# Life Cycle of a Data Science Project

```
start → Business          →  Data Acquisition
         Understanding        & Understanding
                                      ↕
Deployment ← Modelling            Data source
    ↓              ↑               Pipeline
Customer      Feature Eng          Environment
Acceptance    Model training       cleaning wrangly
    ↓         Model eval           Exploration
  End
```

EDA
↳ Data Analysis         → feature Eng → making the data into
↳ Data preprocessing <                  right format to train
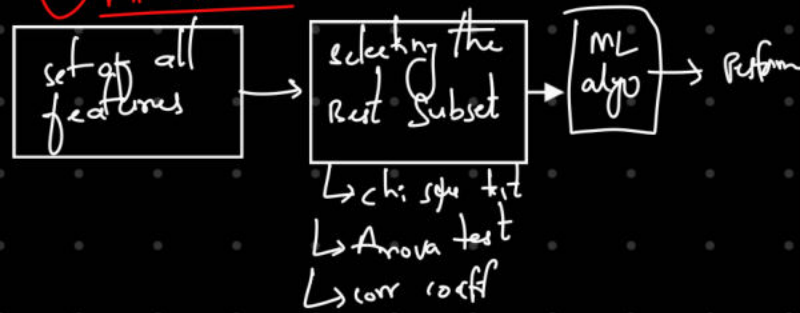                       → feature selection → only select    machine
                                             necessary feature
                                             to train

# feature selection techniques

① Univariate selection
② feature importance
③ Correlation matrix with heatmap.

① Filter method

```
set of all  →  Selecting the  →  ML   → Perform
features       Best Subset        algo
                    ↳ chi squ tst
                    ↳ Anova test
                    ↳ corr coeff
```

② Wrapper method
i> forward Selection :- first A train
                       next A & B if good keep
                       → A, B, C → 9 so like this

ii) Backward elimination :-
A, B, C, D, E → find which has low impact
                on our target variable.
                → remove that one.
                        ↖
by chi·2 test        remove
     ↓                → not imp → remove
  P≤0.05 imp    P>0.05              that

---

③ Recursive feature elimination :-
↳ greedy optimization algo.
→ keeps best performing subset in each iteration.

*→ These 2 techniques works only if the
dataset is small.

③ Embedded Methods

A B C D E → finds the permutation
            & comb of all the features

① Univariate Selection
Statistical test can be used to select those features
that have the (strongest rlnship with the output var)
→ the scikit-learn library provides the SelectKBest
class that can be used with a suite of different
statistical tests to select a specific no.of features
↳ internally it uses some tests eg: chi^2 test.

# Types of encoding

① Nominal Encoding —→ One hot encoding
                    → one hot En with
                      many categorical
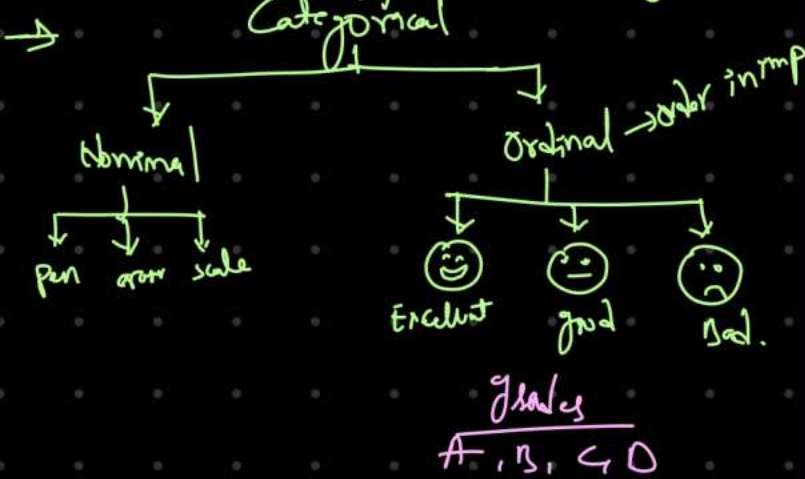                    → Mean Encoding

② ORDINAL ENCODING —→ Label Encoding
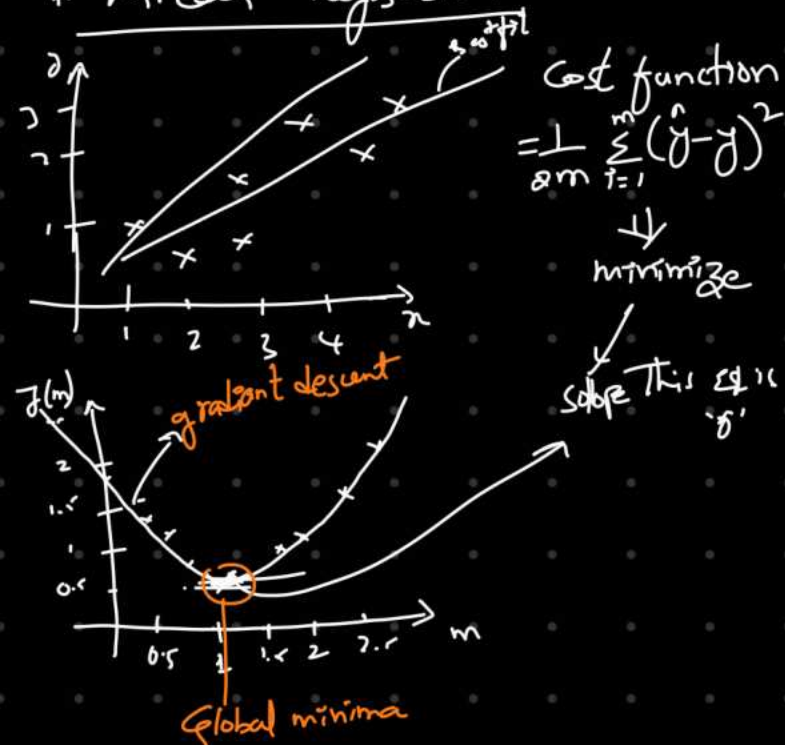                    → Target guided
                      ordinal Encoding.

# Handling missing values in Categorical variable

① Deleting the rows → may miss imp data
② Replace with the most frequent values
                              ↳ imbal dataset
③ Apply classifier algo. to predict → good
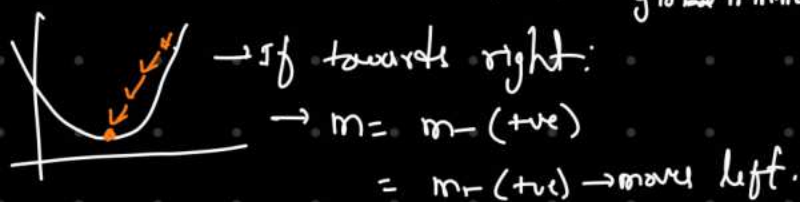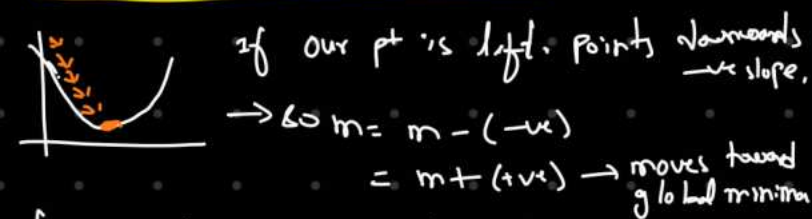④ Apply unsupervised ML → good.

