

Midterm Exam, Fall 2024
ESE 577
SOLUTIONS

Instructor: Jorge Mendez-Mendez

Date: Thursday October 10th, 2024

Do not tear exam booklet apart!

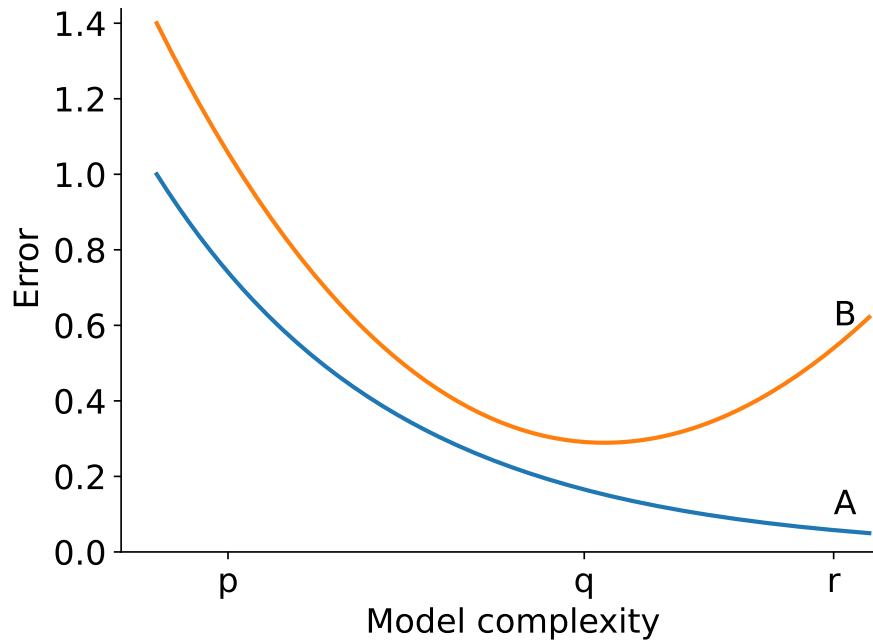
- This is a closed book exam. One sheet (8 1/2 in. by 11 in.) of notes, front and back, are permitted. Calculators are not permitted.
- The total exam time is 2.5 hours.
- The problems are not necessarily in any order of difficulty.
- Record all your answers in the places provided. If you run out of room for an answer, continue on a blank page and mark it clearly.
- If a question seems vague or under-specified to you, make an assumption, write it down, and solve the problem given your assumption.
- If you absolutely *have* to ask a question, come to the front.
- **Write your name on every piece of paper.**

Name: _____ SBU email: _____

Question	Points
1	12
2	19
3	22
4	20
5	9
6	18
Total:	100

Model Evaluation

1. (12 points) The following plot shows the training and validation errors as a function of the model complexity. For each of the following questions about the plots, provide a one-sentence justification for your answer.



- (a) Which of the curves likely corresponds to the training and validation error?

- ☐ Curve A is the training error and curve B is the validation error
- ☐ Curve B is the validation error and curve A is the training error
- ☐ Both curves are the training error
- ☐ Both curves are the validation error

Justification:

ANSWER: Curve B is the validation error and curve A is the training error. The training error goes down as a function of model complexity, while the validation error decreases up to a point, and then the model overfits and the validation error increases.

- (b) What type of fit occurs at the point labeled “p” on the horizontal axis?

- ☐ Underfitting.
- ☐ Proper fit.
- ☐ Overfitting

Justification:

ANSWER: Underfitting. When the model complexity is too low, it is not possible to fit the training data sufficiently well and the model underfits.

- (c) What type of fit occurs at the point labeled “r”?

- ☐ Underfitting.
- ☐ Proper fit.
- ☐ Overfitting

Justification:

ANSWER: Overfitting. When the model complexity is too high, we can fit the training data perfectly, but generalization on the validation data decreases.

- (d) What type of fit occurs at the point labeled “q”?

- ☐ Underfitting.
- ☐ Proper fit.
- ☐ Overfitting

Justification:

ANSWER: Proper fit. When the validation loss reaches its lowest point, we have an adequate fit to the training data that enables the maximum possible generalization.

- (e) Suppose that the curves above correspond to a linear hypothesis trained to minimize the ridge regression objective, given by:

$$J(X, y; \theta) = \frac{1}{n} \left(\sum_{i=1}^n \left(\theta^\top x^{(i)} - y^{(i)} \right)^2 \right) + \lambda \|\theta\|^2 .$$

Which of the following could correspond to the value being plotted in the horizontal axis?

