

Missing Data

① Importing libraries.

* Python :

```
import numpy as np, (mathematical tools - mathematics)
import matplotlib.pyplot as plt
import pandas as pd
    ↗ plot charts in python .
    ↗ import & manage datasets .
```

* R :

We don't need to import libraries. By default
Some usual libraries are imported .

② Importing Dataset

* Python :

- Goto file explorer → select machinelearningA →
→ selection 2-part 1 → dataset.csv(location).
- Now save your python file in that location .
- If will become your home directory .
- Now click "Run" and it will become your working directory .

```
import numpy as np,
import matplotlib.pyplot as plt-
import pandas as pd
# Importing dataset
dataset = pd.read_csv('name.csv')
```

[ctrl + enter - to execute] [Format : '...of']

- * In python you need to distinguish dependent and independent variables.
- * In given dataset :

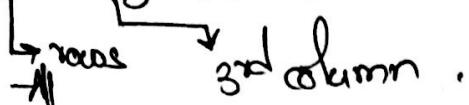
Country, Age, Salary - Independent Variables.
Purchased - Dependent Variable.

~~Importing dataset -~~

```
dataset = pd.read_csv('name.csv')
```

```
X = dataset.iloc[:, :-1].values.
```

```
y = dataset.iloc[:, 3].values.
```


→ rows → 3rd column .

* * R:

~~Importing~~

- * In R → set the directory as working directory,
Same steps as python except instead of 'Run'
click more → Set as working directory.

- * Here we don't need to distinguish variables.

~~Importing dataset~~ .

```
dataset = read.csv('name.csv')
```

(Independent variables → Matrix of features).

Mising Values Dealing :

~~* Python :~~

~~Fig.~~ Importing library that takes care of missing state .

```
from sklearn.preprocessing import Imputer,
```

↳ Syikit learn ~~class~~ → preprocessor library → import class.
Contains libraries → Contains class / methods
object of ~~class~~ → deals with missing data.

Create object of inputer class. [Ctrl + I — for help].

'imputer' = Imputer(missing_values='NaN', strategy='mean',
 axis=0)

↳ impute column, $1 \rightarrow \text{row}$

Fit the values in the matrix. impule column, 1

impules = impulser. \$it (x[, 1:3])

Upper bound excluded.

Now you need to transform the dataset with
impute values .

```
x[, 1:3] = imputed_transform(x[, 1:3])
```

KR:

• Taking care of missing date.

`dataset$Age = ifelse(is.na(dataset$Age), ave(dataset$Age, FUN = function(x) mean(x, na.rm=TRUE)),`

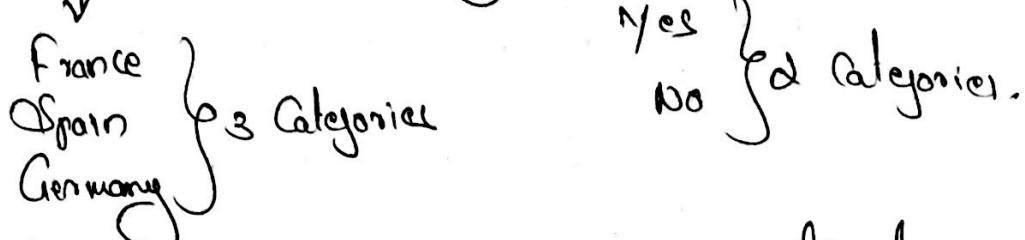
Detecting Age column in dataset . \downarrow dataset #Age)

`ifelse(condition, if true, else)`

Categorical Data

In the given dataset there are 2 categorical variables.

Country → age salary purchased .



In machine learning since we do mathematical calculations, to keep text is a problem. So we need to encode them into numbers, so here Country & Purchased are encoded into numbers.

* Python :

importing the library that encodes .

from sklearn.preprocessing import LabelEncoder
↳ class .

labelencoder - $x = \text{LabelEncoder}()$

$x[:, 0] = \text{labelencoder} - x.\text{fit_transform}(x[:, 0])$

France - 0

Spain - 2

Germany - 1

* Here we have a problem. Since Machine Learning is based on equations. It will think that Germany has higher value than France because (1 > 0). But that is not right because those are categories .

To prevent the system from thinking like that we use dummy encoding.
To do that encoding we use OneHotEncoder class from sklearn.preprocessing import LabelEncoder.

One Hot Encoder
onehotencoder = OneHotEncoder(categorical_features=[0])
 $X = \text{onehotencoder}.fit_transform(X).toarray()$

* R:

dataset\$country = factor(dataset\$country,
levels = c('France', 'Spain', 'Germany'),
labels = c(1, 2, 3))

Splitting the data set into TrainingSet
+ TestSet

& from sklearn.cross-validation import
train-test-split.

x-train, x-test, y-train, y-test = train-test-split(
x, y, test_size=0.2, random_state=①)
(for above results as video).

* Python

R:

=
install.packages('caTools') → install caTools package.
library(caTools) → to select caTools library

* In python we use RandomState to get same results,
here in R we use seed()

set.seed(123)

split = sample.split (dataset & purchased, SplitRatio=0.8)
ctr+1

→
but in python we do it for test set.

* In python we split train-set & test-set for dependable and non-dependable variables separately.

* But in R we only do it for dependable variables.

train.set = subset (dataset, split = TRUE)

test.set = subset (dataset, split = FALSE)

Feature Scaling:

* In the given dataset - Age ranges from 0-50 and Salary from 0-100k.
∴ Age & Salary are not in same scale.

When calculating Euclidean distance Age doesn't affect because it has small values, so machine learning will neglect that.

So to avoid that we do ~~feature scaling~~,
so that both are in same scale.

→ We do feature scaling in two ways

(i) Standardisation — $x_{\text{stand}} = \frac{x - \text{mean}(x)}{\text{Standard deviation}(x)}$

(ii) Normalization

$x_{\text{Norm}} = \frac{x - \text{min}(x)}{\text{max}(x) - \text{min}(x)}$

Python:

From sklearn.preprocessing import StandardScaler

sc-x = StandardScaler(),

x-train = sc-x.fit_transform(x-train)

x-test = sc-x.transform(x-test)

R:

In R, to do scaling.

training-set = scale(training-set)

test-set = scale(test-set)

But by running this we get error i.e.
x is not numeric.

But when we see training-set we see that
all are 'in numeric'.

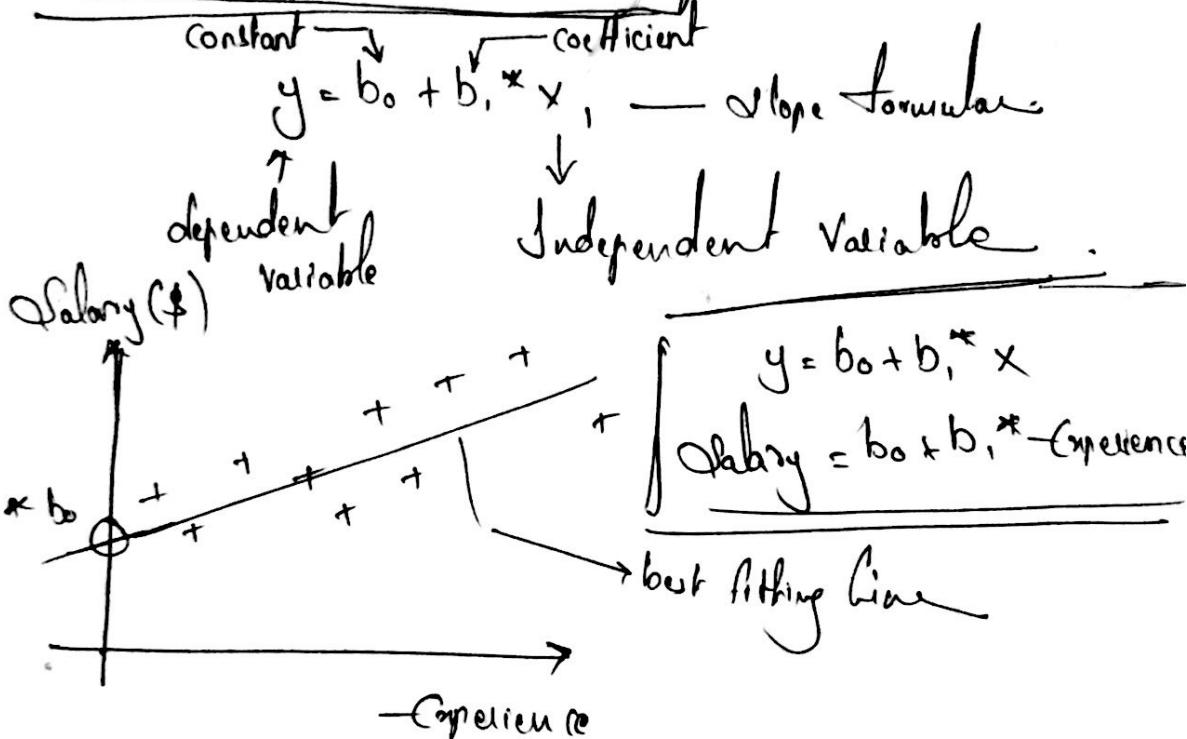
The problem we have done factoring in R,
and factoring means not changing in numeric.

So here we leave the two columns that we
factored and scale the remaining one.

Data Preprocessing template :

ection - 2

Regression - Simple - Linear :



b_0 - means point where line crosses vertical axis.

b_1 - Slope of Line

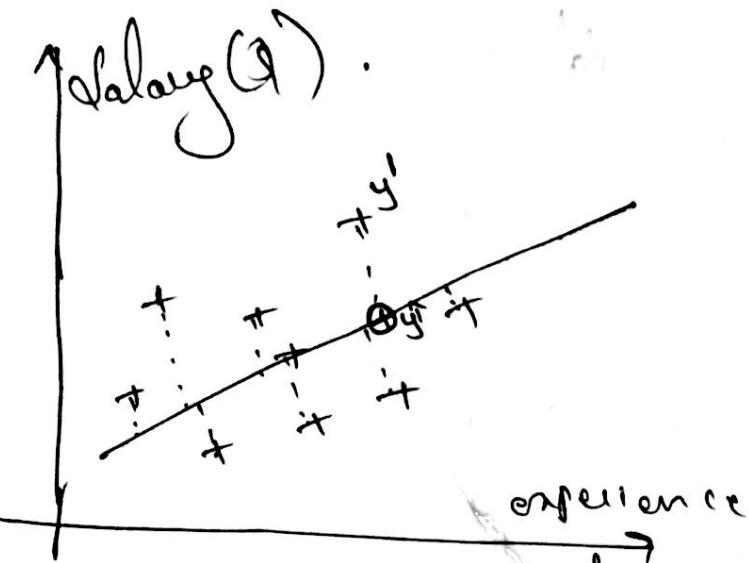
If slope is more, then for more experience more salary.

→ Simple linear Regression is used to find out the best fitting line.

- → Ways of formulae

$$\begin{cases} y = b_0 + b_1 x \\ y = b_1 x \end{cases}$$





The best fitting line is $\underset{\min}{\text{Sum of Squared differences between } (y - y^{\hat{}})}$

$$\text{i.e. } \sum (y - y^{\hat{}})^2 \rightarrow \min.$$

Python :

```
import numpy as np
```

```
np.set_printoptions(threshold=np.nan)
```

```
import matplotlib.pyplot as plt
```

```
import pandas as pd
```

importing the dataset

```
dataset = pd.read_csv('Salary-Data.csv')
```

```
x = dataset.iloc[:, :-1].values
```

```
y = dataset.iloc[:, 1].values
```

splitting into test & training sets .

- from sklearn.cross_validation import train-test-split

```
x_train, x_test, y_train, y_test = train-test-split
```

$$(x, y, \text{test_size} = 1/3, \text{random_state} = 0)$$

* In Simple - Linear Regression, we don't need to do the feature scaling, the library ~~automatically~~ does take care of it automatically.

fitting Simple Linear Regression in train set .

from sklearn.linear_model import LinearRegression

regressor = LinearRegression()

regressor.fit(x-train, y-train)

predicting the test set results .

y-pred = regressor.predict(x-test)

Visualizing the training set results .

To plot the graphs we use "matplotlib" library .

plt.scatter(x-train, y-train, color='red')

plt.plot(x-train, regressor.predict(x-train), color='blue')

plt.title('Salary vs Experience (Training set)')

plt.xlabel('Years of Experience')

plt.ylabel('Salary')

plt.show()

* We don't change x-train in plt.plot of test set because we want to predict x-test values on the regressor line formed on training set

* If we change to x-test it will form a new regressor line

```
# visualising the test set results  
plt.scatter(x-test, y-test, color='red')  
plt.plot(x-train, regressor.predict(x-train), color='blue')  
plt.title('Salary vs Experience')  
plt.xlabel('Experience')  
plt.ylabel('Salary')  
plt.show()
```

* * Q :

Salary - Dependent Variable

- Experience - Independent Variable

importing dataset

```
dataset = read.csv('Salary-Data')
```

Splitting Data & enabling ca-tools

```
library('caTools')
```

```
set.seed(123)
```

```
split = sample.split(dataset$Salary, SplitRatio = 2/3)
```

training-set = subset(split, split == TRUE)

test-set = subset(split, split == FALSE)

#1 Fitting Simple Linear Regression to the Training Set .

regressor = lm(
↓
formula = Salary ~ YearsExperience,
fit-help data = training-set)

Console : > summary(regressor)
↳ gives info .

predicting the test results .

y-pred = predict(regressor, newdata = test-set)
↳ vector of values

visualising the training set results .

ggplot() +

geom-point(aes(x = training-set \$ YearsExp ,
y = " " \$ Salary), colour = 'red')
+

geom-line(aes(x = training-set \$ YearsExperience ,
y = predict(regressor, newdata = training-set),
colour = 'blue')) +

ggtitle('Salary vs Exp ') +

xlab('Year of Exp ') ,

ylab('Salary ')

↳ ylabel .

