

# Computational Linguistics - 1, Assignment 1

Zubair Abid, 20171076

The assignment works on two language data/sets, English and Bengali. The task is to assign (correctly) appropriate POS Tags for different words in either language based on rules.

## 1 English POS Tagging

Tagset used: the Brown Corpus tagset (As given in assignment resources)

### 1.1 Method used for tagging

For the given ruleset, each word was initially tagged with the first given tag, and then corrected manually if needed. Along the way rules were constructed to automatically tag further data. The general procedure followed is Brill's method.

### 1.2 Development data with tags

The correct tag for each word is highlighted in bold text

1. *She has been absent since last Wednesday.*

She/**PPS** has/**HVC** been/**BEN** absent/**JJ/VB** since/**IN/CS/RB** last/**AP/VB/NN/RB** Wednesday/**NR** ././-HL

2. *It doesn't matter what excuse he gives me, I can't forgive him.*

It/**PPS/PPO** doesn't/**DOZ\*** matter/**NN/VB** what/**WDT** excuse/**NN/VB** he/**PPS** gives/**VBZ** me/**PPO** ,/-HL I/**PPSS** can't/**MD\*** forgive/**VB** him/**PPO** ././-HL

3. *I canceled my appointment because of urgent business.*

I/**PPSS** canceled/**VBD/VBN** my/**PP\$** appointment/**NN** because/**CS/RB** of/**IN** urgent/**JJ** business/**NN** ././-HL

4. *What do you do in Japan?*

What/**WDT** do/**DO** you/**PPSS/PPO** do/**DO** in/**IN/RP** Japan/**NP** ?./.-HL./-TL./-NC

5. *The Handmaid's Tale is an awesome piece of dystopian fiction.*

The/**AT** Handmaid's/**NP\$** Tale/**NN** is/**BEZ** an/**AT/CC** awesome/**JJ** piece/**NN/VB** of/**IN** dystopian/**JJ** fiction/**NN** ././-HL

6. *OK. Now what?*

OK/**JJ/RB** ././-HL Now/**RB** what/**WDT** ?./.-HL./-TL./-NC

7. *I was laughed at by everyone.*

I/**PPSS** was/**BEDZ** laughed/**VBD/VBN** at/**IN** by/**IN/RB** everyone/**PN** ./.

8. *There were people everywhere, covering the roads along the route from the BJP headquarters to the Smriti Sthal from side to side, with security personnel maintaining strict vigil to ensure that nothing goes wrong.*

There/**EX/RB** were/**BED** people/**NNS/VB** everywhere/**RB/NN** ,/-HL covering/**VBG/NN** the/**AT** roads/**NNS** along/**IN/RB** the/**AT** route/**NN** from/**IN** the/**AT** BJP/**NP** headquarters/**NN/NNS** to/**TO/IN/NPS/QL** the/**AT** Smriti/**NP** Sthal/**NP** from/**IN** side/**NN/JJ** to/**TO/IN/NPS/QL** side/**NN/JJ** ,/-HL with/**IN/RB** security/**NN** personnel/**NNS** maintaining/**VBG** strict/**JJ** vigil/**NN** to/**TO** ensure/**VB** that/**IN** nothing/**PN** goes/**VBZ** wrong/**JJ/RB/NN** ././-HL

### 1.3 Disambiguation Rules

Some disambiguation rules (based entirely on the given development dataset and tags) are:

1. If there is a JJ after HV\*/BE\*, change to VB
2. NN after P\*/N\* and VB\* → VB
3. Single word sentence JJ → RB
4. If a word ends with -ed and has VBN in options and follows BE\*/HV\*, tag as VBN
5. VB\* after J\* is NN\*
6. Untagged words default to NP, NP\$ with apostrophe

## 1.4 Tagged Test Data

During/IN a/AT visit/VB to/TO the/AT Cleve-  
land/NP Indians/NPS ,/-HL Beane/NP meets/VBZ  
Peter/NP Brand/NP ,/-HL a/AT young/JJ  
Yale/NP economics/NN graduate/NN with/IN rad-  
ical/JJ ideas/NNS about/IN how/WRB to/TO  
assess/VB player/NN value/NN ./ Beane/NP  
tests/VBZ Brand's/NP\$ theory/NN by/RB ask-  
ing/VBG whether/CS he/PPS would/MD have/HV  
drafted/VBN Beane/NP out/IN of/IN high/JJ  
school/NN ./ Though/CS scouts/NNS consid-  
ered/VBD Beane/NP hugely/JJ promising/VBG  
 ,/-HL his/PP\$ career/NN in/IN the/AT Ma-  
jor/NP Leagues/NP was/BEDZ disappointing/JJ ./  
Brand/NP admits/VBZ that/CS ,/-HL based/VBN  
on/IN his/PP\$ method/NN of/IN assessing/VBG  
player/NN value/NN ,/-HL he/PPS would/MD  
not/\* have/HV drafted/VBN him/PPO until/IN  
the/AT ninth/OD round/NN ./ Impressed/NP ,/-  
HL Beane/NP hires/VBZ Brand/NP as/CS his/PP\$  
assistant/NN manager/NN ./

