# Automatic Feature Extraction using Unsupervised Learning
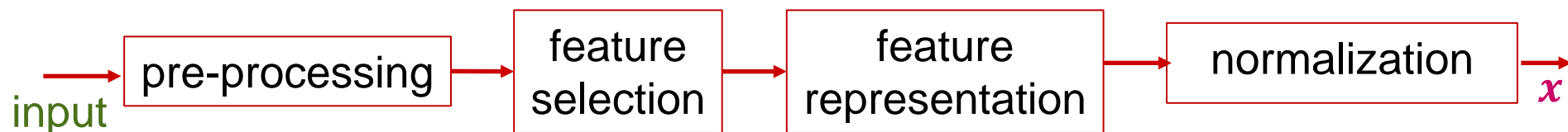
Reading:     19.1, 19.2, 19.3, 19.5, 19.6, 19.7, 19.9

- Recollect:    Several steps for Feature extraction

- Operationally,   Feature extraction  =  Dimensionality reduction

- Example:    signal waveform (infinite dimensions) $\rightarrow$  10-dimensional $\boldsymbol{x}$

- Two ways of doing Feature extraction ...

1. Using Manual methods:    Highly domain-specific ✓

2. Using <u>Automatic</u> dimensionality reduction methods

   — Useful if no domain knowledge

   — Need lots of Training data to overcome lack of domain knowledge

   — Deep learning networks increasingly used to extract features now

input $\rightarrow$ pre-processing $\rightarrow$ feature selection $\rightarrow$ feature representation $\rightarrow$ normalization $\rightarrow$ $\boldsymbol{x}$
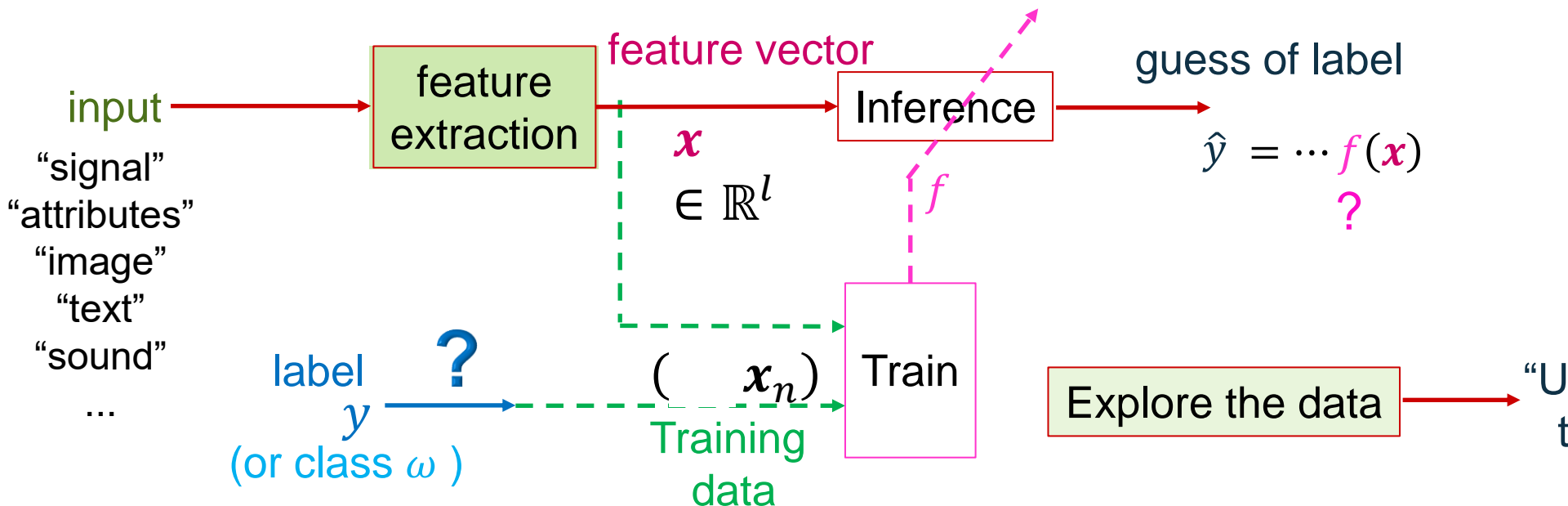
# Automatic feature extraction (Unsupervised Learning)

- Linear Discriminant Analysis (LDA)  and  Kernel LDA

- Principal Component Analysis (PCA)  and  Kernel PCA

- Latent Semantic Analysis (LSA)

- Independent Component Analysis (ICA)

Linear methods

- Multidimensional scaling (MDS),   t-SNE

- Autoencoders

Python `sklearn` has most of these

- Embeddings

So, this is also a lecture on Unsupervised learning

- Clustering  (k-Means)

- Except LDA methods, other methods use Unsupervised learning

- Discuss only at high level,  for motivation/intuition
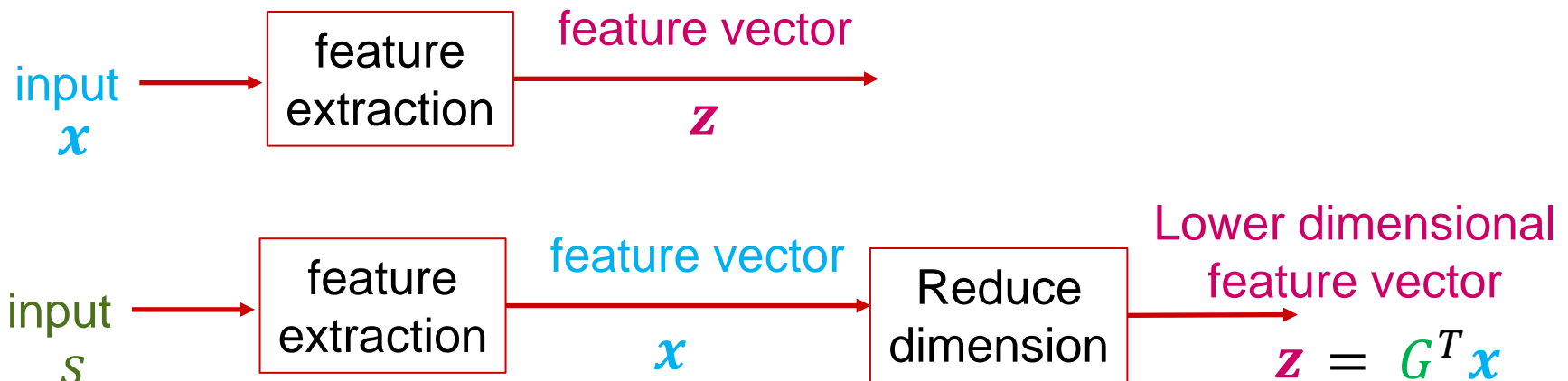
- <u>Unsupervised learning:</u>   $N$ Training points are (    $\boldsymbol{x}_n$)

    feature vector

- No concept of labels here

- $\Rightarrow$   Cannot do Classification or Regression here

- But can 'Explore data'   and do <u>Automatic Feature extraction</u>

- Let's study Feature extraction application of Unsupervised learning



feature vector

guess of label

input

→ feature extraction → Inference →

"signal"
"attributes"
"image"
"text"
"sound"
...

$\boldsymbol{x}$

$\in \mathbb{R}^l$

$\hat{y} = \cdots f(\boldsymbol{x})$

?

$f$

label
$y$
(or class $\omega$ )

?

(    $\boldsymbol{x}_n$)

Training
data

Train

Explore the data →

"U
t

# Feature extraction: Linear methods

$$z_n = G^T x_n$$

- First, look at <u>Linear</u> feature extraction

- Map inputs $x_1, x_2, \ldots, x_n$ to extracted features $z_1, z_2, \ldots, z_n$ <u>linearly</u>

- But linear methods require input to already be a <u>vector</u>

- $\Rightarrow$ Methods often used to reduce dimension of feature <u>vector</u> $x$

- (Explains why we called input as $x$ and extracted feature as $z$)

input $x$ → [ feature extraction ] → feature vector $z$

input $s$ → [ feature extraction ] → feature vector $x$ → [ Reduce dimension ] → Lower dimensional feature vector $z = G^T x$

# LDA,
# Kernel LDA

- Suppose feature vector $x$ has high-dimension $p$ $(x \in \mathbb{R}^p)$

- Binary classification problem with Training data

- Earlier saw a simple <u>binary</u> classifier with linear $f(x) = \theta^T x - c$

- <u>Fisher Linear Discriminant classifier:</u> $\hat{y} = \text{sgn} f(x)$

- Projection direction $\theta$?

$$\text{max FDR} = \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2}$$

1) Means should be well separated
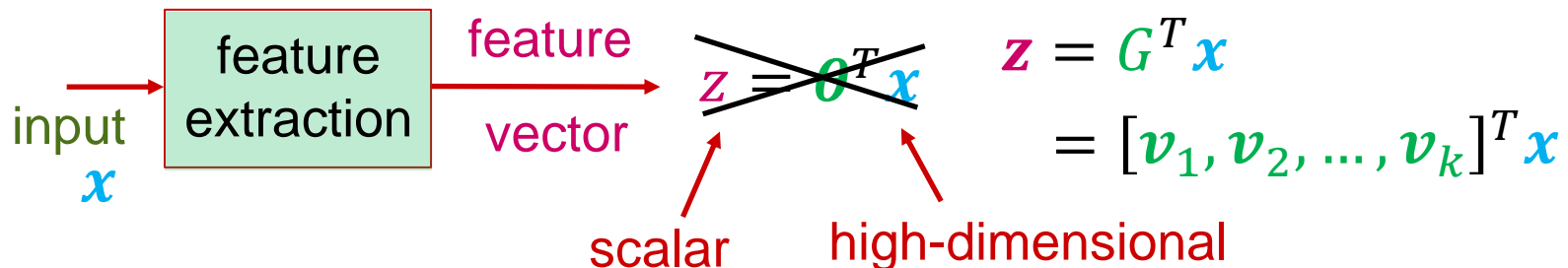
2) Point spreads should be small

Good projection direction

$\theta^T x$ increases

$x_2$

$\hat{y} = +1$  large $\theta^T x$

$\mu_1$

$\theta$

small $\theta^T x$

$c$

$\sigma_1^2$

$\sigma_2^2$  $\mu_2$

$\hat{y} = -1$

Good classifier!

