercept and roughly parallel shape of the three lines in the plot below the 70th percentile suggest that the distribution of concentrations on site is shifted to higher levels than the background distri-

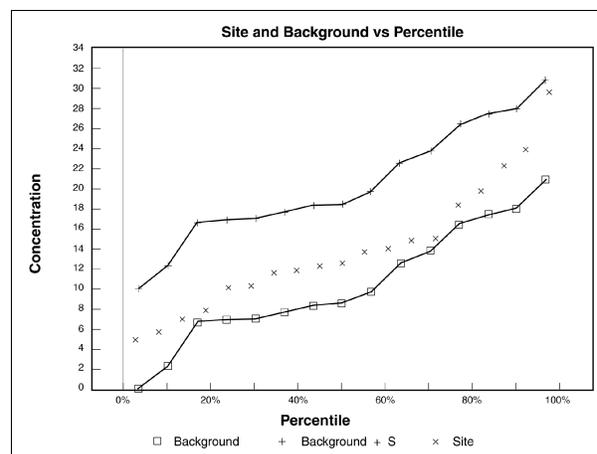


Figure 4.1 Example of a double quantile plot

bution, with a larger shift above the 70th percentile. The upper curve in the figure shows the background distribution augmented by $S = 10$, a hypothetical value for a substantial difference over background. In this example, the entire site distribution lies below the S -augmented background distribution, indicating that the site does not exceed background by more than a substantial difference.

Issues affecting the determination of site-specific values for a substantial difference are discussed in more detail in Appendix A of this guidance.

The formal statistical test procedures presented in Chapter 5 may be used to make decisions that confirm or deny these graphical indications with predetermined error rates. In this and the following exhibits, contaminant concentrations are plotted using a linear scale. If the data are highly variable, it may be necessary to transform the graph by using a logarithmic scale for the concentration axis. Use of the logarithmic transformation does not affect the ranks of the data.

4.2.2 Quantile-Quantile Plots

A *quantile-quantile plot* is useful for comparing two distributions in a single graph. The vertical axis of this plot represents the first distribution of values, and the horizontal axis represents the second distribution. The scales for the concentration axes may be either both linear or both logarithmic. If the two distributions are identical, the quantile-quantile plot will form a straight line at 45 degrees when equal scales are used for the two axes. The slope of this line has a value of one, regardless of the selected scales. Deviations from this line show the differences between the two distributions.

There are two common applications of the quantile-quantile plot. One type is used for parametric applications, and the other for nonparametric comparisons.

- ▶ *Parametric Quantile-Quantile Plot.* In parametric applications of the quantile-quantile plot, the horizontal axis represents the quantiles from

a known distribution, such as the normal distribution. This application is referred to as a *normal probability plot*. If the data follow a normal distribution, the plot will appear as a straight line. Probability plots are useful for determining if the site data or the background data follow a normal or lognormal distribution. More information on the use of the quantile-quantile plot to compare with known parametric distributions is provided in Section 2.3 of EPA QA/G-9.¹

- ▶ *Empirical Quantile-Quantile Plot.* In nonparametric applications, the empirical quantile-quantile plot is used to compare two data sets. In our case, the two data sets are the site distribution and the background distribution. If there are an equal number of data values in the two data sets, it is very easy to construct an empirical quantile-quantile plot. The graph is constructed by plotting each ranked site value against the corresponding background value with the same rank.

The empirical quantile-quantile plot is useful because it provides a direct visual comparison of the two data sets. An example of the quantile-quantile plot is shown in Figure 4.2. If the site and background distributions are identical, the plotted values would lie on a straight line through the origin with slope equal to 1, shown in the figure as the line

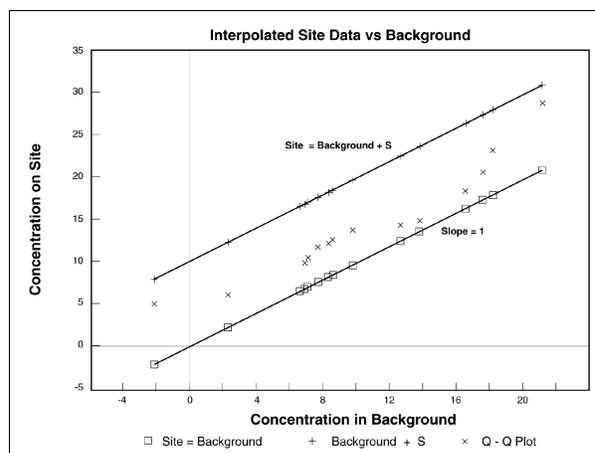


Figure 4.2 Example of an empirical quantile-quantile plot

labeled “Site=Background.” Any deviation from this line shows differences between the two distributions. The points that mark the empirical Q-Q plot in the figure are above the line that indicates equality of the two distributions. This indicates that the distribution of site measurements exceeds the distribution of background measurements. If the site differs from the background data distributions only by an additive difference along the entire distribution, the plotted site values will lie on a straight line with slope 1 that does not pass through the origin. If the site distribution is t units above the background distribution, the straight line will have slope 1 and a y intercept at $+t$.

A hypothetical level of substantial contamination, S , is shown in the upper plot in Figure 4.2 labeled “Background + S .” Note that the median interpolated site value is plotted against the median of the background values at the center of the plot. When this point lies above the equal-distribution line with slope 1, the median interpolated site value is larger than the median background value.

When the size of the data set differs from the size of the background data set, interpolation is used to construct the empirical quantile-quantile plot. Detailed procedures for creating a quantile-quantile plot with unequal sample sizes are provided in Section 2.3.7.4 of EPA QA/G-9.¹

4.2.3 Quantile Difference Plot

The nonparametric *quantile difference plot* is a variant of the empirical quantile-quantile plot. When site data are compared to background data, the quantity of greatest interest is the amount by which the site distribution exceeds the background distribution. This difference can be viewed in the empirical quantile-quantile plot of Figure 4.2 as the difference between two sloped lines, the quantile-quantile plot and the line with slope 1 where site equals background. More resolution for examining the differences between the site and background distributions is obtained by subtracting each background value from its corresponding interpolated site value, then plotting the differences versus their

corresponding background values.

An example of the quantile difference plot is shown in Figure 4.3. In the quantile difference plot, background is represented by the horizontal axis. The distribution of background values is shown by the symbols plotted on this axis. A hypothetical level of substantial contamination of $S = 10$ appears in this plot as a horizontal line, not to be exceeded. In this example, the entire quantile difference plot lies between the background and the substantial difference level, indicating that the site exceeds background by a small amount, but does not exceed background by more than a substantial difference.

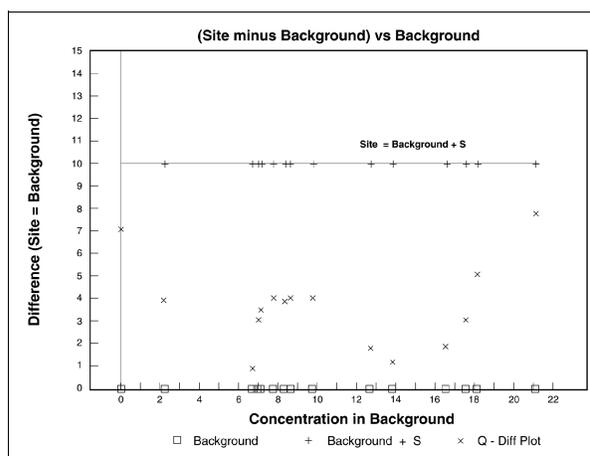


Figure 4.3 Example of a quantile difference plot

The quantile difference plot permits a quick visual evaluation of the amount by which the site exceeds background. In this example, the largest differences occur in the upper half of the distribution. It is clear that the interpolated site values do not exceed background by more than the hypothetical $S = 10$ concentration units. This conclusion is not as obvious using the sloped quantile-quantile plot in Figure 4.2.

Similar warnings exist for use of the quantile difference plot as for the empirical quantile-quantile plot when there are more than twice as many site values as background values. The empirical quantile-quantile plot and the quantile difference plot work best when the site and background data sets are of approximately the same size, and they depend

upon the choice of S.

4.3 Outliers

Outliers are measurements that are unusually larger or smaller than the remaining data. They are not representative of the sample population from which they were drawn, and they distort statistics if used in any calculations. Statistical tests based on parametric methods generally are more sensitive to the existence of outliers in either the site or background data sets than are those based on nonparametric methods.

Outliers can lead to both Type I and Type II errors. They can lead to inconclusive results if the results are highly sensitive to the outliers. There are many plausible reasons for the presence of outliers in a data set:

- ▶ Data entry errors. Data that are extremely high or low should be verified for data entry errors.
- ▶ Missing values and non-detects. It is important that missing value and non-detect codes are not read as real data. For example the number 999 might be a code for missing data, but the computer program used to analyze the data, if not properly designated, could misread this as an extreme value of 999. This is easily remedied.
- ▶ Sampling error. In this case the sample results for the sample that is not from the population of interest should be deleted. However, using data from a population other than the one of concern is not easily recognized. Therefore, this type of error can result in the presence of outliers in the data set.
- ▶ Non-normal population. An outlier might also exist when a sample is from the population of interest, but its distribution has more extreme values than the normal distribution. In this situation, the sample can be retained if a statistical approach is selected for which the outliers do not have undue impact.

Outliers may misrepresent the sample population from which they were taken, and any conclusion drawn that is based on these results may be suspect. Outliers may be true measurements of conditions on the site, or may be due to faulty sample collection, cross contamination, lab equipment failure, or improper data entry. To determine which case applies, the outliers should first be identified. If there is a large number of outliers in the data set, it may be necessary to reassess the area. Outliers in the site data set have different implications from outliers in the background data set. For example, an onsite outlier can indicate a “hot spot,” which indicates that the one spot needs attention. An outlier in the background data set, however, might indicate that one of the background samples was collected in a location that is not truly background. In such a case, an outlier test should be used (along with a qualitative study of where the sample in question was collected) to see if that data point should be discarded from the background set. Additional guidance for handling outliers is provided in EPA QA/G-9, Section 4.4.¹

Data points that are flagged as outliers should be eliminated from the data set if field or laboratory records indicate that the sample location was not a reasonable reference area, or if there was a problem in collecting or analyzing the sample. However, background areas are not necessarily pristine areas. A data point should not be eliminated from the background data set simply because it is the highest value that was observed. The use of nonparametric hypothesis tests for background comparisons greatly reduces the sensitivity of test results to the presence of outliers. Parametric tests based on the lognormal distribution may yield results that are extremely sensitive to the presence of one or more outliers.

Statistical outlier tests give probabilistic evidence that an extreme value does not “fit” with the distribution of the remainder of the data and is therefore a statistical outlier. There are five steps involved in treating extreme values or outliers:

1. Identify extreme values that may be potential outliers;

2. Apply statistical tests;
3. Review statistical outliers and decide on their disposition;
4. Conduct data analyses with and without statistical outliers; and
5. Document the entire process.

More guidance on handling outliers is given in Chapter 5.

4.4 Censored Data (Non-Detects)

Contamination on the site or in background areas may be present at concentrations close to the detection limits. A sample is said to be “censored” when certain values are unknown, although their existence is known. *Type I censoring* occurs when the sample is censored by reference to a fixed value. Non-detect measurements are examples of Type I left censoring. The specific value is unknown, but the existence of a concentration value in the closed interval from 0 to the reporting limit may be inferred. The value may be 0, or a small positive value less than the detection limit. If other measurements collected on the site indicate concentrations above the detection limit, then the likelihood that at least some of the non-detects represent small positive values is increased. Concentration values may be censored at their detection limits or at some arbitrary level based on detection limits.

A detection limit is the smallest concentration of a substance that can be distinguished from zero. Consequently, non-detects may not represent the absence of a chemical but its presence at a concentration below its reliable minimum detection level. Many parametric statistical methods require numerical values for all data points. One approach is to impute a surrogate value for non-detects, commonly assumed to be half the reporting limit. The use of $L / \sqrt{2}$ has also been recommended.⁴ Alternatively, a random value between the reporting limit and zero may be chosen to represent each non-detect for the

purposes of testing assumptions concerning distributions. Both approaches may seriously affect the estimated distribution parameters.⁵

If less than 15 percent of the site and background samples are non-detects, then distributions of both the background and the site sample may be determined by using surrogate values. Probability plots and goodness-of-fit tests may be performed for each data set, first including the non-detects as part of the sample using random values for non-detects, and second, excluding the non-detects from the sample. If the two sets of estimated parameters differ only slightly, then the non-detect problem is of lesser importance. However, if the two sets of estimates differ significantly, then the surrogate value approach should be re-evaluated.

If more than 15 percent and less than 50 percent of the measurements in the background sample set or the site sample set are non-detects, the use of specialized methods for analyzing non-detects is recommended. Section 4.7 of EPA QA/G-9¹ describes in detail several methods for estimating the mean and standard deviation of data sets with non-detects.

If more than 50 percent of the measurements in either the background sample set or the site sample set are non-detects, it may not be possible to compare the means of the two distributions. An alternative approach is to compare the upper percentiles of the two distributions by comparing the proportion of the two populations that is above a fixed level.¹ Comparisons may be made for the upper percentiles of each distribution despite the large number of non-detects.

Nonparametric methods may be used to avoid the necessity of imputing surrogate values for non-detect measurement. Nonparametric methods are often based only on the ranks of the data, and the non-detect values can be assigned unambiguous ranks without the need for assigning surrogate values. Bootstrapping and other nonparametric methods have recently received attention.⁶

CHAPTER NOTES

1. U.S. Environmental Protection Agency (EPA). 2000. *Guidance for Data Quality Assessment: Practical Methods for Data Analysis, EPA QA/G-9, QA00 Version*. EPA 600-R-96-084. Quality Assurance Management Staff, Washington, DC. Available at http://www.epa.gov/quality/qa_docs.html.
 - ▶ See Chapter 2 for more information on the role of preliminary data analysis in the Data Quality Assessment process.
 - ▶ See Section 2.3 for information on the use of a quantile-quantile plot, and Section 2.3.7.4 for detailed procedures for creating a quantile-quantile plot with unequal sizes.
 - ▶ See Section 2.3.9.1 for guidance on preparing a posting plot.
 - ▶ See Section 3.3.2.1 for recommendations on dealing with high proportions of non-detects.
 - ▶ See Chapter 4 for examples of how to test for normality.
 - ▶ See Section 4.4 for guidance on outliers.
 - ▶ See Section 4.7 for methods of estimating mean and standard deviation of data with non-detects.
 - ▶ See Table A-2 for critical values.
2. U.S. Environmental Protection Agency (EPA). 1992. *Supplemental Guidance to RAGS: Calculating the Concentration Term*, Publication 9285.7-08I, Office of Solid Waste and Emergency Response, Washington, DC.
3. Singh, A., A.K. Singh, and G. Flatman. 1994. "Estimation of background levels of contaminants," *Mathematical Geology*, 26:3.
4. See, for example, cleanup regulations in the Model Toxics Control Act, State of Washington, WAC 173-340-708(11)(e).
5. A detailed consideration of non-detects is included in *Statistical Guidance for Ecology Site Managers, Supplement S-6, Analyzing site of background data with below-detection limit or below-PQL values (Censored data sets)*, Washington State Department of Ecology, Olympia, WA, August 1993.
6. U.S. Environmental Protection Agency (EPA). 1997. *The Lognormal Distribution in Environmental Applications*, EPA/600/R-97/006. Office of Research and Development, Environmental Sciences Division, Las Vegas, NV.