Name: _____

- ☐ θ
- ☐ λ
- ☐ $\frac{1}{\lambda}$
- ☐ $y^{(i)}$

Justification:

ANSWER: $\frac{1}{\lambda}$. As the value of lambda increases, the complexity of the model decreases; so as the value of $\frac{1}{\lambda}$ increases, the complexity of the model increases.

Regression or Repetition

2. (19 points) Suppose that we are given a small dataset and we would like to learn the parameters of a linear regressor hypothesis taking the form $h(x) = \theta^\top x + \theta_0$ for fitting the data.

(a) Consider the following dataset \mathcal{D}_1 containing three data points (in feature-label pairs):

x	y
-4	-3
2	3
-1	6

Suppose that we would like to minimize:

$$J_1(\theta, \theta_0; \mathcal{D}_1) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 .$$

- i. For J_1 , we know that there exist θ^* and θ_0^* that minimize it. Can we find θ^* via the analytical solution formula? (No need for justification)

- ☐ Yes.
☐ No.

ANSWER: Yes. With a single feature, there is no possibility of colinear features or $n < d$ (since $d = 1$), so the matrix inverse will be well-defined.

- ii. Suppose we know that for J_1 , one set of minimizing parameters has $\theta_0^* = 3$. What is the corresponding unknown θ^* ?

ANSWER:

$$\begin{aligned} 3J_1 &= (-4\theta + 3 + 3)^2 + (2\theta + 3 - 3)^2 + (-\theta + 3 - 6)^2 \\ &= (4\theta - 6)^2 + (2\theta)^2 + (\theta + 3)^2 \\ &= (16\theta^2 - 48\theta + 36) + (4\theta^2) + (\theta^2 + 6\theta + 9) \\ &= 21\theta^2 - 42\theta + \text{some constant} . \end{aligned}$$

Taking the gradient and equating to 0 we get:

$$\theta^* = \frac{42}{21 \times 2} = 1 .$$

- iii. What is J_1^* , the minimum value achievable of J_1 ?

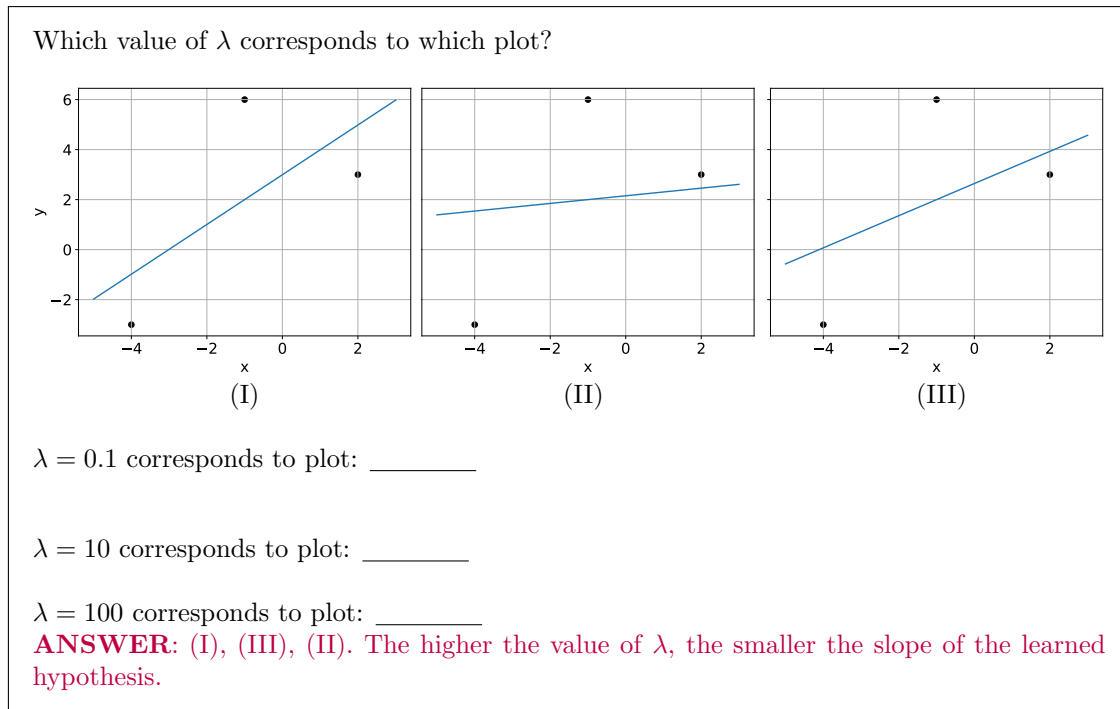
ANSWER: Using our obtained θ^* , we have:

$$\begin{aligned} J_1^* &= \frac{(4 - 6)^2 + (2)^2 + (1 + 3)^2}{3} \\ &= \frac{4 + 4 + 16}{3} = 8 . \end{aligned}$$

