

HW1

Submitted By: Parthiv Borgohain (pb25347)

Note: Some of the questions have been answered in the attached ipynb file. Please refer to that file as well.

Question 1. In this question, we will define tensors in PyTorch and use them to find the derivative of a linear function w.r.t its variables. We then find the mean squared error and conclude with finding the optimum parameters of a linear model.

1.1 For the function $y = f(x) = w \cdot x + b$, where $w=[2,1]$ and $b=3$, find the partial derivatives of y w.r.t the components of x (i.e.: dy/dx_1 and dy/dx_2) at $x = [4,2]$ both on paper and using PyTorch.

Answer:

HW1

Q1 1.1 $y = f(x) = w^T x + b$

Given $b = 3$
 $w = [2, 1]$

$$\therefore \frac{\partial y}{\partial x_1} = \frac{\partial (2x_1 + x_2 + 3)}{\partial x_1}$$
$$= 2$$
$$\therefore \frac{\partial y}{\partial x_2} = \frac{\partial (2x_1 + x_2 + 3)}{\partial x_2}$$
$$= 1.$$

$$\therefore, \frac{\partial y}{\partial x_1} = 2, \frac{\partial y}{\partial x_2} = 1$$

So, the gradients w.r.t x_1 and x_2 are 2 and 1 respectively.

1.2 For a model $y = f(x)$, the predicted and true values are as follows:

$$y_{\text{true}} = [0, 1, 1, 0]$$

$y_{\text{pred}} = [0.1, 0.95, 1.1, 0.2]$

Find the mean squared error both on paper and using PyTorch. In PyTorch, solve it using an inbuilt function and also by defining your own squared error function.

Answer:

Handwritten calculation of Mean Squared Error (MSE):

$$\begin{aligned} \text{1.2} \quad & \text{Given} \\ & y_{\text{true}} = [0, 1, 1, 0] \\ & y_{\text{pred}} = [0.1, 0.95, 1.1, 0.2] \\ \text{So} \quad & \text{MSE} = \frac{1}{4} \left[(0 - 0.1)^2 + (1 - 0.95)^2 + (1 - 1.1)^2 + (0 - 0.2)^2 \right] \\ & = \frac{1}{4} [0.01 + 0.0025 + 0.01 + 0.04] \\ & = \frac{0.0625}{4} \\ & \Rightarrow \text{MSE} = 0.015625 \end{aligned}$$

So, the MSE is 0.015625

1.3 You are given the following dataset (note that it describes an XOR function):

x1	x2	y
0	0	0
0	1	1
1	0	1
1	1	0

Assume that we fit a linear function, $y = w \cdot x + b$ to this dataset with $w = [2, 1]$ and $b = 3$.

- Find the mean squared error (loss) over this dataset for the above weights and bias.
- Which direction should the weights move in to decrease the loss by maximum amount? Find that direction both on paper and using PyTorch.
- For what values of w and b is the loss minimum? Solve on paper only.

Answer:

i)

1.3

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

i)

Gegeben: $y = w \cdot x + b$, $w = [2 \ 3]$, $b = 3$

so, $y = w \cdot x + b$

$$\Rightarrow y = 2x_1 + 3x_2 + 3$$

$$\text{so, } MSE = \frac{1}{4} \left[(3-0)^2 + (4-1)^2 + (5-1)^2 + (6-0)^2 \right]$$

$$= \frac{1}{4} [9 + 9 + 16 + 36]$$

$$\Rightarrow MSE = 17.5$$

ii)

ii)

$$MSE = \frac{1}{4} \left[(0-0-0-h)^2 + (1-0-w_2-h)^2 \right. \\ \left. + (1-w_1-0-h)^2 + (0-w_1-w_2-h)^2 \right]$$

$$\Rightarrow MSE = \frac{h^2 + (-w_2-h)^2 + (-w_2-h)^2 + (w_1+w_2+h)^2}{4}$$

$$\Rightarrow \frac{\partial(MSE)}{\partial w_1} = \frac{1}{4} \left[-2(1-w_1-h) \right. \\ \left. -2(-w_1-w_2-h) \right]$$