# Ordinal Encoding / label encoding.

→ 
Categorical

- Nominal
  - pin
  - order scale
- Ordinal → order in imp
  - 😊 Excellent
  - 😐 good
  - 🙁 Bad.

grades
A, B, C D

---

# Linear Regression



Cost function
$$= \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}-y)^2$$

↓ minimize
↓
slope This eq is 'θ'

gradient descent

Global minima

→ Convergence Theorem :

$$m_{new} = m_{old} - \frac{\partial [J(m)]}{\partial m} \times \alpha$$ → learning rate
0.001 small.

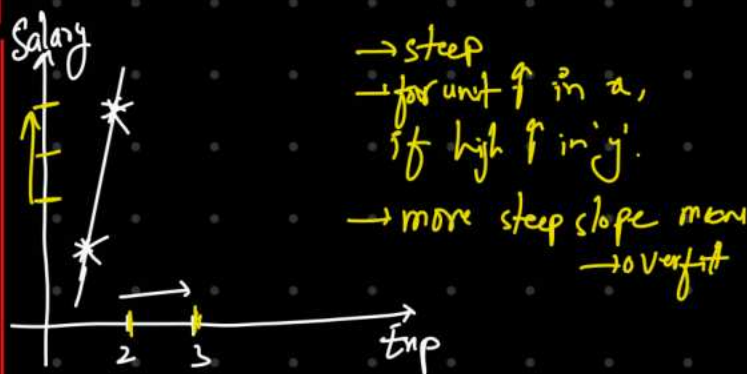If our pt is left, points downwards
→ve slope,
→ so m = m − (−ve)
= m + (+ve) → moves toward global minima

→ If towards right:
→ m = m − (+ve)
= m − (+ve) → moves left.

→ if α is big → ⟍⟋ jumps

---

# RIDGE & LASO REGRESSION



↳ high variance

Cost fn
Sum of residuals
$$\sum_{i=1}^{n} (y-\hat{y})^2 = 0.bbl$$
$$\hat{y} = m x + c$$

Over fitting → High Variance
↓ to
Low variance ⇒ Ridge

→ steep
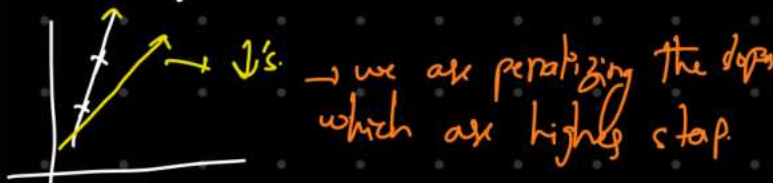→ for unit ↑ in x,
if high ↑ in 'y'.
→ more steep slope means
→ overfit

→ Cost function :
$$\sum_{i=1}^{n} (y-\hat{y})^2 + \lambda \times (slope)^2 = 0$$

↳ After the actual value reach 'θ', but still $\lambda(slope)^2$ is +ve, it even reduces.

$$\underbrace{\sum (y-\hat{y})^2}_{0} + \underbrace{\lambda}_{2} \times \underbrace{(2)^2}_{4}$$

$$0 + 4 \rightarrow \Sigma \text{ will even} \downarrow$$

→ ∫'s → we are penalizing the slope which are highly steep.

→ As we ↑ λ, slope will ↓ closer to θ.

→ lambda value is selected using CV.

→ $y = m_1 x_1 + m_2 x_2 + c$, ⇒ $\lambda \times [m_1^2 + m_2^2]$

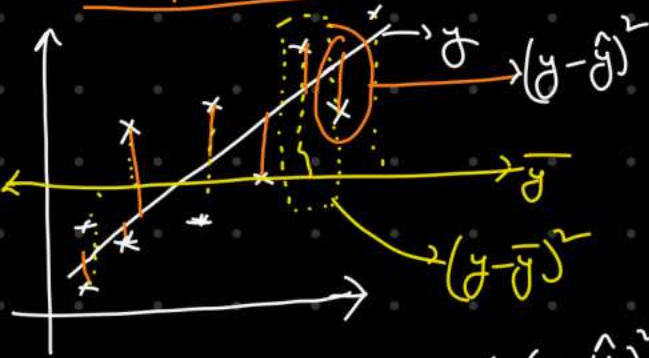⇒ Created Generalized model.

# Lasso : $L_2$ regularization

→ It not only fix overfitting, but also helps in feature selection.

$$\text{Cost fn:} \quad \sum_{i=1}^{n}(y-\hat{y})^2 + \lambda|\text{slope}|$$

→ $y = m_1 a_1 + m_2 a_2 + m_3 a_3 + m_4 a_4 + c_1$

→ $\lambda|m_1 + m_2 + m_3 + m_4|$

If slopes all closer to '0' remove them.

---

# R square and Adjusted R squar



$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} = \frac{\sum(y-\hat{y})^2}{\sum(y-\bar{y})^2}$$

→ $SS_{res}$ → sum of residuals.

→ $R^2$ → goodness of the fit

→ $1$ → good , $0$ → bad

→ If the best fit is worser than avg line, then

$1 - \frac{big}{small}$ → $R^2 = -ve$ value

---

# Adjusted $R^2$

→ $R^2$ always ↑'s as we go on ↑'g features even they are not correlated.

→ Adj. $R^2$ penalizes the features that are not correlated.

$$R^2_{adj} = 1 - \frac{(1-R^2)(N-1)}{N-p-1}$$

$R^2 = \text{sample} R$
$P = \text{no. of predictors}$
$N = \text{sample size}$

→ if ↑ ↑, features, will be added in denomi.

So, $R^2_{adj}$ will ↓.

↳ But if correlated, numeric multiplied by (1+1) also, $R^2$ ↑'s. so overall ↑'s.

→ Adj $R^2$ value always less than or equal to $R^2$ value.

---

→ Hypothesis Testing — Null $H_0$ / Alt $H_1$

① make assumption

② Collect evidence to prove $H_0$, if not fail and reject.



⇒ P-Value : It is the probability for the Null hypothesis to be true. {if $p<0.05$} reject.

$p = 0.01$



If we touch 100 times we only touch that 1 time.

---

Null hypothesis : Treats everything same or eql.