(b) Suppose instead of J_1 , we try to minimize:

$$J_2(\theta, \theta_0; \mathcal{D}_1, \lambda) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 ,$$

with $\lambda = 0.1, 10$, and 100 . Identify the λ used to generate plot (III).



(c) Suppose we add a second feature for each of the three datapoints in \mathcal{D}_1 , to obtain the new dataset \mathcal{D}_2 :

x_1	x_2	y
-4	8	-3
2	-4	3
-1	2	6

Suppose that we would like to minimize:

$$J_3(\theta, \theta_0; \mathcal{D}_2) = \frac{1}{n} \sum_{i=1}^n (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 .$$

- i. For J_3 , we also know that there exist θ^* and θ_0^* that minimize it. Can we find θ^* via the analytical solution formula? If yes, provide such θ^* ; if not, briefly justify why not.

☐ Yes.

☐ No.

$\theta^* =$ **ANSWER:** No. The two features are linearly dependent. Hence $X^\top X$ will not be invertible.

- ii. Compare J_3^* , the minimum value achievable of J_3 , with J_1^* . Which option below is true? Briefly justify your choice.

- ☐ $J_3^* > J_1^*$
- ☐ $J_3^* = J_1^*$
- ☐ $J_3^* < J_1^*$
- ☐ It depends

$\theta^* =$

ANSWER: $J_3^* = J_1^*$. Note that in \mathcal{D}_2 , the 2nd feature is linearly dependent on the 1st feature. Intuitively, this means that the 2nd feature gives us nothing additional to learn from towards a linear hypothesis, i.e. the 2nd feature is redundant. Algebraically, we can write out explicitly J_3 and J_1 to see their connection. Suppose, for the sake of notational cleanliness, that we let the set of parameters for J_1 be α and β :

$$J_1 = [(-4\alpha + \beta - 15)^2 + (2\alpha + \beta + 3)^2 + (-\alpha + \beta)^2] / 3 ,$$

and let the parameters for J_3 be θ and θ_0 :

$$J_3 = [(-4\theta_1 - 8\theta_2 + \theta_0 - 15)^2 + (2\theta_1 + 4\theta_2 + \theta_0 + 3)^2 + (-\theta_1 - 2\theta_2 + \theta_0)^2] / 3 .$$

Then, we realize that by letting $\theta_1 + 2\theta_2 = \alpha$ and $\theta_0 = \beta$, any value achievable by J_1 is achievable by J_3 and vice-versa.

