

CIS 520, Machine Learning, Fall 2018: Homework 1

Due: Sunday, September 16th, 11:59pm, PDF to Canvas

Instructions. Please write up your responses to the following problems clearly and concisely. We encourage you to write up your responses using L^AT_EX; we have provided a L^AT_EX template, available on Canvas, to make this easier. **Submit your answers in PDF form to Canvas. We will not accept paper copies of the homework.**

Collaboration. You are allowed and encouraged to work together. You may discuss the homework to understand the problem and reach a solution in groups up to size **two students**. However, *each student must write down the solution independently, and without referring to written notes from the joint session.* **In addition, each student must write on the problem set the names of the people with whom you collaborated.** You must understand the solution well enough in order to reconstruct it by yourself. (This is for your own benefit: you have to take the exams alone.)

Collaborators:

Type Collaborator Name Here

1 Conditional Independence in Probability Models [18 points]

Consider the following probability model. For a set of points x_1, \dots, x_n , we have k possible generating distributions, f_1, \dots, f_k . (That is, we know each point x_i was generated from one of the f_j ; we just don't know which j .) Let $z_i = \{1, \dots, k\}$ be an indicator variable which indicates that the i 'th data point x_i was generated from f_j if $z_i = j$. Furthermore, we specify that $p(z_i = j) = \pi_j$. Thus, our model specifies the following:

$$(i) \ p(x_i \mid z_i = j) = f_j(x_i), \quad (ii) \ p(z_i = j) = \pi_j.$$

Hint: The key to this problem is applying the basic rules of probability—marginalization, Bayes Rule, the chain rule, and/or conditional independence (not necessarily in that order). Note that some of the answers may follow directly from the definitions given above. You may make the independence assumption: x_1, x_2, \dots, x_n are independent of each other.

1. [4 points] Derive the formula for $p(x_i)$ in terms of (i) and (ii) above.
2. [7 points] Derive the formula for $p(x_1, \dots, x_n)$ in terms of (i) and (ii) above.
3. [7 points] Derive the formula for $p(z_u = v \mid x_1, \dots, x_n)$ in terms of (i) and (ii) above.

2 Non-Normal Norms [12 points]

1. [4 points]

Given the following four vectors,

$$x_1 = [3.1, 1.4, 0.7, 0.1]$$

$$x_2 = [3.0, 2.0, 1.0, 0.5]$$

$$x_3 = [2.9, 2.3, 0.6, 0.1]$$

$$x_4 = [3.1, 4.0, 0.7, 2.0]$$

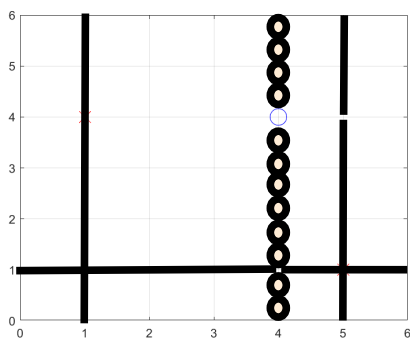
Which point (other than x_1) is closest to x_1 under each of the following norms

- a) L_0
- b) L_1
- c) L_2
- d) L_{inf}

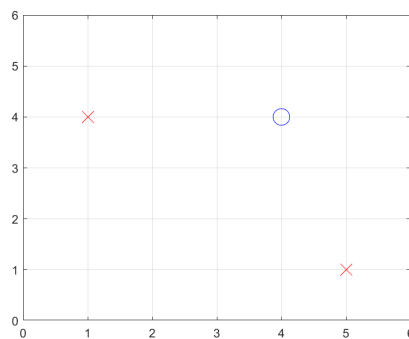
2. [8 points]

Draw the 1-Nearest Neighbor decision boundaries with the given norms and lightly shade the region in which points will be classified as circles:

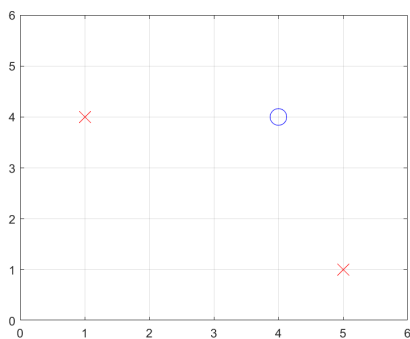
a) L_0



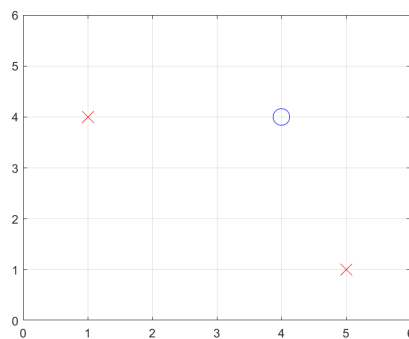
b) L_1



c) L_2



d) L_{inf}



(If there are any classification regions of zero area, draw the line or point of that classification region and label it clearly as belonging to either x or o. In the event of region that is tied, leave it unmarked.)

Hint: you can

- *scan: use your phone or a scanner to take an image of the graph with your solution drawn in and include it in the .pdf you hand in or*
- *use a pdf tool like adobe acrobat to write directly on the pdf or*
- *use any editor you like (e.g. powerpoint)*

3 Decision Trees [30 points]

1. Consider the following set of training examples for the unknown target function $\langle X_1, X_2 \rangle \rightarrow Y$. Each row indicates the values observed, and how many times that set of values was observed. For example, $(+, T, T)$ was observed once, while $(-, T, T)$ was observed 5 times.

Y	X_1	X_2	Count
+	T	T	1
+	T	F	5
+	F	T	3
+	F	F	7
-	T	T	5
-	T	F	2
-	F	T	2
-	F	F	5

Table 1

- (a) **[3 points]** What is the sample entropy $H(Y)$ for this training data (with logarithms base 2)?
 - (b) **[10 points]** What are the information gains $IG(X_1) \equiv H(Y) - H(Y|X_1)$ and $IG(X_2) \equiv H(Y) - H(Y|X_2)$ for this sample of training data?
 - (c) **[10 points]** Draw the decision tree that would be learned by ID3 (without postpruning) from this sample of training data. (Hints: Entropy is zero when one outcome is certain. Be sure to specify the predicted values for the leaves; This is done by a simple majority vote of the labels of examples classified by each leaf in the training data. Feel free to include a photo instead of using latex.)
2. When we discussed learning decision trees in class, we chose the next attribute to split on by choosing the one with maximum information gain, which was defined in terms of entropy. To further our understanding of information gain, we will explore its connection to KL-divergence (introduced above). We can define information gain as the KL-divergence from the observed joint distribution of X and Y to the product of their observed marginals.¹

$$IG(x, y) \equiv KL(p(x, y) || p(x)p(y)) = - \sum_x \sum_y p(x, y) \log \left(\frac{p(x)p(y)}{p(x, y)} \right)$$

When the information gain is high, it indicates that adding a split to the decision tree will give a more accurate model.

- (a) **[3 points]** If variables X and Y are independent, is $IG(x, y) = 0$? If yes, prove it. If no, give a counter example.

¹The negative sign is introduced in this definition because $\log(p/q) = -\log(q/p)$; flipping the fraction will give us KL as defined previously.

- (b) **[4 points]** Show that this definition of information gain is equivalent to the one given in class. That is, show that $IG(x, y) = H[x] - H[x|y] = H[y] - H[y|x]$, starting from the definition in terms of KL-divergence.

4 High Dimensional Hi-Jinx [25 points]

Nearest neighbor classifiers can be very flexible and accurate, but what if we have many irrelevant features? Suppose that we have two classes, $y = 1, 2$ and $P(\mathbf{x}|y)$ is Gaussian. Let's consider what happens to distances between points in the same class and in different classes as the dimension of \mathbf{x} grows but only one of the dimensions is informative. We'll assume for simplicity that for all Gaussians below the variance is the same, σ^2 . Reminder: If $X \sim N(\mu, \sigma^2)$, $\mathbf{E}[X] = \mu$, $\text{Var}[X] = \sigma^2$, and $\mathbf{E}[X^2] = \mu^2 + \sigma^2$. Also, recall that expectation is linear, so it obeys the following three properties:

$$\begin{aligned}\mathbf{E}[X + c] &= \mathbf{E}[X] + c \quad \text{for any constant } c \\ \mathbf{E}[X + Y] &= \mathbf{E}[X] + \mathbf{E}[Y] \\ \mathbf{E}[aX] &= a\mathbf{E}[X] \quad \text{for any constant } a\end{aligned}$$

Also, note that if X and X' are independent, then $\mathbf{E}[XX'] = \mathbf{E}[X]\mathbf{E}[X']$.