CHAPTER 5

COMPARING SITE AND BACKGROUND DATA

This chapter provides guidance on selecting quantitative statistical approaches for comparing site data to background data. Statistical methods allow for specifying (controlling) the probabilities of making decision errors and for extrapolating from a set of measurements to the entire site in a scientifically valid fashion.¹

Several methods are available for comparing background to site data. These can be divided into several major categories: data ranking and plotting, descriptive summaries, simple comparisons, parametric tests, and nonparametric tests. For many of these methods, data users first should determine whether the data are normally distributed, using any of several tests for normality. Data can also be assessed in terms of the whole data set from the site, or with a focus on outliers in the background data set or in the contaminant concentrations at the site (see Chapter 4).

The issue of randomness is an important element of most statistical procedures when sample results are to be extrapolated to the entire site or background sampling area, rather than only representing the areas where measurements were made. The statistical tests discussed in this chapter assume that the data constitute a random sample from the population. If a sample of measurements is to represent the entire site, every sampling point within the area represented by the sample should have a non-zero probability of being selected as part of the sample. If all points have an equal opportunity for selection, the sampling procedure will generate a simple random sample. A random sample implies independence, loosely meaning that the samples are also uncorrelated. If

samples are too closely spaced, then adjacent samples may exhibit a high degree of correlation. This lack of independence is avoided by using a grid sampling technique.

Most procedures presented in this chapter require a simple random sample. Stratification of the site will usually result in differing probabilities of selection within each stratum. A stratified sample is not a simple random sample, and a statistician should be consulted before conducting the analysis. In this context, the statistician would advise on the appropriate calculations to use for estimation and hypothesis testing if a stratified design has been selected.

Judgmental (or “authoritative”²) samples are those collected in areas suspected to have higher contaminant concentrations due to operational or historical knowledge. Judgmental samples may result from sampling conducted for overall site characterization, developing exposure point concentrations, or sampling specifically to delineate areas requiring remediation. Judgmental samples cannot be extrapolated to represent the entire site. In some cases, there is a great deal of bias associated with the collection of judgmental samples. The statistical hypothesis testing procedures recommended in this chapter are based on random samples and should not be used on judgmental samples. If judgmental sampling is used on site, while background measurements are collected randomly, direct comparison of the means of the two data sets is not recommended.

Graphical methods, such as posting plots, may be

Method	Application	Comments
Descriptive Summary <ul style="list-style-type: none"> ▶ □ Mean ▶ □ Median ▶ □ Standard deviation ▶ □ Variance ▶ □ Percentiles 	Preliminary examination of data for comparison with site history and land use activities in the establishment of background. Use as a preliminary screening tool.	Simple and straightforward; less statistical rigor.
Simple Comparisons	Used with very small data sets.	Not recommended
Parametric Tests <ul style="list-style-type: none"> ▶ □ Student t-test ▶ □ Behrens-Fisher Student t-test 	Tests require approximate normality of the estimated means. Use if a larger number of data points are available ($n > 25$). For smaller data sets, examine data for normality or lognormality in distribution. ⁴	Statistically robust and used frequently in parametric data analysis.
Nonparametric Tests <ul style="list-style-type: none"> ▶ □ Wilcoxon Rank Sum Test (also called the “Mann-Whitney Test”) ▶ □ Gehan Test 	Use when data are not normally distributed, as rank-ordered tests make no assumption on distribution.	Statistically robust and used frequently in background estimation.

used to display judgmental data. These displays may reveal likely sources and pathways of contamination. Kriging³ and other spatial smoothing algorithms may be applied to identify areas with suspected high concentrations for conducting the remediation, although the estimated mean concentrations should be recognized for their upward bias.

Depending upon the data and other site-specific considerations, statistical analysis should involve one or a combination of the following methods:

- ▶ Parametric statistical comparison methods involving comparison of one or more parameters of the distribution of site samples with the corresponding parameter of the background distribution, such as the Student t-test; or
- ▶ Nonparametric tests, such as Wilcoxon Rank Sum (WRS) test.

The box at the top of this page lists examples of the statistical tests and applications recommended for establishing background constituent concentrations. These and other useful tests are discussed in more

detail in the following sections.

5.1 Descriptive Summary Statistics

Several statistics can be used to describe data sets. These statistics may be used in many of the tests described later in this chapter. There are two important features of a data set: *central tendency* and *dispersion*.

Estimators of central tendency include the arithmetic mean, median, mode, and geometric mean. The sample mean is an arithmetic average for simple random sampling designs; however for complex sampling designs, such as stratification, the sample mean is a weighted arithmetic average. The sample mean is influenced by extreme values (large or small) and can easily be influenced by non-detects. The sample median value is directly in the middle of the data when the measurements are ranked in order from smallest to largest. More simply, the median is the middlemost value in the data set when the number of data values is odd. When the number of

data points is even, the median is usually defined as the average of the two middlemost values. The median is less affected by the presence of values recorded as being below the detection limit.

The dispersion around the central tendency is described by such items as the range, variance, sample standard deviation, and coefficient of variation. The easiest measure of dispersion is the sample range. For small samples, the range is easy to interpret and may adequately represent the spread of the data. For large samples, the range is not very informative because it only considers and is greatly influenced by extreme values. The sample variance measures the dispersion from the mean of a data set and is affected by extreme values and by a large number of non-detects. The coefficient of variation (CV) is a unitless measure that allows the comparison of dispersion across several sets of data. The CV is often used instead of the standard deviation in environmental applications because the standard deviation is often proportional to the mean. The standard deviation is affected by values below the detection limit, and some method of substituting numerical values for these should be found.⁴

5.2 Simple Comparison Methods

Simple comparison methods rely on descriptive summary statistics, such as comparing means or maximums. These approaches can be used with very small data sets but are highly uncertain.

5.3 Statistical Methods for Comparisons with Background

Many statistical tests and models are only appropriate for data that follow a particular distribution. Statistical tests that rely on knowledge of the form of the population distribution for the data are known as *parametric* tests, because the test is usually phrased in terms of the parameters of the distribution assumed for the data. Two of the most important distributions for tests involving environmental data are the normal distribution and the lognormal distribution. A normal distribution has only two parameters, the mean and

variance. Lognormal distributions also have only two parameters, but there are several common ways to parameterize the lognormal distribution. In this chapter, use of parametric comparison methods like t-tests or ANOVA may require normalization of data by conversion to a log scale.⁵

Tests for the distribution of the data (such as the Shapiro-Wilk test for normality) often fail if there are insufficient data, if the data contain multiple populations, or if there is a high proportion of non-detects in the sample.⁶ Tests for normality lack statistical power for small sample sizes. In this context, “small” may be defined roughly as less than 20 samples, either on site or in background areas. Some standard tests for a particular distribution against all alternatives, such as the Lilliefors form of the Kolmogorov-Smirnov test, require as many as 50 samples. Therefore, for small sample sizes or when the distribution cannot be determined, nonparametric tests should be used to avoid incorrectly assuming the data are normally distributed when there is not enough information to test this assumption.

Statistical tests that do not assume a specific mathematical form for the population distribution are called distribution-free or *nonparametric* statistical tests. Nonparametric tests have good test performance for a wide variety of distributions, and their performance is not unduly affected by outliers. Nonparametric tests can be used for normal or non-normal data sets. If one or both of the data sets fail to meet the test for normality, or if the data sets appear to come from different types of distributions, then nonparametric tests may be the only alternative for the comparison with background. However, for normal data with no outliers or non-detect values, the parametric methods discussed in the next section are somewhat more powerful. Nonparametric tests are discussed in Section 5.3.2.

The relative performance of different testing procedures may be summarized by comparing their *p-values*. The p-value of a statistical test is defined as the smallest value of α at which the null hypothesis would be rejected for the given observations. (The

p-value of the test is sometimes called the critical level, or the significance level, of the test.)

Statistical tests may also be compared based on their *robustness*. Robustness means that the test has good performance for a wide variety of data distributions, and that performance is not unduly affected by outliers. In addition, nonparametric tests used to compare population means and medians generally are unaffected by a reasonable number of non-detect values. There are different circumstances that should be considered:

- ▶ If a parametric test for comparing means is applied to data from a non-normal population and the sample size is large, the parametric test will work well. The central limit theorem ensures that parametric tests for the mean will work because parametric tests for the mean are robust to deviations from normal distributions as long as the sample size is large. Unfortunately, the answer to the question of how large is large enough depends on the nature of the particular distribution. Unless the population distribution is very peculiar, you can safely choose a parametric test for comparing means when there are at least 24 data points in each group.
- ▶ If a nonparametric test for comparing means is applied to data from a normal population and the sample size is large, the nonparametric test will work well. In this case, the p values tend to be a little too large, but the discrepancy is small. In other words, nonparametric tests for comparing means are only slightly less powerful than parametric tests with large samples.

- ▶ If a parametric test is applied to data from a non-normal population and the sample size is small (for example, less than 20 data points), the p value may be inaccurate because the central limit theory does not apply in this case.
- ▶ If a nonparametric test is applied to data from a non-normal population and the sample size is small, the p values tend to be too high. In other words, nonparametric tests may lack statistical power with small samples.

In conclusion, large data sets do not present any problem. In this case the nonparametric tests are powerful and the parametric tests are robust. However, small data sets are challenging. In this case the nonparametric tests are not powerful, and the parametric tests are not robust.

5.3.1 Parametric Tests

Parametric statistical tests, examples of which are listed in the box at the bottom of this page, assume the data have a known distributional form. They may also assume that the data are statistically independent or that there are no spatial trends in the data. Parametric statistical comparison methods, in the context of this guidance, involve comparison of one or more distribution parameters of site samples with corresponding parameters of the background distribution.

Tests for the distribution of the data offer clues on metals detected frequently at higher concentrations. For example, as a general rule, naturally occurring

Parametric Tests		
Test	Purpose	Assumptions
t-test	Test for difference in means	Normality, equal variances
Upper Tolerance Limit (UTL)	Test for outliers	Normality
Extreme Value (Dixon's) Test	Test for one outlier	Normality, not including outlier
Rosner's Test	Test for up to 10 outliers	Normality, sample size 25 or larger
Discordance Test	Test for one outlier	Normality, not including the outlier

aluminum, iron, calcium, and magnesium tend to be normally distributed, while trace metals tend to have lognormal distributions.

Tests of Means

The most common method for background comparisons involves a comparison between means using t-tests or similar parametric methods. If the estimated means do not differ by a statistically significant amount (given a predetermined level of significance such as 0.05), then there is no substantial difference in the mean of the site data as compared to the mean of the background data.

To conduct a t-test, a null hypothesis should first be developed. (See Section 3.1 for developing null hypotheses.) The t-statistic calculated from the data is then compared to a critical value for the test which depends on the level of confidence selected to determine whether or not the null hypothesis should be rejected. Although the t-test is derived based on normality, the conclusion that the data do not follow a normal distribution does not discount the t-test. Generally, the t-test is robust and therefore not sensitive to small deviations from the assumptions of normality.

If the two populations have significantly different variances, the two-sample t-test should not be used for comparing means. Procedures are available to test for equality of variance. Instructions for performing Bartlett's test and Levene's test are presented in EPA QA/G-9, Section 4.5.²

Any t-test should be discussed with a statistician prior to use since there are a number of variations and assumptions that can apply. The Student t-test has good application when comparing background sites to potentially contaminated sites.⁷

Methods such as Cochran's Approximation to the Behrens-Fischer Student t-test may be useful when replicated measurements are available. This statistical comparison method requires that two or more discrete samples be taken at each sampling station. Note that the choice of a specific t-test depends on

site-specific information and other statistical considerations.

Tests of Outliers

There are many parametric tests for outliers, based on deviations from the normal distribution. Three of these tests are explained in detail in EPA QA/G-9,² including Dixon's test, Rosner's test, and the Discordance test shown in the box on the previous page. In addition to these tests, suspected outliers may be identified using a tolerance limit approach. There are parametric and nonparametric forms of tolerance intervals. This section discusses the parametric version.⁸ A nonparametric version of tolerance intervals is presented in Section 5.3.2.

While mean tests explore whether the true means of two populations are significantly different, other tests can be used to indicate whether a single sample is likely to be an outlier in the data set. This type of test can be useful in identifying a "hot spot" that may exceed background, even if the average site concentration does not seem to be different from background. One such test is the tolerance interval. A thorough discussion of normal, Poisson, and nonparametric tolerance limits can be found in Chapter 4 of Gibbons.⁹

A tolerance limit (TL) is a confidence limit on a percentile of the data, rather than a confidence limit on the mean. Tolerance limits provide an interval within which at least a certain proportion of the population lies with a specified probability that the stated interval does indeed "contain" that proportion of the population. An example would be a situation in which you are trying to draw a random sample, and want to know how large the sample size should be so that you can be 95 percent sure that at least 95 percent of the population lies between the smallest and the largest observation in the sample. Similarly, one-sided TLs can be developed. Establishing a TL is recommended for identifying outliers.

For example, a 95 percent one-sided TL for 95 percent coverage represents the value below which 95 percent of the population are expected to fall

(with 95 percent confidence).

5.3.2 Nonparametric Tests

The statistical tests discussed in the previous section rely on the mathematical properties of the population distribution (normal or lognormal) selected for the comparison with background. When the data do not follow the assumed distribution, use of parametric statistical tests may lead to inaccurate comparisons. Additionally, if the data sets contain outliers or non-detect values, an additional level of uncertainty is faced when conducting parametric tests. Since most environmental data sets do contain outliers and non-detect values, it is unlikely that the current widespread use of parametric tests is justified, given that these tests may be adversely affected by outliers and by assumptions made for handling non-detect values.

Nonparametric Tests	
Test	Assumptions
Wilcoxon Rank Sum (WRS)	Both samples are randomly selected from respective populations and mutually independent; distributions are identical (except for possible difference in location parameter).
Gehan Test	Multiple detection limits and non-detect.
Quantile Test	Populations are identical except for differences above a given percentile

Tests that do not assume a specific mathematical form for the underlying distribution are called distribution-free or nonparametric statistical tests. The property of *robustness* is the main advantage of nonparametric statistical tests. Nonparametric tests have good test performance for a wide variety of distributions, and that performance is not unduly affected by outliers.

