# NER USING CRF

## Ayush Kumar Gupta

# Project Overview

As Mentioned in Previous Slide, I had to learn CRF, Feature Functions and had to make a Pipeline from Chunking,POS to NER all through CRF. I also had to do Annotations. I had to test on two different data andcheck their results.

# Contents

# Table Of Contents

# How I Proceeded?

My Goal was to test use CRF Feature Functions for English which includes many Lexical properties for NER, such as Capitalization, Position in a sentence, Articles etc. with a rich premade-gazeteer.

But after Presentation, My Topic changed to Hindi, so I researched about CRF, Feature Functions, How they work, Formulaes, Linear CRF and How to Implement on Hindi.

# Summarise Learnings

- CRF is Conditional Random Field model, which works on its Template and Its Feature Functions, whereas Template provide information about which parts are dependent on which. Function provide probability the of current node in given nodes on which it is dependent.
- Functions are of form $f(X, Y_{i-1}, Y_i, i)$, which have arguments current_val, prev_state, current_state, weight(acc. to length).
- $P(Y/X) = 1/Z * e^{\wedge}Sigma(F(X, Y_{i-1}, Y_i, i))$, where $F(X, Y_{i-1}, Y_i, i) = Sigma(w_j * f_j)$
- **w** is calculated while training the model.
- The Above Formulaes are for Linear CRF++, in which Current State is only dependent on Previous adjacent Stage.

# Problems

My Goal was to test use CRF Feature Functions for English which includes many Lexical properties for NER, such as Capitalization, Position in a sentence, Articles etc. with a rich premade-gazeteer,(Yup Ctrl+C,V).Whereas Hindi had no capitalization, It is Morphologically Rich, Having hard to find good gazeteer, No Particular Word Order, or Articles.

SO I DROPEED USING CRF FEATURE FUNCTIONS FOR NER

# About

## NER

Named Entity Recognition (NER) is a subtask of information extraction that seeks to locate and classify named entities mentioned in unstructured text into predefined categories such as the names of persons, organizations, locations.

## CRF

Conditional Random Fields (CRF) is a type of probabilistic graphical model used for structured prediction, especially in sequence labeling tasks such as named entity recognition (NER), part-of-speech (POS) tagging, and chunking.

# Process

**1**

## NER

Got NER Data from Hugging Face and Trained Model

**2**

## POS

Collected the Data and trained POS Model. Also anotated some POS test.

**3**

## PIPELINE

Connected Pipeline so, If Data is provided,it will return POS ,NER Annotated Data

**4**

## TOKENIZER

To make It complete Text based, one System.

# NAME ENTITY RECOGNITION

STEPS:

IMPLEMENTED USING CRF++

COLLECTED TRAINING DATA WITH 150K LINES OF ANNOTATION


TRAINED MODEL ON PROVIDED DATA


TESTED  ON PROVIDED COLLECTED TEST

TESTED ON NEWS DATA  (MANUALLY)

TESTED ON STORY DATA  (MANUALLY)


MADE PIPELINE TO PROVIDE TEXT->NER

# PARTS OF SPEECH

STEPS:

IMPLEMENTED USING CRF++

COLLECTED TRAINING DATA WITH 100K LINES OF ANNOTATION

TRAINED MODEL ON PROVIDED DATA

TESTED  ON ANNOTATED TEST (MANUALLY)

TESTED ON NEWS DATA  (MANUALLY)

TESTED ON STORY DATA  (MANUALLY)

MADE PIPELINE TO PROVIDE TEXT->NER

# PIPELINE

STEPS:

WROTE JUPYTER SCRIPT TO TRAIN MODEL WITH JUST PROVIDING ANNOTATED DATA.

WROTE JUPYTER SCRIPT TO ANNOTATE DATA WITH JUST PROVIDING TEXT

IF DATA IS ALREADY ANNOTATED, IT TESTS THE PRECISION,RECALL,F1 SCORE OF MODEL ON ANNOTATED DATA

(PIPELINE CAN BE EXTRAPOLATED TO CHUNKING TOO,IF ANNOTATED DATA IS PROVIDED)

# TOKENIZER

STEPS:
WROTE CODE TO CONVERT TEXT INTO SENTENCES
WROTE CODE TO CONVERT SENTENCES INTO WORD IN CONLL FORMAT FOR
TESTING/ANNOTATING DATA


(I HAVENT USED REGEX,IT WAS SIMPLE CODE, USING REGEX I CAN REMOVE
PUNCTUATIONS OR ONE CAN DO AFTER TOKENIZATION,WITH JUST CTRL+F
(ALT+ENTER))

