

# CS 475 Machine Learning: Homework 5

## Graphical Models, Inference, and Structured Prediction

### Introduction

Due: Saturday May 2, 2020, 11:59pm

100 Points Total      Version 1.0

**Make sure to read from start to finish before beginning the assignment.**

## 1 Partner Policy

**You may work with a partner for this assignment.** If you choose to do so, you will only make one submission for the two of you on Gradescope (make sure to include your partner when you submit). We *highly* recommend that you do every part of the assignment together instead of splitting it up. For the programming assignment, you can start a Zoom session with your partner and alternate who shares their screen to facilitate pair programming. For the analytical assignment, you can work on the same Overleaf document and think through the questions together.

## 2 Introduction

Assignments in this course can consist of two parts: Analytical and Programming. **For this assignment, you will have three submissions to Gradescope:** a PDF submission for the analytical section, code for the programming assignment, and **a PDF where you will interpret the results of your code.** There will be three Gradescope assignments; your pair must submit to each assignment.

1. **Analytical:** These questions will ask you to consider questions related to the topics covered by the assignment. You will be able to answer these questions without relying on programming.
2. **Programming:** The goal of the programming assignments in this course is to learn about how machine learning algorithms work through hands-on experience. In each homework assignment you will first implement an algorithm and then run experiments with it. You may also be asked to experiment with these algorithms in a Python Notebook to learn more about how they work.

Assignments are worth various points. Typical assignments are worth 100 each, with point totals indicated in the assignment. Each assignment will contain a version number at the top. While we try to ensure every homework is perfect when we release it, errors do happen. When we correct these, we'll update the version number, post a new PDF and announce the change. Each homework starts at version 1.0.

### 3 Analytical (40 Points)

In addition to completing the analytical questions, your assignment for this homework is to learn  $\text{\LaTeX}$ . All homework writeups must be PDFs compiled from  $\text{\LaTeX}$ . Why learn  $\text{\LaTeX}$ ?

1. It is incredibly useful for writing mathematical expressions.
2. It makes references simple.
3. Many academic papers are written in  $\text{\LaTeX}$ .

The list goes on. Additionally, it makes your assignments much easier to read than if you try to scan them in or complete them in Word.

We realize learning  $\text{\LaTeX}$  can be daunting. Fear not. There are many tutorials on the Web to help you learn. We recommend using `pdflatex`. It's available for nearly every operating system. Additionally, we have provided you with the tex source for this PDF, which means you can start your writeup by erasing much of the content of this writeup and filling in your answers. You can even copy and paste the few mathematical expressions in this assignment for your convenience. As the semester progresses, you'll no doubt become more familiar with  $\text{\LaTeX}$ , and even begin to appreciate using it.

Be sure to check out this cool  $\text{\LaTeX}$  tool for finding symbols. It uses machine learning! <http://detexify.kirelabs.org/classify.html>

For each homework assignment we will provide you with a  $\text{\LaTeX}$  template. You **must use the template**. The template contains detailed directions about how to use it.

At this point, please open the file `2020_Homework5_template.pdf` to read and `2020_Homework5_template.tex` respond to the analytical questions.

### 4 Programming (60 points)

#### 4.1 Python Libraries

We will be using Python 3.7.3. We are *not* using Python 2, and will not accept assignments written for this version. We recommend using a recent release of Python 3.7.x, but any recent Python3 should be fine.

For each assignment, we will tell you which Python libraries you may use. We will do this by providing a `requirements.txt` file. We *strongly* recommend using a *virtual environment* to ensure compliance with the permitted libraries. By strongly, we mean that unless you have a very good reason not to, and you really know what you are doing, you should use a virtual environment. You can also use anaconda environments, which achieve the same goal. The point is that you should ensure you are only using libraries available in the basic Python library or the `requirements.txt` file.

#### 4.2 Virtual Environments

Virtual environments are easy to set up and let you work within an isolated Python environment. In short, you can create a directory that corresponds to a specific Python version with specific packages, and once you activate that environment, you are shielded from the various Python / package versions that may already reside elsewhere on your system. Here is an example:

```
# Create a new virtual environment.
python3 -m venv python3-hw5
# Activate the virtual environment.
source python3-hw5/bin/activate
# Install packages as specified in requirements.txt.
pip install -r requirements.txt
# Optional: Deactivate the virtual environment, returning to your system's setup.
deactivate
```

When we run your code, we will use a virtual environment with *only* the libraries in `requirements.txt`. If you use a library not included in this file, your code will fail when we run it, leading to a large loss of points. By setting up a virtual environment, you can ensure that you do not mistakenly include other libraries, which will in turn help ensure that your code runs on our systems.

Make sure you are using the correct `requirements.txt` file from the current assignment. We may add new libraries to the file in subsequent assignments, or even remove a library (less likely). If you are using the wrong assignments `requirements.txt` file, it may not run when we grade it. For this reason, we suggest creating a new virtual environment for each assignment.

It may happen that you find yourself in need of a library not included in the `requirements.txt` for the assignment. You may request that a library be added by posting to Piazza. This may be useful when there is some helpful functionality in another library that we omitted. However, we are unlikely to include a library if it either solves a major part of the assignment, or includes functionality that isn't really necessary.

