

Entropy of n -possibilities:

$$S = -\sum p_i \log_2 p_i$$

p_i = prob. of i -th ev.

Information $\propto 1/S$ i.e. Info \propto info ↓
if $p_i = 0$ for some $i \Rightarrow p_i \log_2 p_i = 0$

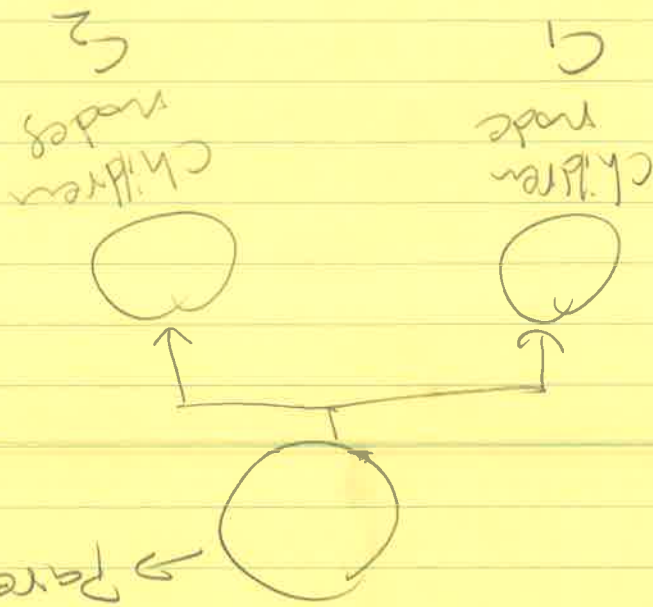
Information gain = Change in entropy

$$= -(S_f - S_i)$$

Basically: upon classifying objects based on the answers of a particular quest. (in dec. tree), if the entropy is S_f of the new of the given state is S_i

then Info is gained if $S_f < S_i$

→ Parent node



$$\text{Info gain} = \text{Spurvent} = -\frac{1}{2}(S_{C1} + S_{C2})$$

avg. of children nodes.

more general:

Info gain = Entropy (parent)

$$-\left[\left(\frac{m}{m+n} \right) S_1 + \left(\frac{n}{m+n} \right) S_2 \right]$$

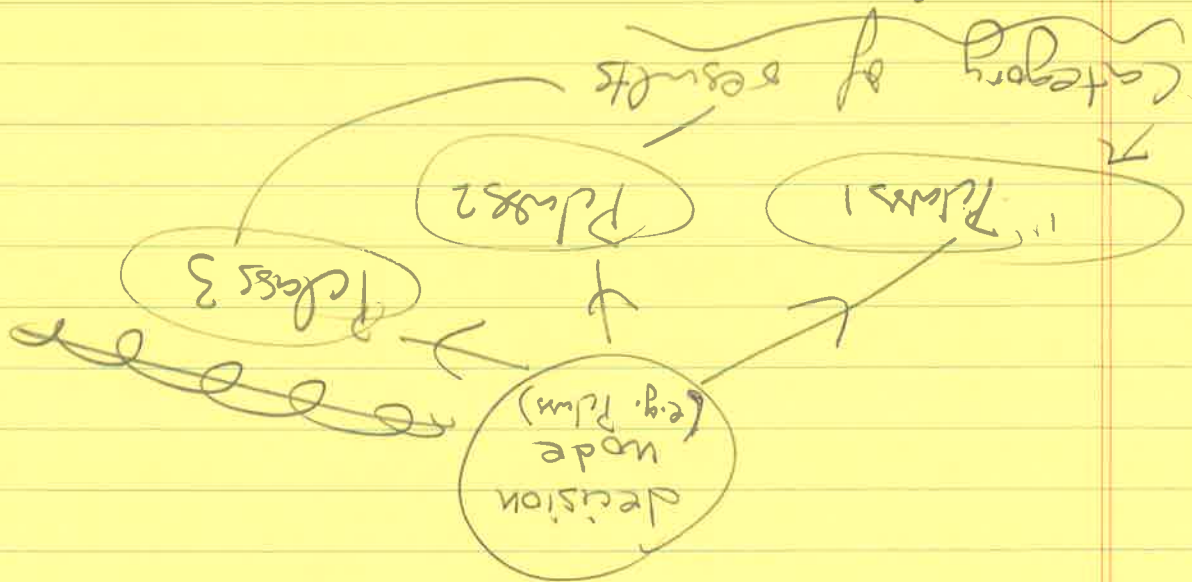
if S_1 contains only m of the training data instances
 1/a. S_2 contains n of the total $(m+n)$ instances.

