

CSCI567 Machine Learning (Fall 2018)

Prof. Haipeng Luo

U of Southern California

Oct 31, 2018

November 14, 2018 1 / 60

Outline

- 1 Review of last lecture
- 2 (Hidden) Markov models
- 3 Principal Component Analysis (PCA)

November 14, 2018 3 / 60

Administration

HW3 solution is available, HW4 is due on Sunday (11/4)

Minor typo in P4 for the formula of multivariate Gaussian density, see Piazza pinned post as well as the updated P4.pdf.

- the comment in gmm.py:

$$p = e^{(-0.5(x-\text{mean}) * (\text{inv}(\text{variance})) * (x-\text{mean})' / \text{sqrt}(c))}$$

should be

$$p = e^{(-0.5(x-\text{mean}) * (\text{inv}(\text{variance})) * (x-\text{mean})') / \text{sqrt}(c)}$$

November 14, 2018 2 / 60

Review of last lecture

Outline

- 1 Review of last lecture
- 2 (Hidden) Markov models
- 3 Principal Component Analysis (PCA)

November 14, 2018 4 / 60

General EM algorithm

Step 0 Initialize $\theta^{(1)}$, $t = 1$

Step 1 (E-Step) update the posterior of latent variables

$$q_n^{(t)}(\cdot) = p(\cdot \mid \mathbf{x}_n; \theta^{(t)})$$

and obtain **Expectation** of complete likelihood

$$Q(\theta; \theta^{(t)}) = \sum_{n=1}^N \mathbb{E}_{z_n \sim q_n^{(t)}} [\ln p(\mathbf{x}_n, z_n; \theta)]$$

Step 2 (M-Step) update the model parameter via **Maximization**

$$\theta^{(t+1)} \leftarrow \underset{\theta}{\operatorname{argmax}} Q(\theta; \theta^{(t)})$$

Step 3 $t \leftarrow t + 1$ and return to Step 1 if not converged

Outline

1 Review of last lecture

2 (Hidden) Markov models

- Markov chain
- Hidden Markov Model
- Inferring HMMs
- Learning HMMs

3 Principal Component Analysis (PCA)

Applying EM to learn GMMs

EM for clustering:

Step 0 Initialize $\omega_k, \mu_k, \Sigma_k$ for each $k \in [K]$

Step 1 (E-Step) update the “soft assignment” (fixing parameters)

$$\gamma_{nk} = p(z_n = k \mid \mathbf{x}_n) \propto \omega_k N(\mathbf{x}_n \mid \mu_k, \Sigma_k)$$

Step 2 (M-Step) update the model parameter (fixing assignments)

$$\omega_k = \frac{\sum_n \gamma_{nk}}{N} \quad \mu_k = \frac{\sum_n \gamma_{nk} \mathbf{x}_n}{\sum_n \gamma_{nk}}$$

$$\Sigma_k = \frac{1}{\sum_n \gamma_{nk}} \sum_n \gamma_{nk} (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T$$

Step 3 return to Step 1 if not converged

Markov Models

Markov models are powerful probabilistic tools to analyze **sequential data**:

- text or speech data
- stock market data
- gene data
- ...

Definition

A **Markov chain** is a stochastic process with **Markov property**: a sequence of random variables Z_1, Z_2, \dots s.t.

$$P(Z_{t+1} \mid Z_{1:t}) = P(Z_{t+1} \mid Z_t) \quad (\text{Markov property})$$

i.e. *the current state only depends on the most recent state* (notation $Z_{1:t}$ denotes the sequence Z_1, \dots, Z_t).

We only consider the following case:

- All Z_t 's take value from the same **discrete** set $\{1, \dots, S\}$
- $P(Z_{t+1} = s' \mid Z_t = s) = a_{s,s'}$, known as **transition probability**
- $P(Z_1 = s) = \pi_s$
- $(\{\pi_s\}, \{a_{s,s'}\}) = (\boldsymbol{\pi}, \mathbf{A})$ are **parameters of the model**

High-order Markov chain

Is the Markov assumption reasonable? Not completely for the language model for example.

