

Contents

Problem 1: XOR Problem	2
Problem 2: L_2 Regularized Linear Regression Via Stochastic Gradient Descent	6
Problem 3: Logistic Sigmoid Identity	
7 a. Prove $\sigma(-x) = 1 - \sigma(x) \forall x$	7
b. Prove $\sigma'(x) = \frac{\partial \sigma}{\partial x}(x) = (\sigma(x)(1 - \sigma(x))) \forall x$	8
Problem 4: Regularization to encourage symmetry	9
Problem 5: Linear-Gaussian Prediction Model	10

Problem 1: XOR Problem

① XOR \rightarrow

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

$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$
↓ ↓ ↓ ↓

give $\rightarrow y = \begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix}^T = (f^*(x))$
(ground truth)

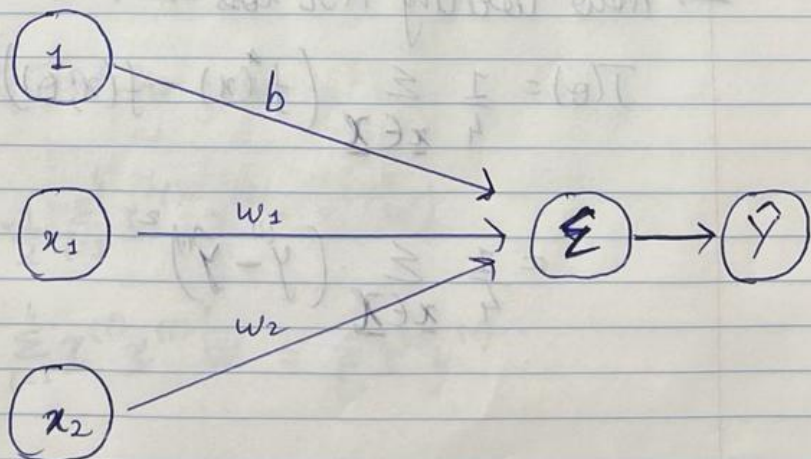
$$\text{MSE} \Rightarrow J(0) = \frac{1}{4} \sum_{x \in X} (f^*(x) - f(x; 0))^2$$

$$f(x; w, b) = x^T w + b$$

Here, $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}^T$ (vector) & $x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$b = b$ (scalar)

\rightarrow 2 layer NN representation \rightarrow



→ After understanding the NN model, modifying the given data to make the matrix equations more consistent.

$$(w)^* W = \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} \quad \text{(weight matrix)} \quad Y \leftarrow w \cdot p$$

$$X = \begin{bmatrix} 0 & 0 & 1 & x_1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

$$Y = \begin{bmatrix} 0 & 1 & 1 & 0 \end{bmatrix}^T \quad \begin{matrix} y^{(1)} & y^{(2)} & y^{(3)} & y^{(4)} \end{matrix}$$

$$x^{(i)} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 = 1 \end{bmatrix}$$

→ Now writing MSE loss →

$$J(\theta) = \frac{1}{4} \sum_{x \in X} (f^*(x) - f(x; \theta))^2$$

$$\hat{J} = \frac{1}{4} \sum_{x \in X} (y^{(i)} - \hat{y}^{(i)})^2$$

$$\hat{y} = \underline{x}^{(i)T} \underline{w} \quad (\text{Since this includes the bias no need for extra } (+b) \text{ term})$$

$$\therefore J(\theta) = \frac{1}{4} \sum_{i=1}^4 (y^{(i)} - \underline{x}^{(i)T} \underline{w})^2$$

→ differentiating with respect to $\nabla_{\underline{w}}$ (i.e. whole θ both w_1 & w_2 & b)

$$\begin{aligned} \nabla J(\theta) &= \nabla_{\underline{w}} \frac{1}{4} \sum_{i=1}^4 (y^{(i)} - \underline{x}^{(i)T} \underline{w})^2 \\ &= \frac{1}{4} \sum_{i=1}^4 \nabla_{\underline{w}} (y^{(i)} - \underline{x}^{(i)T} \underline{w})^2 \\ &= \frac{1}{4} \sum_{i=1}^4 \underline{x}^{(i)} \cdot 2 (y^{(i)} - \underline{x}^{(i)T} \underline{w}) \\ &= -\frac{1}{2} \sum_{i=1}^4 \underline{x}^{(i)} (y^{(i)} - \underline{x}^{(i)T} \underline{w}) \end{aligned}$$

