

Applied
Machine

Learning

With Python

Applied Machine Learning with Python

ML - The study of computer program or algorithms that can learn by example.

- * Algorithms learn rules from labelled examples
- * Training data - The labelled examples used for training of algorithms are training dataset.
- * Test data - The dataset used to evaluate the algorithm whether it is working correctly to unknown dataset which is not in Training.

Supervised learning - Algorithms that uses labelled datasets to predict target values of unknown datasets.

i) Regression - Target values are in continuous forms. There is relationship between input given and target value and that relationship is continuous.

ii) Classification - Target values are discrete values or discrete classes. And target value is need to be identified as whether it is belonging to particular class or not.

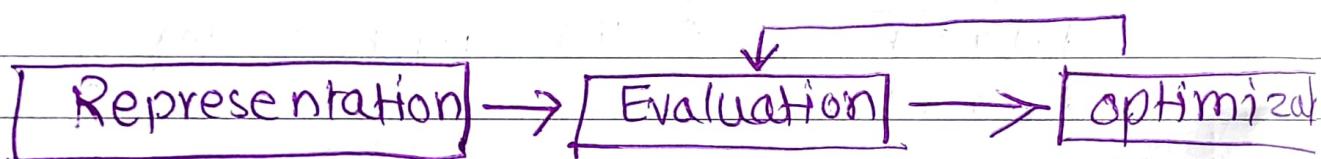
iii) Decision tree - Target values are determined by set of (if else) questions

Unsupervised - Algorithm do not use labelled data and finds structure in unlabelled data directly to make predictions.

i) clustering - Find group of similar instance in the data.

ii) outlier detection - Finding unusual patterns

Basic Machine learning workflow



e.g. choose:

- i) A features representation
- ii) Type of classifier / Algorithm to use

e.g. image pixel with K-nearest neighbour classifier

choose :-

- i) what criteria used to distinguish good v/s bad classifier

e.g. % of correct predictions on test set

choose:

How to search for settings or parameters that gives best classification for this evaluation criteria.

e.g. Try a range of values of 'K' in K-nearest neighbor classifier

Examining the data

- 1) certain pre processing is needed for dataset.
e.g. to see whether some cells in dataset are missing or not, which affects model badly.
- 2) sometimes examining dataset can make you realize that no machine learning algorithm required for problem solving.
- 3) for inspecting the dataset we use visualisation methods like feature plots.

scikit learn -

i) scipy - i) with scikit learn it provides support for sparse matrices a way to store large tables that consists mostly of zeros.

- ii) provide useful scientific tools like -
 - a) statistical distributions,
 - b) optimization of functions.
 - c) mathematical functions.

ii) Numpy - i) provides fundamental data structures used by scikit learn like multidimensional array.
 ii) typically data input to scikit is Numpy array.

iii) Pandas - i) provides key data structures like data frames.
 ii) Also support Reading /writing data in dataset.

iv) Matplotlib - plotting libraries to show spread of data.

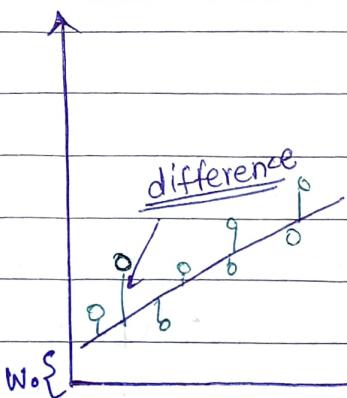
Regression

Takes set of training instances and learn a mapping to target value.

Target value is continuous (floating point)

Real valued target values are required regression model.

Least square linear regression



w_0 and w_1 are parameters defined by many methods correspond to different fit criteria and goals and ways to controls model complexity.

$$RSS = \sum_{i=1}^n [y_i - (w_0 + w_1 x_i)]^2$$

The learning algorithm finds the parameters that optimize an objective function that minimize the error between predicted target value and actual target value.

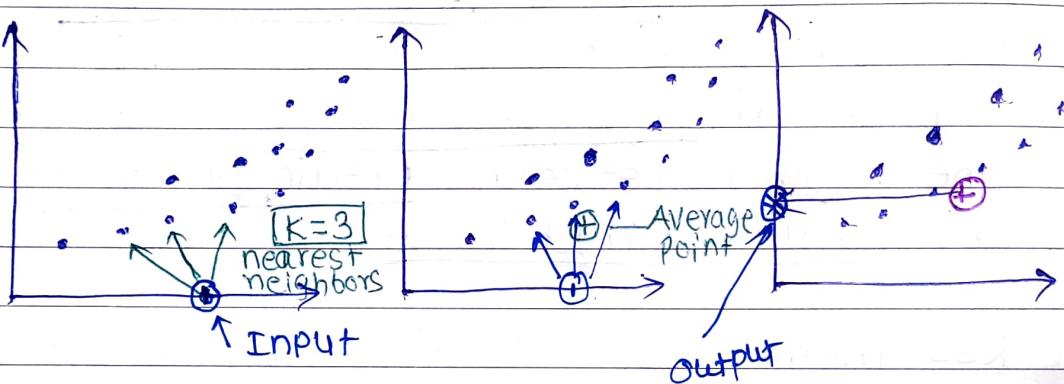
RSS method used to minimize the difference of predicted and actual values.

for Least square Linear Regression line is always straight.

- * The value of w_0 & w_1 ... w_n are finded by least square method i.e. the line which gives lesser RSS is chosen.

KNN & Regression

- 1) KNN Regression models finds the nearest neighbor to input and no. of neighbor based on value of $[k]$ then
- 2) Take average of all this points (neighbors) and according to regression model it defines target value for given input.



R^2 - (r-squared) Regression Score - Measure of determining how well Regression model is Predicting the data .

Score (0 to 1)

0 - corresponds to constant model which gives mean value of all training data given.

1 - corresponds to perfect predictions.

Also known as coefficient of determination.

Key Parameters

Model complexity

n_neighbors - number of nearest neighbors

[default = 5]

Model fitting

metric - distance function, [Default $\Rightarrow p=2$] Euclidean metric

Ridge Regression (L2 Regularization)

Ridge regression learns w_i, w_0 using same least square criterion but adds a penalty for large variations in w parameters

$$(RSS)_{\text{Ridge}}(w, w_0) = \sum_{i=1}^n [y_i - (w_i x + w_0)]^2 + \alpha \sum_{j=1}^n w_j^2$$

Once the parameters are learned, the ridge regression prediction formula is the same as ordinary least squares.

Regulation - method that prevents overfitting by restricting the model, typically to reduce complexity by addition parameters.

Regulation is controlled by α term higher the α higher is the regulation and less complex model and, ~~more~~ less likely to be overfitted.

Lasso regression (L1 Regularization)

minimize the sum of the absolute values of the coefficients

$$\text{RSS}_{\text{Lasso}}(w, w_0) = \sum_{i=1}^n [y_i - (w x_i + w_0)]^2 + \left[\alpha \sum_{j=1}^n |w_j| \right]$$

This has the effect of setting parameters weights in w to zero for the least influential variables. This is called sparse solution kind of feature selection.

