

Assignment 1 - Probability, Linear Algebra, & Computational Programming

Jake Bova

Netid: jb240893

Names of students you worked with on this assignment: Bri, Brock, Joel

Learning Objectives

The purpose of this assignment is to provide a refresher on fundamental concepts that we will use throughout this course and provide an opportunity to develop skills in any of the related skills that may be unfamiliar to you. Through the course of completing this assignment, you will...

- Refresh your knowledge of probability theory including properties of random variables, probability density functions, cumulative distribution functions, and key statistics such as mean and variance.
- Revisit common linear algebra and matrix operations and concepts such as matrix multiplication, inner and outer products, inverses, the Hadamard (element-wise) product, eigenvalues and eigenvectors, orthogonality, and symmetry.
- Practice numerical programming, core to machine learning, by loading and filtering data, plotting data, vectorizing operations, profiling code speed, and debugging and optimizing performance. You will also practice computing probabilities based on simulation.
- Develop or refresh your knowledge of Git version control, which will be a core tool used in the final project of this course
- Apply your skills altogether through an exploratory data analysis to practice data cleaning, data manipulation, interpretation, and communication

We will build on these concepts throughout the course, so use this assignment as a catalyst to deepen your knowledge and seek help with anything unfamiliar.

If some references would be helpful on these topics, I would recommend the following resources:

- [Mathematics for Machine Learning](#) by Deisenroth, Faisal, and Ong
- [Deep Learning](#); Part I: Applied Math and Machine Learning Basics by Goodfellow, Bengio, and Courville

- [The Matrix Calculus You Need For Deep Learning](#) by Parr and Howard
- [Dive Into Deep Learning](#); Appendix: Mathematics for Deep Learning by Weness, Hu, et al.

Note: don't worry if you don't understand everything in the references above - some of these books dive into significant minutia of each of these topics.

Probability and Statistics Theory

Note: for all assignments, write out equations and math using markdown and [LaTeX](#). For this assignment show ALL math work for questions 1-4, meaning that you should include any intermediate steps necessary to understand the logic of your solution

1

[3 points]

Let $f(x) = \begin{cases} 0 & x < 0 \\ \alpha x^2 & 0 \leq x \leq 2 \\ 0 & 2 < x \end{cases}$ For what value of α is $f(x)$ a valid probability density function?

A probability density function must integrate to 1, so we can solve for α by setting the integral of $f(x)$ equal to 1 and solving for α :

$$\int_{-\infty}^{\infty} f(x) dx = 1 \quad (1)$$

$$\int_{-\infty}^0 0 dx + \int_0^2 \alpha x^2 dx + \int_2^{\infty} 0 dx = 1 \quad (2)$$

$$\frac{\alpha}{3} x^3 \Big|_0^2 = 1 \quad (3)$$

$$\frac{\alpha}{3} (2^3 - 0^3) = 1 \quad (4)$$

$$\frac{8}{3} \alpha = 1 \quad (5)$$

$$\alpha = \frac{3}{8} \quad (6)$$

2

[3 points] What is the cumulative distribution function (CDF) that corresponds to the following probability distribution function? Please state the value of the CDF for all possible values of x .

$$f(x) = \begin{cases} \frac{1}{3} & 0 < x < 3 \\ 0 & \text{otherwise} \end{cases}$$

ANSWER

The CDF is the integral of the PDF, so we can solve for the CDF by integrating $f(x)$:

$$F(x) = \int_{-\infty}^x f(x)dx \quad (7)$$

$$= \int_{-\infty}^0 0dx + \int_0^x \frac{1}{3}dx + \int_x^3 0dx \quad (8)$$

$$= \frac{1}{3}x \Big|_0^x \quad (9)$$

$$= \frac{1}{3}x \quad (10)$$

The CDF is: $f(x) = \begin{cases} \frac{1}{3}x & 0 < x < 3 \\ 0 & x \leq 0 \\ 1 & x \geq 3 \end{cases}$

