

FINAL

CS273a Midterm Exam

Introduction to Machine Learning: Fall 2012

Tuesday December 11th, 2012

Your name:

SOLUTIONS

Name of the person in front of you (if any):

Name of the person to your right (if any):

- 50
- Total time is 1:45. 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 to come over.
 - Turn in any scratch paper with your exam.

Problem 1: VC Dimension (10p)

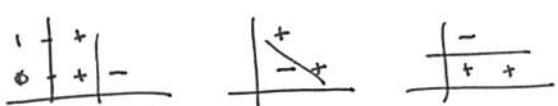
Argue by example / counterexample what is the VC dimension of each of the following classifiers.

- (S_p) (a) A perceptron on two *binary* features

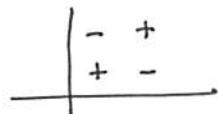
$$\text{Binary features} \Rightarrow x_i \in \{0, 1\}$$

$$VC\ Dim = 3$$

Can shatter 3 points



Cannot shatter 4 points

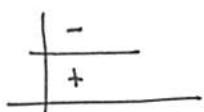


- (S_p) (b) A decision stump on two *binary* features

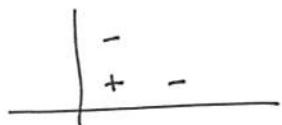
$$\text{Binary features} \Rightarrow x_i \in \{0, 1\}$$

$$VC\ Dim = 2$$

Can shatter 2 points:



cannot shatter 3



(No axis-aligned split can separate)

Problem 2: Decision Trees (12p)

We plan to use a decision tree to predict an outcome y using four features, x_1, \dots, x_3 . We observe six training patterns, each of which we represent as $[x_1, x_2, x_3]$ (so, "010" means $x_1 = 0, x_2 = 1, x_3 = 0$). We observe the training data,

$$y = 0 : [100], [111], [001]$$

$$y = 1 : [110], [110], [011]$$

You may find the following values useful (although you may also leave logs unexpanded):

$$\log_2(1) = 0 \quad \log_2(2) = 1 \quad \log_2(3) = 1.59 \quad \log_2(4) = 2$$

$$\log_2(5) = 2.32 \quad \log_2(6) = 2.59 \quad \log_2(7) = 2.81 \quad \log_2(8) = 3$$

(3p)

- (a) What is the entropy of y ?

$$\frac{3}{6} \log \frac{6}{3} + \frac{3}{6} \log \frac{6}{3} = 1 \text{ bit}$$

(3p)

- (b) Which variable would you split first? Justify your answer.

	x_1		x_2		x_3		By inspection, x_2 has the lower expected entropy \Rightarrow higher info gain
$y=0$	=0	=1	=0	=1	=0	=1	
$y=1$	001	100 111	100 001	111	100	111 001	
	011	110 110	-	110 110	110 110	011	

(3p)

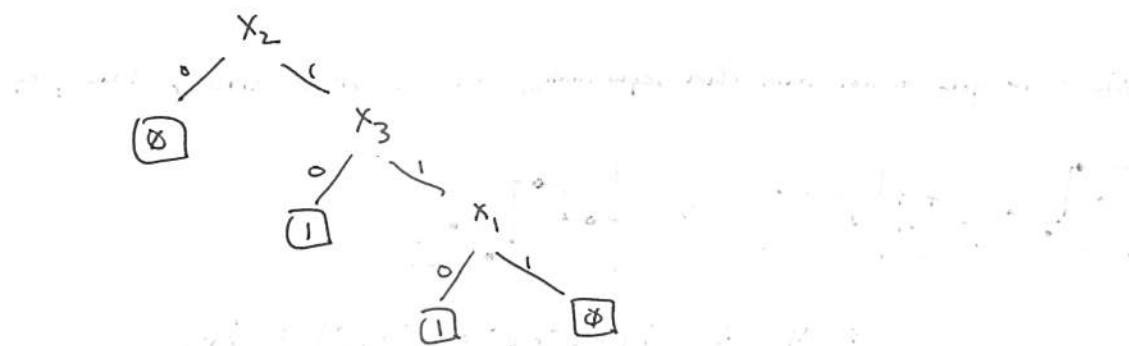
- (c) What is the information gain of the variable you selected in part (b)?

$$H(y|x_1) \geq H(y|x_3) \geq H(y|x_2)$$

$$\begin{aligned} H(y) &= \left[\frac{2}{6} H(0) + \frac{4}{6} H(1) \right] \\ &= 1 - \left[\phi + \frac{2}{3} \left(\frac{1}{4} \log 4 + \frac{3}{4} \log \frac{4}{3} \right) \right] \\ &= 1 - \frac{1}{6} \log 4 - \frac{1}{2} \log \frac{4}{3} = \frac{2}{3} - \frac{1}{2} \log (0.4) \end{aligned}$$

(3p)

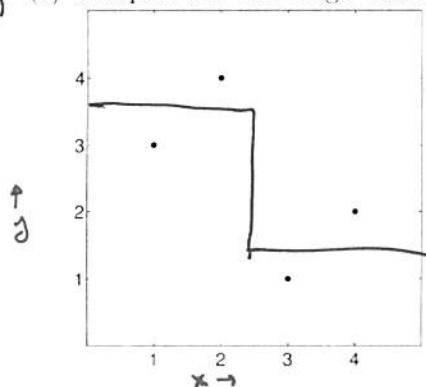
- (d) Draw the rest of the decision tree learned on these data.



Problem 3: Gradient Boosting (9p)

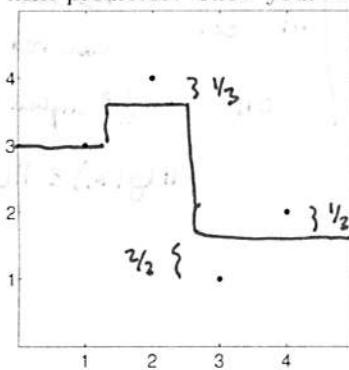
Consider the following data set consisting of four points; for convenience, the data are repeated in each part.

- (3p) (a) Compute the best single decision stump regressor function, to minimize mean squared error.



(this predictor has $MSE = \frac{1}{4} \cdot (4 \cdot 1^2) = 1$)

- (3p) (b) Now, we wish to create a gradient boosted ensemble of decision stumps to minimize MSE. Starting from the decision stump learned in 9a), and using a learning rate of 1, what is the next predictor? Show your work.

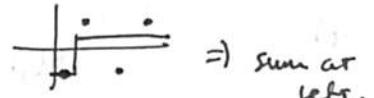


After (a), we have error residual



If we fit a decision stump, we get

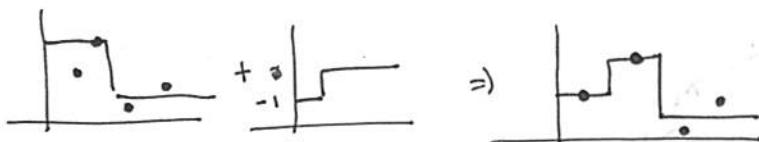
$$MSE = \frac{1}{4} \cdot (0^2 + 1^2 + 1^2 + 0.5^2) = \frac{3}{18}$$



= sum at
left.

- (3p) (c) Is the resulting ensemble the best possible ensemble of two decision stumps? If yes, why? If not, give a better ensemble.

No - if you do not train them sequentially, you can get a better predictor, e.g.



$$\Rightarrow MSE \text{ is } \frac{1}{4} \left(0^2 + 0^2 + 1/2^2 + 1/2^2 \right) = 1/8$$

Problem 4: Naïve Bayes (10p)

We plan to use a naïve Bayes classifier to predict an outcome y using four features, x_1, \dots, x_3 . We observe five training patterns, each of which we represent as $[x_1, x_2, x_3]$ (so, "010" means $x_1 = 0, x_2 = 1, x_3 = 0$). We observe the training data,

$$y = 0 : [000], [111]$$

$$y = 1 : [100], [010], [001]$$

(4p)

- (a) Compute (& show) all the necessary probabilities for a naive Bayes model.

$$p(y) = 3/5$$

$$p(x_1 | y=0) = 1/2$$

$$p(x_1 | y=1) = 1/3$$

$$p(x_2 | y=0) = 1/2$$

$$p(x_2 | y=1) = 1/3$$

$$p(x_3 | y=0) = 1/2$$

$$p(x_3 | y=1) = 1/3$$

(3p)

- (b) Suppose you observe $x = [110]$. What class (value of y) would you predict? Show your work.

$$\begin{aligned} p(y=0) p(x | y=0) &= 2/5 \cdot 1/2 \cdot 1/2 \cdot 1/2 = 1/20 \\ p(y=1) p(x | y=1) &= 3/5 \cdot 1/3 \cdot 1/3 \cdot 2/3 = 2/45 \end{aligned} \quad \Rightarrow \text{predict } y=0.$$

(3p)

- (c) Suppose you observe $x = [100]$. Compute the posterior probability, $p(y|x = 100)$.

$$p(y=0) p(x | y=0) = 2/5 \cdot 1/2 \cdot 1/2 \cdot 1/2 = 1/20 = \frac{9}{5 \cdot 2 \cdot 2 \cdot 3 \cdot 3}$$

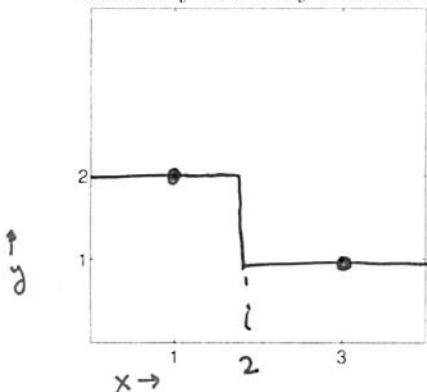
$$p(y=1) p(x | y=1) = 3/5 \cdot 1/3 \cdot 1/3 \cdot 2/3 = 4/45 = \frac{4 \cdot 4}{5 \cdot 3 \cdot 3 \cdot 2 \cdot 2}$$

$$\Rightarrow p(y=1 | x) = \frac{16}{16+9} = \frac{16}{25}$$

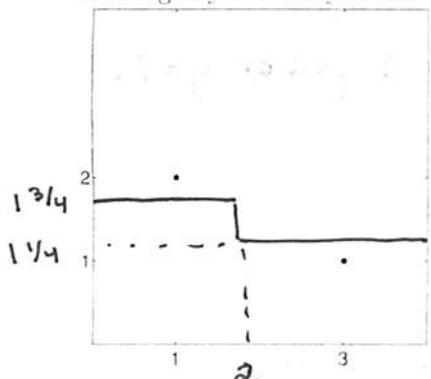
Problem 5: Bagging (9p)

Consider the data set, consisting of two data points, given in each part.

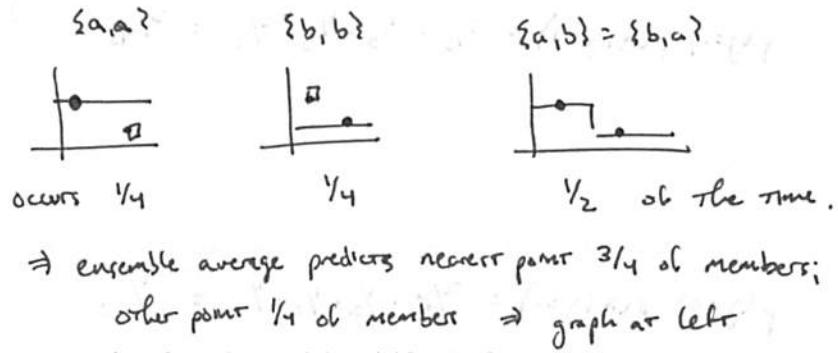
- (3p) (a) Draw the regression function (predicted values for all x) using a nearest-neighbor regressor. Label any necessary values on your graph.



- (4p) (b) Suppose that we create a very large ensemble of *bagged* nearest-neighbor regressors, using data set draws of size two. Compute the regression function of the complete ensemble, again labeling any necessary values.



There are four possible draws \Rightarrow 3 unique data sets:



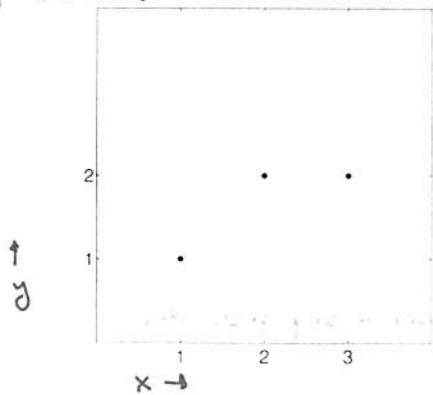
- (c) Is the model in (b) simpler or more complex than the model in (a)? Why?

- (2p) Simpler:
- (1) Bagging tends to reduce complexity / overfitting
 - you can see training error has gone up (& the data are no longer memorized)
 - (2) The function is "simpler" - it is closer to a constant predictor.

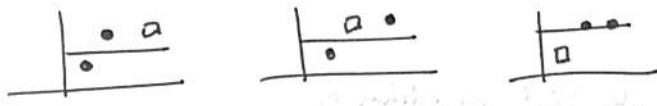
Problem 6: Cross-validation (8p)

Consider the following data points, copied in each part. We wish to perform linear regression to minimize mean squared error.

- (4p) (a) Compute the leave-one-out cross-validation error of a zero-order (constant) predictor.



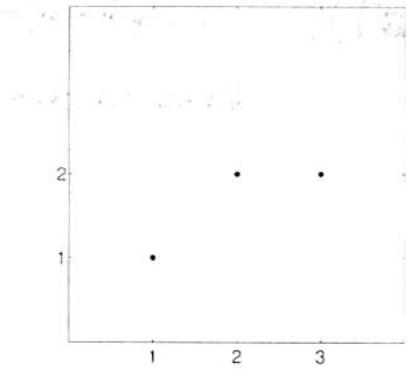
Leave out each data point \Rightarrow



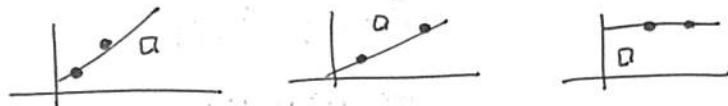
$$(y_2)^2 + (y_2)^2 + (1)^2$$

$$\Rightarrow \frac{1}{3} [(y_2)^2 + (y_2)^2 + 1^2] = \frac{1}{2}.$$

- (4p) (b) Compute the leave-one-out cross-validation error of a first-order (linear) predictor.



Leave out each \Rightarrow



$$1^2 + (\frac{1}{2})^2 + 1^2$$

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

Problem 7: Latent space models (8p)

Suppose that, as in HW5, we wish to model a collection of text documents using a latent space model. For interpretability, we would like our latent representation to be non-negative. As one solution, we use an exponential transform to ensure positive values, giving the model

$$x_j^{(i)} \approx \sum_k \exp(U_{ik}) \exp(V_{kj})$$

Give a stochastic gradient descent algorithm to learn this model, minimizing the mean squared error in the predicted values. Include all necessary details for the implementation.

A simple SGD algorithm is

```

① Initialize  $U, V$  to something at random
② Choose a stopping criterion (e.g. # of steps) and a step size  $\alpha$ .
③ while (!stop) {
    for i=1..m
        for j=1..d
             $E = (x_j^{(i)} - \sum_k \exp(U_{ik}) \exp(V_{kj}))^2$ ; // compute signed error.
             $\tilde{U} = U$ ;  $\tilde{V} = V$ ; // save old values
            for k=1..K:
                 $U_{ik} \leftarrow U_{ik} + \alpha \nabla_{U_{ik}} J$ 
                 $V_{kj} \leftarrow V_{kj} + \alpha \nabla_{V_{kj}} J$ 
            end
        end
    end
}

```

where $J(\cdot) = (x_j^{(i)} - \sum_k \exp(u) \exp(v))^2$ is the squared error loss

and

$$\nabla_{U_{ik}} J = -E \cdot \exp(\tilde{U}_{ik}) \exp(\tilde{V}_{kj})$$

$$\nabla_{V_{kj}} J = -E \cdot \exp(\tilde{U}_{ik}) \exp(\tilde{V}_{kj})$$

are the derivatives with respect
to each element of u, v .



notation should

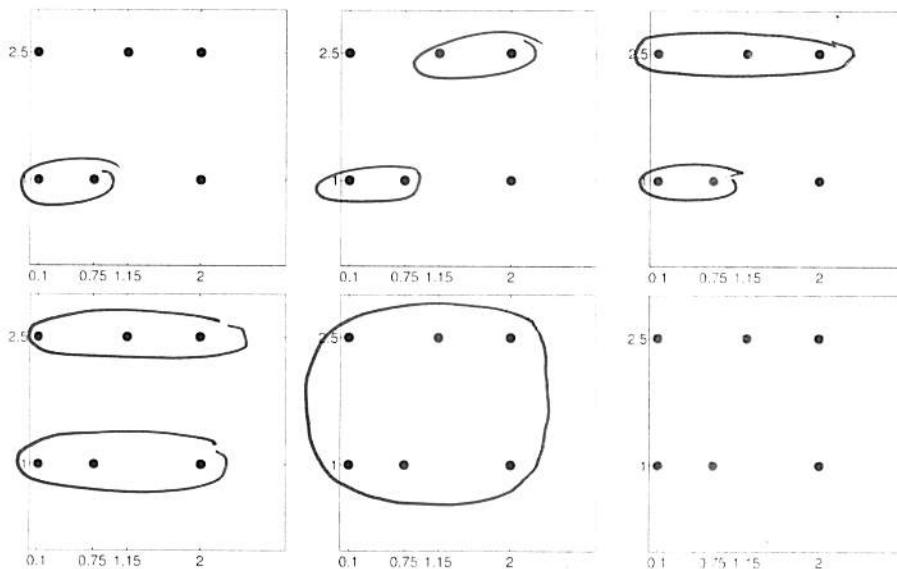
really be

~~$\frac{\partial}{\partial u_{ik}}$~~ $\frac{\partial}{\partial U_{ik}}$ and $\frac{\partial}{\partial V_{kj}}$ instead

Problem 8: Clustering (10p)

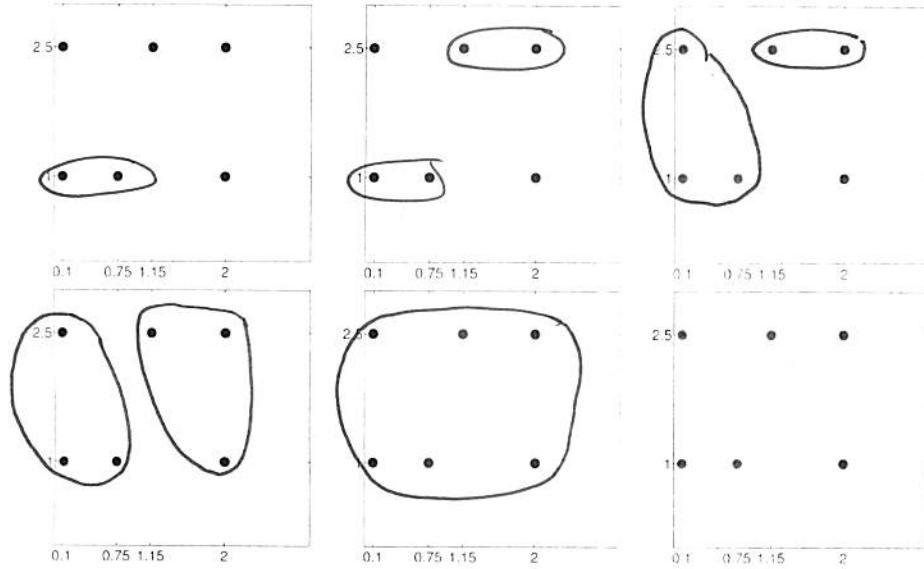
Consider the two-dimensional data points plotted in each panel. In this problem, we will cluster these data using two different algorithms, where each panel is used to show an iteration or step of the algorithm.

