1

# Weight of Evidence (WOE): Quantitative Estimation of Probability of Impact

3

4

5

6  Eric P. Smith

7  Ilya Lipkovich

8  Keying Ye

9

10

11  Dept. of Statistics

12  Virginia Tech

13

14  Corresponding Author

15  Eric P. Smith

16

17  Dept. of Statistics

18  Virginia Tech

19  Blacksburg, VA 24061-0439

20  email:epsmith@vt.edu

21  Phone: 542-231-7929

22  Fax: 540-231-3863

23

24  Updated 12/03/01, 1/14/02, edited 2/5/02

25  Corrected 2/10/02

26

27

28

29

30  Running head: Quantitative weight of evidence for environmental assessment

31

1

2    Abstract

3

4    Environmental decision-making is complex and often based on multiple lines of evidence.

5    Integrating the information from these multiple lines of evidence is rarely a simple process. We

6    present a quantitative approach to the combination of multiple lines of evidence through

7    calculation of weight of evidence, with reference conditions used to define a not impaired state.

8    The approach is risk-based with measurement of risk computed as the probability of impairment.

9    When data on reference conditions are available, there are a variety of methods for calculating

10    this probability.  Statistical theory and the use of odds ratios provide a method for combining the

11    measures of risk from the different lines of evidence.  The approach is illustrated using data from

12    the Great Lakes to predict the risk at potentially contaminated sites.

13

14    Keywords: Bayesian statistics, odds ratio, hazard ranking, combining information, risk

15    assessment, reference conditions

16

1

## Introduction

3

4    Environmental decision-making is often based on multiple sets of information or lines of

5    evidence.  By line of evidence we mean a set of information that pertains to an important aspect

6    of the environment.  For example, in the sediment quality triad (Chapman 1996), there are three

7    lines of evidence, the toxicity line, the biological field line and the chemistry line.  It is difficult

8    to combine the information from these multiple sources into a single measure for decision-

9    making.  Weight of evidence (WOE) is sometimes used as an approach for combining the

10   information, however it is rarely used in a quantitative manner.  This paper discusses a

11   quantitative approach to WOE.  A statistical approach is taken in which the likelihood of the data

12   is calculated under two different scenarios and a decision made based on the ratio of the

13   likelihood.  Our view is that there are two states, and we must decide which of the states is true.

14   Examples of pairs of states common to environmental decision making are (impaired, not

15   impaired), (remediate, don't remediate), (list, don't list), etc. The view we take is that interest is

16   in a single site, we collect a sample at that site and based on the sample, the site is impacted or

17   the site is not impacted.  For practical reasons, we assume the simple case that there is ample

18   information on reference conditions and interest is in evaluating a single new location.  This

19   gives us the ability to obtain a precise estimate of the probability.

20

## Estimating weight of evidence (WOE)

22

23   A quantitative approach to WOE is based on the concept of statistical weight of evidence.

24   This idea dates back to work by Alan Turing in World War II (for a more general discussion of

25   history and concepts of statistical weight of evidence, see Good 1988).  In this approach, there

26   are two states and we must decide which state is more likely given the data.  We can view the

27   outcomes as the site is or is not impacted.  Without observing the information, we may have

28   opinions or insights into the condition of the locations.  This insight might be based on previous

29   data (condition in previous years) or be from sites that are close in space.  This information may

30   be used to form a prior opinion or probability of impairment.  After the data are collected, we

31   process the data to evaluate the site.  This leads to a Bayesian approach in which the data are

1  used to update the prior information.  The lack of prior information suggests a frequentist

2  approach where the data alone are used to make a decision.  The approach may be based on a

3  single line of evidence or multiple lines.  The individual lines of evidence are usually evaluated

4  separately and by combining them we hope to make a stronger inference.

5  Statistical WOE is based on a quantitative evaluation of the data and requires a model

6  that describes the data. In the simplest approach, there are two states and we must decide which

7  state is more likely given the data.  Because it is not always easy to describe impact, an

8  alternative approach is to evaluate the risk of impairment of a site considering the baseline risk of

9  a not impaired site.  If we view the possible outcomes as the site is or is not impacted, the risk is

10  then the probability that the site is impacted or *P(impact)* where we evaluate this probability after

11  information is collected.

12  The odds are a way to evaluate how big the probability is relative to the baseline risk.

13  Odds of impact are defined as

14

15  $$\text{odds(impact)} = \text{P(impact)/P(no impact)} \qquad (1)$$

16

17  Although this problem is analogous to tossing a coin, estimating the probability of impact is not

18  easy since "impact" is not an observable attribute of a sample.  The state (impact or no impact)

19  must be inferred based on information that is collected on both impacted and unimpacted sites.

20  A reasonable approach is based on Bayes rule.  With no data the probabilities would be

21  estimated based on prior information that may come from previous studies.  We more generally

22  would collect data to improve these estimates.  Given data we have to calculate probabilities of

23  impact or no impact.  However, even with data we do not have probabilities of impact, only

24  probabilities associated with observations given a model for impacted sites and not impacted

25  sites.  For example, if there is ample information on sites that are not impacted we may compare

26  our data to that and estimate the probability the data come from that distribution.  We do not

27  have the probability of no impact only the probability that the data come from that distribution.

28  A representation of what we want to calculate is

29

$$\frac{P(impact \mid data)}{P(no\ impact \mid data)}$$

1

2  Given data and information about the different groups we can calculate how likely the data are

3  given they are from one of the groups.  This probability is written as P(data|impact) or P(data|no

4  impact) and must be computed based on a statistical model for the data (possibly different

5  models for each group). For example, we might specify that the model for the impact sites for

6  dissolved oxygen is normal with mean 4 and for no impact sites is normal with mean 7.  If we

7  obtain a sample with dissolved oxygen equal to 6 and we know the variance of the dissolved

8  oxygen measurements, we can calculate how likely the observation is to have come from each

9  group by calculating the value of the density under each model.  Bayes theorem may be used to

10 calculate the ratio in terms of the probability of the data since

11

$$\frac{P(impact \mid data)}{P(no\ impact \mid data)} = \frac{P(data \mid impact)P(impact)}{P(data \mid no\ impact)P(no\ impact)}$$

12

13 where *P(impact)* is the probability the site is an impact site before collecting the measurements

14 and is referred to as the prior probability.  If the prior probabilities are taken to be equal the result

15 is

$$\frac{P(impact \mid data)}{P(no\ impact \mid data)} = \frac{P(data \mid impact)}{P(data \mid no\ impact)}$$

16

17      The quantity on the right hand side of the equation is referred to as the likelihood ratio or

18 a Bayes Factor in this case and measures the likelihood of the data given the site is in the impact

19 class versus the no impact class. (More generally the Bayes Factor also involves parameters that

20 are treated as random and integrated out of calculations, see Kass and Raftery, 1995).  For

21 environmental problems there may not be simple approaches for estimating these quantities.

22 Building a model for the impacted or unimpacted sites requires information on how data are

23 distributed for these types of sites and other factors that might influence the observations.  Sites

24 classified as unimpacted are often viewed as reference sites.  It may not be possible to obtain

25 these sites or there may be covariates that must be considered. Calculation of the likelihood of

26 the data under impact requires a definition of the impact or model of the data that we might

27 expect if the observation came from the impact group.  One would have to have different models

1    for different types of impacts (for example chemical toxicity versus sedimentation).  The models

2    should depend on the strength of the impact, and may vary in space and time.  These models may

3    involve a good deal of work to describe.  One approach to calculating the ratio given a lack of

4    information on impact is to calculate the odds as $p/(1-p)$ where $p$ is the probability of the data

5    given there is an impact.

