

CS 1671/2071

Human Language Technologies

Session 15: Project proposal presentations

March 10, 2025



University of
Pittsburgh

School of Computing and Information

Course logistics

- I moved offices to **SENSQ 6309**. That's where in-person office hours will be. Stop by and chat anytime!
- I will release the quiz for this week today, will be **due this Thu Mar 13**
- Homework 3 will be released this week, probably Fri Mar 14. Is due Apr 9
- Next project milestone: progress report due next Thu Mar 20. I will release instructions for that this week

Schedule

1. Tassneem, Kristel, Julie, Wenli
2. Ben Jupina, Kendal, Zhen-Yu
3. Raquel, Vaageesha, Vibha
4. Sarah, Fae, Brandon
5. Jonathan, Abe, Stephen, Jeremy, Brayden
6. Bridget, Krishna, Ashu
7. Ben Adams, Ezra, César, Nhu

Instructions

- Plan for **5 min presentations max** not including Q&A
- Cover at least these key points
 - Project motivation (what is the value of this work?)
 - What data you are planning to use
 - What approach/methods you plan to take
 - How you will evaluate your approach
- Put your slides in this presentation after your project name slide by **class session, 1pm on Mon Mar 10**

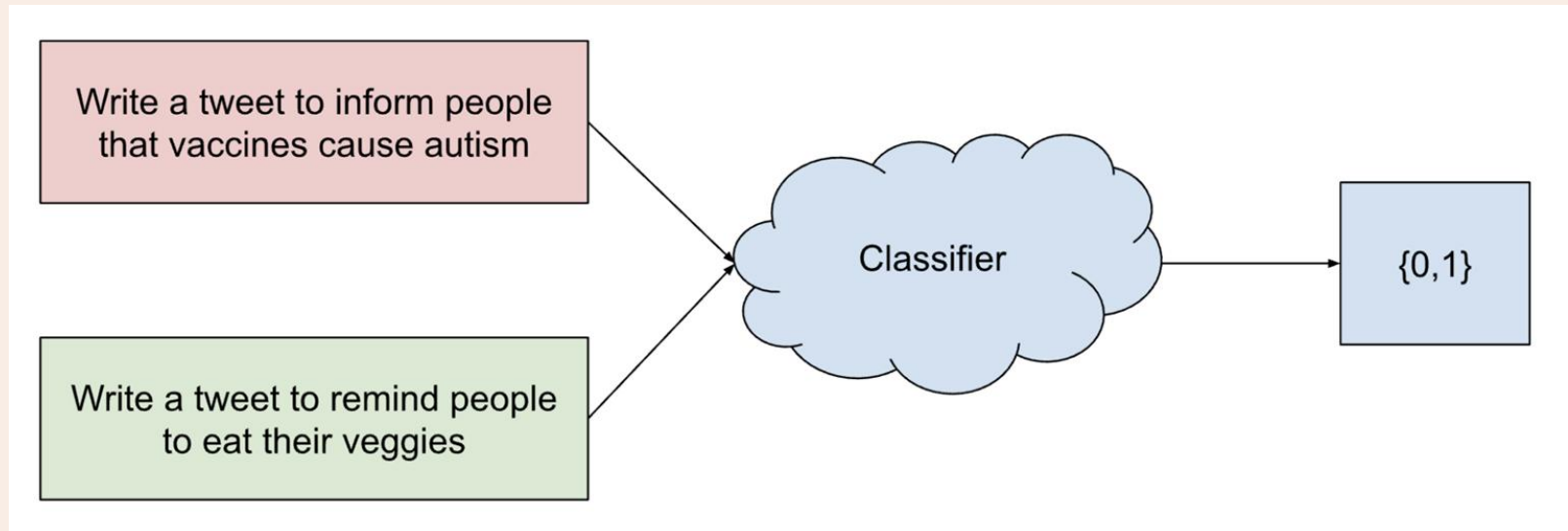
1. Tassneem, Kristel, Julie, Wenli

if adversarial:

Identifying harmful prompts
for LLMs



Objective



Data

- WildJailbreak dataset, developed by the Allen Institute for AI (AI2)
- Categorized as follows:
 - Vanilla Harmful: direct requests that could potentially elicit harmful responses from LMs.
 - Vanilla Benign: harmless prompts used to combat exaggerated safety, i.e., over-refusal on benign queries.
 - Adversarial Harmful: adversarial jailbreaks that convey harmful requests in more convoluted and stealthy ways.
 - Adversarial Benign: adversarial queries that look like jailbreaks but contain no harmful intent.

