

Machine Learning for Data Science (CS4786)

Lecture 4

Canonical Correlation Analysis (CCA)

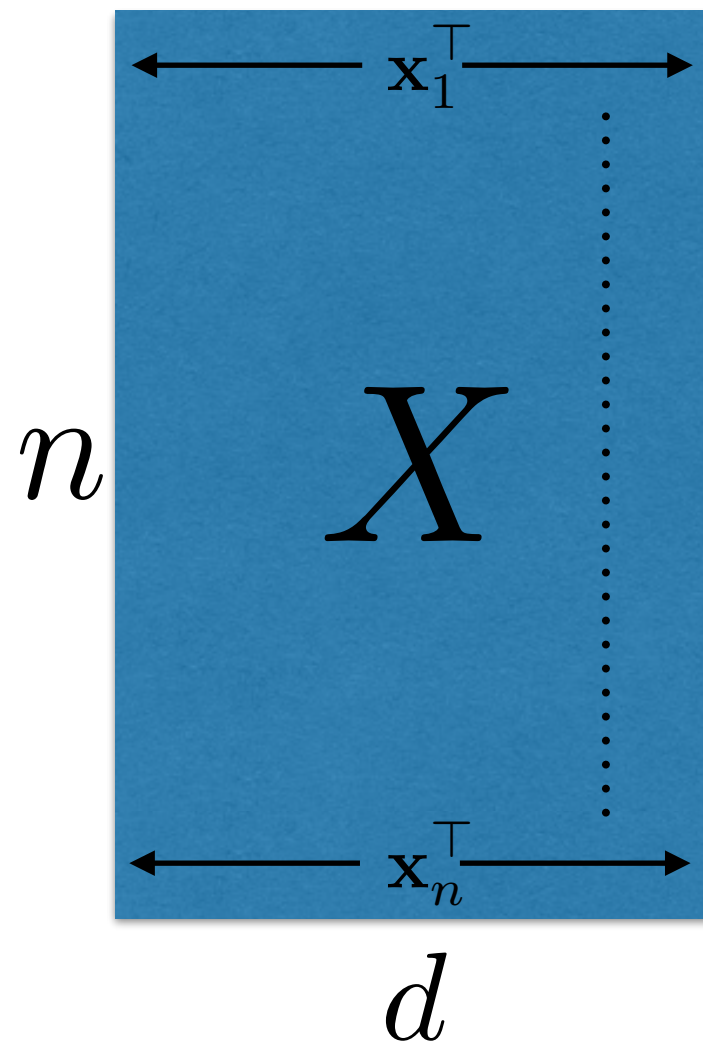
Course Webpage :

<http://www.cs.cornell.edu/Courses/cs4786/2016fa/>

Announcement

- We are grading HW0 and you will be added to cms by monday
- HW1 will be posted tonight on webpage (homework tab)
- HW1 on CCA and PCA (due in a week)

QUIZ



Assume points are centered. Which of the following are equal to the covariance matrix?