6         WOE is a measure of how much an observed feature in the data adds to or subtracts from

7    the evidence of impact.  Numerically it has been defined as (Good 1988)

8

9         *Weight of evidence = 10 log(likelihood ratio).*

10

11   In our applications, this would correspond to the weight of evidence for one line of evidence.  A

12   natural consequence of using logs of ratios is that the weight of evidence from different lines

13   may be added together to get an overall weight of evidence.  For an individual line of evidence

14   values of the ratio are interpreted as follows (Good, 1988):

| Bayes Factor | Weight of Evidence | Strength of Evidence |
|---|---|---|
| <5 | <6.9 | Weak |
| 5-10 | 6.9-10 | Moderate |
| 10-100 | 10-20 | Moderate to strong |
| >100 | >20 | Strong |

15   These rules are guidelines in much the same way that p-values are guidelines and other authors

16   have suggested alternative views (see Kass and Raftery, 1995).  From a hypothesis testing

17   perspective the weight of evidence measures the strength of the evidence against the null

18   hypothesis.

19        Numerical calculation of weight of evidence is not common to statistics.  The reason is

20   that testing of hypotheses and interpretation more common approach of likelihood ratio testing

21   and calculation of Bayes Factors.  Using the natural log scale and twice the log of the Bayes

22   Factor leads to the same scale as likelihood ratio testing in general statistical theory (where –

23   2log(likelihood ratio) is used to test hypotheses) and deviance measures in generalized linear

24   models (McCullagh and Nelder, 1989).  Thus, for applications of WOE in environmental

25   problems a user may choose to summarize results in terms of a WOE measure that is based on

26   the probability of impairment.  Alternatively the user may simply calculate and report the actual

1    probability.   The value of the use of Bayes Factors and odds is that these may be combined

2    easily over the different lines of evidence.

3        The calculation of WOE for multi-attribute environmental studies involves three general

4    stages of analysis.  We assume that the researcher has available the information required for

5    making the decision.  Thus decisions have been made about what information needs to be

6    collected or this information has already been collected.  The three stages of the analysis are the

7    preprocessing stage, the processing stage and the combination of the information over the lines

8    of evidence.

9

10   *Data preprocessing*

11

12       The initial step in the analysis is the preprocessing of the data.  Preprocessing involves

13   selection of the variables to be used in the analysis and scaling or transforming these variables.

14   Variables are selected to provide relevant statistical and scientific information on differences

15   between control and impact.  Scaling and transformation are often used to meet assumptions

16   required for analysis. The assumptions needed depend on the model used to calculate the

17   probability.  Two methods for this calculation are to assume a model (parametric approach) or to

18   calculate the probability using a nonparametric approach.  In the parametric approach we select a

19   probability model for the data.  Here we use a normal or Gaussian model to estimate the risk.

20   Then there are a several assumptions that need to be evaluated for the statistical model to provide

21   a good estimate:

22

23       1.  Normality of the reference data.

24       2.  Independence of samples in the reference set.

25       3.  Homogeneity of variance in the reference set.

26

27   Since a normal distribution is used to calculate probability of impairment the validity of this

28   assumption is required for the estimate of risk to be accurate.  In the normal model, probability is

29   directly related to the standardized distance to the mean.  If the normal model is not reasonable

30   then the estimate may be poor and misleading.  Problems such as skewness and outliers may lead

31   to inaccuracies in the probability estimate.  As environmental data often are not normal, we try to

1  achieve normality via choice of a suitable transformation of the variables or use a method based

2  on the distribution of the data (i.e. logistic regression assumes a binomial distribution). The

3  logarithm is typically used as a transformation with contaminant concentrations. Independence

4  of the reference data is required to provide a good estimate of the variances and covariances in

5  the contaminants. This assumption is best met through choice of the reference locations and

6  sampling occasion. Sites that are spatially close and repeated samples at the same site that are

7  temporally close should be avoided. Homogeneity of the variance in the reference set is required

8  to produce a good estimate of the variances and covariances. An alternative approach would be

9  to allow heterogeneity but include this in the model in some manner. A potential concern here is

10 with multiple sets of reference sites. If the multiple sets are treated as a single set then the

11 estimates of the variances are likely to be smaller than data collected from different sites. Hence

12 detection of impairment is potentially overly sensitive. A useful strategy with multiple samples

13 from a collection of sites is to try to match the test site with similar reference sites rather than to

14 use all of the sites. This might involve forming clusters of reference sites and matching the test

15 site with a cluster or using an auxiliary set of measurements (such as sediment type). Other

16 potential problems include measuring a single site multiple times and the time of sampling.

17 Selecting reference sites is a difficult and important problem.

18 The data transformation need not be a transformation of individual contaminants but may

19 also be on the set of measurements. For example, it is common to analyze composites of

20 variables rather than individual variables. Two common approaches are to use principal

21 components (PC) or correspondence analysis (CA) to form new variables. These two methods

22 are useful when the dimensionality of variable space is high and variables are highly correlated.

23 In the case of the sediment quality triad, the PC transformation would typically be applied to the

24 sediment toxicity and metal chemical variables, while correspondence analysis axes can be used

25 for species composition data. Another possible data transformation could be computing some

26 univariate index. For example, composition data can be represented by diversity measures or an

27 index of biological integrity.

28 When the normal assumption is not valid a possible approach is to use a distance measure

29 and build a nonparametric estimate of the probability of impairment (details are given below).

30 When a nonparametric method is used, there are also assumptions that must be considered. The

31 nonparametric approach that we use involves a measure of the distance from the site to the

reference sites. A probability model is then developed based on how close reference sites are to each other. For this approach to work we have similar assumptions:

     1. The reference site data are from a common distribution
     2. The reference sites are independent
     3. The distance measure is appropriate for detecting change.

If the data for the reference sites come from a common distribution, then a single distance measure will produce reasonable estimates of how similar the test site is to the reference sites. When there are different sets of reference conditions the distance measure would have to be computed with respect to the different distributions or with respect to the set of reference conditions most similar to the test site. Independence of samples implies that equal weight may be given to each of the samples from the reference sites. The choice of distance measure is a critical step as the distance measure defines the measure of impact. One important consideration in the selection of the distance measure is the weight given to the variables used in computing the distance (Smith 1998).

The choice of how to preprocess the information is critical to the analysis as it defines the deviations that are of interest. One should be aware of the limitations associated with these choices. For example, if information on a large number of variables is collected then one has a better chance of detecting a broad scale impact. If the impact is only observed through one variable, the other information becomes of low utility for the detection of impact. Thus there is a need to have a clear idea of what types of impact are to be detected. For example, if a chemical is only toxic to fish then measuring abundance of benthic macroinvertebrates will not be risk informative. It is critical in selecting an approach to be aware of what types of changes will or will not be well detected in the analysis and how likely the analysis is to detect changes of important magnitude. Methods such as power analysis are useful for evaluating variables and their importance in the decision process. Also transformation of the variables is often needed. For example, chemical data are often collected in environmental studies. The user needs to decide if the original or standardized data are used. If standardized the method of standardization needs to be chosen. Options might include an overall standardization,

1    standardization relative to a reference group or standardization in terms of toxic units.  Choice of

2    standardization will change the magnitude of distances between observations.

3

4    *Data processing: Estimating probability of impact*

5

6            In the data processing step, the test site is compared with reference conditions in order to

7    obtain a measure of the degree of impairment.  This may be an indirect or direct evaluation.  For

8    example, with biotic indices (Smoger and Angermeier 1999), there is often a calibration step in

9    which the metrics that make up the index are scaled based on reference conditions.  A direct

10    evaluation is based on comparing the reference and test measurements.  A statistical approach is

11    to compute the probability of impairment through the use of the distance between reference and

12    test measurements.  The use of distance in ecological and environmental impact assessment has a

13    long history that will not be explored here.  For example, distance from control forms the basis

14    of outlier detection methods and control chart approaches that view water quality analysis as a

15    quality control problem.  Distance also is central to many multivariate methods used to assess

16    ecological change such as correspondence analysis (using chi-square distance, Legendre and

17    Legendre, 1998) or multidimensional scaling (Smith et al, 1990, Gray et al, 1990).

18            The estimation of the probability of impact is based on comparison of the site in question

19    with reference sites and for certain statistical models, the probability is directly related to

20    distance from the reference.  For example, with a single measurement and mode of evaluation,

21    the calculation of WOE requires estimation of the distribution and calculating the probability that

22    the observation is from that distribution.  In the simplest case, the probability distribution of the

23    measurement is treated as a normal distribution.  Then the two-tailed probability may be

24    calculated by first computing

25

26                                $$d(x_0, \mu) = (x_0 - \mu) / \sigma$$

27

28    Next, the quantity is used to calculate the probability.  Since the parameters are unknown, they

29    are estimated so one would use

30                                $$\hat{d}(x_0, \mu) = (x_0 - \bar{x}) / s$$

31

1    This idea then suggests a more general approach: first transform the measurement to try to obtain

2    a normal distribution for the reference site data, then use the estimated parameters from the

3    normal distribution to compute the distance and then the probability. This approach is useful in

4    that the direction of deviation may be taken into account. The disadvantage is that when there

5    are several variables, the individual estimates may be correlated and this correlation is not taken

6    into account in the distance measure.

