

MCIS6273 Data Mining (Prof. Maull) / Fall 2017 / HW1

This assignment is worth up to 20 POINTS to your grade total if you complete it on time.

Points Possible	Due Date	Time Commitment (estimated)
20	Saturday, Sep 30 @ Midnight	<i>up to</i> 10 hours

- **GRADING:** Grading will be aligned with the completeness of the objectives.
- **INDEPENDENT WORK:** Copying, cheating, plagiarism and academic dishonesty *are not tolerated* by University or course policy. Please see the syllabus for the full departmental and University statement on the academic code of honor.

OBJECTIVES

- Work with core Pandas and Scikit-Learn concepts in data, data types, representation of data and plotting data
- Explore concepts in statistical inference over real data within Scikit-Learn
- Work with data to understand distance metrics in Scikit-Learn and the impact various metrics have on the outcomes

WHAT TO TURN IN

You are being encouraged to turn the assignment in using the provided Jupyter Notebook. To do so, clone the course repository and modify the `hw1.ipynb` file in the `homework/hw1` directory. If you do not know how to do this, please ask, or visit one of the many tutorials out there on the basics of using Github and cloning repositories.

Turn in a copy of a `.ipynb` file, a PDF or Word Document to Blackboard with the answers to the questions labeled with the § sign.

ASSIGNMENT TASKS

(20%) Work with core Pandas and Scikit-Learn concepts in data, data types, representation of data and plotting data

A great deal of time doing data mining involves understanding the data in a dataset and preprocessing it in preparation for working with it in a real analysis of some sort. We will develop an understanding of:

- how to empirically perform and understand the descriptive statistics of a dataset,
- understand how to reason about data features,
- understand how to use distance metrics,
- understand binarization contexts and techniques, and
- understand how to randomly sample datasets to understand the distribution of data.

We'll be working a new much smaller dataset for this task. Please check the repository for the `bank.csv` data set. There are fewer than 1000 rows to work with in this data. You will need to use Pandas and Scikit-Learn

for all the questions here, and it is recommended you take a look at the optional cheat sheets for lecture 2 in the syllabus:

- Pandas Cookbook Github Repo
- Pandas Cheatsheet @ DataCamp.com

§ In the readings and lecture, we talked about distributions of data. Please produce the frequency distribution (histogram) for the following: (1) rural income, (2) car ownership. You will most certainly need to use the `pandas.Series.plot.hist()` method.

§ Produce the scatter plot for income vs age – it doesn't matter which is on the x and y axis. You will benefit from the `DataFrame.plot.scatter()` method.

(40%) Explore concepts in statistical inference over real data within Scikit-Learn

For questions 1-3: Just answer the question over the entire data as requested.

For questions 4 and 5: Consider this scenario ...

The bank is looking for candidates who might be good for sending mortgage loan information to – that is they are looking for people who are **not homeowners** (i.e. don't own a home), but who have the income and other characteristics of potentially good borrowers.

Your task is to look at all the data for people who have no mortgage and build the case for the profile of the borrower (from this dataset) that they should start calling, sending mail and advertising to.

Imagine you are a real-world analyst and will need to do the following:

- build the dataset that contains non-homeowners (i.e. no mortgage)
- report on the characteristics of non-homeowners based on what is being asked

§ What is the median age of a UNMARRIED, FEMALE, HOME OWNER in the SUBURBAN region?

§ Given this bank dataset what is the joint probability $\Pr(\text{sex} = \text{MALE} \wedge \text{income} \geq 50,000)$? How does this compare with $\Pr(\text{sex} = \text{FEMALE} \wedge \text{income} \geq 50,000)$?

§ What is conditional probability $\Pr(\text{car} = \text{TRUE} | \text{sex} = \text{MALE} \wedge \text{income} \geq 50,000)$?

§ If the lending requirements were a savings account and income greater than \$45,000, what about this data might make it difficult to justify any campaign at all? (**HINT:** most mortgages are 30-years in duration) Provide **concrete evidence** for your claims; you may do this in a variety of ways *including describing the statistic that leads you to your claim*.

§ If the savings account requirement was ignored and the income lowered to a minimum of \$30,000, to whom (MALE or FEMALE) and in which region (RURAL, TOWN, INNER_CITY, SUBURBAN) would the bank likely have the best success?

(40%) Work with data to understand distance metrics in Scikit-Learn and the impact various metrics have on the outcomes

A common thing to understand in a dataset is to determine some number of *nearest neighbors* to a data point. We will see this come back when we get to clustering. For now, let's explore the NearestNeighbor implementation in Scikit-Learn.

You will most certainly need to convert the categorical non-numeric data to binary features, and you may be served well by reading about preprocessing data in Scikit-Learn and also taking a closer look at `sklearn.preprocessing.LabelBinarizer` and `sklearn.preprocessing.OneHotEncoder`. You might also find this resource of value if you read only the first few slides on binarization.

§ What are the 5 nearest neighbors (the indices) to index #7, #21, #40 and #94? Your output will be a list of the tuples of the 5 closest indices and their values (e.g. [(5, 4.56356), (19, 8.83452), (233, 12.23486), ...]). Use the *minkowski (default) distance* first.

§ Use the *cosine similarity metric* and *euclidean distance metric* (e.g. invoke `NearestNeighbor(..., metric='cosine')` and `NearestNeighbor(..., metric='euclidean')`). Produce the same lists for the same indices as in #1. What are the differences in the nearest neighbor lists? Which seem to be the most similar? Provide your thoughts on why there are differences?