

בית ספר "אפי ארזי" למדעי המחשב המרכז הבינתחומי
The Efi Arazi school of computer science
The Interdisciplinary Center

סמסטר ב' תש"פ
Spring 2020

מבחן מועד א בלמידה ממוכנת
Machine Learning Exam A

Lecturer: Prof Zohar Yakhini
Time limit: 2 hours

מרצה: פרופ זהר יכני
משך המבחן: 2 שעות

**You should answer Qn 1 (mandatory),
for 40 pts and two out of the other three
questions (30pts each, 60pts together).**

**In the first page indicate the numbers of
the two questions you answered. If there
is no indication then the first two
solutions will be graded.**

Good Luck!

**יש לענות על שאלה מספר 1 (חובה), 40
נקודות, ולבחור שתיים מתוך שלוש השאלות
הנוספות (30 נקודות לכל אחת, סה"כ 60
נקודות).**

**בעמוד הראשון יש לרשום את מספרי שתי
שאלות הבחירה שעליהן בחרת לענות. אם
לא יהיה רשום תיבדקנה שתי התשובות
הראשונות.**

בהצלחה!

You can use all Moodle course materials as well as student personal notes prepared before the exam. There is no need to print the material. You can use the appropriate files on your PC.

You can use calculators.

Justify all your answers and show your calculations.

All answers should be legibly written and scanned for submission as per IDC's instructions.

ניתן להשתמש בכל החומר הקיים ב Moodle ובנוסף סיכומי סטודנטים שנכתבו לפני תחילת המבחן. אין צורך להדפיס את החומר. אפשר להשתמש במחשב כדי לגשת לקבצים הרלוונטיים.

ניתן להשתמש במחשבון.

יש להצדיק את כל תשובותיך ולהראות את צעדי החישוב.

כל התשובות צריכות להיות כתובות בכתב ברור וקריא ולהיסרק להגשה ע"פ הנחיות IDC.

Question 1 (3 parts, 40 points) – MANDATORY QUESTION

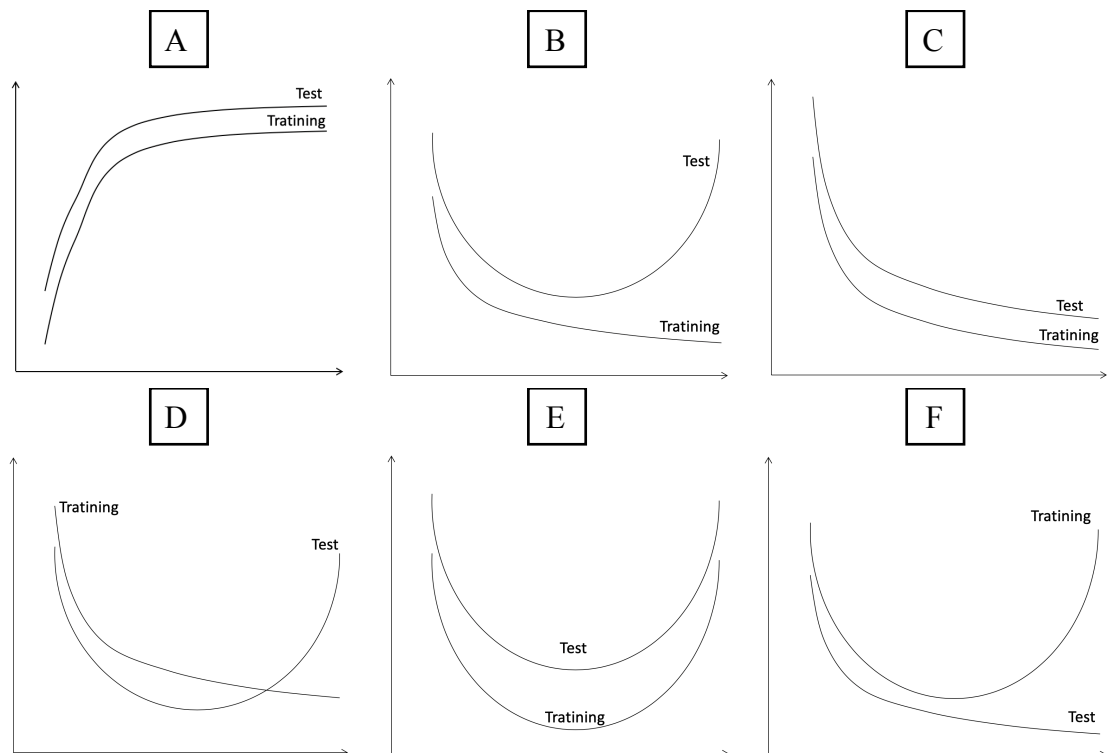
A. You receive 10,000 samples of labeled data, split into training and test, to perform regression to predict house prices from 4 real valued features (x_1, x_2, x_3, x_4) .

