

→ Deep Learning: Learning Theory √ Finite hypothesis

Infinite hypothesis → VC-dimension

$$\checkmark m \geq \frac{1}{\epsilon} (\ln |H| + \ln(\frac{1}{\delta}))$$

for finite hypothesis w/ train error = 0

$$\checkmark m \geq \frac{1}{2\epsilon^2} (\ln |H| + \ln(\frac{1}{\delta}))$$

finite hypothesis w/ train error ≠ 0

$|H| \rightarrow \infty$ Model complexity is computed using
VC dimension of the hypothesis class.

VC-dimension: The maximum number of
samples that can be shattered by a hypothesis
class is called VC-dimension.
" " algorithmic VC

Vapnik-Chervonenkis Dimension

SHATTERER: For all possible labellings of the

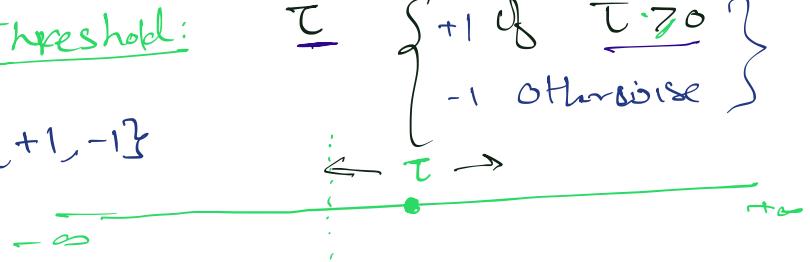
Samples, if there exist a hypothesis (model) ' h ' with zero error (correctly classify all samples) in the hypothesis class ' H ', is called Shattering.

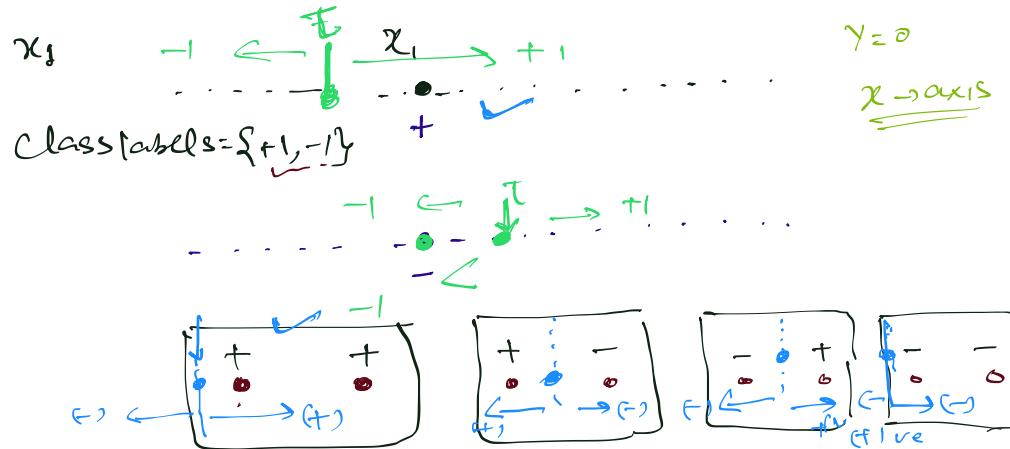
Example:

Let Threshold: τ $\begin{cases} +1 & \text{if } \tau \geq 0 \\ -1 & \text{otherwise} \end{cases}$

Class: $\{+1, -1\}$

Samples:

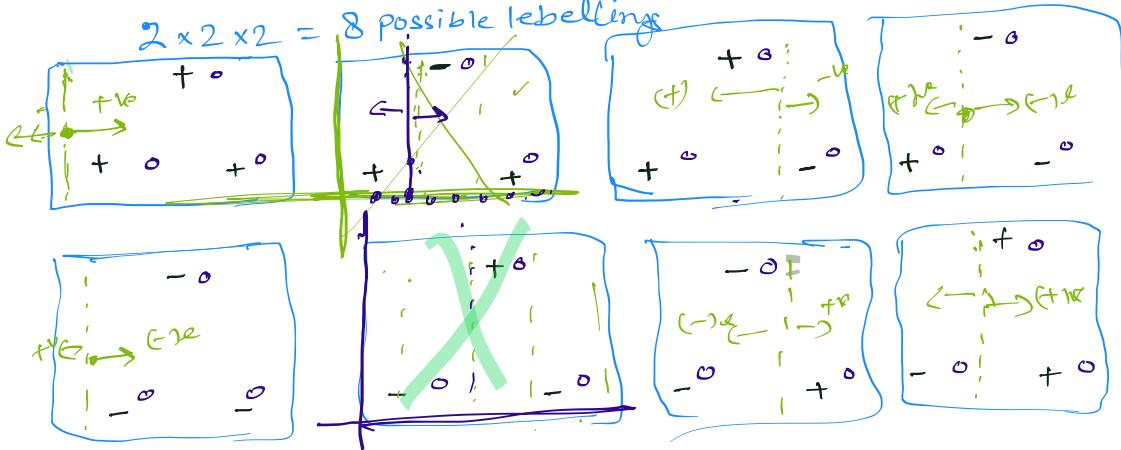




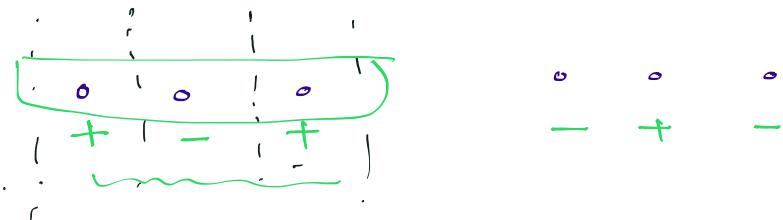
2-Samples, 2 classes:
 $\frac{1}{2} \times \frac{1}{2} = 1$

3-Samples, 2 classes

$$2 \times 2 \times 2 = 8 \text{ possible labellings}$$



threshold does not exist which can shatter.



VC-dimension for threshold hypothesis class is
2. I + can shatter a maximum of
2. samples.

