HOMEWORK 2 SENTENCE CLASSIFICATION FROM SCRATCH¹

CS 678 ADVANCED NATURAL LANGUAGE PROCESSING (SPRING 2023)

https://nlp.cs.gmu.edu/course/cs678-spring23/

OUT: Jan 23, 2023

DATA EXCHANGE: March 3, 2023
MODEL DEPLOYMENT: March 13, 2023
REPORT DUE: March 16, 2023

Based on assignments created by Mohit Iyyer and Graham Neubig.

Your name: Abhishek Jallawaram

Your GID: <u>G01373042</u>

IMPORTANT: The homework is accompanied by a python notebook. After copying this notebook to your Google Drive or One Drive, please paste a link to it below. To get a publicly-accessible link, hit the Share button at the top right, then click "Get shareable link" and copy over the result. Alternatively, you can upload the completed .ipynb file along with your completed .pdf report to Gradescope.

If you fail to do this, you will receive no credit for this homework!

Your Notebook solution:

https://drive.google.com/file/d/1MVyreGXOvHwc2cIfyhtnN5jya56vQZJq/view?usp=sharing

Graded Questions: 100 points Bonus Points: 15 points/student Total Points Available: 115/100

Additional Notes:

- Upload the whole report PDF (including all pages, also if you're scanning it).
- Some questions require writing Python code and computing results, and the rest of them have written answers. For coding problems, you will have to fill out all code blocks that say YOUR CODE HERE.
- For text-based answers, you should replace the text that says "Write your answer here..." with your actual answer.

¹Compiled on Friday 17th March, 2023 at 23:49

1 Introduction

So far in your machine learning classes, you may have experimented with standardized tasks and datasets that were provided and easily accessible. However, in the real world, NLP practitioners often have to solve a problem from scratch, which includes gathering and cleaning data, annotating the data, choosing a model, iterating on the model, and possibly going back to change the data. For this assignment, you'll get to experience this full process.

Please note that you'll be building your own system end-to-end for this project. While we provide starter code (see accompanying notebook), you're welcome to try anything beyond what we provide. You will have to collect your own data and train a model of your choice on the data. We will be releasing an unlabeled test dataset a few days before the assignment deadline, and you will run your already-constructed system over this data and submit the results. We also ask you to follow several experimental best practices, and describe the result in your report.

The full process will require the following steps:

- 1. Understand the task specification
- 2. Collect raw data
- 3. Annotate test and training data for development
- 4. Train and test models using this data
- 5. "Deploy" your System
- 6. Write your report (using this template)

Each of these are outlined in the sections below.

NOTE: This assignment will require time for each step, so make sure to plan ahead!

2 Task Specification

For this assignment, the scenario is that you want to help a social scientist friend, who is studying the *framing* of social issues in media and other documents.

What is framing? Framing refers to the selection of some aspects of a perceived reality and make them more prominent in a communicating text such as a news article and a Tweet. Examples²:

- When some news media emphasize the mental illness of gun shooters over other other aspects of gun violence in covering the issue, this is framing.
- When you choose to purchase a yogurt product that is advertised "90 percent less fat" rather than the one saying "10 percent fat," this is when the framing effect occurs.

Why does it matter? Quoting from one of the seminal studies of the problem?:

In a polarized media environment, partisan media outlets intentionally frame news stories in a way to advance certain political agendas. Even when journalists make their best efforts to pursue objectivity, media framing often favors one side over another in political disputes, thus always resulting in some degree of bias. Hence, a news framing analysis is helpful because it not only tells us whether a news article is left- or right-leaning (or positive or negative), but also reveals how the article is structured to promote a certain side of the political spectrum.

²Taken from http://www.openframing.org/home.html

In communication research, manual identification of media frames is a challenging task due to the large amount of media data in this news-saturated environment. More importantly, there is a high level of complexity in framing analysis that often requires a careful investigation of nuances in news coverage, which is time-consuming.

Hence, a NLP solution that automates media framing identification would immensely help social scientists and other analysts.