7        The possibility of correlated measurements suggests a multivariate distance measure be

8    used. Suppose the measurements at a particular site are given by $\mathbf{x_0}$ where the convention of

9    bold indicates a vector or set of measurements. This vector may be chemical measurements,

10   toxicity information or biological information. The measurements are typically manipulated or

11   transformed in such a way that the distribution of similar vectors for the reference sites follows a

12   normal distribution. Then the squared distance is calculated as the Mahalanobis (Rencher 1995)

13   distance

14

15 $$d^2(\mathbf{x_0}, \boldsymbol{\mu_R}) = (\mathbf{x_0} - \boldsymbol{\mu_R})' \Sigma_R^{-1} (\mathbf{x_0} - \boldsymbol{\mu_R})$$

16

17   where $\boldsymbol{\mu_R}$ is the mean for the reference sites and $\Sigma_R$ is the covariance matrix for the reference

18   sites. The squared distance is estimate by (provided the variance covariance matrix is invertible)

19

20 $$\hat{d}^2(\mathbf{x_0}, \boldsymbol{\mu_R}) = (\mathbf{x_0} - \overline{\mathbf{x}}_{\mathbf{R}})' S_R^{-1} (\mathbf{x_0} - \overline{\mathbf{x}}_{\mathbf{R}})$$

21

22   where the sample mean from the reference sites $\overline{\mathbf{x}}_{\mathbf{R}}$ and the sample covariance matrix $S_R$ are

23   substituted into the distance measure. This distance measure gives less weight to variables that

24   have high variance. Similarly, highly correlated variables do not contribute as much to the

25   variance as uncorrelated variables (Rencher 1995).

26       From the distance measure a probability is calculated. This is easy for multivariate

27   normal data as the probability is determined from the distance to the mean. Other distance

28   measures are also feasible. The advantage to using the Mahalanobis distance is the relationship

29   with the multivariate normal model. From the distance measure, the probability of obtaining a

30   more extreme measurement is easily obtained (as described below). The disadvantage to the

1    distance is that extreme values may occur from extremely good as well as extremely bad sites

2    since the comparison is with the mean.  Another possible measure would be

3

4    $$d^2(\mathbf{x_0}, m_\mathbf{R}) = (\mathbf{x_0} - m_\mathbf{R})' \Sigma_R^{-1} (\mathbf{x_0} - m_\mathbf{R})$$

5

6    where $m_\mathrm{R}$ is the smallest possible value in the reference set (such as the detection limit or zero).

7    Then the distance would not be large for sites with small concentrations of contaminants.  This

8    approach would make sense when the mean does not represent optimal or good conditions.  A

9    model based estimate of probability associated with this distances is difficult to obtain although a

10   model free estimate (nonparametric) is possible given an adequate sample size for reference

11   sites.  When the lower value is zero analytical approaches are possible. The distribution of

12   $d^2 = \mathbf{x_0}' \mathbf{S}_R^{-1} \mathbf{x_0}$ is related to the non-central $T^2$ distribution and the non-central F distribution so

13   probabilities may be evaluated when the observations are normal.

14   The impairment probabilities for each line of evidence can be computed based on the

15   multivariate normal model or, if sample sizes are sufficient, through a data based estimate.  If

16   multivariate normality is assumed, then the probability associated with the distance measure may

17   be calculated using a chi-squared approximation or through Hotelling's $T^2$ distribution (Rencher

18   1995). The squared distance is estimated by

19

20   $$\hat{d}^2(\mathbf{x_0}, \boldsymbol{\mu_\mathbf{R}}) = (\mathbf{x_0} - \overline{\mathbf{x}}_\mathbf{R})' S_R^{-1} (\mathbf{x_0} - \overline{\mathbf{x}}_\mathbf{R})$$

21

22   where $\overline{\mathbf{x}}_\mathbf{R}$ , is the mean for the reference sites and $S_R$ is the within cluster variance-covariance

23   matrix for the reference sites (when no clusters are specified, we of course just use the common

24   mean and variance matrix). The probability of impairment is then defined as one minus the right

25   tail probability of the chi-squared distribution with number of degrees of freedom equal to the

26   number of variables. The probability of impairment is then approximated based on the chi-

27   squared distribution of $X^2 = \dfrac{n}{n+1}(\mathbf{x}_o - \overline{\mathbf{x}}_R)' \Sigma_R^{-1}(\mathbf{x}_o - \overline{\mathbf{x}}_R) \sim \chi_r^2$ where $\Sigma_\mathrm{R}$ is the true covariance

28   matrix for the reference set, and $r$ is number of degrees of freedom equal to the rank of $\Sigma_\mathrm{R}$

29   (typically equal to the number of variables, $p$)

1

2 A more accurate estimate of the probability is obtained using Hotelling $T^2$ distribution. The

3 statistic used is $T_{obs}^2 = \dfrac{n}{n+1}(\mathbf{x}_o - \overline{\mathbf{x}}_R)'\mathbf{S}_R^{-1}(\mathbf{x}_o - \overline{\mathbf{x}}_R)$, where $k$ is number of groups and $n$ is the

4 number of reference sites. To compute the p-values we can use the relationship between $T^2$ and F

5 distribution

6
$$T_{p,n,\alpha}^2 = \frac{p(n-k)}{n-p-k+1}F_{p,n-p-k+1,\alpha}$$

7

8 (Mardia et al. 1979, p. 74). To estimate the probability of impact we first compute

9

10
$$F_{obs} = \frac{n-p-k+1}{p(n-k)}T_{obs}^2$$

11

12 then calculate the probability of a more extreme $F$ value using the $F_{p,n-p-k+1}$ distribution (i.e

13 compute the probability of impairment as $\Pr(\text{impairment}) = \Pr(F_{p,n-p-k+1} > F_{obs})$.

14     An alternative approach is to use a nonparametric estimate of the probability based on the

15 empirical distribution of distances in the reference set. If there is a large enough set of reference

16 sites then an individual reference site may be compared with the other reference sites to build an

17 empirical or data based probability model. The idea is to treat one reference site as a test site and

18 compare it with the other reference sites. By repeating this for all reference sites a distribution of

19 distances is obtained and then the actual test site may be evaluated with that distribution. The

20 distances are defined as above except that the cluster means and within cluster variance matrix

21 are computed with removal of the observation for which the distance is computed. Thus, for the

22 jth observation, we compute

23

24
$$\hat{d}^2(\mathbf{x}_j, \boldsymbol{\mu}_R) = (\mathbf{x}_j - \overline{\mathbf{x}}_{R,-j})' S_{R,-j}^{-1}(\mathbf{x}_j - \overline{\mathbf{x}}_{R,-j})$$

25

26 where the term $-j$ indicates that the observation was removed from the reference set prior to the

27 computation of the parameter estimate. Although the statistically savvy reader will recognize

28 this analysis as computationally intensive, it is possible by employing the fast updating formulas

1    for the variance-covariance matrix and its inverse (for details see Thisted 1988, p. 52). The

2    probability of impairment is computed as the proportion of reference sites whose "leave-one-out"

3    distance is smaller than that of the test site in question.

4

5    *Combining estimates*

6

7    Given estimates of probability of impairment for each line of evidence the problem now

8    becomes combining these together to produce a single weight of evidence estimate. We again

9    assume that the estimate of impact is based on the reference conditions. There are several

10   options available for making the combined estimate. Two approaches involve combining the

11   probabilities for the lines and combining the odds for each of the lines. There are several

12   possibilities for combining information across the different lines by combining the probabilities.

