

# Hecture - Logistic Regression

[https://colab.research.google.com/drive/1c21GqeA5S0do0JZ2H6\\_Uc0DRi8WF1SgD?usp=sharing](https://colab.research.google.com/drive/1c21GqeA5S0do0JZ2H6_Uc0DRi8WF1SgD?usp=sharing)

- First classification model = **logReg** ~~Regression~~
- **linReg**  $\xrightarrow{+}$  **LogReg**

## # Notation

no. of samples -  $m$   
no. of features -  $d$   
no. of classes -  $n$

first axis  
 $X - (m, d)$  ↓  
ith sample -  $X[i], X^{(i)}$   
jth feature -  $X[:, j], X_j$

Output  $\rightarrow y -$

Nominal

One hot encoding

Binary Classification (0,1)

$y \rightarrow (m, 1)$

Multiclass Classif.  $n$  classes

$y \rightarrow (m, n)$

2 classes

1
0
1
0
1
0
0

$m \times 1$

3 classes

0	0	1
0	0	1
1	0	0
0	1	0
1	0	1

$m \times 3$

Class1	001
Class2	010
Class3	100

Binary

Multiclass

Parameters -  $[w_0, w_1, w_2, w_3, \dots, w_d]$

# summary :  $w^T x + w_0$  \*\*

Linear Model - Best line of fit passing through data

Assume d-features.



Wine  $\in$   $d+1$  dimensional space.

$$\text{loss} = \text{MSE} = \frac{1}{m} \sum^m (y^i - \hat{y}^i)^2$$

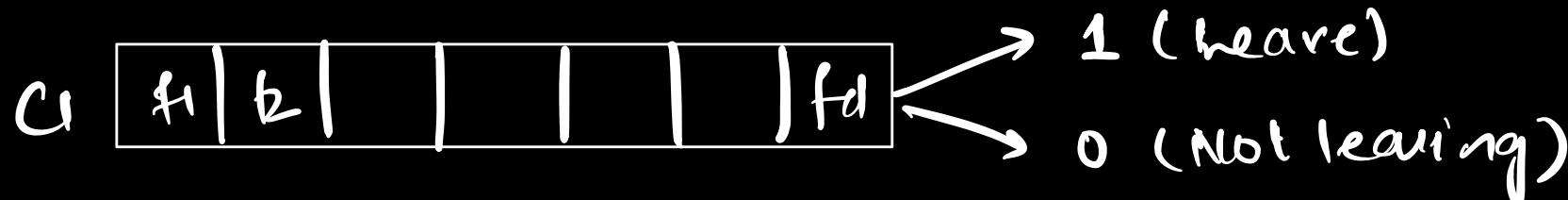
$$y \in (-\infty, \infty), y \in \mathbb{R}$$

## # Intro to Business case - Verizon/Airtel/VI

Customer Churn - leaving subscription.

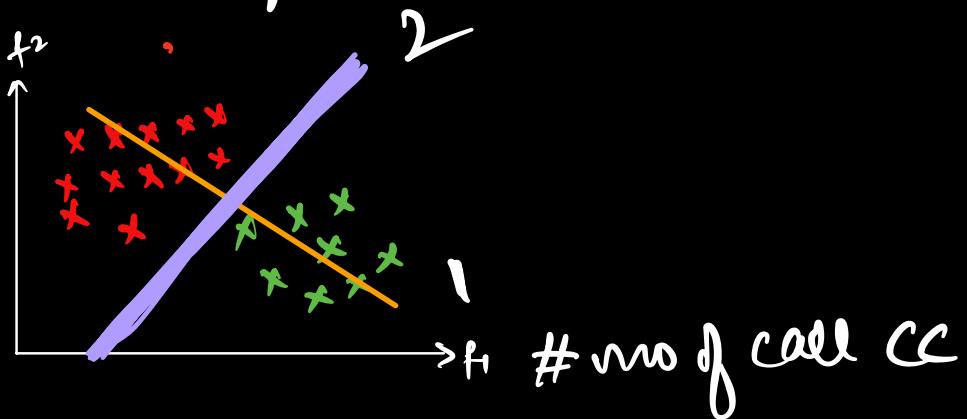
Plan Rate, Network, State/City, # disconnection

Customer service Feedback, # customer care.



