

## Theoretical exercises

### Q2-

1- MSE loss function: *Mean Squared Error* is a loss function which is widely used in regression tasks, it calculates the mean difference between predicted and the actual value, and the domain of numbers is continuous in real number space. The formula of MSE is:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

It calculates the square of the difference between actual value( $Y_i$ ) and predicted value( $\hat{Y}_i$ ) and takes the mean of the sum of squares. It tells us how wrong our model is predicting our data. So in the task of linear regression we are trying to minimize the MSE using an optimizer algorithm which is gradient-descent in this case most of the time.

2- Cross-Entropy: Cross-entropy loss measures the divergence of two distributions. It is widely used in classification tasks, in which the model output is probability of an input belonging to a class, in case of binary classification there are two classes(0 or 1). So the output of the model is the probability of each input, and also a threshold like 0.5 which if  $P(x) > 0.5$  belongs to class 1 and otherwise belongs to 0. The formula of binary cross entropy is:

$$-\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1-y_i) \cdot \log(1-p(y_i))$$

When the output is 0:  $y_i$  is 0 so we have  $(1-y_i)*\log(1-p(y_i))$

When the output is 1:  $(1-y_i)$  is 0 so we have  $y_i * \log(p(y_i))$

The task is minimizing cross entropy loss, which means that we want our prediction distribution to get close enough to the actual distribution.

“For both binary and multi-class classification problems, the main goal is to minimize the cross entropy loss, which in turn maximizes the *likelihood* of assigning the correct class labels to the input data points”

### Conclusion:

According to the above explanation of MSE and Cross-Entropy loss, the main difference between them is MSE is for the models whose outputs are in real number space, like linear regression, and it calculates the *mean difference* between actual value and model prediction. Cross-entropy is for classification tasks in which the output of a model is probability of an instance  $x_i$  belonging to one class, so it calculates the *divergence* of two distributions. Then we use Cross-entropy loss for classification problems.

#### Q4-

The activation functions are used in neural networks for determining which neuron should fire(being active) and which should not, and another purpose of activation functions is to add *non-linearity* to models for dealing with non-linear relations in data for more advanced tasks like images and NLP etc. *Sigmoid* and *ReLU* are two such functions.

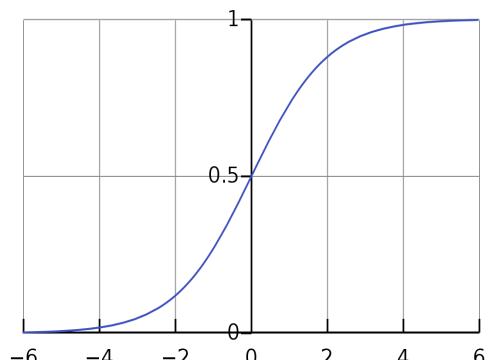
##### **Sigmoid:**

The sigmoid activation function maps weighted sum to a bounded space between 0 and 1. The formula of sigmoid is:

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

Limit of  $f(x)$  to  $+\infty$  is 1 and to  $-\infty$  is 0. The plot of sigmoid is:

It's usually used in classification problem because it maps Continuous input to a probability output between 0 and 1.



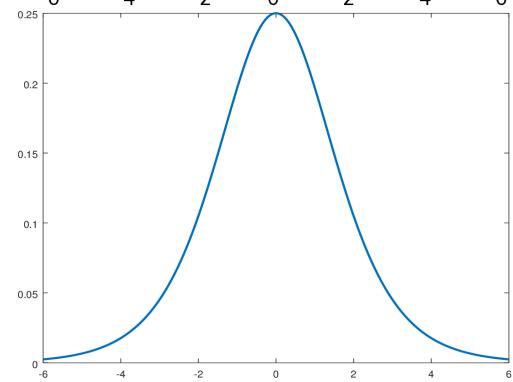
##### **Advantages of Sigmoid:**

1- The output of sigmoid functions is bounded and prevents large or unstable outputs during training.

