

## 6.036: Midterm, Spring 2022

### Solutions

- This is a closed book exam. One page (8 1/2 in. by 11 in. or A4) of notes, front and back, is permitted. Calculators are not permitted.
- The total exam time is 2 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: \_\_\_\_\_ MIT Email: \_\_\_\_\_

Question	Points	Score
1	12	
2	25	
3	22	
4	15	
5	26	
Total:	100	

## 1 How to do ML?

1. (12 points) Please answer each question with a phrase or short sentence.

- (a) If you are training a linear regression model and find that you are getting very low training error but much higher prediction error on a held-out data set from the same source, you should

\_\_\_\_\_ **add or increase the regularization** \_\_\_\_\_ .  
**or change to a hypothesis class with less expressive power**

- (b) If you have trained a classifier on half of your available data and want to estimate how well it will work in actual practice, you should

\_\_\_\_\_ **test it on the other half of the data** \_\_\_\_\_ .

- (c) If you have two different neural-network architectures and you are trying to decide which one to use to perform well given your current data set, you should

\_\_\_\_\_ **use cross-validation to estimate the quality of hypotheses each produces** \_\_\_\_\_ .

- (d) If you are training a logistic regression classifier using gradient descent and find that the prediction accuracy stays around 50% but you look at the weights and they're really big, you should

\_\_\_\_\_ **decrease the step size** \_\_\_\_\_ .

- (e) If you are training a logistic regression classifier using gradient descent and find that the prediction accuracy goes up to 100% but you look at the weights and they're really big, you should

\_\_\_\_\_ **add or increase the regularization** \_\_\_\_\_ .

- (f) If you are training a neural network classifier and find that the loss decreases for a while, but then stops decreasing before it gets to an acceptable accuracy, you should

\_\_\_\_\_ **increase the number of hidden units** \_\_\_\_\_ .

## 2 Media consumption

2. We are trying to help a new public relations company, NetFlacks, try to make predictions about the popularity of various TV shows.

(a) (6 points) Assume that we are trying to predict the number of views of a particular TV show in the next month. What is the best way to encode the following inputs? Specify the number of features that would be used to represent each of these inputs, and how the given input value would be represented. It is fine to leave numerical expressions, if you need them, unevaluated.

i. Genre: non-fiction, comedy, drama, science fiction (each show has a single genre)

Number of features: 4

How would we represent a comedy?

**Solution:** One-hot encoding vector  $[0, 1, 0, 0]$

ii. Famous actors: which members of a list of the 200 most famous actors (*[Arnold Aardvark, ..., Ziggy Ztarduzt]*) appear in it

Number of features: 200

How would we represent a movie starring only Arnold Aardvark and Ziggy Ztarduzt?

**Solution:** Vector  $[1, 0, \dots, 0, 1]$

iii. Production cost:

Number of features: 1

If our dataset had 5 shows with costs \$1M, \$2M, \$3M, \$4M, \$5M, how would we represent the cost of the first of these shows?

**Solution:** Normalizing the production cost, computed by subtracting off the average production costs of all the shows and dividing by the standard deviation. So  $-2/\sqrt{5/2}$ .

Name: \_\_\_\_\_

- (b) (5 points) After the show has been released there are a lot of ratings available. Assume that, for each show, you collect all the published reviews, and compute how many of those reviews gave the show 1 star, how many gave it 2 stars, etc., up to 5 stars. You end up with 5 features for each show:  $(c_1, c_2, c_3, c_4, c_5)$ , where  $c_i$  is an integer number of reviews with  $i$  stars. So, for example, the following shows might have the following features:

- *Big Blockbuster Bonanza* (100, 1000, 1M, 1000, 10)
- *Obscure Omniglot Omnibus* (0, 0, 1, 1, 10)

You want to do linear regression to predict the number of views of the show.

You consider using these features just as they are, but also consider some different encodings.

1. raw counts,  $c_1, \dots, c_5$  (5 features)
2. average  $a = \frac{\sum_{i=1}^5 i \cdot c_i}{\sum_{i=1}^5 c_i}$  (1 feature)  
(Note that this was mis-defined in the original exam.)
3. average  $a$  and the sum,  $s = \sum_{i=1}^5 c_i$  (2 features)
4. the normalized counts  $c_i/s$  (5 features)
5. the normalized counts  $c_i/s$  and the sum  $s$  (6 features)

In order to pick a good encoding you want to be sure it can encode some plausible models of profitability, via linear regression. For each of the rules below, indicate **all** of the encodings that will allow it to be expressed via *linear regression*, or explain why none of them will do.

