

ML in Python.

Machine Learning Algorithms.

Regression:

- Simple linear regression
- Multiple linear regression
- Regression trees

Classification:

- Logistic regression
- KNN
- SVM
- Multiclass prediction
- Decision trees.

Clustering

- K-means.

Machine learning is the subfield of computer science that gives "computers the ability to learn without being explicitly programmed."

- Arthur Samuel.

Major machine learning techniques:

→ Regression/Estimation:

Predicting continuous values

Eg:- CO₂ emission from a car's engine

→ Classification

Predicting the item class/category of a class

Eg:- if a cell is benign or malignant.

→ Clustering
Finding the structure of data,
Segmentation

Eg:- for customer segmentation in
the banking field

→ Associations

Associating frequent co-occurring items

for grocery items that are usually
bought together by a particular
customer

→ anomaly detection

Discovering abnormal and unusual cases

Eg:- for credit card fraud detection

→ Sequence mining

Predicting next events; click stream

(Markov model, HMM)

Eg:- for predicting the next event,
for instance the click stream
in website.

→ Dimension Reduction

Reducing the size of data CP (d)

→ Recommendation systems

Recommending items

Eg:- movies or books recommendations on
platforms

Python libraries for ml

general purpose programming language.

→ NumPy → a math library to work with N-dimensional arrays in Python.

↳ E.g.: - for working with arrays, dictionaries, functions, datatypes and working with images

→ SciPy → collection of numerical algorithms & domain specific toolboxes.

→ Matplotlib → popular plotting package that provides 2D & 3D plotting,

→ Pandas

→ Scikit learn

Data preprocessing → Train / Test split. → Algorithm setup → Model fitting

Model export ← Evaluation ← Prediction

Supervised vs. Unsupervised

Supervised learning

→ we teach the model with knowledge so that it can predict unknown or future instances.

we teach the model by training it with some data from a labelled dataset.

Two types of Supervised learning:

→ Classification

→ Regression

Classification is the process of predicting a discrete class label or category

Regression is the process of predicting a continuous value as opposed in classification

Unsupervised learning:

(unlabelled data)

The model works on its own to discover information.

It learns on datasets and draws conclusions on unlabelled data.

U.L techniques

→ Dimension reduction

→ Density Estimation

→ Market basket analysis

→ Clustering

Clustering:

- is grouping of data points or objects that are somehow similar by
 - Re's common structure
 - generalization
 - anomaly detection.

Regression:

There are 2 types of variables

→ dependent variable & independent variable

↓

State; target or final goal

(Y) should be continuous (not discrete)

↓

Explanatory variable (X).

Simple regression:

one independent variable is used to estimate a dependent variable

→ simple linear or non-linear regression

Eg:- predict CO₂ emission vs. Engine size of vehicles.

Multiple Regression

when more than one independent variable is present.

Eg:- predict CO₂ emission vs. Engine size & cylinders of all cars.

multiple linear or non-linear regression

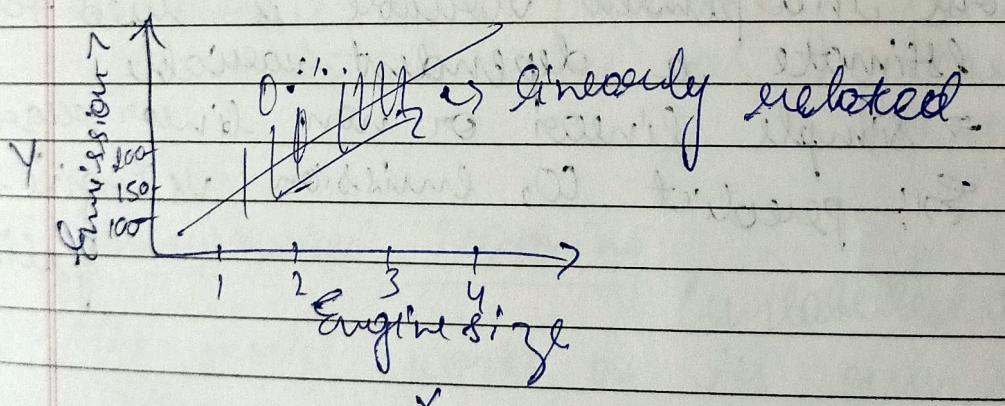
- Application of Regression
- sales forecasting
 - satisfaction analysis
 - price estimation
 - employment income.

