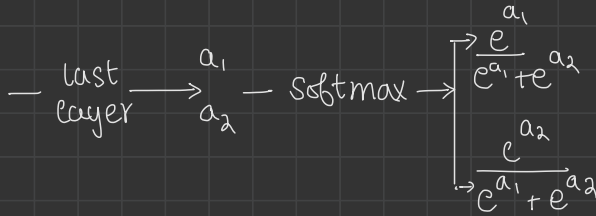


MOCK QUIZ 2

Q1. Prove that sigmoid is a special case of softmax -

i). take a 2 class classification problem:



$$\text{Prob}(\text{class 1}) = \frac{e^{a_1}}{e^{a_1} + e^{a_2}} = \frac{1}{1 + e^{a_2 - a_1}}$$

$$\text{Prob}(\text{class 2}) = \frac{e^{a_2}}{e^{a_1} + e^{a_2}} = \frac{e^{a_2 - a_1}}{1 + e^{a_2 - a_1}}$$

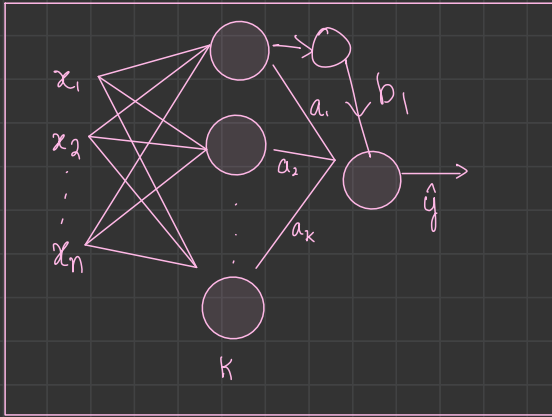
$$\text{let } a_1 - a_2 = x$$

$$\Rightarrow \text{Prob}(\text{class 1}) = \frac{1}{1 + e^{-x}} \rightarrow \text{sigmoid function}$$

$$\text{Prob}(\text{class 2}) = \frac{e^{-x}}{1 + e^{-x}}$$

Q2. A deep NN without non-linear activations is still equivalent to Linear regression

Consider the following neural network with 2 layers:



at the i^{th} neuron in the first layer:

$$x_1 w_{i1}^1 + x_2 w_{i2}^1 + \dots + x_n w_{in}^1 = a_i$$

$$\Rightarrow a_i = \sum_{j=1}^n x_j w_{ij}^1 \rightarrow 1.$$

$$\hat{y} = \sum_{i=1}^k w_i^2 a_i \rightarrow 2.$$

using 1 in 2,

$$\hat{y} = \sum_{i=1}^k w_i^2 \sum_{j=1}^n w_{ij}^1 x_j$$

$$b_i = w_i^4 a_i + b$$

$$\hat{y} = \sum_{i=1}^k w_i^2 b_i$$

$$\hat{y} = \sum_{i=1}^k w_i^2 (w_i^4 a_i + b)$$

$$= \sum w_i^2 w_i^4 a_i + w_i^2 b$$

$$\hat{y} = \sum_{j=1}^n \sum_{i=1}^k w_i^2 w_{ij}^1 x_j = \sum_{j=1}^n w_i^2 w_{ij}^1 \left(\sum_{j=1}^n x_j w_{ij}^1 \right) + w_i^2 b$$

$$= \sum_{j=1}^n x_j \underbrace{\sum_{i=1}^k w_i^2 w_{ij}^1}_{\text{can be written in terms of } w_j^3}$$

$$= \sum_{j=1}^n x_j w_j^3 \quad \therefore \text{where } w_j^3 = \sum_{i=1}^k w_i^2 w_{ij}^1$$

This is equivalent to a linear reg. problem.

We can say that by induction, it will work for multiple layers.