13   The general equation for this over the *k* lines of measurement is calculated as the probability of

14   the union of impairment events (in sediment toxicity *k*=3 corresponding to lines of evidence

15   associated with toxicity, chemistry and biology). A site is declared impaired if it is impaired on

16   any of the modes. Thus,

17

18   $$P(\text{site impaired}) = 1 - \prod_{line=1}^{k} (1 - P(\text{site impaired, line i}))$$

19

20   For our weight of evidence problem,

21           P(site impaired) = 1 – P(site not impaired on any of the 3 lines of evidence)
22                   = 1 – { [1-P(*species composition" impaired"*)]
23                       x [1-P(*toxicity" impairment"*)]
24                       x [1-P(*chemical " impairment"*)]
25

26

27   The formula indicates that the overall probability of site impairment is based on first calculating

28   the probability that the site is not impaired over the different lines, multiplying these together to

29   get an overall probability the site is not impaired, then subtracting from one to get the probability

30   of impairment. Note that with two lines this is the more familiar equation that gives the

31   probability as the union of events (under independence)

P(site impaired) = P(impaired line 1) + P(impaired line 2) – P(impaired both lines)

and is based on computing the individual probabilities and subtracting off the joint probability. There are also some variations on this approach that might be used, similar to that used to estimate joint toxicity using individual estimates of toxicity. The above product approach works if the measurements are independent, which is a difficult assumption. One adjustment is to make bounding assumptions. The most extreme estimate of impairment is to assume no intersection and results in simply a sum i.e.,

$$P_{sum}(\text{site impaired}) = \sum_{lines} P(\text{site impaired, line i})$$

A drawback to the use of this equation is that the probability may exceed one. An average represents an alternative. The least extreme case occurs when there is complete dependence. In this case, one or more of the events is contained inside another event. The estimate of the impairment probability is then greater or equal to the maximum probability across the lines

$$P_{max}(\text{site impaired}) = \max(P_i)$$

where $P_i$ is the probability of impairment on line i.

Another approach is based on combining odds ratios. We can convert the individual probability of impairment from each line into an odds ratio, $O_i = \dfrac{P_i}{1 - P_i}$ and compute the final odds ratio by using the Bayes rule of updating evidence as the product of the odds ratios corresponding to each line, i.e., $O = \prod_{i=1}^{k} \dfrac{P_i}{1 - P_i}$. The advantage of this approach is that it allows one to easily incorporate the prior odds of impairment ($O_{prior}$) of the site by simply multiplying $O = O_{prior} \prod_{i=1}^{k} O_i$. We can account for the effect that the evidence from different sources (lines)

1    may overlap by introducing weights as follow. Let $O = O_{prior} \prod_{i=1}^{k} O_i^{W_i}$ , where the weight $W_i$ for a

2    given line should be the number between 0 and 1, with smaller values reducing the influence of

3    the corresponding odds ratio. The weight factors can be computed as

4    $W_i = 1 - R^2(D_i \mid D_1,.., D_{i-1})$ , where $R^2(D_i \mid D_1,.., D_{i-1})$ is the coefficient of multiple correlation

5    for the regression of the reference distance matrix (extended in a single vector with only upper

6    diagonal elements) associated with $i$-th line regressed on similarly defined distance matrices

7    associated with the lines already accounted for in the combined odds ratio product. Therefore,

8    the weights will account for the relationships among various sources of evidence and effectively

9    down-weight the evidence that does not add much to the information about the disparity between

10   the test and reference sites that was already taken into account. The computations obviously will

11   depend on the order of lines of evidence in the product of odds ratios; therefore we have to use

12   some expert information about relative importance of different lines. For example, with our three

13   lines of comparison, we can use first the species composition data, then toxicity and finally

14   chemistry data.

15        In cases where univariate scores are calculated for each measurement within a line of

16   evidence, these would first have to be combined.  Estimating impairment probability based on

17   multivariate distances makes deviations from the center of reference sites (or the reference

18   clusters) equally unfavorable in any direction. In many cases however, the underlying univariate

19   measures (for example in the case of metals data) have well-defined directions that would cause

20   impairment of the sites. For such a data set we can compute the impairment probability as *one*

21   *minus upper tail probabilities* from the underlying univariate tests based on normal

22   approximation for the quantities $\hat{d}(x_0, \mu) = (x_0 - \bar{x})/s$ . These univariate probabilities are then

23   combined using the intersection principle

24

25          $$P(\text{site impaired}) = 1 - \prod_{i=1}^{m} (1 - P(\text{site impaired, variable } i))$$

26

27   This probability however should be adjusted so as to take into account the effect of multiple

28   testing as follows

29

1 $$P_{adj}(\text{site impaired}) = 1 - \chi^2_{2m}(-2Log(1\text{-}P(\text{site impaired}))) ,$$

2

3 where $\chi^2_{2m}$ is the cumulative distribution function of a chi-squared distribution with *2m* degrees

4 of freedom. This is a standard meta-analytic procedure based on the fact that the null

5 distributions of individual *p*-values from *m* independent tests are independent *Uniform*(0,1)

6 random variables (for details see Hedges and Olkin 1985, p.37).

7

8 **Example**

9

10       As an example we consider data from the Great Lakes that was obtained via the BEAST

11 software (Reynoldson *et al*. 1998; 2000). The reference data consisted of 146 reference samples

12 collected in 1992. Although more data are available on different lines of evidence, these data

13 have information from all three lines of evidence. In addition, there were twenty-five samples

14 taken from Collingwood Harbour that we use as the predictive or test sample. These twenty-five

15 samples were taken at nine locations in 1992, 1995 and 1997 (Table 1). Collingwood Harbour is

16 located in the south shore of Nottawasaga Bay, in the southern extension of Lake Huron's

17 Georgian Bay. The site is of interest as it was identified as an Area of Concern (AOC) by the

18 International Joint Commission but was then de-listed in November 1994, following remediation

19 (for details, see http://www.on.ec.gc.ca/glimr/raps/huron/collingwood/intro.html). Contaminants in the

20 sediment resulted from use of the harbor as a location for ship repair with greatest contamination

21 near the shipyard and in the east and west slips. The nine sampling locations are located as

22 follows: 6703, 6704 and 6705 are located in the harbor with 6703 being farthest from the

23 shipyard and 6705 closest. Sites 6706-6708 are located in the east slip and sites 6709-6711 are

24 located in the west slip. For a map of the locations and additional details on remediation history

25 see http://www.ijc.org/boards/wqb/cases/collingwood/collingwood.html.

26       Several analyses are provided below. The different analyses are compared to illustrate

27 how different assumptions will lead to different estimates of impairment. This is to show the

28 importance of the assumptions and models used for the analysis. First, we consider comparisons

29 of test sites with the entire group of reference sites. No transformations are made. The second

30 model uses transformed chemistry variables. The third model is similar to the second but uses

31 scores from the first two (nontrivial) correspondence analysis axes for the biological data.

1    Additional models are based on principal components reductions, transformations of the toxicity

2    data and empirical measures of probability.  Descriptions of each model are in Table 1.  We first

3    comment on the choices that were made.

4

5    *Chemical Data*

6

7          Graphical displays of the chemical data for the reference sites suggested the data were not

8    normal.  Distributional plots suggested skewness of the distributions and odd observations.  The

9    chemical data were preprocessed using a log transformation for all variables.  Figure 1 displays

10   the scatterplot matrix for the transformed data.  We note that there is evidence of correlations

11   between the chemicals.  Also, there are some locations that have very low or non-detectable

12   concentrations.  This suggests groupings may be useful for analysis.  Two approaches for further

13   processing the data were considered.  The first uses the Mahalanobis distance measure and the

14   multivariate normal distribution.  An alternative is to first combine the chemical data into a

15   single number.  This may be done using toxic units or a principal components analysis.  In the

16   principal components analysis, a linear combination is formed and the combination used as the

17   principal variable for analysis.  The data from the Great Lakes **reference** sites were standardized

18   treating the data as one group, to give each chemical equal weight.  A principal components

19   analysis of the log transformed standardized chemistry data resulted in two components.  The

20   two components explained 75% of the variation in the data.  A plot of the loadings is given in

21   Figure 2 and suggests the first component consists of all the chemicals except arsenic while the

22   second was comprised primarily of arsenic.  Note the loadings are not as high for lead and

23   cadmium.  The model might not accurately describe impairment probability at sites with high

24   levels of these contaminants.  The test sites were scored on the components using the

25   standardization from the reference sites.  The site scores are graphed in Figure 3.  A symbol is

26   used to separate the test sites from the references sites.  In addition, the reference sites were

27   divided into six different groups and the group is also indicated in the graph.  Note that most of

28   the test sites fall to the upper end of the plot suggesting extreme chemical composition.  Also

29   there is a split in the display of points.  This is associated with the arsenic variable that is at the

30   limits of detection for many of the sites.