3

[6 points] For the probability distribution function for the random variable X ,

$$f(x) = \begin{cases} \frac{1}{3} & 0 < x < 3 \\ 0 & \text{otherwise} \end{cases}$$

what is the (a) expected value and (b) variance of X . *Show all work.*

ANSWER

The expected value of X is the mean of $f(x)$:

$$E[X] = \int_{-\infty}^{\infty} x f(x) dx \quad (11)$$

$$= \int_{-\infty}^0 0 dx + \int_0^3 \frac{1}{3} x dx + \int_3^{\infty} 0 dx \quad (12)$$

$$= \frac{1}{3} \frac{x^2}{2} \Big|_0^3 \quad (13)$$

$$= \frac{1}{3} \frac{3^2}{2} \quad (14)$$

$$= \frac{3}{2} \quad (15)$$

The variance of X :

$$Var[X] = E[X^2] - E[X]^2 \quad (16)$$

$$= \int_{-\infty}^{\infty} x^2 f(x) dx - \left(\frac{3}{2}\right)^2 \quad (17)$$

$$= \int_{-\infty}^0 0 dx + \int_0^3 \frac{1}{3} x^2 dx + \int_3^{\infty} 0 dx - \left(\frac{3}{2}\right)^2 \quad (18)$$

$$= \frac{1}{3} \frac{x^3}{3} \Big|_0^3 - \left(\frac{3}{2}\right)^2 \quad (19)$$

$$= \frac{1}{3} \frac{3^3}{3} - \left(\frac{3}{2}\right)^2 \quad (20)$$

$$= 3 - \frac{9}{4} \quad (21)$$

$$= \frac{3}{4} \quad (22)$$

4

[6 points] Consider the following table of data that provides the values of a discrete data vector \mathbf{x} of samples from the random variable X , where each entry in \mathbf{x} is given as x_i .

Table 1. Dataset $N=5$ observations

	x_0	x_1	x_2	x_3	x_4
\mathbf{x}	2	3	10	-1	-1

What is the (a) mean and (b) variance of the data?

Show all work. Your answer should include the definition of mean and variance in the context of discrete data. In this case, use the sample variance since the sample size is quite small

ANSWER

Mean:

$$\bar{x} = \frac{1}{N} \sum_{i=0}^{N-1} x_i \quad (23)$$

$$= \frac{1}{5} (2 + 3 + 10 + -1 + -1) \quad (24)$$

$$= \frac{1}{5} 13 \quad (25)$$

$$= 2.6 \quad (26)$$

Variance:

$$\sigma^2 = \frac{1}{N-1} \sum_{i=0}^N (x_i - \bar{x})^2 \quad (27)$$

$$= \frac{1}{4} ((2 - 2.6)^2 + (3 - 2.6)^2 + (10 - 2.6)^2 + (-1 - 2.6)^2 + (-1 - 2.6)^2) \quad (28)$$

$$= \frac{1}{4} (0.36 + 0.16 + 54.76 + 12.96 + 12.96) \quad (29)$$

$$= \frac{1}{4} 81.2 \quad (30)$$

$$= 20.3 \quad (31)$$

Linear Algebra

5

[5 points] A common task in machine learning is a change of basis: transforming the representation of our data from one space to another. A prime example of this is through the process of dimensionality reduction as in Principle Components Analysis where we often seek to transform our data from one space (of dimension n) to a new space (of dimension m) where $m < n$. Assume we have a sample of data of dimension $n = 4$ (as shown below) and we want to transform it into a dimension of $m = 2$.

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

(a) What are the dimensions of a matrix, \mathbf{A} , that would linearly transform our sample of data, \mathbf{x} , into a space of $m = 2$ through the operation \mathbf{Ax} ?

(b) Express this transformation in terms of the components of \mathbf{x} : x_1, x_2, x_3, x_4 and the matrix \mathbf{A} where each entry in the matrix is denoted as $a_{i,j}$ (e.g. the entry in the first row and second column would be $a_{1,2}$). Your answer will be in the form of a matrix expressing result of the product \mathbf{Ax} .

Note: please write your answers here in LaTeX

ANSWER

a) The dimensions of \mathbf{A} are 2×4 . This is because the number of rows in \mathbf{A} must equal the number of columns in \mathbf{x} , and the number of columns in \mathbf{A} must equal the number of rows in the transformed \mathbf{x} .

b) The transformation is:

$$\mathbf{Ax} = \begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \quad (32)$$

$$= \begin{bmatrix} a_{1,1}x_1 + a_{1,2}x_2 + a_{1,3}x_3 + a_{1,4}x_4 \\ a_{2,1}x_1 + a_{2,2}x_2 + a_{2,3}x_3 + a_{2,4}x_4 \end{bmatrix} \quad (33)$$