Q3. Domain of sigmoid, tanh. Write tanh in terms of sigmoid.

i). Sigmoid(x) = $\frac{1}{1+e^{-x}}$	\rightarrow sigmoid: domain: $(-\infty, \infty)$ Range: $(0, 1)$
ii). tanh(x) = $\frac{e^x - e^{-x}}{e^x + e^{-x}}$	\rightarrow tanh: domain: $(-\infty, \infty)$ Range: $(-1, 1)$

tanh in terms of sigmoid:

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

$$= \frac{1}{1 + e^{-2x}} - \frac{1}{e^{2x} + 1}$$

\therefore division by e^x, e^{-x} respectively

$$= \text{sigmoid}(2x) - \text{sigmoid}(-2x)$$

4. Input size $(32 \times 32 \times 6)$ (image), train for 100 class classification

i) Layer 1 \rightarrow 200 neurons, ReLU

ii) Layer 2 \rightarrow 120 neurons,

iii) 100 class classification.

\therefore Number of params (layer i)

$$= (N_{i-1} + 1) N_i$$

↑ Bias.

Steps in the forward pass:

i). Flatten the image, each image has $32 \times 32 \times 6 = 6144$ data points.

ii). Layer 1: 200 neurons:

each neuron will have 6144 parameters along with a bias term.

Thus each neuron has: 6145 parameters

\rightarrow Layer 1 has $200 \times 6145 = \underline{1,229,000}$ parameters.

iii). ReLU does not have any parameters.

iv). Layer 2: 120 neurons

\rightarrow input to this layer will have a size of 200

\rightarrow Thus, each neuron will have 201 params (with bias)

\rightarrow Layer 2 has $120 \times 201 = \underline{24,120}$ params.

v). Output layer: 100 neurons (100 class classification)

\rightarrow input will be of size 120, each neuron will have 121 params.

\rightarrow Layer 3 has $100 \times 121 = \underline{12,100}$ params.

\Rightarrow total number of params = $12,100 + 24,120 + 1,229,000$

$$\text{total params} = 1,266,220$$

Q5. Derive the vectorized form of gradient descent for logistic reg.

Binary
Classification
Case

$$\hat{y} = \frac{1}{1 + e^{-(x^T \theta + b)}} \quad \because \text{sigmoid.}$$

loss function : $\sigma(x^T \theta + b)$

$$L(\hat{y}, y) = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

$$\text{Now } \frac{dL}{d\theta} = \frac{dL}{d\hat{y}} \times \frac{d\hat{y}}{d\theta} \quad \text{--- 1.}$$

$$\text{i). } \frac{dL}{d\hat{y}} = \frac{-y}{\hat{y}} - \frac{(1-y)(-1)}{1-\hat{y}}$$

$$= \frac{(1-y)}{(1-\hat{y})} - \frac{y}{\hat{y}}$$

$\because \hat{y}$ is not a vector
in this case

$$\frac{dL}{d\hat{y}} = \frac{(1-y)}{(1-\hat{y})} - \frac{y}{\hat{y}} \quad \text{--- 2.}$$

$$\frac{d\sigma(a)}{da} = \sigma(a) \times (1-\sigma(a))$$

$$\because \sigma' = \sigma(1-\sigma)$$

$$\text{ii). } \frac{d\hat{y}}{d\theta} = \frac{d(\sigma(x^T \theta + b))}{d\theta} = \sigma(1-\sigma) \cdot \frac{d(x^T \theta + b)}{d\theta}$$

$$= \hat{y}(1-\hat{y}) \frac{d(\theta^T x + b)}{d\theta} \quad \frac{d(\theta^T x)}{d\theta}$$

$$= \hat{y}(1-\hat{y}) x \quad \text{--- 3.}$$

$$\frac{d(\theta^T x)}{d\theta} = x$$

Sub 3,2 in 1:

$$\begin{aligned}\frac{dL}{d\theta} &= \hat{y}(1-\hat{y})x \times \left\{ \frac{(1-y)}{1-\hat{y}} - \frac{y}{\hat{y}} \right\} \\ &= \cancel{\hat{y}(1-\hat{y})}x \left(\frac{\hat{y}-y\hat{y}-y+y\hat{y}}{\cancel{\hat{y}(1-\hat{y})}} \right)\end{aligned}$$