The parameter α controls regularisation
default value of $\alpha = 1.0$

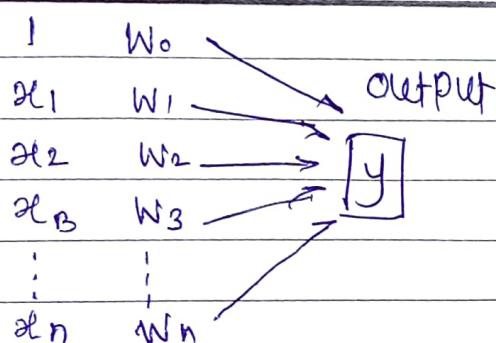
When to use Lasso Vs Ridge

Many small/medium sized effects : use ridge
 Few medium/ large sized effects : use Lasso

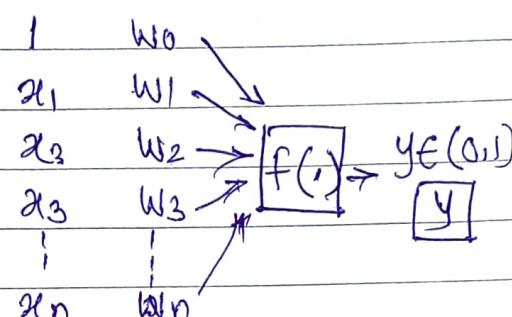
Regulation becomes less important as amount of training data increases.

Logistic Regression

Linear regression



Logistic Regression

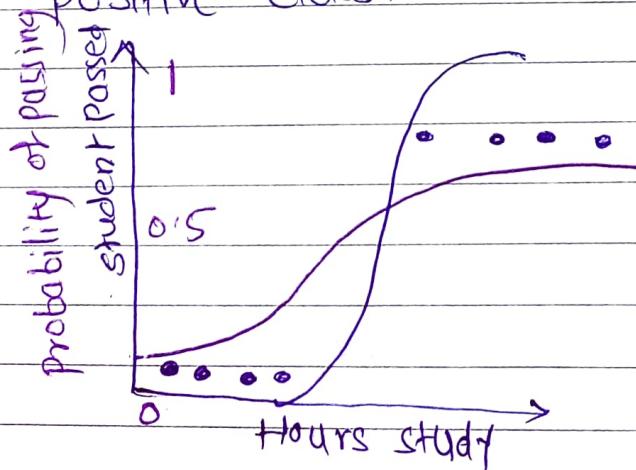


$$\hat{y} = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n$$

$$\hat{y} = \text{logistic}(w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n)$$

$$\hat{y} = \frac{1}{1 + \exp[-(w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n)]}$$

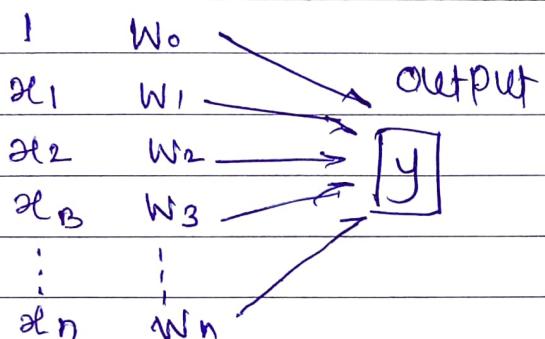
Logistic regression gives output values that falls between (0 to 1) that signifies output or target value falls in positive or negative class. It is the probability the input objects belong to positive class.



Regularization becomes less important as amount of training data increases.

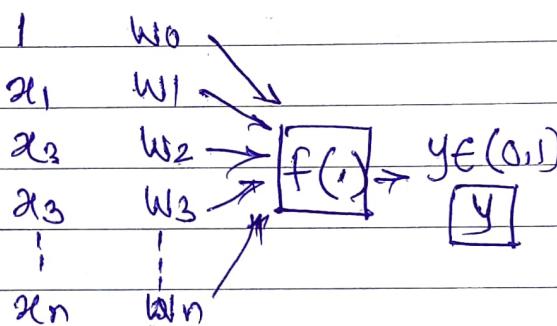
Logistic Regression

Linear regression



$$\hat{y} = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n$$

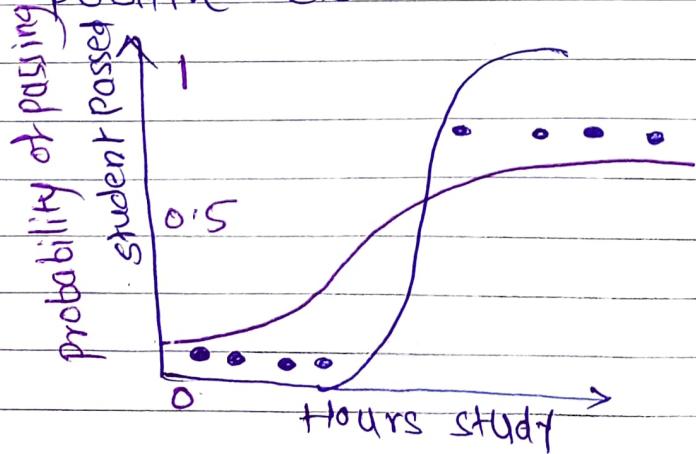
Logistic Regression



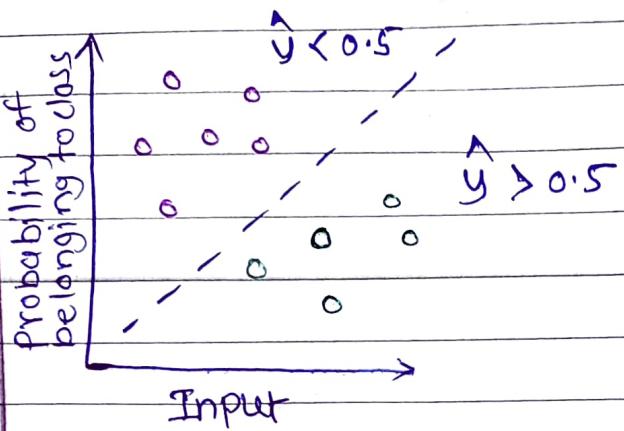
$$\hat{y} = \text{logistic}(w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n)$$

$$\hat{y} = \frac{1}{1 + \exp[-(w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n)]}$$

Logistic regression gives output values that falls between (0 to 1) that signifies output or target value falls in positive or negative class. It is the probability the input objects belong to positive class.