- Binary Classification - 1, 0 → not wave
- 
- Historical data with labels (1|0) - **Supervised**
- $y^i \in \mathbb{R}$ ,  $y^i \in \{0, 1\}$

# charges for extra service



$y$  is represented with colors not with axes.



Class 1	0	1	0
Class 2	1	0	1

1 line of best fit

$d+1$  space

one axis for  $y$

$$y = w^T x + w_0$$

2 Decision Boundary

Best line of separation

$d$  space

# Linear Regression-2

## # Agenda

- LinReg → LogReg - Sigmoid
- Geometric Interpretation
- Likelihood vs Prob - MLE
- Error Fn. LogReg - ~~MSE~~ NLL/ LogLoss
- Gradient Descent for LogReg ✓ - 2.e.X

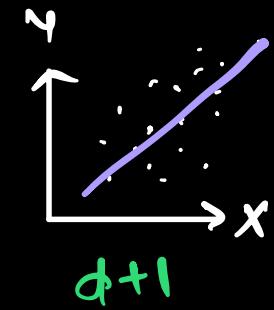
Disclaimer: Content "may" look difficult  
DON'T NEED TO REMEMBER **\*\*\***  
THUMB RULES/RESULTS

# LinReg  $\rightarrow$  Log Reg d+variable

LinReg

$$y = W^T X + W_0$$

$$y = W_0 + W_1 x_1 + W_2 x_2 + W_3 x_3 + \dots + W_d x_d$$



LogReg

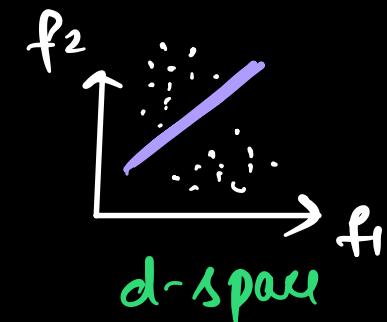
$$\theta = W^T X + W_0$$

$$\theta = W_0 + W_1 x_1 + W_2 x_2 + W_3 x_3 + \dots + W_d x_d$$

d variables.



Equation of a plane



$\Rightarrow$  Separating Hyperplane is present d-space

$$\Rightarrow W^T X + W_0 = 0$$

$$y = \mathbf{w}^T \mathbf{x} + w_0; \quad y \in \mathbb{R} \quad (-\infty, \infty)$$

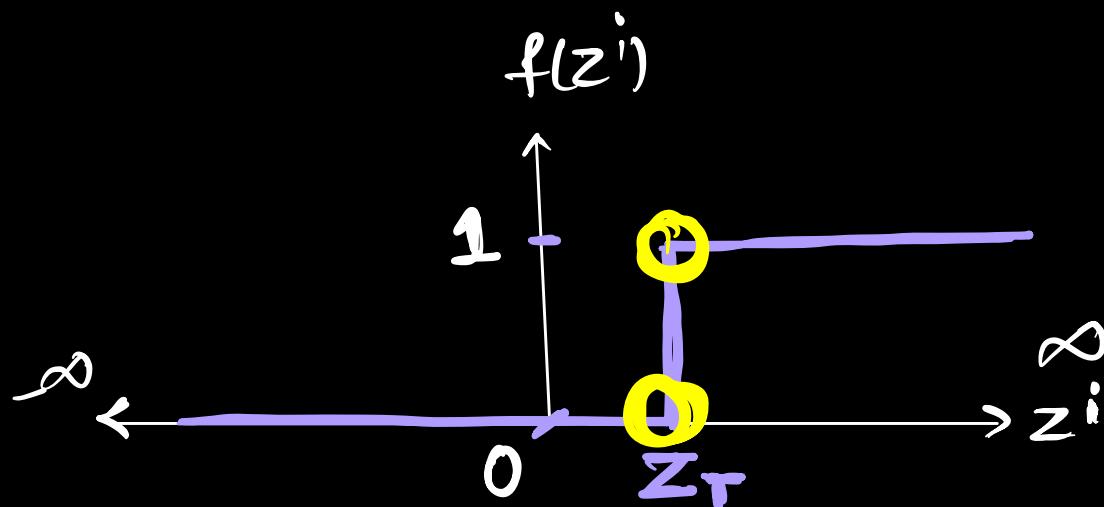
Goal:  $y = f(\mathbf{w}^T \mathbf{x} + w_0) \in \{0, 1\} \cup \{-1, 1\}$

Step 1: Find  $z^i = \mathbf{w}^T \mathbf{x}^i + w_0$

Step 2:  $f(z^i) \rightarrow (0, 1) \text{ or } (-1, 1)$

↑  
find this function.