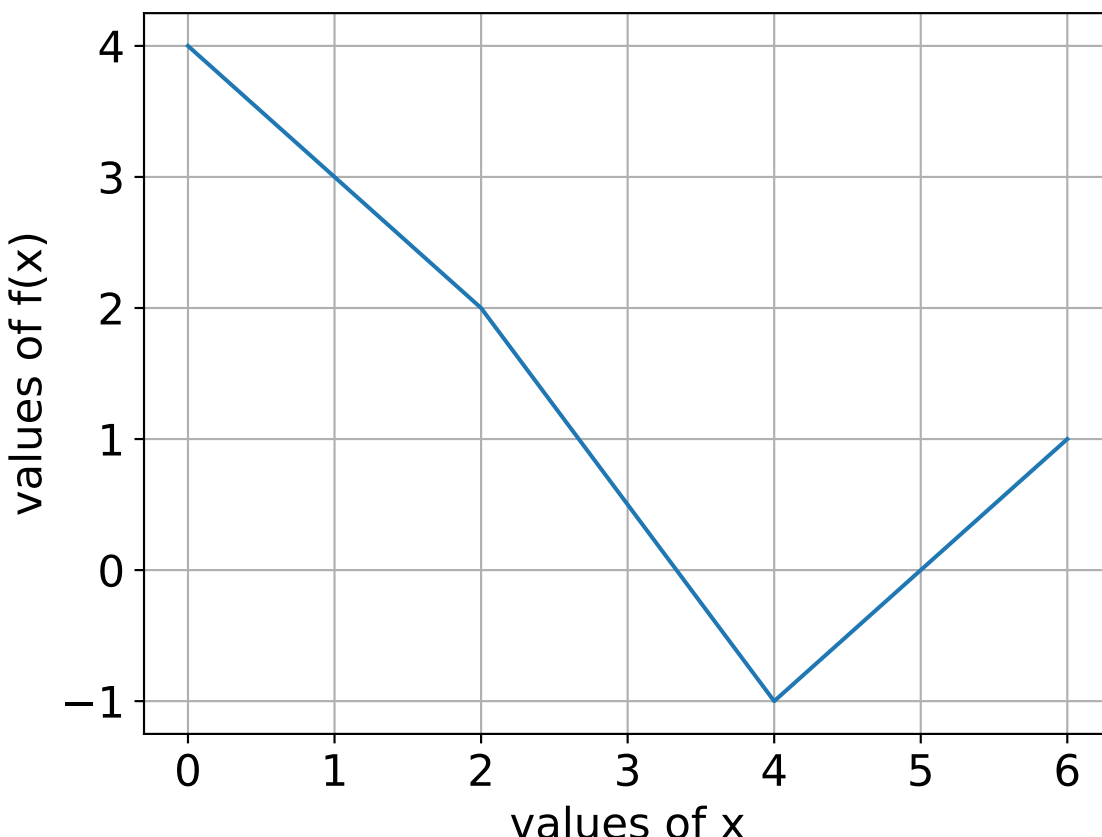
Gradient Descent in Pictures

3. (22 points) John is using standard gradient descent iterations:

$$x^{(k+1)} = x^{(k)} - \eta \nabla f(x^{(k)}) \quad (k = 0, 1, 2, \dots)$$

on a variety of functions f .

- (a) First, John applies gradient descent to a piecewise-linear function $f : \mathbb{R} \mapsto \mathbb{R}$ with the (partial) graph shown in the figure below:



At points $x = 2$ and $x = 4$, John uses $\nabla f(2) = -1.5$ and $\nabla f(4) = 0$.

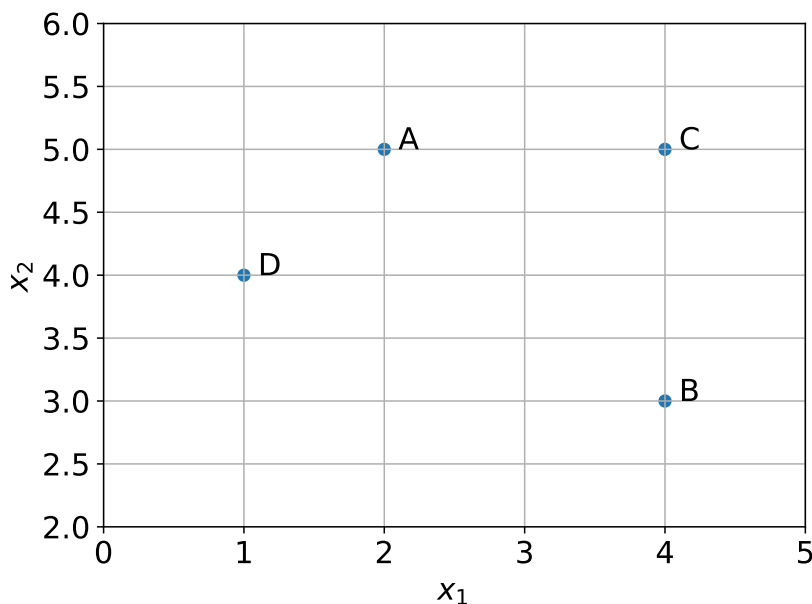
- i. Starting from the initial guess $x^{(0)} = 0.5$, and using a step size $\eta = 2$, what will be the values of $x^{(1)}$, $x^{(2)}$, and $x^{(3)}$?

ANSWER: $x^{(1)} = 2.5$, $x^{(2)} = 5.5$. $x^{(3)} = 3.5$. By inspection of the graph $\nabla f(x^{(0)}) = -1$, hence $x^{(1)} = 0.5 - 2 \cdot (-1) = 2.5$. Next, $\nabla f(x^{(1)}) = -1.5$, hence $x^{(2)} = 2.5 - 2 \cdot (-1.5) = 5.5$. Finally, $\nabla f(x^{(2)}) = 1$, hence $x^{(3)} = 5.5 - 2 \cdot 1 = 3.5$.