- i. The number of reviews predicts the number of viewers in the next month.  
 1.    2.    3.    4.    5.

<b>Solution:</b>
------------------

- ii. The percentage of reviews with more than 3 stars predicts the number of viewers next month.  
 1.    2.    3.    4.    5.

<b>Solution:</b>
------------------

Name: \_\_\_\_\_

(c) (6 points) What is a good choice of *output* encodings if we are trying to make the following predictions with a neural network? Please also specify an appropriate loss function and output activation function (if one is needed).

- i. Number of viewers in the next month  
Output encoding, including number of dimensions

\_\_\_\_\_ **numeric, 1** \_\_\_\_\_

Loss function

\_\_\_\_\_ **squared** \_\_\_\_\_

Output activation

\_\_\_\_\_ **none or linear** \_\_\_\_\_

- ii. Whether each member of the current 6.036 class (with 450 students) will watch it in the next month, predicting all values as outputs of a single neural network.  
Output encoding, including number of dimensions

\_\_\_\_\_ **450 binary values** \_\_\_\_\_

Loss function

\_\_\_\_\_ **average NLL** \_\_\_\_\_

Output activation

\_\_\_\_\_ **450 sigmoids** \_\_\_\_\_

- iii. Whether the exclusive rights to distribute it will be sold to network A, network B, or TwoYoob.  
Output encoding, including number of dimensions

\_\_\_\_\_ **one-hot, 3** \_\_\_\_\_

Loss function

\_\_\_\_\_ **NLLM** \_\_\_\_\_

Output activation

\_\_\_\_\_ **Softmax** \_\_\_\_\_

Name: \_\_\_\_\_

(d) (8 points) You are interested in the effect of the language of a show on its popularity. You consider two encodings of the language:

- A: One-hot encoding of one of 100 modern languages
- B: An integer index into a list of 100 modern languages

You do regression on a data set with two examples in which the language is the only feature. A show in language 2 has had 200 views and one in language 8 has had 800:  $\{(2, 200), (8, 800)\}$

Recall that in linear regression, we do not regularize  $\theta_0$ .

For each of the following learning methods, what prediction would the resulting hypothesis make for language 10 on the list? Provide an approximate numeric output or indicate that it is under-specified.

i. Using strongly regularized linear regression on encoding A?

**Solution:** 500 because  $\theta$  will be driven to 0 by regularization and then  $\theta_0$  will be the average  $y$  value.

ii. Using strongly regularized linear regression on encoding B?

**Solution:** 500 (same reasoning as above)

iii. Using unregularized linear regression on encoding A?

**Solution:** under-specified

iv. Using unregularized linear regression on encoding B?

**Solution:** 1000 because we fit a line to the training data.

### 3 Relugression

3. Let's consider regression in one dimension, so our inputs  $x^{(i)}$  and outputs  $y^{(i)}$  are in  $\mathbb{R}$ .

(a) (4 points) Linny uses regular linear regression. Given the following dataset,

$$\mathcal{D} = \{(1), 1), ((2), 2), ((3), 4), ((3), 2)\}$$

what values of  $\theta$  and  $\theta_0$  optimize the mean squared error of hypotheses of the form  $h(x; \theta, \theta_0) = \theta x + \theta_0$ ?

$$\theta = \underline{\quad \mathbf{1} \quad} \qquad \theta_0 = \underline{\quad \mathbf{0} \quad}$$

(b) (4 points) What property has to be true of your solution above in order for it to be at least a local optimum of the mean squared error objective function

$$J(\theta, \theta_0) = \frac{1}{4} \sum_{(x,y) \in \mathcal{D}} (\theta x + \theta_0 - y)^2 \quad ?$$

Provide conditions on  $\partial J / \partial \theta$  and  $\partial J / \partial \theta_0$

**Solution:**  $\partial J / \partial \theta$  and  $\partial J / \partial \theta_0$  must both be equal to 0. Note that this objective is convex in  $\theta$  and  $\theta_0$  (it's an upward-opening parabola in either variable), and thus a point satisfying these zero-equality conditions is also a global optimum. If our objective were not convex in  $\theta, \theta_0$ ,  $\partial J / \partial \theta = 0$  and  $\partial J / \partial \theta_0 = 0$  could both be true at a *saddle point* rather than a local optimum.