$$\boxed{\frac{\partial(MSE)}{\partial w_1} = \frac{4w_1 + 2w_2 + 4h - 2}{4}}$$

Similarly,

$$\frac{\partial(MSE)}{\partial w_2} = \frac{1}{4} \left[-2(1-w_2-h) -2(w_1+w_2+h) \right]$$

$$\boxed{\frac{\partial(MSE)}{\partial w_2} = \frac{4w_2 + 2w_1 + 4h - 2}{4}}$$

Similarly,

$$\frac{\partial(MSE)}{\partial h} = \frac{1}{4} \left[2h - 2(1-w_2-h) - 2(1-w_1-h) \right. \\ \left. - 2(-w_1-w_2-h) \right]$$

$$\boxed{\Rightarrow \frac{\partial(MSE)}{\partial h} = \frac{4w_2 + 4w_1 + 8h - 4}{4}}$$

\therefore For, $w = [791]$ and $h = 39$

the directions are -

$$\left[\frac{hw_1 + 2w_2 + 4h - 2}{4}, \frac{4w_2 + 2w_1 + 4h - 2}{4}, \right.$$

$$\left. \frac{4w_2 + 4w_1 + 8h - 4}{4} \right]$$

$$\therefore [5, 4, 5, 8]$$

→ This is the
direction vector

The direction vector is $[5, 4, 5, 8]$

iii)

(iii) To find minimum loss, we need to equate the partial derivative to 0

$$\therefore \frac{\partial (MSE)}{\partial w_1} = \frac{4w_1 + 2w_2 + 4b - 2}{4}$$

$$\Rightarrow 4w_1 + 2w_2 + 4b - 2 = 0$$

\rightarrow (a)

Similarly,

$$\frac{\partial (MSE)}{\partial w_2} = \frac{4w_2 + 2w_1 + 4b - 2}{4}$$

$$\Rightarrow 4w_2 + 2w_1 + 4b - 2 = 0$$

\rightarrow (b)

Similarly,

$$\frac{\partial (MSE)}{\partial b} = \frac{1}{4} [4w_2 + 4w_1 + 8b - 4]$$

$$\Rightarrow 4w_2 + 4w_1 + 8b - 4 = 0$$

\rightarrow (c)

\therefore we now need to solve system of equations (a), (b), (c).

Solving it using Python we get -

$$\begin{cases} w_1 = 0 \\ w_2 = 0 \\ b = 1/2 \end{cases}$$

values for which loss is least

So, the values for which the loss is minimum are:

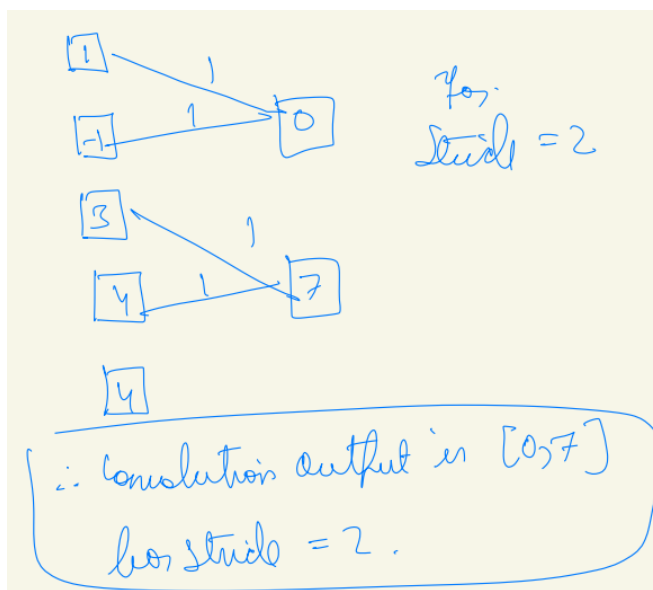
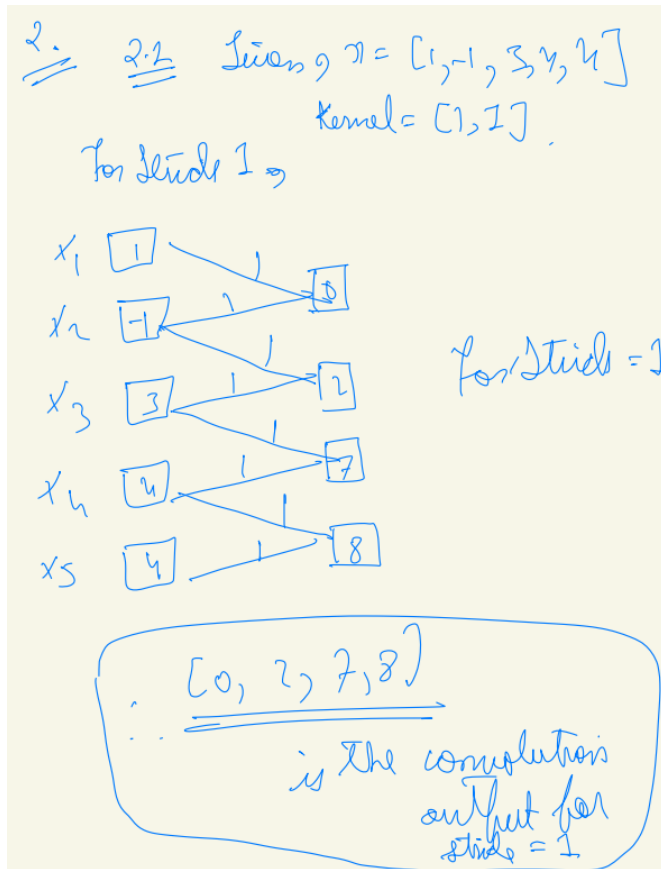
- The weights are $w_1=0, w_2=0$

- The bias is $\frac{1}{2}$ or 0.50

Q2. In this question, we will learn to perform 1-d and 2-d convolutions using different strides on a grayscale image.

2.1 Given a 1-d input $x=[1,-1,3,4,4]$ and a kernel $= [1,1]$, find the 1-d convolution for stride=1 and stride=2. Solve both on paper and using pyTorch.

Answer:



So, the convolution outputs for both the channels are [0,2,7,8] , [0,7]

2.2 You're given the following grayscale image (any matrix in 2 dimensions is a grayscale image):

0.1	-0.6	0.4	0.8
-0.4	0.3	0.9	0.2
0.5	0.2	0.8	-0.7
0.3	0.7	-0.4	0.1

i) You're given 2 2*2 filters (kernels): $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ and $\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$. Find the output when the image is convolved using each filter with stride=1. Solve both on paper and by using pyTorch. Note that you'll get 2 outputs - one for each filter - each output is known as a channel.

ii) What are the dimensions of each channel in (i)? What will be the dimensions of the output when stride=2?

Answer:

2.2

0.1	-0.6	0.4	0.8
-0.4	0.3	0.9	0.2
0.5	0.2	0.8	-0.7
0.3	0.7	-0.4	0.1

i) Filters are $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ and $\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$

For Kernel 1

$$\begin{bmatrix} (0.1 \times 1 + 0.3 \times 1), & (-0.6 + 0.9), & (0.4 + 0.2) \\ (-0.4 + 0.2), & (0.3 + 0.8), & (0.9 - 0.7) \\ (0.5 + 0.7), & (0.2 - 0.4), & (0.8 + 0.1) \end{bmatrix}$$
$$= \begin{bmatrix} 0.4 & 0.3 & 0.6 \\ -0.2 & 1.1 & 0.2 \\ 1.2 & -0.2 & 0.9 \end{bmatrix}$$

