

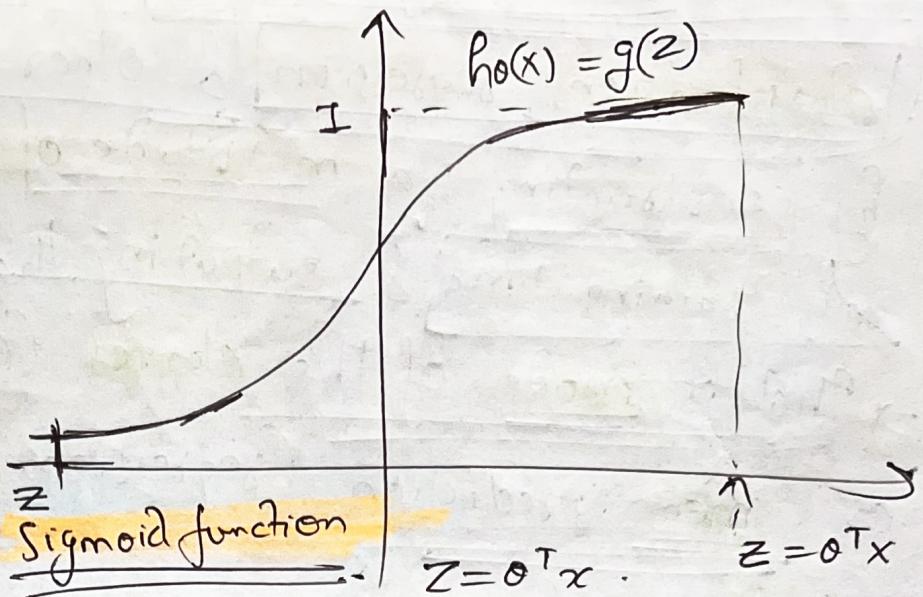
Logistic Regression

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$e^{\frac{1}{\theta^T x}} \rightarrow 0$$

$$\theta^T x \gg 0$$

$$\theta^T x < 0$$



if $y=1$, we want $h_{\theta}(x) = 1$, $\theta^T x \gg 0$

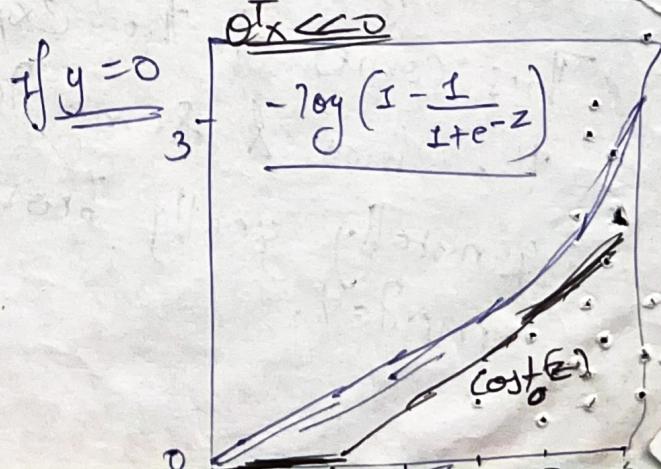
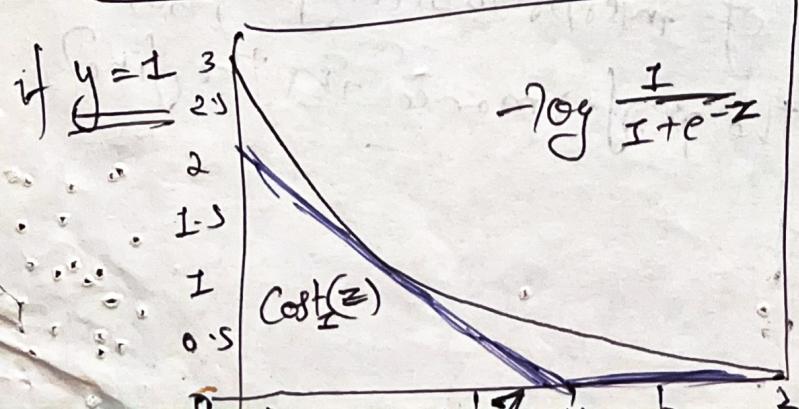
if $y=0$, we want $h_{\theta}(x) \approx 0$, $\theta^T x \ll 0$

Cost function of logistic regression.

$$-\log(h_{\theta}(x)) + (1-h_{\theta}(x)) \log(1-h_{\theta}(x))$$

$$= -y \log h_{\theta}(x) + (1-y) \log(1-h_{\theta}(x))$$

$$= -y \log \frac{1}{1+e^{-\theta^T x}} + (1-y) \log \left(1 - \frac{1}{1+e^{-\theta^T x}}\right)$$



SVM

It's a powerful yet flexible supervised ML algorithms, which are used for classification & regression.

