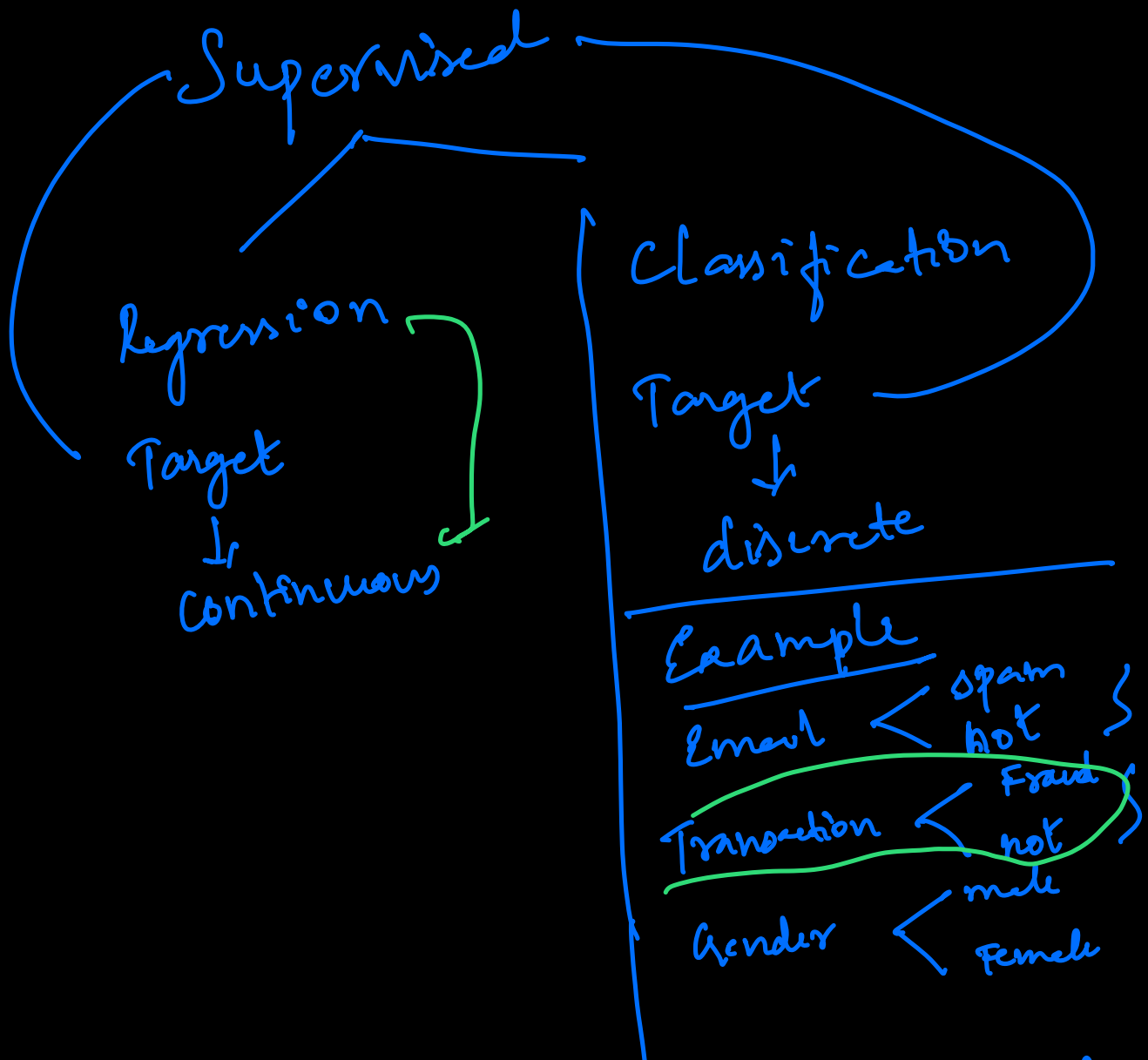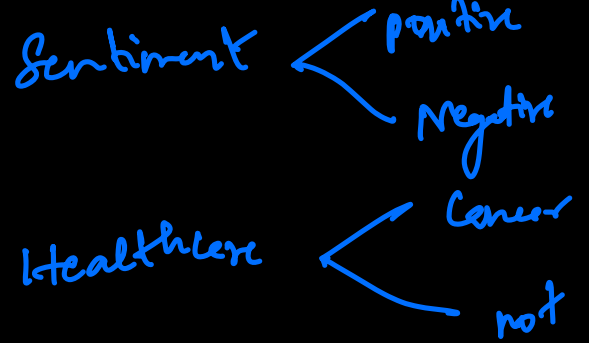So Far

