

# Deep Learning HW #1

## Contents

<b>Problem 1:</b> Python and Numpy Warm-up Exercises .....	2
<b>Problem 2:</b> Linear Regression via Analytical Solution.....	2
a. Age Regressor: .....	2
<b>Problem 3:</b> Probability Distributions .....	3
a. Estimating the Parameters of a Probability Distribution: .....	3
b. Conditional Probability Distributions to Represent the Uncertainty of Functions:.....	4
<b>Problem 4:</b> Proofs and Derivation .....	4
a. Prove $\nabla_x(x^T a) = \nabla_x(a^T x) = a$ .....	4
b. Prove $\nabla_x(x^T Ax) = (A + A^T)x$ .....	7
c. Prove $\nabla_x(x^T Ax) = 2Ax$ .....	9
d. Prove $\nabla_x((Ax + b)^T * (Ax + b)) = 2A^T(Ax + b)$ .....	9

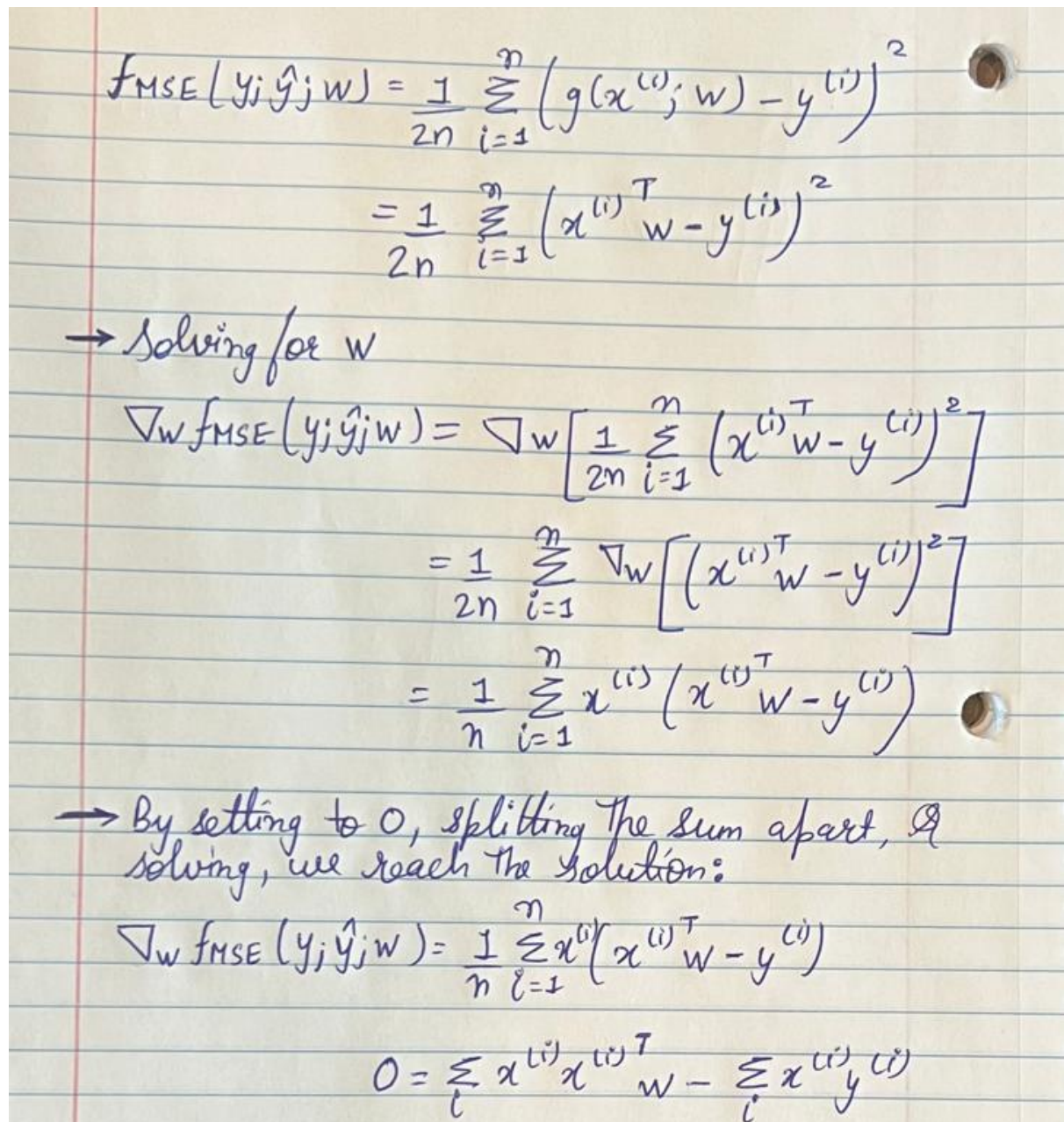
## **Problem 1:** Python and NumPy Warm-up Exercises

