

HOMEWORK 1

M. Neumann

Due: THU 5 SEPT 2019 4PM

Getting Started

Update your SVN repository.

When needed, you will find additional materials for *homework x* in the folder *hw.x*. So, for the current assignment the folder is *hw1*.

SUBMISSION INSTRUCTIONS

WRITTEN:

- all written work needs to be submitted electronically in *pdf format*¹ via GRADESCOPE
- provide the following information on every page of your pdf file:
 - name
 - student ID
- start every problem on a *new page*
- **FOR GROUPS:** make a **group submission** on GRADESCOPE and provide names and student IDs for **all group members** on every page of your pdf file.

CODE:

- No code submission required for this homework.

Find instructions on how to **submit your work** to GRADESCOPE and **add your group member** on the course webpage.

PIAZZA

All course related announcements will be made there. Ask all questions about course materials, logistics, and homework assignments on Piazza using the appropriate tags.

¹ Please, **type your solutions** or use **clear hand-writing**. If we cannot read your answer, we cannot give you credit nor will we be able to meet any regrade requests concerning your writing.

GRADING RESULTS AND REGRADES

Grades will be maintained on Canvas. Grading comments will be provided via GRADESCOPE. We will make a Piazza announcement when the grading is completed. All regrade requests need to be made via GRADESCOPE **within one week** of the grade announcement.

SVN REPOSITORY

Check out your SVN repository. Find instructions on how to checkout, update, and commit to your SVN repository here: http://sites.wustl.edu/neumann/resources/cse427s_resources/

Note: Code will have to be submitted exclusively via your SVN repository. To **check if your submission to the SVN repository was successful** view your repository in a web browser by entering this url (mind browser caching):

https://svn.seas.wustl.edu/repositories/<yourwustlkey>/cse427s_fl19

PROBLEM 1: Big Data Characteristics (50%)

Describe the *characteristics* of the following Big datasets. There may not be one correct answer, so explain your decisions by **including an argument** why each of the characteristics discussed in class *is* or *is not* present.

- (a) **Access Log Data:** all HTTP logs from a popular online vendor's webserver(s) collected over the duration of the past 5 years 2013-2018

The following example shows *four* example entries of this log data:

```
2013-03-15 12:39 - 74.125.226.230 /common/logo.gif 1231ms - 2326
2013-03-15 12:39 - 157.166.255.18 /catalog/cat1.html 891ms - 1211
2013-03-15 12:40 - 65.50.196.141 /common/logo.gif 1992ms - 1198
2013-03-15 12:41 - 64.69.4.150 /common/promoex.jpg 3992ms - 2326
...
```

- (b) **Wikipedia articles:** text (no images, info boxes etc.) of all articles of the English Wikipedia and the links between them

The following example shows (part of) *one* example article with the links indicated in blue:

Big data is [data sets](#) that are so voluminous and complex that traditional [data processing application software](#) are inadequate to deal with them. Big data challenges include [capturing data](#), [data storage](#), [data analysis](#), search, [sharing](#), [transfer](#), [visualization](#), [querying](#), updating and [information privacy](#). There are three dimensions to big data known as Volume, Variety and Velocity. Lately, the term "big data" tends to refer to ...

- (c) **Chemical compounds:** (fixed) data set of all chemical compounds recorded before *today* (e.g. downloaded from the ZINC database). Such datasets are commonly used for drug design.

The following example shows (part of) the representation of *one* example chemical compound:

bonds in form of (bond ID: atom ID, atom ID):

```
(1: 2,1) (2: 14,1) (3: 3,2) (4: 4,3) (5: 12,3) (6: 3,4) (7: 5,4) (8: 6,5)
(9: 5,6) (10: 7,6) (11: 11,6) (12: 6,7) (13: 8,7) (14: 21,7) (15: 7,8) ...
```

bond types as (bond ID: bond type):

```
(1: 47) (2: 47) (3: 47) (4: 50) (5: 47) (6: 47) (7: 47) (8: 50) ...
```

atom types as (atom ID: atom type):

```
(1: 7) (2: 3) (3: 3) (4: 6) (5: 3) (6: 3) (7: 6) (8: 3) (9: 22) ...
```

property:

```
mutagenic
```

Bonds are atomID-atomID pairs indicating that two atoms are connected. Bond types indicate the type of bond (47 = single, 50 = double, ...). Atom types are for example 7=oxygen, 3=chlorine, 6=carbon, 22=hydrogen,... The property of a chemical can be either *mutagenic* or *non-mutagenic*.