31

1 *Biological Data*

2

3       There were several possibilities considered for the analysis of the biological data. First

4 one could use the original or transformed data and then use the distance measure approach.

5 Second, the data could be converted to an index then the index used with the distance approach.

6 A third approach is to compute multivariate axes using correspondence analysis and then

7 compute distances using scores from these axes. In the analysis below, we used the first

8 approach only in model 1 and third approach in the remainder of the models. The reason for

9 discarding the first approach was that the effect of the contaminants is to decrease abundance of

10 the majority of the taxa. Due to large variances, the method is only capable of detecting large

11 increases rather than large decreases. A more detailed description is described below.

12

13 *Toxicity Data*

14

15       Plots of the toxicity data suggested skewness of the measurements. As with the chemical

16 data, the log transformation greatly reduced the skewness. We did not consider combining the

17 toxicity data into a single component, as the correlations between the different tests were not

18 high. Analysis was based on $T^2$ tests or distances with probabilities computed based on

19 empirical distances. We also considered analyses using log transformed variables.

20

21 *Combined Estimates*

22

23       Table 1 compares overall impairment estimates from six different models. What is

24 apparent is that considerable differences in the estimates can occur as a result of different

25 analysis options. For example, for the first sample, the probability is estimated as zero using the

26 first model while the second model estimates probability of impairment as close to 1.0. The

27 difference in estimates is a result of the way the species data are treated. In the first model, the

28 biological data are treated as normal while in the second model, the biological data are first

29 transformed using correspondence analysis and the first two components used in the analysis.

30 The new variables from the correspondence analysis are treated as having a normal distribution

31 and the $T^2$ approach applied. It is instructive to compare the individual estimates of impairment

1    probability from the different lines of evidence for the two models.  For model 1, the individual

2    estimates are 0.99 for toxicity, 0.49 for chemistry and 0.00 for community structure.  Note that

3    even though the estimate for toxicity is practically 1.0, the community structure estimate is close

4    to 0 and dominates the weight of evidence analysis.  For model 2, the estimates are 0.99 for

5    toxicity, 0.43 for chemistry and 0.25 for community structure.  In this case, the toxicity measure

6    dominates the analysis and results in an overall estimate of 0.99. The difference in the overall

7    estimate is mostly a result of the contribution from the biological data.  This example also

8    illustrates the importance of considering the individual components as well as the overall

9    estimate.  The overall estimate for model 2 reflects a different view of the biological data.  The

10   distance measure based on untransformed data (model 1) will detect change if the abundance at

11   the test site is different from the reference mean relative to the variance.  Variance in abundance

12   tends to be relatively large for most taxa hence detection of change is made difficult.  The

13   coefficient of variation for individual taxa abundance ranges from 1.5 to 15 with a mean of 7.5.

14   Since the standard deviation is greater than the mean, it will be difficult to detect a difference

15   unless there is a large change in abundance for a single taxon.  Since it might be expected that

16   the effect of metals contamination is to reduce abundance, the problem is made more difficult as

17   zero is less than one standard deviation away from all the mean abundances.  For example, the

18   high probability for 6704 in 1997 results from high abundance of *Tubificidae*.  The high value

19   for 6711 in 1992 results from a high value of *Caenidae*.  Interestingly the mean abundance at this

20   site is only 0.6.  While this is not high, the mean for that taxon in the reference sites is only 0.11.

21   Sample 6708 in 1992 (which is similar to 6711 in 1992) is not considered impaired even though

22   almost all taxa have zero abundance (only two taxa are present)!  Thus, this modeling approach

23   may be useful for detecting effects that cause an increase in abundance (perhaps due to sewage

24   or deposition) but is poor for detecting change due to toxicity.

25       Correspondence analysis tends to produce better results, as it will produce "average"

26   taxa. Because the primary effect of the metals is reduction in abundance, the effect is more

27   readily detected using the correspondence analysis scores, as the sites with low abundance tend

28   to have smaller scores than the correspondence analysis reference site scores.

29       Another factor that seems to be important is whether or not the variables are transformed.

30   For example, compare samples 6704 in 1997 and 6706 in 1995 for models 3 and 3a.  The

31   difference in models is only a transformation of the toxicity data.  Thus the increase in

1   probability results from higher probability from the toxicity results.  In the first case, the

2   probability associated with the toxicity changes from 0.76 to 1.0 and in the second from 0.62 to

3   0.99.  The increased probability of impairment from the toxicity line of evidence results in a

4   large change in the overall probability of impairment for these two samples.  Similar results hold

5   when chemistry variables are transformed. Other modeling changes do not seem to have as great

6   an effect on the estimates.  For example, model 4 uses a subset of the reference sites for

7   comparisons.  Estimates for this model are not that different from those for model 3.  The use of

8   empirical estimates of probability rather than model based (i.e. based on $T^2$) results in a moderate

9   change for some of the samples.  The use of principal components in the analysis of the metals

10  data can lead to differences.  For example, sample 6706 in 1995 suggests much lower

11  impairment using the component approach than the full variables approach (compare estimates

12  for model 2 and 3).

13         A difference in the estimates from different models suggests the importance of

14  considering how different lines of evidence contribute to the estimate. Probabilities of

15  impairment for the three lines of evidence and the combined estimate are displayed in Figure 4

16  for the twenty-five samples in the test set using model 6. In model 6, the probabilities for metals

17  were calculated by combining probabilities calculated from principal component (two

18  components) distances.  The distances were calculated by first selecting a cluster of reference

19  observations based on covariates and the estimates of probability were based on distances within

20  the reference cluster.  For the toxicity data we calculated the multivariate distance and computed

21  the probability based on the overall empirical distance distribution using all reference distances.

22  Finally, for the biological (community structure) data, correspondence analysis was used and

23  distance and probability evaluated on the first two axes using the empirical distance approach.

24  Figure 4 displays the overall impairment probabilities as well as estimates from individual lines.

25  We note that for most of the sites, the probability of impairment is high.  Exceptions are for 6703

26  in 1995, 6704 in 1992, 6704 in 1995, 6705 in 1995 and 6706 in 1995.  Sites 6707 through 6711

27  appear impaired in all years.  The individual impairment probabilities are highest for metals and

28  sediment toxicity tests.  The 1992 estimates are variable across the different lines of evidence.

29  The biological impairment probabilities are not generally high and perhaps are a result of not

30  comparing abundance with a more closely matched subset of reference sites.  Most of the metal

31  data indicate high levels for 6703 and 6704 but are generally high for the other sites.

1     The results given here are summaries of the different components of each line of

2     evidence.  A further, valuable component of an analysis would be to study which components

3     were important to the individual line of evidence.  For example, the ten toxicity tests are all not

4     equally important to the impairment estimate for toxicity.  One approach would be to remove

5     individual components then to evaluate the effect on the impairment probability.  In this way the

6     components may be evaluated in terms of importance to the probability estimate (see for an

7     example, Smith *et al*. 1990).