# METRICS IN CLASSIFICATION

1) Confusion matrix
2) FPR [Type 1 c-w)
3) FNR [T2 error)
4) Recall [TPR, sensitivity]
5) Precision (tve pred val)
6) Accuracy — for balance
7) $F_\beta$ score
8) Cohen Kappa
9) ROC curve, AUC score

10) PR curve



$A\beta = 0.5$

class labels / Probability

$0.3, 0.4$

1000 records

{ 500 Yes , 500 No
600    400
700    300
800    200 }

for in balanced.

| | 1 | 0 | Actual Value |
|---|---|---|---|
| Predicted Val 1 | TP | (FP) → Type 1 | |
| Predicted Val 0 | (FN) Type 2 | TN | |

$$\Rightarrow FPR = \frac{FP}{FP+TN}$$

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

A ←900   100→ B

$$Acc = \frac{900}{1000} \Rightarrow 90\%$$

1) **Precision** = Out of all the predicted +ve how many of them are actual +ve.

$$Precision = \frac{TP}{TP+FP}$$

→ spam detection :
→ if spam, Predicted not spam (FN)
→ if not spam, predicted spam XX
↓
so FP↓↓↓
⇒ we may miss imp mail

② **Recall** : Out of all +ve values, how many +ve values are predicted.

$$Recall = \frac{TP}{TP+FN}$$
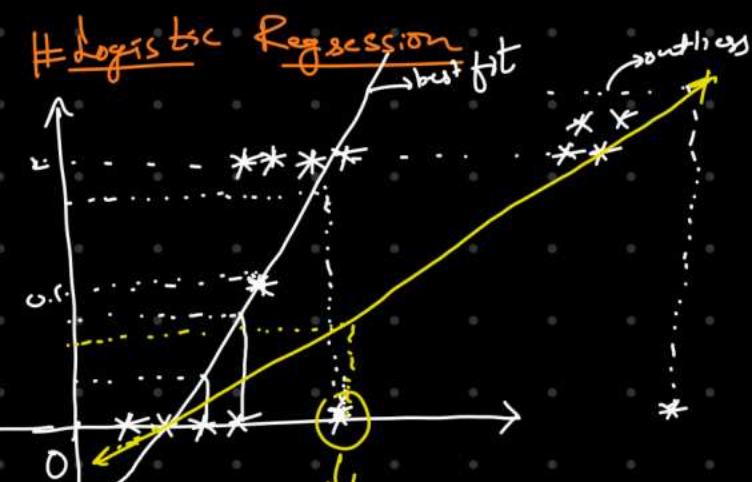
→ cancer → FP
→ if no cancer, predicted cancer, no problem
→ if cancer, predicted no cancer
↓
FN↓↓↓

⇒ So, +ve guy detected no cancer, big problem.
→ so In cancer detection we need to focus to reduce FN↓↓.

③ $F_\beta = (1+\beta^2) \dfrac{Precision \times Recall}{Precision + Recall}$

→ β = 1 → Both Preci, Recall imp.
→ β = 0.5 → FP has higher impact.
→ β = 2 → FN has higher impact.

# Roc & AUC

[0, 0.2, 0.4, 0.6, 0.8, 1]

| y | ŷ_Pred | ŷ(0) | ŷ(0.2) | ŷ(0.4) | ŷ(0.6) | ŷ(0.8) | |
|---|---|---|---|---|---|---|---|
| 1 | 0.8 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 0.96 | 1 | 1 | 1 | 1 | 1 | 0 |
| 1 | 0.4 | 1 | 1 | 0 | 0 | 0 | 0 |
| 1 | 0.3 | 1 | 1 | 0 | 0 | 0 | 0 |
| 0 | 0.2 | 1 | 0 | 0 | 0 | 1 | 0 |
| 1 | 0.7 | 1 | 1 | 1 | 1 | 0 | 0 |



$$TPR = \frac{TP}{TP+FN}$$
$$FPR = \frac{FP}{FP+TN}$$

① $\frac{1}{2+0} = 1$,  $FPR = \frac{2}{2+0} = 1$

② $\frac{4}{4+0} = 1$,  $FPR = \frac{1}{4+1} = 0.5$

③ $TR = \frac{2}{2+2} = 0.5$  $FPR = \frac{1}{1+1} = 0.5$

Ⓕ

⑤ $TR = \frac{0}{0+2} = 0$  $FR = \frac{1}{1+1} = 0.5$

# Logistic Regression
best fit   outliers



Earlier the point is considered as 1, but after adding outliers, it → consider as '0' with lin regession

⇒ logistic regression → Binary classification.

① It will greatly affected by outliers
② also some pts can go >1
③ So, we use logistic regression
→ we will squash the lines at '0' & '1'

# #Logistic Regression

(1) Geometric Intuition
(2) Mathematic Intuition



$$y_i \begin{cases} +ve = +1 \\ -ve = -1 \end{cases}$$

$$y = mx + c$$
$$y = \beta_0 + \beta_1 x$$
$$y = w^T x + b_0$$

$$\boxed{y = w^T x}$$

w is distance from plane



$$\frac{w^T x + b \to 0}{\|w\| \to 1}$$

$$\sum_{i=1}^{n} \boxed{w_i x_i}$$

**Case 1 [red color]**
$$y_i = 1 \quad w^T x_i > 0$$
$$\boxed{y_i \times w^T x} > 0$$
correct class.

**C2 [yellow]**
$$y_i = -1, \quad w^T x_i < 0$$
$$\boxed{y_i \times w^T x_i} > 0$$
Correct classifi

**C3**
$$y_i = -1, \quad w^T x_i > 0$$
$$\boxed{y_i \times w^T x} < 0$$
wrong classifi

→ **Optimization** $\max \boxed{\sum_{i=1}^{n} y_i \, w^T x_i}$

→ Effect of Outliers:



$$\boxed{\begin{array}{l} y = -1 \\ y = +1 \end{array}}$$

$$= 2+2+2+2+2$$
$$+2+2+2+2+2+2$$
$$= 500$$
$$= -480$$



$$-1-2-3-4-5-6$$
$$+1+2+3+4+n-1$$
$$+2$$
$$= \boxed{2} \text{//}$$

→ Sure that the first line is best, but 2nd one getting max for $y_i \times w^T x_i$ because of outlier.
→ So we introduce sigmod to avoid this

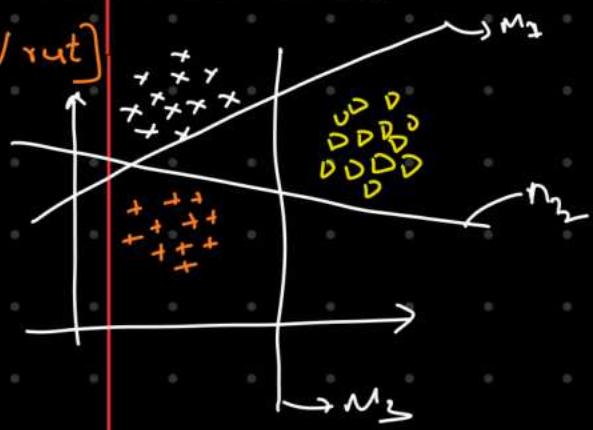→ after finding $\boxed{y_i \times w^T x_i}$, we pass this into another function → Sigmoid function.