Working of SVM (basically a representation of different classes in a hyperplane)

Hyperplane it is generated in an iterative manner by SVM so that the error can be minimized.

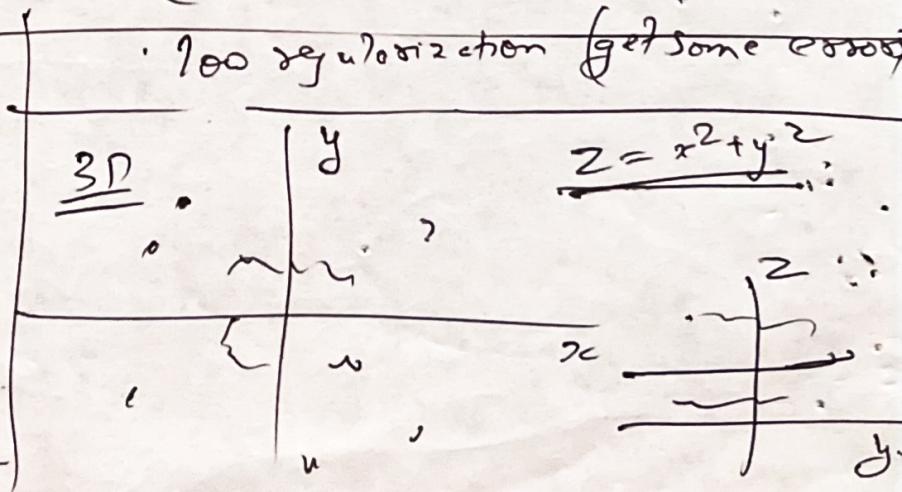
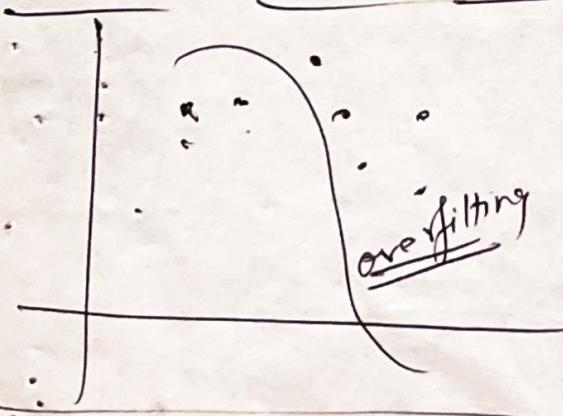
Goal of SVM is to divide the datasets into classes to find a Maximum Marginal Hyperplane (MMH).

Support Vectors: Points close to hyperplane. Separating will be done using these.

Alternative view



3D case (Hyperplane) (for N Dimen)
Hyperplane (gamma & regularization) (c)



Ensembling

When we try to predict the target variable using any ML technique, the main cause of diff. b/w actual & predicted values are: noise, variance & bias.
Ensembling helps to reduce these factors (except noise).

Define: An ensemble is a collection of predictors, which come together to give a final prediction.

The reason we use ensembles is that many different predictors trying to predict some target variables will perform a better job than any single predictor alone.

Ensembling

Bagging

eg. Random Forest

~~A++~~
Handles overfitting
Reduces variance
~~+++~~
Independent Classifiers

Boosting

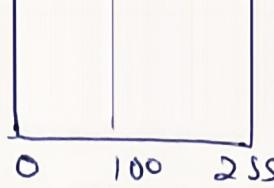
eg. Gradient Boosting

~~A++~~
Can overfit
Reduce bias & variance
~~+++~~
Sequential classifiers

Bagging: is a way to decrease variance in the prediction by generating additional data from the for training from dataset using combinations with repetitions.