A. $\Sigma = \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t^\top \mathbf{x}_t$

B. $\Sigma = \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^\top$

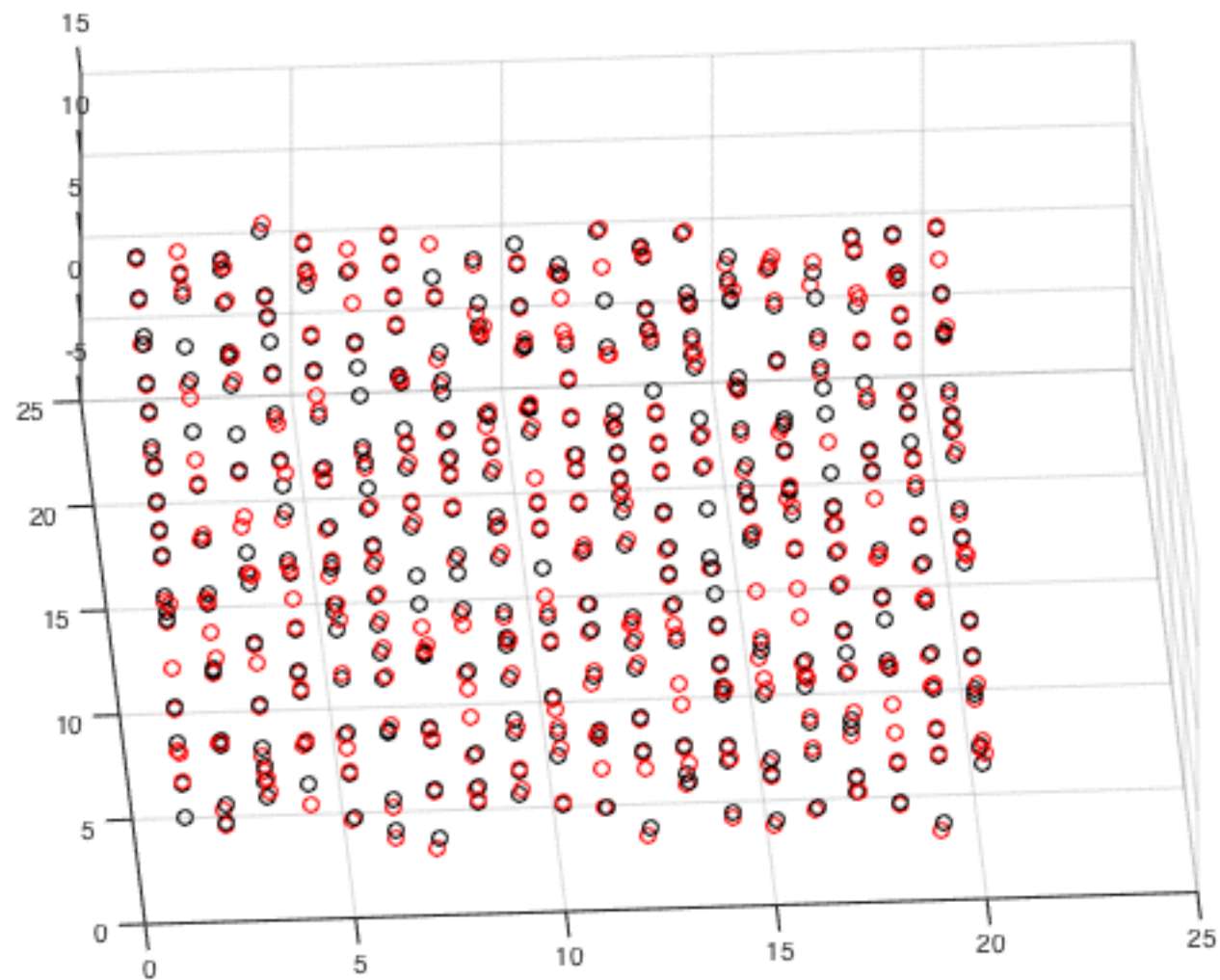
C. $\Sigma = X X^\top$

D. $\Sigma = X X^\top$

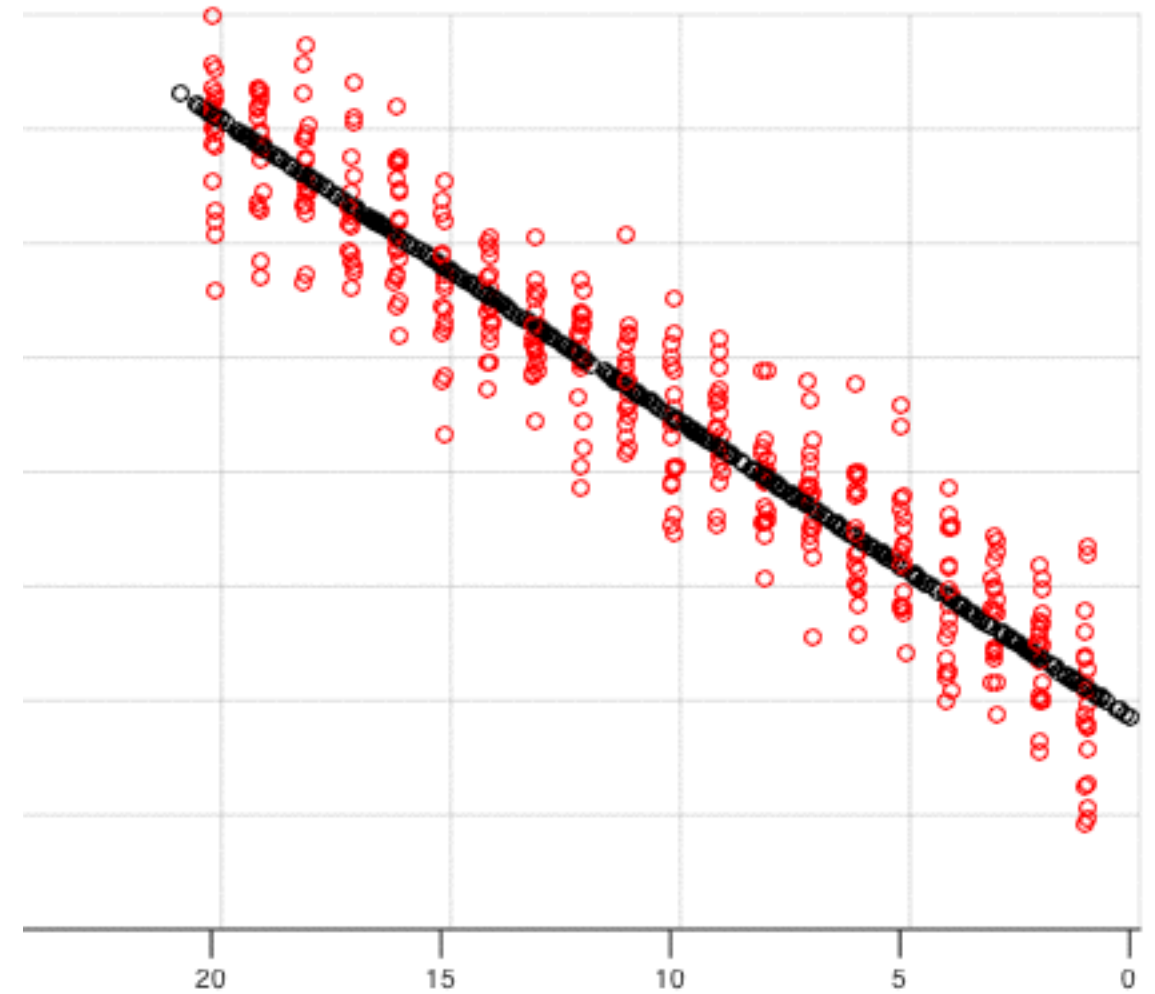
Example: Students in classroom



Maximize Spread



Minimize Reconstruction Error



PRINCIPAL COMPONENT ANALYSIS

1. $\Sigma = \text{COV}(X)$

2. $W = \text{eigs}(\Sigma, K)$

3. $Y = (X - \mu) \times W$

RECONSTRUCTION

4.