(c) (4 points) Rolly read about neural networks and thinks linear regression would be improved by using a ReLU unit. Recall that ReLU is defined as

$$\text{ReLU}(z) = \begin{cases} 0 & \text{if } z < 0 \\ z & \text{otherwise} \end{cases} = \max(0, z)$$

Now we consider a new dataset:

$$\mathcal{D} = \{((0), 0), ((1), 0), ((2), 2), ((3), 4), ((4), 6)\}$$

What values of  $\theta$  and  $\theta_0$  optimize the mean squared error of hypotheses of the form  $h_R(x; \theta, \theta_0) = \text{ReLU}(\theta x + \theta_0)$ ?

$$\theta = \underline{\quad \mathbf{2} \quad} \qquad \theta_0 = \underline{\quad \mathbf{-2} \quad}$$

Name: \_\_\_\_\_

- (d) (4 points) Rolly insists their hypothesis class is bigger than Linny's, in the sense that any data set that can be fit with 0 MSE using a hypothesis in Linny's class can also be fit with 0 MSE using a hypothesis in Rolly's class.

Is Rolly right?

Yes  **No**

Do one of the following:

- Argue that Rolly is right by providing, for any finite data set  $\mathcal{D}$  such that there is a linear hypothesis  $\theta, \theta_0$  with 0 MSE, the parameters of a ReLU hypothesis that also has 0 MSE.

**Solution:** Nope.

- Argue that Rolly is wrong by providing a small dataset that has a 0 MSE linear hypothesis but for which no 0 MSE ReLU hypothesis exists.

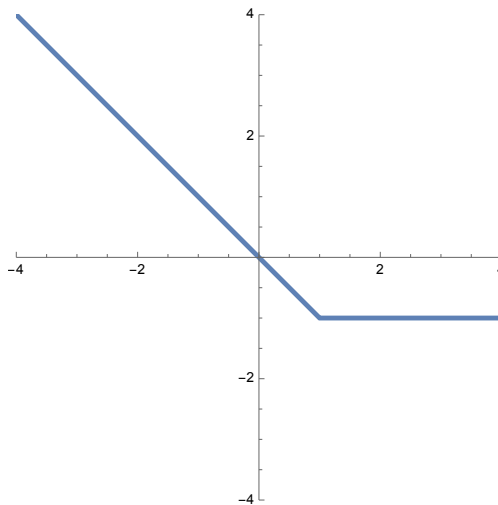
**Solution:**

$$\{(-1, -1), (0, 0), (1, 1)\}$$

- (e) (3 points) Lulu is interested in the following hypothesis class, with parameters  $a, b, c$ , and  $d$ :

$$h(x; a, b, c, d) = a + b \operatorname{ReLU}(cx + d)$$

For the plot below, provide values of  $a, b, c$ , and  $d$  that would generate it. Choose parameter values equal to -1, 0, or 1; or say that it's not possible given this hypothesis class.



$a =$    -1    $b =$    1    $c =$    -1    $d =$    1

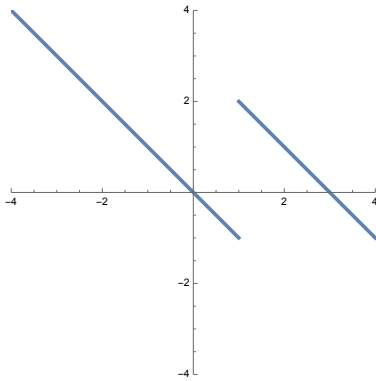


Name: \_\_\_\_\_

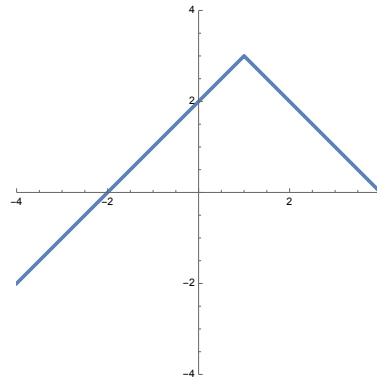
(f) (3 points) Renee wants to be even fancier, and considers hypothesis class with additional parameters  $e$ ,  $f$ , and  $g$ :

$$h(x; a, b, c, d, e, f, g) = a + b \operatorname{ReLU}(cx + d) + e \operatorname{ReLU}(fx + g)$$

