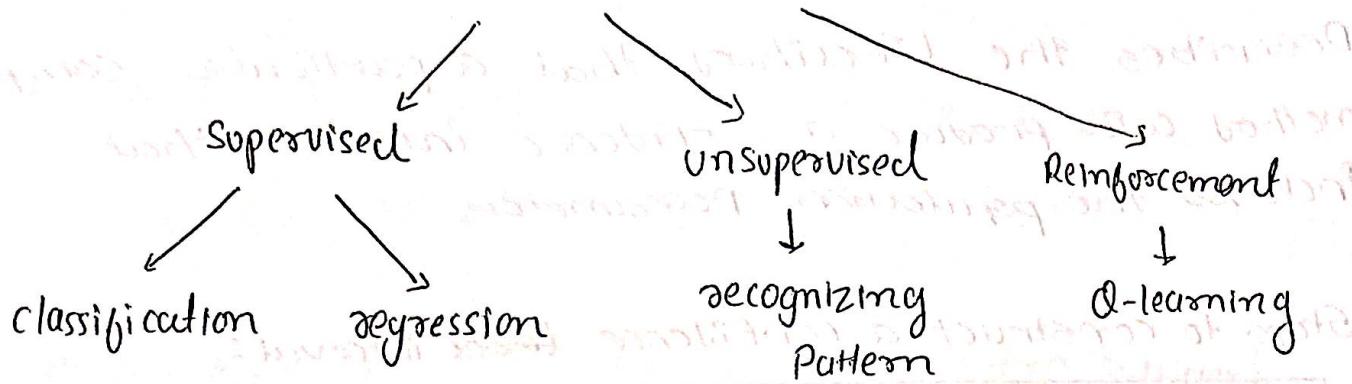
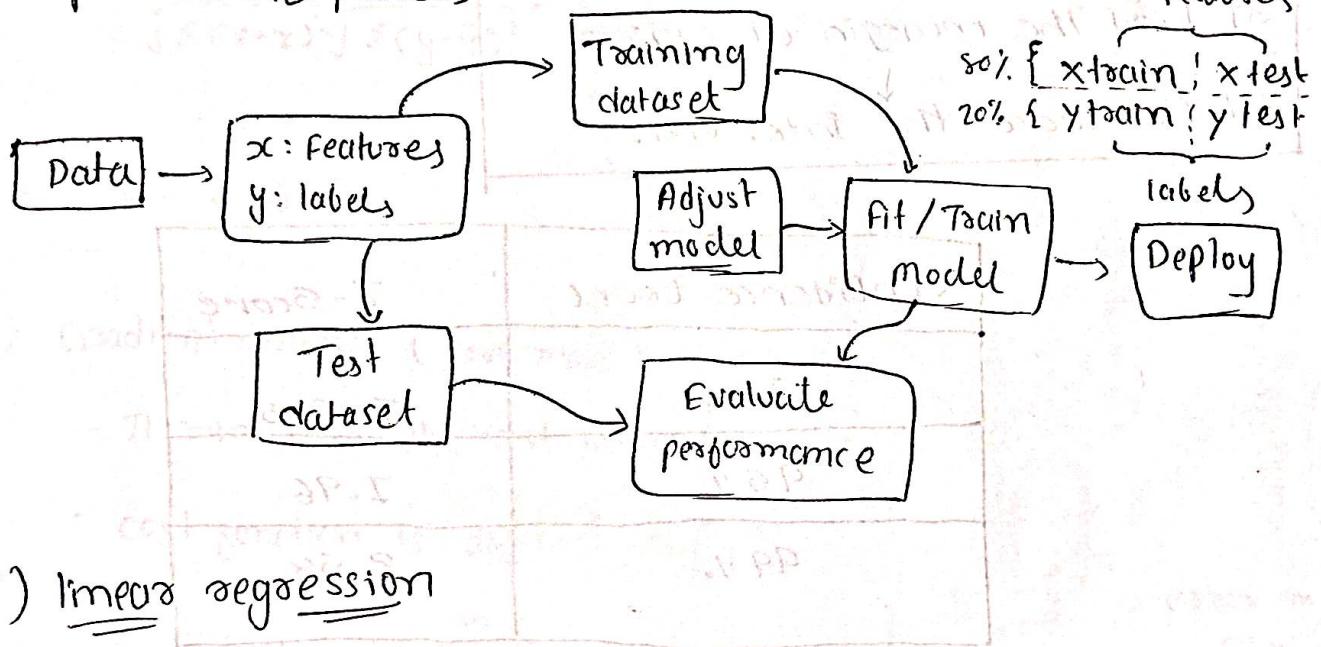


# Machine learning

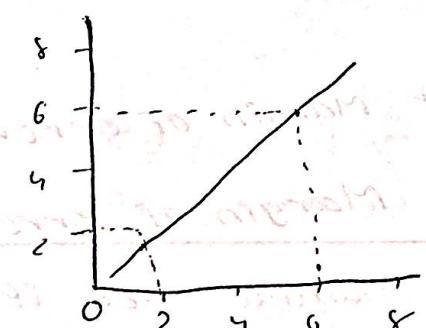


## i) Supervised learning

- supervised ml process



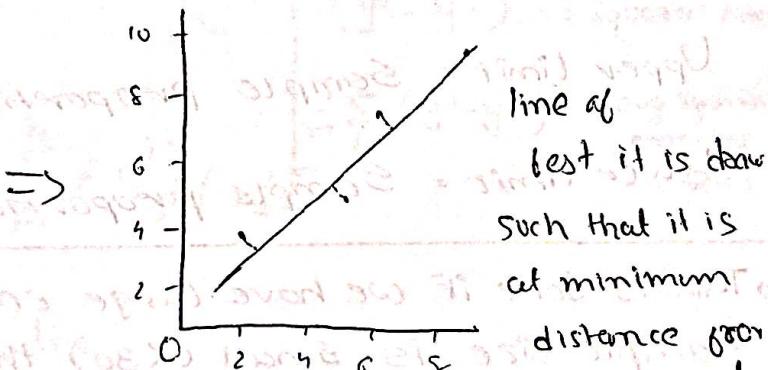
## i) Linear regression



Here we can't predict

any value of  $y$  for

any  $x$  value.



equation of line is,

$$y = mx + c$$

slope  $m$  & intercept  $c$

$$\Rightarrow y = \beta_0 x_0 + \beta_1 x_1 + \dots + \beta_n x_n = \sum_{i=1}^n \beta_i x_i$$

There are 2 methods to calculate slope and intercept or we can say  $\beta$  values

### 1) Ordinary least square (only for 1 feature ' $x$ ')

$$\hat{y} = b_1 x + b_0$$

$$\Rightarrow b_1 = \rho_{xy} = \frac{\sigma_y}{\sigma_x} \quad \text{where } \rho_{xy} = \text{Pearson correlation coefficient}$$

$\sigma_x, \sigma_y = \text{standard Deviation}$

$$b_1 = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}} = \frac{\frac{\sum (y - \bar{y})^2}{n}}{\sqrt{\frac{\sum (x - \bar{x})^2}{n}}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\Rightarrow b_0 = \bar{y} - b_1 \bar{x}$$

### 2) Gradient descent (for more than 1 feature)

- It minimize the cost function.