Nonparametric tests can be used for normal or non-normal data sets. If one or both of the data sets fail to meet the test for normality, or if the data sets appear to come from different types of populations, then

nonparametric tests may be the only alternative for the comparison with background. If the two data sets appear to be from the same family of distributions, use of a specific statistical test that is based on this knowledge is not necessarily required because the nonparametric tests will perform almost as well. However, for normal data with no outliers or non-detect values, the parametric methods discussed in the previous section are somewhat more powerful.

Several nonparametric test procedures, including three listed in the box at left, are available for conducting background comparisons. Nonparametric tests compare the shape and location of the two distributions instead of a statistical parameter (such as mean). Nonparametric tests are currently used by some EPA regions on a case-by-case basis. These methods have varying levels of sensitivity and data requirements and should be considered as the preferred methods whenever data are heavily censored (a high percentage of non-detect values).

Wilcoxon Rank Sum Test for Background Comparisons

The Wilcoxon Rank Sum (WRS)¹⁰ test is an example of a nonparametric test used for determining whether a difference exists between site and background population distributions. The WRS tests whether measurements from one population consistently tend to be larger (or smaller) than those from the other population. This test determines which distribution is higher by comparing the relative ranks of the two data sets when the data from both sources are sorted into a single list. One assumes that any difference between the background and site concentration distributions is due to a shift in the site concentrations to higher values (due to the presence of contamination in addition to background).

Two assumptions underlying this test are: 1) samples from the background and site are independent, identically distributed random samples, and 2) each measurement is independent of every other measurement, regardless of the set of samples from which it came. The test assumes also that the distributions of

the two populations are identical in shape (variance), although the distributions need not be symmetric.

The WRS test has three advantages for background comparisons:

- ▶ The two data sets are not required to be from a known type of distribution. The WRS test does not assume that the data are normally or log-normally distributed, although a normal distribution approximation often is used to determine the critical value for the test for large sample sizes.
- ▶ It allows for non-detect measurements to be present in both data sets (see box below).¹¹ The WRS test can handle a moderate number of non-detect values in either or both data sets by treating them as ties.¹² Theoretically, the WRS test can be used with up to 40 percent or more non-detect measurements in either the background or the site data. If more than 40 percent of the data from either the background or site are non-detect values, the WRS test should not be used.¹³
- ▶ It is robust with respect to outliers because the analysis is conducted in terms of ranks of the data. This limits the influence of outliers because a given data point can be no more extreme than the first or last rank.

Procedures for Non-Detect Values in WRS Test

If there are t non-detect values, they are considered as “ties” and are assigned the average rank for this group. Their average rank is the average of the first t integers, $(t+1)/2$. If more than one detection limit was in use, all observations below the largest detection limit should be treated as non-detect values. Alternatively, the Gehan test may be performed.¹⁴

The WRS test may be applied to either null hypothesis in the two forms of background test discussed in Chapter 3: *no statistically significant difference* or *exceed by more than a substantial difference*. In

either form of background test, the null hypothesis is assumed to be true unless the evidence in the data indicates that it should be rejected in favor of the alternative.

WRS Test Procedure for Background Test Form 1

Null Hypothesis (H_0): The mean of the site distribution is less than or equal to the mean of the background ($\Delta \leq 0$).

Alternative Hypothesis (H_A): The mean of the site distribution exceeds the mean of the background distribution ($\Delta > 0$).

The WRS test for Background Test Form 1 is applied as outlined in the following steps. The lead-contaminated storage yard example from Chapter 3 serves to illustrate the procedure. (Although EPA selected Background Test Form 2 in this example, both forms of the test are evaluated.)

Hypothetical data for the storage yard example is shown in Tables 5.1 and 5.2 for the onsite and background areas, respectively. There is one non-detect measurement (ND) in the data collected on site and five in the background data set. The background non-detects were treated as 0 values when adding S to the background measurements. This is a more conservative approach than using $\frac{1}{2}$ the detection limit or other surrogate or random numbers for the non-detect values.¹⁵ The WRS test is very robust to this small modification as it is unlikely that any reasonable surrogate value will affect significantly the assigned rank of the non-detects in the combined data set.

Table 5.3 demonstrates the WRS test procedure for Background Test Form 1, testing the null hypothesis that there is no statistically significant difference between the site and background distributions. The background measurements ($m = 20$) and the site measurements ($n = 20$) are ranked in a single list in order of increasing size from 1 to N , where $N = m + n = 40$. At the top of the list, all six non-detect values are considered as ties and are assigned an average rank of $3.5 = (6 + 1) \div 2$. (See the box

Data (mg/kg)	Source
ND	Site
34.0	Site
39.5	Site
48.6	Site
54.9	Site
70.9	Site
72.1	Site
81.3	Site
83.2	Site
86.2	Site
88.2	Site
96.1	Site
98.3	Site
104.3	Site
105.6	Site
129.0	Site
139.3	Site
156.9	Site
167.9	Site
208.4	Site

Table 5.1 Site data

Data (mg/kg)	Source
ND	Background
0.1	Background
15.7	Background
46.1	Background
48.1	Background
49.3	Background
53.5	Background
58.0	Background
59.7	Background
68.0	Background
88.5	Background
96.5	Background
115.8	Background
122.9	Background
126.8	Background
147.5	Background

Table 5.2 Background data

Rank	Data (mg/kg)	Source	Ranks for	
			Site	Background
3.5	ND	Site	3.5	
3.5	ND	Background		3.5
3.5	ND	Background		3.5
3.5	ND	Background		3.5
3.5	ND	Background		3.5
3.5	ND	Background		3.5
7	0.1	Background		7
8	15.7	Background		8
9	34.0	Site	9	
10	39.5	Site	10	
11	46.1	Background		11
12	48.1	Background		12
13	48.6	Site	13	
14	49.3	Background		14
15	53.5	Background		15
16	54.9	Site	16	
17	58.0	Background		17
18	59.7	Background		18
19	68.0	Background		19
20	70.9	Site	20	
21	72.1	Site	21	
22	81.3	Site	22	
23	83.2	Site	23	
24	86.2	Site	24	
25	88.2	Site	25	
26	88.5	Background		26
27	96.1	Site	27	
28	96.5	Background		28
29	98.3	Site	29	
30	104.3	Site	30	
31	105.6	Site	31	
32	115.8	Background		32
33	122.9	Background		33
34	126.8	Background		34
35	129.0	Site	35	
36	139.3	Site	36	
37	147.5	Background		37
38	156.9	Site	38	
39	167.9	Site	39	
40	208.4	Site	40	
820		Sum of Ranks	491.5	328.5
			W_s	W_b

Table 5.3 WRS test for Test Form 1
 H_0 : site < background

entitled, "Procedures for Non-Detect Values in WRS Test," on the previous page). The ranks for each area are shown in the two columns at the right of the exhibit. The sum of the ranks of the site measurements ($W_s = 491.5$) and the sum of the ranks of the background measurements ($W_b = 328.5$) are shown at the bottom of the exhibit.¹⁶ The sum of the ranks

of the site measurements ($W_s = 491.5$) is the test statistic used for Background Test Form 1. The sum of the site ranks is used as the test statistic for background test form 1 because EPA is looking for evidence that the site distribution exceeds the background distribution. If W_s is greater than the tabulated critical value for the test, the null hypothesis that there is no significant difference will be rejected.

Most readily available tables for the WRS test only extend up to sample sizes of $n = m = 20$. Critical values for the WRS test when n and m exceed 20 may be calculated from the large sample approximation using this equation:

$$W_{\text{crit}} = m(N+1)/2 + z_{\alpha} [nm(N+1)/12]^{1/2}$$

where $N = n + m$ and z_{α} is the $100(1 - \alpha)^{\text{th}}$ percentile of the standard normal distribution. The first term is the expected value of the sum of ranks W , calculated under the assumption that the null hypothesis is true. The second term is a standard normal variate times the standard deviation of W , under the same assumptions. The first factor in the expectation term m represents the number of ranks that were summed, each having expectation $(N+1)/2$ under the equality assumption included in the null hypothesis.

Table 5.4 shows the critical values for the WRS test for selected values of α for data sets with $n = m = 20$. The critical value for $\alpha = 0.10$ is 458, and the critical value for $\alpha = 0.05$ is 471. Since W_s exceeds the critical values for most commonly used values of α , the null hypothesis is rejected. Hence, the site is distinguishable from background at a confidence

α	Critical Value
0.20	442
0.15	449
0.10	458
0.05	471
0.025	482
0.010	495
0.005	504
0.001	521

Table 5.4 Critical Values for the WRS Test for $n = m = 20$

level of 95 percent. Note that the null hypothesis would not be rejected at $\alpha = 0.01$.

WRS Test Procedure for Background Test Form 2

Null Hypothesis (H_0): The site distribution exceeds the background distribution by more than a substantial difference S ($\Delta > S$).

Alternative Hypothesis (H_A): The site distribution does not exceed the background distribution by more than S ($\Delta \leq S$).

The WRS test for Background Test Form 2 is applied as outlined in the following steps. The lead example will again serve as an illustration of the procedure. In the example from Chapter 3, EPA chose to use Background Test Form 2, with $\alpha = 0.10$ and a substantial difference of $S = 100$ mg/kg. First, the background measurements are adjusted by adding $S = 100$ mg/kg to each measured value. Table 5.5 contains two columns on the right that show the S -adjusted background data for $S = 50$ mg/kg and $S = 100$ mg/kg.

The adjusted background measurements and the measurements from the site in Table 5.5 are ranked in increasing order from 1 to 40. Note that the five adjusted background measurements that were non-detects are tied at 100 mg/kg. They are all assigned the average rank of 16 for that group of tied measurements.

The sum of the ranks of the adjusted measurements from background, $W_b = 544$, is the test statistic for Background Test Form 2. Note that the test statistic for Background Test Form 2 differs from the test statistic for Background Test Form 1. In this case, we are looking for evidence that S plus the background distribution is greater than the site distribution. Earlier, in Background Test Form 1, we were looking for evidence that the site distribution exceeds the (unmodified) background distribution. The critical value for the WRS test (Table 5.4) for $\alpha = 0.10$ is 458. Since W_b is greater than the critical value, the null hypothesis that the site exceeds background by more than a substantial difference of

Rank	Data (mg/kg)	Source	Ranks for	
			Site	Background + 100
1	ND	Site	1	
2	34.0	Site	2	
3	39.5	Site	3	
4	48.6	Site	4	
5	54.9	Site	5	
6	70.9	Site	6	
7	72.1	Site	7	
8	81.3	Site	8	
9	83.2	Site	9	
10	86.2	Site	10	
11	88.2	Site	11	
12	96.1	Site	12	
13	98.3	Site	13	
16	100.0	Background+S		16
16	100.0	Background+S		16
16	100.0	Background+S		16
16	100.0	Background+S		16
16	100.0	Background+S		16
19	100.1	Background+S		19
20	104.3	Site	20	
21	105.6	Site	21	
22	115.7	Background+S		22
23	129.0	Site	23	
24	139.3	Site	24	
25	146.1	Background+S		25
26	148.1	Background+S		26
27	149.3	Background+S		27
28	153.5	Background+S		28
29	156.9	Site	29	
30	158.0	Background+S		30
31	159.7	Background+S		31
32	167.9	Site	32	
33	168.0	Background+S		33
34	188.5	Background+S		34
35	196.5	Background+S		35
36	208.4	Site	36	
37	215.8	Background+S		37
38	222.9	Background+S		38
39	226.8	Background+S		39
40	247.5	Background+S		40
820		Sum of Ranks	276	544
			W_s	W_b

Table 5.5 WRS test for Test Form 2
 H_0 : site > background + 100

100 mg/kg is rejected at the 90 percent confidence level.

Table 5.6 shows the WRS test for the lead example using Background Test Form 2 with a smaller (more conservative) value for a substantial difference, $S =$

Rank	Data (mg/kg)	Source	Ranks for	
			Site	Background + 50
1	ND	Site	1	
2	34.0	Site	2	
3	39.5	Site	3	
4	48.6	Site	4	
7	50.0	Background+S		7
7	50.0	Background+S		7
7	50.0	Background+S		7
7	50.0	Background+S		7
7	50.0	Background+S		7
10	50.1	Background+S		10
11	54.9	Site	11	
12	65.7	Background+S		12
13	70.9	Site	13	
14	72.1	Site	14	
15	81.3	Site	15	
16	83.2	Site	16	
17	86.2	Site	17	
18	88.2	Site	18	
19	96.1	Background+S		19
20	96.1	Site	20	
21	98.1	Background+S		21
22	98.3	Site	22	
23	99.3	Background+S		23
24	103.5	Background+S		24
25	104.3	Site	25	
26	105.6	Site	26	
27	108.0	Background+S		27
28	109.7	Background+S		28
29	118.0	Background+S		29
30	129.0	Site	30	
31	138.5	Background+S		31
32	139.3	Site	32	
33	146.5	Background+S		33
34	156.9	Site	34	
35	165.8	Background+S		35
36	167.9	Site	36	
37	172.9	Background+S		37
38	176.8	Background+S		38
39	197.5	Background+S		39
40	208.4	Site	40	
820		Sum of Ranks	379	441
			W_s	W_b

Table 5.6 WRS test for Test Form 2
 H_0 : site > background + 50

50 mg/kg. The sum of the ranks of the S-adjusted background measurements is $W_b = 441$. After examination of these data, it is clear that the null hypothesis that the site exceeds background by more than 50 mg/kg cannot be rejected at any reasonable level of confidence.

In conclusion, site concentrations in this example are significantly higher than background concentrations. The site distribution may exceed background by 50 mg/kg or more, but it is unlikely that the site distribution is more than 100 mg/kg above background.