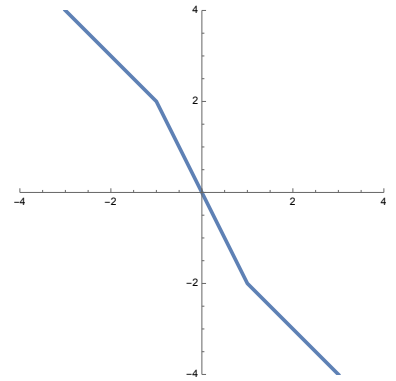
For each plot below, indicate whether it can be expressed using Renee's class.



Yes  No



Yes  No

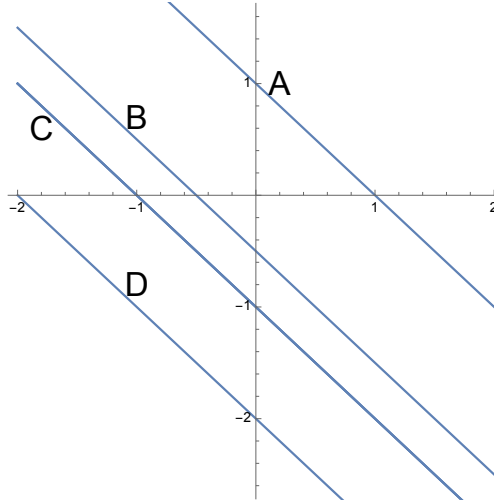


Yes  No

## 4 Class struggle

4. We are interested in doing linear classification of points in a 2-dimensional space.

- (a) (5 points) We started with a classifier, described by  $\theta = [1, 1]$ ,  $\theta_0 = 1$ . All our friends had theories about how to transform it. They drew their candidate separators on a napkin at a restaurant, but forgot to label the or indicate which side was positive vs negative, so we just have drawings of the separators but not the normals.



For each of the transformations below, indicate which separator in the diagram it corresponds to.

- i. Multiply  $\theta$  by 2

A    B    C    D    none

- ii. Multiply  $\theta$  and  $\theta_0$  by 2

A    B    C    D    none

- iii. Multiply  $\theta$  by -1

A    B    C    D    none

- iv. Multiply  $\theta$  and  $\theta_0$  by -1

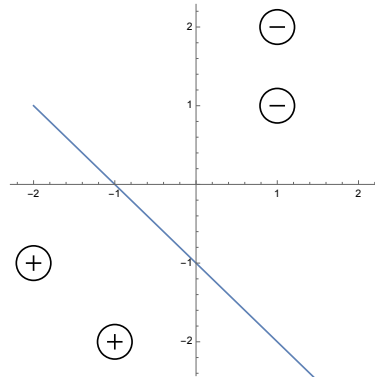
A    B    C    D    none

- v. Add 1 to  $\theta_0$

A    B    C    D    none

Name: \_\_\_\_\_

- (b) (4 points) Consider the following data set, shown with the separator corresponding to the original linear classifier with parameters  $\theta = [1, 1]$ ,  $\theta_0 = 1$ .



Remembering to think about which side of the separator is assigned the positive class, select the transformation below that causes the largest decrease in negative log likelihood (NLL)?

- Multiply  $\theta$  and  $\theta_0$  by 2
  - Multiply  $\theta$  and  $\theta_0$  by  $1/2$
  - Multiply  $\theta$  and  $\theta_0$  by  $-1$
  - Multiply  $\theta$  and  $\theta_0$  by  $-2$**
- (c) (6 points) Now consider a problem with input dimension  $d = 10$ , with separator

$$\theta = [1, 1, 1, 1, 1, 1, 1, 1, 1, 1], \quad \theta_0 = 0$$

We are wondering about the following data points:

- $x^{(1)} = [0, 0, 0, 0, 0, 0, 0, 0, 0, 0]$
- $x^{(2)} = [1, 1, 1, 1, 1, 1, 1, 1, 1, 1]$
- $x^{(3)} = [-1, -1, -1, -1, -1, -1, -1, -1, -1, -1]$
- $x^{(4)} = [1, -1, 1, -1, 1, -1, 1, -1, 1, -1]$
- $x^{(5)} = [1, 0, 0, 0, 0, 0, 0, 0, 0, 0]$

- i. Which point or points are assigned to class 1 (positive) with highest confidence?