6

[14 points] Matrix manipulations and multiplication. Machine learning involves working with many matrices, so this exercise will provide you with the opportunity to practice those skills.

$$\text{Let } \mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 5 \\ 3 & 5 & 6 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} -1 \\ 3 \\ 8 \end{bmatrix}, \mathbf{c} = \begin{bmatrix} 4 \\ -3 \\ 6 \end{bmatrix}, \text{ and } \mathbf{I} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Compute the following **using Python** or indicate that it cannot be computed. Refer to NumPy's tools for handling matrices. While all answers should be computer using Python, your response to whether each item can be computed should refer to underlying linear algebra. There may be circumstances when Python will produce an output, but based on the dimensions of the matrices involved, the linear algebra operation is not possible. **For the case when an operation is invalid, explain why it is not.**

When the quantity can be computed, please provide both the Python code AND the output of that code (this need not be in LaTeX)

1. \mathbf{AA}
2. \mathbf{AA}^T

3. \mathbf{Ab}
4. \mathbf{Ab}^T
5. \mathbf{bA}
6. $\mathbf{b}^T \mathbf{A}$
7. \mathbf{bb}
8. $\mathbf{b}^T \mathbf{b}$
9. \mathbf{bb}^T
10. $\mathbf{b} + \mathbf{c}^T$
11. $\mathbf{b}^T \mathbf{b}^T$
12. $\mathbf{A}^{-1} \mathbf{b}$
13. $\mathbf{A} \circ \mathbf{A}$
14. $\mathbf{b} \circ \mathbf{c}$

Note: The element-wise (or Hadamard) product is the product of each element in one matrix with the corresponding element in another matrix, and is represented by the symbol " \circ ".

ANSWER

to be a valid matrix multiplication, the number of columns in the first matrix must equal the number of rows in the second matrix. i.e. 3×3 times 3×1 , if the 'inner' numbers are the same, they are valid for matrix multiplication, and the size of the resulting matrix is the 'outer' numbers.

```
In [ ]: import numpy as np

A = np.array([[1, 2, 3], [2, 4, 5], [3, 5, 6]])
b = np.array([[ -1], [ 3], [ 8]])
c = np.array([[4], [-3], [6]])
I = np.identity(3)
bt = np.matrix.transpose(b)

# 1 AA
print("1. AA")
print(A @ A) # valid since inner dimensions match 3x3 and 3x3

# 2 AA^T
print("2. AA^T")
print(A @ A.T) # valid since inner dimensions match 3x3 and 3x3

# 3 Ab
print("3. Ab")
print(A @ b) # valid since inner dimensions match 3x3 and 3x1

# 4 Ab^T
print("4. Ab^T")
print("invalid because the inner dimensions do not match 3x3 and 1x3")

# 5 bA
print("5. bA")
print("invalid because the inner dimensions do not match 3x3 and 1x3")
```

```

# 6  $b^T A$ 
print("6.  $b^T A$ ")
print(b.T @ A) # valid since inner dimensions match 1x3 and 3x3

# 7  $bb$ 
print("7.  $bb$ ")
print("invalid because the inner dimensions do not match 1x3 and 1x3")

# 8  $b^T b$ 
print("8.  $b^T b$ ")
print(b.T @ b) # valid since inner dimensions match 1x3 and 3x1

# 9  $bb^T$ 
print("9.  $bb^T$ ")
print(b @ b.T) # valid since inner dimensions match 3x1 and 1x3

# 10  $b + c^T$ 
print("10.  $b + c^T$ ")
print("invalid because the inner dimensions do not match 1x3 and 1x3")

# 11  $b^T b^T$ 
print("11.  $b^T b^T$ ")
print("invalid because the inner dimensions do not match 1x3 and 1x3")

# 12  $A^{-1}b$ 
print("12.  $A^{-1}b$ ")
print(np.linalg.inv(A) @ b) # valid since inner dimensions match 3x3 and 3x1

# 13  $A \circ A$ 
print("13.  $A \circ A$ ")
print(A * A) # element wise multiplication valid since dimensions match 3x3

# 14  $b \circ c$ 
print("14.  $b \circ c$ ")
print(b * c) # element wise multiplication valid since dimensions match 1x3

