

## Data Science Training - Day 3

Multi-Linear Regression:

A	B	C	D	E
RNDSpend	Administration	MarketingSpend	State	Profit
165349.2	136897.8	471784.1	New York	192262
162597.7	151377.59	443898.53	California	191792
153441.51	101145.55	407934.54	Florida	191050
144372.41	118671.85	383199.62	New York	182902
142107.34	91391.77	366168.42	Florida	166188
131876.9	99814.71	362861.36	New York	156991
134615.46	147198.87	127716.82	California	156123
130298.13	145530.06	323876.68	Florida	155753
120542.52	148718.95	311613.29	New York	152212
123334.88	108679.17	304981.62	California	149760
101913.08	110594.11	229160.95	Florida	146122
100671.96	91790.61	249744.55	California	144259
93863.75	127320.38	249839.44	Florida	141586
91992.39	135495.07	252664.93	California	134307

## Multiple Linear Regression

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_n * x_n$$

The dataset won't get fit into the below formula:

$$\frac{\text{sum}(y) * \text{sum}(x^2) - \text{sum}(x)\text{sum}(xy)}{n * \text{sum}(x^2) - \text{sum}(x)^2} \quad b_0 = \text{intercept}$$

$$\frac{n * \text{sum}(xy) - \text{sum}(x) * \text{sum}(y)}{n * \text{sum}(x^2) - \text{sum}(x)^2} \quad b_1 = \text{slope}$$

To handle the state column (category), we need encoding. We will convert it to numeric format.

Encoding/dummy variable creation /One hot encoder:

Encoding takes unique values and alphabetically sorted. It will represent each unique value as a column.

California	<u>d0_C</u>
Florida	<u>d1_F</u>
New York	<u>d2_NY</u>

State	d0_C	d1_F	d2_NY
New York	0	0	1
California	1	0	0
Florida	0	1	0

Machine learns through rows – identifying the pattern instead of columns.

You can remove the dummy variables.

State	d0_C	d1_F	d2_NY
New York	0	0	1
California	1	0	0
Florida	0	1	0

In order to have a balanced learning, dummy variable is best suited. We won't use weighted labels.

```
 numOfDummyVar = numOfUniqVal - 1
```

California	d0_C	0
Florida	d1_F	1
New York	d2_NY	2
Texas		3
		4
		5

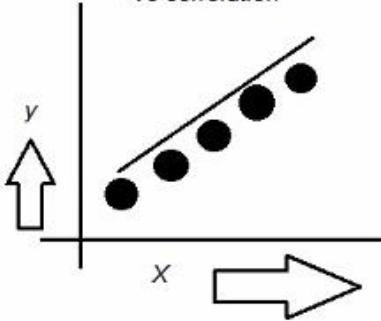
XNew - DataFrame

Index	RNDSpend	Administration	MarketingSpend	State_California	State_Florida	State_New York
0	165349	136898	471784	0	0	1
1	162598	151378	443899	1	0	0
2	153442	101146	407935	0	1	0
3	144372	118672	383200	0	0	1
4	142107	91391.8	366168	0	1	0
5	131877	99814.7	362861	0	0	1
6	134615	147199	127717	1	0	0
7	130298	145530	323877	0	1	0
8	120543	148719	311613	0	0	1
9	123335	108679	304982	1	0	0
10	181913	110594	229161	0	1	0
11	100672	91790.6	249745	1	0	0
12	93863.8	127320	249839	0	1	0
13	91992.4	135495	252665	1	0	0
14	119943	156547	256513	0	1	0
15	114524	122617	261776	0	0	1
16	78013.1	121598	264346	1	0	0
17	94657.2	145078	282574	0	0	1
18	91749.2	114176	294920	0	1	0
19	86419.7	153514	0	0	0	1
20	76253.9	113867	298664	1	0	0
21	78389.5	153773	299737	0	0	1
22	73994.6	122783	303319	0	1	0

Index	RNDSpend	Administration	MarketingSpend	State_Florida	State_New York
0	165349	136898	471784	0	1
1	162598	151378	443899	0	0
2	153442	101146	407935	1	0
3	144372	118672	383200	0	1
4	142107	91391.8	366168	1	0
6	134615	147199	127717	0	0
5	131877	99814.7	362861	0	1
7	130298	145530	323877	1	0
9	123335	108679	304982	0	0
8	120543	148719	311613	0	1
14	119943	156547	256513	1	0
15	114524	122617	261776	0	1
10	101913	110594	229161	1	0
11	100672	91790.6	249745	0	0
17	94657.2	145078	282574	0	1

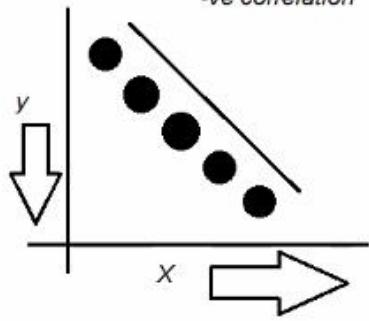
Linear Spread of Data

+ve correlation



Linear Spread of Data

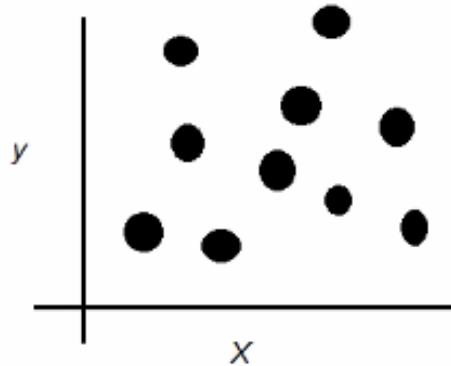
-ve correlation



Linear Regression

moderate to very strong correlation

### *Non-Linear Spread of Data*



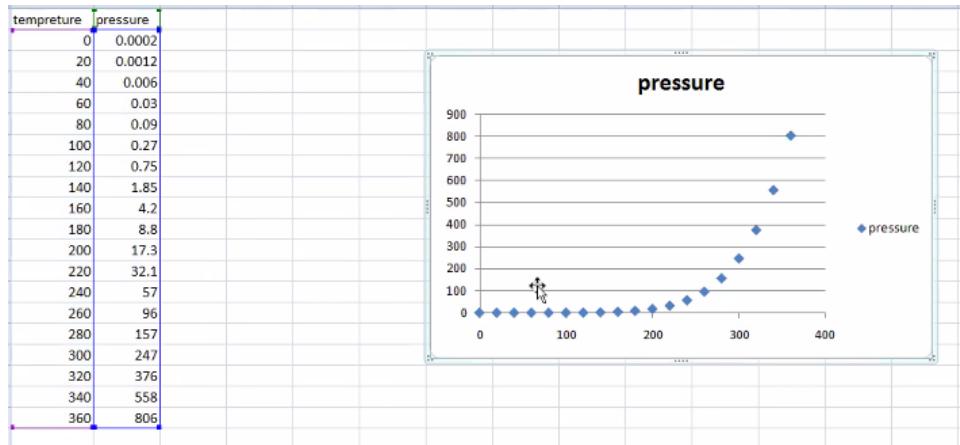
*correlation of most of the X features will be in negligible to no-correlation*

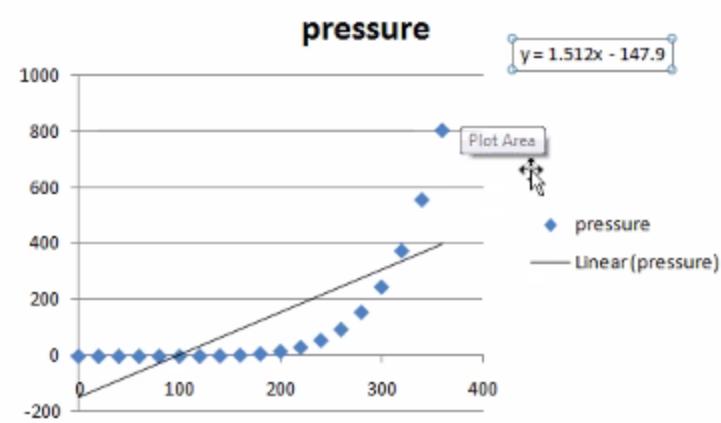
||

*Polynomial Regression  
DecisionTree Reg.  
RandomForest Reg.*

Take the example below:

This is not a linear problem.



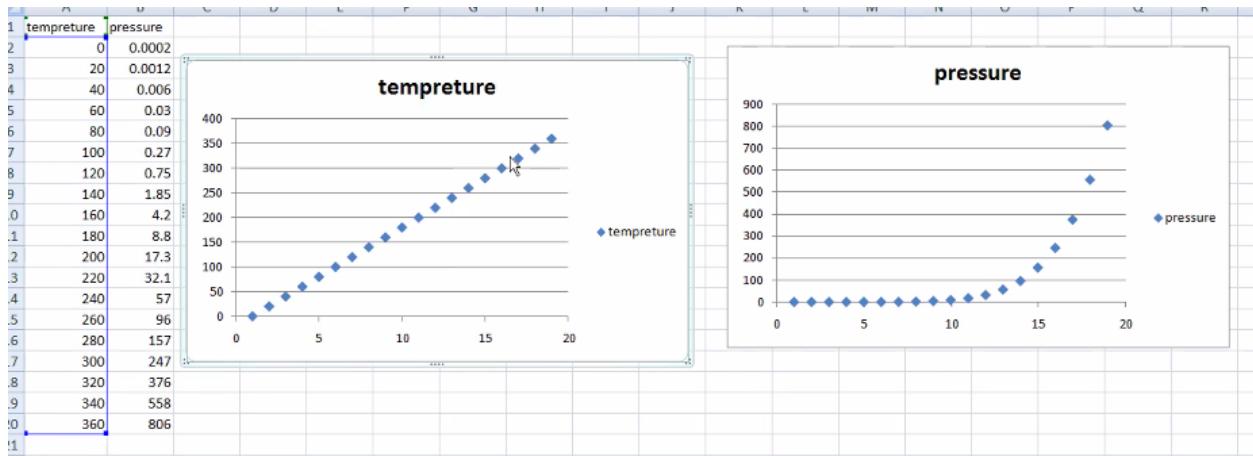
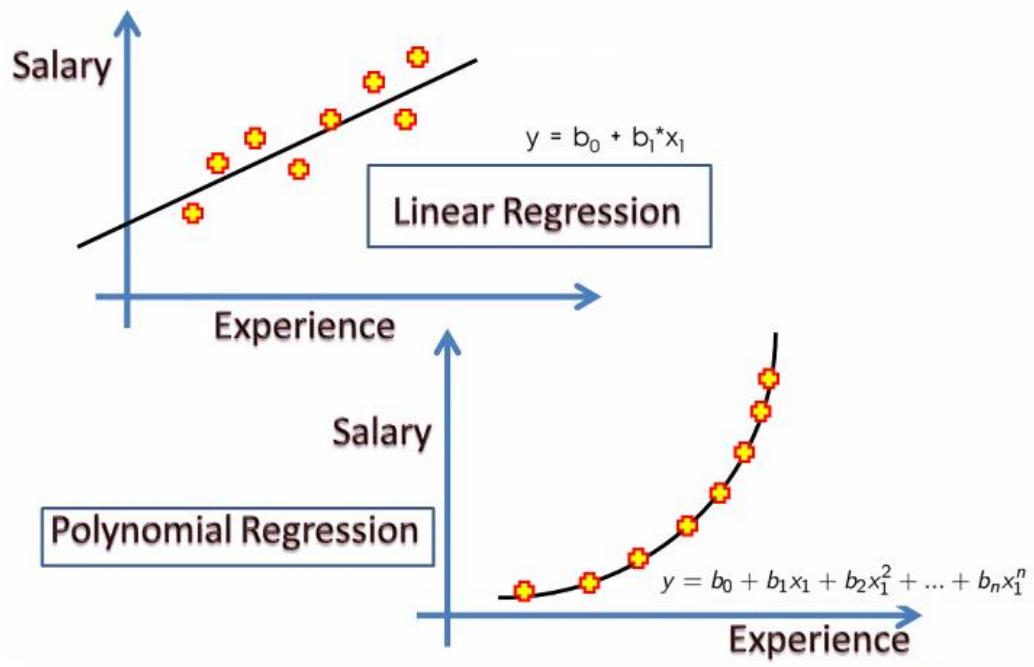


Here linear regression formula falls short. So, in this case you must use polynomial regression

## Polynomial Regression

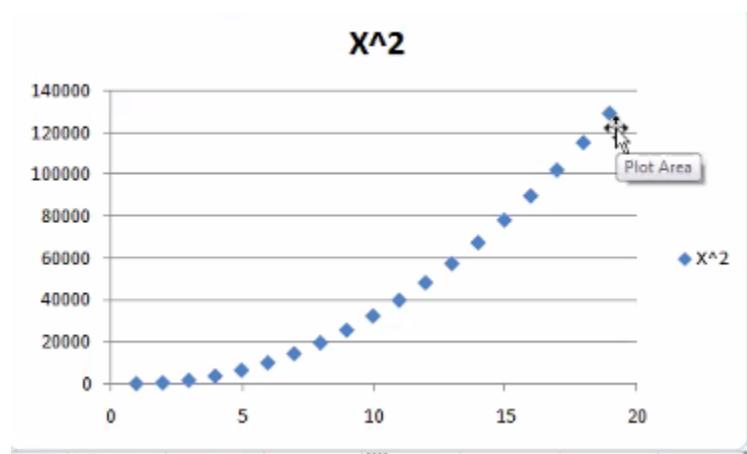
$$y = b_0 + b_1 x_1 + b_2 x_1^2 + \dots + b_n x_1^n$$

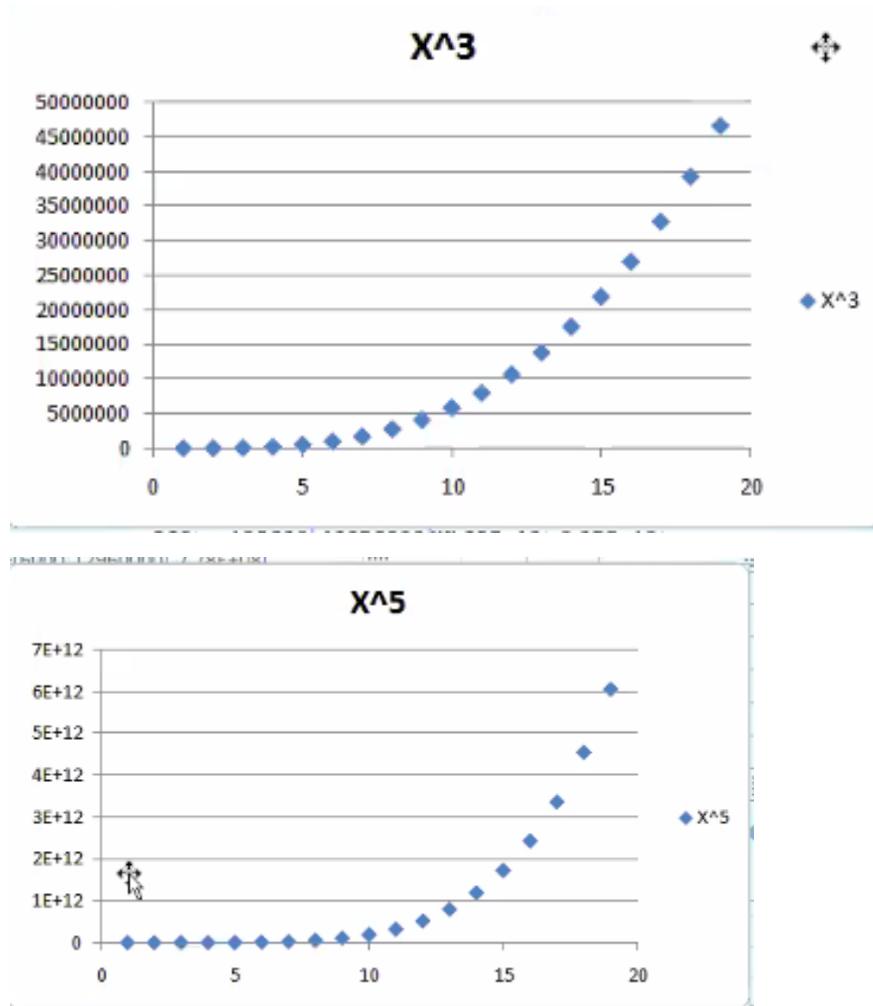
By adding degree to each feature, we can solve it. This will curve the best fit line.



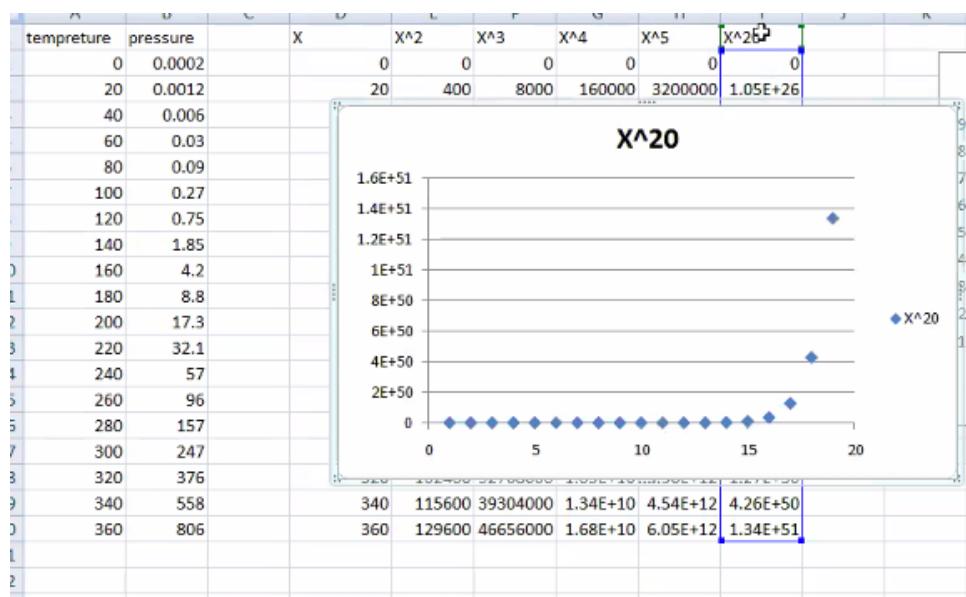
You need to tweet the x value.

tempreture	pressure	X	X^2	X^3	X^4	X^5
0	0.0002		0	0	0	0
20	0.0012		20	400	8000	160000
40	0.006		40	1600	64000	2560000
60	0.03		60	3600	216000	12960000
80	0.09		80	6400	512000	40960000
100	0.27		100	10000	1000000	1E+08
120	0.75		120	14400	1728000	2.07E+08
140	1.85		140	19600	2744000	3.84E+08
160	4.2		160	25600	4096000	6.55E+08
180	8.8		180	32400	5832000	1.05E+09
200	17.3		200	40000	8000000	1.6E+09
220	32.1		220	48400	10648000	2.34E+09
240	57		240	57600	13824000	3.32E+09
260	96		260	67600	17576000	4.57E+09
280	157		280	78400	21952000	6.15E+09
300	247		300	90000	27000000	8.1E+09
320	376		320	102400	32768000	1.05E+10
340	558		340	115600	39304000	1.34E+10
360	806		360	129600	46656000	1.68E+10

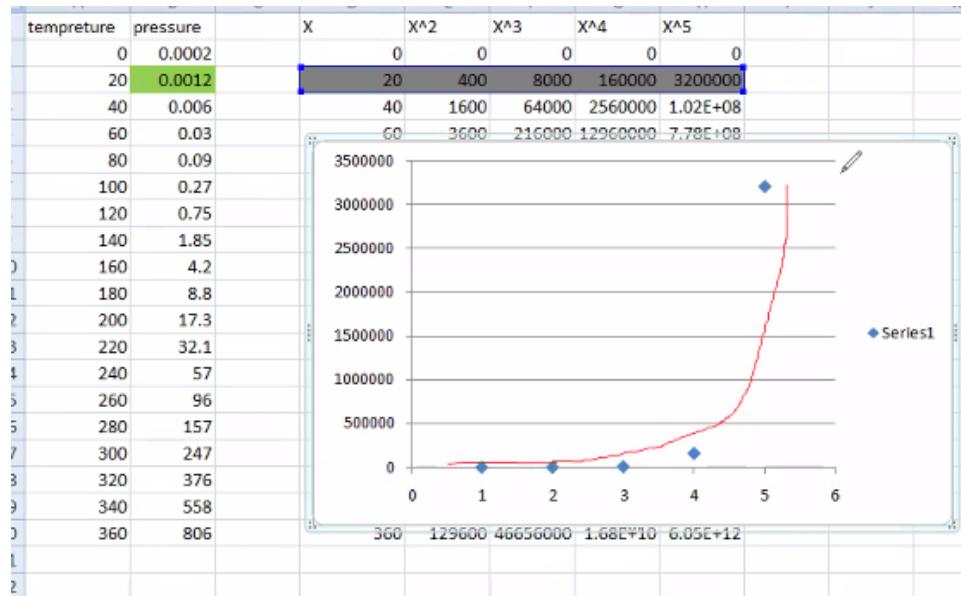




More the degree, better the result. But not too far!



It's a trial-and-error method. Until we try with different degree, we won't know which is best.

