

## Lecture 5: Ecological distance metrics; Principal Coordinates Analysis

### Univariate testing vs. community analysis

- Univariate testing deals with hypotheses concerning individual taxa
  - Is this taxon differentially present/abundant in different samples?
  - Is this taxon correlated with a given continuous variable?
- What if we would like to draw conclusions about the community as a whole?

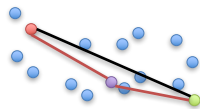
## Useful ideas from modern statistics

- Distances between anything (abundances, presence-absence, graphs, trees);
- Direct hypotheses based on distances;
- Decompositions through iterative structuration;
- Projections;
- Randomization tests, probabilistic simulations.

3

## What is a distance metric?

- Scalar function  $d(.,.)$  of two arguments
- $d(x, y) \geq 0$ , always nonnegative;
- $d(x, x) = 0$ , distance to self is 0;
- $d(x, y) = d(y, x)$ , distance is symmetric;
- $d(x, y) \leq d(x, z) + d(z, y)$ , triangle inequality.

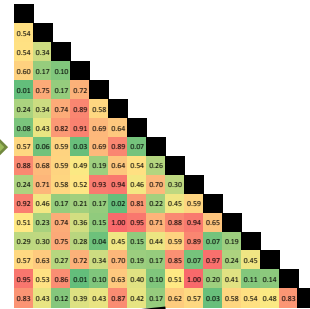


4

Microbiome analysis pipeline:  
 Specimens → Sequence → Abundance Matrix →  
 Distances → Statistics\* → Paper



	S-1	S-2	...	S-16
OTU-1				
OTU-2				
OTU-3				
OTU-4				
OTU-5				
....				
OTU-10,000				

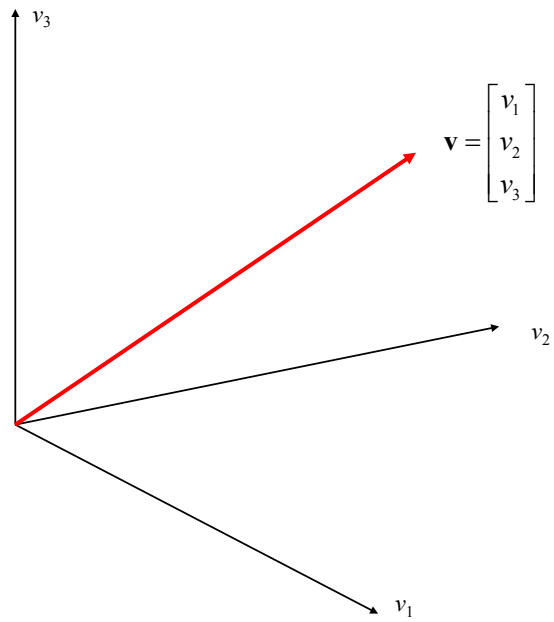


Visualization (Principal coordinates analysis)  
 Statistical hypothesis testing (PERMANOVA)

