

Spring 2021

CSE 465.1



AI

- Lecture 1

- Course Introduction
- Projects
- Questions



Course Objective and Outcome Form

Department of Electrical and Computer Engineering

School of Engineering and Physical Sciences

North South University, Bashundhara, Dhaka-1229, Bangladesh

1. Course Number and Title: CSE 465 Pattern Recognition

2. Number of Credits: 03

3. Type: Core

4. Prerequisites: CSE 373

5. Faculty Name: Dr. Nabeel Mohammed (NbM)

6. Room: SAC 917

7. Office Hours: TBA (Lets schedule as needed)

8. Email: nabeel.mohammed@northsouth.edu + Google Classroom

9. Contact Hours: Lectures – 3 Hours/week

10. Course Summary:

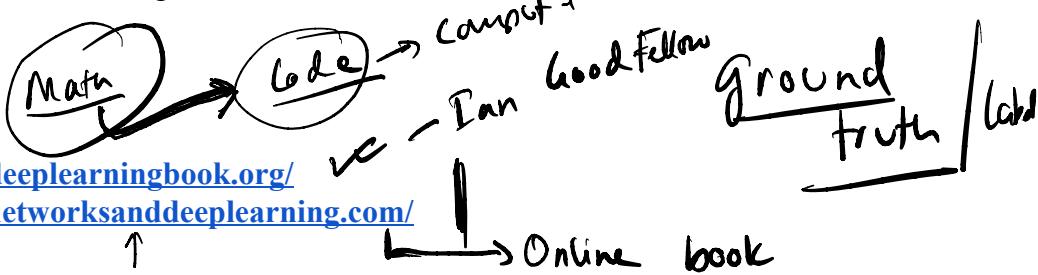
The course will cover topics such as supervised learning (unsupervised learning) mainly in the context of neural networks. We will cover techniques, mainly popularised in the last few years, for neural network training and aim to give students an appreciation for the machine learning process involving this class of models.

11. Resources

Text books:

- <https://www.deeplearningbook.org/>
- <http://neuralnetworksanddeeplearning.com/>

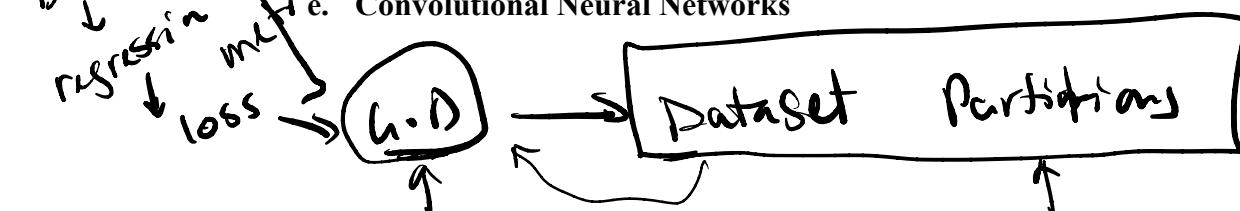
(You will be expected to read a lot of academic papers)



12. Topics:

- a. Supervised/Unsupervised Learning, Data, Dataset partitions, metrics
- b. Regression, loss function, gradient descent
- c. Classification, entropy, KL divergence, cross entropy
- d. Neural Networks
- e. Convolutional Neural Networks

K-fold



- f. Recurrent Neural Networks
- g. Autoencoders
- h. Generative Adversarial Networks
- i. Siamese networks
- j. Adversarial Images
- k. Many other possible topics if time permits

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} ae + bg & af + bh \\ ce + dg & cf + dh \end{bmatrix}$$

Chain Rule Linear Algebra

Calculus
Discipline, Permutation

- Vectors
 - Vector spaces
 - Basis vectors
 - Matrix multiplication
- ↓
What

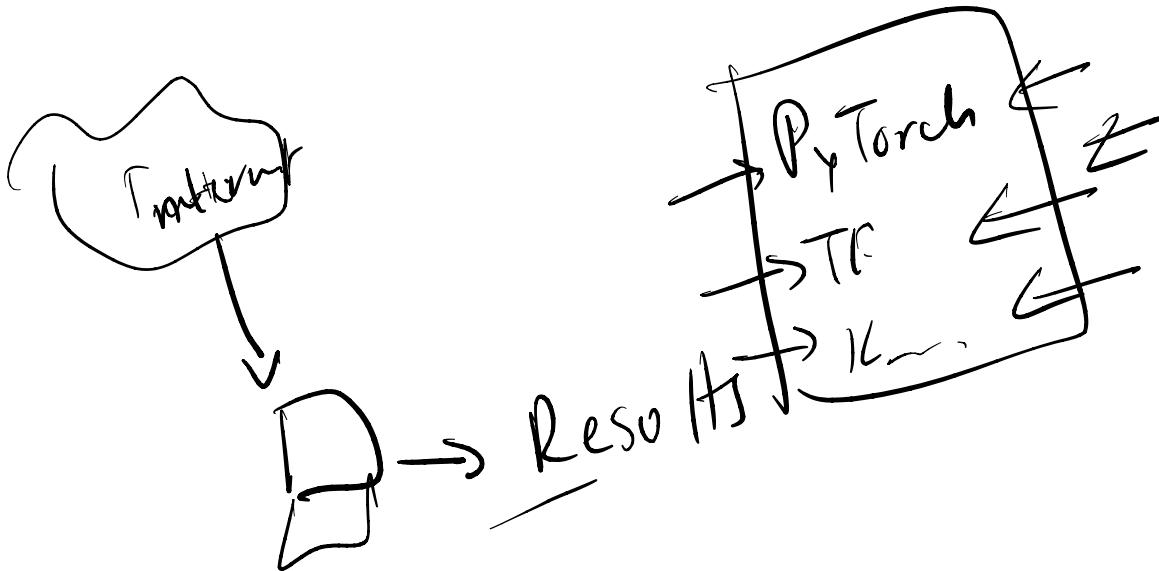
13. Weightage Distribution among Assessment Tools

Assessment Tools	Weightage (%)
Project Work	40
Literature Review of Project Topic (Due week 4)	10
Final Exam	25
Take home assessment (Mid)	25
	100

14. Grading policy: As per NSU grading policy available in
<http://www.northsouth.edu/academic/grading-policy.html>

15. Course Policies:

- a. Be good and follow university policies.



Project Ideas :

- Student should be able to frame a formal Machine Learning problem
- Student should be able to discuss real world & arrange data for model training,
data validation & testing.
} ↑
- Choosing appropriate model, loss & training, validation, testing strategies for model deployment.

→ Spoken Recognition Ø
→ Face Anti-Spoofing Ø

BANNER: A Cost-Sensitive Contextualized Model for Bangla Named Entity Recognition

IMRANUL ASHRAFI[✉], MUNTASIR MOHAMMAD[✉], ARANI SHAWKAT MAUREE[✉],
GALIB MD. AZRAF NIJHUM[✉], REDWANUL KARIM[✉], NABEEL MOHAMMED[✉],
AND SIFAT MOMEN[✉]

Department of Electrical and Computer Engineering, North South University, Dhaka 1229, Bangladesh

Corresponding author: Nabeel Mohammed (nabeel.mohammed@northsouth.edu)

This work was supported by the Special Grant of the Ministry of ICT, Bangladesh, under Grant 56.00.0000.028.20.004.19-248.

ABSTRACT Named Entity Recognition (NER) is a task in Natural Language Processing (NLP) that aims to classify words into a predetermined list of Named Entities (NE). Many architectures have produced good results on high resourced languages like English and Chinese. However, the NER task has not yet achieved much progress for Bangla, a low resource language. In this paper, we perform the NER task on Bangla Language using Word2Vec and contextual Bidirectional Encoder Representations from Transformers (BERT) embeddings. We propose multiple BERT-based deep learning models that use the contextualized embedding from BERT as inputs and a simple statistical approach for class weight cost sensitive learning. The modified cost-sensitive loss function was used to address the class imbalance of the data. In our modified cost-sensitive loss function, we penalize the dominant classes by taking the ratio concerning the maximum sample in a class instead of the whole dataset. This penalty is made so that the learner learns slowly for the dominant class. In addition, we experiment by adding a Conditional Random Field (CRF) layer and incorporating Focal Loss to the training process. We found the best F1 Macro score to be 65.96%, F1 Micro score of 90.64%, and F1 Message Understanding Coreference (MUC) score of 72.04%, which were calculated at Named Entity level. Our experimental results demonstrate that one of the proposed models, which jointly optimizes for the CRF loss and class weighted cost-sensitive loss according to our proposed statistical approach, achieve an improvement of over 8% F1 MUC score on a recently introduced Bangla NER dataset when compared to previously published work.

INDEX TERMS Bangla NER, Bengali NER, BERT, CRF, NLP, focal loss, cost-sensitive learning, contextual embeddings.

I. INTRODUCTION

Named Entity Recognition (NER) is a task in Natural Language Processing (NLP) which seeks to classify words in an unstructured text into particular categories of Named Entities (NE), such as - person, organization, location etc. [18]. NER models can be used in multiple downstream tasks like identifying names of genes and gene products [43], recognition of chemical entities [50] and chatbots [40].

Bangla is a globally spoken language. Hence it is important that globally spoken languages like Bangla are enriched in linguistic knowledge and vocabulary. But current state of Bangla NLP has not achieved much progress in tasks like NER. This is because, unlike English, Bangla is more

complex regarding both usability and vocabulary. Since Bangla is a rich language, many linguistic challenges occur while training models on it.