- cost function is defined as

$$\Rightarrow J(\beta) = \frac{1}{2m} \sum_{i=1}^m (y^i - \hat{y}^i)^2$$

$$\Rightarrow \frac{\partial J(\beta)}{\partial \beta_k} = \frac{1}{m} \sum_{i=1}^m \left( y^i - \sum_{j=0}^n \beta_j x_{ij} \right) (-x_{ik})$$

$y^i - \hat{y}^i \Rightarrow \text{error } m^i$   
 $(y^i - \hat{y}^i)^2 \Rightarrow \text{squared error}$   
 $\frac{1}{m} \sum_{i=1}^m (y^i - \hat{y}^i)^2 \Rightarrow \text{avg squared error of } m \text{ rows}$

$$\Rightarrow \nabla_{\beta} J = \begin{bmatrix} \frac{\partial J}{\partial \beta_0} \\ \vdots \\ \frac{\partial J}{\partial \beta_n} \end{bmatrix}$$

gradient

Derivative  
of cost  
function  
with respect to  
each  $\beta$

## $\Rightarrow$ Generated code

1. import required libraries such as numpy, pandas & matplotlib
2. Read the dataset.
3. from sklearn.model\_selection import train\_test\_split
4. X-train, X-test, y-train, y-test = train\_test\_split(X, y), test\_size=0.3, random\_state=101
5. from sklearn.linear\_model import LinearRegression
6. model = LinearRegression() # Here, we pass hyper parameter of model
7. model.fit(Xtrain, y-train)
8. model.predict(X-test)

## $\Rightarrow$ Evaluating Regression

1. Mean absolute error (MAE)
2. mean squared error (MSE)
3. Root mean square error (RMSE)

1. MAE: ~~absolute difference between y-test & test prediction~~

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

where  $y_i$  = actual value

$\hat{y}_i$  = predicted value

2. MSE

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

3. RMSE

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

8. test-prediction = model.predict(x-test)

9. from sklearn.metrics import mean\_absolute\_error,  
mean\_squared\_error

10. mean\_absolute\_error(y-test, test-prediction)

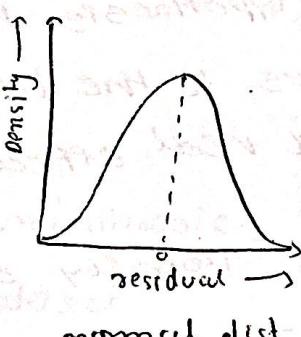
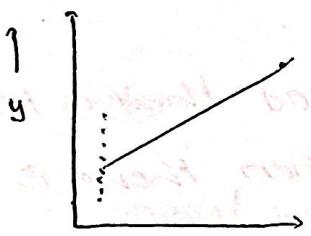
11. # RMSE

np.sqrt(mean\_squared\_error(y-test, test-predictions))

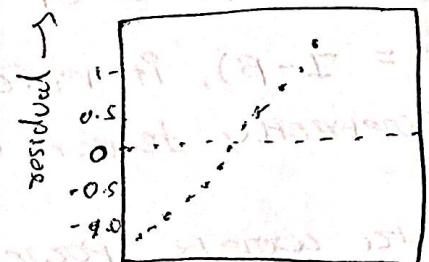
## => Evaluating Residuals

-  $|y_i - \hat{y}_i|$  => Residuals

- plot the below 3 plots



normal dist.  
(must be around  
'0' mean)



residual vs y plot  
(point must be  
random)

12. test-residuals = y-test - test-prediction

- no curve

13. sns.scatterplot(x=y-test, y=test-residuals)

- no line

- no pattern

## => Model Deployment

14 final\_model = LinearRegression()

15 model.fit(x, y)

16 model.coef\_ # return β value for all the features

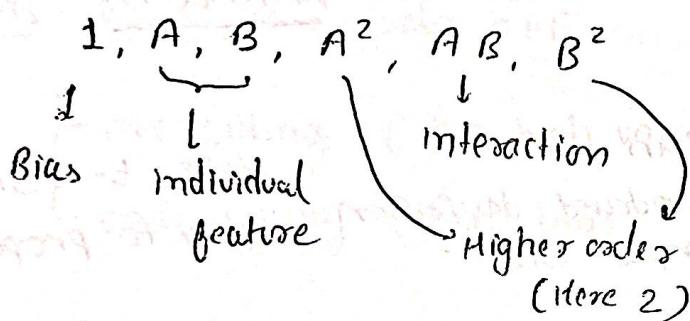
17 from joblib import dump, load

18 dump(final\_model, 'final-sales-model.joblib') # saving

19 loaded\_model = load('final-sales-model.joblib') # reload.

## i. o) Polynomial Regression

- Sometimes A and B feature may not be that useful, but their interaction can be.
- So, we create new features from original feature using scikit learn PolynomialFeatures.
- Converting Two features A and B



$\Rightarrow$  Code :

1. import important libraries
2. load / Read dataset
3. Identify features and label ( $X$  and  $y$ ) [200, 3] Feature labels
4. from sklearn.preprocessing import PolynomialFeatures
5. Polynomial\_converter = PolynomialFeatures (degree=2)
6. polynomial\_converter.fit( $X$ )
7. poly\_feature = polynomial\_converter.transform( $X$ ) [200, 9]
8. poly\_feature = polynomial\_converter.fit\_transform( $X$ )
9. from sklearn.model\_selection import train-test-split
10.  $X\_train, X\_test, y\_train, y\_test = train\_test\_split$  (poly-feature, y, testsize=0.2, random\_state=101)

```

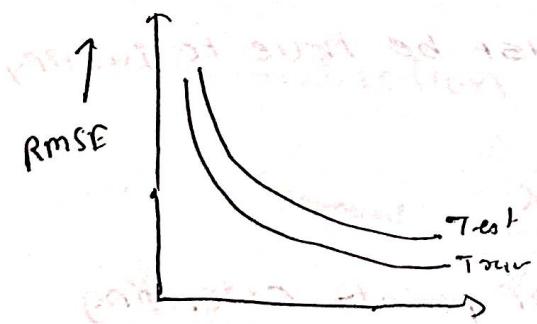
11 import from sklearn.linear_model import LinearRegression
12 model = LinearRegression()
13 model.fit(x_train, y_train)
14 test_predict = model.predict(x_test)
15 from sklearn.metrics import mean_absolute_error,
                           mean_squared_error
16 MAE = mean_absolute_error(y_test, test_prediction)

```

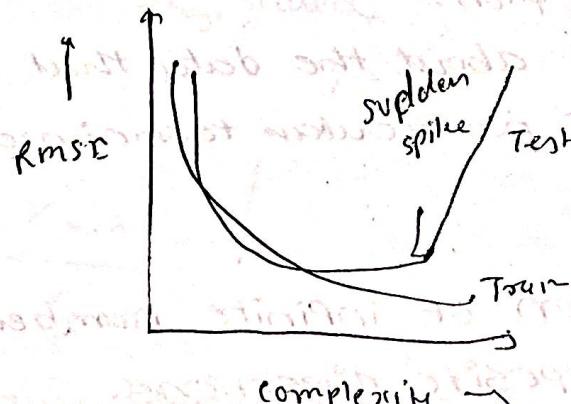
$\Rightarrow$  Bias - Variance Trade off

- overfitting (high variance)
- underfitting (high bias)

- Increasing the degree of polynomial increases the complexity which leads to overfitting problem.
- And lesser complexity than required creates problem of underfitting
- Hence, with very less or very high degree of polynomial RMSE increases



Ideal situation



overfitting at high complexity.

