

Homework 0: Preliminary

Introduction

There is a mathematical component and a programming component to this homework. Please submit your PDF and Python files to Canvas, and push all of your work to your GitHub repository. If a question requires you to make any plots, please include those in the writeup.

This assignment is intended to ensure that you have the background required for CS281, and have studied the mathematical review notes provided in section. You should be able to answer the problems below *without* complicated calculations. All questions are worth $70/6 = 11.\bar{6}$ points unless stated otherwise.

Variance and Covariance

Problem 1

Let X and Y be two independent random variables.

- (a) Show that the independence of X and Y implies that their covariance is zero.
- (b) Zero covariance *does not* imply independence between two random variables. Give an example of this.
- (c) For a scalar constant a , show the following two properties:

$$\begin{aligned}\mathbb{E}(X + aY) &= \mathbb{E}(X) + a\mathbb{E}(Y) \\ \text{var}(X + aY) &= \text{var}(X) + a^2\text{var}(Y)\end{aligned}$$

(a) First we'll show that $X \perp\!\!\!\perp Y \Rightarrow E(XY) = E(X)E(Y)$, assuming X and Y are discrete random variables.

Using the definition of expectation, we get that:

$$\begin{aligned}E(XY) &= \sum_x \sum_y xyP(X = x, Y = y) \\ &= \sum_x \sum_y xyP(X = x)P(Y = y) \\ &= \sum_x xP(X = x) \sum_y yP(Y = y) \\ &= E(X)E(Y)\end{aligned}$$

The same can be shown for independent, continuous random variables. Now we can show (a) easily:

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

(b) We'll use the properties of expectation, variance, covariance, and Bernoulli random variables in the following example.

Let $X \sim \text{Bern}(p)$, $-Y \sim \text{Bern}(p)$

$$\begin{aligned}\text{Cov}(X + Y, X - Y) &= \text{Cov}(X, X) - \text{Cov}(X, Y) + \text{Cov}(Y, X) - \text{Cov}(Y, Y) \\ &= \text{Cov}(X, X) - \text{Cov}(Y, Y) \\ &= \text{Var}(X) - \text{Var}(Y) = 0\end{aligned}$$

$$E(X - Y) = E(X) + E(-Y) = 2p \neq E(X - Y | X + Y = 1) = 1$$

Thus, we have $\text{Cov}(X + Y, X - Y) = 0$ but $X + Y$ and $X - Y$ not independent since

$$E(X - Y) \neq E(X - Y | X + Y = 1)$$

(c) For proof the first property, we'll assume X and Y are discrete random variables and use the definition of expectation as well as properties of independent random variables.

$$\begin{aligned}
E(X) + aE(Y) &= \sum_x xP(X=x) + a \sum_y yP(Y=y) \\
&= \sum_x xP(X=x) \sum_y P(Y=y) + \sum_y ayP(Y=y) \sum_x P(X=x) \\
&= \sum_x \sum_y xP(X=x)P(Y=y) + \sum_x \sum_y ayP(Y=y)P(X=x) \\
&= \sum_x \sum_y (x+ay)P(X=x)P(Y=y) \\
&= \sum_x \sum_y (x+ay)P(X=x, Y=y) = E(X+aY)
\end{aligned}$$

For proof of the second property, we'll use properties of independent random variables, expectation, and variance.

$$\begin{aligned}
Var(X+aY) &= E[(X+aY)^2] - [E(X+aY)]^2 \\
&= E[X^2 + 2aXY + a^2Y^2] - [E(X) + aE(Y)]^2 \\
&= E(X^2) + 2aE(XY) + a^2E(Y^2) - [E(X)]^2 - 2aE(X)E(Y) - a^2[E(Y)]^2 \\
&= Var(X) + a^2(E(Y^2) - [E(Y)]^2) \\
&= Var(X) + a^2Var(Y)
\end{aligned}$$

Densities

Problem 2

Answer the following questions:

- (a) Can a probability density function (pdf) ever take values greater than 1?
- (b) Let X be a univariate normally distributed random variable with mean 0 and variance $1/100$. What is the pdf of X ?
- (c) What is the value of this pdf at 0?
- (d) What is the probability that $X = 0$?
- (e) Explain the discrepancy.