Multiple Linear Regression:

Simple linear regression: $y = b_0 + b_1 x$

Multiple linear regression: $y = b_0 + b_1 x_1 + \dots + b_n x_n$

↓

Dependent Variable (DV)

Independent Variables (IV)

When Should we consider Linear Regression

1. Linearity
 2. Homoscedasticity
 3. Multivariate normality
 4. Independence of errors
 5. Lack of multicollinearity.

Given Dateach

??? — because `astar` is not a dollar value. It
is a string and string cannot be inserted
into an equation,

Since State is a string column we expand our dataset like
dataset (conts)

	State
NY	
CA	
CA	
NY	

Dummy Variables.

D ₁	D ₂
1	0
0	1
0	1
1	0

In NY column for NY - 1, CA - 0

In CA column for NY - 0, CA - 1

Now eq becomes

$$y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + b_4 D_1 \rightarrow \textcircled{1}$$

We shouldn't take all the dummy variables.

Dummy Variable Trap:

If we include the last dummy variable D₂ also in the eq $\textcircled{1}$, the problem is $D_L = 1 - D_1$.

So basically if you include D₂, you are duplicating D₁ again. Then the above equation is used for prediction.

So if you duplicate the dummy variable, then your equation cannot predict properly.

This is called Dummy Variable Trap.

At ~~the~~ any time, $b_0 + b_4 D_1, b_3 D_L$ all three cannot exist at same time.

** $b_0 + b_5 D_2$ is possible.

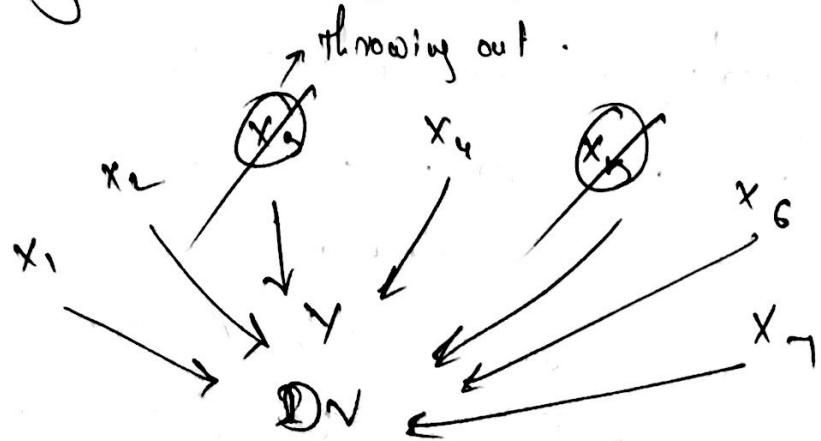
lets say if we have a dummy variable then only include 8.

if we have n dummy variables we should only include (n-1) dummy variables.

*** If we have two sets of dummy variables we should apply (n-1) dummy variables for each set.

*** The phenomenon where one (or) several independent variables predict another in a linear regression is called multi-collinearity.

→ Now in multiple linear regression, We have several I.V's unlike Simple Linear Regress where it only has one D.V.



We will remove unnecessary I.V's while predicting.

- 1) If we remove lot of 2^n's then our model will not predict correctly.
- 2) If we have a thousand 2^n's then we should consider only those 2^n's whose behaviour change reflects on DV.

How To Construct A Model

5 Methods of Building Models:

1. All-in
2. Backward elimination
3. Forward Selection
4. Bidirectional elimination
5. Step-wise regression

i) All-in

(i) Prior knowledge

(ii) You have to

If you know what variables to use (or) company asks you to include all the selected variables then you choose this method.

(iii) If you prepare for backward elimination

(2) Backward Elimination:

Step 1: Select a significance level to stay in the model (eg: $\alpha_L = 0.05$).

Step 2: Fit the full model with all possible predictors. Include all the β 's.

Step 3: Consider the predictor with the highest P-value. If $P > \alpha_L$, go to step 4, otherwise go to FINISH.

Step 4: Remove the predictor.

Step 5: fit model without this variable *

Step 6: Repeat the process until the highest P-value $< \alpha_L$. Then go to finish.

(FIN: Model is ready.)

(3) Forward Selection:

Step 1: Select a significance level to enter the model (eg: $\alpha_L = 0.05$).

Step 2: Fit all simple regression model $y - x_i$.

Select the one with lowest P-value.

Basically here we construct models with only

one DV & DV and calculate p value
for each model and the lowest p value
model is selected.

Step 3: Now selecting the lowest p-value model.
Add another DV to this model and again
calculate p-values for all the models.

Step 4: If the p-value for above model of
2 DV's & 1 DV is less than SL, then
repeat Step 3. Otherwise FIN.

i.e. The model will only be finished when
 $p > SL$

Step 5: If you got, lets say $P(4DV's \text{ or } 1DV) > SL$.
then consider the one you have selected with
 $P(3DV + 1DV) < SL$. ~~Select~~ That's the
final model.

(4) Bi-Directional Elimination:

Step 1: Select a significance level to enter and to stay in the model.

- Eg: $SL_{enter} = 0.05$ or $SL_{stay} = 0.05$

Step 2: perform the next step of forward selection (new variables must have $p < SL_{enter}$ to enter).

Step 3: perform all steps of Backward Elimination (old variables must have $p < SL_{STAY}$ to stay).

You have to iterate from Step 2 to Step 3 & to 2 until you have no variables either to add in Step 2 (or) to remove in Step 3.

Now proceed to Step 4.

Step 4: No new variables can enter and no old variables can exit.



FIN: Your model is ready.

(5) All possible model (Score Comparison)

Step 1: Select a criterion of goodness of fit.

Step 2: Construct all possible regression model. $2^n - 1$.
Total combinations.

Step 3: Select the one with best criterion.



FIN: Your model is ready.

\rightarrow : 10 columns means 1023 models.

In this tutorial we consider Backward
elimination, As it is the fastest one.

Multiple Linear Regression In Python:

~~for~~ here since the given data has a string variable, we encode it.

~~#~~ encoding the data

```
from sklearn.preprocessing import LabelEncoder, OneHotEncoder
```

```
labelencoder_x = LabelEncoder()
```

```
x[:, 3] = labelencoder_x.fit_transform(x[:, 3])
```

```
onehotencoder = OneHotEncoder(categorical_features=[3])
```

```
x = onehotencoder.fit_transform(x).toarray()
```

We also have dummy variable in the above created dataset. We need to avoid the dummy variable trap.

~~#~~ Avoiding the dummy Variable Trap

```
# x = x[:, 1:]
```

Here the first i.e 0th index dummy variable is removed.

~~#~~ Splitting the data

:

fitting multiple linear regression into the - training set .

```
from sklearn.linear_model import LinearRegression  
regression = LinearRegression()
```

regressor.fit(x-train, y-train)

predicting the test results.

y-pred = regressor.predict(x-test)

Finding the highly significant SV in the Given dataset and removing the non-significant ones.

* Here we are going to use Backward Elimination.

Building the optimal model using backward elimination
import statsmodels.formula.api as sm

Usually the above library doesn't take care of b_0 (or) neglects it.

$$y = b_0 + b_1 x_1 + \dots + b_n x_n$$

So we assume there x_0 , where $x_0 = 1$.

So we need to add a column of 1's at 0th index.

* $x = np.append(arr=x, values=np.ones(50).astype(int))$,
 $axis=1$)

The above statement appends array of fifty 1's to last position of x . But we need to append it at first position. So, the trick is

$x = m \cdot \text{append}(\text{arr} = np.ones((50, 1)).\text{astype}(int),$
values = x , axis = 1
either column (or) row = 0

$x_opt = x[:, [0, 1, 2, 3, 4, 5]]$

regressor_OLS = sm.OLS(endog = y, exog = x_opt).fit
OLS dots

regressor_OLS.summary()

 $x_opt = x[:, [0, 1, 3, 4, 5]]$

regressor_OLS = sm.OLS(endog = y, exog = x_opt).fit

regressor_OLS.Summary()

repeat until $p < 0.05$ (SL).

R:

Importing the dataset

dataset = read.csv('30_Startups.csv')

Encoding the state (or) factoring

dataset\$State = factor(dataset\$State,
levels = c('New York', 'California',
'Florida'))
labels = c(1, 2, 3))

fitting multiple linear regression to training set.

regressor = lm(formula = profit ~ R.D.Spend +
Administration + Marketing.Spend + State)
(or)

regressor = lm(formula = profit ~ .,
data = training_set),

* * * All the remaining DV's.

*** In R, if in our dataset we have spaces,
then in R we replace them with `.'

Ex: R.D. Department

*** R automatically takes care of Dummy
Variable + n.

> Summary (regressor)

P - values less than 5% → more significant
P - values greater than 5% → less significant.

Significant codes of P :

P - less b/w - stars

P - 0 - 0.001 - ***

P - 0.001 - 0.01 - **

P - 0.01 - 0.05 - *

P - 0.05 - 0.1 - .

P - 0.1 - 1 - nothing[†]

So, in the above given dataset, the only
P.V significant is R.D.Development.

So, we can rewrite the formula as
~~Formula~~ Profit ~ R.D.Development,

It still gives the same prediction.

Building Optimal Model in R:

* Building the optimal model using backward elimination.

```
regressor = lm(formula = Profit ~ R.D.Spend +  
                Administn + Marketing.Spend + State,  
                data = dataset )
```

✓

You can also take training but usually
recommend to take whole dataset.

```
summary(regressor)
```

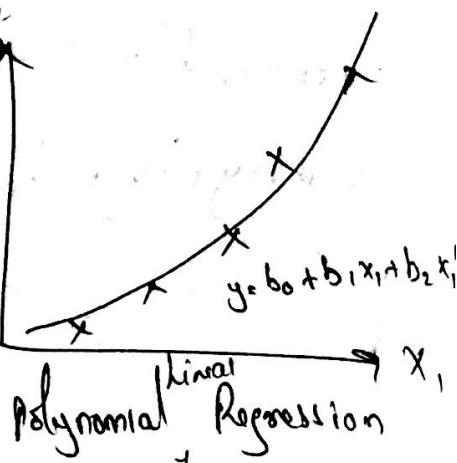
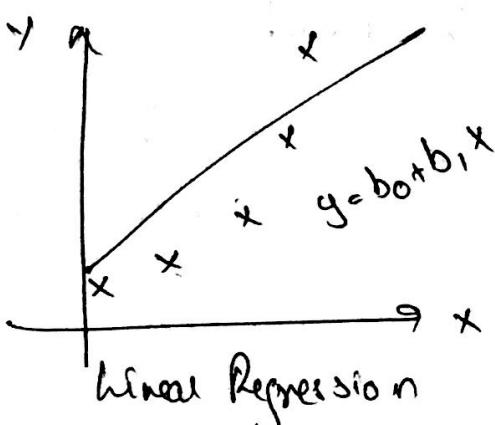
Polynomial Regression :

Intuition :

Simple Linear Reg : $y = b_0 + b_1 x$,

Multiple Linear Reg : $y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$

Polynomial Reg : $y = b_0 + b_1 x_1 + b_2 x_2^2 + \dots + b_n x_n^n$



← here linear doesn't meant x_1, x_2, \dots, x_n . But here linear means coefficients.

Linear : $y = b_0 + b_1 x_1 + \dots + b_n x_n$

Non-linear : $y = \frac{b_0 + b_1 x_1}{b_2 x_2 + b_3 x_3}$

Polynomial Linear Regression is a Special case of
Multiple Linear Regression.

- In machine learning, if you mention only the index like $[:, 1]$, then it will create a vector.
- But our matrix of features should always be in matrices not in vector.
So, $X = \text{dataset}.iloc[:, 1:2].values()$
 \downarrow
upper bound excluded.

- *** If we have very less dataset, then we don't split our dataset into training & test-set.
We will take the entire dataset and train the machine since data is less.
- *** Polynomial regression automatically takes care of feature scaling.

- Here in this tutorial we build linear regression & polynomial regression to compare them.
- Polynomial regression automatically adds column of 1's like we did in Multiple Linear Regression.
These 1's are used for b_0 .

fitting polynomial regression to the training-set
 from sklearn.preprocessing import PolynomialFeatures
 $\text{poly_reg} = \text{PolynomialFeatures(degree=2)}$ creates a col of squared values.
 $\text{X_poly} = \text{poly_reg.fit_transform}(X)$

`lin-reg-2 = LinearRegression()`
`lin-reg-2.fit(transform(x-poly), y)`

Here `x-poly` contains squared values i.e
Polynomial values.

To obtain those polynomial features we use,
Polynomial features library.

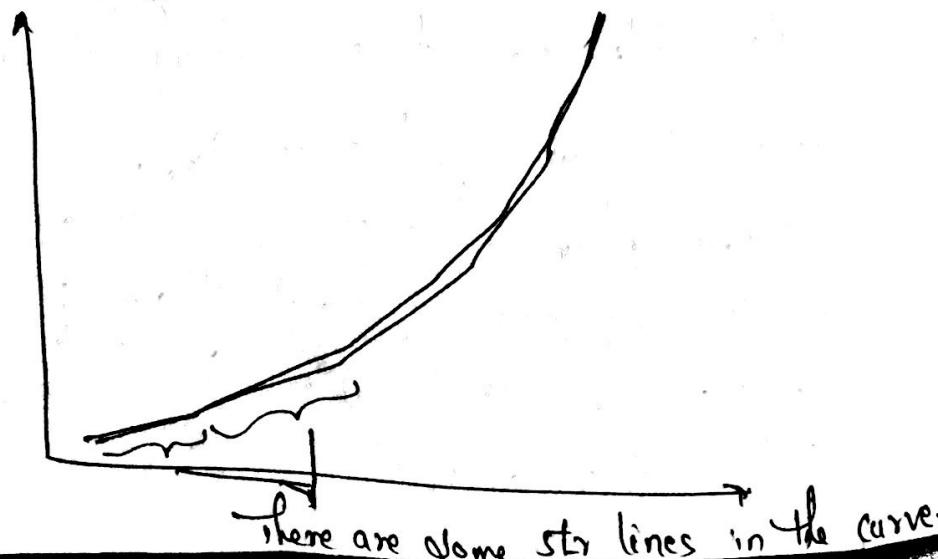
Visualizing the polynomial regression results.