Solution) Step Function + Thresholding



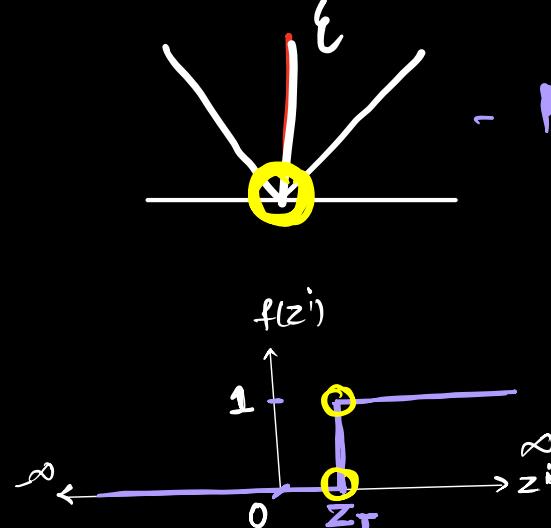
$$z^i = w^T x^i + w_0 \quad - \text{Trivial}$$

If  $z^i > z_T \rightarrow 1$   
elif  $z^i < z_T \rightarrow 0/-1$   
else  $z^i = z_T \rightarrow \text{Unclassified.}$

$$y = \text{SF}(w^T x^i + w_0)$$

## # Problems with SF ?

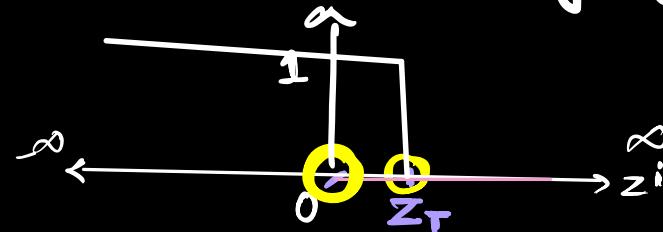
SF has sharp edges



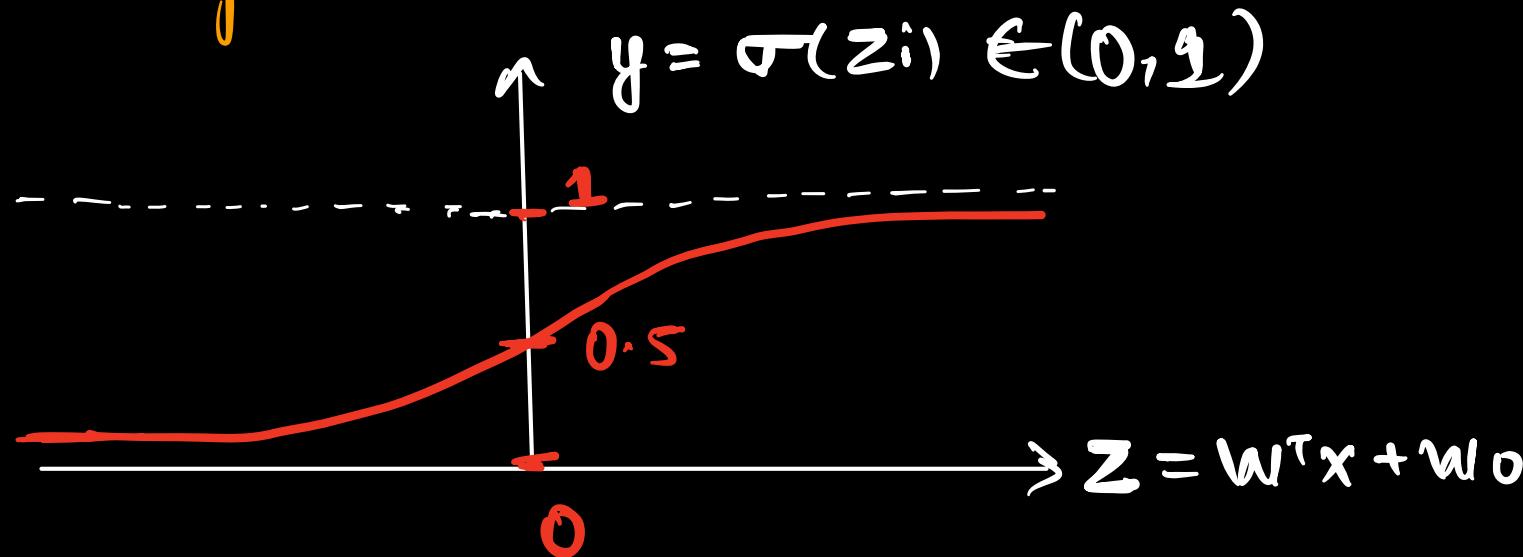
- Not differentiability



$$E = |y - y'|$$



## # Sigmoid Function.



$$y \in (0, 1)$$

- ① Squash  $z \in (-\infty, \infty) \rightarrow (0, 1)$
- ② Differentiability ✓
- ③ Probabilistic Angle to Classification

$$y \in (0, 1) \quad y \in \{0, 1\}$$

~~$y \in \{0, 1\}$~~

$$P(y=1|x^i) = \sigma(z^i)$$

Case 1)  $Z^i \rightarrow \infty$ ;  $\sigma(Z) \rightarrow 1$

Case 2)  $Z^i \rightarrow -\infty$ ;  $\sigma(Z) \rightarrow 0$

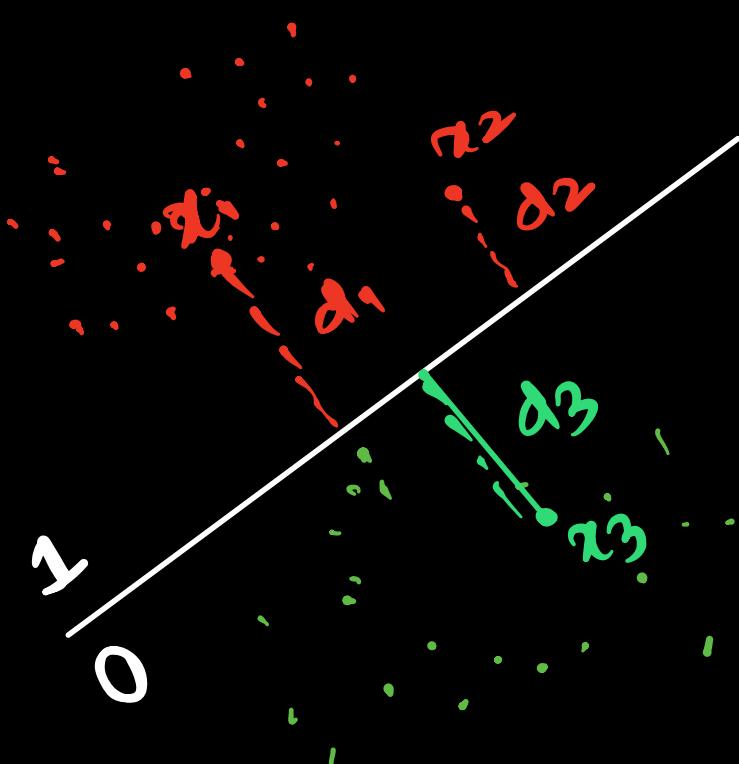
