

STAT 451 Midterm Exam NetID (mine is “jgillett” from “jgillett@wisc.edu”): \_\_\_\_\_

First name Last name (please write clearly so Gradescope’s OCR can read your name):

Indicate your lecture with by filling one circle completely:

- TuTh 8:00-9:15
- TuTh 11:00-12:15

Instructions:

1. Please sit in columns with two empty seats separating columns.
2. Do not open the exam until I say “go.”
3. Put away everything except a pencil or pen, a calculator, and your two one-page (two sides each) notes sheets.
4. For questions with circles in front of the answers, fill in one circle completely. Ok?
  - Yes
  - No
5. Show your work. Correct answers without at least a minimal version of the work normally required may receive no credit. For a question with an answer box “,” write only the answer in the box; work should be outside the box.
6. If you continue writing or do not turn in your exam when I say time is up, you risk a penalty. (The alternative, that you get more time than your peers, is unfair.)
7. If a question is ambiguous, resolve the ambiguity in writing. We will consider grading accordingly. e.g.
  - I think “average” refers to the population mean  $\mu$  (not the sample mean  $\bar{X}$ ).
  - I think “linear regression” refers to OLS, not ridge or lasso.

Question	Points	Earned
Q0 (cover)	2	
Q1	12	
Q2	8	
Q3	12	
Q4	12	
Q5	14	
Q6	12	
Q7	16	
Q8	12	
Total	100	

1. Here are some questions about decision trees.

- (a) Consider a classification decision tree node containing the set of examples  $S = \{(\mathbf{x}, y)\}$  where  $\mathbf{x} = (x_1, x_2, x_3)$ :

$S$			
$x_1$	$x_2$	$x_3$	$y$
2	11	12	1
3	6	14	1
0	8	17	0
4	10	15	1
1	7	13	0
5	9	16	1

- i. The entropy of this node in bits is

**ANSWER:**

The node's  $y$  values are 1, 1, 0, 1, 0, 1, so  $f_{ID3}(S) = \frac{1}{|S|} \sum_{(\mathbf{x}, y) \in S} y = \frac{1}{6}(1 + 1 + 0 + 1 + 0 + 1) = \frac{2}{3}$ .

$H(S) = \frac{2}{3}(-\log_2(\frac{2}{3})) + \frac{1}{3}(-\log_2(\frac{1}{3})) \approx -\frac{2}{3}(-0.585) + \frac{1}{3}(-1.585) \approx 0.918$

- ii. The (feature, threshold) pair  $(j, t)$  that yields the best split for this node is feature

$j =$  $\text{ and threshold } t =$  .

**ANSWER:**

Using feature  $j = 1$  and threshold  $t = 1.5$  (or any  $t \in [1, 2)$ ) splits  $S$  into  $S_- = \{(\mathbf{x}, y) \in S | x^{(j)} \leq t\} = \{0, 0\}$  and its complement  $S_+ = \{(\mathbf{x}, y) \in S | x^{(j)} > t\} = \{1, 1, 1, 1\}$ , each of which has entropy 0.

- (b) Consider a regression decision tree with `max_depth=1` (that is, the root node is split once into two leaves) made from the set of examples  $S = \{(\mathbf{x}, y)\}$  where  $\mathbf{x} = x_1$ :

$S$

$x_1$	$y$
0	10
1	11
2	21
3	22
4	23
5	24

What value does this tree predict for  $x_1 = 4.5$ ?  $\hat{y} =$

**ANSWER: 22.5**

The best split uses feature  $j = 1$  and threshold  $t = 1.5$ , yielding a left subtree containing the first two examples and a right subtree containing the last four. Making a prediction with  $x_1 = 4.5$  would use the right subtree. Its average  $y$  is 22.5, so the tree would predict  $\hat{y} = 22.5$ .

2. Here are questions about feature engineering.

(a) Consider the data -5, 5, 5, 5, 5, which have these summary statistics:

- minimum -5
- mean 3
- median 5
- maximum 5
- (population) standard deviation 4

Do standardization rescaling on feature  $x$ :

(input) $x$	(output) $x_{\text{rescaled}}$
-5	ANSWER: $-2$
5	ANSWER: $\frac{1}{2}$
5	ANSWER: $\frac{1}{2}$
5	ANSWER: $\frac{1}{2}$
5	ANSWER: $\frac{1}{2}$

Or, if we mistakenly use the sample standard deviation  $s = \sqrt{20} \approx 4.47$  instead of the population standard deviation  $\sigma = 4$ , we get this answer (which should receive only a tiny grading penalty, like 0.25 point):

$x$	$x_{\text{rescaled}}$
0	-1.79
5	0.45
5	0.45
5	0.45
5	0.45

- (b) Use one-hot encoding to transform the categorical feature `power_source` into binary features with reasonable names that are in alphabetical order.