Simple Linear Regression

Linear regression is the linear approximation of a linear model used to describe the relationship b/w 2 or more variables.

In simple linear regression → 2 variable
a dependent & an independent variable

Dependent value should be continuous
& can't be discrete value.



$$Y = \alpha_0 + \alpha_1 X_1$$

slope or gradient of fitting line

Response variable
(predicted or dependent)

Intercept

a single predictor (independent)

θ_0 & θ_1 are called the coefficients of linear eqn.

$$\Sigma x_i = x_1 = 5.4$$

$y = 250$ (actual CO_2 emission)

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$\hat{y} = 340$ (predicted emission)

$$\text{Error} = y - \hat{y}$$

$= -90$ (Residual error)

mean squared error $= \frac{1}{n} \sum (y_i - \hat{y}_i)^2$

The objective of linear regression, is to minimize the MSE eqn. We need to find the best (suitable) parameters θ_0 & θ_1 .

$$\theta_1 = \frac{\sum_{i=1}^s (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^s (x_i - \bar{x})^2}$$

$\bar{x} =$

$$\text{Engineering } \bar{x} = (2.0 + 2.4 + \dots) / 9 = 3.03$$

$$\text{CO}_2 \text{ emission } \bar{y} = (796 + 221 + 136 + \dots) / 9 = 226.22$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\theta_1 = \frac{(2.0 - 3.03)(196 - 226.22) + (2.4 - 3.03)}{(2.0 - 3.03)^2 + (2.4 - 3.03)^2}$$

$$\theta_1 = 39.$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\theta_0 = 226.22 - 39 \times 3.03$$

$$\theta_0 = 125.74$$

The true parameters for the line.

$\theta_0 \rightarrow$ bias coefficient.

$$y = 125.74 + 39x,$$

$\text{CO}_2 \text{ emission} = \theta_0 + \theta_1 \cdot \text{Enginesize}$

Posses of linear regression

- Very fast
- No parameter tuning
- Easy to understand.

Model Evaluation in Regression Models

- Train & test on the same dataset
- Train / test split.

Error:

$$\text{Error} = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

Take the CO₂ emission & clean the model & then take few parameters to fit along with their actual values(y), compare it with the predicted values(\hat{y})

$$\text{Error} = (233 - 234) + (255 - 256) + \dots$$

$$\text{Error} = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

- 1) Train & test on the same dataset
 - train the model on the entire dataset, then test it using a portion set of the same dataset.

a portion → high training accuracy & low "out of sample" accuracy of training accuracy → percentage of correct prediction the model makes when using the test dataset
(high TD → results in overfit)

Out of sample accuracy → is the percentage of correct predictions that the model makes on data that the model has not been trained on.

Train & test on same data causes low "out of sample" accuracy → overfit.

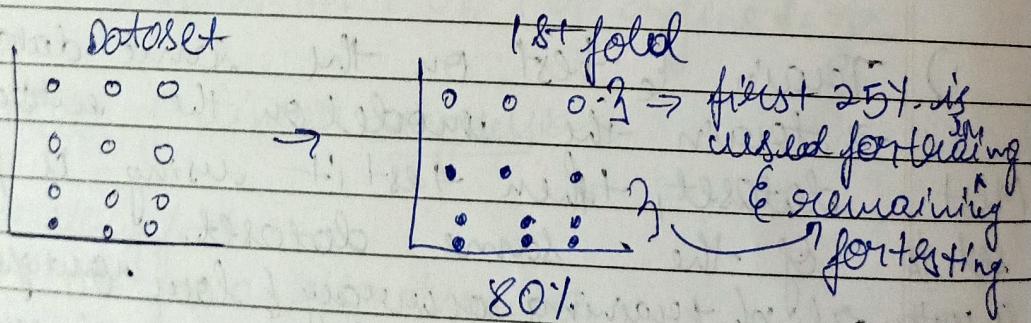
Overfitting - the model is overly trained to the dataset which may ignore noise and produce a non-generalized model.