- ii. John discovers that, starting with $x^{(0)} = 5$, there are many values of $\eta > 0$ for which the gradient descent iterations produce oscillations of period 2 within the range $(0, 6)$ (i.e., $x^{(k+2)} = x^{(k)} \in (0, 6)$ for all $k = 0, 1, 2, \dots$). Find all such values of η .

ANSWER: $\eta \in (3, 5)$. First notice that the gradient at $x^{(0)} = 5$ is $\nabla f(1) = 1$. In order to oscillate with period 2, we need to immediately land at a point with gradient $\nabla f(x^{(1)}) = -1$. We observe that $\nabla f(x) = -1$ for all $x \in (0, 2)$. So we need: $x^{(2)} = 5 - \eta \in (0, 2) \Rightarrow \eta \in (3, 5)$.

- (b) After mastering one-dimensional optimization, John applies gradient descent to a smooth convex function $f : \mathbb{R}^2 \mapsto \mathbb{R}$, with $\eta = 0.1$, resulting in the sequence of points $x^{(0)} = A$, $x^{(1)} = B$, $x^{(2)} = C$, $x^{(3)} = x^{(4)} = D$ shown on the plot below:



- i. Find $\nabla f(x^{(0)})$, $\nabla f(x^{(1)})$, $\nabla f(x^{(2)})$, and $\nabla f(x^{(3)})$.

ANSWER: $x^{(k+1)} = x^{(k)} - \eta \nabla f(x^{(k)}) \Rightarrow \nabla f(x^{(k)}) = \frac{1}{\eta}(x^{(k)} - x^{(k+1)})$, hence:

$$\nabla f(x^{(0)}) = \begin{bmatrix} -20 \\ 20 \end{bmatrix}, \nabla f(x^{(1)}) = \begin{bmatrix} 0 \\ -20 \end{bmatrix}, \nabla f(x^{(2)}) = \begin{bmatrix} 30 \\ 10 \end{bmatrix}, \nabla f(x^{(3)}) = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

- ii. Is this statement true or false: “given the information provided, the point D **must** be a global minimum of function f ”? Briefly justify your choice.

☐ True

☐ False

ANSWER: True. We know that $\nabla f(D) = 0$. Since the problem states that the function is convex, we are guaranteed that any point with $\nabla = 0$ is a global minimum.

CTO Logistic Regression

4. (20 points) You have a training dataset \mathcal{D} with each input $x^{(i)}$ consisting of d binary features $x_1^{(i)}, \dots, x_d^{(i)}$, where all $x_j^{(i)} \in \{0, 1\}$ and a binary label $y^{(i)} \in \{0, 1\}$. After hours struggling with *underfitting* (high training-set error) when using logistic regression to make predictions, you decide to call your friend, the CTO of the famous start-up *ClosedAI*, who gives you a few ideas to get a lower training error. For each of them, specify whether 1) it generally reduces underfitting; 2) it will not make a difference; or 3) it generally worsens underfitting. **Note that we are only talking about training loss, not test loss. Justify each answer.**

- (a) “Combine and conquer”: For each data point, augment the dimensions by adding the *and* of each feature with its neighbor; that is $[x_1^{(i)}, \dots, x_d^{(i)}] \mapsto [x_1^{(i)}, \dots, x_d^{(i)}, x_1^{(i)} \cdot x_2^{(i)}, \dots, x_{d-1}^{(i)} \cdot x_d^{(i)}]$.

- ☐ Reduces underfitting
☐ No change
☐ Worsens underfitting

ANSWER: Reduces underfitting. Because the new features are a non-linear function of the original features, they add expressive power to our hypothesis class.

- (b) “Two is bigger than one”: Use two sigmoids instead of one; that is, instead of parameterizing the solution as $g^{(i)} = \sigma(\theta^\top x^{(i)} + \theta_0)$, parameterize it as $g^{(i)} = \sigma(\sigma(\theta^\top x^{(i)} + \theta_0))$.

- ☐ Reduces underfitting
☐ No change
☐ Worsens underfitting

ANSWER: Worsens underfitting. The output of $\sigma(z)$ is between 0 and 1 for all z . Therefore the output of $\sigma(\sigma(z))$ is between $\sigma(0) = 0.5$ and $\sigma(1) \approx 0.73$, which severely limits its range. In particular, this hypothesis class cannot predict $y^{(i)} = \text{False}$.