→ $\sum_{i=1}^{n} f\left( y_i \times w^T x_i \right)$ → Sigmoid fn toanyform the values btw $\boxed{0,1}$. So removes

the effect of outliers. $\boxed{\text{Sigmoid fn} = \dfrac{1}{1+e^z}}$ //

# #multiclass logistic classification [One V rest]

| $f_1$ | $f_2$ | $f_3$ | OP | $O_1$ | $O_2$ | $O_3$ |
|---|---|---|---|---|---|---|
| $i_1$ | $i_2$ | $i_3$ | $O_1$ | +1 | -1 | -1 |
| $i_4$ | $i_5$ | $i_6$ | $O_2$ | -1 | +1 | -1 |
| $i_7$ | $i_8$ | $i_9$ | $O_3$ | -1 | -1 | +1 |
| $i_{10}$ | $i_{11}$ | $i_{12}$ | $O_1$ | +1 | -1 | -1 |
| $i_{13}$ | $i_{14}$ | $i_{15}$ | $O_2$ | -1 | +1 | -1 |

→ $m_1$ → $m_2$ → $m_3$
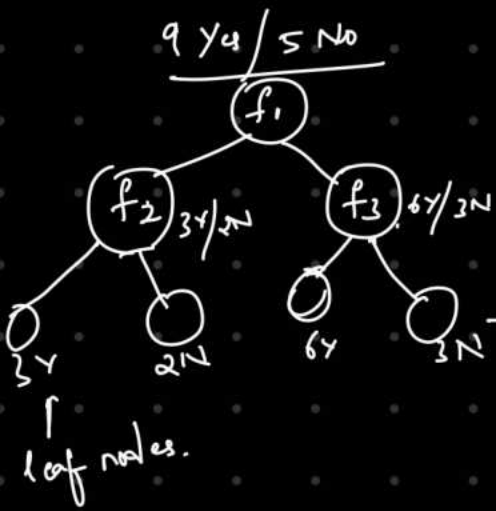


→ $M_1$
→ $M_2$
→ $M_3$

→ OVR
→ new test data
$$[0.20, 0.25, 0.55]$$
$$\downarrow$$
So, $O_3$

# #Decision Tree Entropy

Entropy := Measures the purity of split

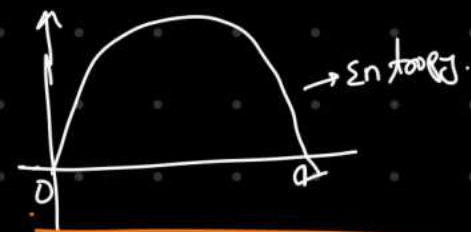$$H(s) = -P_{(+)} \log_2 (P_+) - P_{(-)} \log_2 (P_-)$$

$P_+/P_-$ : % of +ve class / % of -ve class.

$9\ Y_4 / 5\ N_0$



(tree: $f_1$ → $f_2$ $3Y/2N$, $f_3$ $6Y/3N$; $f_2$ → $3Y$, $2N$; $f_3$ → $6Y$, $3N$)

↑
leaf nodes.

→ $f_1\ f_2\ f_3$ → which feature to select at $1^{st}$ node.

→ at 1 node $3Y/3N_0$ → worst split.

→ at 2 node $3X/0N$ → best split.

→ $3Y/2N$

→ $S :-$ subset of Training Ex.

$= -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \left[\frac{3}{5} \log_2\left(\frac{2}{5}\right)\right]$

$= 0.78\%$ → entropy.

→ stops splitting at pure split.

→ Entropy → $0 \to 1$

$\boxed{0 \to best \qquad 1 \to worst}$


→ Entropy.

# information gain [Desicision Tree]

① Entropy



(tree: $f_1$ $9Y/5N$; $f_2$ $6Y/2N$; $f_3$ $3Y/3N$)

$H(s) = -P_+ \log_2(P_+) - P_- \log_2(P_-)$

$= -\frac{6}{8} \log\left(\frac{6}{8}\right) - \frac{2}{8} \log\left(\frac{2}{8}\right)$

$= 0.81$

② Information gain

$Gain(S, A) = H(S) - \sum\limits_{V \in val} \frac{|SV|}{|S|} H(Sv)$

$H(s) = 0.94$



$f_1$ $9Y/5N$ → $H(f_1) = H(s)$

$f_2$ $H(sv)$ $6Y/2N$ $H(f_2) = 0.81$

$f_3$ $2Y/3N$ $H(sv)$ $H(f_3) = 1$

$Gain(S, f_1) = H(s) - \frac{8}{14} H(f_2) - \frac{6}{14} H(f_3)$
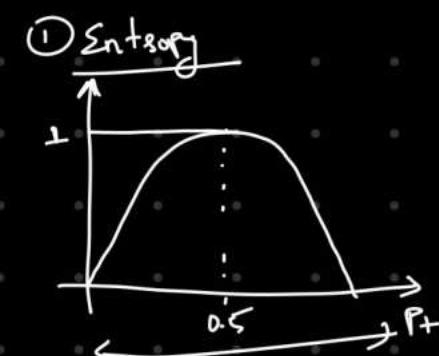
$= 0.94 - \frac{8}{14}(0.81) - \frac{6}{14}(1)$

$= \boxed{0.049}$



which split have the high information value that split will be selected.

A Gini Impurity DT

① Entropy



| $f_1$ | $f_2$ | $f_3$ | OP |
|---|---|---|---|
| $c_1$ | $D_1$ | | Y |
| $c_2$ | $P_2$ | | Y |
| | | | N |
| | | | Y |
| | | | Y |



(tree: $f_0$ $6Y/3N$; $c_1$ $3Y/3N$, $c_2$ $3Y/0N$ — leaf node)
$H(S) = 1 \times t$ $\to H(s) = 0\ Best$

→ ② GINI Impurity

$\boxed{GI = 1 - \sum\limits_{i=1}^{n} (P)^2}$

$GI = 1 - [P_+^2 + P_-^2]$

$= 1 - \left[\left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2\right]$

$= 1 - (0.25 + 0.25)$

$= \boxed{0.5}$


Entro
gini impurity

⊛ gini impurity ranges b/w $0 - 0.5$

⊛ Computationally efficient

⊛ takes less time than entropy

# Decision Tree split for Numerical var

| $f_1$ | o/p |
|-------|-----|
| $x_i =$ 2.3 | Y |
| 3.6 | Y |
| 4 | Y |
| 5.2 | N |
| 6.7 | N |
| 8.4 | N |
| 9.5 | Y |
| 10.2 | N |

① Sorting all the values

② $Th = 2.3$

```
              f₁  4Y/4N
        ≤2.3  /  \  ≥2.3
            /      \
         1Y/0N    3Y/4N
```
↳ Entropy
↳ Information gain

Now for $Th = 3.6$ → $Th = 4$ calculates for all $x_i$ values till 10.2

```
          4Y/4N
     ≤3.6  f₁  ≥3.6
         /    \
      2Y/0N   2Y/4N
```
↳ Entropy
↳ Info gain

Ⓧ After that finds the highest info gain

Ⓧ Then selects that

Ⓧ Disadvantage :- Time complexi

---

# K Nearest Neighbour



cat 1 → 3
cat 2 → ②

2 dimension

Algo :
① K = 5 Neighbour
② We calculate the distance of the neig $K = 5$
③ How many nearest neighbours ∈ cat 1 & cat 2.

① Euclidian Distance

$(x_1, y_1)$ —— $(x_2, y_2)$

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

② Manhattan distance

$(x_2, y_2)$  $(x_1, y_1)$

Dataset ⎧ 900 Yes
         ⎩ 100 No

Ⓧ K-Nearest neibour → biased.

↳ most come from this data pt.

Ⓧ Outliers :-



$k = 5$
→ outliers
→ impacted by outliers.

Ⓧ KNN Regressor

$k = 5$

→ finds the mean of nearest 5 = predict value.