$$\hat{X} = Y \times W^T + \mu$$

WHEN $d \gg n$

- If $d \gg n$ then Σ is large
- But we only need top K eigen vectors.
- Idea: use SVD

$$X - \mu = UDV^T$$

$$\begin{aligned} V^T V &= I \\ U^T U &= I \end{aligned}$$

Then note that, $\Sigma = (X - \mu)^T (X - \mu) = VD^2V$

- Hence, matrix V is the same as matrix W got from eigen decomposition of Σ , eigenvalues are diagonal elements of D^2
- Alternative algorithm:

$$[U, V] = \text{SVD}(X - \mu, K) \quad W = V$$

WHEN TO USE PCA?

- When data naturally lies in a low dimensional linear subspace
- To minimize reconstruction error
- Find directions where data is maximally spread

Canonical Correlation Analysis



+ Age
Gender
Angle

Canonical Correlation Analysis

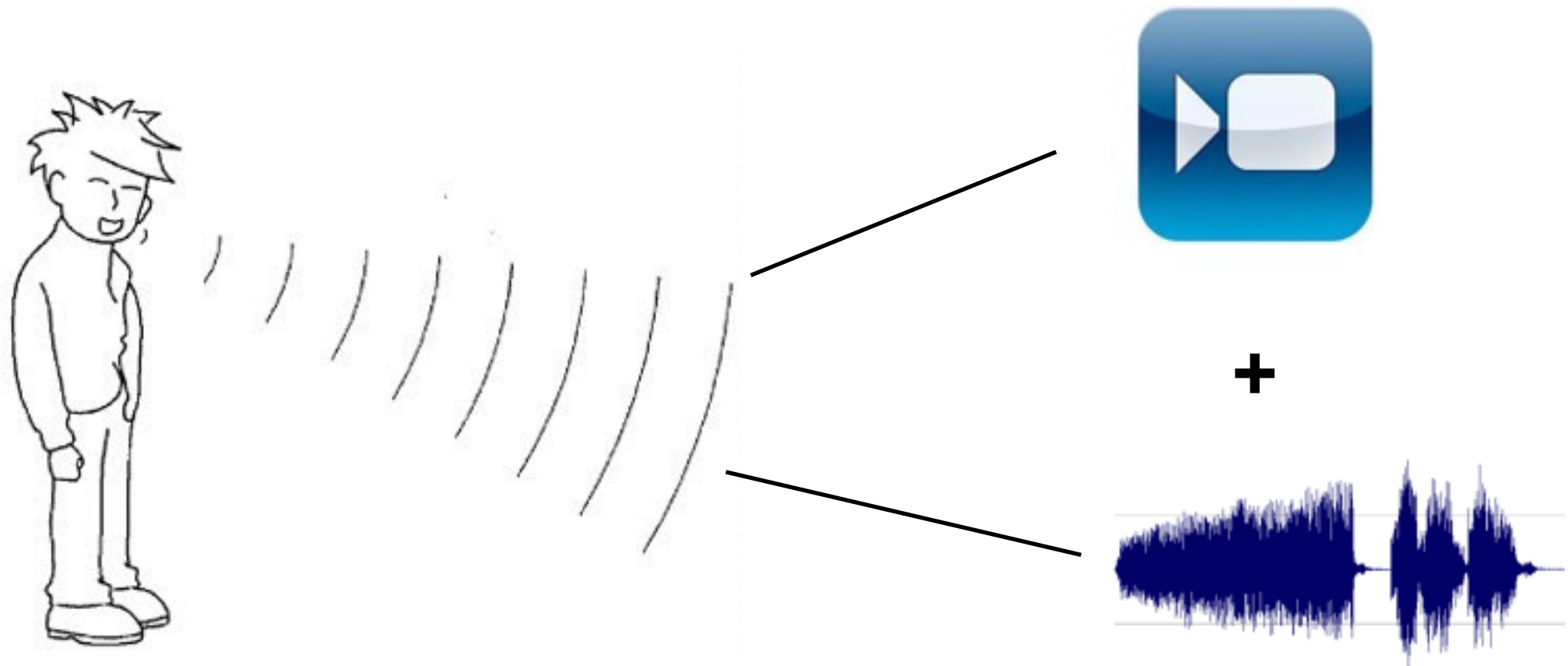


+ Age
Gender
Angle

TWO VIEW DIMENSIONALITY REDUCTION

- Data comes in pairs $(\mathbf{x}_1, \mathbf{x}'_1), \dots, (\mathbf{x}_n, \mathbf{x}'_n)$ where \mathbf{x}_t 's are d dimensional and \mathbf{x}'_t 's are d' dimensional
- Goal: Compress say view one into $\mathbf{y}_1, \dots, \mathbf{y}_n$, that are K dimensional vectors
 - Retain information redundant between the two views
 - Eliminate “noise” specific to only one of the views

