

Exam #1 solution sketches

Instructions. This is a 120-minute test. You may use your notes. You may assume anything that we proved in class or in the homework is true.

Question	Score	Points
1		10
2		10
3		10
4		10
5		10
6		10
Out Of		60

Name: _____

edX Username: _____

1.

- (a) Recall that the misclassification error rate function for a decision tree is $C(a) = 1 - \max(a, 1 - a)$. Graph this function and its negation $-C(a)$ for $a \in [0, 1]$. Is $-C(a)$ a convex function? (You do not need to prove your answer.) Recall that a convex function $f : \mathbb{R} \rightarrow \mathbb{R}$ is one that satisfies

$$\forall x_1, x_2 \in \mathbb{R}, \forall \lambda \in [0, 1] : f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$$

- (b) Fix some training set with Boolean labels S . Recall the algorithm described in lecture for building a decision tree in class. Let T_0 be the empty tree. What is the misclassification error rate of T_0 in terms of the fraction of negative examples in S ? Let T_1 be the tree obtained after executing one iteration of the algorithm described in lecture using the misclassification error rate function $C(a)$ above. So, now T_1 has some variable at the root and two leaves. Prove that the misclassification error rate of T_1 is less than or equal to the misclassification error rate of T_0 . (*Hint.* Use part (a); the proof is short.)

Solution sketch.

- (a) (4 pts.) $-C(a)$ is convex.
- (b) (6 pts.) If we suppose that a λ fraction of the points are on the right side of the root and $1 - \lambda$ on the left, then

$$\begin{aligned} \text{err}(T_1) &= \lambda C(\text{fraction of negative examples on right}) + \\ &\quad (1 - \lambda) C(\text{fraction of negative examples on left}) \\ &\leq C(\text{fraction of negative examples overall}) \\ &= \text{err}(T_0) \end{aligned}$$

by convexity of $-C(a)$.

2. Start with vector $w = (1, 0)$ and run the Perceptron algorithm using data points 1 through 5 below. What is the output after running the algorithm on points 1 through 5? Compute an estimate for its generalization error rate using held-out points 6 through 9 below. Show your work.

	Point	Label
1	$(-1, 1)$	+1
2	$(-0.5, 1.5)$	+1
3	$(1, 1)$	-1
4	$(-1, 0)$	+1
5	$(-1, 2)$	-1
6	$(0.5, 0.5)$	-1
7	$(3, -1)$	+1
8	$(0.4, 0.6)$	+1
9	$(0, 1)$	-1

Solution sketch. (7 pts.) The final weight vector is $(0, -2)$. Note: you needed to make just one pass.

(3 pts.) Generalization error estimate: $1/4$.

3. In this problem we look at PAC-learning using a consistent learner. Recall that a *consistent learner* is one that, given a training data set, outputs a classifier that classifies the entire data set correctly. We work in the Boolean setting, where the domain is $\{0, 1\}^n$, and the labels lie in $\{0, 1\}$.

- (a) Let \mathcal{H} be the concept class of Boolean literals, i.e. functions of the form $h_i(x) = x_i$ or $h_{\neg i}(x) = \neg x_i$. How large is this class?
- (b) Describe a simple and efficient consistent learner for \mathcal{H} . That is, given any finite training set of labeled points $\{(x^1, h(x^1)), \dots, (x^m, h(x^m))\}$ for some $h \in \mathcal{H}$, describe a procedure to come up with a function $h \in \mathcal{H}$ such that h is consistent with all the points in the training set. (Is brute force enough?)
- (c) Write pseudocode describing a simple and efficient PAC-learner for \mathcal{H} that makes use of the consistent learner from part (b). State its sample complexity (i.e. how many training examples it needs) as a function of n , ϵ and δ . (No need for a proof; just use the right theorem from lecture.)
- (d) Let \mathcal{H}' be the concept class of majorities over literals, which are functions of the form $\text{MAJ}(\ell_1, \ell_2, \dots, \ell_n)$ where each literal ℓ_i is either x_i or $\neg x_i$. (Assume n is odd, so there are no ties.) How large is this class? Would the brute force approach yield an efficient consistent learner for \mathcal{H}' ?
- (e) Suppose that we had a consistent learner for \mathcal{H}' . Describe a PAC-learner for \mathcal{H}' , and state its sample complexity as a function of n , ϵ and δ (use the same theorem as for (c)). Compare your answer with the answer for \mathcal{H} in part (c).

Solution sketch.