→ Making decision trees: Pick the split that gives the largest information gain.

(Basically ~~check~~ try splitting by each feature & choose which one gives maximum entropy gain)

Decision trees tend to overfit a lot

In a decision tree, at each node we just split the data according to their result for the query in that node, & then compute the entropy for each category of the results



↳ These are called "leaves" of the nodes. → they become the children of the node we ~~started~~ split.

Hyperplan

max-depth →

underfits in small.

min-samples-leaf →

overfit if

is small

(ie. leaves more

min-samples-leaf

min-samples-split → overfit if min-samples is small

of iterations

Python display + display)

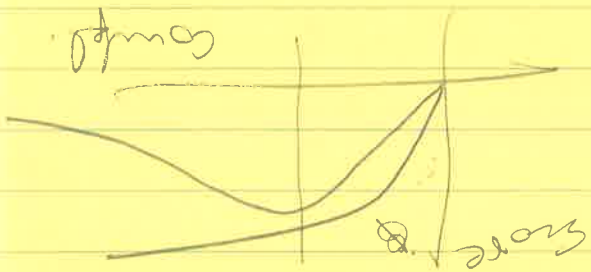
8, 1, 2 % on plot is inline

One-hot encoding
pd.get_dummies

0.820
0.81
0.798
0.793

min-sample-leaf = 2 \rightarrow 10
min-sample-split = 4 \rightarrow 3
complexity \uparrow
complexity \downarrow

to list



Naive Bayes

Bayes thm.

$$P(A|B) =$$

$$\frac{P(A \cap B)}{P(B)}$$

no. of ways A can happen given B has happened.

total ways of B happening = unconditional prob. of B

$$P(A|B) \neq P(B|A)$$

$$P(A|B)P(B) = P(B|A)P(A)$$

$$= P(A \cap B)$$

Phenomenons of false positives:

misdiagnosis

if prob. of error in diag. is much larger than probability of occurrence of the sickness, then

prob. of being actually sick upon testing positive is still low.

c.g. lets make a fictitious disease

that does not exist & prob. of some being afflicted with it is 0.

now devise a test for such a

disease and it is 99% accurate but

even if ~~there~~ someone test pos., their prob. of being sick is 15%.

~~$$P(\text{spam} | \text{each}) = \frac{P(A|BMC)}{P(A|BMC) + P(A|BNC)}$$~~

~~$$P(A|BNC) = \frac{P(A)P(BNC)}{P(A)P(BNC) + P(A)P(BNC)}$$~~

~~| | | | |
|-----------------|---------------|---------------|---------------|
| $\frac{n}{n+1}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |
| $\frac{12}{40}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |
| $\frac{n}{n+1}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |
| $\frac{12}{40}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |
| $\frac{n}{n+1}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |
| $\frac{12}{40}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |
| $\frac{n}{n+1}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |
| $\frac{12}{40}$ | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{2}$ |~~

Naive Bayes

Naive Bayes \Rightarrow Assume that

$$P(A|BNC) = P(A)P(BNC)$$

This is only true if the events are independent and not correlated

for e.g. in day of words teching we assume the probability of each word is independent of other words, when in real world this is not entirely true.

