

# Assignment 1 - Probability, Linear Algebra, Programming, and Git

**Frank Xu**

Netid: hx44

## Probability and Statistics Theory

**1**

$$\text{Let } f(x) = \begin{cases} 0 & x < 0 \\ \alpha x^2 & 0 \leq x \leq 2 \\ 0 & 2 < x \end{cases}$$

For what value of  $\alpha$  is  $f(x)$  a valid probability density function?

*Note: for all assignments, write out all equations and math for all assignments using markdown and LaTeX (<https://tobi.oetiker.ch/lshort/lshort.pdf>) and show all work*

### ANSWER

For a valid probability density function, the definite integral of the function should equals to 1.

Therefore:

$$\begin{aligned} \int f(x)dx &= 0 + \int_0^2 \alpha x^2 dx + 0 = 1 \\ \int_0^2 \alpha x^2 dx &= 1 \\ \frac{1}{3} \alpha x^3 \Big|_0^2 &= 1 \\ \alpha &= \frac{3}{8} \end{aligned}$$

## 2

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 can be calculated from PDF under the following method:

$$F(x) = \int_{-\infty}^x f(t) dt$$

Therefore:

$$\int_0^3 \frac{1}{3} dx = \frac{1}{3} x \Big|_0^3$$
$$F(x) = \begin{cases} 0 & x \leq 0 \\ \frac{1}{3} x & 0 < x < 3 \\ 1 & x \geq 3 \end{cases}$$

## 3

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

(a) Expected value:

$$\mathbf{E}(X) = \int_{-\infty}^{\infty} xf(x)dx = \int_0^3 \frac{1}{3}x dx = \frac{1}{6}x^2 \Big|_0^3 = \frac{3}{2}$$

(b) Variance:

$$\mathbf{V}(X) = \mathbf{E}(X^2) - (\mathbf{E}(X))^2 = \int_0^3 \frac{1}{3}x^2 dx - \left(\frac{3}{2}\right)^2 = 3 - \frac{9}{4} = \frac{3}{4}$$

## 4

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, (b) variance, and the of the data?

Show all work. Your answer should include the definition of mean, median, and variance in the context of discrete data.

### ANSWER

(a) Mean

$$\mu = \frac{\sum x}{n} = \frac{13}{5}$$

(b) Variance

$$V(X) = \frac{\sum (x-\mu)^2}{n-1} = 4.506$$

## 5

Review of counting from probability theory.

(a) How many different 7-place license plates are possible if the first 3 places only contain letters and the last 4 only contain numbers?

(b) How many different batting orders are possible for a baseball team with 9 players?

(c) How many batting orders of 5 players are possible for a team with 9 players total?

(d) Let's assume this class has 26 students and we want to form project teams. How many unique teams of 3 are possible?

Hint: For each problem, determine if order matters, and if it should be calculated with or without replacement.

## ANSWER

(a)  $26^3 \times 10^4 = 175760000$

(b)  $9! = 362880$

(c)  ${}^9P_5 = \frac{9!}{(9-5)!} = 15120$

(d)  ${}^{26}C_3 = 2600$

## Linear Algebra

### 6

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

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

Compute the following or indicate that it cannot be computed:

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

$$1. \mathbf{AA} = \begin{bmatrix} 14 & 25 & 31 \\ 25 & 45 & 56 \\ 31 & 56 & 70 \end{bmatrix}$$

$$2. \mathbf{AA}^T = \mathbf{AA} = \begin{bmatrix} 14 & 25 & 31 \\ 25 & 45 & 56 \\ 31 & 56 & 70 \end{bmatrix}$$

$$3. \mathbf{Ab} = \begin{bmatrix} 29 \\ 50 \\ 60 \end{bmatrix}$$

4.  $\mathbf{Ab}^T$  cannot be computed.

5.  $\mathbf{bA}$  cannot be computed.

$$6. \mathbf{b}^T \mathbf{A} = \begin{bmatrix} 29 & 50 & 60 \end{bmatrix}$$

7.  $\mathbf{bb}$  cannot be computed.

$$8. \mathbf{b}^T \mathbf{b} = \begin{bmatrix} 74 \end{bmatrix}$$

$$9. \mathbf{bb}^T = \begin{bmatrix} 1 & -3 & -8 \\ -3 & 9 & 24 \\ -8 & 24 & 64 \end{bmatrix}$$

10.  $\mathbf{b} + \mathbf{c}^T$  cannot be computed.