- (c) “Ensembling is always good”: Fit a regular logistic regression using the standard sigmoid function with base e and a second logistic regression using a base-2 σ_2 sigmoid function, and return the most confident result (farthest away from 0.5).

- ☐ Reduces underfitting
☐ No change
☐ Worsens underfitting

ANSWER: No change.

$$\sigma_2(z) = \frac{1}{1 + 2^{-z}} = \frac{1}{1 + e^{-\ln 2 z}} = \sigma(\ln 2 z) .$$

Therefore, σ_2 logistic regression and σ logistic regression are equally expressive and training them will yield equal predictions for all datapoints. In consequence, ensembling them as proposed will not change results.

- (d) “Label distillation”: Before fitting the dataset, for every data-point i , add its label as a feature; that is, $[x_1^{(i)}, \dots, x_d^{(i)}] \mapsto [x_1^{(i)}, \dots, x_d^{(i)}, y^{(i)}]$.

- ☐ Reduces underfitting
- ☐ No change
- ☐ Worsens underfitting

ANSWER: Reduces underfitting. Regardless of the original features, we can always correctly predict the label by putting all weight on that new feature, which will directly predict the correct label.

- (e) “Not quite a 2-layer neural network”: Fit a regular logistic regression and obtain a prediction $g^{(i)}$ for each element $x^{(i)}$. Add $g^{(i)}$ as a feature: $[x_1^{(i)}, \dots, x_d^{(i)}, g^{(i)}]$ and fit a new logistic regression to the new dataset.

- ☐ Reduces underfitting
- ☐ No change
- ☐ Worsens underfitting

ANSWER: Reduces underfitting. $g^{(i)}$ is a non-linear function combination of the features, and therefore increases the capacity of the logistic regression. Moreover, $g^{(i)}$ contains information of the labels through the parameters trained during the first logistic regression, which makes it a very informative feature.

Wolfie Madness

5. (9 points) Over the last few years, Wolfie has been seen running through some of the popular lectures on campus. There is no warning as to when or where Wolfie would show up (for example, you might be mid-quiz in 577 and then BAM!, a huge seawolf is sprinting behind your professor in his traditional red uniform).

New student Tyler is a big fan, and is interested in seeing if he can figure out how Wolfie chooses when and where to appear in class.

- (a) First, Tyler interviews former students about what they can remember about the Wolfie encounters, in order to collect some data.

For each of the following, suggest an appropriate feature encoding, and provide the dimension of the feature encoding.

- i. Each interviewee remembers the class they were in when they had an encounter. It was always one of {ESE 557, ESE 534, ESE 506, ESE 358, or ESE 301}.

Encoding:

Encoding dimension:

ANSWER: This is a categorical variable, and there is no relationship between each class. We can use a one-hot vector of size 5, in which only a single class is marked at a time.

- ii. There has long been a rumor on campus that the professors are in on Wolfie's appearances, and they wear distinctive items on the days when Wolfie would show up. By looking at the lecture videos, Tyler can find out if the professor was wearing a bowtie, a hat, or a plaid coat. A professor could wear zero, one, or more of these items.

Encoding:

Encoding dimension:

ANSWER: Since the items can be worn in various combinations, it is a good idea to use a factored encoding: essentially one binary feature for each of bowtie, hat, and plaid coat. A one-hot vector would be less useful, since we would need to observe all possible combinations in the data in order to be able to make predictions on new data points that contain those combinations.

- iii. Tyler can also find out the number of students who were present in lecture for each encounter. This might range from zero students showing up (during busy weeks), to a maximum of 522 students attending (full capacity in the largest lecture hall, Javits-100).

Encoding:

Encoding dimension:

ANSWER: This is a continuous numerical feature with a known range. We ideally want our values to be in an appropriate range (e.g., $[0, 1]$). Both standardization and normalization could be acceptable answers here.

- (b) Using all of this data, Tyler now wants to make predictions about how likely it might've been to see Wolfie.

He wants to design a simple neural network for the task. What should the dimension of Tyler's *output* prediction be (that is, the shape of his output)? what would be a good choice of activation function on the output layer? Briefly justify your answer.

Dimension:

Activation function:

Justification:

ANSWER: Tyler wants a probability, so this should be a scalar number output in the range 0 to 1. "Scalar" or size 1 are acceptable. Tyler wants a single probability, so he should use a sigmoid activation. Note that while there are many cases where softmax is a helpful function for probabilities, and is equivalent to sigmoid when $n = 2$, softmax cannot be used here as the probability on a single output would always be 1.