- $x^{(1)}$      $x^{(2)}$      $x^{(3)}$      $x^{(4)}$      $x^{(5)}$

**Solution:** The confidence of a class 1 assignment is given by  $\sigma(\theta^T x + \theta_0)$ , which is high when  $\theta^T x + \theta_0 = \theta^T x$  is high ( $\sigma$  is monotonic). Note that since  $\theta$  is a vector of ones,  $\theta^T x = \sum_{k=1}^d x_k$ , the sum of the components of  $x$ . Of the  $x^{(i)}$ ,  $x^{(2)}$  has the highest component sum of 10, so  $x^{(2)}$  is assigned to class 1 with highest confidence.

- ii. Which point or points are assigned to class 0 (negative) with highest confidence?

- $x^{(1)}$      $x^{(2)}$      $x^{(3)}$      $x^{(4)}$      $x^{(5)}$

Name: \_\_\_\_\_

**Solution:** In binary classification, maximizing class 0 confidence is equivalent to minimizing class 1 confidence, so we seek the  $x^{(k)}$  with the lowest component sum. This is achieved by  $x^{(3)}$ , which has a component sum of -10.

iii. Which point or points is this classifier maximally uncertain about?

$x^{(1)}$      $x^{(2)}$      $x^{(3)}$      $x^{(4)}$      $x^{(5)}$

**Solution:** Maximal uncertainty is achieved by points which lie exactly on the decision boundary  $\theta^T x + \theta_0 = 0$ . Thus, we seek points with component sum exactly 0, and this is achieved by both  $x^{(1)}$  and  $x^{(4)}$ .

## 5 Circulation

5. Instead of linear classifiers, we're going to think about circular classifiers (CCs). In  $d$  dimensions a CC is parameterized by a point  $\theta \in \mathbb{R}^d$  and a radius  $\theta_0$ . We will classify a point  $x$  as positive if  $\|x - \theta\|_2 \leq \theta_0$  where

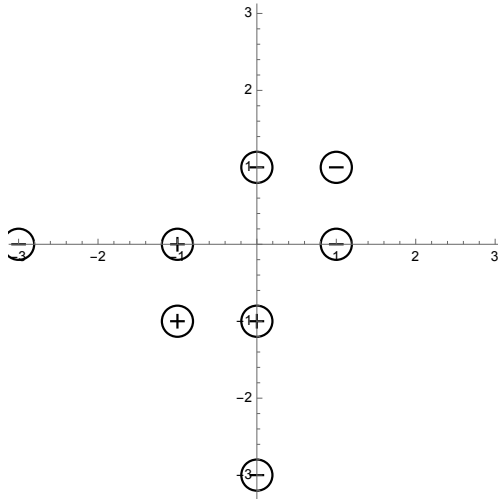
$$\|x - \theta\|_2 = \sqrt{\sum_{j=1}^d (x_j - \theta_j)^2} .$$

- (a) (2 points) Let's start in 1D! Describe the set of points that would be classified as positive by a CC with parameters  $\theta = (2), \theta_0 = 1$ .

**Solution:** Points in the range  $[1, 3]$ .

$$\begin{aligned} \theta &= 2, \theta_0 = 1 \\ \implies \|x - 2\|_2 &\leq 1 \\ \implies -1 &\leq x - 2 \leq 1 \\ \implies 1 &\leq x \leq 3 \end{aligned}$$

- (b) (3 points) Here is a data set in two dimensions. Provide  $\theta$  and  $\theta_0$  for a CC that correctly classifies all the points.



$$\theta = \underline{(-.5, -.5)} \qquad \theta_0 = \underline{\mathbf{1}}$$

**Solution:** There is more than one solution, but one way to derive the above is to first use the three positively labeled points to set-up three equations with three unknowns to get a range for  $\theta$ , and then use the negatively labeled points to get a range on  $\theta_0$ . In particular, the fact that  $(0, -1), (-1, 0), (-1, -1)$  have positive labels imply the