Higher order Markov chains make it more reasonable, e.g.

$$P(Z_{t+1} \mid Z_{1:t}) = P(Z_{t+1} \mid Z_t, Z_{t-1}) \quad (\text{second-order Markov})$$

i.e. the current word only depends on the last two words.

Learning higher order Markov chains is similar, but more expensive.

We only consider standard Markov chains.

Examples

- Example 1 (**Language model**)

States $[S]$ represent a dictionary of words,

$$a_{\text{ice,cream}} = P(Z_{t+1} = \text{cream} \mid Z_t = \text{ice})$$

is an example of the transition probability.

- Example 2 (**Weather**)

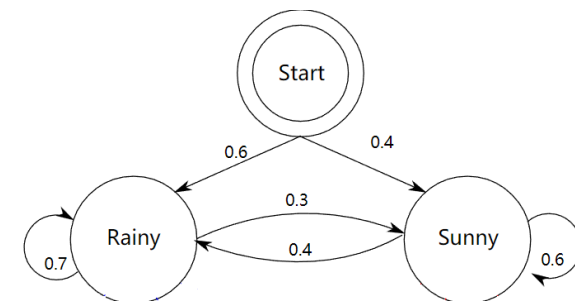
States $[S]$ represent weather at each day

$$a_{\text{sunny,rainy}} = P(Z_{t+1} = \text{rainy} \mid Z_t = \text{sunny})$$

Graph Representation

picture from Wikipedia

It is intuitive to represent a Markov model as a **graph**



Learning from examples

Now suppose we have observed N sequences of examples:

- $z_{1,1}, \dots, z_{1,T}$
- \dots
- $z_{n,1}, \dots, z_{n,T}$
- \dots
- $z_{N,1}, \dots, z_{N,T}$

where

- for simplicity we assume each sequence has the same length T
- lower case $z_{n,t}$ represents the value of the random variable $Z_{n,t}$

From these observations how do we *learn the model parameters* (π, A)?

Finding the MLE

Same story, find the **MLE**. The log-likelihood of a sequence z_1, \dots, z_T is

$$\begin{aligned}
 \ln P(Z_{1:T} = z_{1:T}) &= \sum_{t=1}^T \ln P(Z_t = z_t \mid Z_{1:t-1} = z_{1:t-1}) && \text{(always true)} \\
 &= \sum_{t=1}^T \ln P(Z_t = z_t \mid Z_{t-1} = z_{t-1}) && \text{(Markov property)} \\
 &= \ln \pi_{z_1} + \sum_{t=2}^T \ln a_{z_{t-1}, z_t} \\
 &= \sum_s \mathbb{I}[z_1 = s] \ln \pi_s + \sum_{s,s'} \left(\sum_{t=2}^T \mathbb{I}[z_{t-1} = s, z_t = s'] \right) \ln a_{s,s'}
 \end{aligned}$$

Finding the MLE

So MLE is

$$\begin{aligned}
 \operatorname{argmax}_{\pi, A} \sum_s (\text{\#initial states with value } s) \ln \pi_s \\
 + \sum_{s,s'} (\text{\#transitions from } s \text{ to } s') \ln a_{s,s'}
 \end{aligned}$$

We have seen this many times. The solution is:

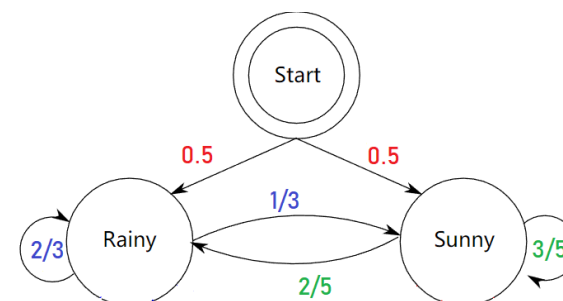
$$\begin{aligned}
 \pi_s &\propto \text{\#initial states with value } s \\
 a_{s,s'} &\propto \text{\#transitions from } s \text{ to } s'
 \end{aligned}$$

Example

Suppose we observed the following 2 sequences of length 5

- sunny, sunny, rainy, rainy, rainy
- rainy, sunny, sunny, sunny, rainy

MLE is the following model



Markov Model with outcomes

Now suppose each state Z_t also “emits” some **outcome** $X_t \in [O]$ based on the following model

$$P(X_t = o \mid Z_t = s) = b_{s,o} \quad (\text{emission probability})$$

independent of anything else.