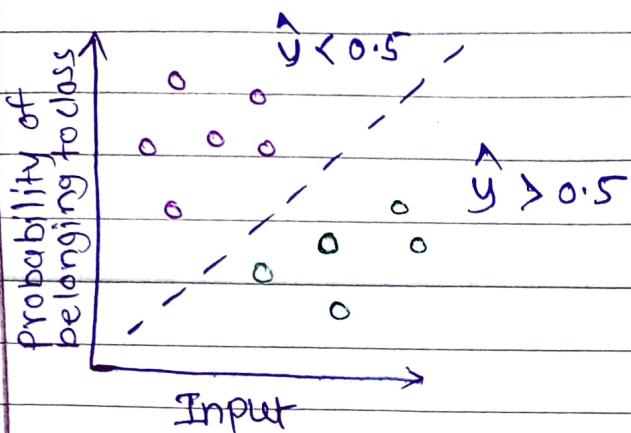


Logistic regression is used for Binary classification.



logistic regression is used for Binary classification.

Logistic regression is used for Binary classification.



logistic regression is used for Binary classification.

Classification

takes set of training instances and learn a mapping to a target value.

for classification target value is discrete class value.

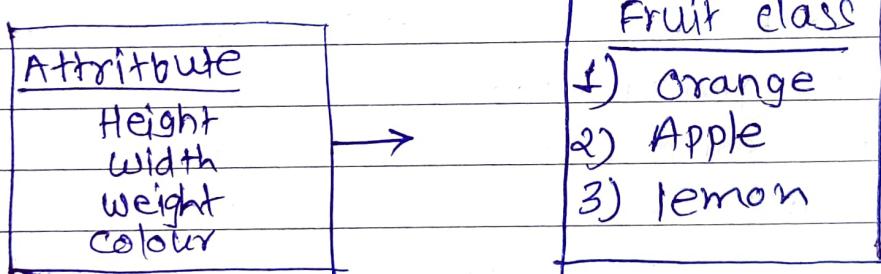
1) Binary

target value is 0 (negative) or 1 (positive)

e.g. fraudulent credit card transaction

2) multiclass

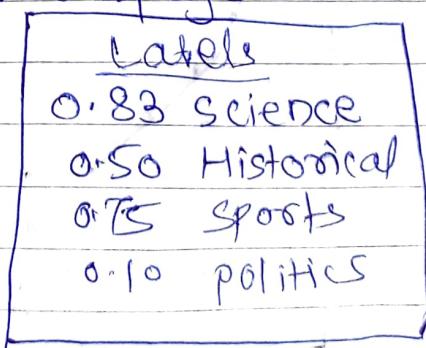
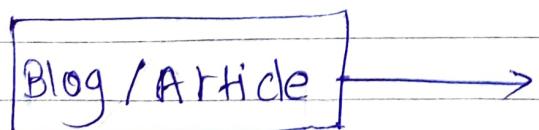
target value is one^{out} of a set of discrete values
e.g. labelling the Name of fruit from given attributes.

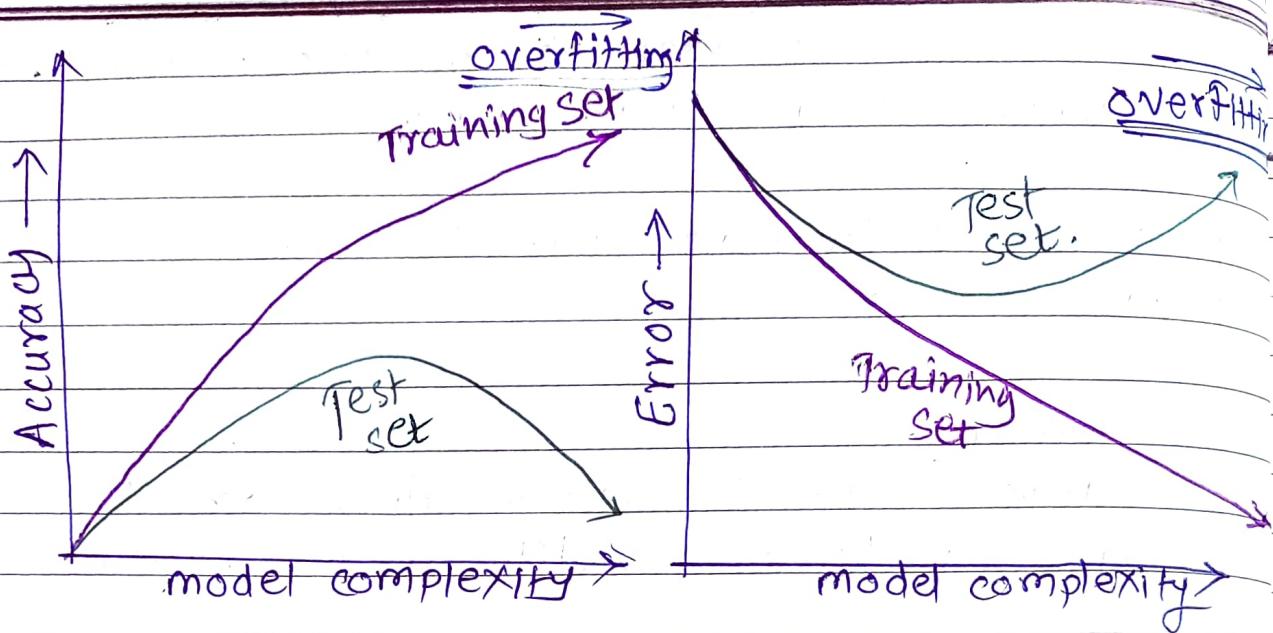


3) multilabel classification

target values are multiple.

e.g. labelling the topic on webpage.





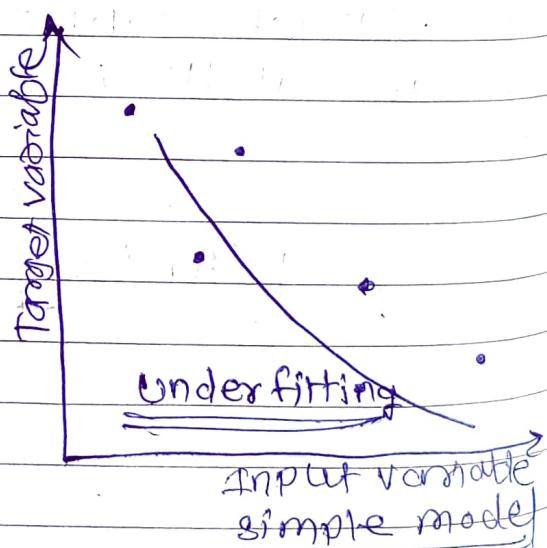
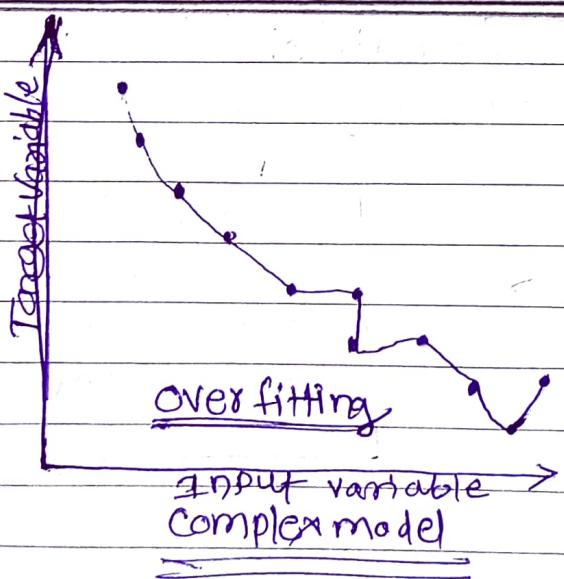
Accuracy and Errors are need to be optimal
If model over fits Accuracy of test set decreases and Error increases:

overfitting

Models that are too specific for training data are not giving good results for test data they are said to be overfitted.

