### Module 1 :

Machine Learning Review

Unsupervised Learning Algorithms



Géraldine Conti, August 2020



## **Discussion Session**

- Review of Notebook 2 (regression)
  - Pipeline for processing numerical attributes
  - Train model
  - linear regression, decision tree, random forest and SVR
  - Compare the models (MSE, MAE, cross-validation)
  - Fine-tune model (grid search, random search
  - Assess the performance on the test set (RMSE, 95% CL)

# Bibliography

- Deep Learning book (Goodfellow, Bengio, Courville)
- Machine Learning @ Stanford (Prof Andrew Ng)
- Hands-On Machine Learning with Scikit-Learn & Tensorflow (Aurélien Géron)



## Learning Objectives







# Clustering

#### - K-means

- Hierarchical clustering



- Partition the dataset into K pre-defined non-overlapping clusters
  - "K" : the number of clusters
  - "means" : data is averaged (ie. finding centroids)
  - Each data point belongs to only one group



K-Means : Introduction

# • Expectation-Maximization (EM) algorithm : general method for finding ML-estimates

- 1) Initialization of the centroids
- 2) Iteration :

K-Means : Optimization 2) Maximization step : Compute the centroids for the clusters by taking the average of all the data points belonging to each cluster



1) Expectation step : Compute the sum of the squared distance between the data points and centroids and assign each data point to the closest cluster

- Stop criteria :
  - Centroids have stabilized wrt some tolerance
  - Number of iterations has been achieved





K-Means : Applications

#### **Customer segmentation**

**Goal** : you own a supermarket mall and you want to find who are your target customers

**Data** : membership card data (ID, age, gender, annual income, spending score)

#### **Document clustering**

**Goal** : you want to analyze Facebook/Twitter/Youtube comments of a particular event

Data : social media text data



Documents	Online	Festival	Book	Flight	Delhi
D1	1	0	1	0	1
D2	2	1	2	1	1
D3	0	0	1	1	1
D4	1	2	0	2	0
D5	3	1	0	0	0
D6	0	1	1	1	2
D7	2	0	1	2	1
D8	1	1	0	1	0
D9	1	0	2	0	0
D10	0	1	1	1	1

#### **Image compression**

**Goal** : to compress a digital image to reduce its cost of storage or transmission

Data : colored image (pixels)







#### sklearn.cluster.KMeans

K-Means :

In practice

- If large dataset (>10k samples) :
  - sklearn.cluster.MiniBatchKMeans

Build a hierarchy of clusters

- Two types :
  - Agglomerative
  - Divisive (not much used)

Hierarchical Clustering : Introduction



Hierarchical Clustering : Optimization

- Agglomerative algorithm: each data point is considered as an individual cluster
  - 1) Calculate the proximity of individual points
  - 2) Iteration :
    - Merge the two closest clusters
    - Update the proximity matrix
- Stop criterium :
  - a single cluster remains
- Divisive algorithm: the opposite



*How can the proximity be calculated ?* 



Hierarchical
Clustering :
Proximity
(Linkage)

	Single linkage algorithm (MIN)	Complete linkage algorithm (MAX)	Group Average *
Proximity	$\begin{array}{c c} & & & & \\ \hline \\ & & & \\ \hline & & \\ \hline & & & \\ \hline \\ \hline$		a
Pros	separate non-elliptical shapes if gap between two clusters is not small	Separates well clusters even if there is noise between clusters	Separates well clusters even if there is noise between clusters
Cons	cannot separate clusters properly if there is noise between clusters	Bias towards globular clusters Tends to break large clusters	Bias towards globular clusters

\* Ward's method is similar, but it calculates the sum of the squares of the distances  $\frac{17}{17}$ 





Hierarchical Clustering : Performance



 Space complexity O(n<sup>2</sup>)

• Time complexity O(n<sup>3</sup>)



hclust (\*, "ward.D")



Hierarchical Clustering: Applications

#### **Animal evolution**

**Goal** : find the phylogenetic tree relating different species together

**Data** : DNA sequences



### **Tracking viruses**

**Goal** : track viral outbreaks and their sources (SARS virus example)

Data : DNA sequences

#### **Find party lines through Twitter**

**Goal** : Cluster US senators into their respective parties

Data : Twitter data





#### Hierarchical Clustering : In Practice

### sklearn.cluster.AgglomerativeClustering

```
# 2D embedding of the digits dataset
print("Computing embedding")
X_red = manifold.SpectralEmbedding(n_components=2).fit_transform(X)
print("Done.")
```