2) Train Test split Evaluation approach,

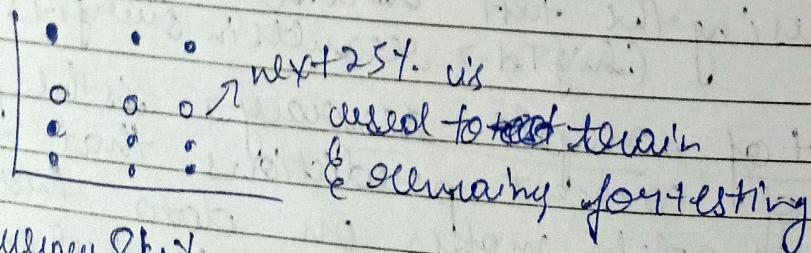
→ involves splitting the dataset into training and testing sets (separately).

- it is mutually exclusive.
- more accurate evaluation on out-of-sample accuracy
- highly dependent on which datasets the data is trained & tested.

3) K-fold Cross-validation.



and fold



accuracy 84%.

3rd fold \rightarrow 82%.

, 11th fold \rightarrow 86%.

$$\text{Accuracy} = \frac{\text{avg}}{\text{11}} \rightarrow 83\%.$$

\rightarrow k-fold cross-validation performs multiple bracketed splits using the same dataset where each split is different.

Then the result \rightarrow it's averaged to produce more consistent out-of-sample accuracy.

Evaluation metrics in Regression models

Error in regression is the difference between the data points & the trend line generated by the algorithm.

$$\text{Mean Absolute Error (MAE)} = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

$$\text{Mean Squared Error (MSE)} = \frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2$$

$$\text{Root MSE (RMSE)} = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

$$\text{Relative Absolute Error (RAE)} = \frac{\sum_{j=1}^n |y_j - \hat{y}_j|}{\sum_{j=1}^n |y_j - \bar{y}|}$$

(or residual sum of squares) mean value of y

$$\text{Relative squared error} = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{\sum_{j=1}^n (y_j - \bar{y})^2}$$

(RSE)

$$R^2 = 1 - RSE$$

represents how close the data values are to the fitted regression line.

higher $R^2 \rightarrow$ better the model fits the data.

Simple Linear Regression:-

Importing packages:-

import pipelite

smoof pipelite - install (['pandas'])

smoof pipelite - install (['matplotlib'])

smoof pipelite - install (['numpy'])

['scikit-learn']

import matplotlib.pyplot as plt

import pandas as pd

: pylab as pl

: numpy as np

: matplotlib inline

Enables drawing of figures in IPython environment

18 download data:
path = "

CSV " deals with HTTP requests.

```
from requests import get
from bs4 import BeautifulSoup
```

Design def download(url, filename):
 response = await pyfetch(url)
 if response.status == 200:
 with open(filename, "wb") as f:
 f.write(await response.bytes())

// Reading the data in.

```
await download(path, "FuelConsumption.csv")
path = "FuelConsumption.csv"
```

```
df = pd.read_csv("FuelConsumption.csv")
```

df.head()

data from in pandas
 ↗ I took a look at the dataset

// Data Exploration

```
df.describe() # to summarise the data
```

```
df = df[['ENGINESIZE', 'CYLINDERS', ...]]
```

df.head(9)

cumulative distribution function
 ↗ finds cumulative probability for given value

To visualize and plot

viz = (df[['Emissions', 'ENGINE SIZE']], df['CO2 EMISSIONS'])

viz.hist()

plt.show()

Now, plot each feature against the Emission, to see linear relationship.

plt.scatter(df['FUEL CONSUMPTION (L/100km)'],

df['CO2 EMISSIONS'], color='blue')

plt.xlabel('Fuel Consumption (L/100km)')

plt.ylabel('Emission')

plt.show()

Split dataset into train and test sets.

80% for training & 20% for testing.

We can create a mask etc select

random rows using np.random.rand() function:

msk = np.random.rand(len(df)) < 0.8

train = df[msk]

test = df[~msk]

Train data distribution

plt.scatter(train['ENGINE SIZE'], train['CO2 EMISSIONS'],
color='blue')

plt.xlabel('ENGINE size')

plt.ylabel('Emission')

plt.show()

Using sklearn package to model data

```
from sklearn import linear_model
regressor = linear_model.LinearRegression()
train_X = np.array([[train_I['Enginesize'], train_I['Cylinders'],
train_Y = np.array([train_I['CO2 emission']])
```

for converting the inputs to numpy arrays

```
regressor.fit(train_X, train_Y)
print("Coefficients: ", regressor.coef_)
print("Intercept: ", regressor.intercept_)
```

parameters of fit line.

```
predictions = regressor.predict(test_X)
```

print("Mean Absolute Error: " + str(np.mean(np.abs(predictions - test_Y))))

Multiple Linear Regression

- Independent variables affectiveness on prediction
- Predict impacts of changes.

