

		Predicted		Total
		Spam P.	legitimate ⁿ	
Actual classes	Positive Spam	TP 40	FN 10	50
	Negative Legitimate	FP 5	TN 45	50
Total		45	55	100

TP (True positive) = 40
Correctly Spam as spam

FN (False Negative) = 10
wrongly predicted as legitimate

FP (False positive) = 5
wrongly predicted legitimate as spam

TN (True Negative) = 45
Correctly legitimate as legitimate

1. Accuracy:

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

$$= \frac{40 + 45}{40 + 45 + 5 + 10} = \frac{85}{100}$$

85% of all predictions were correct = 85% ✓

2. Recall: True Positive rate: Sensitivity

$$\frac{TP}{TP + FN}$$

$$= \frac{40}{40 + 10} = \frac{40}{50} = .8$$

$$= 80\%$$

Model catches 80% of actual spam

3. Specificity

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

$$= \frac{45}{45 + 5} = \frac{45}{50} = 0.9$$

$$= 90\%$$

Model catches 90% of legitimate mails

4. Precision

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

$$= \frac{40}{40 + 5} = \frac{40}{50} = 0.88$$

89%

When model says 'SPAM' it's right 89% of time

5. F1 score

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

Geometric mean

$$\leftarrow \frac{x+y}{2}$$

$$\downarrow$$

$$\textcircled{x}$$

Harmonic mean

$$\Rightarrow \frac{2(xy)}{x+y} \checkmark$$

$$\textcircled{y}$$

$$\frac{x+y}{2} \Rightarrow$$

$$\frac{2xy}{x+y} \downarrow$$

$$x = 10$$

$$y = 2$$

$$\frac{10+2}{2}$$

$$= \frac{12}{2} = 6$$

$$\frac{2(10 \times 2)}{10+2}$$

$$\frac{2(20)}{12} = \frac{40}{12}$$

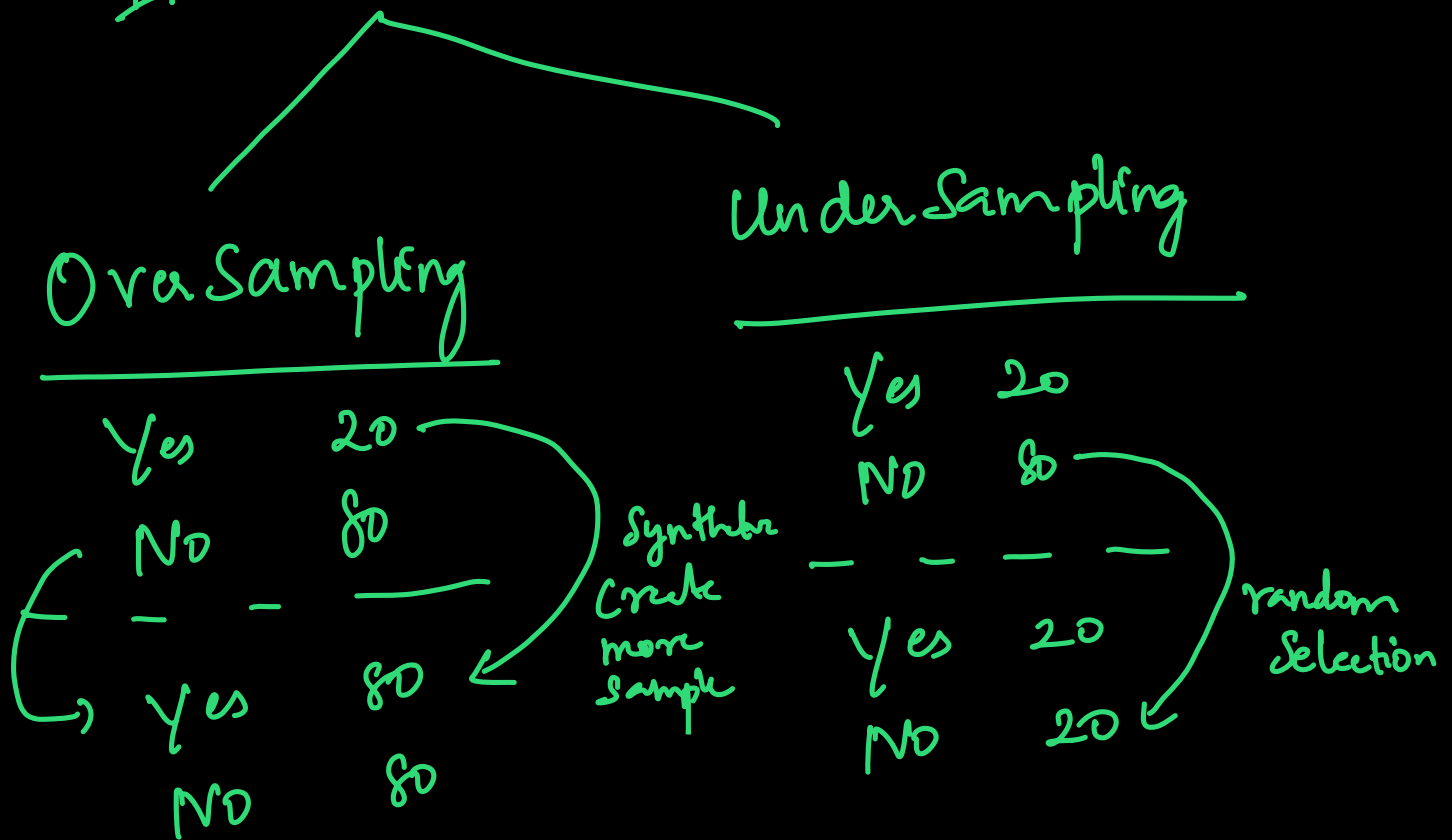
$$= \underline{\underline{3.3}}$$