(a) Yes. For example, let $X \sim Unif\left(0, \frac{1}{10}\right)$. Then $p(x) = 10 \forall x \in \left[0, \frac{1}{10}\right]$. For continuous random variables, the integral over the support needs to be 1. For discrete random variables, the sum over the support needs to be 1.

(b) $X \sim N\left(0, \frac{1}{100}\right) \Rightarrow p(x) = \frac{10}{\sqrt{2\pi}} \exp(-50x^2)$

(c) $p(0) = \frac{10}{\sqrt{2\pi}}$

(d) $P(X = 0) = 0$

(e) The value of the PDF for a random variable X is not the probability that X takes on a particular value x . The value of the PDF at some x that X takes on is the height of the distribution at x . The probability of some event associated with the value of X is the area under the PDF that represents that event. The area under the PDF that represents any one particular value of x that X takes on is 0. This explains why $p(0) \neq P(X = 0) = 0$.

Conditioning and Bayes' rule

Problem 3

Let $\mu \in \mathbb{R}^m$ and $\Sigma, \Sigma' \in \mathbb{R}^{m \times m}$. Let X be an m -dimensional random vector with $X \sim \mathcal{N}(\mu, \Sigma)$, and let Y be a m -dimensional random vector such that $Y|X \sim \mathcal{N}(X, \Sigma')$. Derive the distribution and parameters for each of the following.

- (a) The unconditional distribution of Y .
- (b) The joint distribution for the pair (X, Y) .

Hints:

- You may use without proof (but they are good advanced exercises) the closure properties of multivariate normal distributions. Why is it helpful to know when a distribution is normal?
- Review Eve's and Adam's Laws, linearity properties of expectation and variance, and Law of Total Covariance.

(a) We can see that $Y - X|X \sim N(0, \Sigma')$. This distribution is not dependent on $X \Rightarrow Y - X \sim N(0, \Sigma')$. Since $(Y - X) \perp\!\!\!\perp X$, the sum of independent multivariate normal random vectors is a multivariate normal random vector, and $Y - X + X = Y$, it follows that $Y \sim N(\mu, \Sigma + \Sigma')$.

(b) We are trying to find the joint distribution of $(X, Y)^T$. Since X is MVN and $Y|X$ is MVN, $(X, Y)^T \sim N(\theta, \beta)$ such that $\theta \in \mathbb{R}^{2m}$ and $\beta \in \mathbb{R}^{2m \times 2m}$.

According to Theorem 4.3.1 in Murphy's text, $\theta = (E(X), E(Y))^T = (\mu, \mu)^T$

Also according to Theorem 4.3.1 in Murphy's text, $\beta = \begin{pmatrix} Var(X) & Cov(X, Y) \\ Cov(Y, X) & Var(Y) \end{pmatrix}$.

We know that $Var(X) = \Sigma$ and that $Var(Y) = \Sigma + \Sigma'$, so all we need to do is find $Cov(X, Y)$ and $Cov(Y, X)$.

$$Cov(X, Y) = E(Cov(X, Y)|X) + Cov(E(X|X), E(Y|X)) = 0 + Cov(X, X) = Var(X) = \Sigma$$

$$Cov(Y, X) = E(Cov(Y, X)|X) + Cov(E(Y|X), E(X|X)) = 0 + Cov(X, X) = Var(X) = \Sigma$$

$$\text{Then, } \beta = \begin{pmatrix} \Sigma & \Sigma \\ \Sigma & \Sigma + \Sigma' \end{pmatrix}.$$

$$\text{Therefore, } (X, Y)^T \sim N((\mu, \mu)^T, \begin{pmatrix} \Sigma & \Sigma \\ \Sigma & \Sigma + \Sigma' \end{pmatrix}).$$

I can Ei-gen

Problem 4

Let $\mathbf{X} \in \mathbb{R}^{n \times m}$.

- (a) What is the relationship between the n eigenvalues of $\mathbf{X}\mathbf{X}^T$ and the m eigenvalues of $\mathbf{X}^T\mathbf{X}$?
- (b) Suppose \mathbf{X} is square (i.e., $n = m$) and symmetric. What does this tell you about the eigenvalues of \mathbf{X} ? What are the eigenvalues of $\mathbf{X} + \mathbf{I}$, where \mathbf{I} is the identity matrix?
- (c) Suppose \mathbf{X} is square, symmetric, and invertible. What are the eigenvalues of \mathbf{X}^{-1} ?