$x_1$

# Motivation: LDA

- For the input vector $\boldsymbol{x}$ of high-dimension $p$ ($\boldsymbol{x} \in \mathbb{R}^p$) ...

- The projection of $\boldsymbol{x}$ onto a line is 'Dimensionality reduction' method
  - Scalar (i.e., $1 - D$) easily visualized.   Less memory to store
  - Also, less parameters $\Rightarrow$ Less computations and less over-fitting in inference

- But projection $z$ can be thought of as an extracted feature vector

- So, any Dimensionality reduction $\Rightarrow$ Automatic Feature extraction

- Vector $\boldsymbol{z} \in \mathbb{R}^k$ possible?  Instead of vector $\boldsymbol{\theta}$, use $p \times k$ matrix $G$

- LDA can be defined now.     (Must have #classes $M \geq k + 1$)

feature extraction

feature vector

input $\boldsymbol{x}$

$z = \boldsymbol{\theta}^T \boldsymbol{x}$

scalar     high-dimensional

$\boldsymbol{z} = G^T \boldsymbol{x}$
$= [\boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_k]^T \boldsymbol{x}$

# 1) Fisher's Linear Discriminant Analysis (LDA)

- Fisher's LDA (<u>vector</u> **z** case):  Optimal matrix $G$ ?    Maximize FDR

**Binary classes**

**$M-$classes  $(M \geq 3)$**

$$\frac{1}{2}\frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2}$$

1) Means should be well separated

$\leftarrow$ FDR $\rightarrow$

2) Spreads should be small

$$\frac{\sum_{i=1}^{M}||\boldsymbol{m}_i - \widehat{\boldsymbol{m}}||^2}{\sum_{i=1}^{M}\sigma_i^2}$$

$$= \frac{(\mu_1 - \widehat{\mu})^2 + (\mu_2 - \widehat{\mu})^2}{\sigma_1^2 + \sigma_2^2}$$

$$= \frac{\sum_{i=1}^{2}(\mu_i - \widehat{\mu})^2}{\sum_{i=1}^{2}\sigma_i^2}$$

Mean of mean vectors $\widehat{\boldsymbol{m}} = \frac{1}{M}\sum_{j=1}^{M}\boldsymbol{m}_j$

$\sigma_i^2 = \text{Trace}[S_i]$

Projected <u>vector</u>

$$\boldsymbol{z} = G^T\boldsymbol{x}$$

(Plane, not line)

$\boldsymbol{m}_i$  (mean vector)

$S_i$  (covariance matrix)

$$\widehat{\mu} \doteq \frac{1}{2}(\mu_1 + \mu_2)$$

# Fisher's LDA

```
sklearn.discriminant_analysis.LinearDiscriminantAnalysis
```

- <u>Fisher's LDA</u>: To max FDR, using input features $\boldsymbol{x}$, calculate

1. Sample means and sample covariance matrices of each class

$$\widehat{\boldsymbol{\mu}}_i = \frac{1}{N_i}\sum_{n=1}^{N_i} \boldsymbol{x}_n \qquad \widehat{\Sigma}_i = \frac{1}{N_i-1}\sum_{n=1}^{N_i}(\boldsymbol{x}_n - \widehat{\boldsymbol{\mu}}_i)(\boldsymbol{x}_n - \widehat{\boldsymbol{\mu}}_i)^T$$

2. 'Within-class' covariance matrix $\widehat{\Sigma}_w \doteq \sum_{i=1}^{M} \widehat{\Sigma}_i$

3. Average $\widehat{\boldsymbol{\mu}}$ of class means $\widehat{\boldsymbol{\mu}}_i$: $\qquad \widehat{\boldsymbol{\mu}} \doteq \frac{1}{M}\sum_{i=1}^{M} \widehat{\boldsymbol{\mu}}_i$

4. 'Between-class' covariance matrix $\widehat{\Sigma}_b \doteq \frac{1}{M-1}\sum_{i=1}^{M}(\widehat{\boldsymbol{\mu}}_i - \widehat{\boldsymbol{\mu}})(\widehat{\boldsymbol{\mu}}_i - \widehat{\boldsymbol{\mu}})^T$

5. Max FDR $\rightarrow$ EVD of matrix $B \doteq (\widehat{\Sigma}_w + \lambda I)^{-1}\widehat{\Sigma}_b$     "top-k eigen-vectors"

         Regularize

6. <u>Largest</u> $k$ eigen-values $\rightarrow$ Eigen-vectors collected matrix $G = [\boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_k]$

7. Then, Fisher's LDA $\boldsymbol{z}_n = G^T \boldsymbol{x}_n$

8. Ridge-regularize if $G$ is large matrix    $\text{Reg\_FDR} = \dfrac{\text{Trace}[G^T(\widehat{\Sigma}_b)G]}{\text{Trace}[G^T(\widehat{\Sigma}_w)G] + \underline{\lambda\|G\|}_F^2}$

- Crime level data:   **Communities and Crime Unnormalized Set ⋆**

  — Crime levels in different communities

- Attributes:   145 attributes per community

  1. `population`

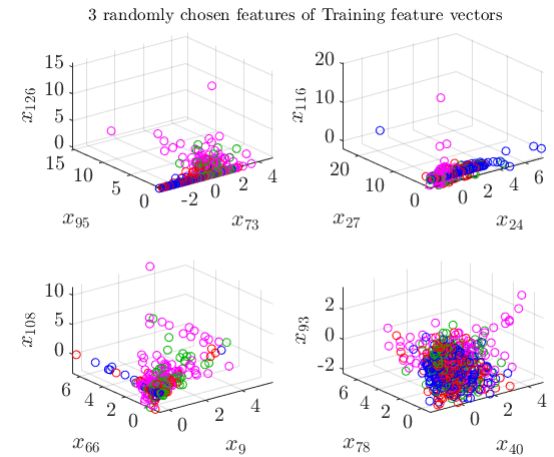  2. `household size`

  3. `% population below poverty`

  4. `...`

  145. `Crime per 100K population`



3 randomly chosen features of Training feature vectors

$$\omega = \begin{cases} \omega_1, & Crime < 25\ percentile \\ \omega_2, & 25 \leq Crime < 50 \\ \omega_3, & 50 \leq Crime < 75 \\ \omega_4, & Crime \geq 75\ percentile \end{cases}$$

- Convert `Crime` into <u>four classes</u>

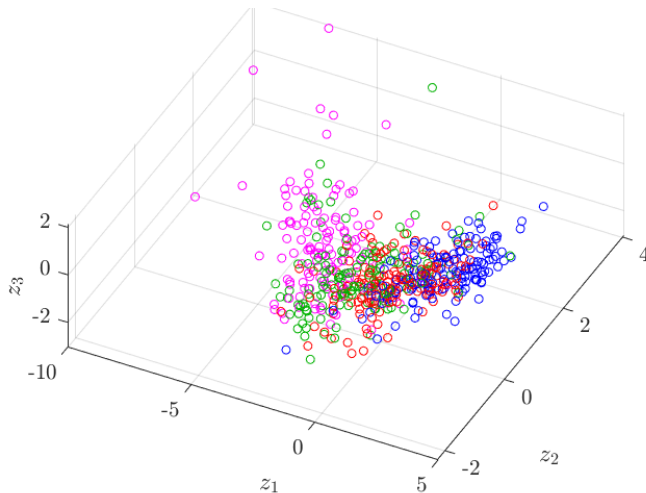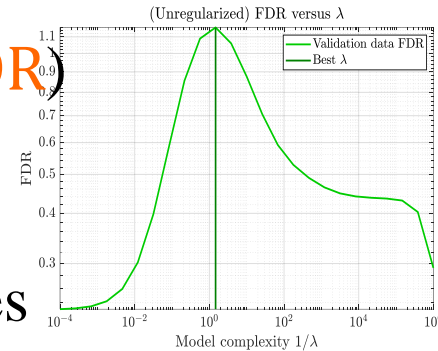- Attributes 1-127 are features $x$.     Any 3 features not informative

- So, visualize using <u>Fisher projected vectors</u> $z_n = G^T x_n \in \mathbb{R}^3$

`http://archive.ics.uci.edu/ml/datasets/Communities+and+Crime+Unnormalized`

# Example:    Fisher's LDA

- Regularization used.   So, normalize by  center and scaling each feature

- Validation of penalty $\lambda$ by plotting FDR (not Reg_FDR)



- Using optimal $G$ ...

- LDA $\mathbf{z}_n$  of Testing data   shows separation of classes

- Projection using a randomly chosen $G$  $\rightarrow$  No separation

-  $\Rightarrow$   LDA does class-aware <u>Automatic feature extraction</u>
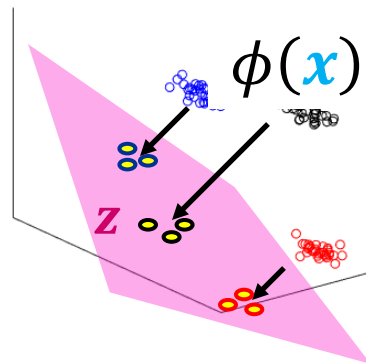
# 2) Kernel LDA

$$\boldsymbol{z}_n = G^T \phi(\boldsymbol{x}_n)$$

■ <u>New idea</u>:  First,  map features $\boldsymbol{x}_n$ to high-dimensional vector $\phi(\boldsymbol{x}_n)$

<u>Then</u> apply LDA to $\phi(\boldsymbol{x}_n)$

■ Why?  High dimensional mapping $\phi(\boldsymbol{x}_n)$  may fit data better

— Recollect:  Common in GLMs

— e.g.,  $\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix} \rightarrow \boldsymbol{\theta}^T \phi(x) = \theta_0 + \theta_1 x + \theta_2 x^2$  (Polynomial)



$\phi(\boldsymbol{x})$

$\boldsymbol{z}$

# Kernel LDA

$$\mathbf{z}_n = G^T \phi(\boldsymbol{x}_n) = \tilde{G}^T \boldsymbol{\kappa}_n$$

- <u>Problem</u>: For very high-dimensional $\phi(\boldsymbol{x}_n)$, computing $G$ is hard