- (4 points) If you perform Linear Regression with all 4 features, what is the dimension of the inferred vector of coefficients $\vec{\theta}$?
- (5 points) If you perform Polynomial Regression using all monomials of degree less or equal 5 (for example: $x_1, x_1^2, x_1^3, x_1^4, x_1^5, x_1^2 x_2^2, x_2 x_3, x_1^2 x_2^2 x_3, x_4^3 x_3^2, x_3^3$). What is the dimension of the inferred vector of coefficients $\vec{\theta}$?
- (5 points) You performed Polynomial Regression using all monomials of degree less or equal 2. You obtained an MSE on both training and test with which you are not happy. Your colleague proposes that you also add $2x_1 x_2$ and $2x_3 x_4$ to your vector of features since these are important quantities in the context of pricing houses. Will this change help the MSE performance? Explain.
- (6 points) Assume that your training / test split was 8,000 / 2,000. You perform Polynomial Regression of increasing degrees.

For each of the following graphs state whether it can conceivably describe the observed performance.

In all graphs below the y-axis represents the MSE and the x-axis represents the degree of the Polynomial Regression.

Explain.



B. (10 points) Let $X = \mathbb{R}^3$. Calculate the VC dimension of the set of all cylinders threaded on the z axis. The cylinder can be both on the positive and negative part of z . Consider data points inside the cylinder as Positives and data points outside of the cylinder classify as Negatives.

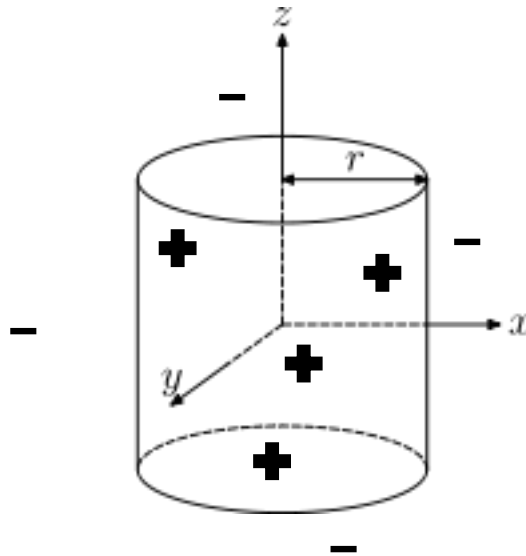
Formally:

For any three numbers $l, u \in \mathbb{R}, r \in \mathbb{R}_+$

we define $h(l, u, r) = \{(x, y, z) \mid x^2 + y^2 \leq r \wedge l \leq z \leq u\}$.

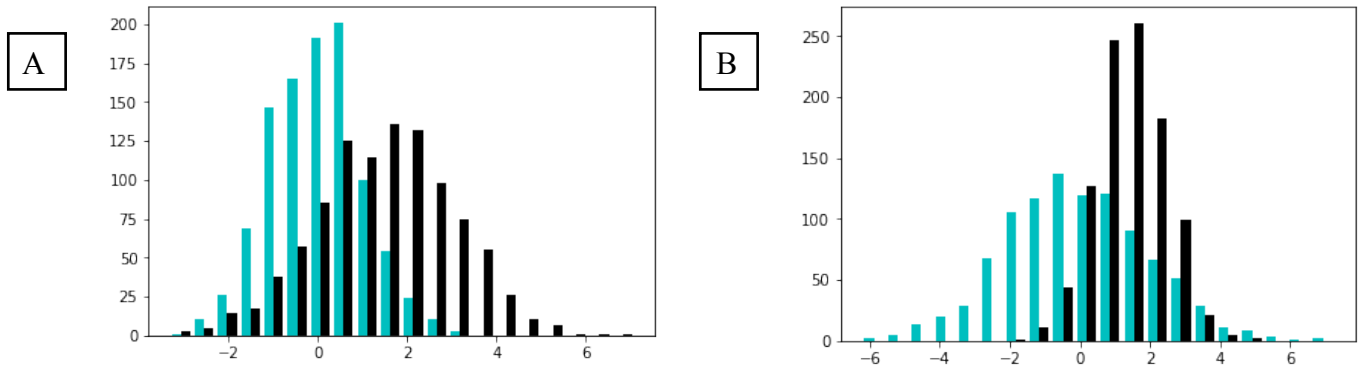
And now we define the hypotheses space as:

$H = \{h(l, u, r) \mid l, u \in \mathbb{R}, r \in \mathbb{R}_+\}$

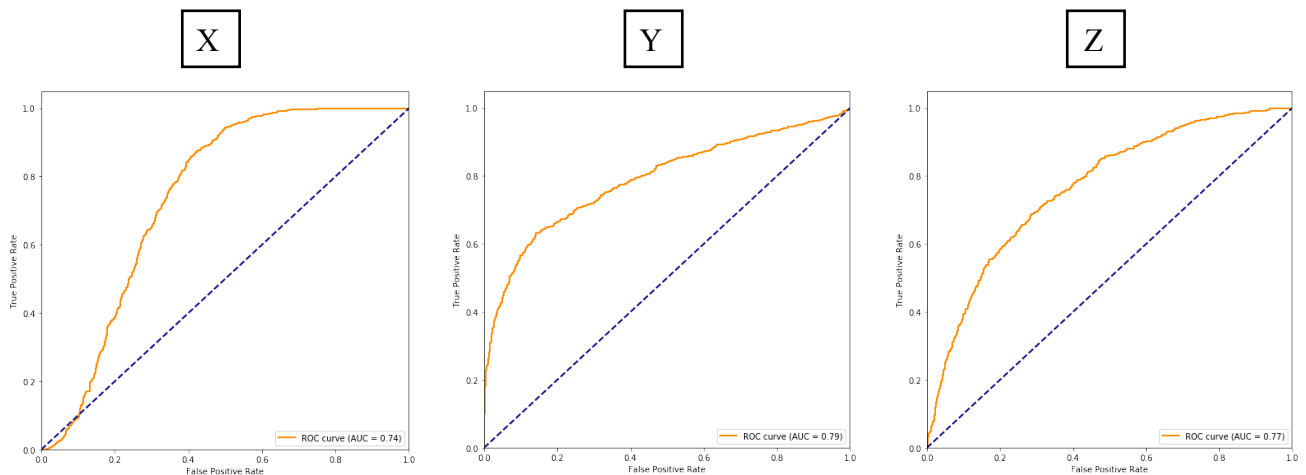


- C. Two companies are proposing a Corona test to a hospital. The hospital asks for your advice as to which test is better. The plots below represent the distributions of the tested quantity for the two companies. The black histograms represent the positive cases. We further depict the ROC curves computed to evaluate their performance.

Distributions



ROC curves



- (4 points) Match two of the three ROC curves above, X, Y and Z, with the two distributions, A and B. Explain.

Recall that the expected benefit of a test, for the hospital, is measured by

$$\pi = \alpha * TPR - FPR$$

- (3 points) Assuming $\alpha = 23$ which company will you select?
- (3 points) Assuming $\alpha = 0.5$ which company will you select?

Question 1 (3 parts, 40 points) – MANDATORY QUESTION – solution

A.

1. The dimension of the inferred vector of coefficients $\vec{\theta}$ will be 5:

$$\vec{\theta} = (\theta_0, \theta_1, \theta_2, \theta_3, \theta_4)$$

2. The dimension of the inferred vector of coefficients $\vec{\theta}$ will be

$$\binom{n+r}{n} = \binom{4+5}{4} = \binom{9}{4} = \frac{9!}{4!5!} = \frac{9*8*7*6}{4*3*2*1} = 126$$

3. This won't change the MSE. The only change will be that the weight of the monomials x_1x_2, x_3x_4 will be split with the new monomials. Adding features that are linear combinations of other features don't add any new information and therefore can't improve the MSE.

4.

- A. No – Training MSE should decrease with the increase in the polynomial degree.
- B. Yes – The test MSE decreasing until the model starting to be too complex for the data = overfit. The training is decreasing as expected.
- C. Yes – This can happen if we still didn't reach overfit conditions.
- D. Yes – Same as B, with the addition that the test MSE can be lower than the training MSE.
- E. No – Training MSE should decrease with the increase in the polynomial degree.
- F. No – Training MSE should decrease with the increase in the polynomial degree.

B. The VC dimension is 3.

$$VC \geq 3:$$

Consider the points

$$x_1 = (0, 0, 1), x_2 = (0, 0, -1), x_3 = (0, 1, 0).$$

$$x_1 = +, x_2 = +, x_3 = + \Rightarrow h(l, u, r) = h(-2, 2, 2)$$

$$x_1 = -, x_2 = -, x_3 = - \Rightarrow h(l, u, r) = h(0, 0, 0)$$

$$x_1 = +, x_2 = +, x_3 = - \Rightarrow h(l, u, r) = h(-2, 2, 0.5)$$

$$x_1 = +, x_2 = -, x_3 = + \Rightarrow h(l, u, r) = h(-0.5, 2, 2)$$

$$x_1 = -, x_2 = +, x_3 = + \Rightarrow h(l, u, r) = h(-2, 0.5, 2)$$

$$x_1 = +, x_2 = -, x_3 = - \Rightarrow h(l, u, r) = h(-0.5, 2, 0.5)$$

$$x_1 = -, x_2 = +, x_3 = - \Rightarrow h(l, u, r) = h(-2, 0.5, 0.5)$$

$$x_1 = -, x_2 = -, x_3 = + \Rightarrow h(l, u, r) = h(-0.5, 0.5, 2)$$

$$VC < 4:$$

Consider 4 points x_1, x_2, x_3, x_4 .

Let l be the minimum of the z values.

Let u be the maximum of the z values.

Let r be the maximum of the $x^2 + y^2$ values.

W.L.O.G let's say:

x_1 has the minimum z value

x_2 has the maximum z value

x_3 has the maximum $x^2 + y^2$ value

The assignment $x_1 = +, x_2 = +, x_3 = +, x_4 = -$ can't be separated by the hypothesis space.

Therefore, for any 4 points we provided an assignment that can't be separated by any $h \in H$. Thus H does not shatter any set of 4 points.

Q.E.D

C.

1. A – Y, thresholds in the right side of the distribution lead to $TP > 0$ with $FP = 0$.

B – X, thresholds in the left side of the distribution lead to $TP = 1$ with $FP < 1$.