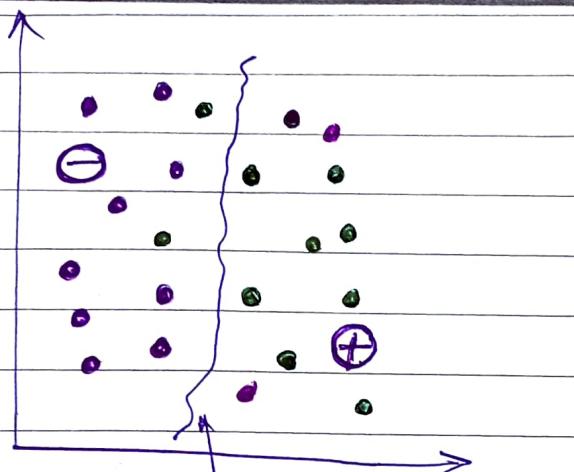
Underfitting

Models that are too simple, not fitting train data well are underfitted.

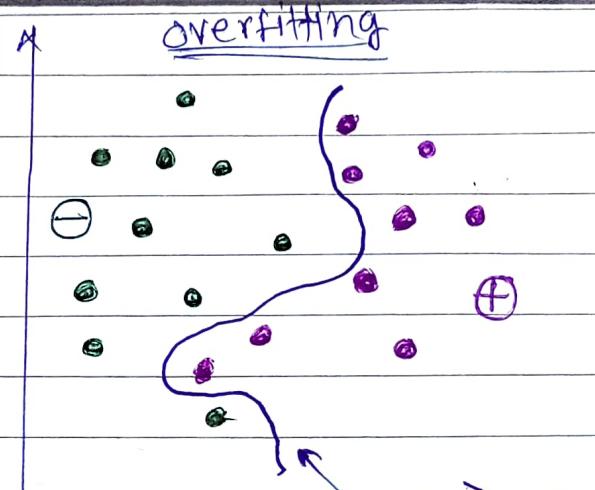


Regression Underfitting and overfitting.

Classification overfitting and underfitting



decision boundary
is not clear.
It allows some
points fails at across
hence not overfitted
and good model.

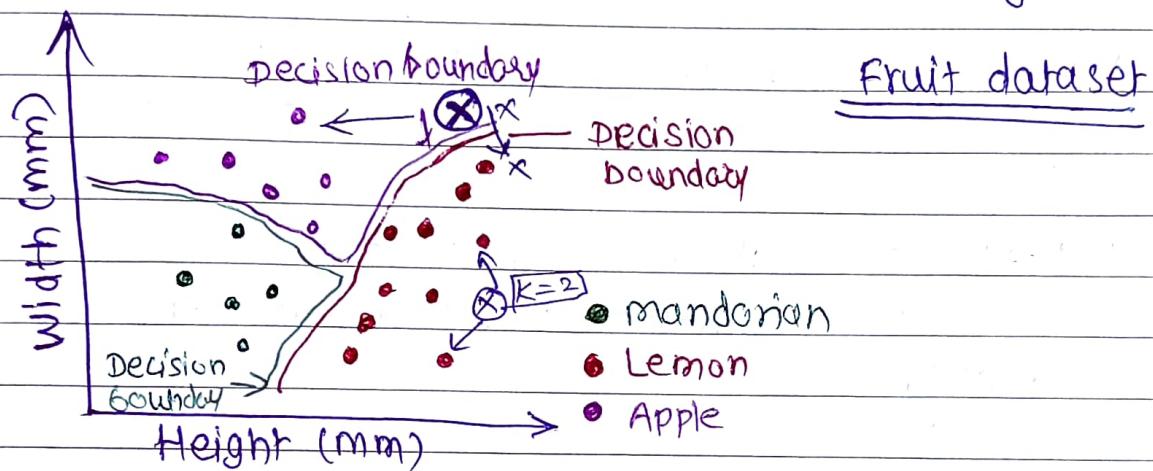


overfitting
decision
boundary
Too
overfitted

Decision boundary defines complexity of classifier model. It should not be too distinct to cover every point, And able to ignore minor variations.

K-nearest neighbor classifier.

K-NN - It is also called as instance based or memory based classifier, because it memorize the instances given to it and make further predictions based on that memory.



How to predict

- when new instance is occurred we find the the nearest point to than instance (query) and label of that point is assigned and predicted.
- For $K > 1$ the two or more according to K no. of points nearest to query point are found.
- Point nearest to same region is mapped and not from different region / class. this classes regions are distinguished with decision boundaries.

Needs four things to be specified.

i) A distance metric

Typically Euclidean metric used by default in scikit with $p=2$ minkowski metric used.

ii) How many nearest neighbors. $[k=?]$

iii) Influence of most nearest neighbors.

optimal weighting functions on the neighbour points.

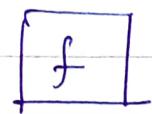
iv) How to aggregate the classes of neighbours

simple majority vote - class with most representatives with nearest neighbor

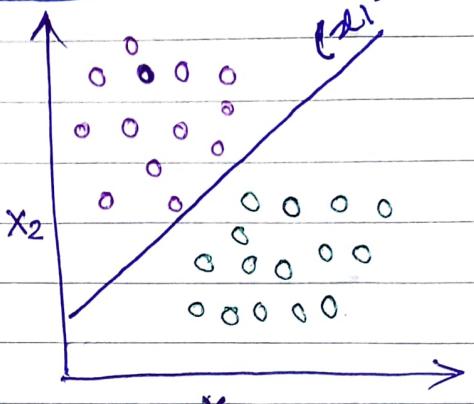
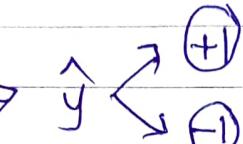
Linear classifiers : support vector machines

Feature
Vector

[Input]



[Output]



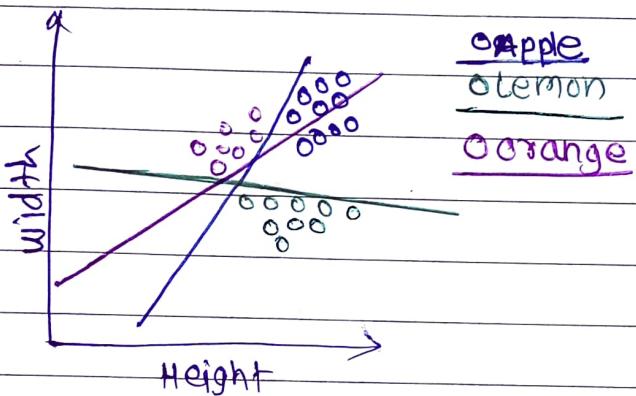
$$f(x, w, w_0) = \text{sign}(w \cdot x + w_0)$$

$$= \text{sign} \left[\sum w(i) \cdot x(i) + w_0 \right]$$

dot product

Multiclass classification

Scikit learn make easy to make multiclass classifier by creating Binary classification model for each class against all other classes.



for more features
space becomes 3D
and so on....