`plt.scatter(x, y, color = 'red')`

`plt.plot(x, lin-reg-2.predict(poly-reg-fit.transform(x)),
color = 'blue')`

Here instead of writing `x-poly`, we wrote the
whole `poly-reg-fit-transform()` because, we
cannot write `x-poly` if its `test-set`, i.e `x-test`.
!

Now the graph will be like



There are some straight lines in the curve

To even reduce those straight lines

We use

for visualising the polynomial reg

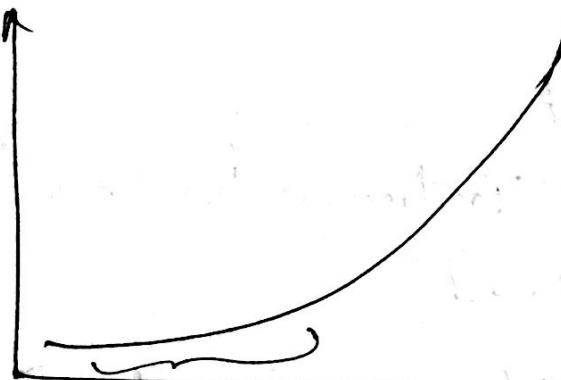
x-grid = np.arange(min(x), max(x), 0.1)

x-grid = X-grid.reshape((len(x-grid), 1)) } We use
plt.scatter(x, y, color = 'red') } these 2
lines.

plt.plot(x-grid, lin-reg2.predict(poly-reg.fit_transform(x-grid)),
, color = 'blue')

plt.show()

Now the graph will be like,



~~Now no straight lines~~

* X-grid contains 90 rows ranging from
1, 1.1, 1.2 . . . 89.9, 90.

* * * To increase the prediction we move the
degree from 2 to 3 .

degree = 3 . (More accuracy) .

* * * To increase even more
degree = 4 (Even more accuracy. fit eye)

*** On Degree = to let predict the results.

$$(\text{lin-reg } (\cancel{\text{poly}})) \cdot \text{predict}(6.5) = 3,30,378$$

far away ←

$$(\text{lin-reg-2} \cdot \text{predict}(\text{poly-reg-fit-transform}(6.5))) \\ = 1,58,862$$

which is very near to 168/k & the
predicted is true.

R:

~~Import~~ the dataset

```
dataset = read.csv('Position_Salaries.csv')
```

dataset = dataset[2:3]

* Feature scaling is taken care automatically.

To add polynomial features to our matrix of features

Fitting polynomial features

dataset \$Level2 = dataset \$Level^2 (Adding squared values as Level)
" \$Level3 = " " " " ^3

8 : 16 : 1

Poly-reg = lm (formula = salary ~ ., data = dataset)

Eff Visualising the polynomial regression

library(ggplot2)

ggplot() +

geom_point(aes(x = dataset \$ level, y = dataset \$ salary),
color = 'red') +

geom_line(aes(x = dataset \$ level, y = predict(poly-reg,
newdata = dataset), color = 'blue')) +

ggtitle()
) +

ggxlab()
) +

ylab()

* * * Predicting the values :

In R, you cannot simply write

y-pred = predict(lin-reg, $\downarrow y$) X

The above one is false

According to Syntax, in y place there should
be a dataset.

So, you need to create a dataset with
single value. To do that we take data.frame

y-pred = predict(lin-reg, data.frame(level = 6.5))

\downarrow
creating a column

→ Predicting linear regression .

predicting new result with polynomial regression

y-pred = predict(poly-reg, data.frame(lev1 = 6.5,
lev2 = 6.5¹2,
lev3 = 6.5¹3,
lev4 = 6.5¹4))

We should take level 2, 3, 4 because our
poly-reg is trained with dataset having level 2, 3, 4

2nd Non-Linear Model SVR

SVR In Python

Support Vector Regression,

In SVR, We need to do feature scaling.

The library in python doesn't care of itself.

Feature scaling

```
from sklearn.preprocessing import StandardScaler
```

```
sc_x = StandardScaler()
```

```
sc_y = StandardScaler()
```

```
x = sc_x.fit_transform(x)
```

```
y = sc_y.fit_transform(y)
```

fitting SVR to our regression.

~~regressor~~ = SVR

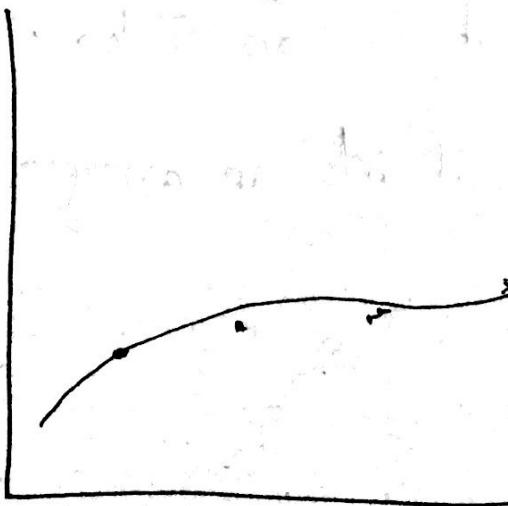
```
from sklearn.svm import SVR
```

```
regressor = SVR(kernel='rbf')
```

```
regressor.fit(x, y)
```

• - outlier.

~~If~~ ~~yo~~



* If you take the graph, the outlier point is CEO point with level 10. So the SVR considers it as an outlier because the distinction b/w each point and to this point is very large. It suddenly jumped to very high. So SVR marks it as an outlier.

* Predicting the result.

y-pred = regressor.predict(sc_x_transform([6.5]))

X

Because the input in sc_x_transform must be an array. So

y-pred = regressor.predict(sc_x_transform(
np.array([[6.5]])))

Why second bracket?

* Because if we only put one bracket it will consider it as a vector.

To make it into an array - take another bracket.

* The above command will return the scaled prediction of salary like '0.05'.
But we want the real salary.

Do,

$y_pred = sc_y.inverse_transform(\text{regressor}.predict(\text{sc_x.transform(np.array([[6.5]]))}))$

Why sc_y , why not sc_x .

Because, we are here predicting y , value and the object related to y . value is sc_y not sc_x .

* The $inverse_transform$ returns the original value of the scaled object.

R:

In R, you need to install package : e1071
`install.packages('e1071')`
`library(e1071)`

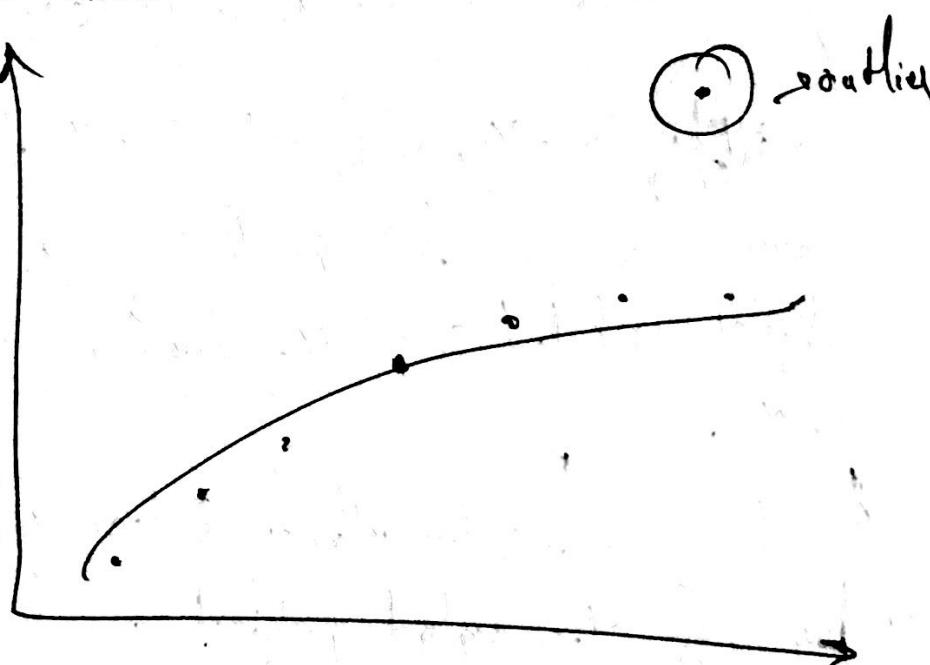
* creating & fitting the regressor to dataset.

`regressor = svm(formula = Salary ~ .,
data = 'Position_salaries.csv',
type = 'eps-regression')`

*** In `type`, if we do regression then
we choose 'eps-regression' (or) if we do
classification then we choose 'c-classification'

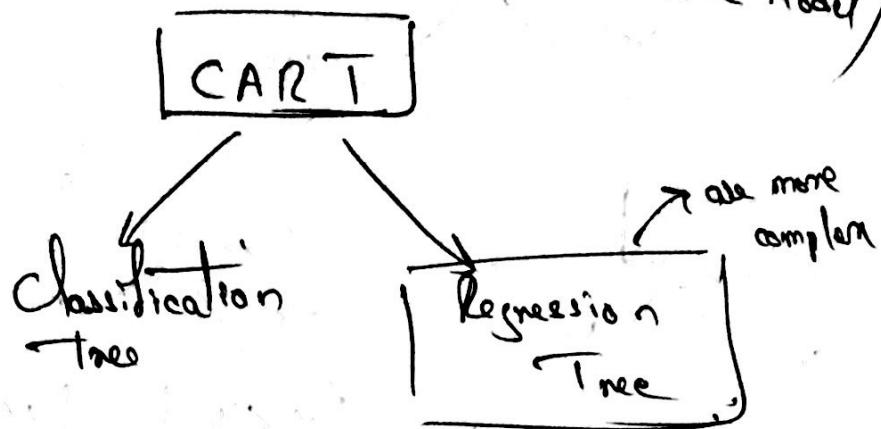
*** Here kernel is by default 'rbf'.

If we don't take 'kernel' argument.

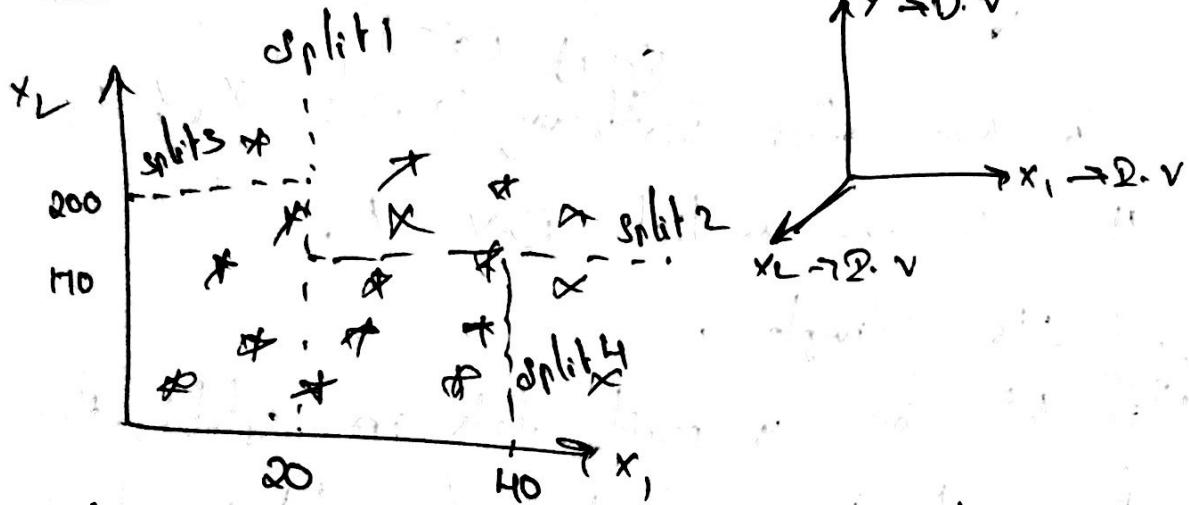




Intuition :

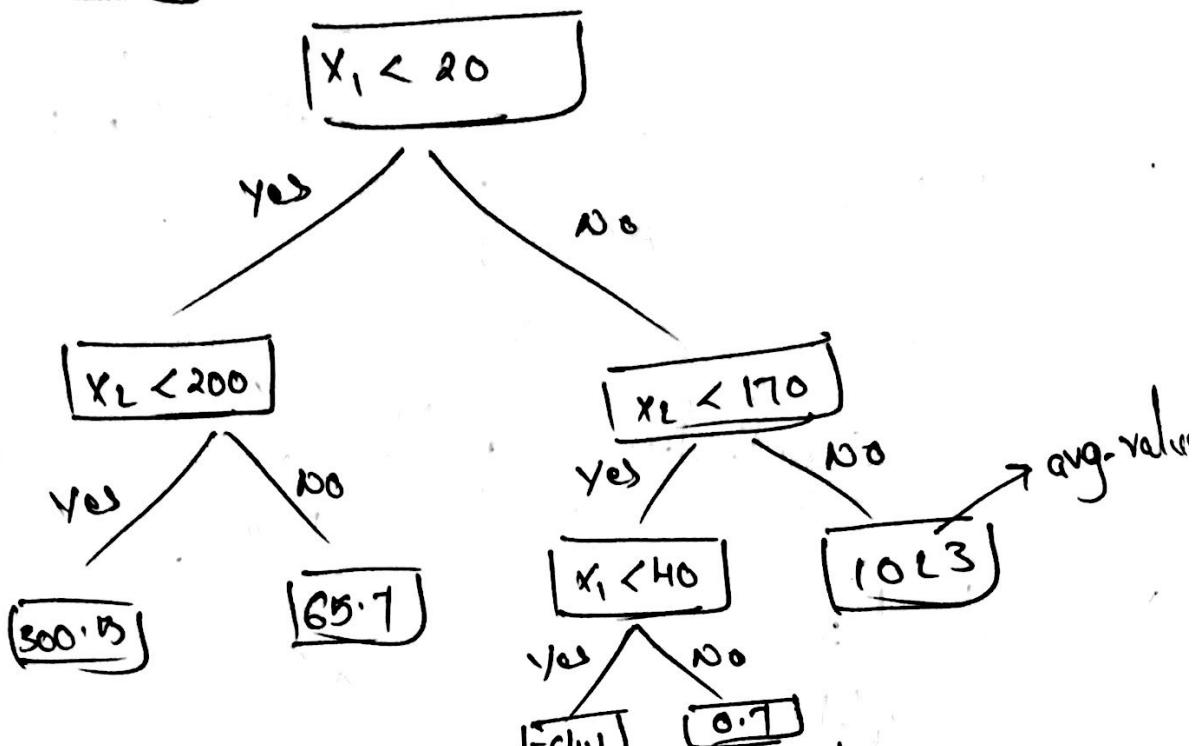


Given Dataset :



The algorithm divides the given dataset into splits according to information entropy. These splits are also known as leaves. The final leaves are known as 'Terminal leaves'. These leaves are optimal.

