

# Homework #0

CSE 446/546: Machine Learning

Profs. Jamie Morgenstern and Ludwig Schmidt

Due: **Wednesday** October 5, 2022 11:59pm

38 points

Collaborators on all problems: Thijs Masmeyer, Nick Andres, Maneeshka Madduri, Amber Chou, Faeze Aminmansov.

## Probability and Statistics

A1. [2 points] (From Murphy Exercise 2.4.) After your yearly checkup, the doctor has bad news and good news. The bad news is that you tested positive for a serious disease, and that the test is 99% accurate (i.e., the probability of testing positive given that you have the disease is 0.99, as is the probability of testing negative given that you don't have the disease). The good news is that this is a rare disease, striking only one in 10,000 people. What are the chances that you actually have the disease?

### Solution:

Let  $T+$  and  $T-$  denote the events {tested positive} and {tested negative}. Also, let  $D+$  and  $D-$  denote the events {have disease} and {have no disease}. By Bayes' rule, we have

$$\begin{aligned} P(D+|T+) &= \frac{P(T+|D+)P(D+)}{P(T+)} = \frac{P(T+|D+)P(D+)}{P(T+|D+)P(D+) + P(T+|D-)P(D-)} = \\ &= \frac{0.99 \cdot 0.0001}{0.99 \cdot 0.0001 + 0.01 \cdot 0.9999} \approx 0.0098. \end{aligned}$$

A2. For any two random variables  $X, Y$  the *covariance* is defined as  $\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$ . You may assume  $X$  and  $Y$  take on a discrete values if you find that is easier to work with.

- a. [1 point] If  $\mathbb{E}[Y | X = x] = x$  show that  $\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])^2]$ .
- b. [1 point] If  $X, Y$  are independent show that  $\text{Cov}(X, Y) = 0$ .

### Solution:

a. Let  $f$  be the PDF of  $X$ . First, note that  $\mathbb{E}[Y] = \mathbb{E}[X]$  and  $\mathbb{E}[XY] = \mathbb{E}[X^2]$ . This is because

$$\mathbb{E}[Y] = \int \mathbb{E}[Y|X = x]f(x)dx = \int xf(x)dx = \mathbb{E}[X]$$

and

$$\mathbb{E}[XY] = \int \mathbb{E}[xY|X = x]f(x)dx = \int x\mathbb{E}[Y|X = x]f(x)dx = \int x^2f(x)dx = \mathbb{E}[X^2].$$

Now, by linearity of expectation,

$$\begin{aligned} \text{Cov}(X, Y) &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \\ &= \text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] \end{aligned}$$

b. Since  $X$  and  $Y$  are independent,  $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$ . So

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] = 0.$$

A3. Let  $X$  and  $Y$  be independent random variables with PDFs given by  $f$  and  $g$ , respectively. Let  $h$  be the PDF of the random variable  $Z = X + Y$ .

- a. [1 point] Show that  $h(z) = \int_{-\infty}^{\infty} f(x)g(z-x)dx$ . (If you are more comfortable with discrete probabilities, you can instead derive an analogous expression for the discrete case, and then you should give a one sentence explanation as to why your expression is analogous to the continuous case.).
- b. [1 point] If  $X$  and  $Y$  are both independent and uniformly distributed on  $[0, 1]$  (i.e.  $f(x) = g(x) = 1$  for  $x \in [0, 1]$  and 0 otherwise) what is  $h$ , the PDF of  $Z = X + Y$ ?

**Solution:**

a. If  $X$  and  $Y$  are discrete, then we want to show  $h(z) = \sum_{x=-\infty}^{\infty} f(x)g(z-x)$ , since sum is like integral. Indeed, by independence

$$h(z) = \mathbb{P}(X+Y=z) = \sum_{x+y=z} \mathbb{P}(X=x \cap Y=y) = \sum_{x+y=z} \mathbb{P}(X=x)\mathbb{P}(Y=y) = \sum_{x+y=z} f(x)g(y) = \sum_{x=-\infty}^{\infty} f(x)g(z-x)$$

b. From above,

$$h(z) = \int_{-\infty}^{\infty} f(x)g(z-x) dx = \int_0^1 g(z-x) dx = \int_{z-1}^z g(u) du.$$

So,  $h(z) = 0$  for  $z < 0$  and  $z > 2$ , and  $h(z) = 1 - |z-1|$  otherwise.

A4. Let  $X_1, X_2, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  be i.i.d random variables. Compute the following:

- [1 point]**  $a \in \mathbb{R}, b \in \mathbb{R}$  such that  $aX_1 + b \sim \mathcal{N}(0, 1)$ .
- [1 point]**  $\mathbb{E}[X_1 + 2X_2], \text{Var}[X_1 + 2X_2]$ .
- [2 points]** Setting  $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$ , the mean and variance of  $\sqrt{n}(\hat{\mu}_n - \mu)$ .