##### **Disadvantages of Sigmoid:**

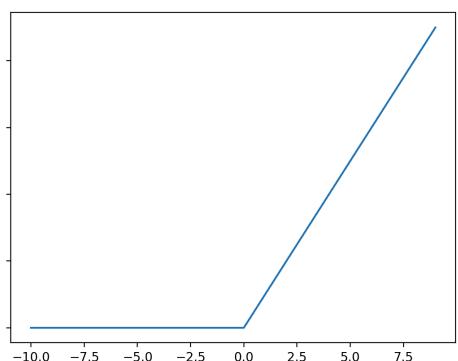
1- Gradient vanishing: “The gradient vanishing problem occurs during training when the gradients of the loss functions become really small as they backpropagate through different layers. This causes slow learning due to less updation of the weights of early layers.”

2- computationally expensive since we have exponential term. If we plot the gradient of sigmoid function, we can visualize The gradient of very negative and very positive numbers Is almost 0 and it cause the slow learning(Gradient vanishing)



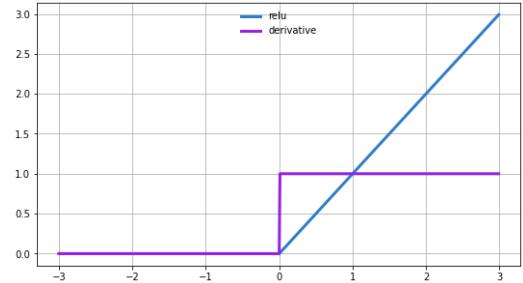
**ReLU:** this activation function treats the weighted sum of neurons as an identity function if it is positive and 0 otherwise.

$$f(x) = \max(0, x)$$



**Advantages of ReLU:** ReLU could overcome both disadvantages of sigmoid function, first we don't have gradient vanishing problem since the gradient of ReLU is 1 or positive numbers and 0 otherwise, and second its more computationally efficient than sigmoid, because it's simply discard negative values and doesn't have exponential term. The ReLU is mostly used in deep learning architecture in hidden layers and it makes faster convergence in networks because the backpropagation becomes very efficient and fast and we don't have a learning slowdown problem.

**Disadvantages of ReLU:** As we said the gradient of ReLU maps all negative numbers to 0, the ReLU prone to causing “dying”. Dying of a neuron means the output is always 0 and it happens when the weighted sum of the neuron becomes negative. In these cases we can use variants of ReLU like “Leaky ReLU”.



### Q3-

In the process of training a neural network, we must have a point to start, so we need to manually initialize the weights. We have several options like setting all weights to 0, *Random initialization* and other techniques which we discuss further.

- 1- One of the issues of neural networks is “symmetry” of weights, which means that all the weights become equal and all neurons learn the same thing about data and it may be stuck in a local minima! so we can't start with setting all weights to 0!!
- 2- Avoiding vanishing and exploding of gradient: if we initialize weights with very small numbers, it may cause gradient vanishing and also if we set them with very large numbers, it causes the gradient to explode!
- 3- Avoiding saturation of activation functions: Some activation functions like sigmoid saturates with very large positive and large negative numbers.
- 4- Better generalization and Faster convergence: Helps the model to start with a reasonable set of parameters close to optimal value and also the optimization algorithm can more quickly lead to optimal solution.

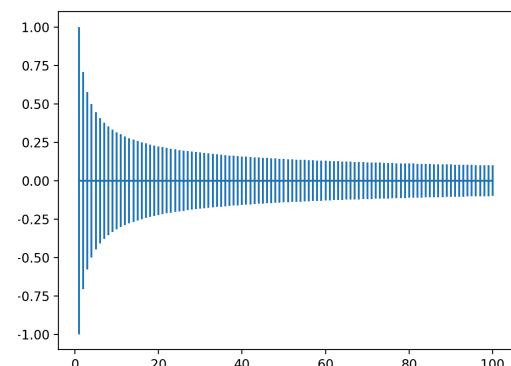
#### Weight initialization techniques:

##### 1- Random initialization:

In this technique we initialize weight randomly from a specific distribution, this technique first introduced to solve the problem of 0 initialization or “Symmetry”.