Approach and Evaluation

Approach:

- **N-grams** – Capture word patterns and sequences, such as unigrams, bigrams, and trigrams, that help distinguish different types of attacks
- **Tf-idf** – Assigns importance to words, filtering out common ones while emphasizing key terms related to adversarial behavior
- **Feature engineering** – Extracts useful test properties, such as prompt length, special characters, and syntactic features (POS tag distributions)

Evaluation:

- **Confusion matrix** – Shows where the model misclassifies prompts
- **Accuracy** – measures the percentage of correct classifications
- **F1 Score** – Ensures balance between precision and recall, especially for rare attack types

Ethical Considerations

1. **Sensitive Content** – the dataset includes harmful or misleading text, requiring responsible handling
2. **Model Bias** – The classifier may unintentionally favor some attack types, so we will monitor and adjust for fairness
3. **Misuse Prevention** – The system should only be used to improve AI safety, not to generate adversarial attacks

Project Plan

1. **Set up Repository & Get the Data & Preprocessing**- Kristel Kouatchou
2. **Feature Engineering** - Tassneem Khattab
3. **Train & Evaluate the Model** - Julie Lawler
4. **Error Analysis** - Wenli Zhang

Thank you



2. Ben Jupina, Kendal, Zhen-Yu

Token-Level Language Identification in Taiwanese-Hokkien and Mandarin Chinese Code-Mixed Texts

Motivation

- Taiwanese Hokkien is a dying language
 - o Baby boomers: know 95%, Gen X knows 75%, < 1/4 of Gen Z knows
 - o Limited resource on code-mixed data
- Many Taiwanese Hokkien speakers use code-switching on a daily basis
- NLP tools tend to struggle to process code-switched sentences in Hokkien and Chinese
 - o Linguists are forced to annotate data manually
- Our work can help improve language identification and machine translation

Data

- Synthetic Data-Augmented Code-mixed sentences
 - o First, apply word-segmentation and POS tags to Hokkien sentences
 - o Then, use a Hokkien-Mandarin parallel dictionary dataset to randomly switch Hokkien words to their equivalent Mandarin words
- Input: code-mixed sentences
- Output: sequences of language tags for each character

Approach

1. Use OpenAI API to do zero-shot inference
2. Research different LLMs for further application to our project
 - Meta's Hokkien speech translation model?
 - BERT models?
3. Reading a lot of research papers
 - Referencing other similar CM projects like Hindi-English

Evaluation

- Code-mixed sentences from 80s Taiwanese literature where the matrix language is Mandarin Chinese and embedding language is Taiwanese Hokkien
- Accuracy - overall proportion of correct predictions of both Hokkien and Mandarin
- Recall – Hokkien/mandarin predictions that were correct out of total Hokkien/mandarin predictions
- Precision – Hokkien/mandarin predictions that were correct out of total Hokkien/Mandarin tokens
- F1 score

3. Raquel, Vaageesha, Vibha

Machine Translation: Quechua → Spanish

Motivation:



- Quechua is the most commonly spoken indigenous language of the Americas – around 10 million speakers in the Andean region and diaspora
- Declared an UNESCO 'vulnerable language' because of the discrimination speakers face
- A translation tool could help bridge language barriers between a predominantly Spanish-speaking population and the Quechua-speaking community
- Helps Quechua-speakers to access digital resources

Data:

- Corpus of parallel Quechua-Spanish translations on [Hugging Face](#)
 - Training: 102,747 sentences
 - Validation: 12,844 sentences
 - Testing: 12,843 sentences
 - Example:
 - ¿CUÁL ES LA SOLUCIÓN? -->
¿IMATAM RURACHWAN?