For example, in the language model, X_t is the speech signal for the underlying word Z_t (very useful for **speech recognition**).

Now the model parameters are $(\{\pi_s\}, \{a_{s,s'}\}, \{b_{s,o}\}) = (\pi, A, B)$.

Joint likelihood

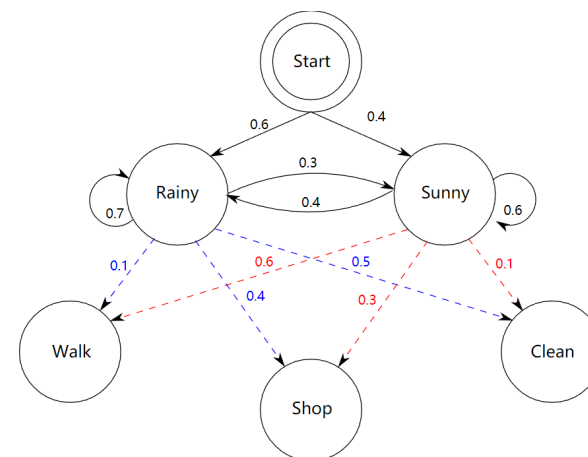
The joint log-likelihood of a **state-outcome sequence** $z_1, x_1, \dots, z_T, x_T$ is

$$\begin{aligned} \ln P(Z_{1:T} = z_{1:T}, X_{1:T} = x_{1:T}) \\ &= \ln P(Z_{1:T} = z_{1:T}) + \ln P(X_{1:T} = x_{1:T} \mid Z_{1:T} = z_{1:T}) \quad (\text{always true}) \\ &= \sum_{t=1}^T \ln P(Z_t = z_t \mid Z_{t-1} = z_{t-1}) + \sum_{t=1}^T \ln P(X_t = x_t \mid Z_t = z_t) \\ &\quad (\text{due to all the independence}) \\ &= \ln \pi_{z_1} + \sum_{t=2}^T \ln a_{z_{t-1}, z_t} + \sum_{t=1}^T \ln b_{z_t, x_t} \end{aligned}$$

Another example

picture from Wikipedia

On each day, we also observe **Bob's activity: walk, shop, or clean**, which only depends on the weather of that day.



Learning the model

If we observe N state-outcome sequences: $z_{n,1}, x_{n,1}, \dots, z_{n,T}, x_{n,T}$ for $n = 1, \dots, N$, the MLE is again very simple (verify yourself):

$$\begin{aligned} \pi_s &\propto \text{\#initial states with value } s \\ a_{s,s'} &\propto \text{\#transitions from } s \text{ to } s' \\ b_{s,o} &\propto \text{\#state-outcome pairs } (s, o) \end{aligned}$$

Learning the model

However, *most often we do not observe the states!* Think about the speech recognition example.

This is called **Hidden Markov Model (HMM)**, widely used in practice

How to learn HMMs? **Roadmap:**

- first discuss how to **infer** when the model is known (key: **dynamic programming**)
- then discuss how to **learn** the model (key: **EM**)

What can we infer for a known HMM?

Knowing the parameter of an HMM, we can infer

- **the transition at some point, given an observation sequence**

$$P(Z_t = s, Z_{t+1} = s' \mid X_{1:T} = x_{1:T})$$

e.g. given Bob's activities for one week, how was the weather like on Wed and Thu?

- **most likely hidden states path, given an observation sequence**

$$\operatorname{argmax}_{z_{1:T}} P(Z_{1:T} = z_{1:T} \mid X_{1:T} = x_{1:T})$$

e.g. given Bob's activities for one week, what's the most likely weather for this week?