$$\boxed{\frac{dL}{d\theta} = x(\hat{y}-y)}$$

With the entire data set

$x \rightarrow$ complete dataset

$$\hat{y} = \sigma(x\theta)$$

vector

$$\text{loss function} = J(\theta) = \frac{1}{m} \sum_{i=1}^m -y_i \log \hat{y}_i - (1-y_i) \log (1-\hat{y}_i)$$

$$\Rightarrow \frac{dL}{d\theta} = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) (x[i, :])^T$$

$$\boxed{\begin{aligned}J(\theta) &= \frac{1}{m} \sum_{i=1}^m L_i(\theta) \\ \frac{dJ(\theta)}{d\theta} &= \frac{1}{m} \sum_{i=1}^m \sum L_i(\theta)\end{aligned}}$$

$$\hat{y} = \frac{1}{1 + e^{-(x^T \theta + b)}} \quad \because \text{sigmoid}$$

loss function :

$$\sigma(x^T \theta + b)$$

$$L(\hat{y}, y) = \frac{-y \log \hat{y} - (1-y) \log(1-\hat{y})}{}$$

$$\text{Now } \frac{dL}{d\theta} = \frac{dL}{d\hat{y}} \times \frac{d\hat{y}}{d\theta} \quad \text{--- 1.}$$

$$\frac{dL}{db} = \frac{dL}{d\hat{y}} \times \frac{d\hat{y}}{db}$$

$$\frac{d\hat{y}}{db} = \frac{d(\sigma(x^T \theta + b))}{db}$$

$$= \hat{y}(1-\hat{y}) \times \frac{d(x^T \theta + b)}{db}$$

$$= \hat{y}(1-\hat{y})$$

$$\begin{aligned} \frac{dL}{db} &= \hat{y}(1-\hat{y}) \times \frac{dL}{d\hat{y}} = \hat{y}(1-\hat{y}) = \cancel{\hat{y}(1-\hat{y})} \times \frac{\hat{y} - y}{\cancel{\hat{y}(1-\hat{y})}} \\ &= \hat{y} - y \end{aligned}$$

$$\frac{dL}{db} = \hat{y} - y$$

$$\begin{aligned} \frac{dL}{db} &= \sum \frac{dL_i}{db} \\ &= \sum_{i=1}^m \hat{y}_i - y_i \end{aligned}$$

→ for the complete dataset

Q. 6.]

$$\text{Entropy} = - \sum_{i=1}^K p_i \log(p_i)$$

no cross terms

Decision
Trees

$$\text{Entropy (for 2 classes)} = - p_1 \log p_1 - p_2 \log p_2$$

$$\text{Cross-Entropy} = - \sum_{i=1}^K y_i \log(\hat{y}_i)$$

cross terms

Logistic
Regression

$$\text{(for 2 classes)} = - y \log \hat{y} - (1-y) \log(1-\hat{y})$$

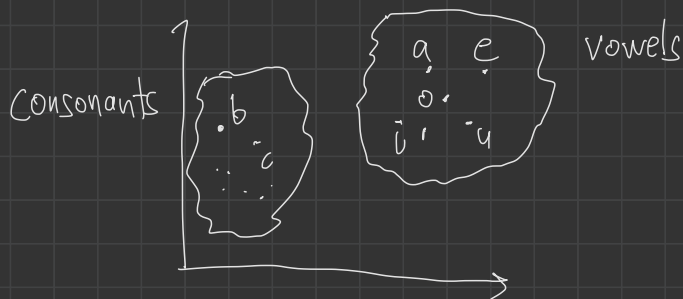
Q.7.] ① Extremely sparse embeddings
for large vocabularies
(wasteful)

$$a = [1 \ 0 \ 0 \ 0 \ \dots \ \dots \ \dots]$$

←-----|v|-----→

② • 1-hot encodings make each
embedding orthogonal (independent
and unrelated)

- We ideally want embeddings
which capture the semantics /
usage patterns too.



How do we solve this ???

