

Lecture 8

Monday, February 7, 2022 9:32 AM

Admin

Principal Component Analysis

① HW due ~~today~~ Thurs.

Reading

- ① Blind source separation
- ② Multi- and megavariate data analysis, 2003, Eriksson.

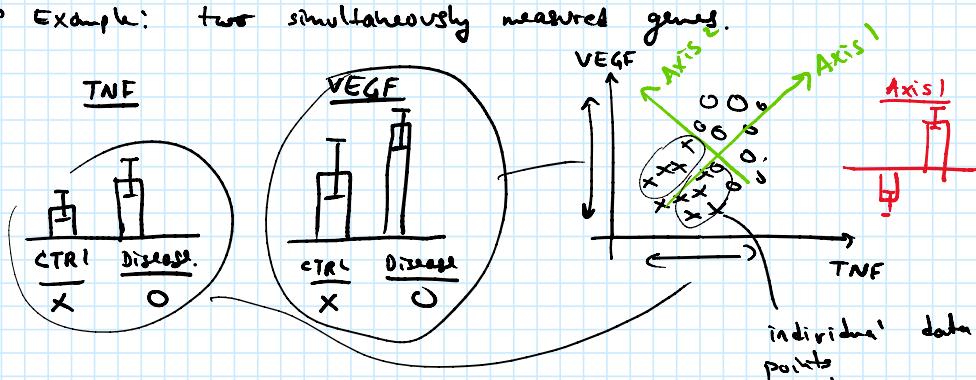
Objectives

- ① Discuss datasets consisting of many variables from which measurements are arising from a lower-dimensional underlying process.
- ② Define principal component analysis (PCA)
- ③ Consider examples

Large datasets arising from lower-dimensional processes

- In many cases we consider datasets in which we collect many datapoints that are highly correlated with one another.
 - multi-dimensional gene expression data from healthy & diseased subjects
 - multiple microphones recording from different parts of a room → "cocktail party problem"
 - multiple ECG recordings from a pregnant woman.
- In each case, the measurements are tightly related to one another, but represent contributions from underlying processes
 - together drive each sensor reading.

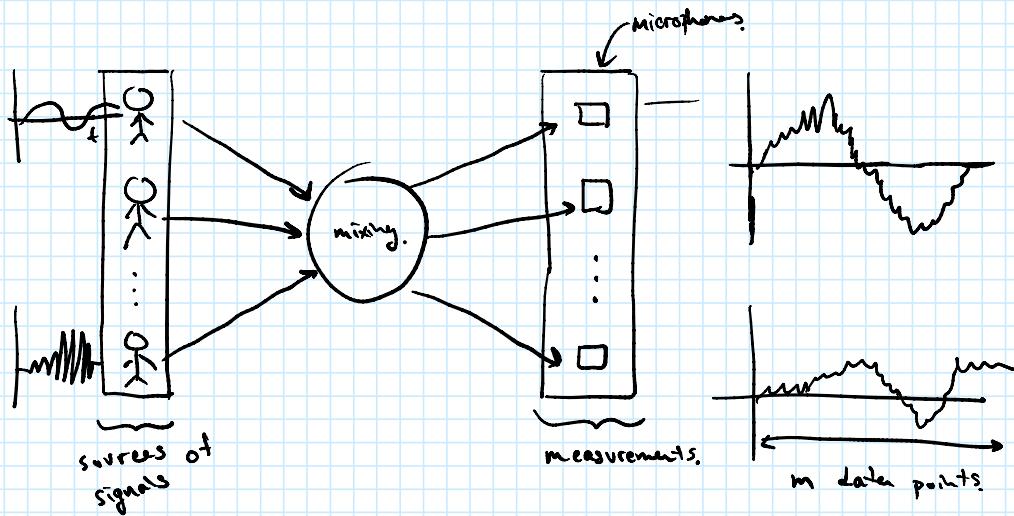
→ Example: two simultaneously measured genes.