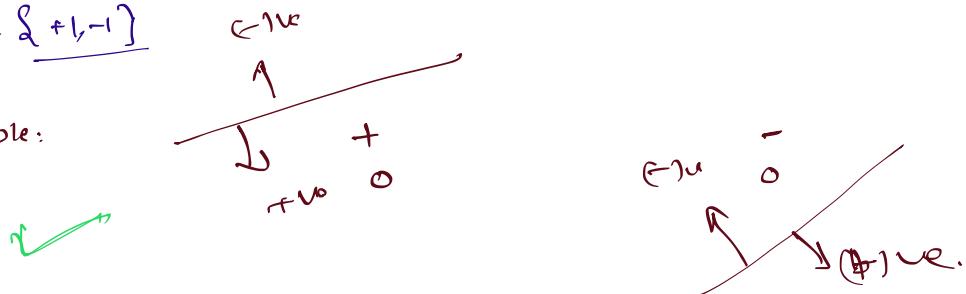
Example 2.

H: Linear models / linear hyperplane.

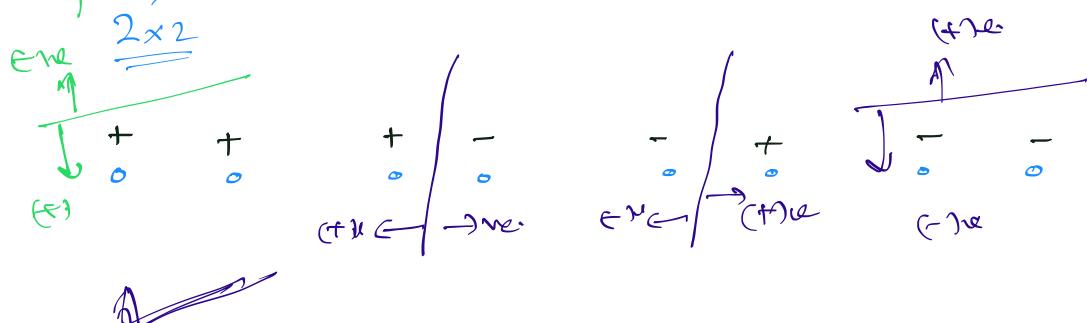
Let example be H: linear model in 2D (line)

$$\text{Class} = \{+1, -1\}$$

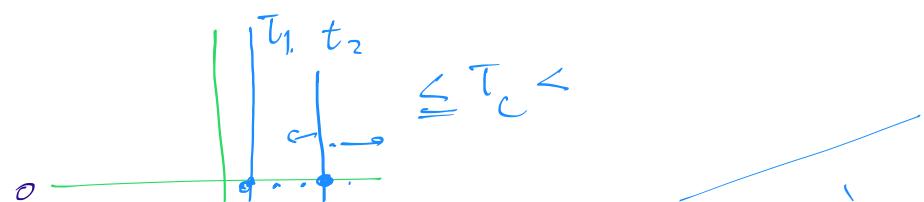
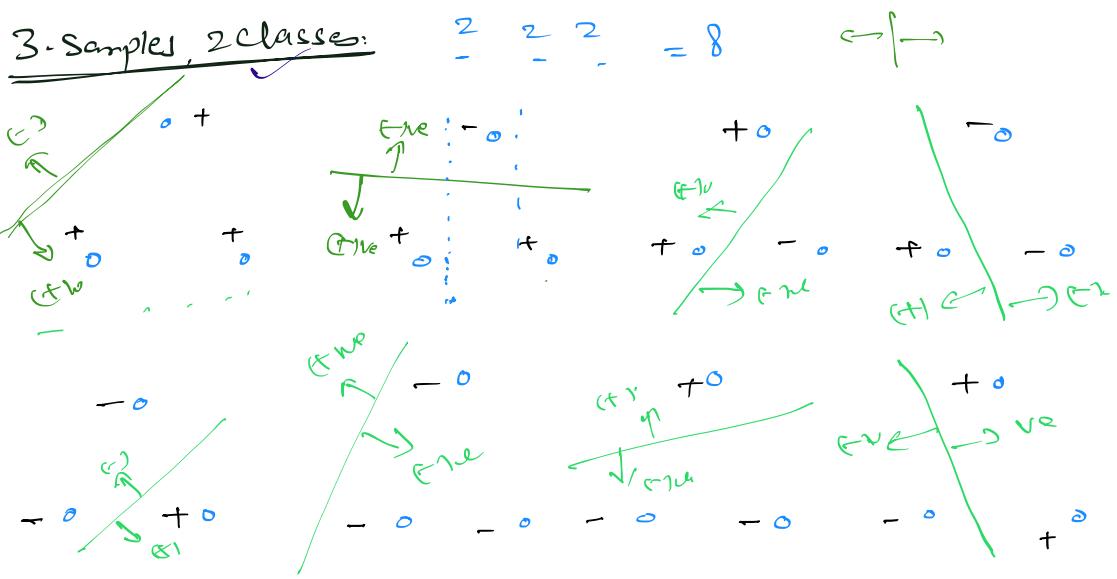
1. Sample:

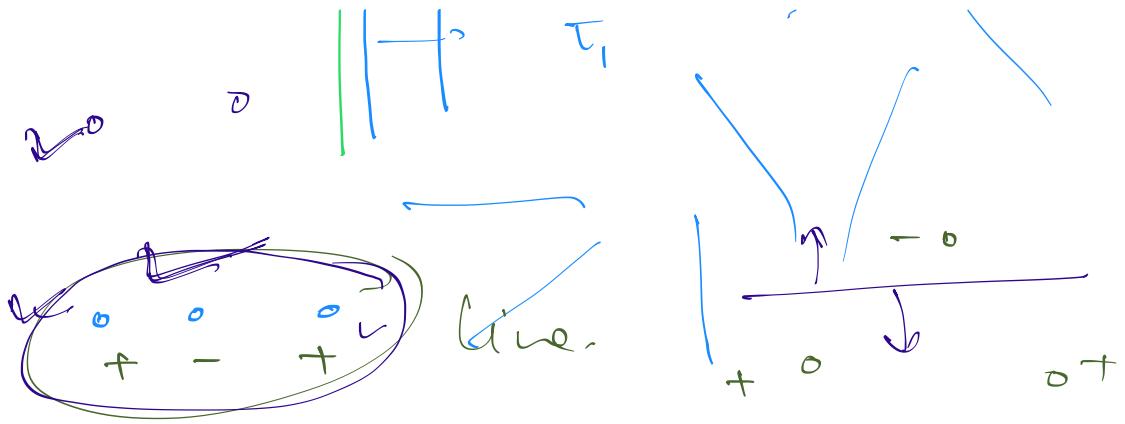


2 Samples, 2 Classes. $\underline{2} \cdot \underline{2} = 4$



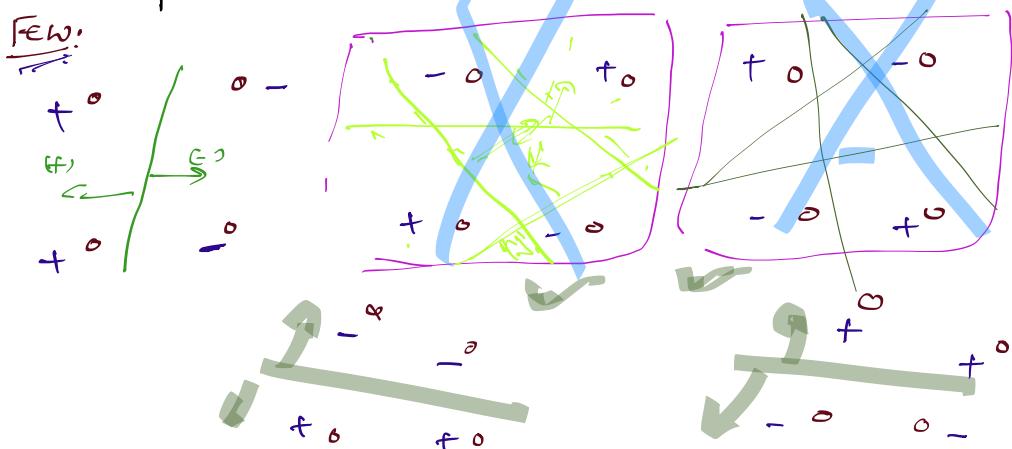
Line can shatter 2 points.





Line can shatter 3 samples.

4 samples, 2 classes : $\underline{2} \times \underline{2} \times \underline{2} \times \underline{2} = 16$ possible labellings.



line can not shatter 4-points

Line can shatter at maximum of 3-points.

VC-dimension of line (or linear hypothesis in 2D)

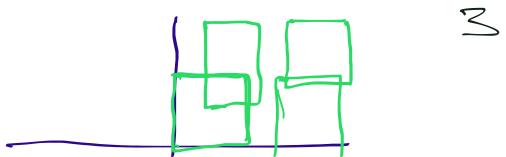
$$\leq 3$$

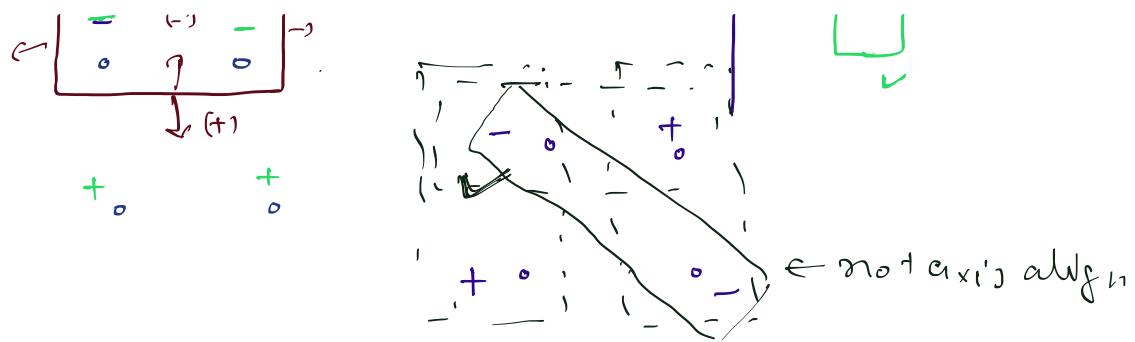
VC-dimension of Linear hyperplane in d-dimensions

is $d+1$

Hypothesis is Axis-aligned Rectangle (2D)

4-sample class = 2
1
2
3
4



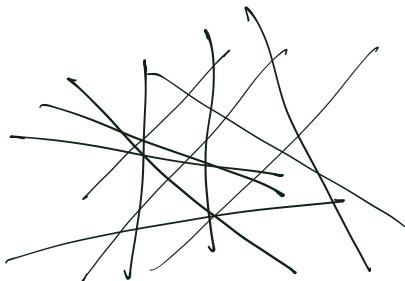


Rectangle in 2D \rightarrow VC - dimension (4)

H = Line in 2D \rightarrow

$$y = mx + c$$

$$\underline{w^T x_i + b = 0}$$



- ✓ b : bias (Distance of the plane from origin of the vector space)
- ✓ w : weight (Vector perpendicular to the plane having x_i)

In a d-dimension - i^{th} sample x_i

$$\mathbf{x}_i = \begin{bmatrix} x_i^1 \\ x_i^2 \\ \vdots \\ x_i^d \end{bmatrix} \in \mathbb{R}^{d \times 1}$$

$$\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \in \mathbb{R}^{d \times 1}$$

$$x_i \in \mathbb{R}^d \quad w \in \underline{\mathbb{R}^d}$$

$V_C(h) \propto \frac{1}{\text{parameters}}$

$$h_{w,b}(x_i) = \underline{\underline{w^T x_i + b}}$$

$$\frac{h(x_i)}{w^T x_i + b} = \pm$$

Linear model in
d-dimensions.

$$m \geq O\left(\frac{VC(H)}{\epsilon}\right)$$

$m \geq O(VC(H) \lg VC(H))$

$$m \geq O\left(VC(H)^2\right)$$

$m \geq O\left(\frac{1}{\epsilon^2}\right)$

$m \geq O\left(\frac{1}{\epsilon}\right)$

$m \geq O\left(\ln |H|\right)$

"

~~Excess~~ [Tolerance] —————

True Error / Generalization Error

$$= \text{Train Error} + \epsilon$$

→ Finite Hyp

Structural Error

i) $m > \frac{1}{\epsilon} (\ln(H) + \ln(\frac{1}{\delta}))$ \uparrow Train Err = 0

ii) $m > \frac{1}{2\epsilon^2} (\ln(H) + \ln(\frac{1}{\delta}))$ \uparrow Train Err = 0

iii) $m \geq \frac{1}{\epsilon^2} \left(C_1 \underline{\text{VC}(H)} + C_2 \cdot \ln \left(\frac{C_3}{\delta} \right) \right)$ $\underbrace{C_1, C_2, C_3}_{\text{Hyp. Constant}}$