11.  $\mathbf{b}^T \mathbf{b}^T$  cannot be computed.

$$12. \mathbf{A}^{-1} \mathbf{b} = \begin{bmatrix} 1 & -3 & 2 \\ -3 & 3 & -1 \\ 2 & -1 & 0 \end{bmatrix} \mathbf{b} = \begin{bmatrix} 6 \\ 4 \\ -5 \end{bmatrix}$$

$$13. \mathbf{A} \circ \mathbf{A} = \begin{bmatrix} 1 & 4 & 9 \\ 4 & 16 & 25 \\ 9 & 25 & 36 \end{bmatrix}$$

$$14. \mathbf{b} \circ \mathbf{c} = \begin{bmatrix} -4 \\ -9 \\ 48 \end{bmatrix}$$

## 6

**Eigenvectors and eigenvalues.** Eigenvectors and eigenvalues are useful for some machine learning algorithms, but the concepts take time to solidly grasp. For an intuitive review of these concepts, explore this [interactive website at Setosa.io](http://setosa.io/ev/eigenvectors-and-eigenvalues/) (<http://setosa.io/ev/eigenvectors-and-eigenvalues/>). Also, the series of linear algebra videos by Grant Sanderson of 3Brown1Blue are excellent and can be viewed on youtube [here](https://www.youtube.com/playlist?list=PLZHQObOWTQDPD3MizzM2xVFitgF8hE_ab) ([https://www.youtube.com/playlist?list=PLZHQObOWTQDPD3MizzM2xVFitgF8hE\\_ab](https://www.youtube.com/playlist?list=PLZHQObOWTQDPD3MizzM2xVFitgF8hE_ab)).

1. Calculate the eigenvalues and corresponding eigenvectors of matrix  $\mathbf{A}$  above, from the last question.
2. Choose one of the eigenvector/eigenvalue pairs,  $\mathbf{v}$  and  $\lambda$ , and show that  $\mathbf{A}\mathbf{v} = \lambda\mathbf{v}$ . Also show that 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 real, symmetric matrices.

## ANSWER

1. *Eigenvalues* :  $\lambda_1 = 11.345, \lambda_2 = 0.171, \lambda_3 = -0.516$

$$\text{Eigenvectors : } v_1 = \begin{bmatrix} -0.328 \\ -0.591 \\ -0.737 \end{bmatrix}, v_2 = \begin{bmatrix} 0.591 \\ -0.737 \\ 0.328 \end{bmatrix}, v_3 = \begin{bmatrix} 0.737 \\ 0.328 \\ -0.591 \end{bmatrix}$$

2. Given the *eigenvalue* : 11.345 and the corresponding *eigenvector* :  $\begin{bmatrix} -0.328 \\ -0.591 \\ -0.737 \end{bmatrix}$ ,

$$\text{left} = \mathbf{A}\mathbf{v} = \begin{bmatrix} -3.721 \\ -6.705 \\ -8.361 \end{bmatrix}, \text{right} = \lambda\mathbf{v} = \begin{bmatrix} -3.721 \\ -6.705 \\ -8.361 \end{bmatrix}$$

$$\text{left} = \mathbf{A}\mathbf{A}\mathbf{v} = \begin{bmatrix} -42.214 \\ -76.067 \\ -94.854 \end{bmatrix}, \text{right} = \lambda^2\mathbf{v} = \begin{bmatrix} -42.217 \\ -76.067 \\ -94.859 \end{bmatrix}$$

Despite the slightly deviation due to calculating accuracy, for both equations, left = right.

$$1. v_1^T v_2 = \begin{bmatrix} -0.328 & -0.591 & -0.737 \end{bmatrix} \begin{bmatrix} 0.591 \\ -0.737 \\ 0.328 \end{bmatrix} = -1.7 \times 10^{-5}$$

$$v_2^T v_3 = -1.7 \times 10^{-5}$$

$$v_3^T v_1 = -1.7 \times 10^{-5}$$

The result is  $-1.7 \times 10^{-5}$  due to accuracy problems. If we apply the eigenvectors with more accuracy, the result will approach more to 0.