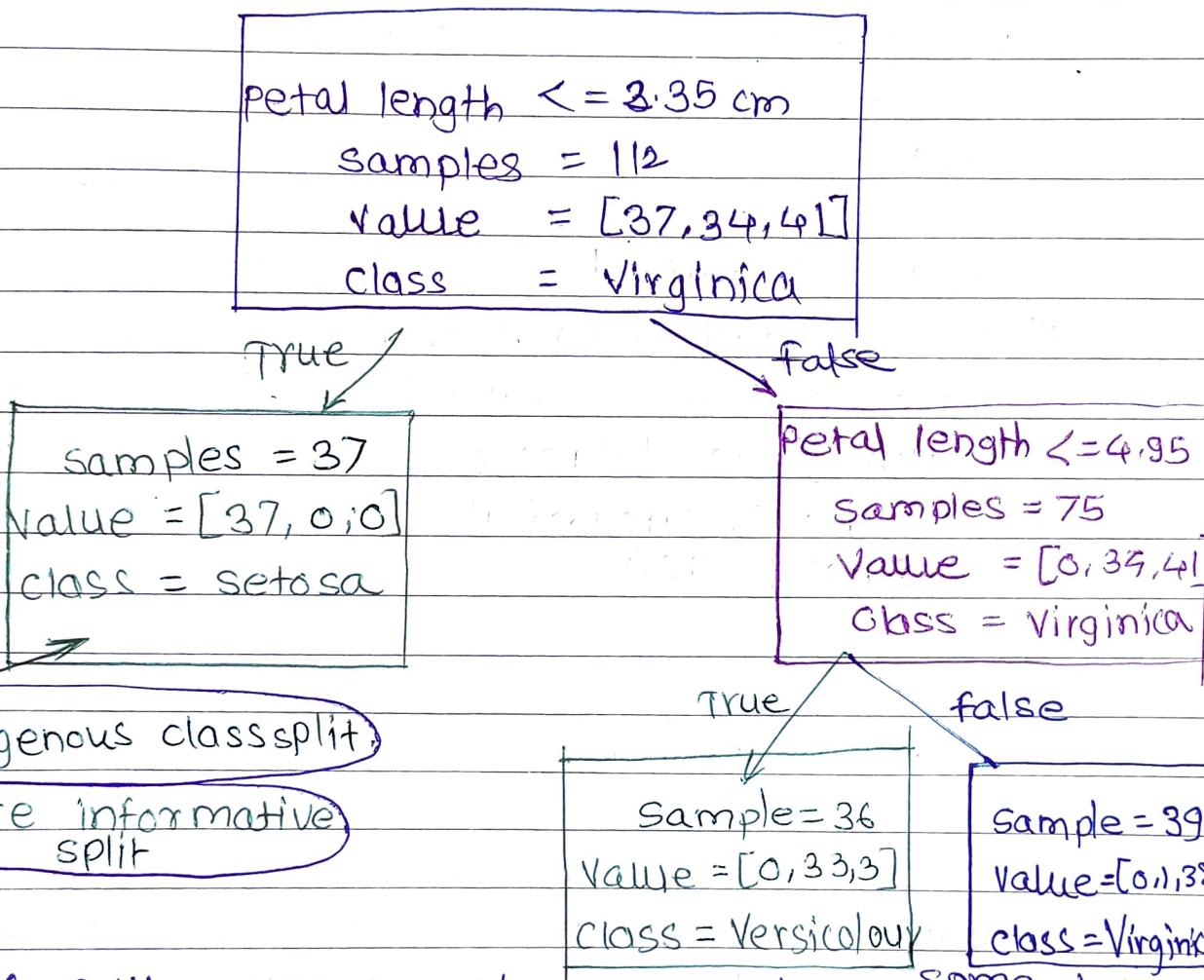
Scikit learn makes Binary classifiers for each class for e.g. - class 1 = Apple then one Binary classifier determines whether object is a apple or not apple.

Similarly all the classes works and classifies object on multiple classes with many parameters.

Decision Trees

Decision Trees - It consists of series of if-else rules which leads to target values

Informativeness of Splits - As split occurs the results must be homogenous as more as possible.



If split occurs as above e.g. then ~~some~~ Node doesn't have homogenous data so further classification cannot happen as certain nodes while further classification required at some nodes.

Decision Trees can be also overfitted which is prevented by restricting

- 1) Depth of Tree
- 2) Max leaf nodes - Total
- 3) Min sample leaf.

of data leaves to avoid splitting

Pros and cons of Decision trees.

Pros

- 1) Better Visualized
- 2) No Normalisation needed
- 3) Work well with mixed features in datasets (continuous, categorical, binary etc.)

