

DS-GA 1003 Machine Learning: Homework 0  
Due 11.59 p.m. EST, February 13, 2024 on Gradescope

(fill in your name here)

We encourage **L<sup>A</sup>T<sub>E</sub>X**-typeset submissions but will accept quality scans of hand-written pages.

## 1 Probability Distributions

For parts (A) to (C), suppose that  $X, Y, Z$  have joint density  $p(X, Y, Z) = p(X)p(Y|X)p(Z|X)$ .

- (A) Use law of total probability to write down the marginal distribution of  $Y$  in terms of  $p(X), p(Y|X), p(Z|X)$ .

*Solution.* Write your solution for each question using the pre-defined **solution** environment. Feel free to use style packages to your convenience, e.g. highlighting parts of your solution that you still need to work on. □

- (B) Use Bayes rule to write down the conditional distribution  $Z|Y$  in terms of  $p(X), p(Y|X), p(Z|X)$ .
- (C) Without further assumptions, which variables are independent? Which are conditionally independent?
- (D) Prove  $\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|X]]$ .

*Note:* Since the inner expectation is a function of only  $X$ , we will generally omit the subscript on the outer expectation since it has to correspond to  $X$ .

- (E) Construct a random variable  $X$ , such that  $\mathbb{P}(X < \infty) = 1$ , but  $\mathbb{E}[X] = \infty$ . Show both properties.
- (F) Construct two continuous random variables  $X, Y$  and a non-constant function  $f$  such that  $f(X, Y)$  is independent of  $X$  and  $f(X, Y)$  is independent of  $Y$ . If impossible, explain why.

## 2 Gradients

- (A) Let  $\mathbf{x} \in \mathbb{R}^2$  be a 2 dimensional real vector where  $\mathbf{x} = [x_1, x_2]$ . Define the scalar-valued function  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  by:

$$f(\mathbf{x}) = \exp[\log(x_1^2) + x_1 x_2]$$

Compute  $\nabla_x f$ .

- (B) Let  $Y \sim \text{Exp}(\lambda)$ , which is the Exponential distribution with parameter  $\lambda$ . Let  $f(y; \lambda)$  denote the evaluation of the PDF at the value  $Y = y$ . Use (univariate) calculus to maximize  $f(2; \lambda)$  with respect to  $\lambda$ . *(We suggest maximizing the log of the density.)*
- (C) The CDF of the Exponential distribution is

$$F(y; \lambda) = 1 - \exp[-\lambda y]$$

Derive the PDF  $f(y; \lambda)$  from  $F(y; \lambda)$ .

### 3 The Gaussian Distribution

In the section below, we use “Gaussian” and “Normal” interchangeably. The univariate Gaussian  $\mathcal{N}(\mu, \sigma^2)$  with mean  $\mu$  and variance  $\sigma^2 > 0$  has PDF

$$p(X = x) = \frac{1}{\sigma\sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{x - \mu}{\sigma} \right)^2 \right]$$

- (A) Given a sample  $X \sim \mathcal{N}(0, 1)$ , specify a function  $f$  (not relying on any other random variables) such that  $f(X) \sim \mathcal{N}(3, 2)$ .
- (B) Given a sample  $X \sim \mathcal{N}(0, 1)$ , name a random variable  $Y$  such that  $X + Y \sim \mathcal{N}(3, 2)$ .
- (C) Let  $\mu$  be a  $D$  dimensional real vector. Let  $\Sigma$  be a  $D \times D$  positive semi-definite matrix. The multivariate Gaussian PDF in  $D$  dimensions with mean  $\mu$  and covariance  $\Sigma$  is:

$$p(X = x) = \det(2\pi\Sigma)^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right]$$

The marginals of each dimension are normal with  $X_i \sim \mathcal{N}(\mu_i, \Sigma_{ii})$ . The 2D case is called the Bivariate Normal. Let  $X = [X_1, X_2]$  be Bivariate Normal  $\mathcal{N}(\mu, \Sigma)$  with

$$\mu = [\mu_1, \mu_2], \quad \Sigma = \begin{bmatrix} \sigma_1^2 & c \\ c & \sigma_2^2 \end{bmatrix}$$

such that  $\Sigma$  is positive semi-definite. Letting  $\rho = \frac{c}{\sigma_1\sigma_2}$ , the 2D case can be written as  $p(X_1 = x_1, X_2 = x_2) =$

$$\frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left[ -\frac{1}{2(1-\rho^2)} \left[ \left( \frac{x_1 - \mu_1}{\sigma_1} \right)^2 - 2\rho \left( \frac{x_1 - \mu_1}{\sigma_1} \right) \left( \frac{x_2 - \mu_2}{\sigma_2} \right) + \left( \frac{x_2 - \mu_2}{\sigma_2} \right)^2 \right] \right]$$

Compute the conditional density  $p(X_1 = x_1 | X_2 = x_2)$ .

*Hint:* Using either form for the 2D Normal PDF, start with Bayes rule and remember that the marginals are Gaussian with  $X_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$ . You may also use the fact that conditionals of Gaussians are Gaussian. Since Gaussians are fully specified by their mean and variance, this means you only need to identify the mean and variance of  $p(X_1 | X_2 = x_2)$ .