The Task Specifically, we will ask you to identify the framing of a given parapgraph/sentence in several languages. Effectively, this is a straight-forward text classification problem in a multilingual setting.

Input: The input to the model will be a text file with one paragraph/sentence per line. The text will already be tokenized using the spacy tokenizer,³ and you should not change the tokenization. An example of the input looks like this (notice the spacing of the punctuation):

Some economists say that immigrants , legal and illegal , produce a net economic gain , while others say that they create a net loss .

Output: The output of your model should be a .tsv file, with one sentence per line, a tab, and then the corresponding label.

Please refer to the provided example data files for a specific example.

The Annotation Standard Social scientists have created a widely accepted list of 15 cross-cutting framing dimensions, such as economics, morality, and politics. These were originally developed by ? and are termed the "Policy Frames Codebook". The paper describing the codebook is publicly available and the codebook is described in depth in Section 3.2.4 You should read carefully and make sure you understand the scheme. The framing dimensions are also listed in the Table to the right.

Capacity and resources

Morality Fairness and equality Legality, constitutionality and jurisprudence Policy prescription and evaluation Crime and punishment Security and defense Health and safety

Quality of life Cultural identity

Public opinion

Political

External regulation and

costs, benefits, or other financial implications availability of physical, human or financial resources, and capacity of current systems

religious or ethical implications

balance or distribution of rights, responsibilities, and resources rights, freedoms, and authority of individuals, corporations, and gov-

discussion of specific policies aimed at addressing problems

effectiveness and implications of laws and their enforcement threats to welfare of the individual, community, or nation health care, sanitation, public safety

threats and opportunities for the individual's wealth, happiness, and well-being

traditions, customs, or values of a social group in relation to a policy attitudes and opinions of the general public, including polling and de-

mographics considerations related to politics and politicians, including lobbying,

elections, and attempts to sway voters international reputation or foreign policy of the U.S

Domains Your social scientist friend is planning on eapplying your system on news data related to two other any coherent group of frames not covered by the above categories

• immigration

social issues:

same-sex marriage

Languages The seed data are all in English. However, your social scientist friend cares a lot about multilinguality, and ideally your system should be able to provide labels for sentences in the following languages: English (eng), Mandarin Chinese (cmn), Hindi (hin), Telugu (tel), Nepali (nep), Bengali (ben), Greek (ell), German (deu), Swahili (swa), as well as two surprise languages that will be revealed later (see Timeline).

³https://spacy.io/api/tokenizer/

⁴https://kilthub.cmu.edu/ndownloader/files/11903384

3 Data Collection [35 pts]

Seed Data To help you get started, we provide 250 sentences in English, annotated with their frames. All sentences are on the immigration domain. You may study them to understand the task and the annotation schema, and you may use them however you wish (e.g. as training or dev data or as a starting point for data augmentation).

First Round of Data Labeling In this first round, you will collect/create a labeled dataset for our task of **at least 150 sentences**. Notice that this is a *minimum* number of data you need to annotate, but annotating more is allowed.⁵

The purpose of this part is to familiarize everyone with the annotation process and the codebook. Also, as any good data scientist would tell you, a little bit of experience performing the task you are trying to automate is invaluable!

It is up to you to find (or create) sentences/paragraphs relevant to the task. We provide some ideas, pointers, and considerations about this part below.

- You could create a list of keywords related to the two domains, search the web for relevant articles, and scrape sentences from there.
- You could perhaps write some your own sentences (but be careful of your own biases!)
- You could perhaps use a writing assistant or a chatbot (e.g. ChatGPT) and create sentences relevant to a domain with a certain framing.
- There are two target domains and 15-ish output classes. How will you make sure that all classes and domains are represented in the data you will annotate? Remember that a neural model will only be able to predict something that it has been trained on!
- Can you use the seed data in any way?
- In which language(s) will you search for, create, or annotate data?
- What would you do if more than one labels seem appropriate for a sentence?
- Be creative!