Hints:

- Make use of singular value decomposition and the properties of orthogonal matrices. Show your work.
- Review and make use of (but do not derive) the spectral theorem.

(a) Suppose λ is a non-zero eigenvalue of $\mathbf{X}^T\mathbf{X}$ and is associated with eigenvector $\mathbf{a} \in \mathbb{R}^m$. Then $\mathbf{X}^T\mathbf{X}\mathbf{a} = \lambda\mathbf{a} \Rightarrow (\mathbf{X}\mathbf{X}^T)\mathbf{X}\mathbf{a} = \lambda\mathbf{X}\mathbf{a} \Rightarrow \mathbf{X}\mathbf{X}^T$ also has λ as an eigenvalue and it is associated with eigenvector $\mathbf{X}\mathbf{a}$. Since the choice of λ was arbitrary, it follows that the non-zero eigenvalues of $\mathbf{X}^T\mathbf{X}$ and $\mathbf{X}\mathbf{X}^T$ are the same.

It follows that in the case that $n = m$, $\mathbf{X}^T\mathbf{X}$ and $\mathbf{X}\mathbf{X}^T$ will have the same eigenvalues.

If $n \neq m$, the larger of $\mathbf{X}^T\mathbf{X}$ and $\mathbf{X}\mathbf{X}^T$ will have more eigenvalues with value 0. This is because $\text{rank}(\mathbf{X}^T\mathbf{X}) = \text{rank}(\mathbf{X}\mathbf{X}^T) \leq \min(n, m)$, and these ranks correspond to the maximum number of non-zero eigenvalues that $\mathbf{X}^T\mathbf{X}$ and $\mathbf{X}\mathbf{X}^T$ can have $\Rightarrow \mathbf{X}^T\mathbf{X}$ has $m - n$ more eigenvalues with value 0 than $\mathbf{X}\mathbf{X}^T$ has if $m > n$ or $n - m$ less eigenvalues with value 0 than $\mathbf{X}\mathbf{X}^T$ has if $m < n$.

(b) The spectral theorem says that if \mathbf{X} is a square, symmetric matrix, all the eigenvalues of \mathbf{X} are real.

Let λ be an eigenvalue of \mathbf{X} associated with eigenvector \mathbf{a} . For the eigenvalues of $\mathbf{X} + \mathbf{I}$, we start with $\mathbf{X}\mathbf{a} = \lambda\mathbf{a} \Rightarrow (\mathbf{X} + \mathbf{I})\mathbf{a} = (\lambda\mathbf{I} + \mathbf{I})\mathbf{a} \Rightarrow (\mathbf{X} + \mathbf{I})\mathbf{a} = (\lambda + 1)\mathbf{a}$. Since the choice of λ was arbitrary, it follows the eigenvalues of $\mathbf{X} + \mathbf{I}$ can be found by adding 1 to each of the eigenvalues of \mathbf{X} .

(c) Suppose λ is an eigenvalue of $\mathbf{X} \in \mathbb{R}^{n \times n}$ (λ must be non-zero since \mathbf{X} is invertible) and is associated with eigenvector $\mathbf{a} \in \mathbb{R}^n$. Then $\mathbf{X}\mathbf{a} = \lambda\mathbf{a} \Rightarrow \mathbf{X}^{-1}\mathbf{X}\mathbf{a} = \lambda\mathbf{X}^{-1}\mathbf{a} \Rightarrow \mathbf{a} = \lambda\mathbf{X}^{-1}\mathbf{a} \Rightarrow \mathbf{X}^{-1}\mathbf{a} = \frac{1}{\lambda}\mathbf{a} \Rightarrow \mathbf{X}^{-1}$ has $\frac{1}{\lambda}$ as an eigenvalue and it is associated with eigenvector \mathbf{a} . Since the choice of λ was arbitrary, it follows the eigenvalues of \mathbf{X} and \mathbf{X}^{-1} are reciprocals of each other.

Vector Calculus

Problem 5

Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^m$ and $\mathbf{A} \in \mathbb{R}^{m \times m}$. Please derive from elementary scalar calculus the following useful properties. Write your final answers in vector notation.

