

CS273a Midterm Exam
Machine Learning & Data Mining: Fall 2013
Thursday November 7th, 2013

Your name:

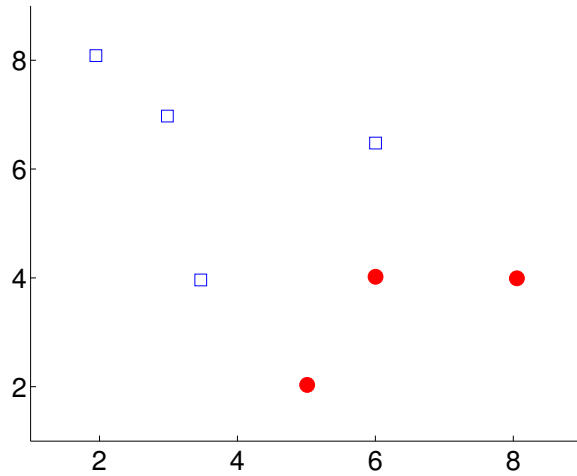
Name of the person in front of you (if any):

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

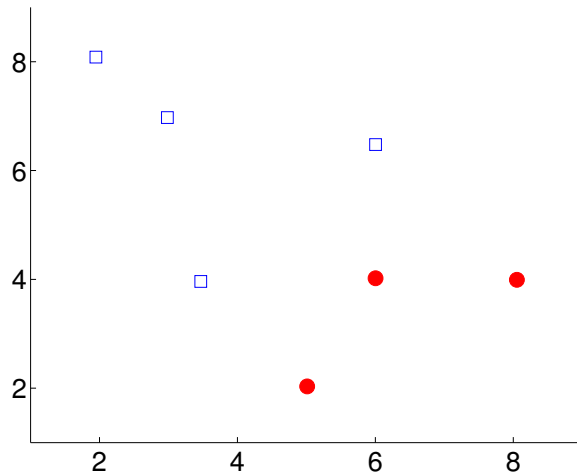
- Total time is 1:15. READ THE EXAM FIRST and organize your time; don't spend too long on any one problem.
- Please **write clearly** and **show all your work**.
- If you need clarification on a problem, please raise your hand and wait for the instructor to come over.
- Turn in any scratch paper with your exam.

Problem 1: (9 points) K-Nearest Neighbor Classification

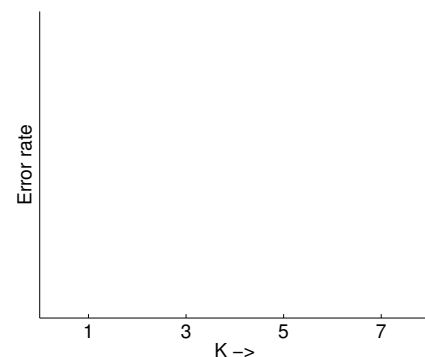
Consider the following set of training data, consisting of two-dimensional real-valued features and a binary class value, for a k-nearest-neighbors classifier. Positive data are shown as circles, negative as squares.



(1) Sketch the decision boundary for $k = 1$. Show your work and justify your answer in a few sentences (2-3).



(2) Sketch the decision boundary for $k = 5$, in the relevant part of the feature space (i.e., near the training data). Again, show your work and justify your answer in a few sentences.



(3) Sketch the basic shape you would expect to see for the error rate on training data, and on test data, as a function of increasing $k = 1, \dots, 7$. For the training error rate, **indicate the values** (error rates) of the endpoints ($k = 1$ and $k = 7$).

Problem 2: (10 points) Under- and Over-fitting

Circle one answer for each:

When training a linear classifier with gradient descent, we decrease the maximum number of iterations performed by the algorithm. This will make it **more** **equally** **less** likely to overfit the data.

Suppose we are using a *weighted* neighbor regression model, where for a test point x we assign the data weights $w^{(i)} = \exp(-c\|x - x^{(i)}\|)$. Increasing c will make it **more** **equally** **less** likely to overfit the data.

When training a decision tree, we had a parameter (**minParent**) that forced us to *never* split a node if there were fewer than **minParent** data in that node. Suppose we decrease **minParent** from its current value. This will make it **more** **equally** **less** likely to overfit the data.

Adding features to a decision tree classifier will make it **more** **equally** **less** likely to overfit the data.

For the next several questions, suppose that we are using a linear classifier, and we currently believe our model to be **overfitting**. We decide to increase the number of training data used by our learner. Choose one answer for each part:

Training error will most likely **increase** **stay the same** **decrease**.

Test error will most likely **increase** **stay the same** **decrease**.

The VC dimension of our learner will most likely **increase** **stay the same** **decrease**.

Now suppose we believe our model is **underfitting**. We again increase the number of training data. Choose one answer for each part:

Training error will most likely **increase** **stay the same** **decrease**.