Building D.Tree :



Qs. Now after Scattering T.N's X_1 & X_2 ,
How do you predict O.V. Y.

* * * D.T.'s Simple .

* * * If lets say y point falls in split 1, then
there will be avg-value of split 1 and this
value will be assigned to y .

* * * This avg-value is calculated using all
points in that split and will be assigned to
any 'y' value that falls into that split .

Importance of splits:

If we ~~do not~~ have not done any split then we ~~should~~ might take avg of whole dataset and assign to every 'y' which is very less accurate.

By splitting the data and assigning that split avg value, we get a very accurate value,

Python

Fitting the DT regression to the dataset.

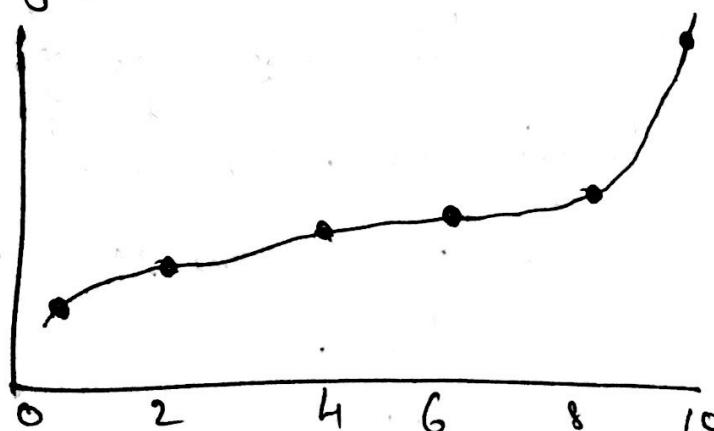
```
from sklearn.tree import DecisionTreeRegressor
```

```
regressor = DecisionTreeRegressor(random_state=0)
```

```
regressor.fit(x, y)
```

to get same result as tutorial.

* If you plot normally as in L.H.L, Poly - then the graph will be like



But in this interval there should only be one pred-value that is avg.

**** So, To clearly see the interval have a constant value for the interval, you may plot the higher-resolution Graph i.e.

visualising with higher resolution

```
x-grid = np.arange(min(x), max(x), 0.01)
```

```
x-grid = x-grid.reshape(len(x-grid), 1)
```

```
plt.scatter(x, y, color='red')
```

```
plt.plot(x-grid, y, color='blue')
```



Install library 'rpart'

```
# install.packages('rpart')
```

```
regression = rpart(formula = Salary ~ .,  
                    data = dataset)
```

If you simply plot then you get a single straight-line. We encountered like this in SVR - python, there we didn't do feature scaling.

But D.T.R takes care of feature scaling. Do - What's Problem here?

Sample
Graph.

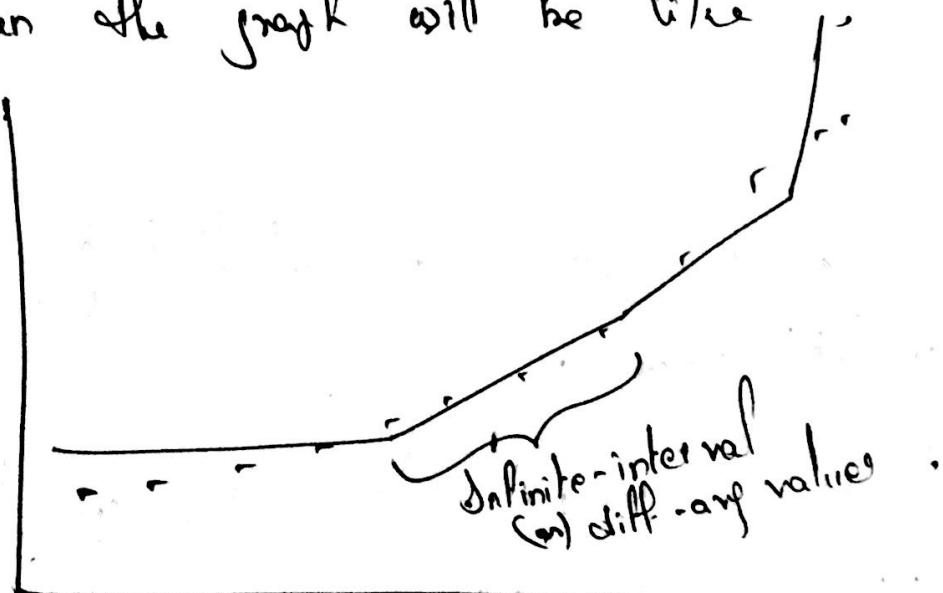


The problem is it took the whole dataset
as a single split.

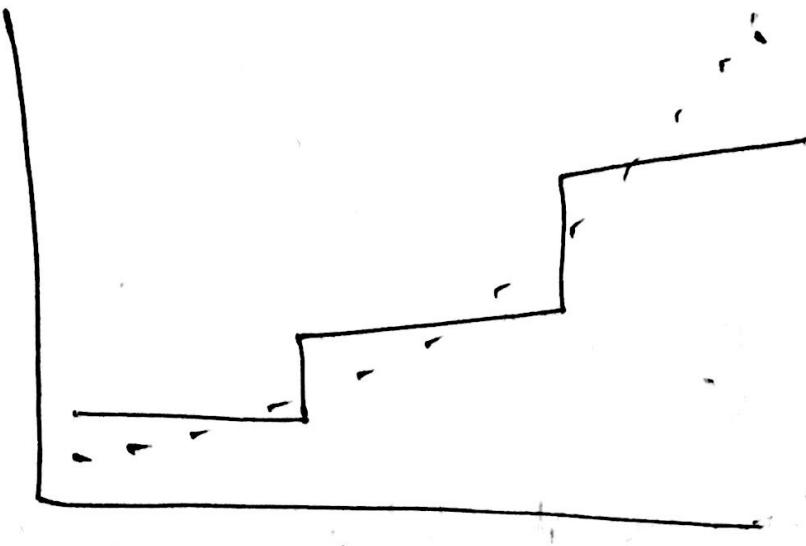
So, To solve this problem, we put a parameter,
i.e control.

regressor = rpart (formula = Salary ~ .,
data = dataset, control = rpart.control(
minsplit = 1))

Then the graph will be like



So, To solve the above problem,
Consider higher-resolution plotting.



Perfect Graph

Visualising using higher resolution.

library(ggplot2)

*-grid ~ seq(min(dataset\$file), max(dataset\$file)
0.01)

geom_point() +

geom_line() +

Here, the y.pred is controlled by
'minsplit' value.

**** Change the 'minsplit'. value then y.pred value
will be changed.

Random - Forest Regression

Non-linear Reg

(i) Multiple Deep Forest trees

Ensemble learning:

Is when you take multiple algorithms (i) same algorithm multiple times to get more powerful than the original one.

Random - Forest is one type of ensemble learning.

Intuition :

Step 1: Pick a random k points from training set.

Step 2: Build the Decision tree associated to these k data points.

Step 3: Choose the number of trees you want to build and repeat steps 1 & 2.

Step 4: For a new data point, make each one of your N trees predict y value. Now calculate the avg of all these y values and assigned the $\text{avg}(y\text{-values})$ as the predicted value.

Python :

```
# fitting random forest regression into the dataset  
from sklearn.ensemble import RandomForestRegressor  
regressor = RandomForestRegressor(n_estimators=10,  
                                     random_state=0)
```

```
regressor.fit(x, y)
```

```
y_pred = regressor.predict(6.5)
```

* When n.est = 10 - The prediction was closed.

* When n.est = 100 - " " " even closed

* " " " 300 - " " " more even, closed and collated with real value

Random Forest has beaten Polynomial Regression Prediction,

R

* If we dataset \$ level (or) dataset \$ salary.
The \$ sign creates a vector.

* To have a data frame we dataset[1] renders

fitting the randomForest Regressor to the dataset

install.packages('randomForest')

library(randomForest)

set.seed(1234)

regressor = randomForest(x = dataset[,
y = dataset\$Salary, ntree = 500])

y.pred = regressor\$predict(data.frame(here =
6.5))

*** The Best M.L models are combination
of M.L Models .

R²-Adjusted:

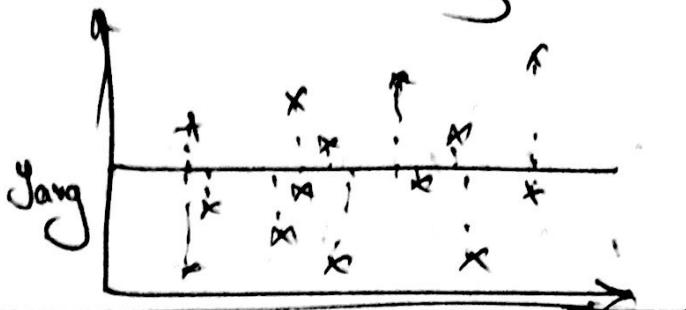
Linear Regression:

$$SS_{res} = \text{sum } (y_i - \hat{y}_i)^2$$

→ sum of squared of residuals.



Now instead take avg line i.e.



$$SS_{tot} = \text{sum } (y_i - \bar{y})^2$$

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

R^2 - tells us how good is our line (res) compare to the avg line.

The closer R^2 is to 1, the better.

SS_{res} should be smaller to get R^2 closer to 1

R^2 can be -ve. The model here is broken.

Normally R^2 lies b/w 0 & 1.

Adjusted R^2 :

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

problem: the more variables you add.

so when adding new variable, this new D.V will affect SS_{res} such that SS_{res} will be decreased and R^2 will be increased.

In the worst case, if the model thinks this new variable will worsen the prediction (\Rightarrow) increase SS_{res} , it will put the coefficient of D.V as

0, making it 0.

$$\text{Ex: } y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 \quad \begin{matrix} \xrightarrow{\text{coefficient}} \\ b_3 \end{matrix} \quad \begin{matrix} \xrightarrow{\text{D.V}} \\ x_3 \end{matrix}$$

$$\boxed{\text{Adjusted } R^2 = 1 - (n - R^2) \frac{n - 1}{n - p - 1}}$$

P = no. of regressors

n = sample size

D.V increase, P increase

Evaluating Regression Models:

Evaluating Regression Models:

- Adj R^2 should be more.

R^2 should be more.

Observe the coefficients. See if they are in the same, changing them will change the value of D.V in the same direction.

* If a value is like

~~7.966e-01~~ means

0.7966

(01)

2.991e-02 means

0.02991

Classification

Regression : The o/p variable takes continuous values.

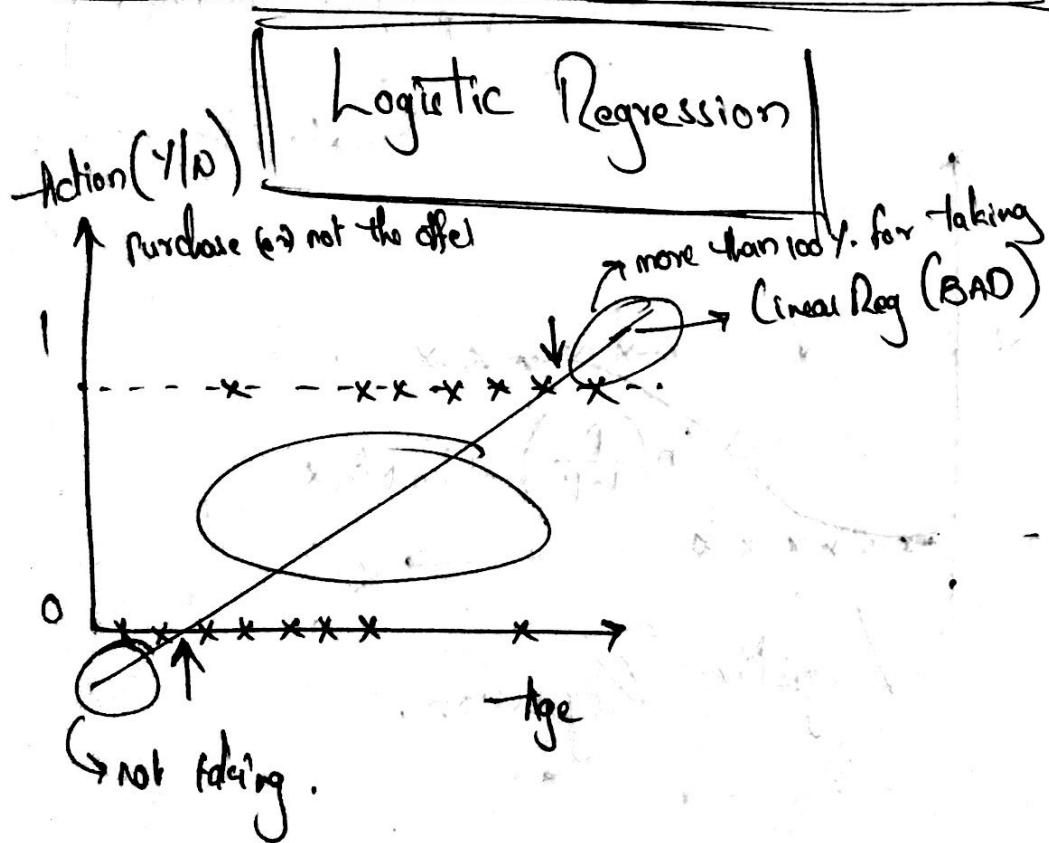
classification : The o/p variable takes class labels.