What can we infer about an HMM?

Knowing the parameter of an HMM, we can infer

- **the probability of observing some sequence**

$$P(X_{1:T} = x_{1:T})$$

e.g. prob. of observing Bob's activities "walk, walk, shop, clean, walk, shop, shop" for one week

- **the state at some point, given an observation sequence**

$$P(Z_t = s \mid X_{1:T} = x_{1:T})$$

e.g. given Bob's activities for one week, how was the weather like on Wed?

Forward and backward messages

The key to infer all these is to compute two things:

- **forward messages:** for each s and t

$$\alpha_s(t) = P(Z_t = s, X_{1:t} = x_{1:t})$$

- **backward messages:** for each s and t

$$\beta_s(t) = P(X_{t+1:T} = x_{t+1:T} \mid Z_t = s)$$

Computing forward messages

Key: *establish a recursive formula*

$$\begin{aligned}
 \alpha_s(t) &= P(Z_t = s, X_{1:t} = x_{1:t}) \\
 &= P(X_t = x_t \mid Z_t = s, X_{1:t-1} = x_{1:t-1}) P(Z_t = s, X_{1:t-1} = x_{1:t-1}) \\
 &= b_{s,x_t} \sum_{s'} P(Z_t = s, Z_{t-1} = s', X_{1:t-1} = x_{1:t-1}) \quad (\text{marginalizing}) \\
 &= b_{s,x_t} \sum_{s'} P(Z_t = s \mid Z_{t-1} = s', X_{1:t-1} = x_{1:t-1}) P(Z_{t-1} = s', X_{1:t-1} = x_{1:t-1}) \\
 &= b_{s,x_t} \sum_{s'} a_{s',s} \alpha_{s'}(t-1) \quad (\text{recursive form!})
 \end{aligned}$$

Base case: $\alpha_s(1) = P(Z_1 = s, X_1 = x_1) = \pi_s b_{s,x_1}$

Computing backward messages

Again establish a recursive formula

$$\begin{aligned}
 \beta_s(t) &= P(X_{t+1:T} = x_{t+1:T} \mid Z_t = s) \\
 &= \sum_{s'} P(X_{t+1:T} = x_{t+1:T}, Z_{t+1} = s' \mid Z_t = s) \quad (\text{marginalizing}) \\
 &= \sum_{s'} P(Z_{t+1} = s' \mid Z_t = s) P(X_{t+1:T} = x_{t+1:T} \mid Z_{t+1} = s', Z_t = s) \\
 &= \sum_{s'} a_{s,s'} P(X_{t+1} = x_{t+1} \mid Z_{t+1} = s') P(X_{t+2:T} = x_{t+2:T} \mid Z_{t+1} = s') \\
 &= \sum_{s'} a_{s,s'} b_{s',x_{t+1}} \beta_{s'}(t+1) \quad (\text{recursive form!})
 \end{aligned}$$

Base case: $\beta_s(T) = 1$

Forward procedure

Forward procedure

For all $s \in [S]$, compute $\alpha_s(1) = \pi_s b_{s,x_1}$.

For $t = 2, \dots, T$

- for each $s \in [S]$, compute

$$\alpha_s(t) = b_{s,x_t} \sum_{s'} a_{s',s} \alpha_{s'}(t-1)$$

It takes $O(S^2T)$ time and $O(ST)$ space.

Backward procedure

Backward procedure

For all $s \in [S]$, set $\beta_s(T) = 1$.

For $t = T-1, \dots, 1$

- for each $s \in [S]$, compute

$$\beta_s(t) = \sum_{s'} a_{s,s'} b_{s',x_{t+1}} \beta_{s'}(t+1)$$

Again it takes $O(S^2T)$ time and $O(ST)$ space.