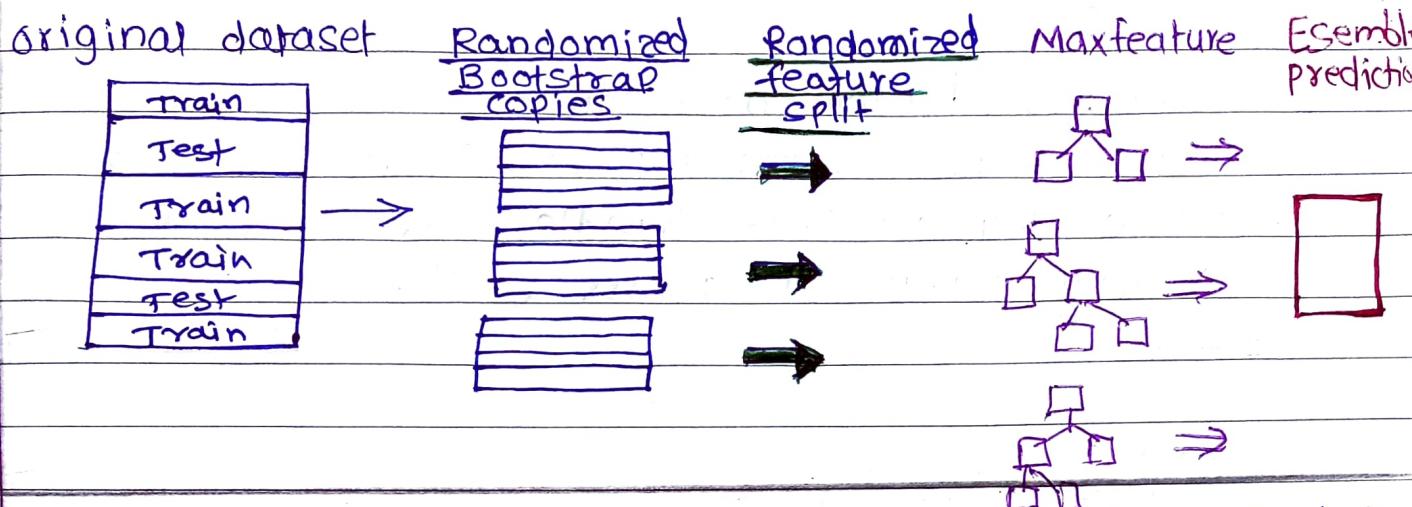
cons

- 1) Tendency to overfit even after tuning
- 2)

Random forests

It is a ensemble of many decision trees which uses Randomized bootstrap copies of data and Randomized feature splits which leads to different kinds of decision trees which are ensembled (brought together) to make predictions.

Bootstrap sample - Samples of training sets are selected randomly from original data hence each sample is slightly different than other.



1) original Dataset is Randomly splitted as training and test datasets this Various copies of original data set is called bootstrap samples.

2) Feature are extracted randomly from each bootstrap copy and

3) separate models are made for each bootstrap

4) Esembling of this trees are made.

5) Predictions are made for every tree in forest and combination of this prediction done as -

i) Regression :- mean of individual tree prediction

ii) Classification:- a) Each tree gives probability of each class
b) probabilities are averaged.

c) predict class with high probability,

max-feature parameter controls the no. of feature to be selected for each tree.

i) For High value of max-features = < no. of total feature gives similar forests with similar trees

ii) For low value of max-features = 1 trees are diverse.

PROS

- i) widely used, Excellent prediction
- ii) don't require careful Normalisation of features
- iii) like decision tree handle mixture of features.
- iv) Easily parallelised across CPU's.

CONS:

- i) models are difficult to be interpreted by Humans
- ii) may not be good for High dimensional tasks like text classifiers, compared to linear mode

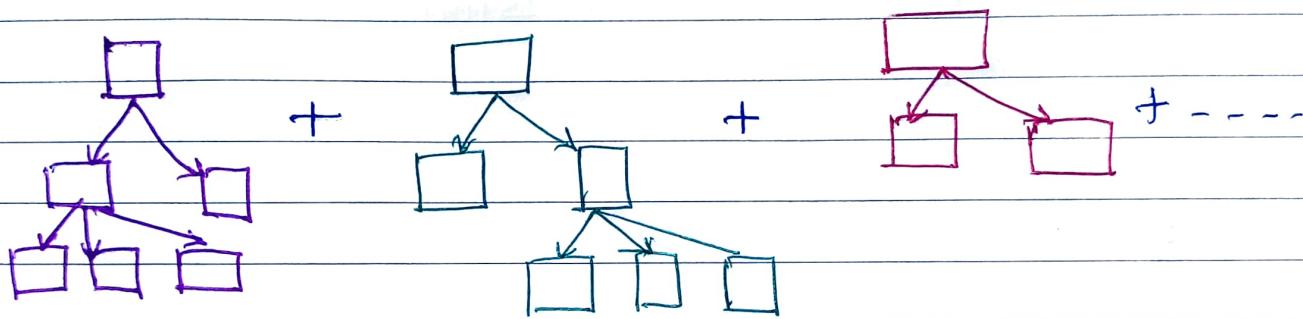
key parameters

- 1) n_estimators : Number of trees to ensemble
should be larger for larger datasets
(default = 10)
- 2) max_features : Influence on diversity of trees
default works better
- 3) max_depth :- controls depth of each tree.
(default = None)
- 4) n_jobs :- How many cores to use in parallel during training.

choose fixed settings for random_state parameter if need reproducible results.

Gradient Boosted decision Trees (GBDT)

No. of shallow trees that fixes the mistakes made by previous trees. hence with many trees in series reduces mistakes made by trees and leads to achieve higher accuracy.



Learning Rate - The parameter which controls how new tree try to force to correct the mistakes of previous tree.

High learning rate \rightarrow High complexity
 Low learning rate \rightarrow Low complexity.

* Pros and cons are similar to Random forests.

Key parameters

- 1) n_estimators - # of small trees to use which is used first always to best exploit memory and CPU.
- 2) max_depth - Typically very less of (3 to 5)
- 3) Learning_rate - controls emphasis on fixing errors

CROSS VALIDATION

Process in which we evaluate the models with different train test splits. we create different models and evaluate them with different sets K fold cross Validation (3 fold)

dataset	model 1	model 2	model 3
fold 1	Train	Train	Test
fold 2	Train	Test	Train
fold 3	Test	Train	Train

choosing different values to train and test may give different Accuracy

Hence we define different models and obtain accuracy as a Average of all folds.

Leave one out cross validation. (N sample in dataset)

Dataset	model 1	model 2	N model
fold 1 = 1 data	(Test) Train 	Train (Test) 	Train Train (Test)
fold N	Train	Train	

only one data used as Test data and all other are Training data for N data items their are N models.

Cross Validation

Process in which we evaluate the models with different train test splits. we create different models and evaluate them with different sets K fold cross Validation (3 fold)

dataset	model 1	model 2	model 3
fold 1	Train	Train	Test
fold 2	Train	Test	Train
fold 3	Test	Train	Train

choosing different values to train and test may give different Accuracy

Hence we define different models and obtain accuracy as a Average of all folds.

Leave one out cross validation. (N sample in dataset)

Dataset	model 1	model 2	N model
fold 1 = 1 data	(Test) Train Train	Train (Test) Train	Train Train (Test)
fold N	Train	Train	Train

only one data used as Test data and all other are Training data for N data items their are N models.

Unsupervised learning

NO target values, labels

Clustering

- 1) Find groups in the data -
- 2) Assign every point in dataset to one of this group -

clustering - Finding a way to divide a dataset into groups. (clusters).

- * Data points within the same cluster should be close or similar in some way;
- * Data points within the different clusters should be far or different in some away.

Hard clustering - Each point belongs exactly one cluster

Soft (fuzzy) clustering - Each point having certain weight, score or probability of membership of certain cluster.