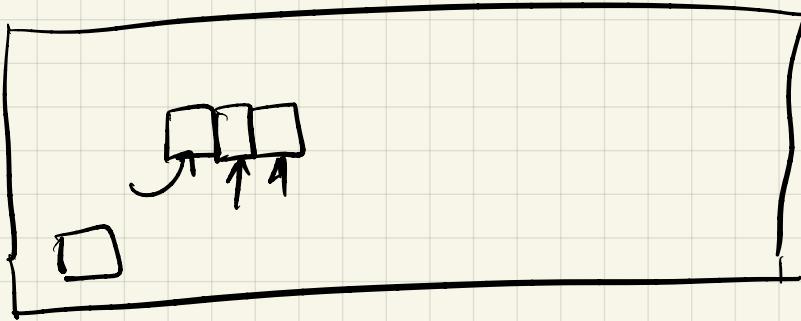
Figure 1 provides examples which illustrate the real challenges of Bangla language more precisely. Unlike English, there are minimal indicators for tags like capitalization in Bangla as stated by [17], [42]. For instance, under *Capitalization* in the figure, the words '*chottogram*' and '*kushtia*' refer to *Chittagong* and *Kushtia*, which are names of places. From the English words, it is clearly understood that *Chittagong* and *Kushtia* are capitalized and names (of locations). Hence, it is possible that they can be tagged as *Location* through NER models trained in English Language. On the other hand, observing the Bangla sentences, it is not clear as to which parts of speech do each word signify simply by looking at the words appearance without prior knowledge of Bangla.

The associate editor coordinating the review of this manuscript and approving it for publication was Mauro Tucci[✉].

ML Engineer → X
↓

S. Engineer (Must be able to do ML)

- Project Ideas
 - Keyboard layout optimization
 - Speaker Recognition
 - Shaka AI ✓
 - Bangla Named Entity Recognition ✓
 - Violence Detection ✓
 - Bangla H. writing recognition
 - Machine Translation → Bangla



QWERTY

(metric)

!!!

Spell Corrections
↓
Self Adaptive
↓
Domain specific

Spell Corrections
↓
NN.

NN

Optimization

0 0 0 0

→ Statistically Learnt

Lecture 2

$$f(x) = x^2 + 2 \quad \underline{\text{is } ax^2 + bx + c, \quad a=1, b=0, c=2}$$

$$f(1) = 1^2 + 2 = 3 \rightarrow (1, 3)$$

Don't know the form of f

We don't know the parameters. $\Theta = ()$

x	1	2	3	4	5
y	2	5	5	10	12

↓ ↓

$\{(1, 2), (2, 5), (3, 5), (4, 10), (5, 12)\}$

$$y = f(x) + n$$

$$\underset{\text{Independent}}{\hat{x}} = \frac{\hat{f}(x)}{\text{Estimated function}}$$

$$\hat{x} = g(x)$$

$$Y = g(x)$$

$$\hat{g}$$

image face

fingerprint

x	1	2	3	4	5	
y	2	5	5	10	12	

$$f(x) = \boxed{ax + b}$$

$$\theta = (a, b)$$

(VN)

Algorithm:

Purpose of this algorithm:

close

x	1	2	3	4	5	
y	2	5	5	10	12	.

$f(x) = x$	1	2	3	4	5	
$\theta = (1, 0)$	1	3	2	6	7	$\sum = 19$

$f(x) = x+1$	2	3	4	5	6	
$\theta = (1, 1)$	0	2	1	5	6	$\sum = 14$

$f(x) = 3x$	3	6	9	12	15	
$\theta = (3, 0)$	-1	-1	-4	-2	-3	$\sum = -11$

$f(x) = 2x$	2	4	6	8	10	
$\theta = (2, 0)$	0	1	-1	2	2	$\sum = 4$

$f(x) =$ magic(x)	2	5	5	10	12	
$\theta = (H, P)$	0	0	0	0	0	$\sum = 0$

close
How close
measured
at

$\frac{1}{(2-1)^2} + \frac{9}{(5-2)^2} + \frac{4}{(5-3)^2} + \frac{36}{(10-4)^2} + \frac{49}{(12-5)^2} = 99$

x	1	2	3	4	5	
y	2	5	5	10	12	.
$f(x) = x$ $\theta = (1, 0)$	1	2	3	4	5	$19(A), 99(S)$
$f(x) = x+1$ $\theta = (1, 1)$	2	3	4	5	6	$14(A), 66(S)$
$f(x) = 3x$ $\theta = (3, 0)$	3	6	9	12	15	$71(A), 31(S)$
$f(x) = 2x$ $\theta = (2, 0)$	2	4	6	8	10	$6(A), 10(S)$
$f(x) = \max(x, 0)$ $\theta = (HP)$	2	3	5	10	12	$0(A), 0(S)$

min distance = 0

possible

total distance = function individual distances

individual distances ≥ 0

$$\sum_{i=1}^n \text{individual distance}$$

⇒ Total Distance = $\sum_{i=1}^n g(y_i, f(x_i))$

g is a distance function

$$g(a, b) = |a - b| \rightarrow \begin{cases} \text{Absolute error} \\ \text{error} \end{cases}$$

$$g(a, b) = (a - b)^2 \rightarrow \begin{cases} \text{Squared error} \\ \text{error} \end{cases}$$

$y, f(x)$

An initial form of Gradient Descent

initialize a to some random value

while (not happy with error)

{

calculate E .

calculate $\frac{dE}{da}$

$$\text{Step} = \begin{cases} -1, & \frac{dE}{da} > 0 \\ 1, & \frac{dE}{da} < 0 \\ 0, & \frac{dE}{da} = 0 \end{cases}$$

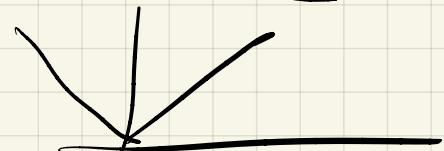
$a = a + \text{Step}$

Data partitions



Check happiness with new error.

Distance function

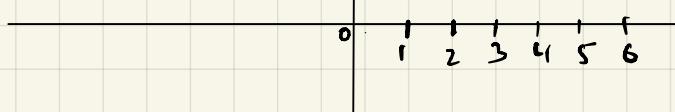


$$f(x) = ax$$

~~Distance / Error~~

$$\theta(a)$$

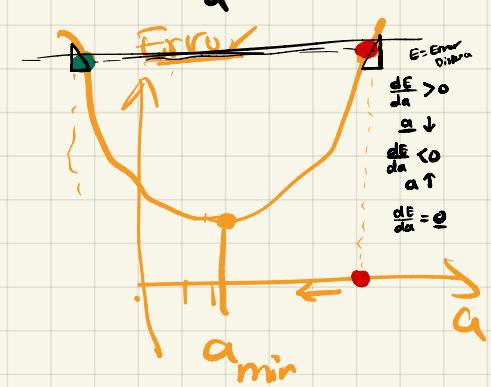
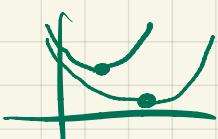
$$|y_i - f(x_i)|$$



a^*

Math

Computer \Leftarrow Reality

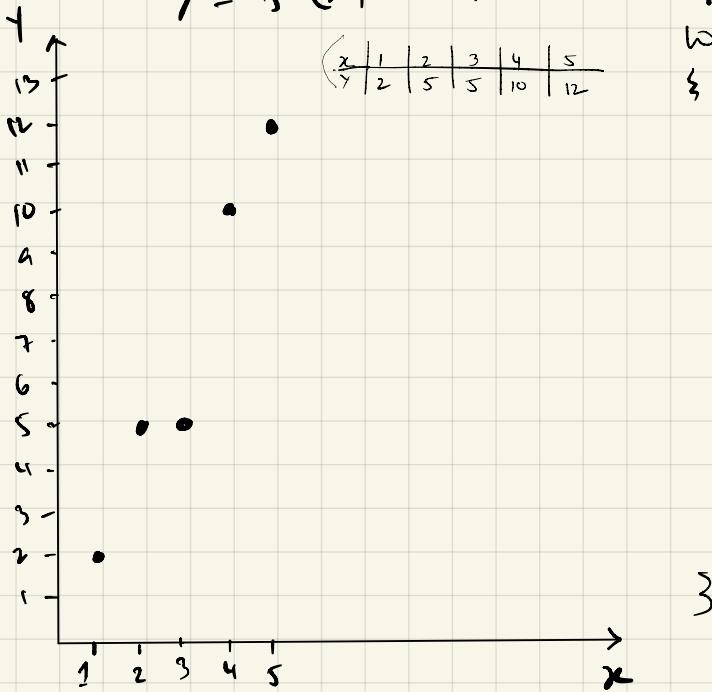


Lecture 3

x	1	2	3	4	5
y	2	5	5	10	12

$$\hat{y} = \hat{f}(x) = ax$$

$$\begin{array}{|c|c|c|c|c|c|} \hline x & 1 & 2 & 3 & 4 & 5 \\ \hline y & 2 & 5 & 5 & 10 & 12 \\ \hline \end{array}$$



initialize a randomly
while (not happy with error)
 {

 calculate error $E \rightarrow \sum_{i=1}^5 (f(x_i) - y_i)^2$
 calculate $\frac{dE}{da}$

 step = $\begin{cases} -1, & \frac{dE}{da} > 0 \\ 1, & \frac{dE}{da} < 0 \\ 0, & \frac{dE}{da} = 0 \end{cases}$
 $a = a + \text{step}$

 check happiness with new error

Shayne Paben

$$a = 5$$

while (not happy error)

$$E =$$

$$\text{Error} = (ax_1 - y_1)^2 +$$

$$(ax_2 - y_2)^2 + (ax_3 - y_3)^2 + (ax_4 - y_4)^2 + (ax_5 - y_5)^2 = \sum_{i=1}^5 (\hat{f}(x_i) - y_i)^2$$