Regression is used to predict o/p values.

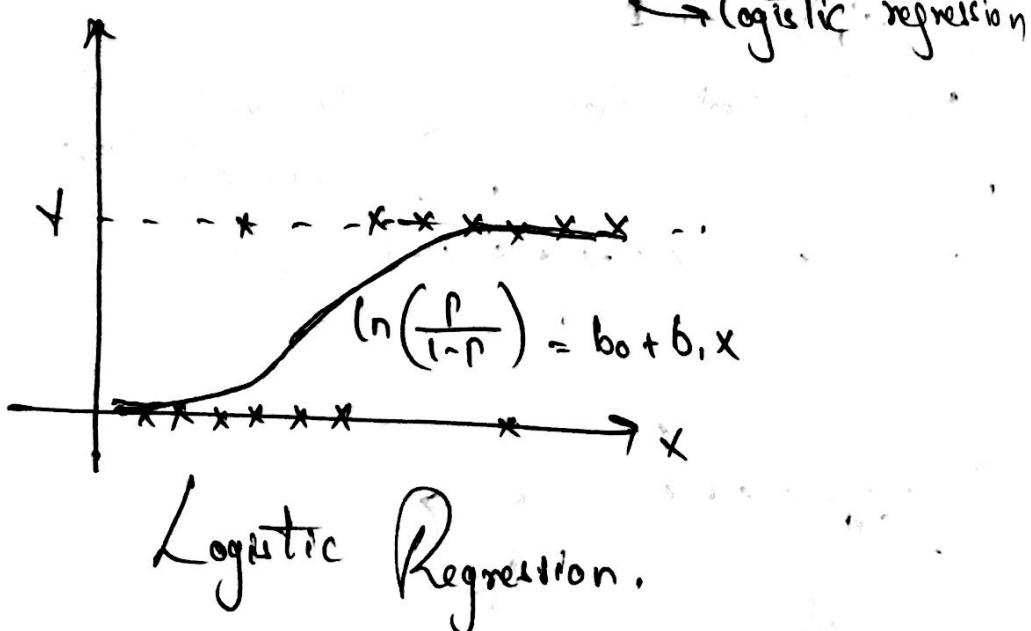
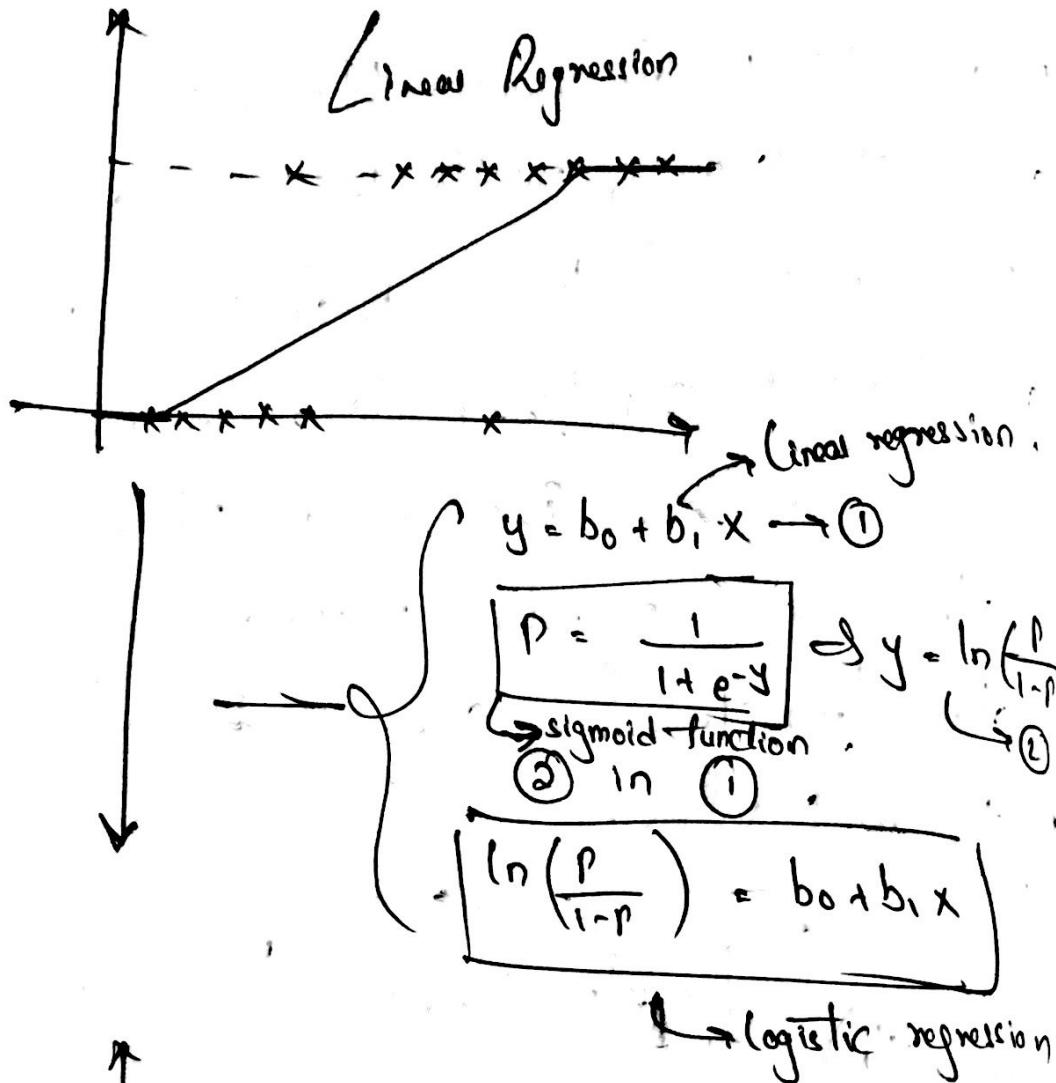
classification is group the o/p into a class.

Regression trees have O.p variables that are continuous values (or) ordered values.

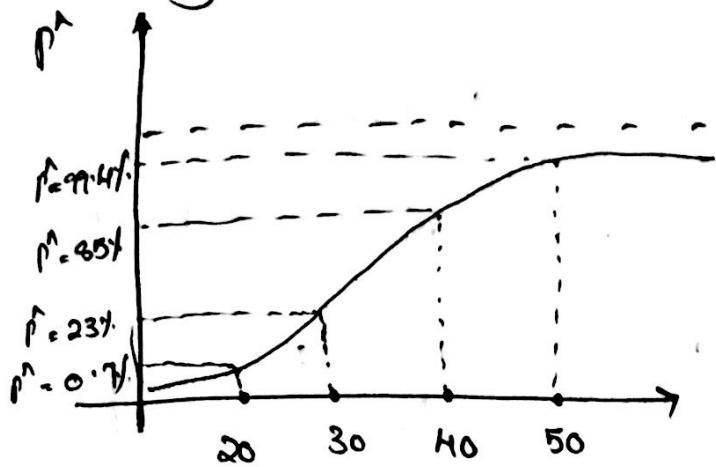
Classification trees have Categorical ~~data~~ and Unordered values.



Modifying the above Diagram

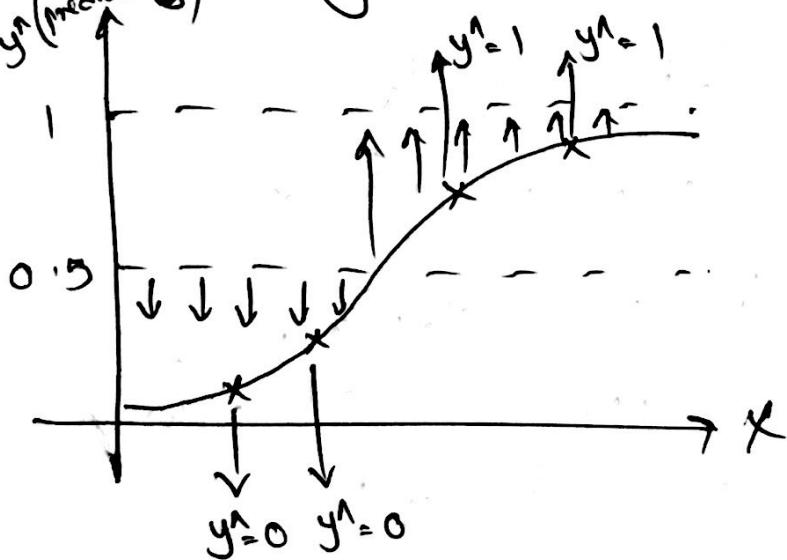


probability $\rightarrow \hat{P}(\hat{P})$



How do you predict?

Consider an arbitrary line (Here it is 0.5)



Python:

Use Data Preprocessing Template

*** In logistic regression We have to do
feature Scaling.

fitting logistic regression to our dataset

```
from sklearn.linear_model import LogisticRegression
```

```
classifier = LogisticRegression(random_state=0)
```

```
classifier.fit(X_train, Y_train)
```

predicting test set results

```
y_pred = classifier.predict(X_test)
```

Making the Confusion Matrix.

```
from sklearn.metrics import confusion_matrix
```

```
cm = confusion_matrix(y_test, y_pred)
```

This confusion matrix has the correct & wrong predicted values of y_pred compared to y_test.

Output of cm in terminal will be like

```
array([[65, 3], [8, 24]])
```

Correct

65 + 24 = correct

8 + 3 = incorrect.

Predicted	Predicted:		
No	Yes		
Actual	65	3	
Yes	8	24	

Graph (Very Important ****):

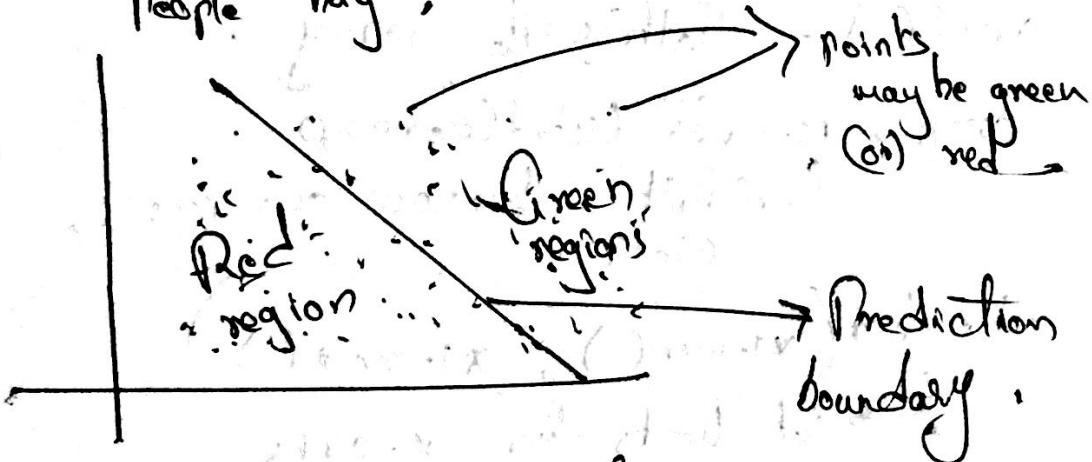
red point - User didn't buy

Green point - User bought

Red region - The region is predicted region where people don't buy.

Green region - The region is predicted region where people buy.

e.g.:



Point - It is the truth.

Region - It is Prediction

*** Logistic Regression is a linear classifier. So the prediction boundary is a straight line.

visualising the training set results .

from matplotlib.colors import ListedColormap :-

x-set, y-set = x-train, y-train

$x_1, x_2 = np.meshgrid(np.arange(start = x-set[:, 0].min() - 1, stop = x-set[:, 0].max() + 1, step = 0.01), \dots)$

plt.contourf(x1, x2, classifier.predict(np.array([x1.ravel(), x2.ravel()]).T)).reshape(x1.shape)

alpha = 0.75, cmap = ListedColormap(['~~blue~~red', 'grey'])

contour - It draws a line and shades colors to either side of line .

ravel() - flattens the array .

cmap - None or ListedColormap

Default \hookrightarrow your own colors

plt.xlim(x1.min(), x1.max())

\hookrightarrow set limits for x-axis

plt.ylim(x2.min(), x2.max())

for i, j in enumerate(np.unique(y-set)):

plt.scatter(x-set[y-set == j, 0], x-set[y-set == j, 1])

C = ListedColormap([('red', 'green'))(i),
label = j)

```
plt.title( )  
plt.xlabel( )  
plt.ylabel( )  
plt.legend()  
plt.show()
```

* Also Visualise the training set.

R:

- * Use Data preprocessing template
- * Change the needed requirement

** Logistic Regression needs feature scaling.

Here we only scale column 1, 2.

Feature Scaling

```
training-set[, 1:2] = scale(training-set[, 1:2])  
test-set[, 1:2] = scale(test-set[, 1:2])
```

Fitting logistic regressor to our training set:

Here we use existing function of R i.e.

'glm' to build logistic regression.

Generalised linear model,

classifier = glm(formula = Purchased ~ .,
family = binomial, data = training-set)
↳ Logistic regression belongs to binomial family.

Predicting Test Results:

Usually we predict in one step. But here we are calculating in 2 steps adding calculation of Probability. logistic

prob-pred = predict(classifier, type = 'response',
newdata = test-set[-3]) Removing D.V
— Here 'type = response' for logistic regression
We choose 'type = response', so that it will return the results in a single vector.

y-pred = ifelse(prob-pred > 0.5, 1, 0)
~~* * * To add [-3]~~ true
card false

* * * To add [-3] i.e. to remove D.V, its your wish.

* * * It is mandatory to add response.

* * * Logistic Regression returns probability that you need to convert.

Build Confusion Matrix i.e. Evaluating:

$cm = \text{table}(\text{test_set[, 3]}, \text{y_pred})$

↳ Create confusion matrix.