Power of the WRS Test

The exact power of the WRS test is difficult to calculate. An approximation by Lehmann¹⁷ is based on the Mann-Whitney form of the WRS test statistic.¹⁸ The Mann-Whitney test statistic is equal to the Wilcoxon rank sum statistic minus its smallest possible value. Thus the Mann-Whitney test statistic is

$$W_{MW} = W_s - n(n+1)/2$$

for Background Test Form 1, and

$$W_{MW}^* = W_b - m(m+1)/2$$

for Background Test Form 2. In each case, the power of the WRS test is calculated using the mean and variance of the approximating normal distribution for the corresponding Mann-Whitney test statistic. Lehmann describes the method for approximating the power of the WRS test used in Background Test Form 1. The U.S. Nuclear Regulatory Commission (NRC)¹⁹ offers a detailed application of the normal approximation for the power of Background Test Form 2 using a different notation for the gray region than is used in this guidance. At the upper end of the gray region, EPA's substantial difference (S) is called by NRC the "Design Concentration Guideline Level" (DCGL). The width of the gray region (EPA's MDD) is defined by NRC as $\Delta = DCGL - LBGR$, where "LBGR" is the lower bound of the gray region. The NRC document also implements Lehmann's power approximation for Background Test Form 1, obtained by letting $LBGR = 0$.²⁰ The NRC document also contains tables for use in evaluating the mean and variance of W_{MW} and W_{MW}^* in Lehmann's approximation for the power of the WRS test. Due to the differences in notation noted above, the NRC tables are tabulated in terms of the design parameter Δ/σ , which corresponds with MDD/σ in this guidance.

Gehan's Form of the WRS Test

The Gehan test is a generalized version of the WRS test.²¹ If there are a large number of non-detect measurements and several different detection levels, Gehan's form of the WRS test is a more powerful test for the background comparison. The Gehan test addresses multiple detection limits using a modified ranking procedure rather than relying on the "all ties get the same rank" approach used in the WRS test.¹⁵ After the modified ranking is completed, the standard WRS test procedure discussed above is applied to determine if the null hypothesis should be rejected. It has been recommended that there should be at least 10 data values in each data set to use this test.

Quantile Test

In many instances, releases of chemicals have impacted only portions of the site. Under such conditions, chemical concentrations in relatively small areas at the site could be elevated relative to the underlying background concentrations. As a result, only a small portion in the upper tail of the distribution of site measurements would be expected to be shifted to higher concentrations than the distribution of background measurements. The quantile test is a nonparametric test that is designed to compare the upper tails of the distributions. The quantile test may detect differences that are not detected by the WRS test. The quantile test is described in detail in Chapter 7 of Gilbert and Simpson.¹

Walsh's Tests for Outliers

Walsh's test is a nonparametric test for determining the presence of outliers in either the background or onsite data sets. This test was developed to detect up to a specified number of outliers, r . The test requires large sample sizes ($n > 60$ for $\alpha \equiv 0.10$; and $n > 220$ for $\alpha \equiv 0.05$). Procedures for conducting this test is discussed in Section 4.4 of EPA QA/G-9.²

Nonparametric Tolerance Intervals

The parametric tolerance intervals presented in

Section 5.2 are derived based on the assumption of a normal distribution. If the data are not normal and are not easily transformed to normality, then non-parametric tolerance intervals may be calculated for the background distribution to provide a tolerance level for screening site data. A readable discussion of the use of nonparametric tolerance intervals is provided by Glick.²²

5.4 Hypothesis Testing

Hypothesis testing was discussed in detail in Section 3. Here, some of this information is reviewed, and additional aspects of such testing are discussed. The emphasis is on classical methods for testing hypotheses, including parametric and nonparametric methods. The Bayesian approach is an alternative to classical methods for hypothesis testing, but is not included in the discussion. The Bayesian approach for comparing the means of two populations is discussed by many authors, including Box and Tiao.²³ Bayesian methods permit the incorporation of prior knowledge and provide the ability to update information as results come in from successive rounds of sampling.

5.4.1 Initial Considerations

For CERCLA sites, use of a null hypothesis and alternative hypothesis is recommended when comparing data sets from potentially impacted areas with background data. For example, a null hypothesis could be “there is no difference between the mean contaminant concentration in samples from potentially impacted areas and background data sets.” The alternative hypothesis would be “there is a difference between mean contaminant concentration in samples from potentially impacted areas and background data sets.” To conduct the comparison, parametric or non-parametric statistical tests are recommended. Use of parametric comparison methods like t-tests or ANOVA may require normalization of data, such as the conversion to a log scale. Depending upon the data and other site-specific considerations, statistical analysis should involve one or a combination of the following methods:

- ▶ A preliminary descriptive analysis involving the comparison of median, mean, and upper range concentrations between sample sets considered site-related and background;
- ▶ Parametric statistical comparison methods involving the comparison of one or more parameters of the distribution of site samples with corresponding parameters of the (assumed or sampled) background distribution, such as Gosset’s Student t-test or Cochran’s Approximation to the Behrens-Fischer Student t-test; or
- ▶ Nonparametric tests, such as the Wilcoxon Rank Sum test (on a case-by-case basis).

Once a test has been selected, the assessor should consider several questions:

- ▶ *What should the null and alternative hypotheses be? What are we testing? What are we trying to support or reject about the site and background?*
- ▶ *Should the test be one-tailed or two-tailed? Should we ask whether the site and background are from the same population, or should we focus on whether one is more contaminated than the other?*
- ▶ *What confidence level should be used? At what “cut-off” point do we accept or reject the hypothesis?*

5.4.2 Examples

It may be easiest to explore these questions by using an example. Suppose we have an area that meets our criteria for local background. The data from this area for Chemical X (mg/kg) are as follows:

66 67 68 68 69 69 69
70 70 70 71 71 71 72
72 72 72 73 74 74 75

These data were collected randomly and are normally distributed. There are 21 measurements ($n = 21$), with an average of 70.6 mg/kg and a standard

deviation of 2.37 mg/kg.

We also have data from an onsite process area. These data for Chemical X (mg/kg) are as follows:

62 63 64 65 66 68 68
69 69 70 71 71 72 72
72 73 74 75 77 78 80

These data were collected randomly and are normally distributed. There are 21 measurements ($n = 21$), with an average of 70.4 mg/kg and a standard deviation of 4.86 mg/kg.

The background and onsite areas appear to be similar, but some of the onsite data exceed the background data. We would like to be able to state with a given level of confidence whether the data are essentially from the same population, or not. If we use the t test to compare the true means of these data sets, we could test the hypothesis that the background mean and the site mean are essentially equal (H_0 , the null hypothesis). If H_0 is not true, then we would support the alternative hypothesis that the means are not equal. This is a two-tailed test, because H_0 could be rejected if the site mean is greater than the background mean or if the site mean is less than the background mean.

Example 1: $H_0: \mu_s = \mu_b$
 $H_A: \mu_s \neq \mu_b$

(Note that this is a two-tailed version of Test Form 1.) Using the equations in EPA QA/G-9, for t, we find that $t = 0.1693$.² At 40 degrees of freedom, for a two-tailed test, our t falls below the t of 0.681, where $\alpha = 0.5$.²⁴ Therefore, if we had chosen an α of 0.01 (99 percent confidence), 0.05 (95 percent confidence), or 0.1 (90 percent confidence), we would not reject our null hypothesis. Only if we were testing at less than 50 percent confidence would we reject H_0 .

When using Test Form 1, the higher the confidence limit, the more likely this test is to find that the site is from the same population as background. Choosing the rejection range for the hypothesis involves balancing both kinds of error. *In general, EPA recom-*

mends a minimum confidence limit of 80 percent and a maximum confidence limit of 95 percent.

Suppose we want to compare our background data set with another onsite process area. The data for Chemical X (mg/kg) are as follows:

56 58 60 62 66 67 68
70 72 73 75 76 81 82
84 85 87 90 91 92 103

These data were collected randomly and are normally distributed. There are 21 measurements ($n = 21$), with an average of 76.1 mg/kg and a standard deviation of 12.68 mg/kg.

Is this area significantly different from background? The arithmetic mean is 76.1 mg/kg, compared to the background mean of 70.6 mg/kg. But is this difference truly significant? After all, the mean of the first process area, 70.4 mg/kg, was different from the background mean. According to the t test, however, we did not find the difference of 0.2 mg/kg to be significant at the 80-99 percent confidence levels. What about the second process area?

Suppose we decide that we are really interested in whether the site is “dirty” (above background). Instead of a 2-tailed test, we could perform a 1-tailed test:

Example 2: $H_0: \mu_s > \mu_b$
 $H_A: \mu_s \leq \mu_b$

(Note that this is Test Form 2 with $S = 0$.) This test is 1-tailed because the rejection region is only on one side of the distribution; that is, we are only interested in whether the site is greater than the background.

To use the normal distribution theory correctly, for a 1-tailed t test, with 40 degrees of freedom, the t of -1.95 is calculated for the background mean minus the site mean. This t falls between the 95 percent and 97.5 percent confidence levels. If we were testing at 80 percent or 95 percent confidence, we would reject H_0 and find that the site is less than or

equal to background—in other words, “clean.” At 99 percent confidence, H_0 could not be rejected. In this case, therefore, a lower confidence limit seems to *increase* the chances of finding that the site is clean, where in our earlier tests we found that a lower confidence limit *decreased* the chances of considering the site clean. Why is this?

The difference is in the setup of the hypotheses. In the first case (example 1), the null hypothesis was that the site and background were from the same population (the site was clean). In the later case, the null hypothesis was that the site mean exceeded the background mean (the site is dirty). In essence, we have shifted the burden of proof. If we are really interested in whether the site is dirty (greater than background), then our last test could have looked at these hypotheses:

Example 3:

$$H_0: \mu_s \leq \mu_b$$

$$H_A: \mu_s > \mu_b$$

(Note that this is a one-tailed version of Test Form 1.) Using the site mean minus the background mean for this test, we derive a t of 1.95. At the 80 percent confidence level, we would reject H_0 and find that the site is dirty. At the 95 percent confidence level and above, we would accept H_0 and find that the site is clean because the data are insufficient to support this higher level of confidence demanded for rejection. Once again, with Test Form 1, a *lower* confidence level results in a *more conservative* approach to environmental protection.

There is another problem, besides burden-of-proof, with Example 2. As discussed in Chapter 3, the null hypothesis that there *is* a substantial difference (Test Form 2, $\Delta > 0$) should only be tested if some minimal difference (S) is specified. This is because the null hypothesis $H_0: \Delta > 0$ (i.e., $H_0: \mu_s > \mu_b$) will be rejected only if the site mean is significantly below the background mean. In a more typical case, the site mean may be almost equal to or slightly below the background mean, and the null hypothesis will only be rejected when a large number of samples is collected to reduce the uncertainty to below the magnitude of the difference in means.

Essentially, Test Form 1 uses the default assumption that the site is “clean” unless it can be shown otherwise; Test Form 2 uses the default assumption that the site is “dirty” unless it can be shown otherwise.

5.4.3 Conclusions

Now we return to our original three questions. Table 5.7 also summarizes this information.

- ▶ *What should the null and alternative hypotheses be?*
- ▶ *Should the test be one-tailed or two-tailed?*
- ▶ *What confidence level should be used?*

To determine whether the site and background are from the same population, these hypotheses can be used in a two-tailed Test Form 1:

$$H_0: \mu_s = \mu_b$$

$$H_A: \mu_s \neq \mu_b$$

For this test, the confidence level should be at least 80 percent but no more than 95 percent. For a more conservative test, use the lower end of the confidence range.

To determine whether the site is significantly greater than background, these hypotheses can be used in a one-tailed Test Form 1:

$$H_0: \mu_s \leq \mu_b$$

$$H_A: \mu_s > \mu_b$$

For this test, the confidence level should be at least 80 percent; for a more conservative test, use the lower end of the confidence range and require adequate power.

If testing the hypotheses in reverse—Test Form 2—to show whether the site is greater than background + S , use a higher confidence level, such as 95 percent, and specify a substantial difference S . (See Appendix A for guidance on choosing S .) To determine whether the site exceeds background by more than S , these hypotheses can be used in a one-tailed Test Form 2:

$$H_0: \mu_s \geq \mu_b + S$$

$$H_A: \mu_s < \mu_b + S$$

For this test, the confidence level should be at least 80 percent; for a more conservative test, use higher levels of the confidence range.

What to test:	H ₀	H _A	Recommended alpha	Rejection criteria
H ₀ : site and background are from the same population; vs. H _A : site and background are from different populations (Two-tailed, Test Form 1)	$\mu_s = \mu_b$	$\mu_s \neq \mu_b$	80-95% confidence ($\alpha = 0.2$ to 0.05) [More conservative: $\alpha = 0.2$]	For 2-sided t test, e.g., reject H ₀ if $t > t_{\alpha/2}$ or if $t < -t_{\alpha/2}$
H ₀ : site is less than or from the same population as background; vs. H _A : site is greater than background (One-tailed, Test Form 1)	$\mu_s \leq \mu_b$	$\mu_s > \mu_b$	80-95% confidence ($\alpha = 0.2$ to 0.05) [More conservative: $\alpha = 0.2$]	For 1-sided t test, e.g., reject H ₀ if $ t > t_{\alpha}$ For 1-sided t test, e.g., reject H ₀ if $t > t_{\alpha}^*$
H ₀ : site is greater than background + S; vs. H _A : site is less than or from the same population as background + S (One-tailed, Test Form 2)	$\mu_s \geq \mu_b + S$	$\mu_s < \mu_b + S$	80-95% confidence ($\alpha = 0.2$ to 0.05) [More conservative: $\alpha = 0.05$]	For 1-sided t test, e.g., reject H ₀ if $ t > t_{\alpha}$ For 1-sided t test, e.g., reject H ₀ if $t < -t_{\alpha}^*$

* Assuming the test statistic, t, is calculated using site mean minus background mean (or background mean + S, for Test Form 2) in the numerator.

Table 5.7 Summary of hypothesis tests

CHAPTER NOTES

1. Gilbert, R.O. & J.C. Simpson. June 1994. *Statistical Methods for Evaluating the Attainment of Cleanup Standards, Volume 3*. EPA 230-R-94-004.
2. U.S. Environmental Protection Agency (EPA). July 2000. *Guidance for Data Quality Assessment: Practical Methods for Data Analysis, EPA QA/G-9, QA00 Version*. Quality Assurance Management Staff, Washington, DC, EPA 600-R-96-084. Available at http://www.epa.gov/quality/qa_docs.html.
 - ▶ See Section 1.3.1 for guidance on “authoritative samples.”
 - ▶ See Section 3.3.1.1 for guidance on how to calculate t.
3. Cressie, N. 1991. *Statistics for Spatial Data*, New York: John Wiley & Sons. See Section 3.2.
4. A variety of methods for addressing non-detects are presented in Section 4.7 of EPA QA/G-9, *op. cit.* Simulation results using $L/\sqrt{2}$ are reported by Hornung, R.W. and Reed, L.D. “Estimation of average concentration in the presence of nondetectable values,” *Applied Occupational and Environmental Hygiene*, 5(1), p.45-51, January, 1990.
5. Although the use of t-tests after logarithmic transformation to approximate normality has a long history, some authors have recommended against using t-tests on log transformed data. See A. K. Singh, A. Singh, and M. Engelhardt, *The Lognormal Distribution in Environmental Applications*, EPA Technology Support Center Issue, U.S. EPA Office of Research and Development, National Exposure Research Laboratory, Las Vegas, NV, EPA/600/R-97/006, December 1997. The authors note that the H-statistic based on the Upper Confidence Limit for the mean of a lognormal has erratic performance with sample sizes smaller than 30.
6. When using parametric statistical tests, a limit of 15 percent non-detect measurements in either data set is suggested in the Navy’s *Procedural Guidance for Statistically Analyzing Environmental Background Data*. Nonparametric statistical methods are recommended if this limit is exceeded.
7. Michigan Department of Environmental Quality Waste Management Division. April 1994. *Guidance Document: Verification of Soil Remediation*. Revision 1. <http://www.deq.state.mi.us/wmd/docs/vsr.html>.
8. Devore, J.L. 2000. *Probability and Statistics for Engineering and the Sciences*, 5th Ed., Duxbury Press, Pacific Grove, California.
9. Gibbons, R.D. 1994. *Statistical Methods for Groundwater Monitoring*, John Wiley & Sons, Inc.
10. The WRS test is also called the Mann-Whitney test, which is mathematically equivalent to the WRS test. Sometimes, the combined name is used: Wilcoxon-Mann-Whitney test.
11. In general, the use of “non-detect” values in data reporting is not recommended. Wherever possible, the actual result of a measurement, together with its uncertainty, should be reported. Estimated concentrations should be reported for data below the detection limit, even if these estimates are negative, because their relative magnitude compared to the rest of the data is of importance.