- (D) Construct a pair of variables  $X, Y$  that have  $\text{Cov}(X, Y) = 0$  but  $X$  is not independent of  $Y$ . Is this possible if  $X, Y$  are jointly Gaussian? Why or why not?

## 4 Monte Carlo Estimators

Let  $X \sim D$  be a random variable and denote  $\mu = \mathbb{E}_{X \sim D}[X]$  and  $\sigma^2 = \mathbf{Var}_{X \sim D}[X]$  as its mean and variance respectively. Assume that  $X$  has finite variance, i.e.  $\sigma^2 < \infty$ . While you do not know  $\mu$  or  $\sigma^2$ , you can collect  $N$  independent samples of  $X$ , which we denote as  $\{X_i\}_{i=1}^N$ .

(A) Is the mean  $\mu$  finite? If yes, why? If not, construct an example of such a random variable  $X$ .

(B) From your  $N$  samples, you can construct a **Monte Carlo estimator** of  $\mu$  as:

$$\hat{\mu}_N = \frac{1}{N} \sum_{i=1}^N X_i$$

Find the mean and variance of  $\hat{\mu}_N$ .

(C) Based on your answer in (B), name a potential advantage and disadvantage of using  $\hat{\mu}_N$  to estimate  $\mu$ .

(D) Assuming that  $\sigma^2 < \infty$  and  $\mu < \infty$ , then prove for any  $k > 0$  the following inequality:

$$\mathbb{P}(|X - \mu| > k) \leq \frac{\sigma^2}{k^2}$$

(E) Using parts (B) and (D), prove for any  $k > 0$  that:

$$\lim_{N \rightarrow \infty} \mathbb{P}(|\hat{\mu}_N - \mu| > k) = 0$$

(F) In your own words, why is the result in (E) useful?

## 5 Kullback-Liebler Divergence

One way to measure the similarity between two distributions  $P, Q$  is the **KL divergence**, which is defined using their densities  $p, q$  as:

$$KL(P||Q) = \int_{x \in \mathbb{R}} p(x) \log \frac{p(x)}{q(x)} dx$$

The KL is non-negative and is 0 if and only if the two distributions are equal. These properties also hold when  $P, Q$  are discrete.

Assume that the densities  $p(x), q(x) > 0$  for all  $x \in \mathbb{R}$ . Prove the following two statements:

- when  $P = Q$ ,  $KL(P||Q) = 0$ .
- when  $P \neq Q$ ,  $KL(P||Q) > 0$  (strict inequality).

*Hint:* Use Jensen's inequality, which states that given a strictly-convex function  $f$  and a (non-constant) random variable  $X$ :

$$f(\mathbb{E}(X)) < \mathbb{E}(f(X))$$

## 6 Setting Up PyTorch

This question is mostly to get you to install PyTorch, one of the two popular machine learning libraries for python (the other being Tensorflow), and to start writing a few lines of sampling code. It should be easy to get started by choosing your system settings on this page <https://pytorch.org/get-started/locally/>. The non-GPU version for your regular laptop is fine for our purposes.

Assuming you have installed the library you should be able to `import torch`. We expect you are familiar with basic usage of Numpy, where `np.array` is the main data structure. In Torch, the equivalent is a `torch.tensor`:

- `x=torch.tensor([[1.0,2.0],[3.0,4.0]])` is a  $2 \times 2$  matrix. You can verify the shape by using `x.shape`.
- `tensors` have lots of convenient methods. Try `x.sum()`, `x.sum(0)`, `x.sum(1)`, `x.mean(0)`, `x.std()`, `x.abs()`, `x.pow(2)` etc... See <https://pytorch.org/docs/stable/index.html> for more.

For this homework question, we want you to teach yourself how to do the following in PyTorch:

1. Draw  $N$  univariate normal samples  $x_i \sim \mathcal{N}(0, \sigma^2)$  for some value of  $\sigma^2$ . For this you will need

`torch.distributions.Normal`

Be sure to give the right arguments (e.g. standard deviation and not variance). Compute the square of each sample and record the average of these squares  $\hat{\mu}_N = \frac{1}{N} \sum_i x_i^2$ .

2. Let's call the estimate  $\hat{\mu}_N$  we obtain in Step 1 as a single "trial". Now perform  $T$  trials for a fixed choice of  $N$ . Denote the mean produced by trial  $t$  as  $\hat{\mu}_{N,t}$  for  $t \in \{1, \dots, T\}$ . Now, compute the mean and standard deviation across trials of  $\hat{\mu}_{N,t}$ . For example, for the mean, you would compute  $\frac{1}{T} \sum_t \hat{\mu}_{N,t}$ .

Now that you can code these two steps:

- (A) Set  $T = 100$  and  $\sigma^2 = 10$ . Perform Steps 1 and 2 for each value of  $N \in \{1, 10, 50, 100, 200, 500, 1000\}$ . Plot the means and variances on a single graph each, i.e. you should have two graphs, one for the means and one for the variances, where the  $x$ -axis is  $\log N$ .
- (B) What do you observe about the mean and variances as  $N$  increases? How do these trends relate to your answers in Question 4?