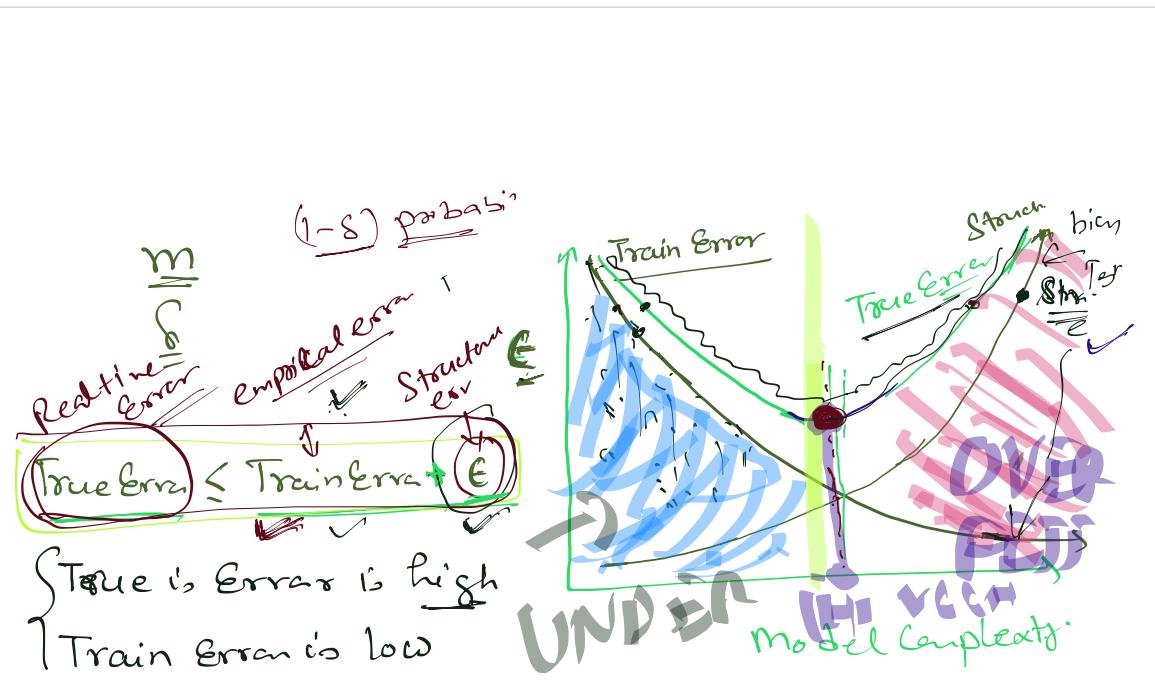
(i) $E = \frac{\ln(M) + \ln(\frac{1}{\delta})}{m}$

$$E \propto (1-\delta)$$

Prob (Conf)

(ii) $E = \frac{\ln(M) + \ln(\frac{1}{\delta})}{2m}$

(iii) $E \approx \frac{C_1 \cdot VC(H) + C_2 \ln(\frac{2}{\delta})}{m}$



~ overfitting (Variance)

VC(H) a |H|

Train Erra \rightarrow high : Under fitting. (bias _{probly})

Good hypothesis : Low overfitting low Under fit.

high Variance \rightarrow overfitting. }

high bias \rightarrow Under fitting } \checkmark

Bias - Variance Trade Off

Model : low bias & low variance is a good v

→ Model is suffering from overfitting (Variance) issue:

- Remedies: → Reducing the model complexity
- ✓ ↳ Regularization
- ✓ ↳ Dropout
- ($E \rightarrow$ reduces) ↳ Change to less complex model
- ✓ • Increase the # samples in the training
- ↳ Augmentation (Data methods)
- ✓ • NN: (early stopping) ↳ High Confidence
- ✓ • "g": $\rightarrow \alpha(1-\delta)$ reduced.
- Explain(?)

$$P(\text{True Error} \leq E) > (1-\delta)$$

(True = γ)

$\delta = 0.01$
 $1-\delta = .99$
 $\delta = 0.05$
 $1-\delta = .95$

=
Confidence interval, p-value
Significance
Null hypothesis

→ Distribution Shift: ~~test~~
Training data, Validation data, Test data

If there is a mismatch (all these three does n't follow same distribution: distribution shift)

↳ Collecting data from different sources

~~Validation~~ ↳ Data splitting during training
vs (upto 1000)

↳ Cross-validation (

Small-data: Leave-one-out (~~Leave-one~~)

medium data: 80-20, 70-30, 60-40, (Train, validate) + K-fold
Cross-val.

Large dataset: 1 to 10% data for validation

↳ low K-fold ($K = 3 \text{ to } 5$)

↳ Random sampling for fixing the hyperparameters.

Hyper-parameters: Set by the user,
(Cross-validation, Splitting)

Validation

Step 1: Split data into 2-parts → Training

Testing

Step 2: Training data → Training + Validation

Step 3: Fit the model w/ hyperparameters with
good validation accuracy. (CV, splitting)

Step 4: Fixed Hyperparameters in the previous step

Train the model ^{Whole} on Training Data and
Evaluate model on Test Data

Some steps

Train error
 $M_1 : 0.1 (10\%)$

Validation
 $M_{10} (12\%)$

Test error
 $M_{100} (10\%)$

Real T. ✓
Test error

Validation
error

Test error

Test error

Test error



Distribution shift is found:

- . Try to Increase the %age of validation
- Random Sample the samples from Test Distribution and Try to finetune the large model on this random samples. (If you are not able to retrain) otherwise Augment the Training Data \rightarrow random Sample from Test distribution and retrain the model and fix the hyper-parameter using validation data
- Normalized
feature

→ Model is suffering from Underfitting (bias)

Issue: (Training error is high)

↳ Increase the model Complexity

- ↳ Somehow try to increase the VC-dim.
 - ↳ Introduces non-linearity
 - ↳ Change to complex model.

↳ Increase the epochs (iterations)

↳ Look into the Data

- ↳ Preprocessing (lossy techniques)
 - ↳ Enhancement

Biased Data

Hand

Winter

Blood Image

Example:

	Training Error	Validation Error	Test Error	Problem
<u>M₁</u>	10%	9%	11%	Low bias, low var. No dist shift
<u>overfitting M₂</u>	8%	20%	18%	Low bias, High variance, ??
<u>M₃</u>	30%	40%	45%	High bias, ??
→ M ₄	13%	14%	20%	Low bias, low variance, Dist.