Name: \_\_\_\_\_

following inequalities (here we have  $\theta = (\theta_1, \theta_2)$ ):

$$\begin{aligned}(0 - \theta_1)^2 + (-1 - \theta_2)^2 &\leq \theta_0^2 \\ (-1 - \theta_1)^2 + (0 - \theta_2)^2 &\leq \theta_0^2 \\ (-1 - \theta_1)^2 + (-1 - \theta_2)^2 &\leq \theta_0^2.\end{aligned}$$

Subtracting the third line from the first line gives:

$$\begin{aligned}\theta_1^2 + (-1 - \theta_1)^2 &\leq 0 \\ \implies \theta_1 &\leq -0.5.\end{aligned}$$

Subtracting the third line from the second line similarly gives  $\theta_2 \leq -0.5$ . Let us now pick  $\theta = (-0.5, -0.5)$ . Substituting this into all three inequalities results in the following condition for  $\theta_0$

$$\begin{aligned}\theta_0^2 &\geq \frac{1}{2} \\ \implies \theta_0 &\geq \frac{1}{\sqrt{2}}.\end{aligned}$$

( $\theta_0$  cannot be negative so we only need to consider positive values). Now we need to use the negatively classified values to refine the bound on  $\theta_0$ . It is clear that we only need to consider the points  $(1, 0)$ ,  $(0, 1)$ , since the other points are further away from the center  $\theta = (-0.5, -0.5)$ . This gives

$$\begin{aligned}(1 - \theta_1)^2 + (0 - \theta_2)^2 &> \theta_0^2 \\ (0 - \theta_1)^2 + (1 - \theta_2)^2 &> \theta_0^2 \\ \implies \theta_0 &< \frac{\sqrt{10}}{2}.\end{aligned}$$

So  $\theta_0 = 1$  will do.

Name: \_\_\_\_\_

- (c) (2 points) Circe wants to optimize CC's using gradient descent. Her first thought is to use 0-1 loss, but she decides against it. Why would 0-1 loss be a poor choice?

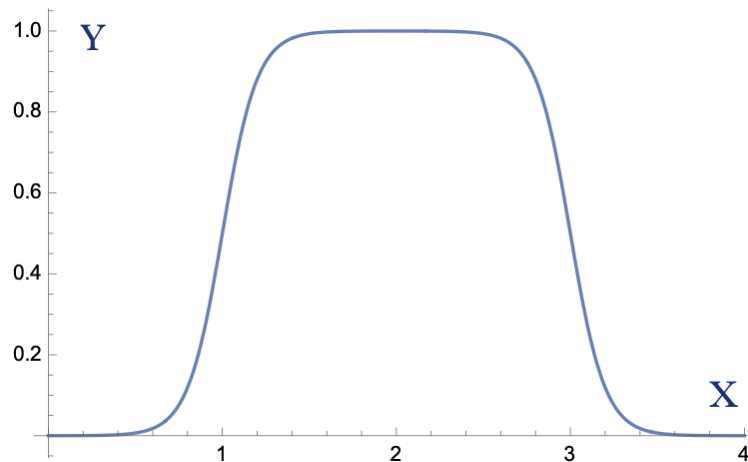
**Solution:** Since the 0-1 loss only takes on 2 values, it is necessarily a piece-wise constant function of  $\theta, \theta_0$ , so the derivative is 0 everywhere, except where it doesn't exist.

Note that the answer is NOT that the 0-1 loss is non-differentiable. ReLU and absolute value (think L1 regularization) are also non-differentiable, but are ubiquitous in the gradient descent-optimized neural net world. The derivative of the 0-1 loss is well-defined almost everywhere - it just gives no information about which direction to update the parameter in.

- (d) (3 points) Circe decides to emulate logistic regression and use hypothesis class

$$h_c(x; \theta, \theta_0) = \sigma(10\theta_0 - 10\|x - \theta\|_2)$$

On the axes below, sketch  $\sigma(10 - 10\|x - 2\|_2)$ . Clearly label the X and Y axes.