Package your annotations in a file named first_data.tsv. The .tsv file should include the following five columns:

- text: the sentence that will be the input to the model
- language: the language of the sentence. If codemixed, use a comma-separated list of languages. Use either language names or language ISO-3 codes.
- label: the framing label of the sentence. Use either the label name (e.g. "economic") or the label ID according to the annotation schema.
- source: the source of the sentence. If scraped, provide the URL. If written by you or someone else, provide the author (or system) name. If created in another way, provide an explanation in your report.
- annotator: this should be your GMUID
- 1. (10 points) Describe the source of your unlabeled data, why you chose it, and what kind of sentence selection process you used (if any) to choose the (minimum) 150 sentences for annotation. If you chose to annotate more data, please justify your choices on the source and amount of data you create. To receive the question credits, you must provide the above-described file with your data.

⁵A good model will likely need more data that this, but be careful to not waste too much of your time annotating data though!

Data Source Description

English articles on immigration and same-sex marriage were web scrapped using Duckduckgo & GoogleSearch along with BeautifulSoup and requests to attain the top 10 websites for each label. Text was scrapped from these websites to collect the required sentences.

Pre-processing was performed to clean the sentences and 170 were selected for the annotator task. Another 1000 sentences were used for the first seed data which were labelled manually during the course of the task.

Web_Scrap_URL.txt has been provided which contains a few sample URL used to collect the sentences.

Second Round of Data Labeling In this second round, you will collect annotations on another 150 sentences from at least two classmates. Everyone in this class will need to both create their own dataset and also serve as an annotator for (at least) two other classmates. In total, each student will annotate a minimum of 450 sentences (150 in the first round, and then 150+150 in the second round). The data exchange between annotators should happen at the latest by March 3.6

The idea now is that you will select another batch of 150 sentences to be annotated, but this time each sentence should be annotated by at least 2 annotators (that are not you).

In particular, you should:

- 1. Create another batch of 150 unlabeled sentences (minimum).
- 2. Find two classmates willing to label 150 sentences each (use the Piazza "search for teammates" thread if you're having issues finding labelers).
- 3. Send the data to your annotators in a way that they can easily annotate the data (e.g., a spreadsheet, Google form, or by using Labelstudio).
- 4. Collect the annotations from your annotators and sanity check the data for basic cleanliness (are all examples annotated? are all labels allowable ones?)

Notes:

- 1. How will you select the data to share for annotation? What language(s) will they be in?
- 2. Plan ahead on your commitments/arrangements for data annotation!
- 2. (10 points) Write down the names and emails of the two classmates⁸ who will be annotating your data below. Once they are finished annotating, create two .tsv files (annotator1.tsv and annotator2.tsv)

⁶This is a **minimum**. You are allowed to annotate more than 300 sentences in this second round, but be mindful of your time.

⁷Or three classmates willing to label 100 sentences each. Or 5 classmates willing to label 60 sentences each. You get it.

⁸Feel free to modify if using more than two annotators.

that contain each annotator's labels. Each file should have three columns with headers text, label, and GMUID respectively. You will include these files in your submission when you're finished with this homework.

Annotator 1

Name: Stephen Dias Email: sdias3@gmu.edu GMUID: G01387625

Annotator 2

Name: Aditya Limbekar Email: alimbekar@gmu.edu GMUID: G01384408

Annotator 3

Name: Navneeth Rao Balguri Email: nbalguri@gmu.edu GMUID: G01398077

Annotator 4

Name: Teja Sajja Email: tsajja@gmu.edu GMUID: G01395837

Annotator 5

Name: Aditya Singh Email: asingh89@gmu.edu GMUID: G01381857

Please contact the course staff with any issues regarding the data exchange. Email us if e.g. you committed to provide data for a classmate but cannot do so, or if a classmate committed to provide data for you but has not done so. We will try to resolve this in the best way possible. But the success of this homework relies on everyone behaving professionally: set and meet deadlines!