- <u>Kernel idea</u>: Not hard if mapping $\phi(\cdot)$ is based on a kernel

  — Kernel? A measure of similarity. E.g., $\kappa(\boldsymbol{x}_m, \boldsymbol{x}_n) = e^{-\|\boldsymbol{x}_m - \boldsymbol{x}_n\|^2 / \sigma^2}$

  (RBF kernel)

  — (RKHS theory of Non-parametric inference later)

- So, <u>Kernel LDA</u>: LDA applied to high-dimensional $\phi(\boldsymbol{x}_n)$

1. Each $\boldsymbol{x}_n \rightarrow$ Kernel vector

   <u>Not</u> $\phi(\boldsymbol{x}_n)$ $\longrightarrow$ $\boldsymbol{\kappa}_n = \begin{pmatrix} \kappa(\boldsymbol{x}_1, \boldsymbol{x}_n) \\ \kappa(\boldsymbol{x}_2, \boldsymbol{x}_n) \\ \vdots \\ \kappa(\boldsymbol{x}_N, \boldsymbol{x}_n) \end{pmatrix}$

   $\dim(\boldsymbol{\kappa}_n) = N$

   sample mean

   $\hat{\boldsymbol{\mu}}_i = \dfrac{1}{N_i} \displaystyle\sum_{n=1}^{N_i} \boldsymbol{\kappa}_n$

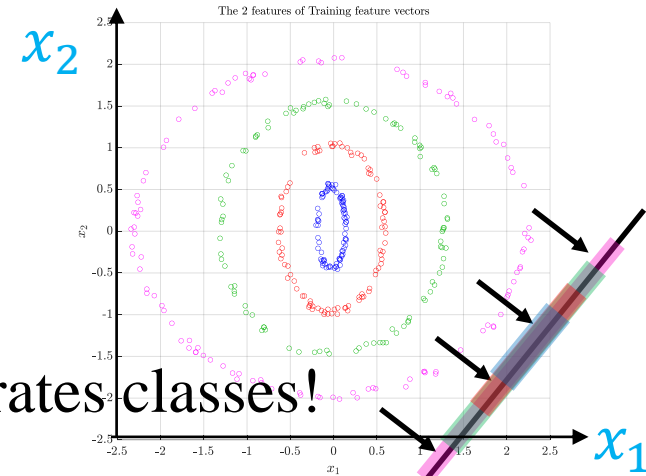   $\vdots$

2. Do <u>LDA calculation</u> on $\boldsymbol{\kappa}_n$ (instead of on $\boldsymbol{x}_n$) to get $\tilde{G}$ matrix

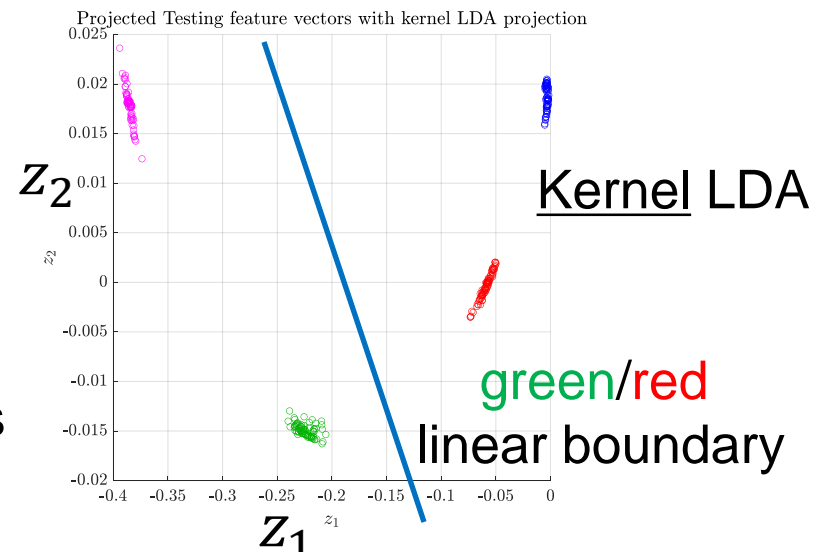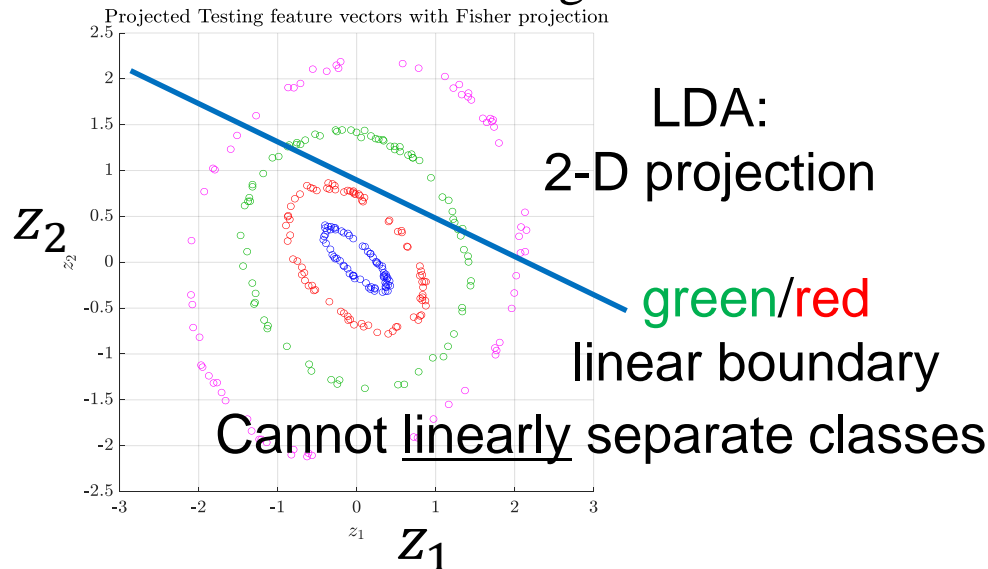3. Projection uses $\tilde{G}$ matrix operating on $\boldsymbol{\kappa}_n$ (instead of on $\phi(\boldsymbol{x}_n)$)

# Example:   Kernel LDA

- <u>Synthetic</u> data with $x \in \mathbb{R}^2$ of 4 classes

- LDA → Cannot <u>linearly</u> separate classes

  — ⇒ LDA features are uninteresting

- <u>Kernel</u> LDA projection into $\mathbb{R}^2$ <u>linearly</u> separates classes!

  — Since it is still <u>linear</u>, but in <u>high-dimensional</u> $\phi(x)$ space

  — These are interesting features!



The 2 features of Training feature vectors

LDA:
1-D projection

LDA:
2-D projection

green/red
linear boundary

Cannot <u>linearly</u> separate classes

Projected Testing feature vectors with Fisher projection

Projected Testing feature vectors with kernel LDA projection

Kernel LDA

green/red
linear boundary

# PCA,
# Kernel PCA,
# LSA

# 3) Principal Component Analysis (PCA)

$$\mathbf{z}_n = G^T \mathbf{x}_n \qquad \text{Class} \qquad \widehat{\Sigma}_i = \frac{1}{N_i - 1} \sum_{n=1}^{\boxed{N_i} = 1} \underbrace{(\mathbf{x}_n - \widehat{\boldsymbol{\mu}}_i)}_{= 0}(\mathbf{x}_n - \widehat{\boldsymbol{\mu}}_i)^T = 0$$

sample covariance

- PCA is for unlabeled data.   Suppose $N$ points of <u>unlabeled</u> $\mathbf{x}_n \in \mathbb{R}^p$

- <u>No classes to separate</u> here.    "Best" matrix $G$ for Fisher LDA?

- Note:    "No classes"  $\leftrightarrow$   Every point $\mathbf{x}_n$ is "its own class"

- So <u>LDA calculation</u> reduces to …

- PCA:   Uses matrix $G = [\boldsymbol{v}_1, \boldsymbol{v}_2, \dots, \boldsymbol{v}_k] = $ Top-$k$ eigenvectors of $\widehat{\Sigma}$

Sample mean of <u>all points</u>
$$\widehat{\boldsymbol{\mu}} \doteq \frac{1}{N} \sum_{n=1}^{N} \mathbf{x}_n \qquad\qquad \widehat{\Sigma} \doteq \frac{1}{N-1} \sum_{i=1}^{N} (\mathbf{x}_i - \widehat{\boldsymbol{\mu}})(\mathbf{x}_i - \widehat{\boldsymbol{\mu}})^T$$

sample covariance matrix of <u>all points</u>

- (Alternative view:   PCA maximizes randomness of projected $\mathbf{z}_n$)

# Example:    PCA

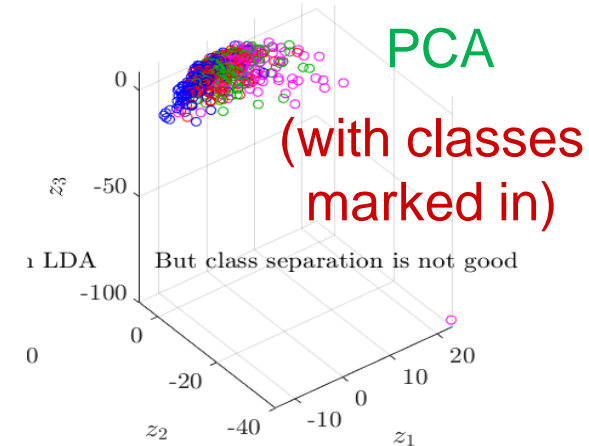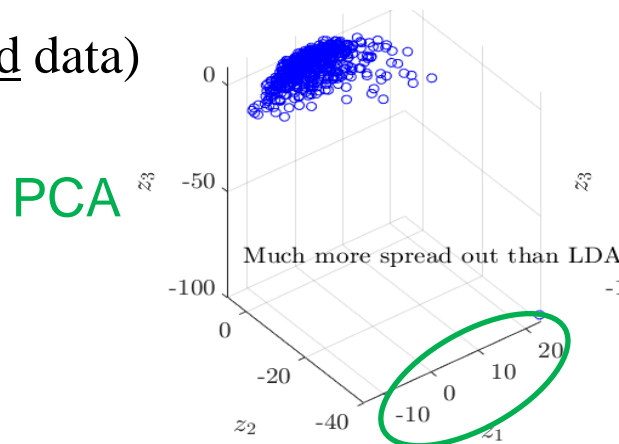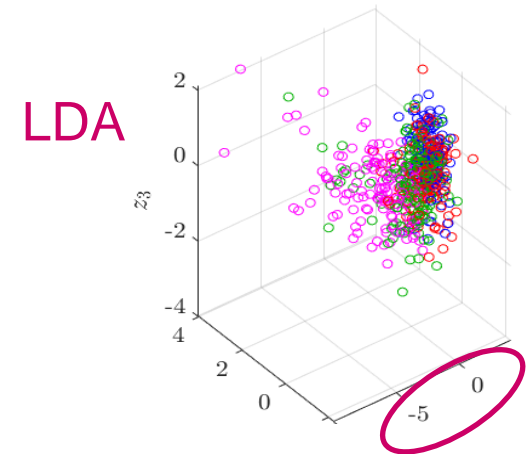- Crime level data:    **Communities and Crime Unnormalized Set**