Test error will most likely **increase** **stay the same** **decrease**.

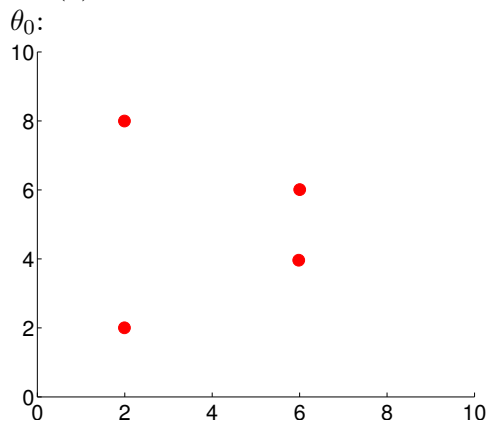
The VC dimension of our learner will most likely **increase** **stay the same** **decrease**.

(Hint: think about the relationship between training and test error rates in each regime, what this will mean for our performance on the newly added data, and how much these new data will influence our model parameters.)

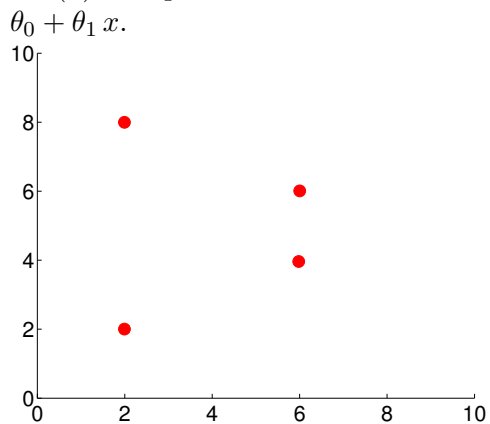
Problem 3: (8 points) Cross-validation and Linear Regression

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, $\hat{y}(x) =$



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



Problem 4: (11 points) Gradient Descent

Suppose that we wish to train the following non-linear regression model:

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

(a) Given data $D = \{(x^{(i)}, y^{(i)})\}$, write down the mean-squared error loss function $J(w)$ for the model, and compute the gradient of $J(w)$ with respect to the weights.

(b) Suppose now that, instead of feature 2, we have a transformed version of feature 1, i.e.,

$$\hat{y}(x) = \exp(w_0 + w_1x_1 + w_2h_2(x_1))$$

where $h_2(x_1) = \log(x_1 + \alpha)$, and we wish to also update α using gradient descent. Give the derivative of J with respect to α .

(c) Suppose that, instead of using an exponentiated linear function as described above, we transformed our target y by taking its log, $\tilde{y} = \log(y)$, and performed a standard linear regression with target \tilde{y} . What loss function is this minimizing? How would its predictions likely differ from the nonlinear regression above? (Hint: comment on the relative magnitude of errors for different values of y .)

Problem 5: (12 points) Decision Trees

We plan to use a decision tree to predict an outcome y using four features, x_1, \dots, x_3 . We observe nine 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,

$y = 0$: [001], [010], [010], [110]

$y = 1$: [000], [011], [101], [111]

(Note that one feature vector is observed twice.)

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

$\log_2(1) = 0$ $\log_2(2) = 1$ $\log_2(3) = 1.59$ $\log_2(4) = 2$ $\log_2(5) = 2.32$

$\log_2(6) = 2.59$ $\log_2(7) = 2.81$ $\log_2(8) = 3$ $\log_2(9) = 3.17$ $\log_2(10) = 3.32$

In 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.

- (a) What is the entropy of y ?

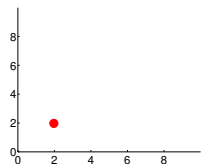
- (b) Which variable would you split first? Justify your answer.

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

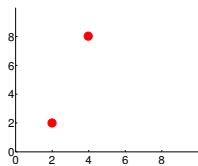
- (d) Draw the rest of the decision tree learned on these data.

Problem 6: (8 points) Shattering & VC Dimension

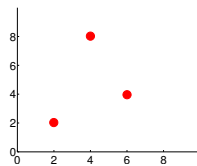
Which of the following examples can be shattered by each of the learners below? (You do not have to formally prove, but justify your answer briefly.)



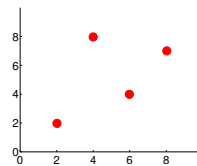
(a)



(b)



(c)



(d)

For the two learners, $T[z]$ is the sign threshold function, $T[z] = +1$ for $z \geq 0$ and $T[z] = -1$ for $z < 0$. The learner parameters a, b, c are real-valued scalars, and each data point has two real-valued input features x_1, x_2 .

(a) $\hat{y} = T[(x_1 - a)^2 + (x_2 - b)^2 + c]$ (note: uses both x_1, x_2)

(b) $\hat{y} = T[ax_1 + b \exp(x_1)]$ (note: uses only x_1 !)