EXAMPLE I: SPEECH RECOGNITION



- Audio might have background sounds uncorrelated with video
- Video might have lighting changes uncorrelated with audio
- Redundant information between two views: the speech

EXAMPLE II: COMBINING FEATURE EXTRACTIONS

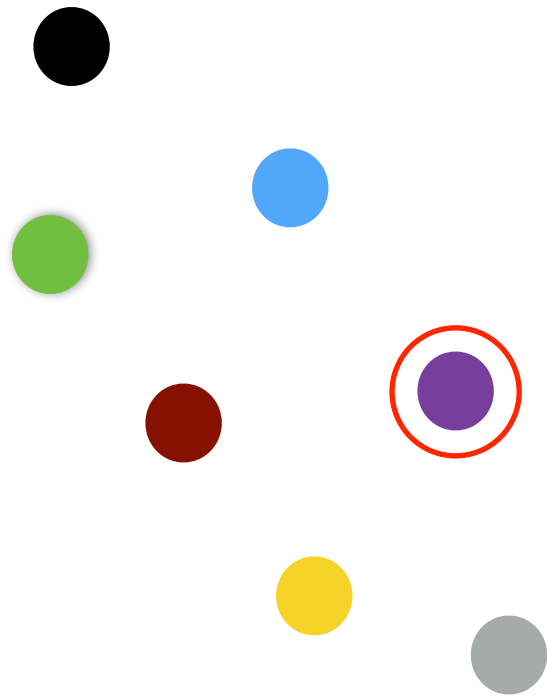
- Method A and Method B are both equally good feature extraction techniques
- Concatenating the two features blindly yields large dimensional feature vector with redundancy
- Applying techniques like CCA extracts the key information between the two methods
- Removes extra unwanted information

How do we get the right direction? (say $K = 1$)

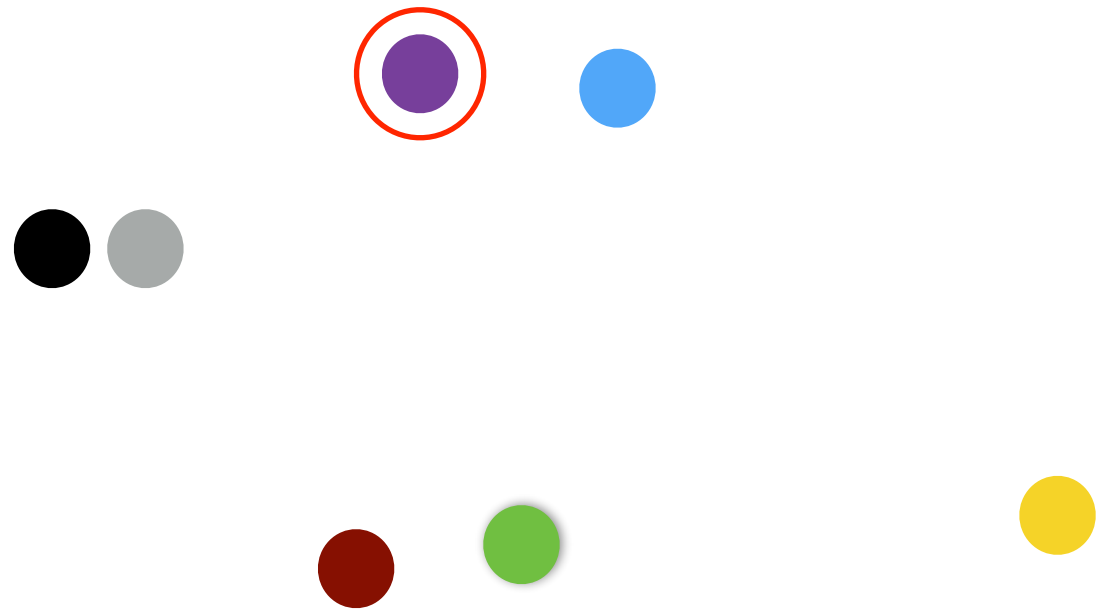


Age
+ Gender
Angle

WHICH DIRECTION TO PICK?