8

9     **Discussion**

10

11     We have presented an approach for estimation of the probability or risk of impairment for

12     a site based on multiple lines of evidence.  Many variations on the approach are possible based

13     on different distances and different summarization of the information.  Decisions about distance

14     measures and summarization are best made in the initial stages of the study. Although the results

15     of the weight of evidence analysis summarize impairment in terms of a single number, this

16     approach is generally restrictive.  A single measure attempts to summarize the multivariate

17     degree of impairment.  It is certainly possible that several different scenarios will lead to the

18     same or similar measure of impairment.  Hence it is important to consider the individual lines of

19     evidence as well as the data themselves.  Graphical display of the data is necessary to check for

20     problems and assess assumptions.  Biological and environmental evaluation of the collected data

21     is also necessary to verify that the evaluation is scientifically correct as well as statistically valid.

22     We envision the above approach to be most useful for comparing and ranking sites.  As

23     Figure 4 illustrates, it is possible to display the estimates of impairment for different lines and a

24     combined estimate for several sites/times in a single display.  The information plotted may be

25     ordered in time or space to look for change in impairment probabilities.  This might be useful for

26     restoration/recovery studies.  The information for different sites may also be compared to

27     identify sites with greatest risk or trends in impairment.  For example, if the sites were located

28     along a toxic gradient, then one would expect to see increases in the toxicity estimate of

29     impairment and this could be displayed on the diagram.

30     Although we have focused on a statistical approach based on estimation of the probability

31     of impairment, other approaches are available for obtaining a combined estimate of impairment.

1    One common approach is an index-oriented method.  In this approach, the numerical values are

2    combined, possibly after a transformation.  One example of this approach is Wildhaber and

3    Schmitt (1996).  They combine data over biological, toxicological and chemical lines.  To do this

4    the chemistry data are standardized by calculating the ratio of the bioavailable component of the

5    contaminant to the chronic toxicity water quality criterion.  The toxicological data are

6    standardized by adjusting the test endpoint for the control endpoint.  Values may then be

7    averaged over each of the lines of evidence.  The biological data are not evaluated in that

8    manner, rather tolerance values are used and a biological index is computed as a tolerance

9    weighted average.  The values are then combined over the three lines of evidence.  A variation

10   on this approach is given in Soucek *et al.* (2000) and Cherry *et al.* (2001).  In these papers the

11   data are not combined within each line of evidence.  Rather, important metrics from each line are

12   selected then the metrics are combined into an overall index.

13        An important aspect of the analysis is the choice of distance measure.  Our approach is to

14   use a distance measure that is directly related to a probability distribution.  It is therefore

15   important that the assumptions of the distribution be checked so the distance measure is directly

16   related to probability.  It is also important that the distance measure reflects impairment in that

17   the farther away from the center, the more impaired the site.  One approach to achieve this is to

18   base distance on ordinated data rather than on actual observations.  For example, when dealing

19   with chemical data an approach is to calculate principal components for the data from the

20   reference sites.  Then decide if the location of the test site is indicative of impairment.  A

21   distance measure that reflects a directional distance would be appropriate.  Distance might be

22   calculated as zero if the site is on the safe side of the distribution and the ordinary distance

23   measure used if the site is located on the not-safe side.  Principal components would be

24   computed based on standardized data to give all the chemicals equal weight in the derivation of

25   the components.

26

27    Acknowledgements

28

1

**References**

2

3

4  Chapman PM. 1996. Presentation and interpretation of Sediment Quality Triad data.
5       Ecotoxicology 5: 327-339

6

7  Cherry DS, Currie RJ, Soucek DJ, *et al.* 2001.  An integrative assessment of a watershed
8       impacted by abandoned mined land discharges.  Envir Poll 111:377-388

9

10  Gilbert RO. 1987. Statistical Methods for Environmental Pollution Monitoring.  Van Nostrand
11       Reinhold, New York, NY, USA

12

13  Good IJ.  1988. Statistical evidence.  In: Kotz S and Johnson NL (eds), Encyclopedia of Statistics
14       vol. 8, pp 651-655. John Wiley and Sons, New York, NY, USA

15

16  Gray JS, Clarke KR, Warwick RM, Hobbs G. 1990. Detection of initial effects of pollution on
17       marine benthos: An example from the Ekofisk and Eldfisk oilfields, North Sea.  Mar Ecol
18       Prog Ser 66:285-299.

19

20  Hedges LV, Olkin I. 1985. Statistical Methods for Meta-Analysis. Academic Press, London, UK

21

22  Kass RE, Raftery AE. 1995. Bayes factors.  Jour Amer Stat Assoc 90:773-795

23

24  Legendre P, Legendre L. 1990. Numerical Ecology, 2$^{nd}$ ed, Elsevier, Amsterdam, Holland

25

26  Mardia KV, Kent JT, Bibby JM. 1979. Multivariate Analysis. Academic Press, London, UK

27

28  McCullagh P, Nelder JA. 1989. Generalized Linear Models, 2$^{nd}$ ed, Chapman and Hall, London,
29       UK

30

31  Rencher A. 1995. Methods of Multivariate Analysis.  John Wiley & Sons, New York, USA

1

2   Reynoldson, TF, Day, KE. 1998.  Biological guidelines for the assessment of sediment quality in

3       the Laurentian Great Lakes.  NWRI Report 98-232, Dept of Fisheries and Oceans

4       Canada, Burlington, ON, Canada

5

6   Reynoldson TB, Day KE, Pascoe T. 2000 The development of BEAST: A predictive approach

7       for assessing sediment quality in the North American Great Lakes.  In: Wright JF,

8       Sutcliffe DW, Furse MT (eds), Assessing the Quality of Freshwaters: RIVPACS and

9       Other Techniques, Chapter 11, pp 165-180.  Freshwater Biological Association,

10      Ambleside, UK

11

12  Smith EP. 1998. Randomization methods and the analysis of multivariate ecological data.

13     Environmetrics 9:37-51

14

15  Smith EP, Pontasch KW, Cairns J Jr. 1990. Community similarity and the analysis of

16     multispecies environmental data: A unified statistical approach.  Water Res 24:507-514

17

18  Smogor RA, Angermeier, PL. 1999. Relations between fish metrics and measures of

19     anthropogenic distrurbance in three IBI regions in Virginia. In: Simon TP (ed), Assessing

20     the Sustainability and Biological Integrity of Water Resources Using Fish Communities,

21     pp.585-610. CRC Press, Boca Raton, FL, USA

22

23  Soucek DJ, Cherry DS, Currie RJ, *et al*. 2000. Laboratory to field validation in an integrative

24     assessment of an acid mine drainage-impacted watershed.  Environ Toxicol Chem

25     19:1036-1043

26

27  Thisted RA. 1988. Elements of Statistical Computing. Chapman and Hall, London, UK

28

29  Wildhaber MW, Schmitt CJ. 1996. Hazard ranking of contaminated sediments based on

30     chemical analysis, laboratory toxicity tests, and benthic community composition:

31     Prioritizing sites for remedial action.  J Great Lakes Res 22:639-652

1
2

Figure 1.  Scatterplot matrix of log transformed chemical reference measurements.  The

following abbreviations are used: LAS=log(arsenic), LCD=log(cadmium), LCR=log(chromium),

