

# CSCI567 Machine Learning (Spring 2021)

Sirisha Rambhatla

University of Southern California

March 31, 2021

1 / 26

Review of last lecture

## Outline

1 Review of last lecture

2 Principal Component Analysis (PCA)

3 / 26

## Outline

1 Review of last lecture

2 Principal Component Analysis (PCA)

2 / 26

Review of last lecture

## Bayes optimal classifier

Suppose  $(\mathbf{x}, y)$  is drawn from a joint distribution  $p$ . The **Bayes optimal classifier** is

$$f^*(\mathbf{x}) = \operatorname{argmax}_{c \in [C]} p(c | \mathbf{x})$$

i.e. **predict the class with the largest conditional probability.**

$p$  is of course unknown, but we can estimate it, which is *exactly a density estimation problem!*

4 / 26

## A “naive” assumption

Naive Bayes assumption:

conditioning on a label, features are independent, which means

$$p(\mathbf{x} \mid y = c) = \prod_{d=1}^D p(x_d \mid y = c)$$

Now for each  $d$  and  $c$  we have a simple **1D density estimation problem!**

Is this a reasonable assumption? Sometimes yes, e.g.

- use  $x = (\text{Height, Vocabulary})$  to predict  $y = \text{Age}$
- Height and Vocabulary are dependent
- but **condition on Age, they are independent!**

More often this assumption is *unrealistic and “naive”*, but still Naive Bayes can work very well even if the assumption is wrong.

## Outline

- 1 Review of last lecture
- 2 Principal Component Analysis (PCA)
  - PCA
  - Kernel PCA

## Dimensionality reduction

**Dimensionality reduction** is yet another important unsupervised learning problem.

**Goal:** reduce the dimensionality of a dataset so

- it is **easier to visualize and discover patterns**
- it **takes less time and space** to process for any applications (classification, regression, clustering, etc)
- **noise is reduced**
- ...

There are many approaches, we focus on a linear method:

**Principal Component Analysis (PCA)**

## Example

[picture from here](#)

Consider the following dataset:

- **17 features**, each represents the **average consumption of some food**
- **4 data points**, each represents some **country**

Alcoholic drinks	375	135	458	475
Beverages	57	47	53	73
Carcase meat	245	267	242	227
Cereals	1472	1494	1462	1582
Cheese	105	66	103	103
Confectionery	54	41	62	64
Fats and oils	193	209	184	235
Fish	147	93	122	160
Fresh fruit	1102	674	957	1137
Fresh potatoes	720	1033	566	874
Fresh Veg	253	143	171	265
Other meat	685	586	750	803
Other Veg	488	355	418	570
Processed potatoes	198	187	220	203
Processed Veg	360	334	337	365
Soft drinks	1374	1508	1572	1256
Sugars	156	139	147	175

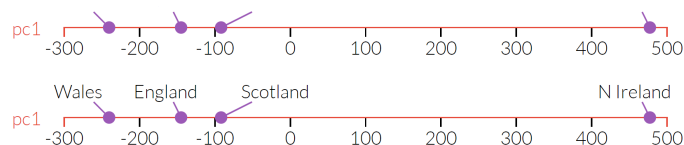
*What can you tell?*

Hard to say anything looking at all these 17 features.

## Example

picture from here

PCA can help us! Plot along the **first principal component** of this dataset:



i.e. we **reduce the dimensionality from 17 to just 1**.

Now one data point is clearly different from the rest!

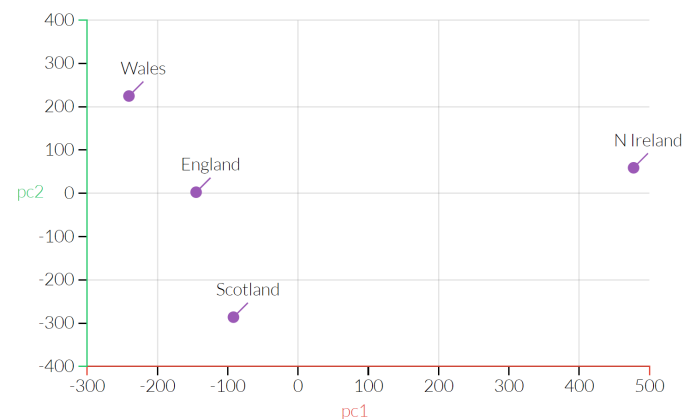
That turns out to be data from **Northern Ireland**, *the only country not on the island of Great Britain out of the 4 samples*.