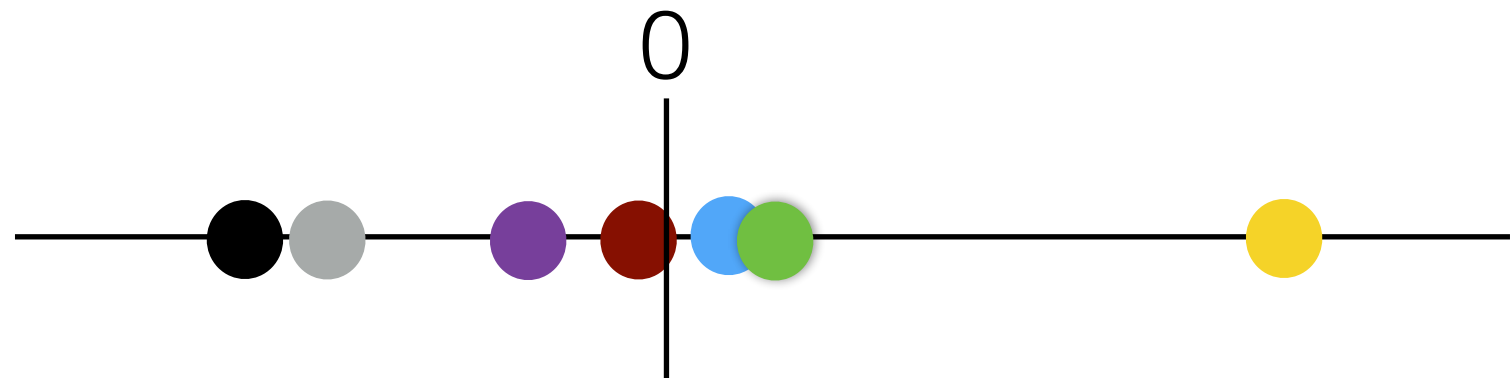
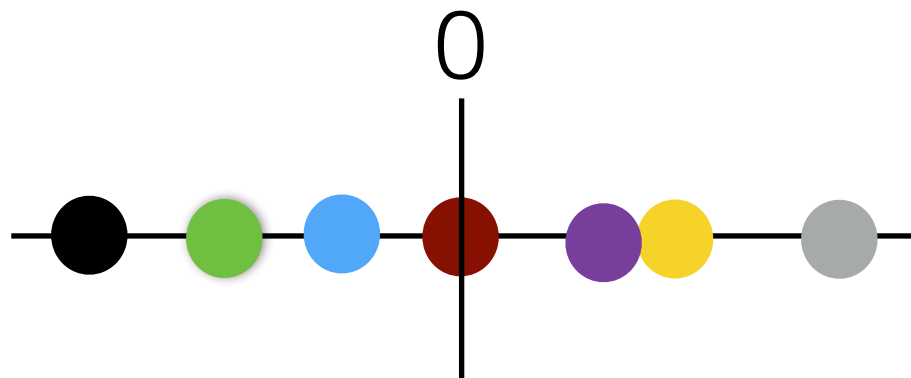
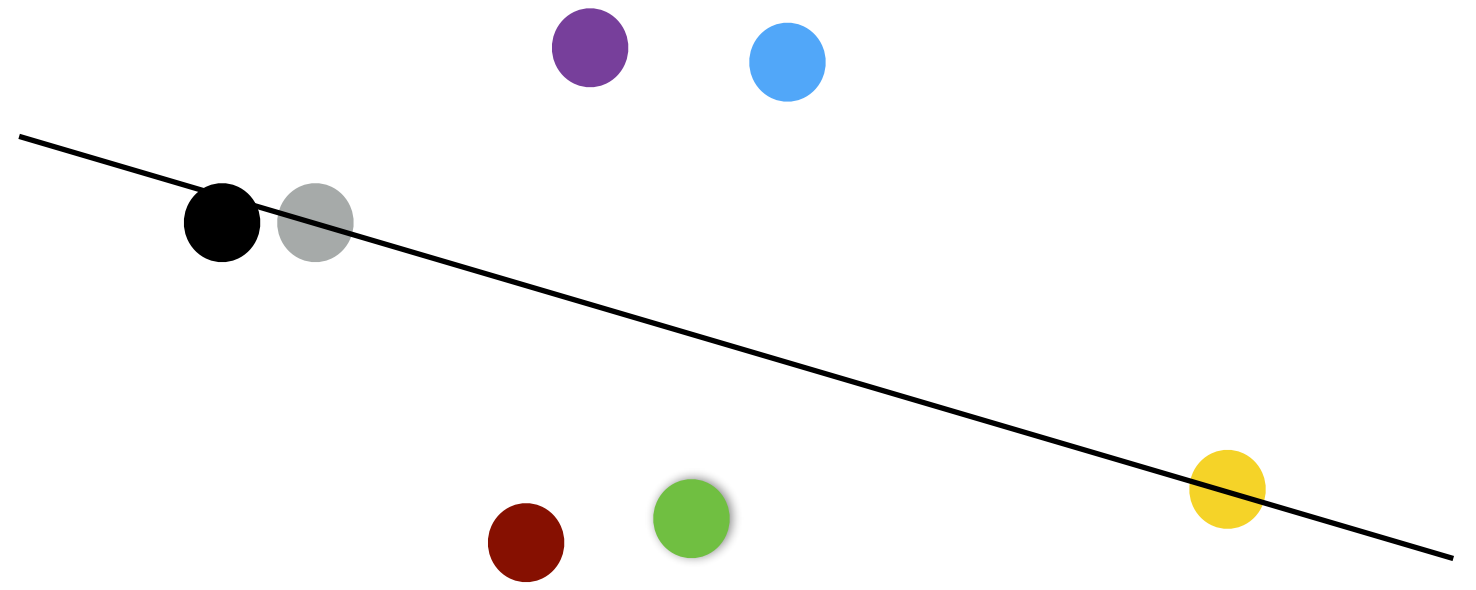
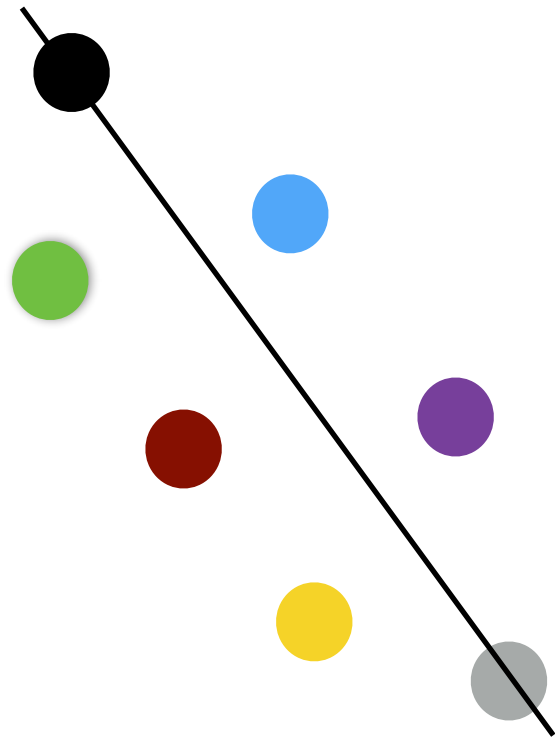
View I