Similarly, for Kernel 2,

$$\text{Output} = \begin{bmatrix} (-0.6+0.4) & (0.4+0.3) & (0.8+0.9) \\ (0.3+0.5) & (0.9+0.2) & (0.2+0.8) \\ (0.2+0.3) & (0.8+0.7) & (-0.7+0.4) \end{bmatrix}$$

$$\Rightarrow \text{Output} = \begin{bmatrix} -1.0 & 0.7 & 1.7 \\ 0.8 & 1.1 & 1.0 \\ 0.5 & 1.5 & -1.1 \end{bmatrix}$$

(ii) The dimensions of each channel in the output in part(i) are 3x3

If stride = 2 \rightarrow then the dimension of each channel will be 2x2

iii) The given image is a 2-d matrix. How can you convolve it so that the output channel has only one dimension?

iv) Perform a 2*2 max-pooling with stride=1 on the image both on paper and using PyTorch. What are the dimensions of the channel after max-pooling?

Answer:

(iii) like convs so using using a kernel of the same size as the image.

e.g for the image given in this question, we will get an output of only 1 dimension

(iv) Max Pooling ^{output} with 2x2 kernel

and stride 1 = $\begin{bmatrix} 0.3 & 0.9 & 0.9 \\ 0.5 & 0.9 & 0.9 \\ 0.7 & 0.8 & 0.8 \end{bmatrix}$

by taking max values

Extra-credit: If the size of the output channel after performing a 2*2 convolution (stride=1) and a 2*2 max-pooling (stride=1) is the same, why do you think we needed max-pooling when we could perform convolution to decrease the size of the image. Think in terms of the advantages that max-pooling may offer over convolution. When do you think a max-pooling operation may not be advantageous (and in fact may hurt the network)?

Answer:

Max pooling is a widely used operation in convolutional neural networks (CNNs) that aids in the reduction of the input's spatial dimensions. This reduction can significantly reduce the model's computational cost and prevent overfitting. Max pooling does not introduce any trainable parameters during training, which means that it does not add any complexity to the model. Conversely, convolutional layers possess trainable parameters known as convolutional kernels, which are used to learn crucial features from the input data. During training, the kernel weights are modified using backpropagation to minimize the loss function. This process can add complexity to the model and increase the computational cost of training.

While both convolution and max pooling are effective dimensionality reduction techniques, they serve distinct purposes. Convolution captures all the information present in the input, whereas max pooling selects the most significant information. Max pooling has its advantages, such as extracting

the most vital features from the input and discarding less important ones. It also introduces some level of invariance to small translations or rotations in the input, which can be beneficial for tasks like object recognition.

However, the use of the max-pooling operation may not be beneficial in certain situations, such as when the input image is small, the task requires precise localization, or the features are already well-extracted. In such cases, using max pooling can result in the loss of crucial information and important details that could be vital to the task at hand. For instance, in image recognition tasks where the input image is small, max pooling may cause the image to become too small and lose vital features. Similarly, in tasks that require precise localization, max pooling may cause the loss of spatial information, making it difficult to pinpoint the exact location of the object of interest.

Question 3 This question deals with the addition of perturbations to an image x to create an adversarial example from that image.

3.1 Identify True/False i) In FGSM method, the parameters of the trained model change.

ii) In FGSM method, the pixel values of input image changes.

iii) FGSM can only be used for undirected adversarial attacks (An undirected adversarial attack is one in which the aim is to only perturb the original image in the direction of maximum loss. The attack does not care about which incorrect class the input example is classified into after perturbation).

iv) PGD always finds points inside the threat model.

Answer:

i) False. In FGSM, the parameters of the trained model are not directly modified. However, the input data is perturbed to generate adversarial examples that can fool the model's predictions. This means that the model's output can be different for these adversarial inputs, even though its parameters remain the same.

ii) True. The Fast Gradient Sign Method (FGSM) involves perturbing the pixel values of an input image by a small amount in the direction of the sign of the gradient of the loss function with respect to the input. This perturbation is designed to create an adversarial example that can mislead the model's predictions.

iii) False. The Fast Gradient Sign Method (FGSM) is a form of untargeted attack in which the perturbation is not specifically designed to target certain features of the input data that are most important for classification.

iv) True. The PGD algorithm ensures that the final point is within the allowable threat model and close to the original input.