- (a) (2 pts.) $|\mathcal{H}| = 2n$.
- (b) (2 pts.) Brute force all $2n$ possibilities to find one that classifies everything correctly. This takes $O(nm)$ time at most.
- (c) (2 pts.) Draw a sufficiently large number m of random training points, and then find a consistent hypothesis as above. We will need $m = \frac{1}{\epsilon} \log \frac{|\mathcal{H}|}{\delta} = \frac{1}{\epsilon} \log \frac{2n}{\delta}$.
- (d) (2 pts.) $|\mathcal{H}'| = 2^n$. The brute force approach would no longer work since it would take exponential time.
- (e) (2 pts.) Again, draw a sufficiently large number m of training points and find a consistent hypothesis. In this case $m = \frac{1}{\epsilon} \log \frac{|\mathcal{H}'|}{\delta} = \frac{1}{\epsilon} \log \frac{2^n}{\delta} = \frac{1}{\epsilon} (n + \log \frac{1}{\delta})$. In (c) the sample complexity was logarithmic in n , whereas here it is linear.

4.

- (a) Suppose we have a data set consisting of three points in \mathbb{R}^2 : $(1, 2), (2, 4), (3, 6)$. How many principal components does this data set have? Write down the first principal component.
- (b) Given the SVD of matrix

$$A = U\Sigma V^T = \begin{bmatrix} 3 & 7 & 11 \\ 6 & -1 & -5 \\ 3 & 10 & 18 \end{bmatrix}$$
$$= \begin{bmatrix} -0.531 & 0.215 & -0.819 \\ 0.162 & 0.975 & 0.150 \\ -0.832 & 0.053 & 0.553 \end{bmatrix} \cdot \begin{bmatrix} 25.0197 & 0 & 0 \\ 0 & 6.916 & 0 \\ 0 & 0 & 0.416 \end{bmatrix} \cdot \begin{bmatrix} -0.124 & -0.487 & -0.864 \\ 0.962 & 0.153 & -0.225 \\ 0.242 & -0.859 & 0.449 \end{bmatrix}.$$

Write down the matrix that is the best rank-2 approximation to A . You don't need to calculate the exact numbers, a formula or expression is enough.

Solution sketch.

- (a) (5 pts.) This is a rank-1 dataset, so there is just one principal component. It is $(\frac{1}{\sqrt{5}}, \frac{2}{\sqrt{5}})$ (since it must be a unit vector).
- (b) (5 pts.) This is given by keeping only the top two singular values.

$$\begin{bmatrix} -0.531 & 0.215 & -0.819 \\ 0.162 & 0.975 & 0.150 \\ -0.832 & 0.053 & 0.553 \end{bmatrix} \cdot \begin{bmatrix} 25.0197 & 0 & 0 \\ 0 & 6.916 & 0 \\ 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} -0.124 & -0.487 & -0.864 \\ 0.962 & 0.153 & -0.225 \\ 0.242 & -0.859 & 0.449 \end{bmatrix}.$$

5.

- (a) In each of the following plots, a training set of data points X in \mathbb{R}^2 labeled either $+$ or $-$ is given, where the original features are the coordinates (x, y) . You can assume that the data is origin-centered (despite what the axes may suggest). For each of the two training sets below, answer the following questions:
- (i) Draw all the principal components (eyeball it).
 - (ii) Can we correctly classify this dataset by using a halfspace after projecting onto one of the principal components? If so, which principal component should we project onto? If not, explain in 1–2 sentences why it is not possible.
- (b) Is it possible to have a data set in \mathbb{R}^2 that is linearly separable by a halfspace in \mathbb{R}^2 but is not linearly separable after projecting onto *either* of the two principal components? If so, give a simple example along the lines of the above data sets. If not, explain in 1–2 sentences why it is not possible.

Solution sketch.

- (a) (6 pts.)
- (i) The principal components are the same for both datasets and are shown below.
 - (ii) Dataset 1 can be correctly classified by a halfspace after projecting onto the first principal component. Dataset 2 cannot be correctly classified by a halfspace after projecting onto either principal component.
- (b) (4 pts.) Yes, this is possible and an example is shown below as Dataset 3. Essentially, it can happen when the true halfspace in 2D is at an angle to both the principal components.