- Can identify a composite axis (Axis 1) that more clearly separates groups.
- ⇒ our approach can better distinguish source of variation in underlying data.

Example — The cocktail party problem

Example - The cocktail party problem



Principal Component Analysis

- problem: multi-dimensional dataset

$$\mathbf{X} = \underset{\substack{\# \text{ of} \\ \# \text{ of} \\ \text{observation.}}}{m \times n} \underset{\# \text{ of variables.}}{\uparrow}$$

- Goal is to represent the data in a lower-dimensional plane

→ first described by Cauchy, later by Pearson.
↳ lines/planes in data be formulated as finding of closest

→ representation in lower-dimensional plane will reveal groups of differences between samples, sensors, that would not otherwise be obvious.

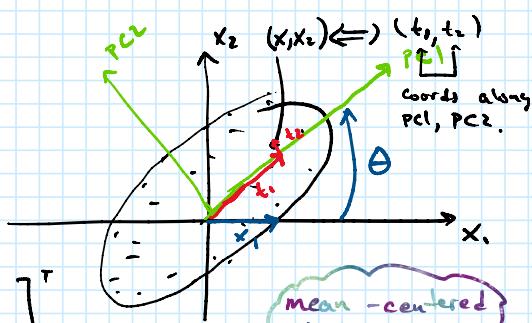
- We can re-express the data in a new coord. system by rotating.

$$x_1 = t_1 \cos \theta + t_2 \sin \theta$$

$$x_2 = -t_1 \sin \theta + t_2 \cos \theta$$

$$\underbrace{[x_1 \ x_2]}_{\text{original data}} = \underbrace{[t_1 \ t_2]}_{\text{PCA coords}} \underbrace{\begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}}_{R^T}$$

orthogonal rotation matrix.



- we can choose our rotation, θ , to:

- 1) maximize variation (variance) in the data explained by t_1 .
- 2) reduce covariance between variables (t_1, t_2, \dots)

- We measure m observations, want to transform them

$$\mathbf{X} = \bar{\mathbf{X}} + \mathbf{T} \mathbf{P}^T + \mathbf{E}$$

↓
 # samples
 $(m \times n)$
 measurement means.
 \uparrow
 $m \times n$.
 $\begin{bmatrix} \text{mean(meas1)} & \dots & \text{mean(meas1)} \\ \vdots & \ddots & \vdots \\ \text{mean(meas1)} & \dots & \text{mean(meas1)} \end{bmatrix}$

↓
 # measurement variables.
 $(m \times n)$
 \mathbf{T}
 \mathbf{P}^T
 $(n \times k)$
 ↓
 # of prin. comp.
 scores matrix.

↓
 # samples
 # of prin. comp.
 $(n \times k)$
 \mathbf{E}
 residual noise. (Gaussian).

- we often pre-treat the data
 - subtract column means $\Rightarrow \bar{\mathbf{X}} = 0$
 - scale by standard dev.

$$\left. \begin{array}{l} \text{- subtract column means} \Rightarrow \bar{\mathbf{X}} = 0 \\ \text{- scale by standard dev.} \end{array} \right\} \text{z-score } \tilde{x}_{ij} = \frac{x_{ij} - \bar{x}_j}{\sigma_j}$$

- How do we compute \mathbf{T}, \mathbf{P}
 - \rightarrow In PCA, we choose $E[t_{ij}, t_{kj}] = \underbrace{\text{cov}(\text{col } i, \text{col } j)}_{\text{two cols.}} = 0, i \neq j$