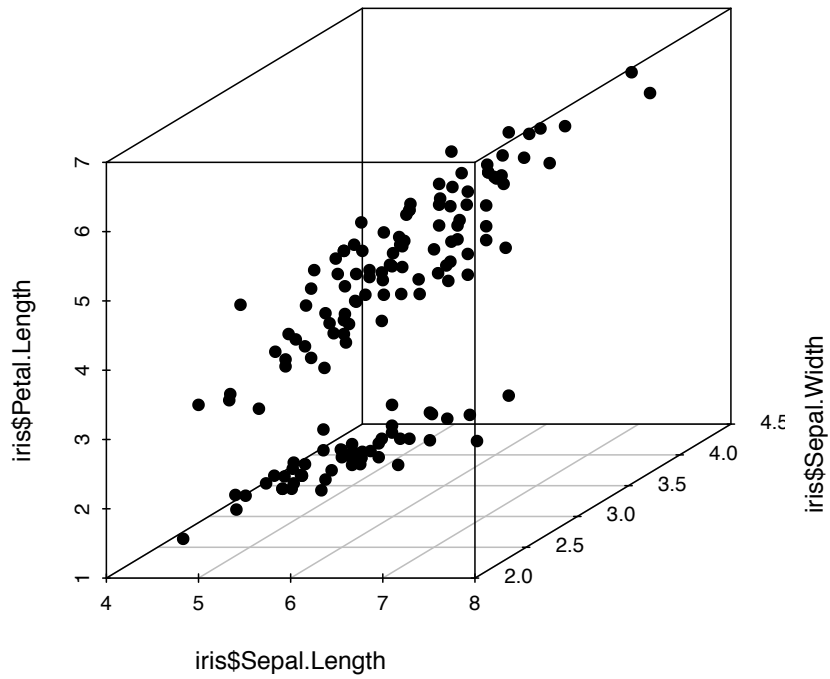
\*This step can sometimes be omitted for lack of necessity.