(input) <code>power_source</code>	(output)
grid	
solar	
generator	
grid	

ANSWER:

(input) <code>power_source</code>	(output)		
	<code>generator</code>	<code>grid</code>	<code>solar</code>
grid	0	1	0
solar	0	0	1
generator	1	0	0
grid	0	1	0

3. Consider the gradient descent algorithm.

- (a) Consider applying gradient descent with step size  $\alpha = 0.1$  to find the  $\mathbf{x}$  that minimizes the function  $f(\mathbf{x}) = f((x^{(1)}, x^{(2)})) = (x^{(1)} - 1)^2 + (x^{(2)} + 2)^2$  starting from  $\mathbf{x}_0 = (0, 0)$ . Find the value  $\mathbf{x}_1$  after one iteration.

$\mathbf{x}_1 =$

ANSWER:

$\nabla f(\mathbf{x}) = (2(x^{(1)} - 1), 2(x^{(2)} + 2))$ , which is  $(-2, 4)$  at  $\mathbf{x} = (0, 0)$ .

Move to  $\mathbf{x}_1 = \mathbf{x}_0 - \alpha \nabla f(\mathbf{x}_0) = (0, 0) - (-0.2, 0.4) = (0.2, -0.4)$ .

(b) Mark each statement as true or false.

- Gradient descent can fail to converge on a convex function if step size  $\alpha$  is such that it gets stuck in a cycle, oscillating between two or several values.  
 True ANSWER:   
 False
- For a non-convex function, gradient descent can fail to converge by descending without bound.  
 True  
 False

ANSWER:

I regret asking this question.

$y = -e^x$  is an example of a non-convex function on which gradient descent can descend without bound, so originally I marked it True.

In fall 2024 a student persuaded me I should have marked it False. I don't remember the reason; it involved me understanding "fail to converge" to mean "fail to converge to the global minimum" when it doesn't mean that.

By spring 2025 I am back to True.

Waffling is bad, so I withdraw the question.

- Gradient descent can fail to converge if it gets stuck in a local minimum.  
 True ANSWER:   
 False
- Gradient descent can fail to converge on a convex function if the step size  $\alpha > 0$  is too large.  
 True ANSWER:   
 False

4. Consider 3-NN (three nearest neighbors) using the Minkowski distance with  $p = 1$ .

(a) Find the distance from  $\mathbf{z} = (2, 2)$  to each of the other points  $\mathbf{x}$ :

$\mathbf{x}$	$y$	Distance from $\mathbf{z}$ to $\mathbf{x}$
$(-1, -1)$	1	ANSWER: 6
$(0, 1)$	0	ANSWER: 3
$(1, 0)$	0	ANSWER: 3
$(2, 3)$	1	ANSWER: 1

(b) How does 3-NN classify  $\mathbf{z}$ ?

$$\hat{y} = \boxed{\phantom{0}}$$

ANSWER: 0

The 3-NN have distances 1, 3, and 3, and  $y$  values 0, 0, and 1. 0 is most frequent.

(c) How does weighted 3-NN classify  $\mathbf{z}$ ?

$$\hat{y} = \boxed{\phantom{1}}$$

ANSWER: 1

The reciprocals of the 3 smallest distances  $\{d_i\}$  are  $\left\{\frac{1}{d_i}\right\}$ :  $\frac{1}{1}$ ,  $\frac{1}{3}$ , and  $\frac{1}{3}$ . The weight on  $\hat{y} = 1$  is 1. The weight on  $\hat{y} = 0$  is  $\frac{1}{3} + \frac{1}{3} = \frac{2}{3}$ .

(d) What  $y$  value does 3-NN regression predict for  $\mathbf{z}$ ?

$$\hat{y} = \boxed{\phantom{\frac{1}{3}}}$$

ANSWER:  $\hat{y} = \frac{1}{3}(0 + 0 + 1) = \frac{1}{3}$

5. Consider finding the linear regression line by hand for the points  $\{(\mathbf{x}, y)\} = \{(x, y)\} = \{(1, 2), (2, 4), (3, 3)\}$ . Match each mathematical quantity on the left with its value on the right. (Hint: Very little arithmetic is required.)



