

# Assignment 2

## Methodology and Assumptions

- Used Add-k Smoothing
- Added Start and End dummy words in each sentence
- Randomised the dataset before splitting
- For any tag, its first two characters are considered and rest ignored. For example - Merger\_NN-HL, then word **Merger** will be considered to be tagged as **NN** and not **NN-HL**.
- **Case folding** is used for word normalisation.
- Words without tag in dataset are tagged as OOV(Out of Vocabulary).
- **SpeedUp - 9x via Multiprocessing + Using the pypy3 compiler instead of python3**
- Time for Bigram - **1.5 min** and for trigram - **30 min**
- Used **Hold-out** method for **trigram**.
- Command to run code - **python3 <code.py> <dataset\_file\_name>**
- Bigram model has **want\_parallelism** variable to disable multiprocessing and it takes 3.5 min around for execution.
- Empty lines are removed from dataset

## Q1 Bigram Model

Following results are based on **weighted-macro**

### Fold 1

**Accuracy** = 0.8820700565019499    **Precision** = 0.9358948348936105  
**Recall** = 0.9364771856118227    **F1- score** = 0.933143581595565

### Fold 2

**Accuracy** = 0.8833315157862018    **Precision** = 0.9365496910984699  
**Recall** = 0.9369760927122053    **F1- score** = 0.9339699478245068

### Fold 3

**Accuracy** = 0.8822880991486078    **Precision** = 0.935943581336039  
**Recall** = 0.9364604342287322    **F1- score** = 0.9333650508698922

### Average of 3-folds

**Accuracy** = 0.8825632238122533    **Precision** = 0.9361293691093732  
**Recall** = 0.9366379041842533    **F1-score** = 0.9334928600966547

**Accuracy** = **0.936** when I used formula

Accuracy = correctly\_tagged / (correct\_tagged + wrong\_tagged)

## Q1 Trigram Model

Following results are based on **weighted-macro**

**Accuracy** = 0.8680384420167648    **Precision** = 0.9294264089360309  
**Recall** = 0.9284344250811633    **F1- score** = 0.9244502903635906

Accuracy = **0.926135833131034** when I used formula

Accuracy = correctly\_tagged / (correct\_tagged + wrong\_tagged)

Q-2) For markov assumption length = 2

$$P(x^n) = \prod_{i=1}^n P(x_i | x_{i-1} x_{i-2})$$

We need to find such a tag sequence  $x^n$   
for a given word sequence  $w^n =$   
 $\langle w_1, w_2, \dots, w_n \rangle$  that  
maximises  $P(x^n | w^n)$

$$\text{Now } P(x^n | w^n) = \frac{P(w^n | x^n) P(x^n)}{P(w^n)}$$

since  $w^n$  is given and it will be same  
in all possible  $x^n$  so we ignore it.

$$P(x^n | w^n) = P(w^n | x^n) P(x^n)$$

$$= \prod_{i=1}^n (w_i | x_i) \cdot P(x_i | x_{i-1} x_{i-2})$$

↓  
likelihood  
assumption

↓  
for trigram  
assumption

$$= \prod_{i=1}^n \frac{\text{count}(w_i, x_i)}{\text{count}(x_i)} \times \frac{\text{count}(x_i x_{i-1} x_{i-2})}{\text{count}(x_i x_{i-1})}$$

For  $P(x_1 | x_0 x_{-1})$  and  $P(x_2 | x_1 x_0)$  to be  
valid, we insert two ~~dummy~~ dummy  
words  $w_{-1}, w_0$  and  $w_{n+1}, w_{n+2}$  at back  
and front of word sequence  $w^n$ .  
and then calculate the required emission  
and transmission matrix.

### Q3

- Words which are uniformly distributed over the tag-set are tagged incorrectly because they have the probability are close to each other for any tag consider and hence they are the ones which are tagged incorrectly.
- Similarly if  $t_i$  and  $t_{i+1}$  have uniformly distribution then they are the ones which get wrong tags due to same reason as above.
- If the tags are not enough for a particular word in a dataset then also wrong tag is given due to less data available.