$$= \frac{2(PR)}{P+R} = \frac{2 \times 0.88 \times 0.8}{0.88 + 0.8}$$

$$= 84.2\%$$

Harmonic mean balancing Precision and recall!

Imbalanced Data



6. Roc Curve

Receiver Operating Characteristics

Curve which shows how the model performs at different threshold settings

email 1 : 0.95 probability of spam
↳ High confidence spam

email 2 : 0.60 Medium confidence spam

email 3 : 0.20 probability of legitimate

Set a threshold $\begin{cases} \text{spam} \\ \text{legitimate} \end{cases}$

Email	Actual class	Spam Probability (fmodel)
A	Spam	0.95
B	Spam	0.85
C	legitimate	0.70
D	spam	0.65
E	legitimate	0.55
F	Spam	0.45
G	legitimate	0.35
H	Spam	0.25
I	legitimate	0.15
J	legitimate	0.05

Handwritten notes:

- A vertical green line is drawn at 0.55, separating the 'Spam' region (left) from the 'legitimate' region (right).
- Red arrows point from the 'Actual class' column to the predicted class based on the threshold: A (Spam), B (Spam), C (legitimate), D (Spam), E (legitimate), F (Spam), G (legitimate), H (Spam), I (legitimate), J (legitimate).
- A green bracket on the right side of the table groups rows A through J, with a 'Spam' label next to it.

$$\text{Threshold} = 0.5$$

$$\text{Predicted Spam} = A \text{ B C D E}$$

$$TP = 3 \quad FP = 2 \quad FN = 2 \quad TN = 3$$

$$TPR = \text{Recall} = \frac{TP}{TP + FN}$$

$$= \frac{3}{3 + 2}$$

$$= \frac{3}{5}$$

$$= \underline{0.60}$$

$$\text{Specificity} =$$

$$\frac{TN}{TN + FP} = \frac{3}{5}$$

$$FPR = 1 - \text{Specificity} = \underline{0.6}$$

$$= 1 - 0.6 \Rightarrow \underline{0.4}$$

Threshold (0.7)