(a)  $\mathbf{X} =$

(A) (B) (C) (D) (E) (F) (G)

G: must be  $3 \times 2$

(b)  $\mathbf{X}^T \mathbf{X} =$

(A) (B) (C) (D) (E) (F) (G)

B: must be  $2 \times 2$

(c)  $\mathbf{y} =$

(A) (B) (C) (D) (E) (F) (G)

D: from data

(A) 3.5

(B)  $\begin{bmatrix} 3 & 6 \\ 6 & 14 \end{bmatrix}$

(d)  $\mathbf{w} =$

(A) (B) (C) (D) (E) (F) (G)

C: must be  $2 \times 1$

(C)  $\begin{bmatrix} 2 \\ 0.5 \end{bmatrix}$

(D)  $\begin{bmatrix} 2 \\ 4 \\ 3 \end{bmatrix}$

(e)  $\hat{\mathbf{y}} =$

(A) (B) (C) (D) (E) (F) (G)

E: must be  $3 \times 1$  and D is taken by  $\mathbf{y}$

(E)  $\begin{bmatrix} 2.5 \\ 3.0 \\ 3.5 \end{bmatrix}$

(F) 2

(G)  $\begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix}$

(f)  $f_{\mathbf{w},b}(3) =$

(A) (B) (C) (D) (E) (F) (G)

A: must be scalar and 2 is taken by  $b$

(g) intercept =

(A) (B) (C) (D) (E) (F) (G)

F: first element of  $\mathbf{w}$

(My answers above are mostly from considering dimensions. With a little more work, the matrix arithmetic can be done to find the answers.)

6. Here are some questions on support vector machines.

- (a) Suppose we have a soft-margin SVM for which  $\mathbf{w} = (2, 3)$  and  $b = -1$ . How does the SVM classify  $(\mathbf{x} = (1, 1), y = 1)$ ?

ANSWER:

$$\mathbf{w}\mathbf{x} + b = (2, 3) \cdot (1, 1) + (-1) = 4 > 0, \text{ so } \hat{y} = 1.$$

- (b) Suppose we have some training data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  (where  $\mathbf{x}_i$  is 2D) in matrices  $\mathbf{X}$  and  $\mathbf{y}$ . We have plotted the data with  $y = -1$  examples red and  $y = 1$  examples blue. Which line of code gives the best model for predicting new examples?

For each question, write the best answer from among these lines labeled “A” through “H”.

A: `clf = svm.SVC(kernel='linear', C=1); clf.fit(X, y)`

B: `clf = svm.SVC(kernel='linear', C=1000); clf.fit(X, y)`

C: `clf = svm.SVC(kernel='rbf', C=1, gamma=1); clf.fit(X, y)`

D: `clf = svm.SVC(kernel='rbf', C=1, gamma=10); clf.fit(X, y)`

E: `clf = svm.SVC(kernel='euclidean', C=1, gamma=2); clf.fit(X, y)`

F: `clf = svm.SVC(kernel='euclidean', C=1000, gamma=2); clf.fit(X, y)`

G: `clf = svm.SVC(kernel='gini', C=1); clf.fit(X, y)`

H: `clf = svm.SVC(kernel='gini', C=1000); clf.fit(X, y)`

- i. (A) (B) (C) (D) (E) (F) (G) (H)

The red points are scattered between  $x_1 = 0$  and  $x_1 = 2\pi$  and roughly along  $x_2 = \sin x_1$ , a wave. The blue points are scattered over the same  $x_1$  interval and roughly along  $x_2 = \sin x_1 + 1$ , a wave 1 higher than the first wave.

ANSWER: C

- ii. (A) (B) (C) (D) (E) (F) (G) (H)

The data consist of two clouds of points, one red and one blue, that are linearly-separable except for a few outliers of each color.

ANSWER: A

- iii. (A) (B) (C) (D) (E) (F) (G) (H)

The data are mixed red and blue points scattered randomly in the disk  $x_1^2 + x_2^2 \leq 1$ .

ANSWER: D

7. Mark each statement True or False.

(a) In linear regression, a reasonable alternative to the cost function *mean squared error* =  $\frac{1}{N} \sum_{i=1}^N [f_{\mathbf{w},b}(\mathbf{x}_i) - y_i]^2$  is *sum of squared error* =  $\sum_{i=1}^N [f_{\mathbf{w},b}(\mathbf{x}_i) - y_i]^2$ .

True ANSWER: ●

The two functions will have the same minimum, as they differ only by the constant  $\frac{1}{N}$ .

False

(b) For the soft-margin SVM with decision boundary  $\mathbf{w}\mathbf{x} + b = 0$  where  $\mathbf{w} = (1, 2)$  and  $b = 3$ , the example  $(\mathbf{x}, y) = ((4, 5), -1)$  has hinge loss 18.

True ANSWER: ●

Hinge loss is  $\max(0, 1 - y_i(\mathbf{w}\mathbf{x} + b)) = \max(0, 1 - (-1)((1 \cdot 4 + 2 \cdot 5) + 3)) = 18$ .

False

(c) For training data  $\{(\mathbf{x}, y)\}$  such that  $\mathbf{x}_i \neq \mathbf{x}_j$  for all  $i$  and  $j$ , we can build a 3NN model that classifies the training examples without error.