## ⇒ Regularization

(20)

- 1) Minimize model complexity
- 2) Penalize the loss function
- 3) Reduce model overfitting

- The 3 main types of regularization

- 1) L1 Reg. (Lasso regression)
- 2) L2 Reg. (Ridge regression)
- 3) combine L1 and L2 (Elastic Net)

- Feature Scaling

- If our data or features are on different scale, we use feature scaling
- 2 main ways of scaling features:

1. Standardization (rescale data to have  $\mu = 0$  &  $\sigma^2 = 1$ )

2. Normalization (rescale all data values betn 0-1)

- 1. Standardization:

$$X_{\text{changed}} = \frac{X - \mu}{\sigma} \quad X = \text{original} \quad \mu = \text{mean} \quad \sigma = \text{standard deviation}$$

- 2. Normalization

$$X_{\text{changed}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

- Fit & transform data classes are used for scaling
- Simply calculates the necessary statistics  
actually scales data one  
return new scaled data

## - Feature scaling process

1. Perform train test split
2. Fit to training feature data
3. Transform training feature data
4. Transform test feature data.

=> code

1. from sklearn.preprocessing import StandardScaler
2. scaler = StandardScaler()
3. scaler.fit(x\_train)
4. x\_train\_scale = scaler.transform(x\_train)
5. x\_test\_scale = scaler.transform(x\_test)

## 1) L2 Regularization or Ridge Regression

- general formula for regression line

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p$$

- These  $\beta$  coefficient were solved by minimizing RSS (Residual sum of squares)

$$\Rightarrow \text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\Rightarrow \sum_{i=1}^n (y_i - \hat{\beta}_0 x_{i1} - \hat{\beta}_1 x_{i2} - \dots - \hat{\beta}_p x_{ip})^2$$

$$= \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + \lambda \sum_{j=1}^p \beta_j^2$$

Penalty term

- Ridge regression minimize RSS + Penalty term

=> code:

1. From sklearn import Ridge Hyperparameters
2. ridge-model = Ridge (alpha=10)
3. ridge-model.fit (x-train, y-train)
4. test-predict = ridge-model.predict (x-test)
5. calculate error MAE & RMSE
6. test-prediction = ridge-cv-model.predict (x-test) error string
- # Ridge regression with cross validation
7. ridge-cv-model = RidgeCV (alphas=(0.1, 1.0, 10.0), scoring='r2')
8. ridge-cv-model.fit (x-train, y-train)
9. ridge-cv-model.alpha - # gives best performing alpha value
10. test-prediction = ridge-cv-model.predict (x-test)
11. MAE = mean\_absolute\_error (y-test, test-predictions)

## 2) L1 Regularization or LASSO regression

- It adds penalty equal to the absolute value of the magnitude of coefficient
- It can yield sparse model where some coefficient can be zero

$$\text{Error} = \sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij})^2 + |\lambda \sum_{j=1}^p \beta_j| \rightarrow \begin{matrix} \text{Penalty} \\ \text{term} \end{matrix}$$

- LASSO - least absolute shrinkage and selection operator

$\Rightarrow$  code

1. from sklearn.linear\_model import LassoCV
2. lasso\_cv\_model = LassoCV (eps=0.001, n\_alphas=100, cv=5)
3. lasso\_cv\_model.fit (x\_train, y\_train)
4. lasso\_cv\_model.alpha\_
5. test\_prediction = lasso\_cv\_model.predict (x\_test)
6. calculate the error

$\Rightarrow$  Elastic net

- minimize  $\beta \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\}$  subject to  $\sum_{j=1}^p |\beta_j| \leq S$
- minimize  $\beta \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\}$  subject to  $\sum_{j=1}^p \beta_j^2 \leq S$

- For 2 features

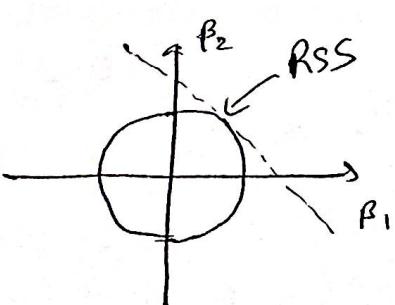
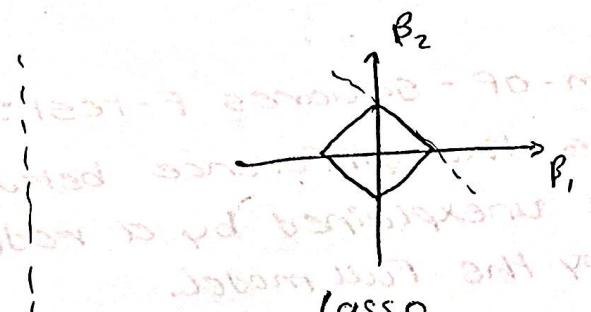
$$\hat{y} = \beta_1 x_1 + \beta_2 x_2$$

Lasso regression

$$|\beta_1| + |\beta_2| \leq S$$

Ridge regression

$$\beta_1^2 + \beta_2^2 \leq S$$



$$\begin{aligned}
 - \text{Error} &= \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda_1 \sum_{j=1}^p \beta_j^2 + \lambda_2 \sum_{j=1}^p |\beta_j| \\
 &= \frac{\sum_{i=1}^n (y_i - x_i \hat{\beta})^2}{2n} + \lambda \left( \frac{1-\alpha}{2} \sum_{j=1}^m \hat{\beta}_j^2 + \alpha \sum_{j=1}^m |\hat{\beta}_j| \right)
 \end{aligned}$$

$\Rightarrow$  code

1. from sklearn.linear\_model import ElasticNetCV
2. Elastic-model = ElasticNetCV (l1\_ratio = [0.1, 0.5, 0.7, 0.9, 0.95],  
 $\epsilon = 0.001$ , n\_alpha = 100,  
 $\text{max\_iters} = 1000000$ )
3. Elastic-model.fit(x-train, y-train)
4. Elastic-model.l1\_ratio\_ # gives the best l1-ratio
5. test\_prediction = elastic-model.predict(x-test)

## \* Feature Engineering

- Extracting features from raw data via data mining techniques.