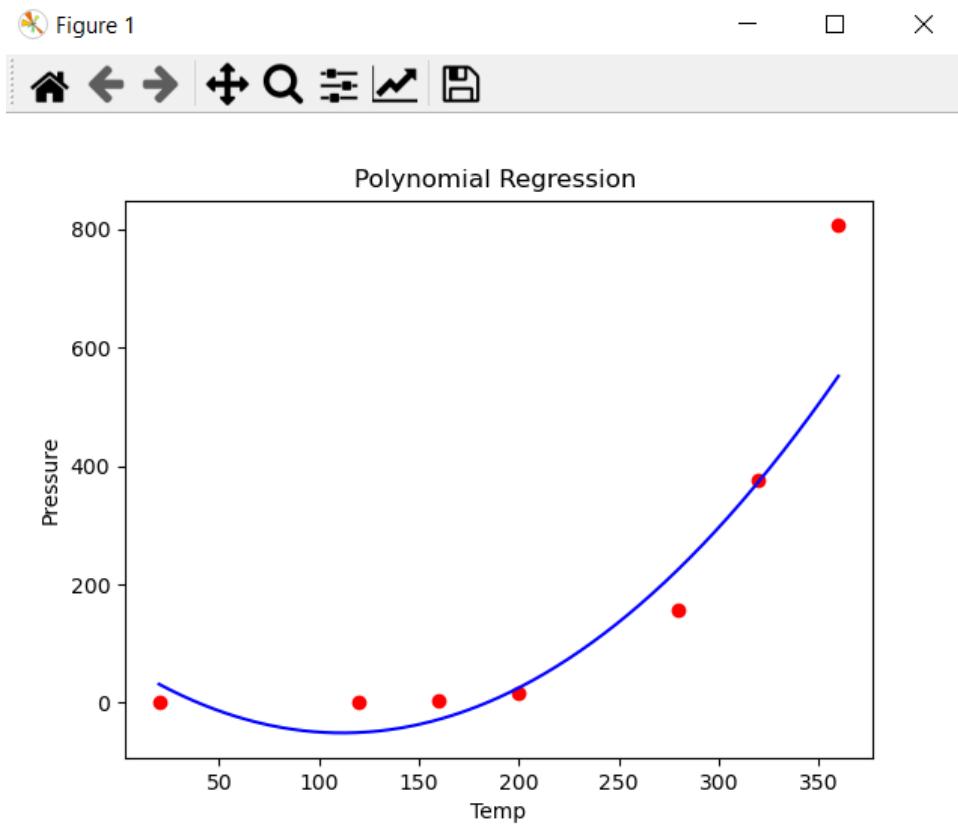


Machine is least interested in column names. Machine is not aware of which column it deals with. It considers it as feature 1,2,3,4, and so on.

tempreture	pressure	x1	x2	x3	x4	x5
0	0.0002	0	0	0	0	0
20	0.0012	20	400	8000	160000	3200000
40	0.006	40	1600	64000	2560000	1.02E+08
60	0.03	60	3600	216000	12960000	7.78E+08
80	0.09	80	6400	512000	40960000	3.28E+09
100	0.27	100	10000	1000000	1E+08	1E+10
120	0.75	120	14400	1728000	2.07E+08	2.49E+10
140	1.85	140	19600	2744000	3.84E+08	5.38E+10
160	4.2	160	25600	4096000	6.55E+08	1.05E+11
180	8.8	180	32400	5832000	1.05E+09	1.89E+11
200	17.3	200	40000	8000000	1.6E+09	3.2E+11
220	32.1	220	48400	10648000	2.34E+09	5.15E+11
240	57	240	57600	13824000	3.32E+09	7.96E+11
260	96	260	67600	17576000	4.57E+09	1.19E+12
280	157	280	78400	21952000	6.15E+09	1.72E+12
300	247	300	90000	27000000	8.1E+09	2.43E+12
320	376	320	102400	32768000	1.05E+10	3.36E+12
340	558	340	115600	39304000	1.34E+10	4.54E+12
360	806	360	129600	46656000	1.68E+10	6.05E+12

Degree 2 output:

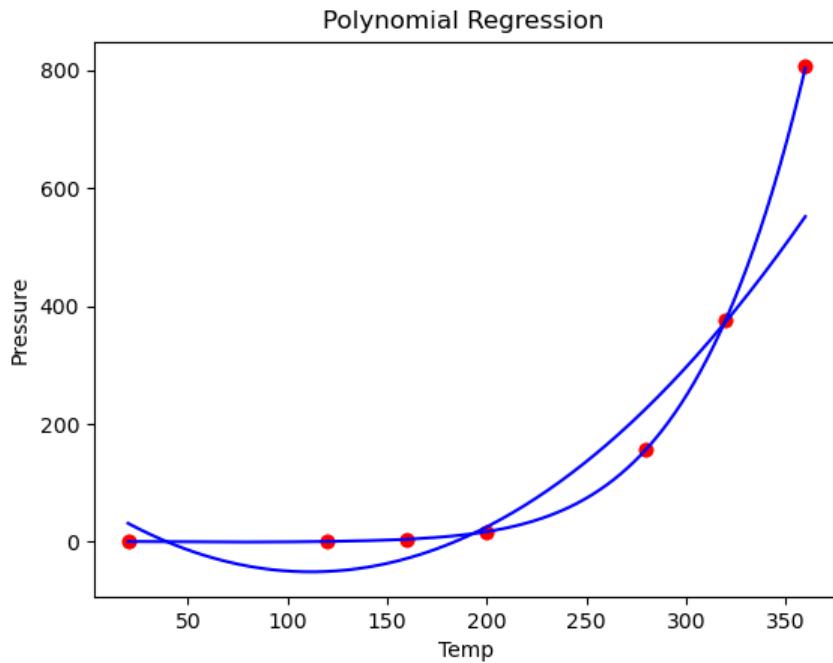
Error is huge here.



Degree 5:

Error is less. It's fits 100%.

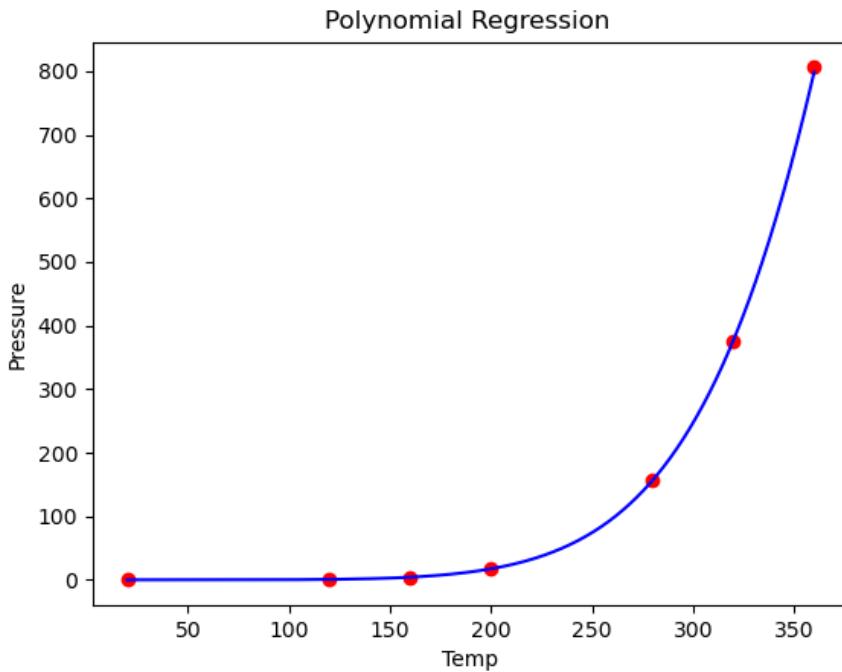
Figure 1



error_2	float64	1	10553.342229990534
error_5	float64	1	0.40559764248728847

Degree 10:

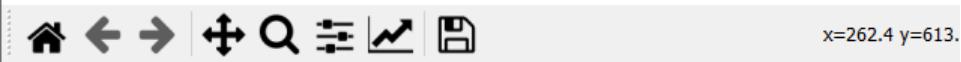
Figure 1



error_2	float64	1	10553.342229990534
error_5	float64	1	0.40559764248728847
error_10	float64	1	5.366908457395362

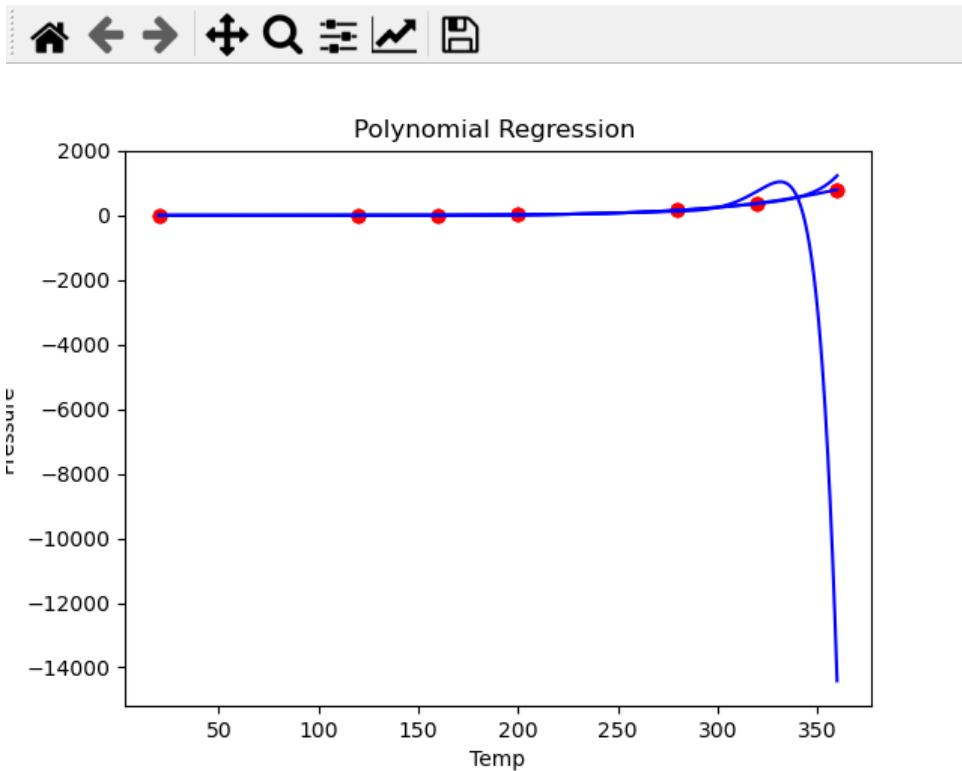
Degree 15:

Figure 1



Degree 20:

Figure 1



			Process
error_2	float64	1	10553.342229990534
error_5	float64	1	0.40559764248728847
error_10	float64	1	5.366908457395362
error_15	float64	1	26561.71204403531
error_20	float64	1	33961969.21293781

What is bias here?

```
from sklearn.preprocessing import PolynomialFeatures
poly_Obj = PolynomialFeatures(degree=20 , include_bias=False)
Xtrain_poly = poly_Obj.fit_transform(X_train)
Xtest_poly = poly_Obj.fit_transform(X_test)
```

We already discussed about the intercept. It helps to get the curve close to the spread of datapoints.

We can manually insert the intercept to keep it close to the spread of the information instead of machine calculating.

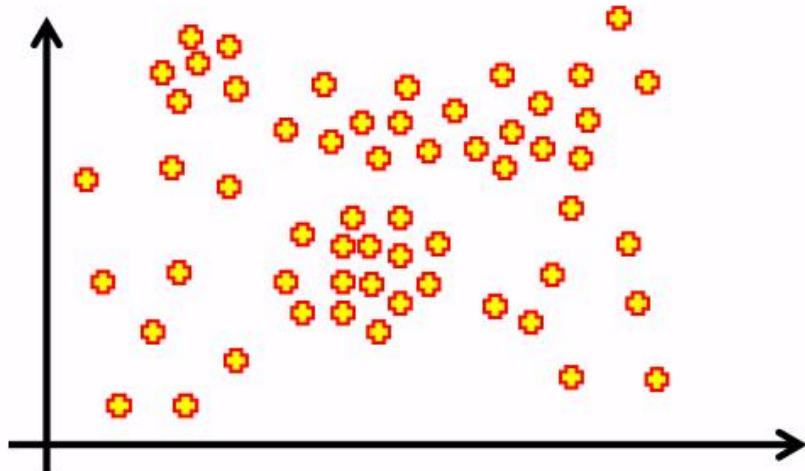
Feeding the intercept manually helps or not – range 0 to 1 will be good

Sometimes polynomial won't help. Go for Decision tree and Random Forest.

#### **Decision tree regression:**

Forms in a tree structure. Breaks down into smaller and smaller subsets.

## Decision Tree Regression

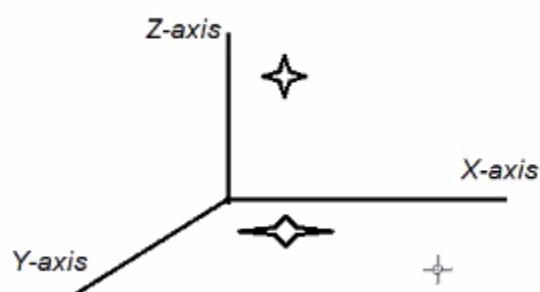
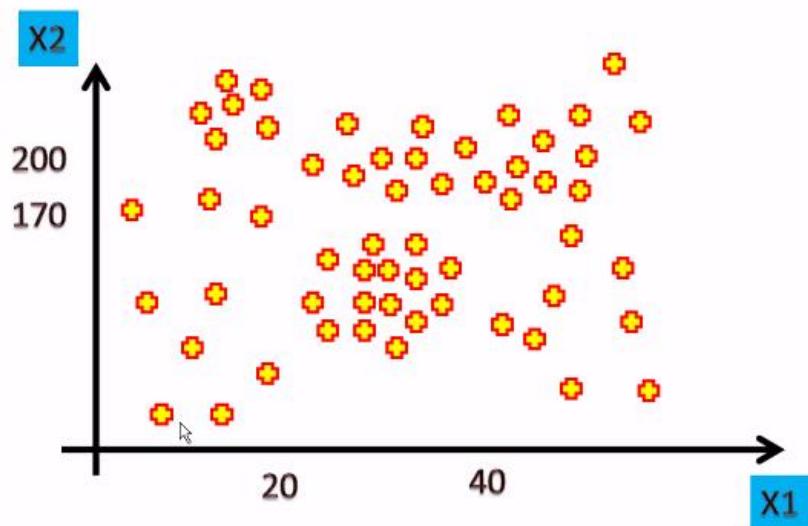


## Decision Tree Regression

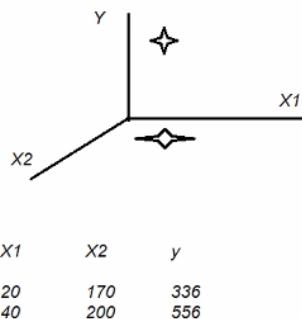
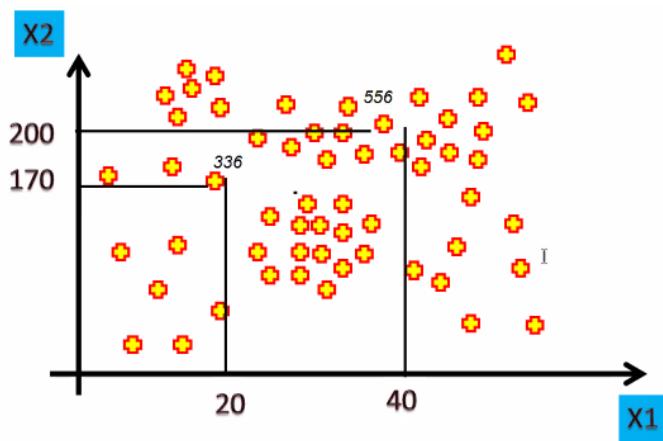
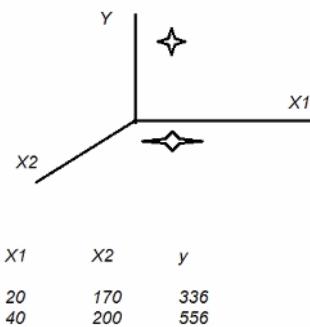
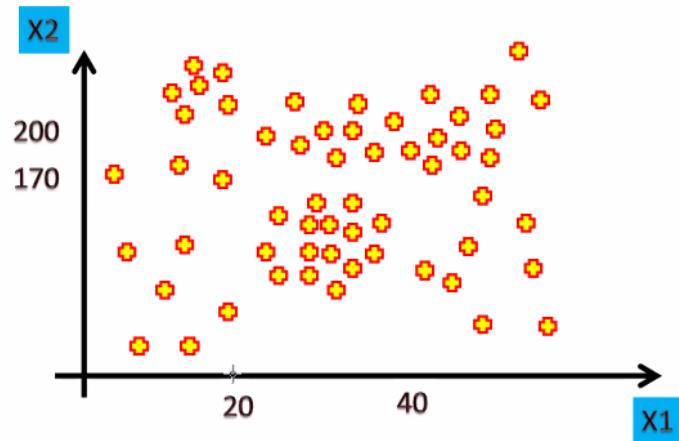
Decision tree builds regression or classification models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed.

3 D Plotting looks like:

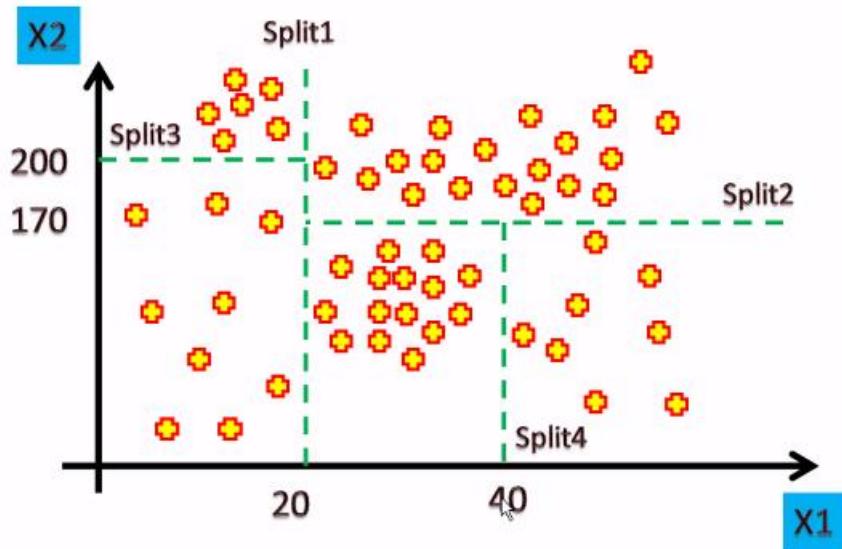
# Decision Tree Regression



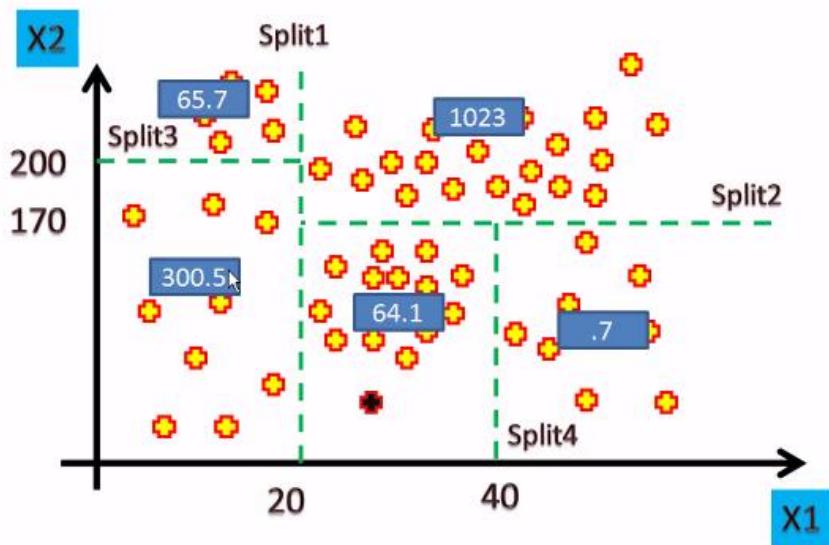
<i>input</i>		<i>output</i>
$X_1$	$X_2$	$y$
20	170	336
40	200	556



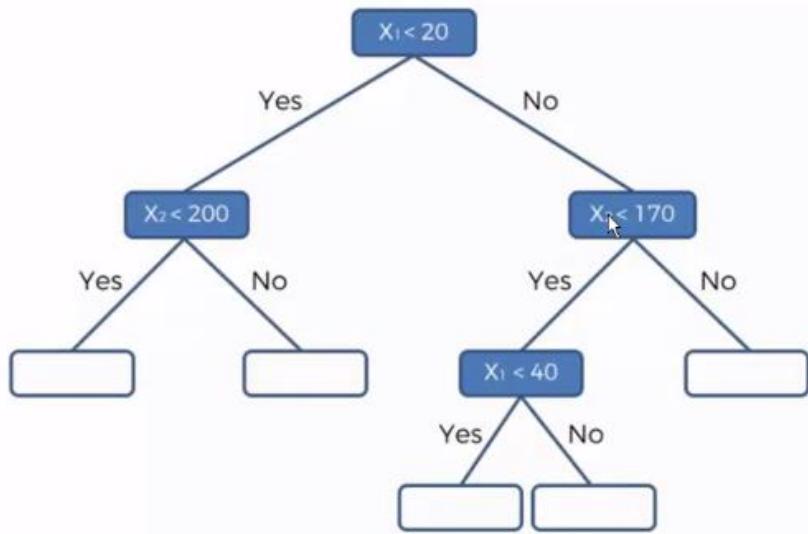
## Decision Tree Regression



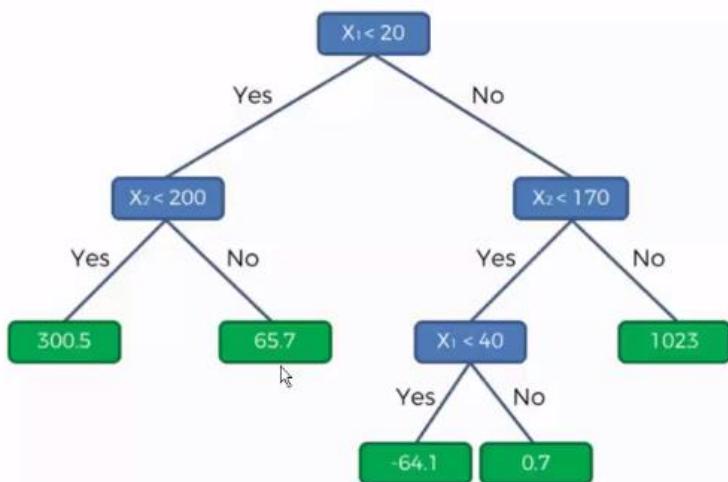
## Decision Tree Regression



# Decision Tree Regression



# Decision Tree Regression

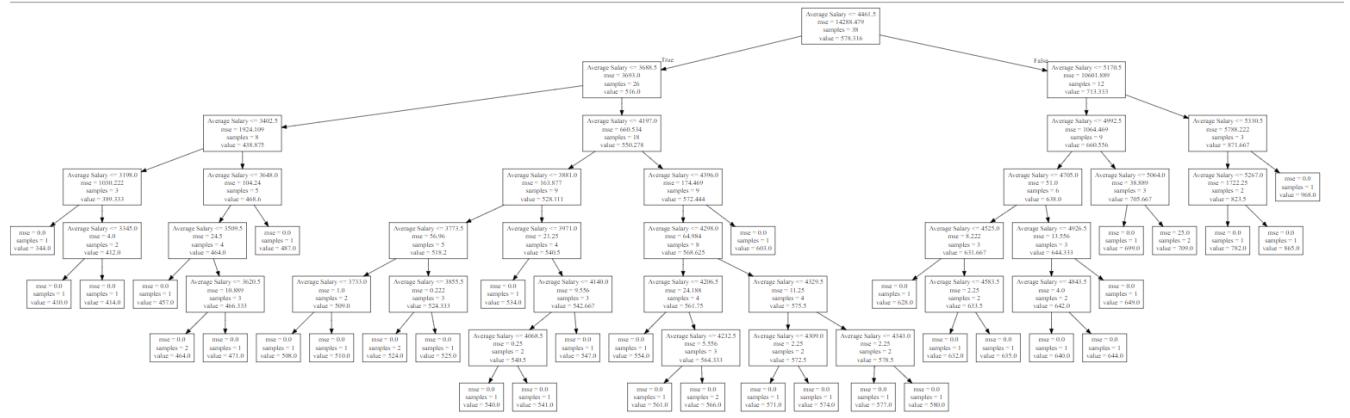


- Iterative in nature and time consuming
- Handle the complexity better than the polynomial regression.
- Good solution
- Performance – takes extra time.
- We can see the actual tree as well
- You will come here only when the polynomial regression is not good.