Naive Bayes estimator

$$P(A | g_1, g_2, \dots, g_n) = \frac{P(g_1, g_2, \dots, g_n | A) P(A)}{P(g_1, \dots, g_n)}$$

$$P(\bar{A} | g_1, g_2, \dots, g_n) = \frac{P(g_1, g_2, \dots, g_n | \bar{A}) P(\bar{A})}{P(g_1, \dots, g_n)}$$

But usually one can not evaluate the unconditional probs. $P(g_1, \dots, g_n)$

But $P(g_1, g_2, \dots, g_n | A), P(g_1, g_2, \dots, g_n | \bar{A})$

$P(A)$ & $P(\bar{A})$ can be obtained from labelled data

$$\Rightarrow P(A | g_1, g_2, \dots, g_n) \propto P(g_1, g_2, \dots, g_n | A) P(A)$$

$$P(\bar{A} | g_1, g_2, \dots, g_n) \propto P(g_1, g_2, \dots, g_n | \bar{A}) P(\bar{A})$$

→ evaluate proportionality constant by using

$$P(A | g_1, g_2, \dots, g_n) + P(\bar{A} | g_1, g_2, \dots, g_n) = 1$$

Use naive Bayes
 $\Rightarrow P(g_1, g_2, \dots, g_n | A) = P(g_1 | A) \times P(g_2 | A) \times \dots \times P(g_n | A)$

Pandas. read-table

↳ deprecated

use read_csv

Bag of words: Only the frequency of ~~each~~ each word appearing in a message is counted, the ordering of words is ignored

sklearn CountVecorizer

str.lower() → in Python string methods
n.lower() for string in x
**kwargs?

string.punctuation → have to import string.

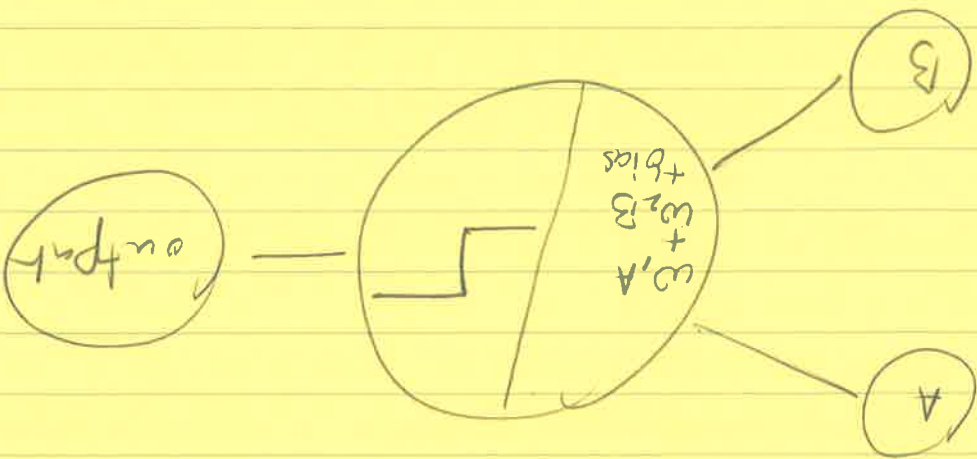
str.maketrans(x,y,z)

Counter() from collections in Python

Multinomial Naive Bayes: discrete values
Gaussian Naive Bayes: continuous values

→ Perceptron

graph → via udacity



python: zip()

~~map~~ how to properly use map & lambda

np.matmul

Pandas.map (to ~~convert~~ map named classes / labels into numericals)

Pandas.apply

$$\begin{array}{r} 340'9 \\ 11208 \\ \hline 45222 \end{array}$$

→ A good bench mark is the naive predictor, that naively predicts the answer for any input data to be the most frequently occurring class in the training database.

↳ Shows what a model without any ~~learn~~ intelligence will look like

→ problem of continuous data not being gaussian → log transform often helps.