- 3 approach to perform feature Engineering

1. Extracting info.

2. Combining info.

3. Transforming info.

[ Integer Encoding (convert categories into 1, 2, 3.. N)  
One-hot Encoding (convert each categories into 0 & 1) ]

## \* Data cleaning

### 1) Dealing with outliers:

- Outlier consideration depends on

1. Range and limit

2. Percentage of Data

- Highly dependent on domain

- Outlier with some methodology

  - Inter Quartile Range

  - Standard Deviation

  - Visualized

- $IQR = Q_3 - Q_1$

- lower-limit =  $Q_1 - 1.5 \cdot IQR$

- upper-limit =  $Q_3 + 1.5 \cdot IQR$

code:

```
q75, q25 = np.percentile(df, [75, 25]) # return array of Q3 & Q1
```

```
IQR = q75 - q25
```

## 2) Missing value imputation: (20)

- Missing value can be NAN, NULL, 0 etc

- code:

1. import libraries

2. import Dataset as df

3. df = ds.drop ("PID", axis=1) # Drop unique column

4. 100 \* df.isnull().sum() / len(df) # percentage of null values

5. def percent\_missing (df):

percent\_nan = 100 \* df.isnull().sum() / len(df)

percent\_num = percent\_nan [percent\_nan > 0], sort\_values()

return percent\_num

6. percent\_missing (df) # gives the % of null values in all columns.

7. sns.barplot (x = percent\_num.index, y = percent\_num)

plt.xticks (rotation = 90) # optional to visualize the null values

## 8. If the percentage of null values is less than 1%.

then, we will drop or fill few rows missing the data feature

9. percent\_num [percent\_num < 1] # gives all the feature whose % is < 1

10. df [ds ['Electrical'].isnull()] # gives the rows which has Electrical as null value

## 11. Now, If the feature column has 'Nan' or '0' for no data registered, we can fill na with respective values else, we have to delete the row

12. df = df.dropna (axis=0, subset = ['Electrical', 'GarageCars'])

13. # NUMERIC COLUMNS --> Fill na 0

num-cols = ['col.name1', 'col.name2']

df[num-cols] = df[num-cols].fillna(0)

# STRING COLUMNS --> fill na 'none'

str-cols = ['col.name1', 'col.name2']

df[str-cols] = df[str-cols].fillna('None')

## 14. Now, If the percentage of null values is greater than 1% then, we will drop or fill columns. If percentage is greater than 50%, we will drop all those columns

15. df = df.drop(['poolQC', 'miscFeature', 'Alley'])

## 16. If column has no categorical is filled with NA or 0, will fillna with 'Na' or '0' for respective column. And for other numeric column, we will fillna with some statistical computation.

17. df['Garage Yr Blt'] = df['Garage Yr Blt'].fillna(0)

18. dd.groupby('Neighborhood')[['LotFrontage']].transform

(lambda value: value.fillna(value.mean()))

↳ std = StandardScaler()

↳ scaled = std.fit\_transform(df)

↳ scaled = scaled.transform(X\_train)

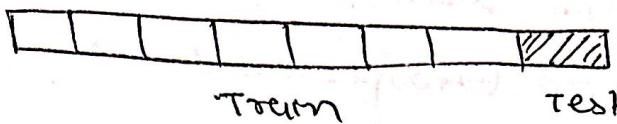
↳ scaled = scaled.fillna(scaled.mean())

↳ scaled = scaled.replace([np.inf, -np.inf], np.nan).replace(np.nan, 0)

↳ scaled = scaled.fillna(scaled.mean())

↳ scaled = scaled.fillna(scaled.mean())

## \* Cross validation



Process general Train / Test split

### $\Rightarrow$ Train / Test split procedure

- clean and adjust data for  $x$  only
- split data in Train / Test split for both  $x$  and  $y$
- fit / Train scales on Training  $x$  Data
- Scale  $x$  Test Data
- Create model
- Fit / Train model on  $x$  Train
- Evaluate model on  $x$  Test data

### $\Rightarrow$ code

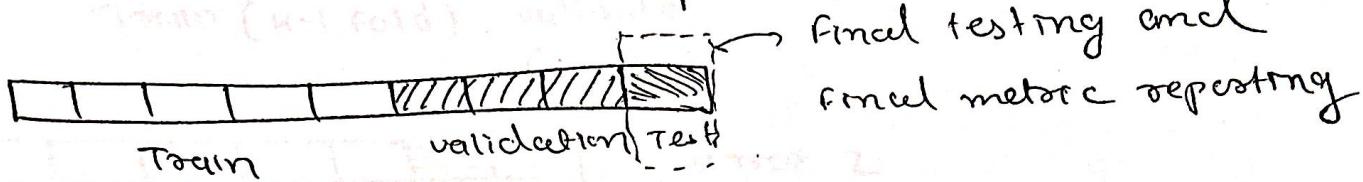
1.  $x = df.drop(['sales'])$
2.  $y = df['sales']$
3.  $x\_train, x\_test, y\_train, y\_test = train-test-split(x, y, test-size=0.3, random-state=42)$
4. from sklearn.preprocessing import StandardScaler
5.  $scaler = StandardScaler()$
6.  $scaler.fit(x\_train)$
7.  $x\_train = scaler.transform(x\_train)$
8.  $x\_test = scaler.transform(x\_test)$
9. from sklearn.linear\_model import Ridge
10.  $model = Ridge(alpha=1)$
11.  $model.fit(x\_train, y\_train)$

12.  $y_{\text{pred}} = \text{model.predict}(X_{\text{test}})$

13. From sklearn.metrics import mean\_squared\_error

14.  $\text{MSE} = \text{mean_squared_error}(y_{\text{test}}, y_{\text{pred}})$

$\Rightarrow$  Train, validation & test split



- we will do this by performing train-test-split() twice

1. Once split off larger training set.

2. remaining data into a validation & test set

- Code:  $\# 70: 15!15$  (similar to general split)

•  $x_{\text{train}}, x_{\text{other}}, y_{\text{train}}, y_{\text{other}} = \text{train-test-split}(x, y, \text{test\_size}=0.3, \text{random\_state}=101)$

•  $x_{\text{eval}}, x_{\text{test}}, y_{\text{eval}}, y_{\text{test}} = \text{train-test-split}(x_{\text{other}}, y_{\text{other}}, \text{test\_size}=0.5, \text{random\_state}=101)$