- Simple is better
- Prediction is just substituting values into the formula and getting it.
- Decision tree is good when we have good amount of data.
- Less data – not good to go with decision tree
- If the points falls under the split region, then the predicted value will be the average of that split region.

Web Graphviz: To visualize the tree structure

<http://www.webgraphviz.com/>



There will be overfitting problem! More splitting is done here.

Splitting has done in a granular level.

### Overfitting:

Machine learns much on the training data and does only good only on training data.

It lacks generalization capabilities

Gives high error/ low accuracy on testing data

Multiple test sets tried against such model will lead to unstable error or accuracy and it won't be able to give appropriate results.

**Can we control the number of splits? Whether to split or not depends on the value of N?**

Value of N is considered for taking this decision.

Min\_samples\_split = n Default set to 2

If we have 2 or more data points in the split we divide it further.

```

from sklearn.tree import DecisionTreeRegressor
regressor = DecisionTreeRegressor()
regressor.fit(X_train,y_train)
# =====#
# Model Testing
# =====#
y_pred = regressor.predict(X_test)
from sklearn.metrics import mean_squared_error
error_DT = mean_squared_error(y_te
# =====#
# Tree plotting
# =====#

```

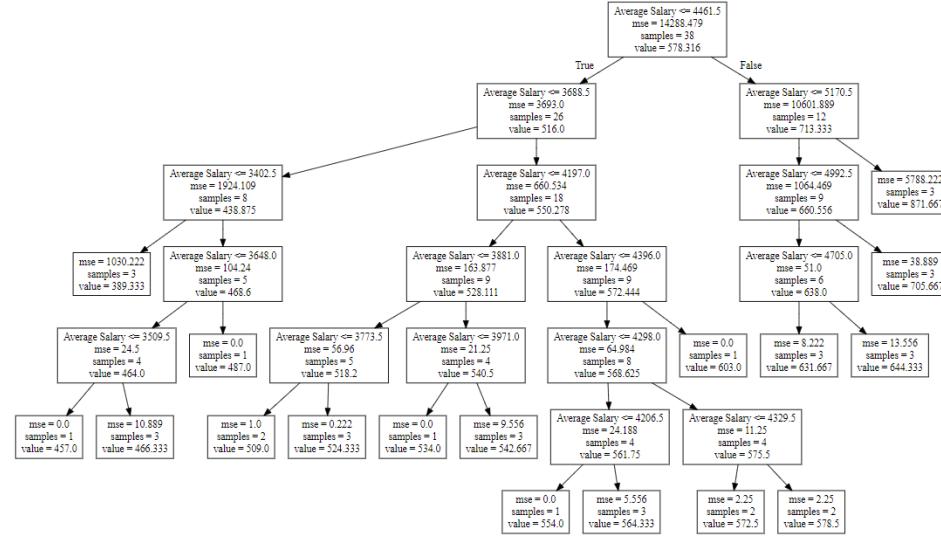
DecisionTreeRegressor(criterion="mse",  
splitter="best",  
max\_depth=None,  
**min\_samples\_split=2,**  
min\_samples\_leaf=1,  
min\_weight\_fraction\_leaf=0.,  
max\_features=None,  
random\_state=None,  
max\_leaf\_nodes=None,  
min\_impurity\_decrease=0.,  
min\_impurity\_split=None,  
presort='deprecated',  
ccp\_alpha=0.0)

### Should we split or confirm it as leaf node?

In the above case, it has spitted more. Results in overfitting!

Now, provide the min sample split as 4

error_DT	float64	1	39.9
error_DT4	float64	1	46.17222222222222



If the split value is more, then it results in **underfitting! (Too crisp and too much generalization)**

Step 1: Train model with min sample slit = n

## Step 2: Test sample      Error

Sample 1	39.0
Sample 2	45.0
Sample 3	47
Sample 4	46.0

underfitting

Not learning Enough  
Too much of generalization the model is learning to do  
testing error will be very high / acc will be very low

Overfitting

Machine learns too much on training data  
Does only good on training data  
it lacks generalization capabilities  
gives high error / low accuracy on testing data  
multiple test sets tried against such model will lead to unstable error or accuracy

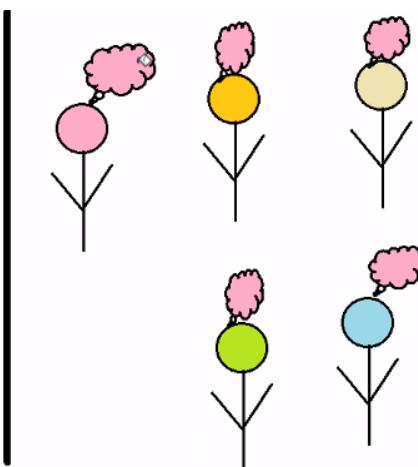
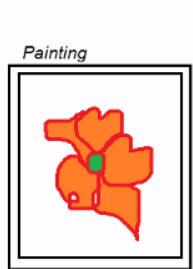
## Good fit

- Enough learning
- Has good generalization capability
- Testing error low / accuracy high, and stable

**Note:** The error should be stable. You should be able to reproduce the same output. It should do good in all machines. So, stability check is a must!

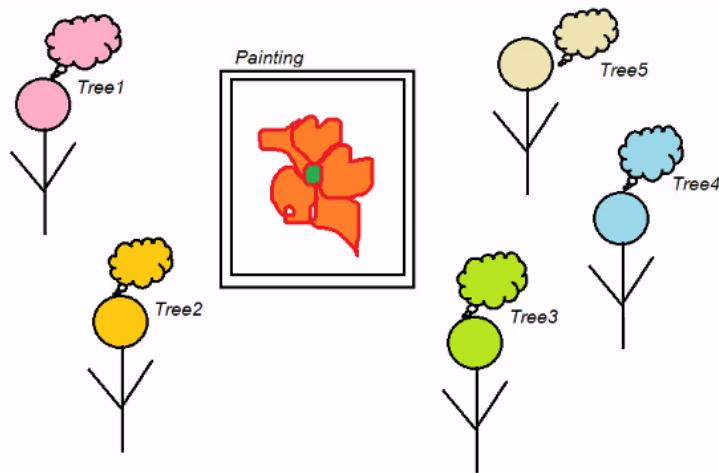
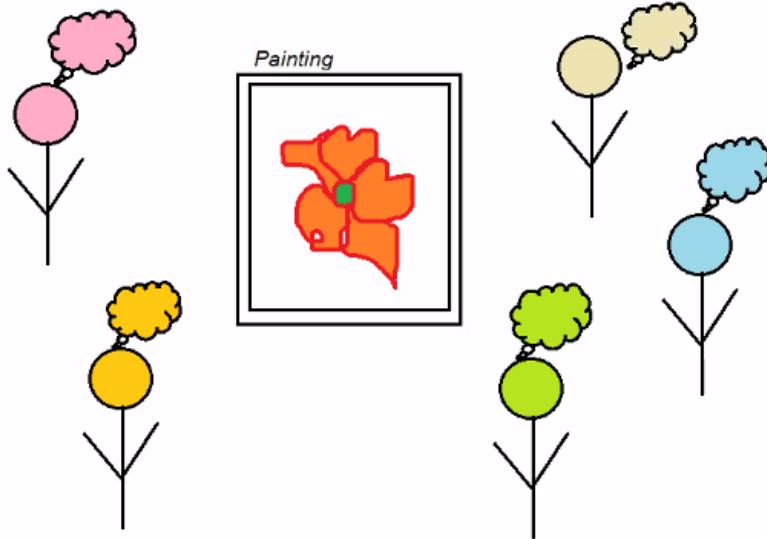
## Random Forest:

One tree is one perspective of looking at the data. That is the problem -> Inappropriate results



Now everyone is going inside the art gallery and having their own perspective. Better idea about the painting. Multiple perspective out of same painting -> Random Forest.

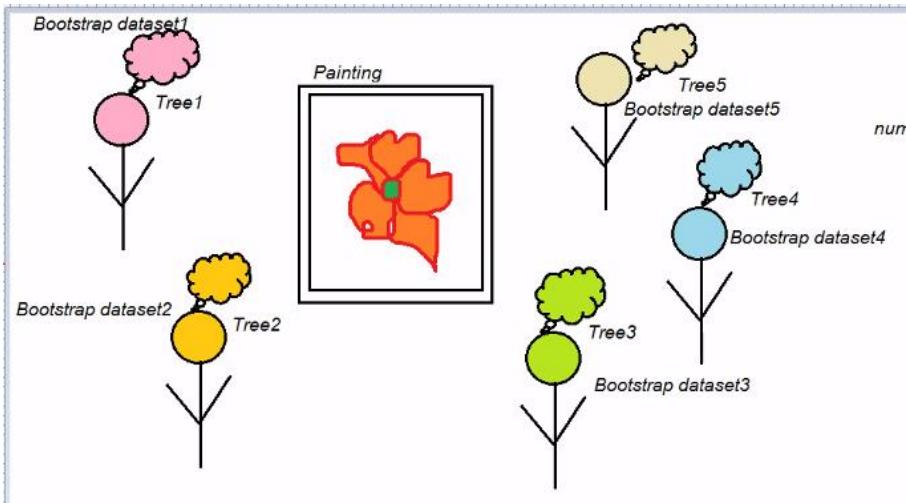
It decides from multiple trees.



Overall Dataset will be same

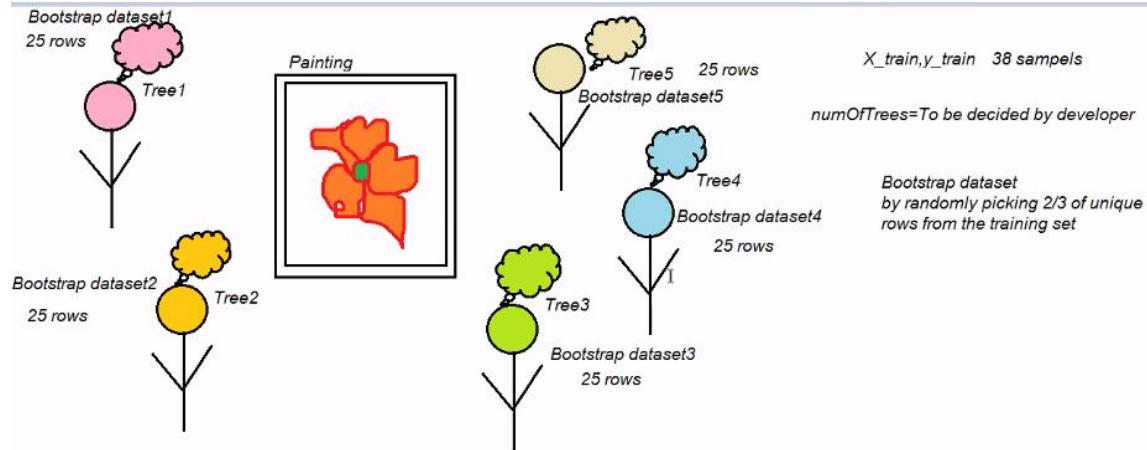
**Number of trees = to be decided by developer**

Each tree will have its own bootstrap dataset



### How is the bootstrap dataset is constructed?

- Method to split the dataset: Bootstrap.
- Without replacement and its pure random
- By randomly picking 2/3 of unique rows from the training dataset.
- Each tree is different as it picks up different bootstrap dataset.



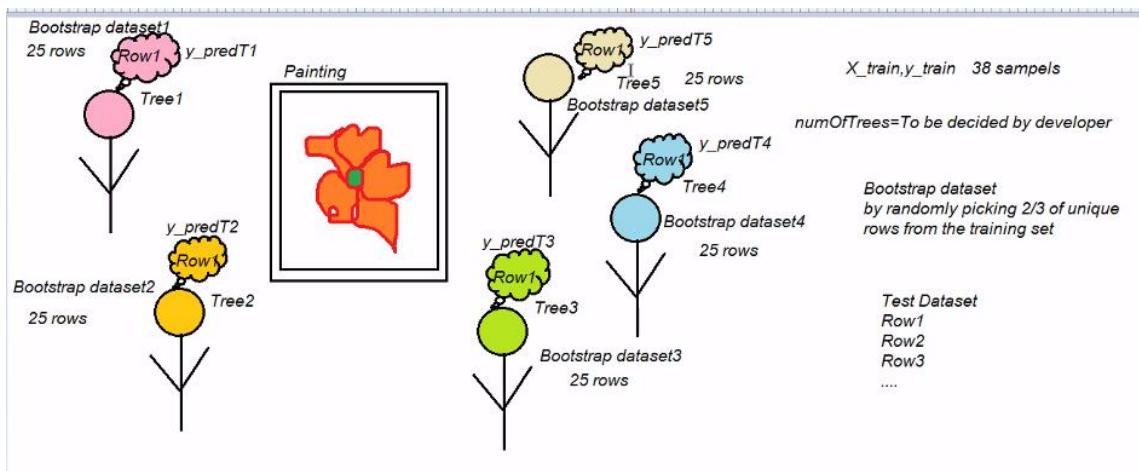
### Test Dataset:

Row 1

Row 2

Row 3

.....



## What am I going to do with all these predictions?

## *Test Dataset*

Row1 = y\_pred = avg(y\_predT1,y\_predT2,y\_predT3,y\_predT4,y\_predT5)

Row2 = y\_pred = avg(y\_predT1,y\_predT2,y\_predT3,y\_predT4,y\_predT5)

Row3 =  $y\_pred = \text{avg}(y\_predT1, y\_predT2, y\_predT3, y\_predT4, y\_predT5)$

卷之三

NumOfTrees = To be decided by developer

RandomForest Tree Creation (X\_train, y\_train) 38 samples

**Step1:** => BootStrap Dataset Creation i.e. Random Row selection from the given training set

=> BootStrap dataset will be SAME SIZE as Training Dataset

e.g. X\_train,y\_train 38 samples

=> But Bootstrap Dataset doesn't use all 38 values in the training set to create the BootStrap dataset  
it just uses 2/3 of the total values in the training set to create the BootStrap dataset

I it just uses 2/3 of the total values in the training set to create the BootStrap dataset  
2/3(training set) -----> 66% rows, Bagging, remainder is called out of Bag Record

=> BootStrap dataset after taking 2/3 unique rows the remainder rows are filledup to match the training set size by randomizing rows of 2/3 taken

### RandomForest (10 trees)

Tree1 BootStrap Dataset1  
Tree2 BootStrap Dataset2  
Tree3 BootStrap Dataset3  
Tree4 BootStrap Dataset4

..... Tree10 BootStrap Dataset10

### Training Sample

```
Row1    pattern1
Row2    pattern2
Row3    pattern1
Row4    pattern1
Row5    pattern2
Row6    pattern1
Row7    pattern3
Row8    pattern2
Row9    pattern2
Row10   pattern1
Row11   pattern2
```

### BootStrap Dataset1

```
Row1    pattern1
Row2    pattern2
Row3    pattern1
Row4    pattern1
Row5    pattern2
Row6    pattern1
Row7    pattern3
Row7    pattern3
Row7    pattern3
Row6    pattern1
Row6    pattern1
```

It's a random pick. For understanding, done as systematic pick.

2/3 volume is the pickup volume.

### Training Sample

Row1	pattern1
Row2	pattern2
Row3	pattern1
Row4	pattern1
Row5	pattern2
Row6	pattern1
Row7	pattern3
Row8	pattern2
Row9	pattern2
Row10	pattern1
Row11	pattern2

### BootStrap Dataset1

Row1	pattern1
Row2	pattern2
Row3	pattern1
Row4	pattern1
Row5	pattern2
Row6	pattern1
Row7	pattern3
Row7	pattern3
Row6	pattern1
Row6	pattern1

Randomly pick up from bootstrap dataset to fill the gap.

Bootstrap and training dataset size should be same.

### BootStrap Dataset1

Row1	pattern1
Row2	pattern2
Row3	pattern1
Row4	pattern1
Row5	pattern2
Row6	pattern1
Row7	pattern3
Row7	pattern3
Row6	pattern1
Row6	pattern1

+

Randomness won't be effective after certain number of trees.

The screenshot shows the Spyder Python IDE interface with two panes. The left pane displays the code for a Random Forest Regression model, while the right pane shows the execution results and a variable explorer.

**Code (Spyder (Python 3.8)):**

```
9 # Model Implementation
10 # Fitting Random Forest Regressor
11 # DT Error 39.9
12 #
13 from sklearn.ensemble import RandomForestRegressor
14 regressor = RandomForestRegressor()
15 regressor.fit(X_train,y_train)
16
17 # =====
18 # Model Testing
19 #
20 y_pred = regressor.predict(X_test)
21
22 from sklearn.metrics import mean_squared_error
23
24
25 error_RF100 = mean_squared_error(y
26 #error_RF200 = mean_squared_error(
27 #error_RF300 = mean_squared_error(
28 #error_RF500 = mean_squared_error(
29
```

**Execution Results (Arguments and Variable Explorer):**

**Arguments:**

```
RandomForestRegressor(n_estimators=100,
                     criterion="mse",
                     max_depth=None,
                     min_samples_split=2,
                     min_samples_leaf=1,
                     min_weight_fraction_leaf=0.,
                     max_features="auto",
                     max_leaf_nodes=None,
                     min_impurity_decrease=0.,
                     min_impurity_split=None,
                     bootstrap=True,
                     oob_score=False,
                     n_jobs=None,
                     random_state=None,
                     verbose=0,
                     warm_start=False,
```

**Variable Explorer:**

Name	Type	Size	Value
dataset	DataFrame	(48, 2)	Column names: Average_Income, Petrol_Consumption
error_RF100	float64	1	23.117580000000002
regressor	ensemble._forest.RandomForestRegressor	100	RandomForestRegressor object of sklearn.ensemble._forest module
X	Array of int64	(48, 2)	[1060, 13331]
X_test	Array of int64	(18, 2)	[3528, 37321]
X_train	Array of int64	(30, 2)	[3721, 4447]
y	Array of int64	(48,)	[344 418 414 ... 782 865 988]
y_pred	Array of float64	(18,)	[381.11 463.74 577.11 596.22 605.79 638.89 626.88 648.67 497.96 504.05
y_test	Array of int64	(18,)	[387 468 577 591 610 640 631 648 467 498]

error_RF100	float64	1	23.11758000000002
error_RF200	float64	1	22.971057500000065
error_RF300	float64	1	23.58203666666742
error_RF500	float64	1	23.001594500000017

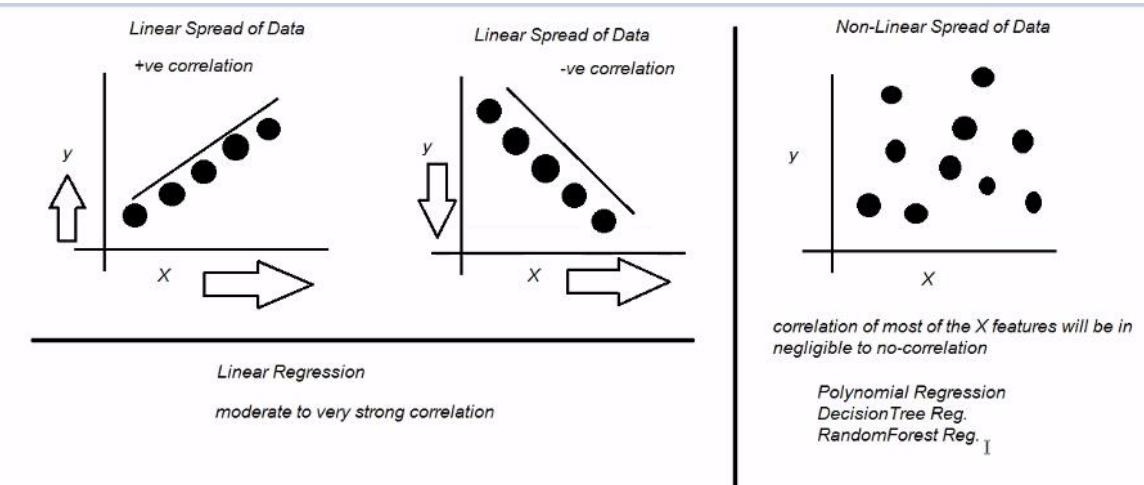
## **Observation:**

Exhausted the error value after 100. You can see the same similarity.

**Try it between 120 to 150:**

error_RF100	float64	1	23.11758000000002
error_RF120	float64	1	21.793486111111214
error_RF150	float64	1	21.89207555555556
error_RF200	float64	1	22.971057500000065
error_RF300	float64	1	23.582036666666742
error_RF500	float64	1	23.001594500000017

## Regression overview:



## Which regression supports what kind of spread of data?

SUPERVISED LEARNING	Input Data & Output Data	Remember & Generalize
Regression	Value to Predict is Continuous	
Linear Regression Polynomial Regression DecisionTree Regressor RandomForest Regressor	Linear spread of data non-Linear spread of data non-Linear spread of data non-Linear spread of data	
Classification	Value to Predict is Discrete	
Logistic Regression Support Vector Machine DecisionTree Classifier RandomForest Classifier Naive Bayes		

## How to decide on which regression to pick up?

If most of the X features (80%) are moderate to very strongly correlated to Y feature  
Then Choose LINEAR REGRESSION

If most of the X features (80%) are no correlation to negligible correlation to Y feature  
Then Choose Poly Reg. , DT , RF

## R2 score:

Tell us how much variation in y that can be explained by all x features.

# R-squared

- R-Squared is the proportion of variation in the dependent (response) variable that has been explained by the model.

# R-squared

- R-Squared is the proportion of variation in the dependent (response) variable that has been explained by the model.

Also known as coefficient of determination, it tells us how much is the variation in the dependent variable (salary) can be explained by the independent variable (Experience)

R2 score of 80 and above is healthy.

```
1 from sklearn.metrics import r2_score
2 score = r2_score(y_test,y_pred)
3
```

Simple regression:

```
y_pred = regressor.predict(X_test)

from sklearn.metrics import mean_squared_error
error = mean_squared_error(y_test,y_pred)

errorUnits = np.sqrt(error) # root mean squared error

from sklearn.metrics import r2_score
score = r2_score(y_test,y_pred)

# =====
# Add-On Visualization
# =====

plt.scatter(X_test,y_test)
plt.plot(X_test,y_pred,c="Red")
plt.show()
```

X_train	float64	(20, 1)	[2.9]
dataset	DataFrame	(30, 2)	Column names: YearsOfExp
error	float64	1	21026037.329511296
errorUnits	float64	1	4585.4157204675885
score	float64	1	0.9749154407708353
y	int64	(30,)	[ 39343 46205 37731 ...
y_pred	float64	(10,)	[ 40835.10590871 123079...