## PRINCIPAL COORDINATES ANALYSIS - MULTIDIMENSIONAL SCALING

Every multivariate sample can be represented as a vector in some vector space



7

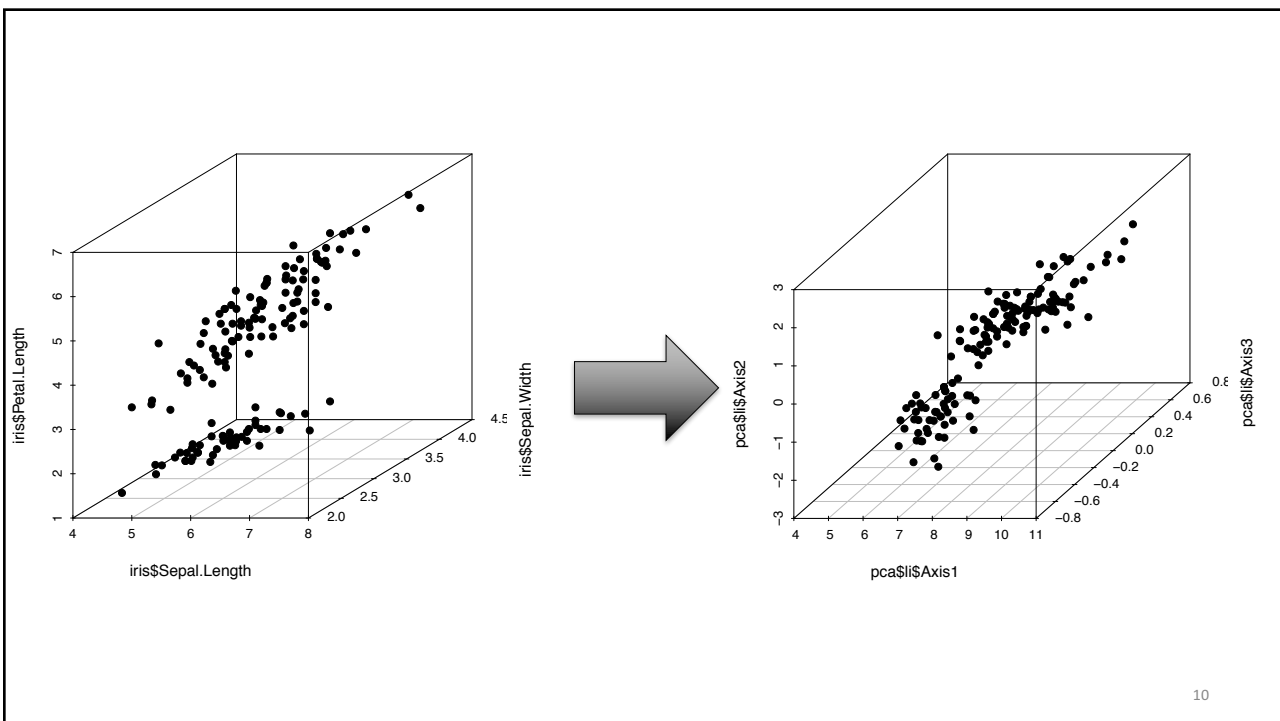


8

# Vector Basis

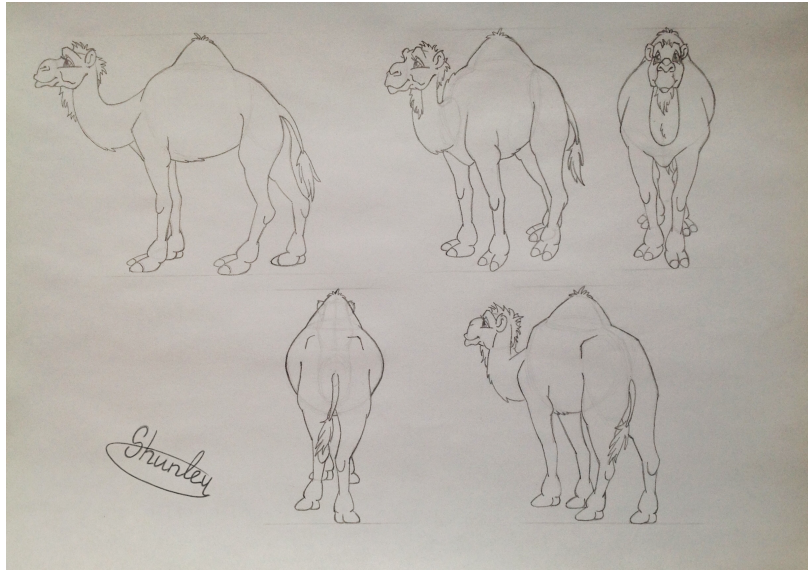
- A basis is a set of linearly independent (dot product is zero) vectors that **span** the vector space.
- **Spanning** the vector space: Any vector in this vector space may be represented as a linear combination of the basis vectors.
- The vectors forming a basis are orthogonal to each other. If all the vectors are of length 1, then the basis is called orthonormal.

9



10

Idea: Change basis so that we can better discover patterns in the data.



11

## Basic idea for analysis of multidimensional data

- Compute distances
- Reduce dimensions
- Embed in Euclidean space
- The general framework behind this process is called **Duality**

**diagram:**  $(\mathbf{X}_{n \times p}, \mathbf{Q}_{p \times p}, \mathbf{D}_{n \times n})$

- $\mathbf{X}_{n \times p}$  (centered) data matrix
- $\mathbf{Q}_{p \times p}$  column weights (weights on variables)
- $\mathbf{D}_{n \times n}$  row weights (weights on observations)

12

## Duality diagram defines the geometry of multivariate analysis

$$\begin{array}{ccc}
 \mathbb{R}^{p^*} & \xrightarrow{X} & \mathbb{R}^n \\
 Q \uparrow & & \downarrow D \\
 \mathbb{R}^p & \xleftarrow{X^t} & \mathbb{R}^{n^*} \\
 & & \uparrow W
 \end{array}$$

- $V = X^T D X$
- $W = X Q X^T$
- Duality:
  - The eigen decomposition of  $VQ$  leads to eigen-decomposition of  $WD$
- *Inertia* is equal to trace (sum of the diagonal elements) of  $VQ$  or  $WD$ .

13

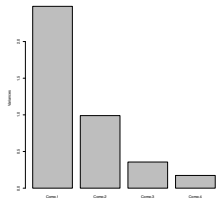
## Principal Component Analysis (PCA)

- Let  $Q = I$  and  $D = 1/n I$  and let  $X$  be centered.
- $VQ = X^T D X Q = 1/n X^T X$ .
- The inertia  $\text{Tr}(VQ) = \text{sum of the variances}$ .
- PCA decomposes the variance of  $X$  into independent components.
- To decompose the inertia means to find the eigen-system of  $VQ$  or equivalently  $WD$  matrices.
- Eigenvalues give the amount of inertia explained in corresponding dimension.
- Eigenvectors give the dimensions of variability.