```



```

1. AA
[[14 25 31]
 [25 45 56]
 [31 56 70]]
2. AA^T
[[14 25 31]
 [25 45 56]
 [31 56 70]]
3. Ab
[[29]
 [50]
 [60]]
4. Ab^T
invalid because the inner dimensions do not match 3x3 and 1x3
5. bA
invalid because the inner dimensions do not match 3x3 and 1x3
6. b^TA
[[29 50 60]]
7. bb
invalid because the inner dimensions do not match 1x3 and 1x3
8. b^Tb
[[74]]
9. bb^T
[[ 1 -3 -8]
 [-3  9 24]
 [-8 24 64]]
10. b + c^T
invalid because the inner dimensions do not match 1x3 and 1x3
11. b^T*b^T
invalid because the inner dimensions do not match 1x3 and 1x3
12. A^-1b
[[ 6.]
 [ 4.]
 [-5.]]
13. A o A
[[ 1  4  9]
 [ 4 16 25]
 [ 9 25 36]]
14. b o c
[[-4]
 [-9]
 [48]]

```

7

[8 points] Eigenvectors and eigenvalues. Eigenvectors and eigenvalues are useful for some machine learning algorithms, but the concepts take time to solidly grasp. They are used extensively in machine learning and in this course we will encounter them in relation to Principal Components Analysis (PCA), clustering algorithms, For an intuitive review of these concepts, explore this [interactive website at Setosa.io](#). Also, the series of linear

algebra videos by Grant Sanderson of 3Brown1Blue are excellent and can be viewed on youtube [here](#). For these questions, numpy may once again be helpful.

1. Calculate the eigenvalues and corresponding eigenvectors of matrix \mathbf{A} above, from

the last question. Let $\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 5 \\ 3 & 5 & 6 \end{bmatrix}$

2. Choose one of the eigenvector/eigenvalue pairs, \mathbf{v} and λ , and show that $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$.

This relationship extends to higher orders: $\mathbf{A}\mathbf{A}\mathbf{v} = \lambda^2\mathbf{v}$

3. Show that the eigenvectors are orthogonal to one another (e.g. their inner product is zero). This is true for eigenvectors from real, symmetric matrices. In three dimensions or less, this means that the eigenvectors are perpendicular to each other. Typically we use the orthogonal basis of our standard x, y, and z, Cartesian coordinates, which allows us, if we combine them linearly, to represent any point in a 3D space. But any three orthogonal vectors can do the same. We will see this property is used in PCA to identify the dimensions of greatest variation in our data when we discuss dimensionality reduction.

ANSWER

```
In [ ]: eigenvalues, eigenvectors = np.linalg.eig(A)
print("1. Eigenvalues of A")
print(eigenvalues)
print("1. Eigenvectors of A")
print(eigenvectors)

# 2. show that Av = lambda*v
print("2. show that Av = lambda*v")
print(A@eigenvectors[:,0])
print(eigenvalues[0]*eigenvectors[:,0])

# 3 show that the eigenvectors are orthogonal
print("3. show that the eigenvectors are orthogonal")
print(eigenvectors[:,0]@eigenvectors[:,1]) # slice the first column of the e
print(eigenvectors[:,0]@eigenvectors[:,2])
print(eigenvectors[:,1]@eigenvectors[:,2])
print("The eigenvectors are orthogonal since the dot product of each eigenve
```

```

1. Eigenvalues of A
[11.34481428 -0.51572947  0.17091519]
1. Eigenvectors of A
[[-0.32798528 -0.73697623  0.59100905]
 [-0.59100905 -0.32798528 -0.73697623]
 [-0.73697623  0.59100905  0.32798528]]
2. show that Av = lambda*v
[-3.72093206 -6.70488789 -8.36085845]
[-3.72093206 -6.70488789 -8.36085845]
3. show that the eigenvectors are orthogonal
-1.6653345369377348e-16
-3.608224830031759e-16
-5.828670879282072e-16
The eigenvectors are orthogonal since the dot product of each eigenvector is
0 (or very close to 0, python moment)

```

Numerical Programming

8

[10 points] Loading data and gathering insights from a real dataset

In data science, we often need to have a sense of the idiosyncrasies of the data, how they relate to the questions we are trying to answer, and to use that information to help us to determine what approach, such as machine learning, we may need to apply to achieve our goal. This exercise provides practice in exploring a dataset and answering question that might arise from applications related to the data.