12. The Gehan test discussed in the next section should be considered if there are many non-detect values with different detection levels.
13. A limit of 50 percent non-detect values is suggested in the Navy's *Procedural Guidance for Statistically Analyzing Environmental Background Data*. A more conservative limit of 40 percent non-detect values is recommended in the *Multi-Agency Radiation Survey and Site Investigation Manual (MARSSIM)*.
14. As a third alternative, Morris DeGroot (1986) recommends that in the case of ties, the WRS could "be carried out twice. In the first test, the smaller ranks in each group of tied observations should be assigned to the x 's and the larger ranks should be assigned to the y 's. In the second test, these assignments should be reversed. If the decision to accept or reject H_0 is different for the two assignments, or if the calculated tail areas are very different, the data must be regarded as inconclusive." *Probability and Statistics*, 2nd Edition. Addison-Wesley, Reading, MA, pp. 517 and 584.
15. Additional information on the treatment of non-detects is given by Newman, *et al.*, (1989) "Estimating the mean and variance for environmental samples with below detection limit observations," *Water Resources Bulletin* 25(4): 905-916.
16. Critical values for the WRS test are available in many published texts and reference books. Two sources are W.J. Conover (1980) *Practical Nonparametric Statistics*, 2nd Ed., John Wiley & Sons, New York; and Zwillinger, D. and S. Kokoska (2000) *CRC Standard Probability and Statistics Tables and Formulae*, Chapman and Hall/CRC Press, Boca Raton, Florida.
17. Lehmann, E.L. and H.J.M. D'Abbrera. 1998. *Nonparametrics: Statistical Methods Based on Ranks*, revised 1st Ed., Prentice Hall, New Jersey. Section 2.3.
18. Many tables of critical values for the WRS test, including Lehmann's, are expressed in terms of the Mann-Whitney form of the test statistic. Care should be exercised in selecting an appropriate table when results of the WRS test are evaluated.
19. Gogolak, C.V., G.E. Powers, and A.M. Huffert. *A Nonparametric Statistical Methodology for the Design and Analysis of Final Status Decommissioning Surveys*, Office of Nuclear Regulatory Research, U.S. Nuclear Regulatory Commission (NRC), NUREG-1505. June 1998 (Rev. 1). See "Scenario A" in Section 10.4.
20. See "Scenario B" in Section 10.5 of Gogolak *et al.*, *op cit.* This scenario, with $LBGR > 0$, introduces a new set of WRS tests not discussed in this guidance. These tests have null hypotheses similar to the null hypothesis used in Background Test Form 1, but with a less stringent definition of "clean" ($H_0: \Delta < LBGR$). The test statistic for rejecting these null hypotheses is obtained by ranking in one column the site measurements, which are first adjusted by subtracting the LBGR from each measurement, together with the (unadjusted) background measurements. The WRS test statistic is defined as the sum of the ranks of the adjusted site measurements.
21. Detailed instructions for conducting the Gehan test are found in *Handbook for Statistical Analysis of Environmental Background Data*, Naval Facilities Engineering Command, SWDIV and EFA WEST, July, 1999, Section 3.7. Also see Appendix C.1 of the draft Supplemental Guidance to RAGS: Region 4 Bulletins – Addition #1, Statistical Tests for Background Comparison at Hazardous Waste Sites, U.S.

- EPA, Waste Management Division, Region 4, Office of Technical Services, November, 1998. An earlier use is found in Millard, W.P., and S.J. Deverel. 1988. "Non-Parametric Statistical Methods for Comparing Two Sites Based on Data with Multiple Non-Detect Limits." *Water Resources Research*, 24:12, pp. 2087-2098.
22. Glick, N. "Breaking records and breaking boards", *American Mathematical Monthly*, Vol. 85, No. 1, pp. 2-26, 1978.
23. *Bayesian Inference in Statistical Analysis*, G.E.P. Box and G.C. Tiao, Addison-Wesley, Reading, MA, 1973.
24. In this context, *degrees of freedom* ($n - 1$) is the number of independent observations ("n") minus the number of independent parameters estimated in computing the variation. The shape of the t-distribution curve depends upon the number of degrees of freedom. Distributions with fewer degrees of freedom have heavier tails.

APPENDIX A

SUPPLEMENTAL INFORMATION FOR DETERMINING “SUBSTANTIAL DIFFERENCE”

“Substantial difference” (S) is the difference in mean concentration in potentially impacted areas over background levels that presents a “substantial risk.”

In situations where regulatory requirements indicate that contamination at or below background concentrations presents an unacceptably high risk, it may not be possible to define a reasonable level for a substantial difference. However, background analysis is important in situations where background chemicals occur at concentrations above risk-based criteria, and the statistical methods presented in this guidance are useful tools for background analysis in these situations.

This guidance does not establish a value for S because it should be developed within the Quality Assurance Project Plan on a case-by-case basis as part of the planning process for site investigations.¹ This Appendix does not establish policy, and is provided as supplemental information on the statistical considerations that support a selection for S .

A.1 Precedents for Selecting a Background Test Form

Hypothesis testing is used to make decisions under conditions of uncertainty. The DQO process provides a way for decision makers to determine the requirements of the test, based on evaluation of the consequences of making a Type 1 error (α) or a Type 2 error (β). Statisticians involved in developing the theory of hypothesis testing have noted the

asymmetry between these two types of errors:

The justification for fixing the Type 1 error to be α (usually small and often taken as .05 or .01) seems to arise from those testing situations where the two hypotheses are formulated in such a way that one type of error is more serious than the other. The hypotheses are stated so that the Type 1 error is more serious, and hence one wants to be certain that it is small.²

These opinions are echoed by Bickel and Doksum,³ who use the symbol H for the null hypothesis (we have used H_0) and K for the alternative (our H_A):

Even when we leave the area of scientific research the relative importance of the errors we commit in hypothesis testing is frequently not the same. ... There is a general convention that, if the labeling of H and K is free, the label H is assigned so that type 1 error is the most important to the experimenter.

These opinions relate to the choice between two complementary hypothesis tests with the difference being the reversal of the null and the alternative. This is a burden-of-proof issue. The comparison of the two background test forms also involves selection of an appropriate value for a “substantial difference.” It is also important to distinguish between the value that characterizes a “substantial difference over background” and the appropriate risk-based “action level” for the chemical of concern.

Existing EPA guidance in the data quality objectives (DQO) process for choosing the null hypothesis has focused on the burden-of-proof, when the contaminant concentration is to be compared to a fixed, risk-based action level, L . The choice of test forms for this type of decision includes either

$$a) H_0: X < L \text{ vs. } H_A: X > L$$

or

$$b) H_0: X > L \text{ vs. } H_A: X < L,$$

where X represents the parameter of interest for the distribution of contaminant concentrations in the potentially impacted areas. Hypothesis test a compares the site concentrations to the action level using a null hypothesis that the site does not exceed the action level and an alternative hypothesis that the site exceeds the action level. Hypothesis test b is the opposite of test a , using a null hypothesis, the site exceeds the action level. Background issues are not addressed directly in this framework.

One way to address background comparisons is to reformulate the hypotheses using the difference (delta— Δ) between the distribution of contaminant concentrations and background:

$$a') H_0: \Delta < S \text{ vs. } H_A: \Delta > S$$

and

$$b') H_0: \Delta > S \text{ vs. } H_A: \Delta < S.$$

In hypothesis tests a' and b' , concentrations in potentially impacted areas and in background locations are compared to determine if there is or is not a substantial difference between the two areas. Test a' uses the null hypothesis that the site does not exceed background by more than a substantial difference, while the opposite test b' uses the null hypothesis that the site exceeds background by more than a substantial difference (S). Approaches for selecting a value for S are addressed in the following section. Note that Test Form b' is the one discussed in Section 3.1.2 (Background Test Form 2).

Background Test Form 1 focuses interest on comparisons using a “substantial” difference of $S = 0$. In this case, the two alternative tests are

$$a'') H_0: \Delta < 0 \text{ vs. } H_A: \Delta > 0$$

and

$$b'') H_0: \Delta > 0 \text{ vs. } H_A: \Delta < 0.$$

Background Test Form 1 (Section 3.1.1) is identical with test a'' . This discussion demonstrates that the two background tests addressed in this paper are not opposite forms of the same test in the same sense that tests a and b are opposite forms of the same test with the same threshold. Since the guidance reviewed in this section compares opposite forms of tests with the same action level, the guidance does not contain a direct recommendation for choosing

The two background test forms differ both in terms of burden of proof and in the choice of a substantial difference:

- ▶ Test Form 1 uses a conservative value for a substantial difference of $S = 0$, but relaxes the burden of proof by selecting the null hypothesis that there is no statistically significant difference.
 - ▶ Test Form 2 requires a stricter burden of proof, but permits a larger value for a substantial difference.
-

between Test Forms 1 and 2. Distinguishing characteristics are listed in the box below.

EPA QA/G-9⁴ (Section 1.2) provides the following guidance on the selection of an appropriate null hypothesis in a choice between Test Forms a and b :

The decision on what should constitute the null hypothesis and what should be the alternative is sometimes difficult to ascertain. In many cases, this problem does not arise because the null and alternative hypotheses are determined by specific regulation. However, when the null hypothesis is not specified by regulation, it is necessary to make this determination. The test of hypothesis procedure prescribes that the null hypothesis is only rejected in favor of the alternative, provided there is overwhelming evidence from the data that the null hypothesis is false. In other words, the null hypothesis is considered to

be true unless the data show conclusively that this is not so. Therefore it is sometimes useful to choose the null and alternative hypotheses in light of the consequences of possibly making an incorrect decision between the null and alternative hypotheses. The true condition that occurs with the more severe decision error (not what would be decided in error based on the data) should be defined as the null hypothesis. For example, consider the two decision errors: "decide a company does not comply with environmental regulations when it truly does" and "decide a company does comply with environmental regulations when it truly does not." If the first decision error is considered [the] more severe decision error, then the true condition of this error, "the company does comply with the regulations" should be defined as the null hypothesis. If the second decision error is considered the more severe decision error, then the true condition of this error, "the company does not comply with the regulations" should be defined as the null hypothesis.

For background comparisons, that guidance may be extrapolated. When deciding between Test Forms *a* and *b*", there are two possible decision errors:

- (i) decide the site exceeds background when it truly does not; and
- (ii) decide the site does not exceed background when it truly does.

Decision error (i) occurs when a "clean" site is wrongly rejected. If decision error (i) is more serious than decision error (ii), and if the choice is between tests *a*" and *b*" with a substantial difference of 0, then Background Test Form 1 (*a*") should be selected.

When deciding between Test Forms *a'* and *b'*, there are two possible decision errors:

- (i) decide the site exceeds background + S when it truly does not; and

- (ii) decide the site does not exceed background + S when it truly does.

Decision error (ii) occurs when a truly contaminated site goes undetected. If decision error (ii) is considered more serious than error (i) and the choice is between tests *a*" and *b*" with a substantial difference of S, then Background Test Form 2 should be selected. Note that this logic does not provide a direct comparison of the two forms of background tests considered here, but does indicate situations when Test Forms 1 or 2 may be recommended over their respective opposites.

Chapter 6 of EPA QA/G-4⁵ is succinct and definitive for deciding between Test Form *a* and *b*:

"Define the null hypothesis (baseline condition) and the alternative hypothesis and assign the terms "false positive" and "false negative" to the appropriate decision error.

"In problems that concern regulatory compliance, human health, or ecological risk, the decision error that has the most adverse potential consequences should be defined as the null hypothesis (baseline condition). In statistical hypothesis testing, the data must conclusively demonstrate that the null hypothesis is false. That is, the data must provide enough information to authoritatively reject the null hypothesis (reject the baseline condition) in favor of the alternative. Therefore, by setting the null hypothesis equal to the true state of nature that exists when the more severe decision error occurs, the decision maker guards against making the more severe decision error by placing the burden of proof on demonstrating that the most adverse consequences will not be likely to occur."

This suggests that environmental concerns are not like the jury trial process, and that the "innocent until proven guilty" assumption is an environmentally risky approach. From this viewpoint, a more protective approach would be to presume guilt, and demand proof of innocence: "guilty until proven

innocent.” Remember that this comparison assumes that opposite forms of the same test (*a* and *b*) are being evaluated. Extrapolation of this logic to the background problem would indicate that Test Form 2 is preferred over its true opposite, but Test Form 1 is not preferred over its opposite.

EPA guidance⁶ adopts a conservative approach by stating that when the results of the investigation are uncertain, erroneously concluding that the sample area does not attain the cleanup standard is preferable to concluding that the sample area attains the cleanup standard when it actually may not. Again the recommended approach favors protection of human health and the environment.

A.2 Options for Establishing the Value of a Substantial Difference

Selection of an appropriate value to represent a substantial difference when testing for differences between concentrations in potentially impacted areas and background areas depends on the intended application of the test and a variety of factors. These factors include site and background variability and appropriate cleanup goals.

The term “substantial difference” (*S*) was defined at the beginning of this Appendix as the difference in mean concentration that presents a “substantial risk.” Alternatively, *S* may represent a selected “not-to-exceed” action level that is appropriate for the decision at hand. The application of either test form for a background comparison requires that an upper bound be established for the magnitude of the difference before the site is determined to exceed background. When using Test Form 1, the power of the test is specified at the right edge of the gray region, which has a width equal to the minimum detectable difference (MDD). In this case, the value of *S* serves as an upper bound for the width of the gray region. When using Test Form 2, *S* is explicitly incorporated into the test procedure. *S* is measured in concentration units above the mean background concentration. The decision to use a specific value for a substantial difference may be based on direct

risk assessment, a generic regulatory value, or other level selected to reflect site-specific conditions.

