

1. AdaBoost (Adaptive Boosting)

- ▷ In-depth Math intuition
- ⇒ Dataset

Salary	Credit	Approval
$\leq 50k$	B	No
$\leq 50k$	G	Yes
$\leq 50k$	G	Yes
$> 50k$	B	No
$> 50k$	G	Yes
$> 50k$	N	Yes
$\leq 50k$	N	No

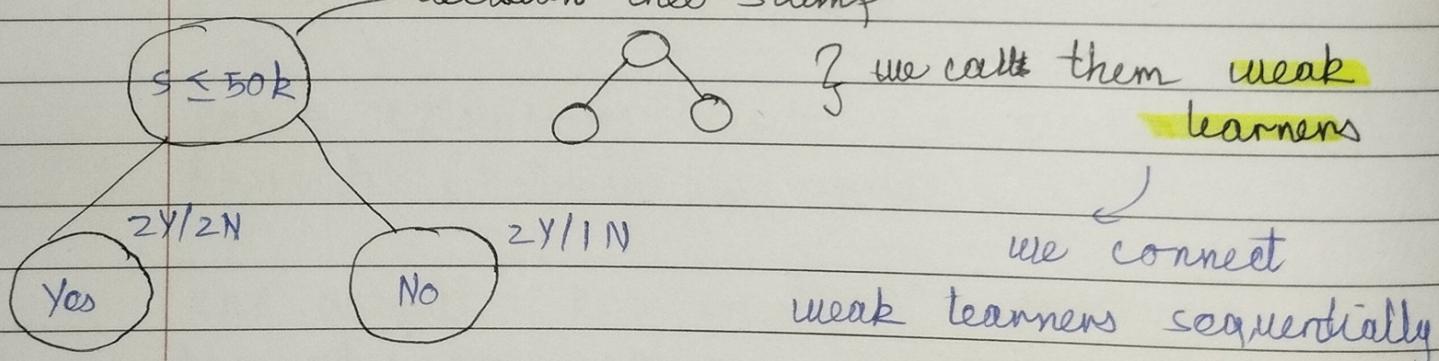
⇒ Step 1: We have to create a decision ~~stump~~ stump and we select the best stump

? What is Decision tree stump

↪ A decision-tree stump is a one-level decision

tree that makes a prediction using a single split on one feature, typically via a simple rule like "if feature \leq threshold then class A else class B"

decision tree stump



- Another decision-tree

↳ Which decision-tree stump we would select?

Best Decision ←
tree stump

based on
purity split

$C = G_1$

$3Y/1N$

$1X/3N$

↳ This depth-1
tree is $C(M_1)$

+ To select the best decision-tree stump we use
1. Entropy

$$\hookrightarrow \text{formula : } H(S) = - \sum_{i=1}^m p_i \log_2(p_i)$$

\hookrightarrow for k classes explicitly

$$\Rightarrow H(S) = - (p_1 \log_2 p_1)$$

$$\Rightarrow H(S) = - (p_1 \log_2 p_1 + p_2 \log_2 p_2 + \dots + p_k \log_2 p_k)$$

2. Gini Index (Impurity)

$$\cdot \quad \Omega_1 = 1 - \sum_{i=1}^m (p_i)^2$$

$$\hookrightarrow \text{expanded version : } \Omega_1 = 1 - (p_1^2 + p_2^2 + \dots + p_k^2)$$

3. Information Gain (this tells us how much uncertainty is reduced if we split the data using a particular feature)

$$\hookrightarrow \text{formula : } \text{Gain}(S|S_i) = H(S) - \frac{\sum |S_v| \times H(S_v)}{|S|}$$

$\therefore H(S)$: Entropy of root node

$|S_v|$: size of each subset v

$|S|$: total size of the dataset

$H(S_v)$: Entropy of each subset

⇒ Step 2:

- We assign a sample weight

Now how do we assign a same sample weight to all instances in a feature?

- $s_w \rightarrow s_{ot} := \text{Divide all by total for 7 nodes}$
 $1/7 \Rightarrow \text{for all seven instances (nodes)}$

(Sample Weight)

Dataset

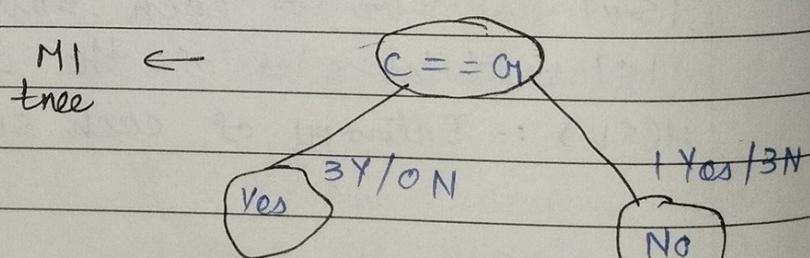
⇒ Salary	Credit	Approval	s_w
$\leq 50k$	B	No	1/7
$\leq 50k$	C ₁	Yes	1/7
$\leq 50k$	C ₁	Yes	1/7
$> 50k$	B	No	1/7
$> 50k$	C ₁	Yes	1/7
$> 50k$	N	Yes	1/7
$\leq 50k$	N	No	1/7