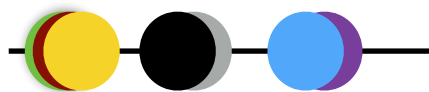
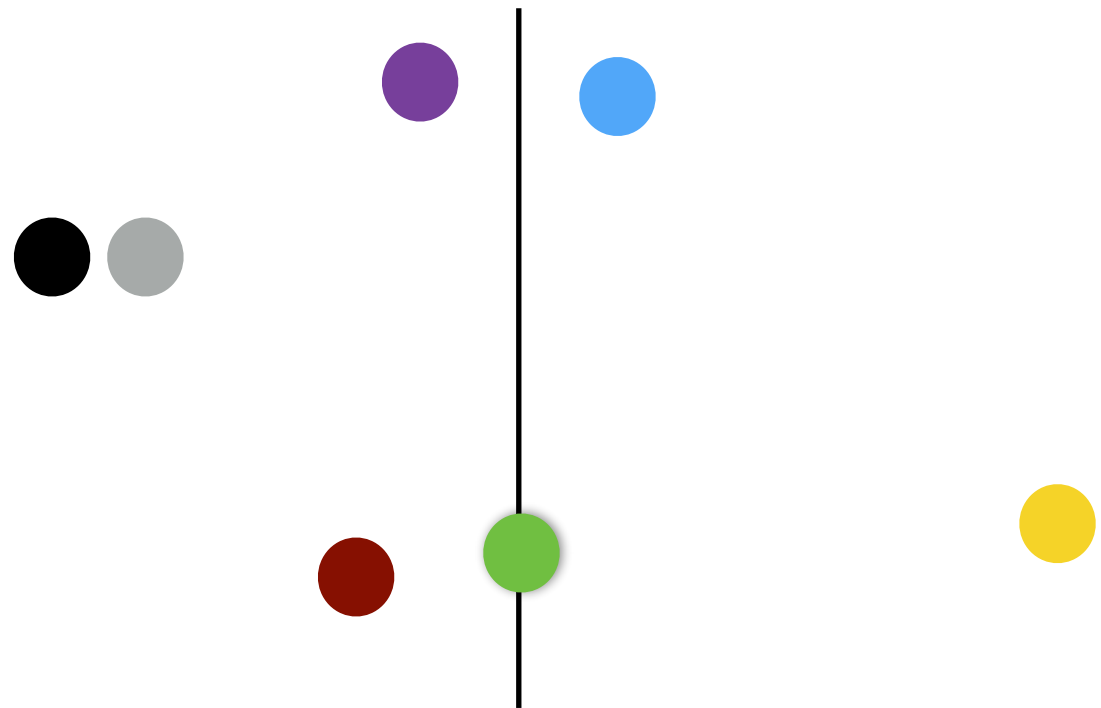
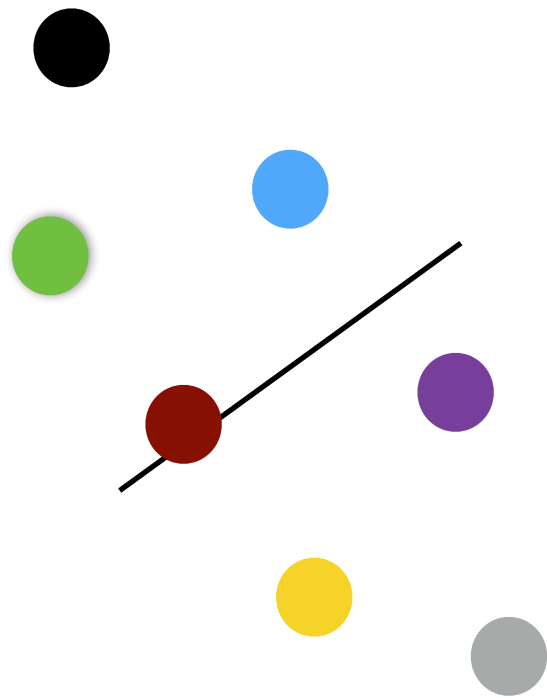
View II

WHICH DIRECTION TO PICK?