## Numerical Programming

## 7

Speed comparison between vectorized and non-vectorized code. Begin by creating an array of 10 million random numbers using the numpy random.randn module. Compute the sum of the squares first in a for loop, then using Numpy's dot module. 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: all code should be well commented, properly formatted, and your answers should be output using the print() function as follows (where the # represents your answers, to a reasonable precision):

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

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

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

### ANSWER

```
In [2]: import numpy as np
import time

# Generate the random samples
arr = np.random.randn(10000000)

# Compute the sum of squares the non-vectorized way (using a for loop)
t1 = time.time()
ans1 = 0
for i in arr:
    ans1 += i**2
t_nv = time.time() - t1

# Compute the sum of squares the vectorized way (using numpy)
t2 = time.time()
ans2 = sum(arr ** 2)
t_v = time.time() - t2

# Print the results
print('Time [sec] (non-vectorized): {}'.format(t_nv))

print('Time [sec] (vectorized): {}'.format(t_v))

print('The vectorized code is {} times faster than the vectorized code
'.format(t_nv / t_v))
```

Time [sec] (non-vectorized): 3.4999070167541504

Time [sec] (vectorized): 0.7863631248474121

The vectorized code is 4.450751702571609 times faster than the vectorized code



## 8

One popular Agile development framework is Scrum (a paradigm recommended for data science projects). It emphasizes the continual evolution of code for projects, becoming progressively better, but starting with a quickly developed minimum viable product. This often means that code written early on is not optimized, and that's a good thing - it's best to get it to work first before optimizing. Imagine that you wrote the following code during a sprint towards getting an end-to-end system working. Vectorize the following code and show the difference in speed between the current implementation and a vectorized version.

The function below computes the function  $f(x, y) = x^2 - 2y^2$  and determines whether this quantity is above or below a given threshold, `thresh=0`. This is done for  $x, y \in \{-4, 4\}$ , over a 2,000-by-2,000 grid covering that domain.

- (a) Vectorize this code and demonstrate (as in the last exercise) the speed increase through vectorization and
- (b) plot the resulting data - both the function  $f(x, y)$  and the thresholded output - using `imshow` ([https://matplotlib.org/api/\\_as\\_gen/matplotlib.pyplot.imshow.html?highlight=matplotlib%20pyplot%20imshow#matplotlib.pyplot.imshow](https://matplotlib.org/api/_as_gen/matplotlib.pyplot.imshow.html?highlight=matplotlib%20pyplot%20imshow#matplotlib.pyplot.imshow)) from `matplotlib`.

*Hint: look at the `numpy meshgrid` (<https://docs.scipy.org/doc/numpy-1.13.0/reference/generated/numpy.meshgrid.html>) documentation*

```
In [4]: import numpy as np
import time
import matplotlib.pyplot as plt

# Initialize variables for this exercise
thresh = 0
points = np.linspace(-4, 4, 2000)
x, y = np.meshgrid(points, points)
t0 = time.time()

# Nonvectorized implementation
f1 = np.zeros(shape=(2000,2000))
for i in range(len(points)):
    for j in range(len(points)):
        if (x[i,j]**2 - 2*y[i,j]**2) > thresh:
            f1[i,j] = 1
    pass
pass
t1 = time.time()

# Vectorized implementation
f = (x**2 - 2*y**2)
f2 = 1 * (f > thresh)
t2 = time.time()

# Print the time for each and the speed increase
print('Time [sec] (non-vectorized): {}'.format(t1-t0))

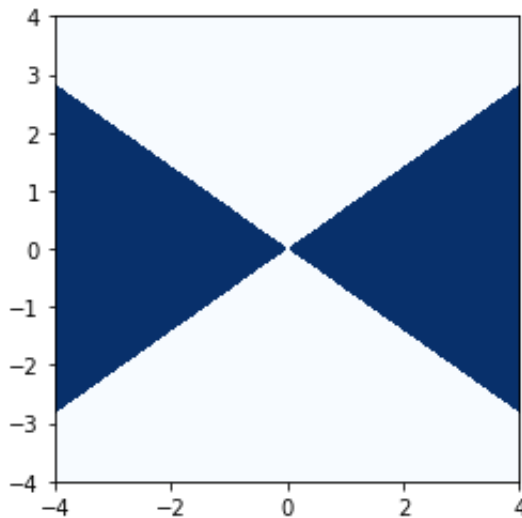
print('Time [sec] (vectorized):      {}'.format(t2-t1))

print('The vectorized code is {} times faster than the vectorized code
'.format((t1-t0)/(t2-t1)))

# Plot the result
plt.imshow(f2, cmap='Blues', interpolation='none', extent=[-4,4,-4,4])
# Darker blue indicates f(x) is above the threshold.
```

```
Time [sec] (non-vectorized): 4.454632997512817
Time [sec] (vectorized):      0.03504300117492676
The vectorized code is 127.11904940094298 times faster than the vect
orized code
```

```
Out[4]: <matplotlib.image.AxesImage at 0x134975f60>
```



## 9

This exercise will walk through some basic numerical programming exercises.