  — 127-dimensional feature vectors $x$

  — LDA gives some class separation here

  — PCA?   Ignore the classes (label)

- PCA spreads out $z_n$ much more than LDA does

  — Because PCA maximizes sample variance of $z_n$

- PCA does not separate classes, unlike LDA

  — (PCA is for <u>unlabeled</u> data)



LDA



PCA

Much more spread out than LDA

PCA

(with classes marked in)

But class separation is not good

$$\widehat{\Sigma} \doteq \frac{1}{N-1} \sum_{n=1}^{N} (\boldsymbol{x}_n - \widehat{\boldsymbol{\mu}})(\boldsymbol{x}_n - \widehat{\boldsymbol{\mu}})^T = V \, L \, V^T \quad \textbf{EVD}$$

Columns $\boldsymbol{v}_i$ are eigen-vectors

Diagonal $L_{ii}$ are eigen-values

- PCA matrix $G$ calculated using EVD of data covariance matrix

- Covariance matrix $\widehat{\Sigma} = \left(\frac{1}{\sqrt{N-1}} X_c\right)^T \left(\frac{1}{\sqrt{N-1}} X_c\right)$

$$= A^T A$$

- Matrix of form such as $\widehat{\Sigma}$ is called a 'Gram matrix'

Matrix of <u>centered</u> feature vectors

$$X_c = \begin{bmatrix} \boldsymbol{x}_1^T - \widehat{\boldsymbol{\mu}}^T \\ \boldsymbol{x}_2^T - \widehat{\boldsymbol{\mu}}^T \\ \vdots \\ \boldsymbol{x}_N^T - \widehat{\boldsymbol{\mu}}^T \end{bmatrix}$$

# Singular Value Decomposition (SVD)

$$A \; = \; U \; D \; V^T \qquad \textbf{SVD}$$

Diagonal elements $D_{ii}$ are 'singular-values'

Columns $\boldsymbol{v}_i$ are 'right singular vectors'

- For any Gram matrix $\hat{\Sigma} = \; A^T A$, EVD can also be calculated as …

```
[U,D,V] = svd(A)        U,D,Vt = numpy.linalg.svd(A)
```

- Singular Value Decomposition (SVD) of matrix $A \; = \; \dfrac{1}{\sqrt{N-1}} X_c$

  — Right singular vectors of $A$ are eigen-vectors of $\hat{\Sigma}$

  — Singular values $D_{ii}$ of $A$ relate to eigen-values $L_{ii}$ of $\hat{\Sigma}$ as $D_{ii} = \sqrt{L_{ii}}$

  — So, for PCA, find largest $k$ <u>singular</u> values $D_{ii}$
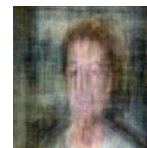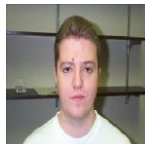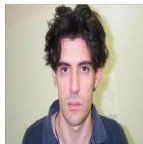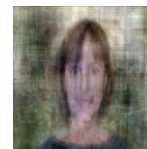    <u>eigen</u>-values $L_{ii}$

  eigen-vectors $\boldsymbol{v}_i$

  — And collect their right singular-vectors $\boldsymbol{v}_i$ into $G = [\boldsymbol{v}_1, \boldsymbol{v}_2, …, \boldsymbol{v}_k]$

  PCA matrix ✓

- Advantage?   Potentially large matrix $\hat{\Sigma}$ not computed anywhere

# Example:    PCA using SVD

- Example:    **Caltech vision images data set (use 435 face images)★**

  — Image dimensions $= 227 \times 227 \times 3 = 154587$

  — ($\hat{\Sigma}$ is $154587 \times 154587$, while $A$ is $435 \times 154587$)   (quite low dimensional)

- So for PCA $G$, use SVD instead of EVD $\rightarrow$   $z_n = G^T x_n \in \mathbb{R}^{50}$

- To visualize $z_n$, map it back to original space $\hat{x}_n = G z_n$ "Eigen-faces"

- Even 50 dimensions in $z_n$ captures quite a bit of image structure



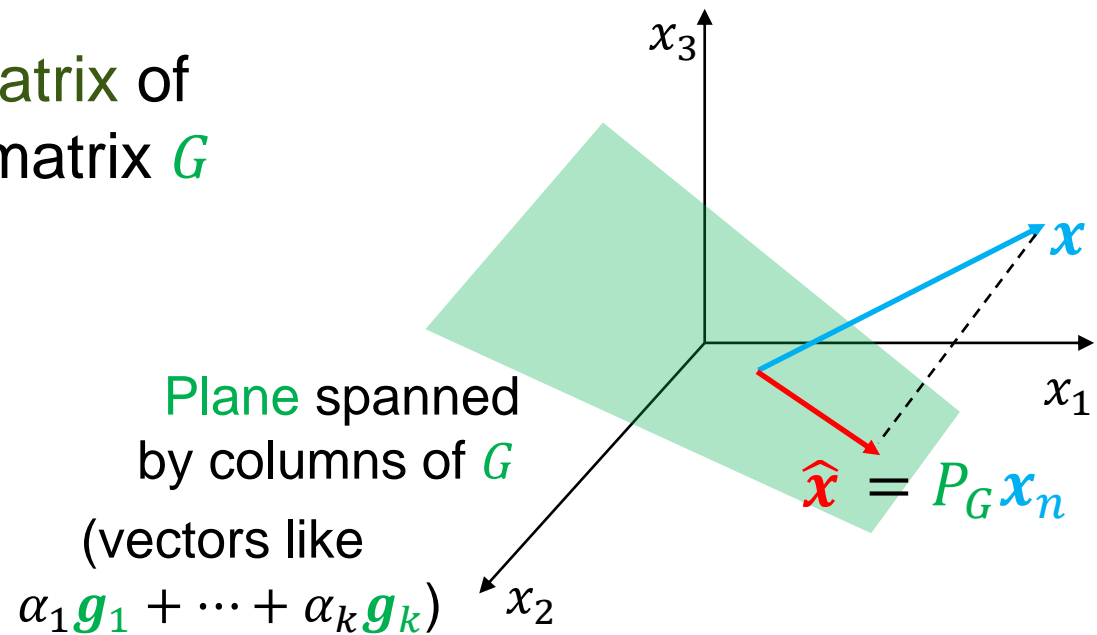http://www.vision.caltech.edu/Image_Datasets/Caltech101/101_ObjectCategories.tar.gz

- Eigen-faces $\widehat{\boldsymbol{x}}_n$:     $\boldsymbol{z}_n$ mapped back to original space

- Essentially, applying a "Projection matrix" to original $\boldsymbol{x}_n$

$$\widehat{\boldsymbol{x}}_n = G \mathbf{z}_n$$

$$= G G^T \boldsymbol{x}_n$$

$$= P_G \boldsymbol{x}_n$$

Projection matrix of
orthonormal matrix $G$

Plane spanned
by columns of $G$
(vectors like
$\alpha_1 \boldsymbol{g}_1 + \cdots + \alpha_k \boldsymbol{g}_k$)

$x_3$

$x_1$

$x_2$

$\boldsymbol{x}$

$\widehat{\boldsymbol{x}} = P_G \boldsymbol{x}_n$

# 4) Kernel PCA

$$\mathbf{z}_n = G^T \phi(\mathbf{x}_n) = \tilde{G}^T \boldsymbol{\kappa}_n$$

- Same motivation as Kernel LDA $\rightarrow$ PCA using transformed $\phi(\mathbf{x}_n)$

- As earlier, computations difficult, unless $\phi(\cdot)$ is based on a kernel

Example (RBF kernel)
$$e^{-\|\mathbf{x}_m - \mathbf{x}_n\|^2/\sigma^2}$$

```
sklearn.decomposition.KernelPCA
```

- Kernel PCA:   To maximize randomness in $\mathbf{z}_n$,
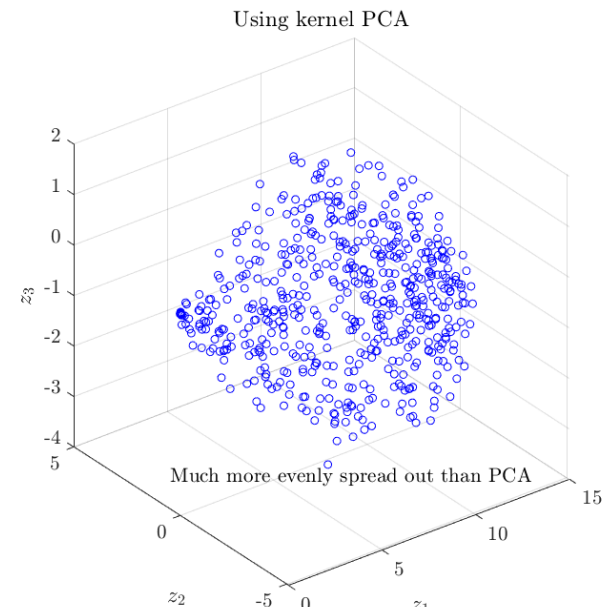
1.  Each $\mathbf{x}_n \rightarrow$ Kernel vector
$$\boldsymbol{\kappa}_n = \begin{pmatrix} \kappa(\mathbf{x}_1, \mathbf{x}_n) \\ \kappa(\mathbf{x}_2, \mathbf{x}_n) \\ \vdots \\ \kappa(\mathbf{x}_N, \mathbf{x}_n) \end{pmatrix}$$