1. Linear Regression
2. Polynomial Regression
3. Ridge
4. Lasso
5. Elastic Net

} Supervised

Regression → Conti

Supervised

Regression
Target
↳
Continuous

Classification
Target
↳
discrete

Example

Email < spam / not }

Transaction < Fraud / not }

Gender < male / Female

# Classification

Binary , Multi label

Sentiment < positive / Negative

Healthcare < Cancer / not

## Logistic Regression

$\quad\hookrightarrow$ misleading

linear regression

Continous Outcome $= \beta_0 + \beta_1 x + \cdots + \varepsilon$

range $= (-\infty, +\infty)$

$\quad 0 \leq\ >$

How? Probability
$\quad\quad 0 - 1$

range $\nearrow$

$(0, 1)$

logistic regression $\nearrow$ $(0, 0.1, 0.2 ; \cdots 1)$

$0 - 0.5 \quad , \quad 0.5 ; 1$

$\quad\quad\downarrow \quad\quad\quad \downarrow$

$\quad( 0 \quad\quad\quad 1 )$

$\quad\quad 0.2 \quad\quad 0.8$

$$y = \beta_0 + \beta_1 x \qquad \text{Linear Regression}$$

$$P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}} \qquad \text{Logistic Regression}$$

1. Input $= \dfrac{(-\infty, +\infty)}{}$

2. $(0, \infty) \quad e^{\cdots} = \checkmark$

exponential $\longleftrightarrow (0, \infty) \longmapsto$

3. $(0, \infty)$

$\uparrow$

$(0, 1) \qquad \dfrac{P}{P+1}$

$P = 561$

$= \dfrac{561}{561 + 1}$

$(0, 1)$

$$y = \beta_0 + \beta_1 x$$

$$P = \frac{\exp(\beta_0 + \beta_1 x)}{\exp(\beta_0 + \beta_1 x) + 1}$$

$$P = \frac{e^{\wedge} y}{1 + e^{\wedge} y}$$

$$q = 1 - P = 1 - \frac{e^{\wedge} y}{1 + e^{\wedge} y}$$

$P$ = probablity

$q = 1 - P$

$$P = \frac{e^y}{1 + e^y}$$

$$P = \frac{\frac{e^y}{e^y}}{\frac{1 + e^y}{e^y}} = \frac{1}{\frac{1}{e^y} + 1}$$

$$\boxed{P = \frac{1}{1 + e^{-y}}}$$

Sigmoid
function

$$p = \cfrac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

$\to$ $p$

$1 - p = q$

Odds $\longrightarrow$ $\cfrac{p}{1-p}$ $=$

$\longrightarrow$ range
$(0, \infty)$

India vs pakistan

$100 \to 80, 20$

$\dfrac{8\emptyset}{2\emptyset} = 4$

$\dfrac{p}{1-p}$

$$p = \cfrac{\dfrac{1}{1 + e^{-y}}}{1 - \dfrac{1}{1 + e^{-y}}}$$

$$= \cfrac{\cfrac{1}{1+e^{-y}}}{\cfrac{1+e^{-y}-1}{1+e^{-y}}}$$

$$= \frac{1}{1+e^{-y}} \times \frac{1+e^{-y}}{e^{-y}}$$

$$= \frac{1}{e^{-y}}$$

$$\boxed{\frac{p}{1-p} = \frac{e^{y}}{}} \quad (0, \infty)$$

Odds

apply log on each side

$$\boxed{\log\left(\frac{p}{1-p}\right) = y}$$

$$\boxed{\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x}$$

$$\frac{1}{1+e^{-y}}$$

$$\underset{logit}{\underline{\log\ odds}} = \underset{(-\infty,\ +\infty)}{linear\ regression}$$