3. (10 points) Now, compute the inter-annotator agreement between your annotators. In the provided code cell in the accompanying notebook, read the data from the two files and compute both the raw agreement (% of examples for which both annotators agreed on the label) and the Cohen's Kappa. Feel free to use implementations in existing libraries. After you're done, report the raw agreement and Cohen's κ scores below:

Raw Agreement	Cohen's Kappa					
81.034483	0.793150					

If you're curious, Cohen suggested the Kappa result be interpreted as follows: values ≤ 0 as indicating no agreement and 0.01–0.20 as none to slight, 0.21–0.40 as fair, 0.41–0.60 as moderate, 0.61–0.80 as substantial, and 0.81–1.00 as almost perfect agreement.

4. (10 points) To produce your final dataset, you will need to **aggregate** the annotations from your annotators. That means that for cases where they disagree, you will need to choose a single label. Use any method you like *other than random label selection* to perform this aggregation (e.g., have the two annotators discuss each disagreement and come to consensus, use Piazza to create a poll on a label, or choose the label you agree with the most). You will need to create a file aggregated_data.tsv, with the same format as your first_data.tsv.

Below, describe your aggregation strategy and why you chose it. To receive the question credits, you

⁹https://en.wikipedia.org/wiki/Cohen%27s_kappa

¹⁰e.g., sklearn: https://scikit-learn.org/stable/modules/generated/sklearn.metrics.cohen_ kappa_score.html

must provide the above-described file with your data.

Aggregation Strategy

The initial raw agreement and Cohen Kappa was 0.1. So 3 more annotators were requested to perform the annotation. The raw agreement and Cohen Kappa for the combinations were calculated to select the best two annotated files. The other 20% of the files which were mismatched were selected using a majority rule between the 5 annotators and the final aggregated file required for classification was generated using this approach. Odd number of annotators were selected to make the voting easier.

5. (5 points) Based on what you learned about the task and your data, you might want to update your labels for the first_data.tsv file (remember, noone has "checked" this, you were the only one who provided annotations for this portion of the data). If you do so, provide the updated file as updated_data.tsv with the same format.

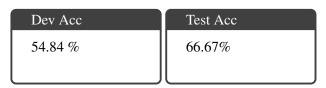
Last, combine the files from the two rounds of annotations into a single file, named final_data.tsv with the same format as the others. To receive the question credits, you must provide this file with your data.

4 Model Training and Testing [25 pts]

The next step is to train a text classification model and try to get an estimation of how well it would perform, based on your data. We provide code to do this in the accompanying notebook. You will need to upload the data to the Colab session (if you use Colab, as recommended, just click the folder icon in the sidebar to the left of the screen).

Follow the steps in the notebook to connect to a GPU, install the Huggingface transformers library, read-in the data, create train/dev/test splits, set hyperparameters (feel free to play with these!), and then fine-tune your model.

6. (5 points) After you're satisfied with your hyperparameters (i.e., you're unable to achieve higher validation accuracy by modifying them further), it's time to evaluate your model on the test set! Run the corresponding cell in the accompanying notebook and report below the validation and test set accuracy you obtain.



7. (10 points) Congratulations! You've now gone through the entire fine-tuning process and created a model for your downstream task. Now, describe your hyperparameter selection process in words. If you based your process on any research papers or websites, please reference them. Why do you think the hyperparameters you ended up choosing worked better than others? Also, is there a significant discrepancy between your test and validation accuracy? Why do you think this is the case?

Modeling Decisions

For hyperparameter selection, I started with the Hugging Face's Transformers library's default parameters for the BERT model and AdamW optimizer. Firstly, according to the task and the volume of the dataset, the parameters were updated..

The learning rate was set to 2e-5, which is a standard setting for fine-tuning BERT on small to medium-sized datasets. The learning rate was considered within the range of 1e-4 to 1e-6. The eps was considered within the range of 1e-7 to 1e-8. The number of epochs was initially considered as 25 then they were adjusted based on the total loss observed for the system.

There was a marginal difference between the accuracy of val and test sets observed when the model was run multiple times. This is expected since the dataset contains 15 labels and during the split the val or the test set may contain labels which might not be available in the train set. This might cause some misclassification as discussed further in error analysi.