#### from sklearn.cluster import AgglomerativeClustering

```
for linkage in ('ward', 'average', 'complete', 'single'):
    clustering = AgglomerativeClustering(linkage=linkage, n_clusters=10)
    t0 = time()
    clustering.fit(X_red)
    print("%s :\t%.2fs" % (linkage, time() - t0))
```



# **Density estimation**

*Observed data from a complex but unknown probability distribution* 

#### - Gaussian Mixture Model (GMM)





• mixing probability c that defines how big the Gaussian will be

• Mixture models are linear combinations of densities :

$$p(x \mid \Theta) = \sum_{i=1}^{K} c_i p(x \mid \Theta_i)$$
  
with  $\sum_{i=1}^{K} c_i = 1$ ,  $\int_{x} p(x \mid \Theta_i) dx = 1$ 

GMM: Formalism

• For Gaussian mixtures :

$$\Theta_i = \{\mu_i, \Sigma_i\}, \quad \Rightarrow \quad p(x \mid \Theta_i) \equiv N(\mu_i, \Sigma_i)$$

• To find the parameters, minimize the <u>-log-likelihood</u> :

 $L(\Theta) = \ln p(X; \Theta)$  $\Theta^* = \arg \max_{\Theta} L(\Theta)$ 

No analytical expression  $\rightarrow$  EM-algorithm

Gaussian Mixture : Performance

- Speed (fastest algorithm for mixture models)
- Agnostic (no biases of the means towards zero/of the cluster sizes to have specific structures that might (or not) apply)

- Singularities
- Number of components



Bayesian Information Criterion (BIC)



#### sklearn.mixture.GaussianMixture





Gaussian Mixture: Applications

#### **Demographic Statistics**

**Goal** : compute the average height of people of different ethnicities (African-American, Caucasian, Asian and Latino)

#### Data : People's heights

#### Finance

Goal : forecast asset returns

Data : timeseries of asset prices

#### **Detect objects in challenging environments**

Goal : detect coloured buoys underwater

Data : video sequence









## **Sequence Prediction**

- HMM

- RNN (module 2)

- LSTM (module 2)



• Invisible Markov chain that we cannot observe. Each state generates a random one out of k observations

Hidden Markov

Model

(HMM)

Number of states m to be decided

observations  $p(x_1, ..., x_m, s_1, ..., s_m; \theta) = \pi(s_1) \prod_{j=2}^m t(s_j | s_{j-1}) \prod_{j=1}^m e(x_j | s_j)$ 

InitialΠprobabilitiesμπ of thefstatest

Transition probabilities t from a state to another

Emission probabilities e of an observation given a state

### • 3 main tasks :

Task1 : Probability of an observation sequence (state

estimation)

Calculate the probability that we feel cold for two consecutive days

HMM Tasks

### Task 2 : Most likely explanation (inference/decoding)

find the most probable hidden states from a series of observations → Viterbi algorithm

In a sequence of 14 days, we know if it was hot or cold. From it, we want to find the most likely weather forecast (snow, rain, sunshine) for these two weeks

# Task 3 : Learning the HMM parameters → Baum-Welch algorithm

- Transition probability matrix
- Observation likelihood matrix



HMM : Applications

### Word recognition

**Goal** : recognize a word by looking at the characters separately

Data : images of characters



#### Stock trading

**Goal** : predict monthly closing prices of the S&P 500 index

Data : monthly prices



• No implementation in sklearn

<u>https://pystruct.github.io/</u>

HMM: In practice

<u>https://larsmans.github.io/seqlearn/</u>

• seqlearn.hmm.MultinomialHMM

<u>https://hmmlearn.readthedocs.io/en/latest/index.html</u>



## Feature Extraction

- PCA
- Kernel PCA
- Manifold Learning

The Curse of Dimensionality

- Many ML problems involve thousands of features for each training instance
  - Extremely slow, much harder to find a good solution
- Pick two points randomly, average distance between them :

Unit square	Unit 3D cube	Unit 1mio-dim hypercube
0.52	0.66	408.25

- High-dimensional datasets are at risk of being very sparse!
- Any new instance will likely be far away from any training instance, making predictions much less reliable

The Curse of Dimensionality: Solutions

- Increase the size of the training set
  - Required size grows exponentially with number of dimensions
- Reduce the number of features considerably
  - extremely useful for data visualization as well



Principal Component Analysis (PCA) : Introduction

- Find a new set of orthogonal dimensions and ranked according to the variance of data along them.
  - more important principal axis occurs first
  - the new dimensions allow to predict/reconstruct the original dimensions



Singular Value Decomposition (SVD)

- Eigen-decomposition :
  - 1) Calculate the covariance matrix *X* of data points.

$$C_x = \frac{1}{n-1} (X - \bar{X})(X - \bar{X})^T \qquad \qquad X^T = Transpose \ of \ X$$