- (5p) (a) Execute the hierarchical agglomerative clustering (linkage) algorithm on these data points, using "single linkage" for the cluster scores. Stop when converged, or after 6 steps, whichever is first. Show each step separately in a panel.



Join nearest cluster.
Cluster distance =
nearest points in the
two clusters.

- (5p) (b) Now execute hierarchical agglomerative clustering on the data points, but use "complete linkage" for the cluster scores. Stop when converged, or after 6 steps, whichever is first. Show each step separately in a panel.



Join nearest two clusters.
Cluster distance =
furthest distance between
two points in the clusters.

From board:

$$\sqrt{(.85)^2 + (1.5)^2} = 1.724$$

$$\sqrt{(.65)^2 + (1.5)^2} = 1.635$$

CS273a Final Exam
Introduction to Machine Learning: Fall 2013
Thursday December 12th, 2013

Your name:

SOLUTIONS

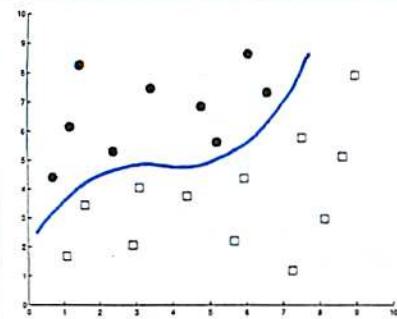
Your UCInetID (all caps):

Your Seat (row and number):

- Total time is 1:50. READ THE EXAM FIRST and organize your time; don't spend too long on any one problem.
- Closed book; one page of (your own) notes
- Please write clearly and show all your work.
- If you need clarification on a problem, please raise your hand and wait for the instructor to come over.
- Turn in any scratch paper with your exam.

Problem 1: Separability

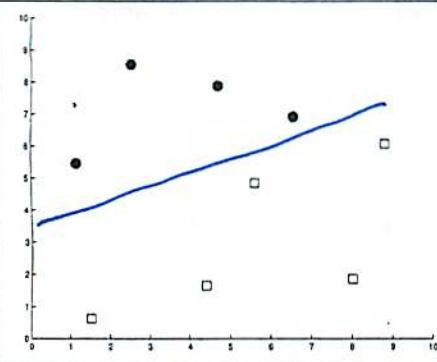
For each of the following examples of training data, sketch a classification boundary that separates the data. State whether or not the data are linearly separable, and if not, give a set of features that would allow the data to be separated.



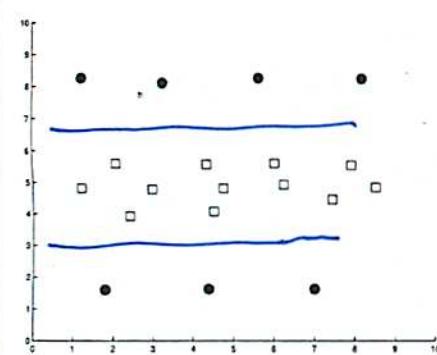
No, not linearly separable.

Based on the curve, $x_2 = a + bx_1 + cx_1^2 + dx_1^3$

The features $[1 \ x_1 \ x_1^2 \ x_1^3 \ x_2]$ should work.



Yes, linearly separable.



No, not linearly separable.

The decision boundary can be obtained as

$$ax_2 + bx_2 + c = 0$$

$\Rightarrow [1 \ x_2 \ x_2^2]$ is enough.

Problem 2:

Select the best choice to complete each statement.

Increasing the number of hidden nodes in a neural network will most likely
(increase decrease not change) the bias.

Decreasing regularization on the weights in logistic regression will most likely
(increase decrease not change) the VC dimension.

Increasing the amount of data will most likely
(increase decrease not change) the variance.

Increasing the regularization on a perceptron will most likely
(increase decrease not change) the bias.

Decreasing the maximum depth of a decision tree will most likely
(increase decrease not change) the VC dimension.

Increasing the depth of a decision tree will most likely
(increase decrease not change) the bias.

Possible argument for no change - if depth is already very large compared to the size of data

The predictions of a k-nearest neighbor classifier
(will will not) be affected by pre-processing to normalize the data.

Reducing the number of features using PCA will most likely
(increase decrease not change) the variance.

Linear regression
(can cannot) be solved using either matrix algebra or gradient descent.

The predictions of a regression tree
(will will not) be affected by pre-processing to normalize the features.

Problem 3: Regression

Suppose that we train a *non-linear* regression model on m data, where our prediction is

$$\hat{y}(x) = a + \exp(bx)$$

for two scalar parameters a, b .

- (a) Write down the formula for the mean-squared error on the training data, and compute its gradient with respect to the parameters.

$$MSE = J(a, b) = \frac{1}{m} \sum_i (y^i - \hat{y}(x^i))^2 = \frac{1}{m} \sum_i (y^i - a - \exp(bx^i))^2$$

$$\nabla J = \begin{bmatrix} \frac{\partial J}{\partial a} & \frac{\partial J}{\partial b} \end{bmatrix}$$

$$\frac{\partial J}{\partial a} = \frac{1}{m} \sum_i (y^i - \hat{y}(x^i)) \cdot (-1). \quad (1)$$

$$\frac{\partial J}{\partial b} = \frac{1}{m} \sum_i (y^i - \hat{y}(x^i)) \cdot (-1) \exp(bx^i) \cdot b. \quad (2)$$

- (b) Give pseudo-code for (batch) gradient descent on this problem. Be sure to specify initialization, the update itself (in enough detail to enable coding), and a stopping condition (again, in enough detail to enable coding).

```

Init a, b : a=0, b=0.      a' = inf, b' = inf.      alpha = stepsize.

while (~done)
    for i=1..m,   y^i = a + exp(bx^i)
        a ← a - alpha * ∂J / ∂a           as in (1)           % take a step (could update stepsize)
        b ← b - alpha * ∂J / ∂b           as in (2)
    done = [(a-a')^2 + (b-b')^2 < ε]       % check for convergence
    a=a'; b=b';                           % save old values

```

- (c) Give at least one advantage of batch gradient descent over stochastic gradient descent. In contrast, when would using stochastic gradient be more appropriate?

Batch - always a descent on J (for sufficiently small α)

→ easy to debug, easy to assess convergence, monotonic.

SGD - more appropriate than batch when m is very large (many data)

- in this case, batch updates will be very slow. (all data processed before each step).

Problem 4: Naïve Bayes

Consider the following table of measured data:

x_1	x_2	x_3	y
1	1	1	0
1	1	0	0
0	1	0	0
1	0	1	0
1	1	1	1
0	1	1	1
0	0	1	1

We will use the three observed features x_1, x_2, x_3 to predict class y . In the case of a tie, we will prefer to predict class $y = 1$.

- (a) Write down the probabilities necessary for a naïve Bayes (NB) classifier:

$$\Pr[y=0] = \Pr[y=1] = 3/7$$

$$\Pr[x_1=1 | y=0] = 3/4 \quad \Pr[x_2=1 | y=0] = 3/4 \quad \Pr[x_3=1 | y=0] = 1/2$$

$$\Pr[x_1=1 | y=1] = 1/3 \quad \Pr[x_2=1 | y=1] = 2/3 \quad \Pr[x_3=1 | y=1] = 1.$$

- (b) Using your NB model, what value of y is predicted given observation $(x_1, x_2, x_3) = (000)$.

Predict $y=0$. ($y=1$ has $\Pr[x_3=0 | y=1] = 0$).

- (c) Using your NB model, what is the probability $p(y=1 | x_1=1, x_2=1, x_3=1)$?

$$= \frac{3/7 \cdot 1/3 \cdot 2/3 \cdot 1}{(1) + 4/7 \cdot 3/4 \cdot 3/4 \cdot 1/2} = \frac{6/9}{6/9 + 9/8} = \frac{48}{48 + 81} = \frac{48}{129}$$

- (d) Using your NB model, what is the probability $p(y=1 | x_1=0)$?

$$= \frac{3/7 \cdot 2/3}{(1) + 4/7 \cdot 1/4} = \frac{2/7}{2/7 + 1/2} = \frac{2/7}{7/14 + 7/14} = \frac{2/7}{14/14} = \frac{2}{14} = \frac{1}{7}.$$

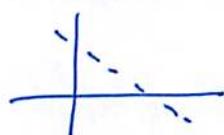
Problem 5: Perceptrons and VC Dimension

In this problem, consider the following perceptron model on two features:

$$\hat{y}(x) = \text{sign}(w_0 + w_1x_1 + w_2x_2)$$

and answer the following questions about the decision boundary and the VC dimension.

- (a) Describe (in words, with diagrams if desired) the decision boundaries that can be realized by this classifier



Any hyperplane - eg, decision boundary is an arbitrary line.

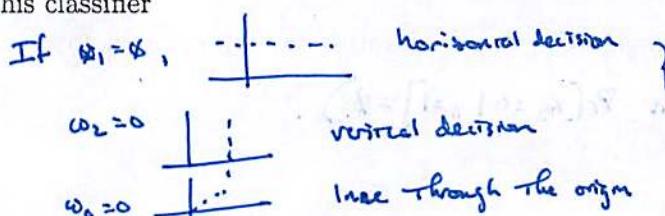
- (b) What is its VC dimension?

3

- VC dim of a perceptron is $d+1$, and $d=2$.

Now suppose that I also enforce an additional condition on the parameters of the model: that at most two of the weights w_i are non-zero (so, at least one weight is zero).

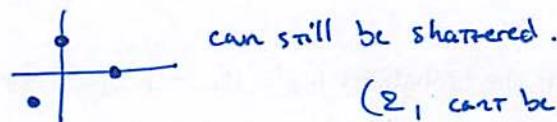
- (c) Describe (in words, with diagrams if desired) the decision boundaries that can be realized by this classifier



Decision boundary can be an arbitrary horizontal or vertical split, or a line through the origin.

- (d) What is its VC dimension?

Still 3:



can still be shattered.

(Σ_1 can't be higher than part (b))

Finally, I enforce that at most one of the weights w_i is non-zero (so, at least two are zero).

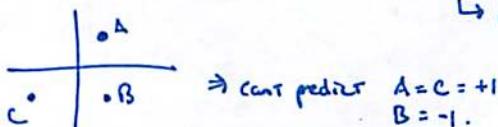
- (e) Describe (in words, with diagrams if desired) the decision boundaries that can be realized by this classifier

$w_0 \neq 0$ - entire plane same decision
 $w_1 \neq 0$ - right vs left half-plane
 $w_2 \neq 0$ - upper vs lower half-plane.

} union of these functions.

- (f) What is its VC dimension?

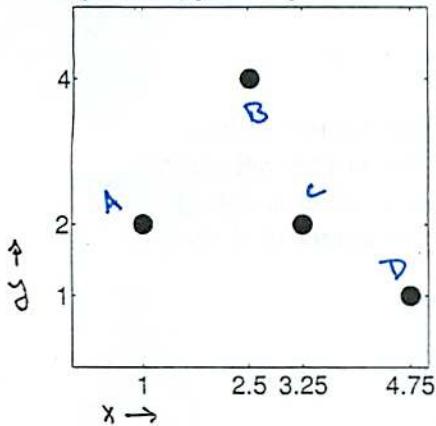
Now 2 - can't shatter above example; not possible to arrange points to get all three ++- patterns...
 ↳ (can't put points on the axes, though).



⇒ can predict $A = C = +1$
 $B = -1$.

Problem 6: Cross-validation

Consider a regression problem for predicting the following data points, using the k-nearest neighbor regression algorithm from class and the homework to minimize mean squared error (MSE). (Note: if you like, you may leave an arithmetic expression, e.g., leave values as “(.6)²”.)



- (a) For $k = 1$, compute the training error on the provided data.

$$\phi.$$

- (b) For $k = 1$, compute the leave-one-out cross-validation error on the data.

nearest nbr:

$$\begin{array}{c} \rightarrow \vec{\text{A}} \\ \rightarrow \vec{\text{B}} \\ \rightarrow \vec{\text{C}} \\ \rightarrow \vec{\text{D}} \end{array} \Rightarrow \frac{1}{4} [2^2 + 2^2 + 2^2 + 1^2] = 13/4$$

- (c) For $k = 3$, compute the training error on the provided data.

neighbors: predict:

$$\begin{array}{ll} \text{A: ABC } & 8/3 \\ \text{B: ABC } & 8/3 \\ \text{C: BCD } & 7/3 \\ \text{D: BCD } & 7/3 \end{array} \Rightarrow \frac{1}{4} [(2/3)^2 + (4/3)^2 + (1/3)^2 + (4/3)^2] = 37/36$$

- (d) For $k = 3$, compute the leave-one-out cross-validation error on the data.

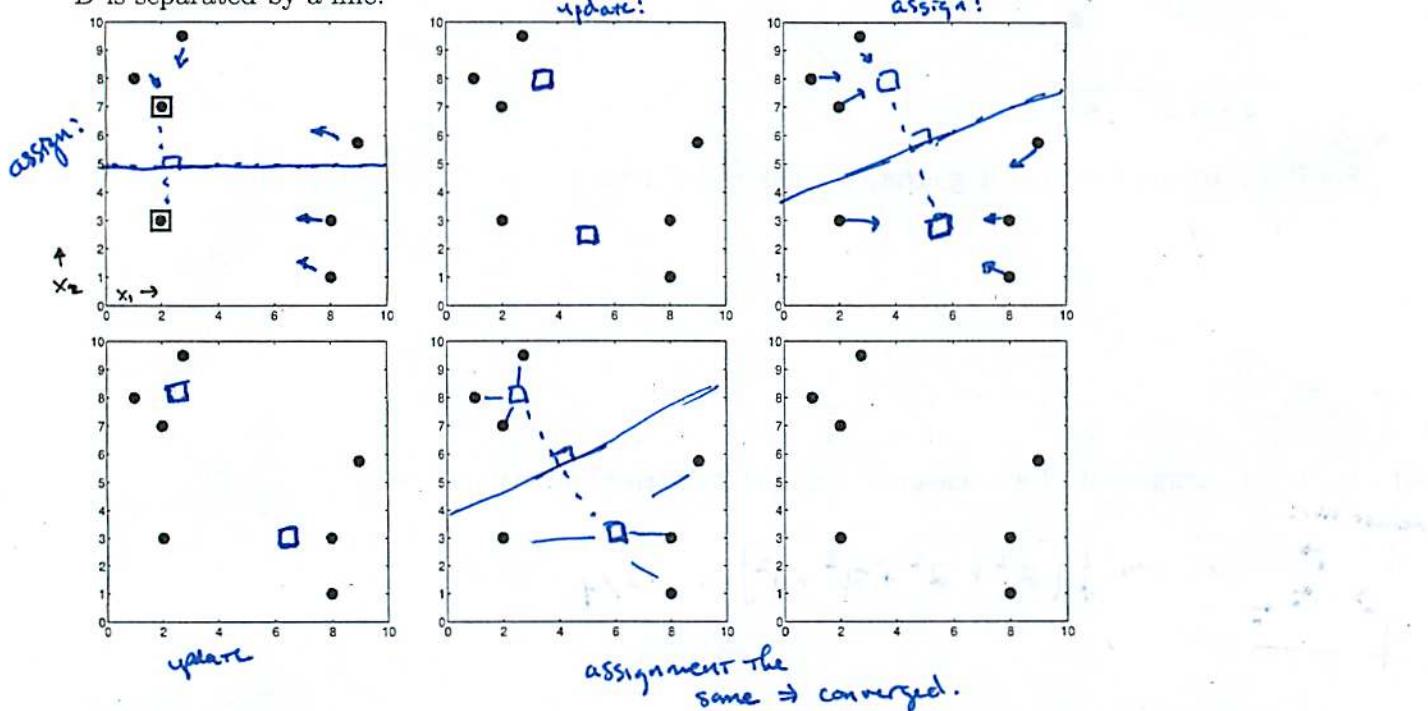
$$\begin{array}{ll} \text{A: BCD } & 7/3 \\ \text{B: ACD } & 5/3 \\ \text{C: ABD } & 7/3 \\ \text{D: ABC } & 8/3 \end{array} \Rightarrow \frac{1}{4} [(1/3)^2 + (7/3)^2 + (1/3)^2 + (5/3)^2] = 76/36$$

Problem 7: Clustering

Consider the two-dimensional data points plotted in each panel. In this problem, we will cluster these data using two different algorithms, where each panel is used to show an iteration or step of the algorithm.

k-means

- (a) Starting from the two cluster centers indicated by squares, perform k-means clustering on the data points. In each panel, indicate (somehow) the data assignment, and in the next panel show the new cluster centers. Stop when converged, or after 6 steps (3 iterations), whichever is first. It may be helpful to recall from our nearest-neighbor classifier that the set of points nearer to A than B is separated by a line.