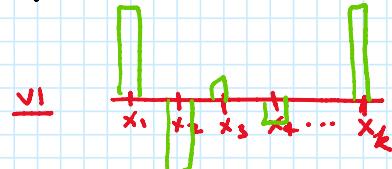
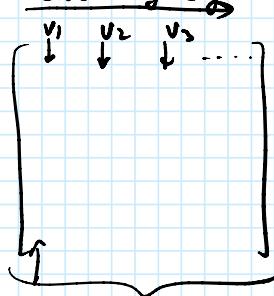
- To do so, compute cov matrix of \mathbf{X}

$$C_{ij} = \frac{1}{m} \sum_{k=1}^m x_{ki} x_{kj} \Rightarrow \underline{\underline{C}} = \frac{1}{m} \mathbf{X}^T \mathbf{X} \leftarrow \text{want to diagonalize.}$$

$$\underline{\underline{C}} \mathbf{V} = \lambda \mathbf{V}$$

↑
 eigenvectors
 components along x_1, x_2, \dots

eigenvalues
descending eigenvalues.



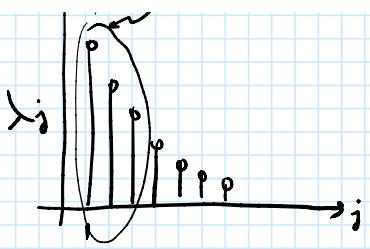
components along each eigen vector are called "loadings".

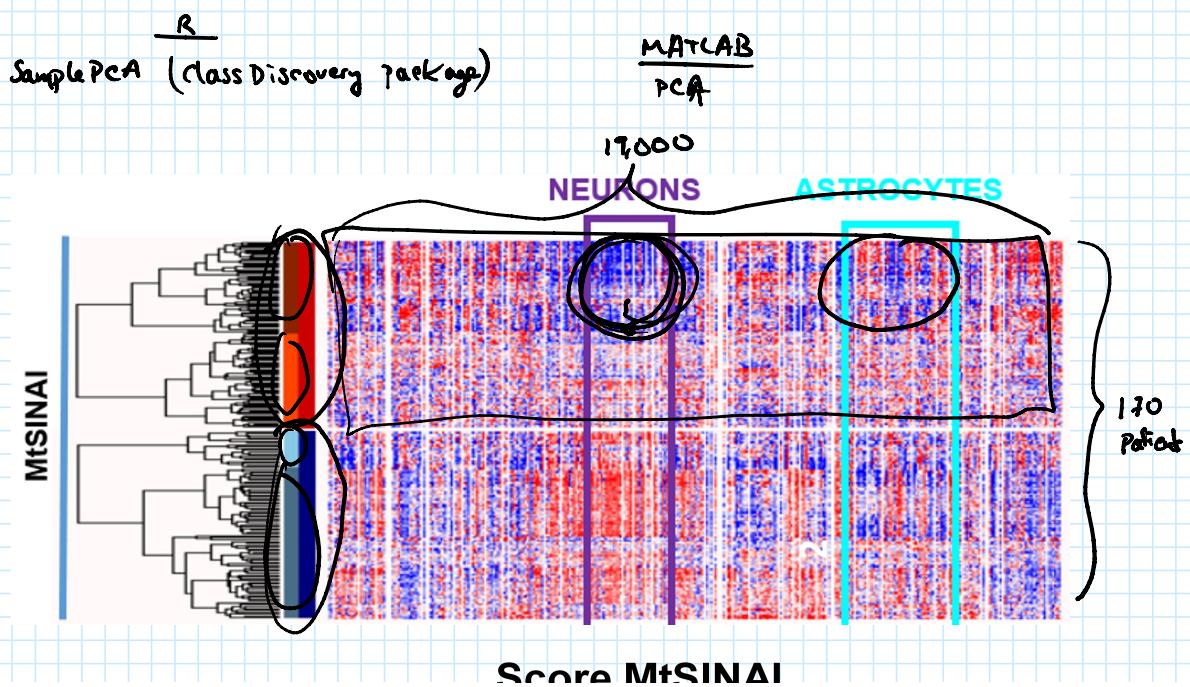
\rightarrow represent the weight of each measurement variable in defining that PC.

- what do eigenvalues mean?

contributes a lot of variance.









Score MtSINAI