In this and future assignments we will allow you to use `numpy` and `scipy`.

“Can I use the command `import XXX`?” For every value of `XXX` the answer is: if you are using a virtual environment setup with the distributed `requirements.txt` file, then you can import and use anything in that environment.

**Please verify you use the exact same command line flags we dictate.** There are always students with a typo in their command line flags, which means their code fails when we run it.

### 4.3 How to Run the Provided Framework

The framework has only one mode: inference. This is defined in `main.py`, which has the following arguments:

- `--data` This is required! Specifies the directory where the documents in the corpus are located.
- `--predictions-file` This is required! Tells the program where to store a model's topic assignments over words in the vocab.
- `--inference` This is required! Specify whether to use sampling or variational inference method. You need to implement these! The inference methods we'll use are `gibbs-sampling`, `sum-product`
- `--num-documents` Specifies the number of documents from the corpus to include within the document-word matrix. By default, 100 documents are included.
- `--top-k` Specifies the the top  $k$  words in each topic to include in predictions. By default, this is set to 10.

In addition to these meta arguments, `main` also takes some hyperparameter arguments, namely:

- `--num-topics` The number of topics to uncover in the data. By default, this is set to 10.
- `--alpha` The Dirichlet parameter for document-portions ( $\theta_d$ ). By default this parameter is set to 0.1.
- `--beta` The Dirichlet parameter for topic-word distribution ( $\phi_k$ ). By default this parameter is set to 0.01.
- `--iterations` The number of iterations to use. By default, this value is 100.

Feel free to explore values for the last four hyperparameters to try to find hyperparameters that do the best for your algorithms. Note that when we run your code, we will use predetermined values for these, to ensure consistency amongst everyone's implementation.

### 4.3.1 Running Inference

The usage for this framework is:

```
python3 main.py --inference inference_method --predictions-file predictions_file --data data_directory
```

The `inference` option indicates which inference method to use. Each assignment will specify the string argument for an algorithm. The `data` option indicates the document directory. Finally, results are saved to the `predictions-file`.

### 4.3.2 Examples

As an example, the following trains a Gibbs Sampler on the documents from political blogger Matthew Yglesias:

```
python3 main.py --inference gibbs-sampling --predictions-file my.gibbs_sampling.topics --data datasets/my/
```

## 4.4 Data

The data are provided are blog post text extracted from raw html texts from five American political blogs. The blog posts are each stored in an xml file.

The blog posts from the five American blogs are organized into the following directories:

- `my`: Matthew Yglesias
- `dk`: Daily Kos
- `cb`: Carpetbagger Report
- `rs`: Red State
- `rwn`: Right Wing News

In each directory there is a file called `common_words.txt` that is used for preparing the document-word matrix.

### 4.4.1 Creating the Document-Word Matrix

The text within each document is pre-processed before creating the document-term matrix. This involves removing stopwords, punctuation, and meta tags. Then the document-word matrix is constructed by computing the frequency of each word in each document. Finally, it is stored in a sparse `coo_matrix`<sup>1</sup> using a 32-bit integer format for space efficiency. See Section 4.1 in the programming assignment for tips on working with `coo_matrix`.

## 4.5 Existing Components

The foundations of the learning framework have been provided for you. You will need to complete this library by filling in code where you see a `TODO` comment. You are free to make changes to the code as needed provided you do not change the behavior of the command lines described above. We emphasize this point: **do not change the existing command line flags, existing filenames, or algorithm names**. We use these command lines to test your code. If you change their behavior, we cannot test your code.

The code we have provided is fairly compact, and you should spend some time to familiarize yourself with it. Here is a short summary to get you started:

- `data.py` – This contains the `build_dtm` function, reads in all the documents from the specified directory and returns the document-word matrix and vocab list. The document-word matrix is stored as a sparse matrix of 32-bit integers (and in particular as a `scipy.sparse.coo_matrix` of `np.intc`), which has `num_docs` rows and `|vocab|` columns. The index of words in the vocab list correspond to their column indices in the document-term matrix.
- `main.py` – This is the main testbed to be run from the command line. It takes care of parsing command line arguments, entering train/test mode, saving models/predictions, etc. Once again, **do not change the names of existing command-line arguments**.
- `models.py` – This contains a `Model` class which you should extend. We have implemented the LDA class for you, which calls the inference methods you will write and infers the words associated with each topic in the `predict` method. There is also an `Inference` class which is extended into the `GibbsSampling` and `SumProduct` classes you will implement. At the very least, you should have `initialize` and `inference` methods, but you are encouraged to add other methods as necessary.

## 4.6 Evaluation

We will evaluate your code using automatic tests, the graphs produced in the jupyter notebook `visualizing_results.py` and your written answers in the jupyter notebook. We have provided the code to produce the necessary graphs and figures. Be sure these graphs are visible in the PDF you submit from `visualizing_results.py`.