We can predict continuous values.

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$g_j = \theta^T X$$

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots]$$

$$X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \end{bmatrix}$$

optimised parameters are the ones which lead to a model with the fewest errors.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - g_i)^2$$

(need to minimise MSE value)

→ How to estimate θ^T

→ ordinary least squares.

→ an optimisation algorithm.

→ gradient descent

Classification:

- Supervised learning approach
- Categorising some unknown items into a discrete set of categories or 'classes'
- The target attribute is a categorised variable
- Classification determines the class label for an unlabelled test case.

Classification algorithms →

- Decision Trees
- Naive Bayes
- Linear Discriminant Analysis
- K-nearest Neighbour
- Logistic Regression
- Neural Networks
- SVM → Support Vector Machines

KNN:

The K-nearest Neighbors algorithm is a classification algorithm that takes a bunch of labelled points & uses them to learn how to label other points.

- A method of classifying cases based on their similarity to other ones.
- Cases that are near each other are neighbors

- algo \rightarrow KNN
- \rightarrow pick a value for K .
 - \rightarrow calc the distance of unknown w.r.t from all w.r.t
 - \rightarrow select the K -observations in the learning data that are nearest to the unknown data point
 - \rightarrow Predict the response of the unknown data point using the most popular response value from the K -nearest neighbors.

Evaluation metrics for classification

Acuracy measurements

\rightarrow Jaccard index

y : Actual labels

\hat{y} : Predicted labels

$$J(y, \hat{y}) = \frac{|y \cap \hat{y}|}{|y \cup \hat{y}|} = \frac{|y \cap \hat{y}|}{|y| + |\hat{y}| - |y \cap \hat{y}|}$$

$$y = [0, 0, 0, 0, 0, 1, 1, 1, 1, 1]$$

$$\hat{y} = [1, 0, 0, 0, 1, 1, 1, 1, 1]$$

$$J(y, \hat{y}) = \frac{8}{10+10-8} = 0.66$$

F1-score:

Precision \rightarrow True Positive
 $(TP / (TP + FP))$

$$\text{Recall} = \frac{TP}{(TP + FN)}$$

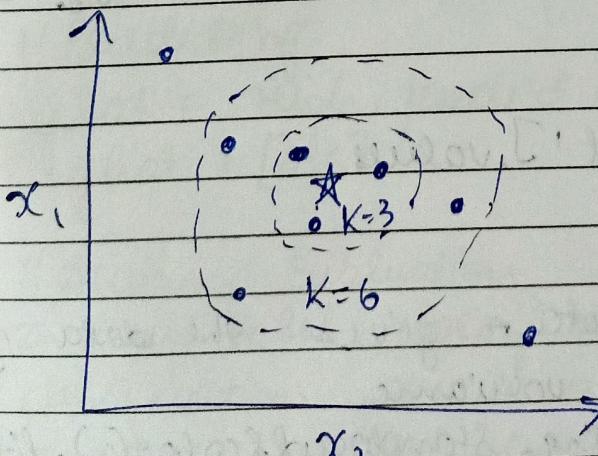
$$\text{F1-score} = \frac{2 \times (\text{precision} \times \text{recall})}{(\text{precision} + \text{recall})}$$

Both Jaccard & F1-score can be used for multi-class classifiers.

Log loss.

Performance of a classifier where the predicted output is a probability value between 0 and 1.

$$\text{Log loss} = -\frac{1}{n} \sum (y \times \log(p) + (1-y) \times \log(1-p))$$



To build a classifier, to predict unknown values (by using K nearest neighbor)

II Load data from CSV file

await download (path, 'telcust1000t.csv')
df = pd.read_csv('telcust1000t.csv')
df.head()

II Data Visualisation & analysis.

df['age'].value_counts()

O/P: 281 Plus Service, 266 -Basic service,
236 Total Service, 217 E-Services customers

df['age'][column n = 'income', limit = 50]

df['age'].columns

To use scikit-learn, we have to convert the pandas data frame to a 'Numpy array.'

X = df[['Region', 'Housing', 'Age', 'Marital',
'Address', 'Income']]
X = X[0:5]

y = df['custcat'][values]
y = y[0:5]

II Normalise data → gives the data zero mean & unit variance.

