

Machine Learning

3rd April 2019

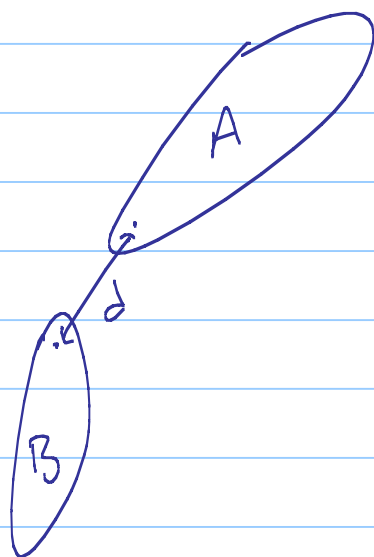
K-means clustering

- Supply the number of clusters. This is a hyperparameter
- The computer iteratively finds the center of k clusters
- At the initial step randomly choose k points as cluster centers, each data point is assigned cluster corresponding to its nearest cluster center. The cluster centers are now moved to the mean-position of their respective clusters.

Hierarchical clustering

- The number of clusters is a priori unfixed.
- We need to supply the value for minimal separation between two distinct clusters. Call this d_{min} .

→ If the closest points between two clusters are less than d_{min} apart, then the clusters are merged.



→ if $d < d_{min}$, then
merge A & B;
else A & B are
left as distinct
clusters.

→ At the initial stage, each data point is identified as an independent cluster. Then the computer iteratively merges the clusters by measuring their separation.

🚩 Numpy :

→ if A is a numpy array then
 $A[-k:]$ gives the last k elements
in the array.

Pandas :

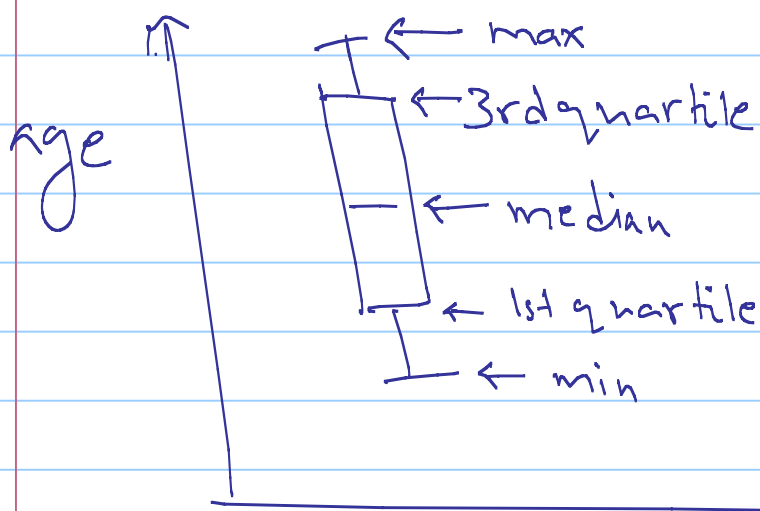
→ if df is a pandas dataframe, with columns : 'age', 'gender', ...

→ find average age for each gender

sol : `df.groupby('gender')['age'].mean()`

Can also get median, mod, std etc

→ `df['age'].plot(kind='box')` is the distribution of age



median : $(\# \text{ data points} < \text{median}) = (\# \text{ data points} > \text{median})$

4th April 2014

→ Make sure to use only the training data and NOT the cross-validation or test data during the preprocessing stages.



This is exactly why using a pipeline that integrates the various preprocessing and machine learning libraries is a good idea, as it automatically takes care of this.



Scoring on classification tasks:
confusion matrix

Prediction Labels ↓	→	
	Pos.	Neg.
Pos.	TP	FN
Neg.	FP	TN

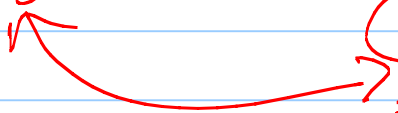
$$\text{Recall} : \frac{TP}{\text{Labelled positive}}$$