Predicted Spam = A, B, C

$$TP = 2 \quad FP = 1 \quad FN = 3 \quad TN = 4$$

$$\text{Recall} = \frac{2}{5} = 0.4$$

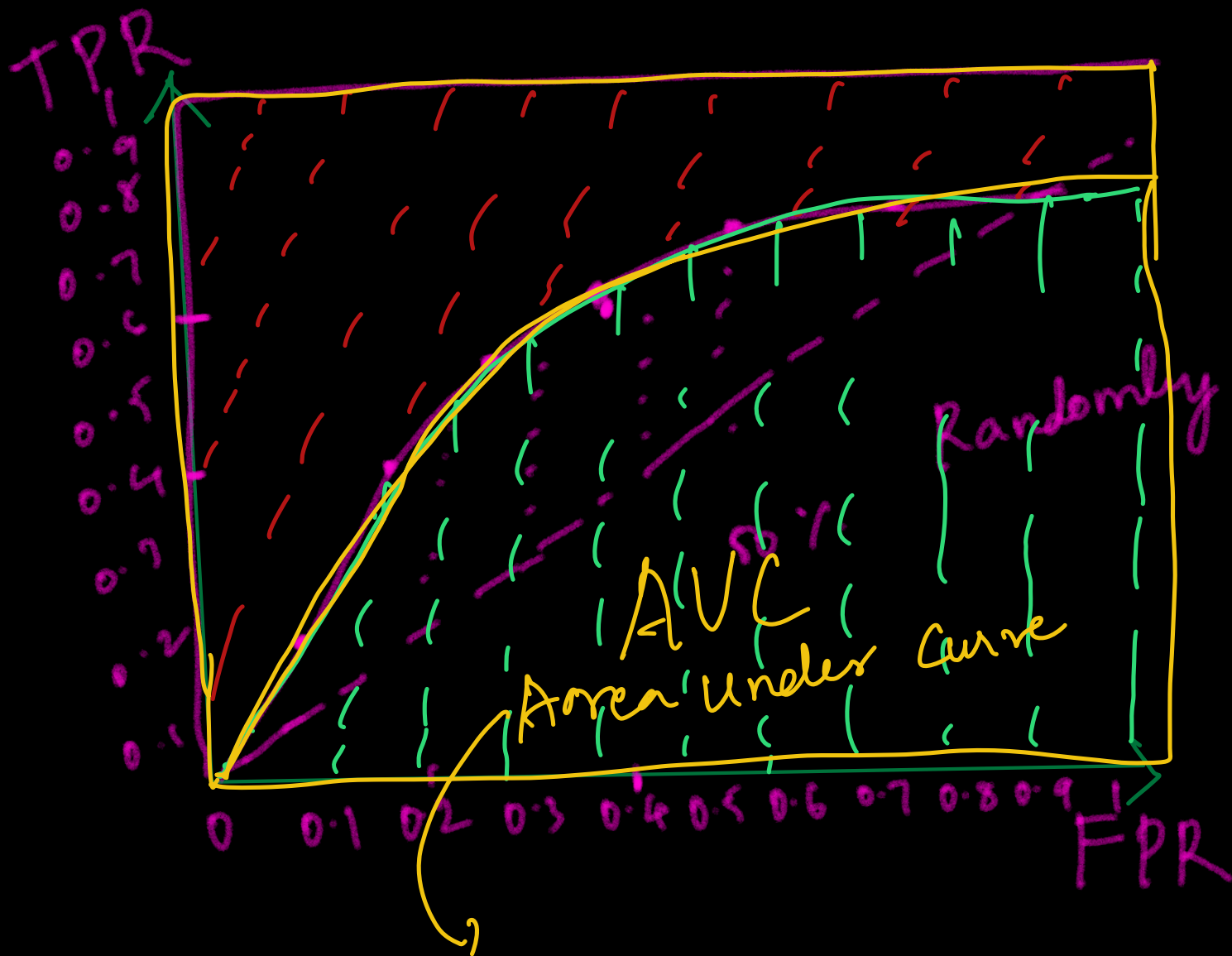
$$\text{Specificity} = \frac{4}{4+1} = \frac{4}{5} = 0.8$$