•  $x_{\text{eval}} = \text{scaler.transform}(x_{\text{eval}})$

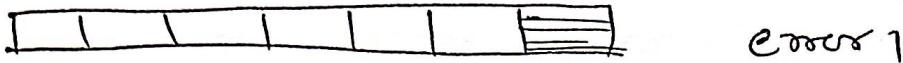
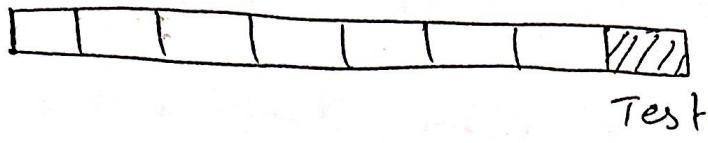
•  $y_{\text{eval}} = \text{regressor.predict}(x_{\text{eval}})$

•  $\text{regressor} = \text{RandomForestRegressor}(\text{n_estimators}=100, \text{random\_state}=101)$

•  $\text{cross_val_score}(\text{regressor}, \text{x}, \text{y})$  (will return array of scores)

•  $\text{GridSearchCV}(\text{regressor}, \text{grid}, \text{cv}=5)$  (will return array of scores)

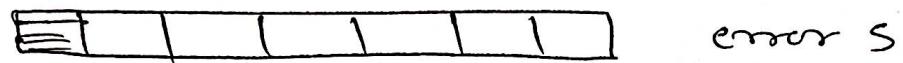
$\Rightarrow$  K-Fold cross validation



Train (K-1 fold) validate (1 fold)



? upto all possible folds.



-----  
avg error

= codes

~~1. df['sales']~~,

1.  $X = df.drop(['sales'], axis=1]$

2.  $y = df['sales']$

3.  $x\_train, x\_test, y\_train, y\_test = train\_test\_split(X, test\_size=0.3, random\_state=101)$

4. from sklearn.preprocessing import StandardScaler

5.  $scaler = StandardScaler()$

6.  $scaler.fit(x\_train)$

7.  $x\_train = scaler.transform(x\_train)$

8.  $x\_test = scaler.transform(y\_test)$  9.  $model = Ridge(alpha=1,$

10. from sklearn.model\_selection import cross\_val\_score

11.  $scores = cross_val_score(model, x\_train, y\_train,$

scoring='neg\_mean\_squared\_error', cv=5)

12. `abs(scores.mean())` # report avg error

13. `model.fit(X-train, y-train)`

14. `y-final-pred = model.predict(X-test)`

15. `mean_squared_error(y-test, y-final-pred)`

- We can also use `cross_val_score()` function similarly as `cross_val_scores()`, it just gives more information like fit-time and score-time.

### \* Grid search

- A grid search is a way of training and validating a model on every possible combination of multiple hyperparameter options.

=> code

1. create X and y

2. Train and Test split

3. scale data

4. from sklearn.linear\_model import ElasticNet

5. base\_elastic\_model = ElasticNet()

6. param\_grid = { 'alpha': [0.1, 1, 5, 10, 50, 100],  
                   'l1\_ratio': [.1, .5, .7, .95, .99, 1] }

7. from sklearn.model\_selection import GridSearchCV

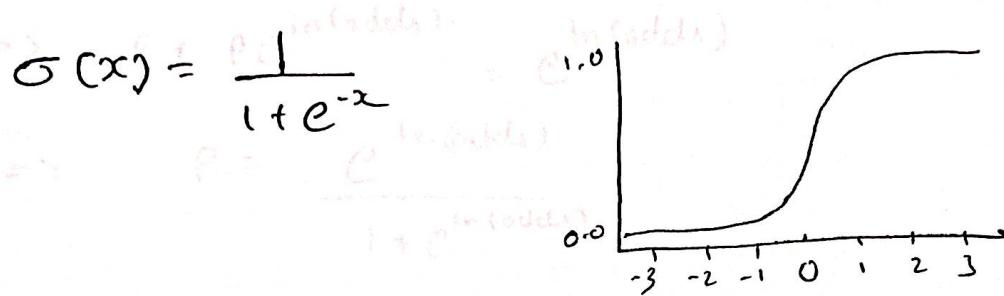
8. grid\_model = GridSearchCV( estimator = base\_elastic\_mod  
                           param\_grid = param\_grid,  
                           scoring = 'neg\_mean\_squared\_err  
                           cv = 5, verbose = 2 )

9. grid\_model.fit(X-train, y-train)

## 2) Logistic Regression

(32)

- It is a classification algorithm. It allows us to predict a categorical label based on historical feature data.
- The categorical column (label) may have two or more discrete class labels.
- logistic function or sigmoid function



$\Rightarrow y_i$  are trying to maximize the likelihood

$$y_i = \frac{1}{1+e^{-\sum_{i=1}^n \beta_i x_i}} \quad [ \because \hat{y}(x) = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n ]$$

$$\sigma(x) = \frac{1}{1+e^{-x}} \quad \therefore \text{putting } x = \hat{y}(x)$$

$$\Rightarrow \hat{y} + \hat{y} e^{-\sum_{i=1}^n \beta_i x_i} = 1$$

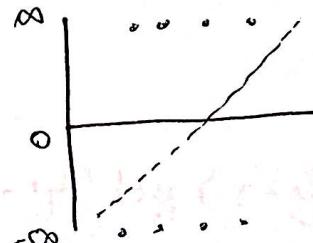
$$\Rightarrow \hat{y} e^{-\sum_{i=1}^n \beta_i x_i} = 1 - \hat{y}$$

$$\Rightarrow e^{-\sum_{i=1}^n \beta_i x_i} = \frac{1 - \hat{y}}{\hat{y}}$$

$$\Rightarrow e^{\sum_{i=1}^n \beta_i x_i} = \frac{\hat{y}}{1 - \hat{y}} \rightarrow \text{odds of log}$$

$$\Rightarrow \ln\left(\frac{\hat{y}}{1 - \hat{y}}\right) = \sum_{i=0}^n \beta_i x_i$$

odds is ratio of probability of event occurring to probability of event not occurring



⇒ 1<sup>st</sup> step for maximum likelihood is to go from log odds back to probability (32)

$$\Rightarrow \ln\left(\frac{P}{1-P}\right) = \text{ln}(\text{odds})$$

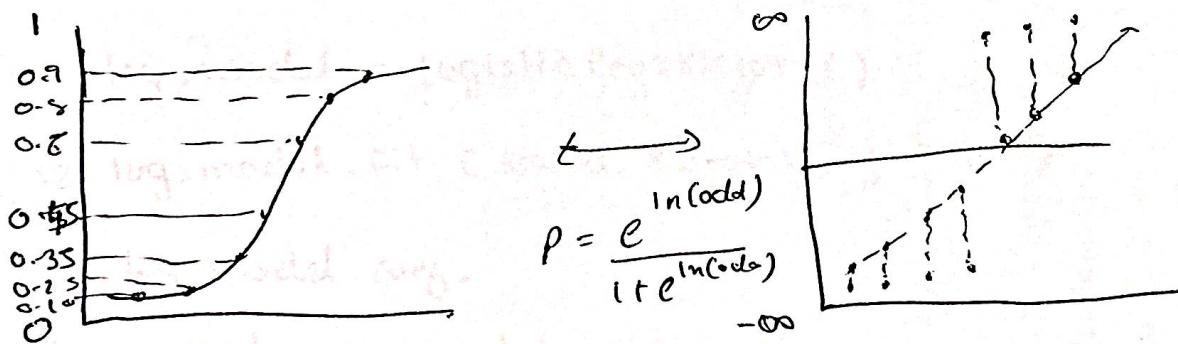
$$\Rightarrow \frac{P}{1-P} = e^{\text{ln}(\text{odds})}$$

$$\Rightarrow P = \frac{e^{\text{ln}(\text{odds})}}{1 + e^{\text{ln}(\text{odds})}}$$

$$\Rightarrow P + PC^{\text{ln}(\text{odds})} = e^{\text{ln}(\text{odds})}$$

$$\Rightarrow P = \frac{C^{\text{ln}(\text{odds})}}{1 + e^{\text{ln}(\text{odds})}}$$

