

Pos tagging with HMM & SVM



1

Dataset



Brown Corpus

- 500 samples of English-language text from works published in the United States in 1961
 - 1'100'000 million words
 - 56'000 different words
- Universal Tagset
 - Twelve universal part-of speech categories

Tag	Meaning	English Examples
ADJ	adjective	<i>new, good, high, special, big, local</i>
ADP	adposition	<i>on, of, at, with, by, into, under</i>
ADV	adverb	<i>really, already, still, early, now</i>
CONJ	conjunction	<i>and, or, but, if, while, although</i>
DET	determiner, article	<i>the, a, some, most, every, no, which</i>
NOUN	noun	<i>year, home, costs, time, Africa</i>
NUM	numeral	<i>twenty-four, fourth, 1991, 14:24</i>
PRT	particle	<i>at, on, out, over per, that, up, with</i>
PRON	pronoun	<i>he, their, her, its, my, I, us</i>
VERB	verb	<i>is, say, told, given, playing, would</i>
.	punctuation marks	<i>. , ; !</i>
X	other	<i>ersatz, esprit, dunno, gr8, univeristy</i>

2

HMM Pos Tagger



HMM

Considering an Hidden Markov Model,

- Tags are the **hidden** events, since they are unobservable.
- Words are the **observable** events.

Markov Assumption: $P(q_i | q_1 \dots q_{i-1}) = P(q_i | q_{i-1})$

Decoding: Viterbi Algorithm

- Determining which sequence of hidden state(tags) is the underlying source of some observations(words)



How does it work?

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} \frac{P(w_1^n | t_1^n) P(t_1^n)}{P(w_1^n)}$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(w_1^n | t_1^n) P(t_1^n)$$



$$P(w_1^n | t_1^n) \approx \prod_{i=1}^n P(w_i | t_i)$$

$$P(t_1^n) \approx \prod_{i=1}^n P(t_i | t_{i-1})$$

$$\hat{t}_1^n = \operatorname{argmax}_{t_1^n} P(t_1^n | w_1^n) \approx \operatorname{argmax}_{t_1^n} \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$$

$P(w_i | t_i)$ = *Emission probability*

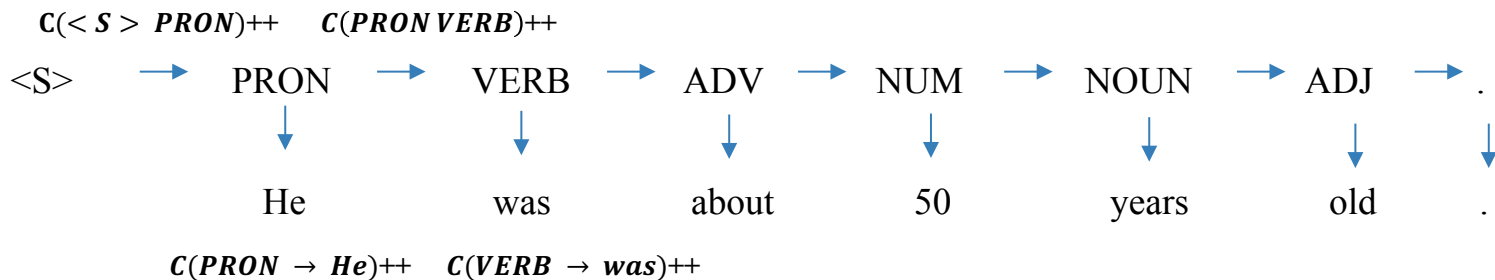
$P(t_i | t_{i-1})$ = *Transition probability*



Learning Markov Models

Count the number of occurrences in the corpus

Sentence: "He was about 50 years old."