##### 2- Xavier:

In this method we calculate weights with random numbers from a Uniform distribution with mean = 0 and range =  $[-1 / \sqrt{n_{\text{input}}}, 1 / \sqrt{n_{\text{input}}}]$  and  $\text{std} = \sqrt{2 / (n_{\text{input}} + n_{\text{output}})}$ . which  $n_{\text{input}}$  is the number of input neurons and  $n_{\text{output}}$  for output neurons. It is mainly used with sigmoid and Tanh activation functions. So the spread of weights is dependent on the number of input neurons, which we can see in the plot. So with a greater number of input neurons, We have smaller weights. “The rationale behind this scaling is to prevent the activations and gradients from becoming too large or too small as they propagate through the network.”



If the weights are too large, it can lead to exploding gradients, causing training instability. Conversely, if the weights are too small, it can lead to vanishing gradients, making it difficult for the network to learn effectively.”

### **He method:**

This method is used with the ReLU activation function. We initialize weights by calculating random numbers from a Gaussian distribution with mean = 0 and std =  $\sqrt{2 / n_{\text{input}}}$ . When using ReLU activation functions, it's important to initialize the weights properly to prevent the issue of dead neurons, where neurons in the network can become stuck in a state where they always output zero due to negative weights. If weights are initialized too small, it's possible for a large portion of neurons to never activate during training, effectively rendering them useless. Compared with the Xavier method, which considers both input and output neurons, the He method only considers inputs, because ReLU ignores the negative values and almost halves the gradient, it must be initialized with weights with higher variance to counteract this effect.

### **Potential issues with symmetric weights:**

Redundancy and lack of diversity, as the output of all neurons in the layer is the same! It may cause Convergence issues, since we can get stuck in suboptimal points. And the solution for these issues is using a proper weight initialization technique which we discussed above.

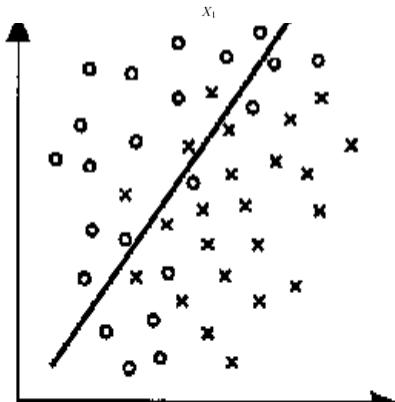
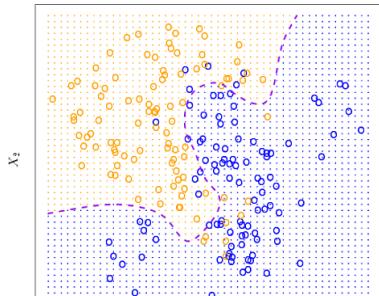
## **Q1-**

Let us first start with explaining the *Logistic regression* and *Perceptron* individually.

**Note:** We assume a Perceptron is a single neuron with sign(step) function as activation function.

### **1- Logistic regression:**

Logistic regression is a machine learning algorithm which is used mostly in Binary classification tasks, it takes the net input and calculates the probability of an instance belonging to a class which means the output is a continuous number between 0 and 1. The activation function of the model is sigmoid,(introduced in Q4). The loss function which is used in this model is Binary cross entropy(introduced in Q2), which calculates the divergence of two distributions, which is quite suitable for the task. So by minimizing the BCE or maximizing likelihood, we can train the model on the data, which is done by an optimization algorithm(like SGD or Adam etc). Since the sigmoid function is a non-linear activation function, the decision boundary of Logistic regression is non-linear, so it could solve non-linear separable data as well(but not very complex!). We can see in the figure below. Also the logistic regression could solve the multi-class classification with appropriate modification(one-vs-all, softmax). And due to the probabilistic nature of the output and loss module, it's less sensitive to outliers(if the skewness of the data isn't very high).