Background comparisons may be conducted at various stages of site characterization and remediation cycle. In the characterization stage, areas with some likelihood of contamination may be compared to background areas to determine if contamination is present in excess of background levels. For example, the goal at this stage may be to determine the areal extent of contamination on a large site. The site is divided into sub-units that are compared to background to determine if contamination is present in the sub-unit. At this stage, Background Test Form 1 is useful for determining if the difference between the site mean and the background mean is significantly greater than zero. An upper bound for the MDD of the test is set by determining a value of the substantial difference *S* which will represent a threshold value for identifying possibly contaminated sub-units.

Later in the site evaluation process, background comparisons may be used to determine if a sub-unit with known contamination has been sufficiently remediated. At this stage, Background Test Form 2 is useful to demonstrate that the remediation was successful. If the goal of the remediation is to reduce contamination to near-background levels, than an appropriate value of *S* is selected that will represent the maximum amount by which a remediated sub-unit may exceed background.

A.2.1 Proportion of Mean Background Concentration

One choice for selecting a value of *S* is to use a specified proportion of typical mean background concentrations for the contaminant of concern:

$$S = rM_b$$

where M_b is the mean background concentration and *r* is the specified proportion. This choice may be appropriate for determining if contamination exists in a sub-unit, or if a sub-unit has been remediated successfully. There is no theoretical reason for restric-

ting r to proportions less than 1, if background concentrations are far below the level that presents a substantial risk. Values of r near 1 may require a high number of samples, because the MDD for the test should be less than S .

The required sample size is determined by MDD/σ , where σ is the standard deviation of the concentrations in potentially impacted areas. Even if the area has little or no contamination, then σ will be approximately as large as the background standard deviation, which is usually at least as large as the background mean. Hence, if r is less than 1, then it is very likely that MDD/σ also is less than 1. If there is contamination in the potentially impacted area, then MDD/σ will be much less than 1.

A.2.2 A Selected Percentile of the Background Distribution

Due to the high variability in background concentrations of many chemicals, defining S as a fraction of the mean background concentration may not be appropriate. Another choice for a value to represent a substantial difference is to use a specified percentile of the distribution of background concentrations for the contaminant of concern:

$$S = (B_p - M_b)$$

where B_p is the p^{th} percentile of the background distribution and M_b is the mean background concentration. Values of p less than 0.85 may require a high number of samples, because the MDD for the test should be less than S . This is because the 85th percentile is approximately 1 standard deviation above the background mean. When there is little or no contamination on the site, S is approximately equal to σ , and hence, MDD/σ usually will be near 1. If there is contamination, then MDD/σ will be much less than 1.

A.2.3 Proportion of Background Variability

A third choice for selecting a value to represent a substantial difference is to use a specified proportion of variance of background concentrations for

the contaminant of concern:

$$S = r\sigma_b$$

where σ_b is the standard deviation of background concentrations and r is the specified proportion. This choice for a substantial difference is closely related to the use of a percentile of the background distribution discussed in Section A.3.2.

Areas with relatively high mean background concentrations generally also have high variance of background. Values of r less than 1 may require a high number of samples, for the reasons noted in Section A.2.2.

A.2.4 Proportion of Preliminary Remediation Goal

The concept of calculating risk-based soil concentrations to serve as reference points for establishing site-specific cleanup levels was introduced in RAGS.⁷ If a preliminary remediation goal (PRG) is available for the contaminant of concern, the value of S may be based on a proportion of the PRG:

$$S = r \cdot \text{PRG}$$

A proportion less than 1 may be required, because the total risk will be the sum of the incremental risk due to S plus the risk due to background concentrations of the contaminant. If the PRG is less than the mean or standard deviation of background, a high number of samples may be required for conclusive test results.

A.2.5 Proportion of Soil Screening Level

If a PRG is not available for the contaminant of concern, a risk-based value of S may be based on the soil screening level (SSL) for the contaminant.⁸

$$S = r \cdot \text{SSL}$$

SSLs are based on a 10^{-6} individual risk for carcinogens and a hazard quotient of 1 for noncarcinogens. SSLs were established to identify the lower bound

of the range of risks of interest in decision making, and are not cleanup goals. SSL target risks should be adjusted to reflect established cleanup level targets. Again, a proportion less than 1 may be required, because the total individual risk will be the sum of the incremental risk due to S plus the risk due to background concentrations of the contaminant. If the (adjusted) SSL is less than the mean or standard deviation of background, a high number of samples may be required for the background comparison.

A.3 Statistical Tests and Confidence Intervals for Background Comparisons

This section provides supplementary material on the use of hypothesis tests and confidence intervals for conducting background comparisons. The science of statistics is often divided into two parts: estimation theory and hypothesis testing. Estimation theory includes the calculation of confidence intervals as estimates for population parameters, while hypothesis testing focuses on the use of statistical tests to accept or reject hypotheses concerning these parameters. Although only the use of hypothesis tests has been discussed in the main text, the one-to-one correspondence between hypothesis tests for Δ conducted at level α and the estimated 100(1- α) percent confidence interval for Δ permits the use of either method to conduct a background comparison. While the emphasis of this section is technical in nature, mathematical proofs of results have been omitted.

When using Test Form 1, a one-sided, level- α hypothesis test of the null hypothesis $\Delta \leq 0$ will only reject the null hypothesis if we conclude that Δ is significantly greater than zero by comparing the test statistic to the tabulated critical value. The critical value is selected to ensure that the probability the test statistic will exceed the critical value by chance alone is less than α . A similar conclusion is reached when the lower bound of the one-sided, 100(1- α) percent confidence interval for Δ is greater than zero. There are two ways to reach the same conclusion that Δ is significantly greater than zero. A two-

sided confidence interval for Δ is often more useful than a one-sided confidence interval to summarize the information about Δ that is contained in the data. In this case, a two-sided, 100(1- α) percent confidence interval for Δ will correspond to a one-sided, level- $\alpha/2$ hypothesis test for Δ .

A.3.1 Comparisons Based on the t-Test

Background comparisons based on the t-test rely on the assumption of a normal distribution for the data, or for a transformation of the data. Hypotheses are tested using the t-statistic, which follows the Student t-distribution. Similar results are obtained by estimating a confidence interval for $\Delta = \mu_y - \mu_x$, where μ_y is the mean concentration in the potentially impacted areas and μ_x is the mean background concentration.

NORMAL THEORY, CASE 1: EQUAL BUT UNKNOWN VARIANCES⁹

For simplicity, we first assume that the site data (Y_1, \dots, Y_n) and background data (X_1, \dots, X_m) are independent random samples from normal distributions with the same variance, σ^2 , but with different means, μ_y and μ_x , respectively:

$$Y_j \sim N [\mu_y, \sigma^2]$$

and

$$X_j \sim N [\mu_x, \sigma^2].$$

In this case, the test statistic for the two-sample t-test is based on the difference in the estimated means, M_y and M_x , where

$$M_y = \Sigma Y_j / n \sim N [\mu_y, \sigma^2/n]$$

and

$$M_x = \Sigma X_j / m \sim N [\mu_x, \sigma^2/m].$$

A pooled estimate for σ^2 , the common variance of the distributions, is

$$s_p^2 = [\Sigma (Y_j - M_y)^2 + \Sigma (X_j - M_x)^2] / (n + m - 2).$$

The test statistic for conducting a t-test using Background Test Form 1 is

$$t_1 = (M_y - M_x) / s^*$$

where

$$s^* = s_p(1/n + 1/m)^{1/2}.$$

In Background Test Form 1, the test statistic t_1 has the standardized Student-t distribution with $n+m-2$ degrees of freedom if $\mu_y = \mu_x$ ($\Delta = 0$). Let $t_{1-\alpha}$ represent the $100(1-\alpha)$ th quantile of the Student t-distribution with $n+m-2$ degrees of freedom. The value $t_{1-\alpha}$ is the critical value for the test. If the test statistic t_1 exceeds the critical value $t_{1-\alpha}$, the null hypothesis in Background Test Form 1 ($H_0: \Delta \leq 0$) may be rejected with $100(1-\alpha)$ percent confidence.

The test statistic for conducting a two-sample t-test using Background Test Form 2 is

$$t_2 = (M_x + S - M_y) / s^*$$

where the quantity S is a substantial difference. The test statistic t_2 has a standard Student-t distribution with $n+m-2$ degrees of freedom when $\mu_S = \mu_B + S$. If the test statistic t_2 exceeds the critical value $t_{1-\alpha}$, then the null hypothesis in Background Test Form 2 ($H_0: \Delta > S$) may be rejected with $100(1-\alpha)$ percent confidence.

A $100(1-\alpha)$ percent confidence interval for Δ is an interval denoted as (Δ_1, Δ_2) that satisfies the requirement

$$\Pr\{\Delta_1 \leq \Delta \leq \Delta_2\} \geq 1 - \alpha.$$

Here Δ_1 represents the lower limit of the confidence interval, and Δ_2 represents the upper limit of the confidence interval. Although one-sided hypothesis tests were considered on the previous page, the desired confidence interval is two-sided and symmetric, meaning that there is a probability of $\alpha/2$ that Δ will be below this interval and a probability of $\alpha/2$ that it will be above this interval.

If the lower limit of a $100(1-\alpha)$ percent confidence interval for Δ is greater than zero, then the mean in the potentially impacted area is significantly greater than the background mean. This means that a one-sided, level- $\alpha/2$ test of the null hypothesis $H_0: \Delta \leq 0$

(Test Form 1) will reject the null hypothesis. Similarly, if the upper limit of a $100(1-\alpha)$ percent confidence interval for Δ is less than S , then the difference between the mean in the potentially impacted area and the background mean is significantly less than a substantial difference. This means that a one-sided, level- $\alpha/2$ test of the null hypothesis $H_0: \Delta > S$ (Test Form 2) will reject the null hypothesis.

A symmetric confidence interval for the difference $\Delta = \mu_y - \mu_x$ is constructed using $t_{1-\alpha/2}$, which represents the $100(1-\alpha/2)$ th quantile of the Student t-distribution with $n+m-2$ degrees of freedom. A $100(1-\alpha)$ percent confidence interval for Δ has the form (Δ_1, Δ_2) , where the lower bound is

$$\Delta_1 = (M_y - M_x) - t_{1-\alpha/2} s^*$$

and the upper bound is

$$\Delta_2 = (M_y - M_x) + t_{1-\alpha/2} s^*.$$

Although the distribution of the test statistic for the two-sample Student t-test is derived based on the assumption of normal distributions with equal variances, the test is robust and has demonstrated good performance when the variances are unequal, and when the population distributions are not normal. However, the estimates M_y , M_x and s_p^2 are sensitive to outliers in either data set. If either or both data sets contain non-detects, then the test will be sensitive to most common methods of handling these values. Confidence intervals derived using the two-sample test statistic are expected to have similar properties.

NORMAL THEORY, CASE 2: UNEQUAL, UNKNOWN VARIANCES¹⁰

Now assume that the site data (Y_1, \dots, Y_n) and background data (X_1, \dots, X_m) are independent random samples from normal distributions with different means, μ_y and μ_x , and different variances, σ_y^2 and σ_x^2 , respectively:

$$Y_j \sim N[\mu_y, \sigma_y^2]$$

and

$$X_j \sim N[\mu_x, \sigma_x^2].$$

Estimates for the sample variances are

$$s_y^2 = \Sigma(Y_j - M_y)^2 / (n - 1)$$

and

$$s_x^2 = \Sigma(X_j - M_x)^2 / (m - 1).$$

An estimate of the approximate degrees of freedom is

$$v = \tau^2/b$$

where

$$\tau = s_y^2/n + s_x^2/m$$

and

$$b = (s_y^2/n)^2 / (n - 1) + (s_x^2/m)^2 / (m - 1).$$

A symmetric confidence interval for the difference $\Delta = \mu_y - \mu_x$ is constructed using the Student t-distribution with v^* degrees of freedom, where v^* is the closest positive integer to v . Let $t_{1-\alpha/2}$ represent the $100(1-\alpha/2)^{\text{th}}$ quantile of this t-distribution with v^* degrees of freedom. An approximate $100(1-\alpha)$ percent confidence interval for Δ has the form (Δ_1, Δ_2) , where the lower bound is

$$\Delta_1 = (M_y - M_x) - t_{1-\alpha/2}\tau^{1/2}$$

and the upper bound is

$$\Delta_2 = (M_y - M_x) + t_{1-\alpha/2}\tau^{1/2}$$

A.3.2 Comparisons Based on the Wilcoxon Rank Sum Test

The Wilcoxon Rank Sum (WRS)¹¹ test is a nonparametric test for testing whether there is a difference between the site and background population distributions. The WRS test examines whether measurements from one population tend to be consistently larger (or smaller) than those from the other population. The test determines which is the higher distribution by comparing the relative ranks of the two data sets when the data from both sources are sorted into a single list. One assumes that any difference between the site and background concentration distributions represents a shift of the site

concentrations to higher values due to the presence of contamination in addition to background. The WRS test is most effective when contamination is spread throughout a site.

Two assumptions underlying the WRS test are:

- 1) Samples from the background and site are independent, identically distributed random samples; and
- 2) Each measurement is independent of every other measurement, regardless of the set of samples from which it came.

The WRS test assumes that the distributions of the two populations are identical in shape (variance), although the distributions need not be symmetric.

The WRS test has three advantages over the t-test for background comparisons:

- 1) The two data sets are not required to be from a known type of distribution. The WRS test does not assume that the data are normally or log-normally distributed, although a normal distribution approximation often is used to determine the critical value for the test for large sample sizes.
- 2) The WRS test is robust with respect to outliers because the analysis is conducted in terms of ranks of the data. This limits the influence of outliers because a given data point can be no more extreme than the first or last rank.
- 3) The WRS test allows for non-detect measurements to be present in both data sets. The WRS test can handle a moderate number of non-detect values in either or both data sets by treating them as ties.¹²

Theoretically, the WRS test can be used with up to 40 percent or more non-detect measurements in either the background or the site data. Such a high proportion of non-detects indicates that there will be a large number of ties. In this case, the simple

expediency of assigning all ties the same ranks may not be adequate. More specific procedures have been developed to address data sets with a large number of ties.¹³ If more than 40 percent of the data from either the background or site are non-detect values, the WRS test should not be used.

The WRS test may be applied to both forms of background test, *no statistically significant difference* or *exceed by more than a substantial difference*. In either form of background test, the null hypothesis is assumed to be true unless the evidence in the data indicates that it should be rejected in favor of the alternative.