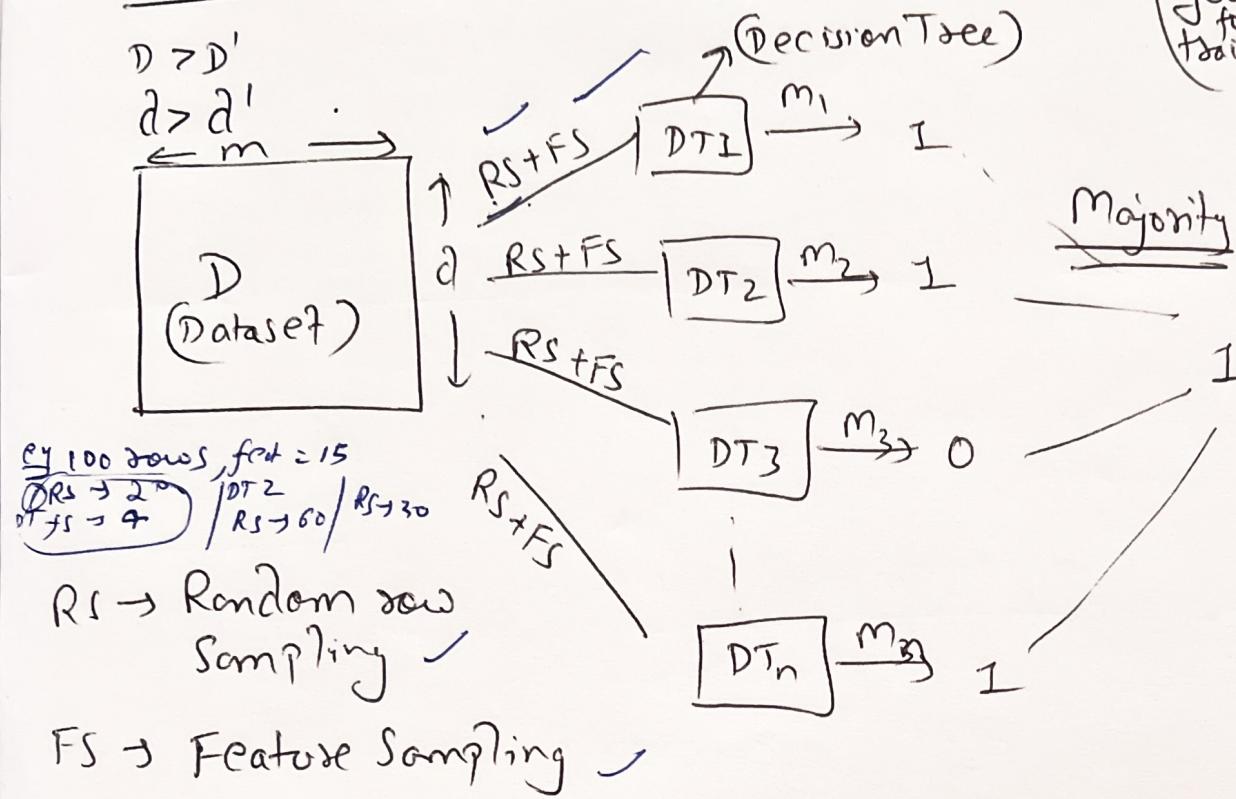
Boosting: It is an ensemble technique in which the predictors are not made independently, but sequentially. This technique employs the logic in which the subsequent predictors learn from the mistakes of the previous predictors.

Bagging: It is simple ensembling technique in which we build many independent predictors/modals & combine them using averaging technique.



MNIST (Handwritten digits)

'Random Forest' (Bagging)



$m, d \rightarrow$ no. of cols, no. of rows

Note: In case of Classification Problem, the output is Majority vote or mode.

Whereas,
 f_n

Regression Problem, output is mean or median.

Note

(will be good for training set)
 In DT
 ① low Bias
 ② High var.
 ↓
 (not good for test data)

Whereas;

In RF;
 high var → gets converted to low var.

Random Forest or random decision ~~trees~~ forests

are an ensembling learning method for

classification, regression

& other tasks that operate by constructing a multitude of decision trees

at training time & outputting the class.

That is mode of the classes (classification)

or mean prediction (regression) of the individual

trees.

Random Forest corrects for decision trees' habits of overfitting to their training sets.

Note:

In general, trees that are grown very deep tend to learn highly irregular patterns: they overfit their training set i.e. have low bias & high variance.