$$p = probablity \longrightarrow (1, 0)$$

$$\frac{p}{1-p} = odds \longrightarrow (0, \infty)$$

$$\log\left(\frac{p}{1-p}\right) = logits\ or\ log\text{-}odds \longrightarrow (-\infty, +\infty)$$

$$\boxed{\log\left(\frac{p}{1-p}\right) = y}$$

apply $exp$ on both

$$\frac{p}{1-p} = e^y$$

$$p = (1-p)\, e^y$$

$$p = e^y - pe^y$$

$$p + pe^y = e^y$$

$$p(1 + e^y) = e^y$$

$$p = \frac{e^y}{1 + e^y} \qquad \div\ e^y \text{ on num}$$
$$\text{denom}$$

$$p = \frac{\dfrac{e^y}{e^y}}{\dfrac{1 + e^y}{e^y}} = \boxed{\dfrac{1}{1 + e^{-y}}}$$

Sigmoid function

$$\text{Logistic Regression} = \frac{1}{1 + e^{-(\underline{\text{Linear Regression}})}}$$

Classification

# Metrics

Yes
   correctly — Yes ✓
           ✓
   wrong — No
       No ✓

No
   correct — No ✓
   wrong
      Yes ✓

| | Prediction class | |
|---|---|---|
| | Yes | No |
| Actual Class **Yes** | ✓ | ✗ |
| **No** | ✗ | ✓ |

**Prediction class**

| Actual Class | Yes *(positive)* | No. *(Negative)* |
|---|---|---|
| **Yes** *Positive* | ✓ Truly Predicting as positive | ✗ Falsely Predicting as Negative → *wrongly Predicting as NO.* |
| **No** *Negative* | ✗ Falsely Predicting as positive → *wrongly predicting as `Yes'* | Truly Predicting as Negative ✓ |

| Actual class | | Predicting Class | |
|---|---|---|---|
| | | Yes | No |
| | Yes | TP | FN |
| | No. | FP | TN |

**Confusion Matrix**

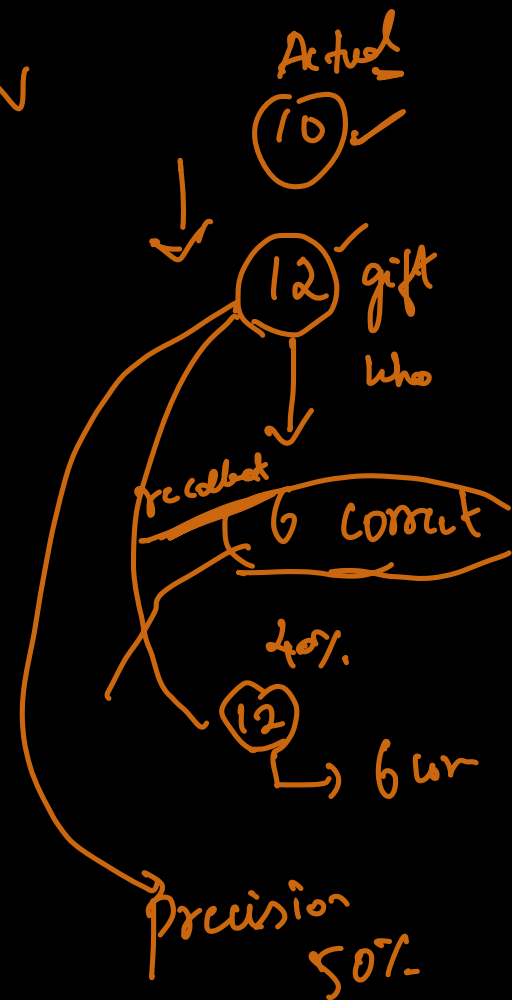|  |  | Predicting | Class |
|---|---|---|---|
|  |  | Yes | No |
| Actual class | Yes | TP ✓ | FN |
|  | No | FP | TN |

Actual Whole Positive — Recall (or) Sensitivity

Actual whole Negative — Specificity

Predicted whole positive — Precision

Predicted whole Negative

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

$$Recall = \text{Based on actual data}$$

(True positive rate)
(True Negative rate)

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

On your Actual data How much you recalled correctly

Actual ⑩

⑫ gift who

recalled 6 correct

40%

⑫ → 6 un

Precision 50%

Precision = Based on Predicted data

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

---

$$10{,}000 \begin{cases} 9990 \rightarrow \text{Non Cancer } \checkmark \\ 10 \rightarrow \text{Cancer } \checkmark \end{cases} \text{Training}$$