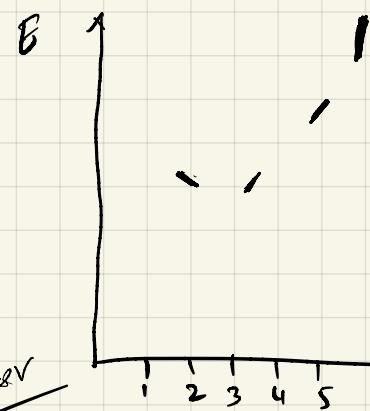
$$(ax_1 - y_1)^2 + (ax_2 - y_2)^2 + (ax_3 - y_3)^2 + (ax_4 - y_4)^2 + (ax_5 - y_5)^2 = \sum_{i=1}^5 (ax_i - y_i)^2$$

a	E	dE/da
5	403	296
4	162	186
3	31	76
2	10	-34

$$\frac{dE}{da} = 2(ax_1 - y_1) \cdot x_1 +$$

$$2(ax_2 - y_2) \cdot x_2 + 2(ax_3 - y_3) \cdot x_3 + 2(ax_4 - y_4) \cdot x_4 + 2(ax_5 - y_5) \cdot x_5 = \sum_{i=1}^5 (ax_i - y_i) \cdot x_i$$

a	E	dE/da
5	403	296
4	162	186
3	31	76
2	10	-34



$$f(x) = ax^2 + bx + c$$

constant

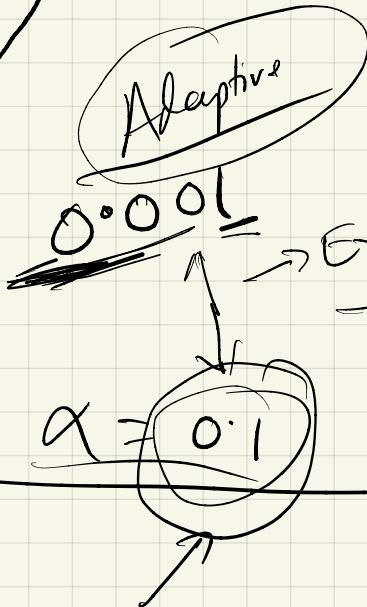
initialize a randomly

calculate $E = \frac{1}{5} \sum_{i=1}^5 (ax_i - y_i)^2$
While (not happy with E)

calculate $\frac{dE}{da}$ learning rate

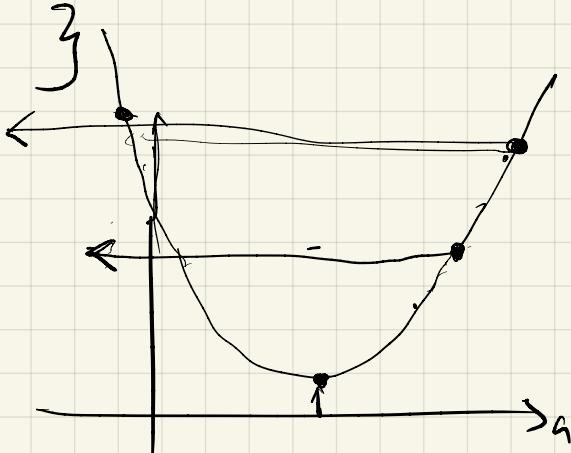
$$\text{Step} = \alpha * \frac{dE}{da} \quad \downarrow \text{Hyper parameter}$$

$$a = a - \text{Step}$$



calculate $E = \frac{1}{5} \sum_{i=1}^5 (ax_i - y_i)^2$

x	1	2	3	4	5
y	2	5	5	10	12



a	E	dE/da	Step
5	80.6	59.2	5.92
-0.92	115.64	-71.04	-7.104
6.184			0.1

X	$f(x)$	$(y - f(x))^2$
1	-0.92	$(2 - (-0.92))^2 = 2.92^2 = 8.52$
2	-1.89	$(5 - (-1.89))^2 = 6.89^2 = 46.78$
3	-2.76	$(5 - (-2.76))^2 = 7.76^2 = 60.27$
4	-3.68	$(10 - (-3.68))^2 = 13.68^2 = 187$
5	-4.6	$(12 - (-4.6))^2 = 275.56$

$$\frac{C}{da} =$$

$$\frac{578.23}{5}$$

Lecture 4

\$\$\$

Generalizes

x	1	2	3	4	5	6	7	8
y	2	5	5	10	12	19	25.5	33
2.312	2.31	4.62	6.93	9.24	11.55	13.86	16.17	18.48

$$f(x) = ax$$

Validation

a E	E
5 80.6	86.75
4 32.4	10.75
3 6.2	34.08
2 2	15.6
= 2.5 1.35	8.3
2.25 0.9875	11.6
2.31 0.9491	10.8

0-1
L2
10
10

100's

test =

x	9	10	11
y	42	50	63
42	35	40	44

Error = 165.67

Lecture 5

$y = f(x) + \eta$, $\hat{y} = \hat{f}(x)$, \hat{f} is parameterized by θ
 $y, x, \hat{y} \rightarrow \underline{\text{scalars}}$

y, x, \hat{y} (can be) is vectors.

$$\hat{Y} = ax, \quad x \in \mathbb{R}^3, \quad Y \in \mathbb{R}^4$$

$a \in \mathbb{R}^{4 \times 3}$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ \hat{y}_4 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \\ a_{41} & a_{42} & a_{43} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$\hat{y}_i = a_{i1}x_1 + a_{i2}x_2 + a_{i3}x_3$$

$$\hat{Y}_i = a_{i1}x_1 + a_{i2}x_2 + a_{i3}x_3$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \quad \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

$$= \frac{1}{4} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + (\hat{y}_3 - y_3)^2 + (\hat{y}_4 - y_4)^2$$

For 1 training

Training data

Input X	\vec{x}_1	\vec{x}_2	\vec{x}_3	...	\vec{x}_K
Actual Y	\vec{y}_1	\vec{y}_2	\vec{y}_3	...	\vec{y}_K
Pred \hat{Y}	$\hat{\vec{y}}_1$	$\hat{\vec{y}}_2$	$\hat{\vec{y}}_3$...	$\hat{\vec{y}}_K$

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

$$\hat{\vec{x}}_i = \begin{bmatrix} x_1^i \\ x_2^i \\ x_3^i \end{bmatrix} \Leftarrow$$

$$+ (\hat{y}_4^i - y_4^i)^2$$

$$\hat{\vec{y}}_i = \begin{bmatrix} \hat{y}_1^i \\ \hat{y}_2^i \\ \hat{y}_3^i \\ \hat{y}_4^i \end{bmatrix} \Rightarrow$$

$$\begin{bmatrix} \hat{y}_1^i \\ \hat{y}_2^i \\ \hat{y}_3^i \\ \hat{y}_4^i \end{bmatrix} = \begin{bmatrix} \hat{\vec{y}}_1^i \\ \hat{\vec{y}}_2^i \\ \hat{\vec{y}}_3^i \\ \hat{\vec{y}}_4^i \end{bmatrix} = \begin{bmatrix} \hat{\vec{x}}_1^i \\ \hat{\vec{x}}_2^i \\ \hat{\vec{x}}_3^i \\ \hat{\vec{x}}_4^i \end{bmatrix}$$

$$\frac{1}{K} \sum_{j=1}^K \frac{1}{4} [(\hat{y}_1^j - y_1^j)^2 + (\hat{y}_2^j - y_2^j)^2 + (\hat{y}_3^j - y_3^j)^2 + (\hat{y}_4^j - y_4^j)^2]$$

$$= \frac{1}{DK} \sum_{j=1}^K \sum_{d=1}^D (\hat{y}_d^j - y_d^j)^2 \quad \text{where } \hat{y}_d^j, y_d^j \in \mathbb{R}^D$$

$$= \frac{1}{DK} \sum_{j=1}^K \sum_{d=1}^D (\hat{y}_d^j - y_d^j)^2 \quad \text{when } \hat{y}_d^j, y_d^j \in \mathbb{R}^P$$

$$= \frac{1}{DK} \sum_{j=1}^K m(\hat{f}(x_j), y_j)$$

$$= \frac{1}{DK} \sum_{j=1}^K m(Ax_j, y_j)$$

$$E = \frac{1}{DK} \sum_{j=1}^K \sum_{d=1}^D (a_{d1}x_1^j + a_{d2}x_2^j + a_{d3}x_3^j - y_d^j)^2$$

$$\hat{y}^j =$$

initialize a Random A
calculate E

while (not happy)
{

for each a_{pq} ($1 \leq p \leq 4, 1 \leq q \leq 3$)
Step = $\alpha \cdot \frac{dE}{da_{pq}}$

$$a_{pq} = a_{pq} - \text{Step}$$

calculate E

$$\vec{x} \in \mathbb{R}^3$$

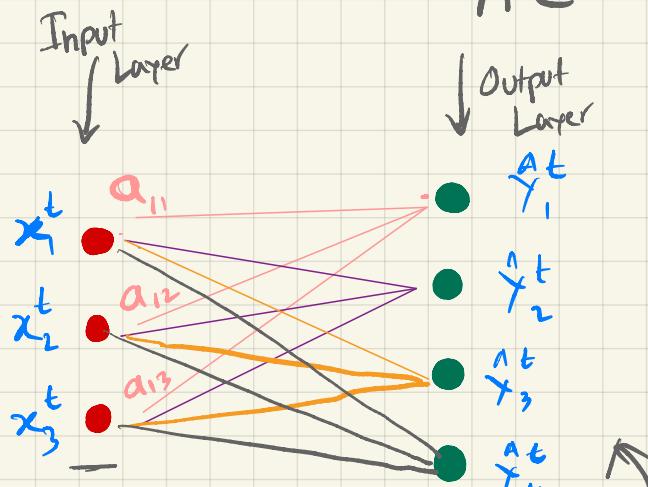
$$C\vec{x} \in \mathbb{R}^{100}, C \in \mathbb{R}^{100 \times 3}$$

$$\vec{y} \in \mathbb{R}^4, B \in \mathbb{R}^{4 \times 100}$$

$$\begin{bmatrix} 4 \times 100 \times 3 \\ 4 \times 100 \\ 4 \times 3 \end{bmatrix}$$