- (d) Both data sets described in (a) and (b) are represented as *text*.
- What is the main difference between those data sets?
 - What does that imply for data analysis tasks such as information extraction or pattern recognition?
- (e) For the example data sets described in (a)-(c), what are the **data points** and what **data types/data structures** do we use to represent those data points programmatically?

PROBLEM 2: Bonferroni's Principle (25%)

"Be careful with what you mine in Big data – it could be random!"

This question generalizes the example of "evil-doers" visiting hotels, as in Section 1.2.3 of the MMDS book. Suppose (as described there) that there are one billion people being monitored for 1000 days. Each person has a 1% probability of visiting a hotel on any given day, and hotels hold 100 people each, so there are 100,000 hotels. However, our test for evil-doers is different. We consider a group of p people evil-doers if they all stayed at the same hotel on d different days. **Derive** the formula for the (approximated) expected number of false accusations f (that is, the expected number of sets of p people that will be suspected of evil-doing), assuming that in fact there are no evil-doers, but all people behave at random, following the conditions stated in this problem (1% probability of visiting a hotel, etc.).

Note: You may assume that d and p are sufficiently small and thus, $\binom{1000}{d} \approx \frac{1000^d}{d!}$, and similarly for p . Show your work to receive maximum credit.

Hint: you can use this table to check your formula (the values for f are rounded):

d	p	f
2	2	2.5×10^5
2	3	$0.83 \approx 1$
3	2	10^{-1}
3	3	3×10^{-14}

PROBLEM 3: The Unreasonable Effectiveness of Data (25%)

Read the article "The Unreasonable Effectiveness of Data" by Alon Halevy, Peter Norvig, and Fernando Pereira. Based on this article, answer the following questions.

- (a) Compare the *amount* of data needed when using unlabeled versus labeled/annotated data. Give an explanation why they are the same or why they are different.
- (b) Summarize the *data-based approach* described in the article. How do models look like in this approach?
- (c) What are the limits of this approach? Discuss the *amount* of data needed dependent on the *size of the feature space* (amount of features).

Reflection (Bonus Problem for 5% up to a max. of 100%)

Reflect on your homework experience! Write a paragraph of at least 50 words to express your experiences and feelings when working on this assignment. Try to answer at least 2 of the following questions:

- What did you like/dislike about the assignment and why?
- What is the most important thing you learned and why do you think so?
- What surprised you, and why?
- Assuming you could start over again (with working on the assignment), what would you do differently and why?

Why?

Hopefully, we will be able to use this data for sentiment analysis at the end of the semester! Sentiment Analysis tries to assess the *emotion* or *mood* in a text document. Essentially it can be computed by comparing the set of words in each document to an existing dictionary of positive words, negative words, and neutral words. At the end of the semester, given that we have enough data for each homework, you will perform sentiment analysis to see which assignments you and your peers regarded as "positive" and which as "negative".

Submission Instructions

Store your reflection in the `hw1_reflection.txt` file provided in the `hw1` folder in your SVN repository and commit it.

This file should only include the reflection, **no other personal information** such as name, wustlkey, etc. reflections are not graded based on the content, but solely for completion.

To submit your reflection `cd` into the `hw1` folder and run:

```
$ svn commit -m 'hw1 reflection submission' .
```

Hint: You can **check your submission** to the SVN repository by viewing https://svn.seas.wustl.edu/repositories/<yourwustlkey>/cse427s_f119 in a web browser.

Finally provide a **star rating** for the homework assignment similar as if it were a product you could review on Amazon. Why? To be able to evaluate your sentiment analysis approach you will need the ground truth. So, in addition to your textual review, take 1 minute to **submit a star rating for hw1 via this link**: https://wustl.az1.qualtrics.com/jfe/form/SV_bIRKP1xobp08yIB.

Grading - no group work!

You can only earn bonus points if you write a *meaningful reflection* of **at least 50 words** answering at least 2 of the prompted questions and provide the corresponding **star rating**. You will **not** be graded on what your reflection says and the number of stars you assign, but rather solely the completion of it.

Bonus points are given to the **owner of the repository only**. No group work!.