2.  Do PCA calculation on $\boldsymbol{\kappa}_n$ (instead of on $\mathbf{x}_n$) to get $\tilde{G}$ matrix

3.  Projection $\mathbf{z}_n$ uses $\tilde{G}$ matrix operating on $\boldsymbol{\kappa}_n$ (instead of on $\phi(\mathbf{x}_n)$)

# Example:  Kernel PCA

- <u>Crime level data</u>:  **Communities and Crime Unnormalized Set**

- Map  Input vector $x_n \in \mathbb{R}^{127}$ to feature vector  $z_n \in \mathbb{R}^3$

- PCA  does spread out points  (it is variance maximizing)

  — But notice large cluster of points close together

- Kernel PCA  spreads out points much more evenly

# 5) Latent Semantic Analysis (LSA)

- Example: Application of LSA to Text Analysis

  — Dataset of $N$ text documents (e.g., newspaper articles)

  — Each document → 'Bag of words' $x$

$$x = \begin{bmatrix} 3 \\ \vdots \\ 0 \\ 2 \end{bmatrix}$$

← 'prize' occurs 3 times

← 'party' is absent

$n^{th}$ document $x_n$

noise

$h_2$ "politics-type" document

$h_1$ "sports-type" document

$\hat{x}_n$ noiseless document

- <u>Latent factor Model</u>: Assume each document combines $k$ '<u>prototypes</u>'

$$x_n = \hat{x}_n + Gaussian\ noise \qquad \text{where}$$

$$\hat{x}_n = \sum_{i=1}^{k} h_i\, c_{i,n}$$

$i^{th}$ coefficient of $n^{th}$ document $x_n$

$i^{th}$ latent factor vector ("prototype document")

# Latent Semantic Analysis (LSA)

$$x_n = \widehat{x}_n + \text{Gaussian noise} \qquad \text{where}$$

$i^{th}$ coefficient of $n^{th}$ document $x_n$

$$\widehat{x}_n = \sum_{i=1}^{k} h_i \, c_{i,n}$$

$i^{th}$ latent factor vector

- Orthonormal Latent factors $h_1, h_2, \ldots, h_k$ <u>common</u> to all documents

- But, coefficient vector $c_n$ makes document $x_n$ unique

- <u>LSA problem</u>:   Using only Training data $x_1 \ldots, x_N$ (documents)

  Find MLE of all factors $(h_i)$ and coefficients $(c_n)$

$$\widehat{h_i}, \widehat{c_n} = \underset{h_i, c_n}{\arg\min} \sum_{n=1}^{N} \underline{\|x_n - \widehat{x}_n\|^2}$$

Quadratic NLL (*Gaussian noise*)

coefficient vector $c_n \doteq \begin{bmatrix} c_{1,n} \\ c_{2,n} \\ \vdots \\ c_{k,n} \end{bmatrix}$ of $x_n$

$$\checkmark \quad \widehat{h_i}, \widehat{c_n} = \arg\min_{h_i, c_n} \sum_{n=1}^{N} \| x_n - \widehat{x}_n \|^2$$

- Closed form solution is possible here.    Called ...

- Latent Semantic Analysis (LSA):

  — "Document-Term" feature matrix $X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} = \begin{bmatrix} \text{document 1 bag-of-words} \\ \text{document 2 bag-of-words} \\ \vdots \\ \text{document } N \text{ bag-of-words} \end{bmatrix}$

  $$= \ U \ D \ V^T \quad \text{(SVD)}$$

  — PCA matrix   using SVD method    $G = [v_1, v_2, \ldots, v_k]$    (top$-k$ right-singular vectors)

  — Projection of $x_n \ \rightarrow \ z_n = \ G^T x_n$    (feature vectors)

  — Then,  a) MLE of coefficient vector  $\widehat{c_n} = z_n$

  — Then,  b) MLE of latent factor  $\widehat{h_i} = v_i$   (right-singular vector)

# Example:   LSA for Text Documents

- Example:   **BBC news articles dataset***

  — 1490 BBC news articles in 5 categories

- Document-Term matrix using tf-idf

- Compute LSA $G$, then project $\boldsymbol{z}_n = G^T \boldsymbol{x}_n \in \mathbb{R}^k$

  — (For Text, generally choose large $k$, say 50)

  — Here, choose $k = 3$ for easy visualization

| | ArticleId | Text | Category |
|---|---|---|---|
| 0 | 1833 | worldcom ex-boss launches defence lawyers defe... | business |
| 1 | 154 | german business confidence slides german busin... | business |
| 2 | 1101 | bbc poll indicates economic gloom citizens in ... | business |
| 3 | 1976 | lifestyle governs mobile choice faster bett... | tech |
| 4 | 917 | enron bosses in $168m payout eighteen former e... | business |
| ... | ... | ... | ... |
| 1485 | 857 | double eviction from big brother model caprice... | entertainment |
| 1486 | 325 | dj double act revamp chart show dj duo jk and ... | entertainment |
| 1487 | 1590 | weak dollar hits reuters revenues at media gro... | business |
| 1488 | 1587 | apple ipod family expands market apple has exp... | tech |
| 1489 | 538 | santy worm makes unwelcome visit thousands of ... | tech |

**3078 dictionary terms**

**1490 documents**

Document-Term matrix $X$

| | 000 | 10 | 100 | 100m | 11 | 12 | 120 | 13 | 14 | 15 | ... | year | years | yen | yes | yet | york | young | younger | yukos | zealand |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.000000 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.000000 | 0.033765 | 0.0 | 0.0 | 0.000000 | 0.057611 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 1 | 0.000000 | 0.041944 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.025758 | 0.000000 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 2 | 0.028581 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.042893 | 0.044012 | 0.080704 | ... | 0.000000 | 0.025569 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 3 | 0.023337 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.071872 | 0.000000 | ... | 0.000000 | 0.020877 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.094044 | 0.0 | 0.0 |
| 4 | 0.000000 | 0.038268 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.000000 | 0.032200 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 1485 | 0.056301 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.071659 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.036759 | 0.000000 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 1486 | 0.000000 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.016108 | 0.022071 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 1487 | 0.000000 | 0.049054 | 0.0 | 0.0 | 0.12762 | 0.058725 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.210869 | 0.000000 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 1488 | 0.000000 | 0.053066 | 0.0 | 0.0 | 0.00000 | 0.031764 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.016294 | 0.000000 | 0.0 | 0.0 | 0.032494 | 0.000000 | 0.036562 | 0.000000 | 0.0 | 0.0 |
| 1489 | 0.035532 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.000000 | 0.000000 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |

`https://www.kaggle.com/c/learn-ai-bbc`

- Note: LSA did not use labels. Yet, documents get clustered nicely!

- ⇒ LSA "discovered" underlying 'prototype documents'

- PCA or LSA → $k$ selected using "elbow" of singular values plot

Singular value $D_{ii}$ of matrix $X$

$k = 6$ dimensions better choice?

LSA-projected vectors $\boldsymbol{z}_n \in \mathbb{R}^3$

Tech and Politics are far

Tech and business are close

Document-Term matrix $X$

| | 000 | 10 | 100 | 100m | 11 | 12 | 120 | 13 | 14 | 15 | ... | year | years | yen | yes | yet | york | young | younger | yukos | zealand |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.000000 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.000000 | 0.033765 | 0.0 | 0.0 | 0.000000 | 0.057611 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 1 | 0.000000 | 0.041944 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.025758 | 0.000000 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 2 | 0.028581 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.042893 | 0.044012 | 0.080704 | ... | 0.000000 | 0.025569 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 3 | 0.023337 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.071872 | 0.000000 | ... | 0.000000 | 0.020877 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.094044 | 0.0 | 0.0 |
| 4 | 0.000000 | 0.038268 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.000000 | 0.032200 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| ... | | | | | | | | | | | ... | | | | | | | | | | |
| 1485 | 0.056301 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.071659 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.036759 | 0.000000 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 1486 | 0.000000 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.016108 | 0.022071 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 1487 | 0.000000 | 0.049054 | 0.0 | 0.0 | 0.12762 | 0.058725 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.210869 | 0.000000 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |
| 1488 | 0.000000 | 0.053066 | 0.0 | 0.0 | 0.00000 | 0.031764 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.016294 | 0.000000 | 0.0 | 0.0 | 0.032494 | 0.000000 | 0.036562 | 0.000000 | 0.0 | 0.0 |
| 1489 | 0.035532 | 0.000000 | 0.0 | 0.0 | 0.00000 | 0.000000 | 0.0 | 0.000000 | 0.000000 | 0.000000 | ... | 0.000000 | 0.000000 | 0.0 | 0.0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 |

politics
sport
entertainment
tech
business

# Independent Component Analysis (ICA)

- Cocktail party problem :    Illustrates ICA

- $k$ independent audio signals $c_{i,n}$  →   Assume at least $k$ microphones

- Goal:   From mixed samples vector $x_n$,  estimate separate speeches $c_n$

$= h_1^T c_n + noise$              $x_n$   (Mixed speech samples vector at time $n$)

$= \sum_{i=1}^{k} h_{1,i} c_{i,n} + noise$

Vectorize

$x_{p,n} = h_p^T c_n + noise$

$x_{1,n}$

$x_{2,n}$

$\ldots, c_{k,n}, \ldots$

$\ldots, c_{2,n}, \ldots$

$\ldots, c_{1,n}, \ldots$ speech

$$x_{i,n} = \mathbf{h}_i^T \mathbf{c}_n + noise, \qquad i = 1,2,\ldots,p$$

$$p \geq k \text{ microphones}$$

vectorize

$$H = \begin{bmatrix} \mathbf{h}_1^T \\ \vdots \\ \mathbf{h}_p^T \end{bmatrix} \text{ mixing matrix}$$

$$\mathbf{x}_n = H\mathbf{c}_n + Gaussian\ noise$$

- This is usual Normal Discriminative model (with vector label $\mathbf{c}$)