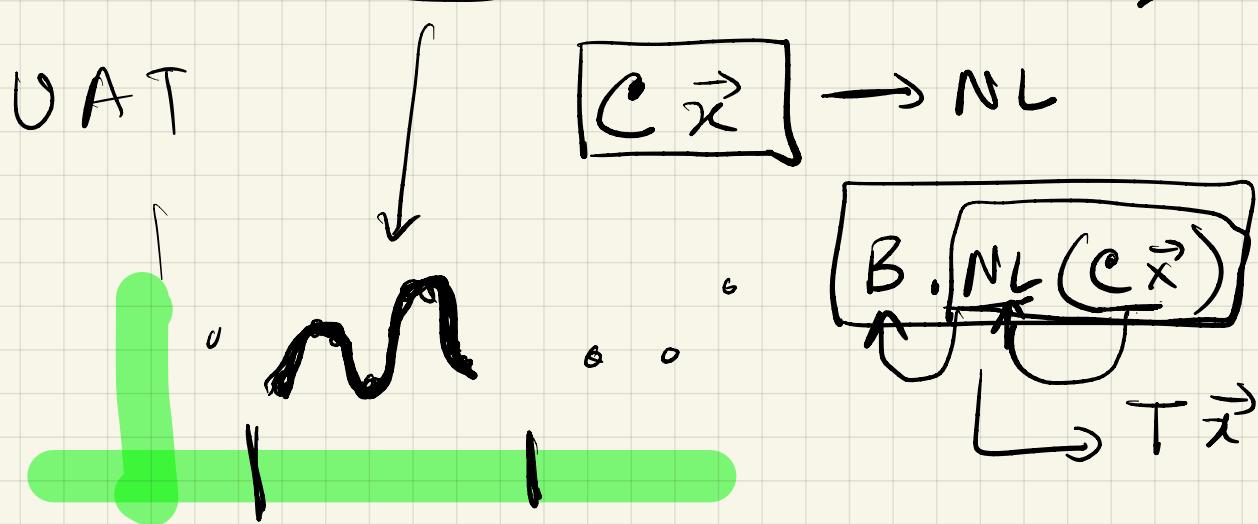
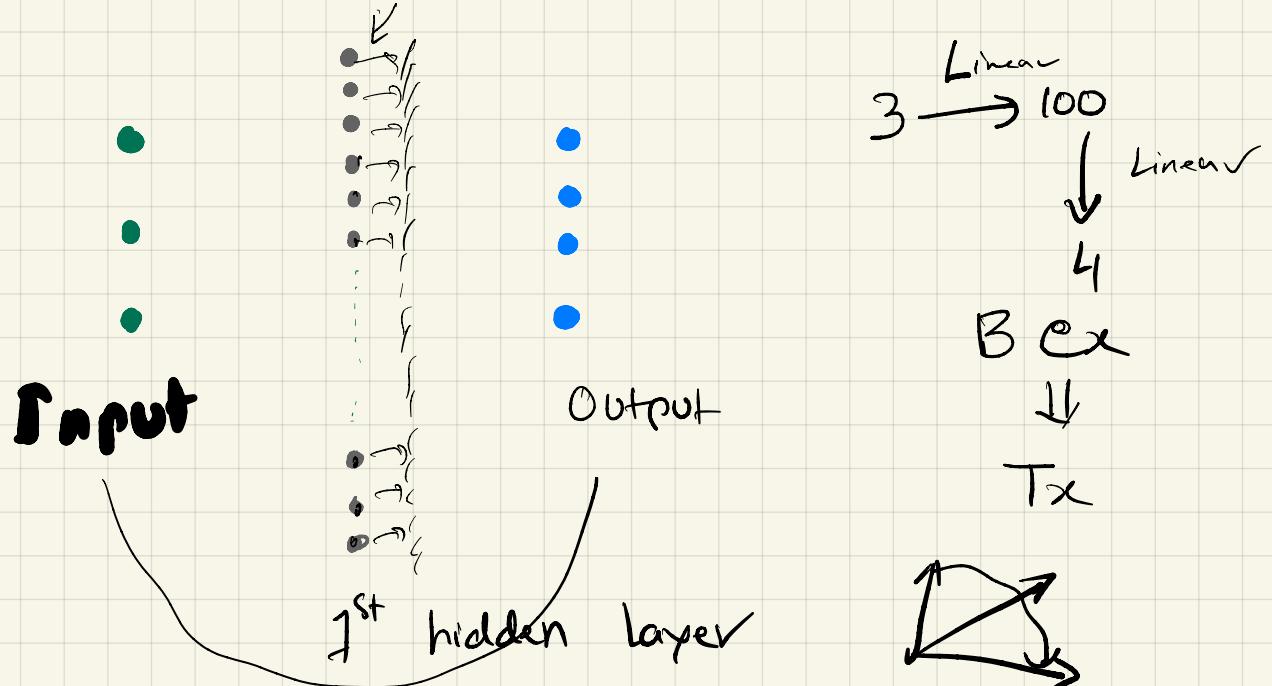
$$\begin{aligned} \vec{y} &= A\vec{x} \\ \vec{x} &= A^{-1}\vec{y} \\ \vec{x} &= A^{-1}B\vec{y} \end{aligned}$$

$$A \in \mathbb{R}^{100 \times 3}$$



Neural Network

$$\begin{aligned} \vec{y} &= A\vec{x} \\ &= B(C\vec{x}) \end{aligned}$$



E must be differentiable

$$\hat{Y} = \frac{B(c \vec{x})^2}{\text{or } B(c \vec{x})^3} \xrightarrow{\substack{\text{Elementwise} \\ \text{Square}}} \text{UAT}$$

or $B \cdot \log(c \vec{x})$

Back Propagation

$$\hat{Y} = A \cdot (B \vec{x})^2$$

$X \in \mathbb{R}^3$

$\vec{x} \in \mathbb{R}^4$

$B \in \mathbb{R}^{5 \times 3}$

$A \in \mathbb{R}^{4 \times 5}$

Lecture 6

Back propagation

Context

Init Θ randomly
etc calculate E

Θ is a ~~large~~ vector
of all parameters

while (not happy)
|
Step = $\alpha \frac{dE}{d\Theta}$ → partial derivatives.
 $\Theta = \Theta - \text{Step}$
Calculate E

$$\begin{matrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{matrix} \quad \begin{matrix} x \in \mathbb{R}^3 \\ H_1 = \mathbb{R}^4 \\ \hat{x} = \mathbb{R}^2 \end{matrix} \quad E = \frac{1}{D} \sum_i (y_i - \hat{x}_i)^2$$

$$H_1 = \underbrace{\text{Sq}(Ax)}_{\rightarrow \mathbb{R}^4}$$

Sq - is elementwise square.

$$\hat{x} = \underbrace{B H_1}_{\rightarrow \mathbb{R}^2}$$

$$A \in \mathbb{R}^{4 \times 3}$$

$$B \in \mathbb{R}^{2 \times 4}$$

$$H_1 = \text{Sq} \left(\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \\ a_{41} & a_{42} & a_{43} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \right)$$

$$= \begin{bmatrix} (a_{11} \cdot x_1 + a_{12} \cdot x_2 + a_{13} \cdot x_3)^2 \\ (a_{21} \cdot x_1 + a_{22} \cdot x_2 + a_{23} \cdot x_3)^2 \\ (a_{31} \cdot x_1 + a_{32} \cdot x_2 + a_{33} \cdot x_3)^2 \\ (a_{41} \cdot x_1 + a_{42} \cdot x_2 + a_{43} \cdot x_3)^2 \end{bmatrix}$$

$$H_1 = \left[\begin{array}{c} (a_{11}x_1 + a_{12}x_2 + a_{13}x_3)^2 \\ (a_{21}x_1 + a_{22}x_2 + a_{23}x_3)^2 \\ (a_{31}x_1 + a_{32}x_2 + a_{33}x_3)^2 \\ (a_{41}x_1 + a_{42}x_2 + a_{43}x_3)^2 \end{array} \right]$$

$$\hat{Y} = \left[\begin{array}{c} b_{11} & b_{12} & b_{13} & b_{14} \\ b_{21} & b_{22} & b_{23} & b_{24} \end{array} \right] \quad \left[\begin{array}{c} (a_{11}x_1 + a_{12}x_2 + a_{13}x_3)^2 \rightarrow h_{11} \\ (a_{21}x_1 + a_{22}x_2 + a_{23}x_3)^2 \rightarrow h_{12} \\ (a_{31}x_1 + a_{32}x_2 + a_{33}x_3)^2 \rightarrow h_{13} \\ (a_{41}x_1 + a_{42}x_2 + a_{43}x_3)^2 \rightarrow h_{14} \end{array} \right]$$

$$= \frac{b_{11} \cdot h_{11} + b_{12} \cdot h_{12} + b_{13} \cdot h_{13} + b_{14} \cdot h_{14}}{b_{21} \cdot h_{11} + b_{22} \cdot h_{12} + b_{23} \cdot h_{13} + b_{24} \cdot h_{14}} \rightarrow \hat{Y}_1$$

$$E = \frac{1}{2} \left((\hat{Y}_1 - Y_1)^2 + (\hat{Y}_2 - Y_2)^2 \right)$$

$$12 + 8 = 20$$

$$\frac{dE}{da_{11}}$$

Big
Annoying

$$E = k_{11}(a_{11})$$

$$E = k_{12}(a_{12})$$

$$E = \frac{1}{2}((\hat{x}_1 - x_1)^2 + (\hat{x}_2 - x_2)^2)$$

$$\begin{aligned}\frac{dE}{d\hat{x}_1} &= Q \downarrow \nabla(\hat{x}_1 - x_1) \\ &= (\hat{x}_1 - x_1)\end{aligned}$$