Equating $\nabla J(\theta)$ to zero

$$0 = -\frac{1}{2} \sum_{i=1}^4 \underline{x}^{(i)} (y^{(i)} - \underline{x}^{(i)T} \underline{w})$$

$$\sum_{i=1}^4 \underline{x}^{(i)} \underline{x}^{(i)T} \underline{w} = \sum_{i=1}^4 \underline{x}^{(i)} y^{(i)}$$

$$\therefore W = \left(\sum_{i=1}^4 \underline{x}^{(i)} \underline{x}^{(i)T} \right)^{-1} \sum_i \underline{x}^{(i)} y^{(i)}$$

→ This equation can also be represented in terms of the design matrix \underline{X} & \underline{Y}

$$\therefore W = (\underline{X}\underline{X}^T)^{-1} \underline{X}\underline{Y}$$

$$\therefore W =$$

$$(\underline{X}\underline{X}^T)^{-1} = \left(\begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix} \right)^{-1}$$

$$= \left(\begin{bmatrix} 2 & 1 & 2 \\ 1 & 2 & 2 \\ 2 & 2 & 4 \end{bmatrix} \right)^{-1}$$

$$= \begin{bmatrix} 1 & 0 & -0.5 \\ 0 & 1 & -0.5 \\ -0.5 & -0.5 & 0.75 \end{bmatrix} \quad \left(\begin{array}{l} \text{calculated} \\ \text{using} \\ \text{MATLAB} \end{array} \right)$$

$$\underline{X}\underline{Y} = \begin{bmatrix} 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}$$

$$\therefore W = (\underline{X}\underline{X}^T)^{-1} \underline{X}\underline{Y}$$

$$= \begin{bmatrix} 1 & 0 & -0.5 \\ 0 & 1 & -0.5 \\ -0.5 & -0.5 & 0.75 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix}$$

$$= \begin{bmatrix} 1 + (-0.5) \times 2 \\ 1 + (-0.5) \times 2 \\ (-0.5) + (-0.5) + (2 \times 0.75) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0.5 \end{bmatrix}$$

$$\therefore \begin{bmatrix} w_1 \\ w_2 \\ b \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0.5 \end{bmatrix}$$

$$\therefore w_1 = 0 \quad w_2 = 0 \quad b = 0.5$$

Problem 2: L_2 Regularized Linear Regression Via Stochastic Gradient Descent

- We developed a 2-layer neural network using the data given in homework 1. The file is saved as **homework2_590146168_249812226.py**.
- The model is trained using Linear regression via stochastic gradient descent method.
- There are 4 Hyperparameters that can be altered:
 - Mini-Batch Size
 - Learning Rate
 - Number of Epochs
 - L_2 Regularization Strength
- The training data is split into Training (80%) and Validation (20%) dataset.
- We have selected a set of 10 values for each hyperparameter i.e., 10^4 combinations.
- The weights are initialized randomly and this same weight is used for training the model for each set of hyperparameters to keep the results fair.
- Once training is done, we calculate the cost on validation set. If the cost is minimum, the best hyperparameter is updated. This process is done till all hyperparameters are checked.
- It took approximately 22 hours to check all 10^4 combinations.
- At the end we obtain the best set of hyperparameters. Using these hyperparameters we train the model on validation + training dataset.
- Now this trained model is used to calculate the cost (unregularized MSE) for the test dataset and the result is reported below:

```
C:\Users\Fenil\anaconda3\python.exe "C:/THIS PC/WPI/SPRING 2022/CS 541 Deep Learning/h
-----
Performing Grid Search for 10e4 combinations of Hyperparameters:
-----
Grid Search Completed
Results after performing Grid Search:
Best Hyperparameters:
  Epochs= 750
  Alpha= 1
  Learning Rate= 0.00075
  Mini Batch Size= 32
Cost on Validation Set with Best Hyperparameters= 80.66432176996726

-----
Training on Training + Validation Dataset:
-----
Training Completed

-----
Performance Evaluation
-----
Cost on Test Dataset= 85.14249556730823
```