**Solution:**

- Set  $a = \frac{1}{\sigma^2}$  and  $b = -\frac{\mu}{\sigma^2}$ , so that  $aX_1 \sim \mathcal{N}(\mu/\sigma^2, 1)$ , and  $aX_1 + b \sim \mathcal{N}(0, 1)$ .
- By linearity of expectation,

$$\mathbb{E}[X_1 + 2X_2] = \mathbb{E}[X_1] + 2\mathbb{E}[X_2] = 3\mu.$$

By independence,

$$\text{Var}[X_1 + 2X_2] = \text{Var}[X_1] + 2^2 \text{Var}[X_2] = 5\sigma^2.$$

c. By linearity of expectation,

$$\begin{aligned} \mathbb{E}[\sqrt{n}(\hat{\mu}_n - \mu)] &= \sqrt{n} (\mathbb{E}[\hat{\mu}_n] - \mu) = \sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] - \mu \right) \\ &= \sqrt{n} (\mathbb{E}[X_1] - \mu) = 0 \end{aligned}$$

By independence,

$$\begin{aligned} \text{Var}[\sqrt{n}(\hat{\mu}_n - \mu)] &= n (\text{Var}[\hat{\mu}_n]) = n \left( \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] \right) \\ &= \sigma^2. \end{aligned}$$

- **Part a:**  $a, b$ , and the corresponding calculations
- **Part b:**  $\mathbb{E}[X_1 + 2X_2], \text{Var}[X_1 + 2X_2]$
- **Part c:**  $\mathbb{E}[\sqrt{n}(\hat{\mu}_n - \mu)], \text{Var}[\sqrt{n}(\hat{\mu}_n - \mu)]$
- **Parts a-c** Corresponding calculations

**Linear Algebra and Vector Calculus**

A5. Let  $A = \begin{bmatrix} 1 & 2 & 1 \\ 1 & 0 & 3 \\ 1 & 1 & 2 \end{bmatrix}$  and  $B = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 0 & 1 \\ 1 & 1 & 2 \end{bmatrix}$ . For each matrix  $A$  and  $B$ :

- [2 points]** What is its rank?
- [2 points]** What is a (minimal size) basis for its column span?

**Solution:**

- a. Both are of rank 2.  $A_2 = 3A_1 - A_3$ , and  $B_2 = B_3 - B_1$ .  
 b. For  $A$  it's  $\{A_1, A_3\}$ . For  $B$  it's  $\{B_1, B_3\}$ . This is because  $A_2$  is a linear combination of  $A_1$  and  $A_3$ , and same for  $B$ .

A6. Let  $A = \begin{bmatrix} 0 & 2 & 4 \\ 2 & 4 & 2 \\ 3 & 3 & 1 \end{bmatrix}$ ,  $b = [-2 \quad -2 \quad -4]^\top$ , and  $c = [1 \quad 1 \quad 1]^\top$ .

- a. [1 point] What is  $Ac$ ?  
 b. [2 points] What is the solution to the linear system  $Ax = b$ ?

**Solution:**

Both are basic computations.

- a.  $Ac = [6, 8, 7]^\top$ .  
 b.  $x = [-2, 1, -1]^\top$ .

A7. For possibly non-symmetric  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$  and  $c \in \mathbb{R}$ , let  $f(x, y) = x^\top \mathbf{A}x + y^\top \mathbf{B}y + c$ . Define

$$\nabla_z f(x, y) = \left[ \frac{\partial f}{\partial z_1}(x, y) \quad \frac{\partial f}{\partial z_2}(x, y) \quad \dots \quad \frac{\partial f}{\partial z_n}(x, y) \right]^\top \in \mathbb{R}^n.$$

- a. [2 points] Explicitly write out the function  $f(x, y)$  in terms of the components  $A_{i,j}$  and  $B_{i,j}$  using appropriate summations over the indices.  
 b. [2 points] What is  $\nabla_x f(x, y)$  in terms of the summations over indices *and* vector notation?  
 c. [2 points] What is  $\nabla_y f(x, y)$  in terms of the summations over indices *and* vector notation?

