# CS 178 Midterm Exam
## Machine Learning and Data Mining: Fall 2018
## Monday November 5th, 2018

**Your name:**

Solutions

**Row/Seat Number:**

Lectum

**Your ID #(e.g., 123456789)**

**UCINetID (e.g.ucinetid@uci.edu)**

PanTeater @ uci.edu

- Please put your name and ID **on every page.**

- Total time is 50 minutes. READ THE EXAM FIRST and organize your time; don't spend too long on any one problem.

- Please **write clearly** and **show all your work**.

- If you need clarification on a problem, please raise your hand and wait for the instructor or TA to come over.

- You may use **one** sheet containing handwritten notes for reference, and a (basic) calculator.

- Turn in your notes and any scratch paper with your exam.

## Problems

**Total,** *(42 points.)*

## Bayes Classifiers, *(10 points.)*

Consider the table of measured data given at right. We will use the
two observed features $x_1$, $x_2$ to predict the class $y$. Each feature
can take on one of three values, $x_i \in \{a, b, c\}$.
In the case of a tie, we will prefer to predict class $y = 0$.

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| a | c | 0 |
| a | b | 0 |
| b | a | 0 |
| b | b | 0 |
| a | b | 1 |
| b | c | 1 |
| b | c | 1 |
| c | c | 1 |

(1) Write down the probabilities learned by a naïve Bayes classifier: *(4 points.)*

$p(y = 0)$ :  $\frac{1}{2}$          $p(y = 1)$ :  $\frac{1}{2}$

$p(x_1 = a \mid y = 0)$ :  $\frac{1}{2}$          $p(x_1 = a \mid y = 1)$ :  $\frac{1}{4}$

$p(x_1 = b \mid y = 0)$ :  $\frac{1}{2}$          $p(x_1 = b \mid y = 1)$ :  $\frac{1}{2}$

$p(x_1 = c \mid y = 0)$ :  $\emptyset$          $p(x_1 = c \mid y = 1)$ :  $\frac{1}{4}$

$p(x_2 = a \mid y = 0)$ :  $\frac{1}{4}$          $p(x_2 = a \mid y = 1)$ :  $\emptyset$

$p(x_2 = b \mid y = 0)$ :  $\frac{1}{2}$          $p(x_2 = b \mid y = 1)$ :  $\frac{1}{4}$

$p(x_2 = c \mid y = 0)$ :  $\frac{1}{4}$          $p(x_2 = c \mid y = 1)$ :  $\frac{3}{4}$

(2) Using your naïve Bayes model, compute: *(3 points.)*
$p(y = 1 \mid x_1 = b, x_2 = c)$ :                    $p(y = 0 \mid x_1 = b, x_2 = c)$ :

Calculate  $p(y=0, x=bc) = Q_0 = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{4}$   and   $p(y=1, x=bc) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{3}{4}$

Then  $p(y=0 \mid bc) = \frac{Q_0}{Q_1 + Q_0} = \frac{1}{1+3} = \frac{1}{4}$   and   $p(y=1 \mid bc) = \frac{Q_1}{Q_1 + Q_0} = \frac{3}{4}$.

(3) Compute the probabilities $p(y = 1 \mid x_1 = b, x_2 = c)$ and $p(y = 0 \mid x_1 = b, x_2 = c)$ for a joint
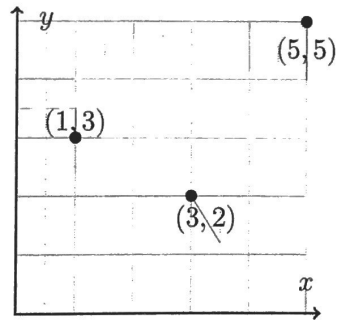Bayes model trained on the same data. *(3 points.)*

By inspection,  $p(y=1 \mid bc) = 1$
$p(y=0 \mid bc) = \emptyset$.

3

## Linear and Nearest Neighbor Regression, *(12 points.)*

Consider the data points shown at right, for a regression problem
to predict $y$ given a scalar feature $x$.

In the case of ties, we prefer to use the
data to the left (smaller values of $x$).



(1) Compute **training** MSE of a 1-nearest neighbor predictor. *(3 points.)*

$\emptyset$ .

(2) Compute the **leave-one-out** cross-validation error (MSE) of a 1-nearest neighbor predictor.
*(3 points.)*

Leave out (1,3) → predict 2    error 1

(3,2)         3  →        1     → $\frac{1}{3}(1^2 + 1^2 + 3^2) = \frac{11}{3}$ .

(5,5)         2           3        MSE

(3) Compute the **leave-one-out** cross-validation error (MSE) of a 2-nearest neighbor predictor.
*(3 points.)*

Leave out    predict      error      MSE

(1,3)         $3\frac{1}{2}$    $\frac{1}{2}$

(3,2)  =)     4      →      2     → $\frac{1}{3}((\frac{1}{2})^2 + 2^2 + (2\frac{1}{2})^2) = \frac{7}{2}$

(5,5)         $2\frac{1}{2}$    $2\frac{1}{2}$

(4) Compute the **leave-one-out** cross-validation error (MSE) of a linear regressor, e.g., a model
of the form $f(x) = \theta_0 + \theta_1 x$. *(3 points.)*

Leave out  ↓ predict  → error     MSE

(1,3)         −1         4

(3,2)         4          2     → $\frac{1}{3}[4^2 + 2^2 + 4^2] = 12$ .