Usefulness of SVM: Ability to find optimal & stable separating hyper planes, and loose over-fitting tendency

in a large feature space

(11, 14)
(11, 15, 16)
(9, 15, 16)
(16, 14, 6)

(18, 6, 18)

Decision tree has too many hyperparameters
This can make it difficult to compare
complexity of two models if too as
more of their hyperparameters are diff.
e.g. one has larger max-depth but larger leaf size
other has smaller max-depth but also smaller leaf size

sklearn.preprocessing: Polynomial features.

sklearn - linear Regressor

Require almost no data preprocessing
Decision trees are useful less prone to
sparse data & hence useful in cases
where unbalanced class distribution → eg
Cancer vs not cancer.

→ Limitation: tend to overfit

Scale logarithmically with size of data
→ can handle large data sets.

Random Forest: Bootstrapping
a sample

Panda: values?

Python: clone()?!

~~* train columns: values.~~
~~→ Pandas~~

Passing F_3 with a particular value for β

$$\text{score} = \text{make_score}(\text{fit_score}, \beta = 0.5)$$

Clustering

1) k -means clustering:

How to choose the correct no. of clusters?

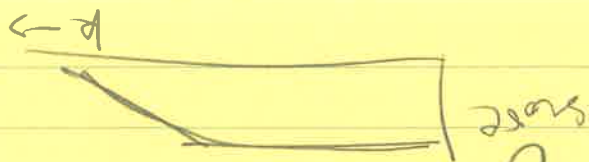
Answer: Elbow method!

or

compute mean squared distance of each ~~point~~ cluster and plot it against the number of clusters



Afterwards look at silhouette score



education, education-num, occupation, week-class

age, sex, capital-gain

age, education-num, capital-gain

Mar. Stat, education-num, age, capital-gain, hours-per-week

sex

39.7 & 45.7

45%
43%
@vAF

KMeans:

Silhouette ~~score~~ score: tells us how similar a point is to its own cluster as opposed to ~~a~~ diff some other cluster

How is this diff. from pandas.concat?

pandas.merge

pandas.to_csv
to convert

pandas.pivot-table

to sparse
DataFrame

Other clustering algorithms

→ Hierarchical clustering

→ DBSCAN: robust against noise

→ does not require no. of clusters

→ clusters densely packed data points & labels rest as noise

→ can be used

for anomaly / fraud detection
Yes, sample points are clusters and 1st noise
2 & min sample
s.t. must have 1st noise

Single link Hierarchical clustering

→ Dendrograms.

→ Very powerful to get insights about clustering

→ also contains lib. for hierarchical clustering methods

have import from scipy.cluster.hierarchy

Agglomerative clustering ~~is~~ also.

↳ included in sklearn

↳ Complete Link clustering

↳ distance between the

clusters is ~~measured~~ give by the distance between two

maximally separated points in the clusters

↳ than the two "nearest" clusters

are merged.

→ Average Link clustering:

↳ Distance between clusters

is the average distance

between two points from

the clusters.

→ Ward's method

↳ Tries to ensure that merging

clusters leads to least increase

in variance in the cluster after

merging

Cluster Validation Index

→ Since ~~correlation~~ clustering is based on distance, if different features have different scales, then their contribution to the distance will get weighted accordingly. So it best to normalize all of them to same scale

Note that the objective is to make sure the variance in each feature is ~~the~~ given equal weight.

Seaborn clustermap.

→ DBSCAN

→ Gaussian Mixture Models (GMM)

→ ~~Did~~ Used in Speaker Verification
(see paper by Douglas Reynolds, Thomas F. Quatieri, Robert B. Dunn)

Also in Biometric Verification.

~~Cluster indices~~

→ Cluster Validation Index

1) External indices (adjusted Rand score)

2) Internal indices (e.g. Silhouette

co-efficient)

Should not be used for DBSCAN

Also not

the most optimal score if the data points is not the form of dense circular cluster

if it will not give good score or be good for scoring

ie data is in the form of two co-circular clusters

Use DBCV

→ Feature Scaling :

Projecting data on the direction of maximum variance
↓
we'll have minimum loss of information

The principal components of a data are always chosen to be perpendicular to each other.
Rule of thumb: Try to use as many components as you can, such that 80% of variance gets explained.
PCA can be used for noise reduction.