Deep neural networks

6. (18 points) Nori thinks about ReLU units and wonders whether there's a better alternative for an activation function, and decides to explore the LURe function, defined as:

$$f_{\text{LURe}} = \min(z, 0) \text{ .}$$

- (a) What is the derivative of this function $\frac{df_{\text{LURe}}(z)}{dz}$?

ANSWER:

$$\frac{df_{\text{LURe}}(z)}{dz} = \begin{cases} 1 & \text{if } z < 0 \\ 0 & \text{otherwise} \end{cases}$$

Or any variants of this. Note that the derivative is **never** -1 !

- (b) Nori's friend Ori thinks this is cool and suggests making a neural network with two activation functions per layer, so that:

$$a^l = f_{\text{LURe}}(f_{\text{ReLU}}(z^l)) \text{ .}$$

Explain what effect this will have on the network.

ANSWER: This will effectively make the activation function $f(z) = 0$ for all z , destroying all expressivity of the NN.

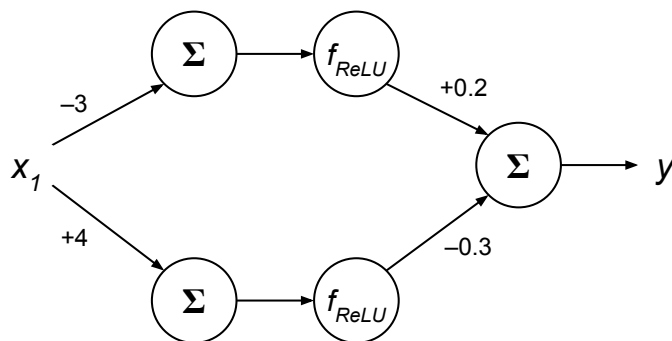
- (c) Nori's other friend Dori thinks we should try this trick with two ReLUs, so that:

$$a^l = f_{\text{ReLU}}(f_{\text{ReLU}}(z^l)) \text{ .}$$

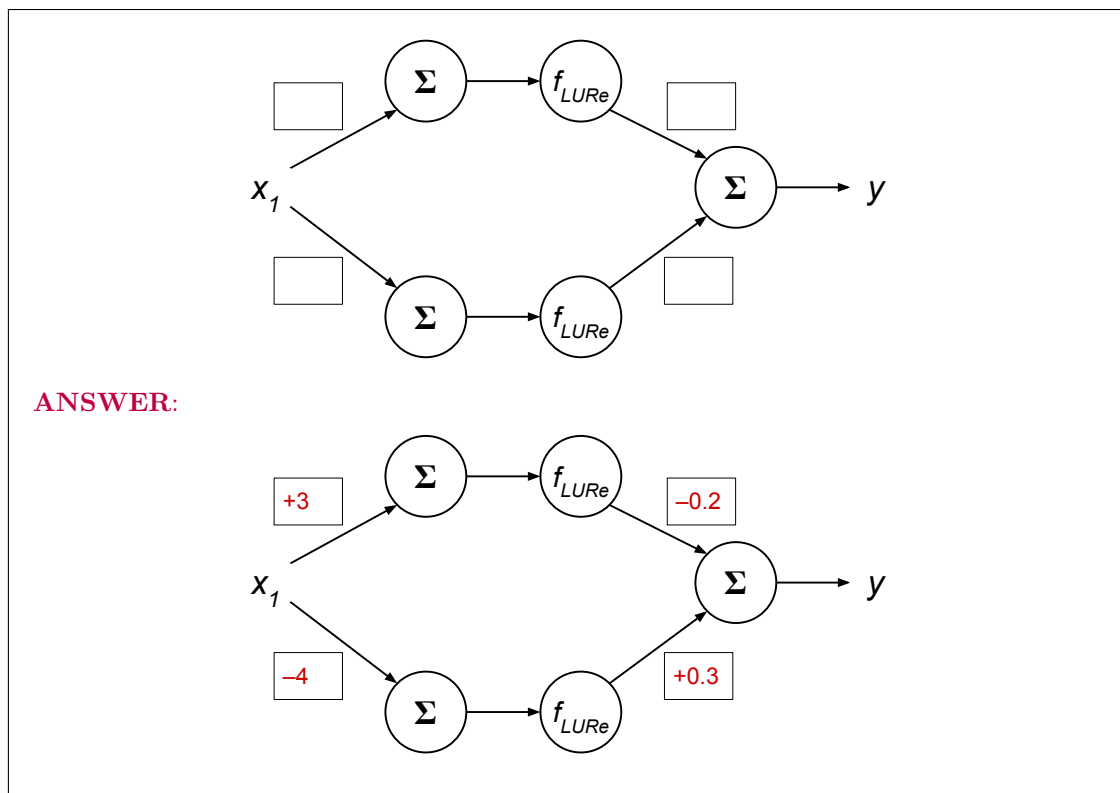
Explain what effect this will have on the network.