x, y

$$\frac{dE}{dx_1} = (\hat{x}_1 - x_1) \checkmark$$

$$\frac{dE}{dx_2} = (\hat{x}_2 - x_2) \checkmark$$

$$\frac{dE}{db_{11}} = \frac{dE}{d\hat{x}_1} \cdot \frac{d\hat{x}_1}{db_{11}}$$

$$\frac{dE}{db_{12}} = \frac{dE}{d\hat{x}_1} \cdot \frac{d\hat{x}_1}{db_{12}}$$

$$\frac{dE}{db_{13}} = \frac{dE}{d\hat{x}_1} \cdot \frac{d\hat{x}_1}{db_{13}}$$

$$\frac{dE}{db_{14}} = \frac{dE}{d\hat{x}_1} \cdot \frac{d\hat{x}_1}{db_{14}}$$

$$\hat{x} = \frac{B \operatorname{Sq}(Ax)}{I}$$

[$\begin{matrix} \hat{x}_1 \\ \hat{x}_2 \end{matrix}$]

$$\left[\frac{dE}{d\hat{x}_i} \right] \quad \left[\frac{dE}{d\hat{x}_l} \right]$$

$$\frac{dE}{dh_{11}} = \frac{dE}{dx_1} \cdot \frac{dx_1}{dh_{11}} + \frac{dE}{dx_2} \cdot \frac{dx_2}{dh_{11}}$$

$$= - \cdot b_{11} + - \cdot b_{21}$$

RNN

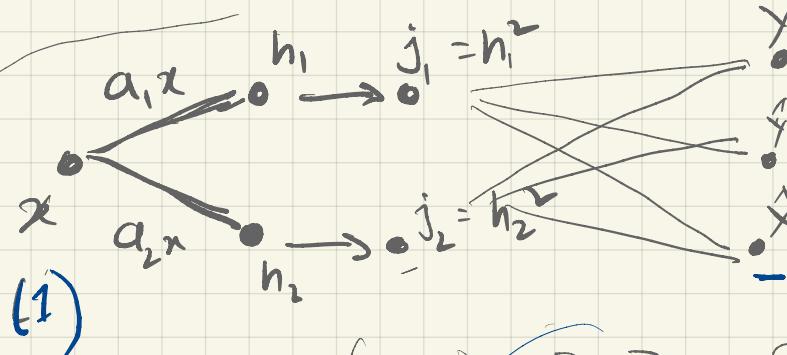
$$\frac{dE}{da_{11}} = \frac{dE}{dh_{11}} \cdot \frac{dh_{11}}{da_{11}}$$

$$h_{11} = (a_{11} \cdot x_1 + a_{12} \cdot x_2 + a_{13} \cdot x_3)$$

$$= \underbrace{[2(a_{11}x_1 + a_{12}x_2 + a_{13}x_3)]}_{x_1}$$

$$\frac{dE}{da_{11}} = \frac{dE}{dh_{11}} \cdot \frac{dh_{11}}{da_{11}}$$

$$\frac{dE}{dh_{11}} \cdot \frac{d\hat{x}_1}{dx_1}$$



$$(1) \quad \begin{aligned} & \frac{dE}{da_{11}} = \frac{dE}{dh_{11}} \cdot \frac{dh_{11}}{da_{11}} \\ & \frac{dE}{dh_{11}} = \frac{dE}{d(x_1 - \hat{x}_1)} \cdot \frac{d(x_1 - \hat{x}_1)}{dx_1} \\ & \frac{dE}{dx_1} = \frac{dE}{d(x_1 - \hat{x}_1)} \cdot \frac{d(x_1 - \hat{x}_1)}{dx_1} \end{aligned}$$

$$\begin{aligned} a_1 &= 22 \\ a_2 &= 66 \\ k_{11} &= 6 \\ k_{12} &= 5 \\ k_{21} &= 0.9 \\ k_{22} &= 11 \\ k_{31} &= 5.3 \\ k_{32} &= -4 \\ x_1 &= 1000 \\ x_2 &= -1000 \\ x_3 &= 0 \end{aligned}$$

$$\begin{aligned} H &= \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} x = \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} \\ J &= S_q(H) \\ \hat{Y} &= \begin{bmatrix} k_{11} & k_{12} \\ k_{21} & k_{22} \\ k_{31} & k_{32} \end{bmatrix} \begin{bmatrix} j_1 \\ j_2 \end{bmatrix} = \begin{bmatrix} \hat{x}_1 \\ \hat{x}_2 \\ \hat{x}_3 \end{bmatrix} \end{aligned}$$

$$\frac{dE}{da_{11}} = x = 1 \quad \begin{aligned} \frac{d\hat{x}_1}{dh_{11}} &= 2h_1 \\ &= 44 \end{aligned} \quad = 24684$$

$$\frac{dE}{da_{22}} = x = 1 \quad \begin{aligned} \frac{d\hat{x}_2}{dh_{22}} &= 2h_2 \\ &= 132 \end{aligned} \quad = -14858.8$$

$$R = q_1 + q_2 + q_3 \approx 2.4 \times 10^9$$

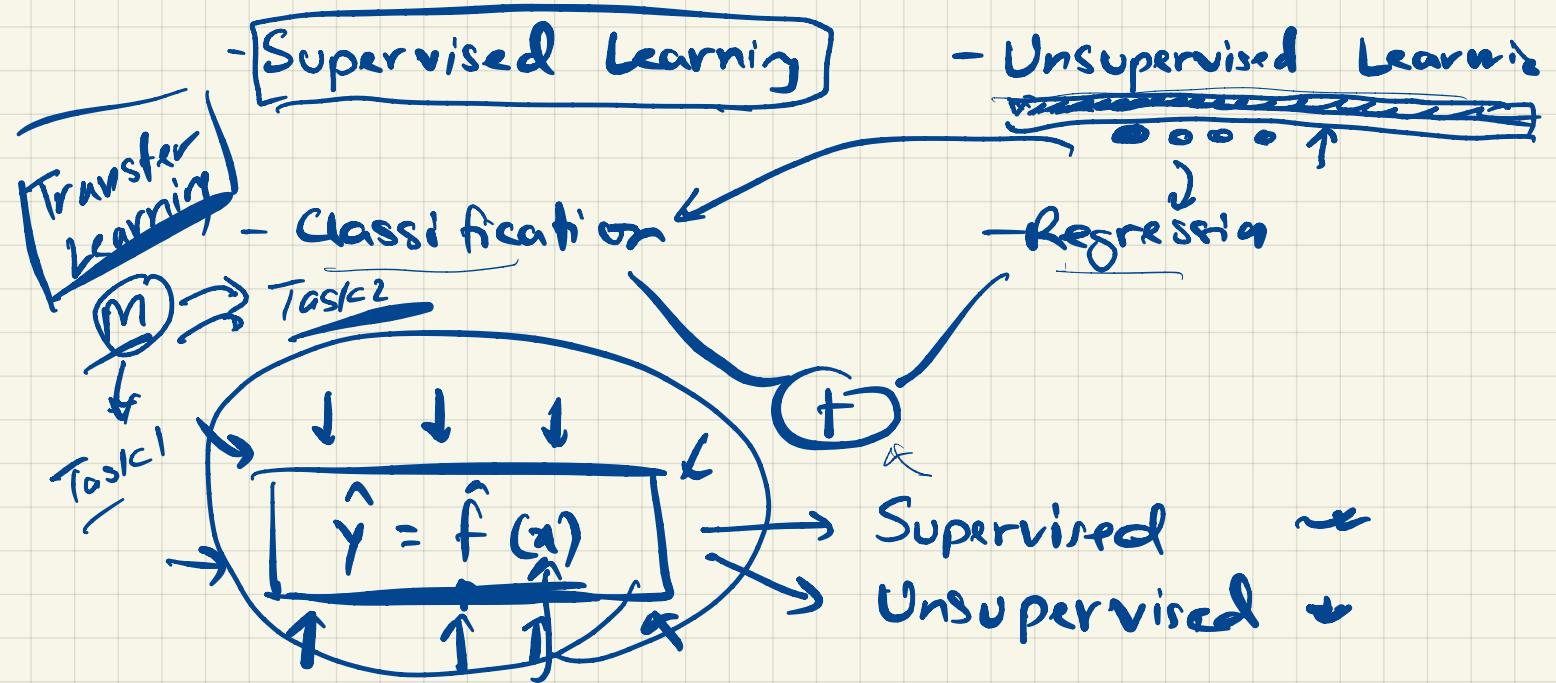
$$E = R/3 = 8 \times 10^8$$

$$\begin{aligned} \frac{dE}{dR} &= \frac{1}{3} \\ \frac{dE}{dq_1} &= \frac{dE}{dR} \cdot \frac{dR}{dq_1} \\ &= \frac{1}{3} = \frac{dE}{dq_1} = \frac{dE}{dq_2} = \frac{dE}{dq_3} \end{aligned}$$

$$\begin{aligned} q_1 &= p_1^2 \\ &= 560.93 \cdot 10^6 \\ \frac{dE}{dp_1} &= \frac{1}{2} \frac{dp_1}{dq_1} = \frac{1}{2} \frac{dp_1}{dq_2} = \frac{1}{2} \frac{dp_1}{dq_3} \end{aligned}$$

$$\begin{aligned} q_2 &= p_2^2 \\ &= 2.4 \times 10^9 \\ q_3 &= p_3^2 = 220.78 \times 10^6 \end{aligned}$$

Lecture 7



(x_1, y_1)

:

(x_k, y_k)

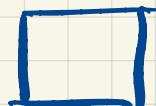
$$y = f(x) + \eta$$

$$\hat{y} = \hat{f}_\theta(x)$$

$$\lambda = E(\hat{y}, y) \leftarrow$$

$$\frac{dE}{d\theta}$$

Gradient Descent



→ Cat / Dog

Image

Cat

Human

(x_1, y_1)