True

False ANSWER: ●

For example, if there are only  $N = 3$  examples with  $y_1 = 0, y_2 = 0, y_3 = 1$ , all three examples will be classified as the majority  $y$  value, namely 0.

(d) If we train a hard-margin linear SVM on linearly separable data, then discard training examples which are support vectors, and then train a new SVM on the remaining examples, the first SVM will have a wider “road” than the second.

True

False ANSWER: ●

The first “road” will be narrower than the second.

(e) A linear SVM with decision boundary  $(1, 2, 2) \cdot \mathbf{x} - 2 = 0$  has a smaller margin between +1 and -1 support vectors than one with boundary  $(1, 4, 8) \cdot \mathbf{x} + 3 = 0$ .

True

False ANSWER: ●

The margin for the first SVM is  $\frac{2}{\|\mathbf{w}\|} = \frac{2}{\sqrt{1^2+2^2+2^2}} = \frac{2}{3}$ , while the margin for the second is  $\frac{2}{\|\mathbf{w}\|} = \frac{2}{\sqrt{1^2+4^2+8^2}} = \frac{2}{9}$ .

(f) Every decision tree regression function is a step function.

Hint: A *step function* is a function that is constant over each of one or more intervals.

True ANSWER: ●

False

(g) Every  $k$ -NN regression function is a step function.

Hint: A *step function* is a function that is constant over each of one or more intervals.

True ANSWER: ●

False

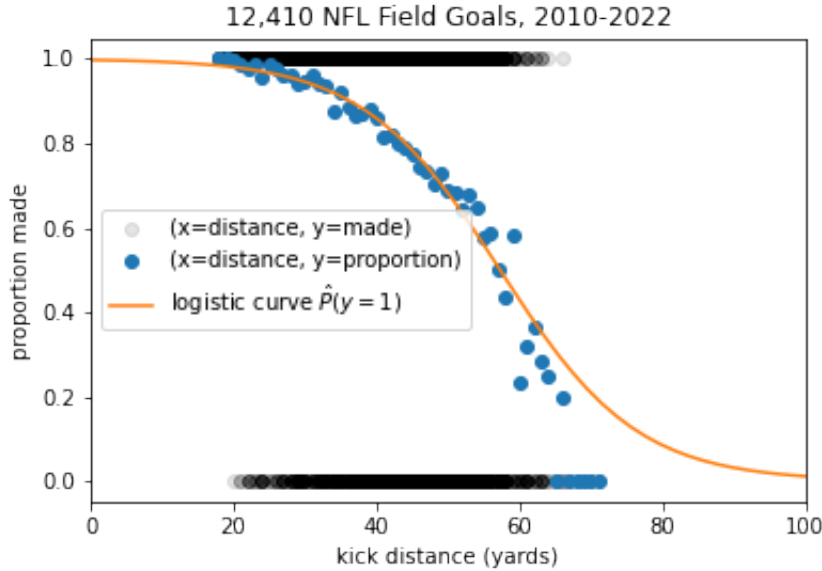
(h) In logistic regression, we use the natural log function to facilitate finding a closed-form expression for the coefficients  $\mathbf{w}$  and  $b$  in terms of the data.

True

False ANSWER: ●  
We do not have such a closed-form expression.

8. Consider a logistic regression model with  $\mathbf{w} = (1, 2)$  and  $b = 0$ .

- (a) From the logistic regression model represented in the figure, estimate the likelihood of an NFL field kicker making two field goals in a row, one from 20 yards and one from 60. We may suppose these attempts are independent and make other reasonable simplifying assumptions.



$P(\text{NFL kicker makes the two field goals}) \approx \boxed{\phantom{0.4}}$ . (Give an estimate.)

ANSWER:

The likelihood is

$$L_{\mathbf{w},b} = \prod_{i=1}^N l_{\mathbf{w},b}(\mathbf{x}_i, y_i) = P(y = 1|x = 20) \cdot P(y = 1|x = 60) \approx 1 \cdot 0.4 = 0.4.$$

Any answer in  $[0.35, 0.45]$  is ok too.

- (b) In calculating the coefficients for a logistic regression model, why do we minimize negative log likelihood instead of maximizing likelihood? Mark each statement as a true or false.
- A product of probabilities can overflow in fixed-precision computer arithmetic.
    - True
    - False ANSWER: ●
  - A product of probabilities can underflow in fixed-precision computer arithmetic.
    - True ANSWER: ●
    - False
  - The natural log of a product is naturally expressed as a sum, and differentiating a sum is easier than differentiating a product.
    - True ANSWER: ●
    - False

iv. The natural log is strictly increasing, so maximizing the likelihood is the same as minimizing the negative log likelihood.

True ANSWER: ●

False