**Solution:**

$$\begin{aligned} h_c(x; 2, 1) &= \sigma(10 - 10\|x - 2\|_2) \\ &= \begin{cases} \sigma(10 - 10(x - 2)) & \text{if } x > 2 \\ \sigma(10 - 10(2 - x)) & \text{otherwise} \end{cases} \\ &= \begin{cases} \sigma(30 - 10x) & \text{if } x > 2 \\ \sigma(10x - 10) & \text{otherwise} \end{cases} \\ &= \begin{cases} \sigma(-10(x - 3)) & \text{if } x > 2 \\ \sigma(10(x - 1)) & \text{otherwise} \end{cases} \end{aligned} \tag{1}$$

which we know passes through points  $(3, 0.5)$  and  $(1, 0.5)$

Name: \_\_\_\_\_

(e) (6 points) Let's try to optimize negative log likelihood for this hypothesis class. Just using your intuition, for each case, indicate whether a small gradient step on the specified parameter would increase or decrease its value in order to improve the objective value.

i. Let  $\theta = 5$ ,  $\theta_0 = 2$ ,  $x = 4$ ,  $y = 0$ .

$\theta$  should  **increase**  decrease  stay the same

$\theta_0$  should  increase  **decrease**  stay the same

ii. Let  $\theta = 5$ ,  $\theta_0 = 2$ ,  $x = 0$ ,  $y = 1$ .

$\theta$  should  increase  **decrease**  stay the same

$\theta_0$  should  **increase**  decrease  stay the same

iii. Let  $\theta = 5$ ,  $\theta_0 = 2$ ,  $x = 5$ ,  $y = 1$ .

$\theta$  should  increase  decrease  **stay the same**

$\theta_0$  should  increase  **decrease**  stay the same

**Solution:** Part i:

The sigmoid function for this combination of  $x, \theta, \theta_0$  is  $\sigma(10\theta_0 - 10 * \|4 - \theta\|_2) = \sigma(10\theta_0 - 10 * (\theta - 4)) = \sigma(10\theta_0 + 40 - 10\theta)$ . So, to make the sigmoid closer to  $y = 0$  via gradient descent, the term  $z$  inside  $\sigma(z)$  should be smaller, meaning  $\theta$  should increase and  $\theta_0$  should decrease.

Another way to think about this:

$$10\|x - \theta\|_2 \leq 10\theta_0$$

$$10\|4 - 5\|_2 = 10 \leq 20$$

Currently,  $x$  is being classified as positive since the above relation holds. In a small step of gradient descent,  $\theta$  would increase to make  $\|x - \theta\|_2$  a larger value, and  $\theta_0$  would decrease to make  $\|x - \theta\|_2$  relatively larger compared to  $\theta_0$ . Both of these changes make the classification closer to negative ( $y = 0$ ).

By similar reasoning, for part ii:

$\sigma(10\theta_0 - 10 * \|0 - \theta\|_2) = \sigma(10\theta_0 - 10(\theta - 0)) = \sigma(10\theta_0 - 10\theta)$ . To make this term inside the sigmoid larger so that the classification is closer to  $y = 1$ ,  $\theta$  should decrease and  $\theta_0$  should decrease.

Thinking about this the other way:

$$10\|0 - 5\|_2 = 10 * 5 \not\leq 10 * 2$$

The point is currently incorrectly being classified as negative, so  $\theta$  should decrease to make  $\|x - \theta\|_2$  smaller and  $\theta_0$  should increase to make  $\|x - \theta\|_2$  greater than  $\theta_0$ , in order to be classified as positive.



Name: \_\_\_\_\_

Finally, part iii:

$\sigma(10\theta_0 - 10\|5 - \theta\|_2) = \sigma(10\theta_0 - 10(\theta - 5)) = \sigma(10\theta_0 - 10(5 - \theta))$ . Since  $\theta$  can be written either way in this case, there is not a clear answer for how it should change according to this formula. Let's think about the classification:

$$10\|5 - 5\|_2 = 0 \leq 10 * 1$$

This is a correct classification, since it is a positive classification and  $y = 1$ . The sigmoid value  $\sigma(20 - 0)$  is already very close to 1.  $\theta$  will not change because we can see that the term including  $\theta$  in the sigmoid is a 0. Since the gradient of the loss with respect to  $\theta$  depends on a factor of  $x - \theta$ , the derivative is 0 at this point and  $\theta$  can not change via gradient descent for this particular combination of  $\theta$  and  $x$ .