The WRS test for Background Test Form 1 is applied as outlined in the following steps. The site and background measurements are ranked in a single list in increasing order from 1 to N, where $N = m + n$. All tied values are assigned the average of the ranks for that group of measurements. All non-detect values are considered as ties and are assigned an average rank (if there are a total of t non-detects, they all are assigned rank $(t+1)/2$, which is the average of the first t integers).

The sum of the ranks of the site measurements (W_y) and the sum of the ranks of the background measurements (W_x) are sufficient statistics for the test, where $W_y + W_x = N(N + 1)/2$. The sum of the ranks of the site measurements (W_y) is the test statistic used for Background Test Form 1. To conduct the test, W_y is compared with w_{α} , which is the critical value for a level- α WRS test for the appropriate values of n and m .¹⁴ If W_y exceeds the critical value for the test, the null hypothesis that there is no statistically significant difference ($\Delta \leq 0$) may be rejected with $100(1-\alpha)$ percent confidence.

The WRS test for Background Test Form 2 is applied as outlined in the following steps. First, the background measurements are adjusted by adding the substantial difference S to each measured value.¹⁵ Second, the S-adjusted background data and the site data are ranked in a single list in increasing

order from 1 to N. Finally, all tied values are assigned the average of the ranks for that group of measurements.

The sum of the ranks of the S-adjusted background measurements (W_{x+S}) is the test statistic for Background Test Form 2. If W_{x+S} is greater than the critical value for the test, w_{α} , the null hypothesis that the site exceeds background by more than a substantial difference ($\Delta > S$) is rejected at the $100(1-\alpha)$ percent confidence level.

Nonparametric confidence intervals for Δ are derived based on the Mann-Whitney form of the WRS test (Section 5.3.2). The Mann-Whitney test statistics are computed from the set of all possible differences between the site and background data sets:

$$\{Y_i - X_j, I = 1, \dots, n; j = 1, \dots, m\}.$$

There are n times m possible differences in this set, so a computer program may be required to perform the necessary calculations. Let the symbol Z_r ($r = 1, \dots, nm$) represent the r^{th} -ranked difference in the ordered set of all possible differences between the site and background data. A symmetric nonparametric confidence interval for Δ is constructed using the k^{th} -smallest ranked difference (Z_k) and the k^{th} -largest ranked difference (Z_{nm-k+1}) in the set of all possible differences, where k depends on n , m , and α .¹⁶ Thus, a $100 \times (1-\alpha)$ percent confidence interval for Δ is a closed interval of the form

$$(\Delta_1, \Delta_2) = (Z_k, Z_{nm-k+1})$$

with

$$k = w_{\alpha/2} - n(n + 1)/2.$$

Here, as noted above for the WRS test, $w_{\alpha/2}$ is the tabulated critical value for a level- $\alpha/2$ WRS test for the appropriate values of n and m . This confidence interval satisfies the requirement

$$\Pr\{ \Delta_1 \leq \Delta \leq \Delta_2 \} \geq 1 - \alpha.$$

APPENDIX A NOTES

1. U.S. Environmental Protection Agency (EPA). 2001. *Requirements for Quality Assurance Project Plans, EPA QA/R-5*. <http://www.epa.gov/QUALITY/qapps.html>).
2. Mood, A.M., Graybill, F. A. And Boes, D. C., *Introduction to the theory of statistics*, 3rd Ed., McGraw Hill, Boston, Mass., 1974, p. 411.
3. Bickel, P.J., and Doksom, K. A., *Mathematical Statistics: Basic Ideas and Selected Topics*, Holden Day, San Francisco, 1977, p. 168.
4. U.S. Environmental Protection Agency (EPA). July 2000. *Guidance for Data Quality Assessment: Practical Methods for Data Analysis, EPA QA/G-9, QA00 Version*. Quality Assurance Management Staff, Washington, DC. EPA 600-R-96-084. Available at http://www.epa.gov/quality/qa_docs.html.
5. U.S. Environmental Protection Agency (EPA). 1994. *Guidance for the Data Quality Objectives Process, EPA QA/G-4*, EPA 600-R-96-065. Washington DC.
6. U.S. Environmental Protection Agency (EPA). 1989. *Statistical Methods for Evaluating the Attainment of Cleanup Standards Volume 3*, subtitled *Reference-Based Standards for Soils and Solid Media*, EPA 230-02-89-042. Washington DC.
7. U.S. Environmental Protection Agency (EPA). 1989. *Risk Assessment Guidance for Superfund Vol. I, Human Health Evaluation Manual (Part A)*. Office of Emergency and Remedial Response, Washington, DC. EPA 540-1-89-002.
8. U.S. Environmental Protection Agency (EPA). 1996. *Soil Screening Guidance: Technical Background Document*, EPA 540-R-95-128.
9. Zwillinger, D. and S. Kokoska. 2000. *CRC Standard Probability and Statistics Tables and Formulae*, Chapman and Hall/CRC Press, New York, Section 9.6.2.
10. Zwillinger and Kokoska, *Op. Cit.*, Section 9.6.3.
11. The WRS test is also called the Mann-Whitney test, which is mathematically equivalent to the WRS test. Sometimes, the combined name is used: Wilcoxon-Mann-Whitney test. See Section 5.3.2.
12. The Gehan form of the WRS test should be considered if there are many non-detect values with different detection levels.
13. If there are many ties, see instructions for applying the WRS test in Conover, W.J., *Practical Nonparametric Statistics, 2nd Ed.*, John Wiley & Sons, Inc., New York, NY, 1980.
14. Critical values for the WRS test are available in many published texts and reference books. Two sources are Conover, W.J., *Practical Nonparametric Statistics, 2nd Ed.*, John Wiley & Sons, NY, 1980; and *CRC Standard Probability and Statistics Tables and Formulae*, D. Zwillinger and S. Kokoska, Chapman and Hall/CRC Press, Boca Raton, Florida, 2000.

15. Conover, *Practical Nonparametric Statistics, 2nd Ed., Op. Cit.*, p. 223, Equation 8.
16. Conover, *Practical Nonparametric Statistics, 2nd Ed., Op. Cit.*, p. 224.

APPENDIX B

POLICY CONSIDERATIONS FOR THE APPLICATION OF BACKGROUND DATA IN RISK ASSESSMENT AND REMEDY SELECTION

Role of Background in the CERCLA Cleanup Program

**U.S. Environmental Protection Agency
Office of Solid Waste and Emergency Response
Office of Emergency and Remedial Response**

April 26, 2002

OSWER 9285.6-07P

Table of Contents

Purpose B-3

History B-3

Definitions of Terms B-4

Consideration of Background in Risk Assessment B-5

Consideration of Background in Risk Management B-6

Consideration of Background in Risk Communication B-7

Hypothetical Case Examples B-7

References B-10

Purpose

This document clarifies the U.S. Environmental Protection Agency (EPA) preferred approach for the consideration of background constituent concentrations of hazardous substances, pollutants, and contaminants in certain steps of the remedy selection process, such as risk assessment and risk management, at Comprehensive Environmental Response, Compensation, and Liability Act (CERCLA or “Superfund”) sites. To the extent practicable, this document may also be applicable to sites addressed under removal actions and time-critical actions. In general, the presence of high background concentrations of hazardous substances, pollutants, and contaminants found at a site is a factor that should be considered in risk assessment and risk management.

The primary goal of the CERCLA program is to protect human health and the environment from current and potential threats posed by uncontrolled releases of hazardous substances, pollutants, and contaminants. Contamination at a CERCLA site may originate from releases attributable to the CERCLA site in question, as well as contamination that originated from other sources, including natural and/or anthropogenic sources not attributable to the specific site releases under investigation (EPA, 1995a). In some cases, the same hazardous substance, pollutant, and contaminant associated with a release is also a background constituent. These constituents should be included in the risk assessment, particularly when their concentrations exceed risk-based concentrations. In cases where background levels are high or present health risks, this information may be important to the public. Background information is important to risk managers because the CERCLA program, generally, does not clean up to concentrations below natural or anthropogenic background levels.

A comprehensive investigation of all background substances found in the environment usually will not be necessary at a CERCLA site. For example, radon background samples normally would not be collected at a chemically contaminated site unless radon, or its precursor (radium, Ra-226) was part of the CERCLA release. Also, EPA normally would not analyze background samples for Ra-226 at a cesium (Cs-137) site, or dioxin at a lead site where dioxin was not the subject of a CERCLA release into the environment.

This document provides guidance to EPA Regions concerning how the Agency intends to exercise its discretion in implementing one aspect of the CERCLA remedy selection process. The guidance is designed to implement national policy on these issues.

Some of the statutory provisions described in this document contain legally binding requirements. However, this document does not substitute for those provisions or regulations, nor is it a regulation itself. Thus, it cannot impose legally-binding requirements on EPA, States, or the regulated community, and may not apply to a particular situation based upon the circumstances. Any decisions regarding a particular remedy selection decision will be made based on the statute and regulations, and EPA decision makers retain the discretion to adopt approaches on a case-by-case basis that differ from this guidance where appropriate. EPA may change this guidance in the future.

History

Background issues are discussed in a number of EPA documents.¹ A need for CERCLA-specific guidance

¹ ▶ *Risk Assessment Guidance for Superfund Volume I, Human Health Evaluation Manual [RAGS]* (EPA, 1989).
(continued...)

was identified during risk assessment reform discussions with stakeholders in 1997. An issue that is often raised at CERCLA sites is whether a reliable representation of background is established (EPA, 1989). To assist Regions with this issue, EPA developed a peer-reviewed practical guide to sampling and statistical analysis of background concentrations in soil at CERCLA sites (EPA, 2001b).

EPA has developed this policy to respond to questions about the general application of background concentration during the CERCLA remedial investigation process.² This policy encourages national consistency and responds to the Agency's goals for risk characterization and communication of risks to the public as expressed in other EPA policy and guidance, including:

- ▶ *Policy for Risk Characterization* which provides principles for fully, openly, and clearly characterizing risks (EPA, 1995b); and
- ▶ *Cumulative Risk Assessment Guidance* which encourages programs to better advise citizens about the environmental and public health risks they face (EPA, 1997c).

Definitions of Terms

For the purposes of this policy, the following definitions are used.

Background refers to constituents or locations that are not influenced by the releases from a site, and is usually described as naturally occurring or anthropogenic (EPA, 1989; EPA, 1995a):

- 1) *Anthropogenic* – natural and human-made substances present in the environment as a result of human activities (not specifically related to the CERCLA release in question); and,
- 2) *Naturally occurring* – substances present in the environment in forms that have not been influenced by human activity.

Chemicals (or constituents) of concern (COCs) are the hazardous substances, pollutants, and contaminants that, at the end of the risk assessment, are found to be the *risk drivers* or those that may actually pose

¹ (...continued)

- ▶ Preamble to the National Oil and Hazardous Substances Pollution Contingency Plan (NCP, 1990a).
- ▶ *Role of the Baseline Risk Assessment in Superfund Remedy Selection Decisions* (EPA, 1991).
- ▶ *Determination of Background Concentrations of Inorganics in Soils and Sediments at Hazardous Waste Sites* (EPA, 1995a).
- ▶ *Soil Screening Guidance: User's Guide* (EPA, 1996).
- ▶ *Ecological Risk Assessment Guidance for Superfund* (EPA, 1997a).
- ▶ *Rules of Thumb for Superfund Remedy Selection* (EPA, 1997b).
- ▶ *Soil Screening Guidance for Radionuclides: User's Guide* (EPA, 2000).
- ▶ *ECO Update. The Role of Screening-Level Risk Assessments and Refining Contaminants of Concern in Baseline Ecological Risk Assessments* (EPA, 2001a).

² The process of determining when risks warrant remedial actions and the degree of cleanup for specific hazardous substances, pollutants, and contaminants involves many factors that are not addressed in this document. Additional guidance is provided in the EPA (1991) *Role of the Baseline Risk Assessment in Superfund Remedy Selection Decisions*.

unacceptable human or ecological risks.³ The COCs typically drive the need for a remedial action (EPA, 1999a).

Chemicals (or constituents) of potential concern (COPCs) generally comprise the hazardous substances, pollutants, and contaminants that are investigated during the baseline risk assessment. The list of COPCs may include all of the constituents whose data are of sufficient quality for use in the quantitative risk assessment, or a subset thereof (EPA, 1989).

Screening is a common approach used by risk assessors to refine the list of COPCs to those hazardous substances, pollutants, and contaminants that may pose substantial risks to health and the environment. Screening involves a comparison of site media concentrations with site-specific risk-based values.⁴

Consideration of Background in Risk Assessment

A baseline risk assessment generally is conducted to characterize the current and potential threats to human health and the environment that may be posed by hazardous substances, pollutants, and contaminants at a site. EPA's 1989 *Risk Assessment Guidance for Superfund* (RAGS) provides general guidance for selecting COPCs, and considering background concentrations. In RAGS, EPA cautioned that eliminating COPCs based on background (either because concentrations are below background levels or attributable to background sources) could result in the loss of important risk information for those potentially exposed, even though cleanup may or may not eliminate a source of risks caused by background levels. In light of more recent guidance for risk-based screening (EPA, 1996; EPA, 2000) and risk characterization (EPA, 1995c), this policy recommends a baseline risk assessment approach that retains constituents that exceed risk-based screening concentrations. This approach involves addressing site-specific background issues at the end of the risk assessment, in the risk characterization. Specifically, the COPCs with high background concentrations should be discussed in the risk characterization, and if data are available, the contribution of background to site concentrations should be distinguished.⁵ COPCs that have both release-related and background-related sources should be included in the risk assessment. When concentrations of naturally occurring elements at a site exceed risk-based screening levels, that information should be discussed qualitatively in the risk characterization. To summarize:

- ▶ The COPCs retained in the quantitative risk assessment should include those hazardous substances, pollutants, and contaminants with concentrations that exceed risk-based screening levels.

³ Guidance for determining if site risks are unacceptable is discussed in the EPA (1991) *Role of the Baseline Risk Assessment in Superfund Remedy Selection Decisions*. As stated in the EPA (1991) memorandum, "EPA uses the general 10^{-4} to 10^{-6} risk range as a "target range" within which the Agency strives to manage risks as part of a Superfund cleanup." The risk used in this decision generally is the "cumulative site risk" to an individual using reasonable maximum exposure (RME) assumptions for either current or future land use and includes all exposure pathways which the same person may consistently face. See also EPA (1989) RAGS, Section 8.3.

⁴ Risk-based values or concentrations are generally based on a cancer risk of one-in-a-million (1×10^{-6}) or a hazard quotient of 1.0 for noncarcinogens (EPA, 1996) or screening-level ecological risk values (EPA, 1997a; EPA, 2001a). COPCs with concentrations below the screening levels might be excluded from the risk assessment unless there are other pathways or conditions that are not addressed by the screening values (EPA, 1996).