Case 3)  $Z^i \rightarrow 0$ ,  $\sigma(Z) = 0.5$

Case 4)  $Z^i = +ve$ ,  $\sigma(Z) > 0.5$

Case 5)  $Z^i = -ve$ ,  $\sigma(Z) < 0.5$ .

$$\sigma(Z) = \frac{1}{1 + e^{-Z}}$$

# # Geometric Interpretation .



$\pi: w^T x + w_0 = 0$  Decision Boundary

$$d_1 = \frac{w^T x^1 + w_0}{\|w\|}$$

$$d_2 = \frac{w^T x^2 + w_0}{\|w\|}$$

$$d_3 = \frac{w^T x^3 + w_0}{\|w\|}$$

Just a constant

Numerator

$$d \propto w^T x + w_0 = z$$

→  $|z|$  tells us how far is  $x$  from  $\pi$

→  $+/- \theta_z$  tells us which class does it belong

$x_1 \in \text{Class 1}$   
 $x_2 \in \text{Class 1}$   
 $x_3 \in \text{Class 0}$

# # MLE · (Maximum Likelihood Estimation)

$$p^i \quad P(y^i=1|x^i) = \sigma(\underbrace{w^\top x^i + w_0}_{z^i}) \in (0,1) \quad \text{--- ①}$$

$$1-p^i \quad P(y^i=0|x^i) = 1 - P(y^i=1|x^i) \quad \text{--- ②}$$

$\uparrow \quad = 1 - p$

Observed Data.