- (b) Write down the cost function optimized by the k-means algorithm, explaining your notation.

Let μ_c be the center of cluster c .

z_i be the cluster assignment of data point i .

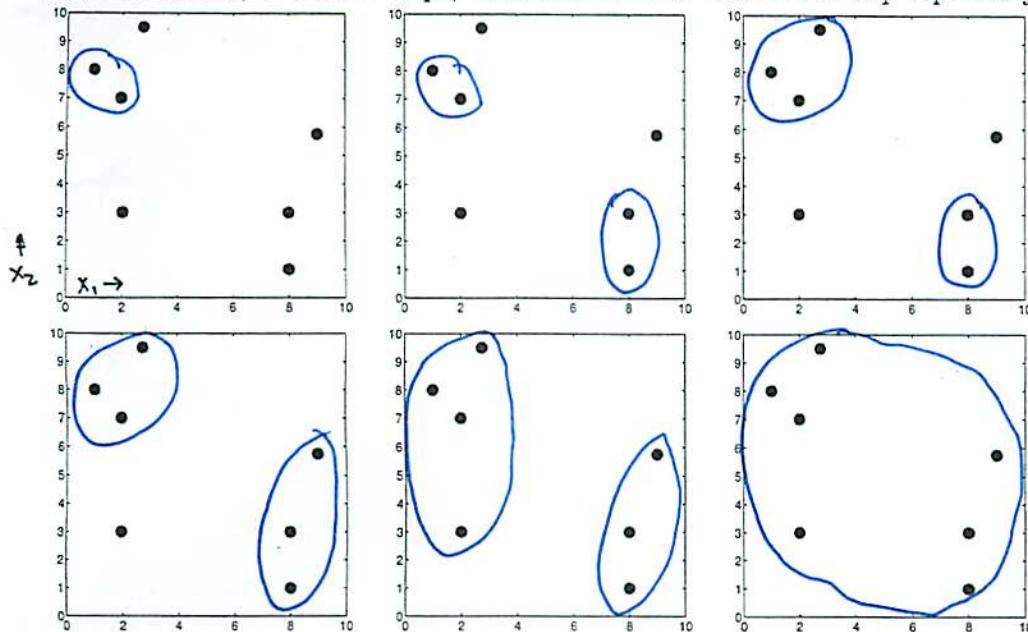
$$J = \sum_{c=1}^m \sum_{i=1}^n \|x^{(i)} - \mu_{z_i}\|^2$$

where $\|\cdot\|^2$,

is the Euclidean distance / length squared.

Linkage

- (a) Now execute the hierarchical agglomerative clustering (linkage) algorithm on these data points, using "complete linkage" (maximum distance) for the cluster scores. Stop when the algorithm would terminate, or after 6 steps, whichever is first. Show each step separately in a panel.



Complete linkage
⇒ join nearest pair of clusters, where distance is given by farthest points.

- (b) What is the algorithmic (computational) complexity of the hierarchical clustering algorithm? Briefly justify your answer.

$$O(n^2)$$

Initial step calculates $\binom{n}{2} = O(n^2)$ pairwise distances

Each of n iterations (one join per step)

requires updating new cluster's distance to remaining clusters

$$\Rightarrow \text{constant } (n-1) + (n-2) + (n-3) + \dots + (1) = O(n^2)$$

$$\text{and, } O(n^2) + O(n^2) = O(n^2)$$

CS273 Final Exam
Introduction to Machine Learning: Winter 2015
Tuesday March 17th, 2015

Your name:

SOLUTIONS

Your UCINetID (e.g., myname@uci.edu):

Your seat (row and number):

- Total time is 1 hour 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 to come over.
- Turn in any scratch paper with your exam.

(This page intentionally left blank)

Problem 1: (6 points) Multiple Choice

For the following questions, assume that we have m data points $y^{(i)}, x^{(i)}, i = 1 \dots m$, each with n features, $x^{(i)} = [x_1^{(i)} \dots x_n^{(i)}]$.

Circle one answer for each:

True or false: Early stopping can be used to reduce overfitting in neural networks.

True or false: The SVM learning algorithm will find the *globally optimal* model with respect to its objective function.

True or false: Increasing k in a k-nearest-neighbor classifier will decrease the bias.

True or false: Increasing the depth of a decision tree classifier will decrease the bias.

True or false: The VC dimension of a perceptron classifier is smaller than the VC dimension of a linear SVM.

True or false: If there exists a set of h instances that cannot be shattered by $f(x)$, then the VC dimension of f is less than h .

Problem 2: (4 points) Short Answer

Give one advantage of the dual (kernel) form of support vector machines over the primal (linear) form, and one advantage of the primal form over the dual.

Dual: may be easier to specify kernel similarity than features for linear classifiers
many kernels work in high / infinite dimensional feature spaces
better if # of features is very high (maybe)
 ∞ -d.m. kernel \Rightarrow nonparametric predictor, often gets better as $m \rightarrow \infty$ (#dara)

Primal: more efficient if $m \gg n$. (lots of data), in computation & model storage
Comp. & model storage are fixed (dont grow with m)
This also means it is often more efficient at test time (if $m \gg n$)
Easy to apply standard algorithms like SGD.

Problem 3: (12 points) Decision Trees

Consider the table of measured data given at right. We will use a decision tree to predict the outcome y using the three features, x_1, \dots, x_3 . In the case of ties, we prefer to use the feature with the smaller index (x_1 over x_2 , etc.) and prefer to predict class 1 over class 0. You may find the following values useful (although you may also leave logs unexpanded):

$$\begin{aligned} \log_2(1) &= 0 & \log_2(2) &= 1 & \log_2(3) &= 1.59 & \log_2(4) &= 2 \\ \log_2(5) &= 2.32 & \log_2(6) &= 2.59 & \log_2(7) &= 2.81 & \log_2(8) &= 3 \end{aligned}$$

x_1	x_2	x_3	y
1	0	0	1
1	1	1	1
0	0	1	1
1	1	0	0
1	1	0	0
0	1	1	0

- (a) What is the entropy of y ?

$$1 \text{ bit}$$

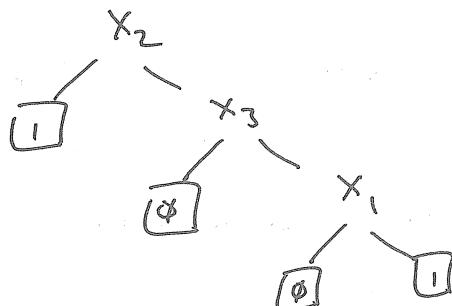
- (b) Which variable would you split first? Justify your answer.

$x_1 =$	$y = 0$	$y = 1$	$x_2 =$	$y = 0$	$y = 1$	$x_3 =$	$y = 0$	$y = 1$
0	011	001	0	X	100 001	0	110 110 011	101
1	110 110	100 111	1	110 110 011	111	1	011	111 001

- (c) What is the information gain of the variable you selected in part (b)?

$$\begin{aligned} I(G) &= 1 - [Y_0 \cdot 0 + Y_1 \cdot (\frac{3}{4} \log \frac{4}{3} + \frac{1}{4} \log 4)] \\ &= 1 - Y_0 \log Y_0 - Y_1 \log Y_1 \\ &= Y_0 \log 3 - Y_1 \log 4 \approx 0.4591 \text{ bits} \end{aligned}$$

- (d) Draw the rest of the decision tree learned on these data.



x_3 next - get 2 right (vs 1)
(can also calculate I_G if desired,
but clear by inspection)

Problem 4: (9 points) Gradient Descent & Latent Space Models

Suppose that, as in lecture, we wish to model a collection of text documents using a latent space model such as Latent Semantic Indexing. However, for interpretability, we would like our latent representation to be non-negative, so that each “direction” can only make a set of words more likely to appear, not less likely. As one solution, we use an exponential transform to ensure positive values, giving the model

$$x_j^{(i)} \approx \sum_k \exp(u_{ik}) \exp(v_{jk})$$

(so, data point i is a nonnegative linear combination u_i of nonnegative directions v_j). We wish to train our model (U, V) to minimize squared error from the observed data, and will train it using gradient descent.

- (a) Write down an expression for J_{ij} , the mean squared error (MSE) of our model for element $x_j^{(i)}$, and for J , the overall MSE of our model.

$$J_{ij} = (x_j^{(i)} - \sum_k \exp(u_{ik}) \exp(v_{jk}))^2$$

$$J = \frac{1}{m} \sum_{ij} J_{ij} \quad (\text{also ok to normalize by } \frac{1}{mn})$$

- (b) Compute the derivatives of J_{ij} with respect to u_{ik} and v_{jk} .

$$\frac{\partial J_{ij}}{\partial u_{ik}} = 2 \left(x_j^{(i)} - \sum_k \exp(u_{ik}) \exp(v_{jk}) \right) \cdot (-\exp(u_{ik}) \exp(v_{jk}))$$

$$\frac{\partial J_{ij}}{\partial v_{jk}} = 2 \left(x_j^{(i)} - \sum_k \exp(u_{ik}) \exp(v_{jk}) \right) \cdot (-\exp(u_{ik}) \exp(v_{jk})) \quad (\text{same!})$$

- (c) Briefly describe how we could apply (e.g., give pseudocode for) stochastic gradient descent to learn u and v .

Init u, v to something (random usually)

Set stopping condition (eg # iterations) & step size α .

while (! done)

 for each i in random order

 for each j in random order

 (or do all j , depending on interpretation)

$$\hat{u} = u_j, \hat{v} = v_j$$

 for $k = 1..K$

$$u_{ijk} = u_{ik} - \alpha \frac{\partial J}{\partial u_{ik}}$$

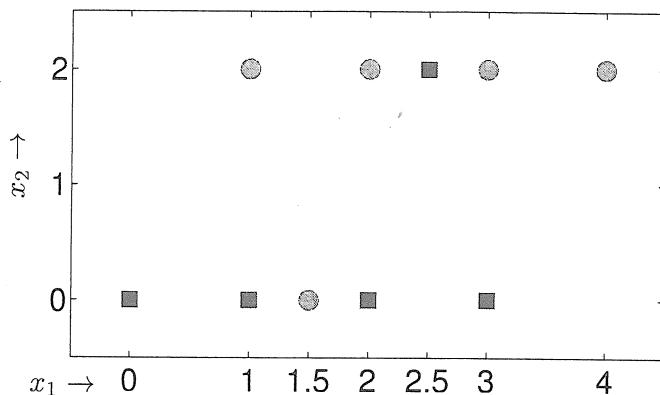
$$v_{jkn} = v_{jn} - \alpha \frac{\partial J}{\partial v_{jn}}$$

(use \hat{u}, \hat{v} in ∂J
to avoid overwriting)

Problem 5: (8 points) Cross-validation

Suppose that we learn a classifier on the following binary classification data. There are two real-valued features, x_1 and x_2 , and a binary class $y \in \{0, 1\}$.

x_1	x_2	y
0	0	0
1	0	0
1.5	0	1
2	0	0
3	0	0
1	2	1
2	2	1
2.5	2	0
3	2	1
4	2	1



We decide to learn a decision tree as described in class. As in class, when the decision tree splits on the real-valued features, it puts the split threshold halfway between the data points on either side of the highest-scoring split. For example, if we first split on x_2 , the algorithm would choose to split at $x_2 = 1$, which is halfway between the data at $x_2 = 0$ and $x_2 = 2$. In the case of ties, we prefer to predict class 0.

- (a) What is the training error rate of a decision *stump* (decision tree with max depth 1, or two leaf nodes) trained on these data?

$$\text{Split on } x_2 = 1$$

$$\Rightarrow 2/10 = 1/5 \text{ error}$$

- (b) What is the training error rate of a full decision tree (no maximum depth or pruning) trained on these data?

$$\text{Split until all correct}$$

$$\Rightarrow 0 \text{ error}$$

- (c) What is the leave-one-out cross-validation error rate of a decision *stump* (decision tree with max depth 1) trained on these data?

$$\text{All } x_2 \text{'s split on } x_2 = 1$$

$$\Rightarrow 2/10 = 1/5 \text{ error}$$

- (d) What is the leave-one-out cross-validation error rate of a full decision tree (no maximum depth) trained on these data?

$$\text{All } x_2 \text{'s split on } x_2 = 1, \text{ then at midpoints:}$$

$$\Rightarrow \text{wrong at } x_2 = 0, x_1 = 1, 1.5, 2$$

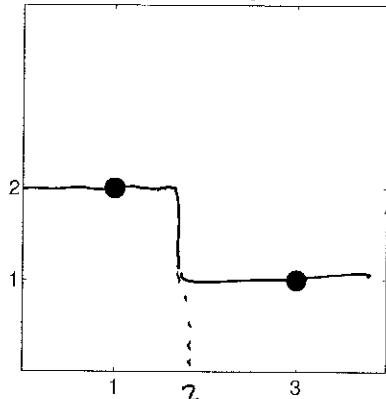
$$x_2 = 2, x_1 = 2, 2.5, 3$$

$$6 \Rightarrow 6/10 = 3/5 \text{ error.}$$

Problem 6: (8 points) Bagging

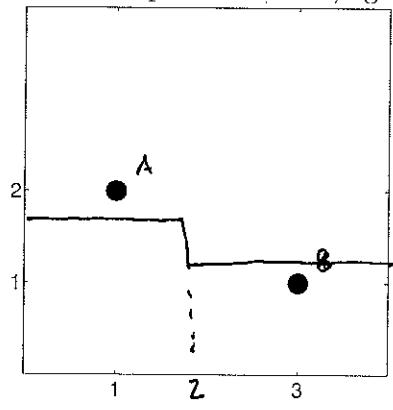
Consider a data set, consisting of two data points, plotted in each part.

- (a) Draw the regression function (predicted values for all x) using a nearest-neighbor regressor. Label any necessary values on your graph.

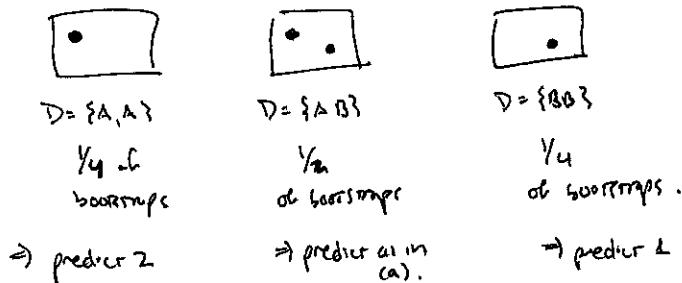


$$f(x) = \begin{cases} 2 & x < 2 \\ 1 & \text{ow.} \end{cases}$$

- (b) Suppose that we create a very large ensemble of *bagged* nearest-neighbor regressors, using data set draws of size $m = 2$ during the bootstrap sampling. Compute the regression function of the complete ensemble, again labeling any necessary values.



Bootstrap sampling draws data sets with replacement, which means we will get one of three possible sets, with probability:



The bagged ensemble averages the predictions of its members $\Rightarrow Y_4 \cdot 2 + Y_4 \cdot 1 + Y_2 \cdot \begin{cases} 2 & x < 2 \\ 1 & \text{ow.} \end{cases}$

$$= \begin{cases} 1 \frac{3}{4} & x < 2 \\ 1 \frac{1}{4} & \text{ow.} \end{cases}$$

- (c) How does the training error in model (b) compare to the training error in (a)? (Note: you don't need to have answered (b) to answer this part.)

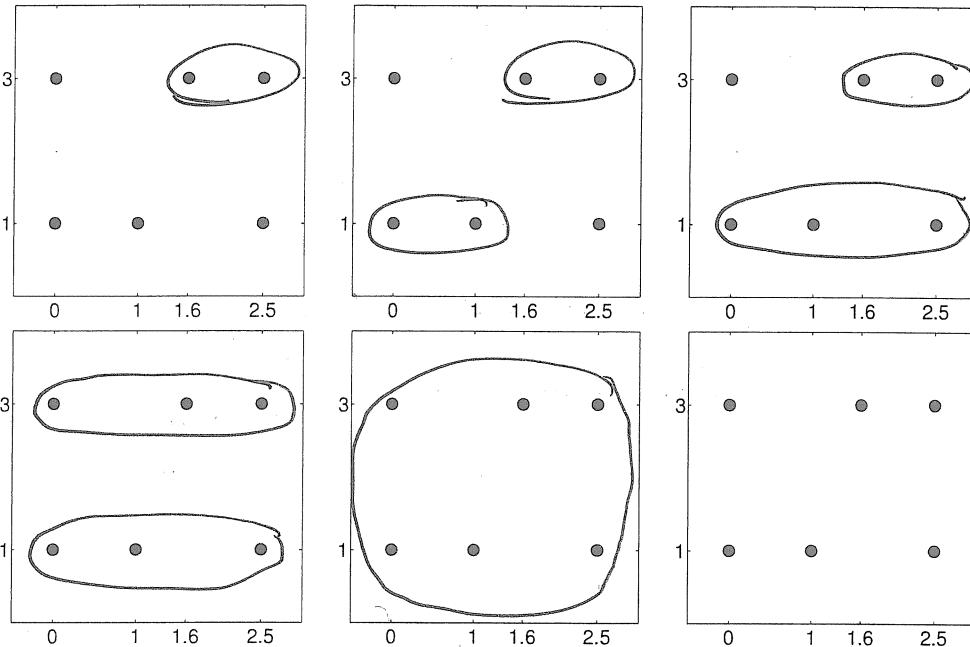
Training error has increased (MSE was 0, now $(\frac{1}{4})^2$)

- our bagged model is overfitting less than the original; it is unable to "memorize" the data points.

