

Short HW1 - Preparing for the course

Useful python libraries, Probability, and Linear algebra

Instructions

General

- **First, don't panic!**
 - This assignment seems longer than it actually is.
 - In the first part, you are mostly required to run *existing* code and complete short python commands here and there.
 - In the two other parts you need to answer overall 4 analytic questions.
 - **Note:** The other 4 *short* assignments will be shorter and will *not* require programming.
- **Individually or in pairs?** Individually only.
- **Where to ask?** In the [Piazza forum \(http://piazza.com/technion.ac.il/spring2023/236756\)](http://piazza.com/technion.ac.il/spring2023/236756).
- **How to submit?** In the webcourse.
- **What to submit?** A pdf file with the completed jupyter notebook (including the code, plots and other outputs) and the answers to the probability/algebra questions (Hebrew or English are both fine).
Or two separate pdf files in a zip file. All submitted files should contain **your ID number** in their names.
- **When to submit?** Tuesday 04.04.2023 at 23:59.

Specific

- First part: get familiar with popular python libraries useful for machine learning and data science. We will use these libraries heavily throughout the major programming assignments.
 - You should read the instructions and run the code blocks sequentially.
In 10 places you are required to complete missing python commands or answer short questions (look for the **TODO** comments, or notations like **(T3)** etc.). Try to understand the flow of this document and the code you run.
 - Start by loading the provided jupyter notebook file (*Short_HW1.ipynb*) to [Google Colab \(https://colab.research.google.com/\)](https://colab.research.google.com/), which is a very convenient online tool for running python scripts combined with text, visual plots, and more.
 - Alternatively, you can [install jupyter \(https://jupyter.org/install\)](https://jupyter.org/install) locally on your computer and run the provided notebook there.
- Second and third parts: questions on probability and linear algebra to refresh your memory and prepare for the rest of this course.
The questions are mostly analytic but also require completing and running simple code blocks in the jupyter notebook.
 - Forgot your linear algebra? Try watching [Essence of LA \(https://www.youtube.com/playlist?list=PLZHQObOWTQDPD3MizzM2xVFitgF8hE_ab\)](https://www.youtube.com/playlist?list=PLZHQObOWTQDPD3MizzM2xVFitgF8hE_ab) or reading [The Matrix Cookbook \(http://www2.imm.dtu.dk/pubdb/views/edoc_download.php/3274/pdf/imm3274.pdf\)](http://www2.imm.dtu.dk/pubdb/views/edoc_download.php/3274/pdf/imm3274.pdf).
 - Forgot your probability? Try reading [Probability Theory Review for Machine Learning \(https://see.stanford.edu/materials/aimlcs229/cs229-prob.pdf\)](https://see.stanford.edu/materials/aimlcs229/cs229-prob.pdf).
 - Correction: In 3.2 it says that $X \perp Y \implies \text{Var}(X + Y) = \text{Var}(X)\text{Var}(Y)$ but it should say $X \perp Y \implies \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$.

Important: How to submit the notebook's output?

You should only submit PDF file(s). In the print dialog of your browser, you can choose to `Save as PDF`. However, notice that some of the outputs may be cropped (become invisible), which can harm your grade.

To prevent this from happening, tune the "scale" of the printed file, to fit in the *entire* output. For instance, in Chrome you should lower the value in `More settings->Scale->Custom` to contain the entire output (50%~ often work well).

Good luck!

What is pandas?

Python library for Data manipulation and Analysis

- Provide expressive data structures designed to make working with “relational” or “labeled” data both easy and intuitive.
- Aims to be the fundamental high-level building block for doing practical, real world data analysis in Python.
- Built on top of NumPy and is intended to integrate well within a scientific computing.
- Inspired by R and Excel.

Pandas is well suited for many different kinds of data:

- **Tabular data** with heterogeneously-typed columns, as in an SQL table or Excel spreadsheet
- Ordered and unordered (not necessarily fixed-frequency) **time series data**.
- **Arbitrary matrix data** (homogeneously typed or heterogeneous) with row and column labels
- Any other form of observational / statistical data sets (can be unlabeled)

Two primary data structures

- **Series** (1-dimensional) – Similar to a column in Excel's spreadsheet
- **Data Frame** (2-dimensional) – Similar to R's data frame

A few of the things that Pandas does well

- Easy handling of **missing data** (represented as NaN)
- Automatic and explicit **data alignment**
- Read and Analyze **CSV**, Excel Sheets Easily
- Operations
- Filtering, Group By, Merging, Slicing and Dicing, Pivoting and Reshaping
- Plotting graphs