→ finds best 'k' value by elbow method
→ classifier = KNNeighClas (n_neighbors = p) → loop in in range(1, 40) → 50

error

best k value

---

# Ensemble Techniques :-

```
        Bootstrap
        aggregation
       /          \
   Bagging      Boosting
① Random Forest  ① Ada boost
                 ② Gradient boost
                 ③ Xg boost
```

Bagging    row sampling with replacement

```
                           voting
m≤n      D'm →  [ ]  M₁ → 1
Test data  D''m → [ ] M₂ → 0
[D] →    D'''m → [ ] M₃ → 1    → ①
 n       ...      [ ] Mₙ → 1
Bootstrapping              Aggregate
```

# Random forest



$D' < D$
$n < m$
$d' < d$

## DT
① Low Bias    high variance.
② RF → Row sampling → f sampling
↳ so RF becomes expert so we will get low variance.
→ for regression → mean of all OP.

majority
2

---

# Ada boost { Boosting Technique }



Incorrect records

incorrect record

→ Boosting

1 features → 3 — stumps → DT with 2 depth

## Ada boost    step①    $w = \frac{1}{n} = \frac{1}{5}$

① $f_1$  $f_2$  $f_3$   O/p   sample wt.

| | f_1 f_2 f_3 | O/p | sample wt. |
|---|---|---|---|
| 1 | | | 1/5 |
| 2 | | | 1/5 |
| 3 | | | 1/5 |
| 4 | | | 1/5 |
| 5 | | | 1/5 |

incorrect record

↳ Entropy / gini impurity
↳ information gain

↑ if this has less entropy more info gain
(★) this is 1st weak learner.

## step 2



4 correct    incorr

Calculate total error.
$T\varepsilon = \frac{1}{5}$ → sum of all error weights.

## step 3 Performance of stump.

Performance of stump
$= \frac{1}{2} \log_e \left( \frac{1 - T\varepsilon}{T\varepsilon} \right)$

$= \frac{1}{2} \log \left( \frac{1 - \frac{1}{5}}{\frac{1}{5}} \right)$

$= 0.896$

(eg).

| | initial | updated wt | normalized | |
|---|---|---|---|---|
| 1 | 1/5 | 0.05 | 0.01 | 0 - 0.01 |
| 2 | 1/5 | 0.349 | 0.513 | 0.01 - 0.5 |
| 3 | 1/5 | 0.05 | 0.07 | 0.58 - 0.6 |
| 4 | 1/5 | 0.05 | 0.07 | 0.65 - 0.72 |
| 5 | 1/5 | 0.05 | 0.07 | |
| | | sum | | |

## Step 4  update weights  [↑incorr, ↓ right]

Performance = 0.896

New sample wt = $wt \times e^{\text{performance}}$  for incorr.

for correct = $wt \times e^{-\text{perfor}}$

$= \frac{1}{5} \times e^{+0.896} = 0.349$

$= \frac{1}{5} \times e^{-0.896} = 0.05$

—DIN ext step is to create new dataset based on normalized wt & buckets.
$f_1 \ f_2 \ f_3$ OP

we will select mostly only the records with high error by algo.
& create new dataset.

Again same steps :

| K_n |
| 1/n |
| 1/n |
| 1/n |
| 1/n |



Entropy → 2nd weak learner

↑ new dataset

normalized & bucket → update wts ← performance of stump

tell no. of DT selected.
finally error is less

finally
weak 1   weak 2
(on testing)

$\frac{\text{3}}{\text{y}} \ 2N$

finally

test data →



↓
1

majority.

---

# Kmeans clustering:


→ cluster 1
→ clu-2
↳ Kmeans

① Algorithm
② metrics
③ ⟶ euclidean
     ↳ manhattan
④ Elbow method
   Select the `k` value.

① K-value → centroids k=2
② Initialize the centroid randomly
③ Select the group & find the mean.

# Hierarichal clustering ___ —unsupervised.



Dendogram.

↳ we need to select the longest vertical line without any line cutting it.
↳ ✶ draw a line, no. of cuts = no. of clusters.
↳ so, 2 clusters.

---
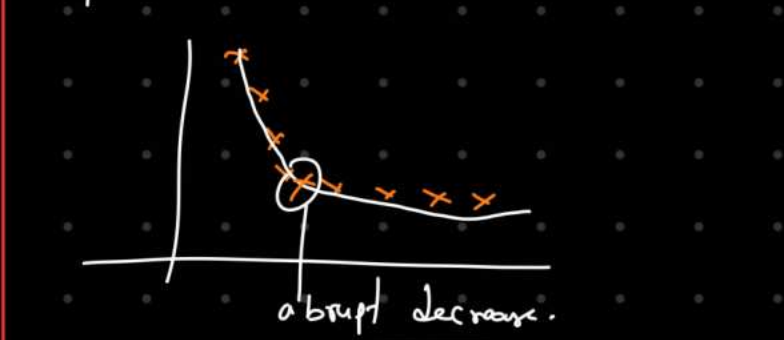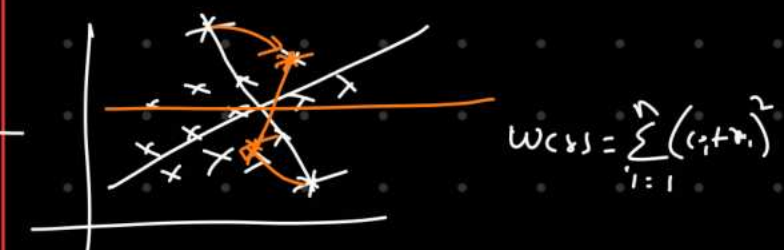


The left image depicts a more traditional clustering method that does not account for multi-dimensionality. Whereas the right image shows how DBSCAN can contort the data into different shapes and dimensions in order to find similar clusters.

---



find the distance
→ Then mean of both cluster
→ move centroid to there.

$$WCSS = \sum_{i=1}^{n}(c_i + n_i)^2$$

abrupt decrease.

# Density Based Special clustering (applications with Noise [DBSCAN] — unsupervised



core points.

Border pt

Noise
Eps    ε—radius.

→ Red core pts ⇒ min 4 pts inside a circle.

→ Yellow Border pts ⇒ still part of clusters, no 4 pts, But atleast 1 Border pt inside circle.

→ Blue Noise; no 4 pts, no Border pt

Advantages ↓
① great at creating clusters of high & low density.
② great at handling outliers.

⇒ Disadvantages ↓
① Doesn't work well when cluster of similar densities.
② struggles with very high dimensional data.

# Silhouette [← clustering]

④ $C_1$



→ finds the distance within the cluster pts.

→ $a_i = \dfrac{\sum\limits_{j \in C_i ; i \neq j} d(i,j)}{|C_i| - 1}$

② $C_1$   $C_2$



→ finds the dist. b/w one cluster pt to all pts of other cluster

→ $b_i = \min\limits_{k \neq i} \dfrac{\sum\limits_{j \in C_k} d(i,j)}{|C_k|}$

③ $a_i <<< b_i$ for correct clustering.
if $b_i > a_i$ → wrong clustering.

④ finally $S_i = \dfrac{b(i) - a(i)}{\max\{a(i), b(i)\}}$ , if $|C_i| > 1$

$S_i$ is b/w $-1$ to $1$ //

---

# Principal Component Analysis (PCA)

→ PCA is n't exactly afully ML algo, but instead an unsupervised learning algo.