8. (10 points) Next, perform an *error analysis* on your model. This is good practice for your final project. Write some code in the provided code cell to print out the text of up to five test set examples that your model gets **wrong**. If your model gets more than five test examples wrong, randomly choose five of them to analyze. If your model gets fewer than five examples wrong, please design five test examples that fool your model (i.e., *adversarial examples*). Then, provide a qualitative analysis of these examples. See if you can figure out any reasons for errors that you observe, or if you have any informed guesses (e.g., common linguistic properties of these particular examples). Does this analysis suggest any possible future steps to improve your classifier?

Error Analysis Discussion

Based on the accuracy and misclassified sentences for both the validation and test sets, it indicates that the model is experiencing issues in accurately categorizing sentences with complex syntax and semantics. Misclassified sentences, which have longer and more complex sentence structures. It's also likely that the model is experiencing issues in accurately distinguishing words with various meanings, particularly in context.

Increase the size of the training data set, can provide the model with more examples of complex sentence structures and enhance its ability to accurately categorize them. Pre-processing the web scrapped sentences and tokenizing them properly can reduce the misclassification in our scenario.

5 Data Augmentation [20 pts]

So far we have worked with whatever limited data you and your classmates were able to annotate in a small amount of time. But remember that in the end we want to deploy the model on a hidden test set that will include examples from 11 languages.

One approach is to use whatever data you have already created and perform *data augmentation*. There are several approaches for this, and we invite you to get creative with this! To help you get started, we provide some ideas below, but don't feel obliged to use any or all of them – you do need to do *something* though.¹¹

Some ideas to get you started:

- Use automatic **translation** to translate your data in another language. For example, if you have annotated data in English, then you can translate the sentences into e.g. French. You can then use the same labels for the French sentences as the ones the corresponding English sentences had. Tada now you have annotated data in French! We provide a helper notebook with some functions implemented that will allow you to easily query Google Translate and obtain translations.

 Food for thought: Which languages would you translate into? The target test languages? Others? What if training on some e.g. Spanish data would also help with some of our target languages? But what is training on some e.g. Spanish data interferes with learning and you end up performing worse on some target languages? In case you're curious, both of these are possible! The first is called "cross-lingual transfer learning", and the second is called "cross-lingual negative interference".
- Use **paraphrasing**. You can paraphrase the sentences that you currently have annotated presumably the framing (and hence the label) shouldn't change! There's a multitude of models and approaches for doing this:
 - A simple example for English paraphrasing with the Pegasus model is here. 12
 - Another way of paraphrasing is through round-trip translation. Let's say you have an English sentence. You can translate it into French, and then translate that French sentence back into English. It is very likely that you will now have a paraphrase of the original English sentence! The choice of the pivot language (e.g. French or Chinese or somethine else) could also lead to a different paraphrase.
- Label propagation. Instead of generating paraphrases, you might be able to find data that are similar or comparable to the ones you have already annotated. This would allow you to *propagate* the labels from your annotated examples to the unannotated ones. Or, you could make some assumptions about sentences with some careful data selection that will allow for label propagation. For example, imagine you have found long news articles with several sentences on the immigration topics, but you have only annotated a sample of these sentences could you perhaps propagate the annotations to the rest of a paragraph's sentences or to all sentences in the article? Or maybe you could only annotate the article's title and propagate the label to all sentences in the article. Or you could try something else!
- The above are just some of the ideas that the instructors could come up with get creative!

Compile all your augmented data in a file named augmented_data.tsv with a format similar to the other datasets. Please provide an indication of the augmentation process used for each example in the source column (e.g. you may use "eng—swa" to indication translation, "eng—swa—eng" to indicate round-trip translation, "paraphrasing-[model]" for paraphrasing, etc).

¹¹We suspect that some combination of these ideas will work well.

 $^{^{12}}$ https://www.thepythoncode.com/article/paraphrase-text-using-transformers-in-python]

Load your augmented dataset to the notebook, and now train and evaluate your model. Does your model performance improve on your test set? (make sure to keep the same train/dev/test set on your original data – you should only augment your train (and perhaps dev) data, not the test data!)

9. (20 points) Describe your data augmentation approach. Did it result in improvements on your test set? Why do you think this happened?