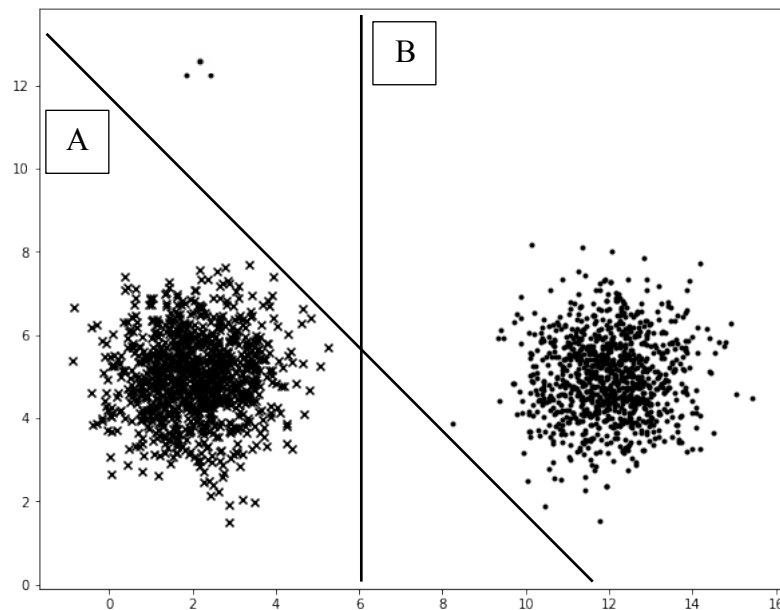
Z is a symmetric ROC curve and therefore doesn't match to any of the two distributions.

2. With $\alpha = 23$ we will prefer ROC curve X which belong to company B.

3. With $\alpha = 0.5$ we will prefer ROC curve Y which belong to company A.

Question 2 (6 parts 30 points)

- A. (5 points) Consider the data in the figure below and the two different linear decision boundaries (A & B) as depicted. One of them was obtained by running Logistic Regression on the data and the other one was obtained by running Naïve Bayes with normal distributions. Which is which? Explain.



The following pertains to parts B through F.

In rolling two dice, the first with 6 sides and the second with 4 sides, we have the following distributions of results for two casino houses – Casino A and Casino B:

Casino A

Die1 \ Die2	1	2	3	4
1	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$
2	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$
3	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$
4	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$
5	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$
6	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$

Casino B

Die1 \ Die2	1	2	3	4
1	$\frac{1}{12}$	$\frac{1}{24}$	0	$\frac{1}{24}$
2	$\frac{1}{24}$	$\frac{1}{12}$	$\frac{1}{24}$	0
3	0	$\frac{1}{24}$	$\frac{1}{12}$	$\frac{1}{24}$
4	$\frac{1}{24}$	0	$\frac{1}{24}$	$\frac{1}{12}$
5	$\frac{1}{12}$	0	$\frac{1}{12}$	0
6	0	$\frac{1}{12}$	0	$\frac{1}{12}$

The prior probability for playing in Casino A is $\frac{2}{5}$.

Given a game outcome (2 numbers), we want to classify whether the game was played in Casino A or B.

- B. (5 points) We observe the following game outcome: 1st die is 6, 2nd die is 2. Which casino will a Naïve Bayes classifier predict? Show your calculations.
- C. (5 points) Given the same game result as in Part B, which casino will a Full Bayes classifier predict? Show your calculations.
- D. (5 points) What is the minimal prior we need to assign to Casino A in order for the Full Bayes classifier to predict A regardless of the game outcome?
That is, find the minimal number π that satisfies:
(Prior of A $> \pi$) \implies (Full Bayes always decides A).
Show/explain your calculations.
- E. (5 points) Assume that the prior of A is the number π that you found in Part D. You can now change two entries in the joint distribution matrix of Casino B. Perform a change that produces a joint distribution that will lead the Full Bayes classifier to select Casino B, when observing the results (6,2), like in Part B.
- F. (5 points) What is the minimal prior we need to assign to Casino A in order for the task of part E to be impossible? Explain

Question 2 (6 parts 30 points) – solution

- A. Denote the left Gaussian as class A_1 and the right Gaussian as A_2 . Since Naïve Bayes assumes conditional independence,
 $p(x, y|A_i) = p(x|A_i) \cdot p(y|A_i)$. then on the y-axis (vertical axis), $p(y|A_1) \approx p(y|A_2)$ and hence that feature will not help determine the class
 $(\operatorname{argmax}_i p(A_i) \cdot p(x|A_i) \cdot p(y|A_i))$. Hence separation is done only using the x-axis which is why the Naïve Bayes is the B boundary (the 3 points in the top left have negligible influence on the Gaussian parameters). Logistic regression is a discriminative algorithm that looks for a separation line that best fits the data. Logistic regression therefore will produce the hyperplane that perfectly separates the data – boundary A.
 Regularization can change this but we didn't speak about regularization in class.

The following pertains to parts B through F.