$cm = \begin{matrix} & 0 & 1 \\ 0 & 57 & 7 \\ 1 & 10 & 26 \end{matrix}$

$\left. \begin{matrix} & 57 & 7 \\ & 10 & 26 \end{matrix} \right\} 100$

Visualising The Training Set Results:

Install 'ElemStatLearn'

install.packages('ElemStatLearn')

library(ElemStatLearn)

set = ~~test~~ set training.set

$x_1 = \text{seq}(\min(\text{set[, 1]}), \max(\text{set[, 1]}),$
by = 0.01)

$x_2 = \text{seq}(\min(\text{set[, 2]}), \max(\text{set[, 2]}),$
by = 0.01)

grid.set = expand.grid(x_1, x_2)

colnames(grid.set) = c('Age', 'EstimatedSalary')

*** Setting column names for newly created

grid.set which contains x_1 & x_2 as its

columns. We then predict for this ~~grid~~

grid-set giving colors 'springgreen3' or 'tomato'

prob-set = predict(classifier, type = 'response',
newdata = grid-set)

y-grid > ifelse(prob.set > 0.5, 1, 0)

plot(cet[, -3]), plotting points + boundary

main = 'Logistic Regression (Training set)'

xlab = 'Age', ylab = 'Estimated Salary',

xlim = range(x1), ylim = range(x2))

contour(x1, x2, matrix(as.numeric(y-grid),
length(x1), length(x2)), add = TRUE)

Drawing the line - Logistic Reg line

points(grid-set, pch = '.', col = ifelse(y-grid == 1,
'springgreen3', 'tomato'))

coloring the plane

points(~~grid-set~~, pch = 19, bg = ifelse(cet[, 3]
== 1, 'green4', 'red3'))

Should the points be in circular-type, rectangle, etc

else

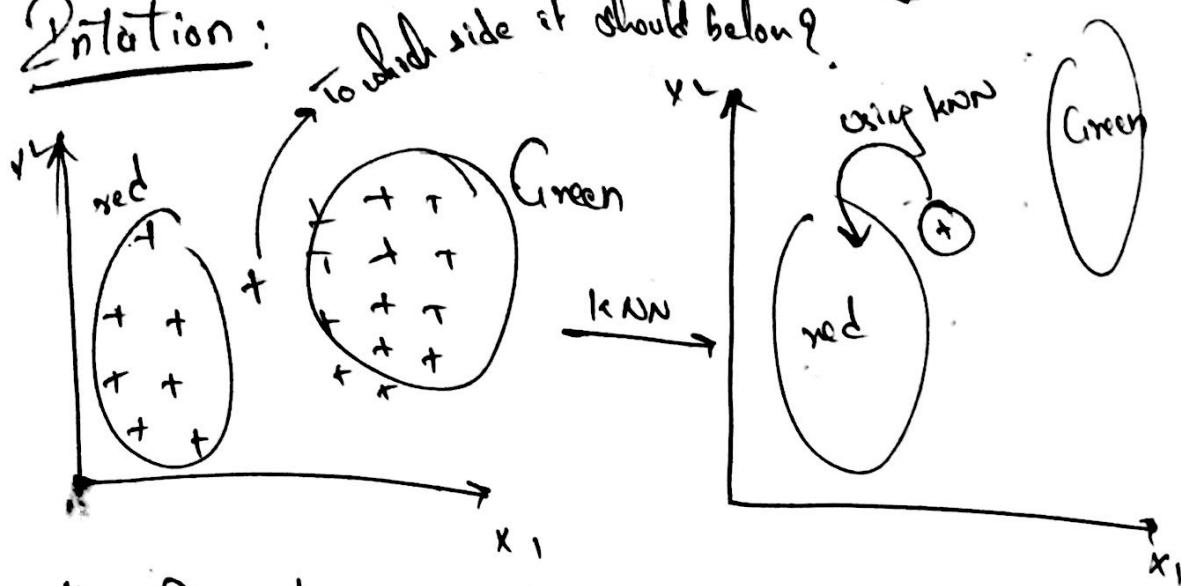


rect

K - Nearest Neighbours

Non-Linear
Classifier.

Intuition: To which side it should belong.



How Does KNN Work:

Step 1: Choose no of k neighbours. By default $k = 3$.

Step 2:

Take k nearest neighbours of the new data point, according to the Euclidean distance.

Step 3:

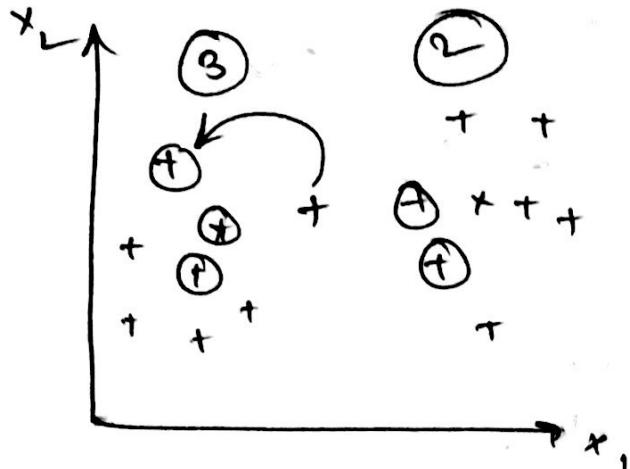
Among these neighbours, count the number of data points in each category.

Step 4:

Assign the new data point to the category where you counted the most neighbours.

Euclidean distance b/w two points $p_1(x_1, y_1)$
& $p_2(x_2, y_2)$

$$= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$



Python :

Use classification - template.py .

fitting the kNN to our training set .

```
from sklearn.preprocessing import KNeighborsClassifier  
classifier = KNeighborsClassifier(n_neighbors=5  
                                    metric='minkowski', p=2)
```

* * * To choose Euclidean distance for kNN,

* * * You have to take metric = 'minkowski' &
 $p=2$.

classifier.fit(x_train, y_train) .

~~# fit predict x-test~~

Build confusion-matrix

0	1	
0	64	4
1	3	29

64+29 - correctly predicted.

7 - incorrectly predicted.

Comparing Incorrectly predicted results.

* Logistic Reg (Incorrect) - ~~17~~ 17

* KNN (Incorrect) - 7 (Much Better)

Visualise the Graph.

*** You might think that the graph should change for training-set & test-set. It doesn't because, the graph is build with every pixel point with difference of 0.01. So, the test-points are also in that array containing points with diff of 0.01. So when you take test-set it will just scatter points.

R:

Use the template.

fitting knn to the training set & predicting the test-set results.

$y.\text{pred} = \text{knn}(\text{train} = \text{training.set} \setminus \text{test}[, -3],$
 $\text{test} = \text{test.set}[, -3], \text{cl} = \text{training.set}[, 3]$
 $k = 5)$

↓
no. of neighbours

telling D.V's.

* * Note: Usually every time, we first create regressor & then predict the test-set. But in knn, we do both at same time.

Making confusion matrix.

$\text{cm} = \text{table}(\text{test.set}[-3], \text{y.pred})$

Visualising

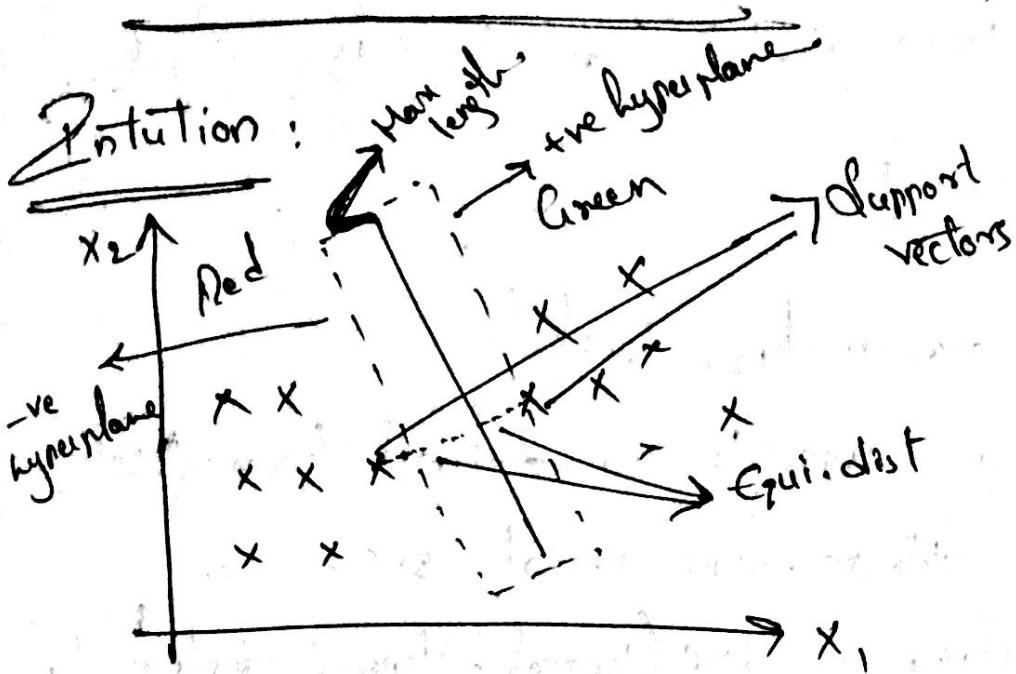
- Everything is same except only one change.
Usually

$y.\text{grid} = \text{predict}(\text{classified}, \text{newdata} = \text{grid.set})$

Change →

$y.\text{grid} = \text{knn}(\text{train} = \text{training.set}[, -3], \text{test} = \text{grid.set}, \text{cl} = \text{training.set}[, 3],$
 $k = 5)$

Support Vector Machine (SVM)



What is popular about SVM?

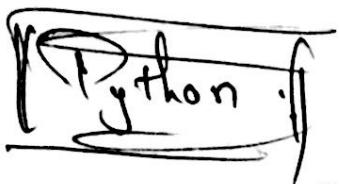
Problem: Machine must be able to identify apples & oranges.

Usually: Usually we take apples that ~~are~~ that are more likely apples (100% apples) & oranges and train machine on them & then predict

SVM: But in SVM we consider apples that are less likely apples but apples (~~or~~) oranges and train machine on them.

* * * * These less likely apples (~~or~~) oranges are called Support Vectors.

* * * This SVM may seem like risky because
We train machine on boundary objects.



One classification-template

fitting svm to our ~~dataset~~ training-set.

```
from sklearn.svm import SVC
```

```
classifier = SVC(kernel='linear', random_state
```

```
= 0)
```

```
classifier.fit(x_train, y_train)
```

CM:	0	1
0	66	2
1	8	24

incorrect - 10

correct - 90

Its similar to logistic Regression.

visualise as in classification-template.

R:

Use the classification template

fitting our svm to the dataset

library(e1071)

classifier = svm(formula = Purchased ~ . ,
data = training-set, type = 'C-classification',
kernel = 'linear')

predicting test-set results

y-pred = predict(classifier, newdata = test-set).

CM:

	0	1
0	57	7
1	13	23

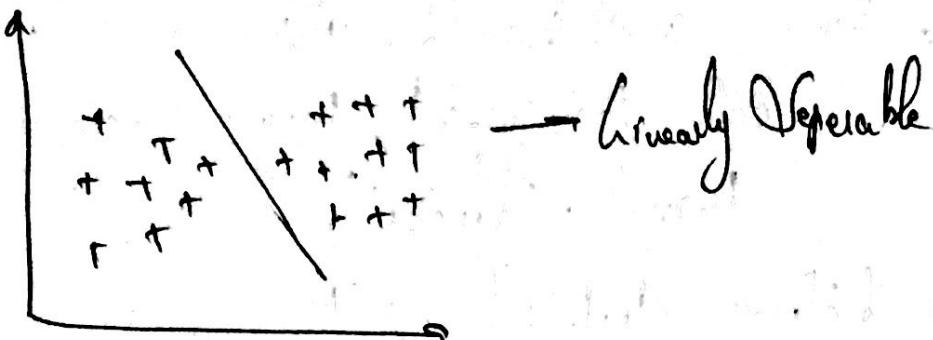
~~Correct~~

~~Incorrect~~

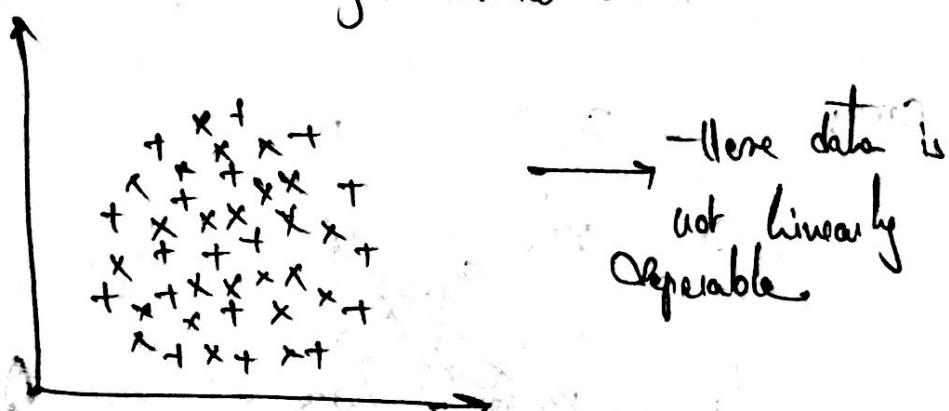
Kernel SVM

Intuition :

Usually you can draw a line when



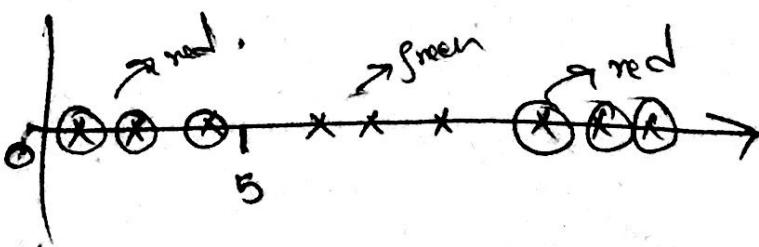
then how will you draw when



So, here we will take a higher dimension and make the data → to linearly separable and then do kernel . svm ,

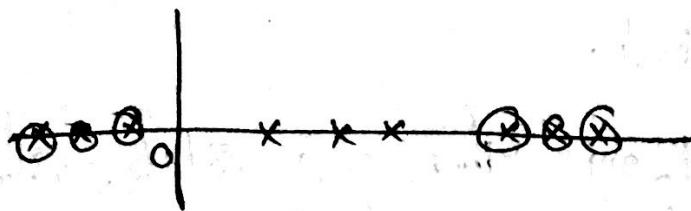
Higher - Dimensional Space

Consider a single dimension with 9 points .

