

(1 point for easily legible writing on this cover sheet.)

NetID: \_\_\_\_\_ (mine is `jgillett` from `jgillett@wisc.edu`)

Last name: \_\_\_\_\_ First name: \_\_\_\_\_

First name and last name: \_\_\_\_\_ (I am asking for these separately above and together here to use with optical character recognition software.)

Mark your lecture with “X”:

\_\_\_\_\_ TuTh 1:00-2:15

STAT 451 Midterm Exam 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 one-page (two sides) notes sheet.
4. Show your work. Correct answers without at least a minimal version of the work normally required may receive no credit.
5. The exam ends when I call time. If you continue writing after I call time, you risk a penalty. (The alternative, that you get more time than your peers, is unfair.)
6. 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)	1	
Q1	10	
Q2	5	
Q3	15	
Q4	15	
Q5	10	
Q6	12	
Q7	10	
Q8	22	
Total	100	



4. 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	0
1	6	14	0
0	8	17	1
3	10	15	0
4	7	16	1
5	9	13	0

i. The entropy of this node in bits is \_\_\_\_\_.

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

- (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	12
3	13
4	23
5	24

What value does this tree predict for  $x_1 = 4.5$ ? \_\_\_\_\_

5. Here are questions about feature engineering.

(a) Consider the data 0, 5, 5, 5, 5, which have these summary statistics:

- minimum 0
- mean 4
- median 5
- maximum 5
- (population) standard deviation 2

Do standardization rescaling on feature **x**:

(input) <b>x</b>	(output) <b>x_rescaled</b>
0	
5	
5	
5	
5	

(b) Use one-hot encoding to transform the categorical feature **activity** into binary features with reasonable names.

(input) <b>activity</b>	(output)
swim	
bike	
run	
bike	

6. Consider 3-NN (three nearest neighbors) with the negative cosine similarity distance.

(a) Find the distance from  $\mathbf{z} = (3, 4)$  to each of the other points  $\mathbf{x}$ :

$\mathbf{x}$	$y$	Distance from $\mathbf{z}$ to $\mathbf{x}$
$(-4, 3)$	1	
$(0, 1)$	0	
$(1, 0)$	0	
$(6, 8)$	1	

Hint: Notice the “negative” in “negative cosine similarity.”

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

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

7. In linear regression we minimize the *mean squared error* (MSE).

(a) Find the MSE for the points  $(0, 0)$  and  $(1, 2)$  relative to the line  $\hat{y} = f_{\mathbf{w}, b}(\mathbf{x}) = \mathbf{w}\mathbf{x} + b$ , where  $\mathbf{w} = 2$  and  $b = 3$ .

(b) For the best-fitting line for these data,  $\mathbf{w} = \underline{\hspace{2cm}}$  and  $b = \underline{\hspace{2cm}}$ .

8. Mark each statement true or false by circling the appropriate choice.
- (a) TRUE / FALSE A support vector machine on 2D  $\mathbf{x}$  with decision boundary  $(3, 4) \cdot \mathbf{x} + 5 = 0$  has a smaller margin between  $+1$  and  $-1$  support vectors than one with boundary  $(5, 12) \cdot \mathbf{x} + 13 = 0$ .
  - (b) TRUE / FALSE Logistic regression sends an example  $\mathbf{x}$  through a linear function to get a real number, which it sends through an exponential function to get a positive number, which it sends through a logistic function to predict a probability between 0 and 1. (Then it optionally compares that probability to a threshold to predict  $\hat{y} = 0$  or  $\hat{y} = 1$ .)
  - (c) TRUE / FALSE The linear regression model is not sensitive to the signs of labels  $\{y_i\}$  because it minimizes mean *squared* error.
  - (d) TRUE / FALSE A decision tree node containing the examples  $\{(\mathbf{x}, y)\} = \{((1, 2), 0), ((3, 4), 1), ((5, 6), 0)\}$  has lower entropy than the one containing the examples  $\{((7, 8), 0), (9, 10), 1)\}$ .
  - (e) TRUE / FALSE Weighted  $k$ NN evaluates a new example  $\mathbf{x}$  using  $\mathbf{x}$ 's  $k$  nearest neighbors weighted with weights proportional to the corresponding distances to  $\mathbf{x}$ .
  - (f) TRUE / FALSE If we train an SVM on linearly separable data, then discard all support vectors, and then train a new SVM on the remaining examples, the first SVM will have a larger margin between  $+1$  and  $-1$  examples than the second.
  - (g) TRUE / FALSE The hinge loss function of a soft-margin SVM gives a nonzero value for any  $\mathbf{x}$  such that  $\mathbf{w}\mathbf{x} + b \geq 0$ .
  - (h) TRUE / FALSE If stochastic gradient descent (SGD) and gradient descent (GD) are both run in the same amount of time on a well-behaved function of some parameters over a large  $N$  of training examples each having a small  $D$  of features, SGD can use more iterations with a smaller learning rate  $\alpha$  than GD.
  - (i) TRUE / FALSE  $1 + 1 = 2$
  - (j) TRUE / FALSE Ridge regression tends to set most coefficients to zero.
  - (k) TRUE / FALSE Training data are used to set parameters, validation data are used to choose hyperparameter settings and/or models, and test data are used to evaluate the chosen model.