Problem 3: Logistic Sigmoid Identity

a. Prove $\sigma(-x) = 1 - \sigma(x) \forall x$

$$\begin{aligned} \text{L.H.S} &= \frac{1}{1+e^{-(-x)}} = \frac{1}{1+e^x} \\ \text{R.H.S} &= 1 - \frac{1}{1+e^{-x}} = \frac{1+e^{-x}-1}{1+e^{-x}} \\ &= \frac{e^{-x}}{1+e^{-x}} \\ &= \boxed{\frac{1}{1+e^x}} \quad \left(\text{dividing numerator \& denominator by } e^{-x} \right) \\ \boxed{\text{L.H.S} &= \text{R.H.S}} \end{aligned}$$

b. Prove $\sigma'(x) = \frac{\partial \sigma}{\partial x}(x) = (1 - \sigma(x)) \forall x$

(b) Prove that $\sigma'(x) = \frac{\partial \sigma}{\partial x}(x) = \sigma(x)(1 - \sigma(x)) \forall x$

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

$$\sigma'(x) = \frac{\partial}{\partial x} \left(\frac{1}{1+e^{-x}} \right) =$$

$$= \frac{-1}{(1+e^{-x})^2} \cdot (-e^{-x}) = \boxed{\frac{e^{-x}}{(1+e^{-x})^2}}$$

$$= \left(\frac{1}{1+e^{-x}} \right) \cdot \left(\frac{e^{-x}}{1+e^{-x}} \right)$$

$$= \left(\frac{1}{1+e^{-x}} \right) \left(\frac{1}{1+e^{+x}} \right) = \underline{\text{L.H.S}}$$

$$= \sigma(x)(1 - \sigma(x))$$

$$\therefore \underline{\underline{\text{L.H.S} = \text{R.H.S}}}$$

$$\left\{ \left(\frac{1}{1+e^{+x}} = (1 - \sigma(x)) \right) \right. \\ \left. \text{from last proof} \right\}$$

Problem 4: Regularization to encourage symmetry

④ L_2 regularization \rightarrow

$$J_{MSE}(w) = \frac{1}{2n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2 + \underbrace{\frac{\alpha}{2n} w^T w}_{\text{regularization term}}$$

\rightarrow our goal - discourage the weights from becoming too asymmetric.

$$\text{let } w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

\rightarrow according to our goal we should add $(w_2 - w_1)^2$ to the $J_{MSE}(w)$ for regularization. $(w_2 - w_1)^2$ ensures the weights are as similar to each other as possible.
- The squared term ensures regularization is non-negative.

\therefore Now, we are advised to use $\frac{\alpha}{2n} w^T I w$ for regularization.

\therefore let us assume a square matrix $S_{2 \times 2}$ such that

$$S = \begin{bmatrix} s_1 & s_2 \\ s_3 & s_4 \end{bmatrix}$$

S follows the following property,

$$(w_2 - w_1)^2 = \begin{bmatrix} w_1 & w_2 \end{bmatrix} \begin{bmatrix} s_1 & s_2 \\ s_3 & s_4 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$(w_2 - w_1)^2 = s_1 w_1^2 + w_1 w_2 (s_2 + s_3) + s_4 w_2^2$$

\rightarrow now expanding the L.H.S

$$w_2^2 - 2w_1 w_2 + w_1^2 = s_1 w_1^2 + w_1 w_2 (s_2 + s_3) + s_4 w_2^2$$

\rightarrow so from comparison we get to know that

$$s_1 = 1$$

$$s_4 = 1$$

$$s_2 + s_3 = -2 \Rightarrow s_2 = -1 \text{ \& } s_3 = -1$$

$$\therefore S = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

Problem 5: Linear-Gaussian Prediction Model

$$(5) \quad P(y|x) = \mathcal{N}(y|x^T w, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-x^T w)^2}{2\sigma^2}\right)$$