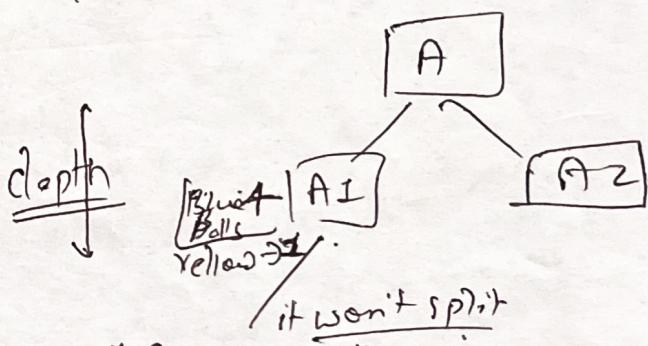
Random Forest are a way of averaging multiple deep decision trees, trained on different parts of the same training data set, with the goal of reducing the variance.

This comes at the expense of a small increase in the bias & some loss of interpretability, but generally greatly boosts the performance in the final model.

Random Forest (Hyperparameters)

Stratification ??

- ① n_estimators (no. of decision trees)
default 10
- ② max_depth (Specifies max-depth of each tree).
default None
- ③ min_samples_leaf (if specifies the min. no. of samples required to be at a leaf node).
Default = 1
- ④ min_samples_split (min. samples required to split on internal leaf node)
Default = 2



Code:

```
from sklearn.model_selection import GridSearchCV.
```

```
from sklearn.ensemble import RandomForestClassifier.
```

```
param_grid = [ 'max_depth': [30, 50, 80],  
              'max_features': -- ]
```

Model

$\text{rf} = \text{RandomForestClassifier}$
 $\text{grid_search} = \text{GridSearchCV}(\text{estimator} = \text{rf},$
 $\text{param_grid} = \text{param_grid}, \text{cv} = 3, \text{n_jobs} = -1, \text{verbose} = 2)$

Fit to model

```
grid_search.fit(X, Y)
```

$X, Y \rightarrow (\text{training dep. \& indep. variable})$

```
grid_search.best_params_
```

max_depth = 80

max_features = 2

n_estimators = 100

min_samples_split = 8
-Leaf = None

$\text{rfc} = \text{RandomForestClassifier}()$

$\text{model_rf} = \text{rfc} \cdot \text{fit}(X, Y)$

$\text{Y_pred_rf} = \text{model_rf} \cdot \text{predict}(X)$

$KFold = \text{StratifiedKFold}(\text{n_splits} = 10)$

$\text{result} = \text{cross_val_score}(\text{rfc}, X, Y, cv = KFold)$

Random Forest

ex) of parameter of random forest.
↳ are variable & thresholds used
(to split each node learned
during training)

ex) of hyperparameter of RF.

① no. of decision trees in the forest
(n-estimator)

no. of features considered
while splitting a node.

(max-features)

max number of levels in a decision tree
(max-depth)

min number of data points placed in a node before split

(min-samples-split)

min number of data points allowed in a leaf

(min-samples-leaf)

For ex if min-samples-split = 5 & these are 7 samples
& on internal node, then split is allowed -
2 leaf's split resulting 2 leaves, one with sample of size 3 &
another with 6 size leaf if min-samples-leaf = 2
then won't be allowed

Parameter

Model parameters
are learned during
training - such as
slope, intercept
in linear regression

Hyperparameters

set by Data-scientist
before training

Decision Tree

Splitting measures

(Starting from
root node &

further down by
splitting the nodes)

Gini Index

Entropy (Information gain)

6 Gini Index / Gini Impurity.

measure the degree or probability of a particular variable being wrongly classified when it is randomly chosen.

But what is impurity?

If all elements belong to same class than it is pure.

Note The degree of Gini Index varies from

0 to 1, where 0 denotes that all

elements belong to a certain class or if there exists only one class.

Note

$\& \Sigma$ denotes that the elements are randomly distributed across various class.

Note

A Gini index of 0.5 denotes equally distributed elements into some classes.

Formula of Gini Index.

$$\text{Gini} = 1 - \sum_{i=1}^n (p_i)^2$$

where, p_i is the probability of an object classified to a particular class.

Note

We prefer attribute/feature with

lowest Gini Index of root node.

Formula for Entropy

$$E(S) = \sum_{i=1}^c -p_i \log(p_i)$$

Not preferred due to
→ 'log' function of it
increases computational complexity.

where; P denotes the probability & $E(S)$ denotes entropy.

Note) we prefer attribute with least entropy
at root node.

- Note UDF
User Defined functions make it
possible for user to write its own
transformation.

from pyspark.sql.functions import udf.