Data. The data for this problem can be found in the `data` subfolder in the `assignments` folder on [github](#). The filename is `a1_egrid2016.xlsx`. This dataset is the Environmental Protection Agency's (EPA) [Emissions & Generation Resource Integrated Database \(eGRID\)](#) containing information about all power plants in the United States, the amount of generation they produce, what fuel they use, the location of the plant, and many more quantities. We'll be using a subset of those data.

The fields we'll be using include:

field	description
SEQPLT16	eGRID2016 Plant file sequence number (the index)
PSTATABB	Plant state abbreviation
PNAME	Plant name
LAT	Plant latitude
LON	Plant longitude

field	description
PLPRMFL	Plant primary fuel
CAPFAC	Plant capacity factor
NAMEPCAP	Plant nameplate capacity (Megawatts MW)
PLNGENAN	Plant annual net generation (Megawatt-hours MWh)
PLCO2EQA	Plant annual CO2 equivalent emissions (tons)

For more details on the data, you can refer to the [eGrid technical documents](#). For example, you may want to review page 45 and the section "Plant Primary Fuel (PLPRMFL)", which gives the full names of the fuel types including WND for wind, NG for natural gas, BIT for Bituminous coal, etc.

There also are a couple of "gotchas" to watch out for with this dataset:

- The headers are on the second row and you'll want to ignore the first row (they're more detailed descriptions of the headers).
- NaN values represent blanks in the data. These will appear regularly in real-world data, so getting experience working with these sorts of missing values will be important.

Your objective. For this dataset, your goal is to answer the following questions about electricity generation in the United States:

(a) Which plant has generated the most energy (measured in MWh)?

(b) What is the name of the northern-most power plant in the United States?

(c) What is the state where the northern-most power plant in the United States is located?

(d) Plot a bar plot showing the amount of energy produced by each fuel type across all plants.

(e) From the plot in (d), which fuel for generation produces the most energy (MWh) in the United States?

ANSWER

```
In [ ]: import pandas as pd

# get the file at data/a1_egrid2016.xlsx
plants_df = pd.read_excel('data/a1_egrid2016.xlsx', sheet_name='PLNT16')
```

```
In [ ]: plants_df.dtypes
```

```
Out[ ]: SEQPLT16      int64
        PSTATABB      object
        PNAME         object
        LAT           float64
        LON           float64
        PLPRMFL       object
        CAPFAC        float64
        NAMEPCAP       float64
        PLNGENAN       float64
        PLC02EQA       float64
        dtype: object
```

```
In [ ]: # *(a)* Which plant has generated the most energy (measured in MWh)?
        plants_df.loc[plants_df['PLNGENAN'].idxmax()].PNAME
```

```
Out[ ]: 'Palo Verde'
```

```
In [ ]: # *(b)* What is the name of the northern-most power plant in the United States?
        plants_df.loc[plants_df['LAT'].idxmax()].PNAME
```

