

Michael W. Palmer, Daniel J. McGlinn, Lars Westerberg, and Per Milberg. 2008. Indices for detecting differences in species composition: some simplifications of RDA and CCA. *Ecology* 89:1769–1771.

Appendix B: Simplification of CCA eigenvalue.

Notation (modified after Legendre and Legendre 1998):

As with Appendix A, we assume for convenience that we are dealing with fixed plots

sampled twice, although the results pertain to split plot designs and balanced

independent samples with no loss of generality.

Samples are indexed by the letter i : 1, 2, ... n

Pairs of samples or ‘plots’ are indexed by the letter k : 1, 2, ... $n/2$

Species are indexed by the letter j : 1, 2, ... r

Elements y_{ij} of the response matrix \mathbf{Y} ($n \times r$) denote the abundance in sample i (rows) of species j (columns)

y_{+j} = sample totals (sum across rows to get column totals)

y_{i+} = species totals (sum across columns to get row totals)

y_{++} = grand total for sample by species matrix

The explanatory variable \mathbf{x} is a column vector of length n and is a dummy variable where

x_i equals 0 or 1 for times 1 and 2, respectively.

The response matrix \mathbf{Y} is divided into two $\frac{n}{2} \times r$ submatrices \mathbf{G} and \mathbf{H} , where \mathbf{G} contains

the samples corresponding to $x_i = 0$ and where \mathbf{H} contains the samples

corresponding to $x_i = 1$.

p is the proportion of the total abundance in samples where $x_i = 1$, i.e. $p = h_{++}/y_{++}$.

S stands for a dispersion matrix except without dividing by degrees of freedom (see

Legendre and Legendre 1998, p596).

D stands for a diagonal matrix of species totals, y_{i+}

Proof:

The CCA eigenvector equation (Legendre and Legendre 1998, p596) is:

$$(\mathbf{S}_{\bar{\mathbf{Q}}\mathbf{D}^{1/2}\mathbf{x}^*}\mathbf{S}_{\mathbf{x}^*\mathbf{D}\mathbf{x}^*}^{-1}\mathbf{S}_{\mathbf{x}^*\mathbf{D}^{1/2}\bar{\mathbf{Q}}} - \lambda_a\mathbf{I})\mathbf{u}_a = 0 \quad (\text{B.1})$$

This equation is similar to equation A.1 for the RDA test statistic, but important differences exist in how the component matrices are scaled and weighted. First, **Y** must be scaled by its contribution to χ^2 (Legendre and Legendre 1998) into:

$$\bar{\mathbf{Q}} = [\bar{q}_{ij}] = \frac{y_{ij}y_{++} - y_{i+}y_{+j}}{y_{++}\sqrt{y_{i+}y_{+j}}} \quad (\text{B.2})$$

Without loss of generality, we define the sum of all abundances as 1, i.e. $y_{++} = 1$. Then:

$$\bar{\mathbf{Q}} = [\bar{q}_{ij}] = \left[\frac{y_{ij} - y_{i+}y_{+j}}{\sqrt{y_{i+}y_{+j}}} \right] \quad (\text{B.3})$$

The explanatory variable must also be transformed into a weighted standardized variable by subtracting the weighted average and dividing by the maximum likelihood estimator of the standard deviation:

$$\bar{x} = \frac{\sum_i y_{i+}x_i}{y_{++}} = \frac{\sum_i (g_{i+})(0) + \sum_i (h_{i+})(1)}{y_{++}} = \frac{h_{++}}{y_{++}} = \frac{h_{++}}{1} = p \quad (\text{B.4})$$

$$\text{var}(\mathbf{x}) = \sigma_{\mathbf{x}}^2 = \frac{\sum_i [y_{i+}(x_i - \bar{x})^2]}{y_{++}} = \frac{\sum_i [g_{i+}(\bar{x}^2)] + \sum_i [h_{i+}(1 - \bar{x})^2]}{1} \quad (\text{B.5})$$

$$\text{var}(\mathbf{x}) = (p)^2 \sum_i g_{i+} + (1-p)^2 \sum_i h_{i+} = (p)^2(1-p) + (1-p)^2 p = p(1-p) \quad (\text{B.6})$$

$$\text{sd}(\mathbf{x}) = \sigma_{\mathbf{x}} = \sqrt{p(1-p)} \quad (\text{B.7})$$

Therefore the standardized value \mathbf{x}^* is:

$$x_i^* = \frac{x_i - \bar{x}}{\sigma_{\mathbf{x}}} = \begin{cases} \frac{-p}{\sqrt{p(1-p)}}, x_i = 0 \\ \frac{1-p}{\sqrt{p(1-p)}}, x_i = 1 \end{cases} \quad (\text{B.8})$$

Also because CCA relies upon weighted multiple regression instead of conventional multiple regression as in RDA, the square root of a diagonal matrix of weights \mathbf{D} must be applied to \mathbf{x}^* everywhere it occurs in B.1 (Legendre and Legendre 1998, p. 595).

Therefore the inverse of the variance of the weighted standardized \mathbf{x} matrix is a scalar equal to:

$$\mathbf{S}_{\mathbf{x}^* \mathbf{D} \mathbf{x}^*}^{-1} = \frac{1}{\sigma_{\mathbf{x}^*}^2} = \frac{1}{1} = 1 \quad (\text{B.9})$$

Furthermore, in our case, $\mathbf{S}_{\bar{\mathbf{Q}} \mathbf{D}^{1/2} \mathbf{x}^*}$ and $\mathbf{S}_{\mathbf{x}^* \mathbf{D}^{1/2} \bar{\mathbf{Q}}}$ are column and row vectors respectively.