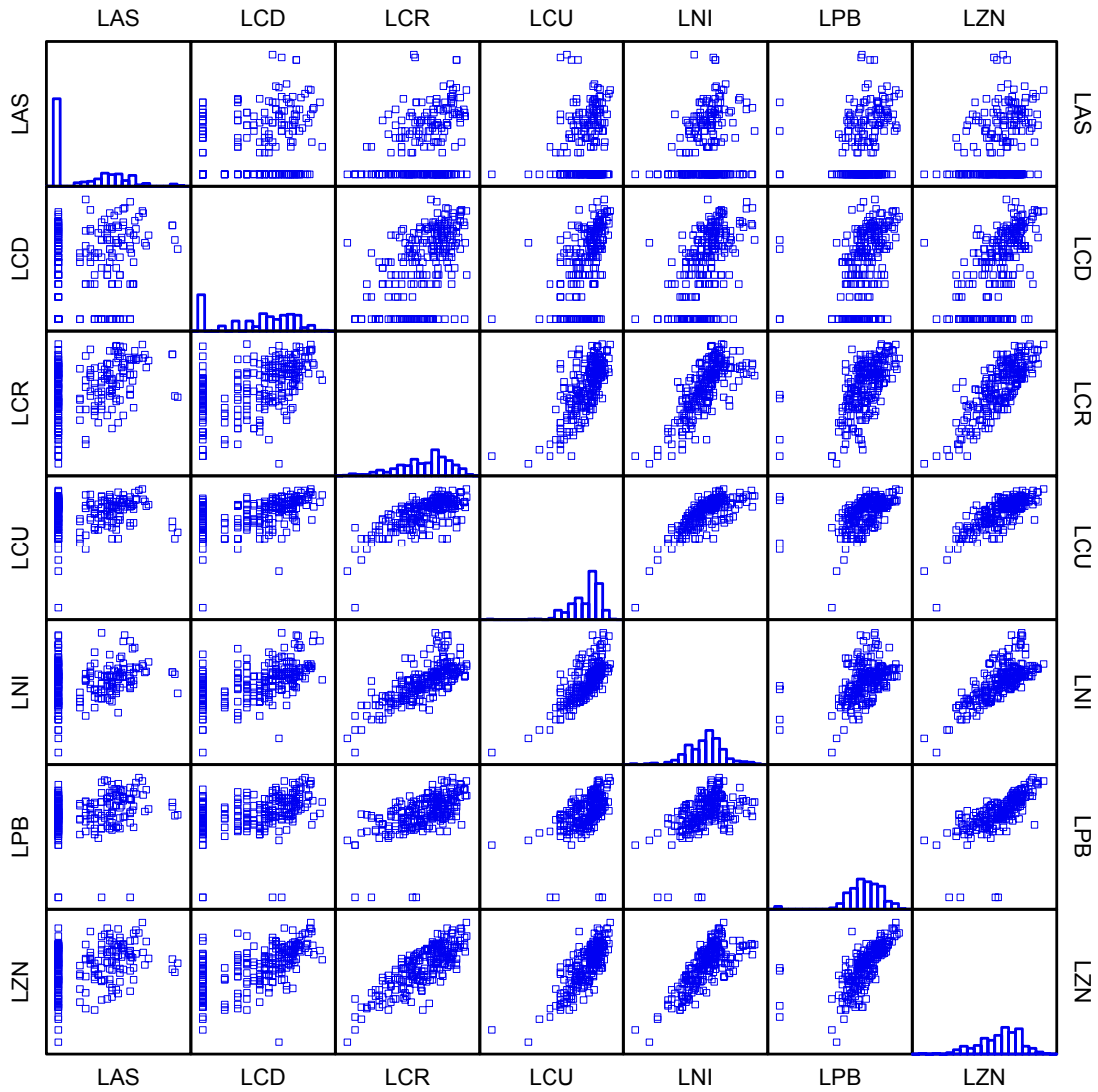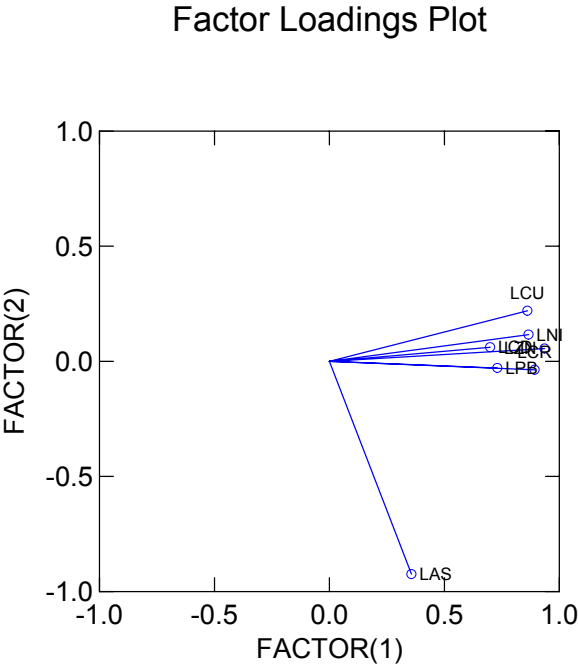LCU=log(copper), LNI=log(nickel), LPB=log(lead) and LZN=log(zinc).

1    Figure 2. Plot of loading for the principal components for the metal data.  The coordinates (loadings) are:

2    LAS (0.36, 0.92), LCD (0.70, -0.06), LCR (0.89, 0.04), LCU (0.86, -0.22), LNI (0.87, -0.11), LPB (0.73,

3    0.03) and LZN (0.93, -0.05).

4

5



Factor Loadings Plot

6

7

8

1

2    Figure 3. Plot of site scores for reference and test sites.  Symbols are coded based on whether

3    they are reference or test (r=reference, test=test) and the number after the "r" indicates the

4    reference cluster.

5



6

1     Figure 4. Plot of probabilities for three lines of evidence and overall estimate for sites in test data set. The center corresponds to zero and tick marks

2     represent tenths. The first four numbers in the site/sample label indicate the site number while the last 2 provide the year of sampling.

3



Legend:
- Odds Ratio Estimate
- Sediment Toxicity (BioAssayQuery, original values, test: Empirical distances)
- Metal Chemicals (Habitat, log transformed values, test: Empirical distances)
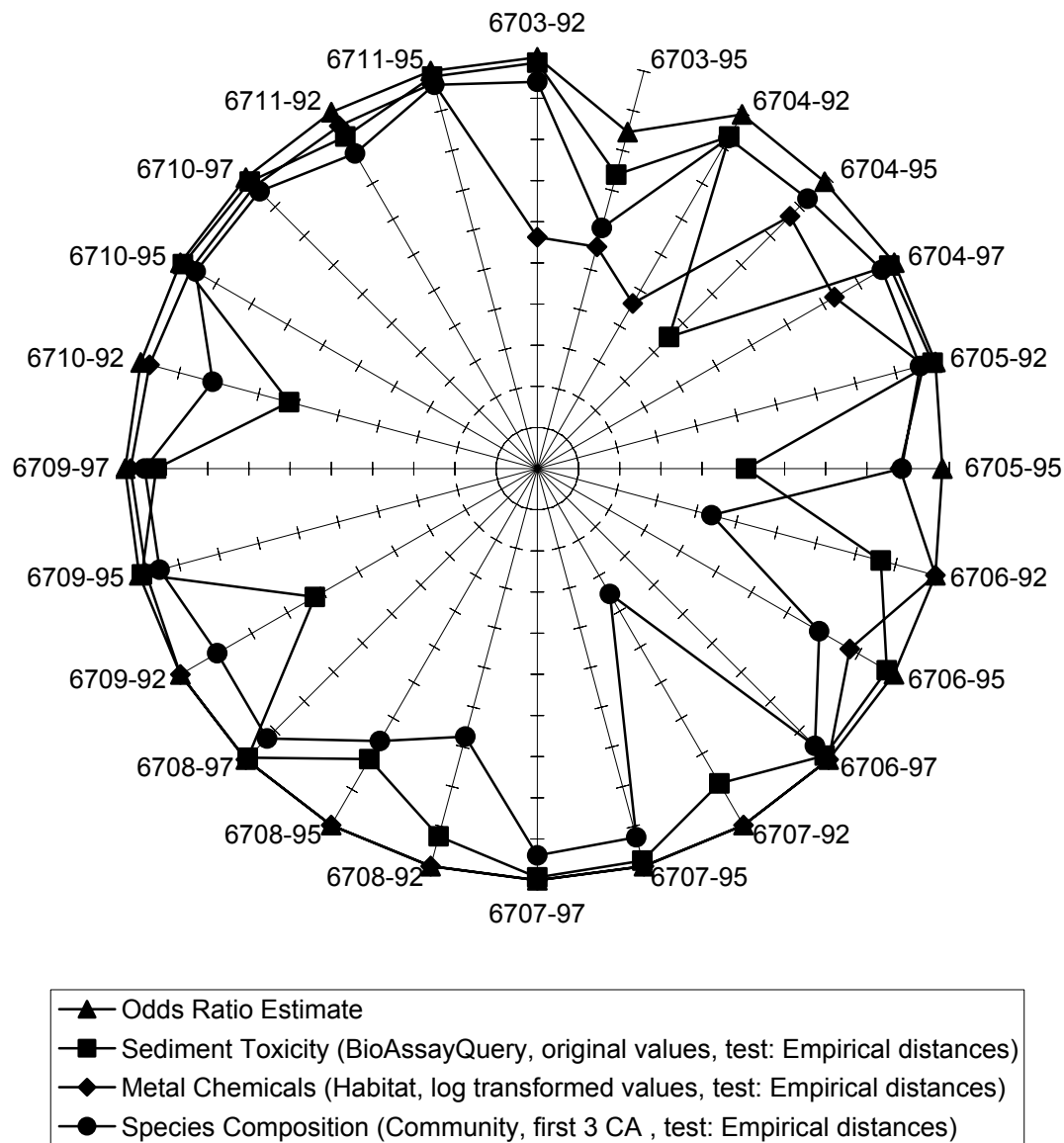- Species Composition (Community, first 3 CA , test: Empirical distances)

1 Table 1. Comparison of impairment probabilities using different analysis options. Hotelling $T^2$

2 used in procedure 1-4, empirical distance approach used for models 5 and 6.