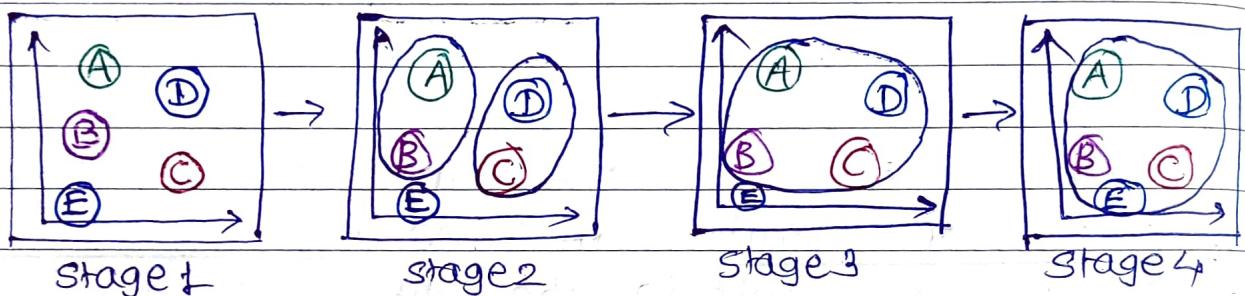
K-means clustering -

- 1) Pick a number of clusters 'k' you want to find then pick a random point for clusters.
- 2) A - Assign every point in cluster to single point selected randomly, which is nearest one.

B - Update each cluster center by replacing with mean of all points of same cluster.

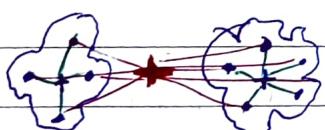
g) Repeat the (2) process until cluster converges to stable solution.

Agglomerative clustering

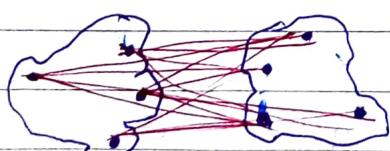


Linkage criteria - in scikit learn

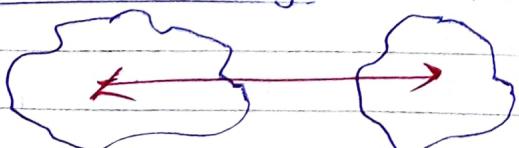
1) Ward's method - Least Increase in total variance around cluster centroids.



2) Average linkage - Average distance between clusters.



3) complete linkage - Max distance between clusters.

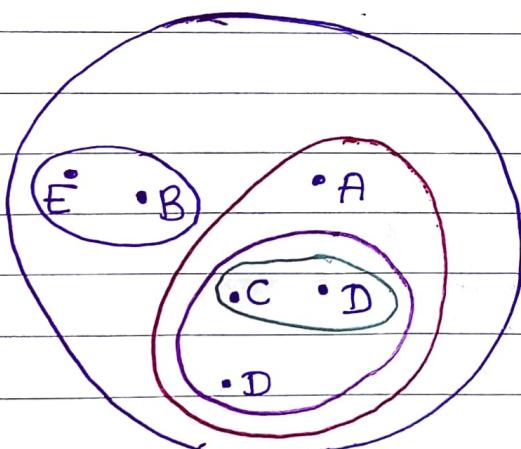


* Average linkage works well, in most of the Applications

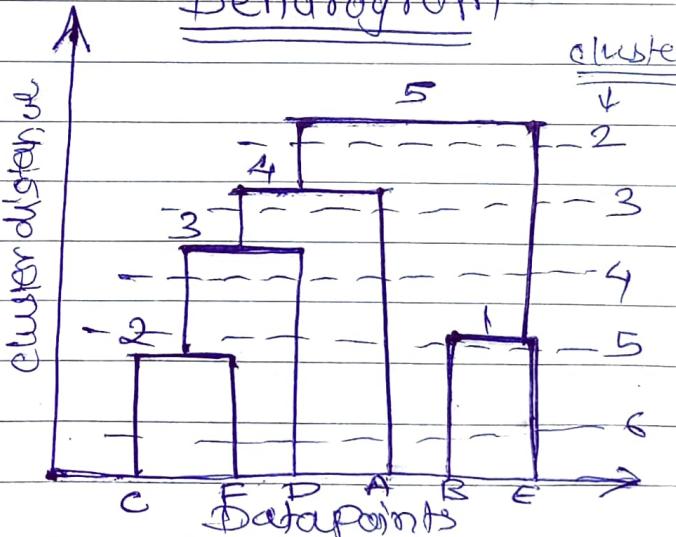
Agglomerative clustering

- 1) Initially Each instance is separate cluster then
- 2) Each cluster get linked to other clusters which are certain distances according to linkage types specified.
- 3) This process of combining clusters is repeated until no. clusters shrinks to specific number of clusters.

Hierarchical clustering



Dendrogram



- 1) One cluster is a part of/instance of another cluster.
- 2) Distance metric is made by finding distance between two points form a matrix of such distance metrics.
- 3) The two object of clusters merged together as single cluster, and now again process is repeated.

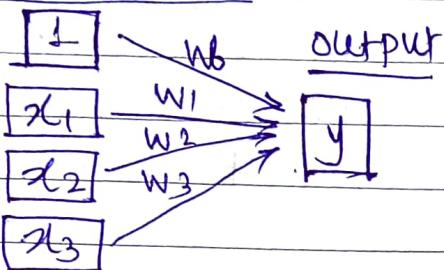
* Dendrogram plots point vs distance of points on x and y axes.

Neural Networks for regression

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:						

Linear Regression

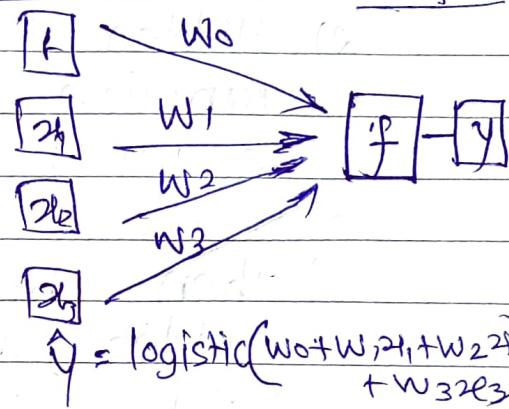
Input feature



$$\hat{y} = w_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2 + \hat{w}_3 x_3$$

logistic Regression

Input feature



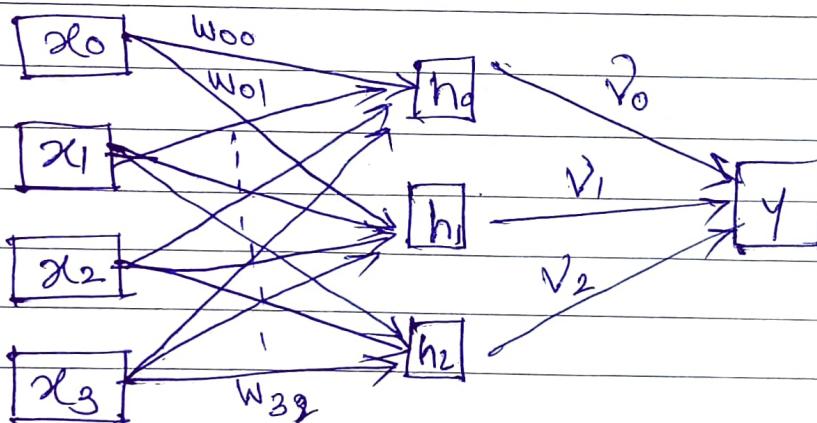
$$\hat{y} = \text{logistic}(w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3)$$

MLP) Multi layer perceptron with one Hidden layer.