Pandas is very useful for interactive data exploration at the data preparation stage of a project

The official guide to Pandas can be found [here \(http://pandas-docs.github.io/pandas-docs-travis/10min.html\)](http://pandas-docs.github.io/pandas-docs-travis/10min.html)

Pandas Objects

```
In [ ]: import pandas as pd
import numpy as np
```

Series is like a column in a spreadsheet.

```
In [ ]: s = pd.Series([1,3.2,np.nan,'string'])
s
```

DataFrame is like a spreadsheet – a dictionary of Series objects

```
In [ ]: data = [['ABC', -3.5, 0.01], ['ABC', -2.3, 0.12], ['DEF', 1.8, 0.03],
['DEF', 3.7, 0.01], ['GHI', 0.04, 0.43], ['GHI', -0.1, 0.67]]

df = pd.DataFrame(data, columns=['gene', 'log2FC', 'pval'])

df
```

Input and Output

How do you get data into and out of Pandas as spreadsheets?

- Pandas can work with XLS or XLSX files.
- Can also work with CSV (comma separated values) file
- CSV stores plain text in a tabular form
- CSV files may have a header
- You can use a variety of different field delimiters (rather than a 'comma'). Check which delimiter your file is using before import!

Import to Pandas

```
df = pd.read_csv('data.csv', sep='\t', header=0)
```

For Excel files, it's the same thing but with read_excel

Export to text file

```
df.to_csv('data.csv', sep='\t', header=True, index=False)
```

The values of header and index depend on if you want to print the column and/or row names

Case Study – Analyzing Titanic Passengers Data

```
In [ ]: import matplotlib.pyplot as plt
        %matplotlib inline
        import numpy as np
        import pandas as pd
        import os

        #set your working_dir
        working_dir = os.path.join(os.getcwd(), 'titanic')

        url_base = 'https://github.com/Currie32/Titanic-Kaggle-Competition/raw/master/{}.csv'
        train_url = url_base.format('train')
        test_url = url_base.format('test')

        # For .read_csv, always use header=0 when you know row 0 is the header row
        train = pd.read_csv(train_url, header=0)
        test = pd.read_csv(test_url, header=0)
        # You can also load a csv file from a local file rather than a URL
```

(T1) Use `pandas.DataFrame.head` (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.head.html>) to display the top 6 rows of the `train` table

```
In [ ]: # TODO: print the top 6 rows of the table
```

VARIABLE DESCRIPTIONS:

Survived - 0 = No; 1 = Yes

Age - Passenger's age

Pclass - Passenger Class (1 = 1st; 2 = 2nd; 3 = 3rd)

SibSp - Number of Siblings/Spouses Aboard

Parch - Number of Parents/Children Aboard

Ticket - Ticket Number

Fare - Passenger Fare

Cabin - Cabin ID

Embarked - Port of Embarkation (C = Cherbourg; Q = Queenstown; S = Southampton)

```
In [ ]: train.columns
```

Understanding the data (Summarizations)

```
In [ ]: train.info()
```

```
In [ ]: train.shape
```

```
In [ ]: # Count values of 'Survived'
train.Survived.value_counts()
```

```
In [ ]: # Calculate the mean fare price
train.Fare.mean()
```

```
In [ ]: # General statistics of the dataframe
train.describe()
```

Selection examples

Selecting columns

```
In [ ]: # Selection is very similar to standard Python selection
df1 = train[["Name", "Sex", "Age", "Survived"]]
df1.head()
```

Selecting rows

```
In [ ]: df1[10:15]
```

Filtering Examples

Filtering with one condition

```
In [ ]: # Filtering allows you to create masks given some conditions
df1.Sex == 'female'
```

```
In [ ]: onlyFemale = df1[df1.Sex == 'female']
onlyFemale.head()
```

Filtering with multiple conditions

(T2) Alter the following command so `adultFemales` will contain only females whose age is 18 and above. You need to filter using a **single** mask with multiple conditions (google it!), i.e., without creating any temporary dataframes.

Additionally, update the `survivalRate` variable to show the correct rate.

```
In [ ]: # TODO: update the mask
adultFemales = df1[(df1.Sex == 'female')]

# TODO: Update the survival rate
survivalRate = 0
print("The survival rate of adult females was: {:.2f}%".format(survivalRate * 100))
```

Aggregating

Pandas allows you to aggregate and display different views of your data.

```
In [ ]: df2 = train.groupby(['Pclass', 'Sex']).Fare.agg(np.mean)
df2
```

```
In [ ]: pd.pivot_table(train, index=['Pclass'], values=['Survived'], aggfunc='count')
```

The following table shows the survival rates for each combination of passenger class and sex.