PCA direction



WHICH DIRECTION TO PICK?



Direction has large covariance

How do we pick the right direction to project to?

MAXIMIZING CORRELATION COEFFICIENT

- Say \mathbf{w}_1 and \mathbf{v}_1 are the directions we choose to project in views 1 and 2 respectively we want these directions to maximize,

$$\frac{1}{n} \sum_{t=1}^n \left(\mathbf{y}_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[1] \right) \cdot \left(\mathbf{y}'_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}'_t[1] \right)$$

where $\mathbf{y}_t[1] = \mathbf{w}_1^\top \mathbf{x}_t$ and $\mathbf{y}'_t[1] = \mathbf{v}_1^\top \mathbf{x}'_t$

MAXIMIZING CORRELATION COEFFICIENT

- Say \mathbf{w}_1 and \mathbf{v}_1 are the directions we choose to project in views 1 and 2 respectively we want these directions to maximize,

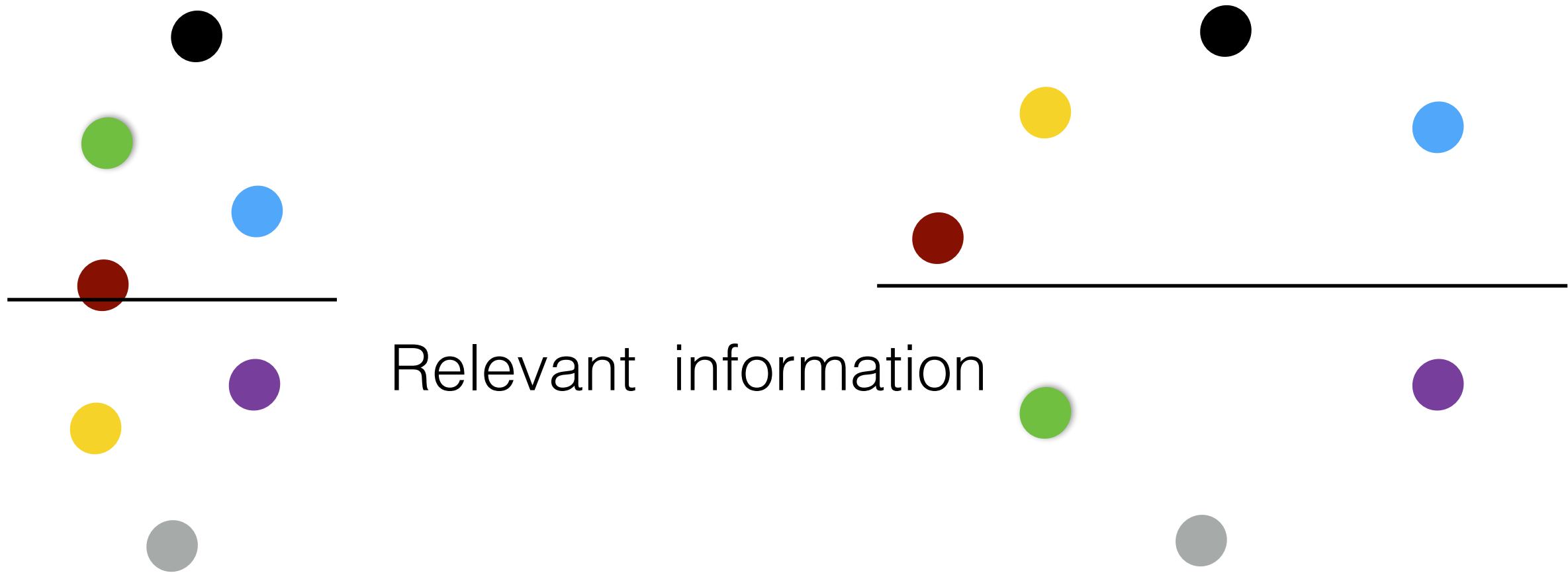
$$\frac{1}{n} \sum_{t=1}^n \left(\mathbf{y}_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[1] \right) \cdot \left(\mathbf{y}'_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}'_t[1] \right)$$

$$\text{s.t. } \frac{1}{n} \sum_{t=1}^n \left(\mathbf{y}_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[1] \right)^2 = \frac{1}{n} \sum_{t=1}^n \left(\mathbf{y}'_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}'_t[1] \right)^2 = 1$$

where $\mathbf{y}_t[1] = \mathbf{w}_1^\top \mathbf{x}_t$ and $\mathbf{y}'_t[1] = \mathbf{v}_1^\top \mathbf{x}'_t$

What is the problem
with the above?