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. Plot the histogram of the data in  $\mathbf{y}$  on a (new) plot 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 [5]: # 1
x = np.random.normal(2, 1, 10000)

# 2
mu = np.mean(x)
sigma = np.std(x)
print('mu = {:.4}, sigma = {:.4}'.format(mu, sigma))

# 3
plt.figure()
plt.hist(x, bins = 30)

# 4
P90 = np.percentile(x, 90)
print(P90)

# 5
P99 = np.percentile(x, 99)
print(P99)

# 6
y = np.random.normal(0, 3, 10000)

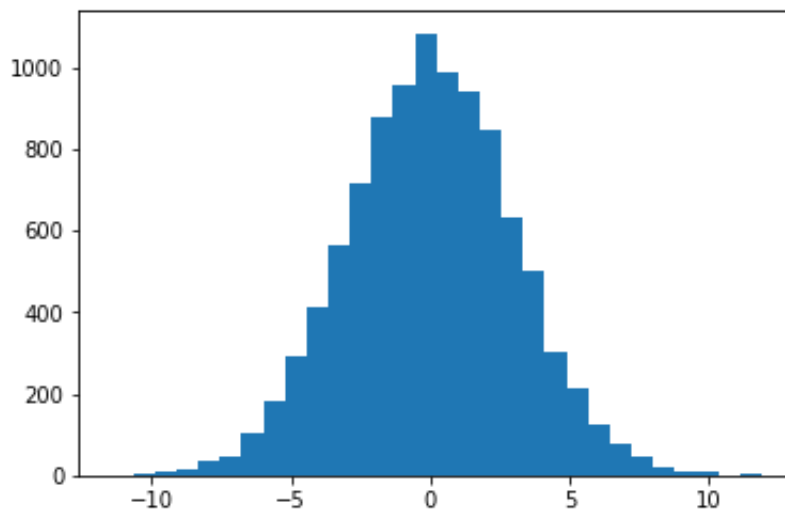
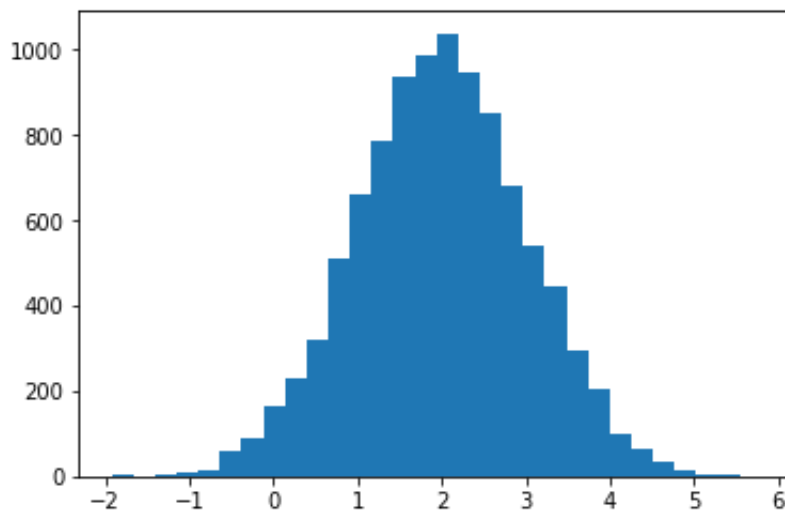
# 7
plt.figure()
plt.hist(y, bins = 30)