3.2 Solve this question assuming that you're the attacker. You're given a linear classifier which classifies each input into a dog or a cat:

$$f(x) = Pr(x = cat) = \text{sigmoid}(w'x)$$

Given $w = [1, 1]$.

i) It is given that $x = [2, 1]$ represents a cat. What is $Pr(x = dog)$ and $Pr(x = cat)$?

Answer:

Q3.2.

$$(i) \quad h(x) = P_1(x = \text{cat}) = \text{sigmoid}(w^T x)$$

Given $w = [1, 1]$.

$x_1 = [2, 1]$ represents a cat.

$$\Rightarrow P_1(x = \text{cat}) = \frac{1}{1 + e^{-w^T x}}$$

$$\Rightarrow P_1(x_1 = \text{cat}) = \frac{1}{1 + e^{-(1 \cdot 2 + 1 \cdot 1)}}$$

$$\Rightarrow P_1(x_1 = \text{cat}) = \frac{1}{1 + e^{-3}} = 95.26\%$$

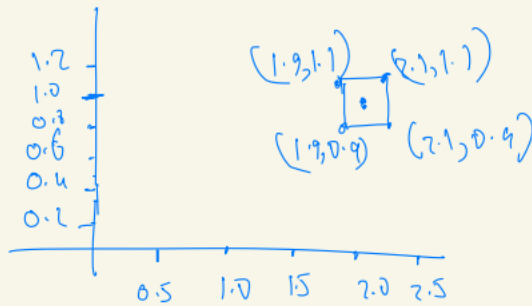
$$\therefore P_1(x_1 = \text{dog}) = 1 - 0.9526 = 4.74\%$$

ii) Now you want to change the components of x (i.e, change from $[2, 1]$ to something else) so that the probability of dog increases. However, you can only change x acc. to the threat model $\|x - x_1\|_1 \leq 0.1$. The threat model specifies the region in which the input x is allowed to vary. 1 Plot the threat model on a 2-d graph and specify the co-ordinates of the corners of the quadrilateral thus formed.

Answer:

(ii) Given, $x_1 = [2, 1]$

$$\|x - x_1\|_{\text{inf}} \leq 0.1$$



\therefore Centre of the square is (2, 1)

\therefore coordinates of the corners are

(1.9, 0.9), (2.1, 0.9), (2.1, 1.1), (1.9, 1.1)

iii) Write down the optimization equation and constraints for finding the optimum x (optimum x is that x which assigns maximum probability to x being a dog). Denote optimum x by x_{optim}

Answer:

(ii) This is the optimization problem -

$$\text{Max } \sigma(w^T x)$$

$$\text{such that } \|x - x_1\|_{\text{inf}} \leq 0.1$$

\therefore we have -

$$\text{Max}_x \sigma(w^T x)$$

$$\text{s.t. } \max(|x - x_1|, |y - y_1|) \leq 0.1$$

As seen in the plot in part (ii),

optimum value of x_1

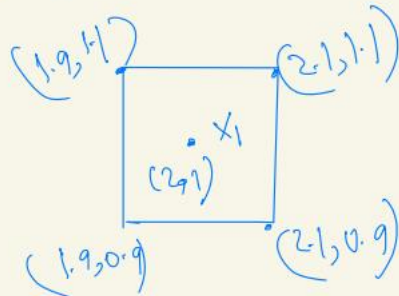
$$x_{\text{optimum}} = (1.9, 0.9)$$

has linear model

iv) Find an FGSM attack using these values of step size, α : 0.01, 0.1, 0.25, 0.5 (finding FGSM attack means finding the updated values of x). 1 Which of these FGSM attack lies in the threat model?