# RESULTS

# Given Test Data

```
 1  गुजरात    B-ORG    N_NNP    B-ORG
 2  हाईकोर्ट   I-ORG    N_NN     I-ORG
 3  ने  0      PSP 0
 4  मंगलवार   0    N_NNP    0
 5  को  0      PSP 0
 6  कांग्रेस   B-ORG    V_VM     B-ORG
 7  नेता 0     N_NN     0
 8  अहमद      B-PER    CC_CCD   B-PER
 9  पटेल I-PER    N_NN     I-PER
10  पर  0      PSP 0
11  5   0      QT_QTC   0
12  हजार 0     QT_QTC   0
13  रुपये 0    N_NN     0
14  का  0      PSP 0
15  जुर्माना   0    N_NN     0
16  लगा 0      V_VM     0
17  दिया 0     V_VAUX   0
18  जब  0      PR_PRL   0
19  उनके 0     PR_PRP   0
20  वकील      0    N_NN     0
21  ने  0      PSP 0
22  2017       0    N_NNP    0
23  में  0      PSP 0
24  राज्यसभा   B-ORG    JJ   B-ORG
25  चुनाव 0    N_NN     0
26  से  0      PSP 0
27  पहले 0     N_NST    0
28  दिये 0     V_VM     0
29  गए  0      V_VAUX   0
30  विधायकों   0    N_NN     0
31  के  0      PSP 0
32  बयानों 0   N_NN     0
33  की  0      PSP 0
34  सीडी 0     N_NN     0
35  पेश 0      JJ   0
36  करने 0     V_VM     0
```

## NER:
## F1 Score:0.97
## Precision:0.95

```
Matching count: 18980
Total entities: 19893
Precision: 0.9541044588548736
Recall: 1
F1 Score: 0.9765132611324054
```

I wonder O are so many so should we consider them in precision or not

## POS:
## F1 Score:0.93
## Precision:0.88

# News Test Data

**Tabs:** TrainingCode.ipynb | TestingCode.ipynb | output_NER.txt

```
1   मुंबई   N_NNP     B-LOC
2   दंगा   QT_QTC    O
3   मामले  N_NN      O
4   में     PSP O
5   4      QT_QTC    O
6   नवंबर  N_NNP     O
7   2022         N_NN    O
8   को     PSP O
9   राज्य   N_NN      O
10  सरकार        N_NN      O
11  को     PSP O
12  कुछ    QT_QTF    O
13  निर्देश  N_NN      O
14  दिए    V_VM      O
15  थे      V_VAUX    O
16  लेकिन        CC_CCD    O
17  इनका   PR_PRP    O
18  पालन   N_NN      O
19  नहीं    RP_NEG    O
20  किया   V_VM      O
21  गया    V_VAUX    O
22
23  इस     PR_PRP    O
24  पर     PSP O
25  सुप्रीम       N_NNP     B-ORG
26  कोर्ट   N_NN      I-ORG
27  ने     PSP O
28  नाराजगी      N_NNP     O
29  जताई   V_VM      O
30  है      V_VM      O
31
32  अदालत       N_NNP     O
33  ने     PSP O
34  19    N_NNP     O
35  जुलाई  N_NNP     O
36  से     PSP O
```

**NER:**
F1 Score:0.97
Precision:0.95

This One Feels True having true Precision

**POS:**
F1 Score:0.85
Precision:0.91

Tabs: TrainingCode.ipynb | TestingCode.ipynb | output_NER.txt

```
1   भुवाली        NNP  N_NN      0
2   की     PSP  PSP  0
3   इस     DMR  DM_DMD  0
4   छोटी-सी        QTF  JJ  0
5   कॉटेज          NN   N_NN  0
6   में     PSP  PSP  0
7   लेटा-लेटा  JJ   ECH  JJ  0
8   मैं     PRP  RP_RPD  0
9   सामने  NST  N_NST      0
10  के     PSP  PSP  0
11  पहाड़   NN   N_NN      0
12  देखता VM   JJ   0
13  हूँ      VAUX       V_VM      0
14  ।      PUNC         RD_PUNC  0
15
16  पानी-भरे    JJ   N_NNP      0
17  सूखे-सूखे    JJ   ECH  V_VM      B-LOC
18  बादलों NN   N_NN      0
19  के     PSP  PSP  0
20  घेरे    NN   PSP  0
21  देखता VM   PSP  0
22  हूँ      VAUX       V_VM      0
23  ।      PUNC         RD_PUNC  0
24
25  बिना   NEG  N_NNP      0
26  आँखों  NN   N_NN      0
27  के     PSP  PSP  0
28  झटक-झटक       RB   ECHO       V_VM      0
29  जाती   VB   V_VAUX  0
30  धुंध    NN   N_NN      0
31  के     PSP  PSP  0
32  निष्फल        RB   JJ  0
33  प्रयास  JJ   N_NN      0
34  देखता VB   N_NNP      0
35  हूँ      AUX  N_NN      0
36  और     CCD  CC_CCD  0
37  फिर    RB   0
```

## NER:
**F1 Score:0.98**
**Precision:0.97**

But the truth is there are only o tags, and most NER tags are wrong or missed. so Precision is around 0.1

## POS:
**F1 Score:0.80**
**Precision:0.67**

| | |
|---|---|
| Number(Correct_Tags) | 352 |
| Number(Wrong_Tags) | 157 |
| Precision | 0.6692307692 |
| Recall | 1 |
| F1 Score | 0.801843318 |

# What else?

With Chunking Training Data, One can easily ass Chunk Annotator in Pipeline.
With Regex One make Tokenizer better

BIO Tagging is with NER Tagging so no need to do another time

**With Some more Time and Effort,**
**One can make complete Pipeline which takes text and annotate it**

The Problem is Precision, It depends on whether Training Data include such type or not
can be fixed by increasing amount of Training Data

# SOURCES

Github Repo: https://github.com/AKGIIITH/NER_CRF_Hindi
[Read README.md]

# THANK YOU