

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

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.
10. CS at a non initial position is IN

2 Bengali POS Tagging

Tagset used: BIS Tagset

2.1 Method used for tagging

2.2 Development data with tags

1. আমার খুব খিদে পেয়েছে
আমার/PR_PRP খুব/RP_INTF খিদে/N_NN পেয়েছে/
2. ভেজা মেঝেতে দৌড়োতে হয় না বাবা
ভেজা/V_VM_VNG মেঝেতে/N_NN দৌড়োতে/ 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

2.4 Tagged Test Data

1. আমার আজ দুপুর থেকেই বেশ ঘুম পেয়েছে
2. এই কাজ জিনিষটা আমার একদম ভালো লাগে না
3. পরীক্ষাটা একদম জঘন্ন হতে চলছে
4. আমার বাড়িতে একদিন আয়, বেশ মজা হবে