X = preprocessing.StandardScaler().fit(X).
transform(X, astype='float')
X = X[0:5]

1) Train, Test, split

from sklearn.model_selection import train_test_split

X-train, X-test, y-train, y-test = train-test-split(X, y, test_size=0.2, random_state=10)

print('Train set:', X-train.shape, y-train.shape)
print('Test set:', X-test.shape, y-test.shape)

II Classification

II K nearest neighbor(KNN)

from sklearn.neighbors import KNeighborsClassifier

II Training

K=4

neigh = KNeighborsClassifier(n_neighbors=K).fit(X-train, y-train)

neigh

II Predicting

y-hat = neigh.predict(X-test)

y-hat = [0, 5]

II Accuracy Evaluation

In multi-label classification, accuracy classification score is a function that computes subset accuracy = jaccard-score between true labels and predicted labels seen matched in the test set.

from sklearn import metrics
point ("Train set accuracy", metrics.accuracy_score(y_train, neigh.predict(X_train)),
score(y_train, neigh.predict(X_train))]
point ("Test set accuracy", metrics.accuracy_score(y_test, yhat))

To calculate the diversity of KNN for different values of K.

$$K_S = 10$$

$$\text{mean-acc} = \text{np. zeros}(K_S - 1)$$

$$\text{std-acc} = \text{np. zeros}(K_S - 1)$$

for n in range(1, K_S):

 neigh = KNeighborsClassifier(n_neighbors=n),
 fit(X_train, y_train)

 yhat = neigh.predict(X_test)

 mean-acc[n-1] = metrics.accuracy_score(y_test, yhat)

 std-acc[n-1] = np.std(yhat == y_test),
 np.sqrt(np.var(yhat, ddof=1))

mean-acc

Decision Tree

- Each internal node corresponds to a test
- Each branch corresponds to a result of the test
- Each leaf node designs a classification algo

- Choose an attribute from your dataset.
 - Calculate the significance of attribute in splitting of data
 - Split data based on the value of the best attribute
- Repeat again

Building Decision Tree!

- More Predictions
- Less Entropy
- Less Impurity
- Pure node.

Entropy → measure of randomness or uncertainty.

Lower the entropy, the less uniform the distribution, the purer the node.

$$\text{Entropy} = - p(A) \log_2(p(A)) - p(B) \log_2(p(B))$$

probability

Information gain is the information that can increase the level of certainty after splitting

Information Gain = (Entropy before split) - (weighted. entropy after split)

from sklearn.tree import DecisionTreeClassifier
import sklearn.tree as tree

// Download data

path = ' .csv'

path = "diabetes.csv"

// Read the data

my_data = pd.read_csv("diabetes.csv", delimiter=',')
my_data[0:5]

// Data processing

using my_data as糖尿病.csv

// X as the Feature Matrix

// y as the response vector.

X = my_data[['age', 'sex', 'BP', ...]].values
X[0:5]

from sklearn import preprocessing

le_sex = preprocessing.LabelEncoder()

le_sex.fit(['F', 'M'])

X[:, 1] = le_sex.transform(X[:, 1])

le_BP = preprocessing.LabelEncoder()

le_BP.fit(['LOW', 'NORMAL', 'HIGH'])

X[:, 2] = le_BP.transform(X[:, 2])

X[:, 3] = le_chol.transform(X[:, 3])

X[0:5]

$y = \text{my_data}["Pclass"]$
 $y[6:5]$

II. Setting up the DecisionTree

from sklearn.model_selection import
 train_test_split

$X_{trainset}, X_{testset}, y_{trainset}, y_{testset} =$
 $\text{train_test_split}(X, y, \text{test_size} = 0.3,$
 $\text{random_state} = 3)$

II Modelling.

deugTree = DecisionTreeClassifier(criterion="entropy",
 max_depth=4)

deugTree.

deugTree.fit(X_trainset, y_trainset)

II Prediction

predTree = deugTree.predict(X_testset)

print(predTree[0:5])

print(y_testset[0:5])

II Evaluation

from sklearn import metrics
 import matplotlib.pyplot as plt
 print("DecisionTree's Accuracy:",
 metrics.accuracy_score(y_testset, predTree))

II Visualise

tree.plot_tree(deugTree)
 plt.show()

- Logistic Regression
- is a classification algorithm for categorical variables (like 0/1 binary)
 - When is logistic regression suitable
 - if data is binary.
 - if $y \in \{0, 1\}$.
 - if you need probabilistic results.
 - When you need a linear decision boundary.
 - When we need to understand the impact of a feature.