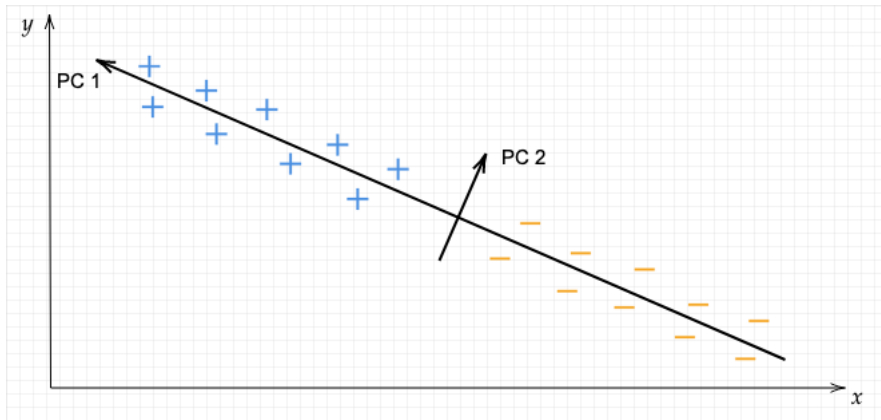


Figure 1: Dataset 1, with principal components

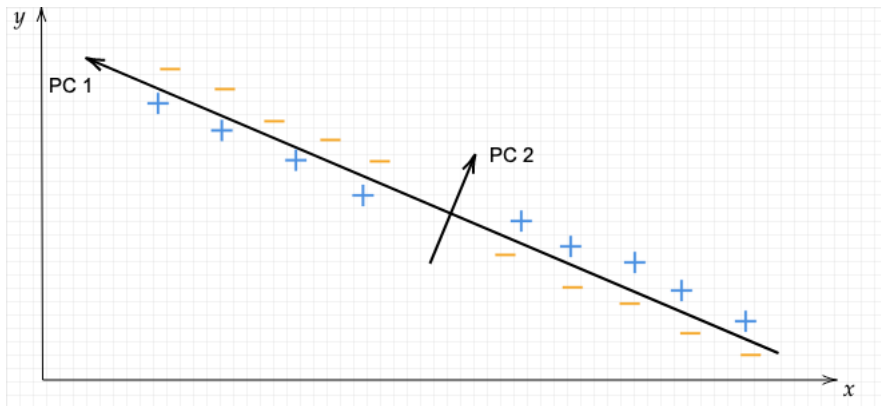


Figure 2: Dataset 2, with principal components

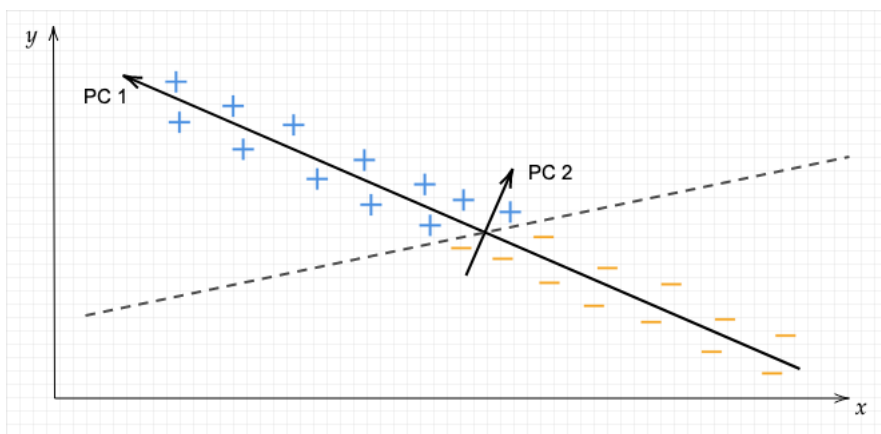


Figure 3: Dataset 3, with principal components and labeling halfspace

6. Regression problems.

- (a) You are given a data set $S = (x_1, y_1), \dots, (x_t, y_t)$ where each x_i and y_i are real numbers. You perform simple linear regression to obtain the line $\beta_0 + \beta_1 x$. Now re-scale the x_i 's so that $x'_i = \alpha x_i$ for some real number α . Perform simple linear regression again. How do the coefficients β_0, β_1 change for the new line, quantitatively? You may reason by drawing a picture or using formulas for these coefficients from class.
- (b) For each of the following scenarios, state whether or not we can use linear regression, and give a short reason.
- (i) We have training data (x, y) (where $x \in \mathbb{R}^2, y \in \mathbb{R}$) satisfying $y = \alpha x_1 + \beta x_2$, and we want to learn the model parameters α, β . (That is, we have training data of the above form for various different x .)
 - (ii) We have training data (x, y) (where $x \in \mathbb{R}^2, y \in \mathbb{R}$) satisfying $y = \alpha x_1 x_2^3$, and we want to learn the model parameter α .
 - (iii) We have training data (x, y) (where $x \in \mathbb{R}^2, y \in \mathbb{R}$) satisfying $y = 2^\alpha x_1^\beta$, and we want to learn the model parameters α, β .

Solution sketch.

- (a) (4 pts.) Only the slope changes: $\beta'_0 = \beta_0, \beta'_1 = \beta_1/\alpha$.
- (b) (6 pts.)
- (i) Yes, because the dependence on the weights α, β is linear.
 - (ii) Yes, because the dependence on α is linear.
 - (iii) No, because the dependence on α, β is *not* linear.