Variable explorer File explorer Help

IPython console

Console 1/A Console 2/A Console 3/A

```
...: errorUnits = np.sqrt(error) # root mean
squared error
...
Slope / b1 [9345.94244312]
Intercept / b0 26816.19224403119
```

## Multiple regression:

```

18 # =====
19
20 y_pred = regressor.predict(X_test)
21
22 from sklearn.metrics import mean_squared_error
23 error = mean_squared_error(y_test,y_pred)
24 errorUnits = np.sqrt(error)
25
26 from sklearn.metrics import r2_score
27 score = r2_score(y_test,y_pred)
28
29
30
31

```

Python console output:

error	float64	1	83502864.03257737
errorUnits	float64	1	9137.990152794944
score	float64	1	0.9347068473282425
y	float64	(50,)	[192261.83 191792.86 19105.
v_pred	float64	(10,)	[103015.20159796 132582.21

Pearsonr score:

```

datasetEval = pd.read_csv("02Companies.csv")

datasetEval = pd.get_dummies(data=datasetEval, columns=["State"], drop_first=True)

from scipy.stats import pearsonr

pearsonr(datasetEval['RND Spend'], datasetEval['Profit'])
pearsonr(datasetEval['Administration'], datasetEval['Profit'])
pearsonr(datasetEval['Marketing Spend'], datasetEval['Profit'])
pearsonr(datasetEval['State_Florida'], datasetEval['Profit'])
pearsonr(datasetEval['State_New_York'], datasetEval['Profit'])

```

Taking all features vs taking specific features:

```

MultLinearReg.py | SimpleLinearReg.py | MultLinearReg_Remplement.py
12 # =====
13
14 y_pred = regressor.predict(X_test)
15
16 from sklearn.metrics import mean_squared_error
17 error = mean_squared_error(y_test,y_pred)    # old 83502864.03257737
18                      # new 67220275.37568115
19 errorUnits = np.sqrt(error)      # old 9137.990152794944
20                      # new 8198.797190788484
21
22 from sklearn.metrics import r2_score
23 score = r2_score(y_test,y_pred) # old 0.9347068473282425
24                      # new 0.947438644768489
25
26
27
28

```

Comparison of Linear Regression and Multiple regression:

error_LR	float64	1	556.0375904395161
error_RF	float64	1	21.8920755555556

## Classification:

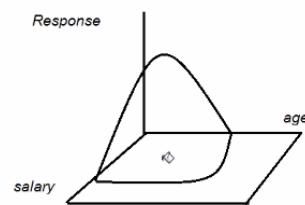
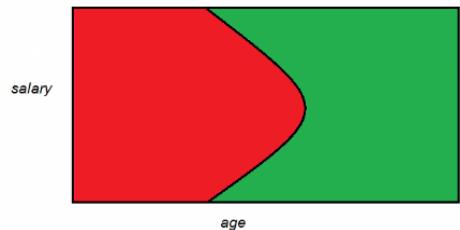
Use cases:

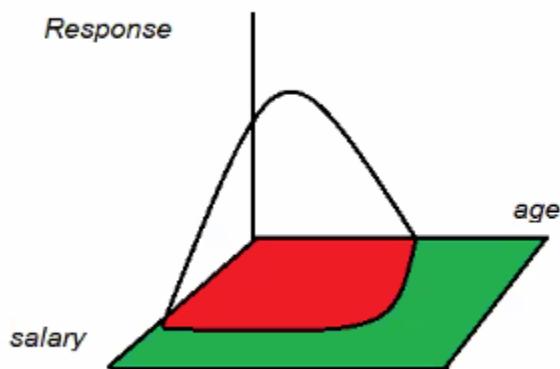
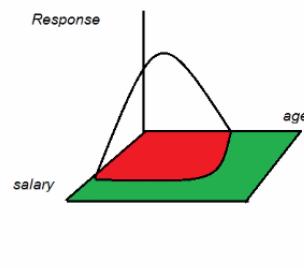
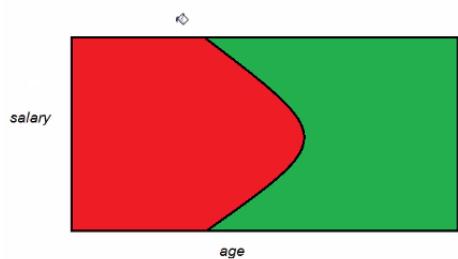
User ID	Gender	Age	EstimatedSalary	Purchased
15624510	Male	19	19000	0
15810944	Male	35	20000	0
15668575	Female	26	43000	0
15603246	Female	27	57000	0
15804002	Male	19	76000	0
15728773	Male	27	58000	0
15598044	Female	27	84000	0
15694829	Female	32	150000	1
15600575	Male	25	33000	0
15727311	Female	35	65000	0
15570769	Female	26	80000	0
15606274	Female	26	52000	0
15746139	Male	20	86000	0
15704987	Male	32	18000	0
15628972	Male	18	82000	0
156007000	Male	20	80000	0

Discrete Value Prediction, Measure of model performance is Accuracy percentage

Age	ActualResponse	PredictedResponse	6/9*100=66% accuracy score
20	0	0	
30	1	1	
21	0	0	
35	1	0	
40	1	1	
45	1	1	
50	0	1	
18	0	0	
19	1	0	

How am I going to generalize here?





For Classification, accuracy is in the range 0 to 1. We can't take MSE here.

Consider multiple test sets and get accuracy.

Do we have a good training set?

K-fold validation – It is used for Evaluation.

Statistical calculation.

**Classification      Value to Predict is Discrete**

- Logistic Regression
- Support Vector Machine
- DecisionTree Classifier
- RandomForest Classifier
- Naive Bayes

**SVM: Support Vector Machine**

Dataset:

dataset - DataFrame

Index	User ID	Gender	Age	EstimatedSalary	Purchased
0	15624510	Male	19	19000	0
1	15810944	Male	35	20000	0
2	15668575	Female	26	43000	0
3	15603246	Female	27	57000	0
4	15804002	Male	19	76000	0
5	15728773	Male	27	58000	0
6	15598044	Female	27	84000	0
7	15694829	Female	32	150000	1
8	15600575	Male	25	33000	0
...	156207244	Female	25	65000	0

Scaling:

Age	Salary	BankBalance	CreditScore	OwnHouse
30	24000	50000	650	1

Age	Salary	BankBalance	CreditScore	OwnHouse
30	24000	50000	650	1

MORE  
Less

Bank balance: Contributes more in calculation

Own House: Contributes less in calculation

Each has its own scale

Age	Salary	BankBalance	CreditScore	OwnHouse
30	24000	50000	650	1

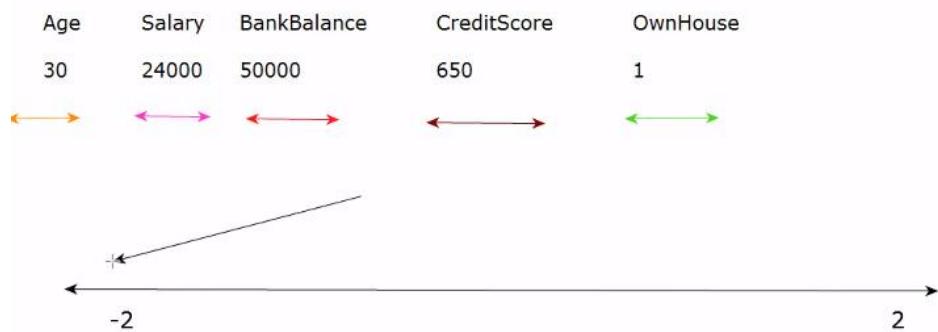
Machine doesn't understand its scale.

So, get everything under one scale

**Take example:** Common scale – range -2 to 2

To take balanced data

To avoid the imbalance



There are scalar available to perform this operations.

Standard scaler – good if your dataset doesn't have many outlier

Min and max scalar – good to use when your dataset has more outlier

Max absolute scalar

Robust scalar

Most popular standard, and min and max scalar

Standardization StandardScaler      *(if the dataset does not have too many outliers)*

$$X_{\text{stand}} = \frac{X - \text{mean}(X)}{\text{standard deviation } (x)}$$

MinMaxScaler      *(if the dataset has too many outliers)*

$$X_{sc} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}.$$

+

We won't have scaling in regression as we didn't face imbalance situation there.

```
In [6]: np.min(X_train)  
Out[6]: -1.9931891594584856
```

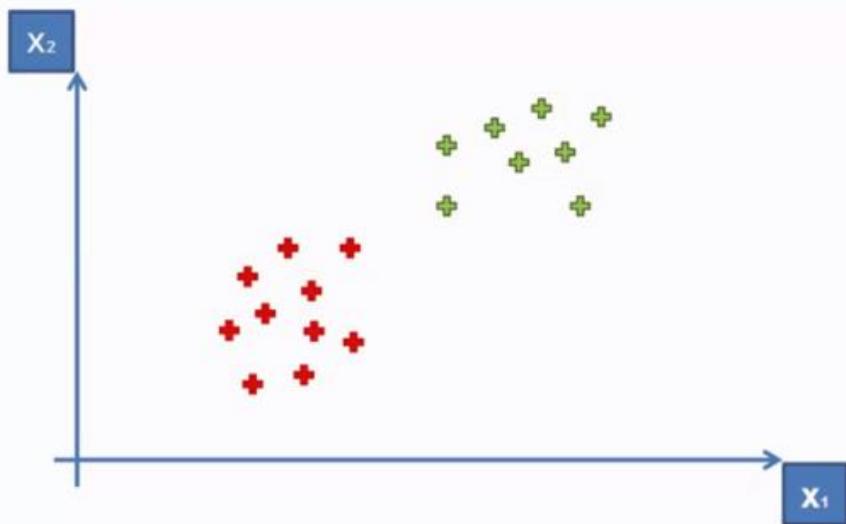
```
In [7]: np.max(X_train)  
Out[7]: 2.3315320031817324
```

```
In [8]: |
```

Pretty balanced! Not too low or high.

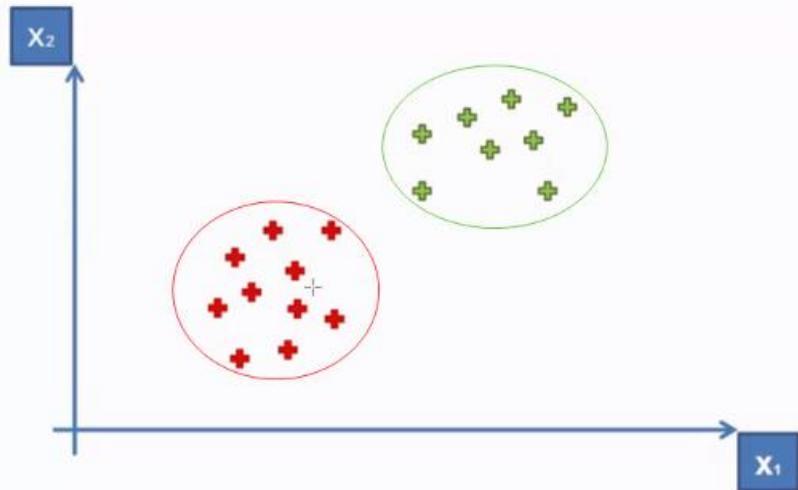
SVM:

# Support Vector Machine



Learn by looking into the similarity and dissimilarity of the data.

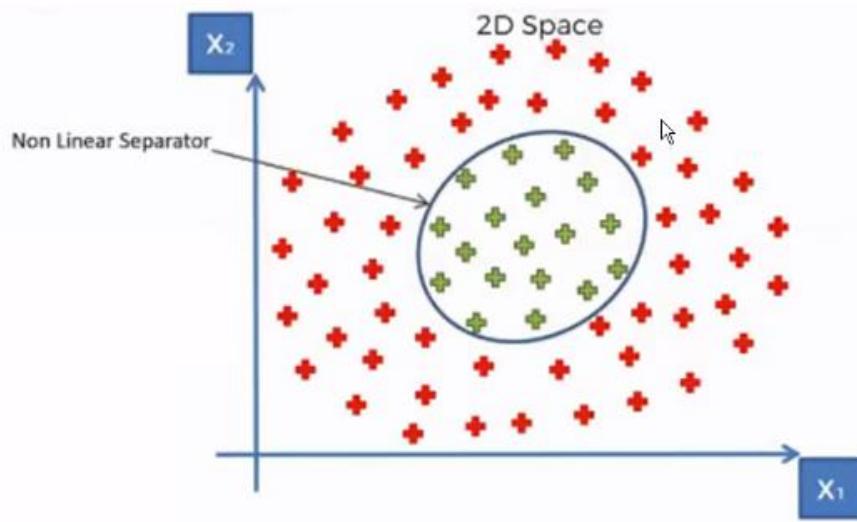
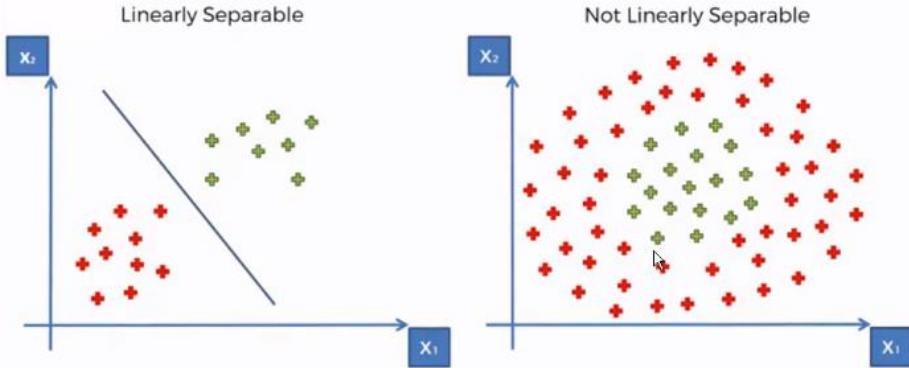
# Support Vector Machine



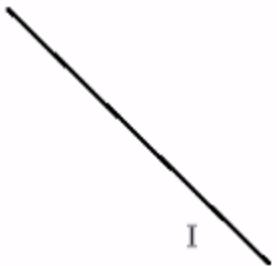
At boundaries, you can see data will not be similar.

We will draw the hyperplane

# Kernel SVM



**Simple problem:**



*kernel=linear*

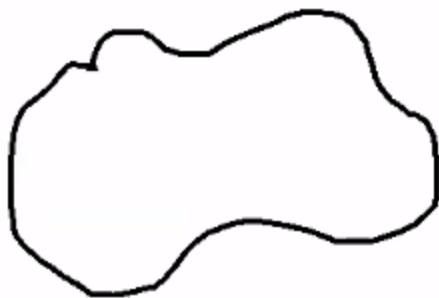
For complex problem: poly kernel



*kernel=poly* +  
*degree = ?*

Didn't serve long as data is more complex in real world

Radial basis function



*kernel=rbf*

*radial basis function*

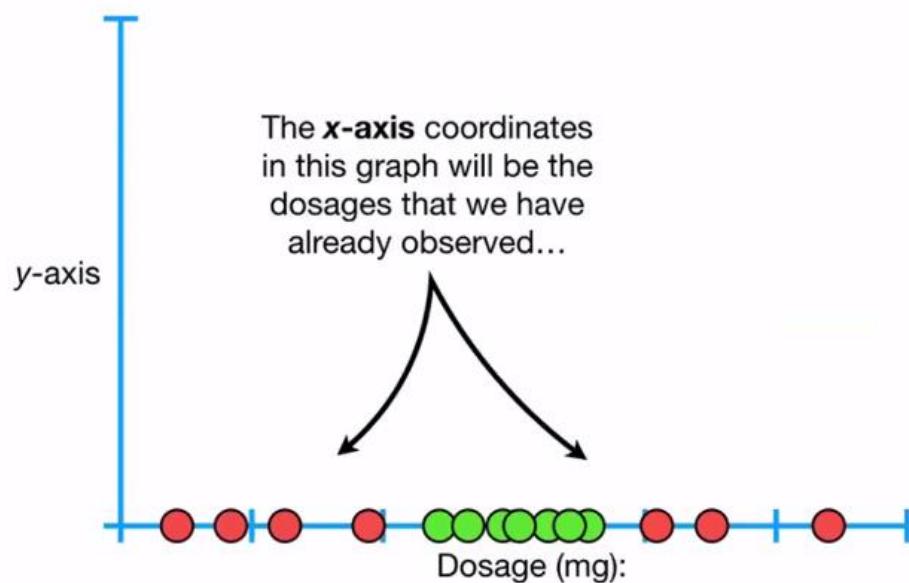
What is the math behind them?

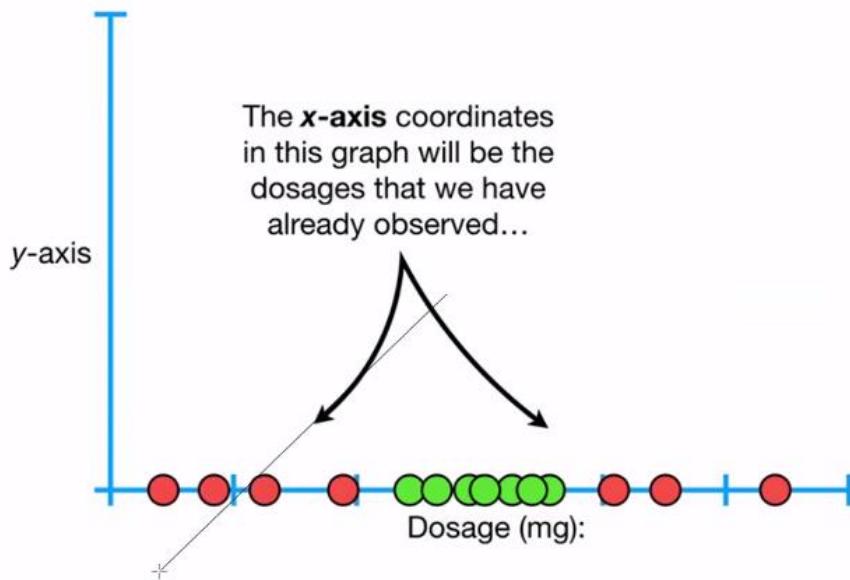
Example: Doctor prescribed the medicines with different dosages

Green – cured

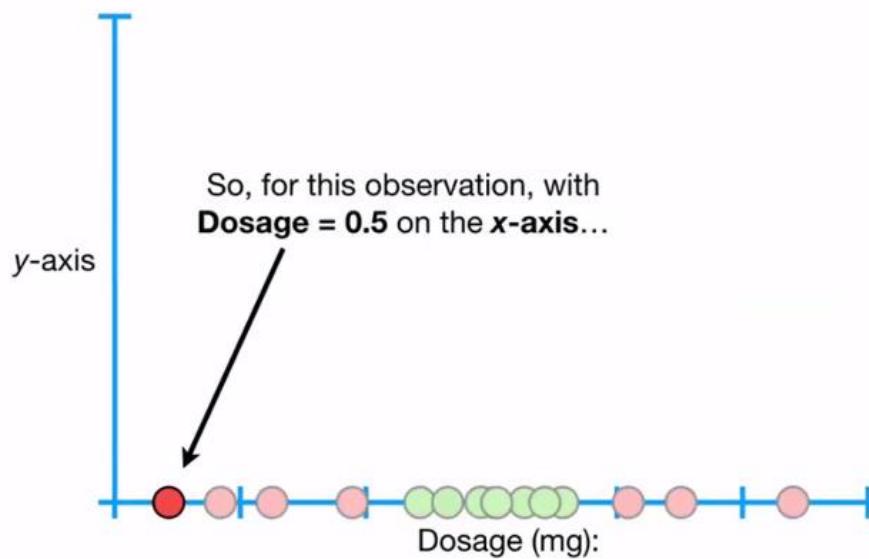
Red – not cured.

We need to draw the best fit line

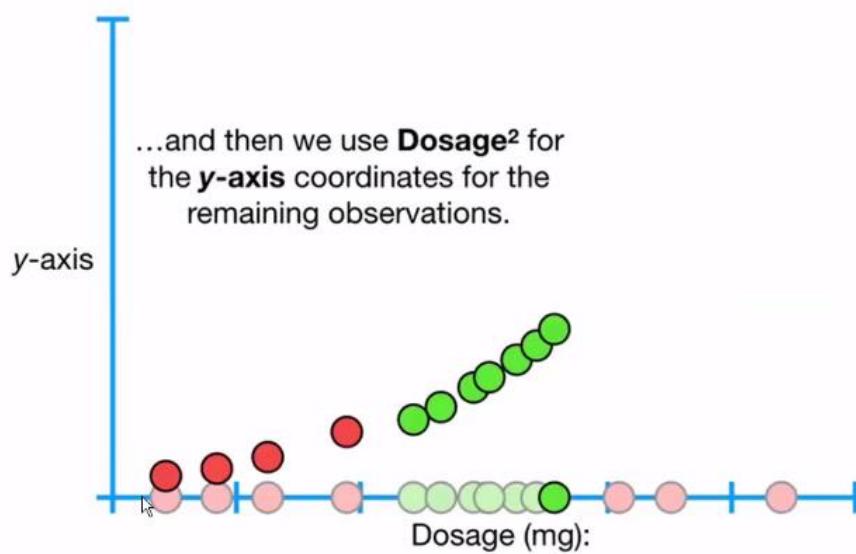




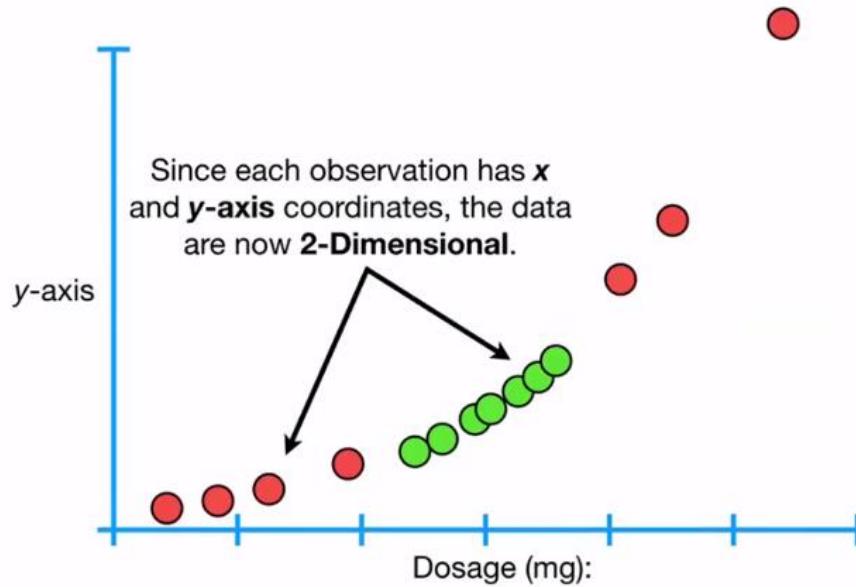
Pull into higher dimension

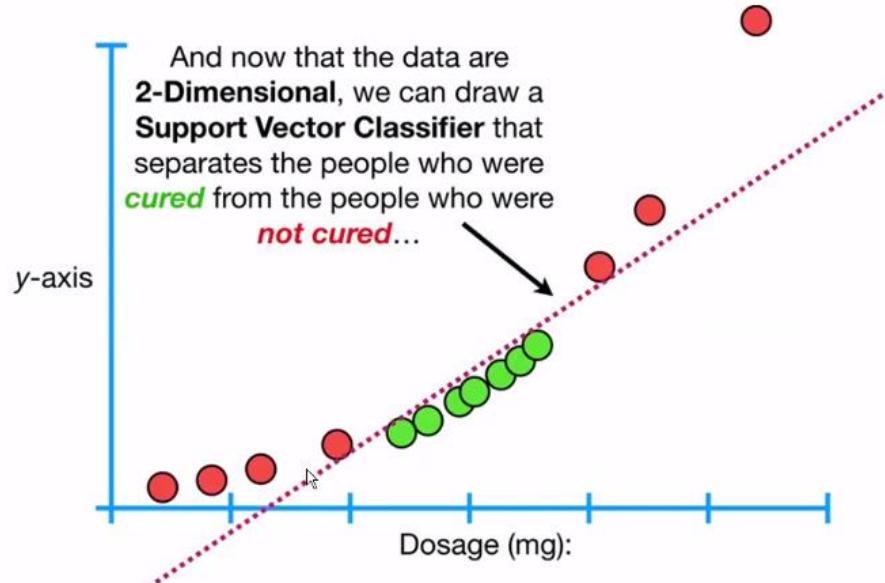


...and then we use **Dosage<sup>2</sup>** for the **y-axis** coordinates for the remaining observations.