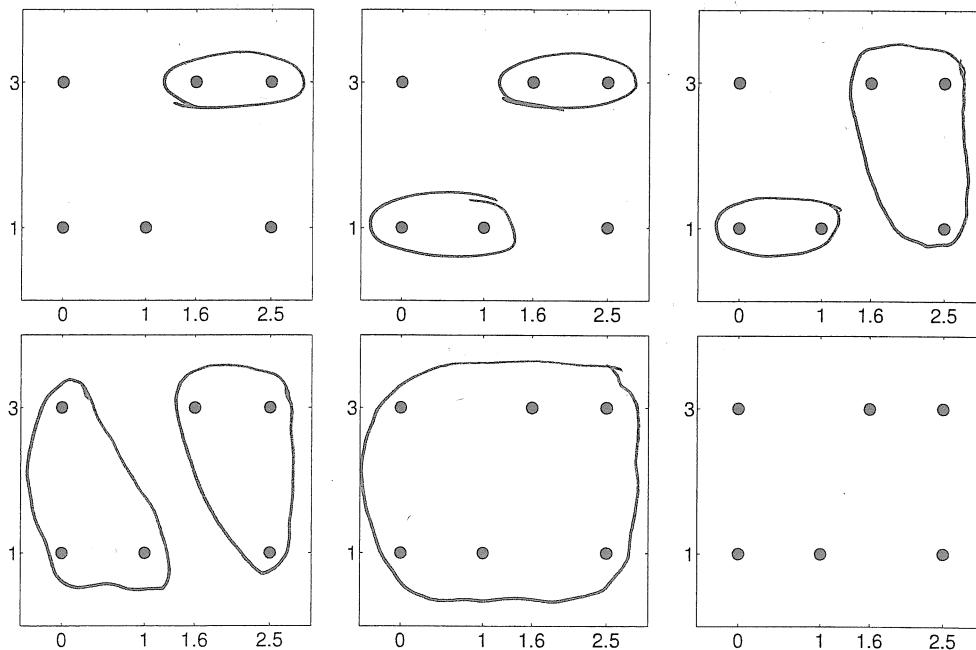
Problem 7: (8 points) Clustering

Consider the two-dimensional data points plotted in each panel. In this problem, we will cluster these data using two different algorithms, where each panel is used to show an iteration or step of the algorithm.

- (a) Execute the hierarchical agglomerative clustering (linkage) algorithm on these data points, using “single linkage” (minimum distance) for the cluster scores. Stop when converged, or after 6 steps, whichever is first. Show each step separately in a panel.

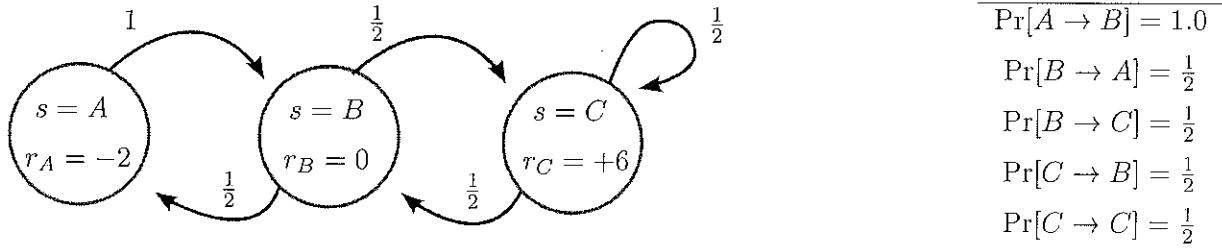


- (b) Now execute hierarchical agglomerative clustering on the data points, but use “complete linkage” (maximum distance) for the cluster scores. Stop when converged, or after 6 steps, whichever is first. Show each step separately in a panel.



Problem 8: (8 points) Markov models

Consider the Markov model shown here:



where the transition probabilities are shown next to each arc and at right, and the rewards r_s associated with each state s are shown inside the circles. Assume a future discounting factor of $\gamma = \frac{1}{2}$.

- (a) Compute $J^1(s)$, the expected discounted reward for state sequences of length 1 starting in each state s .

$$J^1 = \begin{array}{ccc} \underline{A} & \underline{B} & \underline{C} \\ -2 & 0 & 6 \end{array}$$

- (b) Compute $J^2(s)$, the expected discounted reward for state sequences of length 2 starting in each state s .

$$\begin{aligned} J^2 = & \begin{array}{ccc} \underline{A} & \underline{B} & \underline{C} \\ -2 + \frac{1}{2} \cdot 1 \cdot 0 & 0 + \frac{1}{2} \cdot \frac{1}{2} \cdot (-2) & 6 + \frac{1}{2} \cdot \frac{1}{2} \cdot 6 \\ = -2 & + \frac{1}{2} \cdot \frac{1}{2} \cdot 6 & + \frac{1}{2} \cdot \frac{1}{2} \cdot 6 \\ & = 1 & = 7\frac{1}{2}. \end{array} \end{aligned}$$

- (c) Compute $J^3(s)$, the expected discounted reward for state sequences of length 3 starting in each state s .

$$\begin{aligned} J^3 = & \begin{array}{ccc} \underline{A} & \underline{B} & \underline{C} \\ -2 + \frac{1}{2} \cdot 1 \cdot (-2) & 0 + \frac{1}{2} \cdot \frac{1}{2} \cdot (-2) & 6 + \frac{1}{2} \cdot \frac{1}{2} \cdot 1 \\ \uparrow & + \frac{1}{2} \cdot \frac{1}{2} \cdot 6 & + \frac{1}{2} \cdot \frac{1}{2} \cdot 6 \\ = -1\frac{1}{2} & = \frac{15}{8} - \frac{1}{2} & = 6 + \frac{1}{2} + \frac{15}{8} \\ & = \frac{11}{8} & = 8\frac{1}{8}. \end{array} \end{aligned}$$

CS178 Final Exam
Machine Learning & Data Mining: Winter 2016
Thursday March 17th, 2016

Your name: *Solutions*

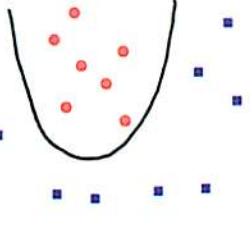
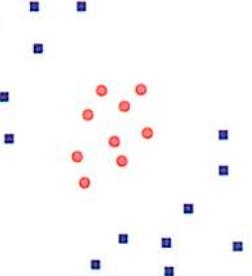
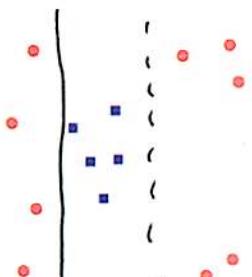
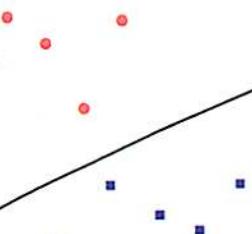
Your ID Number and UCINetID
(e.g., 12345678 / myname@uci.edu): *31415926 / pmreuter@uci.edu*.

Your seat (row and number): *leconm*

- Total time is 1 hour 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 to come over.
- You may use one sheet of your own, handwritten notes for reference, and a calculator.
- Turn in any scratch paper with your exam

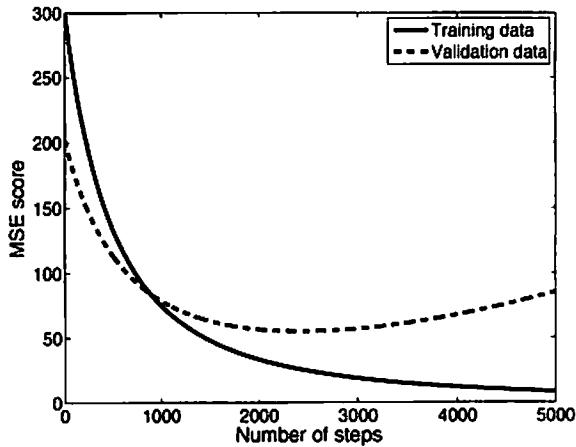
Problem 1: (8 points) Separability & Classifiers

For each of the following examples of training data and classifiers, state whether there exists a set of parameters that can separate the data and justify your answer briefly (~1 sentence).

	Linear classifier with quadratic features: <i>Yes - sketched boundary is</i> $x_2 = ax_1^2 + bx_1 + c$
	Depth-two decision tree: <i>No - no 2 level axis-aligned split.</i>
	Depth-two decision tree: <i>Yes, e.g.:</i> <pre> graph TD Root[x1 > a] --> L[x1 > b] Root --> R[x1 > b] L --> L1[+1] R --> R1[-1] R --> R2[+1] </pre>
	Linear perceptron classifier: <i>Yes; linear decision boundary.</i>

Problem 2: (9 points) Training & Test Error

Consider the following plot, which shows the training set error and the validation test set error for a neural network model as it is trained, i.e., the horizontal axis indicates the number of iterations of training (gradient steps). Note that the training error decreases monotonically, while the test error does not.



- (a) Explain what is happening and why; suggest a possible solution.

overfitting - the model is becoming too tuned to the training data, degrading validation performance.

Use early stopping, regularization, or some other form of complexity control.

Now suppose that we were to re-train the model with 10 times as much data, while keeping all other aspects (initialization, etc.) the same.

- (b) Would you expect the training curve to be different? If so, sketch how it might change.

No / Similar - maybe slightly higher, esp. at the far right.

- (c) Would you expect the validation (test) curve to be different? If so sketch how it might change.

Yes - flatter, probably no or less upturn at the far right



Problem 3: (12 points) Decision Trees

Consider the table of measured data given at right. (Note that some data points are repeated.) We will use a decision tree to predict the outcome y using the three features, x_1, \dots, x_3 . In the case of ties, we prefer to use the feature with the smaller index (x_1 over x_2 , etc.) and prefer to predict class 1 over class 0. You may find the following values useful (although you may also leave logs unexpanded):

$$\begin{aligned} \log_2(1) &= 0 & \log_2(2) &= 1 & \log_2(3) &= 1.59 & \log_2(4) &= 2 \\ \log_2(5) &= 2.32 & \log_2(6) &= 2.59 & \log_2(7) &= 2.81 & \log_2(8) &= 3 \end{aligned}$$

- (a) What is the entropy of y ? $p(y) = 4/7$

$$H = \frac{4}{7} \log \frac{7}{4} + \frac{3}{7} \log \frac{7}{3} \approx .985 \text{ bits}$$

x_1	x_2	x_3	y
0	0	1	1
0	1	0	1
1	1	1	1
1	1	1	1
0	0	0	0
0	0	0	0
1	1	0	0

- (b) Which variable would you split first? Justify your answer.

x_3 clearly has lower entropy after split (highest info gain)

$$\begin{array}{lll} x_0: 0 \geq 1100 & x_2: 0 \geq 1000 & x_3: 0 \geq 1000 \\ & 1 \geq 110 & 1 \geq 111 \end{array}$$

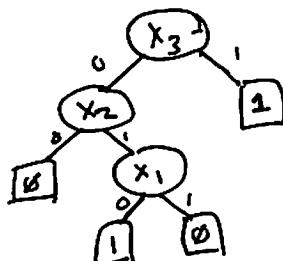
- (c) What is the information gain of the variable you selected in part

$$(b)? H(\frac{4}{7}) - \frac{4}{7} H(\frac{4}{4}) - \frac{3}{7} \cdot \frac{1}{2}.$$

$$\approx .985 - .463 - .18.$$

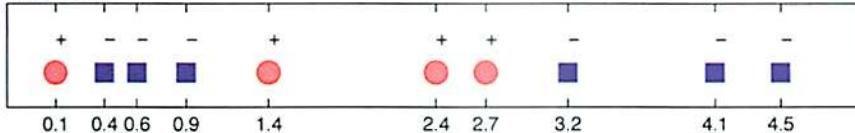
$$\approx .522 \text{ bits.}$$

- (d) Draw the rest of the decision tree learned on these data.



Problem 4: (12 points) Classification in 1D

We observe a collection of training data with one feature, “ x ” and a class label $y \in \{-, +\}$, shown here; class + is indicated by circles and - by squares, and also labeled with text for redundancy. Answer each of the following questions. Express error rates as the fraction of data points incorrectly classified.



- (a) What is the best training error rate we can achieve on these data from a linear classifier on the original input features ($f(x) = \text{sign}(ax + b)$)? Explain briefly (sketch + 1-2 sentences): how it is achieved.

3/10

$$\begin{array}{c} + \\ 0 | \square \square \square \end{array} \quad \begin{array}{c} \square \square \square \\ 0 \end{array} \quad \begin{array}{c} \square \square \square \\ 0 \end{array}$$

or

$$\begin{array}{c} (+) \quad (-) \\ \square \square \square \quad \square \square \square \end{array}$$

set the decision boundary, $ax+b=0$,
at $x=.3$ or $x=3$, for example.

- (b) What is the best training error we can achieve from a linear classifier with quadratic features, e.g., $f(x) = \text{sign}(ax^2 + bx + c)$? Explain briefly how it's achieved.

1/10

$$\begin{array}{c} ax^2+bx+c \\ 0 \quad \square \square \square \quad \square \square \square \end{array}$$

$$ax^2+bx+c = -(x-1)(x-3).$$

(for example)

- (c) What is the best training error we can achieve from a decision tree classifier? Explain briefly how it's achieved.

0/10

keep splitting until all the training data
are correct.

- (d) What is the best training error we can achieve from a two-layer neural network (multi-layer perceptron) with input features “ x ” and “1”? Explain briefly how it's achieved (e.g., # of hidden nodes & what they look like).

0/10

Use say 3 hidden nodes, activating at $x=.3$, $x=1$, and $x=3$.

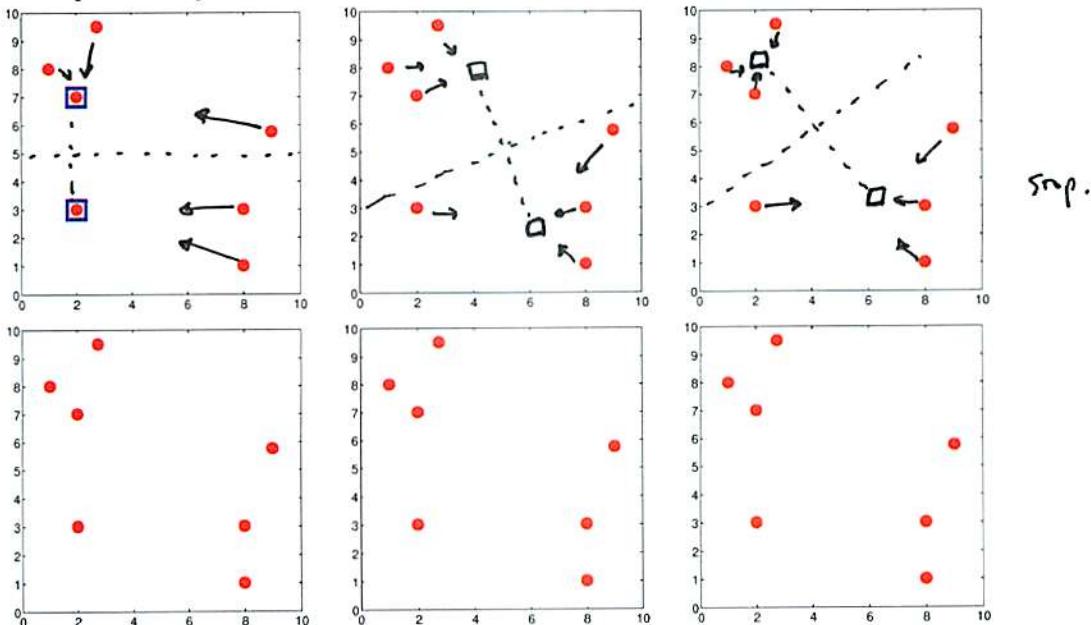
then, a linear combination of these values can separate the data.

Problem 5: (12 points) Clustering

Consider the two-dimensional data points plotted in each panel. In this problem, we will cluster these data using two different algorithms, where each panel is used to show an iteration or step of the algorithm.

k-means

(a) Starting from the two cluster centers indicated by squares, perform k-means clustering on the data points. In each panel, indicate (somehow) the data assignment, and in the next panel show the new cluster centers. Stop when converged, or after 6 steps (3 iterations), whichever is first. It may be helpful to recall from our nearest-neighbor classifier that the set of points nearer to A than B is separated by a line.