9 / 26

## Example

picture from here

PCA can find the **second (and more) principal component** of the data too:



10 / 26

## High level idea

*How does PCA find these principal components (PC)?*



The first PC is in fact **the direction with the most variance**, i.e. the direction where the data is most spread out.

11 / 26

## Finding the first PC

More formally, we want to find a direction  $\mathbf{v} \in \mathbb{R}^D$  with  $\|\mathbf{v}\|_2 = 1$ , so that the **projection of the dataset on this direction has the most variance**, i.e.

$$\max_{\mathbf{v}: \|\mathbf{v}\|_2=1} \sum_{n=1}^N \left( \mathbf{x}_n^T \mathbf{v} - \frac{1}{N} \sum_m \mathbf{x}_m^T \mathbf{v} \right)^2$$

- $\mathbf{x}_n^T \mathbf{v}$  is exactly the **projection of  $\mathbf{x}_n$  onto the direction  $\mathbf{v}$**
- if we **pre-center the data**, i.e. let  $\mathbf{x}'_n = \mathbf{x}_n - \frac{1}{N} \sum_m \mathbf{x}_m$ , then the objective simply becomes

$$\max_{\mathbf{v}: \|\mathbf{v}\|_2=1} \sum_{n=1}^N \left( \mathbf{x}'_n^T \mathbf{v} \right)^2 = \max_{\mathbf{v}: \|\mathbf{v}\|_2=1} \mathbf{v}^T \left( \sum_{n=1}^N \mathbf{x}'_n \mathbf{x}'_n^T \right) \mathbf{v}$$

- we will simply assume  $\{\mathbf{x}_n\}$  is centered (to avoid notation  $\mathbf{x}'_n$ )

12 / 26

## Finding the first PC

With  $\mathbf{X} \in \mathbb{R}^{N \times D}$  being the data matrix, we want

$$\max_{\mathbf{v}: \|\mathbf{v}\|_2=1} \mathbf{v}^T (\mathbf{X}^T \mathbf{X}) \mathbf{v}$$

The Lagrangian is

$$\mathbf{v}^T (\mathbf{X}^T \mathbf{X}) \mathbf{v} - \lambda (\|\mathbf{v}\|_2^2 - 1)$$

The stationary condition implies  $\mathbf{X}^T \mathbf{X} \mathbf{v} = \lambda \mathbf{v}$ , which means  $\mathbf{v}$  is *exactly an eigenvector!* And the objective becomes

$$\mathbf{v}^T (\mathbf{X}^T \mathbf{X}) \mathbf{v} = \lambda \mathbf{v}^T \mathbf{v} = \lambda$$

To maximize this, we want the **eigenvector with the largest eigenvalue**

**Conclusion:** the first PC is the top eigenvector of the covariance matrix

## Finding the other PCs

If  $\mathbf{v}_1$  is the first PC, then the **second PC** is found via

$$\max_{\mathbf{v}_2: \|\mathbf{v}_2\|_2=1, \mathbf{v}_1^T \mathbf{v}_2=0} \mathbf{v}_2^T (\mathbf{X}^T \mathbf{X}) \mathbf{v}_2$$

i.e. **the direction that maximizes the variance among all other dimensions**

This is just the **second top eigenvector of the covariance matrix!**

**Conclusion:** the  $d$ -th principal component is the  $d$ -th eigenvector (sorted by the eigenvalue from largest to smallest).