Let us first add the marginal probabilities (we will use this table for the next sections as well):

Casino A					
Die2 \ Die1	1	2	3	4	
1	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{4}{24}$
2	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{4}{24}$
3	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{4}{24}$
4	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{4}{24}$
5	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{4}{24}$
6	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{1}{24}$	$\frac{4}{24}$
	$\frac{6}{24}$	$\frac{6}{24}$	$\frac{6}{24}$	$\frac{6}{24}$	1

Casino B					
Die2 \ Die1	1	2	3	4	
1	$\frac{1}{12}$	$\frac{1}{24}$	0	$\frac{1}{24}$	$\frac{4}{24}$
2	$\frac{1}{24}$	$\frac{1}{12}$	$\frac{1}{24}$	0	$\frac{4}{24}$
3	0	$\frac{1}{24}$	$\frac{1}{12}$	$\frac{1}{24}$	$\frac{4}{24}$
4	$\frac{1}{24}$	0	$\frac{1}{24}$	$\frac{1}{12}$	$\frac{4}{24}$
5	$\frac{1}{12}$	0	$\frac{1}{12}$	0	$\frac{4}{24}$
6	0	$\frac{1}{12}$	0	$\frac{1}{12}$	$\frac{4}{24}$
	$\frac{6}{24}$	$\frac{6}{24}$	$\frac{6}{24}$	$\frac{6}{24}$	1

- B. Naïve Bayes chooses $\operatorname{argmax}_i p(A_i) \cdot p(y|A_i) \cdot p(y|A_i)$.

$$p(A) \cdot p(\text{Die1} = 6|A) \cdot p(\text{Die2} = 2|A) = \frac{2}{5} \cdot \frac{4}{24} \cdot \frac{6}{24} = \frac{2}{120}$$

$$p(B) \cdot p(\text{Die1} = 6|B) \cdot p(\text{Die2} = 2|B) = \frac{3}{5} \cdot \frac{4}{24} \cdot \frac{6}{24} = \frac{3}{120}$$

Since $\frac{3}{120} > \frac{2}{120}$ Naive Bayes will predict class B.

- C. Full Bayes predicts the class by $\operatorname{argmax}_i p(A_i) \cdot p(x, y|A_i)$.

$$p(A) \cdot p(\text{Die1} = 6, \text{Die2} = 2|A) = \frac{2}{5} \cdot \frac{1}{24} = \frac{1}{60}$$

$$p(B) \cdot p(\text{Die1} = 6, \text{Die2} = 2|B) = \frac{3}{5} \cdot \frac{1}{12} = \frac{3}{60}$$

Since $\frac{3}{60} > \frac{1}{60}$ Full Bayes will also predict class B.

D. We wish that:

$$p(A) \cdot p(\text{Die1} = 6, \text{Die2} = 2|A) > p(B) \cdot p(\text{Die1} = 6, \text{Die2} = 2|B)$$

$$p(A) \cdot \frac{1}{24} > p(B) \cdot \frac{1}{12} \Rightarrow p(A) > p(B) \cdot 2$$

Since $p(B) = 1 - P(A)$ we substitute $p(B)$ in the inequality:

$$p(A) > (1 - p(A)) \cdot 2 \Rightarrow p(A) > \frac{2}{3}$$

E. We will choose $p(A) = \frac{2}{3}$.

We need to change two entries in the same row/column to maintain the marginal distributions valid.

Consider the following change:

Casino B					
Die1 \ Die2	1	2	3	4	
1	$\frac{1}{12}$	$\frac{1}{24}$	0	$\frac{1}{24}$	$\frac{4}{24}$
2	$\frac{1}{24}$	$\frac{1}{12}$	$\frac{1}{24}$	0	$\frac{4}{24}$
3	0	$\frac{1}{24}$	$\frac{1}{12}$	$\frac{1}{24}$	$\frac{4}{24}$
4	$\frac{1}{24}$	0	$\frac{1}{24}$	$\frac{1}{12}$	$\frac{4}{24}$
5	$\frac{1}{12}$	0	$\frac{1}{12}$	0	$\frac{4}{24}$
6	0	$\frac{2}{12}$	0	$\frac{0}{12}$	$\frac{4}{24}$
	$\frac{6}{24}$	$\frac{8}{24}$	$\frac{6}{24}$	$\frac{4}{24}$	1

$$p(A) \cdot p(\text{Die1} = 6, \text{Die2} = 2|A) = \frac{2}{3} \cdot \frac{1}{24} = \frac{1}{36}$$

$$p(B) \cdot p(\text{Die1} = 6, \text{Die2} = 2|B) = \frac{1}{3} \cdot \frac{2}{12} = \frac{2}{36}$$

$\frac{1}{36} < \frac{2}{36} \Rightarrow$ The Full Bayes will predict B.