~~Shift Issue~~ | Shift Issue

→ These remedies doesn't guarantee the perfect model:

. In practice it works

. ~~medium size~~
• Moderate size data is needed > 10k
 $\epsilon \propto \sqrt{\frac{1}{m}}$: Precise samples

• Experience Matters (Expert)

+ Comparing the models: Statistical Test are required
errors: $(\bar{U}_m, \pm I_m), (U_m, \pm I_m) - \dots - ()$

Model with low mean err and low confidence interval of err is a GOOD model

i^{th} model: Repeat the training ' n ' times to different needs.

$$m_{m_i} = \frac{1}{n} \cdot \sum_{j=1}^n \frac{\text{valid/test.}(e_j)}{e_j}$$

e₁ ↗ 1
 e₂ ↗ 2
 :
 e_n ↗ n
 Testen ✓

→ Computation Graph : Pictorial Rep
→ Pictorial Representation of a Computation

$$\Rightarrow f(x) = x^2 = \underline{\underline{x*x}}$$

$f(x,y) = \frac{\downarrow}{\cancel{x*y}}$

$\frac{\partial f}{\partial x} = 2x$

$\frac{\partial f}{\partial y} = 0$

$\frac{\partial f}{\partial x} = x$

$\frac{\partial f}{\partial y} = 1$

$\frac{\partial f}{\partial x} = x + x = \underline{\underline{2x}}$

$f(x, y) = \underline{x} + \underline{y}$

$$\frac{\partial f}{\partial f} = 1$$

Chain Rule

$f(x) = \log(x)$

$$\frac{\partial f}{\partial f} = 1$$

Chain Rule

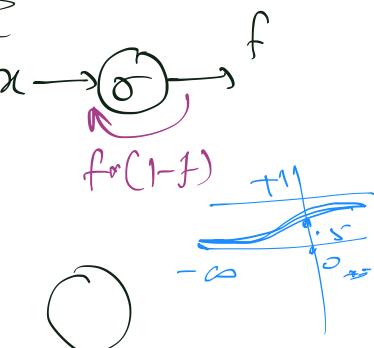
$$\sigma(x) = \frac{1}{1+e^{-x}} = \text{sigmoid}$$

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1-\sigma(x))$$

$$f(x) = \frac{1}{x} \checkmark$$

$$f(x,y) = x \div y \text{ (or } x/y) \Rightarrow$$

$$\frac{\partial f}{\partial x} \cdot \frac{\partial x}{\partial z} + \frac{1}{y} \Rightarrow \frac{1}{y}$$



nom
denom

$$\frac{\partial f}{\partial \text{ nominator}} = \frac{1}{\text{denominator}}$$

$$\frac{\partial f}{\partial \text{den.}} = \frac{-\text{nominator}}{(\text{denominator})^2} \checkmark$$

$$\frac{\partial f}{\partial y} = \frac{y \cdot \frac{\partial x}{\partial y} - x \cdot \frac{\partial y}{\partial y}}{(y)^2} = -\frac{x}{y^2}$$

U "

dem x $\xrightarrow{\frac{1}{y}}$ $\frac{x}{y}$ \xrightarrow{f} $\frac{\partial f}{\partial f} = 1$ ✓

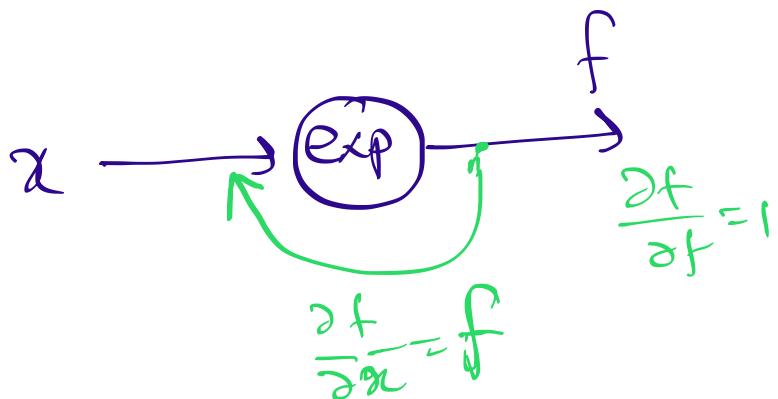
dem y $\xrightarrow{-\frac{x}{y^2}}$ f

$f(x,y) = \underline{x-y}$

$\frac{\partial f}{\partial x} = 1 ; \frac{\partial f}{\partial y} = -1$

$x \xrightarrow{f} -$ $\frac{\partial f}{\partial f} = 1$

$$f(x) = \underline{e^x} ; \quad \frac{\partial f}{\partial x} = \underline{e^x}$$



$$f(x) = \sin(x)$$

$$\frac{\partial f}{\partial x} = \cos(x)$$

$$f'(x) = \cos(x)$$

$$\frac{\partial f}{\partial x} = -\sin(x)$$

$$\frac{\cos^2 x}{\cos^2 x + \frac{\sin^2 x}{2}} + \frac{\sin^2 x}{2}$$

$$f(x) = \tan(x) = \frac{\sin(x)}{\cos(x)}$$

$$\frac{\partial f}{\partial x} = \frac{\sec^2(x)}{\sqrt{\cos(x)}} = \frac{1 + \tan^2 x}{\sqrt{\cos(x)}}$$

$$= \frac{\cos(x) \cdot \frac{\partial \sin(x)}{\partial x} - \sin(x) \cdot \frac{\partial \cos(x)}{\partial x}}{(\cos(x))^2}$$

