

## NLP - ex2 - Practical Part

Guy Kornblit, 308224948; Mohamad Salama, 318983384

### 1. Data:

After loading the raw data, we considered only the prefix of each tag, before a '+' or '-' sign, resulting:

	# Sentences	# Words	# Wordforms	# Unique Tags
Train	4161	90,538	13,576	98
Test	462	10,016	2,826	79

In addition, we measured 2,008/2,826 wordforms and 78/79 tags from test that appeared in the train set.

### 2. Models Comparison.

In all the following HMM methods, we manually inserted smoothing with  $1e-20$  for all zero probabilities (transitions and emissions).

Model	Error rate for Known Word	Error Rate for Unknown Words	Total Error Rate
Baseline MLE	0.194	0.75	0.258
Base HMM	0.122	0.123	0.122
HMM-Add-1 Smoothing	0.130	0.134	0.130
HMM with Pseudo-words	0.127	0.123	0.126
HMM-Add-1 Smoothing with Pseudo-words	0.135	0.129	0.134

Generally, we can see that HMM variations perform rather similarly, and outperforms the MLE baseline model.

There is a major difference between baseline to the other models with respect to the unknown words error rate, that because we handle unknown words entirely different using the Viterbi algorithm - we just assume that the emission probability for that word is 1 for the default POS tag 'NN' and zero otherwise, same as baseline.

But, considering the `back_pointers_table[k]` where `k` is the location of the unknown word, we updated for every POS tag in the row the previous tag that maximizes the route and the transition. i.e

$$\operatorname{argmax}_t \{ \text{transitions} * \pi[k-1].\max(\text{axis} = 1) \}$$

Thus, we kept a probability model for these unknown words instead of tagging them in a deterministic way.

### 3. Pseudo words creation:

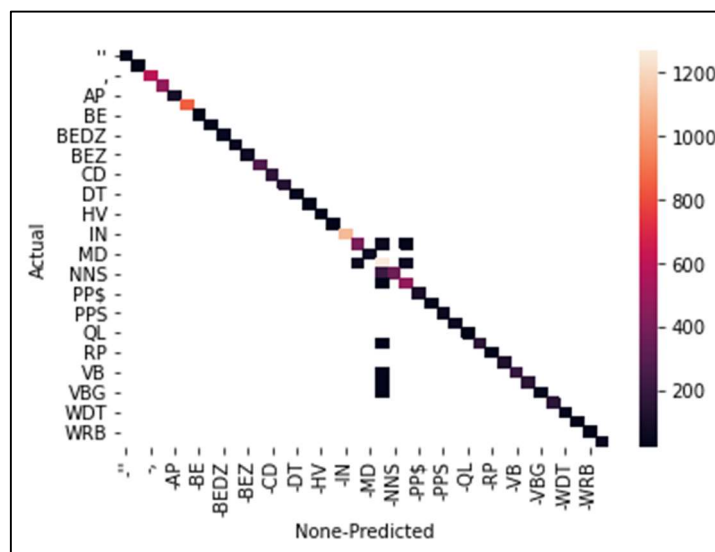
We chose threshold=2 to categorize low frequency training words by the following statistics:

```
for t in [1,2,3,5, 20, 50]:
    print("Threshold ", t, "below: ", (word_freq.freq <= t).sum(), ". above: ", (word_freq.freq > t).sum())

Threshold 1 below: 7398 . above: 6178
Threshold 2 below: 9510 . above: 4066
Threshold 3 below: 10537 . above: 3039
Threshold 5 below: 11610 . above: 1966
Threshold 20 below: 13101 . above: 475
Threshold 50 below: 13432 . above: 144
```

### 4. Error analysis:

Due to the sparsity of the confusion matrix, we decided to plot it only for the top 50 values of joint occurrences of tags. In addition, we show the number of error and correct predictions for every actual label.





מספר ז'נר  
308224948

מספר סלמ  
318983384

\*

\*

\* 2827

\*

\*

1 שלב

מספר הכתובת אפשר להעביר

1) מכתב מס' (1) נקבע שגור הת'ו'ים  $y_1, y_2, \dots, y_n$

$$P(y_1, \dots, y_n) = \prod_{i=1}^n P(y_i | y_{i-1}, y_{i-2}, \dots, y_{i-n})$$

ולו כיוון הנוסחה שמ'צ'ת את כתביו כסדר מרוק מ'ס'  $M$

2) מכתב מס' (2) נקבע שגור ה'ל'ס  $x_1, x_2, \dots, x_n$  ות'ו'ים  $y_1, y_2, \dots, y_n$

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | y_i) = \prod_{i=1}^n e(x_i | y_i)$$

כעת, ע'י  $y$  לא תל'ה כ'א, ו'ם ההסתברות של  $x_i$  ל'היות  $y_i$  לא תל'ה כ'א ו'ם, א'ז ההסתברות המשותפת של מ'ס'  $y_i$  של ע'רת הת'ו'ים תה'ה מ'כ'ת ההסתברות כלומר:

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = \prod_{i=1}^n P(y_i | y_{i-1}, y_{i-2}, \dots, y_{i-n}) \cdot \prod_{i=1}^n e(x_i | y_i) \\ = \prod_{i=1}^n P(y_i | y_{i-1}, y_{i-2}, \dots, y_{i-n}) \cdot e(x_i | y_i)$$

ולו כיוון ה'ל'ס של HMM שגור ע'רת מ'ק'ו'ת מ'ס'  $M$ ,  $M$



emission

	H	L
A	0.2	0.3
C	0.3	0.2
G	0.3	0.2
T	0.2	0.3

transsion		
	H	L
H	0.5	0.5
L	0.4	0.6

2 of the

input  $\Rightarrow S = ACCGTGCA$

קבוצה  $k=1, \dots, 8$  של  $n$  נקודות  $\{H_i, L_i\}$   $\forall i \in \{1, \dots, n\}$

$$\pi(k, v) = \max_{u \in S_{k-1}} \{ \pi(k-1, u) \cdot q(v|u) \cdot c(x_k|v) \}$$

K	$X_k$	BP-H	H	BP L	L
1	A	H	0.1	H	0.12
2	C	L	0.018	L	0.0144
3	C	H	$2.7 \cdot 10^{-3}$	L	$1.728 \cdot 10^{-3}$
4	G	H	$4.05 \cdot 10^{-4}$	H	$2.16 \cdot 10^{-4}$
5	T	H	$4.05 \cdot 10^{-5}$	H	$4.86 \cdot 10^{-5}$
6	G	L	$7.29 \cdot 10^{-6}$	L	$5.832 \cdot 10^{-6}$
7	C	H	$1.0935 \cdot 10^{-6}$	L	$6.9984 \cdot 10^{-7}$
8	A	H	$1.0935 \cdot 10^{-7}$	H	$1.3122 \cdot 10^{-7}$

כל נהפך עת'א'ים הנכונים הם:

5- L H H H L H H L

והסתברות של  $\leq$  פנינת הרצף של  $f$  היא

$$S = 1.2 \cdot e^{-31}$$

אחצה לא את זה ש' דנאסחה

הנחת נוספת היא שהפרמטרים  $\theta$  אינם תלויים בנתונים  $x$  ו- $y$ .  

$$P(x_1, x_2, \dots, x_8, y_1, \dots, y_8) = \prod_{i=1}^8 P(y_i | x_i) \prod_{i=1}^8 P(x_i | y_i)$$

$X = (A, C, C, G, T, G, C, A)$   $Y = \{L, H, H, H, L, H, H, L\}$  70%

$$y_0 = 1$$