**Solution:**

a. First,

$$x^\top \mathbf{A}x = \sum_{i=1}^n x_i \left( \sum_{j=1}^n A_{i,j} x_j \right) = \sum_{i=1}^n \sum_{j=1}^n x_i A_{i,j} x_j = \sum_{i \in [n], j \in [n]} A_{i,j} x_i x_j$$

Similarly,

$$y^\top \mathbf{B}y = \sum_{i \in [n], j \in [n]} B_{i,j} y_i y_j$$

So

$$f(x, y) = c + \sum_{i \in [n], j \in [n]} A_{i,j} x_i x_j + B_{i,j} y_i y_j$$

b. By linearity of derivative,

$$\nabla_x f(x, y) = \sum_{i \in [n], j \in [n]} A_{i,j} \nabla_x (x_i x_j) + B_{i,j} y_i \nabla_x (x_j).$$

with  $\frac{\partial x_i x_j}{\partial x_k} = 0$  for  $k \neq i, j$ . So the  $k^{\text{th}}$  component of  $\nabla_x f(x, y)$  is

$$[\nabla_x f(x, y)]_k = (2A_{k,k} x_k + B_{k,k} y_k) + \sum_{j \neq k} A_{k,j} x_j + \sum_{i \neq k} A_{i,k} x_i + B_{i,k} y_k,$$

where the first term comes from  $i = k, j = k$ , the second from  $i = k$ , and the third from  $j = k$ . Simplifying,

$$[\nabla_x f(x, y)]_k = (2A_{k,k}x_k + B_{k,k}y_k) + \sum_{j \neq k} A_{k,j}x_j + \sum_{i \neq k} A_{i,k}x_i + B_{i,j}y_k = \sum_{i=1}^k (A_{k,i} + A_{i,k})x_i + B_{i,k}y_k.$$

So in vector notation,  $\nabla_x f(x, y) = x^\top (A + A^\top) + y^\top B$ .

c. Since  $x^\top Ax$  does not depend on  $y$ ,

$$\nabla_y f(x, y) = (\nabla_y y)^\top Bx = Bx.$$

In terms of indices, the  $k^{\text{th}}$  component of  $\nabla_y f(x, y)$  is

$$[\nabla_y f(x, y)]_k = \sum_{i=1}^n B_{k,i}x_i.$$

A8. Show the following:

- [2 points]** Let  $g : \mathbb{R} \rightarrow \mathbb{R}$  and  $v, w \in \mathbb{R}^n$  such that  $g(v_i) = w_i$ . Find an expression for  $g$  such that  $\text{diag}(v)^{-1} = \text{diag}(w)$ .
- [2 points]** Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be orthonormal and  $x \in \mathbb{R}^n$ . An orthonormal matrix is a square matrix whose columns and rows are orthonormal vectors, such that  $\mathbf{A}\mathbf{A}^\top = \mathbf{A}^\top\mathbf{A} = \mathbf{I}$  where  $\mathbf{I}$  is the identity matrix. Show that  $\|\mathbf{A}x\|_2^2 = \|x\|_2^2$ .
- [2 points]** Let  $\mathbf{B} \in \mathbb{R}^{n \times n}$  be invertible and symmetric. A symmetric matrix is a square matrix satisfying  $\mathbf{B} = \mathbf{B}^\top$ . Show that  $\mathbf{B}^{-1}$  is also symmetric.
- [2 points]** Let  $\mathbf{C} \in \mathbb{R}^{n \times n}$  be positive semi-definite (PSD). A positive semi-definite matrix is a symmetric matrix satisfying  $x^\top \mathbf{C}x \geq 0$  for any vector  $x \in \mathbb{R}^n$ . Show that its eigenvalues are non-negative.

**Solution:**

- $g(x) = 1/x$ . This is because  $\text{diag}(v)^{-1} = \text{diag}(1/v_1, \dots, 1/v_n)$ .
- Just compute

$$\begin{aligned} \|\mathbf{A}x\|_2^2 &= \sum_{i=1}^n (\mathbf{A}x)_i^2 = \sum_{i=1}^n \left( \sum_{j=1}^n \mathbf{A}_{i,j}x_j \right)^2 = \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \mathbf{A}_{i,j} \mathbf{A}_{i,k} x_j x_k \\ &= \sum_{j=1}^n \sum_{k=1}^n x_j x_k \sum_{i=1}^n \mathbf{A}_{i,j} \mathbf{A}_{i,k} = \sum_{j=1}^n \sum_{k=1}^n x_j x_k (\mathbf{A}^\top \mathbf{A})_{j,k} = \sum_{j=1}^n x_j x_j = \|x\|_2^2. \end{aligned}$$

The second to last equality holds because  $\mathbf{A}$  is orthonormal.

c. Taking transpose of

$$\mathbf{I} = \mathbf{B}\mathbf{B}^{-1},$$

we have

$$\mathbf{I} = (\mathbf{B}^{-1})^\top \mathbf{B}^\top = (\mathbf{B}^{-1})^\top \mathbf{B}.$$

So  $(\mathbf{B}^{-1})^\top$  is an inverse of  $\mathbf{B}$ . Since inverses are unique, this means that  $(\mathbf{B}^{-1})^\top$  is the inverse of  $\mathbf{B}$ , i.e.  $(\mathbf{B}^{-1})^\top = \mathbf{B}^{-1}$ .

d. Let  $x$  be an eigenvector of  $\mathbf{C}$  with eigenvalue  $\lambda$ . Then

$$0 \leq x^\top \mathbf{C}x = x^\top \lambda x = \lambda \|x\|_2^2.$$

Since  $\|x\|_2^2 \geq 0$ , this implies that  $\lambda \geq 0$  as well.