(T3) Add a column showing the mean **age** for such a combination.

```
In [ ]: # TODO: Also show the mean age per group
pd.pivot_table(train, index=['Pclass', 'Sex'], values=['Survived'], aggfunc='mean')
```

(T4) Use [this \(https://stackoverflow.com/questions/21441259/pandas-groupby-range-of-values\)](https://stackoverflow.com/questions/21441259/pandas-groupby-range-of-values) question on stackoverflow, to find the mean survival rate for ages 0-10, 10-20, etc.).

Hint: the first row should roughly look like this:

	Age	Survived
Age		
(0, 10]	4.268281	0.593750

```
In [ ]: # TODO: find the mean survival rate per age group
ageGroups = np.arange(0, 81, 10)
survivalPerAgeGroup = None

survivalPerAgeGroup
```

```
In [ ]: type(train.groupby(pd.cut(train.Age, ageGroups)).Survived.mean())
```

Filling missing data (data imputation)

Note that some passenger do not have age data.

```
In [ ]: print("{} out of {} passengers do not have a recorded age".format(df1[df1.Age.isna()].shape[0], df1.shape[0]))
```

```
In [ ]: df1[df1.Age.isna()].head()
```

Let's see the statistics of the column **before** the imputation.

```
In [ ]: df1.Age.describe()
```

Read about `pandas.Series.fillna` (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.Series.fillna.html?highlight=fillna#pandas.Series.fillna>).

(T5) Replace the missing ages `df1` with the general age *median*, and insert the result into variable `filledDf` (the original `df1` should be left unchanged).

```
In [ ]: # TODO : Fill the missing values
filledDf = df1
```

```
In [ ]: print("{} out of {} passengers do not have a recorded age".format(filledDf[filledDf.Age.isna()].shape[0], filledDf.shape[0]))
```

Let's see the statistics of the column **after** the imputation.

```
In [ ]: filledDf.Age.describe()
```

(T6) Answer below: which statistics changed, and which did not? Why? (explain briefly, no need to be very formal.)

Answer: TODO

Plotting

Basic plotting in pandas is pretty straightforward

```
In [ ]: new_plot = pd.crosstab([train.Pclass, train.Sex], train.Survived, normalize="index")
new_plot.plot(kind='bar', stacked=True, grid=False, figsize=(10,6))
plt.yticks(np.linspace(0,1,21))
plt.grid()
```

(T7) Answer below: which group (class × sex) had the best survival rate? Which had the worst?

Answer: TODO

What is Matplotlib

A 2D plotting library which produces publication quality figures.

- Can be used in python scripts, the python and IPython shell, web application servers, and more ...
- Can be used to generate plots, histograms, power spectra, bar charts, errorcharts, scatterplots, etc.
- For simple plotting, pyplot provides a MATLAB-like interface
- For power users, a full control via OO interface or via a set of functions

There are several Matplotlib add-on toolkits

- Projection and mapping toolkits [basemap](http://matplotlib.org/basemap/) (<http://matplotlib.org/basemap/>) and [cartopy](http://scitools.org.uk/cartopy/) (<http://scitools.org.uk/cartopy/>).
- Interactive plots in web browsers using [Bokeh](http://bokeh.pydata.org/en/latest/) (<http://bokeh.pydata.org/en/latest/>).
- Higher level interface with updated visualizations [Seaborn](http://seaborn.pydata.org/index.html) (<http://seaborn.pydata.org/index.html>).

Matplotlib is available at www.matplotlib.org (www.matplotlib.org).

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

Line Plots

The following code plots the survival rate per age group (computed above, before the imputation).

(T8) Use the [matplotlib documentation](https://matplotlib.org/) (<https://matplotlib.org/>) to add a grid and suitable axis labels to the following plot.

```
In [ ]: plt.plot(survivalPerAgeGroup.Age, survivalPerAgeGroup.Survived)
_ = plt.title("Survival per age group")
# TODO : Update the plot as required.
```

```
In [ ]: survivalPerAgeGroup
```

Scatter plots

(T9) Alter the [matplotlib.pyplot.scatter](https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.scatter.html) (https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.scatter.html) command, so that the scattered dots will be green , and their size will be 10 .