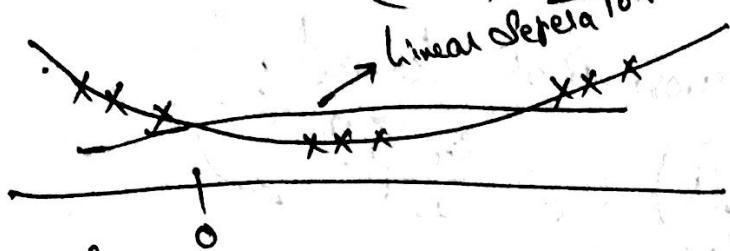


consider 3rd point at 5.

① then $f = x - 5$ will be



② Now do $f = (x - 5)^2$



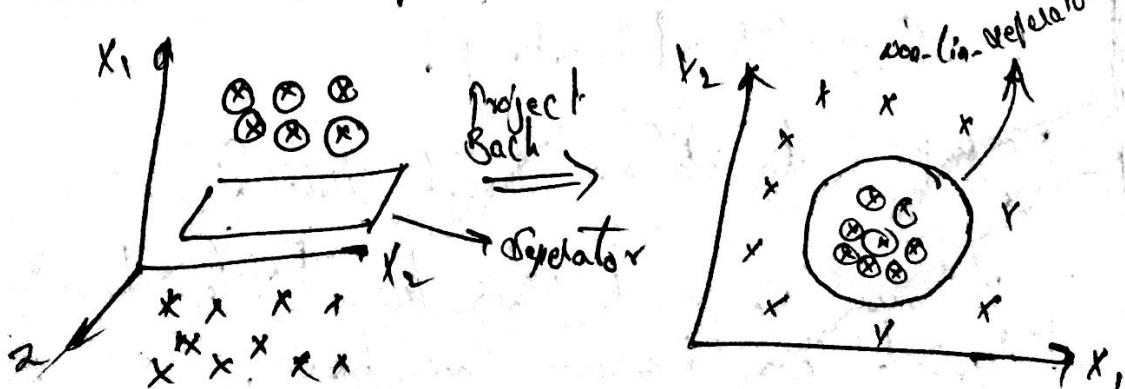
This happened possible by adding another dimension.

Similarly for 2D, you add another dimension and make it 3D.

$$f(x_1, x_2) \xrightarrow{\text{Mapping function}} (x_1, x_2, z)$$

③ Now after having linear separator we project it back to original dimension.

i.e. When you dim \leftarrow to 2D and make it 3D take linear separator as plane.



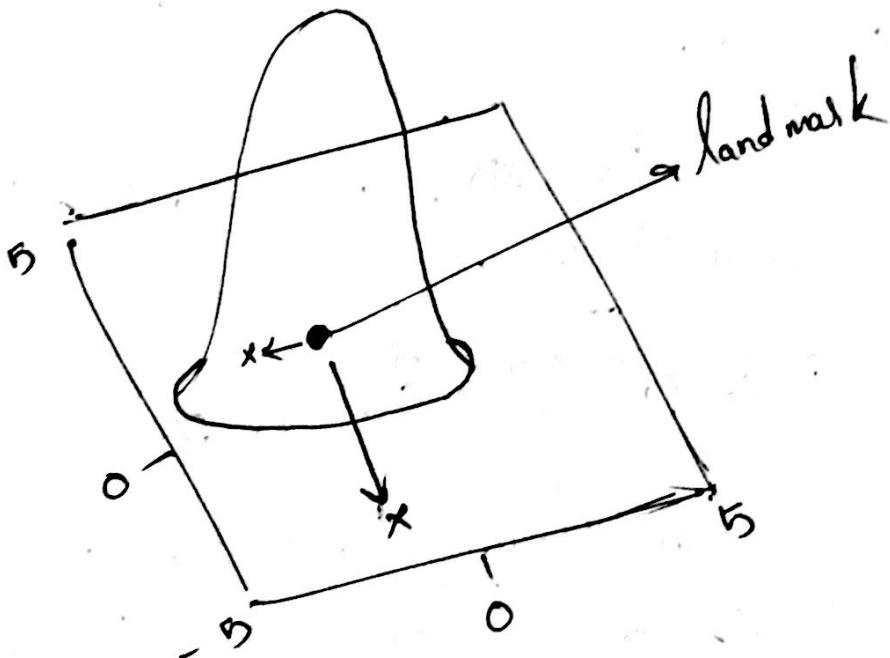
*** Mapping to higher Dimensional Space is highly Compute intensive. So this approach isn't the best.

But in Kernel-SVM there is a catch, without going into higher Dimensional Space.

The Gaussian RBF Kernel:

$$k(\vec{x}, \vec{l}_i) = e^{-\frac{\|\vec{x} - \vec{l}_i\|^2}{2\sigma^2}}$$

point landmark

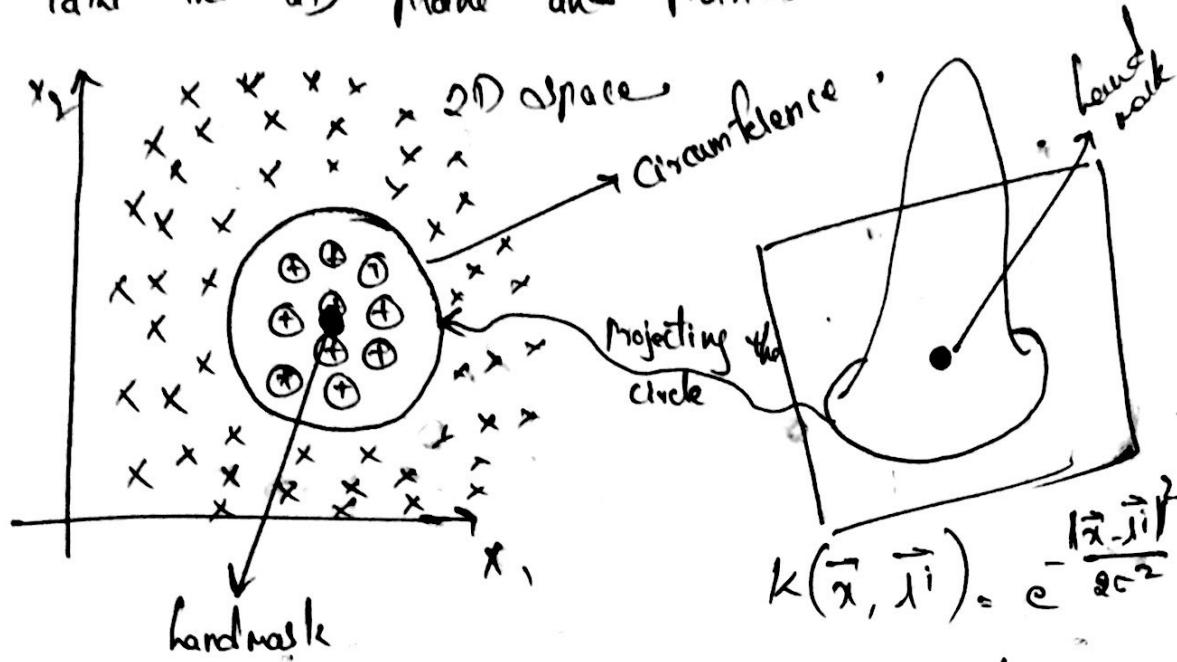


** If your point is closer to landmark then the k will be nearer to 1 and vice versa (By nature in formula).

**** Why do we need this?

Because we need this kernel function to separate our decision boundary.

Take the 2D plane and Points



*** Any point that will fall out of the Circumference then the value of that point will be nearer to 0.

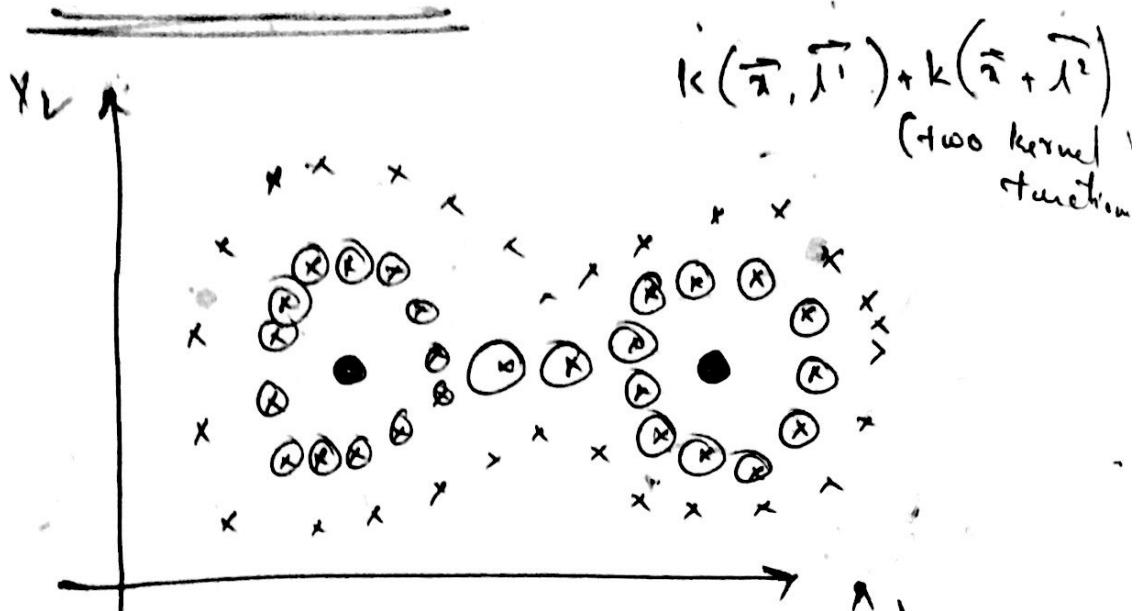
*** Any point that will fall out of the circumference, its value will be 1.

This is how you separate Green points and Red points.

*** Role of σ :

σ Defines the circumference. Increase the σ , circumference increases. So to get best results, You need to find best ' σ '.

Complex Example :



Now take two kernel functions:

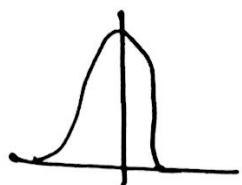
$$k(\vec{x}, \vec{j}_1) + k(\vec{x} + \vec{l}_1) \rightarrow ①$$

Substitute the point in ① in both kernel
If it's close to 0 - red point (x)
close to 1 - Green point (o)

Types of Kernel Functions:

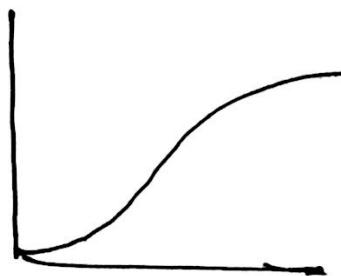
In all these types, the essence is
same. You choose a landmark and
from there you calculate the distance
using the formulae.

Types :



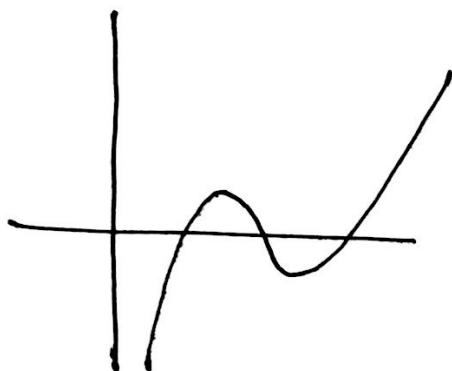
Gaussian
RBF kernel

$$k(\vec{x}, \vec{y}) = e^{-\frac{\|\vec{x} - \vec{y}\|^2}{2\sigma^2}}$$



Sigmoid
kernel

$$k(x, y) = \tanh(\gamma \cdot x^T y + r)$$



Polynomial
kernel

$$k(x, y) =$$

$$(\gamma \cdot x^T y + r)^d, \gamma > 0$$

Note : You ^{do not} need to go deep in every type of kernel, just know the types.

To just see how these kernels exists in 3D:

Visit : [mlkernels.readthedocs.io/en/latest/kernelfunctions.html](#)

Python :

fitting the regression to our training set.

```
from sklearn.svm import SVC
```

```
classifier = SVC(kernel='rbf', random_state=0)
```

```
classifier.fit(training-set X_train, Y_train)
```

Rest is same

Confusion Matrix:

	0	1
0	64	4
1	3	29

After using kernel-SVM, it makes the data linearly separable

R:

You have to install package 'e1071'(or) kernlab.

```
install.packages('e1071')
```

```
library(e1071)
```

```
classifier = svm(formula = Purchased ~ .,
```

```
data = training-set, type = 'C-classification'
```

```
kernel = 'radial')
```

~~These 'radial' is Gaussian kernel~~

Remaining is Laine

CM

	0	1
0	38	6
1	4	32

Naive - Bayes Theorem :

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Bayes Theorem}$$

P - Probability.

Let's take :

Machine 1 : 30 wrenches/hr ~~P(Mach)~~

Machine 2 : 20 wrenches/hr

Out of all produced parts :

We can see 1% are defective.

Out of all the defective parts :

Machine 1 - 50%.

Machine 2 - 50%.

Question :

What is the probability that a part produced by machine 2 is defective.