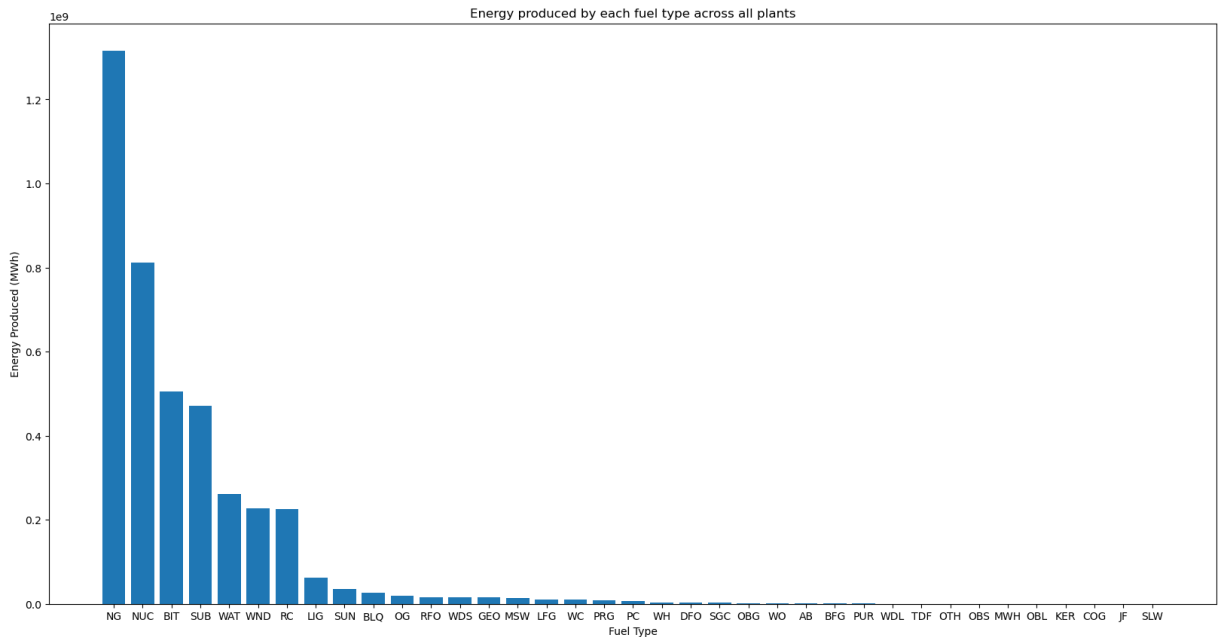
```
Out[ ]: 'Barrow'
```

```
In [ ]: # *(c)* What is the state where the northern-most power plant in the United States is located?
        plants_df.loc[plants_df['LAT'].idxmax()].PSTATABB
```

```
Out[ ]: 'AK'
```

```
In [ ]: # *(d)* Plot a bar plot showing the amount of energy produced by each fuel type across all plants
        import matplotlib.pyplot as plt

        energyprod_by_fuel = plants_df.groupby('PLPRMFL')['PLNGENAN'].sum()
        energyprod_by_fuel.sort_values(ascending=False, inplace=True)
        plt.bar(energyprod_by_fuel.index, energyprod_by_fuel.values)
        plt.title('Energy produced by each fuel type across all plants')
        plt.xlabel('Fuel Type')
        plt.ylabel('Energy Produced (MWh)')
        # change width of plot
        plt.rcParams['figure.figsize'] = [10, 5]
        plt.show()
```



```
In [ ]: # **(e)** From the plot in (d), which fuel for generation produces the most
print('Appears to be NG (natural gas)')
```

Appears to be NG (natural gas)

9

[6 points] *Vectorization.* When we first learn to code and think about iterating over an array, we often use loops. If implemented correctly, that does the trick. In machine learning, we iterate over so much data that those loops can lead to significant slow downs if they are not computationally efficient. In Python, vectorizing code and relying on matrix operations with efficient tools like numpy is typically the faster approach. Of course, numpy relies on loops to complete the computation, but this is at a lower level of programming (typically in C), and therefore is much more efficient. This exercise will explore the benefits of vectorization. Since many machine learning techniques rely on matrix operations, it's helpful to begin thinking about implementing algorithms using vector forms.

Begin by creating an array of 10 million random numbers using the numpy `random.randn` module. Compute the sum of the squares of those random numbers first in a for loop, then using Numpy's `dot` module to perform an inner (dot) product. Time how long it takes to compute each and report the results and report the output. How many times faster is the vectorized code than the for loop approach? (Note - your results may vary from run to run).

Your output should use the `print()` function as follows (where the # symbols represent your answers, to a reasonable precision of 4-5 significant figures):

Time [sec] (non-vectorized): #####

Time [sec] (vectorized): #####

The vectorized code is ##### times faster than the nonvectorized code

ANSWER

```
In [ ]: import numpy as np
import time

rand_nums = np.random.randn(10000000)

# sum of squares in for loop
mean = np.mean(rand_nums)
start_for_loop = time.time()
sum_of_squares = 0
for i in rand_nums:
    sum_of_squares += pow(i - mean, 2)
end_for_loop = time.time()

# inner dot
start_vectorized = time.time()
sum_squares_dot = np.dot(rand_nums, rand_nums)
end_vectorized = time.time()

forloop_runtime = round(end_for_loop - start_for_loop, 5)
vectorized_runtime = round(end_vectorized - start_vectorized, 5)

print('Time [sec] (non-vectorized): {}'.format(forloop_runtime))
print('Time [sec] (vectorized): {}'.format(vectorized_runtime))

runtime_diff_x_times_faster = round(forloop_runtime / vectorized_runtime, 5)

print('The vectorized code is {} times faster than the non-vectorized code'.

print('{} ..... {}'.format(sum_of_squares, sum_squares_dot))
```

Time [sec] (non-vectorized): 4.39721

Time [sec] (vectorized): 0.0028

The vectorized code is 1570.43214 times faster than the non-vectorized code
10002030.168880885 10002030.243680751

10

[10 points] This exercise will walk through some basic numerical programming and probabilistic thinking exercises, two skills which are frequently use in machine learning for answering questions from our data.