Building a model for customer churn

$$X \in \mathbb{R}^{m \times n}$$

$$y \in \{0, 1\}$$

$$\hat{y} = P(y=1|x)$$

label / actual
vector

Sigmoid function

→ Logistic function.

$$\sigma(\Theta^T X) = \frac{1}{1 + e^{-\Theta^T X}}$$

$$\sigma(\Theta^T X) = 0$$

$$[0, 1]$$

O/P of the model:

$$\rightarrow P(y=1|x)$$

$$\rightarrow P(y=0|x) = 1 - P(y=1|x)$$

Logistic Regression

→ is a classification algorithm for categorical variables (like 0/1 binary)

- When is logistic regression suitable
- if data is binary.
 - if you need probabilistic results.
 - when you need a linear decision boundary
 - when we need to understand the impact of a feature.

Building a model for customer churn

$$X \in \mathbb{R}^{m \times n}$$

$$y \in \{0, 1\}$$

$$\hat{y} = P(y=1|x)$$

label / actual
vector

Sigmoid function

→ Logistic function

$$\sigma(\theta^T X) = \frac{1}{1 + e^{-\theta^T X}}$$

$$\sigma(\theta^T X) = 0$$

$$[0, 1]$$

O/P of the model:

$$\rightarrow P(Y=1|X)$$

$$\rightarrow P(Y=0|X) = 1 - P(Y=1|X)$$

$$\rightarrow P(\hat{y} = 0 \mid x) \rightarrow P(y=1 \mid x)$$

Training process [for customer churn example]

→ Initialise θ

$$\theta = [1, 2]$$

→ calculate $\hat{y} = \sigma(\theta^T x)$ for customer

$$\hat{y} = \sigma([1, 2] \times [1, 5]) \\ = 0.7$$

→ compare the O/p of \hat{y} with actual O/p of chrt, if & second it as error

$$\text{Error} = 1 - 0.7 = 0.3$$

→ calculate the error of all customers.

$$\text{Cost} = J(\theta)$$

→ Change the θ to

reduce the cost

θ_{new}

→ Repeat

logistic Regression Training:

$$\sigma(\theta^T x) \rightarrow P(y=1 \mid x)$$

→ change the weight \rightarrow Reduce the cost

→ cost function $\text{Cost}(f(\hat{y}, y)) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(\hat{y}_i, y_i)$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(\hat{y}_i, y_i)$$

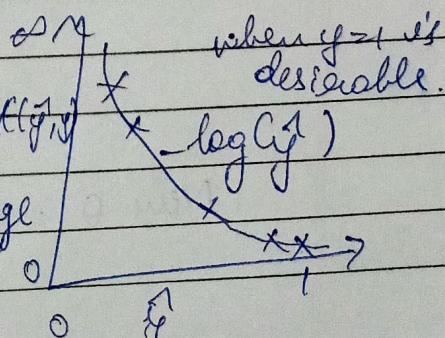
Plotting the cost function model.

model \hat{y}

actual value $y=1 or 0$

\rightarrow If $y=1, \hat{y}=1 \rightarrow \text{cost} \approx 0$

\rightarrow If $y=1, \hat{y}=0 \rightarrow \text{cost} \approx \text{large}$



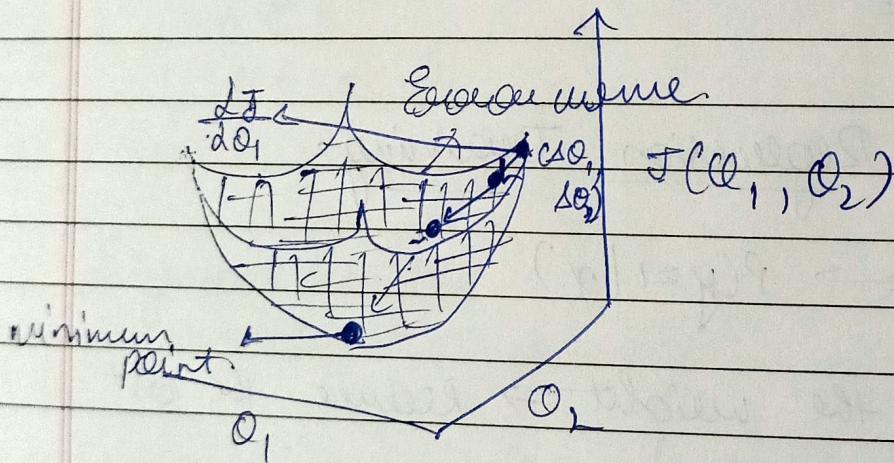
$$\text{cost}(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y=1 \\ -\log(1-\hat{y}) & \text{if } y=0 \end{cases}$$

$$\Rightarrow J(\theta) = -\frac{1}{m} \sum_{i=1}^m y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

Minimized the cost function of the model.

