

LING 575K HW1

Due 11PM on Apr 5, 2021

1 Setting Up Environment on Patas [10 pts]

Install Anaconda. This is a necessary component for running the code for most assignments in this course. We have setup a conda environment for the assignment, but you will need to install anaconda in order to use that environment. These are “free points”. From your home directory, please execute the following steps:

1. `wget https://repo.anaconda.com/archive/Anaconda3-2020.11-Linux-x86_64.sh`
2. `sh Anaconda3-2020.11-Linux-x86_64.sh`

The first two lines of `run_hw9.sh` show how to now activate the environment that we have supplied with all of the necessary libraries.

2 Implementing a Vocabulary Class [30 pts]

3 Data Statement for Stanford Sentiment Treebank [35 pts]

For natural language processing applications, data plays a crucial role, since it largely shapes the resulting models and systems that are used in deployment. Emily Bender and Batya Friedman have recently been developing and advocating for the practice of “data statements”: explicit documentation of the nature and origins of datasets used in NLP. For more information, please consult their paper:

- “Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science”, *Transactions of the ACL*

For the first half of this course, we will make heavy use of the Stanford Sentiment Treebank. Before beginning to use this data in modeling, you will attempt to write a data statement for this dataset to the best of your ability given the documentation available. Afterwards, there will be space to reflect on what was missing. Please see:

- The paper presenting the data (§3 in particular): <https://www.aclweb.org/anthology/D13-1170/>
- The website for the dataset: <https://nlp.stanford.edu/sentiment/index.html>

Please answer, in a couple of sentences, the following questions.

Q1: Curation Rationale Which texts were included and what were the goals in selecting texts?

Q2: Language Variety From what language variety are the texts? You may provide a prose description and/or a BCP47 code.

Q3: Speaker Demographics Who were the producers of the texts? Information may include (if available): age, gender, native language, socioeconomic status, number of speakers.

Q4: Annotator Demographics Who were the annotators of the data? Information may include (if available): age, gender, native language, socioeconomic status, number of annotators.

Q5: Speech Situation What were the conditions of text production? Details may include whether it was written or spoken, whether it was spontaneous or not, and who the intended audience was.

Q6: Text Characteristics What are the genre and topic of the texts?

Q7: Reflections Which of these questions were hard to answer with the information provided about the dataset? Why might it be helpful to have that information more explicitly documented?

4 Submission Instructions

Submission: In your submission, include the following:

- `readme.(txt|pdf)` that includes your answers to Q1-Q6. No need to submit anything for Q7.
- Since this assignment does not require programming, there is no need to submit `hw.tar.gz`, and no need to run `check_hwX.sh` script.