F. We need:

$$p(A) \cdot p(\text{Die1} = 6, \text{Die2} = 2|A) > p(B) \cdot p(\text{Die1} = 6, \text{Die2} = 2|B)$$

$$p(A) \cdot \frac{1}{24} > (1 - p(A)) \cdot p(\text{Die1} = 6, \text{Die2} = 2|B)$$

$$p(A) \cdot \left(\frac{1}{24} + p(\text{Die1} = 6, \text{Die2} = 2|B) \right) > p(\text{Die1} = 6, \text{Die2} = 2|B)$$

$$p(A) > \frac{p(\text{Die1} = 6, \text{Die2} = 2|B)}{\frac{1}{24} + p(\text{Die1} = 6, \text{Die2} = 2|B)}$$

For all possible $p(\text{Die1} = 6, \text{Die2} = 2|B)$. The biggest probability we can add to entry (6,2) is $\frac{1}{12}$ (this is the max probability in the other entries and we wish to keep the distribution valid). Resulting $p(\text{Die1} = 6, \text{Die2} = 2|B) = \frac{2}{12}$.

$$\text{So } p(A) > \frac{\frac{2}{12}}{\frac{1}{24} + \frac{2}{12}} = \frac{4}{5}$$

Question 3 (5 parts 30 points)

A. (10 points) We want to cluster a set of S instances into k groups. Below is a pseudo code of an algorithm called k -medians with L_1 as the distance measure:

Initialize c_1, \dots, c_k by randomly.

Loop:

Assign all n instances to their closest c_i , with L_1 as the distance metric, and create k clusters S_1, \dots, S_k

For each cluster S_i ($1 \leq i \leq k$) define a new c_i :

$c_i = \text{median}(S_i)$

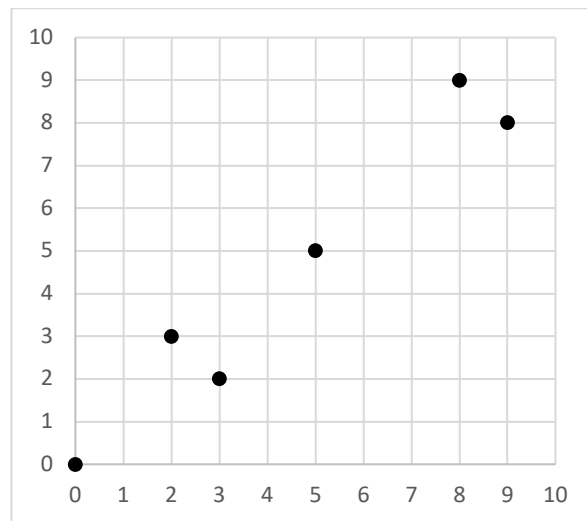
Until no change in c_1, \dots, c_k

Return c_1, \dots, c_k

* Note: the median of a set of vectors $S = v_1, \dots, v_n$ is a vector of the same dimension. Each entry of this vector is obtained by computing the median of the corresponding entries in v_1, \dots, v_n .

Consider a set of instances, S , of size 6 with 2 features as shown in the figure and the table. Run the k -medians algorithm with $k = 2$ and with the starting centers at $c_1 = (7, 9)$ and $c_2 = (7, 8)$.

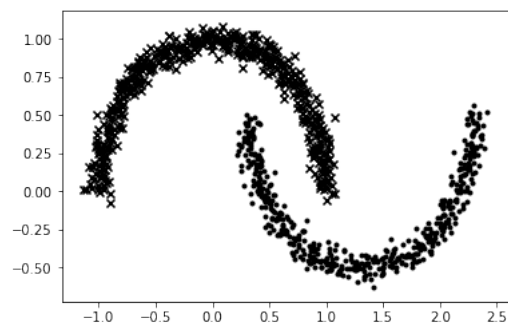
At each iteration write down the new centers and which center each point is assigned to. No need to show all intermediate calculations.



x	y
0	0
2	3
3	2
5	5
8	9
9	8

B. (5 points) Consider the following cluster structure. The 'xs' (the upper half circle) belong to Cluster 1 and the 'dots' (the lower half circle) belong to Cluster 2 (see picture).

Can it be the final assignment of the execution of k -medians? Explain.



- C. (5 points) Give data for which k-medians, as described above, converges to a different clustering assignment solution than the one resulting from running k-means, starting at the same initial points, with $k=2$. Justify your answer.
- D. (5 points) We ran k-means with $k=8$ on the following image for color quantization:



We ran the algorithm twice.
We got the following two images:



Explain how it is possible for us to get two different images.

- E. (5 points) How many parameters do you need to learn in order to infer a multivariate Gaussian mixture model in \mathbb{R}^2 with $k = 3$?

Question 3 (5 parts 30 points) – solution

A. Iteration 1:

$$c_1: \{(8,9)\} \rightarrow c_1^{new} = (8,9)$$

$$c_2: \{(0,0), (2,3), (3,2), (5,5), (9,8)\} \rightarrow c_2^{new} = (3,3)$$

Iteration 2:

$$c_1: \{(8,9), (9,8)\} \rightarrow c_1^{new} = (8.5, 8.5)$$

$$c_2: \{(0,0), (2,3), (3,2), (5,5)\} \rightarrow c_2^{new} = (2.5, 2.5)$$

- B. It cannot be the final assignment. The centers of each cluster are the medians of all the points assigned to that cluster. Although the density of the points in each cluster is unknown, it must be that the centers are contained within the convex hull of each cluster. For any position of the centers in the convex hulls, at least one point from cluster (1) will be assigned to cluster (2) in the next iteration (or vice versa).
- C. There are more than one (∞) answers to this question.
Consider the following 1D data: (1,2,4,5,100) and the initial centers $c_1 = 2, c_2 = 5$. For k-medians, it is clear that the first center will be updated to $c_1 = 1.5$ and c_2 will remain the same and this will be the final assignment. However, for k-means, the outlier will pull the second center such that $c_2 = 100$ and $c_1 = 3$.
- D. k-means is sensitive to the initialization of the centers. Different initializations might converge to different solutions. Each solution is a different local minimum.
- E. Multivariate Gaussian mixture model with $k = 3$ has 3 w's, one for each multivariate gaussian. Multivariate gaussian in \mathbb{R}^2 has 2 variances, one for each dimension, and one covariance. In addition, each multivariate gaussian in \mathbb{R}^2 has 2 means.
This means that multivariate Gaussian mixture model in \mathbb{R}^2 with $k = 3$ has 18 parameters. (We actually need to infer 17 because we know the last w from the other two).

Question 4 (3 parts 30 points)

A. (10 points) Prove:

If $\alpha, \beta > 0$ and K, L are kernels then $\alpha K + \beta L$ is also kernel.

B. Consider the instance space $X = \mathbb{R}^2$ with some probability distribution π .

Consider the concept space C that consists of right-angle isosceles triangles whose head vertex is positioned on the $x = y$ line and whose base is to the right of its head vertex.

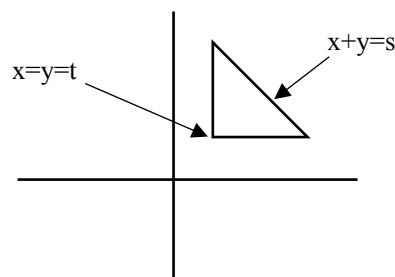
Formally, for any $t \geq 0$ and $s > t$ define

$$c(s, t) = \{(x, y) \mid x \geq t, y \geq t, x + y \leq s\}$$

and then:

$$C = \{c(s, t) \mid s > t \geq 0\}.$$

Let $H = C$.



1. (6 pts) Propose a consistent learning algorithm L that takes as input labeled points $D^m = \{(x_1, y_1), (x_2, y_2) \dots, (x_m, y_m)\} \subseteq \mathbb{R}^2$ and returns $h \in H$.
2. (7 pts) Prove that C is PAC learnable from H by computing a sufficiency bound on the sample complexity.
3. (7 pts) For $\varepsilon = 0.05$ and $\delta = 0.01$ compute a sufficiency bound on the size, m , of a set D^m of independently drawn training samples, that guarantees that for any $c \in C$ we have:
$$\text{Prob}(\text{Err}(c, L(D^m)) > \varepsilon) < \delta$$

Question 4 (3 parts 30 points) – solution

A. Since K, L are kernels, there exists two mappings ψ_1, ψ_2 respectively such that

$$\forall x, x', K(x, x') = \langle \psi_1(x), \psi_1(x') \rangle$$

And

$$\forall x, x', L(x, x') = \langle \psi_2(x), \psi_2(x') \rangle$$

Let us show that there exists a mapping ψ , such that

$$\forall x, x', (\alpha K + \beta L)(x, x') = \langle \psi(x), \psi(x') \rangle$$

Consider the following mapping: $\psi(x) = (\sqrt{\alpha} \cdot \psi_1(x), \sqrt{\beta} \cdot \psi_2(x))$

Since ψ_1, ψ_2 are mappings, ψ is a mapping.

$$\begin{aligned} \forall x, x', \langle \psi(x), \psi(x') \rangle &= \langle (\sqrt{\alpha} \cdot \psi_1(x), \sqrt{\beta} \cdot \psi_2(x)), (\sqrt{\alpha} \cdot \psi_1(x'), \sqrt{\beta} \cdot \psi_2(x')) \rangle \\ &= \langle \sqrt{\alpha} \cdot \psi_1(x), \sqrt{\alpha} \cdot \psi_1(x') \rangle + \langle \sqrt{\beta} \cdot \psi_2(x), \sqrt{\beta} \cdot \psi_2(x') \rangle \\ &= \alpha \cdot \langle \psi_1(x), \psi_1(x') \rangle + \beta \cdot \langle \psi_2(x), \psi_2(x') \rangle \\ &= \alpha \cdot K(x, x') + \beta \cdot L(x, x') = (\alpha K + \beta L)(x, x') \end{aligned}$$

And hence $\alpha K + \beta L$ is a kernel.

B.

1. The algorithm will produce a hypothesis which is the smallest relevant triangle that contains all the positive points. This can be done in $O(m)$ as follows:

Let $\Delta = \Delta^m = (x_i, y_i)_{i=1}^m$ be a set of points in the plane, labeled positive and negative (see Fig1).

Our algorithm seeks to return a hypothesis $h \in H$.

Let $(x_i, y_i)_{i=1}^{m^{(+)}}$ be all positively labeled data points.

