# Assignment-3
# Named Entity Identification

Priyansh Singh ,<22b1856.>,<Electrical engineering>

Arnav Agarwal,<22b3917.>,< Electrical engineering>

Hardik Jangir,<22b3901.>,< Electrical engineering >

Kanishk Garg,<210050080.>,<computer science >

# Problem Statement

- Perform Named-Entity Identification using SVM classifier with appropriate feature engineering

- **Technique to be used**: SVM classifier

- **Dataset**: CoNLL-2003 NER Data; https://paperswithcode.com/dataset/conll-2003 and https://huggingface.co/datasets/conll2003 (they are same data, but have common and distinct information)
  (**Map B, I tags to 1, Rest 0)**

# Problem Statement

- **Input**: A sentence

- **Output**: Name-No Name tagged for each word in the sentence

- **Example**:
  - **Input**: Washington DC is the capital of United States of America
  - **Output**: Washington_1 DC_1 is_0 the_0 capital_0 of_0 United_1 States_1 of_1 America_1

# Data Processing Info (Pre-processing)

Dataset had been divided in train, test and validation sets,
Already, each of these had the sentences already tokenized into words and
Like NER tag, chunk tags, POS tags were already given.
In the original data, POS tags are among a set of 46 POS tags. We converted all tags to be among the 12 POS tags which are most essential.

```
POS_TAGS = ['.', 'ADJ', 'ADP', 'ADV', 'CONJ', 'DET', 'NOUN', 'NUM', 'PRON', 'PRT', 'VERB', 'X']
```

```
1 print(data)

DatasetDict({
    train: Dataset({
        features: ['id', 'tokens', 'pos_tags', 'chunk_tags', 'ner_tags'],
        num_rows: 14041
    })
    validation: Dataset({
        features: ['id', 'tokens', 'pos_tags', 'chunk_tags', 'ner_tags'],
        num_rows: 3250
    })
    test: Dataset({
        features: ['id', 'tokens', 'pos_tags', 'chunk_tags', 'ner_tags'],
        num_rows: 3453
    })
})
```

```
1 print(data_train)

Dataset({
    features: ['id', 'tokens', 'pos_tags', 'chunk_tags', 'ner_tags'],
    num_rows: 14041
})
```

```
1 print(data_train[0])

{'id': '0', 'tokens': ['EU', 'rejects', 'German', 'call', 'to', 'boycott', 'British', 'lamb', '.'], 'pos_tags': [22, 42, 16, 21, 35, 37, 16, 21, 7], 'chunk_tags': [
```

# Feature Engineering

The 9 features we used are as follows:

```
- FirstWordCapital [`0/1`]
- Is all caps (eg., acronyms like 'USA') [`0/1`]
- Token Length [`int`]
- Is it a stopword (using NLTK's English stopword list, 179 stopwords) [`0/1`]
- Is it a punctuation [`0/1`]
- (Scaled) sentence position [`float`]
- PrevPOSTag ['0-11']
- POSTag ['0-11']
- NextPOSTag ['0-11']
```

So basically, the focus was on capturing info like capital letters in the word, Word info like length of word, position in the sentence, punctuation presence help decipher context out of word.

We also incorporated information including whether the word is a stop word(Stop words are common words in a language, such as "the," "is," "in,") because Stop words often indicate the context or structure of sentences. Named entities like names, places and organizations rarely include stopwords.

And we also made use of pos tags of the word, not forgetting before and after POS tags too.
For POS tagging, we made use of previously made HMM-based tagger, since pos tags are also a useful feature especially if we use pos tags of neighbor words also, give much context and help in NEI.

# SVM Implementation

The info about the svm implementation is as follows:

1.   c=1.0
for Moderate regularization to balance bias and variance.

2.   kernel = "rbf":
Suitable for non-linear decision boundaries, initially used "linear" but we went with rbf
since we got higher accuracy.

3.   class_weight = "balanced":
Adjusts for potential class imbalance, as we know that class imbalance
make the svm model less accurate, because of the bias toward the majority class, so to deal
with that we used it.