Finally, you may leave summations in your answers below. There's no need to convert to matrix notation.

- [5 points]** (Intra-class distance) Consider two points from class 1 in one dimension: $X \sim N(\mu_1, \sigma^2)$ and $X' \sim N(\mu_1, \sigma^2)$. What is the expected squared distance between them? ($\mathbf{E}[(X - X')^2] = \dots$)
- [5 points]** (Inter-class distance) Now consider two points from different classes in one dimension. Now X and X' have different means: $X \sim N(\mu_1, \sigma^2)$ and $X' \sim N(\mu_2, \sigma^2)$. What is the expected squared distance between them? ($\mathbf{E}[(X - X')^2] = \dots$)
- [5 points]** (Intra-class distance, m-dimensions) Again, consider two points from class 1 but in m dimensions. For each dimension j , we have a different mean μ_{1j} : j th dimension of X is $X_j \sim N(\mu_{1j}, \sigma^2)$ and j th dimension of X' is $X'_j \sim N(\mu_{1j}, \sigma^2)$. What is the expected squared distance between them? ($\mathbf{E}[\sum_{j=1}^m (X_j - X'_j)^2] = \dots$)
- [5 points]** (Inter-class distance, m-dimensions) Finally, consider two points from different classes but in m dimensions. That is, j th dimension of X is $X_j \sim N(\mu_{1j}, \sigma^2)$ and j th dimension of X' is $X'_j \sim N(\mu_{2j}, \sigma^2)$. What is the expected squared distance between them? ($\mathbf{E}[\sum_{j=1}^m (X_j - X'_j)^2] = \dots$)
- [5 points]** Suppose that only one dimension is informative about class values, that is $\mu_{11} \neq \mu_{21}$, but all others have the same mean $\mu_{1j} = \mu_{2j}$ for $j = 2, \dots, m$. Write down the ratio of expected intra-class distance divided by inter-class distance under this assumption. Briefly explain the significance of this ratio. As m increases towards ∞ , what value does this ratio approach? What does this limit imply about the performance of a NN classifier in this case?

5 Fitting distributions with KL divergence [15 points]

In this problem, you will use *Kullback-Leibler divergence* (KL-divergence) to measure the difference between two probability distributions. KL-divergence is an important concept in information theory and machine

learning. For more on these concepts, refer to Section 1.6 in Bishop. Please note that for most purposes (as is the case for the remainder of the course) the log is base 2.

The KL-divergence from a distribution $p(x)$ to a distribution $q(x)$ can be thought of as a distance measure from P to Q (this is just for intuition, though you should check why this is not a formal distance metric):

$$\begin{aligned} KL(p(x)||q(x)) &= \mathbf{E}_p \left[\log \frac{p(x)}{q(x)} \right] \\ &= \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx && \text{For continuous } p \text{ and } q \\ &= \sum p(x) \log \frac{p(x)}{q(x)} && \text{For discrete } p \text{ and } q \end{aligned}$$

From an information theory perspective, the KL-divergence specifies the number of additional bits required on average to transmit values of x if the values are distributed with respect to $p(x)$ but we encode them assuming the distribution $q(x)$. If $p(x) = q(x)$, then $KL(p||q) = 0$. Otherwise, $KL(p||q) > 0$. The smaller the KL-divergence, the more similar the two distributions.

1. **[8 points]** Provide the formula the Kullback-Leibler divergence $KL(p(x)||q(x))$ between two univariate Gaussians distributions:

$$p(x) = \mathcal{N}(\mu_1, \sigma^2), \quad q(x) = \mathcal{N}(\mu_2, 1).$$

Write your answer in terms of expectation over p ; i.e., fill in f and g such that $KL(p(x)||q(x)) = \mathbf{E}_p[f(x, \mu_1, \mu_2, \sigma)] + g(\sigma)$.

2. **[7 points]** For fixed μ_2 and σ , what value of μ_1 minimizes $KL(p(x)||q(x))$? At the minimum, what is the value of $KL(p(x)||q(x))$? Your answers should depend only on μ_2 and/or σ . It's not necessary to determine first and second order conditions.