Ground Truth

Combine ① and ②:

$$(p^i)^{y^i} (1-p^i)^{1-y^i}$$

$$\text{If } y^i=1, \underline{\underline{1}} = p^i \quad \text{--- ①}$$

$$\text{If } y^i=0, \underline{\underline{1}} = 1-p^i \quad \text{--- ②}$$

common eq.

Independent Event A, B  $P(A \cap B) = P(A) \cdot P(B)$

$$P(A, B, C) = P(A) \cdot P(B) \cdot P(C)$$

$$P(A, B, C, \dots, Z) = P(A) \cdot P(B) \cdots P(Z)$$

m samples, independent

$$P(Y^1=1, Y^2=1, Y^3=1, \dots, Y^m=1 | X) = P_1^1 \cdot P_2^2 \cdot P_3^3 \cdot P_4^4 \cdots P_m^m$$

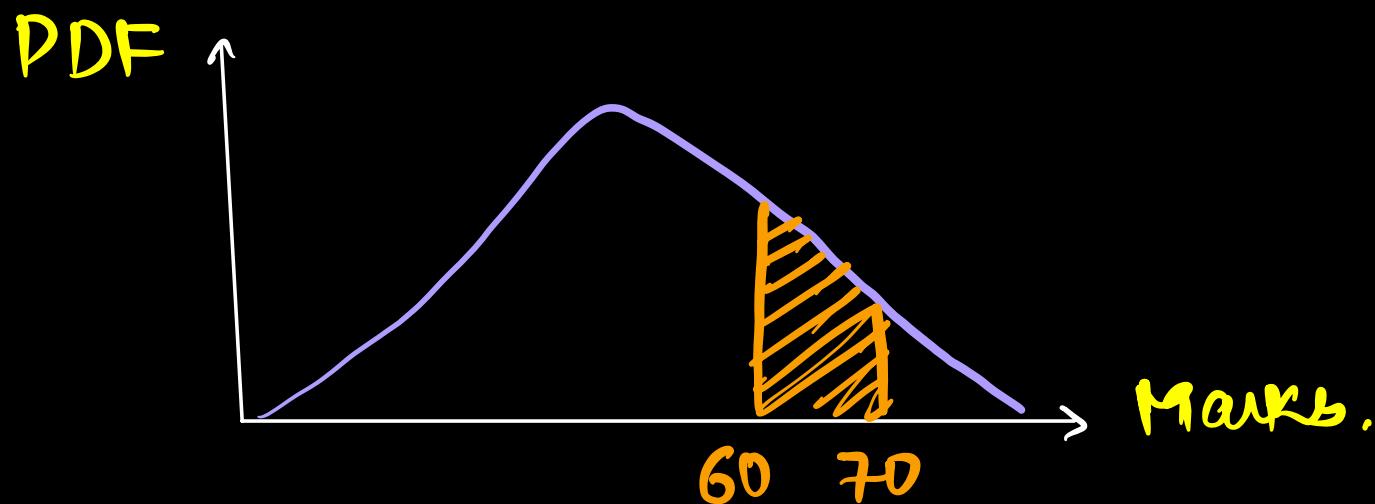
$$\sum_{\text{Multiplication}} = \prod_{i=1}^m [(p^i)^{y^i} \times (1-p^i)^{(1-y^i)}]$$

= PDF

= Joint Distribution of data

= Likelihood.