Input features

Hidden layer

output



$$\hat{y} = v_0 h_0 + v_1 h_1 + v_2 h_2$$

$$h_i = \tanh(w_{0i}x_0 + w_{1i}x_1 + w_{2i}x_2 + w_{3i}x_3)$$

Activation function

* There can more than one hidden layers.

- 1) By addition of addition process layer called hidden layer in Logistic Regression represented by $\boxed{h_0} \boxed{h_1} \boxed{h_2}$ which called hidden boxes
- 2) Hidden units in hidden layers computes non linear function for each feature with weighted sum in input, resulting in output values v_0, v_1, v_2
- 3) Then the MLP computes weighted sum of hidden unit outputs to form final output value (\hat{Y})
- 4) Non linear function that hidden unit applies is called Activation functions.

Data leakage

- 1) When data we are using to train contains information about what you are trying to predict.
- 2) Introducing information about targets during training that would not legitimately be available during actual use.

E.g. - Including the label to be predicted as a feature
- Including test data with training data ↴
data to be predicted as which fruit has one ↴
feature with its actual label.

e.g. Predict if user on financial site will open an account or not? and asking account id. which could be only filled after opening account.

Model Evaluation

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:						

Represent - Train - Evaluate - Refine cycle

Representation

Extract and
select object
features

Train models

Fit the estimator
to the data



Feature and
model Refinement



Evaluation

Evaluation method are need to be choose wisely that match the goal of application.

We compute selected Evaluation metric for the multiple different models. and then select the model with best evaluation metric.

1) Accuracy with imbalanced class

Just looking at accuracy for evaluating the model is not sufficient sometimes.

e.g. suppose we have two classes

- 1) Relevant - Positive class
- 2) Irrelevant - Negative class

out of 1000 randomly selected items ,on average only 1 item is Relevant (Positive) and all other are Irrelevant (Negative) then even if we only define a dummy classifier that too gives accuracy of [99.9%]

Dummy classifier - scikit provides a dummy classifier class which produce model which don't train model on training dataset just

gives results based on majority of class as depending upon whether majority of instances belongs to Negative or positive class.

Dummy classifier is not used for actual application.

Hence Accuracy cannot be always reliable to be used as Evaluation metric of model.

$$\text{Accuracy} = \frac{\# \text{ correct Predictions}}{\# \text{ instances}}$$

If Accuracy of model is same as dummy classifier Accuracy is close to Null accuracy line.

- 1) Because of large imbalance of class
- 2) missing important features or parameters

2) confusion metrics - Accuracy is not always para we also need to look for the consequences of mistakes made by models.

e.g. model that predict Brain Tumour.

		Positive	FN (4)	TP (7)
True	Positive	T N (107)	FP (26)	
	Negative	Negative	Positive	
	Predicted			

Here even if person have tumor and model predicts it as negative that is not acceptable.

$$\text{Accuracy} = \frac{\# TP + \# TN}{\# TP + \# TN + \# FP + \# FN} = \frac{107 + 7}{107 + 7 + 26 + 4} =$$

$$\text{Error} = \frac{\# FN + \# FP}{\# FN + \# FP + \# TP + \# TN} = \frac{4+26}{4+26+107+7} =$$

$$\boxed{\text{Accuracy} = 1 - \text{Error}}$$

Precision -

$$\text{Precision} = \frac{TP}{TP+FP} = \frac{7}{7+26}$$

Recall -

$$\text{Recall} = \frac{TP}{TP+FN} = \frac{7}{7+4}$$

Recall and Precision trades are needed as high or low depending upon the application
e.g. Tumour detector require high recall.

F-score - Generalises F-1-score for combining precision and recall into a single number.

$$F_1 = \frac{2 \cdot \text{Precision} \cdot \text{recall}}{\text{Precision} + \text{recall}} = \frac{2TP}{2TP + FN + FP}$$

$$F_\beta = (1+\beta^2) \cdot \frac{\text{Precision} \cdot \text{recall}}{(\beta^2 \cdot \text{Precision}) + \text{recall}} = \frac{(1+\beta^2) \cdot TP}{(1+\beta^2)TP + \beta FN + FP}$$

β Allows users to adjustment of metric to control the emphasis on recall vs precision.

$\beta = 0.2$ - Recall oriented users

False negative hurt performance more than false positive.

e.g. Tumour detection.

$\beta = 0.5$ - Precision oriented users

False negative hurt less than false positive.

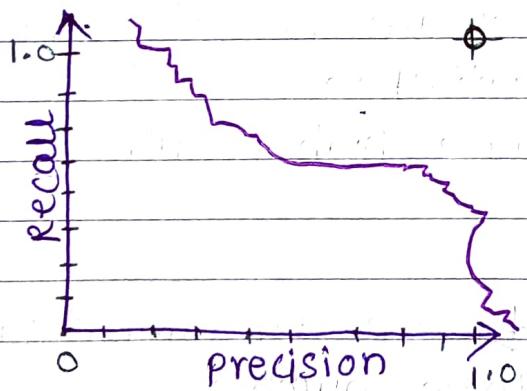
Hence apart from accuracy various parameters can be used as evaluation metric for model depending upon the goal of model.

Precision Recall curve

Ideal point - Top right corner

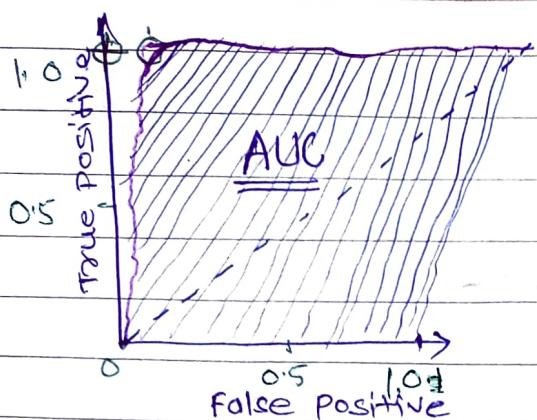
Precision = 1

Recall = 1



ROC curve - (

Ideal point - Top left corner



Ideal

False positive = 0

True positive = 1

AUC should be High

AUC = 0 Bad classifier

AUC = 1 Good classifier

Multiclass Evaluation

Regression Evaluation

Typically used R^2 score for evaluating regression.

Alternative metrics are -

- 1) mean absolute error - absolute diff. of target & pred
- 2) mean squared error - squared diff. of target & pred
- 3) median absolute error - median of error distribution

To determine whether R^2 score is good or not we take a help of dummy regression model and compare R^2 score of both models.