1. Synthesize $n = 10^4$ normally distributed data points with mean $\mu = 2$ and a standard deviation of $\sigma = 1$. Call these observations from a random variable X , and call the vector of observations that you generate, \mathbf{x} .
2. Calculate the mean and standard deviation of \mathbf{x} to validate (1) and provide the result to a precision of four significant figures.
3. Plot a histogram of the data in \mathbf{x} with 30 bins
4. What is the 90th percentile of \mathbf{x} ? The 90th percentile is the value below which 90% of observations can be found.
5. What is the 99th percentile of \mathbf{x} ?
6. Now synthesize $n = 10^4$ normally distributed data points with mean $\mu = 0$ and a standard deviation of $\sigma = 3$. Call these observations from a random variable Y , and call the vector of observations that you generate, \mathbf{y} .
7. Create a new figure and plot the histogram of the data in \mathbf{y} on the same axes with the histogram of \mathbf{x} , so that both histograms can be seen and compared.
8. Using the observations from \mathbf{x} and \mathbf{y} , estimate $E[XY]$

ANSWER

```
In [ ]: # 1
import matplotlib.pyplot as plt
x = np.random.normal(2, 1, 10000)

# 2
mean = round(np.mean(x), 4)
sd = round(np.std(x), 4)
print('mean: {}, sd: {}'.format(mean, sd))
```

mean: 1.9961, sd: 0.9964

```
In [ ]: # 4
print('90th percentile of x: {}'.format(np.percentile(x, 90)))
# 5
print('99th percentile of x: {}'.format(np.percentile(x, 99)))
```