Supervised
Learning

(x_2, y_2)

Ground Truth

Unsupervised



Image
(B-W)

→ | m | →



Color

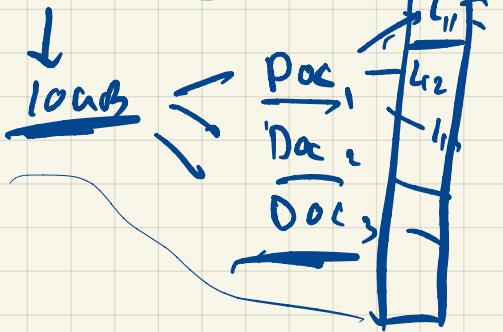
(x_1, y_1)

(x_k, y_k)

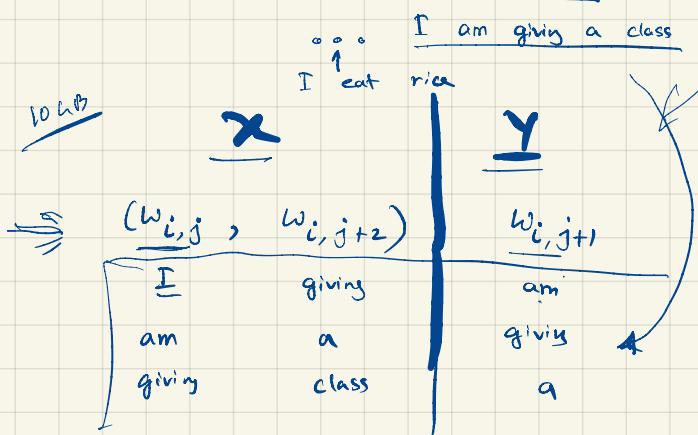
10,000
Buy

10 GB → Text

Web Scraper



$$S_i = (w_{i,1}, \dots, w_{i,m})$$



Unsupervised learning

$$t_1, t_2, t_3 \rightarrow [m] \rightarrow \hat{t}_2 = y^* \\ \hat{t}_2 = f(y^*)$$

Classification

vs

Regression

Discrete ↗
[::]

Continuous ↗
[::]

↗

$$= B \text{NL}(A x)$$

↗

Interpretation

I → M → Cat / Dog

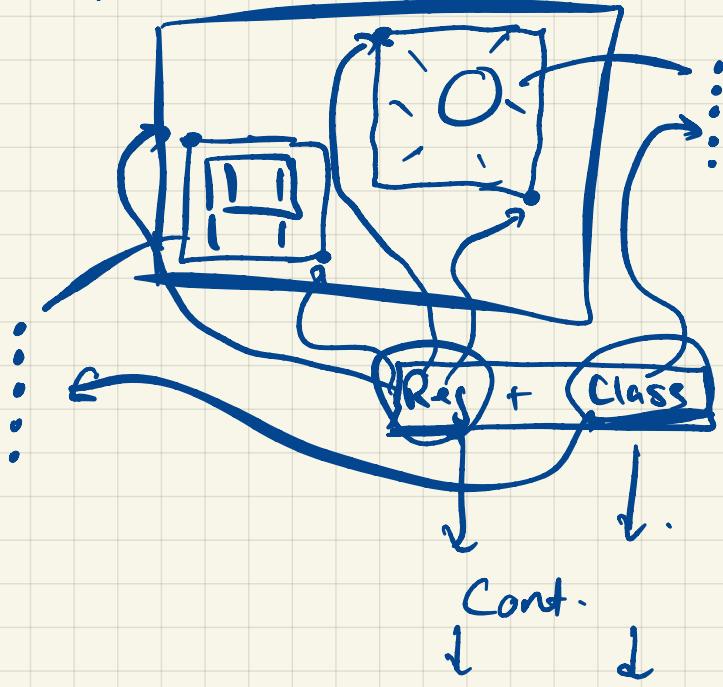
I → M → S_cat, S_dog

S_cat, S_dog, S_late, S_ugly

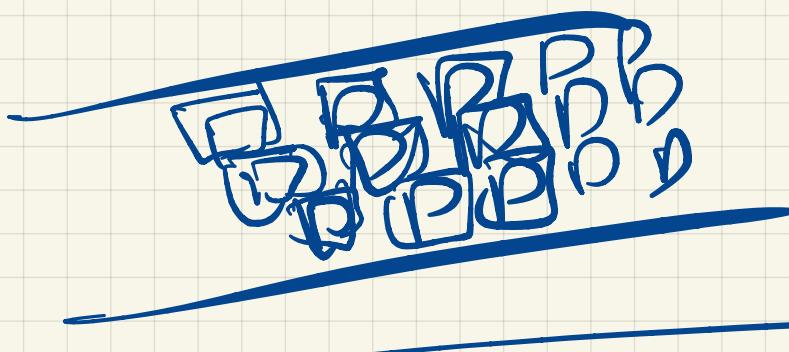
I → M →

- Table →
- Furniture →
- Car →
- Toyota →
- Horse →

Object Detection



Yolo
↳
Faster R-CNN
Efficient Det



- Entropy (Information Theory)

$P \rightarrow I.T$
↓
M.S.E → classification → One Hot Encoding
↓
P. Dist. →
→ {Entropy} ←

Lecture 8

Classification

PyTorch
Friday

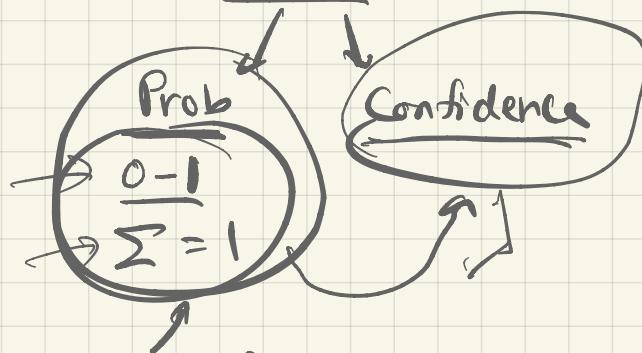
2130 pm

- 1 out of many classes

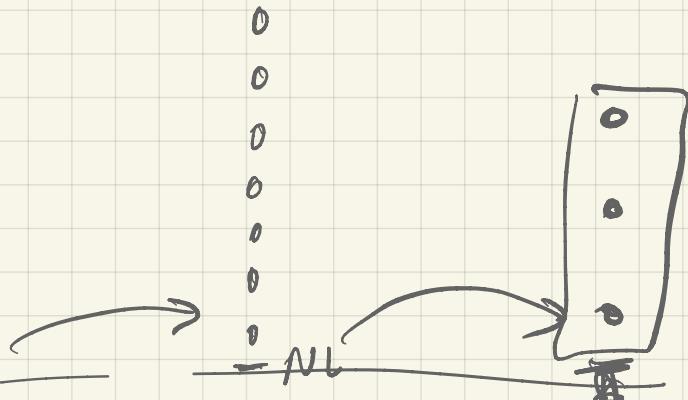
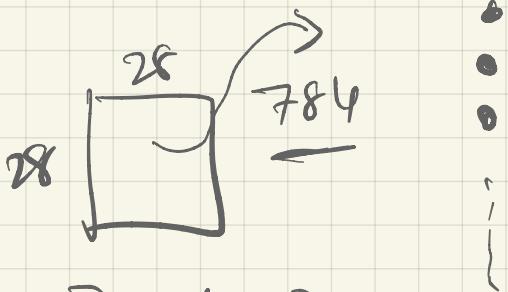


- Asking the model to predict scores for each class.

cat	dog	bird
321	500	<u>786</u>



- # of classes is fixed.



Input	Label	Transformed
x_1	Cat	$[1 \ 0 \ 0]^T$
x_2	Dog	$[0 \ 1 \ 0]^T$
x_3	Bird	$[0 \ 0 \ 1]^T$

One Hot representation

wordvec

7

$$\hat{Y}_i = \text{Impe} \quad A \cdot \text{NL}(B\vec{x}_i)$$

NL = Element Wise Sq.

$$L = \text{MS.E}(\hat{Y}_i, Y_i) \quad P.D$$

$$\begin{bmatrix} 0 & 10 \\ 0 & 01 \end{bmatrix} \quad \begin{bmatrix} 100 \\ 100 \end{bmatrix}$$

$[3.5, 0.000001, 2 \text{ Billion}]^T$ One hot encoding

$$\vec{x}_i \in \mathbb{R}^3$$

Each value 0-1
 $\sum = 1$

$$\hat{Y}_i = S(A \cdot \text{NL}(B\vec{x}_i))$$

$$S \left(\begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix} \right) = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix}$$

$$0 \leq p_j \leq 1 \quad \sum_{j=1}^3 p_j = 1$$

$$\begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix} \left[\begin{bmatrix} -200 \\ 10 \\ 50 \cdot 2 \end{bmatrix} \right]$$

$$p_j = \frac{e^{t_j}}{\sum_{k=1}^3 e^{t_k}} \quad (\text{Softmax})$$

$$t_1 \rightarrow \frac{e^{t_1}}{e^{t_1} + e^{t_2} + e^{t_3}}$$

Shopne Paisi

$$\begin{matrix} 0 & \rightarrow & x_2 \\ 0 & \rightarrow & x_2 \end{matrix}$$

P.D

(Empirically found)

$$\rightarrow \text{Softmax} \rightarrow \hat{x}_i$$

78M

$$\Rightarrow \text{MSE } (\hat{x}_i, x_i) \Leftarrow$$

$$\begin{matrix} \cancel{\frac{0.2}{+}} \\ \cancel{\frac{0.3}{-}} \\ \cancel{\frac{0.5}{+}} \\ \cancel{\frac{0}{-}} \\ \cancel{\frac{0}{+}} \end{matrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \begin{matrix} \overset{①}{\frac{1}{0}} \rightarrow [0] \\ \overset{②}{\frac{0}{0}} \rightarrow [0.09] \\ \overset{③}{\frac{0}{1}} \rightarrow [0.04] \end{matrix} \quad \begin{matrix} x_i \\ \hat{x}_i \end{matrix} \quad \text{very large}$$

$$\begin{aligned} \text{MSE}(\hat{x}_i, x_i) &= \frac{1}{3} [(0.5 - 1)^2 + (0.3 - 0)^2 + (0.2 - 0)^2] \\ &= \frac{1}{3} [(-0.5)^2 + (0.3)^2 + (0.2)^2] \\ &= \frac{1}{3} [(0.25) + (0.09) + (0.04)] \\ &= \frac{1}{3} [0.38] = \underline{0.1267} \end{aligned}$$

Measurement of distance between two

P.D.s..