- $\Rightarrow$ Solved by Linear Statistical Regression if mixing matrix $H$ known
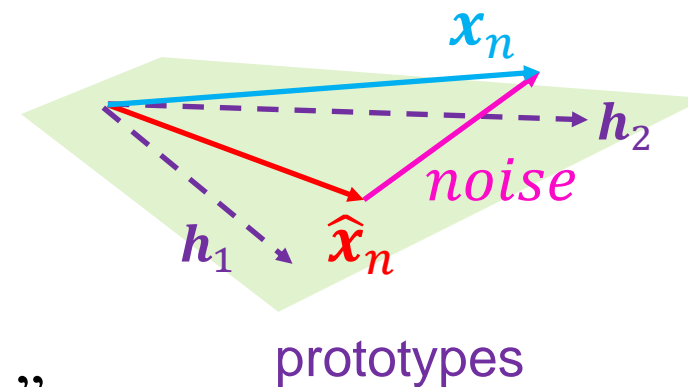
- But mixing matrix $H$ is <u>unknown</u>, so use ICA

- Each EEG signal $x_{i,n}$ is linear mixture of several signals $c_{m,n}$

- <u>Goal</u>:    Isolate brain waves from other artifacts

- (i.e., Pre-processing step in Feature extraction)

brain waves,

eye-blinking,

heart contractions,

…

$x_{1,n}$

$n$

EEG
signals

$x_{2,n}$

$x_{3,n}$

Multichannel
EEG

EEG
signals
input

ICA

pre-processing

$\cdots$

$x$

brain
waves

# 6) Independent Component Analysis (ICA)

- ICA model:  $\boldsymbol{x}_n = H\boldsymbol{c}_n + $ *Gaussian noise*

  $$= \widehat{\boldsymbol{x}}_n + \text{ } \textit{Gaussian noise}$$

- But this is just LSA model!

- But LSA assumed ~~factor matrix $H$ is orthonormal~~.   Instead  ...  **No**

- <u>Independent Component Analysis (ICA)</u>:   Assume

  — Elements $c_{i,n}$ of $\boldsymbol{c}_n$ are  <u>independent</u> signals "components"

  — $c_{i,n}$  are <u>non-Gaussian</u> variables

- <u>Goal of ICA</u>:   Estimate $H$ and  $\boldsymbol{c}_n$

  $H$ unknown        $c_{i,n}$

- ICA is example of "Blind source separation"

$\boldsymbol{x}_n$

$\boldsymbol{h}_2$

*noise*

$\boldsymbol{h}_1$    $\widehat{\boldsymbol{x}}_n$

prototypes

# ICA algorithms

$$x_n = Hc_n + noise$$

- There are many <u>heuristic</u> ICA algorithms

- General idea: Estimate $c_n$ from $x_n$ <u>linearly</u> as $z_n = G^T x_n$

- Optimal "unmixing matrix" $G \approx H^{-1}$

  $(\approx c_{i,n}$ wanted)

- But $H$ unknown. So, find $G$ that makes elements $z_i$ of $z_n$

  "as independent and non-Gaussian" as possible, just like $c_{i,n}$
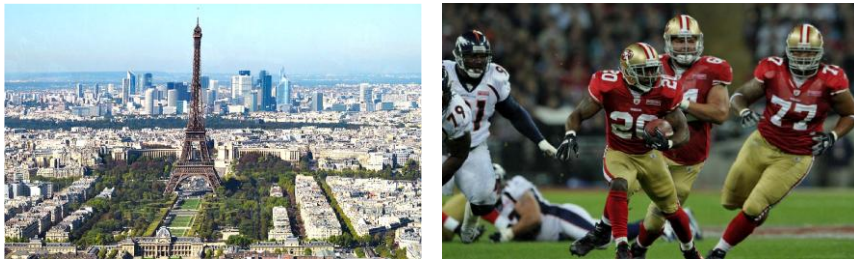
- Example (<u>FastICA algorithm</u>): Minimizes "Entropy"

  ```
  sklearn.decomposition.FastICA()
  ```

Two Mixed images

$x_n = [x_{1,n}, x_{2,n}]^T$ is $n^{th}$ pixel pair

ICA's optimal $G$ matrix always has a scaling ambiguity



Optimally unmixed images (FastICA)

$$z_n = [z_{1,n}, z_{2,n}]^T$$



Unmixed and then scaled

$$D\, z_n = [\lambda_1 z_{1,n}, \quad \lambda_2\, z_{2,n}]^T$$

Underlying true images

$$c_n = [c_{1,n}, c_{2,n}]^T$$

# Non-linear methods

- We have explored Linear feature extraction $\mathbf{z}_n = \cancel{G^T \mathbf{x}_n} \ g(\mathbf{x}_n)$

- Auto-encoder: Non-linear Feature extraction using <u>self-prediction</u>

  — Has an Encoder: $\mathbf{x} \rightarrow \mathbf{z} \in \mathbb{R}^k$ (low dimensional)

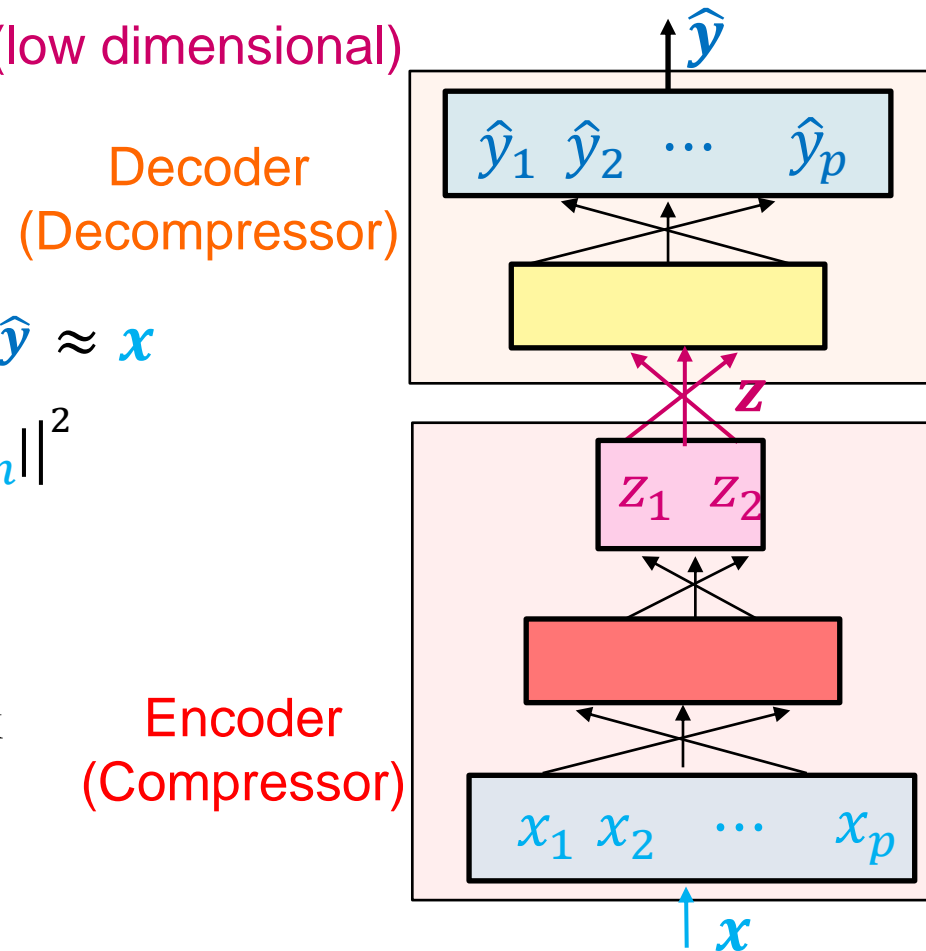  — Has a Decoder: $\mathbf{z} \rightarrow \widehat{\mathbf{y}}$

  — Train Encoder and Decoder to make $\widehat{\mathbf{y}} \approx \mathbf{x}$

  $$\text{minimize} \ \sum_{n=1}^{N} || \widehat{\mathbf{y}}_n - \mathbf{x}_n ||^2$$
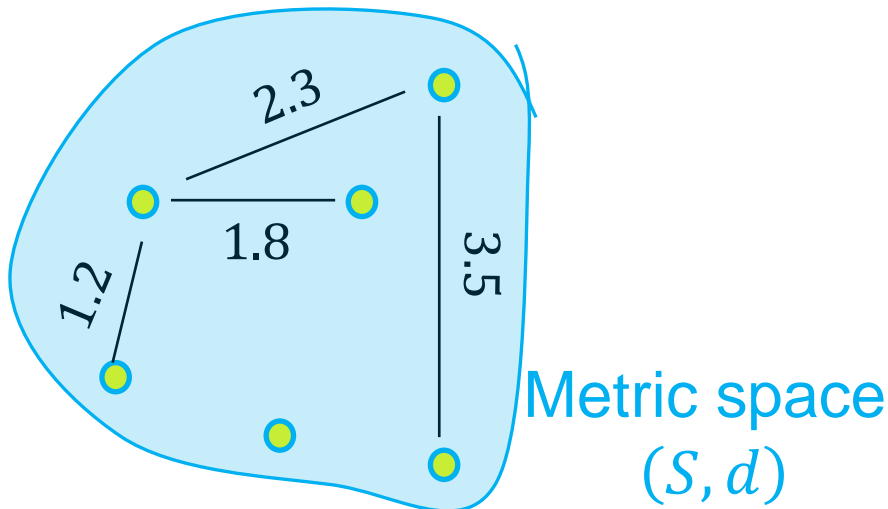
  — i.e., Self-prediction of input $\mathbf{x}$

  — Nowadays uses Deep neural network

  — (Why useful? Details later)



Decoder (Decompressor)

$\widehat{y}_1 \ \widehat{y}_2 \ \cdots \ \widehat{y}_p$

$\mathbf{z}$

$z_1 \ z_2$

Encoder (Compressor)

$x_1 \ x_2 \ \cdots \ x_p$

$\mathbf{x}$

# Metric space

- Until now, "input" was <u>vector $\boldsymbol{x}$</u>.     Now, more general inputs …

- Set  $S = \{s_1, s_2, \ldots, s_N\}$  of  $N$ inputs   (Abstract ideas, not vectors)