$$\Leftarrow = \frac{(\cos^2(x) + \sin^2(x))}{\cos(x)}$$

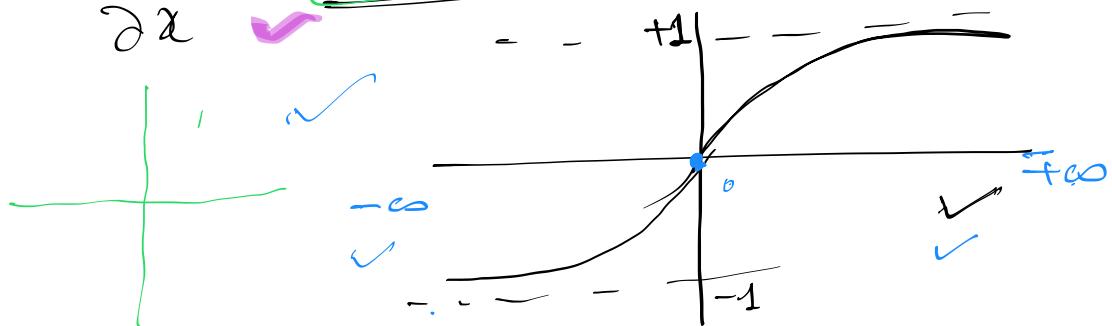
$\cos u$

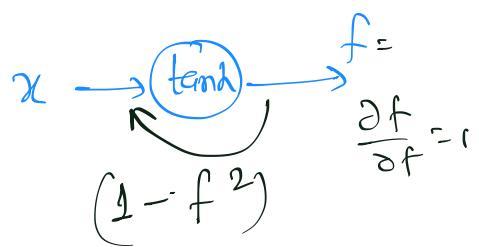
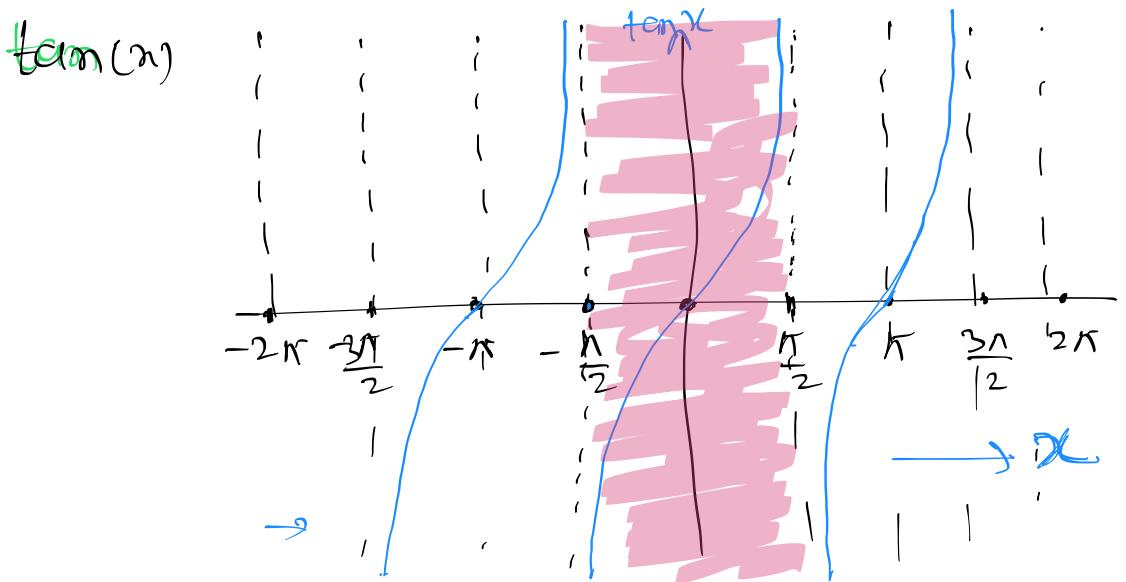
$\cos ny$

$(\cos u) -$

$$f(x) = \tanh(x) \quad (\text{hyperbolic tangent})$$

$$\frac{\partial f}{\partial x} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$





Revn

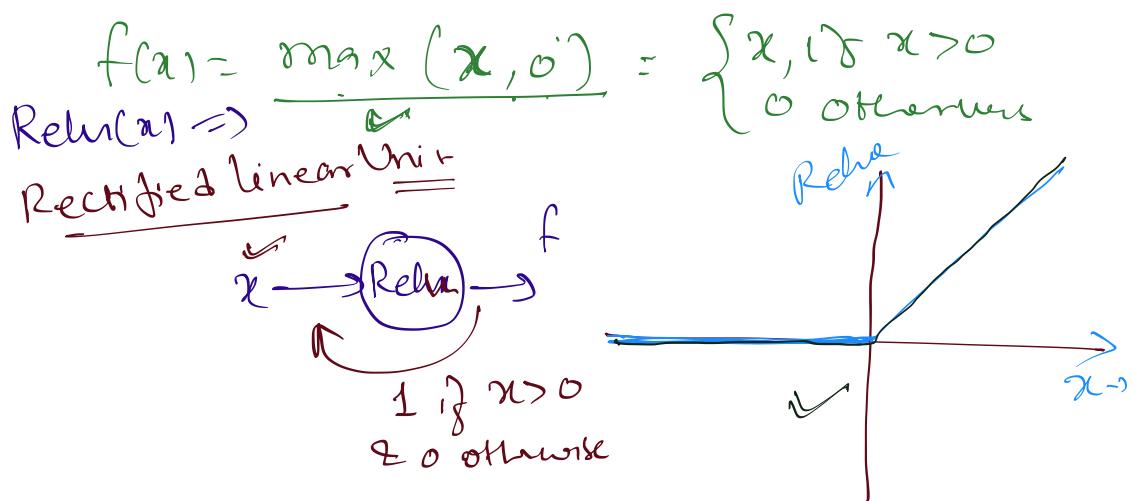
$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

$\frac{\partial f}{\partial x} \leftarrow 1 \text{ if } x > 0$

$f \rightarrow \frac{\partial f}{\partial x} = 1$

$$\overline{x} = \begin{cases} 0 & \text{otherwise} \end{cases}$$

$1 \text{ if } \ln p > 0$
otherwise '0'



$$\sigma(x) = \frac{1}{1+e^{-x}} ; \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Exercise: $\tanh(x) = \frac{1}{1+\sigma(x)}$

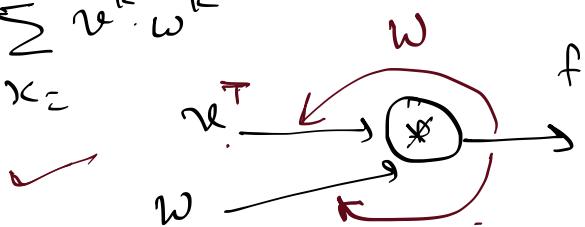
Vectors. v, w $v \in \mathbb{R}^d, w \in \mathbb{R}^d$

$$f(v, w) = v^T \cdot w$$

$$[v^T] [w]$$

$$= v^1 \cdot w^1 + v^2 \cdot w^2 + \dots + v^d \cdot w^d$$

$$= \sum_{k=1}^d v^k \cdot w^k$$

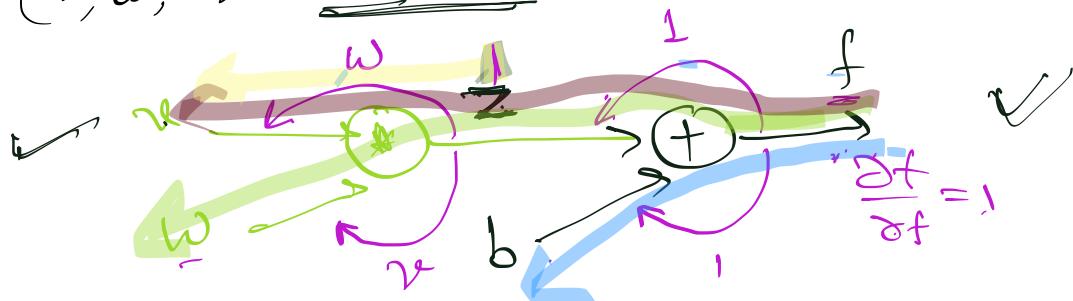


v^\top

$$v, w \in \mathbb{R}^d$$
$$\underline{b \in \mathbb{R}^l}$$

$$v = \begin{bmatrix} v^1 \\ v^2 \\ \vdots \\ v^d \end{bmatrix}; w = \begin{bmatrix} w^1 \\ w^2 \\ \vdots \\ w^d \end{bmatrix}$$