14

# Example PCA

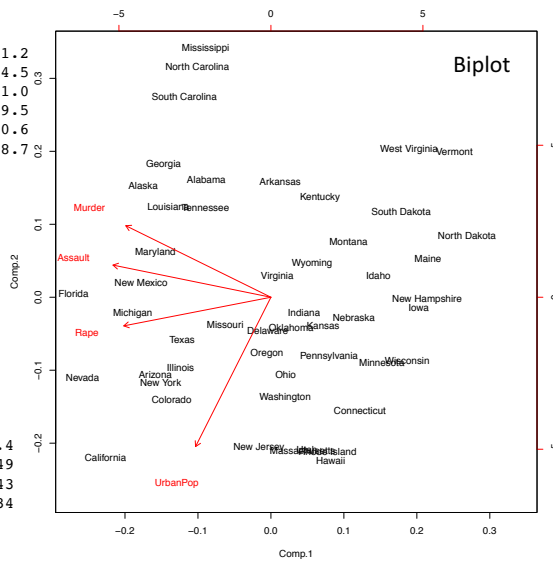
```
head(USArrests)
  Murder Assault UrbanPop Rape
Alabama   13.2    236     58 21.2
Alaska   10.0    263     48 44.5
Arizona   8.1    294     80 31.0
Arkansas  8.8    190     50 19.5
California 9.0    276     91 40.6
Colorado  7.9    204     78 38.7
```



Screeplot: plot of inertia

Loadings:

	Comp.1	Comp.2	Comp.3	Comp.4
Murder	-0.536	0.418	-0.341	0.649
Assault	-0.583	0.188	-0.268	-0.743
UrbanPop	-0.278	-0.873	-0.378	0.134
Rape	-0.543	-0.167	0.818	



15

## Centering

- Let  $Y$  be not centered data matrix with  $n$  observations (rows) and  $p$  variables (columns)
- Let  $H = (I - 1/n \mathbf{1}\mathbf{1}')$
- Then  $X = HY$  is centered

16



## From Euclidean distances to PCA to PCoA

- Note that if  $\mathbf{D}$  is a Euclidean distance, then
- $\mathbf{X} \mathbf{X}' = 1/n \mathbf{H} \mathbf{D}^{(2)} \mathbf{H}$ .
- PCoA is a generalization of PCA in that knowledge of  $\mathbf{X}$  is not required, all you need to represent the points is  $\mathbf{D}$ , the inter-point distance matrix.

17

## Representation of (arbitrary) distances in Euclidean space

- The idea is to use singular value decomposition (SVD) on the centered interpoint distance matrix to extract Euclidean dimensions
- SVD:  $\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}$ , where  $\mathbf{S}$  is diagonal matrix with diagonal elements  $s_1, s_2, \dots, s_n$ , and  $\mathbf{U}$  and  $\mathbf{V}$  are unit matrices (i.e. their determinant is 1 and they span (they are form orthonormal basis) their corresponding spaces)

18

## Definition

Let  $A$  be an  $n \times n$  matrix

Let  $\vec{x}$  and  $\lambda$  be such that

$$A\vec{x} = \lambda\vec{x} \quad \text{with } \vec{x} \neq \vec{0}$$

then  $\lambda$  is called an *eigenvalue* of  $A$  and

and  $\vec{x}$  is called an *eigenvector* of  $A$  and

## Note:

$$(A - \lambda I)\vec{x} = \vec{0}$$

If  $|A - \lambda I| \neq 0$  then  $\vec{x} = (A - \lambda I)^{-1} \vec{0} = \vec{0}$

thus  $|A - \lambda I| = 0$

is the condition for an eigenvalue.