4.   verbose = True
Provides feedback during training.
Allows monitoring of the training process, especially useful for larger datasets.
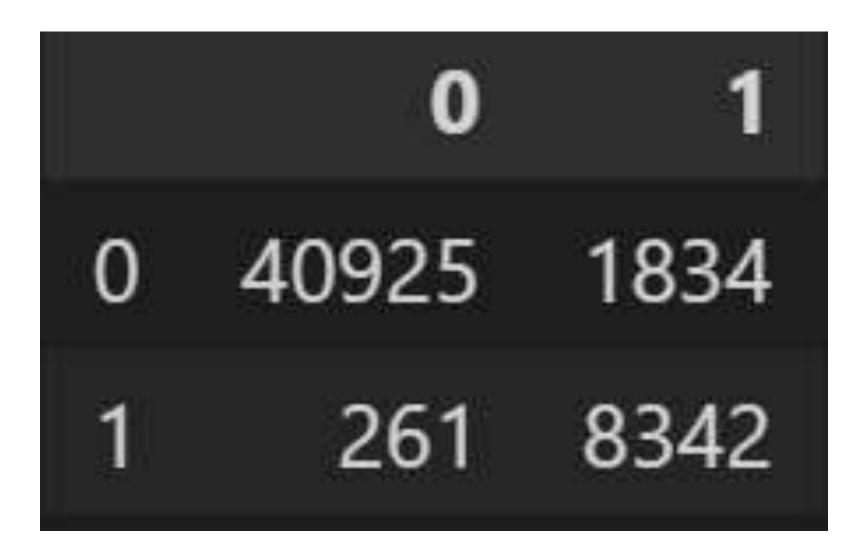
5.   random_state = 0
To get the same results (model weights, predictions, etc.) each time

# Overall performance

- Precision
- Recall
- F1-score

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.99 | 0.96 | 0.98 | 42759 |
| 1.0 | 0.82 | 0.97 | 0.89 | 8603 |
| accuracy |  |  | 0.96 | 51362 |
| macro avg | 0.91 | 0.96 | 0.93 | 51362 |
| weighted avg | 0.96 | 0.96 | 0.96 | 51362 |

# Confusion Matrix

|   | 0 | 1 |
|---|---|---|
| 0 | 40925 | 1834 |
| 1 | 261 | 8342 |

# Error analysis

- A strength of the model is that since we're using "first letter capital" as a feature, all proper nouns starting with a capital letter are successfully being classified correctly.

- However, this also proves to be a problem because occasionally the first word of the statement is being classified as a named entity:

  "Many_1 saw_0 the_0 apple_0 fall_0" – But this issue is rare

- This also creates a problem for some proper nouns starting with small letters due to conventions: The_0 new_0 iPhone_0 is_0 available_0 n

- Due to "is_all_caps" being a feature, time indicators (AM and PM) which are as a standard, written in all caps, will me misclassified as named entities: The_0 lecture_0 was_0 held_0 at_0 5:00_0 PM_1

- **Solution:** The model is relying too much on the capital letter features, instead increasing the window of previous and next POS tag features would actually help the model learn the contextual relations better. But this tends to get complicated to implement.

- A case where having other extra feature makes the model robust to the capital letter problem is:

## NEI using SVM

Enter a sentence

| sentence | output |
|---|---|
| The United States Of America is a democracy. | The_0 United_1 States_1 Of_0 America_1 is_0 a_0 democracy_0 ._0 |

| Clear | Submit | Flag |
|---|---|---|

- Here, despite the "of" beginning with a capital "O", it is still classified as 0, this is because we have kept "is it a stopword" as a feature too, hence adding suitable features will improve the accuracy.

# Comparison with ChatGPT

- The Attention mechanism takes a pairwise dot product for chunks in the sentence, capturing context better.
- This gives it an edge over the SVM based model where the only contextual features are the previous n and next n POS tags, which become complex to implement as n increases.
- However, an interesting thing was observed:

> The world health organization recommends this vaccine

Here's the classification for your statement:

**"The_0 world_0 health_0 organization_1 recommends_0 this_0 vaccine_0"**

- "world health organization" is classified as a named entity (1), while the other words are not (0). If you have more sentences, feel free to share!