Using forward and backward messages

With forward and backward messages, we can easily infer many things, e.g.

$$\begin{aligned}\gamma_s(t) &= P(Z_t = s \mid X_{1:T} = x_{1:T}) \\ &\propto P(Z_t = s, X_{1:T} = x_{1:T}) \\ &= P(Z_t = s, X_{1:t} = x_{1:t})P(X_{t+1:T} = x_{t+1:T} \mid Z_t = s, X_{1:t} = x_{1:t}) \\ &= \alpha_s(t)\beta_s(t)\end{aligned}$$

What constant are we omitting in “ \propto ”? It is exactly

$$P(X_{1:T} = x_{1:T}) = \sum_s \alpha_s(t)\beta_s(t),$$

the probability of observing the sequence $x_{1:T}$.

This is true for any t ; a good way to check correctness of your code.

Finding the most likely path

Though can't use forward and backward messages directly to find the most likely path, it is **very similar to the forward procedure**. Key: compute

$$\delta_s(t) = \max_{z_{1:t-1}} P(Z_t = s, Z_{1:t-1} = z_{1:t-1}, X_{1:t} = x_{1:t})$$

the probability of the most likely path for time $1 : t$ ending at state s

Using forward and backward messages

Another example: the conditional probability of transition s to s' at time t

$$\begin{aligned}\xi_{s,s'}(t) &= P(Z_t = s, Z_{t+1} = s' \mid X_{1:T} = x_{1:T}) \\ &\propto P(Z_t = s, Z_{t+1} = s', X_{1:T} = x_{1:T}) \\ &= P(Z_t = s, X_{1:t} = x_{1:t})P(Z_{t+1} = s', X_{t+1:T} = x_{t+1:T} \mid Z_t = s, X_{1:t} = x_{1:t}) \\ &= \alpha_s(t)P(Z_{t+1} = s' \mid Z_t = s)P(X_{t+1:T} = x_{t+1:T} \mid Z_{t+1} = s') \\ &= \alpha_s(t)a_{s,s'}P(X_{t+1} = x_{t+1} \mid Z_{t+1} = s')P(X_{t+2:T} = x_{t+2:T} \mid Z_{t+1} = s') \\ &= \alpha_s(t)a_{s,s'}b_{s',x_{t+1}}\beta_{s'}(t+1)\end{aligned}$$

The **normalization constant** is in fact again $P(X_{1:T} = x_{1:T})$

Computing $\delta_s(t)$

Observe

$$\begin{aligned}\delta_s(t) &= \max_{z_{1:t-1}} P(Z_t = s, Z_{1:t-1} = z_{1:t-1}, X_{1:t} = x_{1:t}) \\ &= \max_{s'} \max_{z_{1:t-2}} P(Z_t = s, Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t} = x_{1:t}) \\ &= \max_{s'} P(Z_t = s \mid Z_{t-1} = s')P(X_t = x_t \mid Z_t = s) \cdot \\ &\quad \max_{z_{1:t-2}} P(Z_{t-1} = s', Z_{1:t-2} = z_{1:t-2}, X_{1:t-1} = x_{1:t-1}) \\ &= b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1) \quad (\text{recursive form!})\end{aligned}$$

Base case: $\delta_s(1) = P(Z_1 = s, X_1 = x_1) = \pi_s b_{s,x_1}$

Exactly the same as forward messages except replacing “sum” by “max”!

Viterbi Algorithm (!)

Viterbi Algorithm

For each $s \in [S]$, compute $\delta_s(1) = \pi_s b_{s,x_1}$.

For each $t = 2, \dots, T$,

- for each $s \in [S]$, compute

$$\delta_s(t) = b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1)$$

$$\Delta_s(t) = \operatorname{argmax}_{s'} a_{s',s} \delta_{s'}(t-1)$$

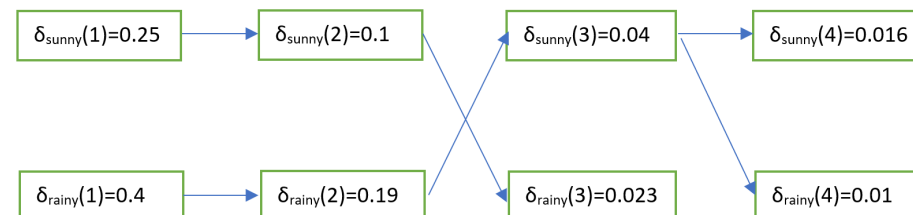
Backtracking: let $z_T^* = \operatorname{argmax}_s \delta_s(T)$.