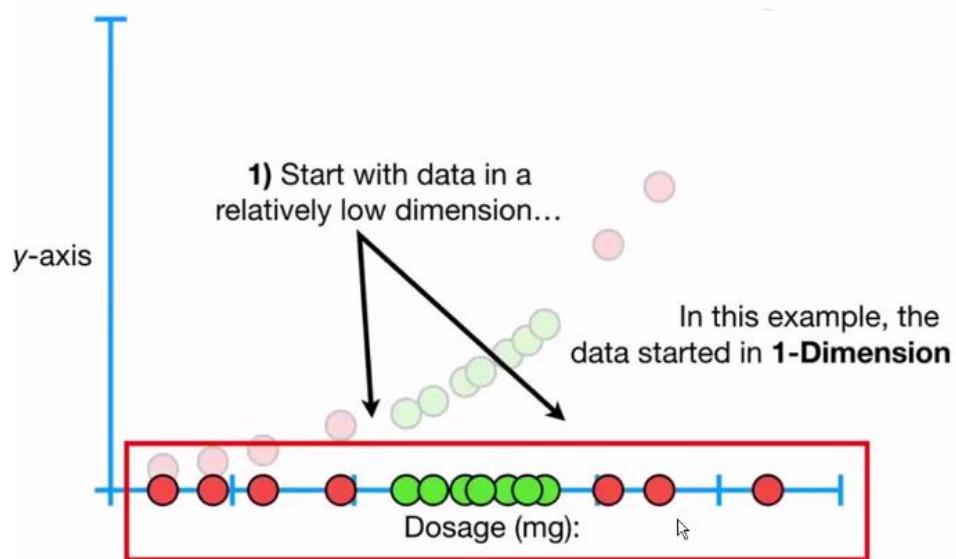


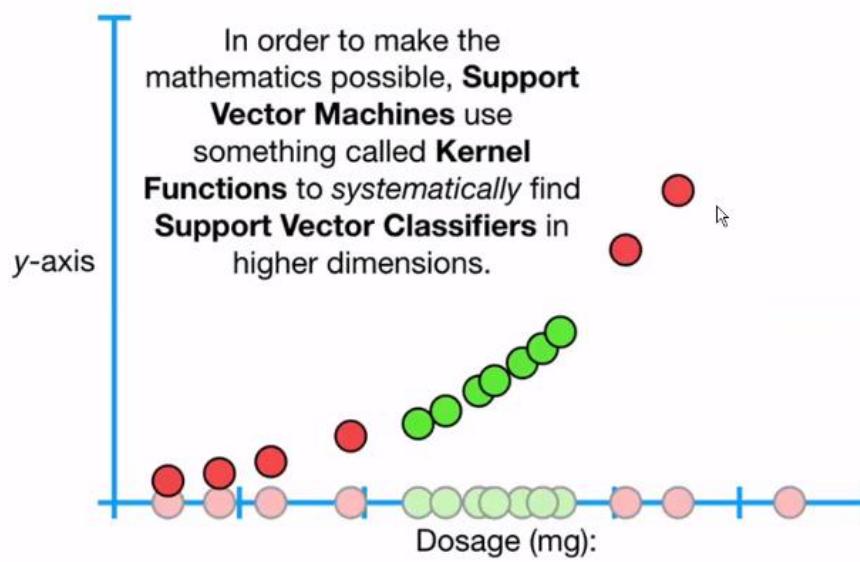
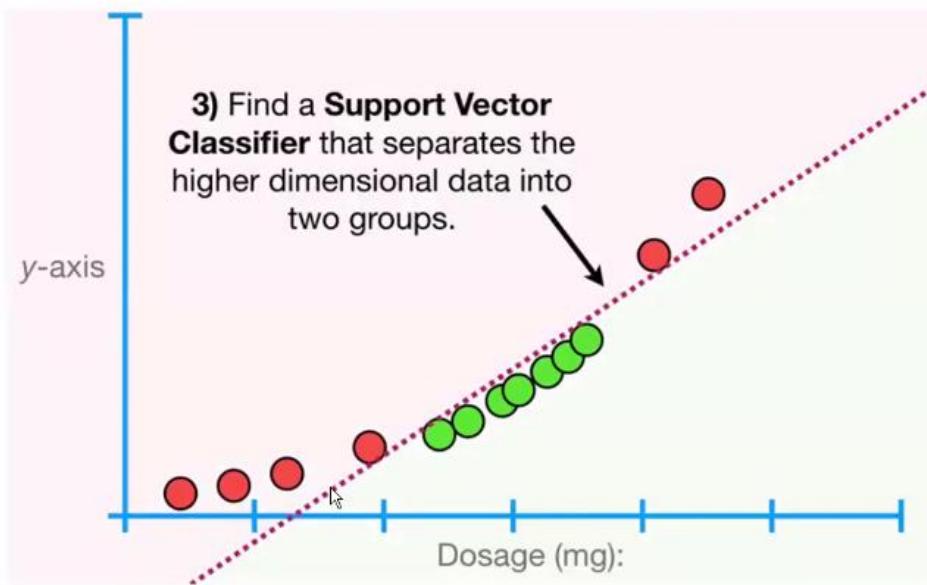
Since each observation has **x** and **y-axis** coordinates, the data are now **2-Dimensional**.

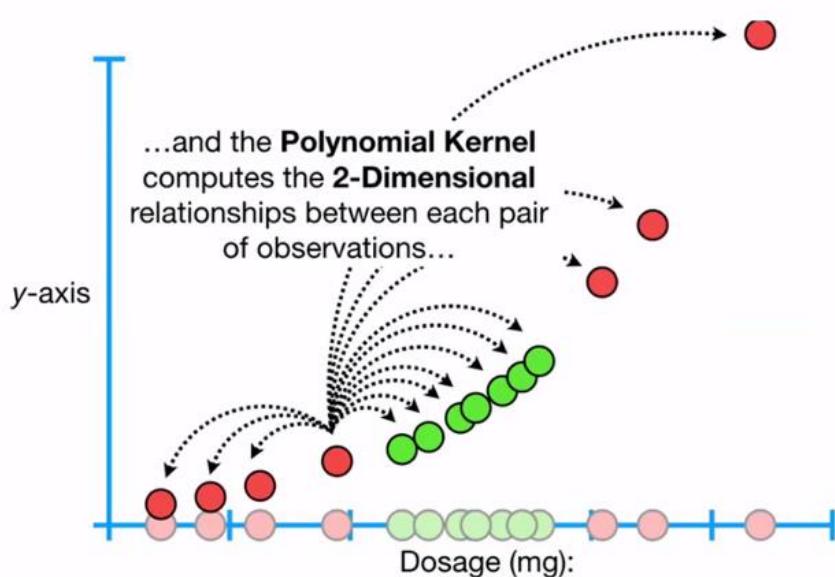




Taking the points which are very boarder of it. We will try to draw the lines (hyperplane)







Another very commonly used **Kernel** is the **Radial Kernel**, also known as the **Radial Basis Function (RBF) Kernel**.

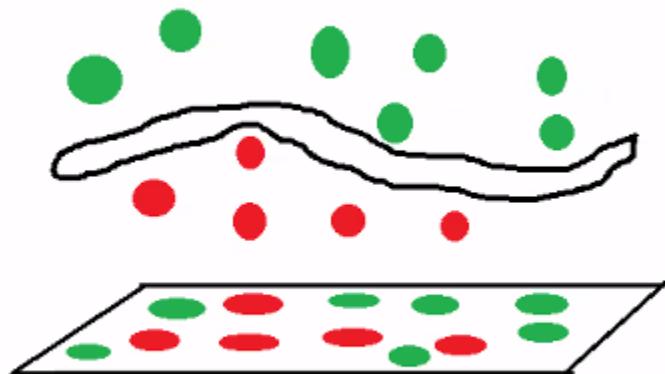
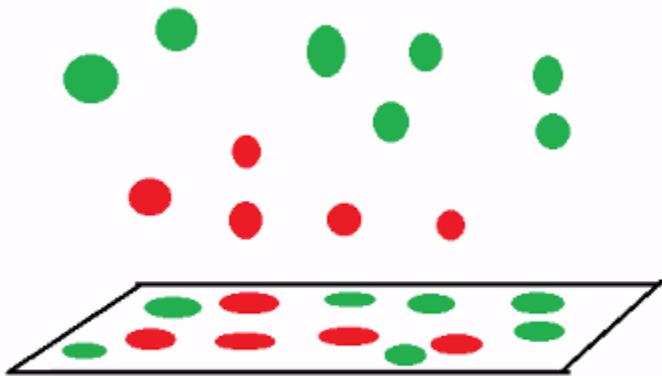
**Radial Kernel finds Support Vector Classifiers in *infinite dimensions***

Datapoints:

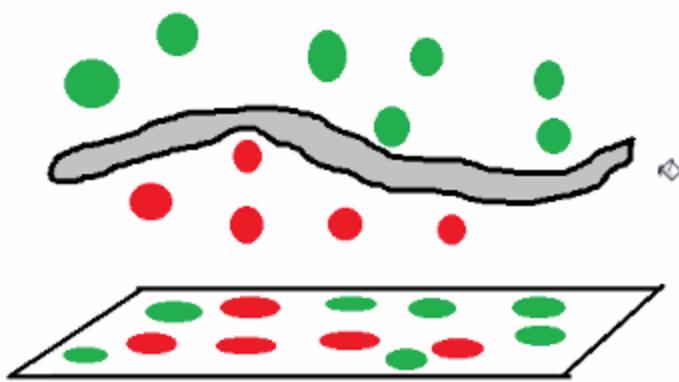


Pull them into higher dimension space:

Floating in the air



Radial basis function: (RDF)



Default kernel is RBF.

y_pred - NumPy object array	
	0
0	0
1	0
2	0
3	0
4	0
5	0
6	0
7	1
8	0
9	0
10	0

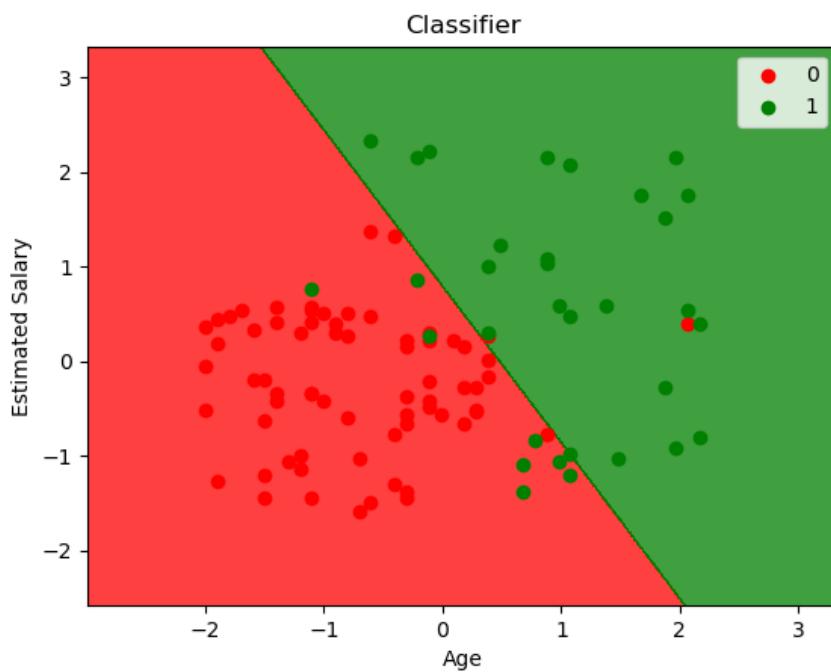
y_test - NumPy object array	
	0
0	0
1	0
2	0
3	0
4	0
5	0
6	0
7	1
8	0
9	0
10	0

Accuracy:

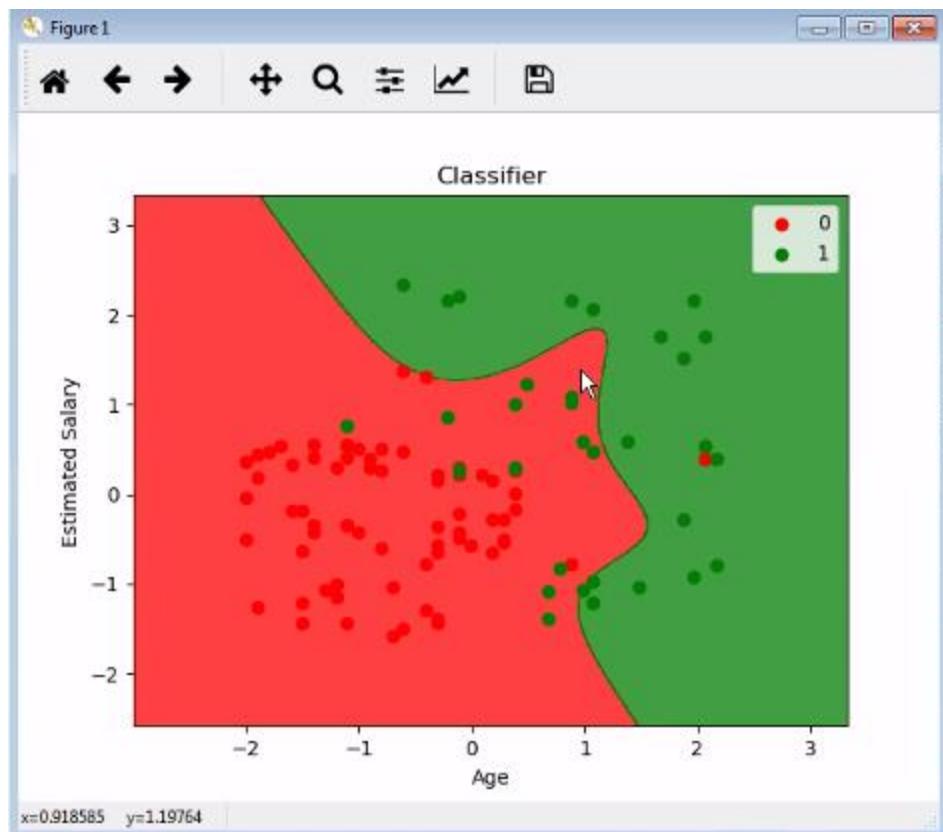
acc	float64	1	0.9
classifier	svm._classes.SVC	1	SVC object of sklearn.svm._classes module

### Visualization:

Linear:



Polynomial with degree 5:



RBF:

Figure 1



Classifier

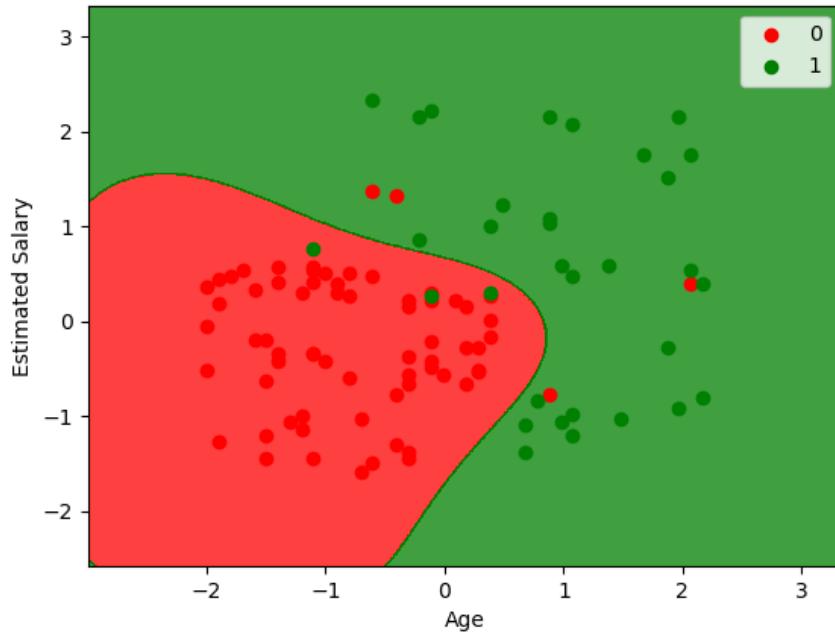
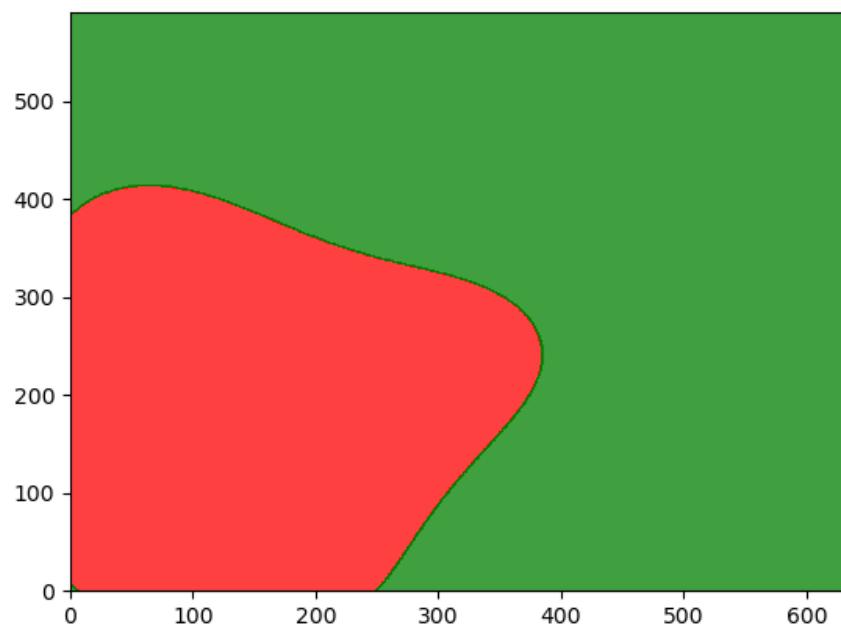
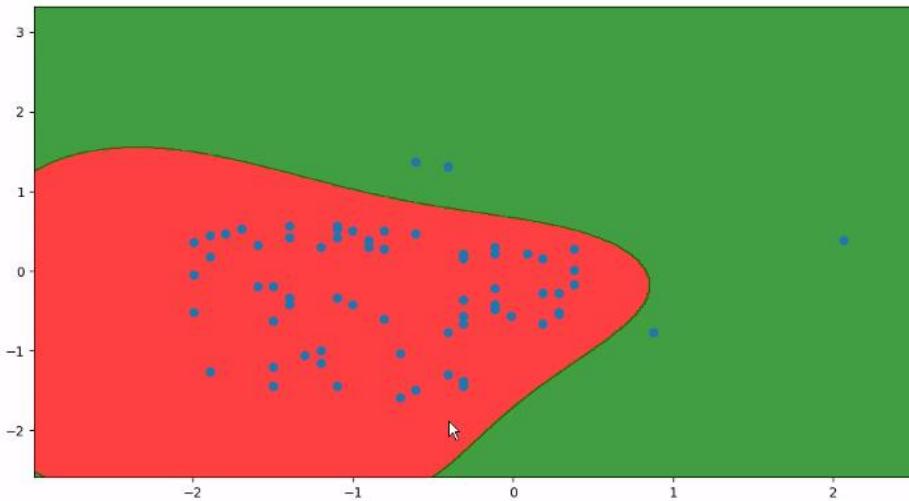


Figure 1

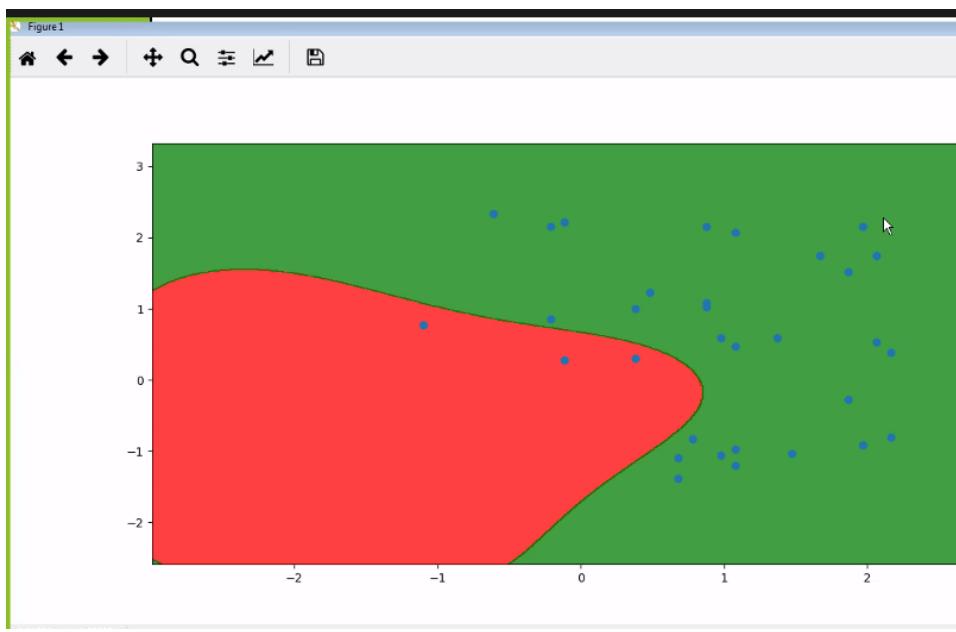


### False cases:

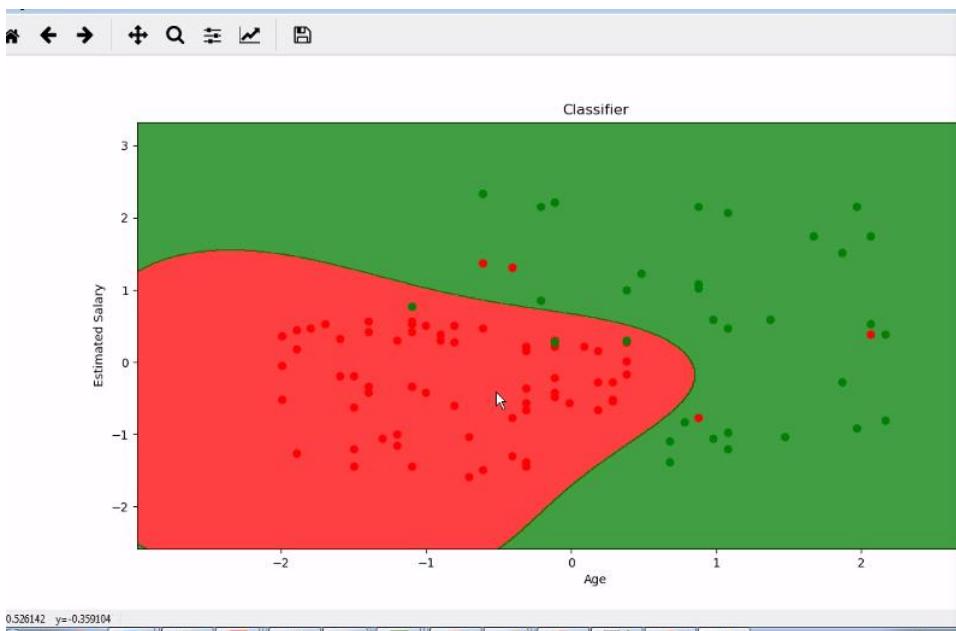


```
1 X1, X2 = np.meshgrid(np.arange(start=X_set[:,0].min()-1,stop=X_set[:,1].max()+1,step=0.01),
2                         np.arange(start=X_set[:,1].min()-1,stop=X_set[:,1].max()+1,step=0.01))
3
4 plt.contourf(X1,X2,
5               classifier.predict(np.array([X1.ravel(),X2.ravel()]).T).reshape(X1.shape),
6               alpha = 0.75 ,
7
8 plt.scatter(X_test[y_test==0],X_test[y_test==1])
9 plt.show()
10
11 plt.xlim(X1.min(),X1.max())
12 plt.ylim(X2.min(),X2.max())
13
14 for i, j in enumerate(np.unique(y_set)):
15     plt.scatter(X_set[y_set == j, 0],
16                 X_set[y_set == j, 1],
17                 c=ListedColormap(['red','green'])(i), label = j)
18
19 plt.title('Classifier')
20 plt.xlabel('Age')
21 plt.ylabel('Estimated Salary')
22 plt.legend()
23 plt.show()
```