# Programming

**These problems are available in a .zip file, with some starter code. All coding questions in this class will have starter code. Before attempting these problems, you will need to set up a Conda environment that you will use for every assignment in the course. Unzip the HW0-A.zip file and read the instructions in the README file to get started.**

A9. For  $\nabla_x f(x, y)$  as solved for in Problem 7:

- [1 point]** Using native Python, implement the summation form.
- [1 point]** Using NumPy, implement the vector form.
- [1 point]** Report the difference in wall-clock time for parts a-b, and discuss reasons for the observed difference.

## What to Submit:

- **Part c:** Difference in wall-clock time for parts a-b
- **Part c:** Explanation for difference (1-2 sentences)
- **Code** on Gradescope through coding submission

A10. Two random variables  $X$  and  $Y$  have equal distributions if their CDFs,  $F_X$  and  $F_Y$ , respectively, are equal, i.e. for all  $x$ ,  $|F_X(x) - F_Y(x)| = 0$ . The central limit theorem says that the sum of  $k$  independent, zero-mean, variance  $1/k$  random variables converges to a (standard) Normal distribution as  $k$  tends to infinity. We will study this phenomenon empirically (you will use the Python packages Numpy and Matplotlib). Each of the following subproblems includes a description of how the plots were generated; these have been coded for you. The code is available in the .zip file. In this problem, you will add to our implementation to explore **matplotlib** library, and how the solution depends on  $n$  and  $k$ .

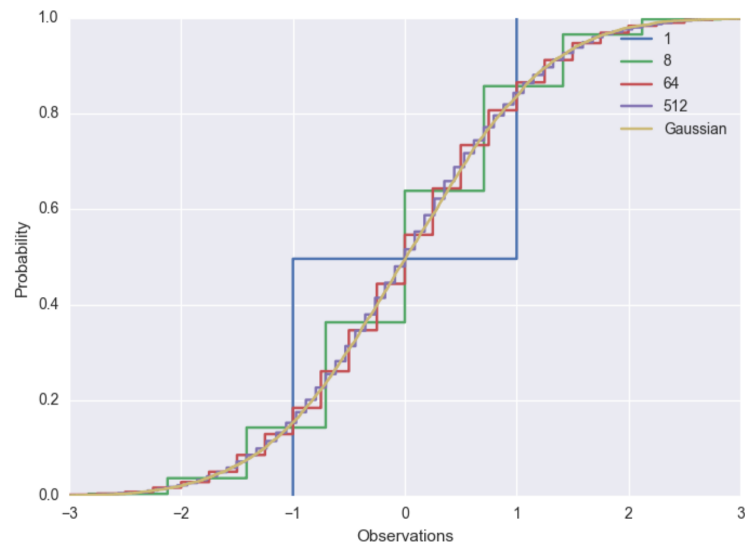
- [2 points]** For  $i = 1, \dots, n$  let  $Z_i \sim \mathcal{N}(0, 1)$ . Let  $F(x)$  denote the true CDF from which each  $Z_i$  is drawn (i.e., Gaussian). Define  $\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{Z_i \leq x\}$  and we will choose  $n$  large enough such that, for all  $x \in \mathbb{R}$ ,

$$\sqrt{\mathbb{E} \left[ \left( \hat{F}_n(x) - F(x) \right)^2 \right]} \leq 0.0025 .$$

Plot  $\hat{F}_n(x)$  from  $-3$  to  $3$ .

- [2 points]** Define  $Y^{(k)} = \frac{1}{\sqrt{k}} \sum_{i=1}^k B_i$  where each  $B_i$  is equal to  $-1$  and  $1$  with equal probability and the  $B_i$ 's are independent. We know that  $\frac{1}{\sqrt{k}} B_i$  is zero-mean and has variance  $1/k$ . For each  $k \in \{1, 8, 64, 512\}$  we will generate  $n$  (same as in part a) independent copies  $Y^{(k)}$  and plot their empirical CDF on the same plot as part a.

Be sure to always label your axes. Your plot should look something like the following (up to styling) (Tip: checkout **seaborn** for instantly better looking plots.)



**What to Submit:**

- **Part a:** Value for  $n$  (Hint: You will need to print it) **Part a:** In 1-2 sentences: How does empirical CDF change with  $k$ ?
- **Parts a and b:** Plot of  $\hat{F}_n(x) \in [-3, 3]$
- **Code** on Gradescope through coding submission