# 8
E_xy = np.mean(x)*np.mean(y)
print('E[xy] = {}'.format(E_xy))
```

```

mu = 1.997, sigma = 1.008
3.310118483220631
4.331604504343274
E[xy] = -0.008680591079317249

```



$X \sim \text{Gauss}(\mu = 2, \sigma = 1)$

## 10

Estimate the integral of the function  $f(x)$  on the interval  $0 \leq x < 2.5$  assuming we only know the following points from  $f$ :

Table 1. Dataset containing  $n=5$  observations

$x_i$	0.0	0.5	1.0	1.5	2.0
$y_i$	6	7	8	4	1

## ANSWER

Suppose the  $y$  value between  $x_i$  and  $x_{i+1}$  is  $y_i$ .

$$\int xf(x)dx = \sum \int (x_{i+1} - x_i)y_i dx = 0.5 \times (6 + 7 + 8 + 4 + 1) = 13$$

## Version Control via Git

### 11

Complete the [Atlassian Git tutorial \(https://www.atlassian.com/git/tutorials/what-is-version-control\)](https://www.atlassian.com/git/tutorials/what-is-version-control), specifically the following sections. Try each concept that's presented. For this tutorial, instead of using BitBucket, use Github. Create a github account here if you don't already have one: <https://github.com/> (<https://github.com/>)

1. [What is version control \(https://www.atlassian.com/git/tutorials/what-is-version-control\)](https://www.atlassian.com/git/tutorials/what-is-version-control)
2. [What is Git \(https://www.atlassian.com/git/tutorials/what-is-git\)](https://www.atlassian.com/git/tutorials/what-is-git)
3. [Install Git \(https://www.atlassian.com/git/tutorials/install-git\)](https://www.atlassian.com/git/tutorials/install-git)
4. [Setting up a repository \(https://www.atlassian.com/git/tutorials/install-git\)](https://www.atlassian.com/git/tutorials/install-git)
5. [Saving changes \(https://www.atlassian.com/git/tutorials/saving-changes\)](https://www.atlassian.com/git/tutorials/saving-changes)
6. [Inspecting a repository \(https://www.atlassian.com/git/tutorials/inspecting-a-repository\)](https://www.atlassian.com/git/tutorials/inspecting-a-repository)
7. [Undoing changes \(https://www.atlassian.com/git/tutorials/undoing-changes\)](https://www.atlassian.com/git/tutorials/undoing-changes)
8. [Rewriting history \(https://www.atlassian.com/git/tutorials/rewriting-history\)](https://www.atlassian.com/git/tutorials/rewriting-history)
9. [Syncing \(https://www.atlassian.com/git/tutorials/syncing\)](https://www.atlassian.com/git/tutorials/syncing)
10. [Making a pull request \(https://www.atlassian.com/git/tutorials/making-a-pull-request\)](https://www.atlassian.com/git/tutorials/making-a-pull-request)
11. [Using branches \(https://www.atlassian.com/git/tutorials/using-branches\)](https://www.atlassian.com/git/tutorials/using-branches)
12. [Comparing workflows \(https://www.atlassian.com/git/tutorials/comparing-workflows\)](https://www.atlassian.com/git/tutorials/comparing-workflows)

For your answer, affirm that you either completed the tutorial or have previous experience with all of the concepts above. Do this by typing your name below and selecting the situation that applies from the two options in brackets.

## ANSWER

I, **[Frank Xu]**, affirm that I have **[completed the above tutorial]**

# 12

Using Github to create a static HTML website:

1. Create a branch in your `machine-learning-course` repo called "gh-pages" and checkout that branch (this will provide an example of how to create a simple static website using [Github Pages](https://pages.github.com/) (<https://pages.github.com/>))
2. Create a file called "index.html" with the contents "Hello World" and add, commit, and push it to that branch.
3. Submit the following: (a) a link to your github repository and (b) a link to your new "Hello World" website. The latter should be at the address [https://\[USERNAME\].github.io/ECE590-assignment0](https://[USERNAME].github.io/ECE590-assignment0) ([https://\[USERNAME\].github.io/ECE590-assignment0](https://[USERNAME].github.io/ECE590-assignment0)) (where [USERNAME] is your github username).

## ANSWER

<https://github.com/Frank-Xu-Huaze/machine-learning-course> (<https://github.com/Frank-Xu-Huaze/machine-learning-course>)

<https://github.com/Frank-Xu-Huaze/machine-learning-course/blob/gh-pages/index.html>  
(<https://github.com/Frank-Xu-Huaze/machine-learning-course/blob/gh-pages/index.html>)

# Exploratory Data Analysis

## 13

Here you'll bring together some of the individual skills that you demonstrated above and create a Jupyter notebook based blog post on data analysis.

1. Find a dataset that interests you and relates to a question or problem that you find intriguing
2. Using a Jupyter notebook, describe the dataset, the source of the data, and the reason the dataset was of interest.
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.
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.
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 as if your target audience was the readership of a major news organization - boil down your findings in a way that is accessible, but still accurate.
6. Create a public repository on your github account titled "machine-learning-course". In it, create a readme file that contains the heading "ECE590: Introductory Machine Learning for Data Science". Add, commit, and push that Jupyter notebook to the master branch. Provide the link to the that post here.

## ANSWER

[https://github.com/Frank-Xu-Huaze/machine-learning-course/blob/master/EDA/Exploratory\\_Data\\_Analysis.ipynb](https://github.com/Frank-Xu-Huaze/machine-learning-course/blob/master/EDA/Exploratory_Data_Analysis.ipynb) ([https://github.com/Frank-Xu-Huaze/machine-learning-course/blob/master/EDA/Exploratory\\_Data\\_Analysis.ipynb](https://github.com/Frank-Xu-Huaze/machine-learning-course/blob/master/EDA/Exploratory_Data_Analysis.ipynb))