## PCA

**Input:** a dataset represented as  $\mathbf{X}$ , #components  $p$

**Step 1** **Center the data** by subtracting the mean

**Step 2** **Find the top  $p$  eigenvectors (with unit norm) of the covariance matrix**  $\mathbf{X}^T \mathbf{X}$ , denote it by  $\mathbf{V} \in \mathbb{R}^{D \times p}$

**Step 3** Construct the new compressed dataset  $\mathbf{XV} \in \mathbb{R}^{N \times p}$

## How many PCs do we want?

One common rule: pick  $p$  large enough so it **covers about 90% of the spectrum**, i.e.

$$\frac{\sum_{d=1}^p \lambda_d}{\sum_{d=1}^D \lambda_d} \geq 90\%$$

where  $\lambda_1 \geq \dots \geq \lambda_N$  are sorted eigenvalues.

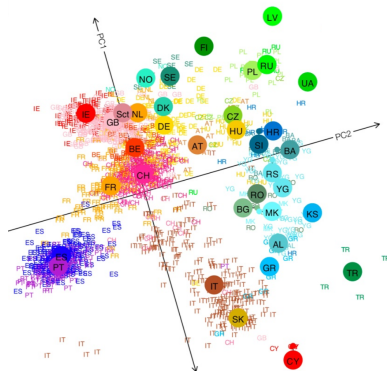
Note:  $\sum_{d=1}^D \lambda_d = \text{Tr}(\mathbf{X}^T \mathbf{X})$ , so no need to actually find all eigenvalues.

For **visualization**, also often pick  $p = 1$  or  $p = 2$ .

## Another visualization example

A famous study of **genetic map**

- dataset: **genomes of 1,387 Europeans**
- First 2 PCs shown below; *looks remarkably like the geographic map*

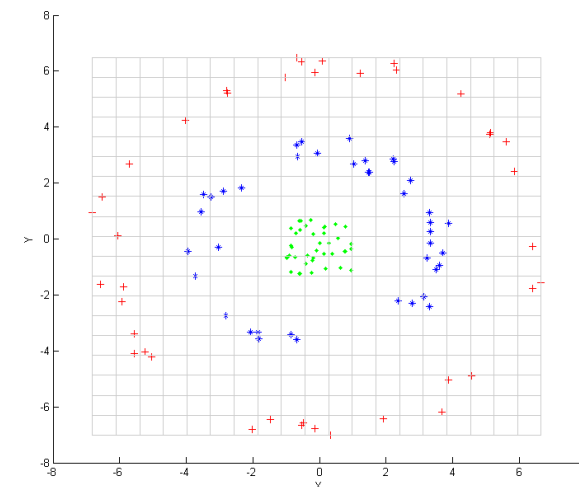


17 / 26

## Does PCA always work?

picture from Wikipedia

PCA is a **linear method** (recall the new dataset is  $XV$ ), it does not do much when **every direction has similar variance**.



18 / 26

## KPCA: high level idea

Similar to learning a linear classifier, when we encounter such data, *we can apply kernel methods*.

### Kernel PCA (KPCA):

- first map the data to a more complicated space via  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M$
- then apply regular PCA to reduce the dimensionality

Sounds a bit counter-intuitive, but the key is this gives a **nonlinear method**.

*How to implement KPCA efficiently without actually working in  $\mathbb{R}^M$ ?*

19 / 26

## KPCA: finding the PCs

Suppose  $v \in \mathbb{R}^M$  is the first PC for the nonlinearly-transformed data  $\Phi \in \mathbb{R}^{N \times M}$  (centered). Then let

$$v = \frac{1}{\lambda} \Phi^T \Phi v = \Phi^T \alpha$$

for some  $\alpha \in \mathbb{R}^N$ , i.e. it's a **linear combination of data**.

Plugging into  $\Phi^T \Phi v = \lambda v$  gives

$$\Phi^T \Phi \Phi^T \alpha = \lambda \Phi^T \alpha$$

and thus with the Gram matrix  $K = \Phi \Phi^T$ ,

$$\Phi^T (K \alpha - \lambda \alpha) = 0.$$

*So  $\alpha$  is an eigenvector of  $K$ !*

**Conclusion:** KPCA is just finding top eigenvectors of the Gram matrix

20 / 26

## One issue: scaling

Should we scale  $\alpha$  s.t  $\|\alpha\|_2 = 1$ ?