\Rightarrow expected value of y is $x^T w$

\Rightarrow the variance of y is constant (σ^2) for all possible x .

collection of dataset $\mathcal{D} = \{(x^{(i)}, y^{(i)})\}_{i=1}^n \rightarrow$ MLE for w & σ

\Rightarrow MLE - optimizing on w & σ^2

$$P(\mathcal{D}|w, \sigma^2) = \prod_{i=1}^n P(y^{(i)}|x^{(i)}, w, \sigma^2) \quad (\text{considering the conditional independence})$$

\rightarrow maximizing $\log P(\mathcal{D}|w, \sigma^2)$

$$\log P(\mathcal{D}|w, \sigma^2) = \log \prod_{i=1}^n P(y^{(i)}|x^{(i)}, w, \sigma^2)$$

$$= \sum_{i=1}^n \log P(y^{(i)}|x^{(i)}, w, \sigma^2)$$

$$= \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(y^{(i)} - x^{(i)T} w)^2}{2\sigma^2}\right) \right)$$

$$= \sum_{i=1}^n \left(\log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{(y^{(i)} - x^{(i)T} w)^2}{2\sigma^2} \right)$$

$$= \sum_{i=1}^n \left(\log \frac{1}{\sqrt{2\pi\sigma^2}} \right) - \sum_{i=1}^n \left(\frac{(y^{(i)} - x^{(i)T} w)^2}{2\sigma^2} \right) \quad \text{--- (1)}$$

→ Partially differentiating w.r.t to W

$$= 0 - \sum_{i=1}^n \frac{2(-x^{(i)})}{2\sigma^2} (y^{(i)} - x^{(i)T}W)$$

$$= 0 + \sum_{i=1}^n \frac{x^{(i)}(y^{(i)} - x^{(i)T}W)}{\sigma^2}$$

— equating it to zero

$$\sum_{i=1}^n x^{(i)} y^{(i)} = \sum_{i=1}^n x^{(i)} x^{(i)T} W = 0$$

$$\sum_{i=1}^n x^{(i)} y^{(i)} = \sum_{i=1}^n x^{(i)} x^{(i)T} W$$

$$W = \left(\sum_{i=1}^n x^{(i)} x^{(i)T} \right)^{-1} \left(\sum_{i=1}^n x^{(i)} y^{(i)} \right)$$

→ Partially differentiating eqn (1) w.r.t to σ^2

rewriting eqn (1) → $\sum_{i=1}^n \left(\log \frac{1}{\sqrt{2\pi\sigma^2}} \right) - \sum_{i=1}^n \left(\frac{(y^{(i)} - x^{(i)T}W)^2}{2\sigma^2} \right)$

$$\sum_{i=1}^n \left(\frac{-1}{2} \log(2\pi\sigma^2) \right) - \sum_{i=1}^n \left(\frac{(y^{(i)} - x^{(i)T}W)^2}{2\sigma^2} \right)$$

— now differentiating w.r.t σ^2

$$\sum_{i=1}^n \left(\frac{-1}{2} \times \frac{1}{2\pi\sigma^2} \right) - \sum_{i=1}^n \frac{1}{2} \left(\frac{(y^{(i)} - x^{(i)T}W)^2}{\sigma^4} \right) = 0$$

Now,

$$\sum_{i=1}^n \frac{1}{2} \frac{1}{\sigma^2} = \sum_{i=1}^n \frac{1}{2} \frac{(y^{(i)} - x^{(i)T}W)^2}{\sigma^4}$$

$$n = \sum_{i=1}^n \frac{(y^{(i)} - x^{(i)T}W)^2}{\sigma^2}$$

$$\therefore \sigma^2 = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - x^{(i)T}W)^2$$

End Of Assignment