Please refer to ***homework1\_template.py*** for solutions of 1(a) to 1(n).

## **Problem 2:** Linear Regression via Analytical Solution

### **a. Age Regressor:**

- We have implemented linear regression via analytical solution using linear algebraic operations in NumPy. Find the methods *linear\_regression* and *train\_age\_regressor* in ***homework1\_template.py***.
  - Note: Before running the script, please include the datasets in the same folder.
- Analytical Solution to obtain the weights ( $w$ ):

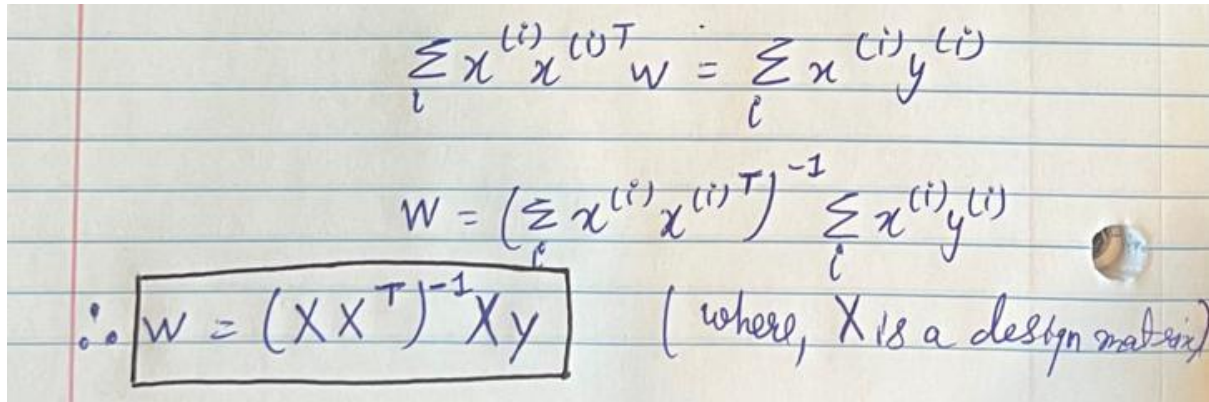

$$\begin{aligned} f_{\text{MSE}}(y; \hat{y}; w) &= \frac{1}{2n} \sum_{i=1}^n \left( g(x^{(i)}; w) - y^{(i)} \right)^2 \\ &= \frac{1}{2n} \sum_{i=1}^n \left( x^{(i)T} w - y^{(i)} \right)^2 \end{aligned}$$

→ Solving for  $w$

$$\begin{aligned} \nabla_w f_{\text{MSE}}(y; \hat{y}; w) &= \nabla_w \left[ \frac{1}{2n} \sum_{i=1}^n \left( x^{(i)T} w - y^{(i)} \right)^2 \right] \\ &= \frac{1}{2n} \sum_{i=1}^n \nabla_w \left[ \left( x^{(i)T} w - y^{(i)} \right)^2 \right] \\ &= \frac{1}{n} \sum_{i=1}^n x^{(i)} \left( x^{(i)T} w - y^{(i)} \right) \end{aligned}$$