⁵ Technical guidance should be consulted for sampling and analysis of background concentration data (EPA, 2001b).

- ▶ The Risk Characterization should include a discussion of elevated background concentrations of COPCs and their contribution to site risks.
- ▶ Naturally occurring elements that are not CERCLA hazardous substances, pollutants, and contaminants, but exceed risk-based screening levels should be discussed in the risk characterization.

This general approach is preferred in order to:

- ▶ Encourage national consistency in this area;
- ▶ Present a more thorough picture of risks associated with hazardous substances, pollutants, and contaminants at a site; and
- ▶ Prevent the inadvertent omission of potentially release-related hazardous substances, pollutants, and contaminants from the risk assessment.

This approach is consistent with the *Policy for Risk Characterization* which provides principles for fully, openly, and clearly characterizing risks (EPA, 1995b). Risks identified during the baseline risk assessment should be clearly presented and communicated for risk managers and for the public. Risk characterization is one of many factors in determining appropriate CERCLA risk management actions (EPA, 1991; EPA, 1995b).

Consideration of Background in Risk Management

Where background concentrations are high relative to the concentrations of released hazardous substances, pollutants, and contaminants, a comparison of site and background concentrations may help risk managers make decisions concerning appropriate remedial actions. The contribution of background concentrations to risks associated with CERCLA releases may be important for refining specific cleanup levels for COCs that warrant remedial action.⁶

Generally, under CERCLA, cleanup levels are not set at concentrations below natural background levels. Similarly, for anthropogenic contaminant concentrations, the CERCLA program normally does not set cleanup levels below anthropogenic background concentrations (EPA, 1996; EPA, 1997b; EPA, 2000). The reasons for this approach include cost-effectiveness, technical practicability, and the potential for recontamination of remediated areas by surrounding areas with elevated background concentrations. In cases where area-wide contamination may pose risks, but is beyond the authority provided under CERCLA, EPA may be able to help identify other programs or regulatory authorities that are able to address the sources of area-wide contamination, particularly anthropogenic (EPA, 1996; EPA, 1997b; EPA, 2000). In some cases, as part of a response to address CERCLA releases of hazardous substances, pollutants, and contaminants, EPA may also address some of the background contamination that is present on a site due to area-wide contamination.

⁶ For example, in cases where a risk-based cleanup goal for a COC is below background concentrations, the cleanup level may be established based on background.

The determination of appropriate CERCLA response actions and chemical-specific cleanup levels includes the consideration of nine criteria as provided in the National Oil and Hazardous Substances Pollution Contingency Plan (NCP, 1990b). In cases where applicable or relevant and appropriate requirements (ARARs) regarding cleanup to background levels apply to a CERCLA action, the response action generally should be carried out in the manner prescribed by the ARAR. In the case where a law or regulation is determined to be an ARAR and it requires cleanup to background levels, the ARAR will normally apply and be incorporated into the Record of Decision, unless the ARAR is waived.

Consideration of Background in Risk Communication

EPA strives for transparency in decision-making (EPA, 1995c) and encourages programs to better advise citizens about the environmental and public health risks they face (EPA, 1997c). The presence of high background concentrations of COPCs may pose challenges for risk communication. For example, the discussion of background may raise the expectation that EPA will address those risks under CERCLA. The knowledge that background substances may pose health or environmental risks could compound public concerns in some situations.

On the other hand, knowledge of background risks could help some community members place CERCLA risks in perspective. Also, the information about site and background risks can be helpful for both risk managers who make an appropriate CERCLA decision, and for members of the public who should know about environmental risk factors that come to light during the remedial investigation process.

As a general policy matter, EPA strives for early and frequent outreach to communities in order to share information and encourage involvement (EPA, 2001c). EPA has made a clear commitment to fully, openly, and clearly characterize and communicate risks (EPA, 1995b; EPA, 1995c). There is no one-size-fits-all technique that can help explain risks associated with CERCLA releases or with background levels, or the basis of risk management decisions. Approaches will depend on the site, the issues, and the level of community interest. Early on in the process, Regions should clarify their understanding of stakeholder expectations and clearly explain the relevant constraints and limitations of the CERCLA remedial process (EPA, 1999b; EPA, 2001c).

In some cases where area-wide contamination may pose a risk, but is beyond the authority of the CERCLA program, communication of potential risks to the public may be most effective when coordinated with public health agencies. Examples of situations where Regions might coordinate risk communication with local, state or federal health officials are sites where widespread lead contamination or high levels of naturally occurring radiation have been found, but are not the subject of a CERCLA release into the environment. Public health agency officials may combine education and outreach efforts to inform residents about ways to reduce exposures and risks.

Hypothetical Case Examples

Three general hypothetical case examples are given to show how background may be considered in risk assessment and risk management at CERCLA sites:

Case 1 presents an example of a chemical site with widespread background contamination.

Case 2 presents an example of a radiation site with both natural- and release-related sources.

Case 3 presents an example of a site with hazardous substances, pollutants, and contaminants from both natural- and release-related sources.

In these examples, it is presumed that adequate samples are collected from appropriate background reference locations and evaluated using appropriate statistical methods. It is presumed that background is not used to screen out substances from the risk assessment. For simplicity, only one pathway⁷ is used for hypothetical human health risk assessments.⁸

Based on the presumptions above, the basic concepts these examples are designed to highlight are:

- ▶ Background issues should be discussed in the risk characterization portion of the baseline risk assessment in order to inform risk management decisions;
- ▶ Information about unacceptable risks should be communicated to public; and
- ▶ Other factors, such as the nine criteria provided in the NCP, should be considered by the risk manager in making final decisions.

Hypothetical Case 1

The ABC Industrial Site risk assessment included all COPCs that exceed site-specific risk-based concentrations for soil pathways. The results of the risk assessment identified the following COPCs with risks above or at the high end of the 10^{-4} to 10^{-6} risk range: arsenic, dieldrin, and 4,4-DDT. The hazard quotients were below 1.0.

Arsenic is a potential background substance—it is a common naturally occurring element—but is also a hazardous substance that was released at this site. The available site characterization data indicate that soil arsenic concentrations may be naturally occurring or consistent with background concentrations. Dieldrin and DDT are present at high concentrations that contribute to an unacceptable site risk. However, only dieldrin is known to be associated with the CERCLA site activities and releases. Since there are no known historical uses of DDT at this CERCLA site, the RPM suspects that the DDT in soil originated from area-wide agricultural pesticide applications in this part of the state. Based on this information, the RPM requests additional sampling of background locations for arsenic and DDT analysis. A statistical comparison of sampling data for arsenic and 4,4-DDT in on-site samples and background samples indicates that site concentrations for DDT are consistent with background concentrations. Local and regional data support the conclusion that DDT is an area-wide contaminant. The additional data indicate that arsenic concentrations

⁷ At most CERCLA sites, risks for the reasonably maximum exposed individual typically are combined across several exposure pathways to estimate the total risks at a CERCLA site. This is done only for the pathways which the same individual would be likely to face consistently (EPA, 1989). Depending on the particular CERCLA site, risks could be calculated for the entire area of the site or for separate units (see Section 4.5 of RAGS (EPA, 1989)). More technical guidance for characterizing background concentrations and comparing data sets is provided in EPA (2001b) and other technical references cited previously in this document.

⁸ Guidance on the consideration of background concentrations during screening level ecological risk assessments is provided in EPA (2001a).

on the site are above background concentrations. Therefore, the arsenic risks cannot be attributed solely to background.

In this example, arsenic and dieldrin are the soil COCs for which cleanup goals should be derived. The risk characterization should present information about DDT as an area-wide background contaminant that is unrelated to releases at this site, and the Agency should explain whether or not it will be addressed. The RPM should consider whether other regulatory programs or authorities are able to address the area-wide DDT contamination in a coordinated response effort. If available, the location(s) of additional information on pesticide use in this part of the state should be provided for concerned citizens.

Hypothetical Case 2

At ABC Radium Production Site, site characterization data indicate that radium (Ra-226) and inorganics are present in soil. Arsenic concentrations exceed screening levels but are assumed to be within naturally occurring levels. To confirm this assumption, the RPM evaluates site-specific background samples for comparison to site concentrations. The site-specific background analysis confirms that arsenic concentrations collected on the site are consistent with background concentrations in soils. There are no known regional anthropogenic sources of arsenic (such as smelters or pesticide manufacturers). Arsenic, in this case, is considered to be a naturally occurring substance and is excluded from further consideration in the quantification of site risks. However, the finding of natural background arsenic at concentrations that may pose health risks should be discussed in the text of the risk characterization.

The risk assessment indicates that Ra-226 exceeds the high end of the acceptable risk range of 10^{-4} to 10^{-6} . It is commonly known that Ra-226 occurs naturally in the environment. Samples collected in an appropriate background location near this site indicate that Ra-226 levels from natural sources are lower than the site levels, but are associated with a risk at the upper end of the risk range (10^{-4}).

In this example, only Ra-226 should be a COC for which a cleanup goal should be derived. The risk characterization, however, should include a discussion of natural background levels of both arsenic and Ra-226.

Hypothetical Case 3

XYZ Site contains buried chemical wastes, but some anecdotal accounts indicate that radium may have been used. Preliminary site characterization data show that arsenic, manganese, and Ra-226 concentrations exceed the site-specific, risk-based concentrations. A comparison of arsenic and manganese concentrations in groundwater samples collected from upgradient background locations indicates that only manganese site concentrations are consistent with background levels and considered to be naturally occurring. Naturally occurring manganese is not considered further in the quantification of risks, but is included in a qualitative discussion of risks in the risk characterization.

The RPM decides to analyze for Ra-226 both at the site and in background locations because it is commonly known that Ra-226 occurs naturally in the environment. Samples are collected in an appropriate background location near this site. The samples indicate that Ra-226 levels at this site are not different from naturally occurring levels. Therefore, Ra-226 is not a COPC for further consideration in the quantification of risks. Subsequent site investigation data confirms the use of chemicals, but not radionuclides.

In this example, only arsenic risks are quantified in the risk assessment. The baseline risk for groundwater indicates that arsenic poses an unacceptable risk. The risk characterization should include a discussion of the natural Ra-226 and manganese concentrations because the levels exceeded risk-based concentrations. Site characterization data indicate that site disposal activities caused naturally occurring arsenic in soil to be mobilized and leach to groundwater. Arsenic, therefore, is the subject of a CERCLA release into the environment and a cleanup goal for it should be derived.

References

- U.S. Environmental Protection Agency (EPA). 1989. *Risk Assessment Guidance for Superfund (RAGS): Volume I: Human Health Evaluation Manual (HHEM), (Part A), Interim Final*, Office of Emergency and Remedial Response, Washington, DC. EPA/540/1-89/002, OSWER 9285.70-02B.
- U.S. Environmental Protection Agency (EPA). 1991. *Role of the Baseline Risk Assessment in Superfund Remedy Selection Decisions*, Office of Emergency and Remedial Response, Washington, DC. OSWER 9355.0-30.
- U.S. Environmental Protection Agency (EPA). 1995a. *Engineering Forum Issue Paper. Determination of Background Concentrations of Inorganics in Soils and Sediments at Hazardous Waste Sites*, R.P Breckenridge and A.B. Crockett, Office of Research and Development, Office of Solid Waste and Emergency Response, Washington, DC. EPA/540/S-96/500.
- U.S. Environmental Protection Agency (EPA). 1995b. *Policy for Risk Characterization at the U. S. Environmental Protection Agency*, Science Policy Council, Washington, DC. <http://www.epa.gov/OSP/spc/rcpolicy.htm>.
- U.S. Environmental Protection Agency (EPA). 1995c. *Risk Characterization Handbook*, Science Policy Council, Washington, DC. EPA 100-B-00-002.
- U.S. Environmental Protection Agency (EPA). 1996. *Soil Screening Guidance: User's Guide*. Office of Emergency and Remedial Response, Washington, DC. EPA/540-R-96/018, OSWER 9355.4-23.
- U.S. Environmental Protection Agency (EPA). 1997a. *Ecological Risk Assessment Guidance for Superfund: Process for Designing and Conducting Ecological Risk Assessments, Interim Final*, EPA/540-R-97-006, OSWER 9285.7-25.
- U.S. Environmental Protection Agency (EPA). 1997b. *Rules of Thumb for Superfund Remedy Selection*. Office of Emergency and Remedial Response, Washington, DC. EPA 540-R-97-013, OSWER 9355.0-69.
- U.S. Environmental Protection Agency (EPA). 1997c. *Cumulative Risk Assessment Guidance-Phase I Planning and Scoping*, Science Policy Council, Washington, DC. <http://www.epa.gov/osp/spc/cumrisk2.htm>.
- U.S. Environmental Protection Agency (EPA). 1999a. *A Guide to Preparing Superfund Proposed Plans, Records of Decision, and Other Remedy Selection Decision Documents*, Office of Emergency and

Remedial Response, Washington, DC. OSWER 9200.1-23.P.

U.S. Environmental Protection Agency (EPA). 1999b. *Risk Assessment Guidance for Superfund: Volume 1 – Human Health Evaluation Manual Supplement to Part A: Community Involvement in Superfund Risk Assessments*, Office of Emergency and Remedial Response, Washington, DC. OSWER 9285.7-01E-P.

U.S. Environmental Protection Agency (EPA). 2000. *Soil Screening Guidance for Radionuclides: User's Guide*, Office of Radiation and Indoor Air, OSWER 9355.4-16A.

U.S. Environmental Protection Agency (EPA). 2001a. *ECO Update, The Role of Screening-Level Risk Assessments and Refining Contaminants of Concern in Baseline Ecological Risk Assessments*, OSWER 9345.0-14.

U.S. Environmental Protection Agency (EPA). 2001b. *Guidance for Characterizing Background Chemicals in Soil at Superfund Sites*, External Review Draft, Office of Emergency and Remedial Response, OSWER. 9285.7-41. [Replaced by *Guidance for Comparing Background and Chemical Concentrations in Soil for CERCLA Sites*, EPA 540-R-01-003, September 2002.]

U.S. Environmental Protection Agency (EPA). 2001c. *Early and Meaningful Community Involvement*, Office of Emergency and Remedial Response, OSWER 9230-0-9.

NCP, 1990a. Preamble to the National Oil and Hazardous Substances Pollution Contingency Plan (NCP), 40 CFR Part 300, 53 *Federal Register* 51394 and 55 *Federal Register* 8666.

NCP, 1990b. National Oil and Hazardous Substances Pollution Contingency Plan (NCP), 40 CFR Part 300, 55 *Federal Register* 8666, March 8, 1990.