For each $t = T, \dots, 2$: set $z_{t-1}^* = \Delta_{z_t^*}(t)$.

Output the most likely path z_1^*, \dots, z_T^* .

Example

Arrows represent the “argmax”, i.e. $\Delta_s(t)$.



The most likely path is “rainy, rainy, sunny, sunny”.

Learning the parameters of an HMM

All previous inferences depend on **knowing the parameters** $(\pi, \mathbf{A}, \mathbf{B})$.

How do we learn the parameters based on N observation sequences $x_{n,1}, \dots, x_{n,T}$ for $n = 1, \dots, N$?

MLE is **intractable due to the hidden variables** $Z_{n,t}$'s (similar to GMMs)

Need to apply **EM** again! Known as the **Baum–Welch algorithm**.

Applying EM: E-Step

Recall in the E-Step we fix the parameters and find the **posterior distributions q of the hidden states** (for each sample n), which leads to the complete log-likelihood:

$$\begin{aligned} & \mathbb{E}_{z_{1:T} \sim q} [\ln(Z_{1:T} = z_{1:T}, X_{1:T} = x_{1:T})] \\ &= \mathbb{E}_{z_{1:T} \sim q} \left[\ln \pi_{z_1} + \sum_{t=1}^{T-1} \ln a_{z_t, z_{t+1}} + \sum_{t=1}^T \ln b_{z_t, x_t} \right] \\ &= \sum_s \gamma_s(1) \ln \pi_s + \sum_{t=1}^{T-1} \sum_{s,s'} \xi_{s,s'}(t) \ln a_{s,s'} + \sum_{t=1}^T \sum_s \gamma_s(t) \ln b_{s,x_t} \end{aligned}$$

We have discussed how to compute

$$\begin{aligned} \gamma_s(t) &= P(Z_t = s \mid X_{1:T} = x_{1:T}) \\ \xi_{s,s'}(t) &= P(Z_t = s, Z_{t+1} = s' \mid X_{1:T} = x_{1:T}) \end{aligned}$$

Applying EM: M-Step

The maximizer of complete log-likelihood is simply doing **weighted counting** (compared to the unweighted counting on Slide 20):

$$\begin{aligned}\pi_s &\propto \sum_n \gamma_s^{(n)}(1) = \mathbb{E}_q [\text{\#initial states with value } s] \\ a_{s,s'} &\propto \sum_n \sum_{t=1}^{T-1} \xi_{s,s'}^{(n)}(t) = \mathbb{E}_q [\text{\#transitions from } s \text{ to } s'] \\ b_{s,o} &\propto \sum_n \sum_{t: x_t=o} \gamma_s^{(n)}(t) = \mathbb{E}_q [\text{\#state-outcome pairs } (s, o)]\end{aligned}$$

where

$$\begin{aligned}\gamma_s^{(n)}(t) &= P(Z_{n,t} = s \mid X_{n,1:T} = x_{n,1:T}) \\ \xi_{s,s'}^{(n)}(t) &= P(Z_{n,t} = s, Z_{n,t+1} = s' \mid X_{n,1:T} = x_{n,1:T})\end{aligned}$$

Summary

Very important models: **Markov chains**, **hidden Markov models**

Several algorithms:

- forward and backward procedures
- inferring HMMs based on forward and backward messages
- Viterbi algorithm
- Baum–Welch algorithm

Baum–Welch algorithm

Step 0 Initialize the parameters (π, A, B)

Step 1 (E-Step) Fixing the parameters, **compute forward and backward messages for all sample sequences**, then use these to compute $\gamma_s^{(n)}(t)$ and $\xi_{s,s'}^{(n)}(t)$ for each n, t, s, s' (see Slides 29 and 30).