Eigenspace : PCA for face recog.

(sp. learn choose

PCA with sub-samples

= "randomized"

When feature is too large for PCA

Random Projection: (based on Johnson-Lindenstrauss lemma)

→ find component along randomly chosen directions
→ might not work for low dimensional feature spaces
→ looks quite with when feature space is large.

Basically
projection is a
matrix of
random
no.

→ Gaussian Random Projection
→ Sparse Random Projection.

→ No. of samples must be greater than the no. of components we seek.

ICA: Independent Component Analysis

→ hopes/assumes that features are admixtures of independent sources & tries to extract these independent sources.

→ for e.g. 3 or more independent sources of sound in a recording.

→ This problem is called "blind source separation"

(Fast ICA)

Writing a dataset?

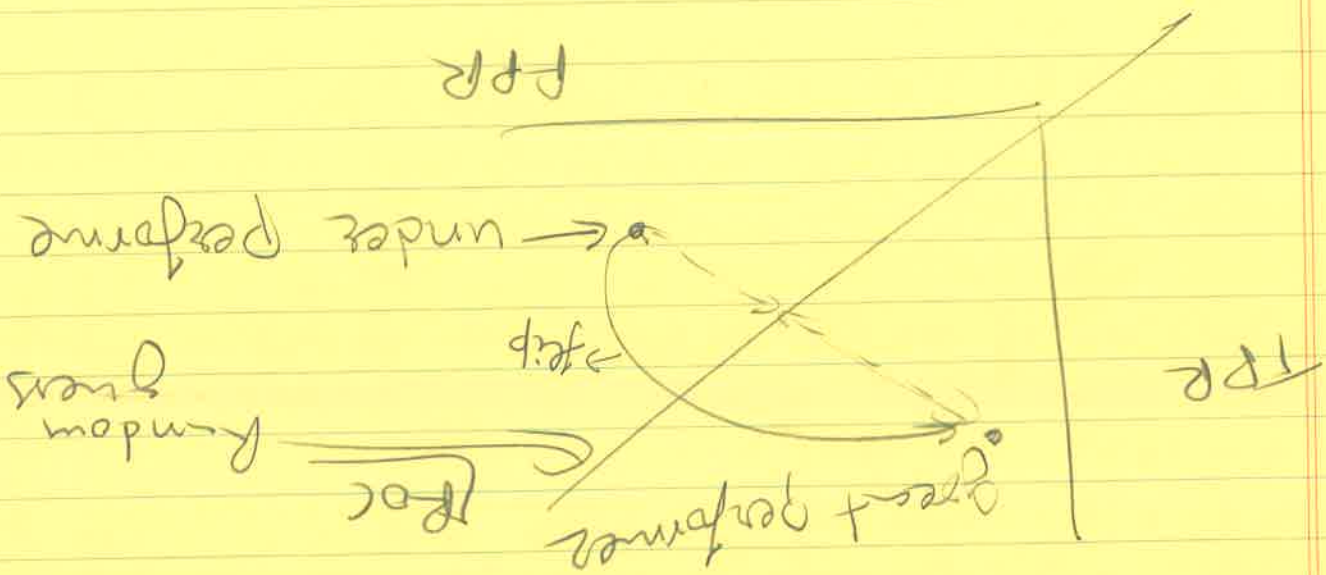
→ Read & understand central limit thm.!

The independent components must have ~~Non-Gaussian~~ distribution (V. imp)

Imp

Note that in classification (for 2-class classification) random guesses will give 50% accuracy.

A classifier that underperforms a random guess can be easily mirrored into a good classifier by its ~~pre~~ predictions. Making the opposite of its ~~pre~~ predictions.



Python zip
np.ravel, np.flatten