⇒ we are trying to maximize the likelihood



$$\text{so likelihood} = 0.9 \times 0.8 \times 0.4 \times (1 - 0.45) \times (1 - 0.35) \times (1 - 0.25) \times (1 - 0.10)$$

maximize likelihood

this by try different line of best fit.

⇒ In terms of cost function, we seek to minimize the following (log loss):

$$J(x) = -\frac{1}{m} \sum_{i=1}^m y^i \log(\hat{y}^i) + (1 - y^i) \log(1 - \hat{y}^i)$$

(34)

$\Rightarrow$  Code

1. import important libraries
2.  $x = df.\text{drop}(\text{'test\_result'}, \text{axis}=1)$
3.  $y = df[\text{'test\_result'}]$
4. from sklearn.model\_selection import train\_test\_split
5. from sklearn.preprocessing import StandardScaler
6.  $X\_train, X\_test, y\_train, y\_test = \text{train\_test\_split}(X, y,$   
 $\quad \quad \quad \text{test\_size=0.1,}$   
 $\quad \quad \quad \text{random\_state=101})$
7.  $\text{scaler} = \text{StandardScaler}()$
8.  $\text{scaled\_x\_train} = \text{scaler}.\text{fit\_transform}(X\_train)$
9.  $\text{scaled\_x\_test} = \text{scaler}.\text{transform}(X\_test)$
10. from sklearn.linear\_model import LogisticRegression()
11.  $\text{log\_model} = \text{LogisticRegression}()$
12.  $\text{log\_model}.\text{fit}(\text{scaled\_x\_train}, y\_train)$
13.  $\text{log\_model}.\text{coef\_}$
14.  $y\_pred = \text{log\_model}.\text{predict}(\text{scaled\_x\_test})$

$\Rightarrow$  classification evalution metrics

- Confusion matrix

		Actual	
		Infected	Healthy
Predicted	Infected	True +ve	False -ve
	Healthy	False -ve	True +ve

$$\cdot \text{Accuracy} = \frac{(TP+TN)}{\text{Total}}$$

$$\cdot \text{Recall} = \frac{TP}{\text{Total actual +ve}}$$

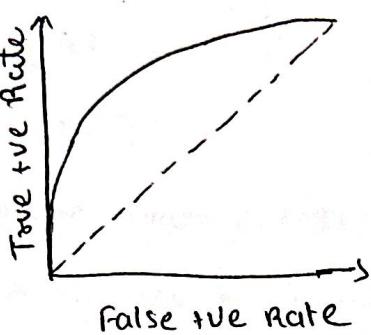
$$\cdot \text{Precision} = \frac{TP}{\text{Total predicted +ve}}$$

$$\cdot \text{F1 Score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- ROC curve : Receiver Operator characteristic curve

(35)

- we gladly accept more false +ve to reduce false negatives



=> code:

1. from sklearn.metrics import accuracy\_score, confusion\_matrix, classification\_report.
2. y-pred = log\_model.predict(scaled\_x-test)
3. accuracy\_score(y-test, y-pred)
4. confusion\_matrix(y-test, y-pred)
5. from sklearn.metrics import plot\_confusion\_matrix
6. plot\_confusion\_matrix(log\_model, scaled\_x-test, y-test)
7. class\_report(classification\_report(y-test, y-pred)) # gives precision, recall, f1-score and support.
8. precision\_score() # Similarly as confusion matrix  
recall\_score() # " "
9. from sklearn.metrics import plot\_precision\_recall\_curve, plot\_roc\_curve
10. plot\_roc\_curve(log\_model, scaled\_x-test, y-test)
11. plot\_precision\_recall\_curve(log\_model, scaled\_x-test, y-test)
12. log\_model.predict\_proba(scaled\_x-test) # get probabilities of belong class

## \* Multiclass Logistic Regression

(36)

1. import libraries
2. create X and y
3. create train / test split
4. fit transform train } scale data
5. transform test
6. from sklearn.linear\_model import LogisticRegression
7. from sklearn.model\_selection import GridSearchCV
8. log-model = LogisticRegression (solver='saga', multi-class='ovr', max\_iter=5000)
9. penalty = ['L1', 'L2', 'elasticnet']  
l1\_ratio = np.linspace(0, 1, 20)  
 $C = \text{np.logspace}(0, 10, 20)$   
param\_grid = {'penalty': penalty, 'l1\_ratio': l1\_ratio, 'C': C}
10. grid-model = GridSearchCV(log-model, param\_grid=param\_grid)
11. grid-model.fit(scaled-X-train, y-train)
12. Evaluate using metrics.

### 3) KNN

- It assigns a label to new data based on the distance between the old data and new data.
- K is the no. of distance point, we consider to label new data point.
- The label which has major no. of shortest distance the new point gets assigned to that label.
- Eg for  $K=5$ , 3 distance are shortest to class A data point  
2 distance are shortest to class B data point  
then the new data point belongs to class A.
- We want K value that minimizes error  
 $\text{error} = 1 - \text{Accuracy}$
- KNN considerations for distance metric:
  - minkowski
  - Euclidean
  - manhattan
  - chebyshew
- Things to keep in mind for KNN algorithm.
  - choosing optimal K value.
  - scaling features.

## $\Rightarrow$ code

1. Import important libraries
2. Import dataset
3. Create X and y
4. Create train-test-split
5. Scale data
6. From sklearn.neighbors import KNeighborsClassifier
7. knn-model = KNeighborsClassifier()
8. knn-model-fit (scaled\_x-train, y-train)
9. y-pred = knn-model.predict (scaled\_x-tester)
- 10 Evaluate errors # confusion metrics and classification report.

## $\Rightarrow$ Pipeline

- Pipeline object can set up a sequence of repeated operations such as a scalar and a model.