(b) Write down the cost function optimized by the k-means algorithm, explaining your notation.

$$J = \sum_i \|x^{(i)} - \mu_{z_i}\|_2^2$$

$$\|v\|_2^2 = \sum_j v_j^2 \text{ is Euclidean length squared}$$

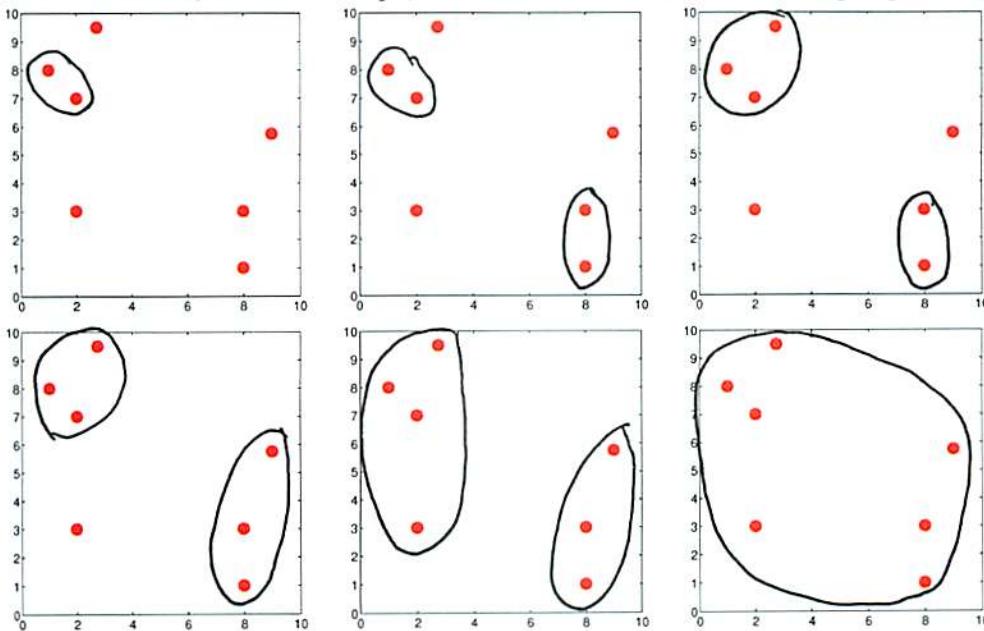
$x^{(i)}$ = data point i .

z_i = ^{cluster}assignment of $x^{(i)}$

μ_c = cluster center of cluster c .

Linkage

(a) Now execute the hierarchical agglomerative clustering (linkage) algorithm on these data points, using “complete linkage” (maximum distance) for the cluster scores. Stop when the algorithm would terminate, or after 6 steps, whichever is first. Show each step separately in a panel.



(b) What is the algorithmic (computational) complexity of the hierarchical clustering algorithm? Briefly justify your answer.

$$O(m^2 \log m)$$

Compute $O(m^2)$ cluster distances for the original m points.

$$\text{Manually sorted} \Rightarrow O(m^2 \log(m^2)) = O(m^2 \log m).$$

Each iteration i (m times):

merge closest cluster pair (constant time)

update $(m-i)$ cluster distances in sorted somehow $\Rightarrow O((m-i) \log m)$

$$\Rightarrow \text{total } O(m^2 \log m).$$

Problem 6: (9 points) Bayes Classifiers and Naïve Bayes

Consider the table of measured data given at right. We will use the three observed features x_1, x_2, x_3 to predict the class y . In the case of a tie, we will prefer to predict class $y = 0$.

- (a) Write down the probabilities necessary for a naïve Bayes classifier:

$$\begin{aligned} p(y=1) &= 4/7 & p(y=0) &= 1 - p(y=1) = 3/7 \\ p(x_1=1 | y=1) &= 3/4 & p(x_1=1 | y=0) &= 1/3 \\ p(x_2=1 | y=1) &= 1/2 & p(x_2=1 | y=0) &= 2/3 \\ p(x_3=1 | y=1) &= 1/2 & p(x_3=1 | y=0) &= 1/3 \end{aligned}$$

x_1	x_2	x_3	y
0	0	0	0
0	0	0	1
0	1	1	0
1	1	0	0
1	1	0	1
1	0	1	1
1	1	1	1

- (b) Using your naïve Bayes model, what value of y is predicted given observation $(x_1, x_2, x_3) = (000)$?

$$\begin{array}{lll} 3/7 \cdot 2/3 \cdot 1/3 \cdot 2/3 & \text{vs} & 4/7 \cdot 1/4 \cdot 1/2 \cdot 1/2 \\ 4/9 & \text{vs} & 1/4 \Rightarrow \text{predict } y=0. \end{array}$$

- (c) What is the class probability $p(y = 1 | x_1 = 0, x_2 = 1, x_3 = 1)$?

$$p(y=1 | 011) = \frac{4/7 \cdot 1/4 \cdot 1/2 \cdot 1/2}{4/7 \cdot 1/4 \cdot 1/2 + 3/7 \cdot 2/3 \cdot 1/3} = \frac{4 \cdot 1 \cdot 1 \cdot 1}{4 \cdot 2 + 3 \cdot 4 \cdot 1} = \frac{4}{9+16} = \frac{4}{25} = 16/25.$$

Problem 7: (9 points) VC Dimension

Consider the following classifiers $f(x)$, defined on data with two real-valued features $x = (x_1, x_2)$ and predicting a binary class $y \in \{-1, +1\}$. Answer the following questions about their VC dimension, by showing that it is at least as large as the value you give.

- (a) What is the VC dimension of the linear classifier,

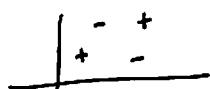
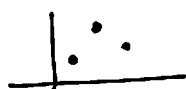
$$f(x) = \text{sign}(a + bx_1 + cx_2),$$

with parameters (a, b, c) ?

3 - perceptron in d features has vc dim d+1

Recall:

can shatter: but cannot separate



- (b) What is the VC dimension of a decision stump,

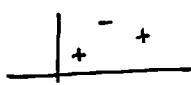
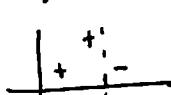
$$f(x) = a \text{ if } x_i < t \text{ else } b,$$

with parameters (a, b, t, i) ?

3 -



Not required, but again, cannot shatter

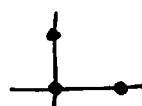


No single axis split.

etc.

- (c) Suppose that the data points were forced to have *binary valued* features, e.g., $x_i \in \{0, 1\}$, rather than real values. Would this change your answer for (a) or (b)?

For (a) - No. Just place the 3 points at 2 still shatters.



For (b) - Yes. For point placement like ↑,

we cannot separate with a single, axis aligned split.

CS273A Final Exam
Introduction to Machine Learning: Winter 2019
Thursday, March 21st, 2019

Your name:

SOLUTIONS

Row/Seat Number:

lemon

Your ID # (e.g., 123456789)

1314159265

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

panteater@uci

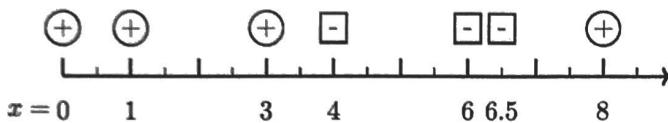
- Please put your name and ID on every page.
- Total time is 1 hour 50 minutes. READ THE EXAM FIRST and organize your time; don't spend too long on any one problem.
- Please write clearly and show all your work.
- If you need clarification on a problem, please raise your hand and wait for the instructor or TA to come over.
- You may use one sheet containing handwritten notes for reference, and a basic calculator; no other electronics allowed.

Problems

1	K-Nearest Neighbors, (8 points.)	3
2	Linear Regression, (8 points.)	5
3	Multiple Choice, (14 points.)	7
4	Decision Trees, (10 points.)	9
5	Bayes Classifiers, (10 points.)	11
6	Clustering, (10 points.)	13
7	Markov Processes, (9 points.)	15
8	VC-Dimension, (10 points.)	17
Total, (79 points.)		

K-Nearest Neighbors, (8 points.)

Consider the following dataset with five points shown below, for a binary classification task ($y = +, -$) with a scalar feature x . In case of ties, prefer the negative class. Put final answers in the box.



- (1) Compute the **training** error of a 1-nearest neighbor classifier. (2 points.)

0

- (2) Compute the **leave-one-out** cross-validation error of 1-nearest neighbor classifier. (2 points.)

✓ ✓ ✗ ✗ ✓ ✓ ✗

3/7

- (3) Compute the **training** error of 2-nearest neighbor classifier. (2 points.)

✓ ✓ ✗ ✓ ✓ ✓ ✗

2/7

- (4) Compute the **leave-one-out** cross-validation error of 2-nearest neighbor classifier. (2 points.)

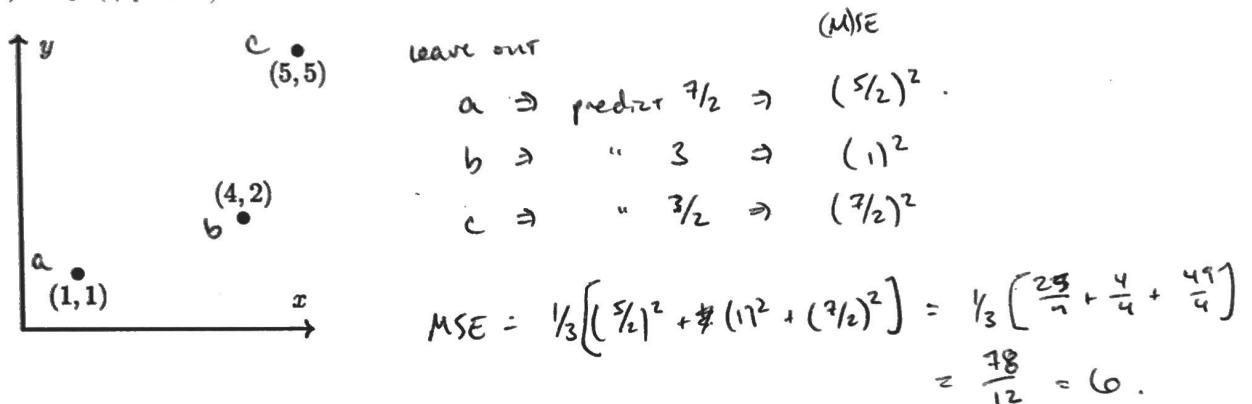
✓ ✓ ✗ ✓ ✓ ✓ ✗

2/7

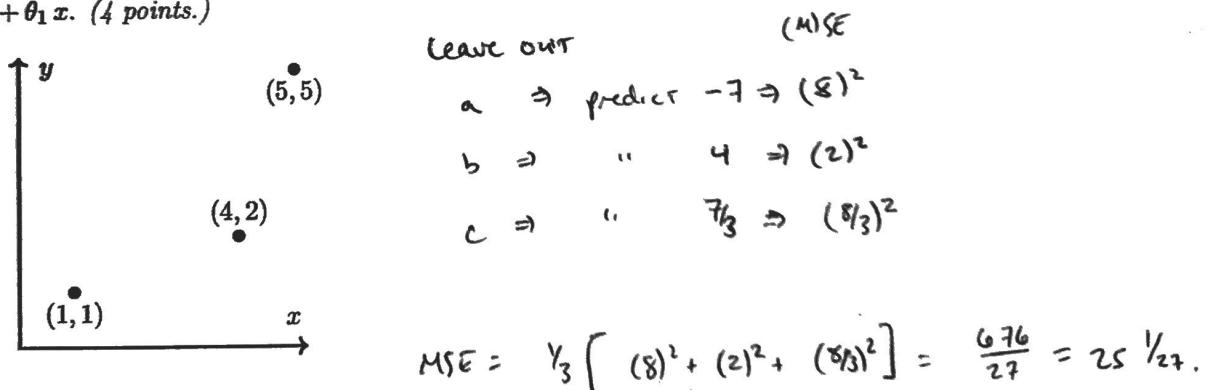
Linear Regression, (8 points.)

Consider the following data points, copied in each part. We wish to perform linear regression to minimize the mean squared error of our predictions.

- (a) Compute the leave-one-out cross-validation error of a zero-order (constant) predictor, $\hat{y}(x) = \theta_0$. (4 points.)



- (b) Compute the leave-one-out cross-validation error of a first-order (linear) predictor, $\hat{y}(x) = \theta_0 + \theta_1 x$. (4 points.)



Multiple Choice, (14 points.)

For each of the following statements, choose whether the statement is true or false.

Statement	True	False
Universal function approximators		
With enough hidden nodes and layers, a neural network can approximate any function.	<input checked="" type="checkbox"/>	<input type="checkbox"/>
With enough hidden nodes and a <i>single</i> layer, a neural network can approximate any function.	<input checked="" type="checkbox"/>	<input type="checkbox"/>
With enough data, a Gaussian Bayes classifier can approximate any function.	<input type="checkbox"/>	<input checked="" type="checkbox"/>
With enough depth, a decision tree can approximate any function.	<input checked="" type="checkbox"/>	<input type="checkbox"/>
Bagging can be used to make any simpler classifier arbitrarily complex.	<input type="checkbox"/>	<input checked="" type="checkbox"/>
Optimization		
Using backpropagation to train a neural network will avoid getting stuck in local optima.	<input type="checkbox"/>	<input checked="" type="checkbox"/>
Stochastic gradient descent is often preferred over batch when the number of data points m is very large.	<input checked="" type="checkbox"/>	<input type="checkbox"/>
Linear regression can be solved using either matrix algebra or gradient descent.	<input checked="" type="checkbox"/>	<input type="checkbox"/>
The K-Means algorithm is guaranteed to converge in finite number of steps.	<input checked="" type="checkbox"/>	<input type="checkbox"/>
The number of steps for agglomerative clustering to complete does not depend on the linkage function.	<input checked="" type="checkbox"/>	<input type="checkbox"/>

Order the following four learners in order of *increasing* complexity (likelihood of overfitting):

(a) Perceptron, $\hat{y} = T(\theta \cdot x)$

(b) Ensemble (mixture) of three perceptrons,

$$\hat{y} = \alpha_1 T(\theta_1 \cdot x) + \alpha_2 T(\theta_2 \cdot x) + \alpha_3 T(\theta_3 \cdot x)$$

trained by jointly optimizing $\{\alpha_1, \alpha_2, \alpha_3, \theta_1, \theta_2, \theta_3\}$

(c) Ensemble of three perceptrons trained by AdaBoost

(d) Ensemble of three perceptrons trained by Bootstrap Aggregation

(Simplest) **a** **b** **c** (Most Complex)

Decision Trees, (10 points.)

Consider the table of measured data given at right. We will use a decision tree to predict the outcome y using the three features, x_1, \dots, x_3 . In the case of ties, we prefer to use the feature with the smaller index (x_1 over x_2 , etc.) and prefer to predict class 1 over class 0. You may find the following values useful (although you may also leave logs unexpanded):

$$\begin{aligned} \log_2(1) &= 0 & \log_2(2) &= 1 & \log_2(3) &= 1.59 & \log_2(4) &= 2 \\ \log_2(5) &= 2.32 & \log_2(6) &= 2.59 & \log_2(7) &= 2.81 & \log_2(8) &= 3 \end{aligned}$$

x_1	x_2	x_3	y
0	0	0	1
1	0	1	1
1	0	1	1
1	1	1	0
0	1	0	0
1	0	0	0

- (1) What is the entropy of y ? (2 points.)

$$P(y=1) = 1/2 \Rightarrow H(y) = 1 \text{ bit.}$$

- (2) Which variable would you split first? Justify your answer. (2 points.)

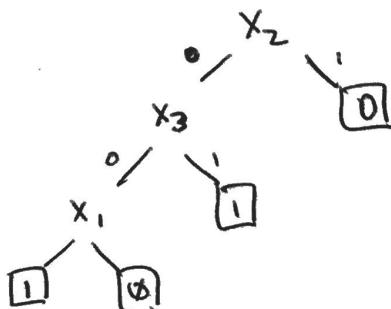
$x_1 = 0$	$y = 0$	$y = 1$	$x_2 = 0$	$y = 0$	$y = 1$	$x_3 = 0$	$y = 0$	$y = 1$
0	010	000	0	100	000	0	010	000
1	100	101	1	101	101	1	100	101

- (3) What is the information gain of the variable you selected in part (2)? (3 points.)

$$\begin{aligned} IG_1 &= 1 - \left[\frac{1}{3} H(0) + \frac{2}{3} H(\frac{1}{2}) \right] \\ &= 1 - \frac{2}{3} \left[\frac{3/4 \log \frac{4}{3}}{0.41} + \frac{1/4 \log 4}{2} \right] \end{aligned}$$

$$\approx 1 - .64 \approx .36 \text{ bits.}$$

- (4) Draw the rest of the decision tree learned on these data. (3 points.)



Name: _____

ID#: _____

Bayes Classifiers, (10 points.)

Consider the table of measured data given at right. We will use the three observed features x_1, x_2, x_3 to predict the class y . In the case of a tie, we will prefer to predict class $y = 0$.

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

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

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

$$p(x_1 = 1|y=0) : \frac{3}{4}$$

$$p(x_1 = 1|y=1) : \frac{2}{4} \frac{1}{2}$$

$$p(x_1 = 0|y=0) : \frac{1}{4}$$

$$p(x_1 = 0|y=1) : \frac{3}{4}$$

$$p(x_2 = 1|y=0) : \frac{1}{4}$$

$$p(x_2 = 1|y=1) : \frac{3}{4}$$

$$p(x_2 = 0|y=0) : \frac{3}{4}$$

$$p(x_2 = 0|y=1) : \frac{1}{4}$$

$$p(x_3 = 1|y=0) : \frac{1}{2}$$

$$p(x_3 = 1|y=1) : \frac{1}{2}$$

$$p(x_3 = 0|y=0) : \frac{1}{2}$$

$$p(x_3 = 0|y=1) : \frac{1}{2}$$

- (2) Using your naïve Bayes model, compute: (3 points.)

$$p(y=1|x_1=0, x_2=1, x_3=1) : \frac{1}{10} \quad p(y=0|x_1=0, x_2=1, x_3=1) : \frac{9}{10}$$

$$p(y=1|x_1=0, x_2=1, x_3=1) = \frac{\frac{1}{2} \cdot \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{2}}{\frac{1}{2} \cdot \frac{3}{4} \cdot \frac{1}{4} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{3}{4} \cdot \frac{3}{4} \cdot \frac{1}{2}} = \frac{1}{1+9} = \frac{1}{10}$$

- (3) Compute the probabilities $p(y=1|x_1=0, x_2=1, x_3=1)$ and $p(y=0|x_1=0, x_2=1, x_3=1)$ for a joint Bayes model trained on the same data. (3 points.)

$$p(y=1|x_1=0, x_2=1, x_3=1) = 1$$

by inspection - two examples w/ $y=1$,

$$p(y=0|x_1=0, x_2=1, x_3=1) = 0$$

no examples w/ $y=0$.

