

# COMP550 – Natural Language processing

## Assignment 2

Badreddine SAYAH

### Question 1

1. The Viterbi algorithm returns to the globally optimal state sequence for hidden Markov models.

$$\operatorname{argmax}_{\vec{Q}} P(\vec{Q}, \vec{O} | O)$$

**True.**

**Proof: by induction**

$V(j)$  denotes the highest probability.

$$V_t(j) = \operatorname{argmax}_{\vec{Q}} P(\vec{Q}, \vec{O} | O)$$

By definition:

$$V_t(j) = \max_{1 \leq i \leq N} [V_{t-1}(i) a_{ij}] b_j(O_t)$$

doing term rearrangement and apply induction hypothesis:

$$B_j(O_t) \max_{1 \leq i \leq N} (a_{ij} \times P(q_{t-1} = Q_i, Q_1, O_1, \dots, O_{t-1}))$$

We can replace the transition  $(a_{ij})$  (by its equation), move  $b_j(O_t)$  inside the max, since it's a constant. Then we can use one-to-one Markov assumption and apply chain rule, we obtain:

$$\begin{aligned} & \max_{1 \leq i \leq N} \max_{\vec{Q}} b_j(O_t) P(q_{t-1} = Q_i, q_t = Q_i, O_1, \dots, O_{t-1}) \\ & \max_{\vec{Q}} b_j(O_t) P(q_{t-1} = Q_i, q_t = Q_i, O_{t-1}) \end{aligned}$$

By replacing  $b_j(O_t)$  The emission probability of  $O_t$  from the state  $S_j$  and again using Markov assumption: we got

$$\begin{aligned} & \max_{\vec{Q}} P(O_t | q_{t-1}, q_t = Q_j, O_{1..t}) \times P(q_{t-1}, q_t = Q_j, O_{t-1}) \\ & \max_{\vec{Q}} P(q_{t-1}, q_t = Q_j, O_{1..t}) \end{aligned}$$

Using the arguments of the maxima; **the argmax** is stored at each step giving the best state sequence for HMM.

## Question 2

### Non terminal categories:

|                |                                       |
|----------------|---------------------------------------|
| <b>S</b>       | <b>sentence/clause</b>                |
| <b>NP</b>      | <b>noun phrase</b>                    |
| <b>VP</b>      | <b>verb phrase</b>                    |
| <b>N</b>       | <b>noun</b>                           |
| <b>PN</b>      | <b>pronoun</b>                        |
| <b>V</b>       | <b>verb</b>                           |
| <b>DT</b>      | <b>determiner</b>                     |
| <b>A</b>       | <b>adjective</b>                      |
| <b>N-Pl</b>    | <b>noun plural</b>                    |
| <b>N-M-Sg</b>  | <b>noun male singular</b>             |
| <b>N-F-Sg</b>  | <b>noun female singular</b>           |
| <b>N-M-Pl</b>  | <b>noun male plural</b>               |
| <b>N-F-Pl</b>  | <b>noun female plural</b>             |
| <b>N-Prop</b>  | <b>proper noun</b>                    |
| <b>PR-1Sg</b>  | <b>pronoun first person singular</b>  |
| <b>PR-2Sg</b>  | <b>pronoun second person singular</b> |
| <b>PR-3Sg</b>  | <b>pronoun third person singular</b>  |
| <b>PR -1Pl</b> | <b>pronoun first person plural</b>    |
| <b>PR -2Pl</b> | <b>pronoun second person plural</b>   |
| <b>PR -3Pl</b> | <b>pronoun third person plural</b>    |
| <b>V-1Sg</b>   | <b>verb first person singular</b>     |
| <b>V-2Sg</b>   | <b>verb second person singular</b>    |
| <b>V-3Sg</b>   | <b>verb third person singular</b>     |
| <b>V-1Pl</b>   | <b>verb first person plural</b>       |
| <b>V-2Pl</b>   | <b>verb second person plural</b>      |
| <b>V-3Pl</b>   | <b>verb third person plural</b>       |
| <b>DT-Pl</b>   | <b>determiner plural</b>              |
| <b>DT-M-Sg</b> | <b>determiner male singular</b>       |