Step 2 (M-Step) Update parameters:

$$\pi_s \propto \sum_n \gamma_s^{(n)}(1), \quad a_{s,s'} \propto \sum_n \sum_{t=1}^{T-1} \xi_{s,s'}^{(n)}(t), \quad b_{s,o} \propto \sum_n \sum_{t: x_t=o} \gamma_s^{(n)}(t)$$

Step 3 Return to Step 1 if not converged

Outline

- 1 Review of last lecture
- 2 (Hidden) Markov models
- 3 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)

November 14, 2018 41 / 60

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
Carcass 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	1504	1572	1256
Sugars	156	139	147	175

What can you tell?

Hard to say anything looking at all these 17 features.

November 14, 2018 42 / 60

Example

picture from here

PCA can help us! 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!

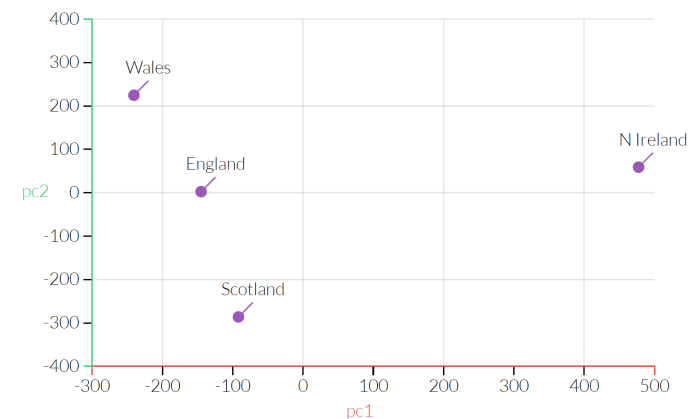
That turns out to be data from **Northern Ireland**, *the only country not on the island of Great Britain out of the 4 samples*.

November 14, 2018 43 / 60

Example

picture from here

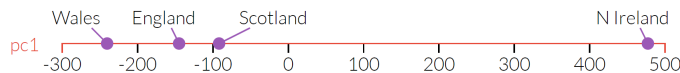
PCA can find the **second (and more) principal component** of the data too:



November 14, 2018 44 / 60

High level idea

How does PCA find these principal components (PC)?



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

Finding the first PC

With $\mathbf{X} \in \mathbb{R}^{N \times D}$ being the data matrix (as in Lec 2), 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 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)

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 we want

Step 1 Center the data by subtracting the mean

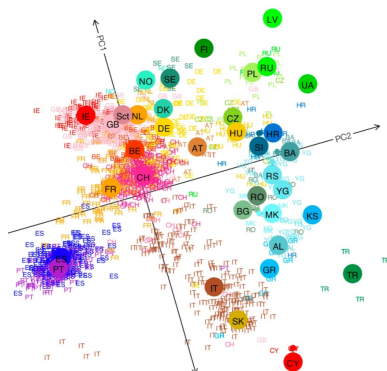
Step 2 Find the top p (unit norm) eigenvectors of the covariance matrix $\mathbf{X}^T \mathbf{X}$, denote it by $\mathbf{V} \in \mathbb{R}^{D \times p}$

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

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*



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.

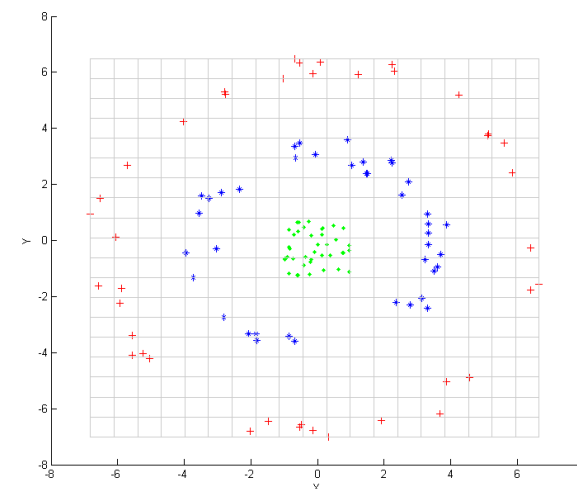
Note: $\sum_{d=1}^D \lambda_d = \text{Tr}(\mathbf{X}^T \mathbf{X})$, so no need to actually find all eigenvalues.