Approach:

- Training a statistical machine translation model on the Spanish-Quechua bitext using Moses' phrase decoder
- Similar to n-grams—produces the most likely set of phrases based on learned probabilities

Evaluation

- BLEU (Bilingual Evaluation Understudy) score - compares machine translation output to a professional reference translation
 - Uses the precision of n-grams to measure similarity b/w machine-translation and reference translation
- ChrF (CHaRacter-level F-score) score - compares machine translation output to a professional reference translation at the character level
 - Calculates F-score (harmonic mean of precision and recall)

4. Sarah, Fae, Brandon

H A T E

SPEECH

Fae, Sarah, Brandon



CONTENTS

What are we doing

Why we chose this project

Our datasets

Our approach / methods

How we will evaluate our approach



WHAT IS

our project?

- Analyze and compare hate speech trends across different countries.
- Use the model to predict where a particular piece of hate speech originated from.

Motivation

- I chose this project because I thought it seemed super interesting. Hate speech has so many layers to it, it can often be so covert adding to its complexity. I'm excited to see how it differs globally -Fae
- Online hate speech is an ongoing problem with many layers and I feel that it would be very insightful to see what trends and patterns exist and differ between countries. -Brandon
- I have a tendency to view hate speech from a very Western perspective, so I am interested to see how it differs across languages - Sarah

Approach/Methods

- We will be using datasets from <https://hatespeechdata.com/>. We wanted to look at hate speech from a global context so we will look at data from a variety of languages and geographic areas. We are differentiating areas mostly based on language rather than country. Here are two examples of our datasets:

- Toxic Language Dataset for Brazilian Portuguese (ToLD-Br)
- Brazilian portuguese
- 21,000 human annotated samples
- Generated from twitter and annotated by 'demographically diverse' volunteers.

- Multi-Label Hate Speech and Abusive Language Detection in Indonesian Twitter
- It was published in 2019
- 13,169 human annotated samples.
- Indonesian
- Built with a twitter crawl between March 20th, 2018 and September 10th, 2018, and annotated by paid crowdsourced labor.

Evaluation

- As it is difficult to auto-generate performance metrics for this type of problem, we will take a random sample of ~ 100 rows of our dataset and manually check if we predicted the target identities correctly. For non-English text we will translate the text with Google Translate.

5. Jonathan, Abe, Stephen, Jeremy, Brayden

By JJABS

SIMPLYFY

Project Motivation

Adults illiterate: **21%**

< 6th grade level: **54%**

Losses: **\$2.2 trillion**

Overall, by simplifying text:

- Less text = easier to read
- Caters to lower attention spans
- More reading = more knowledge

Data Used

- **ASSET**
 - **For simplifying sentences**
 - **Human generated**
- **Example (from Asset)**
 - **Input: In architectural decoration Small pieces of colored and iridescent shells have been used to create mosaics and inlays, which have been used to decorate walls, furniture, and boxes.**
 - **Output: Small pieces of colored shells make mosaics that decorate walls, furniture, and boxes**

Approach/Methods

- **Moses to create statistical machine translation model**
 - **Try and feed multiple sentence data into Moses**
 - **Moses is built to do one sentence translated to one different sentence**
 - **ASSET's data is one complex sentence to multiple simple sentences**
 - **If things break, email Michael**
- **Moses uses Perl scripts, so we'll call those in Python**

Evaluating the approach

Character f-score

Calculates the similarity between machine translation output and a reference translation using character n-grams

BLEU score

Bilingual Evaluation Understanding, gives a score 0 to 1 on translation of text

SARI

System Output Against References and Input is lexical simplicity metric that measures "how good" are the words added, deleted, and kept

Concerns

- Bias in the data because the data is based off of historical information
- Moses not taking one complex sentence to many simple sentence examples.



6. Bridget, Krishna, Ashu



Motivation for Project



- ❑ **Challenge:** U.S. legal decisions rely heavily on precedent, making efficient retrieval crucial for judicial accuracy.
- ❑ **Problem:** The vastness and complexity of legal precedents complicate accurate and timely retrieval.
- ❑ **Objective:** Develop an NLP-based information retrieval system to efficiently identify and rank relevant legal precedents for new legal cases.