K.L. Divergence a, b

\$\\$

{ entropy }

EIE

C.S. E

C.S.

(Sfm1)

(NBM)

Nabeel

Yusuf

Dr. Sifat Momen

0.3

0.333

Cake

0.2

0.3

0.333

Burger

0.5

0.3

0.333

Pizza

0.3

0.4

0.333

$$L(\text{cake}) = -\log 0.2$$

$$L(\text{Burger}) = -\log 0.5$$

$$L(\text{Pizza}) = -\log 0.3$$

$$\text{Average } L = \frac{-\log 0.2 + (-\log 0.5) + (-\log 0.3)}{3}$$

(Sfm 1)

Dr. Sifat Momen

<u>Cake</u>	<u>0.2</u>
<u>Burger</u>	<u>0.5</u>
<u>Pizza</u>	<u>0.3</u>

$$L(\text{Cake}) = -\log 0.2 = \cancel{\cancel{0.2}} \quad 2.32$$

$$L(\text{Burger}) = -\log (0.5) = \cancel{\cancel{0.5}} \quad 1$$

$$L(\text{Pizza}) = -\log (0.3) = 1.73$$

$$\left. \begin{aligned} & \frac{1}{3}C + \frac{1}{3}B + \frac{1}{3}P \\ & \frac{2.32 + 1 + 1.73}{3} \\ & = 1.68 \end{aligned} \right\}$$

<u>C</u>	<u>C</u>	<u>B</u>	<u>B</u>	<u>B</u>	<u>B</u>	<u>P</u>	<u>P</u>	<u>P</u>
2.32	2.32	1	1	1	1	1.73	1.73	1.73

$$\Theta \cdot \frac{2C + 5B + 3P}{10} = \underline{0.2C} + \underline{0.5B} + \underline{0.3P}$$

$$\boxed{H(\text{Sifat})} = -0.2 \log 0.2 + (-0.5 \log 0.5) + (-0.3 \log 0.3)$$

$$= 1.48$$

Lecture 9

$$\rightarrow -\sum_{i=1}^3 s_i \log_2 s_i$$

(SFM1)
Dr. Sifat Momen

<u>Cake</u>	0.2	<u>s₁</u>
<u>Burger</u>	0.5	<u>s₂</u>
<u>Pizza</u>	0.3	<u>s₃</u>

(NBM)
Nabeel

0.3	<u>n₁</u>
0.3	<u>n₂</u>
0.4	<u>n₃</u>

$$H_S(N) =$$

$$\sum_{i=1}^3 n_i \log_2 s_i$$

<u>S_c</u>	$-\log_2 0.2$ $= 2.32$
<u>S_B</u>	$-\log_2 0.5$ $= 1$
<u>S_P</u>	$-\log_2 0.3$ $= 1.73$

<u>N_c</u>	$-\log_2 0.3$ $= 1.73$
<u>N_B</u>	$-\log_2 0.3$ $= 1.73$
<u>N_P</u>	$-\log_2 0.4$ $= 1.32$

$$H(S) = 0.2 * 2.32 + 0.5 * 1 + 0.3 * 1.73 \\ = 1.48$$

$$H(N) = 0.3 * 1.73 + 0.3 * 1.73 + 0.4 * 1.32 \\ = 1.571$$

$$H_S(N) = 0.3 * 2.32 + 0.3 * 1 + 0.4 * 1.73$$

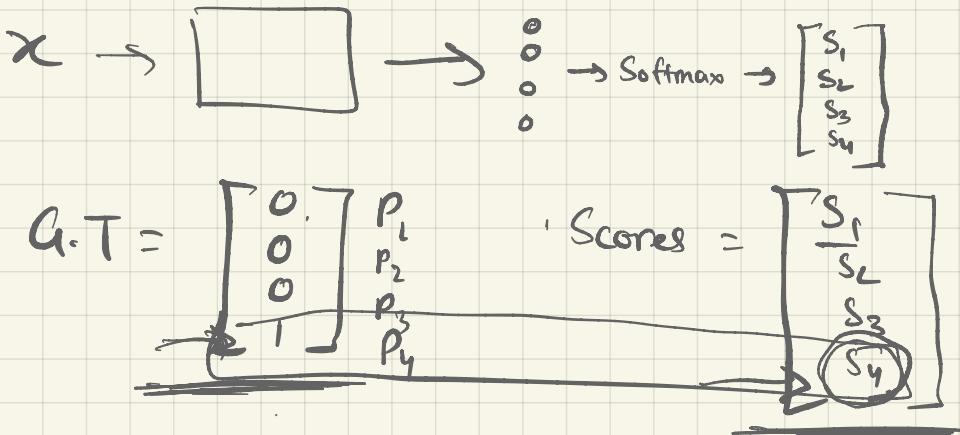
$$\text{cross entropy} = 1.688 \quad \text{cost} = H_S(N) - H(N)$$

$$\text{cost}$$

$$KL_{S, N} = -\sum_{i=1}^3 n_i \log_2 s_i - \left(-\sum_{i=1}^3 n_i \log_2 n_i \right)$$

$$KL_{N, S} = -\sum_{i=1}^3 s_i \log_2 n_i - \left(-\sum_{i=1}^3 s_i \log_2 s_i \right)$$

$$\frac{x}{x_i} \left| \begin{array}{cccc} x \\ [0 \ 0 \ 0 \ 1]^T = y_i \end{array} \right.$$



$$K.L = \sum_{i=1}^4 p_i \log_2 s_i - \left(-\sum_{i=1}^4 p_i \log_2 p_i \right)$$

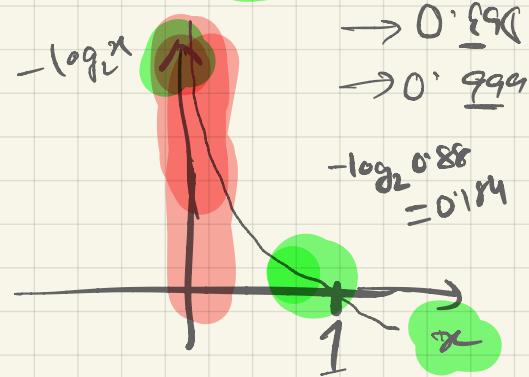
$$= -\log_2 s_4 + \log_2 1$$

$$= -\log_2 s_4$$

Cross entropy loss
 my \rightarrow log likelihood loss
 new \rightarrow log loss, log loss

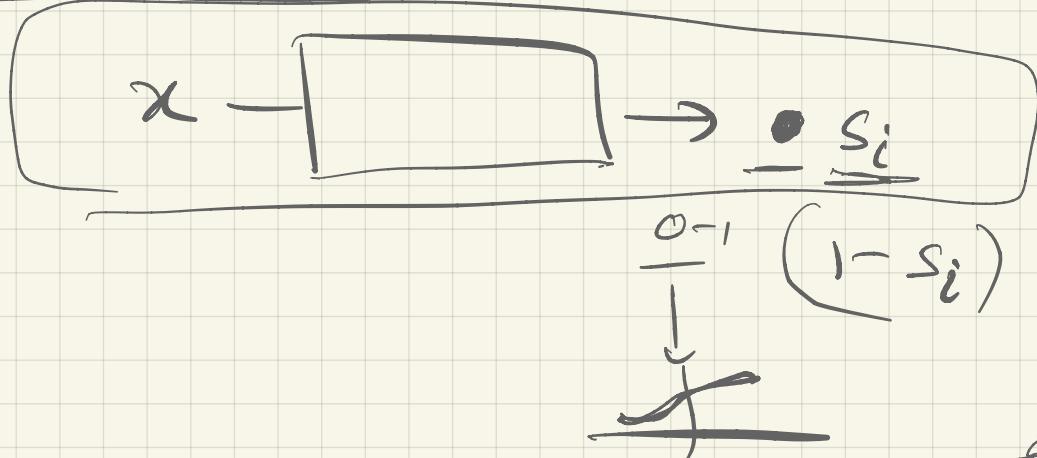
$$[] []$$

$$\frac{C.T}{S_i} \frac{1}{0.001} \frac{-\log_2 0.001}{0.184} = 9.965$$

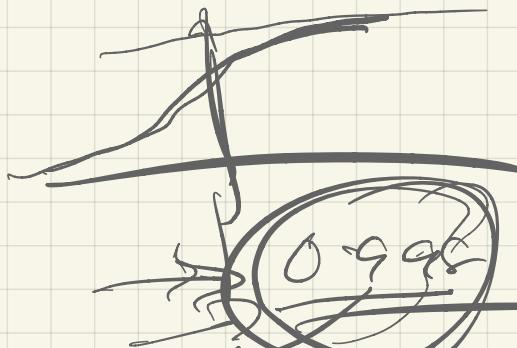


$\alpha \cdot \frac{dE}{d\theta}$

Binary Classification



G.T		Prediction
x_i	$P_i(1/0)$	s_i

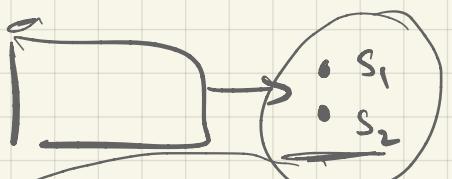


$$\text{Loss} = -P_i \log_2 s_i - (1 - P_i) \log_2 (1 - s_i)$$

x	P_i	s_i	Loss	
.	1	0.99	Small (Zero)	0
)	1	0.02	Bigs	α
)	0	0.02	Small (Zero)	0
4 3	$\frac{e^4}{e^4 + e^3}$	0.9	Bigs (Zero)	α