  — Set of images,   Set of documents,   Set of signals,   …

- Suppose distance function $d(s_i, s_j)$ measures distance between $s_i, s_j$

- <u>Metric space</u> $(S, d)$:    Set $S$  with $d(s_i, s_j)$ assigned  to <u>each</u> pair $s_i, s_j$

  — Example:    Euclidean space $\mathbb{R}^p$  is a Metric space with $d(\boldsymbol{s}_1, \boldsymbol{s}_2) \doteq ||\boldsymbol{s}_1 - \boldsymbol{s}_2||$



Metric space
$(S, d)$

Distance function $d(\cdot, \cdot)$
measures 'dissimilarity' of points
$$d(s_1, s_2) = 0 \quad \Leftrightarrow \quad s_1 = s_2$$
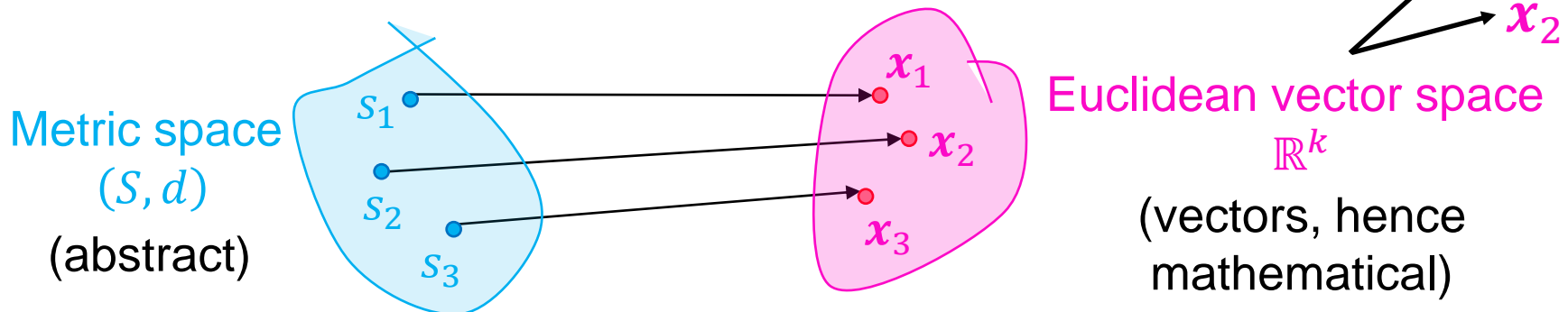$$d(s_1, s_2) = d(s_2, s_1)$$
$$d(s_1, s_3) \leq d(s_1, s_2) + d(s_2, s_3)$$
triangle-inequality

# Embedding

$$x_n = g(\,s_n\,)$$

- Until now, Feature extraction assumed <u>vector</u> input $x$

- Now, from (possibly abstract) input $s$, get feature vector $x$

- <u>Embedding</u> of Metric space $(S, d)$ : For each <u>abstract</u> input $s_n$ in it

  map it to feature <u>vector</u> $x_n \in \mathbb{R}^k$

  — Why? For Automatic feature extraction of abstract ideas $s_n$, or

  — Visualization (if dimension $k = 1,2,3$)

  — Extracted feature vector $x_n$ is called "Embedding" of $s_n$

$x_1$

$x_2$

Metric space
$(S, d)$
(abstract)

$s_1$

$s_2$

$s_3$

$x_1$

$x_2$

$x_3$

Euclidean vector space
$\mathbb{R}^k$

(vectors, hence
mathematical)

- Multi-Dimensional Scaling (MDS): Embedding $\boldsymbol{x}$ found by

  retaining the distance geometry between pairs of points

- Heuristic algorithms to solve. e.g., Metric MDS (see handout)
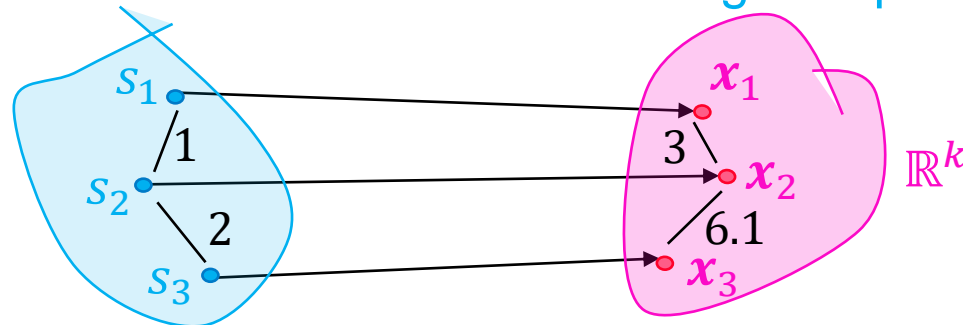
  Python: sklearn.manifold.MDS

$$\min_{\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_N} \sum_{i=1}^{N} \sum_{j=1}^{N} \left( d(s_i, s_j) - \|\boldsymbol{x}_i - \boldsymbol{x}_j\| \right)^2$$

Choose embeddings
that minimize … the total
squared error between ... distance in original space … and distance
in new space

Metric space
$(S, d)$

$s_1$   $x_1$

1   3

$s_2$   $x_2$   $\mathbb{R}^k$

2   6.1

$s_3$   $x_3$

# Metric MDS algorithm

- Metric MDS algorithm:    Heuristic to embed $(S, d)$ in $\mathbb{R}^k$,

  1. Matrix of <u>squared-distances</u>

  $$D_2 \doteq \begin{bmatrix} \left(d(s_1, s_1)\right)^2 & \cdots & \left(d(s_1, s_N)\right)^2 \\ \vdots & \ddots & \vdots \\ \left(d(s_N, s_1)\right)^2 & \cdots & \left(d(s_N, s_N)\right)^2 \end{bmatrix}$$

  2. 'Double-centered' matrix

  $$B = -\frac{1}{2} C \, D_2 \, C$$

  3. EVD of matrix $B = V \Lambda V^T$

  $$C = I_N - \frac{1}{N}\begin{bmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{bmatrix}$$

  4. Choose largest $k$ eigen-values $\lambda_1, \lambda_2, \ldots, \lambda_k$ and eigen-vectors $\boldsymbol{v}_1, \boldsymbol{v}_2, \ldots, \boldsymbol{v}_k$

  5. Calculate feature vector matrix $\rightarrow$ The $N$ rows of $X$ are the embeddings $\boldsymbol{x}_n$

  $$X = \begin{bmatrix} \lambda_1 \boldsymbol{v}_1 & \lambda_2 \boldsymbol{v}_2 & \ldots & \lambda_k \boldsymbol{v}_k \end{bmatrix} = \begin{bmatrix} \boldsymbol{x}_1^T \\ \boldsymbol{x}_2^T \\ \vdots \\ \boldsymbol{x}_N^T \end{bmatrix}$$
  $(k$ columns in $X)$

- Heuristic <u>exactly</u> solves MDS if $(S, d)$ is already <u>Euclidean space</u>

- Another popular Embedding is …

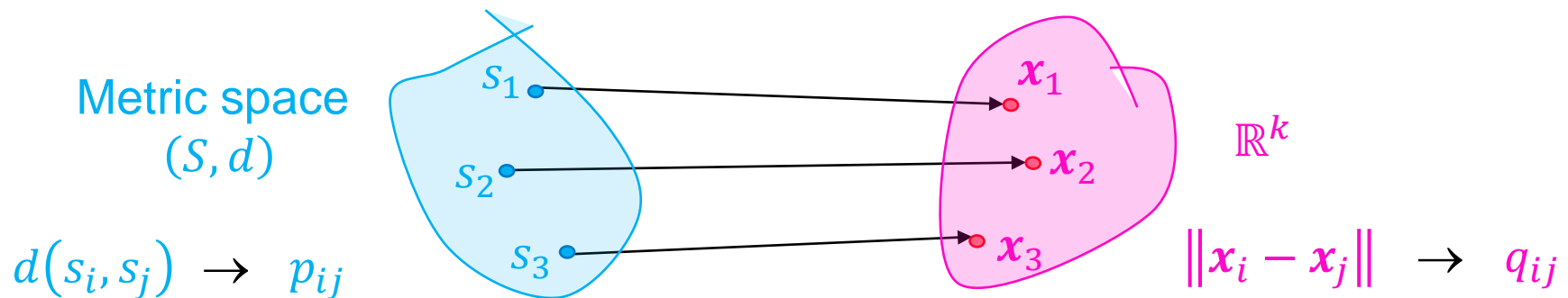- <u>t-distributed Stochastic Neighbor Embedding</u> (t-SNE):

    Similar to MDS, but uses probability models to define distances

    — Algorithm tries to make  $q_{ij} \approx p_{ij}$   $\Rightarrow$   $\|x_i - x_j\| \approx d(s_i, s_j)$

    — See handout for some details

    Python:    sklearn.manifold.TSNE

- (Later:  Embeddings using <u>Deep neural networks</u> becoming popular)



Metric space
$(S, d)$

$s_1$

$s_2$

$s_3$

$x_1$

$x_2$

$x_3$

$\mathbb{R}^k$

$d(s_i, s_j) \rightarrow p_{ij}$

$\|x_i - x_j\| \rightarrow q_{ij}$

# t-SNE

- <u>t-distributed Stochastic Neighbor Embedding</u> (t-SNE):

  — Define conditional probability $p_{j|i} \doteq \begin{cases} \frac{1}{Z_i} \exp\left( -\frac{1}{2\sigma^2}\left( d(s_i, s_j) \right)^2 \right), & j \neq i \\ 0, & j = i \end{cases}$

  — Intuition: $p_{j|i}$ is Gaussian-based probability of "point $s_i$ picking point $s_j$"

  — Define $p_{ij} = \frac{1}{2N}(p_{j|i} + p_{i|j})$

  — In Euclidean space $\mathbb{R}^k$, define $q_{ij} \doteq \begin{cases} \frac{1}{Z}\left( 1 + \|x_i - x_j\|^2 \right)^{-1}, & j \neq i \\ 0, & j = i \end{cases}$

  — Intuition: Like $p_{ij}$, but using a heavy-tailed "Student-t pdf"