$10{,}000 \rightarrow$ Non Cancer $\rightarrow$ Model Prediction

L, (Accuracy) of this model

99.9 % $\checkmark$

L, Too good

$\rightarrow$ Worst

Threating $\rightarrow \dfrac{\text{Murders}}{10 \text{ People}}$

# Precision ✓

Predicting

|  | Yes | No |  |
|---|---|---|---|
| **Yes** | (1) | 9 | 10 |
| **No** | 0 | (9,990) TN | 9,990 |
|  | 1 → | 9,999 | 10,000 |

Actual

$$Acc = \frac{0 + 9,990}{0 + 9990 + 10 + 0} = \frac{9990}{10000}$$
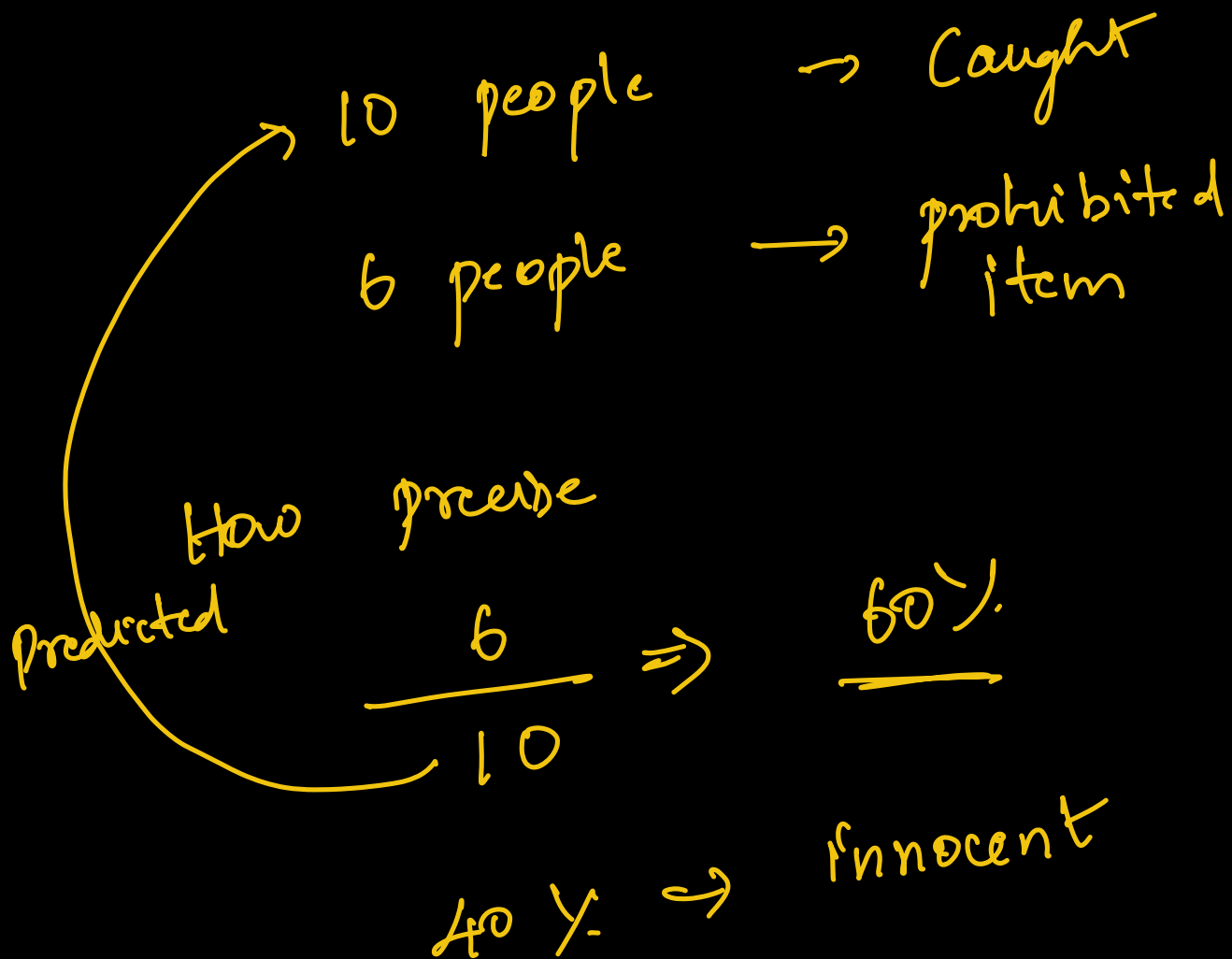
$$\Rightarrow 99.9\% \checkmark$$