$$f(v, w, b) = \underline{\underline{v^\top w + b}}$$



$$\left. \begin{array}{l} \frac{\partial f}{\partial w} = 1 \times 1 \times v = v \\ \frac{\partial f}{\partial v} = 1 \times 1 \times w = w \\ \frac{\partial f}{\partial b} = 1 \times 1 = 1 \end{array} \right\} \quad \frac{\partial \vec{z}}{\partial \vec{x}} = \vec{w}$$

$$\begin{aligned} \frac{\partial z(v^T w + b)}{\partial v} &= \underbrace{\frac{\partial f}{\partial f} \cdot \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial w}}_{\text{Chain Rule}} = 1 \times 1 \times v = \underline{v+b} \\ \frac{\partial (v^T w + b)}{\partial v} &= \underbrace{\frac{\partial f}{\partial f} \cdot \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial v}}_{\text{Chain Rule}} = 1 \times 1 \times w \quad \checkmark \\ \frac{\partial (v^T w + b)}{\partial b} &= \underbrace{\frac{\partial f}{\partial f} \cdot \frac{\partial b}{\partial b}}_{\text{Chain Rule}} = 1 \end{aligned}$$

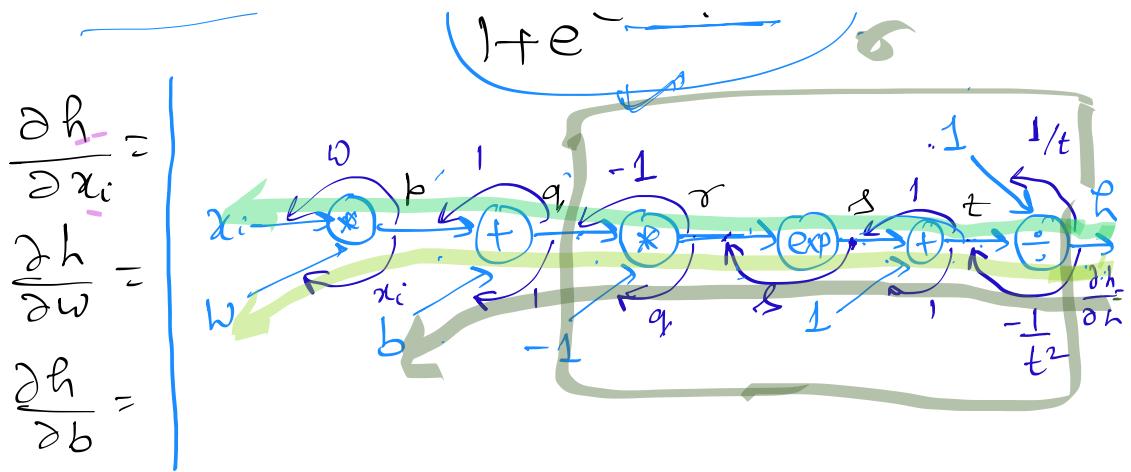
$$x_i \in \mathbb{R}^d; \quad w \in \mathbb{R}^d, \quad b \in \mathbb{R}$$

$$h(w, x_i, b) = \underbrace{w^T x_i + b}_{\text{circled}}$$

$\frac{\partial h}{\partial w} = ? \quad 1 \times 1 \times x_i \quad x_i$
 $\frac{\partial h}{\partial x_i} = ? \quad 1 \times 1 \times w \quad w$
 $\frac{\partial h}{\partial b} = ? \quad 1 \times 1 = 1$

$$h(w, x_i, b) = \frac{1}{-(w^T x_i + b)}$$

$\frac{\partial e^{-x_i}}{\partial x_i} = \underline{e^{-x_i}}$



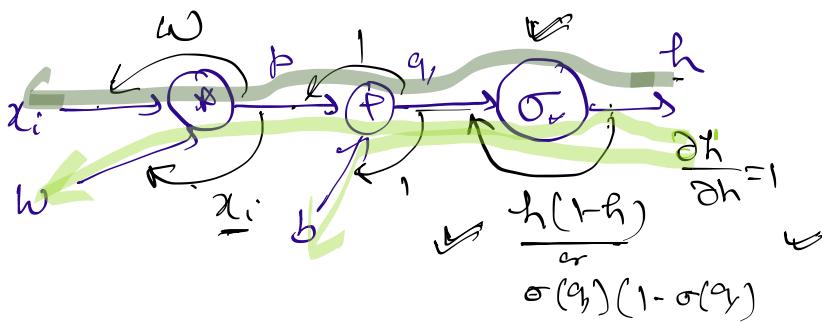
$$\frac{\partial h}{\partial x_i} = 1 * \left(-\frac{1}{t^2}\right) \cdot 1 \times 8 \times (-1) \times 1 \times w = \frac{8 \times w}{t^2}$$

$$\frac{\partial h}{\partial w} = 1 \times \left(-\frac{1}{t^2}\right) \alpha 1 \times 8 \times (-1) \times 1 \times x_i = \frac{8 \times x_i}{t^2}$$

$$\frac{\partial h}{\partial b} = 1 \times \left(-\frac{1}{t^2}\right) \alpha 1 \times 8 \times (-1) \times 1 = \frac{8}{t^2}$$

$$\rightarrow h(\omega, x_i^o, b) = \underline{\sigma} \left(\underline{\underline{\omega^T x_i^o + b}} \right) \quad | \sigma(a) = \frac{1}{1+e^{-a}}$$