⇒ Step 3: Sum of Total Errors & Performance of stump

Dataset

Best Decision - tree stump

Salary	Credit	Approval	s_w
≤ 50	B	No	1/7
≤ 50	C ₁	Yes	1/7
≤ 50	C ₁	Yes	1/7
> 50	B	No	1/7
> 50	C ₁	Yes	1/7
> 50	N	Yes	1/7
≤ 50	N	No	1/7



⇒ credit is N ⇒ so it goes on the leaf in that its yes or no!

→ We know its Yes but depth-1 tree always uses the maximum for prediction. Since we have more No(3) than Yes(1) in No leaf. It will wrong predict it as No

What a stump does ?

- ↳ A decision stump is a depth-1 tree : it chooses one feature and a threshold (or a category test) to split the training set into two leaves, then assigns each leaf the maximum (majority) label of the samples that fall into that leaf. This is why it's called as "Weak learner"

① **Very Important** ***

↳ Decision Tree Stump Prediction Rule

Example Breakdown :-

Scenario :- A new data point falls into a leaf node with

1 "Yes"

3 "No"

Prediction :- The stump predicts "No" because

$$\text{No}(3) > \text{Yes}(1)$$

Majority Wins

* How it Works :-

- 1) Step 1 :- Stump splits data based on one feature - Feature : Age > 30 ?
 - 1 - Yes branch $\rightarrow [1 \text{ Yes}, 3 \text{ No}] \rightarrow$ Predicts "No"
 - 1 - No branch $\rightarrow [5 \text{ Yes}, 2 \text{ No}] \rightarrow$ Predicts "Yes"

→ ~~Step~~ Key Point

- ↳ Even if there's 1 "Yes" in the group, the stump ignores it and goes with the maximum count (3 "No") \Rightarrow This is called as **Plurality Vote**

When it predicts wrong (depth-1 tree) we calculate sum of total errors

To calculate this we need to find out how many records are predicted wrong

If predicted wrong is 1

then sum of TE = $\frac{1}{7}$

(

The formula :- used to calculate sum of Total errors

Total Error = \sum (weights of misclassified samples)

In simple words :- Add up the weights of all wrong predictions

If two instances are predicted wrong and its weight is 0.30 and 0.25

$$\text{Total Error} = 0.30 + 0.25 = 0.55$$

① Performance of stump

To calculate (Performance of stump), we use this formula, which AdaBoost uses internally

$$\text{Performance of stump} = \frac{1}{2} \ln \left[\frac{1 - TE}{TE} \right]$$

→ Calculation of Performance of stump using Total Error

$$\begin{aligned}
 \text{Performance of stump} &= \frac{1}{2} \ln \left[\frac{1 - TE}{TE} \right] \\
 &= \frac{1}{2} \ln \left[\frac{1 - 1/7}{1/7} \right] \\
 &= \frac{1}{2} \ln [6] \approx 0.896
 \end{aligned}$$

$\therefore \ln$ = Natural logarithm

$$\ln(1) = 0$$

$$\ln(2) = 0.693$$

$$\ln(3) = 1.099$$

$$\ln(10) = 2.303$$

Different input → Different output

? When we use Performance of stump?

$$\text{Performance of stump} = 0.896$$

Boosting Algorithm (AdaBoost) :- Function

$$f = \alpha_1(M_1) + \alpha_2(M_2) + \dots + \alpha_n(M_n)$$



$$\alpha_1 = 0.896 \Rightarrow \text{Weights}$$

→ Step 4:- Update the weights

→ Step 4 :- Update the weights for correctly and incorrectly classified points

Dataset

Salary	Credit	Approval	SN
$\leq 50k$	B	No	1/7
$\leq 50k$	G	Yes	1/7
$\leq 50k$	G	Yes	1/7
$> 50k$	B	No	1/7
$> 50k$	G	Yes	1/7
$> 50k$	N	Yes	1/7
$\leq 50k$	N	No	1/7

→ We know that this one is wrong predicted :-
(So we update weights)

→ Why update weights ?

The Big Idea :

Force the next stump to focus on samples the previous stump got Wrong

Simple Example :-

You have 5 students taking a test :-

Test 1 Results :

- Students A, B, C → Passed
- Students D, E → Failed

What do you do for Test 2 ?