## # Prob via Likelihood



AUC b/w 60 and 70 of a PDF

$$P(60 < X < 70 | \mu, \sigma)$$

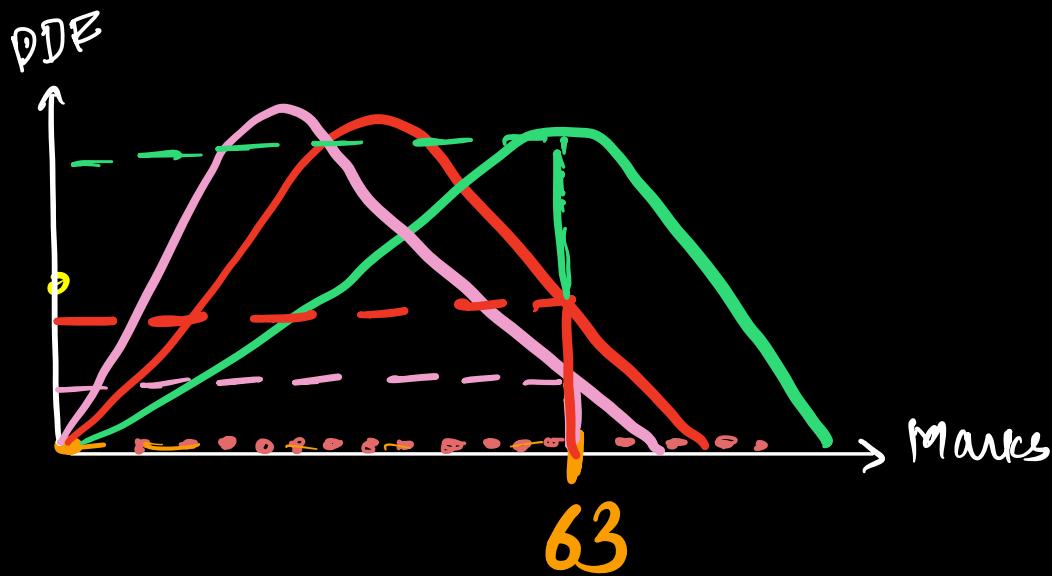
↑ Distribution is given

BUT in our case,

- We have Data
- We don't know the distribution.

$$= P(\mu, \sigma | \text{Data})$$

$$\approx P(\mu, \sigma | m_1 = 63)$$



You only one data pt.  $m_1 = 63$ , which distribution  
is it most likely coming from. - Green

You have  $m$  data pt  $\rightarrow$  Which dist. Is  
most likely?

$P(\text{Distribution} | \text{Data})$



\*\*\*

Likelihood.

Maximum Likelihood Estimation

---

$$\text{Likelihood} = \prod_{i=1}^m (p^i)^{y_i} \cdot (1-p^i)^{1-y_i}$$

↑  
maximize =

$$\log(\text{Likelihood}) = \log \left( \prod_{i=1}^m (p^i)^{y_i} \cdot (1-p^i)^{1-y_i} \right)$$

$$\log(A \cdot B \cdot C) = \log(A) + \log(B) + \log(C)$$

$$= \sum_{i=1}^m \log \left( (p^i)^{y^i} \cdot (1-p^i)^{1-y^i} \right)$$

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

$$\log A^B = B \log A$$

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

Maximise

Class 1

Class 0

$$\text{Case 1) } y^i = 1 \rightarrow \log p^i$$

$$\text{Case 2) } y^i = 0 \rightarrow \log(1-p^i)$$

Loss function for training a model has to be minimized.

$$NLL = - \sum_{i=1}^m y^i \cdot \log p^i + (1-y^i) \log(1-p^i)$$

Minimize

Negative Log likelihood.

③      ②      ①

- step1) Calculate Likelihood of m samples occurring
- step2) Take log both sides
- step3) Take negative of

$$\text{log loss} \left\{ \begin{array}{l} -\log p^i \text{ if } y^i = 1 \\ -\log(1-p^i) \text{ if } y^i = 0 \end{array} \right.$$

$p^i \rightarrow$  Output of log Reg model.

