

CS5012 Language and Computation

Practical 1 – Report POS tagging with HMM

Introduction

The aim of practical is to POS tag words of the sentences using HMM algorithms, Viterbi, Beam Search, Forward and backward probability.

To use these algorithms, it was required to train the model by getting the smoothed emission and transition probabilities using tagged Brown corpus and NLTKs Witten-Bell smoothing function.

For this practical OOP approach is used using Python 3.7 (Worked in lab machines). There are 3 class python files and 1 main file which runs them.

Training

This part is done in the Probabilities.py file, which contains class Probabilities. Constructor saves sentences emitted from the corpus. It also contains uniqueTags list which contains all tags which may appear in corpus. Emission and transition probabilities are created and stored in dictionary variable, so they can be accessed via getEmissionProbabilities and getTransitionProbabilities. Bigrams for the transition probabilities are created manually, without nltk's ngrams method.

Algorithm 1

Viterbi algorithm is implemented using 2d list which stored all probabilities. Another list backpointer is used to store the tags where these probabilities are coming from.

```

self: <ViterbiTagging.ViterbiTagger object at 0x7f847d222048>
backpointer: [['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', ...], ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', ...], ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', ...]]
00: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
01: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
02: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'NOUN', ...]
03: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
04: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
05: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
06: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'NOUN', ...]
07: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
08: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
09: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
10: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
11: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
12: ['<s>', 'DET', 'NOUN', 'ADP', 'DET', 'NUM', 'VERB', 'ADV', 'VERB', 'VERB', 'ADP', 'DET', 'ADJ', 'ADJ', ...]
__len__: 13
probability: <Probabilities.Probabilities object at 0x7f84b061cc50>
resultingTag: []
sentence: ['The', 'importance', 'of', 'this', '5', 'can', 'largely', 'be', 'explained', 'by', 'the', 'natural', 'mathematical', 'properties', ...]
tagsPossible: ['DET', 'ADV', 'PRON', 'ADP', '.', 'NOUN', 'VERB', 'ADJ', 'NUM', 'PRT', '<s>', 'CONJ', 'X']
viterbi: [[0.015397347444342062, 2.2712956750385473e-12, 2.4715689768241237e-16, 5.558889726993048e-10, 8.20003720808096e-20, 1.2028555195808311e-23, 7.290317331653548e-26, 5.165017682937777e-30, 9.411224231
00: [0.015397347444342062, 2.2712956750385473e-12, 2.4715689768241237e-16, 5.558889726993048e-10, 8.20003720808096e-20, 1.2028555195808311e-23, 7.290317331653548e-26, 5.165017682937777e-30, 9.411224231
01: [7.498034477072187e-08, 2.8116697345579937e-10, 1.687018467055395e-14, 5.707453693880601e-16, 1.0150944543941589e-17, 1.5651759214782174e-21, 2.7973121751777108e-21, 2.6036359163932684e-28, 2.14997
02: [8.287088815804835e-09, 9.43410298166718e-12, 9.906851590574286e-16, 1.5717947379336276e-16, 3.405985240926656e-19, 1.917149111914125e-23, 5.124535011243418e-26, 7.399215649528267e-30, 6.61537020605
03: [6.983887604152491e-09, 8.02555665817095e-12, 4.327087795474674e-08, 4.741852408562528e-17, 2.8974599994824754e-19, 4.1153259819511336e-22, 1.7223356404210816e-25, 2.2354817269779916e-29, 2.2233993
04: [2.9456679176147674e-10, 9.006823156895383e-13, 9.167526672066917e-16, 2.2326664663344344e-18, 3.25172481173742e-20, 5.838713184397746e-23, 5.8344864120869075e-27, 1.770985475474889e-30, 5.53186142
05: [4.394223812645027e-07, 7.678013396499295e-07, 3.6275287493127084e-13, 2.7817654353434697e-14, 7.657413622791397e-16, 2.6153111536591403e-19, 4.4708175640085494e-24, 2.5091664306102502e-28, 5.77147
06: [5.673008805332937e-08, 1.3101906872212654e-09, 1.5233626995162174e-13, 2.3075771374489537e-15, 4.730168996915474e-17, 1.8007212673137156e-17, 4.817549241317629e-24, 2.324589020458354e-23, 1.682588
07: [7.583622712631774e-08, 6.931751958582617e-09, 2.1521594216882683e-14, 6.5814299657015146e-15, 2.5025638274339295e-16, 6.23355324600146e-21, 1.8862636537944924e-24, 7.549484594137776e-28, 2.4350174
08: [2.69810706237989e-08, 3.2120360390989437e-10, 1.1563336856782448e-14, 2.6359086742705336e-15, 5.1561098901832446e-14, 1.8411642911431322e-21, 3.5270665625346366e-25, 6.741190056270255e-29, 4.55316
09: [6.285007946614414e-09, 4.91700512417035e-12, 2.732326815672331e-15, 1.080414468495031e-16, 1.7751816909455156e-19, 6.862610663856121e-23, 2.335406990992056e-25, 1.6788730776784052e-29, 3.01482608
10: [1.1985976900931345e-17, 7.796740915310235e-20, 1.6685304919245348e-24, 2.070307383751048e-25, 2.814849969114043e-27, 1.8924932441589383e-30, 7.215161556097402e-35, 3.8442483170926234e-38, 9.314204
11: [1.4876330269580675e-09, 2.1151444547512292e-13, 1.087749145312989e-15, 1.5121296811113998e-18, 7.636285940240624e-21, 4.6206574914640466e-23, 5.841997453766276e-27, 1.0348652818578934e-30, 7.54159
12: [1.864778991331042e-09, 8.733637948450753e-11, 1.3418675587389211e-15, 5.946375157915282e-17, 3.153097815340753e-18, 5.435654097508849e-23, 1.2434112279439794e-26, 2.2082996275846488e-30, 1.6951457

```

Using recursion, the tags are taken from the backpointer list. By training first 10000 words, and next 500 words for testing viterbi gave 0.9500389644125032 accuracy.

Algorithm 2

Used the same approach. However, Beam search has additional functions called keepMaximum. It takes K maximum from each column and makes other values – 1. While calculating next stage probabilities – 1 probabilities are ignored. The reason why 0 is not used, is underflow. It would not be wise to make artificial 0s, which could lead to the fake underflows.

The accuracy of beam search is:

beam K=(1) Accuracy: 0.9332409732444368

beam K=(2) Accuracy: 0.9489133258290761

For $K = 13$ it gives the same probability as Viterbi. As it keeps all values from the probability table, which makes it operate like viterbi algorithm

Algorithm 3

Algorithm 3, forward and backward is implemented. However, its accuracy is 0.9279591306606633. Forward and backward probability tables are calculated and in the third stage these values were used and multiplied by emission and transition probabilities of each position in table. The difference of forward probability table here, was summation instead of the maximisation.

Run

To run the application, the main.py file should be run. Using
`python3 Main.py`