Data Augmentation

My data augmentation strategy was to translate the original data into 11 different languages, therefore enhancing the size of the training set. This was accomplished by employing google translator model. The translated data was then integrated to the original dataset to form a larger training set. When I tested my model on the test set, however, I discovered that the accuracy had actually dropped as compared to when I had only trained on the original data.

After some thought, I believe there might be various causes for this. Firstly, the translations may not have captured the complexities of the original data properly, resulting in noise and confusion in the enhanced dataset. Secondly, the translated data may have introduced biases or irregularities that were not evident in the original data, thereby affecting the performance of my model.

6 Model Deployment [18 pts]

The final "deployment" of your model will consist of running your (best) model over a private test set (text only) and submitting your results to us. You should try to finish building your systems before this set is released, and basically not rely on these data for model training or testing. The test set will be released shortly (3-4 days) before the final submission deadline.

TODO: Can we use Explainaboard for this?

Several Bonus points are available (up to 15 per person) based on your models' performance on this hidden test set. See "Grading" section.

7 Timeline and Effort [2 pts]

Also, please complete below your effort for this assignment, taking all stages into account. Report the approximate time spent total, in hours (integer).

Effort (in hours)
100 hrs

Suggested Timeline You have a total of 4 weeks for this assignment (although nothing stops you from starting this assignment early). We suggest that you keep track of your process and fill in the required parts of the report as you go. **Do NOT leave this for the last minute!** Our suggested timeline:

- week 1: understand the annotation schema, study the seed data, find data and perform the first round of annotation
- week 2: prepare data for second round of annotation, perform second round of annotation (for other students), try out existing code
- week 3: collect data, resolve conflicts, run first iteration of models, explore data augmentation
- week 4: [test set published, surprise languages revealed] explore model hyperparameters, deploy final model (aka run on test set), finalize report.

8 Grading

Beyond the 100 total credits available for completing all above steps in the assignment, we will assign bonus points according to the following guide. Remember than 15 extra points in this assignment correspond to 3 credits in your overall class grade, so this could be the difference between e.g. an A-and an A!

Bonus points will be applied greedily. The maximum bonus for each student is 15 credits.

Bonus points scheme:

- 1. +8 for the best performing system over the blind test set, averaged across all languages, +6 for the second best, +4 for the third best
- 2. +7 for the most *equitable system* across all languages, as measured by the standard deviation across the averages across languages (lower std is better). The system has to be in the top 50% of all submissions to be eligible. +5 for the second most equitable system, +3 for third best.
- 3. +5 for the best performing system on the surprise test languages, if it's different than the top system in criterion 1 above. +4 for the second best performing system in the surprise test languages, +2 for the third best.
- 4. +4 for the best performing system in each non-surprise test language, if it's different than the top system in criterion 1 above. +2 for the second best performing system in each test language. +60

points total

- 5. +2 for each system that is Pareto optimal, in terms of effort reported and downstream performance
- 6. +5 for the submission that provides the most hand-annotated data, +3 for second most, +2 for third most
- 7. +4 for the submission that provides hand-annotated data in most languages, +2 for second most (to be shared with up to 3 submissions. If more, this bonus does not apply)
- 8. +4 for the submission that provides the largest final augmented dataset, +3 for second largest, +2 for third

9 Collaboration Questions

Please answer the following:

1.	Did you	receive	any help	whatsoever	from	anyone	in s	solving	this	assignment	(beyond	data	annota-
	tion)?												

Yes / No.

- If you answered 'yes', give full details:
- (e.g. "Jane Doe explained to me what is asked in Question 3.4")
- 2. Did you give any help whatsoever to anyone in solving this assignment (beyond data annotation)? Yes / No.
 - If you answered 'yes', give full details:
 - (e.g. "I pointed Joe Smith to section 2.3 since he didn't know how to proceed with Question 2")
- 3. Did you find or come across code that implements any part of this assignment? Yes / No. (See policy on "found code")
 - If you answered 'yes', give full details:
 - (book & page, URL & location within the page, etc.).