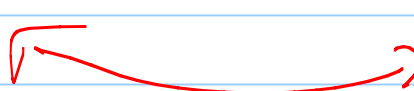
$$= \frac{TP}{TP + FN}$$

$$\text{Precision} : \frac{TP}{\text{Predicted pos.}}$$

$$= \frac{TP}{TP + FP}$$

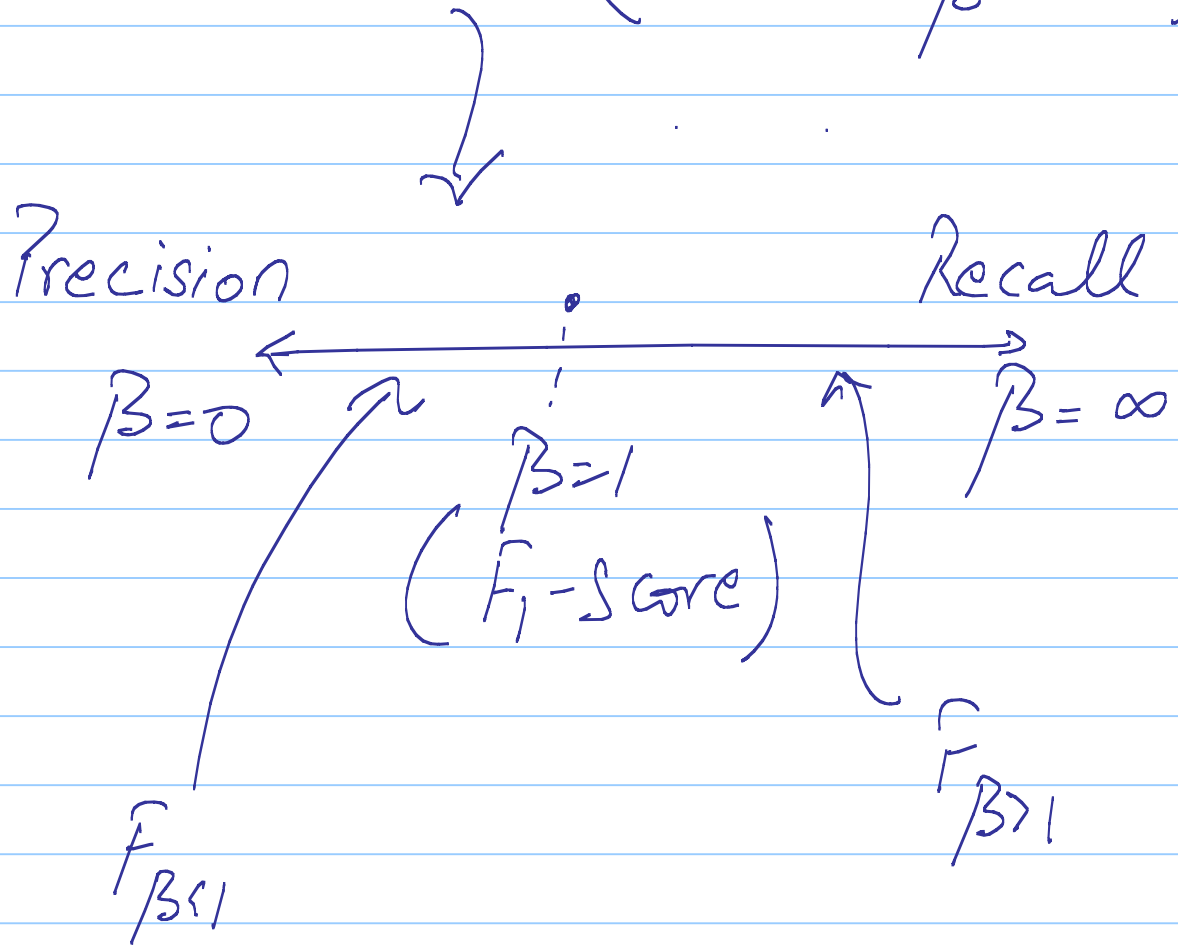
Mnemonic:

$$\text{Recall} = \frac{TP}{\text{Real pos.}}$$


$$\text{Precision} = \frac{TP}{\text{Predicted pos.}}$$


$$F_1 \text{-score} = \frac{\text{Prec.} \times \text{Rec.}}{\left(\frac{\text{Prec.} + \text{Rec.}}{2} \right)}$$

$$F_{\beta} \text{- Score} = \frac{\text{Prec.} \times \text{Rec.}}{\left(\frac{\beta^2 \text{Prec} + \text{Rec}}{1 + \beta^2} \right)}$$



$F_{\beta < 1} \rightarrow$ more imp. to Precision
 (For e.g. we'd rather have a spam email not get filtered than have an imp. email sent to spam. \Rightarrow prefer precision over recall)

$\hat{F}_{\beta > 1} \rightarrow$ more imp. to Recall
(e.g. we'd rather diagnose a healthy patient as mildly sick, then diagnose a sick patient as healthy \Rightarrow want less FN
 \Rightarrow want high recall)

ROC curve

\rightarrow Receiver operating characteristic
 \hookrightarrow curve of signal vs noise

Signal := True pos. rate

$$TPR = \frac{TP}{\text{Real pos.}} = \frac{TP}{TP + FN}$$

Noise = False pos. rate

$$FPR = \frac{FP}{\text{Real neg.}} = \frac{FP}{FP + TN}$$

→ Used to gauge the ability of a binary classifier's ability to distinguish between the positive class and the negative class.

→ To plot ROC we do the following : For each instance, the classifier will output a number.

We then decide a threshold & if that number is greater than the threshold, the instance is predicted to in the positive class.

OR equivalently
the prob. of being in the positive class

ROC curve is obtained by varying the threshold from very low to very high values and then

computing (TPR, FPR) pairs
for each value of threshold.

plotting the list

$$[(\text{TPR}, \text{FPR})_1, (\text{TPR}, \text{FPR})_2, \dots]$$

gives us the ROC curve.

$$(\text{TPR}, \text{FPR})_i = \text{TPR \& FPR for } i\text{-th threshold value}$$

→ If the classifier is able to distinguish well between pos. class & negative class, then area-under the curve (AUC) of ROC will be ~ 1 .

→ For a classifier than randomly guesses the pos. & neg. classes, its ROC-AUC will be ~ 0.5 .

Numpy as np

- `np.clip(array, min, max)`
 - `np.arange(start, stop, step)`
 - let 'a' be a numpy array
 \Rightarrow a.flatten() returns the result of flattening 'a'
 - `np.linspace(start, stop, num)`
 - if 'a' is a numpy array,
`a.sort(axis=axis)`
 \hookrightarrow returns the result of sorting a along the specified axis.
- \rightarrow Can also do this by np.sort(a)

- `np.argsort()`

→ returns the indices that would sort an array.

- `np.argwhere(some_condition)`

↳ returns the indices of the arrays elements where the given condition is true

🚩 `matplotlib.pyplot` as `plt`

- use `plt.imshow()` to plot images.