### True cases:



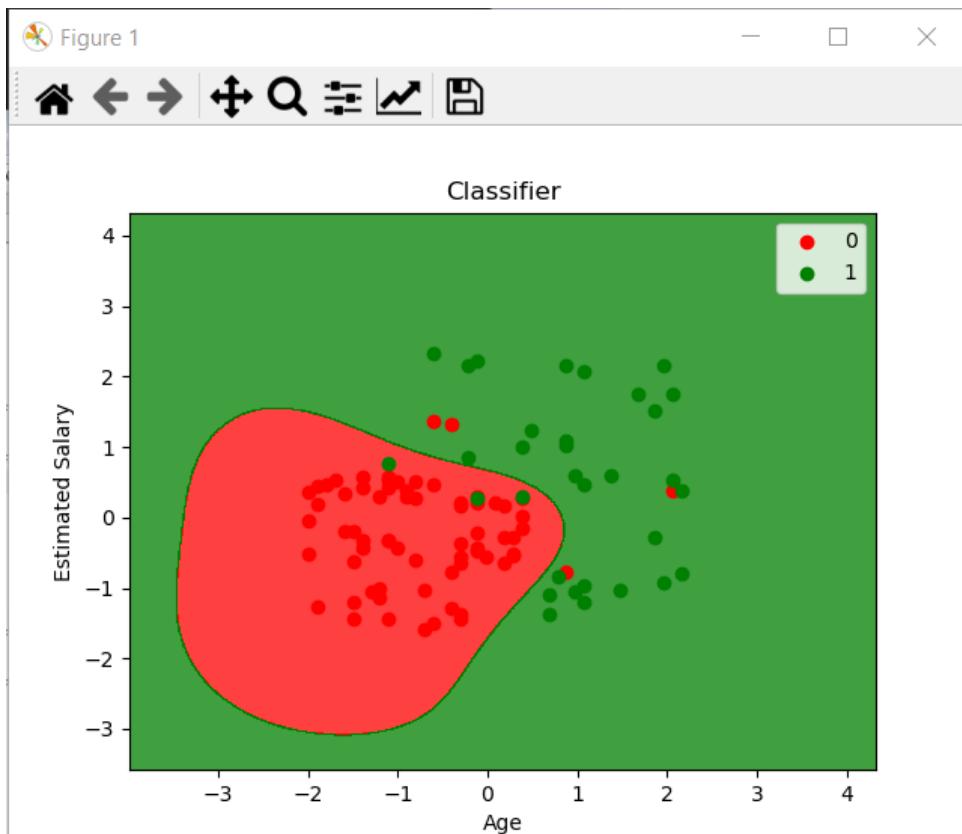
All cases:

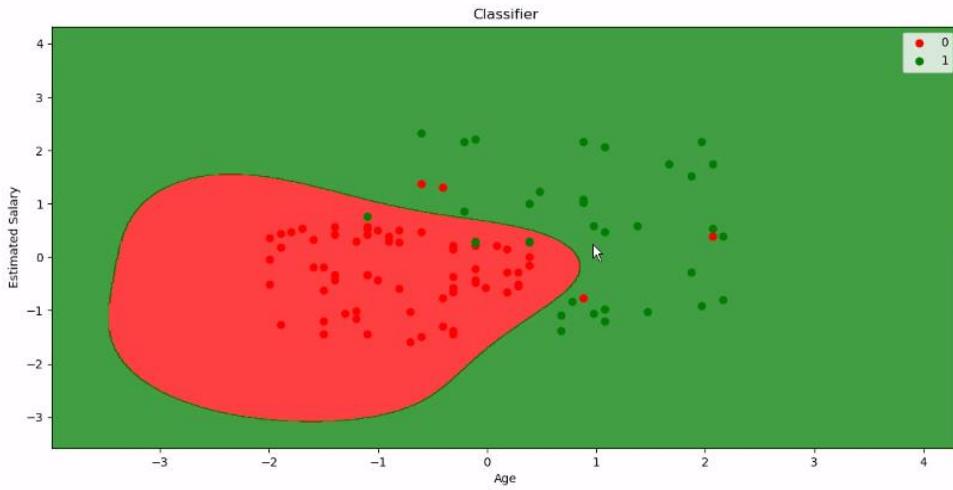


Confusion matrix:

		Predicted		
		NO	YES	
		Predicted		
Actual	NO	64	4	True Negative = 64 Correct Prediction
	YES	3	29	False Positive = 4 False Negative = 3 Incorrect Prediction

Increased the scale:

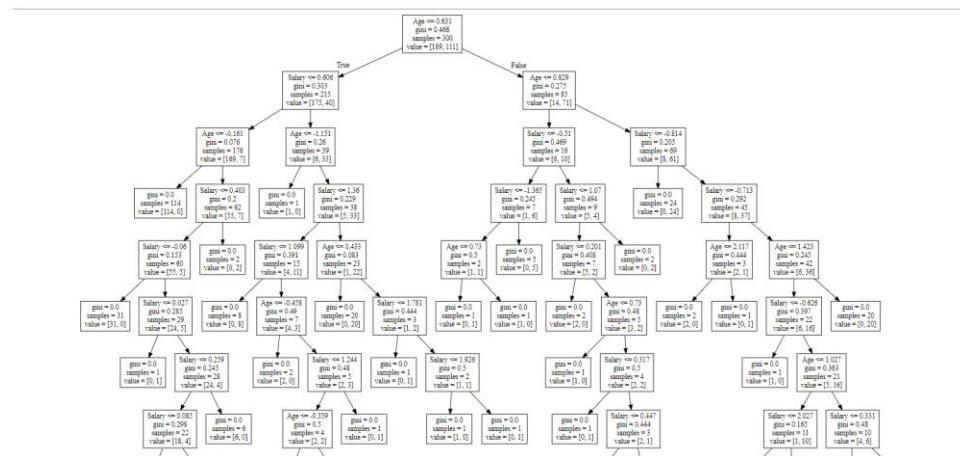




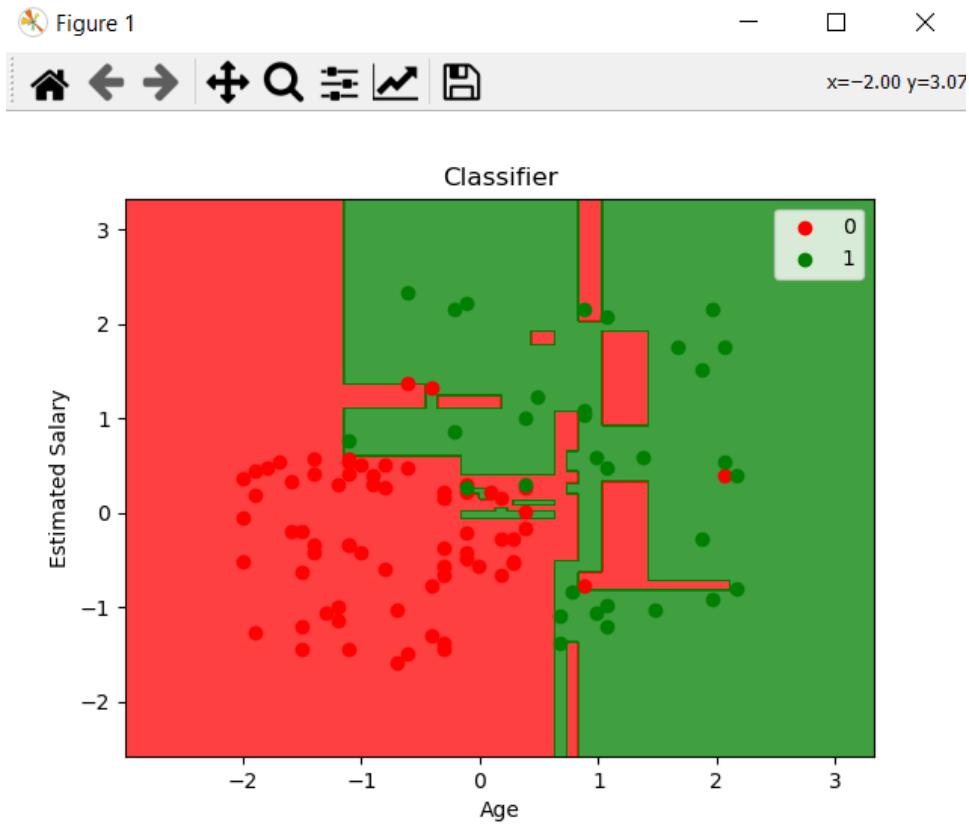
```
3
4 from matplotlib.colors import ListedColormap
5
6 X_set, y_set = X_test,y_test
7
8 X1, X2 = np.meshgrid(np.arange(start=X_set[:,0].min()-1,stop=X_set[:,1].max()+1,step=0.01),
9                      np.arange(start=X_set[:,1].min()-1,stop=X_set[:,1].max()+2,step=0.01))
10
11 plt.contourf(X1,X2,
12               classifier.predict(np.array([X1.ravel(),X2.ravel()]).T).reshape(X1.shape),
13               alpha = 0.75 , cmap = ListedColormap(['red','green']))
14
15 plt.xlim(X1.min(),X1.max())
16 plt.ylim(X2.min(),X2.max())
17
18 for i, j in enumerate(np.unique(y_set)):
```

## Decision Tree Classifier:

Impurity calculation tells how well it perform no and yes response – Gini index



Visualization:

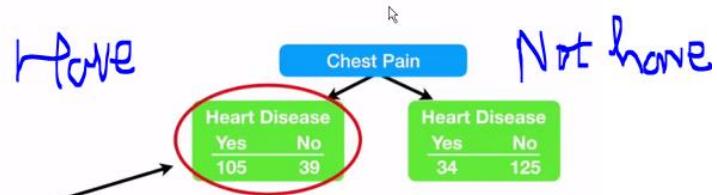


Boundary is not generalized. To control this, we will have minimum split => max depth

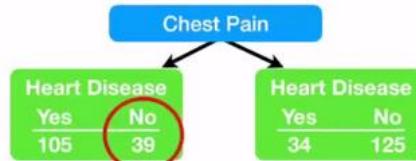
Control the depth here

```
# -----
# Model Implementation
# """
# DecisionTree Classifier
# """
# -----
# from sklearn.tree import DecisionTreeClassifier
# classifier = DecisionTreeClassifier(max_depth=2)
# classifier.fit(X_train,y_train)
# # Model Testing, confusion matrix, accuracy score
# y_pred = classifier.predict(X_test)
# from sklearn.metrics import confusion_matrix,accuracy_score
```

Gini impurity calculation:

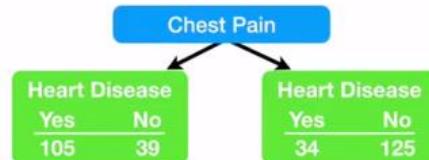


For this leaf, the Gini impurity =  $1 - (\text{probability of "yes"})^2 - (\text{probability of "no"})^2$



For this leaf, the Gini impurity =  $1 - (\text{probability of "yes"})^2 - (\text{probability of "no"})^2$

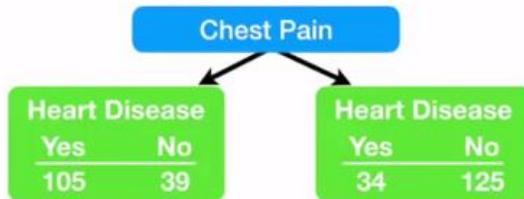
$$= 1 - \left( \frac{105}{105 + 39} \right)^2 - \left( \frac{39}{105 + 39} \right)^2$$



For this leaf, the Gini impurity =  $1 - (\text{probability of "yes"})^2 - (\text{probability of "no"})^2$

$$= 1 - \left( \frac{105}{105 + 39} \right)^2 - \left( \frac{39}{105 + 39} \right)^2$$

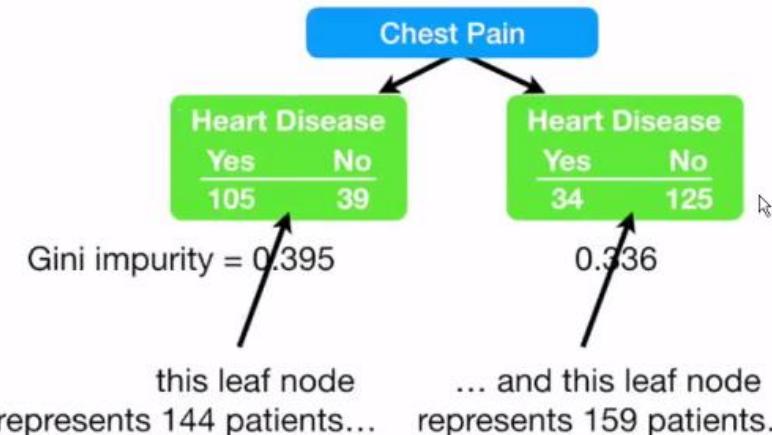
$$= 0.395$$



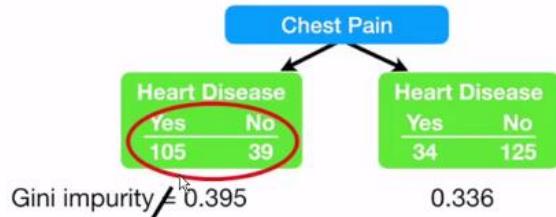
$= 1 - (\text{the probability of "yes"})^2 - (\text{the probability of "no"})^2$

$$= 1 - \left( \frac{34}{34 + 125} \right)^2 - \left( \frac{125}{34 + 125} \right)^2$$

$$= 0.336$$

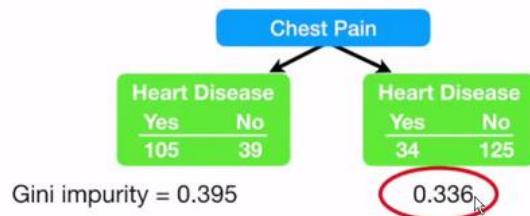


...the leaf nodes do not represent the same number of patients.



Gini impurity for Chest Pain = weighted average of Gini impurities for the leaf nodes

$$= \left( \frac{144}{144 + 159} \right) 0.395$$



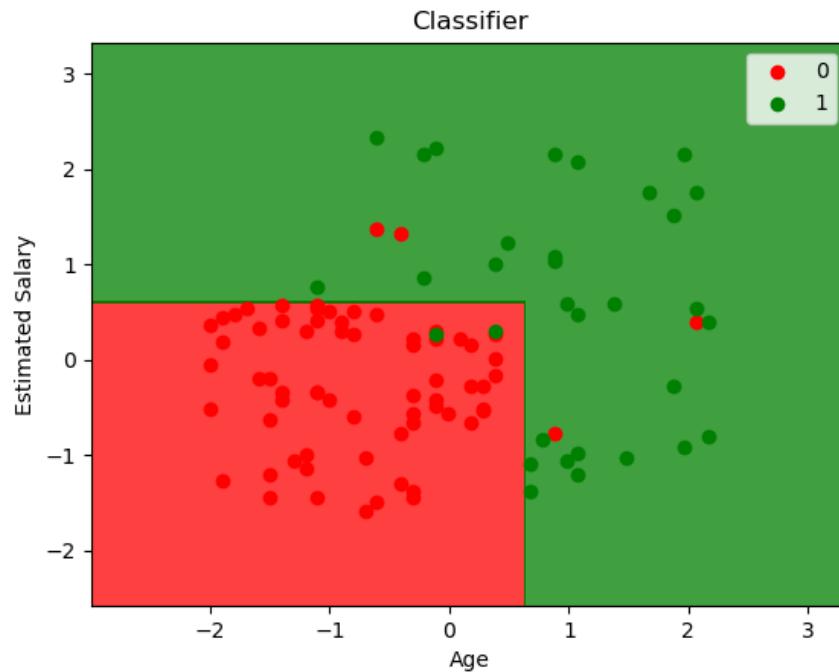
Gini impurity for Chest Pain = weighted average of Gini impurities for the leaf nodes

$$= \left( \frac{144}{144 + 159} \right) 0.395 + \left( \frac{159}{144 + 159} \right) 0.336$$

$$= 0.364$$

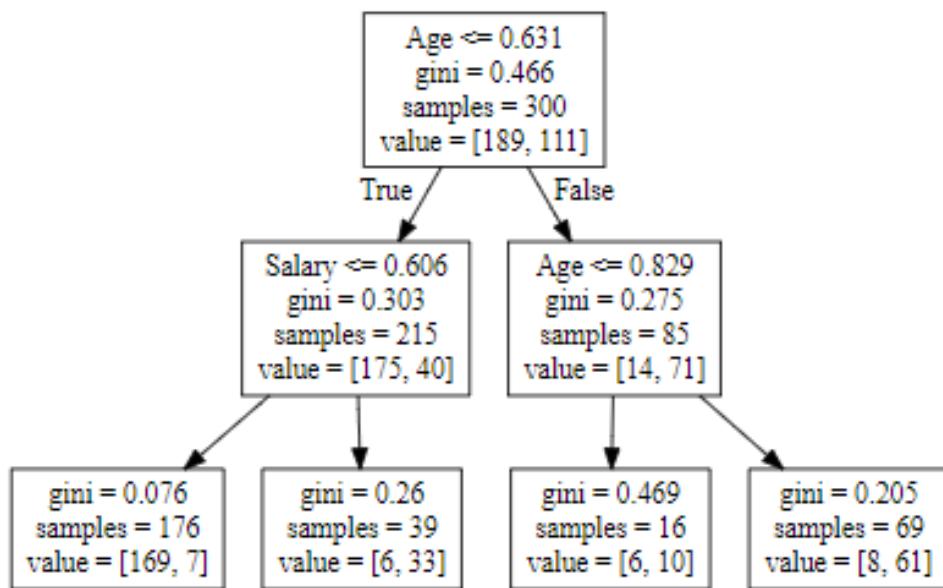
Increase the max depth = 2:

Figure 1



Boundary seems to be generalized.

Good decision tree:



Classification      Value to Predict is Discrete

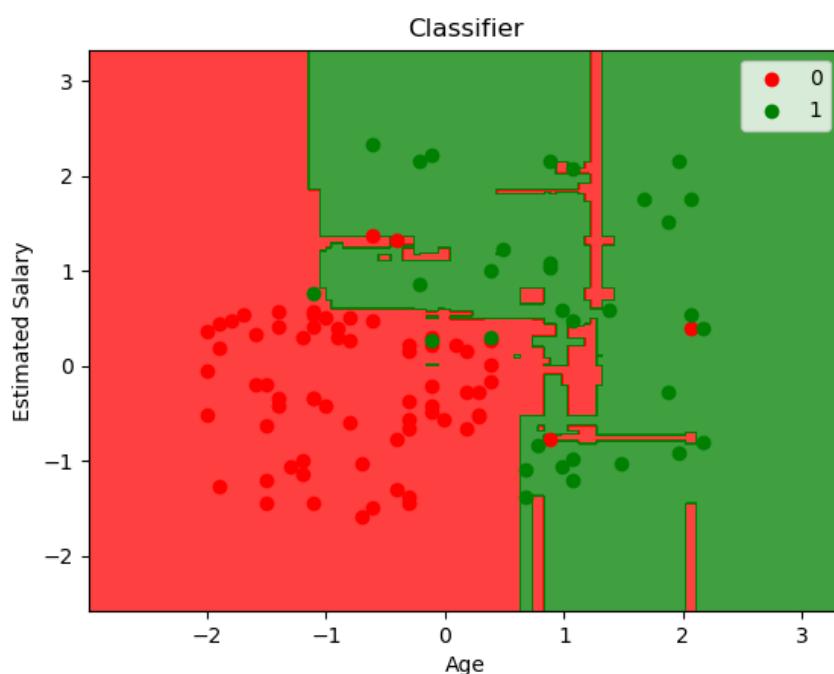
Logistic Regression	93%
Support Vector Machine	94%
DecisionTree Classifier	
RandomForest Classifier	max_depth=2
Naive Bayes	I

Random Forest Classifier:

Name	Type	Size	Value
acc	float64	1	0.92
classifier	ensemble._forest.RandomForestClassifier	10	RandomForestClassif...

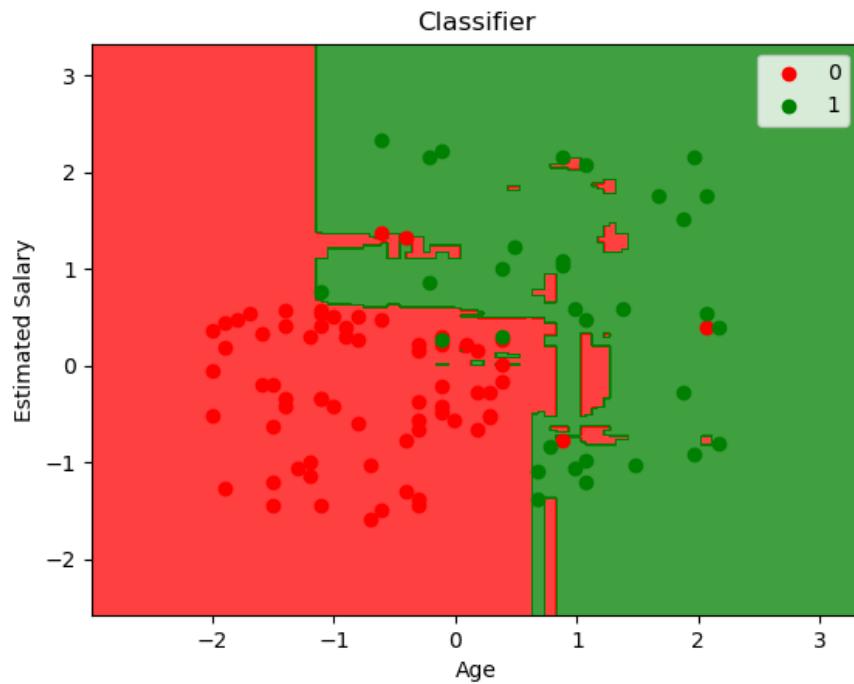
Visualization:

Overfitting issue!