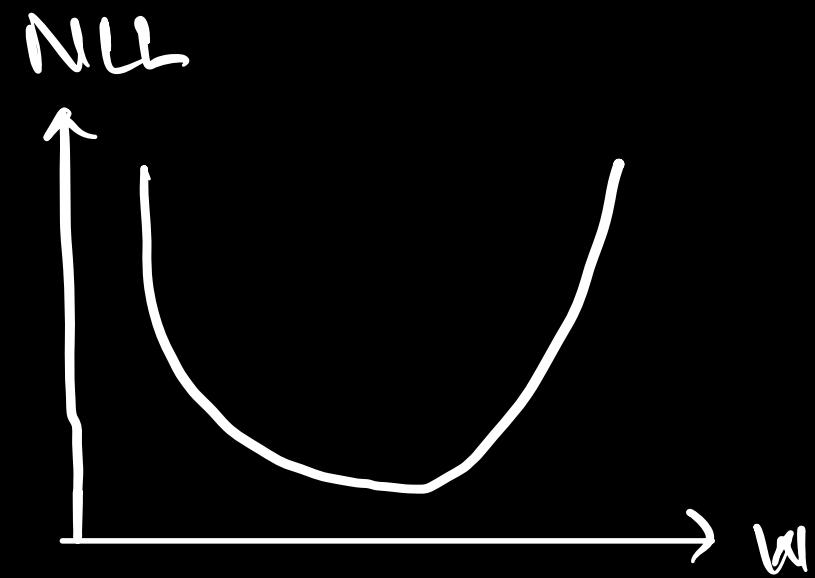
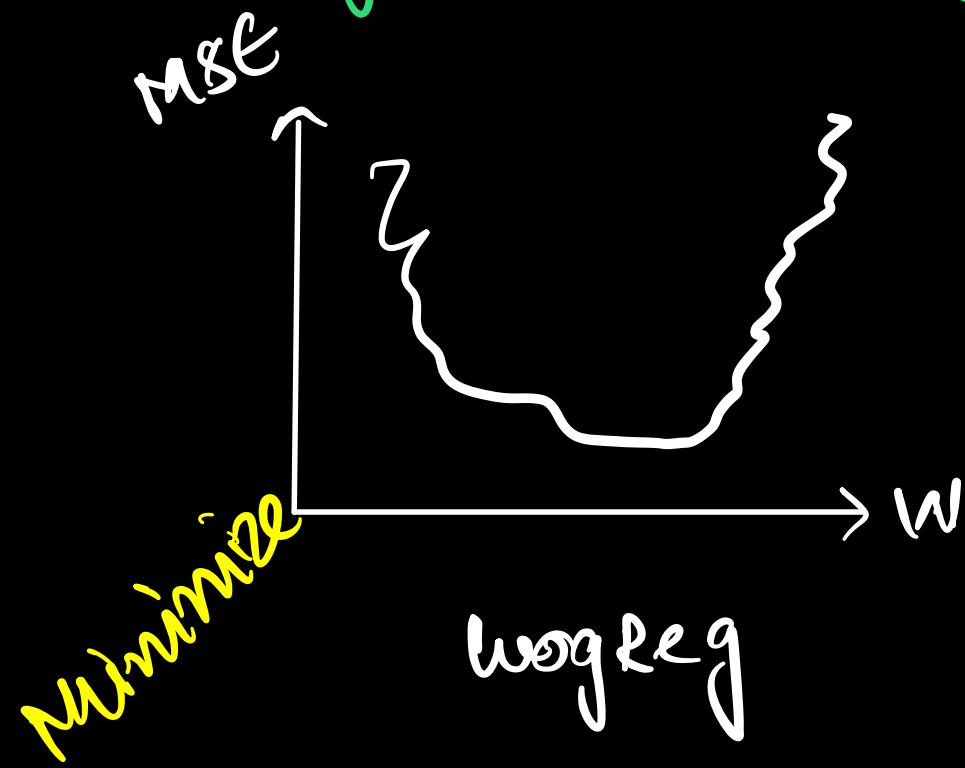
$$\text{log loss} \left\{ \begin{array}{l} -\log \hat{y}^i \text{ if } y^i = 1 \\ -\log(1-\hat{y}^i) \text{ if } y^i = 0 \end{array} \right.$$

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

↑  
minimize.

Q why not simply MSE?