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

Does PCA always work?

picture from Wikipedia

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



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 ?

One issue: scaling

Should we scale α 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 α so that its L2 norm is $1/\sqrt{\lambda}$** , where λ it's the corresponding eigenvalue.

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

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

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 α is an eigenvector of K !

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

Another issue: centering

Should we still pre-center X ?

No. Centering X does not mean Φ is centered!

Remember all we need is Gram matrix. *What is the Gram matrix after Φ is centered?*

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

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

KPCA

Input: a dataset \mathbf{X} , #components p we want, a **Kernel function** k

Step 1 Compute the Gram matrix \mathbf{K} and the **centered Gram matrix**

$$\bar{\mathbf{K}} = \mathbf{K} - \mathbf{E}\mathbf{K} - \mathbf{K}\mathbf{E} + \mathbf{E}\mathbf{K}\mathbf{E}$$

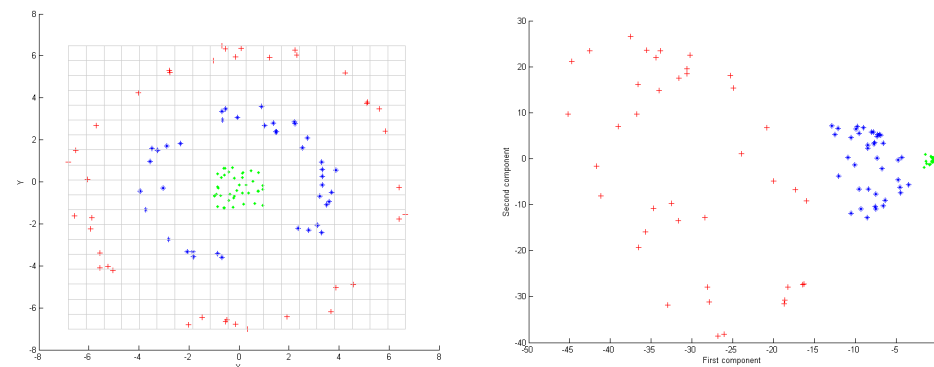
Step 2 Find the top p eigenvectors of $\bar{\mathbf{K}}$ with the appropriate scaling, denote it by $\mathbf{A} \in \mathbb{R}^{N \times p}$

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

Example

picture from Wikipedia

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



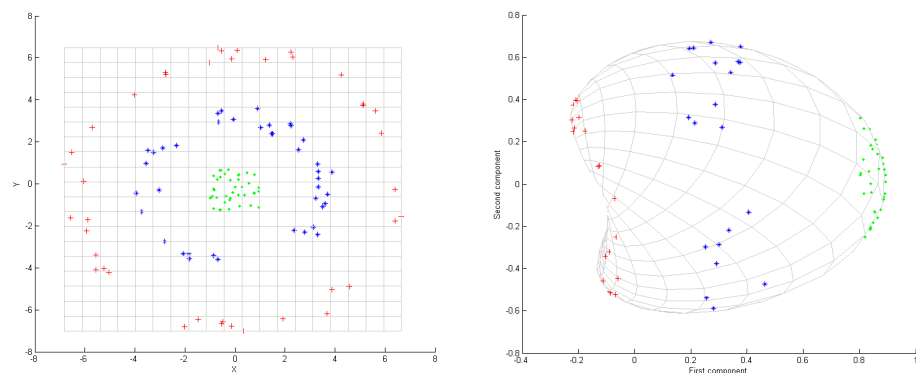
November 14, 2018 57 / 60

November 14, 2018 58 / 60

Example

picture from Wikipedia

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



November 14, 2018 59 / 60

Denoising via PCA

Original data



Data corrupted with Gaussian noise



Result after linear PCA



Result after kernel PCA, Gaussian kernel



November 14, 2018 60 / 60