(5,5)         1          4

## Multiple Choice, *(10 points.)*

Here, assume that we have $m$ data points $y^{(i)}$, $x^{(i)}$, $i = 1 \ldots m$, each with $n$ features, $x^{(i)} = [x_1^{(i)} \ldots x_n^{(i)}]$. For each of the choices below, will it likely increase, decrease, or have no effect on overfitting (circle your choice)? If you think it is equally likely to go either way, pick *No Effect*.
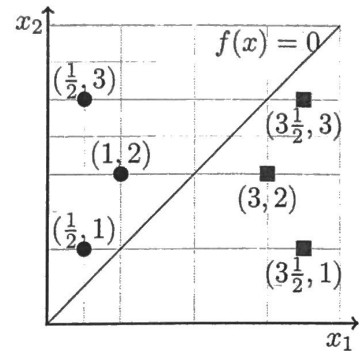
1   Gathering more labeled training data                                    **(Reduce)**   Increase   No Effect

2   For a 3-nearest neighbor classifier, use $2 \times m$ training          Reduce   **(Increase)**   No Effect
    data by copying (duplicating) each data point.
        *Effectively decreases K=3 to K=1*

3   For a 3-nearest neighbor classifier, use $2 \times n$ features          Reduce   Increase   **(No Effect)**
    per data point by copying (duplicating) the features.
        *Doesn't change distances*

4   Increasing $k$ for a k-nearest neighbor classifier                      **(Reduce)**   Increase   No Effect

5   For a linear regressor, use $2 \times m$ training data by              **(Reduce)**   Increase   No Effect
    adding $m$ all-zero ($x$ and $y$) data points.
        *Regularizes a bit - soln likely to pass through (0,0) ⇒ $\Theta_0 \approx \emptyset$.*

6   For a linear regressor, use $2 \times n$ features per data             Reduce   Increase   **(No Effect)**
    point by adding $n$ all-zero features to each.
        *These features cannot be used by the model*

7   For a linear regressor, use $2 \times n$ features per data             Reduce   **(Increase)**   No Effect
    point by adding $n$ random values to each.

8   Adding another layer to a Neural Network                               Reduce   **(Increase)**   No Effect

9   Changing the activation function of hidden nodes                       Reduce   Increase   **(No Effect)**

10  Switching from soft to hard margin SVMs                                Reduce   **(Increase)**   No Effect

## Support Vector Machines, *(10 points.)*

Suppose we are learning a linear support vector machine with two real-valued features $x_1$, $x_2$ and binary target $y \in \{-1, +1\}$. We observe training data (pictured at right):

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0.5 | 1 | -1 |
| 1 | 2 | -1 |
| 0.5 | 3 | -1 |
| 3 | 2 | +1 |
| 3.5 | 1 | +1 |
| 3.5 | 3 | +1 |

*(figure at right, with points labeled $(\frac{1}{2},3)$, $(1,2)$, $(\frac{1}{2},1)$, $(3\frac{1}{2},3)$, $(3,2)$, $(3\frac{1}{2},1)$, and line $f(x)=0$)*

Our linear classifier takes the form

$$f(x; w_1, w_2, b) = \text{sign}(w_1 x_1 + w_2 x_2 + b).$$

(1) For given line $x_1 = x_2$ that perfectly separates the data, list the support vectors. *(2 points.)*

By inspection, SVs are $(\frac{1}{2},1)$ and $(3\frac{1}{2},3)$
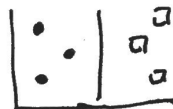
(closer to decision boundary)

(2) Derive the parameter values $w_1, w_2, b$ of this $f(x)$ using the support vectors. What is the length of the margin? *(3 points.)*

$(0,0)$ on bdry $\Rightarrow \quad w_1 \cdot \emptyset + w_2 \cdot \emptyset + b = \emptyset \quad \Rightarrow \quad b = \emptyset.$

SVs $\Rightarrow$

$w_1 \cdot \frac{1}{2} + w_2 \cdot 1 + b = -1$

$w_1 \cdot 3\frac{1}{2} + w_2 \cdot 3 + b = +1$

$\Rightarrow \quad w_1 = 2$

$w_2 = -2$

Margin $= \frac{\sqrt{2}}{2}$ by inspection, or formula $\frac{2}{\sqrt{w^T \cdot w}} = \frac{2}{\sqrt{8}} = \frac{1}{\sqrt{2}}$
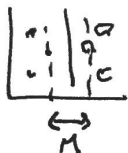
(3) Consider the *best* linear-SVM classifier; one that separates the data and has the largest margin. Sketch the boundary in the above figure, and list the support vectors here. *(2 points.)*

SVs: $(1,2)$ and $(3,2)$.

(4) Derive the parameter values $w_1, w_2, b$ of this $f(x)$ using these support vectors. What is the length of the margin? *(3 points.)*

Then,

$w_1 \cdot 2 + w_2 \cdot (\text{anything}) + b = \emptyset.$

$w_1 \cdot 1 + w_2 \cdot (\text{"}) + b = -1$

$w_1 \cdot 3 + w_2 \cdot (\text{"}) + b = +1$

$\Rightarrow$

$w_1 = 1$

$w_2 = \emptyset.$

$b = -2$

Margin $M = 2$ by inspection, or $\frac{2}{\sqrt{w^T w}} = \frac{2}{1} = 2.$

9