$$FPR = 1 - \text{Specificity} = 1 - 0.8 = 0.2$$

Threshold 0.3

$$\text{Recall} = 0.8$$

$$FPR = 0.6$$



$AUC = 1 \rightarrow$ Perfect model

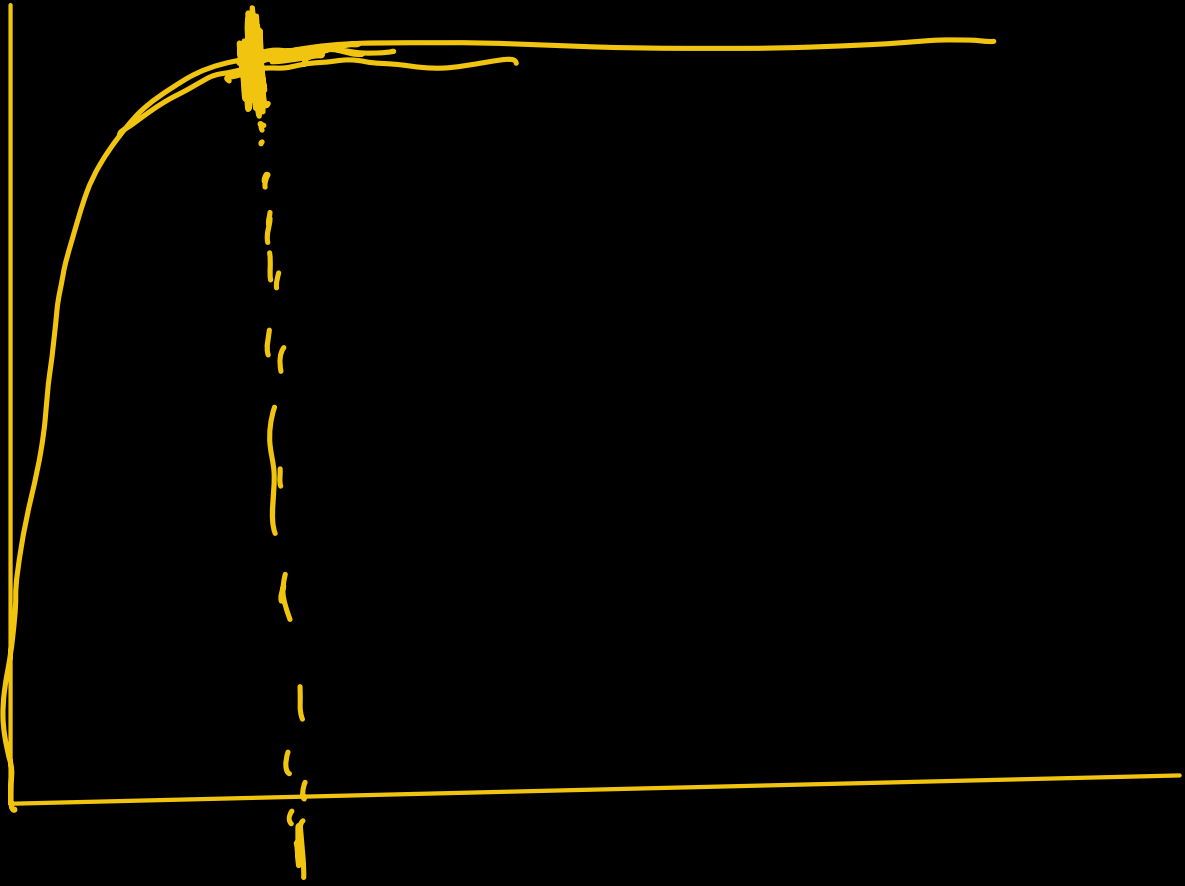
$AUC = 0.9 - 1.0 =$ Excellent

$AUC = 0.8 - 0.9 =$ Good

$AUC = 0.7 - 0.8 =$ Fair model

$AUC = 0.5 \Rightarrow$ No better
than tossing a coin

AUC $< 0.5 \Rightarrow$ Worse than random



Gradient descent

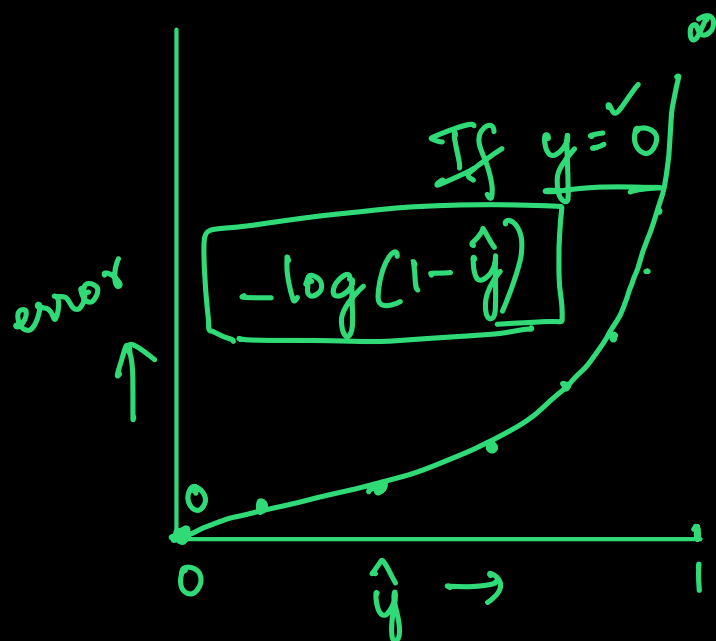
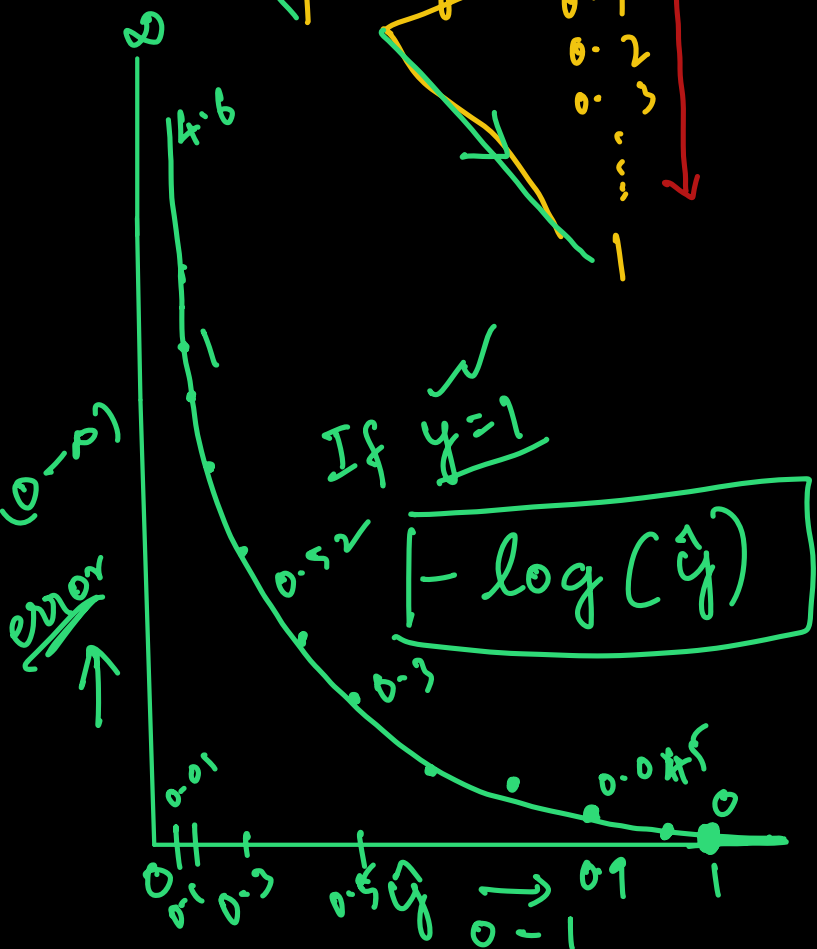