  — Find the embeddings $x_i$ that minimize Kullback-Leibler distance

  $$\min_{x_1, x_2, \ldots, x_N} \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

- **`Senators in 116ᵗʰ Senate*`**: Voting record of 102 Senators

- <u>Task</u>: Embed Senators into $\mathbb{R}^2$ based on voting record (ideology)

- <u>Metric space</u> $(S, d)$: $S = \{s_1, s_2, \ldots, s_{102}\}$ is set of Senators

  —Distance $d$? Define $d(s_i, s_j)$ as number of times Senators $s_i, s_j$ voted differently

  —$d(s_i, s_j)$ is called 'Hamming distance' → So, $(S, d)$ is indeed Metric space

- Metric MDS to embed this metric space into $\mathbb{R}^2$ ($k = 2$)

- ⇒ Each Senator is now vector $\boldsymbol{x_n} \in \mathbb{R}^2$. Plot them!

| Roll call | Senator | Vote |
|---|---|---|
| 1 | Collins | Nay |
| 1 | Manchin | Yea |
| ⋮ | ⋮ | ⋮ |
| 2 | Collins | Yea |
| 2 | Manchin | Yea |
| ⋮ | ⋮ | ⋮ |



`https://voteview.com/data`
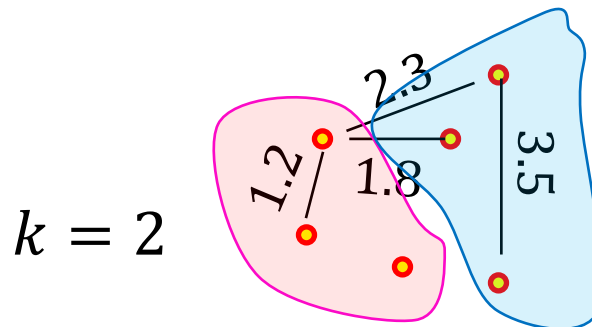
- t-SNE popular alternative to PCA to visualize unlabeled <u>vectors</u> $\boldsymbol{s}_n$ ...

- Here, input Metric space is <u>Euclidean space</u> $\rightarrow$ $d(\boldsymbol{s}_i, \boldsymbol{s}_j) = \|\boldsymbol{s}_i - \boldsymbol{s}_j\|$

- Recollect audio demo : Beethoven piece and Mozart piece

- 20-dim MFCC feature vectors $\boldsymbol{s}_n$

- Ignore $\boldsymbol{s}_n$ labels now $\rightarrow$ Get t-SNE embeddings $\boldsymbol{x}_n \in \mathbb{R}^3$

- Mozart and Beethoven embeddings clustered



Time domain

MFCC
$\boldsymbol{s}_n$

segment $n$

# Clustering

# 9) k-Means clustering

- Suppose unlabeled Training data $s_1, s_2, \ldots, s_N$

- <u>Clustering</u>: Group data into $k$ clusters of nearby points

- Application?

  — E.g., Pre-processing step in Feature extraction

  — (Recollect Bag-of-<u>Visual</u>-Words: Dictionary creation using <u>clustering</u>)

- <u>k-Means clustering</u>: A low-complexity clustering algorithm

- <u>Assumption</u>: Data $s_n$ lives in a Metric space $(S, d)$

  — Example: Euclidean space $\mathbb{R}^p$ with $d(\boldsymbol{s}_i, \boldsymbol{s}_j) = \|\boldsymbol{s}_i - \boldsymbol{s}_j\|$

$k = 2$

**Carnegie Mellon University**

# k-Means clustering

- k-Means clustering:   Choose number of clusters $k$   (say 2)

1. Initialize:   Randomly select  $k$ of the $s_n$  as cluster-heads $c$ (set $\mathcal{C}$)

2. Iterate:      Run these two steps

   a)  Cluster:   Assign each data point $s_n$  to  <u>closest</u> cluster head $c$

   b)  New cluster-heads:  Each cluster  $\rightarrow$  Choose point <u>closest</u> to its points as  $c_i$

3. Terminate:   On convergence (cluster-heads don't change)



$s_n$

So, we have new cluster-heads $c_i$

b)

a)

So, we have $k$ clusters

# k-Means clustering

- <u>k-Means clustering</u>:   Choose number of clusters $k$   (say 2)

1. Initialize:   Randomly select  $k$ of the $s_n$ as cluster-heads $c$ (set $\mathcal{C}$)

2. Iterate:     Run these two steps

   a)  Cluster:     Assign each data point $s_n$ to <u>closest</u> cluster head $c$

   $$\widehat{c_n} \;=\; \underset{c\,\in\,\mathcal{C}}{\text{argmin}}\ \ d(s_n, c)$$

   b)  New cluster-heads:   Each cluster  $\rightarrow$  Choose point <u>closest</u> to its points as  $c_i$

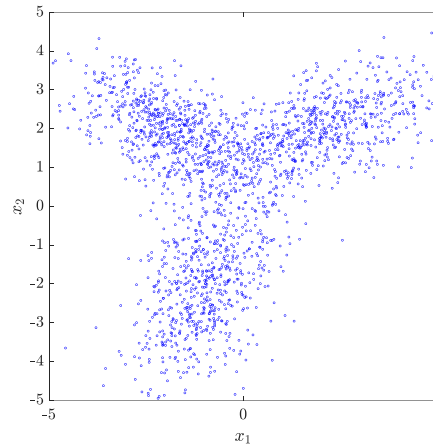   $$c_i \;=\; \underset{s\,\in\,S}{\text{argmin}} \sum_{\substack{s_n: \\ s_n \text{ is in Cluster } i}} d(s_n, s)$$
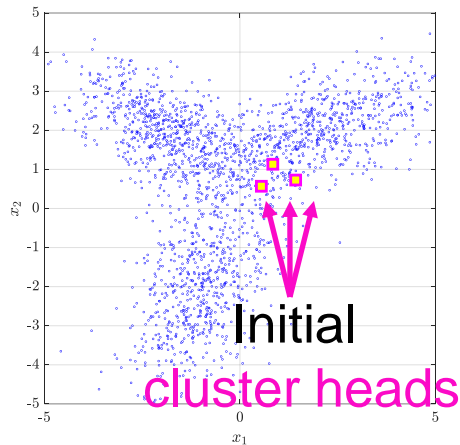
   In Euclidean space

   $$c_i \;=\; \frac{1}{n_i}\Big(s_{k_1} + \cdots + s_{k_{n_i}}\Big)$$

3. Terminate:   On convergence (cluster-heads don't change)
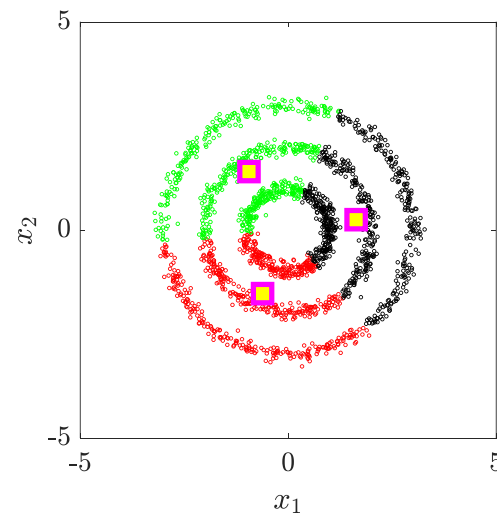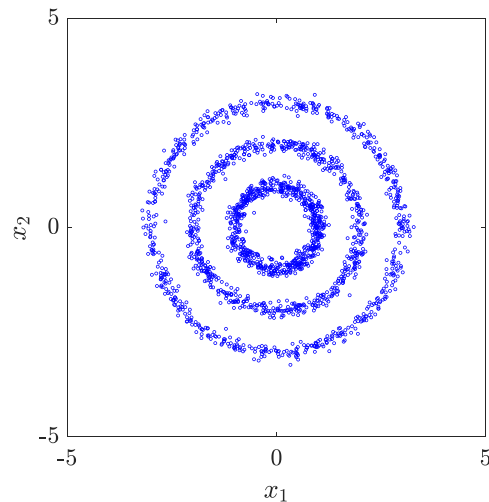
- Choose Euclidean distance $||\boldsymbol{s}_1 - \boldsymbol{s}_2||_2$ and cluster with $k = 3$



Initial cluster heads

$k = 3$ clusters obtained

Converged!

a) $k$ clusters

b) New cluster heads

$k = 3$ cluster-heads obtained $\boldsymbol{c}_1, \boldsymbol{c}_2, \boldsymbol{c}_3$

# Example: k-Means clustering
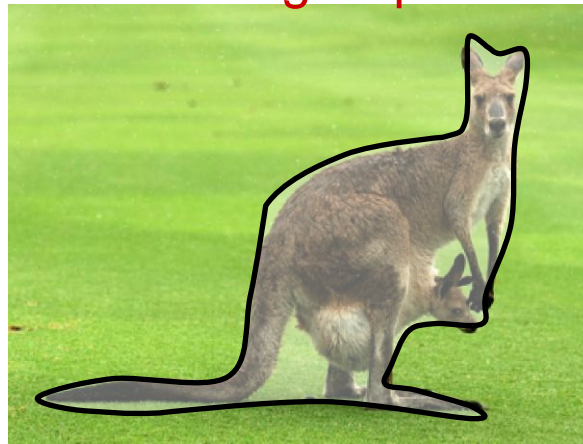
- Metric space choice is important to get "good" clustering

- E.g., Clustering below is "bad" if Metric space is Euclidean space $\mathbb{R}^2$

 — (Need metric space that better captures idea of 'distance' here)

 — (e.g., 'Graph metric space' perhaps)

# Clustering for Feature extraction

- Clustering algorithm → $k$ clusters and cluster heads $c_i$. Use?

1. Dictionary creation: Bag-of-Visual-Words for images (recollect)

2. Segmentation: Partition image into 'objects' (Pre-processing)

   — Each pixel is vector $x_n$ = $\begin{bmatrix} Red\ value \\ Green\ value \\ Blue\ value \end{bmatrix}$

   — Cluster the pixels (vectors $x_n$)

   — Each cluster gives 'mask' → Object to be investigated/classified

Cluster image's pixels          $k$ = 3 clusters