(4)
3

$$\frac{e^3}{e^4 + e^3} = \alpha$$



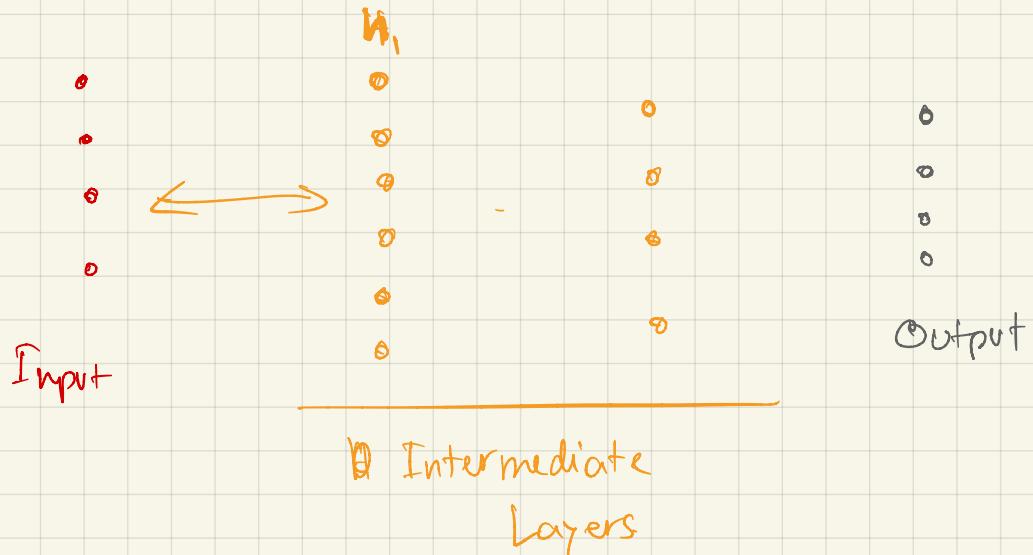
$$\begin{bmatrix} 0 \\ 1 \end{bmatrix} \Bigg/ \begin{bmatrix} 1 \\ 0 \end{bmatrix} \frac{P_i}{P_{i-1}}$$

$$P_2 = (1 - P_1)$$

$$\begin{aligned} \alpha &= -P_1 \log_2 s_1 - P_2 \log_2 s_2) \quad S_2 = (1 - S_1) \\ &= -P_1 \log_2 s_1 - (1 - P_1) \log_2 S_2 \\ &= \boxed{-P_1 \log_2 s_1 - (1 - P_1) \log_2 (1 - s_1)} \end{aligned}$$

Binary Cross Entropy

Lecture 10



$$H_1 = A \cdot x$$

$$H_1^{NL} = \sigma_1(H_1)$$

$$H_2 = B \cdot H_1^{NL}$$

$$H_2^{NL} = \sigma_2(H_2)$$

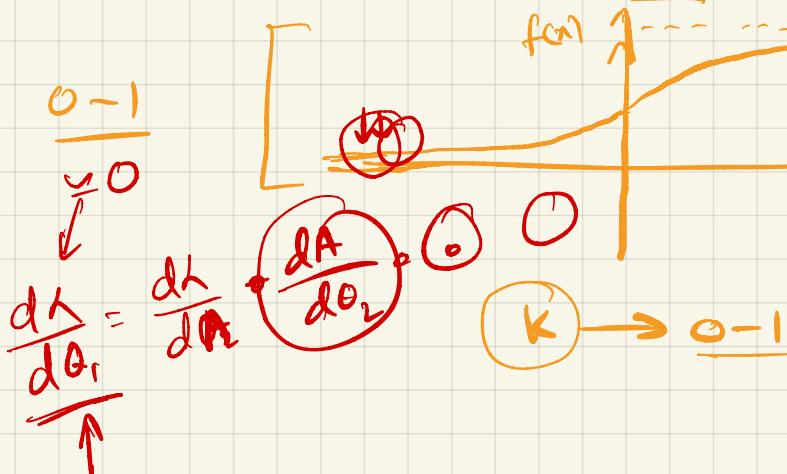
$$O_{\text{Score}} = \frac{C \cdot H_2^{NL}}{\downarrow}$$

$$O_{\text{prob}} = \text{Softmax}(O_{\text{Score}})$$

① Sigmoid

$$f(x) = \frac{1}{(1 + e^{-x})}$$

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



$$\frac{dL_{\text{Final}}}{dO_1} = \frac{\lambda_1 + \lambda_2 + \lambda_3}{c}$$

$$\frac{L_{\text{Final}}}{dO_1} = \frac{(1, 0, 0) + (0, 1, 0) + (0, 0, 1)}{3} = \frac{(1, 1, 1)}{3}$$

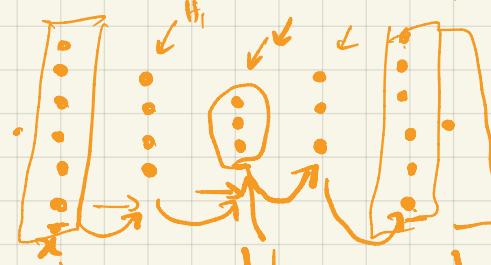
$$\frac{L_{\text{Final}}}{dO_1} = \frac{(1, 1, 1)}{3}$$

Autoencoder



$$\hat{x} = f(x)$$

$$\min ||\hat{x} - x||_1$$



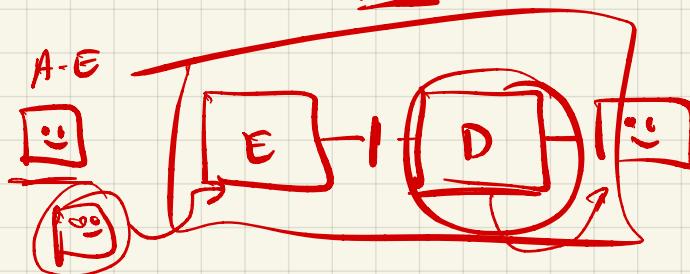
② Color
0-255
0-1

$$x = \begin{bmatrix} 0.2 \\ 0.4 \\ 0.5 \\ 0.8 \end{bmatrix}, \hat{x} = \begin{bmatrix} 0.768 \\ 0.768 \end{bmatrix} = \frac{1}{768}$$

mse

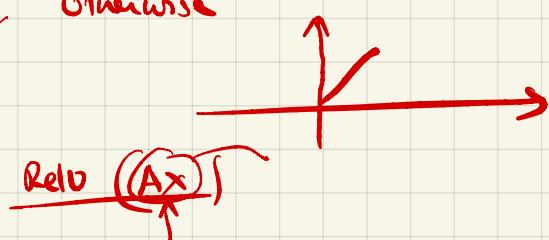
$$\begin{aligned} H_1 &= A(x) \\ H_1^{NL} &= \sigma_1(H_1) \\ H_2 &= B H_1^{NL} \\ H_2^{NL} &= \sigma_2(H_2) \\ H_3 &= C H_2^{NL} \\ H_3^{NL} &= \sigma_3(H_3) \\ H_4 &= D H_3^{NL} \\ H_4^{NL} &= \text{Sigmoid}(H_4) \end{aligned}$$

Empirically



③

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



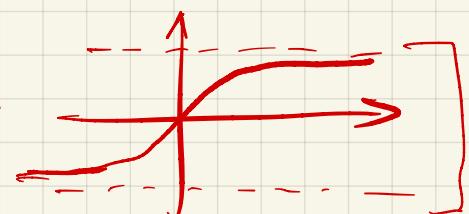
④ Leaky ReLU

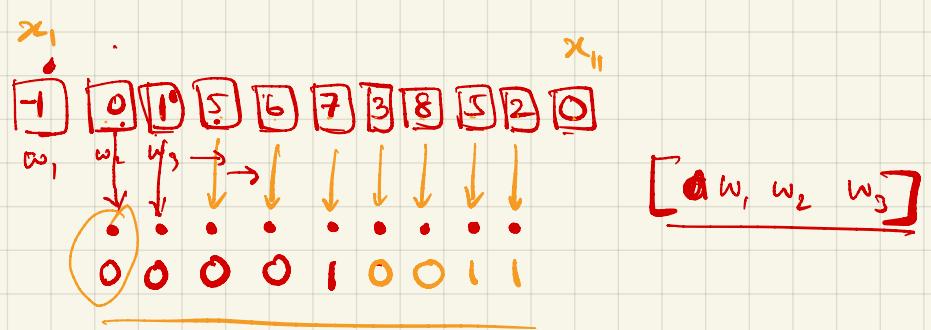
$$\text{LReLU} = \begin{cases} x, & x > 0 \\ \alpha x, & x \leq 0 \end{cases}$$



⑤ tanh(x) =

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





Sigmoid $(A \cdot X)$ \circ $\left[\begin{array}{c} R_1^T \\ R_2^T \\ R_3^T \end{array} \right]$ ↑ column
rows
99 parameters

$$X = \begin{bmatrix} -1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$R_1 = \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$R_1 \cdot X = -2$$

$$\text{Sigmoid } (-2) = 0.119 = 0$$

$$R_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ -1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$R_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ -1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$