→ By setting to 0, splitting the sum apart, & solving, we reach the solution:

$$\begin{aligned} \nabla_w f_{\text{MSE}}(y; \hat{y}; w) &= \frac{1}{n} \sum_{i=1}^n x^{(i)} \left( x^{(i)T} w - y^{(i)} \right) \\ 0 &= \sum_i x^{(i)} x^{(i)T} w - \sum_i x^{(i)} y^{(i)} \end{aligned}$$



$$\sum_i x^{(i)} x^{(i)T} W = \sum_i x^{(i)} y^{(i)}$$

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

$$\therefore W = (X X^T)^{-1} X y \quad (\text{where, } X \text{ is a design matrix})$$

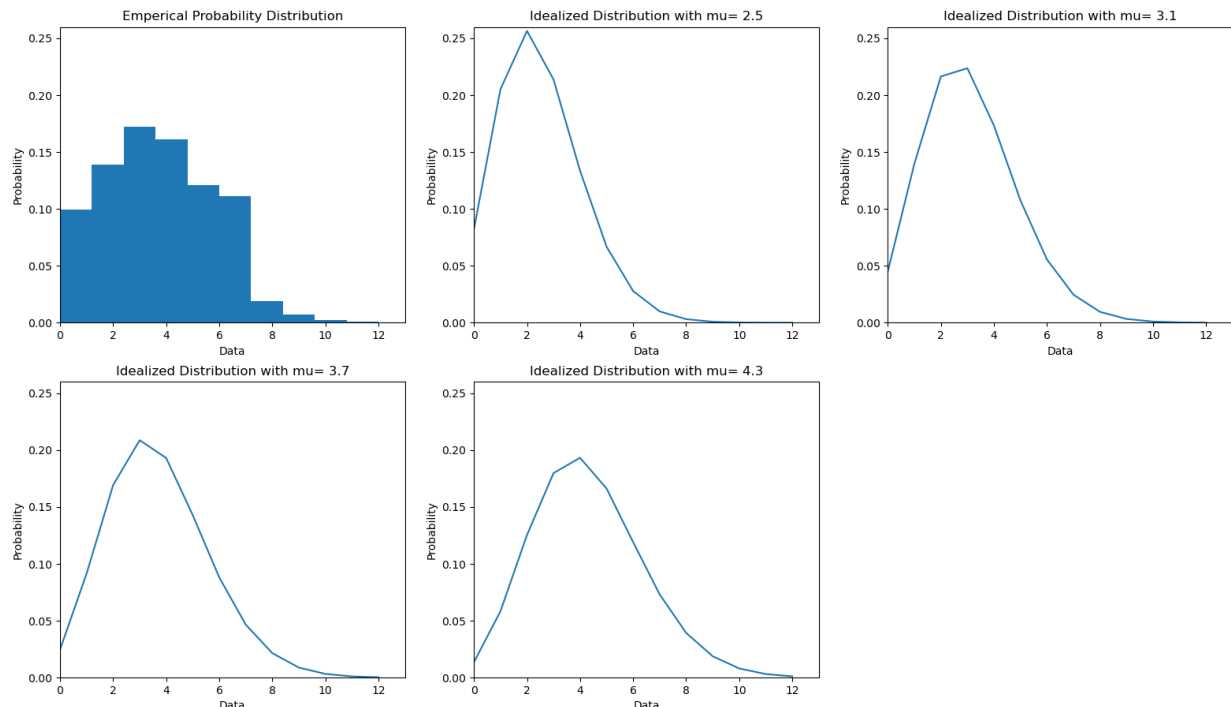
- After optimizing the weights on training set, we then computed the cost  $f_{MSE}$  on both testing and training dataset. The output is reported below:

```
Training Loss: 50.46582283534778
Testing Loss: 268.7869643405511
```