Also, add a grid and suitable axis labels.

```
In [ ]: # TODO : Update the plot as required.
plt.figure(figsize=(10,6))
plt.scatter(train.Age, train.Fare)
```

(T10) Answer below: approximately how old are the two highest paying passengers?

Answer: TODO

Probability refresher

Q1 - Variance of empirical mean

Let X_1, \dots, X_m be i.i.d random variables with mean $\mathbb{E}[X_i] = \mu$ and variance $\text{Var}(X_i) = \sigma^2$.

We would like to "guess", or more formally, estimate (הערכה), the mean μ from the observations x_1, \dots, x_m .

We use the empirical mean $\bar{X} = \frac{1}{m} \sum_i X_i$ as an estimator for the unknown mean μ . Notice that \bar{X} is itself a random variable.

Note: The instantiation of \bar{X} is usually denoted by $\hat{\mu} = \frac{1}{m} \sum_i x_i$, but this is currently out of scope.

1. Express analytically the expectation of \bar{X} .

Answer: $\mathbb{E}[\bar{X}] = \mathbb{E}[\frac{1}{m} \sum_i X_i] = \text{TODO}$.

2. Express analytically the variance of \bar{X} .

Answer: $\text{Var}[\bar{X}] = \text{TODO}$.

You will now verify the expression you wrote for the variance.

We assume $\forall i : X_i \sim \mathcal{N}(0, 1)$.

We compute the empirical mean's variances for sample sizes $m = 1, \dots, 35$.

For each sample size m , we sample m normal variables and compute their empirical mean. We repeat this step 50 times, and compute the variance of the empirical means (for each m).

3. Complete the code blocks below according to the instructions and verify that your analytic function of the empirical mean's variance against as a function of m suits the empirical findings.

```
In [ ]: all_sample_sizes = range(1, 36)
        repeats_per_size = 50

        allVariances = []

        for m in all_sample_sizes:
            empiricalMeans = []

            for _ in range(repeats_per_size):
                # Random m examples and compute their empirical mean
                X = np.random.randn(m)
                empiricalMeans.append(np.mean(X))

            # TODO: Using numpy, compute the variance of the empirical means that are in
            # the `empiricalMeans` list (you can google the numpy function for variance)
            variance = None

            allVariances.append(variance)
```

Complete the following computation of the analytic variance (according to the your answers above). You can try to use simple arithmetic operations between an `np.array` and a scalar, and see what happens! (for instance, `2 * np.array(all_sample_sizes)`.)

```
In [ ]: # TODO: compute the analytic variance
        # (the current command wrongfully sets the variance of an empirical mean
        # of a sample with m variables simply as 2*m)
        analyticVariance = 2 * np.array(all_sample_sizes).astype(float)
```

The following code plots the results from the above code. **Do not** edit it, only run it and make sure that the figures make sense.

```
In [ ]: fig, axes = plt.subplots(1,2, figsize=(15,5))
axes[0].plot(all_sample_sizes, analyticVariance, label="Analytic", linewidth=4)
axes[0].plot(all_sample_sizes, allVariances, label="Empiric", linewidth=3)
axes[0].grid()
axes[0].legend(fontsize=14)
axes[0].set_title("Regular scale", fontsize=14)
axes[0].set_xlabel("Sample size (m)", fontsize=12)
axes[0].set_ylabel("Variance", fontsize=12)

axes[1].semilogy(all_sample_sizes, analyticVariance, label="Analytic", linewidth=4)
axes[1].semilogy(all_sample_sizes, allVariances, label="Empiric", linewidth=3)
axes[1].grid()
axes[1].legend(fontsize=14)
axes[1].set_title("Log scale", fontsize=14)
axes[1].set_xlabel("Sample size (m)", fontsize=12)
axes[1].set_ylabel("Variance", fontsize=12)

_ = plt.suptitle("Empirical mean's variance vs. Sample size",
                fontsize=16, fontweight="bold")

plt.tight_layout()
```

Reminder - Hoeffding's Inequality

Let $\theta_1, \dots, \theta_m$ be i.i.d random variables with mean $\mathbb{E}[\theta_i] = \mu$.

Additionally, assume all variables are bound in $[a, b]$ such that $\Pr[a \leq \theta_i \leq b] = 1$.