3     Procedure 1: No transformations

4     Procedure 2: Log transformation for metals, correspondence analysis (2 axes) for biological

5        data,

6     Procedure 3: Principal components on log transformed metals (2 axes), correspondence

7        analysis for biological data, no transformation for toxicity data

8     Procedure 4: Same as 3 except that log transformed toxicity data used

9     Procedure 5: Same as 3 except metals data are compared with a subset of the controls

10     Procedure 6: Same as 4 using empirical distances to determine probabilities

11     Procedure 7: Same as 5 but using log transformed toxicity data and 3 correspondence

12        analysis axes.

13

| Site-year | Procedure 1 | Procedure 2 | Procedure 3 | Procedure 4 | Procedure 5 | Procedure 6 | Procedure 7 |
|---|---|---|---|---|---|---|---|
| 6703-92 | 0.0000 | 0.9928 | 0.9910 | 0.9999 | 0.9909 | 0.8602 | 0.9993 |
| 6703-95 | 0.0000 | 0.6670 | 0.6290 | 0.4040 | 0.6412 | 0.6558 | 0.8461 |
| 6704-92 | 0.1535 | 0.8469 | 0.8609 | 0.9712 | 0.8631 | 0.6707 | 0.9933 |
| 6704-97 | 0.9997 | 0.7485 | 0.3586 | 0.9998 | 0.3849 | 0.3527 | 0.9856 |
| 6704-95 | 0.0000 | 0.5664 | 0.0126 | 0.0064 | 0.0146 | 0.0210 | 0.9999 |
| 6705-92 | 0.0000 | 1.0000 | 0.9997 | 1.0000 | 0.9996 | 0.9887 | 1.0000 |
| 6705-95 | 0.0000 | 0.5570 | 0.1745 | 0.1745 | 0.1624 | 0.3282 | 0.9829 |
| 6706-92 | 0.0000 | 1.0000 | 0.9998 | 1.0000 | 0.9998 | 0.9989 | 1.0000 |
| 6706-97 | 0.0005 | 1.0000 | 0.9974 | 1.0000 | 0.9972 | 0.9882 | 0.9992 |
| 6706-95 | 0.0000 | 0.8196 | 0.0184 | 0.9449 | 0.0095 | 0.0068 | 1.0000 |
| 6707-92 | 0.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 0.9995 | 1.0000 |
| 6707-97 | 0.0000 | 1.0000 | 0.9953 | 1.0000 | 0.9948 | 0.9902 | 1.0000 |
| 6707-95 | 0.0000 | 1.0000 | 0.9983 | 1.0000 | 0.9981 | 0.9951 | 1.0000 |
| 6708-92 | 0.0000 | 1.0000 | 0.9998 | 1.0000 | 0.9998 | 0.9992 | 1.0000 |
| 6708-97 | 0.9976 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 0.9998 | 1.0000 |
| 6708-95 | 0.0001 | 1.0000 | 0.9999 | 0.9999 | 0.9999 | 0.9984 | 1.0000 |
| 6709-92 | 0.0000 | 1.0000 | 0.9009 | 0.9148 | 0.8937 | 0.9468 | 1.0000 |
| 6709-97 | 0.0000 | 1.0000 | 0.9762 | 0.9957 | 0.9735 | 0.9653 | 1.0000 |
| 6709-95 | 0.0000 | 1.0000 | 0.9732 | 1.0000 | 0.9699 | 0.9407 | 0.9999 |
| 6710-92 | 0.0000 | 0.9972 | 0.6238 | 0.7503 | 0.5986 | 0.9019 | 0.9964 |
| 6710-97 | 0.9500 | 0.9999 | 0.9844 | 0.9999 | 0.9827 | 0.9377 | 1.0000 |
| 6710-95 | 0.0002 | 0.9998 | 0.9363 | 1.0000 | 0.9231 | 0.7486 | 1.0000 |
| 6711-92 | 1.0000 | 0.9875 | 0.6084 | 0.9936 | 0.5837 | 0.7876 | 0.9996 |
| 6711-97 | 0.9977 | 0.9999 | 0.7001 | 0.5816 | 0.6761 | 0.8147 | 1.0000 |
| 6711-95 | 0.0000 | 0.9997 | 0.9298 | 0.9996 | 0.9165 | 0.7990 | 0.9996 |

14

15

1    Table 2. Prediction of impairment from different lines of evidence using empirical distances with

2    log transformed metals, log transformed toxicity data and two correspondence analysis axes

3    (Procedure 7 from Table 1).

| Site-year | Odds Ratio Estimate | Sediment Toxicity (Toxicity, log transformed values, test: Empirical distances) | Metal Chemicals (log transformed values, test: Empirical distances) | Species Composition (Community, first 3 CA , test: Empirical distances) |
|---|---|---|---|---|
| 6703-92 | 0.9993 | 0.9863 | 0.5622 | 0.9399 |
| 6703-95 | 0.8461 | 0.7397 | 0.5579 | 0.6052 |
| 6704-92 | 0.9933 | 0.9315 | 0.4635 | 0.9270 |
| 6704-97 | 0.9856 | 0.4521 | 0.8670 | 0.9270 |
| 6704-95 | 0.9999 | 0.9863 | 0.8326 | 0.9657 |
| 6705-92 | 1.0000 | 0.9932 | 0.9700 | 0.9614 |
| 6705-95 | 0.9829 | 0.5068 | 0.8798 | 0.8841 |
| 6706-92 | 1.0000 | 0.8630 | 1.0000 | 0.4378 |
| 6706-97 | 0.9992 | 0.9795 | 0.8755 | 0.7897 |
| 6706-95 | 1.0000 | 0.9863 | 1.0000 | 0.9528 |
| 6707-92 | 1.0000 | 0.8836 | 1.0000 | 0.3519 |
| 6707-97 | 1.0000 | 0.9863 | 1.0000 | 0.9270 |
| 6707-95 | 1.0000 | 0.9932 | 1.0000 | 0.9399 |
| 6708-92 | 1.0000 | 0.9247 | 1.0000 | 0.6738 |
| 6708-97 | 1.0000 | 0.8151 | 1.0000 | 0.7639 |
| 6708-95 | 1.0000 | 0.9932 | 1.0000 | 0.9270 |
| 6709-92 | 1.0000 | 0.6233 | 1.0000 | 0.8970 |
| 6709-97 | 1.0000 | 0.9932 | 0.9828 | 0.9485 |
| 6709-95 | 0.9999 | 0.9247 | 0.9871 | 0.9528 |
| 671-92 | 0.9964 | 0.6233 | 0.9742 | 0.8155 |
| 671-97 | 1.0000 | 0.9932 | 0.9700 | 0.9571 |
| 671-95 | 1.0000 | 0.9863 | 0.9700 | 0.9528 |
| 6711-92 | 0.9996 | 0.9315 | 0.9614 | 0.8841 |
| 6711-97 | 1.0000 | 0.9863 | 0.9700 | 0.9657 |
| 6711-95 | 0.9996 | 0.6849 | 0.9785 | 0.9657 |

4