Divide by the context to get probability

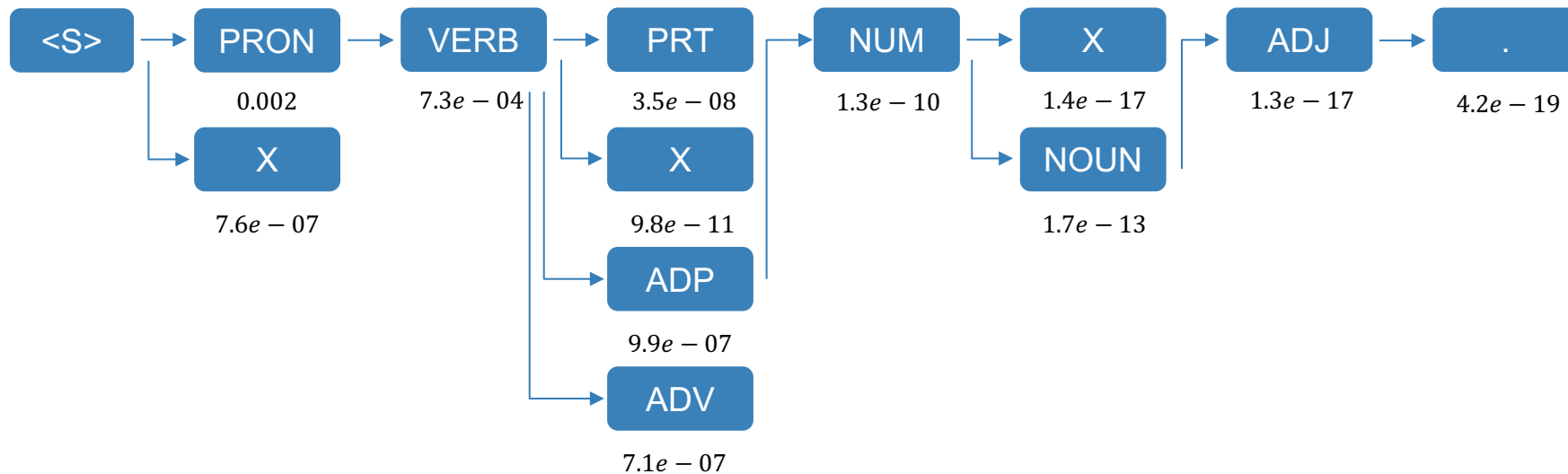
$$P_T(VERB | PRON) = \frac{C(PRON VERB)}{C(PRON)}$$

$$P_E(was | VERB) = \frac{C(VERB \rightarrow was)}{C(VERB)}$$



Viterbi decoding

He was about 50 years old .





Unknown words

Words are considered as unknown if they do not appear in the train set.
Indeed, the emission probability could not be considered.

So, the tag associated to this words only depends on the transition probability.

In the Viterbi step associated to the unknown word, will be only considered $P(t_i|t_{i-1})$ instead of $P(w_i|t_i)P(t_i|t_{i-1})$.



93.2%

Accuracy

Train set length: 51606 sentences

Test set length: 5734

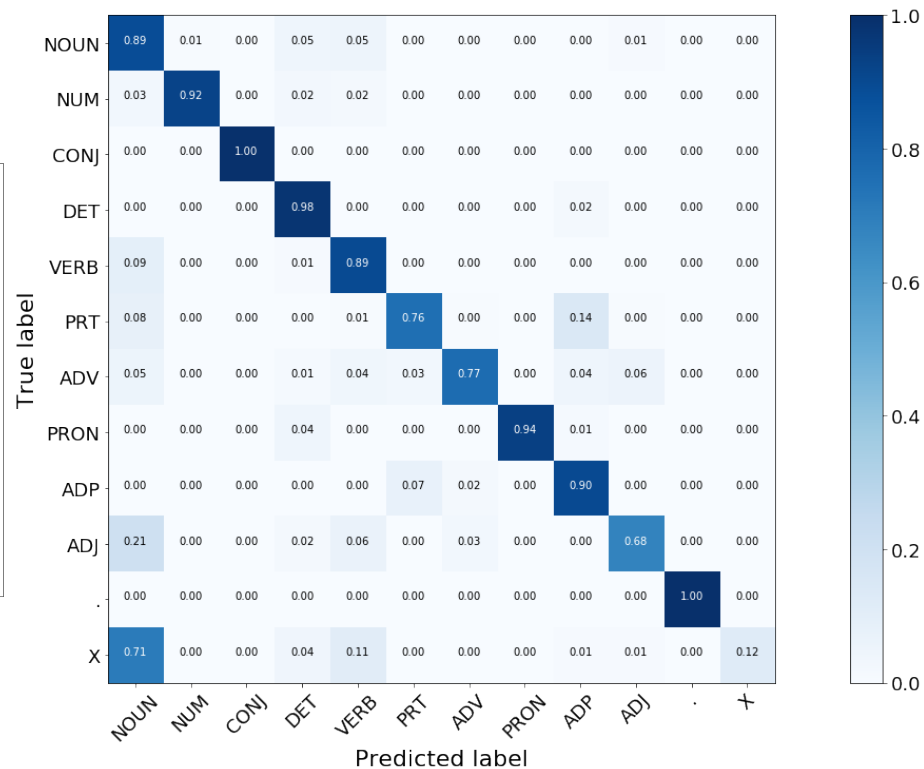
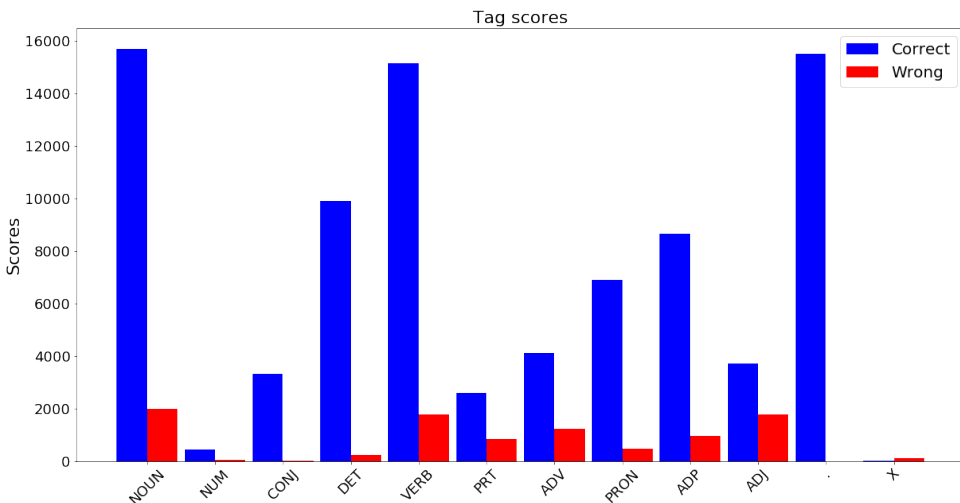
Precision: 93.3%

Recall: 93.2%

F-Measure: 93.2%



Results





90.1%

Accuracy

Train set length: 5000 sentences

Test set length: 5734

Precision: 90.3%

F-Measure: 90%

3

SVM Pos Tagger



How does it work?

SVM training:

- Each training sample is the so-called “word context”.
- The “word context” is a set of features defined on/around the interested word (previous/successive words, tags, suffixes)
- Kernel: RBF
 - Also different kernels has been tried: Linear, sigmoidal and polynomials kernels provided really bad result in term of Accuracy



Problem SVM features

The greatest problem has been the mapping between words and vector (also known as Word Embedding) for which different solution has been tried:

- One hot
 - Each word is encoded as a 1-hot vector (where all the elements are 0 except one, which is 1).
 - Too high-dimensional space, for which the training time was not feasible
- Incremental
 - An incremental number assigned to each word.
 - Meaningless space representation
- Word2Vec
 - Produce a vector space, with each unique word in the corpus being assigned a corresponding vector in the space (10 dimensions).



Features Vector

$[Word_{-1} \ Word \ Word_{+1} \ Tag_{-2} \ Tag_{-1} \ isCapital \ suffixMatch]$

$Word, Word_{-1}, Word_{+1}$: current word and the two words nearby.

Tag_{-2}, Tag_{-1} : The two previous tags

$isCapital$: Binary feature, 1 if the word starts with a capital letter, 0 otherwise

$suffixMatch$: Binary feature, 1 if the word suffix match a pattern, 0 otherwise



88.1%

Accuracy

Train set length: 5000 sentences

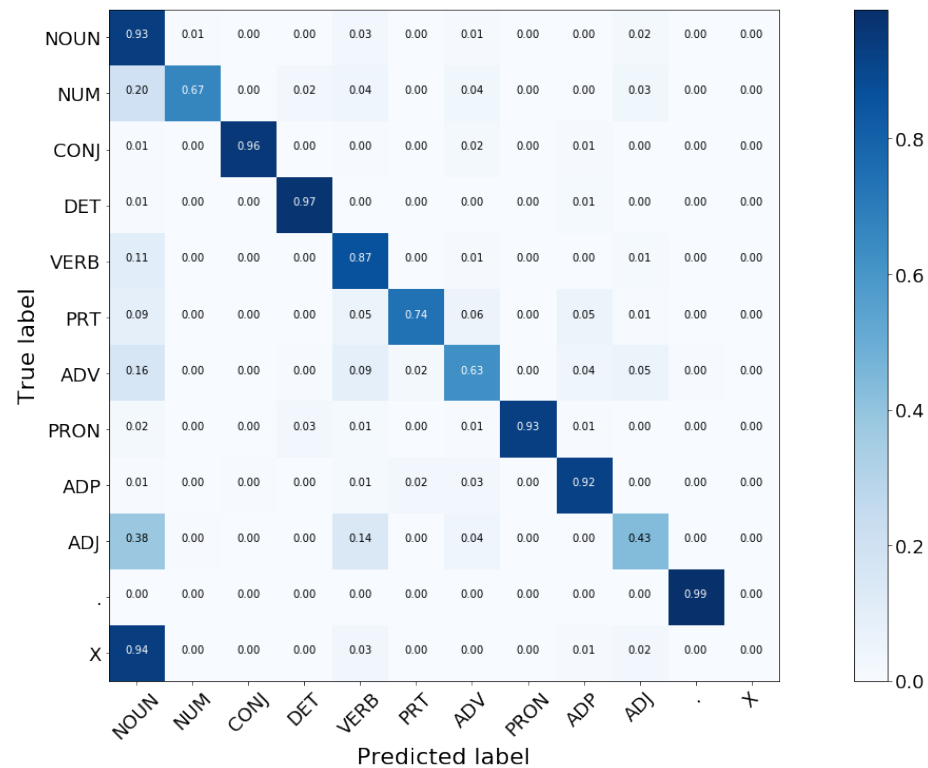
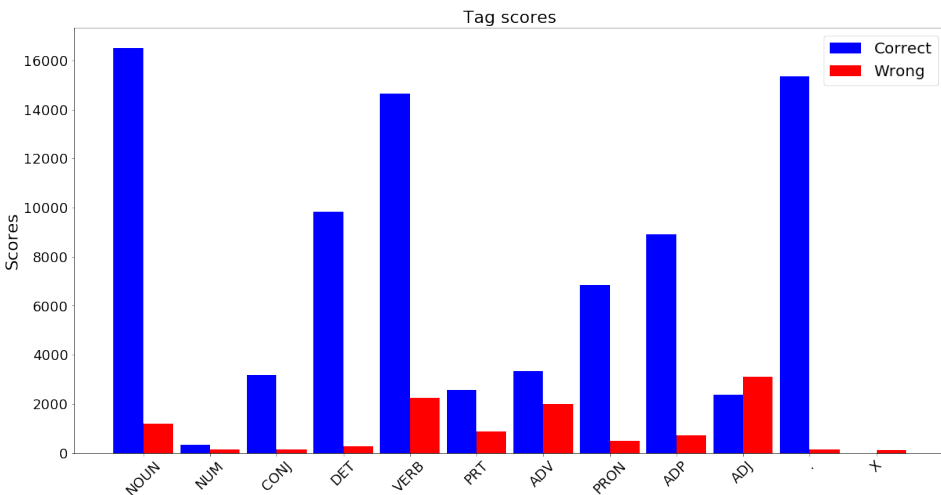
Test set length: 5734

Precision: 88.2%

F-Measure: 87.7%



Results



4

Comparison



Comparison

PROs		CONs	
HMM	SVM	HMM	SVM
<ul style="list-style-type: none">• Interpretability• Computationally low cost	<ul style="list-style-type: none">• Addition of right features lead to better results• Space representation provide a better unknowns handling	<ul style="list-style-type: none">• Limitation due to its structure• Unknown words handling• Out of the ordinary sentences bad predicted	<ul style="list-style-type: none">• Interpretability• Computationally expensive



Thanks!

D'Amicis - Romeo

Cognitive Robotic Project – A.A. 2016-2017