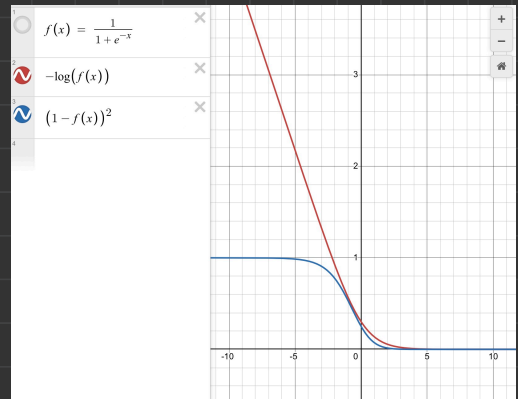
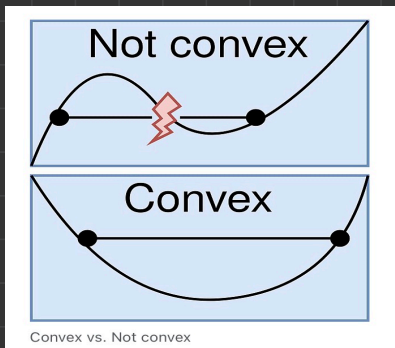
$$Q.8.] J(\theta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\hat{y}_i = \sigma(\theta^T x)$$

Why don't we use Squared Error?

① Squared Error is not convex.

(Difficult for GD to reach global optima.
May get stuck in local optima.)



We generally want to optimize convex functions

Using Gradient Descent.

→ For a function of 1 variable, if the 2nd derivative is non-negative in the domain, the function is convex in the domain

→ Quick test for Convexity - A line joining any two points on the curve must lie above the curve.

② MSE doesn't penalize much even for a perfect mismatch, Eg. $y_i = 1$ $\hat{y}_i = 0$

$$\text{MSE} = (y_i - \hat{y}_i)^2 = (1 - 0)^2 = 1$$

$$\begin{aligned} \text{Log Loss} &= -1 \log 0 - 0 \log 1 \\ &= -\log 0 \quad (\text{tends to infinity}) \end{aligned}$$

Q.9.] For N datapoints and K classes,

$$\text{Multi-class Cross-entropy} = - \sum_{i=1}^N \sum_{k=1}^K y_i^k \log \hat{y}_i^k$$

Eg. ① $\begin{bmatrix} 0.8 \\ 0.1 \\ 0.1 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ ② $\begin{bmatrix} 0.3 \\ 0.3 \\ 0.4 \end{bmatrix}$ $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$

\hat{y}_1 y_1 \hat{y}_2 y_2

$$\begin{aligned} \text{Loss} &= (-1 \log(0.8) - 0 \log(0.1) - 0 \log(0.1)) \\ &\quad + (-0 \log(0.3) - 1 \log(0.3) - 0 \log(0.4)) \\ &= -(\log(0.8) + \log(0.3)) = -\log 0.24 \\ &= \log\left(\frac{25}{6}\right) \end{aligned}$$

Q.10.] $CE(p, y) = \begin{cases} -\log(p) & \text{if } y=1 \\ -\log(1-p) & \text{if } y=0 \end{cases}$

Note:

$y=1$ is the majority class (we have a lot of it's samples)

$$p_t = \begin{cases} p & \text{if } y=1 \\ (1-p) & \text{if } y=0 \end{cases}$$

$$CE(p, y) = CE(p_t) = -\log(p_t)$$

Consider $FL(p_t) = -[1-p_t]^\gamma \log(p_t) ; \gamma \geq 0$

→ Why is $CE(p_t)$ not well-suited for imbalanced datasets (Eg. Cancer detection)??

• Biased towards majority class

→ Can Focal Loss ($FL(p_t)$) help here??

How??

→ In practice, $\alpha_t = \begin{cases} \alpha & \text{if } y=1 \\ 1-\alpha & \text{if } y=0 \end{cases}$

' α ' is a hyperparameter.
' γ ' is a hyperparameter. } can be tuned

$$\rightarrow FL(p_t) = -\alpha_t (1-p_t)^\gamma \log(p_t)$$

(used in practice.)

Answer. When p_t is high (easily classified samples from the majority class), the $(1-p_t)^\gamma$ down-scales the loss contributed by such majority class samples, but when p_t is low, there is lesser down-scaling \Rightarrow ^{more} contribution from minority class.