Change it to 100 trees:

Figure 1



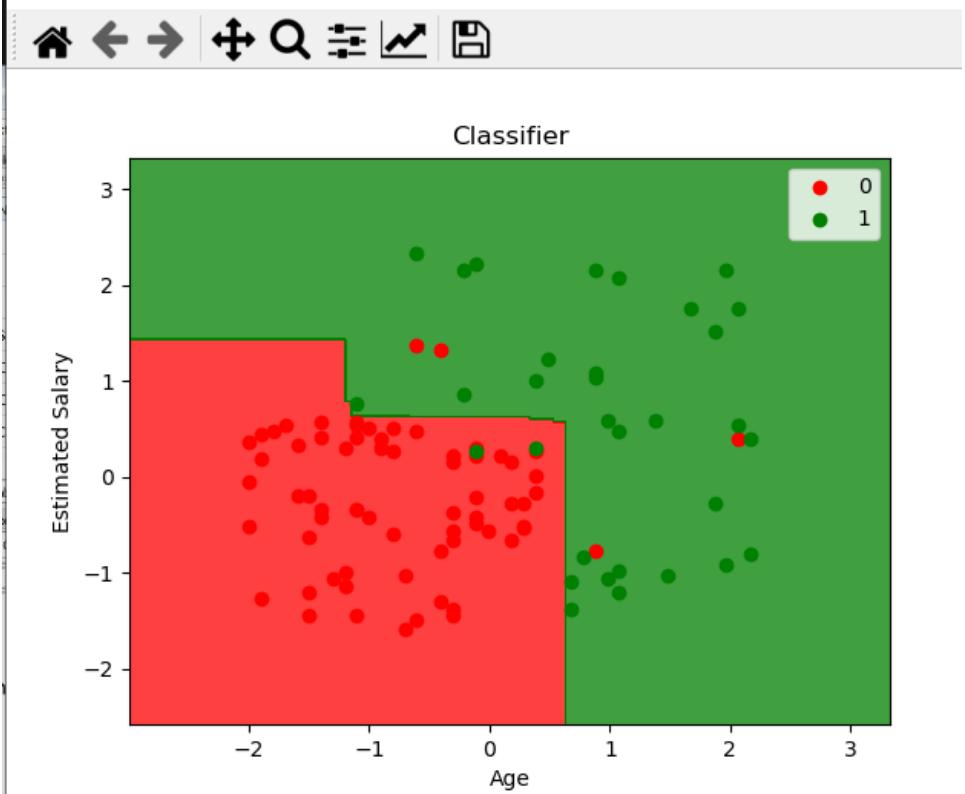
Name	Type	Size	Value
acc	float64	1	0.92
classifier	ensemble._forest.RandomForestClassifier	100	RandomForestClass.

Not helping! Same accuracy

Try to include max\_depth=2:

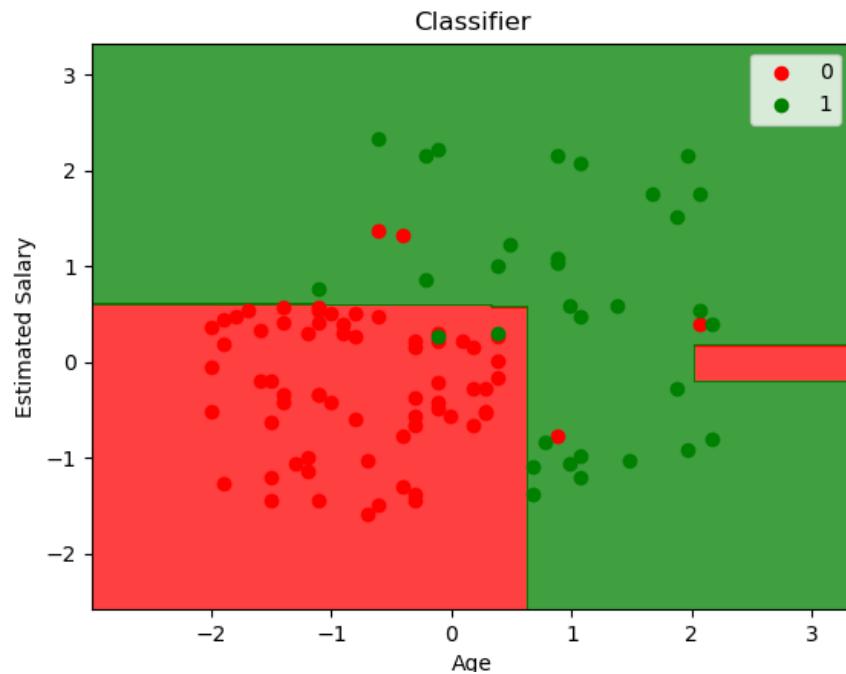
With 100 trees:

Figure 1



With 10 trees:

Figure 1



Name	Type	Size	Value
acc	float64	1	0.94
classifier	ensemble._forest.RandomForestClassifier	100	RandomForestClass..

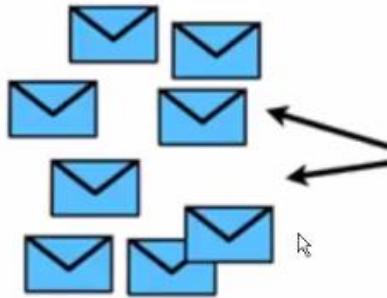
Name	Type	Size	Value
acc	float64	1	0.94
classifier	ensemble._forest.RandomForestClassifier	10	RandomForestClass..

Classification      Value to Predict is Discrete

Logistic Regression	
Support Vector Machine	93%
DecisionTree Classifier	94%
RandomForest Classifier	94%      max_depth=2
Naive Bayes	max_depth=2 ,Trees=10

### Naïve Bayes:

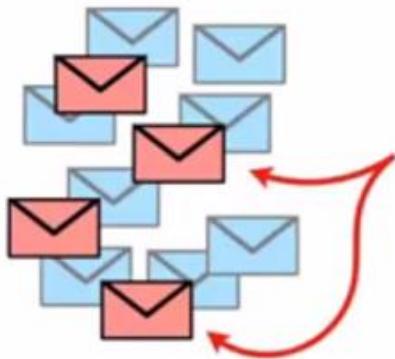
Learned based on probabilities.



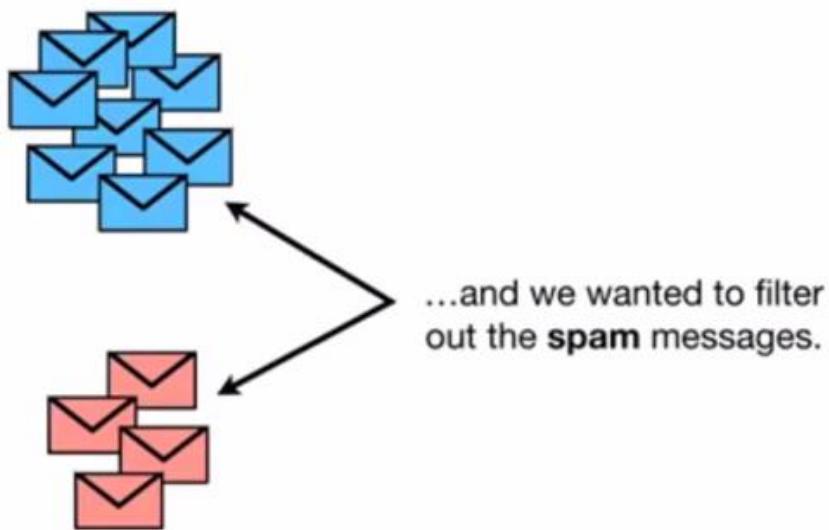
Imagine we received  
**normal messages** from  
friends and family...

Create spam filter

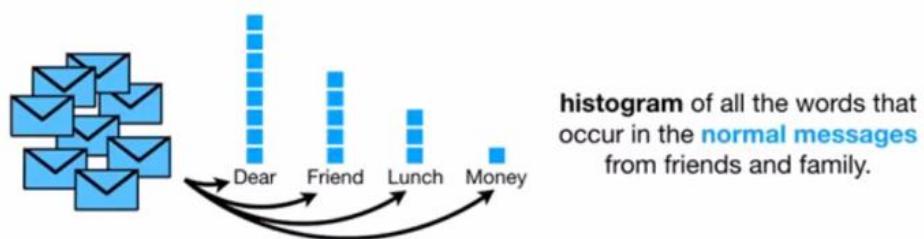
Segreate normal and spam messages.



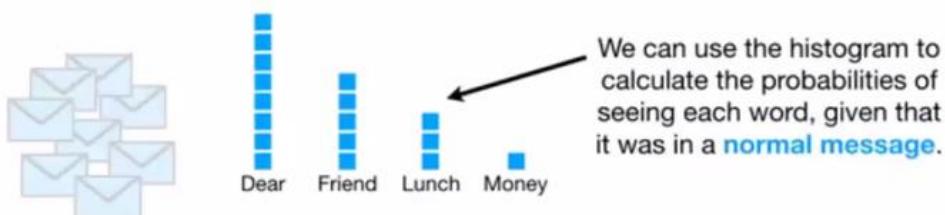
...and we also received  
**spam** (unwanted  
messages that are usually  
scams or unsolicited  
advertisements)...

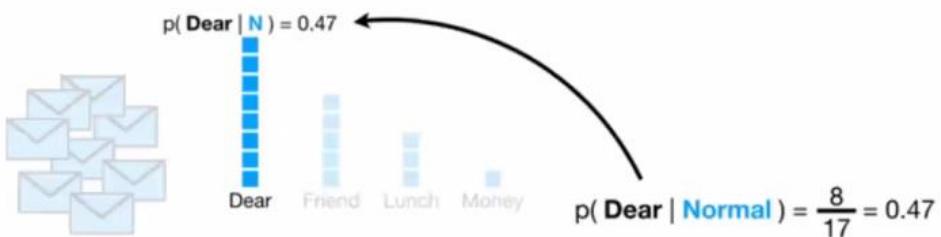
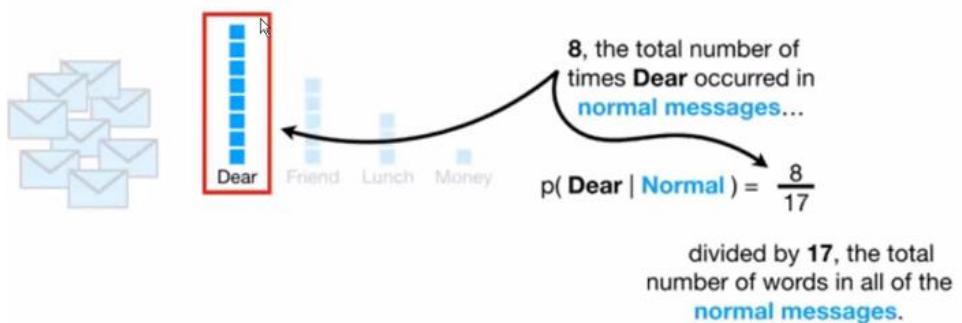
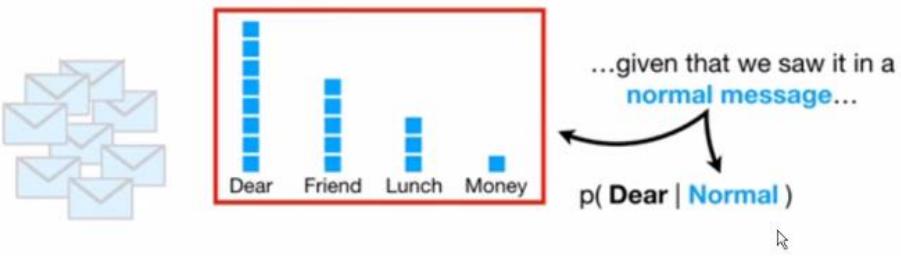


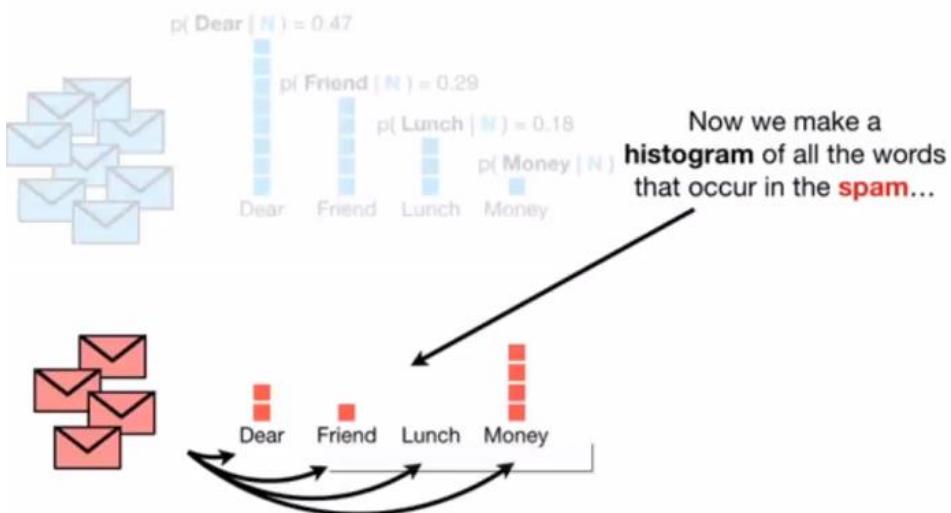
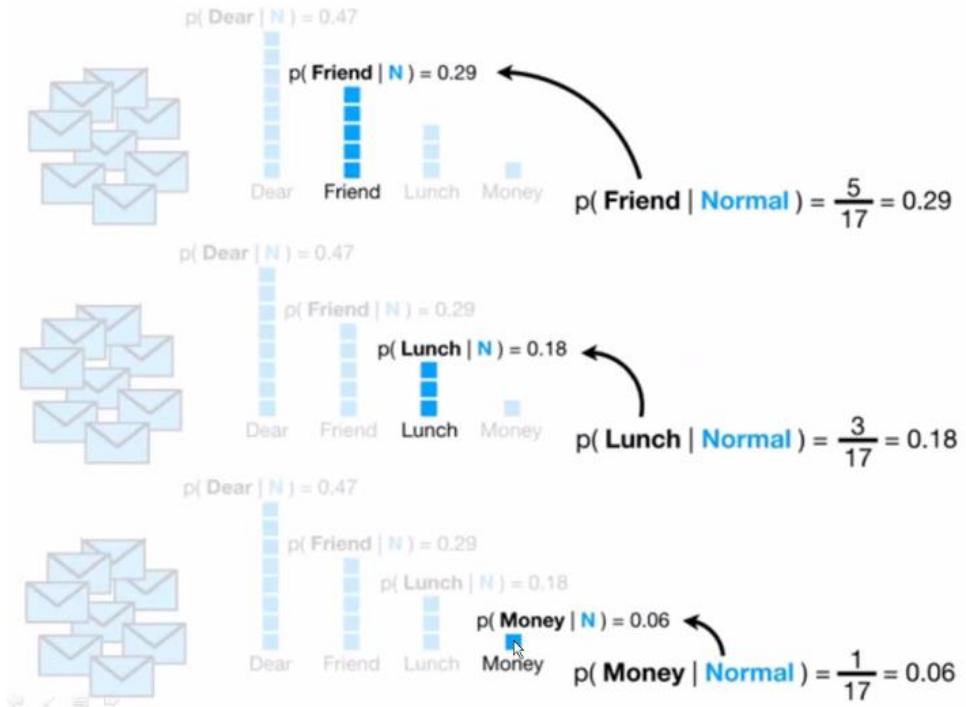
Find probabilities for each words:

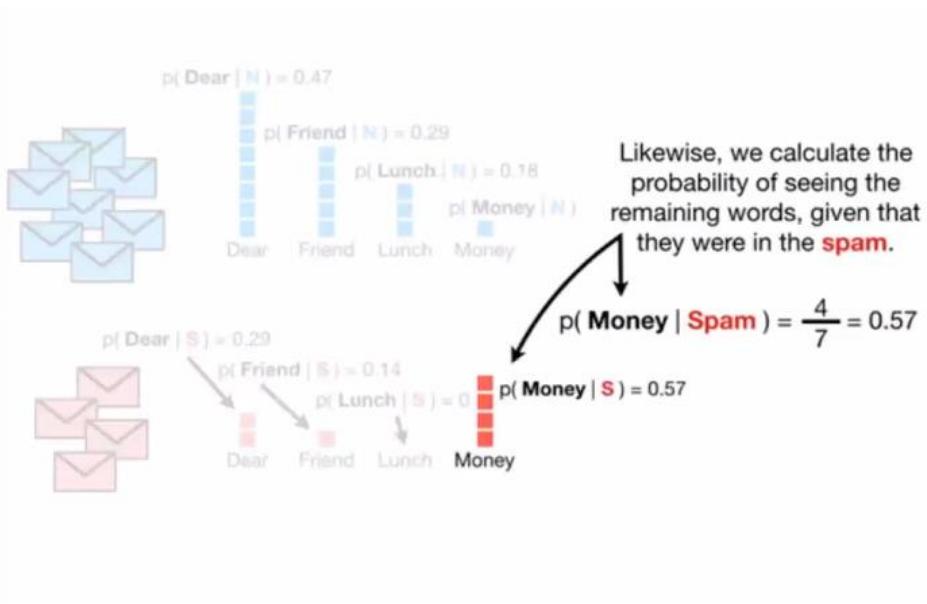
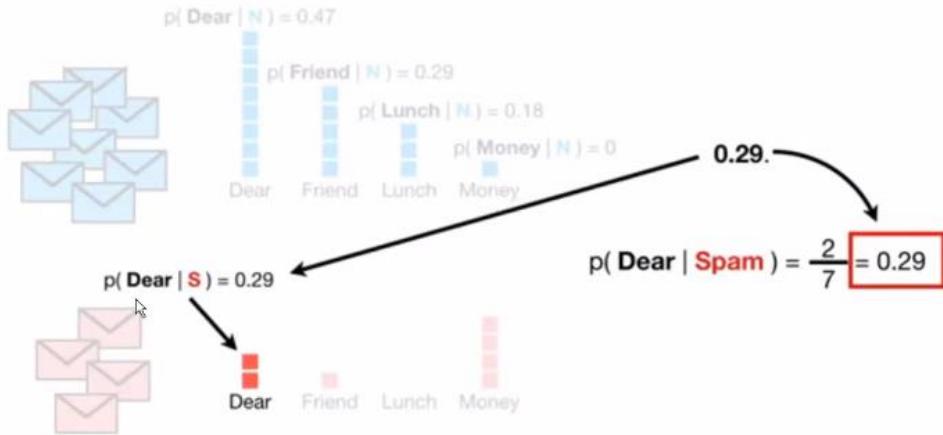


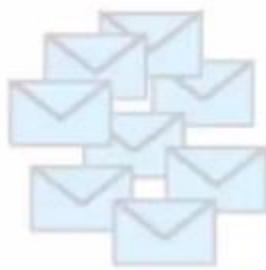
Find the probabilities of words given the visibilities of normal messages











$$\begin{aligned} p(\text{Dear} | \text{N}) &= 0.47 \\ p(\text{Friend} | \text{N}) &= 0.29 \\ p(\text{Lunch} | \text{N}) &= 0.18 \\ p(\text{Money} | \text{N}) &= 0.06 \end{aligned}$$



$$\begin{aligned} p(\text{Dear} | \text{S}) &= 0.29 \\ p(\text{Friend} | \text{S}) &= 0.14 \\ p(\text{Lunch} | \text{S}) &= 0.00 \\ p(\text{Money} | \text{S}) &= 0.57 \end{aligned}$$

We Receive a new Message

Dear Friend



$$\begin{aligned} p(\text{Dear} | \text{N}) &= 0.47 \\ p(\text{Friend} | \text{N}) &= 0.29 \\ p(\text{Lunch} | \text{N}) &= 0.18 \\ p(\text{Money} | \text{N}) &= 0.06 \end{aligned}$$

?

?

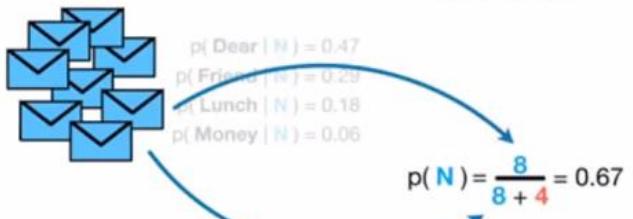
And we want to  
decide if is a **normal**  
**message or spam.**



$$\begin{aligned} p(\text{Dear} | \text{S}) &= 0.29 \\ p(\text{Friend} | \text{S}) &= 0.14 \\ p(\text{Lunch} | \text{S}) &= 0.00 \\ p(\text{Money} | \text{S}) &= 0.57 \end{aligned}$$

?

Dear Friend



since **8** of  
the **12** messages are  
**normal messages**, our  
initial guess will be **0.67**.



Dear Friend





$p(N) = 0.67$

$p(\text{Dear} | N) = 0.47$   
 $p(\text{Friend} | N) = 0.29$   
 $p(\text{Lunch} | N) = 0.18$   
 $p(\text{Money} | N) = 0.06$

$$p(N) \times p(\text{Dear} | N) \times p(\text{Friend} | N)$$



$p(\text{Dear} | S) = 0.29$   
 $p(\text{Friend} | S) = 0.14$   
 $p(\text{Lunch} | S) = 0.00$   
 $p(\text{Money} | S) = 0.57$



$p(N) = 0.67$

$p(\text{Dear} | N) = 0.47$   
 $p(\text{Friend} | N) = 0.29$   
 $p(\text{Lunch} | N) = 0.18$   
 $p(\text{Money} | N) = 0.06$

$$p(N) \times p(\text{Dear} | N) \times p(\text{Friend} | N)$$



$p(\text{Dear} | S) = 0.29$   
 $p(\text{Friend} | S) = 0.14$   
 $p(\text{Lunch} | S) = 0.00$   
 $p(\text{Money} | S) = 0.57$

Now we multiply that initial guess by probability that the word **Dear** occurs in a **normal message**...

...and the probability that the word **Friend** occurs in a **normal message**.

Now we multiply that initial guess by probability that the word **Dear** occurs in a **normal message**...

...and the probability that the word **Friend** occurs in a **normal message**.