WHY NOT MAXIMIZE COVARIANCE



$$\text{Say } \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t[2] \cdot \mathbf{x}'_t[2] > 0$$

Scaling up this coordinate we can blow up covariance

MAXIMIZING CORRELATION COEFFICIENT

- Say \mathbf{w}_1 and \mathbf{v}_1 are the directions we choose to project in views 1 and 2 respectively we want these directions to maximize,

$$\frac{\frac{1}{n} \sum_{t=1}^n (\mathbf{y}_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[1]) \cdot (\mathbf{y}'_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}'_t[1])}{\sqrt{\frac{1}{n} \sum_{t=1}^n (\mathbf{y}_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[1])^2} \sqrt{\frac{1}{n} \sum_{t=1}^n (\mathbf{y}'_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}'_t[1])^2}}$$

BASIC IDEA OF CCA

- Normalize variance in chosen direction to be constant (say 1)
- Then maximize covariance
- This is same as maximizing “correlation coefficient”

COVARIANCE VS CORRELATION

- $\text{Covariance}(A, B) = \mathbb{E}[(A - \mathbb{E}[A]) \cdot (B - \mathbb{E}[B])]$

Depends on the scale of A and B . If B is rescaled, covariance shifts.

- $\text{Corelation}(A, B) = \frac{\mathbb{E}[(A - \mathbb{E}[A]) \cdot (B - \mathbb{E}[B])]}{\sqrt{\text{Var}(A)}\sqrt{\text{Var}(B)}}$

Scale free.

MAXIMIZING CORRELATION COEFFICIENT

- Say \mathbf{w}_1 and \mathbf{v}_1 are the directions we choose to project in views 1 and 2 respectively we want these directions to maximize,

$$\frac{1}{n} \sum_{t=1}^n \left(\mathbf{y}_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[1] \right) \cdot \left(\mathbf{y}'_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}'_t[1] \right)$$

$$\text{s.t. } \frac{1}{n} \sum_{t=1}^n \left(\mathbf{y}_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[1] \right)^2 = \frac{1}{n} \sum_{t=1}^n \left(\mathbf{y}'_t[1] - \frac{1}{n} \sum_{t=1}^n \mathbf{y}'_t[1] \right)^2 = 1$$

where $\mathbf{y}_t[1] = \mathbf{w}_1^\top \mathbf{x}_t$ and $\mathbf{y}'_t[1] = \mathbf{v}_1^\top \mathbf{x}'_t$

CANONICAL CORRELATION ANALYSIS

- Hence we want to solve for projection vectors \mathbf{w}_1 and \mathbf{v}_1 that

$$\text{maximize } \frac{1}{n} \sum_{t=1}^n \mathbf{w}_1^\top (\mathbf{x}_t - \boldsymbol{\mu}) \cdot \mathbf{v}_1^\top (\mathbf{x}'_t - \boldsymbol{\mu}')$$

$$\text{subject to } \frac{1}{n} \sum_{t=1}^n (\mathbf{w}_1^\top (\mathbf{x}_t - \boldsymbol{\mu}))^2 = \frac{1}{n} \sum_{t=1}^n (\mathbf{v}_1^\top (\mathbf{x}'_t - \boldsymbol{\mu}'))^2 = 1$$

where $\boldsymbol{\mu} = \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t$ and $\boldsymbol{\mu}' = \frac{1}{n} \sum_{t=1}^n \mathbf{x}'_t$

CANONICAL CORRELATION ANALYSIS

- Hence we want to solve for projection vectors \mathbf{w}_1 and \mathbf{v}_1 that

$$\text{maximize } \frac{1}{n} \sum_{t=1}^n \mathbf{w}_1^\top (\mathbf{x}_t - \boldsymbol{\mu}) \cdot \mathbf{v}_1^\top (\mathbf{x}'_t - \boldsymbol{\mu}')$$

$$\text{subject to } \frac{1}{n} \sum_{t=1}^n (\mathbf{w}_1^\top (\mathbf{x}_t - \boldsymbol{\mu}))^2 = \frac{1}{n} \sum_{t=1}^n (\mathbf{v}_1^\top (\mathbf{x}'_t - \boldsymbol{\mu}'))^2 = 1$$

$$\text{where } \boldsymbol{\mu} = \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t \text{ and } \boldsymbol{\mu}' = \frac{1}{n} \sum_{t=1}^n \mathbf{x}'_t$$

CANONICAL CORRELATION ANALYSIS

- Hence we want to solve for projection vectors \mathbf{w}_1 and \mathbf{v}_1 that

$$\text{maximize } \mathbf{w}_1^\top \Sigma_{1,2} \mathbf{v}_1$$

$$\text{subject to } \mathbf{w}_1^\top \Sigma_{1,1} \mathbf{w}_1 = \mathbf{v}_1^\top \Sigma_{2,2} \mathbf{v}_1 = 1$$

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} = \text{COV} \left(\begin{pmatrix} X & X' \end{pmatrix} \right)$$

SOLUTION

$$W_1 = \text{eigs}\left(\Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}, K\right)$$

$$W_2 = \text{eigs}\left(\Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}, K\right)$$

CCA ALGORITHM

$$1. \quad X = \begin{pmatrix} n & \begin{matrix} X_1 \\ d_1 \end{matrix}, & \begin{matrix} X_2 \\ d_2 \end{matrix} \end{pmatrix}$$

$$2. \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} = \text{COV} \left(\begin{matrix} X \end{matrix} \right)$$

$$3. \quad W_1 = \text{eigs} \left(\begin{matrix} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \end{matrix}, K \right)$$

$$4. \quad Y_1 = \begin{matrix} X_1 - \mu_1 \\ \end{matrix} \times W_1$$