Find:

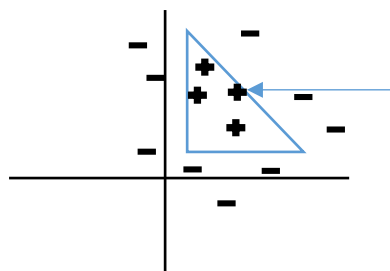
- 1) $a := \min_{1 \leq i \leq m^{(+)}} (x_i)$
- 2) $b := \min_{1 \leq i \leq m^{(+)}} (y_i)$
- 3) $s := \max_{1 \leq i \leq m^{(+)}} (x_i + y_i)$

Given that the head vertex is positioned on the $x = y$ line, the vertices of the hypothesis triangle $h = L(\Delta)$ will be:

$$\gamma = (\min(a, b), \min(a, b))$$

$$\alpha = (\min(a, b), s - \min(a, b))$$

$$\beta = (s - \min(a, b), \min(a, b))$$



The hypothesis triangle

Fig1

2. Consider $c \in C$ and let $\Delta^m(c) = (x_i(c), y_i(c))_{i=1}^m$ be training data generated from c without errors and by drawing m independent points according to a probability distribution π on \mathbb{R}^2 . We will denote the probability distribution thus induced on $(\mathbb{R}^2)^m$ by π^m .

Given $\varepsilon > 0$ and $\delta > 0$ we will now compute a number $m(\varepsilon, \delta)$ so that (Eq.1)

$$m \geq m(\varepsilon, \delta) \Rightarrow e(\Delta^m(c)) = \pi^m(\text{err}_\pi(L(\Delta^m(c)), c) > \varepsilon) \leq \delta$$

Note that $L(\Delta^m(c))$ is the hypothesis h , or the triangle, produced by L when considering data $\Delta^m(c)$ (as we describe in the previous section). $e(\Delta^m(c))$ is a random variable that depends on the stochastic behavior of $\Delta^m(c)$. It is exactly this behavior that we will want to characterize.

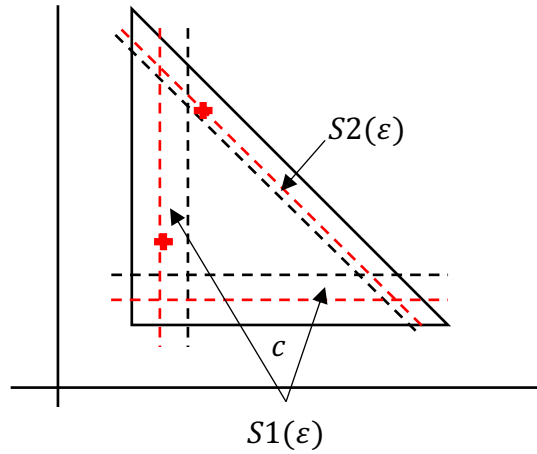


Fig2

Consider the strips parallel to the edges of the triangle c as in Fig2. Note that $S1$ has an L shape where the horizontal and vertical strips have equal size. This is due to the fact that the head vertex is positioned on the $x = y$ line.

These are defined to satisfy:

$$\pi(S1(\varepsilon)) = \pi(S2(\varepsilon)) = \frac{\varepsilon}{2}$$

Now, note that

$$\begin{aligned} \{\Delta^m(c) : \text{err}_\pi(L(\Delta^m(c)), c) > \varepsilon\} \subseteq \\ \{\Delta^m(c) : \Delta^m(c) \cap S1(\varepsilon) = \emptyset\} \cup \\ \{\Delta^m(c) : \Delta^m(c) \cap S2(\varepsilon) = \emptyset\} \end{aligned}$$

This is because if $\Delta^m(c)$ visits both strips (note that negative points can not visit these strips as there are no errors in the training labels) then, according to our construction, the difference between c and $L(\Delta^m(c))$ will have $\pi \leq \pi(S1(\varepsilon) \cup S2(\varepsilon)) < \varepsilon$. See the red triangle example in Fig2.

In terms of probability we therefore get:

$$\begin{aligned} \pi^m(\text{err}_\pi(L(\Delta^m(c)), c) > \varepsilon) &\leq \\ \pi^m(\Delta^m(c) \cap s1(\varepsilon) = \emptyset) + \pi^m(\Delta^m(c) \cap s2(\varepsilon) = \emptyset) &\leq \\ 2\left(1 - \frac{\varepsilon}{2}\right)^m \end{aligned}$$

We now select $m(\varepsilon, \delta) = \frac{2}{\varepsilon} \left(\ln(2) + \ln\left(\frac{1}{\delta}\right) \right)$ to get (Eq.1) to hold.

Q.E.D

3. We will substitute $\varepsilon = 0.05$ and $\delta = 0.01$ in the equation from the previous section:

$$\begin{aligned} m &> \frac{2}{\varepsilon} \left(\ln(2) + \ln\left(\frac{1}{\delta}\right) \right) \\ m &> \frac{2}{0.05} \left(\ln(2) + \ln\left(\frac{1}{0.01}\right) \right) \end{aligned}$$

$$m > 211.9$$

GOOD LUCK!