CS 188 Fall 2021

Introduction to Artificial Intelligence

Challenge Q7 HW10

- **Due:** Tuesday 11/30 at 10:59pm.
- Policy: Can be solved in groups (acknowledge collaborators) but must be submitted individually.
- Make sure to show all your work and justify your answers.
- Note: This is a typical exam-level question. On the exam, you would be under time pressure, and have to complete this question on your own. We strongly encourage you to first try this on your own to help you understand where you currently stand. Then feel free to have some discussion about the question with other students and/or staff, before independently writing up your solution.
- Your submission on Gradescope should be a PDF that matches this template. Each page of the PDF should align with the corresponding page of the template (page 1 has name/collaborators, question begins on page 2.). **Do not reorder, split, combine, or add extra pages**. The intention is that you print out the template, write on the page in pen/pencil, and then scan or take pictures of the pages to make your submission. You may also fill out this template digitally (e.g. using a tablet.)

First name	Nameera	
Last name	FAISAL AKHTAR	
SID	3684244256	
Collaborators	SARA IMAM	

For staff use only:

O7.	Neural Networks and Optimization	/17

Q7. [17 pts] Neural Networks and Optimization

7.1) (2pts) Which of the following models can represent exactly the XOR function for some choice of parameters? The truth table of the XOR gate $X \bigoplus Y$ is given below:

\boldsymbol{x}	y	$x \oplus y$	ax+by+c	ax2+by2+c
0	0	0	С	С
0	1	1	b+c	b+c
1	0	1	a + c	a + c
1	1	0	a+6+C	a+6+C

When no specific details are mentioned (for instance number of layers, size of layers etc.), that option is correct if there exists **some** architecture that can represent the function correctly.

- A. A logistic regression model that uses only x and y as features.
- B. A logistic regression model that uses only x^2 and y^2 as features.
- C. A neural network with no hidden layers.

D. A neural network with one hidden layer.

A) Logistic regression has a threshold value T_i and all $\sigma(x) \angle T \Rightarrow x = 0$ $\sigma(x) \ge T \Rightarrow x = 1$.

Also, ϵ is monotonically increasing. This means that in order of according x, we cannot go from 0 to 1 back to 0. So this doesn't work.

- C) If correct, $\kappa(c) = \kappa(a+b+c) = 0$. We know $c \le a+b+c$ and β χ^2, y^2 have some effect as χ, y for 0/1 variable. κ monotonic \Rightarrow all κ in [c,a+b+c] must have $\kappa(x) = 0$. This means $\kappa(a+c) = 0$, which is false, blc $\kappa(a+c) = 0$ has $\kappa(a+c) = 0$ and $\kappa(a+c) = 0$.
- Always works.

 7.2) (2pts) Which of the following models can represent the function $\neg(x \lor y)$? The truth table for that function is given below:

\boldsymbol{x}	y	$\neg(x\vee y)$	ax+by+c	ax2+by2+c
0	0	1	С	С
0	1	0	b+c	6+c
1	0	0	a + c	a+c
1	1	0	a+6+C	a+6+C

So this works.

Again, if no specific details are mentioned (for instance number of layers, size of layers etc.), that option is correct if there exists **some** architecture that can represent the function correctly.

- A. A logistic regression model that uses only x and y as features.
- B. A logistic regression model that uses only x^2 and y^2 as features.
- C. A neural network with no hidden layers.
- D. A neural network with one hidden layer.

C) (ASE 1: $a \ge b$ If correct, $\kappa(b+c) = \kappa(a+b+c) = 0$. We know $b+c \le a+b+c$ and $\kappa(a+c) = 0$. We know $\kappa(a+c) = 0$. This means $\kappa(a+c) = 0$, which is true. Works.

B) $\kappa(a+c) = 0$, which is true. Works by some logic as above.

D) Always works.

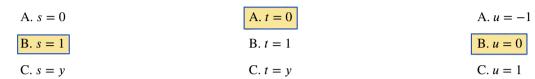
Let's consider a neural network with no hidden layers (one-layer network) that predicts the value $f(x) \in \mathbb{R}$ from an input $x \in \mathbb{R}$, i.e. f(x) = g(Wx + b) with $g(y) = \max(y, 0)$ denoting the ReLU activation function. The initial value for the weight and the bias is W = -1 and b = 0.

7.3) (1pts) Let the derivative of the ReLU function be given as:

C. $g(y) = \sigma(y)$ (sigmoid)



What are the correct values for *s*, *t* and *u*? Choose one answer from each column.



7.4) (3pts) Compute the following partial derivatives using W = -1 and b = 0 as the parameter values.

$$\frac{\partial f}{\partial W}\Big|_{x=-1} = \frac{\delta}{\delta W} \left. g(W_X + b) \right|_{X=-1} = \left. g'(W_X + b) \cdot x \right|_{X=-1} = \left. g'(-x) \cdot x \right|_{X=-1} = \left. g$$

7.5) (3pts) For positive inputs x > 0, what are the upper and lower bounds of the values that the gradient of f with respect to W can take? Assume again that W = -1 and b = 0. $f_z \circ Q(W \times + b)$

$$\frac{\partial f}{\partial W} \leq 0$$

$$= g(-x+0)$$

$$= max(-x,0)$$

$$for x > 0,$$

$$max(-x,0) = 0$$

$$\therefore \frac{\partial f}{\partial W} \leq 0 \text{ and } \frac{\partial f}{\partial W} \geq 0$$

7.6) (2pts) Assume that the objective function is squared loss and that we are using gradient descent to obtain the optimal values for the parameters. The initial values are W=-1 and b=0. If your dataset contains only samples with x>0, then which of the following activation functions will result in the loss function decreasing with the number of gradient descent iterations? More than one option could be correct.

As shown above, for W=-1 and b=0, the derivative of the ReW is always 0.

B.
$$g(y) = ReLU(y)$$
 (rectifier)

B. $g(y) = tanh(y)$ (hyperbolic tangent)

Assume we are training a deep neural network and that we achieve a low training loss. However, the loss on the held-out dataset is rather large.

7.7) (2pts) Which of the following solutions could decrease the loss on the held-out dataset?

Low training loss, high test loss -> overfitting on training!

A. Increasing the number of layers.

Decrease overfitting by removing some layers, so we don't capture

B. Decreasing the number of layers.

"noise" in the data.

7.8) (2pts) Which of the following solutions could decrease the loss on the held-out dataset?

A. Increasing the size of each hidden layer.

Low training loss, high test loss -> overfitting on braining!

B. Decreasing the size of each hidden layer.

Decrease overfitting by removing parameters from each layer, which decreases the size of each hidden layer. This ensures we don't capture

"noise" in the data.