Data

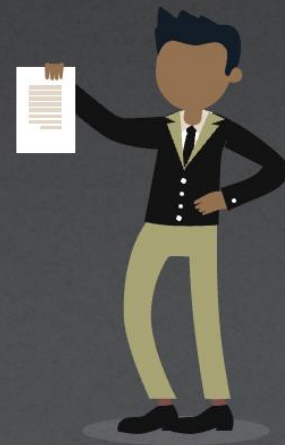
- ❑ Using the Supreme Court Database (SCDB) from Penn State and WashU, containing structured case data from 1946-2023.
- ❑ Data includes case names, dates, docket numbers, legal issues, decisions, rulings, and cited precedents.
- ❑ Data preprocessing steps: tokenization, stopword removal, lemmatization/stemming, case normalization, and n-gram extraction.





Approach and Methods

- ❑ Traditional Retrieval System (BM25): Implemented using Pyserini, focusing on probabilistic retrieval based on TF-IDF.
- ❑ Retrieval-Augmented Generation (RAG): Utilizing a GPT-based LLM for retrieving and summarizing relevant legal precedents.
- ❑ Software used: Pyserini, NLTK, Pandas, NumPy, Scikit-learn, and GPT-based LLM.





Evaluating Our Approach

- **Precision@k**: Measures how many of the top-k retrieved documents are relevant, highlighting immediate retrieval accuracy.
- **Recall@k**: Evaluates the proportion of relevant documents retrieved out of all relevant documents available, assessing completeness.
- **Mean Average Precision (MAP)**: Averages precision scores across multiple queries, reflecting overall retrieval performance.
- **Mean Reciprocal Rank (MRR)**: Highlights the retrieval speed by measuring how quickly the first relevant document is retrieved.



7. Ben Adams, Ezra, César, Nhu



The Good and the Dad

Ben Adams, Ezra Cheifetz, César Guerra-Solano, Nhu
Nguyen

Humor is Underexplored in NLP

- Great significance in society and pop culture
- Current work focuses on humor as a whole
 - Difficulty formalizing humor
 - Effect of noise of different humor types
 - LM abilities fall a bit short :(
 - They've done great at “un-funning” humor though!

Humor Has Structure

- Prevalence of highly-structured humor
 - Question-answer format jokes
 - Common within “dad jokes”
- Structured can be used to aid:
 - Understanding tasks
 - **Generation tasks!**
- Can we use this structure to generate dad jokes?

The Data

- *Dad Jokes* dataset from Kaggle
- 91728 Dad Jokes
- Somewhat unclean set, will require pruning

A steak pun is a rare medium well done.

They say that breakfast is the most important meal of the day. Well, not if it's poisoned. Then the antidote becomes the most important.

What do you get if you cross an angry sheep with a moody cow? An animal that's in a baaaaad mooood.

An apple a day keeps the doctor away. At least it does if you throw it hard enough.

What sounds like a sneeze and is made of leather? A shoe.

Our Focus

- Q&A style dad jokes
- About 7000 in this format

What do you call a kangaroo's lazy joey? A pouch potato.

Why do dragons sleep during the day? Because they like to fight knights.

What are the strongest days of the week? Saturday and Sunday. All the others are weekdays.

What do you call a cow with no legs? Ground beef.

Why did police arrest the turkey? They suspected fowl play.

Methods/Approaches

1. Clean dataset to remove jokes that are misshapen or found offensive by authors
2. Utilize RNNs to achieve baseline question to answer generation
3. Combine our dad-joke dataset with non-joke Q&A text to create a Dad Joke binary classifier
4. Compare RNN generations to Dad Joke binary classifier to benchmark quality of generation
5. As needed, refine RNN to improve generation
6. Integrate the RNN model into a front-end, such as a web app (stretch goal)

Evaluation

Binary Classifier

- F1-Score
- Accuracy

Evaluates

RNN

- Perplexity