Answer:

(iv) Threat Model:



As given, $\alpha = 0.01$.

$$x' = x_1 + \alpha \cdot \text{signum}(\text{grad}(h)) = \begin{bmatrix} 2 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.01 \\ 0.01 \end{bmatrix} = \begin{bmatrix} 2.01 \\ 1.01 \end{bmatrix}$$

$$\therefore \|x_1 - x'\|_{\infty} = 0.01$$

→ This clearly lies inside
the threat model

Nun, given $\alpha = 0.1$ -

$$x = x_1 + \alpha \cdot \text{Signum}(\text{grad}(h)) = \begin{bmatrix} 2 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.1 \end{bmatrix} = \begin{bmatrix} 2.1 \\ 1.1 \end{bmatrix}$$

$$\|x_1 - x\|_{\text{inf}} = 0.1 \leq 0.1$$

\Rightarrow the FGSM attack clearly lies inside the threat model

Nun, given $\alpha = 0.25$ -

$$x = x_1 + \alpha \cdot \text{Signum}(\text{grad}(h)) = \begin{bmatrix} 2 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.25 \\ 0.25 \end{bmatrix} = \begin{bmatrix} 2.25 \\ 1.25 \end{bmatrix}$$

$$\therefore \|x_1 - x\|_{\text{inf}} = 0.25 > 0.1$$

\therefore the FGSM attack lies outside the threat model

Nun, given, $\alpha = 0.5$ -

$$\therefore x = x_1 + \alpha \cdot \text{Signum}(\text{grad}(h)) = \begin{bmatrix} 2 \\ 1 \end{bmatrix} + \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 2.5 \\ 1.5 \end{bmatrix}$$

$$\therefore \|x_1 - x\|_{\text{inf}} = 0.5 > 0.1$$

\swarrow
This attack lies outside the threat model

v) Find a PGD attack using $\alpha: 0.01$.

Answer:

②) for $\alpha = 0.001$,

$$x_b = x_1 + \alpha (\text{grad}(h)) = \begin{bmatrix} 2 \\ 1 \end{bmatrix} + 0.001 \begin{bmatrix} 2 \\ 1 \end{bmatrix} \\ = \begin{bmatrix} 2.002 \\ 1.001 \end{bmatrix}$$

Clearly, x_b lies within

$$\epsilon = 0.1$$

$$\therefore \text{PGD attack} = x_{\text{attack}} = \begin{bmatrix} 2.002, 1.001 \end{bmatrix}$$

for $\alpha = 0.1$

$$x_b = x_1 + \alpha (\text{grad}(h)) \\ = \begin{bmatrix} 2 \\ 1 \end{bmatrix} + 0.1 \begin{bmatrix} 2 \\ 1 \end{bmatrix} \\ = \begin{bmatrix} 2.2 \\ 1.1 \end{bmatrix}$$

$$x_{\text{-attack}} = \arg \min \|x - x_b\|$$

\swarrow
 lies outside threat model
 \longleftarrow

$$\|x\|_2 = 19$$

$$x_b = x_1 + \alpha (\text{grad}(h))$$

$$= \begin{bmatrix} 4 \\ 1 \\ 2 \end{bmatrix}$$

\swarrow
 clearly lies outside
 threat model
 \longleftarrow

vi) What are your observations from the results in (iv) and (v)?

Answer:

When using the actual values of gradient descent in generating adversarial examples, it can be observed that all the points for all given step sizes lie within the allowed threat model. This indicates that the resulting images are not significantly different from the original, non-adversarial images. In contrast, when we use only the signs of the gradient descent, it can lead to a stronger and more effective adversarial attack. This can result in generating new images that are much closer to the targeted image and can significantly mislead the model's predictions. However, when the model is only dealing with two classes, the difference between an untargeted and a targeted attack might not have a significant impact on the overall results. Nonetheless, it is important to consider and evaluate the potential impact of such attacks in the context of the specific task or application.