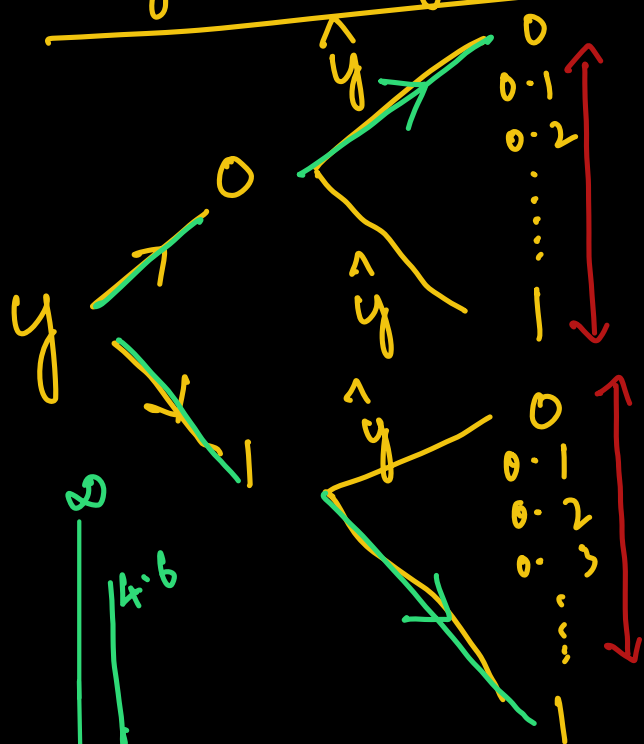
Cost function

Logistic Regression (Classification)

linear Regress

MSE

$$\frac{\sum (y - \hat{y})^2}{n}$$



$$\text{Cost function} = \begin{cases} -\log(\hat{y}) & \text{if } \underline{y=1} \\ -\log(1-\hat{y}) & \text{if } y=0 \end{cases}$$

$$\frac{1}{m} \sum_{i=1}^m y_i [-\log(\hat{y})] + (1-y_i) [-\log(1-\hat{y})]$$

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

$y=0$

$y=1$

Log Loss