**No.** Recall we want  $v = \Phi^T \alpha$  to have unit L2 norm, so

$$v^T v = \alpha^T \Phi \Phi^T \alpha = \lambda \|\alpha\|_2^2 = 1$$

In other words, we in fact need to scale  $\alpha$  so that its L2 norm is  $1/\sqrt{\lambda}$ , where  $\lambda$  it's the corresponding eigenvalue.

21 / 26

## Another issue: centering

Should we still pre-center  $X$ ?

**No.** Centering  $X$  does not mean  $\Phi$  is centered!

Remember all we need is Gram matrix. *What is the Gram matrix after  $\Phi$  is centered?*

Let  $E \in \mathbb{R}^{N \times N}$  be the matrix with all entries being  $\frac{1}{N}$ ,

$$\begin{aligned} \bar{K} &= (\Phi - E\Phi)(\Phi - E\Phi)^T \\ &= \Phi\Phi^T - E\Phi\Phi^T - \Phi\Phi^T E + E\Phi\Phi^T E \\ &= K - EK - KE + EKE \end{aligned}$$

22 / 26

## KPCA

**Input:** a dataset  $X$ , #components  $p$ , a **Kernel function**  $k$

**Step 1** Compute the Gram matrix  $K$  and the **centered Gram matrix**

$$\bar{K} = K - EK - KE + EKE$$

**Step 2** Find the top  $p$  eigenvectors of  $\bar{K}$  with the appropriate scaling, denote it by  $A \in \mathbb{R}^{N \times p}$

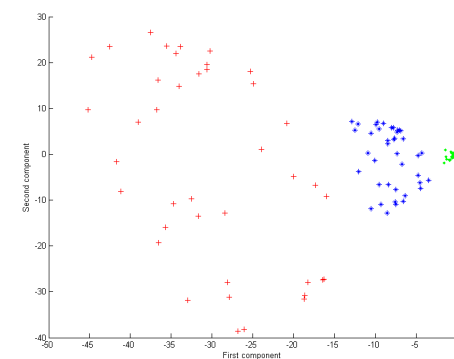
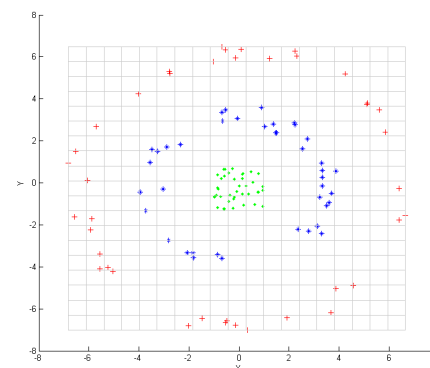
**Step 3** Construct the new dataset  $(\Phi - E\Phi)(\Phi - E\Phi)^T A = \bar{K} A$

23 / 26

## Example

picture from Wikipedia

Applying kernel  $k(x, x') = (x^T x' + 1)^2$ :

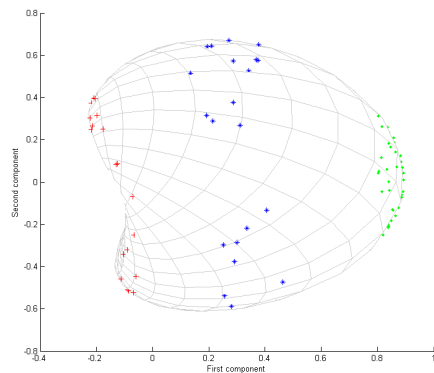
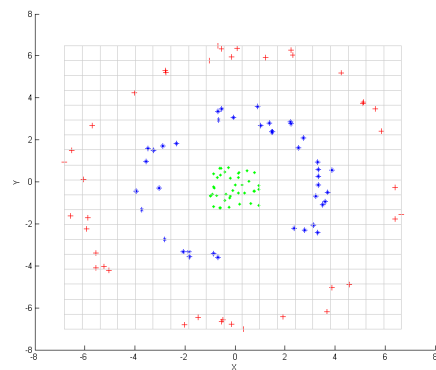


24 / 26

## Example

picture from Wikipedia

Applying Gaussian kernel  $k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}\right)$ :



## Denoising via PCA

Original data



Data corrupted with Gaussian noise



Result after linear PCA



Result after kernel PCA, Gaussian kernel