$p(N) = 0.67$

$p(\text{Dear} | N) = 0.47$   
 $p(\text{Friend} | N) = 0.29$   
 $p(\text{Lunch} | N) = 0.18$   
 $p(\text{Money} | N) = 0.06$

$$p(N) \times p(\text{Dear} | N) \times p(\text{Friend} | N)$$



$p(\text{Dear} | S) = 0.29$   
 $p(\text{Friend} | S) = 0.14$   
 $p(\text{Lunch} | S) = 0.00$   
 $p(\text{Money} | S) = 0.57$

Now we multiply that initial guess by probability that the word **Dear** occurs in a **normal message**...

...and the probability that the word **Friend** occurs in a **normal message**.

$$0.67 \times 0.47 \times 0.29 = 0.09$$

In a simple way, we can think of **0.09** as the score that **Dear Friend** gets if it is a **Normal Message**.



$p(N) = 0.67$

$p(\text{Dear} | N) = 0.47$   
 $p(\text{Friend} | N) = 0.29$   
 $p(\text{Lunch} | N) = 0.18$   
 $p(\text{Money} | N) = 0.06$

Now we multiply that initial guess by probability that the word **Dear** occurs in a **normal message**...

...and the probability that the word **Friend** occurs in a **normal message**.

$$p(N) \times p(\text{Dear} | N) \times p(\text{Friend} | N)$$

$$0.67 \times 0.47 \times 0.29 = 0.09 \quad p(N | \text{Dear Friend})$$



$p(\text{Dear} | S) = 0.29$   
 $p(\text{Friend} | S) = 0.14$   
 $p(\text{Lunch} | S) = 0.00$   
 $p(\text{Money} | S) = 0.57$

In a simple way, we can think of **0.09** as the score that **Dear Friend** gets if it is a **Normal Message**.

$$p(N | \text{Dear Friend}) \approx 0.09 \quad \text{Dear Friend}$$



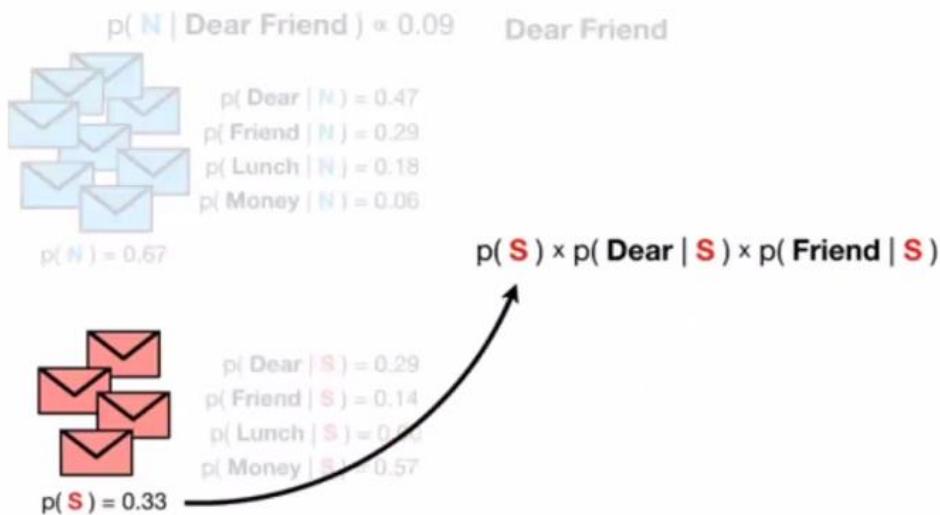
$p(N) = 0.67$

$p(\text{Dear} | N) = 0.47$   
 $p(\text{Friend} | N) = 0.29$   
 $p(\text{Lunch} | N) = 0.18$   
 $p(\text{Money} | N) = 0.06$

$$p(S) = \frac{4}{4 + 8} = 0.33$$

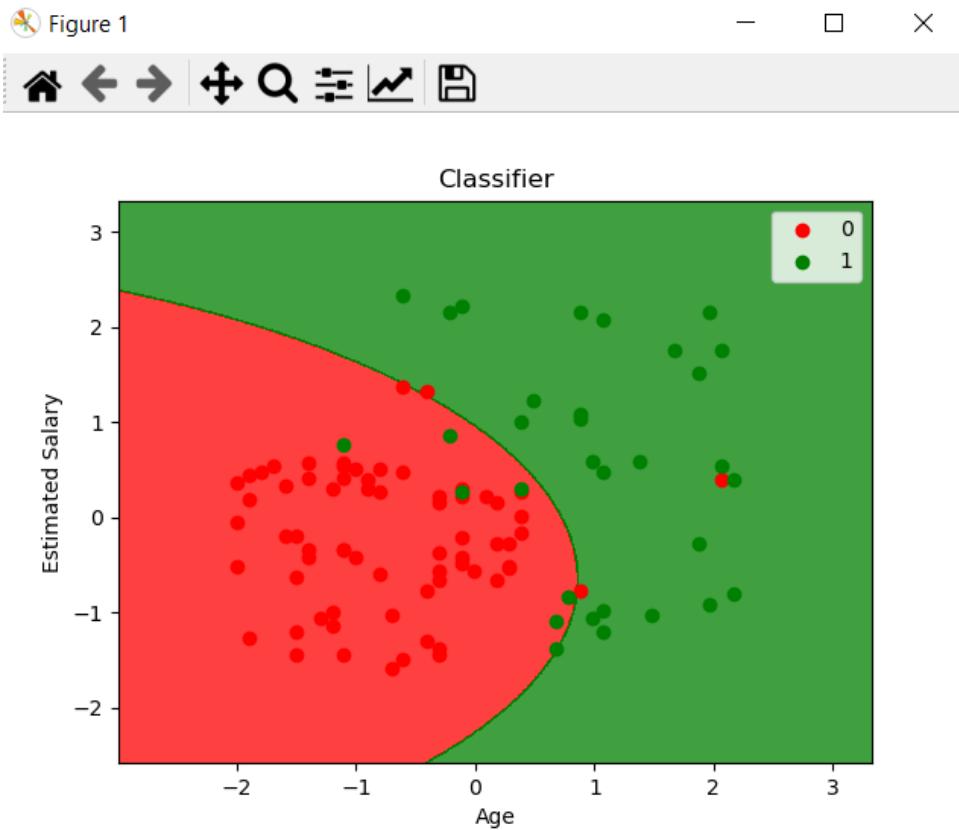


$p(\text{Dear} | S) = 0.29$   
 $p(\text{Friend} | S) = 0.14$   
 $p(\text{Lunch} | S) = 0.00$   
 $p(\text{Money} | S) = 0.57$





### Visualization:



### Accuracy:

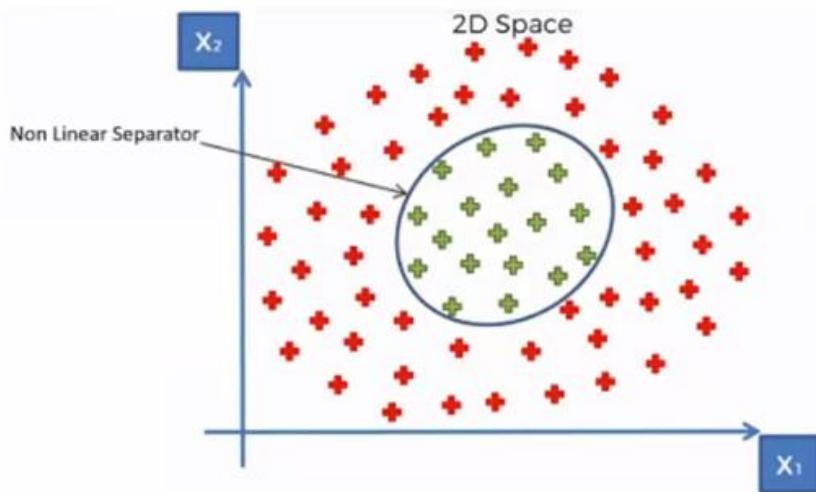
Name	Type	Size	Value
acc	float64	1	0.9
classifier	naive_bayes.GaussianNB	1	GaussianNB object c

Classification      Value to Predict is Discrete

Logistic Regression		
Support Vector Machine	93%	
DecisionTree Classifier	94%	max_depth=2
RandomForest Classifier	94%	max_depth=2 ,Trees=10
Naive Bayes	90%	

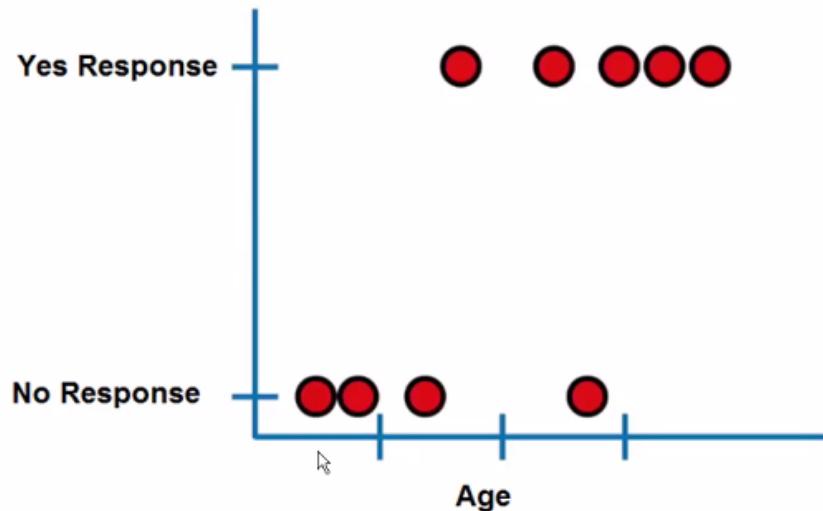
### Logistic Regression:

How the formula is going to give results to classify the problem

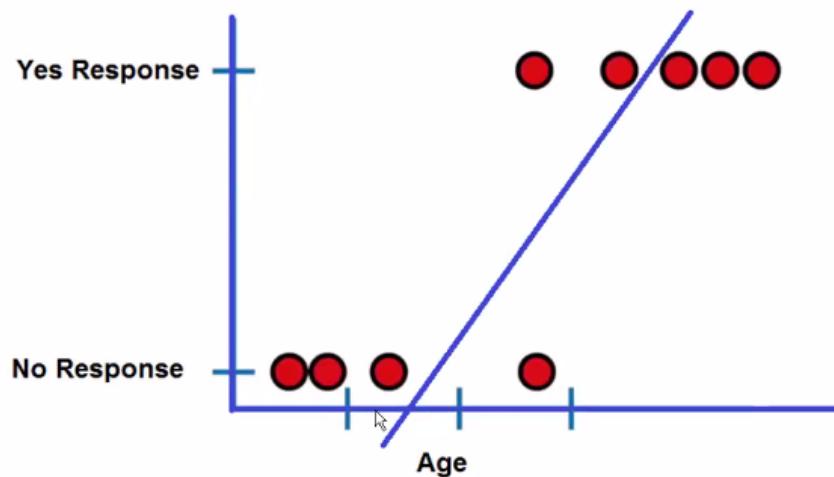


How the straight line fit => Regression

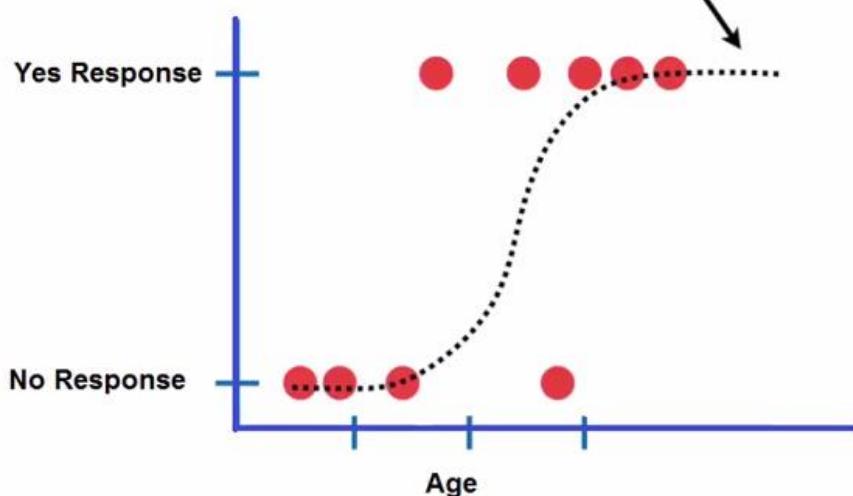
# Logistic Regression



# Logistic Regression



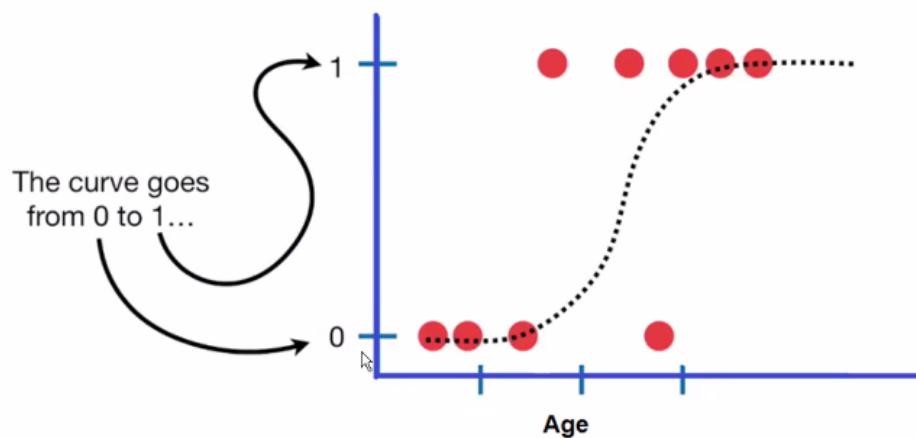
instead of fitting a line to the data, logistic regression fits an "S" shaped "logistic function".



Bring Logistic function

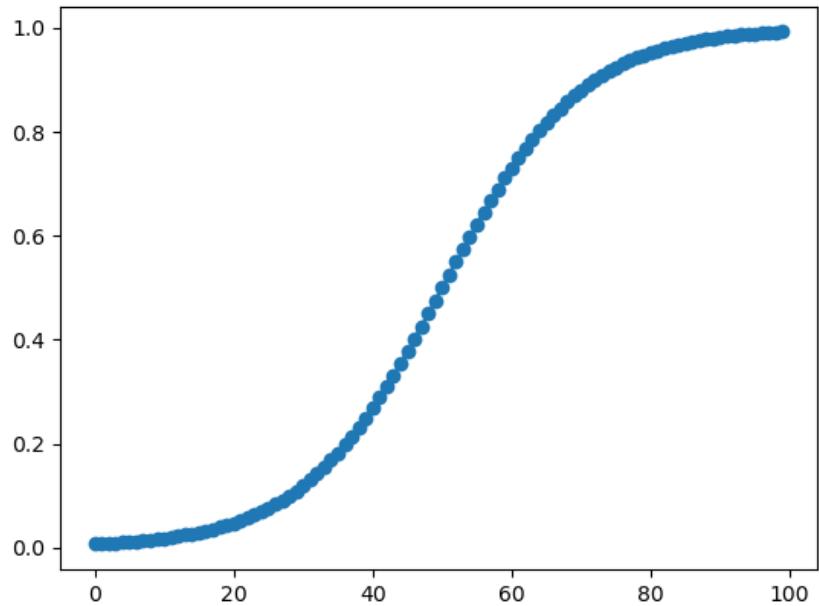
## Logistic Regression

curve tells you the probability



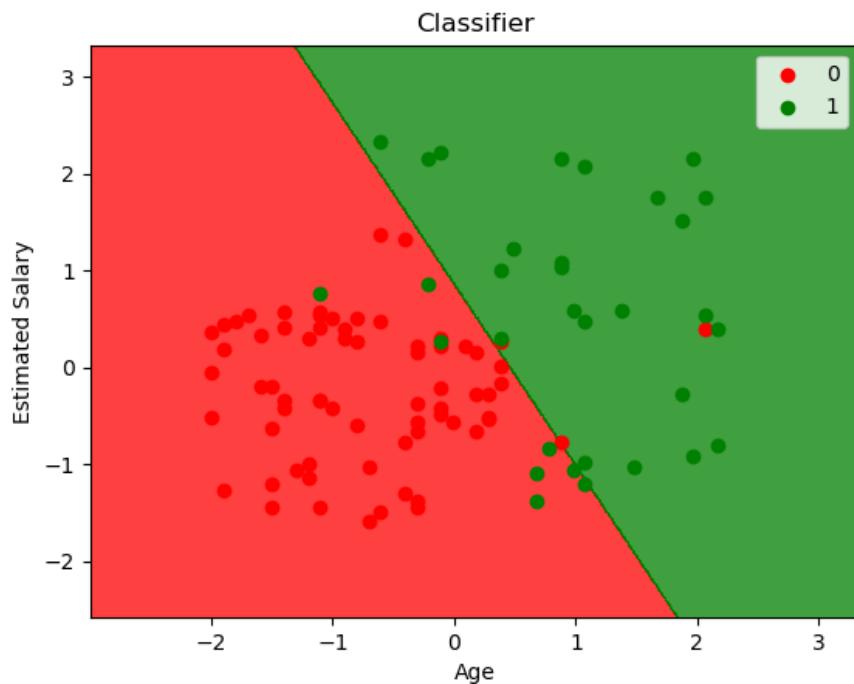
Logistic function is also called as sigmoid functions

Figure 1



**Visualization:**

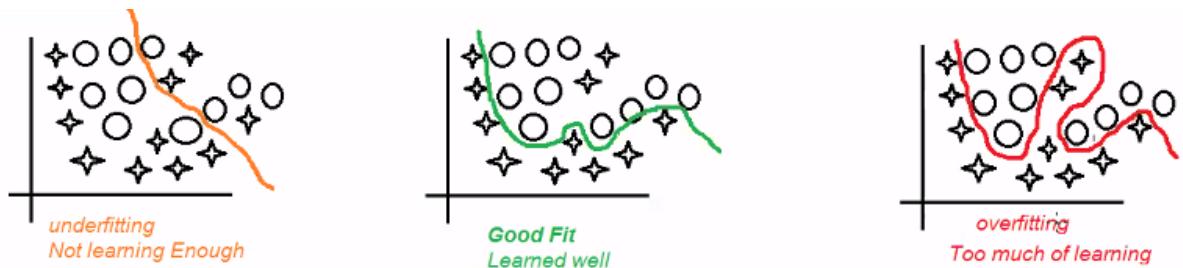
Figure 1



Accuracy:

Name	Type	Size	Value
acc	float64	1	0.89
classifier	linear_model._logistic.LogisticRegression	1	LogisticRegress...

Classification	Value to Predict is Discrete
Logistic Regression	89%
Support Vector Machine	93%
DecisionTree Classifier	94%
RandomForest Classifier	94% max_depth=2 max_depth=2 ,Trees=10
Naive Bayes	90%



*Training Set*

*we train the model on training data*

*Validation Set*

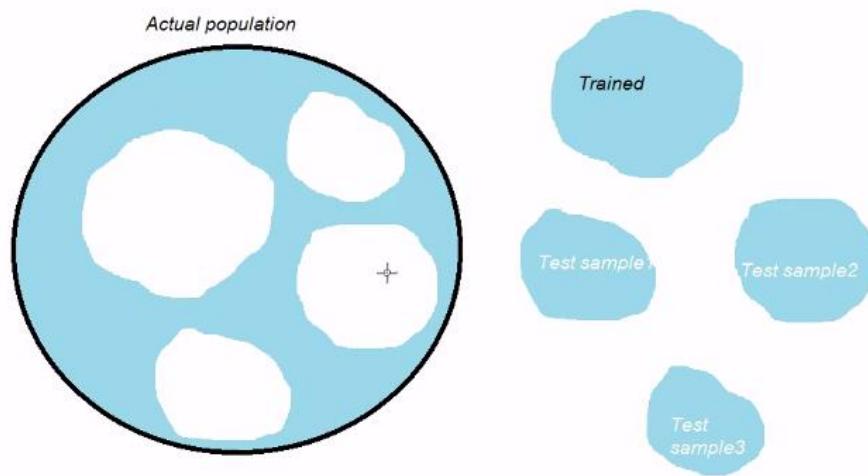
*After training the model, we check how well it performs on the validation set*

*Test Set*

*When we have final Model (model that has performed well on both training and validation). We evaluate it on test set to go and get unbiased estimation of its performance*

### Stability check:

Test the model with various test dataset



### K-Fold Cross Validation:

X\_train, y\_train – divide it into K-Folds

K = 5 to 10 folds

100 rows – dataset

Fold1 20 rows

Fold2 20 rows

Fold3 20 rows

Fold4 20 rows

Fold5 20 rows

Same as bootstrapping – to have equal folds.

Cross validation – 90%-10% Train – Test

#### K-Fold Cross Validation

X\_train,y\_train divide it into K-Folds k=5  
100 rows

Fold1 20 rows

Fold2 20 rows

Fold3 20 rows

Fold4 20 rows

Fold5 20 rows

#### Cross Validation - 90%-10% Train-Test

Training	Testing	Accuracy
----------	---------	----------

Fold1-2-3-4	Fold-5	93%
-------------	--------	-----

Fold-1-2-3-5	Fold-4	80%
--------------	--------	-----

Fold-1-2-4-5	Fold-3	75%
--------------	--------	-----

Fold-1-3-4-5	Fold-2	98%
--------------	--------	-----

Fold-2-3-4-5	Fold-1	90%
--------------	--------	-----

High variance

Accuracy variation > 1%

Accuracy variation == 1%

Overfitting issues

Good Fit

The screenshot shows a Jupyter Notebook interface with two code cells and their corresponding outputs.

**Code Cell 1:**

```
# -*- coding: utf-8 -*-
"""
Created on Wed Mar  9 16:43:49 2022
@author: TSE
"""

import numpy as np
accuracyScore1 = np.array([0.93, 0.80, 0.75, 0.98, 0.90])
accuracyScore1.std()
```

**Code Cell 2:**

```
# -*- coding: utf-8 -*-
"""
Created on Wed Mar  9 16:43:49 2022
@author: TSE
"""

import numpy as np
accuracyScore1 = np.array([0.93, 0.80, 0.75, 0.98, 0.90])
accuracyScore1.std() # 0.08471127433818948
accuracyScore2 = np.array([0.93, 0.92, 0.93, 0.94, 0.93])
accuracyScore2.std() # 0.006324555320336729
```

**Variable Explorer:**

Name	Type	Size	Value
accuracyScore1	float64	(5,)	[0.93 0.80 0.75 0.98 0.9 ]
accuracyScore2	float64	(5,)	[0.93 0.92 0.93 0.94 0.93]

**Python Console:**

```
In [2]: import numpy as np
...: accuracyScore1 =
np.array([0.93, 0.80, 0.75, 0.98, 0.90])

In [3]: accuracyScore1.std()
Out[3]: 0.08471127433818948

In [4]:
```

```
In [2]: import numpy as np
...: accuracyScore1 =
np.array([0.93, 0.80, 0.75, 0.98, 0.90])

In [3]: accuracyScore1.std()
Out[3]: 0.08471127433818948

In [4]: accuracyScore2 =
np.array([0.93, 0.92, 0.93, 0.94, 0.93])
...:
...: accuracyScore2.std()
Out[4]: 0.006324555320336729

In [5]:
```