The j elements of $\mathbf{S}_{\bar{\mathbf{Q}} \mathbf{D}^{1/2} \mathbf{x}^*}$ are equal to the j elements of $\mathbf{S}_{\mathbf{x}^* \mathbf{D}^{1/2} \bar{\mathbf{Q}}}$, i.e.:

$$(\mathbf{S}_{\bar{\mathbf{Q}} \mathbf{D}^{1/2} \mathbf{x}^*})' = \mathbf{S}_{\mathbf{x}^* \mathbf{D}^{1/2} \bar{\mathbf{Q}}} = \sum_i [\bar{q}_{ij} (x_i^* \sqrt{y_{i+}})] \text{ where } j = 1, 2, \dots, r \quad (\text{B.10})$$

Expression B.10 simplifies further by expanding its components and partitioning matrix \mathbf{Y} into the two sets of independent samples: \mathbf{G} and \mathbf{H} . Inserting equation B.3 and simplifying yields:

$$\mathbf{S}_{\bar{\mathbf{Q}} \mathbf{D}^{1/2} \mathbf{x}^*} = \sum_i \left[\frac{y_{ij} - y_{i+} y_{+j}}{\sqrt{y_{i+} y_{+j}}} \cdot (x_i^* \sqrt{y_{i+}}) \right] = \sum_i \left[\frac{y_{ij} - y_{i+} y_{+j}}{\sqrt{y_{+j}}} \cdot x_i^* \right] \quad (\text{B.11})$$

$$\mathbf{S}_{\bar{\mathbf{Q}}'\mathbf{D}^{1/2}\mathbf{x}^*} = \frac{\sum_i (y_{ij} - y_{i+}y_{+j}) \cdot x_i^*}{\sqrt{y_{+j}}} \quad (\text{B.12})$$

Simplifying the numerator and partitioning \mathbf{Y} gives:

$$\sum_i (y_{ij} - y_{i+}y_{+j}) \cdot x_i^* = \sum_k (g_{kj} - g_{k+}y_{+j}) \cdot \left[\frac{-p}{\sqrt{p(1-p)}} \right] + \sum_k (h_{kj} - h_{k+}y_{+j}) \cdot \left[\frac{1-p}{\sqrt{p(1-p)}} \right] \quad (\text{B.13})$$

$$= \frac{\sum_k (g_{kj} - g_{k+}y_{+j}) \cdot (-p) + \sum_k (h_{kj} - h_{k+}y_{+j}) \cdot (1-p)}{\sqrt{p(1-p)}} \quad (\text{B.14})$$

$$= \frac{\sum_k \{(h_{kj} - h_{k+}y_{+j}) - p[(g_{kj} - g_{k+}y_{+j}) + (h_{kj} - h_{k+}y_{+j})]\}}{\sqrt{p(1-p)}} \quad (\text{B.15})$$

$$= \frac{\sum_k h_{kj} - y_{+j} \sum_k h_{k+} - p \left[\sum_k g_{kj} - y_{+j} \sum_k g_{k+} + \sum_k h_{kj} - y_{+j} \sum_k h_{k+} \right]}{\sqrt{p(1-p)}} \quad (\text{B.16})$$

$$= \frac{h_{+j} - y_{+j}h_{++} - p[g_{+j} + h_{+j} - y_{+j}(g_{++} + h_{++})]}{\sqrt{p(1-p)}} \quad (\text{B.17})$$

$$= \frac{h_{+j} - y_{+j}h_{++} + y_{+j}h_{++} - p(g_{+j} + h_{+j})}{\sqrt{p(1-p)}} \quad (\text{B.18})$$

$$= \frac{h_{+j} - p(g_{+j} + h_{+j})}{\sqrt{p(1-p)}} = \frac{h_{+j} - py_{+j}}{\sqrt{p(1-p)}} \quad (\text{B.19})$$

Combining equation B.12 and B.19 gives:

$$\mathbf{S}_{\bar{\mathbf{Q}}'\mathbf{D}^{1/2}\mathbf{x}^*} = \frac{h_{+j} - py_{+j}}{\sqrt{p(1-p)y_{+j}}} \quad (\text{B.20})$$

The CCA eigenvector equation (B.1) becomes $\lambda = \mathbf{S}_{\mathbf{x}^*\mathbf{D}^{1/2}\bar{\mathbf{Q}}} \mathbf{S}_{\bar{\mathbf{Q}}'\mathbf{D}^{1/2}\mathbf{x}^*} \mathbf{S}_{\mathbf{x}^*\mathbf{D}\mathbf{x}^*}^{-1}$ after rearranging

(as in appendix A) and combining this expression with B.9 and B.20 we see that:

$$\lambda = \sum_j \left\{ \frac{h_{+j} - py_{+j}}{\sqrt{p(1-p)y_{+j}}} \right\}^2 = \frac{1}{p(1-p)} \sum_j \frac{(h_{+j} - py_{+j})^2}{y_{+j}} \quad (\text{B.21})$$

To express B.21 in terms of ‘raw’ y , we divide each y -term by the grand total, i.e. y_{++} :

$$\lambda = \frac{1}{p(1-p)} \sum_j \frac{\left\{ \frac{1}{y_{++}} h_{+j} - p \frac{1}{y_{++}} y_{+j} \right\}^2}{\frac{1}{y_{++}} y_{+j}} = \frac{1}{y_{++} p(1-p)} \sum_j \frac{(h_{+j} - py_{+j})^2}{y_{+j}} \quad (\text{B.22})$$

LITERATURE CITED

Legendre, P., and L. Legendre. 1998. Numerical Ecology. Second English edition.

Elsevier, Amsterdam, The Netherlands.