### **Perceptron:**

It's a single layer neural network which takes the weighted sum of input and feeds it to a sign(step) function to make a binary output 0 or 1(with a threshold). So at the first step, the

output space is linear! So it can not solve the nonlinear problems. The loss function used in this algorithm is typically MSE(Mean Squared Error introduced in Q2), which measures the square error of the model output with actual value, which is sensitive to outliers! And it learns with an optimization technique like SGD or Adam.

### **Differences and use cases:**

Aspect	Perceptron	Logistic Regression
Model Type	Single-layer neural network with a binary activation function.	Linear regression model with a logistic (sigmoid) activation function.
Activation Function	Uses a step function (binary threshold) as the activation function.	Employs the logistic (sigmoid) function as the activation function.
Output	Produces binary output (0 or 1) directly based on the weighted sum of inputs.	Outputs continuous values between 0 and 1, representing probabilities.
Learning Algorithm	Uses a simple update rule to adjust weights based on misclassified examples (Perceptron learning rule).	Utilizes gradient descent or other optimization algorithms to find the best weights by minimizing a cost function (typically cross-entropy).
Decision Boundary	Forms a linear decision boundary, suitable for linearly separable data.	Can create non-linear decision boundaries by using higher-order features or polynomial terms.
Classification Task	Typically used for binary classification tasks.	Can handle both binary and multi-class classification problems with appropriate modifications (One-vs-All, Softmax).
Convergence	Guaranteed to converge if the data is linearly separable.	Converges for most datasets, even when they aren't perfectly linearly separable.

Robustness	Sensitive to outliers and noisy data due to its update rule.	More robust to outliers and noise due to probabilistic output.
Probabilistic Output	Does not provide probabilistic outputs; it's a deterministic model.	Provides probabilistic outputs, which can be interpreted as class probabilities.
Use Cases	Simple and suitable for basic classification tasks with linearly separable data.	Widely used in various classification problems, including those with non-linear decision boundaries.

## Q6-

Q6-

A) Feed-forward and calculating loss:

1- calculate net input of 1st layer:

input  $h_1$ :  $i_1 \times w_1 + i_2 \times w_2 + b_1 = 0.05 \times 0.15 + 0.1 \times 0.25 + 0.35 = 0.88$

input  $h_2$ :  $i_1 \times w_3 + i_2 \times w_4 + b_1 = 0.05 \times 0.2 + 0.1 \times 0.3 + 0.25 = 0.79$

2- calculate output of first layer:

out  $h_1: \sigma(h_1) = \frac{1}{1+e^{-0.88}} = 0.9938 \approx 0.99$

out  $h_2: \sigma(h_2) = \frac{1}{1+e^{-0.79}} = 0.9962 \approx 0.99$

3- calculate input of output layer:

input  $o_1 = \text{out } h_1 \times w_5 + \text{out } h_2 \times w_6 + b_2 = 0.99 \times 0.4 + 0.6 \times 0.5 + 0.6 = 1.13$

input  $o_2 = \text{out } h_1 \times w_7 + \text{out } h_2 \times w_8 + b_2 = 0.99 \times 0.4 + 0.6 \times 0.5 + 0.6 = 1.13$

4- calculate output of output layer:

out  $o_1 = \sigma(o_1) = \frac{1}{1+e^{-1.13}} = 0.7558 \approx 0.76$

out  $o_2 = \sigma(o_2) = \frac{1}{1+e^{-1.13}} = 0.7607 \approx 0.77$

5- calculating total error:

Error of  $o_1 = E_{o_1} = \frac{1}{2} (\text{target}_{o_1} - \text{out}_{o_1})^2$  and similar for  $o_2$