Then, for any $\epsilon > 0$, the empirical mean $\bar{\theta}(m) = \frac{1}{m} \sum_i \theta_i$ holds:

$$\Pr\left[\left|\bar{\theta}(m) - \mu\right| > \epsilon\right] \leq 2 \exp\left\{-\frac{2m\epsilon^2}{(b-a)^2}\right\}.$$

Q2 - Identical coins and the Hoeffding bound

We toss $m \in \mathbb{N}$ identical coins, each coin 50 times.

All coins have the same *unknown* probability of showing "heads", denoted by $p \in (0, 1)$.

Let θ_i be the (observed) number of times the i -th coin showed "heads".

1. What is the distribution of each θ_i ?

Answer: $\theta_i \sim \text{TODO}$.

2. What is the mean $\mu = \mathbb{E}[\theta_i]$?

Answer: $\mathbb{E}[\theta_i] = \text{TODO}$.

3. We would like to use the empirical mean defined above as an estimator $\bar{\theta}(m)$ for μ .

Use Hoeffding's inequality to compute the *smallest* sample size $m \in \mathbb{N}$ that can guarantee an error of $\epsilon = 1$ with confidence 0.99 (notice that we wish to estimate μ , not p).

That is, find the smallest m that holds $\Pr\left[\left|\bar{\theta}(m) - \mu\right| > 1\right] \leq 0.01$.

Answer:

TODO

4. The following code simulates tossing $m = 10^4$ coins, each 50 times. For each coin, we use the empirical mean as the estimator and save it in the `all_estimators` array. The (unknown) probability of each coin is 0.65.

Complete the missing part so that for each coin, an array of 50 binary observations will be randomized according to the probability p .

```
In [ ]: m = 10**4
tosses = 50
p = 0.65
all_estimators = []

# Repeat for n coins
for coin in range(m):
    # TODO: Use Google to find a suitable numpy.random function that creates
    # a binary array of size (tosses,), where each element is 1
    # with probability p, and 0 with probability (1-p).
    observations = None

    # Compute and save the empirical mean
    estimator = np.mean(observations)
    all_estimators.append(estimator)
```


5 . The following code plots the histogram of the estimators (empirical means). Run it. What type of distribution is obtained (no need to specify the exact parameters of the distribution)? Explain **briefly** what theorem from probability explains this behavior (and why).

Answer: TODO

```
In [ ]: import seaborn as sns
sns.histplot(all_estimators, bins=tosses, kde=True)
plt.grid()
```

Numerical linear algebra refresher

Reminder - Positive semi-definite matrices

A symmetric real matrix $A \in \mathbb{R}^{n \times n}$ is called positive semi-definite (PSD) iff:

$$\forall x \in \mathbb{R}^n \setminus \{0_n\} : x^\top A x \geq 0.$$

If the matrix holds the above inequality *strictly*, the matrix is called positive definite (PD).

Q3 - PSD matrices

1. Let $A \succ 0_{n \times n}$ be a symmetric PD matrix in $\mathbb{R}^{n \times n}$.

Recall that all eigenvalues of real symmetric matrices are real.

Prove that all the eigenvalues of A are strictly positive.

Answer:

TODO

2. Let $A, B \in \mathbb{R}^{n \times n}$ be two symmetric PSD matrices.

Prove or refute: the matrix $(2A - B)$ is also PSD.

Answer:

TODO

Q4 - Gradients

Define $f : \mathbb{R}^d \rightarrow \mathbb{R}$, where $f(w) = w^\top x + b$, for some vector $x \in \mathbb{R}^d$ and a scalar $b \in \mathbb{R}$.

Recall: the gradient vector is defined as $\nabla_w f = \left[\frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_d} \right]^\top \in \mathbb{R}^d$.

1. Prove that $\nabla_w f = x$.

Recall/read the definition of the [Hessian matrix](https://en.wikipedia.org/wiki/Hessian_matrix#Definitions_and_properties) (https://en.wikipedia.org/wiki/Hessian_matrix#Definitions_and_properties) $\nabla_w^2 f \in \mathbb{R}^{d \times d}$.

- 2 . Find the Hessian matrix $\nabla_w^2 f$ of the function f defined in this question.
- 3 . Is the matrix you found positive semi-definite? Explain.

Now, define $g : \mathbb{R}^d \rightarrow \mathbb{R}$, where $\lambda > 0$ and $g(w) = \lambda \|w\|^2$.

- 4 . Find the gradient vector $\nabla_w g$.
- 5 . Find the Hessian matrix $\nabla_w^2 g$.
- 6 . Is the matrix you found positive semi-definite? is it positive definite? Explain.

Answers:

TODO

```
In [ ]:
```