→ Using Gradient Descent.

A technique to use the derivative of a cost function to change the parameter values, in order to minimize the cost error.



$$\frac{\partial J}{\partial \theta_1} = -\frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i) x_i$$

$$\nabla J = \begin{pmatrix} \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_K} \end{pmatrix} \rightarrow \text{gradient vector.}$$

$$\text{New } \theta = \text{old } \theta - \eta \nabla J$$

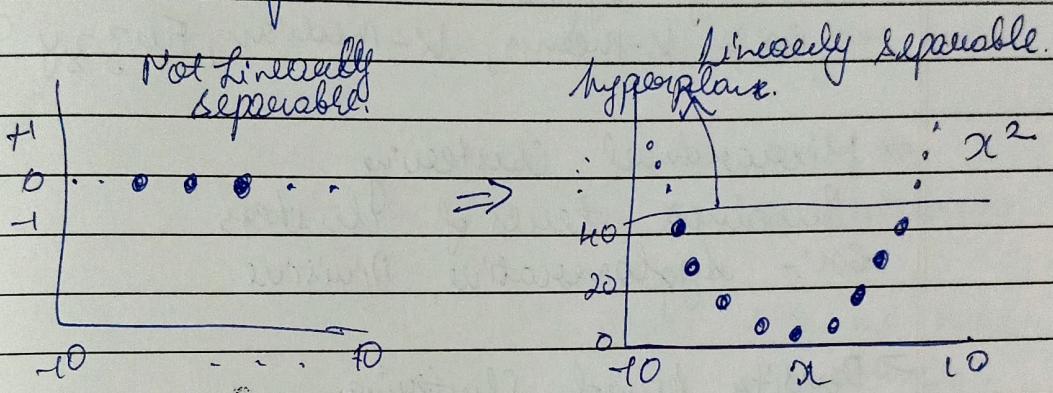
Support Vector Machines

SVM is a supervised algorithm that classifies cases by finding a separator.

- first maps data to a high dimensional feature space
- second, finding a separator.

The SVM algo finds an optimal hyperplane that categorises new examples.

Data transformation.



$$f(x) = [x, x^2]$$

Mapping data to higher dimensional space → is Kernelizing

- linear
- poly
- RBF
- sigmoid.

- they are accurate in high-dim spaces
- memory efficient
- no probability estimation.

Application → image recognition, text assignment, Detecting spam, sentiment analysis, outlier detection, regression

Clustering

→ finding clusters in a dataset; unsupervised

cluster → a group of objects that are similar
to other objects in the cluster & dissimilar
to data points in other clusters.

Clustering Algorithms

→ Partitioned based Clustering

Relatively efficient

Ex:- K-means, K-median, Fuzzy C-means

→ Hierarchical clustering

Produces trees of clusters

Ex:- Agglomerative, Divisive.

→ Density based clustering

Produces arbitrarily shaped clusters

Ex:- DBSCAN.

K-means clustering

- intra cluster distances are minimised
- inter cluster distances are maximised

$$\text{Dis}(x_1, x_2) = \sqrt{\sum_{i=0}^h (x_{1i} - x_{2i})^2}$$

(Minkowski distance)

- initialize K ($K=3$) centroids randomly
- distance calculation \Rightarrow (distance matrix is formed)
- assign each point to the closest centroid.

CSE: the sum of the squared differences between each point & its centroid

$$SSE = \sum_i^n (x_i - c_j)^2$$

→ compute the new centroids for each cluster.

(that is, each cluster center will be updated to be the mean for data points in its cluster).

→ Repeat the steps until there are no more changes.

K-means Accuracy.

- External approach
compare the clusters with the ground truth, if it is available.
- Internal approach
measure the distance b/w data points within a cluster.

Elbow method is used for finding the right K !