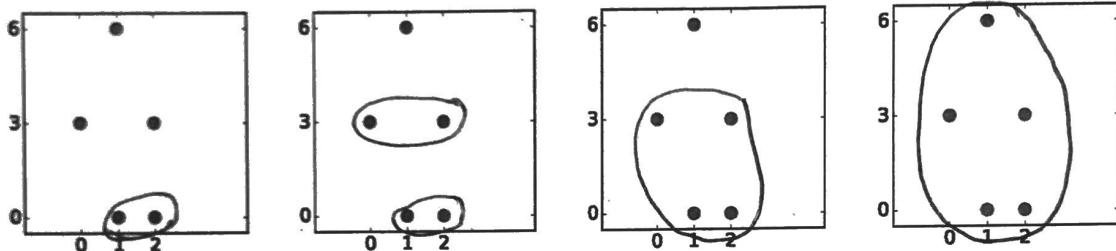
x_1	x_2	x_3	y
0	0	1	0
1	1	1	0
1	0	0	0
1	0	0	0
0	0	0	1
0	1	1	1
0	1	1	1
1	1	0	1

Clustering, (10 points.)

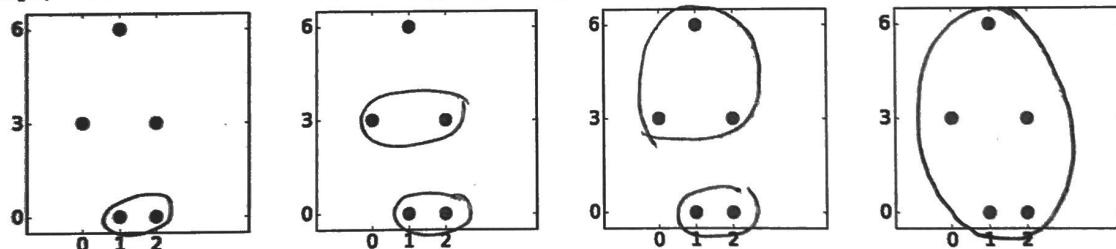
Consider the two-dimensional data points plotted in each panel. In this problem, we will cluster these data.

Linkage

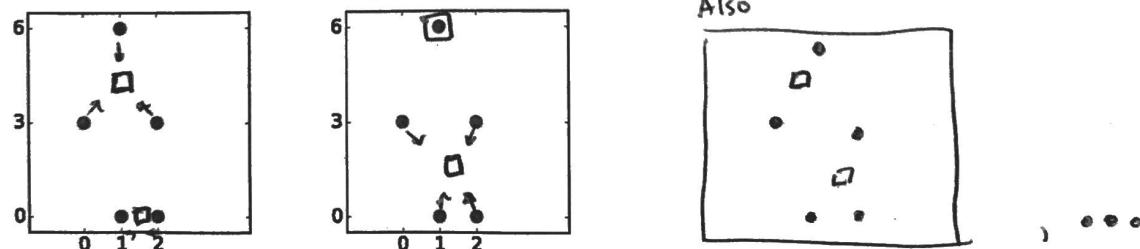
- (a) Execute the hierarchical agglomerative clustering (linkage) algorithm on these data points, using “single linkage” (minimum distance) for the cluster scores. Stop when the algorithm would terminate, or after 4 steps, whichever is first. Show each step separately in a panel.



- (b) Now repeat your agglomerative clustering algorithm, this time using “complete linkage” (maximum distance) for the cluster scores. Stop when the algorithm would terminate, or after 4 steps, whichever is first. Show each step separately in a panel.

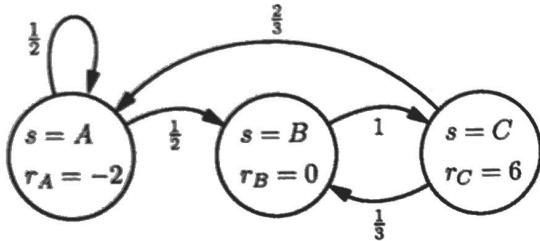
**k-means**

In the panels below, sketch two different clustering solutions that could be found by (would be converged solutions of) k -means, with $k = 2$. Show the final (converged) grouping of data (assignments) and final cluster locations. (You do not have to show the process, only the final clustering.) This illustrates the existence of local optima in k -means.



Markov Processes, (9 points.)

Consider the Markov model shown here:



$$\begin{aligned}\Pr[A \rightarrow A] &= 0.5 \\ \Pr[A \rightarrow B] &= 0.5 \\ \Pr[B \rightarrow C] &= 1.0 \\ \Pr[C \rightarrow B] &= 0.33 \\ \Pr[C \rightarrow A] &= 0.66\end{aligned}$$

where the transition probabilities are shown next to each arc and at right, and the rewards r_s associated with each state s are shown inside the circles. Assume a future discounting factor of $\gamma = \frac{1}{2}$.

- (1) Compute $J^1(s)$, the expected discounted sum of rewards for state sequences of length 1 (e.g., $[A]$) starting in each state s . (3 points.)

$$\begin{array}{ccc} \overline{J_1(s)} & & \\ \hline A & -2 & \\ \hline B & 0 & \\ \hline C & 6 & \end{array}$$

- (2) Compute $J^2(s)$, the expected discounted sum of rewards for state sequences of length 2 (e.g., $[C \rightarrow B]$) starting in each state s . (3 points.)

$$\begin{array}{ccc} \overline{J_2(s)} & & \\ \hline A & -2 + \frac{1}{2} J_1(A) + \frac{1}{2} \cdot \frac{1}{2} J_1(B) & \\ & = -2 \frac{1}{2} & \\ B & 0 + \frac{1}{2} \cdot 1 \cdot J_1(C) & \\ & = 0 & \\ C & 6 + \frac{1}{2} \cdot \frac{2}{3} \cdot J_1(A) + \frac{1}{2} \cdot \frac{1}{3} J_1(B) & \\ & = 5 \frac{1}{3} & \end{array}$$

- (3) Compute $J^3(s)$, the expected discounted sum of rewards for state sequences of length 3 (e.g., $[B \rightarrow A \rightarrow C]$) starting in each state s . (3 points.)

$$\begin{array}{ccc} \overline{J^3(s)} & & \\ \hline A & -2 + \frac{1}{2} \cdot \frac{1}{2} J_2(A) + \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} J_2(B) & \\ & = -2 \frac{1}{8} & \\ B & 0 + \frac{1}{2} \cdot 1 \cdot J_2(C) & \\ & = 2 \frac{2}{3} & \\ C & 6 + \frac{1}{2} \cdot \frac{2}{3} \cdot J_2(A) + \frac{1}{2} \cdot \frac{1}{3} J_2(B) & \\ & = 5 \frac{2}{3} & \end{array}$$

VC-Dimension, (10 points.)

Consider a family of classifiers on two-dimensional data $x = (x_1, x_2)$ that use the angle of the vector x at the origin. We classify a point as positive ($\hat{y} = +1$) if it lies between two angles θ_a and θ_b (parameters of the classifier), moving counterclockwise.

More precisely, we have

$$\hat{y}(x; \theta_a, \theta_b) = \begin{cases} +1 & \theta_a \leq \arctan\left(\frac{x_2}{x_1}\right) + 2k\pi \leq \theta_b \\ -1 & \text{otherwise} \end{cases} \quad \text{for some } k \in \mathbb{Z}$$

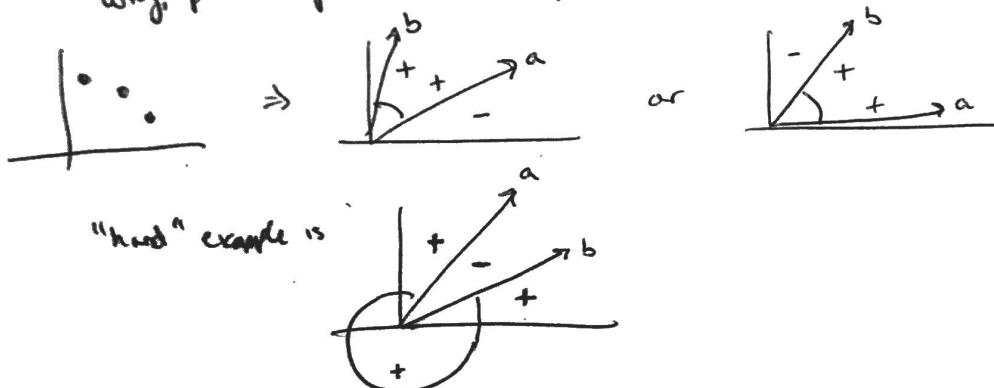
where θ_a, θ_b are parameters to be learned from the data.

- (1) What is the VC dimension H of this classifier? (2 points.)

3

- (2) Demonstrate that H is at least the value you gave in Part (a). (4 points.)

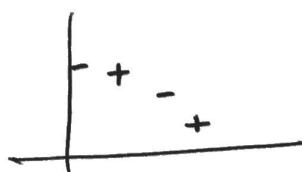
wlog, put the point in the 1st quadrant



- (3) Show by counterexample that H is no larger than the value you gave in Part (a). (4 points.)

An impossible pattern is

(in any position around the circle).



We can only transition once & back; "+" region must form a contiguous slice of the circle.

CS273A Final Exam
Introduction to Machine Learning: Winter 2020
Tuesday March 17th, 2020

Your name:

SOLUTIONS

Your ID #(e.g., 123456789)

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

- Due to the ongoing health emergency, this exam is take home and may be submitted in person or on Canvas (scanned).
- Total time is 2 hours 15 minutes for either in-person delivery (to the classroom) or submission to Canvas.
- READ THE EXAM FIRST and organize your time; don't spend too long on any one problem.
- Please put your name and ID on every page.
- Please write clearly and show all your work.
- If you need clarification on a problem, please post privately on Piazza and we will try to answer it.

Problems

1	Bayes Classifiers, (12 points.)	3
2	Cross-Validation, (9 points.)	5
3	Decision Trees, (10 points.)	7
4	Dimensionality Reduction, (10 points.)	9
5	Hierarchical Clustering, (11 points.)	11
6	K-Means Clustering, (10 points.)	13
7	Markov Processes, (12 points.)	15

Total, (74 points.)

Name: _____ ID#: _____

Bayes Classifiers, (12 points.)

In this problem you will use Bayes Rule: $p(y|x) = p(x|y)p(y)/p(x)$ to perform classification. Suppose we observe some training data with two binary features x_1, x_2 and a binary class y . After learning the model, you are also given some validation data.

Table 1: Training Data

x_1	x_2	y
0	0	0
0	1	0
0	1	1
0	1	1
1	0	1
1	0	1
1	1	0
1	1	0

Table 2: Validation Data

x_1	x_2	y
0	0	1
0	1	0
1	0	0
1	1	0

In the case of any ties, we will prefer to predict class 0.

- (1) Give the predictions of a joint Bayes classifier on the validation data. What is the validation error rate? (Put final answers in boxes.) (4 points.)

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

Error Rate:

$\boxed{3/4}$

- (2) Give the required probabilities to define a naïve Bayes classifier. (4 points.)

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

$$p(x_1=1 | y=0) = \frac{1}{2} \quad p(x_1=1 | y=1) = \frac{1}{2}$$

$$p(x_2=1 | y=0) = \frac{3}{4} \quad p(x_2=1 | y=1) = \frac{1}{2}$$

- (3) Give the predictions of a naïve Bayes classifier on the validation data. What is the validation error rate? (Put final answers in boxes.) (4 points.)

x_1	x_2	y
0	0	1
0	1	0
1	0	1
1	1	0

Error Rate:

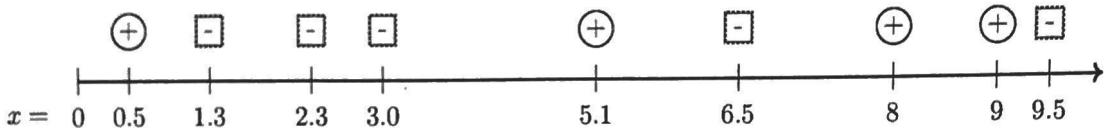
$\boxed{\frac{1}{4}}$

Name: _____

ID#: _____

Cross-Validation, (9 points.)

Consider the following dataset with *nine* points shown below, for a binary classification task ($y = +, -$) with a scalar feature x . In case of ties, prefer the negative class. Put final answers in the box.



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

$\times \times \quad \checkmark \checkmark \quad \times \times \quad \checkmark \times \times$

$\boxed{2/3}$

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

$\times \checkmark \checkmark \checkmark \quad \times \quad \times \times \times \times$

$\boxed{2/3}$

- (3) Compute the **leave-one-out** cross-validation error rate of an 8-nearest neighbor classifier. (3 points.)

All \oplus wrong

All \ominus right

$\boxed{4/9}$

Name: _____ ID#: _____

Decision Trees, (10 points.)

Consider the table of measured data given at right. We will use a decision tree to predict the outcome y using the three features, x_1, \dots, x_3 . In the case of ties, we prefer to use the feature with the smaller index (x_1 over x_2 , etc.) and prefer to predict class 1 over class 0. You may find the following values useful (although you may also leave logs unexpanded):

$$\begin{aligned} \log_2(1) &= 0 & \log_2(2) &= 1 & \log_2(3) &= 1.59 & \log_2(4) &= 2 \\ \log_2(5) &= 2.32 & \log_2(6) &= 2.59 & \log_2(7) &= 2.81 & \log_2(8) &= 3 \end{aligned}$$

- (1) What is the entropy of y ? (2 points.)

$$H(3/8) = -\frac{3}{8} \log \frac{3}{8} - \frac{5}{8} \log \frac{5}{8} \approx .95 \text{ bits}$$

x_1	x_2	x_3	y
0	0	0	0
1	1	0	1
0	1	0	0
0	1	0	0
0	1	0	0
1	0	0	0
1	1	0	1
1	0	1	1

- (2) Which variable would you split first? Justify your answer. (2 points.)

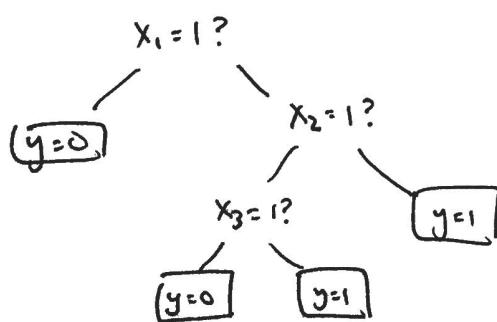
x_1	x_2	x_3
=0 0000	001	000011
=1 0111	00011	1

$\Rightarrow \boxed{x_1}$

- (3) What is the information gain of the variable you selected in part (2)? (3 points.)

$$IG = (.95) - \frac{4}{8} H(\frac{1}{2}) - \frac{4}{8} H(\frac{3}{4}) \approx .5488$$

- (4) Draw the rest of the decision tree learned on these data. (3 points.)



$IG_{x_1=1}:$

110	1
110	1
100	0
101	1

$\Rightarrow \cancel{x_1 = 1}$ $x_2?$ $x_3?$

=0	01	110
=1	11	1

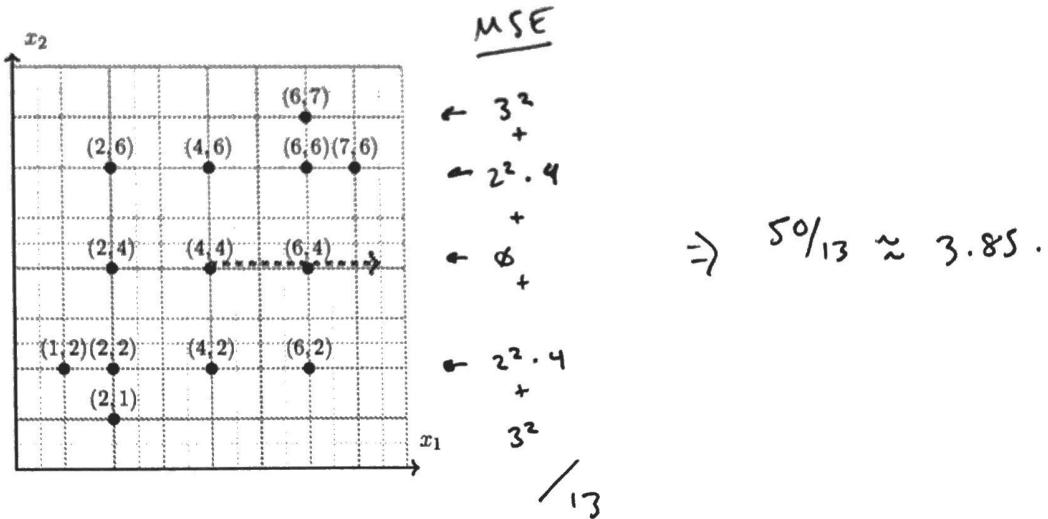
$\Rightarrow x_2 \text{ better.}$

Name: _____

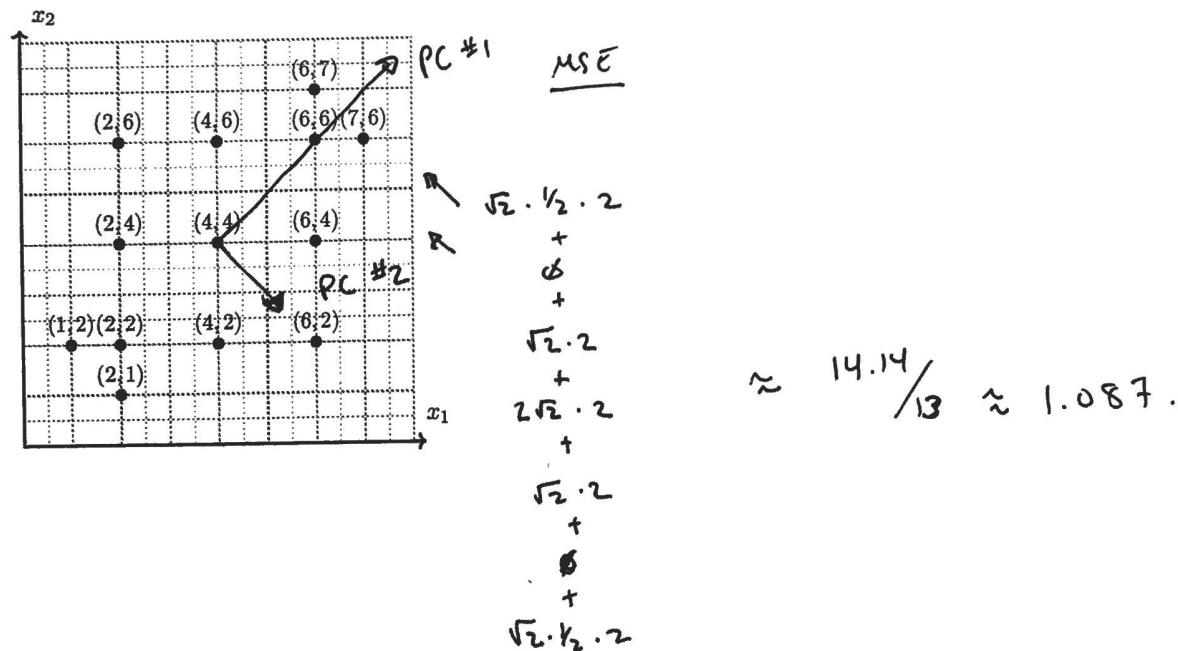
ID#: _____

Dimensionality Reduction, (10 points.)