$$Precision = \frac{1}{1 + 0} = 1.00\%$$

$$Recall = \frac{1}{1 + 9} = \frac{1}{10} = 0.1$$

$$= 10\%$$

# Security

On a day

→ 10 people → Caught

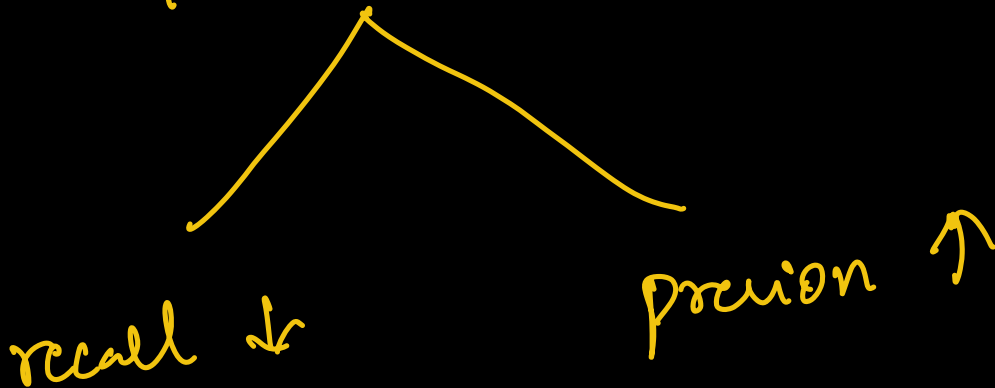6 people → prohibited item

How precise

Predicted $\dfrac{6}{10} \Rightarrow \underline{60\%}$

40% → innocent

20 people have prohibited

$\dfrac{6}{20} = 30\%$

actual

recall

Stop everyone

recall ↑          Precison ↓

Trade off

Stop very very sure

recall ↓          precion ↑

# Model Screens for a Cancer

## High Recall Priority (Don't miss the cancer patient)

Some _healty patient will_ be told to do more test

But you catch almost everyone who actually have cancer

## High Precision Priority (Don't scare healty people)

Only flag patients when I have more confident

Fever healty people get worried unnecessarily

But you might miss - some early stage cancer people

Imagine 1000 emails

$700 \rightarrow$ legitimate

$300 \rightarrow$ spam

Scenario A (Catch all spam)

$TP = 300$   $TPR = 100\%$

$FP = \underline{100}\checkmark$   $\dfrac{600^{TN}}{700} = \cancel{0\%}$ TNR
                                                  $85\%$

| Actual Class | | Predicted Class | |
|---|---|---|---|
| | | Spam | Not Spam |
| | Spam | (300) | 0 | $= 300$ |
| | Not spam | 100 | (600) | 700 |
| | | 400 | 600 | |
| | | | | 1000 |

$$\frac{600}{100 + 600} = \frac{600}{700} = 0.85$$

85% TNR

Perfecly catching the spam but lost 100 important emails

Scenario B: (Protect important mails)

TN = 690

$$TNR = \frac{690}{700} = 98\%$$

$FP = 10$

$TP = 210$

|  | Predicter Spam | Not Spam | |
|---|---|---|---|
| Spam | 210 | 90 | 300 |
| Not Spam | 10 | 690 | 700 |
| | 220 | 780 | 1000 |

Actual

$$TNR = \frac{690}{690 + 10} = 98\%$$

$$TPR = \frac{210}{300} = 70\%$$

90 maols missed from
Spam

10 important classified as
Spam

# Metric of Classification :-

1. Confusion matrix ~
2. Accuracy ✓
3. Recall or Sensitivity ~
4. Specificity ✓
5. Precision ✓
6. F1-score ✓
7. ROC-Curve
8. AUC-ROC