### **Problem 3: Probability Distributions**

#### **a. Estimating the Parameters of a Probability Distribution:**

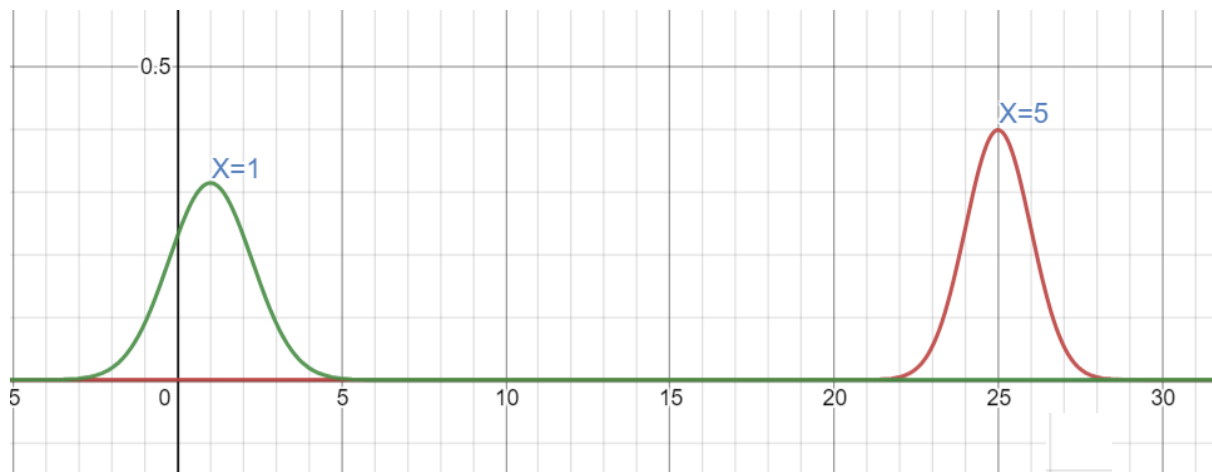
- We have plotted the empirical probability of the data in **PoissonX.npy**.
- Next, we plotted the probability distribution of a Poisson random variable with the following rate parameters: 2.5, 3.1, 3.7, and 4.3.
- The plots are as shown below:



- Based on observations the parameter value of  $\mu = 3.7$  is most consistent with the data.

**b. Conditional Probability Distributions to Represent the Uncertainty of Functions:**

- i. The corresponding value of  $y$  tend to be larger for the values of  $x$  with large magnitude as mean of the given normal distribution varies according to  $x^2$ . We can observe in the graph below that at  $x=5$  the mean of distribution is 25 therefore indicating that  $y$  would be large.
- ii. The uncertainty in the corresponding value of  $y$  tend to be larger for the values of  $x$  with small magnitude. We can observe in the graph below that at  $x=1$  the spread is more when compared to  $x=5$ . More spread indicates that the uncertainty is higher for lower value of  $x$ .



**Problem 4: Proofs and Derivation**

**a. Prove  $\nabla_x(x^T a) = \nabla_x(a^T x) = a$**

→ let  $x$  be a column vector  $x = [x_1, x_2, \dots, x_n]^T$   
 let  $a$  be a column vector  $a = [a_1, a_2, \dots, a_n]^T$   
 now,

$$x^T a = [x_1, x_2, x_3, \dots, x_n] \times \begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_n \end{bmatrix}$$

$$\Rightarrow x^T a = (a_1 x_1 + a_2 x_2 + a_3 x_3 + \dots + a_n x_n) \quad \text{--- (1)}$$



Similarly;

$$a^T x = [a_1, a_2, a_3, \dots, a_n] x \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$

$$\Rightarrow a^T x = (a_1 x_1 + a_2 x_2 + a_3 x_3 + \dots + a_n x_n) \quad \text{--- (2)}$$

now, we know that;