- (1) For the following points in two dimensions, consider performing linear dimensionality reduction along the given vector (dashed line). What is the reconstruction error, in MSE, when using this vector? (4 points.)



- (2) On the figure below, draw the directions of the first two principal components. (2 points.)
- (3) What is the reconstruction error (MSE) of these points when only the first principal component is used to reconstruct each point? (4 points.)

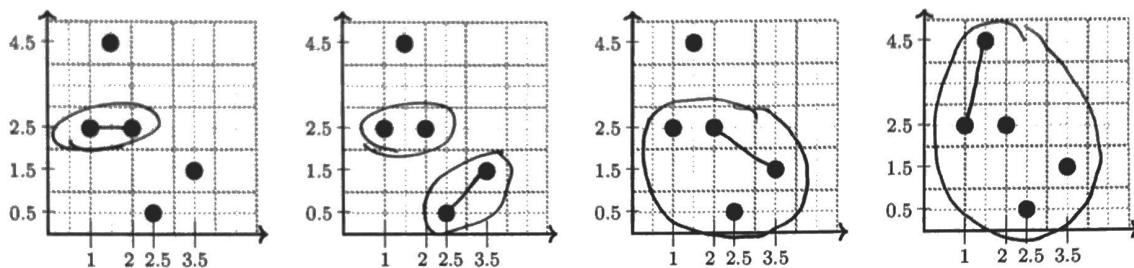


Hierarchical Clustering, (11 points.)

Consider the two-dimensional data points plotted in each panel. In this problem, we will cluster these data.

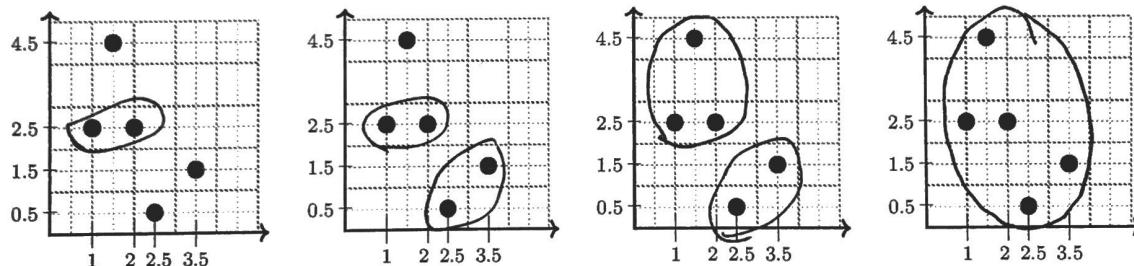
Linkage

- (a) Execute the hierarchical agglomerative clustering (linkage) algorithm on these data points, using "single linkage" (minimum distance) for the cluster scores. Stop when the algorithm would terminate, or after 4 steps, whichever is first. Show each step separately in a panel. (4 points.)



line shows closest pair.

- (b) Now repeat your agglomerative clustering algorithm, this time using "complete linkage" (maximum distance) for the cluster scores. Stop when the algorithm would terminate, or after 4 steps, whichever is first. Show each step separately in a panel. (4 points.)



- (c) What is the (big-O) computational complexity of the hierarchical agglomerative clustering algorithm? Justify your answer briefly in 1-2 sentences. (3 points.)

$$\mathcal{O}(n^2) - \text{compute all distances}$$

$$\mathcal{O}(n^2 \log n^2) - \text{sort}$$

M iterations:

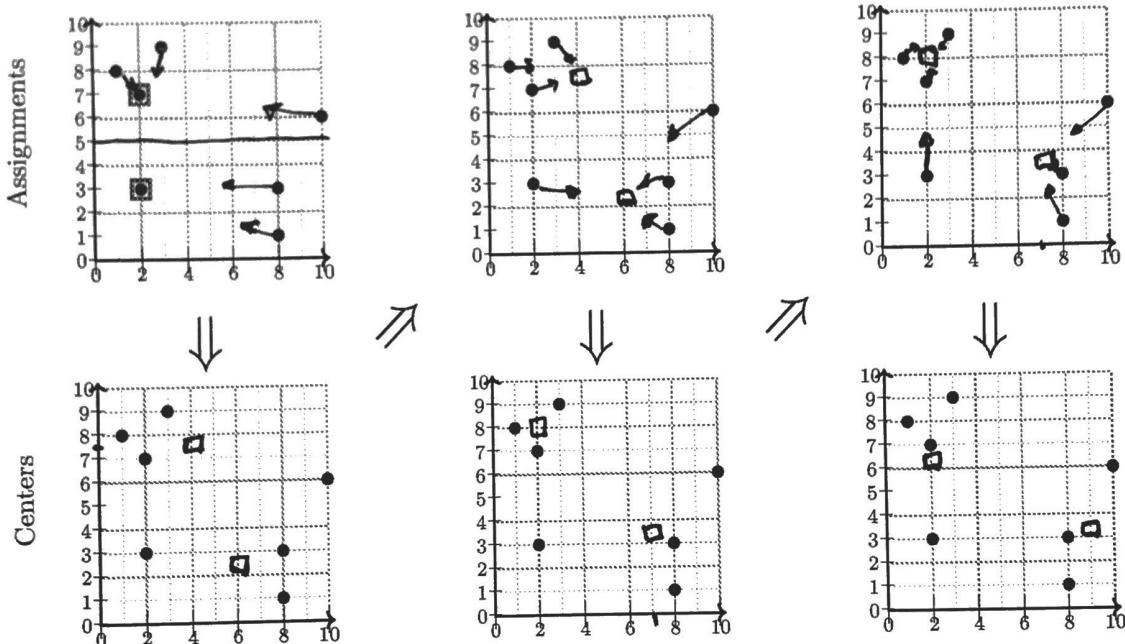
$$\mathcal{O}(m) \text{ distance recomputations} + \mathcal{O}(m \log m^2) \text{ re-insertions}$$

$$\Rightarrow \mathcal{O}(m^2 + m^2 \log m^2 + m^2 + m^2 \log m^2) = \mathcal{O}(m^2 \log m^2) = \mathcal{O}(m^2 \log m).$$

K-Means Clustering, (10 points.)

Consider the 2-D data points plotted in each panel. In this problem, we will cluster these data using the k -means algorithm, where each panel is used to show a single step of the algorithm.

- (a) Starting from the two cluster centers indicated by squares, perform k -means clustering on the data. In the top panels, indicate the assignment of the data, and then in the panel below show the new cluster centers, so each pair of panels shows an iteration of k -means. Stop when converged, or after 6 steps (3 iterations), whichever is first. It may be helpful to recall from our nearest neighbor classifier that the set of points nearer to A than B is separated by a line. (6 points.)



- (b) Write down the cost function optimized by the k -means algorithm, in terms of the data locations $x^{(i)}$, cluster centers μ_c , and cluster assignments $z^{(i)}$. (2 points.)

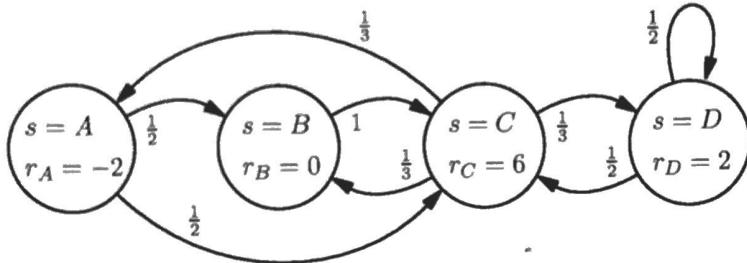
$$J(\mu, z) = \sum_i \|x^{(i)} - \mu_{z^{(i)}}\|^2 = \sum_i \sum_j (x_j^{(i)} - \mu_{z^{(i)} j})^2$$

- (c) What is the (big-O) computational complexity of each iteration of k -means (naïve computation), in terms of the data size m and number of clusters k ? (2 points.)

$O(m \cdot k)$ distance computations to assign
 $O(m)$ to compute new means $\Rightarrow O(mk)$ per iteration.

Markov Processes, (12 points.)

Consider the Markov reward process model shown here:



$$\begin{aligned}\Pr[A \rightarrow B] &= 0.5 \\ \Pr[A \rightarrow C] &= 0.5 \\ \Pr[B \rightarrow C] &= 1.0 \\ \Pr[C \rightarrow A] &= 0.33 \\ \Pr[C \rightarrow B] &= 0.33 \\ \Pr[C \rightarrow D] &= 0.33 \\ \Pr[D \rightarrow C] &= 0.5 \\ \Pr[D \rightarrow D] &= 0.5\end{aligned}$$

where the transition probabilities are shown next to each arc and at right, and the rewards r_s associated with each state s are shown inside the circles. We will use dynamic programming to (start) computing the expected discounted sum of rewards. Assume a future discounting factor of $\gamma = \frac{1}{2}$.

- (1) Compute $J^1(s)$, the expected discounted sum of rewards for state sequences of length 1 (e.g., $[A]$) starting in each state s . (4 points.)

$$J^1(A) =$$

-2

$$J^1(B) =$$

0

$$J^1(C) =$$

6

$$J^1(D) =$$

2

- (2) Compute $J^2(s)$, the expected discounted sum of rewards for state sequences of length 2 (e.g., $[C \rightarrow B]$) starting in each state s . (4 points.)

$$J^2(A) =$$

$-\frac{1}{2}$

$$J^2(B) =$$

3

$$J^2(C) =$$

6

$$J^2(D) =$$

4

$$\begin{aligned}-2 + \gamma_2 \cdot \gamma_2 \cdot 0 \\ + \gamma_2 \cdot \gamma_2 \cdot 6\end{aligned}$$

$$0 + \gamma_2 \cdot 1 \cdot 0$$

$$\begin{aligned}6 + \gamma_2 \cdot \gamma_3 \cdot 0 \\ + \gamma_2 \cdot \gamma_3 \cdot (-2) \\ + \gamma_2 \cdot \gamma_3 \cdot 2\end{aligned}$$

$$\begin{aligned}2 + \gamma_2 \cdot \gamma_2 \cdot 2 \\ + \gamma_2 \cdot \gamma_2 \cdot 4\end{aligned}$$

- (3) Compute $J^3(s)$, the expected discounted sum of rewards for state sequences of length 3 (e.g., $[D \rightarrow C \rightarrow A]$) starting in each state s . (4 points.)

$$J^3(A) =$$

$\frac{1}{4}$

$$J^3(B) =$$

3

$$J^3(C) =$$

$7 \frac{1}{12}$

$$J^3(D) =$$

$4 \frac{1}{2}$

$$\begin{aligned}-2 + \gamma_2 \cdot \gamma_2 \cdot 3 \\ + \gamma_2 \cdot \gamma_2 \cdot 6\end{aligned}$$

$$0 + \gamma_2 \cdot 1 \cdot 6$$

$$\begin{aligned}6 + \gamma_2 \cdot \gamma_3 \cdot (-2) \\ + \gamma_2 \cdot \gamma_3 \cdot 3 \\ + \gamma_2 \cdot \gamma_3 \cdot 4\end{aligned}$$

$$\begin{aligned}2 + \gamma_2 \cdot \gamma_2 \cdot 4 \\ + \gamma_2 \cdot \gamma_2 \cdot 6\end{aligned}$$

CS 273a Final Exam
Intro to Machine Learning: Fall 2023
Monday December 11th, 2023

Your name:

Peter Anteater

Row/Seat Number:

Your ID #(e.g., 123456789)

314159265

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

panteater@uci.edu

- Please put your name and ID **on every page**.
- Total time is 1 hour 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**.
- Please write your final answer **in the provided box** where possible.
- You may use **one** sheet containing handwritten notes for reference, and a **basic** calculator; no other electronics allowed.
- If you need clarification on a problem, please raise your hand and wait for the instructor or TA to come over.

Problems

1	Decision Trees, (<i>10 points.</i>)	3
2	Neural Networks, (<i>9 points.</i>)	5
3	True/False, (<i>10 points.</i>)	7
4	K-Nearest Neighbors, (<i>12 points.</i>)	9
5	Hierarchical Clustering, (<i>12 points.</i>)	11
6	K-Means Clustering, (<i>9 points.</i>)	13
7	Markov Processes, (<i>12 points.</i>)	15

Total, (*74 points.*)

This page is intentionally blank.

While you may use it as scratch paper, it may not be scanned for grading, so please do not include final solutions.

Problem 1 Decision Trees, (10 points.)

Consider the table of measured data given at right. We will use a decision tree to predict the outcome y using the three features, x_1, \dots, x_3 . In the case of ties, we prefer to use the feature with the smaller index (x_1 over x_2 , etc.) and prefer to predict class 1 over class 0. You may find the following values useful:

$$\begin{aligned} \log_2(1) &= 0 & \log_2(2) &= 1 & \log_2(3) &= 1.59 & \log_2(4) &= 2 \\ \log_2(5) &= 2.32 & \log_2(6) &= 2.59 & \log_2(7) &= 2.81 & \log_2(8) &= 3 \\ \log_2(9) &= 3.17 & \log_2(10) &= 3.32 \end{aligned}$$

Note also that $\log(a/b) = \log(a) - \log(b)$.

x_1	x_2	x_3	y
1	1	0	1
0	1	0	0
1	0	1	1
0	1	0	0
1	1	0	1
1	1	0	1
1	1	1	0
1	0	1	0
0	1	1	1

- (1) What is the entropy of y ? (2 points.)

$$\begin{aligned} H(S/9) &= -\frac{5}{9} \log \frac{5}{9} - \frac{4}{9} \log \frac{4}{9} \\ &= \frac{5}{9} (.45) + \frac{4}{9} (.12) \\ &\approx .49 + .52 = .99 \end{aligned}$$

.99 bits

- (2) Which variable would you split first? Justify your answer. (2 points.)

x_1	x_2	x_3
$x_1 = 0?$	001	10
$x_1 = 1?$	111100	1001101
	1001	

$\Rightarrow x_1$, lower entropy / higher IG.

x_1

- (3) What is the information gain of the variable you selected in part (2)? (2 points.)

$$\begin{aligned} IG &= H(S/9) - \frac{2}{3} H(\frac{1}{3}) - \frac{1}{3} H(\frac{1}{3}) \\ &= .99 - .92 = .08 \\ &\quad (.073 \text{ by calculator}) \end{aligned}$$

.073 ≈ .08

- (4) Draw the full decision tree learned on these data. (4 points.)

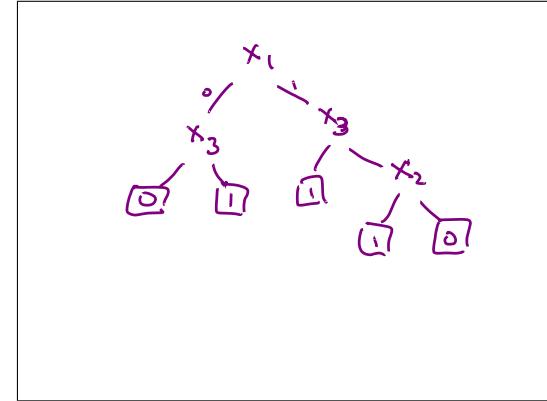
$x_1 = 0?$	$x_1 = 1?$	
$010\ 0$	$110\ 1$	$\Rightarrow x_2?$
$010\ 0$	$101\ 1$	$0 \Rightarrow 01$
$011\ 1$	$110\ 1$	$1 \Rightarrow 110$
$\Rightarrow x_3 \text{ next}$	$110\ 1$	$0 \Rightarrow 111 \checkmark$
	$111\ 0$	$1 \Rightarrow 100$
	$101\ 0$	\checkmark
	$111\ 0$	

$x_1 = 1, x_2 = ?$

$1011 \Rightarrow \text{tie, predict 1}$

1010

1110



This page is intentionally blank.