## 1.5 Additions to disambiguation rules

Some of the tags assigned by the previous ruleset  
turned out to be faulty. Based on the observations,  
some new rules were added to tag the data better:

7. QL followed by a TO tag is WRB
8. RP followed by a IN\* is IN
9. CS at a non initial position is IN

## 2 Bengali POS Tagging

Tagset used: BIS Tagset

### 2.1 Method used for tagging

A method similar to the one used for English tagging.  
Since the tags were unstructured, initially all tags are  
applied to any word and then disambiguation is done.

### 2.2 Development data with tags

1. আমার খুব খিদে পেয়েছে  
আমার/PR\_PRP খুব/RP\_INTF খিদে/N\_NN  
পেয়েছে/V\_VM
2. ভেজা মেঝেতে দৌড়োতে হয় না বাবা  
ভেজা/V\_VM\_VNG মেঝেতে/N\_NN দৌড়োতে/V\_VM  
হয়/V\_VM\_VF/V\_AUX না/RP\_NEG/RP\_RPD  
বাবা/N\_NN
3. পরীক্ষায় ভালো না করলে ঘর থেকে বার করে দেব

পরীক্ষায়/N\_NN ভালো/RB, JJ না/RP\_NEG/RP\_RPD  
করলে/V\_VM\_VNF ঘর/N\_NN থেকে/PSP বার/N\_  
NN, RB করে/V\_AUX/PSP/V\_VM/N\_NN/V\_  
VM\_VNF দেব/V\_VM\_VNF/N\_NNP

4. খাবারে বিরিয়ানি না দিলে বিয়েতে কেউ আসবে না  
খাবারে/ বিরিয়ানি/N\_NNP না/RP\_NEG/RP\_RPD  
দিলে/V\_VM\_VNF/V\_AUX বিয়েতে/N\_NN  
কেউ/PR\_PRI আসবে/V\_VM\_VNF/V\_AUX  
না/RP\_NEG/RP\_RPD

### 2.3 Disambiguation Rules

Some disambiguation rules (based entirely on the given  
development dataset and tags) are:

1. V\_\_ \* without immediately preceding V\_\_ \* is V\_\_  
VM\_\_ \*
2. V\_\_ \* after a V\_\_ VM is V\_\_ AUX IF option avail-  
able
3. In preceding, if no option available flip tags
4. না immediately before or after a V\_\_ is RP\_\_ NEG
5. RB, JJ/N\_\_ \* immediately before RP\_\_ NEG/V\_\_ \*  
is RB
6. PSP, \* immediately after N\_\_ \* is PSP

### 2.4 Tagged Test Data

1. আমার আজ দুপুর থেকেই বেশ ঘুম পেয়েছে  
আমার/PR\_PRP আজ/N\_NN দুপুর/N\_NN  
থেকেই/PSP/V\_VM\_VINF বেশ/RP\_INTF/QT\_  
QTF ঘুম/N\_NN পেয়েছে/V\_VM
2. এই কাজ জিনিষটা আমার একদম ভাল লাগে না  
এই/DM\_DMD কাজ/N\_NN জিনিষটা/N\_NN  
আমার/PR\_PRP একদম/RB ভাল/JJ লাগে/V\_VM\_VF  
না/RP\_NEG
3. পরীক্ষাটা একদম জঘন্ন হতে চলছে  
পরীক্ষাটা/N\_NN একদম/RB জঘন্ন/JJ হতে/V\_VM\_  
VNF/V\_VM/V\_AUX/PSP চলছে/V\_VM\_VF
4. আমার বাড়িতে একদিন আয়, বেশ মজা হবে  
আমার/PR\_PRP বাড়িতে/N\_NN একদিন/N\_NN  
আয়/N\_NN ,/, বেশ/RP\_INTF/QT\_QTF মজা/N\_  
NN হবে/V\_VM\_VF/V\_AUX

### 2.5 Additions to disambiguation rules

Not many discrepancies were found with existing rules,  
but a new rule was added.

1. RP\_\_ INTF, QT\_\_ QTF before a N\_\_ \* is RP\_\_  
INTF