- (a) What is the gradient with respect to \mathbf{x} of $\mathbf{x}^T \mathbf{y}$?
- (b) What is the gradient with respect to \mathbf{x} of $\mathbf{x}^T \mathbf{x}$?
- (c) What is the gradient with respect to \mathbf{x} of $\mathbf{x}^T \mathbf{A} \mathbf{x}$?

(a)

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{y}) = \frac{\partial}{\partial \mathbf{x}} \left(\sum_{i=1}^m x_i y_i \right) = \sum_{i=1}^m \frac{\partial x_i y_i}{\partial \mathbf{x}} = \sum_{i=1}^m \left[\frac{\partial x_i y_i}{\partial x_1}, \frac{\partial x_i y_i}{\partial x_2} \dots \frac{\partial x_i y_i}{\partial x_m} \right]^T = [y_1, y_2 \dots y_m]^T = \mathbf{y}$$

(b)

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} \left(\sum_{i=1}^m x_i^2 \right) = \sum_{i=1}^m \frac{\partial x_i^2}{\partial \mathbf{x}} = \sum_{i=1}^m \left[\frac{\partial x_i^2}{\partial x_1}, \frac{\partial x_i^2}{\partial x_2} \dots \frac{\partial x_i^2}{\partial x_m} \right]^T = [2x_1, 2x_2, \dots, 2x_m]^T = 2\mathbf{x}$$

(c)

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{A} \mathbf{x}) &= \frac{\partial}{\partial \mathbf{x}} \left(x_1 \sum_{i=1}^m x_i a_{i1} + x_2 \sum_{i=1}^m x_i a_{i2} + \dots + x_m \sum_{i=1}^m x_i a_{im} \right) \\ &= \sum_{i=1}^m \frac{\partial x_1 x_i a_{i1}}{\partial \mathbf{x}} + \dots + \sum_{i=1}^m \frac{\partial x_m x_i a_{im}}{\partial \mathbf{x}} \\ &= \sum_{i=1}^m \left[\frac{\partial x_1 x_i a_{i1}}{\partial x_1} \dots \frac{\partial x_1 x_i a_{i1}}{\partial x_m} \right]^T + \dots + \sum_{i=1}^m \left[\frac{\partial x_m x_i a_{im}}{\partial x_1} \dots \frac{\partial x_m x_i a_{im}}{\partial x_m} \right]^T \\ &= \left[2x_1 a_{11} + \sum_{i=2}^m x_i a_{i1}, \dots, x_1 a_{m1} \right]^T + \dots + \left[x_m a_{1m}, \dots, 2x_m a_{mm} + \sum_{i=1}^{m-1} x_i a_{im} \right]^T = (\mathbf{A} + \mathbf{A}^T) \mathbf{x} \end{aligned}$$

We can see that this is true by doing a manual calculation using a small value for m .

Gradient Check

Problem 6

Often after finishing an analytic derivation of a gradient, you will need to implement it in code. However, there may be mistakes - either in the derivation or in the implementation. This is particularly the case for gradients of multivariate functions.

One way to check your work is to numerically estimate the gradient and check it on a variety of inputs. For this problem we consider the simplest case of a univariate function and its derivative. For example, consider a function $f(x) : \mathbb{R} \rightarrow \mathbb{R}$:

$$\frac{df}{dx} = \lim_{\epsilon \rightarrow 0} \frac{f(x + \epsilon) - f(x - \epsilon)}{2\epsilon}$$

A common check is to evaluate the right-hand side for a small value of ϵ , and check that the result is similar to your analytic result.

In this problem, you will implement the analytic and numerical derivatives of the function

$$f(x) = \cos(x) + x^2 + e^x.$$

1. Implement `f` in Python (feel free to use whatever `numpy` or `scipy` functions you need):

```
def f(x):
```

2. Analytically derive the derivative of that function, and implement it in Python:

```
def grad_f(x):
```

3. Now, implement a gradient check (the numerical approximation to the derivative), and by plotting, show that the numerical approximation approaches the analytic as `epsilon` $\rightarrow 0$ for a few values of x :

```
def grad_check(x, epsilon):
```

$$\frac{df}{dx} = -\sin(x) + 2x + e^x$$

The plots on the following page show that the numerical approximations of the gradient approach the analytic values of the gradient as $\epsilon \rightarrow 0$ for four values of x .

