

Crash course - Data Science

Extraction of knowledge from large volume of data that might be structured or unstructured.

Steps in Data science

1. Descriptive statistics.

clustering and basic data summary.

2. Inference.

Making conclusions about populations from samples.

Estimation

confidence intervals

hypothesis test

variability

3. Prediction / classification

Cross Validation

4. Experimental Design

Randomization / sample testing

To get better generalizability of the result population

Machine learning

1) Unsupervised learning

Supervised learning

- Take input and return a prediction.

Nothing to train the

algorithm.

Predictions based on

most information data observed

\Rightarrow observed outcome
and predict when it's
not observed.

- Evaluate performance based on prediction performance.

\Rightarrow Traditional statistics

Sample generalized to superpopulation
 \Rightarrow predictions.

Software Engineering

\Rightarrow Procedure for standardizing / building models
or procedure.

\Rightarrow Interface between I/p & O/p
I/p or O/p are not that important
Interface is import.

Random Forest Regression

- Ensemble learning

Same algorithm multiple times and put them together

- Take prediction from N Decision trees and make them predict together.

→ Average - stable (change in one tree won't impact number of trees)

Average
of all
predictions.

- More accurate prediction.

Regression trees.

→ Algorithm will create a split in information entropy plot

if Split is increasing amount of information we have about the point.

less than 5% of data in each leaf and thus algorithm will stop.

→ Taking average of values from each leaf after splits.

And assigning dependent value as the average.

→ Multiple trees hence more accurate values.

Information Entropy

$$H = \sum_{i=1}^n p_i \times \log_2 \left(\frac{1}{p_i} \right)$$

outputs of data science experiment.

1. Reports
2. Presentations
3. Interactive web pages
4. Apps

- least square error function

$$\sum_{t=1}^n (y_t - \hat{y}_t)^2$$

Gradient descent

To construct
random forest

- ID3 algorithm.

Measure the purity of split.

- Information Gain

$$= H(S) - \sum_{\text{old}} |S_V| \frac{H(S_V)}{|S|} \sum_{\text{new}}$$

Most clarity with highest gain.

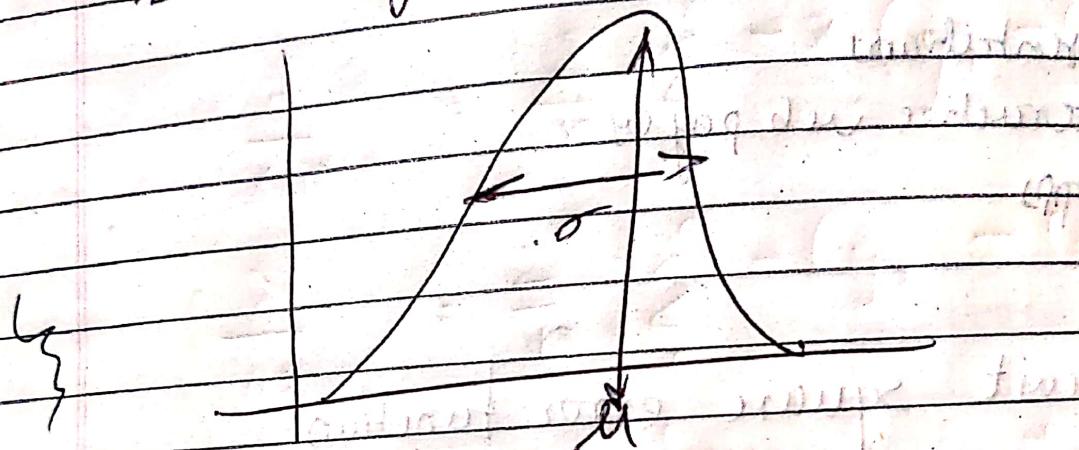
Gain

Date: 11/6

Gain Ratio =

$$\text{Split Entropy}(S_1) = -\sum \frac{(S_i)}{|S|} \log \frac{(S_i)}{|S|}$$

Linear Algebra



$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left(\frac{-(x-\mu)^2}{2\sigma^2} \right)$$

Approximation of goodness

↓
prediction.

- Add attribute as Vector.

(common way of measury information gain: decrease in G)

Gain impurity

$$I(C) = \sum_{k=1}^K P_k (1 - P_k)$$

$$P_k = \frac{1}{|S|} \sum_{i \in S} \mathbb{I}[y_i = c_k]$$

↑
Set of
observations

$$S_{\theta}^k = \{ i : x_i^j > \tau, \forall i \in S \}$$

Value of j^{th} feature for i^{th} observation

* Optimization for the split

$$\beta^* = \arg \max n_S I(S) - n_L I(S_L^U) - n_R I(S_R^R)$$

- Binary Space partitioning trees.

- ABT

Gradient Descent & Cost function

> Prediction function.

Finding a function with the use of

$1/p$ & $\partial p/\partial \theta$ values.

Best fit line of scatter plot

Mean square error $\rightarrow \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

number of data points \rightarrow To reduce -ve problem

Cost function.

$$= \frac{1}{n} \sum_{i=1}^n (y_i - (mx_i + b))^2$$

→ By Lj iteration you can find a proper value of function.

- learning rate in conjunction with the slope to reach minimum point on cost function

$$\text{Cost fun}^n = \frac{1}{n} \sum_{i=1}^n (y_i - (m x_i + b))^2$$

$$\frac{d}{dm} = \frac{2}{n} \sum_{i=1}^n x_i (y_i - (m x_i + b))$$

$$\frac{d}{db} = \frac{2}{n} \sum_{i=1}^n (y_i - (m x_i + b))$$

$m = m - \text{learning rate} * \frac{d}{dm}$

to get next point

np array \leftarrow faster than usual.

- Algoithms to perform split
linear discriminant analysis (LDA)
canonical correlation analysis (CCA)
logistic regression.

Supervised learning - overfit noise
reduce diversity.

- Cancer statuary.

Hyper parameter is a parameter whose value is used to control the deny process.

CCF canonical correlation forest

- Clustering

Flat

- Group number decided

Hierarchical

- Number of groups decided

Given feature set

→ machine decide the clustering.

→ k means. — Flat

→ Mean shift — Hierarchical

- 1) Find centroid.
2) find feature distance from the centroid
3) take mean of all the centroids
And get new centroids
4) keep repeating process until centroids
are not moving anymore.

This supervised learning