2) Calculate eigenvectors and corresponding eigenvalues

 $\begin{bmatrix} Covariance & matrix \end{bmatrix} \cdot \begin{bmatrix} Eigenvector \end{bmatrix} = \begin{bmatrix} eigenvalue \end{bmatrix} \cdot \begin{bmatrix} Eigenvector \end{bmatrix}$ 

- 3) Sort the eigenvectors according to their eigenvalues in decreasing order
- 4) Choose first k eigenvectors and that will be the new k dimensions
- 5) Transform the original n dimensional data points into k dimensions
- Stop criterium :
  - The reconstruction/projection error should be minimized





#### PCA: Applications

#### **Genetic cartography**

**Goal** : characterize genetic variation in European indiviuals

Data : Individual's genotypes



#### **Movie recommendation**

**Goal** : build a system that can make personalized movie recommendations to users

**Data** : historic user-community ratings

	Movie #1 Toy Story	Movie #2 Monsters Inc.	Movie #3 Saw	Movie #4 Ring	Movie #5 Hitch
User #1	4	5	1	NULL	4
User #2	5	NULL	1	1	NULL
User #3	5	4	3	NULL	3
User #4	5	4	1	1	NULL
User #5	5	5	NULL	NULL	3

### sklearn.decomposition.PCA

PCA : In practice

```
>>> import numpy as np
>>> from sklearn.decomposition import PCA
>>> X = np.array([[-1, -1], [-2, -1], [-3, -2], [1, 1], [2, 1], [3, 2]])
>>> pca = PCA(n_components=2)
>>> pca.fit(X)
PCA(n_components=2)
>>> print(pca.explained_variance_ratio_)
[0.9924... 0.0075...]
>>> print(pca.singular_values_)
[6.30061... 0.54980...]
```

• Non-linear extension of PCA

2nd component



$$x_i \to \varphi(xi)$$

Kernel PCA

- Kernel trick : the vectors  $\varphi(xi)$  appear only within scalar products (no mapping needed)
  - kernel function k(x,y) replaces the scalar product  $(\varphi(x) \cdot \varphi(y))$









### Kernel PCA: Applications

#### **Image denoising**

**Goal** : denoise images

Data : 256-dim handwritten digits



#### **Novelty Detection**

**Goal** : detect novel patterns in new data

Data : breast-cancer data



#### Kernel PCA : In Practice

>>> from sklearn.datasets import load\_digits >>> from sklearn.decomposition import KernelPCA >>> X, \_ = load\_digits(return\_X\_y=True) >>> transformer = KernelPCA(n\_components=7, kernel='linear') >>> X\_transformed = transformer.fit\_transform(X) >>> X\_transformed.shape (1797, 7)

### Non-linear method to visualize the structure of a high-dimensional dataset by reducing its dimension.

- N = number of training data points
- D = input dimension

Manifold Learning : Introduction



• Use cases :

- Clustering
- Dimensionality Reduction
- Semi-Supervised (labels)



#### Manifold Learning with 1000 points, 10 neighbors



- Locally Linear Embedding (LLE), Modified LLE (MLLE), Hessian LLE
- Local Tangent Space Alignment (LTSA)
- Multi-Dimensional Scaling (MDS), Isomap (extension of MDS)
- Spectral Embedding (SE)

Manifold

Learning :

Algorithms

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



N: number of training data points D : input dimension k : number of nearest neighbors d: output dimension

**Nearest Neighbors Search :** identify k closest neighbors of training instance x<sup>i</sup>

Weighted graph construction : transform the raw input data into graph representation (data points are treated as nodes of a graph) **Similarity measurement :** between points in the high dimensional space (Gaussian distribution)



keep similar instances close and dissimilar instances apart (step 2 similar to step 1, except that we use a Student t-distribution in the low dim space)

**LLE** : regularization problem if k>D (MLLE, Hessian LLE



**Partial eigenvalue decomposition** to map the training instances to a d-dimensional space

Map probabilities of the high and low dimensionality spaces as best as possible (KL divergence)

#### unroll twisted manifolds when there is **not too much noise**

Manifold Learning with 1000 points, 10 neighbors LLE (0.098 sec) Hessian LLE (0.21 sec) Modified LLE (0.18 sec) LTSA (0.13 sec) Use Cases 1 0<sup>1<sup>2</sup></sup> -1 0 SE (0.069 sec) t-SNE (6.8 sec) Isomap (0.38 sec) MDS (1.4 sec) visualization in a low-dimensional space of 2/3 dimensions visualise **similarity** or dissimilarity between samples



https://b.socrative.com/login/student/

Room : CONTI6128