$E_{o_1} = \frac{1}{2} (0.99 - 0.76)^2 = 0.1089$

$E_{o_2} = \frac{1}{2} (0.99 - 0.77)^2 = 0.04$

$E_{\text{total}} = E_{o_1} + E_{o_2} = 0.1089 + 0.04 = 0.1489$

### B) - Back Propagation:

1. Updating  $w_5, w_6, w_7, w_8$ :  $\alpha = 0.3$

$$\frac{\partial E_t}{\partial w_i} = \frac{\partial E_t}{\partial \text{out}_{o_j}} \frac{\partial \text{out}_{o_j}}{\partial \text{in}_{o_j}} \cdot \frac{\partial \text{in}_{o_j}}{\partial w_i} \quad \text{which } o_j \text{ is the neuron } w_i \text{ connected to.}$$

$$\text{for } w_5: \frac{\partial E_t}{\partial w_5} = \frac{\partial E_t}{\partial \text{out}_{o_1}} \frac{\partial \text{out}_{o_1}}{\partial \text{out}_{o_1}} \frac{\partial \text{in}_{o_1}}{\partial w_5}$$

$$\frac{\partial E_t}{\partial w_5} \rightarrow -(target_{o_1} - \text{out}_{o_1}) = -(0.101 - 0.76) = 0.175$$

$$\frac{\partial \text{out}_{o_1}}{\partial w_5} \rightarrow \text{out}_{o_1} = \sigma(0.76) \Rightarrow \text{out}_{o_1}(1 - \text{out}_{o_1}) = 0.76(0.24) = 0.18$$

$$\frac{\partial \text{in}_{o_1}}{\partial w_5} \rightarrow \text{in}_{o_1} = w_5 \text{out}_{h_1} + w_6 \text{out}_{h_2} + b_2 \Rightarrow \frac{\partial \text{in}_{o_1}}{\partial w_5} = \text{out}_{h_1} = 0.59$$

$$\frac{\partial E_t}{\partial w_5} = 0.175 \times 0.18 \times 0.59 = 0.07$$

$$w'_5 = w_5 - \alpha \frac{\partial E_t}{\partial w_5} = 0.4 - 0.07 = 0.38$$

+ for avoiding all this extra calculation for other weights we define a delta factor:

$$\delta_{o_j} = \frac{\partial E_t}{\partial \text{out}_{o_j}} \frac{\partial \text{out}_{o_j}}{\partial \text{net}_{o_j}} \Rightarrow \delta_{o_j} = -(target_{o_j} - \text{out}_{o_j}) \text{out}_{o_j}(1 - \text{out}_{o_j})$$

$$\frac{\partial E_t}{\partial w_i} = \delta_{o_j} \text{out}_{h_k} \quad (o_j \text{ is } \underset{\text{neuron}}{\cancel{\text{out}}} \text{ one in forward layer } h_k \text{ in "back" layer})$$

$$w'_6 = \delta_{o_2} \text{out}_{h_1} = -(0.99 - 0.77) \times 0.77(0.23) \times 0.59 = -0.02$$

$$w'_6 = 0.45 - (0.3 \times (-0.02)) = 0.44$$

$$w'_7 = \delta_{o_1} \text{out}_{h_2} = -(0.01 - 0.76) \times 0.76(0.24) \times 0.6 = \cancel{0.08} 0.08$$

$$w'_7 = 0.50 - (0.3 \times 0.08) = 0.48$$

$$w'_8 = \delta_{o_2} \text{out}_{h_2} = 0.55 + 0.06 = 0.556$$

2- Updating  $w_1, w_2, w_3, w_4$ ,

using delta rule to make it easier:

$$\frac{\delta E_L}{\delta w_i} = \left( \sum_0 \frac{\delta E_L}{\delta o_{out}} \times \frac{\delta o_{out}}{\delta h_i} \times \frac{\delta h_i}{\delta w_i} \right) \alpha \frac{\delta o_{out} h_k}{\delta h_k} \times \frac{\delta h_k}{\delta w_i}$$

$$= (\sum_0 \delta_0 \times w_{h_0}) \alpha o_{out} h_k (1 - o_{out} h_k) \times i_i;$$

$$w'_1 := w_1 - \alpha \left( \sum_0 \delta_0 w_{h_0} \right) \alpha o_{out} h_k (1 - o_{out} h_k) \times i_i; \quad \delta h_0 = 0/009$$