- Don't give equal attention to everyone
- Give more attention to D and E (they need help!)
- Give less attention to A, B, C (they're doing fine)

Same in AdaBoost :

- Stump 1 Results :
 - Samples 1, 2, 3 → Correct
 - Samples 4, 5 → Wrong

For Stump 2 :

- Increase weight :- for Wrong Prediction
(sample 4, 5) → "Focus here"
- Decrease weight :- for Correct Prediction
(sample 1, 2, 3) → "These are easy"

* Why ?

↳ so each new stump fixes the mistakes of the previous stump !

Weight update formula :-

↳ for correctly classified points
= weight $\times e^{-\text{performance of stump}}$
or

$$\text{New weight} = \text{old weight} \times e^{-c}$$

↳ for wrong predictions

$$\text{New weight} = \text{old weight} \times e^{\text{performance of stump}}$$

or

$$\text{New weight} = \text{old weight} \times e^c$$

Where :

$$e = 2.718 \text{ (Euler's number)}$$

$$\text{alpha}(c) = \text{performance of stump}$$

Eg for correctly classified points

- new weight = old weight $\times e^{-\alpha}$
 $= \frac{1}{7} \times e^{-(0.896)}$
 $= 0.058$

Eg for misclassified points

- new weight = old weight $\times e^{\alpha}$
 $= \frac{1}{7} \times e^{(0.896)}$
 $= 0.349$

Dataset (this is how it looks after weight updation)

Updated weights (Dataset)

salary	credit	Approval	sw	Updated Weights
<= 50k	B	No	1/7 ↓	0.058
<= 50k	a	Yes	1/7 ↓	0.058
<= 50k	o	Yes	1/7 ↓	0.058
> 50k	B	No	1/7 ↓	0.058
> 50k	o	Yes	1/7 ↓	0.058
> 50k	N	Yes	1/7 ↑	0.349
<= 50k	N	No	1/7 ↓	0.058

this is our
misclassified point

Now how we gonna tell the next model - that
this is wrong prediction point and focus
more over here.

Answer :- B in Assignment

⇒ step 5 :- Normalize weights computation and assigning bins

weights

salary	credit	Approval	SW	Updated weights	Normalized,
$\leq 50k$	B	No	$1/7 \downarrow$	0.058	0.083
$\geq 50k$	G	Yes	$1/7 \downarrow$	0.058	0.083
$\leq 50k$	G	Yes	$1/7 \downarrow$	0.058	0.083
$> 50k$	B	No	$1/7 \downarrow$	0.058	0.083
$> 50k$	G	Yes	$1/7 \downarrow$	0.058	0.083
$> 50k$	N	Yes	$1/7 \uparrow$	0.349	0.500
$\leq 50k$	N	No	$1/7 \downarrow$	0.058	0.083
Total			1	0.697	1

- * Normalize weights means :- making sure all weights add up to 1 (or 100%).
- (After after updating weights, they might sum to something weird like 2.5 or 0.8
 \therefore We want : sum of all weights = 1)

* The Formula

$$\text{Normalized Weight} = \text{Weight} / \text{sum of all weights}$$

- * Bin Assignment ?
- Bin Assignment is how we select samples for training the next stump based on their weights
- It forces the next stump to train more on high-weight (misclassified) samples.

Bin Assignment = Groups

Salaney	credit	Approval	SW	Update weights	Normalize	Bin Assig
$\leq 50k$	B	No	1/7	0.058	0.083	0.000 - 0.083
$\leq 50k$	G	Yes	1/7	0.058	0.083	0.083 - 0.166
$\leq 50k$	G	Yes	1/7	0.058	0.083	0.166 - 0.249
$> 50k$	B	No	1/7	0.058	0.083	0.249 - 0.372
$> 50k$	G	Yes	1/7	0.058	0.083	0.372 - 0.455
$> 50k$	N	Yes	1/7	0.349	0.500	0.455 - 0.955
$\leq 50k$	N	No	1/7	0.058	0.083	0.955 - 1.038
Total			1	0.697	1	1

misclassified points

- ↳ Now How the Algorithm automatically picks misclassified points.
- ⇒ Lets break this down super simple

1. The computer generates random numbers between 0 & 1
- random_number = random.uniform(0, 1)
like 0.234, 0.789, etc.
2. How does it know which sample to pick automatically?
→ step-by-step process (super simple):
setup - one time only:
- ↳ the algorithm creates cumulative ranges (bin) based on weights

Your normalized weights

weights = [0.083, 0.083, 0.083, 0.083, 0.083, 0.500, 0.083]

Algorithm automatically calculates cumulative sum
 $\text{cumulative} = [0, 0.083, 0.166, 0.249, 0.372,$
size $0.455, 0.955, 1.0]$

8.37.