## $\Rightarrow$ code

1. scalar = StandardScaler()
2. knn = KNeighborsClassifier()
3. knn.get\_params().keys() # strange code for parameters
4. operations = [( 'scalar', scalar ), ( 'knn', knn )]
5. from sklearn.pipeline import Pipeline
6. pipe = Pipeline (operations)
7. from sklearn.model\_selection import GridSearchCV
8. k-values = list (range (1, 20))
9. param\_grid = { 'KNeighborsClassifier': [ 'Knn-n\_neighbors': k-values ] }

## 4) Support Vector Machines

## Maximum margin classifier

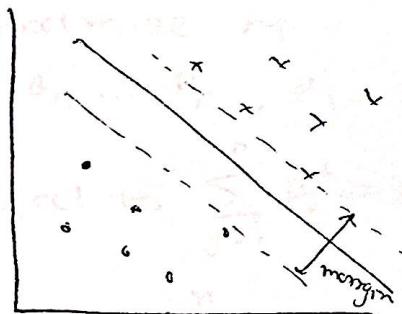
1

Support vector classifier (SVC) (add boxes to allow some mis-classification)

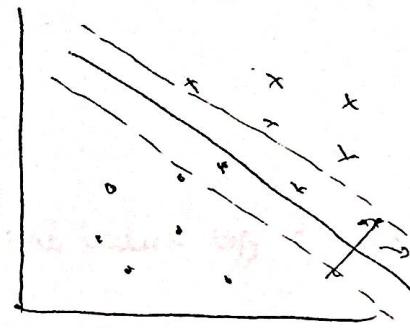
4

## Support vector machines (SVM)

- For N dimensional space, a hyperplane is  $N-1$  dimensional
  - 1-D Hyperplane is a single point
  - 2-D " " " line
  - 3-D " " " plane



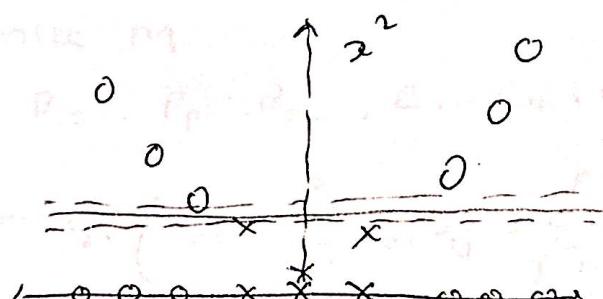
maximum margin



~~SVC~~

of margins  
that allow  
some miss.  
-classification

- SVMs use kernel to project the data into higher dimensions in order to use hyperplane.



⇒ Mathematics of ~~Linear~~ ~~Linear~~

$$x_1 = \begin{pmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1P} \end{pmatrix}, \dots, x_n = \begin{pmatrix} x_{n1} \\ x_{n2} \\ \vdots \\ x_{nP} \end{pmatrix}$$

- maximize  $M$

$$\beta_0, \beta_1, \dots, \beta_P, M$$

$$\text{subject to } \sum_{j=1}^P \beta_j^2 = 1$$

$$- y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_P x_{iP}) \geq M \quad \forall i=1 \dots n$$

⇒ Support vector classifier

- maximize  $M$

$$\beta_0, \beta_1, \dots, \beta_P, \epsilon_1, \dots, \epsilon_n, M$$

$$\text{subject to } \sum_{j=1}^P \beta_j^2 = 1$$

$$- \epsilon_i \geq 0, \sum_{l=1}^n \epsilon_l \leq C \quad \begin{array}{l} \text{Higher the value of } C \text{ less bias} \\ (\text{MMC}) \end{array}$$

$$- y_i (\beta_0 + \beta_1 x_{i1} + \dots + \beta_P x_{iP}) \geq M(1 - \epsilon_i)$$

⇒ support vector machines

$$x_1, x_1^2, x_2, x_2^2, \dots, x_P, x_P^2$$

maximize  $M$

$$\beta_0, \beta_{11}, \beta_{12}, \dots, \beta_{P1}, \beta_{P2}, \epsilon_1, \dots, \epsilon_n, M$$

$$\text{subject to } y_i \left( \beta_0 + \sum_{j=1}^P \beta_{j1} x_{ij} + \sum_{j=1}^P \beta_{j2} x_{ij}^2 \right) \geq M(1 - \epsilon_i)$$

$$\sum_{i=1}^n \epsilon_i \leq C, \epsilon_i \geq 0, \sum_{j=1}^P \sum_{k=1}^2 \beta_{jk}^2 = 1$$

=> To deal with very large feature space we use Kernel trick, which uses the inner product of vectors (dot product)

- Dot product

$$\langle \mathbf{a}, \mathbf{b} \rangle = \sum_{i=1}^r a_i b_i$$

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + a_2 b_2$$

$$\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}| |\mathbf{b}| \cos \theta$$

- SVC rewritten as

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n [\alpha_i] \langle \mathbf{x}, \mathbf{x}_i \rangle$$

↓

$$f(\mathbf{x}) = \beta_0 + \sum_{i \in S} \alpha_i \langle \mathbf{x}, \mathbf{x}_i \rangle$$

- Kernel function

$$\underbrace{\langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle}_{\{ \langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle \}} = \sum_{j=1}^r x_{ij} x_{i'j} \quad f(\mathbf{x}) = \beta_0 + \sum_{i \in S} \alpha_i \underbrace{\langle \mathbf{x}, \mathbf{x}_i \rangle}_{\langle \mathbf{x}, \mathbf{x}_i \rangle}$$

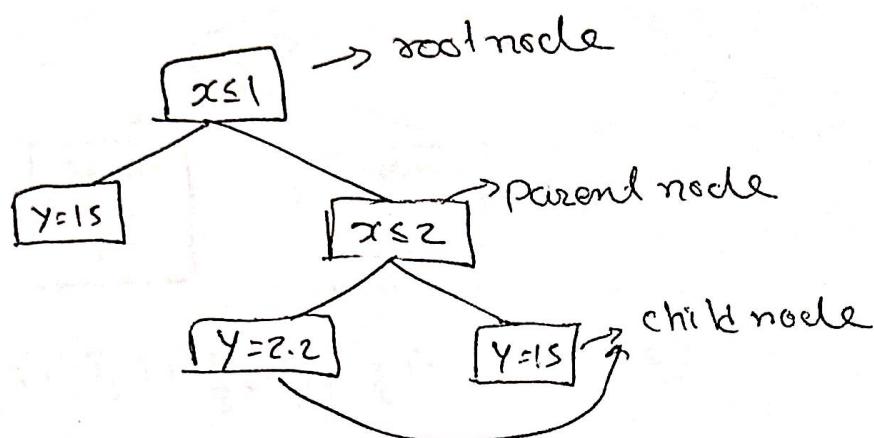
$\downarrow$

$$\underbrace{\langle \mathbf{a}, \mathbf{b} \rangle}_{\{ \langle \mathbf{a}, \mathbf{b} \rangle \}} = \sum_{i=1}^r a_i b_i$$

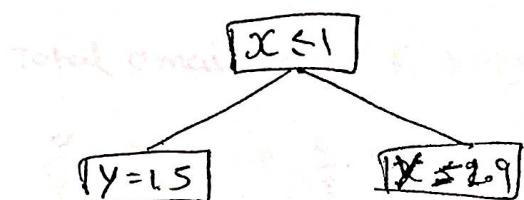
⇒ code:

1. from sklearn.svm import SVC
  2. create X and y
  3. model = SVC (kernel='linear', C=1000) # C is inversely proportional to the margin.  
model.fit(X, y)
  4. from svm\_margins import plot\_svm\_boundary
  5. plot\_svm\_boundary (model, X, y)
- # A low C make decision surface smooth, i.e makes soft margin. whereas high C value makes maximum margin. low C value is better.