$P(\text{Mach}_1)$ producing parts = $30/50 = 0.6$

$P(\text{Mach}_2)$ producing parts = $20/50 = 0.4$

$$P(\text{Defective}) = 1\%.$$

$$P(\text{Machine 1} | \text{Defect}) = 50\%.$$

$$P(\text{Machine 2} | \text{Defect}) = 50\%.$$

$$P(\text{Defect} | \text{Machine 2}) = ?$$

Using Bayes Theorem:

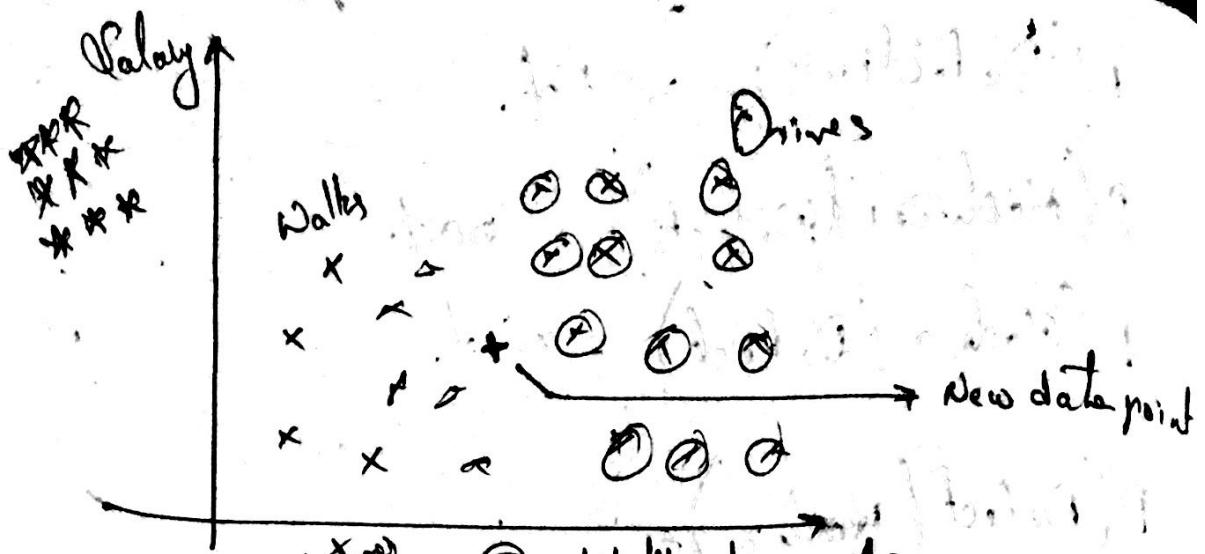
$$\frac{P(\text{Machine 2} | \text{Defect}) * P(\text{Defect})}{P(\text{Machine 2})}$$

$$= \frac{0.5 * 0.1}{0.4} = \frac{0.5}{4} = 0.0125$$

$$= 1.25\%.$$

Probability:

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$



$$P(\text{Walkes} | x) = \frac{P(x | \text{Walkes}) * P(\text{Walkes})}{P(x)}$$

#1 Posterior probability
 #2 Marginal likelihood
 #3 Likelihood
 #4 Prior probability

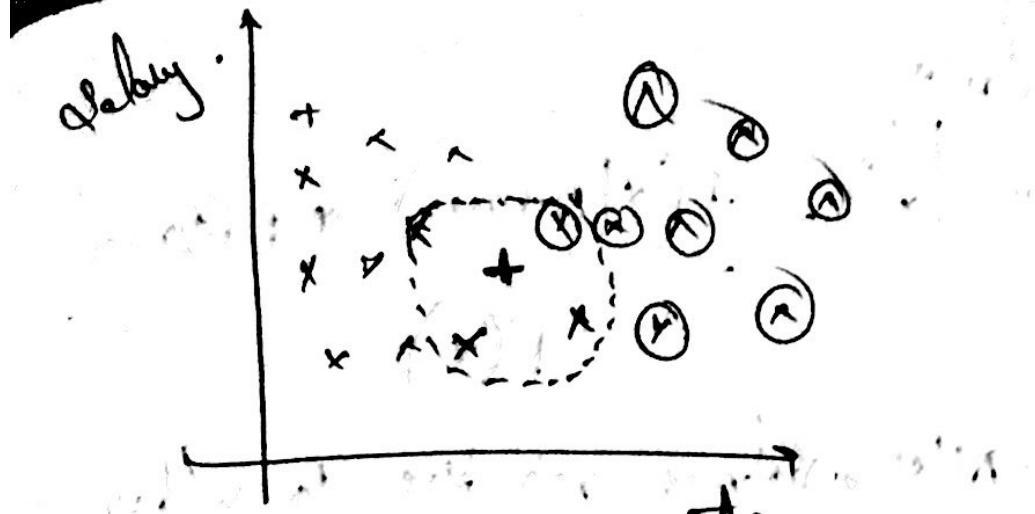
$$P(\text{Drives} | x) = \frac{P(x | \text{Drives}) * P(\text{Drives})}{P(x)}$$

$$P(\text{Walkes} | x) \text{ vs } P(\text{Drives} | x)$$

#1 $P(\text{Walkes}) = \frac{\text{No. of Walkers}}{\text{Total observations}}$

#2 How do you calculate $P(x)$?

P.T.O



1. Take a radius of your own choice and draw a circle around the new point.
2. Take all the points inside the circle except the new one.

$$P(x) = \frac{\text{No. of similar observations}}{\text{Total observations}}$$

$P(x)$ is ~~set~~ by adding new points, the probability of the new point falling in the radius.

#3

$$P(x | \text{Walker})$$

In the circle in the above diag, since we are calculating for walking, the probability of ~~set~~ the newly added point to be walking in the circle

$$P(x | \text{Walker}) = \frac{\text{No. of Walkers in circle}}{\text{Total no of Walkers}}$$

Similarly,

$$P(\text{Drives} | x) = \frac{P(x | \text{Drives}) * P(\text{Drives})}{P(x)}$$

After solving according to video

$$P(\text{Walks} | x) = 0.75$$

$$P(\text{Drives} | x) = 0.25$$

Since $0.75 > 0.25$

that person with features x , will most likely walk.

* * *

1. Why Naive?
2. $P(x)$
3. What happens when there is more than 2 features?
 1. A) Naive Bayes theorem takes independent assumptions. These assumptions are sometimes wrong. That is why it is called "Naive".

~~In~~ In our example, consider Age & Velocity.

Naive Bayes assumes that Age & Velocity are independent. But whereas they are not.

For Naive Bayes they have to be independent.

2.1) $P(x)$

Go Back in calculating $P(x)$. (Formulas)

3.1) If more than 2 features like walking & Driving. Then you have to calculate atleast 2.

In our example we have Walking & Driving. If we have walking probability then we have $P(\text{Driving})$. We do not need to calculate $P(\text{Driving})$ because

$$P(\text{Walking}) + P(\text{Driving}) = 1$$

Python! ~~error~~ There will be no Arguments in NB (ch 4)

```
from sklearn.naive-bayes import GaussianNB  
classifier = GaussianNB()  
classifier.fit(x-train, y-train)
```

cm:

	0	1
0	65	3
1	7	25

R:

```
library(e1071)
```

```
classifier = naiveBayes(x=training-set[, -3],  
y=training-set[, 3])
```

predicting test-set

y-pred = predict(classifier, newdata = test-set)

* After running the above one, in console

You can see factor levels = 0.

* That is because naive Bayes() class doesn't identify (i) cannot automatically encode ~~the factors~~.

You have to manually encode. ~~the factors~~ feature scaling
Do here without encoding. ~~feature scaling~~ factoring doesn't work. (Encoding = factoring).

Encoding the target feature as factor.

dataset \$ Purchased = factor(dataset \$ Purchased,
levels = c(0, 1))

cm :	0	1
0	57	17
1	7	29

① Decision Tree Classification

Pintution:

Splits are taken according to maximize the category. In background it is very deep history.

* These are very old methods. They were ~~were~~ verge on dying.

* So they came with upgrades:

- 1) Random Forest
- 2) Gradient Boosting
- 3) etc

* Decision Trees are used in other methods that creates ~~are~~ very powerful algorithms that are used in facial recognition, xbox connect etc (Wii).

Microsoft has used Random Forest.

Python :

at fitting the regressor to our dataset

```
from sklearn.tree import DecisionTreeClassifier
```

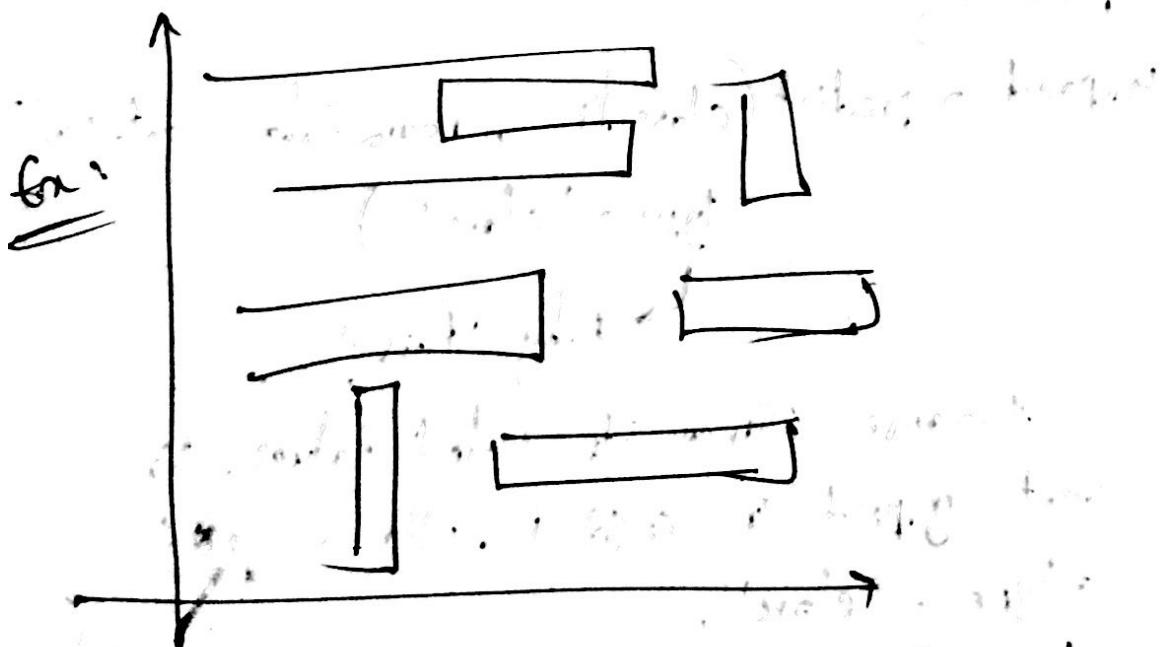
```
classifier = DecisionTreeClassifier(criterion =  
    'entropy', random_state  
    = 0)
```

criterion (ctr + Σ i : to see intro)

↳ entropy - for info gain

↳ gini - for impurity

```
classifier.fit(x_train, y_train)
```



It tries to catch every observation in the training set. It's called overfitting.

cm :

	0	1
0	62	6
1	3	29

R :)

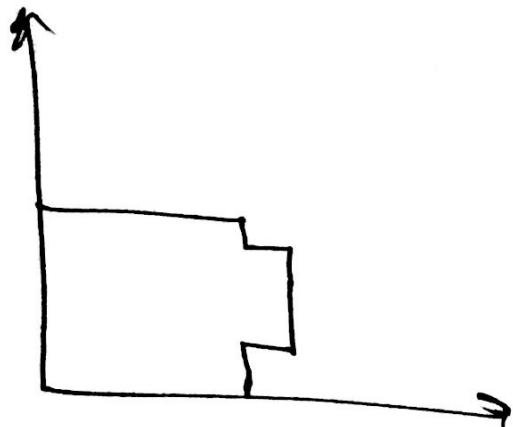
fitting the classifier library (rpart)
classifier = rpart (formula = Purchased ~ .,
data = training_set)

predicting

y_pred = predict (classifier, newdata = test_set,
type = 'class')
↳ Why this?

Because Since its scaled values, You
want y_pred to 0 or 1. So we add
type = 'class'.

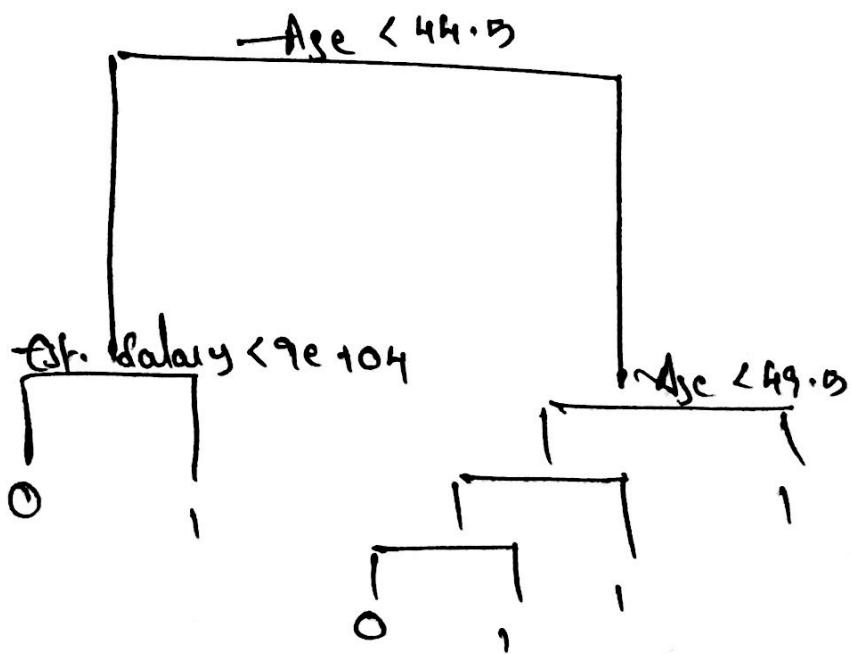
* * * *
* * * Here there is no overfit. So, the graph will be different than python.



→ plotting the decision tree

plot(classifier)

test(classifier)



* * * Since here feature & salary is not required, if you don't use feature & salary they you will have original values in the graph.