PCA1

captures most info
→ lost less variance

loss of info

→ so we select PCA1.
→ 2D → LD.



PCA1

---

# Types of cross Validation

1000 records
70% train
30% test  } train test split random state = 0

randomly split the

85% are

if again change random state
accuracy changes.

---

① K fold CV



200 | 800 train

200 | test

⋮ 5 times

$\dfrac{1000}{5} = 200$

$K = 5$

* Disadvantages :

① If in test dataset, may contain all data of 1 type if binary classification.
② so leads to biased.

---

② Stratified cross Validation    $K = 5$



* → It will make sure that no. of instances of each class will contain in a good proportion in training and testing datasets

③ Time series CV : for stock market etc.

Day 1
Day 2
Day 3
⋮
Day 7

| Day1 | Day2 | D3 | D4 | D5 | | D6 |  O/P
|---|---|---|---|---|---|---|

$D_2$ | D3 | D4 | D5 | D6 | | $D_7$

D3 | D4 | D5 | D6 | D7 | | D8

---

# Bayes Theorem :

① Conditional Probability
② Independent events
③ Dependent events

$$P(A) = \frac{2}{5} \qquad \boxed{P(B/A) = \frac{1}{4}}$$

$$P(A \cap B) = \frac{2}{5} \times \frac{1}{4} = \frac{1}{10}$$

$$\boxed{P(B/A) = \frac{P(A \cap B)}{P(A)}} \qquad \boxed{P(B/A) = \frac{1}{4}}$$

$$= \frac{P(A) \times P(B)}{P(A)} = \frac{\frac{2}{5} \times \frac{1}{4}}{\frac{2}{5}} = \frac{\frac{1}{10}}{\frac{2}{5}} = \frac{1}{4} \,\|$$

---

# Bayes Theorem :-

$$P(A/B) = \frac{P(A \cap B)}{P(B)} \quad ; \quad P(B/A) = \frac{P(B \cap A)}{P(A)}$$

$$P(A \cap B) = P(A/B) \times P(B) \quad ; \quad P(B \cap A) = P(B/A) \times P(A)$$

But $P(A \cap B) = P(B \cap A)$

So $\Rightarrow$  $P(A/B) \times P(B) = P(B/A) \times P(A)$

$$\Rightarrow \boxed{P(A/B) = \frac{P(B/A) \times P(A)}{P(B)}}$$

---

# Naive Bayes Classifier

Data set :- $x = \{x_1, x_2, x_3, x_4 \ldots \ldots x_n\} \; \{y\}$

$$\begin{array}{cccc} f_1 & f_2 & f_3 & y \\ x_1 & x_2 & x_3 & y_1 \end{array}$$

$$P\left(y/x_1, x_2 \ldots x_n\right) = \frac{P(x_1/y) P(x_2/y) \ldots P(x_n/y) \times P(y)}{P(x_1) P(x_2) \ldots P(x_n)}$$

$$= \frac{P(y) \times \prod_{i=1}^{n} P(x_i/y)}{\boxed{P(x_1) P(x_2) \ldots P(x_n)}} \;\; \text{const}$$

$$P(y/x_1, x_2, x_3 \ldots x_n) \propto P(y) \prod_{i=1}^{n} P(x_i/y)$$

Yes = 0.3
No = 0.3

$$\boxed{y = \arg\max_y P(y) \prod_{i=1}^{n} P(x_i/y_i)}$$

---

outlook

| outlook | Yes | No | p(y) | p(N) |
|---------|-----|-----|------|------|
| sunny | 2 | 3 | 2/9 | 3/5 |
| over cast | 4 | 0 | 4/9 | 0/5 |
| Rainy | 3 | 2 | 3/9 | 2/5 |
| Total | 9 | 5 | 80% | 100% |

temp

| temp | Yes | No | p(y) | p(N) |
|------|-----|-----|------|------|
| HOT | 2 | 2 | 2/9 | 2/5 |
| MILD | 4 | 2 | 4/9 | 2/5 |
| cool | 3 | 1 | 3/9 | 1/5 |
| Total | 9 | 5 | 100% | 100% |

| | | p(y) & p(N) |
|------|-----|------|
| Yes | 9 | 9/14 |
| No | 5 | 5/14 |
| Total | 14 | 100% |

$$P(Y_u/today) = \frac{P(sunny/y_u) \times P(HOT/y_u) \times P(y_u)}{P(today)}$$

$$\to today \left(\overset{x_1}{sunny}, \overset{x_2}{Hot}\right)$$

$$P(y/today) \propto \frac{2}{9} \times \frac{2}{9} \times \frac{9}{14} = 0.031$$

$$P(N/today) \propto \frac{3}{5} \times \frac{2}{5} \times \frac{5}{14} = 0.0857$$

$$P(Yes) = \frac{0.031}{0.031 + 0.0857} \Big\} \text{normalizing.}$$

$$\approx 0.22\% \Big\} \text{Normalizing.}$$

$$P(N) = 0.72$$

---

# For classification of text

$$P(A/B) = \frac{P(B/A) \times P(A)}{P(B)}$$

$\to$ stop words
$\to$ stemming
$\to$ BOW
$\to$ TF-IDF
$\to$ NLP

---

| $f_1$ | $f_2$ | $f_3$ | $f_4$ | ... | Sentence $\begin{cases} Good \\ Bad \end{cases}$ |
|-------|-------|-------|-------|-----|----------|
| The | fowl | Debi | Beed | Off | |
| 4 | 1 | 7 | 0 | 1 | $P(y = yes/sentence)$ |
| 1 | 1 | 0 | 1 | 0 | sentence $[x_1, x_2, x_3 \ldots x_n]$ |
| 0 | 1 | 0 | 1 | 0 | |
| 0 | 1 | 1 | 0 | 1 | $P(y/x_1, x_2, x_3, x_4 \ldots x_n)$ |
| 0 | 0 | 0 | 1 | 0 | |

$$\propto P(y) \prod_{i=1}^{n} P(x_1/y) P(x_2/y) \ldots P(x_n/y)$$

$$= P(y = yes) \times P(The/y) P(2nd/y) \times P(Deb/y)$$

$$= \frac{4}{5} \times \frac{1}{2} \times \frac{2}{4} \times \frac{2}{2}$$

$$= 2/5 = 0.21$$

$$P(N/x_1, x_2, x_3, x_4) = P(x_1/N) \times P(x_2/N) P(x_3/N) \times P(N)$$

$$= \frac{3}{5} \times \frac{1}{2} \times \frac{2}{4} \times \approx 0.03$$

$$\text{Normalize} \quad \frac{0.01}{0.01 + 0.03} = \frac{0.01}{0.04} \approx 1/4 = 0.25$$

$$N = 0.75$$

# Support Vector Machine :

1. Support Vectors
2. Hyper plane
3. Marginal distance
4. Linear separable
5. Non-linear separable

hyper plane $h_1$

D1 marginal distance   $h_2$

$D_2$

$D_2 << D_1$

SVM kernels
$\hookrightarrow 2D \rightarrow$ higher dim

$y = \begin{bmatrix} -1 \\ 0 \end{bmatrix} [+0]$

$= -4 \ (-ve).$

$\rightarrow w^T x + b = 0.$

$(; slope = m = -1$

$y = w^T x + 0$

$= \begin{bmatrix} -1 \\ 0 \end{bmatrix} [-4 \ 0]$

$(-4, 0)$

$= 4$ //

All pts below plane +ve.