## Decision Tree



↓ pruning (to short the length of tree to avoid overfitting)



⇒ Gini impurity

- It is measurement of how 'pure' the information in a dataset is.

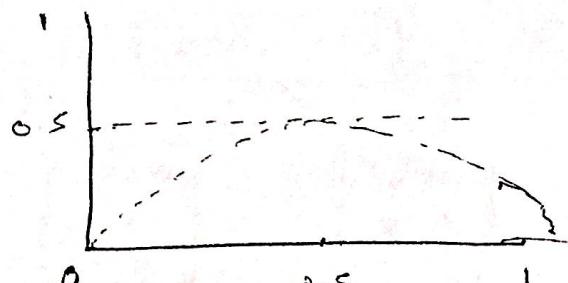
$$G(Q) = \sum_{c \in C} p_c(1 - p_c)$$

↖ probability of class c  
↖ classes  
↙ dataset

- e.g.

X	0
X	0

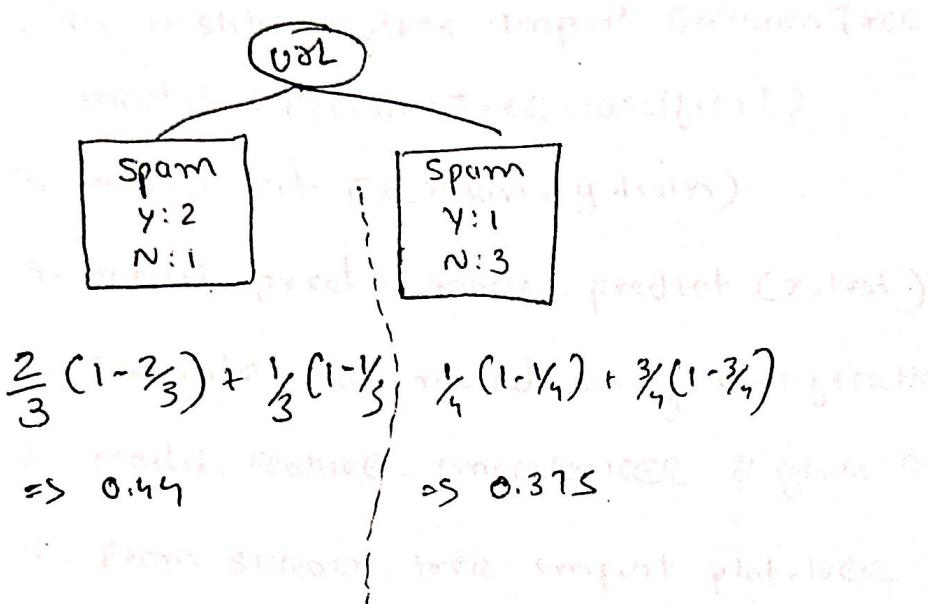
$P(\text{class A}) = \frac{2}{4}(1 - \frac{2}{4}) = 0.25$   
 $P(\text{class B}) = \frac{2}{4}(1 - \frac{2}{4}) = 0.25$



$$\therefore \text{gini impurity} = 0.25 + 0.25 = 0.50$$

- We want to minimize the gini impurity at leaf nodes.

$\Rightarrow$  Gini impurity for binary categorical feature



- Total emails = total 3 left emails & 4 Right emails

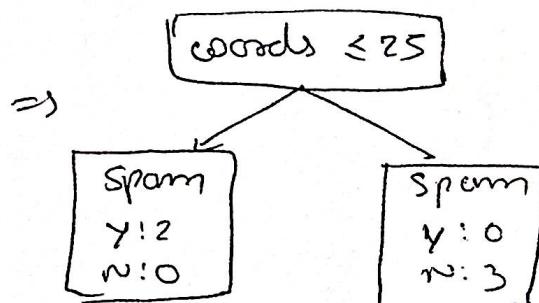
$$\frac{3}{7} \times 0.44 + \frac{4}{7} \times 0.375$$

$$\Rightarrow \underline{\underline{0.375}}$$

$\Rightarrow$  Gini impurity for continuous feature

X-coords      Y-Spam

[10]	-----	yes	Gini = 0.3
[25]	-----	yes	Gini = 0
[35]	-----	no	Gini = 0.25
[45]	-----	no	Gini = 0.4
50	-----	no	



- Choose the best split value such that we get lowest gini.

## Random Forest (RF cannot be overfitted)

- Random forest is a decision tree made on different subset of features.
- Unique Random Forest hyperparameters:
  1. no. of Estimators - How many DT to use in forest?
  2. no. of Features - How many feature to include in each
  3. Bootstrap samples - allows bootstrap sampling of each training subset
  4. Out of Bag Error - calculate oob error during training.

1. no. of Estimators: 68 - 128 DT are good, we can use as many DT as we want. There is no overfitting. Trees just gets duplicated. start with "100".

2. no of Feature: According to original publication subset of  $\log_2(n+1)$  random features is subset ( $n$ ) total feature.

Other methods of choosing are:

-  $\sqrt{n}$

-  $n/3$

3. Bootstrap samples: Allowing to choose same row more than once while training DT subset.

- It affects the overall result of DT.

4. OOB error: we fit the rows that are not considered by bootstrapping during training, and get the  $y$  and calculate the OOB error.

- OOB doesn't affect the model.

$\Rightarrow$  code!

1. from sklearn.ensemble import RandomForestClassifier.
2. model = RandomForestClassifier(n\_estimators=10,  
max\_features='auto',  
random\_state=101)
3. model.fit(X\_train, y\_train)
4. model\_pred = model.predict(X\_test)
5. Evaluate.
6. Performe Grid search cv to choose best Hyper parameter
7. # choosing correct no. of trees  
 $\text{test\_err} = []$   
 For n in range(1, 40):  
 model = RandomForestClassifier(n\_estimators=n,  
 max\_features='auto')  
 model.fit(X\_train, y\_train)  
 test\_pred = model.predict(X\_test)  
 test\_error.append(1 - accuracy\_score(test\_pred, y\_test))  
 plt.plot(range(1, 40), test\_error, label='TestError').  
 plt.legend()