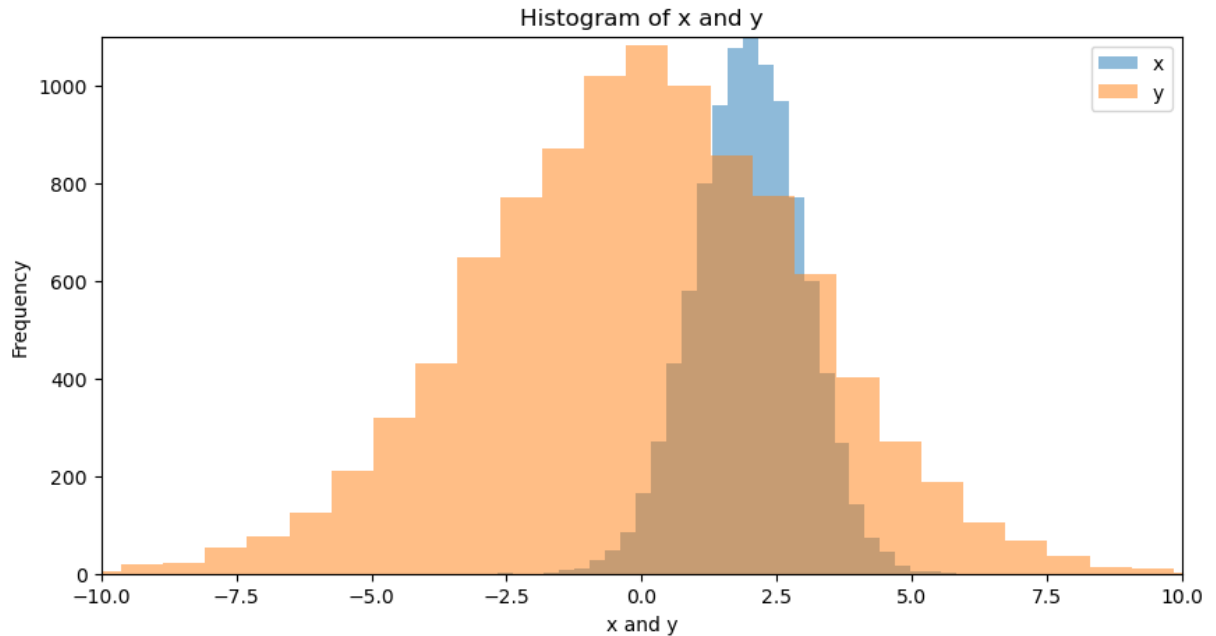
90th percentile of x: 3.270910998735618

99th percentile of x: 4.298753500269296

```
In [ ]: # 6
y = np.random.normal(0, 3, 10000)
# 3,7
# histogram of x and y
plt.hist(x, bins=30, alpha=0.5, label='x')
plt.hist(y, bins=30, alpha=0.5, label='y')
plt.title('Histogram of x and y')
plt.xlabel('x and y')
plt.ylabel('Frequency')
plt.legend(loc='upper right')
plt.rcParams['figure.figsize'] = [10, 5]
# change axes
axes = plt.gca()
axes.set_xlim([-10,10])
```



```
axes.set_ylim([0,1100])  
plt.show()
```



```
In [ ]: # 8  
print('E[XY] = {}'.format(np.mean(x*y))) # normal dist can be multiplied tog  
E[XY] = 0.007567326553118761
```

Version Control via Git

11

[4 points] Git is efficient for collaboration, and expectation in industry, and one of the best ways to share results in academia. You can even use some Git repositories (e.g. Github) as hosts for website, such as with the [course website](#). As a data scientist with experience in machine learning, Git is expected. We will interact with Git repositories (a.k.a. repos) throughout this course, and your project will require the use of git repos for collaboration.

Complete the [Atlassian Git tutorial](#), specifically the following listed sections. Try each concept that's presented. For this tutorial, instead of using BitBucket as your remote repository host, you may use your preferred platform such as [Github](#) or [Duke's Gitlab](#).

1. [What is version control](#)
2. [What is Git](#)
3. [Install Git](#)
4. [Setting up a repository](#)

5. [Saving changes](#)
6. [Inspecting a repository](#)
7. [Undoing changes](#)
8. [Rewriting history](#)
9. [Syncing](#)
10. [Making a pull request](#)
11. [Using branches](#)
12. [Comparing workflows](#)

I also have created two videos on the topic to help you understand some of these concepts: [Git basics](#) and a [step-by-step tutorial](#).

For your answer, affirm that you *either* completed the tutorials above OR have previous experience with ALL of the concepts above. Confirm this by typing your name below and selecting the situation that applies from the two options in brackets.

ANSWER

*I, **[your name here]**, affirm that I have **[completed the above tutorial / I have previous experience that covers all the content in this tutorial]***

I Jake Bova, affirm that I have completed the above tutorial / have previous experience that covers all the content in this tutorial

Exploratory Data Analysis

12

[15 points] Here you'll bring together some of the individual skills that you demonstrated above and create a Jupyter notebook based blog post on your exploratory data analysis. Your goal is to identify a question or problem and to work towards solving it or providing additional information or evidence (data) related to it through your data analysis. Below, we walk through a process to follow for your analysis. Additionally, you can find an [example of a well-done exploratory data analysis here from past years](#).

1. Find a dataset that interests you and relates to a question or problem that you find intriguing.
2. Describe the dataset, the source of the data, and the reason the dataset was of interest. Include a description of the features, data size, data creator and year of creation (if available), etc. What question are you hoping to answer through exploring the dataset?

3. Check the data and see if they need to be cleaned: are there missing values? Are there clearly erroneous values? Do two tables need to be merged together? Clean the data so it can be visualized. If the data are clean, state how you know they are clean (what did you check?).
4. Plot the data, demonstrating interesting features that you discover. Are there any relationships between variables that were surprising or patterns that emerged? Please exercise creativity and curiosity in your plots. You should have at least a ~3 plots exploring the data in different ways.
5. What insights are you able to take away from exploring the data? Is there a reason why analyzing the dataset you chose is particularly interesting or important? Summarize this for a general audience (imagine your publishing a blog post online) - boil down your findings in a way that is accessible, but still accurate.

Here your analysis will be evaluated based on:

1. Motivation: was the purpose of the choice of data clearly articulated? Why was the dataset chosen and what was the goal of the analysis?
2. Data cleaning: were any issues with the data investigated and, if found, were they resolved?
3. Quality of data exploration: were at least 4 unique plots (minimum) included and did those plots demonstrate interesting aspects of the data? Was there a clear purpose and takeaway from EACH plot?
4. Interpretation: Were the insights revealed through the analysis and their potential implications clearly explained? Was there an overall conclusion to the analysis?

ANSWER

SEE ATTACHED Q12.pdf