- (f) (5 points) Now, let's do it more formally and for a general  $d$ -dimensional input  $x \in \mathbb{R}^d$ . Let  $g$  be the output of the hypothesis,

$$g = \sigma(10\theta_0 - 10\|x - \theta\|_2)$$

and recall that

$$\frac{\partial \mathcal{L}_{nll}(g, y)}{\partial \theta} = \frac{\partial \mathcal{L}_{nll}(\sigma(z), y)}{\partial \theta} = (\sigma(z) - y) \frac{\partial z}{\partial \theta} .$$

Note also that

$$\frac{\partial}{\partial \theta} \|x - \theta\|_2 = -\frac{x - \theta}{\|x - \theta\|_2} .$$

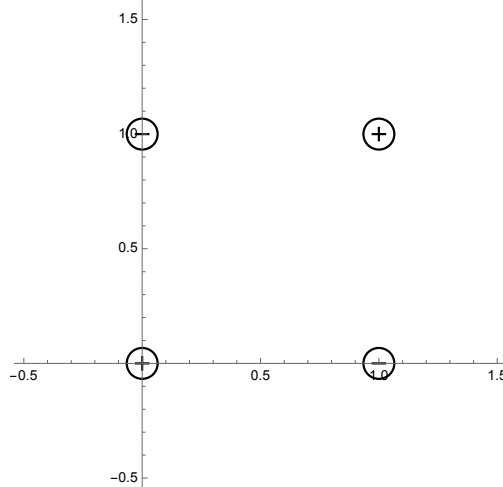
Given current values  $\theta, \theta_0$  and a single data point  $(x, y)$ , write an expression for a gradient-descent update to  $\theta$  in terms of  $\theta, \theta_0, x$ , and  $g$  and step-size  $\eta$ .

**Solution:**

$$\begin{aligned} \theta &:= \theta - \eta \frac{\partial \mathcal{L}_{nll}(g, y)}{\partial \theta} \\ &:= \theta - \eta(g - y) \frac{\partial z}{\partial \theta} \\ &:= \theta - \eta(g - y) \cdot -10 \cdot -\frac{x - \theta}{\|x - \theta\|_2} \\ &:= \theta - 10\eta(g - y) \frac{x - \theta}{\|x - \theta\|_2} \end{aligned}$$

Name: \_\_\_\_\_

- (g) (5 points) Consider a very simple neural network, with two CC units in the hidden layer, and a single output with a sigmoid. Our goal is to use this network to correctly classify the familiar negative XOR data set:



Indicate values for the weights in this network that will result in a correct classification of all four points.

$$h(x) = \sigma(w_0^2 + w_1^2 h_c(x; (w_{11}^1, w_{21}^1), w_{01}^1) + w_2^2 h_c(x; (w_{12}^1, w_{22}^1), w_{02}^1))$$

$$w_{01}^1 = \underline{\mathbf{.5}} \quad w_{11}^1 = \underline{\mathbf{1}} \quad w_{21}^1 = \underline{\mathbf{1}}$$

$$w_{02}^1 = \underline{\mathbf{.5}} \quad w_{12}^1 = \underline{\mathbf{0}} \quad w_{22}^1 = \underline{\mathbf{0}}$$

$$w_0^2 = \underline{\mathbf{-10}} \quad w_1^2 = \underline{\mathbf{20}} \quad w_2^2 = \underline{\mathbf{20}}$$

**Solution:** The weight values specified above place two circles of radius 0.5, one centered at (0, 0) and the other at (1, 1). For the data points (0, 0) or (1, 1), the values of the two  $h_c$  functions, defined in part d as

$$h_c(x; \theta, \theta_0) \equiv \sigma(10\theta_0 - 10\|x - \theta\|_2),$$

are either greater than 0.95 and less than 0.05 or vice-versa. For the data points (0, 1) and (1, 0), the values of the  $h_c$  functions are both less than 0.05. Scaling the results by 20, then summing, and then offsetting by 10, yields  $\approx -10$  for the negative data points and  $\approx +10$  for the positive data points. Applying  $\sigma$  to the scaled and offset result for each data point gives nearly 1 for the plus data points and nearly zero for the minus data points.