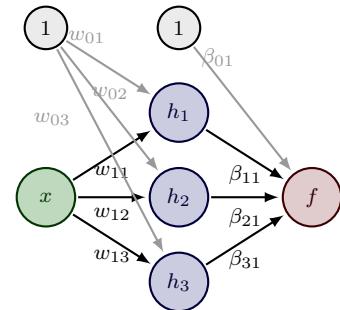
While you may use it as scratch paper, it may not be scanned for grading, so please do not include final solutions.

Problem 2 Neural Networks, (9 points.)

Consider a small neural network designed to classify a scalar feature x as one of $y \in \{-1, +1\}$. We have three hidden nodes h_1, h_2, h_3 and a single output node f_1 .

You are given the weights W of the hidden layer,

$$W = \begin{bmatrix} w_{01} & w_{11} \\ w_{02} & w_{12} \\ w_{03} & w_{13} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 3 & -2 \\ -3 & 3 \end{bmatrix}$$



and the weights B of the output layer,

$$B = [\beta_{01} \quad \beta_{11} \quad \beta_{21} \quad \beta_{31}] = [-1 \quad 3 \quad -1 \quad -1].$$

(For example, w_{12} is the weight connecting x_1 to h_2 ; w_{02} is the constant (bias) term for h_2 , etc.)

The network uses a ReLU activation function, $a(z) = \max(0, z)$, for the hidden layer, and a logistic sigmoid activation function,

$$\sigma(z) = \frac{1}{1 + \exp(-z)} = \frac{1}{1 + \frac{1}{\exp(z)}} = \frac{\exp(z)}{\exp(z) + 1},$$

for the output layer (the value of which, as usual, corresponds to the model's probability that the class is $+1$). You may find the following values useful:

$$\exp(0) = 1 \quad \exp(1) = 2.72 \quad \exp(2) = 7.39 \quad \exp(3) = 20.09 \quad \exp(4) = 54.60 \quad \exp(5) = 148.40$$

- (1) What class is predicted by the model given the input $x_1 = 2$? (3 points.)

$$\begin{aligned} x_1 = 2 \Rightarrow r_1 &= 0+2=2 & h_1 &= 2 \\ r_2 &= 3-4=-1 \Rightarrow h_2 &= 0 & \Rightarrow S = -1+6-0-3 \\ r_3 &= -3+6=3 & h_3 &= 3 \end{aligned}$$

$$\sigma(2) = \frac{\exp(2)}{\exp(2)+1} = \frac{7.4}{8.4} = 0.88$$

- (2) What is the model's estimated probability $p(y = +1 | x_1 = 2)$? (3 points.)

$$\hat{y} = +1$$

$$0.88$$

- (3) Suppose our input is $x_1 = 0$: what is the probability $p(y = +1 | x_1 = 0)$? (3 points.)

$$\begin{aligned} x_1 = 0 \Rightarrow r_1 &= 0+0=0 & h_1 &= 0 \\ r_2 &= 3-0=3 \Rightarrow h_2 &= 3 & \Rightarrow S = -1+0-3-0 \\ r_3 &= -3+0=-3 & h_3 &= 0 \end{aligned}$$

$$.018$$

$$\sigma(-4) = \frac{1}{1+\exp(4)} = \frac{1}{55.6} \approx .018$$

This page is intentionally blank.

While you may use it as scratch paper, it may not be scanned for grading, so please do not include final solutions.

Problem 3 True/False, (10 points.)

For each of the scenarios below, circle one of “true” or “false” to indicate whether you agree with the statement.

True or **false**: Changing our classifier from a nearest centroid classifier to a Gaussian Bayes classifier is likely to reduce our model’s bias.

True or **false**: When training a 1-nearest neighbor classifier, if we decide to double the amount of data available to the learner, we would expect the bias to decrease.

True or **false**: When learning a decision tree model, if we restrict our learner to only use odd-numbered features, we will reduce our variance.

True or **false**: In a random forest model, we use the best-fitting tree in our ensemble for prediction.

True or **false**: In a random forest model, we usually learn using a random subset of features at each node in order to speed up the learning process (reduce computation).

True or **false**: If, when training a random forest model, we decide to use the same subset of features across all levels of a given tree (i.e., the entire tree uses only those features), we will increase our model’s bias.

True or **false**: When clustering using k -means, we can select the number of clusters k using a hold-out (validation) data set.

True or **false**: Computing the expected discounted sum of rewards in a Markov Reward Process can be solved using either matrix algebra or dynamic programming.

True or **false**: The larger the value of gamma, the less far-sighted (less interested in future rewards) the RL agent becomes.

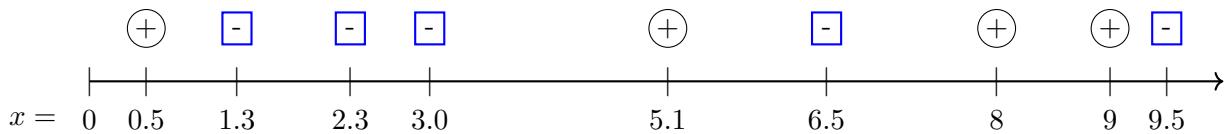
True or **false**: Once the optimal state-action value function is computed, the reinforcement learning problem can be considered solved.

This page is intentionally blank.

While you may use it as scratch paper, it may not be scanned for grading, so please do not include final solutions.

Problem 4 K-Nearest Neighbors, (12 points.)

Consider the following dataset with *nine* points shown below, for a binary classification task ($y = +, -$) with a scalar feature x . In case of ties, **prefer the negative class**. Put final answers in the box.



- (1) Compute the **training** error of a 1-nearest neighbor classifier. (3 points.)

$\frac{6}{9}$

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

$\times \times \quad \checkmark \checkmark \quad \times \times \quad \checkmark \times \times$

$\frac{6}{9}$

- (3) Compute the **training** error of 2-nearest neighbor classifier. (3 points.)

$\times \checkmark \quad \checkmark \checkmark \quad \times \checkmark \quad \checkmark \times \checkmark$

$\frac{3}{9}$

- (4) Compute the **leave-one-out** cross-validation error of 2-nearest neighbor classifier. (3 points.)

$\times \checkmark \quad \checkmark \checkmark \quad \times \times \quad \times \times \times$

$\frac{6}{9}$

This page is intentionally blank.

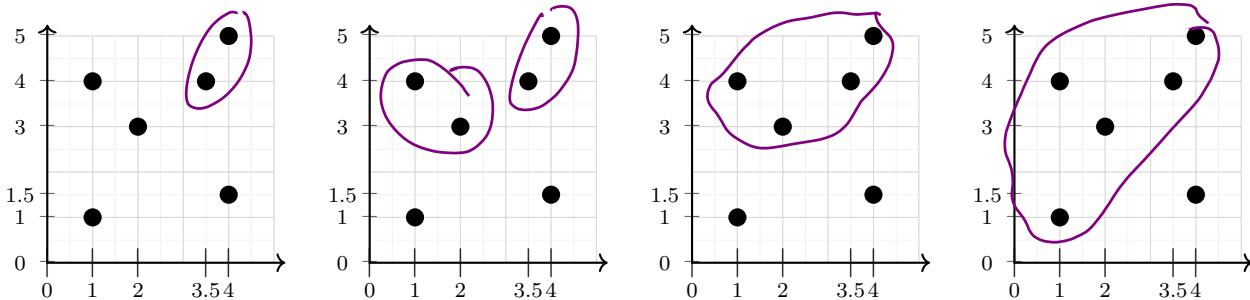
While you may use it as scratch paper, it may not be scanned for grading, so please do not include final solutions.

Problem 5 Hierarchical Clustering, (12 points.)

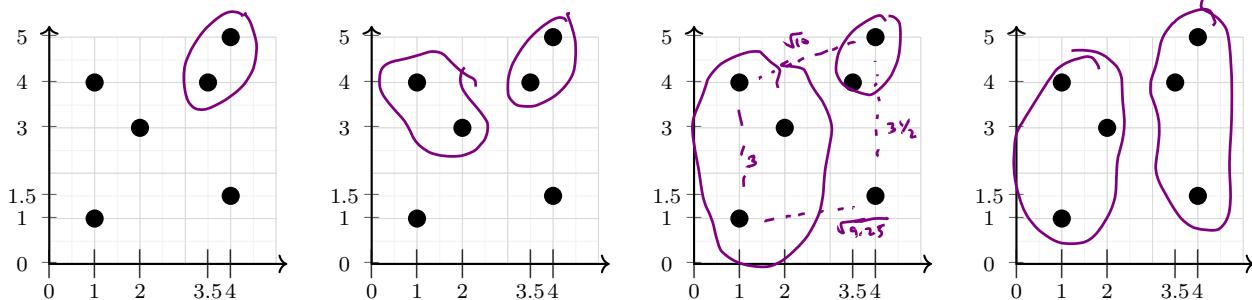
Consider the two-dimensional data points plotted in each panel. In this problem, we will cluster these data.

Linkage

- (a) Execute the hierarchical agglomerative clustering (linkage) algorithm on these data points, using “single linkage” (minimum distance) for the cluster scores. Stop when the algorithm would terminate, or after 4 steps, whichever is first. Show each step separately in a panel. (6 points.)



- (b) Now repeat your agglomerative clustering algorithm, this time using “complete linkage” (maximum distance) for the cluster scores. Stop when the algorithm would terminate, or after 4 steps, whichever is first. Show each step separately in a panel. (6 points.)



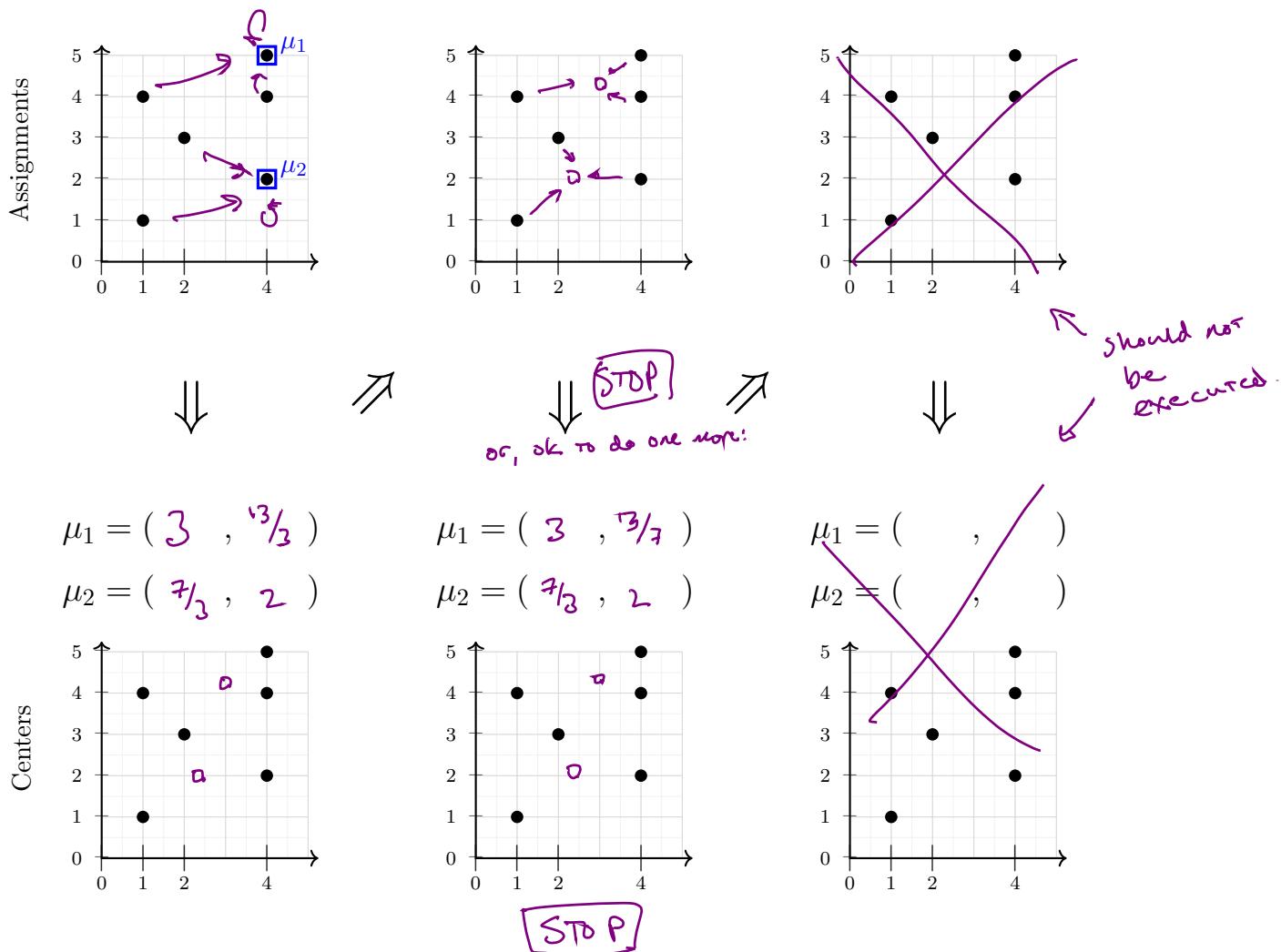
This page is intentionally blank.

While you may use it as scratch paper, it may not be scanned for grading, so please do not include final solutions.

Problem 6 K-Means Clustering, (9 points.)

Consider the 2-D data points plotted in each panel. In this problem, we will cluster these data using the k -means algorithm, where each panel is used to show a single step of the algorithm.

Starting from the two cluster centers indicated by squares, perform k-means clustering on the data. In the top panels, indicate the assignment of the data, and then in the panel below give the values of the updated cluster centers and sketch their location in the plot, so each **column** of panels shows an iteration of k -means. Stop when converged, or after 6 steps (3 iterations), whichever is first. In the case of any ties, we will prefer to assign to cluster 1. It may be helpful to recall from our nearest neighbor classifier that the set of points nearer to A than B is separated by a line.

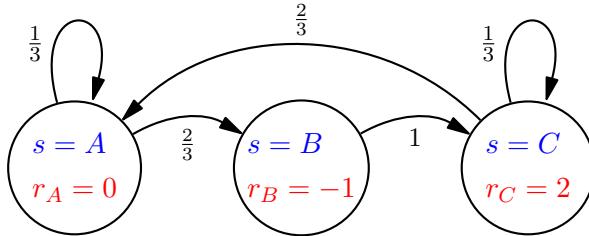


This page is intentionally blank.

While you may use it as scratch paper, it may not be scanned for grading, so please do not include final solutions.

Problem 7 Markov Processes, (12 points.)

Consider the Markov reward process model shown here:



$$\begin{aligned}\Pr[A \rightarrow A] &= 0.33 \\ \Pr[A \rightarrow B] &= 0.67 \\ \Pr[B \rightarrow C] &= 1.0 \\ \Pr[C \rightarrow A] &= 0.67 \\ \Pr[C \rightarrow C] &= 0.33\end{aligned}$$

where the transition probabilities are shown next to each arc and at right, and the rewards r_s associated with each state s are shown inside the circles. We will use dynamic programming to (start) computing the expected discounted sum of rewards. Assume a future discounting factor of $\gamma = \frac{1}{2}$.

- (1) Compute $J^1(s)$, the expected discounted sum of rewards for state sequences of length 1 (e.g., $[A]$) starting in each state s . (4 points.)

$$J^1(A) = \boxed{0}$$

$$J^1(B) = \boxed{-1}$$

$$J^1(C) = \boxed{2}$$

- (2) Compute $J^2(s)$, the expected discounted sum of rewards for state sequences of length 2 (e.g., $[C \rightarrow A]$) starting in each state s . (4 points.)

$$J^2(A) = \boxed{-\frac{1}{3} \approx -0.33}$$

$$0 + \frac{1}{2} \left(\frac{1}{3} \cdot 0 + \frac{2}{3} \cdot (-1) \right) = -\frac{1}{3}$$

$$J^2(B) = \boxed{0}$$

$$-1 + \frac{1}{2} (1 \cdot 2) = 0$$

$$J^2(C) = \boxed{\frac{7}{3} \approx 2.33}$$

$$2 + \frac{1}{2} \left(\frac{1}{3} \cdot 2 + \frac{2}{3} \cdot 0 \right) = \frac{7}{3}$$

- (3) Compute $J^3(s)$, the expected discounted sum of rewards for state sequences of length 3 (e.g., $[C \rightarrow A \rightarrow A]$) starting in each state s . (4 points.)

$$J^3(A) = \boxed{-\frac{2}{9} \approx -0.222}$$

$$0 + \frac{1}{2} \left(\frac{1}{3} \cdot -\frac{1}{3} + \frac{2}{3} \cdot 0 \right) = -\frac{2}{9}$$

$$J^3(B) = \boxed{\frac{1}{6} \approx 0.166}$$

$$-1 + \frac{1}{2} \left(1 \cdot \frac{7}{3} \right) = \frac{1}{6}$$

$$J^3(C) = \boxed{\frac{41}{18} \approx 2.278}$$

$$2 + \frac{1}{2} \left(\frac{1}{3} \cdot \frac{7}{3} + \frac{2}{3} \left(-\frac{1}{3} \right) \right) = \frac{41}{18}$$

This page is intentionally blank.

While you may use it as scratch paper, it may not be scanned for grading, so please do not include final solutions.

Name: ID#:

*This page is intentionally blank.
It may be used for scratch work or to provide additional space for solutions if necessary.*

This page is intentionally blank.

While you may use it as scratch paper, it may not be scanned for grading, so please do not include final solutions.

Name:

ID#:

*This page is intentionally blank.
It may be used for scratch work or to provide additional space for solutions if necessary.*

This page is intentionally blank.

While you may use it as scratch paper, it may not be scanned for grading, so please do not include final solutions.