$w^T x_1 + b = -1$
$(-) \ w^T x_2 + b = 1$
$\overline{w^T (x_2 - x_1) = -2}$

$\frac{w^T}{||w||}(x_2 - x_1) = \frac{2}{||w||}$

$w^T x + b = -1$
$w^T x + b = 0.$
$w^T x + b = +1$

$\hat{w}, b^* \quad max\left(\frac{2}{||w||}\right) \quad st \quad y \begin{cases} +1 & w^T x + b \geq 1 \\ -1 & w^T x + b \leq -1 \end{cases}$

optimizers.

$\boxed{y_p * w^T x_i + b_i \geq 1}$

$\boxed{w^*, b^* = min \ \frac{||w||}{2} + (9 \sum_{q=1}^{n} \xi_i}$

Regularization    too many Errors    value of the error.

# SVM Kernels :

1. Soft margin
2. Hard margin

1. polynomial kernels
2. RBF Kernels
3. Sigmoid kernels.

Transformation from lower dim to higher dim by mathematical formulae

$y$

hyperplane.

$x \rightarrow x^2$

$y = f(x)$
$y = x^2$

$1D \rightarrow 2D$

1. Polynomial kernels:

50%

$x_1 \quad x_2 \quad y_0 \quad y = f(x)$

$f(x_1, x_2) = (x_1^T \cdot x_2 + 1)$

$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot [x_1 \ x_2]$

$= \begin{bmatrix} x_1^2 & x_1 x_2 \\ x_1 x_2 & x_2^2 \end{bmatrix}$

$x_1 \quad x_2 \quad y_0 \Rightarrow x_1 \quad x_2 \quad y_0 \quad x_1^2 \quad x_1 x_2 \quad x_2^2$

2d        5d

SVC   kernel = "poly"
         = "rbf"

hyper plane

# Gradient Boosting (actu - pred)

| Exp | Deg | Sal (act) | $\hat{y}$ (pred) | $R_1$ | $R_2$ | $R_3$ |
|-----|-----|-----|-----|-----|-----|-----|
| 2 | BE | 50k | 75 | -15 | -23 | |
| 3 | master | 70k | 75 | -5 | -3 | |
| 5 | master | 80k | 75 | 5 | 3 | |
| 6 | Degree | 100k | 75 | 25 | 2 | |

① step ① Base model 1 → 1 for doing avg.

$$\frac{50 + 70 + 80 + 100}{4} = 75$$

② compute residuals
(1) error, pseudo Residuals.

③ Construct DT $\{x_i, R_i\}$

$$= 75 + (-23)$$
$$= 52 \quad \text{close to } 50k \quad [\text{overfitting}]$$

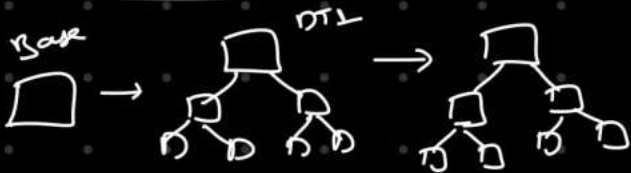so add learning rate. $\alpha \quad 0 \text{ to } 1$

$$= 75 + \alpha(-23)$$
$$= 75 + 0.1(-23)$$
$$= 75 - 2.3 \Rightarrow 72.7 \text{ //}$$

$$\Rightarrow f(x) = h_0(x) + \alpha_1 h_1(x) + \alpha_2 h_2(x) + h_3 \cdots$$
$$+ \alpha_n h_n(x)$$

$$\boxed{F(x) = \sum_{i=1}^{n} \alpha_i h_i(x)}$$

Base $\square \rightarrow$ DTL



all the above is eg purpose
below is actual pseudo algo.

## Pseudo algo

| Exp | Deg | sal |
|-----|-----|-----|
| 2 | BE | 50 |
| 3 | PHD | 70 |
| 4 | mtel | 60 |

① Initialize model with const value.

$$\boxed{F_0(x) = \arg\min\left(\sum_{i=1}^{n} L(y, r)\right)}$$

$$Loss = \sum_{i=1}^{n} \frac{1}{2}(y - \hat{y})^2 \Rightarrow \frac{1}{2}(50 - \hat{y})^2 + \frac{1}{2}(70 - \hat{y})^2 + \frac{1}{2}(60 - \hat{y})^2$$

$$\Rightarrow \frac{\partial}{\partial Loss} \text{ for min} \Rightarrow) \frac{2}{2}(50 - \hat{y}) + (70 - \hat{y}) - (60 - \hat{y})$$

$$\Rightarrow) \frac{180}{3} = 10 \Rightarrow \boxed{\hat{y} = 60}$$

→ initialize base model with this value.

---

② Iterate $M = 1$ to $M$, no. of trees

Compute pseudo residuals

$$r_{jm} = -\left[\frac{\partial L(y, f(x_i))}{\partial F(x_i)}\right] \text{ for } j = 1 \text{ to } m$$
$$F(x) = F_{m-1}(x)$$

$$Loss = \frac{1}{2}(y - \hat{y})^2$$

$$\frac{\partial L}{\partial \hat{y}} = (y - \hat{y})(-1)$$

$$\boxed{\frac{-\partial L}{\partial \hat{y}} = (y - \hat{y})} = r_{jm} = -\left[\frac{\partial L(y, f(x_i))}{\partial F(x_i)}\right] \quad y (60)$$

$$\frac{-\partial L}{\partial \hat{y}} = \frac{\partial L(y, \hat{y})}{\partial \hat{y}} = (y - \hat{y}) \quad \text{we need to compute This}$$

| y | $\hat{y}$ |
|---|---|
| 50 | 10 |
| 70 | 10 |
| 60 | 10 |

$$r_{jm} \Rightarrow \Rightarrow \text{model no.} \begin{cases} r_{11} = -10 \\ r_{21} = 10 \\ r_{31} = 0 \end{cases}$$

$(y - \hat{y})$ $\quad x_i \rightarrow$ indep $\quad r_{im} \leftarrow$ dep

$$\Rightarrow \text{Fit a Base learner } h_m(x) \quad \{(x_i, r_{im})\}$$
$$ip \quad op$$

③ $\boxed{r_m = \arg\min \sum_{i=1}^{n} L(y_i, f_{m-1}(x_i) + r)}$

$$= \sum_{j=1}^{n} \frac{1}{2}(y_i - (60 + \hat{y}))^2$$

$$\rightarrow \frac{\partial}{\partial L} \rightarrow \text{find } \hat{y} \rightarrow \text{repeat step } ① \& ②$$

④ Update model $\boxed{F_m(x) = F_{m-1}(x) + r_m h_m(x)}$



$$60 + (0.1)(-10)$$
$$= 60 - 1$$
$$= 59 \text{ //}$$

Learning rate
$0 - 1$

2nd record $60 + 0.1 \times 10$
$$= 61$$

# XGBoost

| salary | Credit | Approved | Res |
|--------|--------|----------|------|
| <:50k | B | 0 | -0.5 |
| <:50k | G | 1 | 0.5 |
| <:50k | G | 1 | 0.5 |
| >50k | B | 0 | -0.5 |
| >50k | G | 1 | 0.5 |
| >50k | N | 1 | 0.5 |
| <:50k | N | 0 | -0.5 |

→ output is ± oo, so $\frac{1+0}{2} = 0.5$ → Base model probability.