ANSWER: It will act exactly the same as having one ReLU.

- (d) Nori finds a neural network trained by Smaug in his treasure pile that takes a single-dimensional input (so $d = 1$) and looks like this:



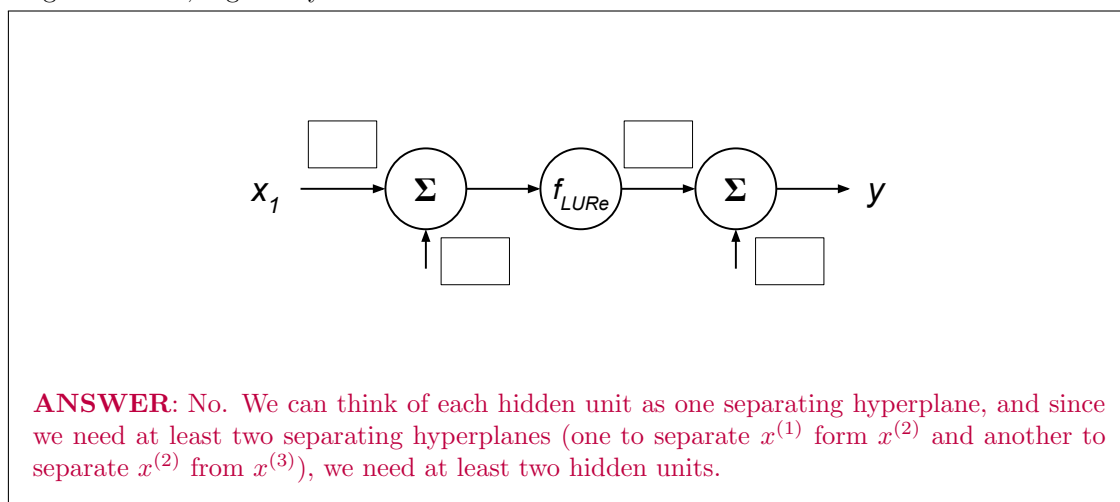
He sees that it computes $\hat{y} = -0.1 \cdot f_{\text{ReLU}}(-5x) + 0.4 \cdot f_{\text{ReLU}}(5x)$ and is very curious to see if he can replace those ReLU activation units with his own LURes. Please help him find another neural network that computes exactly the same function as the one above (that is, maps any input x to the same output as the original one). Provide a set of weights that achieves this in the boxes on the diagram below.



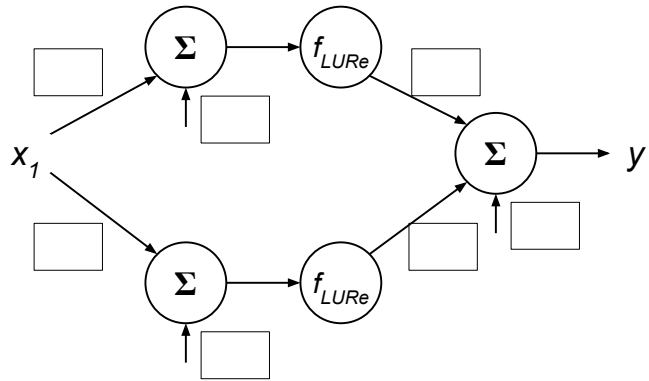
(e) Ori, Dori, and Nori have this dataset with their heights and whether they like beer:

$$\begin{aligned} x^{(1)} &= [1] & y^{(1)} &= 1 \\ x^{(2)} &= [2] & y^{(2)} &= 0 \\ x^{(3)} &= [3] & y^{(3)} &= 1 \end{aligned}$$

They make a two-layer neural network, as shown below. Are there weights and biases for this network that will predict their data correctly? If so, specify them in the boxes on the network diagram. If not, argue why not.



(f) Thorin suggests they add one more unit, so that they have the architecture below. Are there weights and biases for this network that will predict their data correctly? If so, specify them in the boxes on the network diagram. If not, argue why not.



ANSWER:

