

Assignment 2: Advanced Text Processing

Oreen Yousuf

December 27, 2020

1 POS-Tagging

The best version of my tagger was "t3-tagger", which had an Accuracy rate of 93.09 percent. The following are errors made by the tagger:

Sentence 1: "From the AP comes this story: [...]" - "AP" was mistagged as NOUN instead of PROPN. This can be chalked up to anything having the possibility of being a proper noun. This case can be considered as genuinely ambiguous. This instance would not be considered as mistagged by the gold standard ewt-dev-wt.txt file.

Sentence 2: "Bush nominated Jennifer M. Anderson for a 15-year term as associate judge of the [...]" - " - " was 'mistagged' as SYM instead of PUNCT. I do think that the PUNCT tagging is incorrect in this case. The hyphen (-) used here is just a connector for the full phrase "15-year". This mistagging in the gold standard of the ewt-dev-wt.txt file could've been avoided if the phrase "15-year" was tokenized as a single token rather than broken up into three individual tokens.

Sentence 3: "Bush nominated Jennifer M. Anderson for a 15-year term as associate judge of the [...]" - "associate" was mistagged as NOUN instead of ADJ. This instance is another case of ambiguity, as "associate" has three parts-of-speech; verb, noun, and adjective. This instance was utilizing the adjective form of associate, to mean "connected with" and is often seen as a way to rank something.

Sentence 4: "Bush also nominated A. Noel Anketel Kramer for a 15-year term as associate judge of the District of Columbia Court of Appeals, replacing John [...]". - "Appeals" was mistagged as NOUN instead of PROPN. The word "appeal" can be used as a noun or verb. Here the tagger believes it's being used as a noun meaning "a sincere or earnest request or plea", however it's used as an official proper noun for a specific judiciary body.

Sentence 5: "The sheikh in wheel-chair has been attacked with a F-16-launched bomb." - "F", in "F-16", was 'mistagged' as a PROPN instead of NOUN. However, just as in Sentence 2, I believe the tag of PROPN from the t3-tagger is in fact the correct tag, as the token in question is that of a specific item/noun. There isn't much ambiguity to discuss for this example.

Sentence 6: "He could be killed years ago and the israelians have all the reasons [...]". - "israelians" was mistagged as NOUN instead of PROPN. The incorrect tagging of this as a noun could be attributed to the fact that "israelians" is not capitalized in the text, as it should be (I believe). Furthermore, because it's referring to a specific group/nationality it should qualify as a proper noun.

Sentence 7: "...is the spiritual leader of Hamas, but they didn't." - "did" was mistagged as AUX instead of VERB. This one can be ambiguous for it's linguistic nature. For it to be an auxiliary verb, a verb must

follow "did" or "didn't" in its base form. And because a base-form verb (and nothing, really) follows it, it's acting as a verb and not an auxiliary verb in this instance.

My native languages are Harari and English, and there isn't a tagset for Harari as it's a very small language. But the biggest, immediately related language of Amharic has one that we can use to make a relative comparison. Unique properties for the research done on this include: pronouns acting like nouns, prepositions functioning also as conjunctions, adjectives being categorized as a subclass of nouns (this is actually common in most Semitic languages too). Obstacles for conducting POS tagset creations for Amharic are namely frequently mis-tagged words and word classes not necessarily corresponding to prosodic words. The latter of which is the most pressing issue. Certain prosodic words under the microscope contain more than 1 grammatical word that belong to different word classes. Additionally, because the driving force behind Amharic tagset creation is dividing prosodic words into morphological segments (word classes), a handful of special tags had to be created (i.e. PRONPC, PRONC, PRONP, NPC, NC, NP). Below are the 30 tags in the created tagset:

| Basic class | Definition of the tag | Code of the tag |
|--------------|--|-----------------|
| Noun | Verbal/ infinitival Noun, formed from any verb form such as active, passive, and repetitive, by attaching the prefix m(ä)- | VN |
| | Any noun including verbal noun attached with a preposition | NP |
| | Any noun including verbal noun attached with conjunction | NC |
| | Any noun including verbal noun with a proclitic preposition and an enclitic conjunction | NPC |
| | Any other noun; simple or derived | N |
| Pronoun | Pronoun attached with preposition | PRONP |
| | Pronoun attached with conjunction | PRONC |
| | Pronoun with a proclitic preposition and an enclitic conjunction | PRONPC |
| | Any other Pronoun | PRON |
| Verb | Auxiliary verb | AUX |
| | Relative verb | VREL |
| | Any Verb including relative verbs and auxiliaries attached with preposition | VP |
| | Any Verb including relative verbs and auxiliaries attached with conjunction | VC |
| | Any Verb including relative verbs and auxiliaries with a proclitic preposition and an enclitic conjunction | VPC |
| | Verb (all other) | V |
| Adjective | Adjective attached with preposition | ADJP |
| | Adjective attached with conjunctions | ADJC |
| | Adjective with a proclitic preposition and an enclitic conjunction | ADJPC |
| | Any other Adjective | ADJ |
| Preposition | Preposition | PREP |
| Conjunction | Conjunction | CONJ |
| Adverb | Adverb | ADV |
| Numeral | Cardinal | NUMCR |
| | Ordinal | NUMOR |
| | Numeral (cardinal or ordinal) attached with preposition | NUMP |
| | Numeral (cardinal or ordinal) attached with conjunction | NUMC |
| | Numeral (cardinal or ordinal) with a proclitic preposition and an enclitic conjunction | NUMPC |
| Interjection | Interjections | INT |
| Punctuation | Punctuation | PUNC |
| Unclassified | Unclassified | UNC |

Tokenizing highly optimizes NLP models in the end. This is achieved by parts-of-speech being allocated to all elements and tags being utilized for punctuation. This "simple" implementation can lead to punctuation, symbols, special characters, etc. being peeled or not peeled off of words. Denoting end-of-sentence punctuation characters (such as exclamation marks, question marks and periods) from word-internal punctuation (like e.g., i.e., Ph.D, and etc.) is greatly valued in part-of-speech tagging. Next, because tokenization is deployed to partition out a text into words and symbols, it becomes easier later on in the process to fundamentally understand background information in the original text in addition to neatly breaking up somewhat debated on aspects such as contractions, like 'n't', and 's', from their root/stem words. From this somewhat structured attempt at creating concrete tokens, you are more easily able to explicate raw text and data with the compiled list of each and every token that has been tagged.

Shortly, as we intake a word sequence we can then deduce the corresponding hidden tags from said word sequence. This enables us to obtain the key signal sequence and forecast tags, which is a bedrock of utilizing Hidden Markov Models (HMM) for Part of Speech (POS) tagging. HMM models are likelihood sequence models that pair off a list of observed incidents with specific tag sequences. As we learned in the Hidden Markov Models lab, We can use the phrase "emitted signals" to correlate

to the observation events and we can use the phrase "hidden state of the model" for all possible word sequences for said emitted signals. When we use Hidden Markov Models for Part of Speech tagging, these previously described emitted signals will be the sequence of words and Part of Speech tags will be the hidden state. Once the probabilistic calculations of all possible word tagged sequences are done, we can then choose the likely one and start to find the transition probability (the probability of state i transferring to state j) and emission probability (the probabilities of observations o being generated from state q).

2 Lemmatization

The following are error analyses from my lemmatizer python file that reached 96 percent:

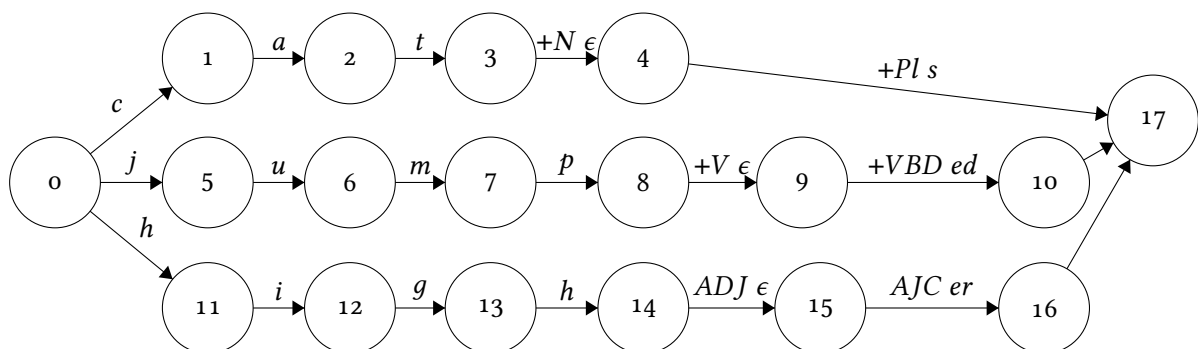
The ADJ lemma of "other" should correctly be "other" but it was seen as "oth". We can rectify this with an if-statement subcondition as ' if word=="other": return "other" '.

The VERB lemma of "went" should correctly be "go", but it was seen as "went". We can rectify this with an if-statement subcondition for this special case as ' if word=="went": return "go" '.

The VERB lemma of "made" should correctly be "make", but it was seen as "made". We can rectify this with an if-statement subcondition, also for this special case as ' if word=="made": return "make" '.

The VERB lemma of "r", broken off from the word "your", should be correctly seen as "be", but was seen as "r". When seeing this I thought again of a subcondition to fix this with an if-statement returning "be", but in my code I already accounted for "yours" or "your" to return "you".

The VERB lemma "submitted" should correctly be "submit". This was a common occurrence with verb inflected in the past form. We can rectify this by seeing if it ends with 2 repeated consonants after removing the "ed", and if so we can remove an additional consonant.



The inputs of the finite state transmitter (FST) are cats/NOUN, jumped/VERB, higher/ADJ and the outputs are cat +Pl -s, jump +VBD -ed, and high AJC -er.

What's happening in the finite state transducer is that it's traversing through the previously mentioned input and "taping" one state onto the next, linking them. It's starting with an input such as "cat" (or rather just the 'c' originated from the whole word input) and jumps to different states depending on the input, while producing output in accordance to its transition sequence. "C" to "a" to "t", etc. and then onto the plural indicator at state 4. This completes the journey of the first word "cats", and the same traversal is repeated for "jumped" and "higher", each transitioning from their initial first character/state inputs and linking to any modifying state such as "VBD" or "AJC" for jumped and higher, respectively. Termination/completion is achieved at the final state (which is represented here by state 17).

In Semitic languages like Harari, Amharic or Arabic, every verb can have 2-3 dozen distinct inflectional forms that also have many special cases. Additionally, because of the numerous plural forms, to be able to parse them all we'd have to save all the morphological variants for each verb. There can still be overlooked data within the corpus even if we create an incredibly large training corpus for morphologically rich languages such as these because it is not feasible to account for all forms of words.

3 VG: Hidden Markov Models

VG: Hidden Markov Models

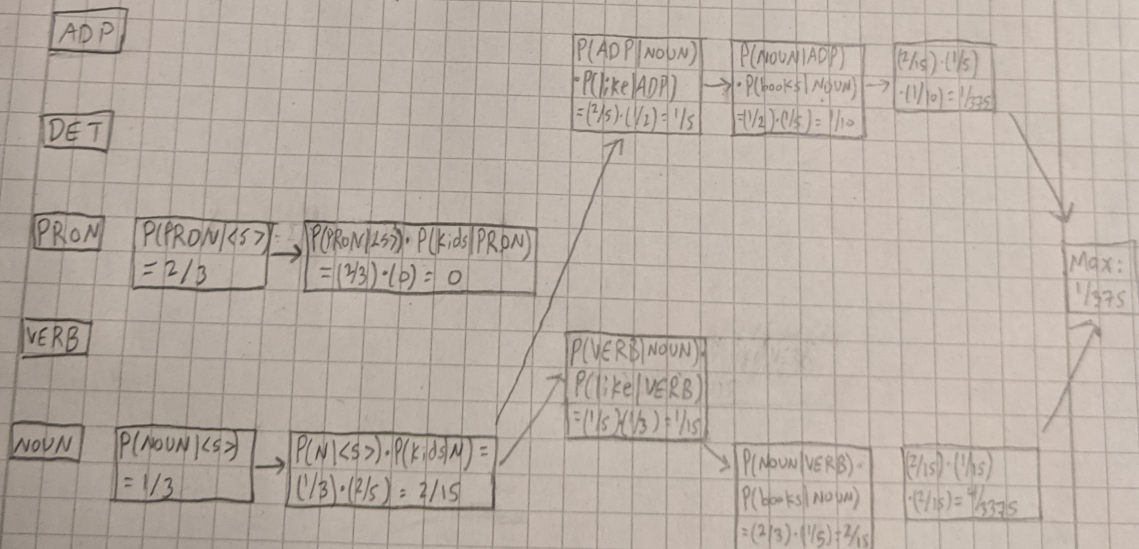
Transition Probability:

| Transition P | Pronoun | Verb | Noun | ADP | DET | Initial | <S> |
|--------------|---------|------|------|-----|-----|---------|-----|
| PRON | 0 | 1/2 | 0 | 0 | 1/4 | PRON | 2/3 |
| VERB | 1/3 | 0 | 2/3 | 0 | 0 | VERB | 0 |
| NOUN | 0 | 1/5 | 0 | 2/5 | 0 | NOUN | 1/3 |
| ADP | 1/2 | 0 | 1/2 | 0 | 0 | ADP | 0 |
| DET | 0 | 0 | 1/2 | 0 | 1/2 | DET | 0 |

Emission Probability:

| Emission P | she | books | trips | for | kids | you | all | the | time | like | these |
|------------|-----|-------|-------|-----|------|-----|-----|-----|------|------|-------|
| PRON | 1/2 | 0 | 0 | 0 | 0 | 1/4 | 0 | 0 | 0 | 0 | 1/4 |
| VERB | 0 | 1/3 | 0 | 0 | 1/3 | 0 | 0 | 0 | 0 | 1/3 | 0 |
| NOUN | 0 | 1/5 | 1/5 | 0 | 2/5 | 0 | 0 | 0 | 1/5 | 0 | 0 |
| ADP | 0 | 0 | 0 | 1/2 | 0 | 0 | 0 | 0 | 0 | 1/2 | 0 |
| DET | 0 | 0 | 0 | 0 | 0 | 0 | 1/2 | 1/2 | 0 | 0 | 0 |

Viterbi trellis of Test:



In the end, the most probabilistic tag-sequence is "kids/noun like/adp books/noun".
 Probability = $(2/15) \cdot (1/5) \cdot (1/10) = 1/375 = 0.0026$