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

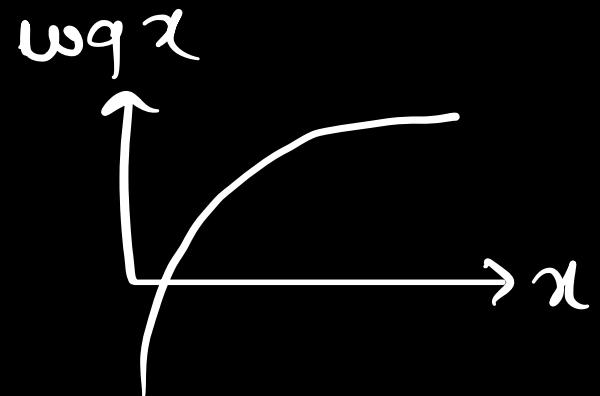
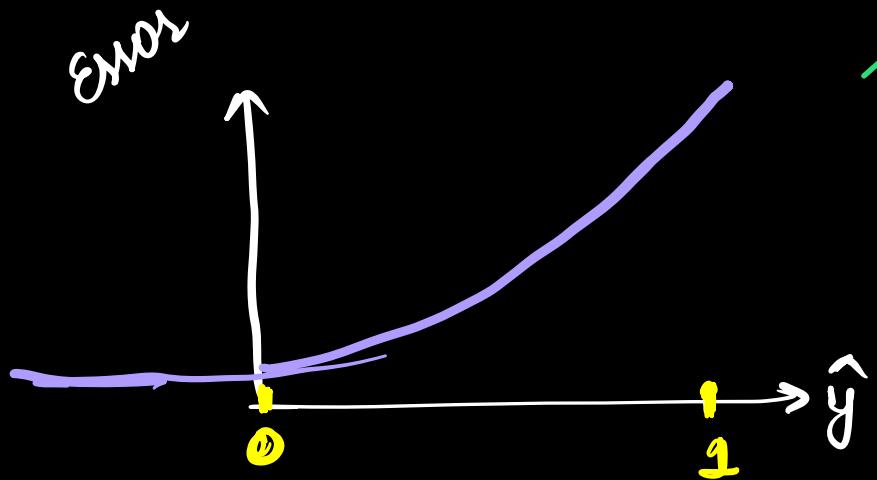


minimize  
log reg

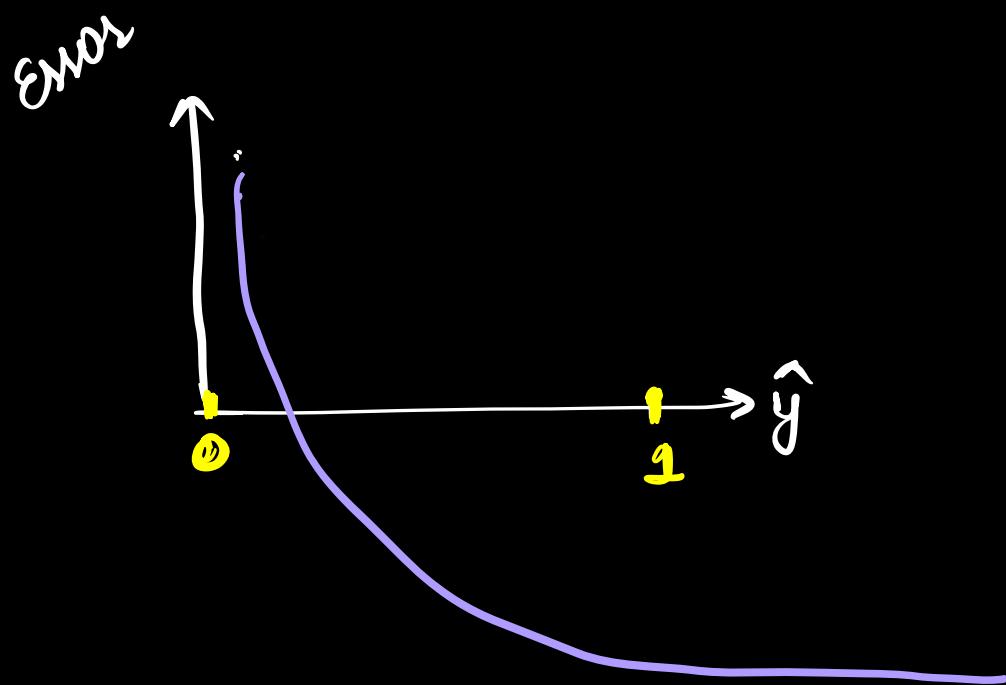
log reg  
☆☆☆

$$NLL = - \sum_i^m y^i \log \hat{y}^i + (1-y^i) \log (1-\hat{y}^i)$$

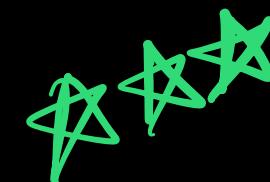
case1) GT=0 ;  $y^i=0$  ;  $NLL = -\log(1-\hat{y})$



Case 2)  $GT = 1, \hat{y}^i = 1, NLL = -\log(\hat{y}^i)$



Thing to remember:



$$NLL = -\sum_i^m y^i \log \hat{y}^i + (1-y^i) \log(1-\hat{y}^i)$$