## Diagonalization

**Theorem** If the matrix  $A$  is symmetric with distinct eigenvalues,  $\lambda_1, \dots, \lambda_n$ , with corresponding eigenvectors  $\vec{x}_1, \dots, \vec{x}_n$

Assume  $\vec{x}_i' \vec{x}_i = 1$

then  $A = \lambda_1 \vec{x}_1 \vec{x}_1' + \dots + \lambda_n \vec{x}_n \vec{x}_n'$

$$= [\vec{x}_1, \dots, \vec{x}_n] \begin{bmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_n \end{bmatrix} \begin{bmatrix} \vec{x}_1' \\ \vdots \\ \vec{x}_n' \end{bmatrix} = PDP'$$

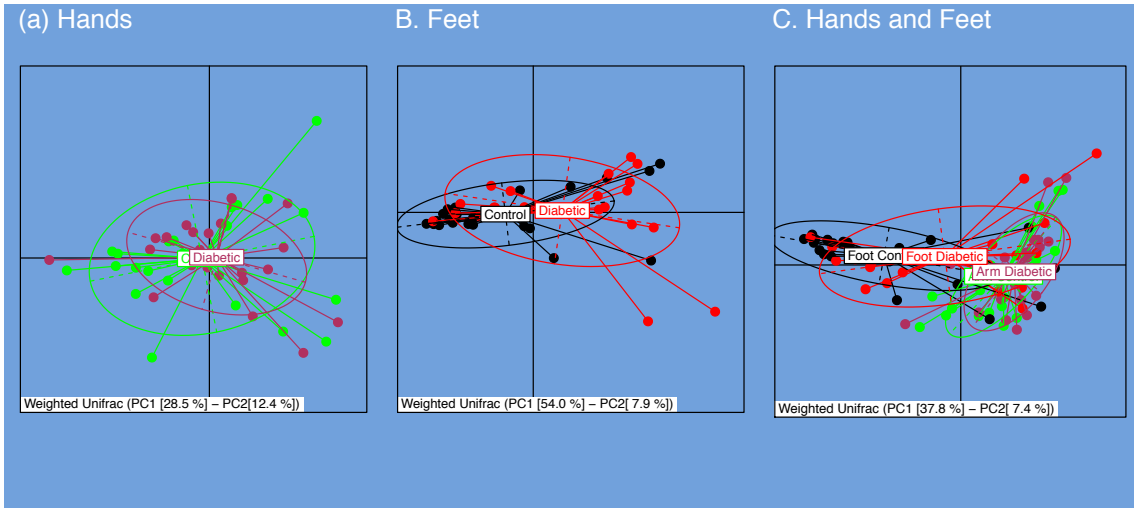
21

## PCoA details

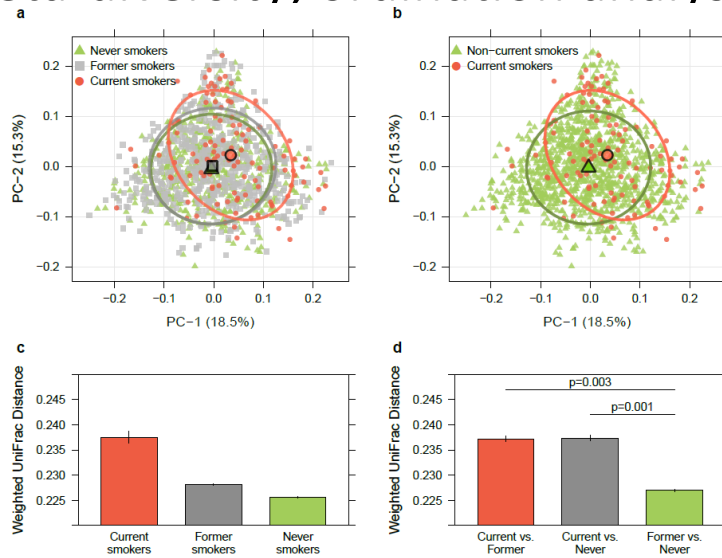
- Algorithm starting from  $\mathbf{D}$  inter-point distances:
  - Center the rows and columns of the matrix of square (element-by-element) distances:  $\mathbf{S} = -1/2\mathbf{H}\mathbf{D}^{(2)}\mathbf{H}$
  - Compute SVD by diagonalizing  $\mathbf{S}$ ,  $\mathbf{S} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$
  - Extract Euclidean representations:  $\underline{\mathbf{X}} = \mathbf{U}\mathbf{\Lambda}^{1/2}$
- The relative values of diagonal elements of  $\mathbf{\Lambda}$  gives the proportion of variability explained by each of the axes.
- The values of  $\mathbf{\Lambda}$  should always be looked at in deciding how many dimensions to retain.