$$= w_1 - 0/3 \left( (\delta_{01} \alpha w_1 + \delta_{02} \alpha w_2) \alpha 0/59 (0/41) \alpha 0/09 \right) =$$

$$\frac{\delta o_1 = 0/13 \quad \delta o_2 = 0/03}{\delta h_2 = 0/09} \rightarrow w_1 - (0/3 \alpha 0/0008) = 0/5 - 0/0002 = 0/498$$

$$w'_2 := w_2 - \alpha \left( (\delta_{01} w_1 + \delta_{02} w_2) \alpha o_{out} h_2 (1 - o_{out} h_2) \times i_2 \right) =$$

$$w_2 = w_2 - (0/3 \alpha 0/0004) = 0/2 - 0/0001 = 0/1999$$

$$w'_3 := w_3 - \alpha \left( (\delta_{01} w_1 + \delta_{02} w_2) \alpha o_{out} h_3 (1 - o_{out} h_3) \times i_3 \right)$$

$$= w_3 - (0/3 \alpha 0/0009) = 0/25 - 0/0002 = 0/2498$$

$$w'_4 = w_4 - \alpha \left( (\delta_{01} w_1 + \delta_{02} w_2) \alpha o_{out} h_4 (1 - o_{out} h_4) \times i_4 \right) =$$

$$w_4 = 0/3 - (0/3 \alpha 0/0009) = 0/3 - 0/0002 = 0/2998$$

3- Update biases:

by using delta rule :  $b' = b - \alpha \delta_k^l$  layer  $k$  neuron  $k$ .

$$b'_{201} = b_{201} - \delta_{01} = 0/5 - (0/13 \times 0/8) = 0/57$$

$$b'_{202} = b_{202} - \delta_{02} = 0/5 - ((-0/03) \times 0/9) = 0/59$$

$$b'_{1h_1} = b_{1h_1} - \delta_{h_1} = 0/35 - (0/3 \times 0/09) = 0/34$$

$$b'_{1h_2} = b_{1h_2} - \delta_{h_2} = 0/35 - (0/3 \times 0/09) = 0/34$$

## Q5-

In the realm of neuron connectivity, two contrasting concepts prevail: Sparse and Dense connectivity.

Within the domain of neural networks, a dense layer manifests as a network layer wherein every neuron establishes connections with every other neuron in the preceding layer.

Conversely, in a sparse layer, each neuron selectively connects only to a subset of neurons from the previous layer.

The primary allure of sparse layers lies in their computational efficiency, coupled with their adeptness at curbing overfitting and reducing overall complexity. These layers tend to excel in domains such as Natural Language Processing (NLP), where input data, often text, is inherently noisy.

*Pruning* stands as a viable option to transition a densely connected neural network into a sparser form.

**Optimal Brain Damage (OBD):** OBD aims to minimize the network's complexity by identifying and eliminating redundant weights. It involves evaluating the second-order derivatives of the loss function, typically represented as the Hessian matrix, with respect to each weight. Mathematically, this can be expressed as:

$$\left[ \frac{\partial^2 \mathcal{L}}{\partial w_i \partial w_j} \right]$$

where  $(\mathcal{L})$  denotes the loss function and  $(w_i)$  and  $(w_j)$  represent individual weights.

Weights with second derivatives close to zero are deemed insignificant and are pruned from the network.

**Sensitivity-Based Pruning:** Sensitivity-based pruning involves quantifying the sensitivity of the network's output to changes in individual weights. This is typically accomplished through gradient calculations using techniques like backpropagation. Mathematically, the sensitivity of the output  $y$  with respect to a weight :

$$\left[ \frac{\partial y}{\partial w_i} \right]$$

**Magnitude-Based Pruning:** Magnitude-based pruning ranks weights based on their magnitudes and selectively removes those below a predefined threshold. Mathematically, this process involves comparing the absolute values of individual weights:

$$[|w_i| < \text{threshold}]$$

where weights below the threshold are pruned, reducing the network's overall complexity.