Approved → Sigmoid prob.

① Construct tree with Root

② Calculate Similarity weight $= \frac{\sum (Residuals)^2}{\sum (pr(1-prob)+\lambda)}$

③ Calculate Gain

$$[-0.5, 0.5, 0.5, -0.5, 0.5, 0.5, -0.5]$$

Salary   sw: 1.43

≤50k                > 50k

sw: 0.               sw: 0.33

$[-0.5, 0.5, 0.5, -0.5]$        $[-0.5, 0.5, 0.5]$

always create binary tree

Similarity wt:
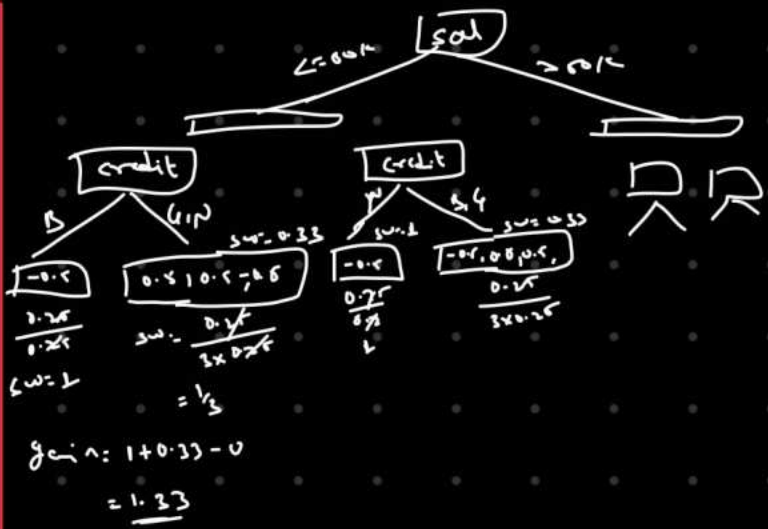$= \frac{(-0.5+0.5+0.5-0.5)-10}{2pr(1+pr)+\lambda}$

$= 0$.

similarity wt:
$= \frac{(0.5)^2}{0.5(1-0.5)+0.5(1-0.5)+0.5(1+0.5)}$

$= \frac{0.25}{3 \times 0.25 \times 0.25} = \frac{1}{3}$

→ sw for Root:

$sw = \frac{(0.5)^2}{7 \times 0.25 \times 0.25} = \frac{1}{7} \approx 0.143$

⇒ Gain = sw left + sw right − sw of Root

$= 0 + 0.33 - 0.14$

$= 0.21$

---

[sal]

<=50k            >50k

credit               credit

B    (in)            N    G

sw:-0.5              sw:-0.33

$[-0.5]$  $[0.5|0.5;+5]$   $[-0.5]$  $[-0.5, 0.5, 0.5,$

$\frac{0.25}{1.25}$  sw:- $\frac{0.y}{3\times0.25}$   $\frac{0.25}{0.75}$  $\frac{0.25}{3\times0.25}$

sw=1                 = 1/3                1

gain: 1+0.33-0

= 1.33

* for post pruning → cover value → pr(1-prob)
= 0.25

if gain < 0.25, cut that DT. then..

* $\log(odds) = \log\left(\frac{p}{1-p}\right)$

$= \log\left(\frac{0.25}{0.25}\right) = \log 1 = 0.y$

→ output is '0' for base model.

$\sigma\left(0 + 0.1\left(\frac{1}{3}\right)\right)$
leaning rate ↑  sw.

$\sigma(0.1) \Rightarrow \frac{1}{1+e^{-0.1}} = 0.6$.

↑
new probability compute for all row

| sal | Cred | App | Res | New P | new R |
|-----|------|-----|-----|-------|-------|
|     |      |     |     | 0.6   |       |

App → new pr

Then for new Residual compute the same process of DT → split → sw → gain → new R.

# XGBoost for Regressor:

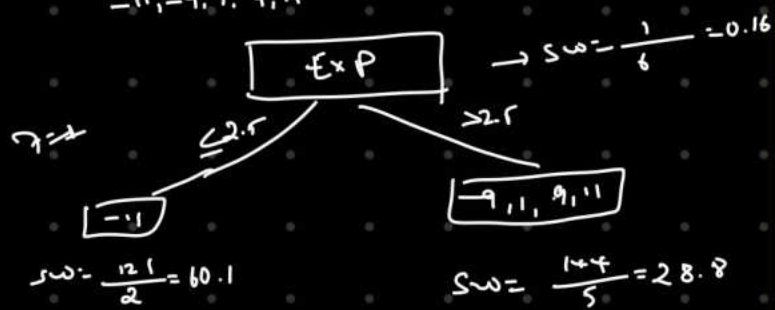| Exp | Gap | Sal | Res1 | O/P |
|-----|-----|-----|------|-----|
| 0.5 | Y | 40k | -11 | |
| 2.5 | Y | 42k | -9 | |
| 3.0 | N | 52k | 1 | |
| 3.5 | N | 60k | 9 | |
| 4 | Y | 62k | 11 | |

→ Res1 ⇒ Sal − Base model ∴ Avg(sal) $\frac{40+42+...}{5} = 51k$

→ Base model = 51k.

* Similarity wt = $\frac{\sum(Residuals)^2}{no. of residuals}$

51k [Base Model] +2

$-11, -7, 1, 9, 11$

[Exp] → $sw = \frac{1}{6} = 0.16$

$\gamma = ?$    $\leq 2.5$     $> 2.5$

[-11]       [-9, 1, 9, 11]

$sw = \frac{121}{2} = 60.1$     $sw = \frac{1+4}{5} = 28.8$

Gain = $60.1 + 28.8 - 0.16$

$= 88.9 - 0.16 \approx \underline{82.84}$

↳ coz calculate gain is only to confirm which split we should consider.

↳ after confirming → calculate output → oₚ

gₐᵣ         $\times 2$    No
$\gamma_1$ ↗↘        $\nearrow$ ↘    $1, 9$

X not go this split    $-9, 11$

only 2 value     $op = \frac{2}{2} = 1$    $op = 5$

s → root node.
   no split.

$op = \underline{-11}$

↳ take values and calculate residuals based
on this tree.

now.

$oupd = 51 + (0.6)^\alpha \left[ -11 \right]$ for 1st x^ord.

↳ versents based on which split they belong.

$= 51 - 5.5$

$= 45.5 \quad \rightarrow 1stmn-$

Sᵣₚ gₐₗ | sₐₗ | ₚₐ ₂ | New oₚ | Rₙ₂
      Y val — new op.

↓
again repeat
↓ DT 2.

⟹ Base Model $+ \alpha_1 (T_1) + \alpha_2 (T_2) + \cdots + \alpha_n (T_n)$