22

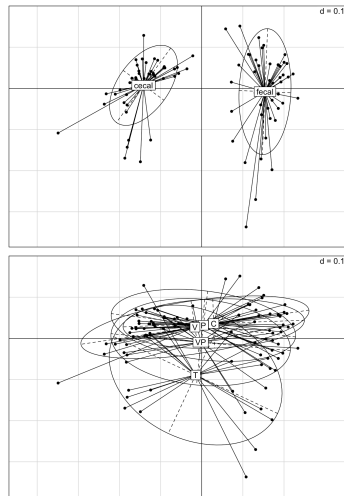
Differentiation of microbiota between diabetic and non-diabetic subjects and across body sites



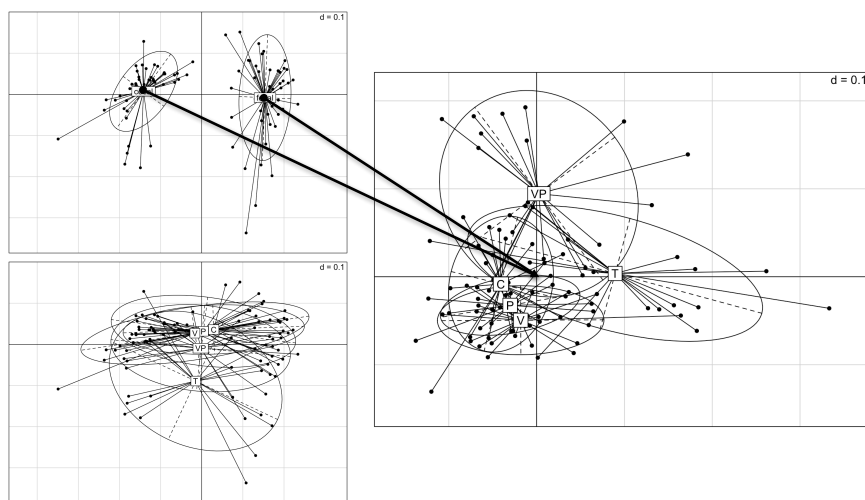
Beta-diversity; ordination analysis



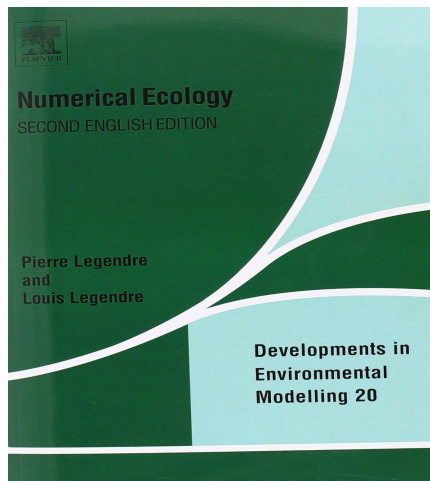
## Within class analysis



## Within class analysis



## Suggested reading



- Susan Holmes “Multivariate Data Analysis: The French Way”, IMS Lecture Notes–Monograph Series, 2006.