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

→ now, let us calculate,  $\frac{\partial (x^T a)}{\partial x_1}$

$$\frac{\partial (x^T a)}{\partial x_1} = a_1 + 0 + 0 + \dots + 0$$

Similarly

$$\frac{\partial (x^T a)}{\partial x_2} = 0 + a_2 + \dots + 0 \quad \& \quad \frac{\partial (x^T a)}{\partial x_n} = a_n$$

$$\therefore \nabla_x (x^T a) = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = a$$

similarly;

$$\frac{\partial (a^T x)}{\partial x_1} = a_1 + 0 + 0 \dots 0$$

$$\frac{\partial (a^T x)}{\partial x_2} = 0 + a_2 + 0 \dots 0$$

$$\frac{\partial (a^T x)}{\partial x_n} = a_n$$

$$\therefore \nabla_x (a^T x) = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = a$$

$$\therefore \boxed{\nabla_x (x^T a) = \nabla_x (a^T x) = a}$$



b. Prove  $\nabla_x(x^T A x) = (A + A^T)x$

→ Prove that:  $\nabla_x(x^T A x) = (A + A^T)x$

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & \dots & \dots & a_{nn} \end{bmatrix}$$

let us represent  $A$  as  $\rightarrow \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$  or  $[b_1 \ b_2 \ \dots \ b_n]$

where,  $a_1 \dots a_n$  are the rows of  $A$ .  
 where,  $b_1 \dots b_n$  are the columns of  $A$ .

so, we can represent  $Ax$  as  $\begin{bmatrix} a_1 x \\ a_2 x \\ \vdots \\ a_n x \end{bmatrix}$

→ now;

$$\begin{aligned} x^T A x &= x^T [a_1 x \ a_2 x \ \dots \ a_n x]^T \\ &= x_1 a_1 x + x_2 a_2 x \ \dots \ x_n a_n x \end{aligned}$$

where,  $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$

→ now, let us take derivative for one component  $x_k$ .

$$\frac{\partial (x^T A x)}{\partial x_k}$$



$$\frac{\partial (x^T A x)}{\partial x_k} = a_k x + a_{1k} x_1 + a_{2k} x_2 + \dots - \dots - a_{nk} x_n$$

$$= a_k x + \sum_{i=1}^n a_{ik} x_i$$

(rows  $\times$   $x$ )      (columns  $\times x$ )

$$= a_k x + (b_k)^T x$$

$$\therefore \nabla_x (x^T A x) = \begin{bmatrix} a_1 x + (b_1)^T x \\ a_2 x + (b_2)^T x \\ \vdots \\ a_n x + (b_n)^T x \end{bmatrix}$$

$$\therefore \begin{bmatrix} a_1 x \\ a_2 x \\ \vdots \\ a_n x \end{bmatrix} + \begin{bmatrix} (b_1)^T x \\ (b_2)^T x \\ \vdots \\ (b_n)^T x \end{bmatrix}$$

$$\therefore \boxed{\nabla_x (x^T A x) = A x + A^T x} \\ \boxed{= (A + A^T) x}$$



c. Prove  $\nabla_x(x^T A x) = 2Ax$

Given  $A$  is symmetric  $n \times n$  matrix  
 $\therefore A = A^T$   
 $\therefore$  from (b) we know  $\nabla_x(x^T A x) = (A + A^T)x$   
 $\therefore A = A^T$   
 $\therefore \nabla_x(x^T A x) = 2A^T x = 2Ax$   
 Hence Proved //

d. Prove  $\nabla_x((Ax + b)^T * (Ax + b)) = 2A^T(Ax + b)$

LHS  $\therefore \nabla_x \left[ (x^T A^T + b^T)(Ax + b) \right]$   
 $\Rightarrow \nabla_x \left[ x^T (A^T A)x + x^T (A^T b) + (b^T A)x + b^T b \right]$   
 Let us consider  $A^T A = K$  ( $n \times n$  matrix)  
 $A^T b = a \rightarrow$  vector ( $n \times 1$ )  
 $b^T A = a^T$   
 $\Rightarrow \nabla_x \left[ x^T K x + x^T a + a^T x + b^T b \right]$   
 $\therefore \nabla_x b^T b = 0, \nabla_x(x^T A x) = 2Ax, \nabla_x(x^T a) = \nabla_x(a^T x) = a$   
 $\Rightarrow \nabla_x 2Kx + a + a + 0$   
 $\Rightarrow 2A^T A x + 2A^T b \Rightarrow 2A^T (Ax + b) = \text{RHS} //$

**End Of Assignment**