8.37.

Now it has ranges :-

8.37.

8.37.

8.37.

50%.

8.37.

8.37.

Sample 1 :- $0.000 - 0.083$ Sample 2 :- $0.083 - 0.166$ Sample 3 :- $0.166 - 0.249$ Sample 4 :- $0.249 - 0.372$ Sample 5 :- $0.372 - 0.455$ Sample 6 :- $0.455 - 0.955 \leftarrow \text{HUGE!}$ Sample 7 :- $0.955 - 1.000$

The Loop (Automatic Selection) :

lets say we need 100 samples for training :-
selected_samples = []

for i in range(100): # Loop runs 100 times

step 1 :- Generate random number (AUTOMATICALLY!)
rand = random.uniform(0, 1)

step 2 :- Check which bin this falls into (AUTOMATICALLY)
if $0.000 \leq \text{rand} < 0.083$:

selected_samples.append(1) # Pick sample 1

elif $0.083 \leq \text{rand} < 0.166$:

selected_samples.append(2) # Pick sample 2

elif $0.166 \leq \text{rand} < 0.249$:

selected_samples.append(3) # Pick sample 3

elif $0.249 \leq \text{rand} < 0.372$:

selected_samples.append(4) # Pick sample 4

elif $0.372 \leq \text{rand} < 0.455$:

selected_samples.append(5) # Pick sample 5

elif $0.455 \leq \text{rand} < 0.955$: # ← HUGE RANGE!

selected_samples.append(6) # Pick sample 6

elif $0.955 \leq \text{rand} \leq 1.000$:

selected_samples.append(7) # Pick sample 7

⇒ Step 5 :- Select data points to the Next Model (Stump)

Real Example (100 iterations):

— Let me show you what happens:

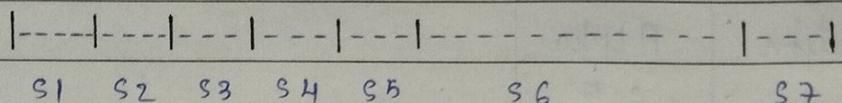
Iterations	Random num	Falls in Range	Sample picked
1	0.234	0.166 - 0.249	3
2	0.678	0.455 - 0.955	6
3	0.012	0.000 - 0.083	1
4	0.521	0.455 - 0.955	6
5	0.789	0.455 - 0.955	6
6	0.145	0.083 - 0.166	2
7	0.601	0.455 - 0.955	6
:	:	:	:
100	0.822	0.455 - 0.955	6

After 100 picks :

- Sample 6 picked - 50 times
- Sample 1 picked - 8 times
- Sample 2 picked - 8 times

- Why Sample G gets Picked more?
- Because its bin is BIGGER!

Imagine throwing a dart at this board:



- ↳ If you throw randomly, where will you hit MOST?
- Sample G's area! It's HUGE

* Simple Analogy:

Lottery Machine:

- Put 1000 balls in a machine
- 83 balls labeled "Sample 1"
- 83 balls labeled "Sample 2"
- 500 balls labeled "Sample G"
- 83 balls labeled "Sample F"
- Spin and pick randomly → You'll get "Sample G" most often!
- same logic! The computer does this automatically

⇒ Step 7 :- Final Prediction in AdaBoost

Simple Example :-

4 stumps vote on a new sample :-

Stump	Prediction	Alpha
1	Yes	0.5
2	Yes	0.3
3	No	0.4
4	No	0.2

① Direct Comparison :-

$$\text{Yes score} = 0.5 + 0.3 = 0.8 \Rightarrow \text{Winner}$$

$$\text{No score} = 0.4 + 0.2 = 0.6$$

Predict : Yes !

② Sign Function :-

Formula :

$$\rightarrow \text{Final Score} = (\alpha_1 \times \text{pred}_1) + (\alpha_2 \times \text{pred}_2) + \dots + (\alpha_n \times \text{pred}_n)$$

$$\begin{aligned}\text{Final} &= \text{sign} (0.5(\text{Yes}) + 0.3(\text{Yes}) + 0.4(\text{No}) + \\ &\quad 0.2(\text{No})) \\ &= 0.5(+1) + 0.3(+1) + 0.4(-1) + 0.2(-1) \\ &= 0.5 + 0.3 - 0.4 - 0.2 \\ &= 0.2 (\text{positive})\end{aligned}$$

Sign Function

If Final score > 0 → Predict "Yes"

If Final score < 0 → Predict "No"

⇒ final score (0.2) > 0 → Yes