## 4.7 Code Readability and Style

In general, you will not be graded for code style. However, your code should be readable, which means minimal comments and clear naming / organization. If your code works

---

<sup>1</sup>Documentation: [https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.coo\\_matrix.html](https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.coo_matrix.html)

perfectly then you will get full credit. However, if it does not we will look at your code to determine how to allocate partial credit. If we cannot read your code or understand it, then it is very difficult to assign partial credit. Therefore, it is in your own interests to make sure that your code is reasonably readable and clear.

## 4.8 Code Structure

Your code must support the command line options and the example commands listed in the assignment. Aside from this, you are free to change the internal structure of the code, write new classes, change methods, add exception handling, etc. However, once again, do not change the names of the files or command-line arguments that have been provided. We suggest you remember the need for clarity in your code organization.

## 4.9 Grading Programming

The programming section of your assignment will be graded using an automated grading program. Your code will be run using the provided command line options, as well as other variations on these options (different parameters, data sets, etc.) The grader will consider the following aspects of your code.

1. **Exceptions:** Does your code run without crashing?
2. **Output:** Some assignments will ask you to write some data to the console. Make sure you follow the provided output instructions exactly.
3. **Accuracy:** If your code works correctly, then it should achieve a certain accuracy on each data set. While there are small difference that can arise, a correctly working implementation will get the right answer.
4. **Speed/Memory:** As far as grading is concerned, efficiency largely doesn't matter, except where lack of efficiency severely slows your code (so slow that we assume it is broken) or the lack of efficiency demonstrates a lack of understanding of the algorithm. For example, if your code runs in two minutes and everyone else runs in 2 seconds, you'll lose points. Alternatively, if you require 2 gigs of memory, and everyone else needs 10 MB, you'll lose points. In general, this happens not because you did not optimize your code, but when you've implemented something incorrectly.

## 4.10 Knowing Your Code Works

How do you know your code really works? That is a very difficult problem to solve. Here are a few tips:

1. Use Piazza. While **you cannot share code**, you can share results. We encourage you to post your results on topic assignments using different inference algorithms. A common result will quickly emerge that you can measure against.
2. Output intermediate steps. Looking at final predictions that are wrong tells you little. Instead, print output as you go and check it to make sure it looks right. This can also be helpful when sharing information on Piazza.
3. Debug. Find a Python debugger that you like and use it. This can be very helpful.

### 4.11 Debugging

The most common question we receive is “how do I debug my code?” The truth is that machine learning algorithms are very hard to debug because the behavior of the algorithm is unknown. In these assignments, you won’t know ahead of time what accuracy is expected for your algorithm on a dataset. This is the reality of machine learning development, though in this class you have the advantage of your classmates, who may post the output of their code to the bulletin board. While debugging machine learning code is therefore harder, the same principles of debugging apply. Write tests for different parts of your code to make sure it works as expected. Test it out on the easy datasets to verify it works and, when it doesn’t, debug those datasets carefully. Work out on paper the correct answer and make sure your code matches. Don’t be afraid of writing your own data for specific algorithms as needed to test out different methods. This process is part of learning machine learning algorithms and a reality of developing machine learning software.

At this point, please open the file `homework5_coding.pdf` to read the programming assignment. You will modify the Python code provided.

## 5 What to Submit

In each assignment you may submit two different things.

1. **Submit your code (.py files) to gradescope.com as a zip file. Your code must be uploaded as code.zip with your code in the root directory.** By ‘in the root directory,’ we mean that the zip should contain `*.py` at the root (`./*.py`) and not in any sort of substructure (for example `hw5/*.py`). One simple way to achieve this is to zip using the command line, where you include files directly (e.g., `*.py`) rather than specifying a folder (e.g., `hw5`):

```
zip code.zip *.py
```

A common mistake is to use a program that automatically places your code in a subfolder. It is your job to make sure you have zipped your code correctly.

We will run your code using the exact command lines described earlier, so make sure it works ahead of time, and make sure that it doesn’t crash when you run it on the test data. A common mistake is to change the command line flags. If you do this, your code will not run.

Remember to submit all of the source code, including what we have provided to you. We will include `requirements.txt` but nothing else.

2. **Submit your evaluation of your code to Gradescope.** Save your notebook `visualizing_results.ipynb` as a PDF and upload it to Gradescope.
3. **Submit your writeup to Gradescope. Your writeup must be compiled from L<sup>A</sup>T<sub>E</sub>X and uploaded as a PDF.** The writeup should contain all of the answers to the analytical questions asked in the assignment. Make sure to include your name in the writeup PDF and to use the provided L<sup>A</sup>T<sub>E</sub>X template for your answers following the distributed template. You will submit this to the assignment called “Homework 5: Graphical Models: Analytical”.

You will need to create an account on [gradescope.com](https://www.gradescope.com) and signup for this class. The course is <https://www.gradescope.com/courses/70713>. Use entry code MK8J8N. **You must either use the email account associated with your JHED, or specify your JHED as your student ID.** See this video for instructions on how to upload a homework assignment: [https://www.youtube.com/watch?v=KMPoby5g\\_nE](https://www.youtube.com/watch?v=KMPoby5g_nE).

## 6 Questions?

Remember to submit questions about the assignment to the appropriate group on Piazza: <https://piazza.com/jhu/spring2020/601475>.