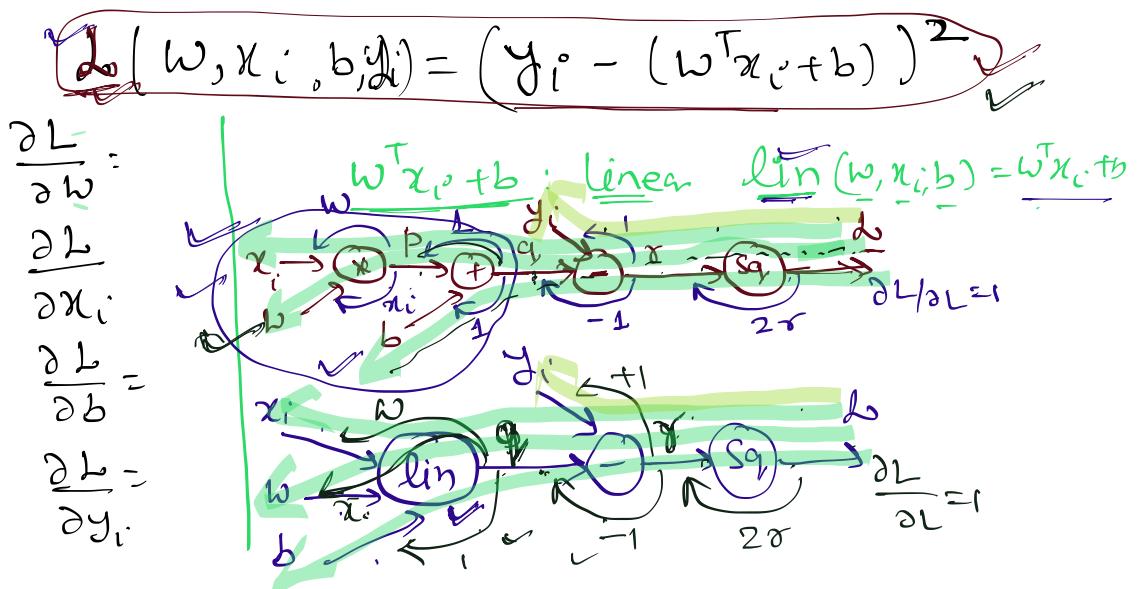
$$\sigma'(a) = \underline{\sigma(a)(1-\sigma(a))}$$



$$\frac{\partial h}{\partial x_i} = 1 \times \sigma(a_1) \cdot (1-\sigma(a_1)) \times 1 \times \omega \quad \}$$

$$\frac{\partial h}{\partial \omega} = 1 \times \sigma(a_1) \cdot (1-\sigma(a_1)) \times 1 \times x_i$$

$$\frac{\partial L}{\partial b} = 1 \times \sigma(g_i) \cdot (1 - \sigma(g_i)) \times 1$$



$$\frac{\partial L}{\partial w} = 1 \times 2r \times (-1) \times \underbrace{(+) \times x_i}_\text{2nd}$$

$$\frac{\partial L}{\partial x_i} = 1 \times 2r \times (-1) \times 1 \times w$$

$$\frac{\partial L}{\partial b} = 1 \times 2r \times (-1) \times 1$$

$$\frac{\partial L}{\partial y_0} = 1 \times 2r \times 1$$

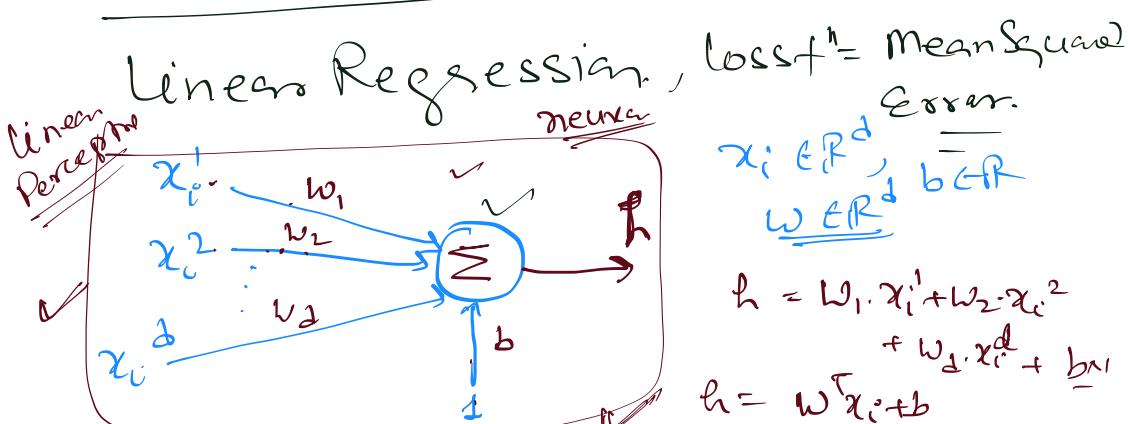
$$\frac{\partial L}{\partial w} = 1 \times 2r \times (-1) \times \underbrace{x_i}_\text{2nd}$$

$$1 \times 2r \times (-1) \times w$$

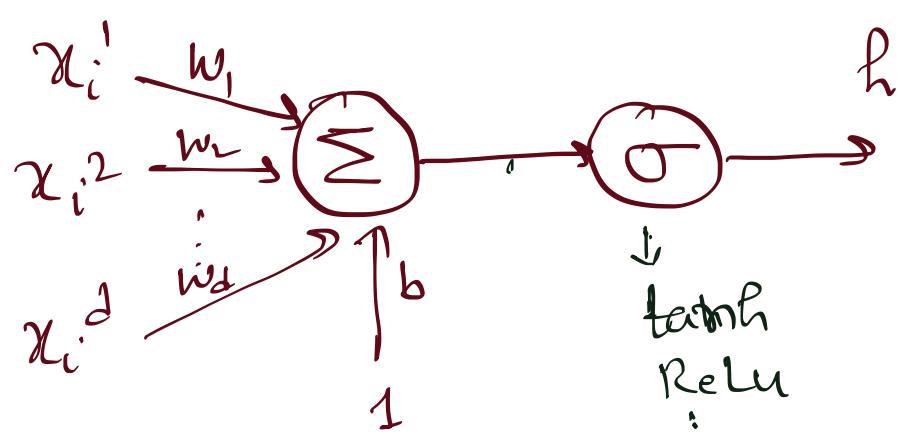
$$1 \times 2r \times (-1) \times 1$$

$$\frac{\partial L}{\partial b} = 1 \times 2r \times 1$$

$$h(w, x_i, b) = \underline{w^T x_i + b}$$



Logistic Regression



$$L(w, x_i, b, y_i) = -y_i \log(\sigma(w^T x_i + b))$$

Binary Cross Entropy

$$= -(1-y_i) \cdot \log(1-\sigma(w^T x_i + b))$$

$$\frac{\partial L}{\partial w} = \checkmark$$

Exercise

$$\frac{\partial L}{\partial x_i} = \checkmark$$

$$\frac{\partial L}{\partial b} = \checkmark$$

$$\frac{\partial L}{\partial y_i} = \checkmark$$