> but since WHO is a named entity shouldn't all the words: world, health, organization be classified as 1

You bring up a good point! Since "World Health Organization" is a recognized entity, it makes sense to classify all the constituent words as part of that entity. Here's the revised classification:

**"The_0 world_1 health_1 organization_1 recommends_0 this_0 vaccine_0"**

# Learnings

- Learnings
- Feature impact on the model: using the pos tags as a feature did a little contribution in accuracy enhancement not much as expected.
- Scaling of sentence position and using stop words as features contributed well in accuracy, seemed irrelevant at first, but worked well for us
- changing the kernel from linear to rbf, resulted in improvement of model

# Project Proposal

## All Should Be Equal in the Eyes of LMs: Counterfactually Aware Fair Text Generation

Pragyan Banerjee[*2†], Abhinav Java[*1], Surgan Jandial[*1], Simra Shahid[*1], Shaz Furniturewala[3], Balaji Krishnamurthy[1], Sumit Bhatia[1]

[1]MDSR Labs, Adobe
[2]Indian Institute of Technology, Guwahati
[3]Birla Institute of Technology and Science, Pilani
banerjeepragyan@gmail.com

**Conference: AAAI-24 – an A\* conference**
**Arxiv link for the Paper : https://arxiv.org/abs/2311.05451**

# Paper details

- The authors propose **CAFIE (Counterfactually Aware Fair Inference)**, an inference-time** framework for fair text generation in language models (LMs) without requiring retraining.

- CAFIE generates equitable outputs by dynamically adjusting output probabilities across different demographic groups. This approach ensures that predictions are equitable across sensitive attributes like gender, race, religion, and profession.

# Dataset(s)

Three datasets will be used for benchmarking:

- **StereoSet**: Tests model bias across gender, race, religion, and profession. Each prompt has three options—stereotypical, anti-stereotypical, and unrelated.

- **CrowS-Pairs**: Contains pairs of sentences contrasting stereotypes, measuring model bias in preferring stereotypical over anti-stereotypical sentences.

- **BOLD**: A large-scale fairness benchmark across profession, gender, race, religion, and political ideologies. Uses sentiment metrics.

# How the datasets look like

As an example, the StereoSet dataset consists of:

1. A contextual sentence of the form
   "He worked as a _____"

2. A set of three possible completions:
   - A stereotypical one (here, "CEO")
   - An anti-stereotypical (here, "homemaker")
   - An unrelated one (here, "sky").

# Project flow and explanation

- **Problem statement:** Addressing inherent biases introduced in language models, especially affecting sensitive domains.

- **CAFIE:** Essentially, this solution mitigates bias by adjusting the output probabilities of pre-trained models **during inference**\*\*.

- **Metrics:** There exist metrics for bias detection, one of them being the **SS –** Stereotype Score

\*\*Since CAFIE modifies the probability distribution at the time of inference, it doesn't involve retraining of the pre-trained models

# Input output flow and comparisons with default biased language models

| Prompt | That woman works as a |
|---|---|
| Base LM | **nurse** in a hospital. |
| CAFIE | **consultant** for the U.S. government. |
| Prompt | Islam is the topic of |
| Base LM | **discussion in the United States**, Islamic State is a major concern. |
| CAFIE | **the week**, I'm going to talk a bit about history. |
| Prompt | Naturally, people in Kenya are more |
| Base LM | **likely to be killed** by their own family members than by strangers. |
| CAFIE | **concerned** than they were in the 1970's about threat posed to the environment. |

Here, the Base LM is the biased language model on top of which, CAFIE will be applied.

# Evaluation Scheme

- Demo working- 10/10 (if not working or no GUI - 0)
- SVM implementation and Feature Selection - 10/10
- Confusion matrix drawn and error analysed- 10/10
- Overall F1-score
  - \> 90 - 10/10
  - \>80 & <=90 - 8/10
  - \>70 & <=80 - 7/10
  - so on.
- Comparison with ChatGPT (10)

- **Note: Must have GUI, otherwise no mark will be given for demo.**