|                    |   |
|--------------------|---|
| <b>DT-F-Sg</b>     | <b>determiner female singular</b>                 |
| <b>A-M-Sg-Pre</b>  | <b>adjective male singular preceding noun</b>     |
| <b>A-M-Pl-Pre</b>  | <b>adjective male plural preceding noun</b>       |
| <b>A-F-Sg-Pre</b>  | <b>adjective female singular preceding noun</b>   |
| <b>A-F-Pl-Pre</b>  | <b>adjective female plural preceding noun</b>     |
| <b>A-M-Sg-Post</b> | <b>adjective male singular post noun</b>          |
| <b>A-M-Pl-Post</b> | <b>adjective male plural post noun</b>            |
| <b>A-F-Sg-Post</b> | <b>adjective female singular post noun</b>        |
| <b>A-F-Pl-Post</b> | <b>adjective female plural post noun</b>          |
| <b>DO-1Sg</b>      | <b>direct object first person singular</b>        |
| <b>DO-2Sg</b>      | <b>direct object second person singular</b>       |
| <b>DO-M-3Sg</b>    | <b>direct object third person male singular</b>   |
| <b>DO-F-3Sg</b>    | <b>direct object third person female singular</b> |
| <b>DO-1Pl</b>      | <b>direct object first person plural</b>          |
| <b>DO-2Pl</b>      | <b>direct object second person plural</b>         |
| <b>DO-M-3Pl</b>    | <b>direct object third person male plural</b>     |
| <b>DO-F-3Pl</b>    | <b>direct object third person female plural</b>   |

## **2.1 Some advantages of modelling French grammar with a CFG's**

- CFG's are able of modelling multiple sentences/phrases in a language using same rules with good precision.
- The rule set is quite compact and readable, moreover we can easily add new features (rules)
- Computationally tractable: we can use CFG's in a program that can 'decide' whether a sentence or phrase is grammatical or not.

## **2.2 Some disadvantages of modelling French grammar with a CFG's**

The main disadvantage is overgeneration, we can have many rules to one terminal word.

The syntax analysis in French grammar is not enough, like (25) *"les noirs chats mangent le poisson"*

It is unlikely that the cats are eating the fish.

## **2.3 some aspects of French that your CFG does not handle**

This CFG does not handle negation in all its forms in the French grammar, also indirect objects (like: *il lui a pr ter sa voiture/ he lent him his car*), multiple types of determiner and the verb tenses “imparfait, plus que parfait, future...etc” or the conjugation in different tenses .

### Question 3 : Decipherment with HMMs

The table below, show the results in term of accuracy for training Hidden Markov model using NLTK module.

First, we can see that training HMM using MLE estimator gives poor results. 9.18% on cipher 1, slightly better on cipher2 even if the latter is more complex cipher.

Essentially by “fixing” the issue for unseen events where the MLE fails; using Laplace smoothing method dramatically increases the accuracies to 97.66 for cipher1, and 83.12 for the second cipher.

The plain text was improved by getting additional text from Project Gutenberg (over than 1M characters). Thus, retrieving the transitions probability improves the results on the Caesar cipher (94.02%). On the second cipher we observed an increase of four points when we consider the state sequences of the optimal path through the HMM.

By contrast, accuracy on cipher 3 do not move at all. In front of limited data available we’ve tried to pull more text from multiple sources (more than 6 million characters), finetuning gamma parameters of Lidstone estimation or initializing differently the HMM probabilities. We suppose that using another smoothing method like Kneser-ney smoothing can improve significantly the results. We may also consider random restarts for decipherment with HMM as in [1].

|                 | Accuracy<br>Built-in test function |                      |                  |  | Accuracy<br>Best path simple (Viterbi) |                      |                  |  |
|-----------------|------------------------------------|----------------------|------------------|--|--|----------------------|------------------|--|
|                 | MLE                                | Laplace<br>smoothing | Text<br>improve. | Text<br>improve.<br>+ Laplace<br>smoothing | MLE                                    | Laplace<br>smoothing | Text<br>improve. | Text<br>improve.<br>+ Laplace<br>smoothing |
| <b>Cipher 1</b> | <b>9.18</b>                        | <b>97.66</b>         | <b>94.02</b>     | <b>97.66</b>                               | 94.12                                  | 97.66                | 94.02            | 94.02                                      |
| <b>Cipher 2</b> | <b>14.20</b>                       | <b>83.12</b>         | <b>14.98</b>     | <b>83.12</b>                               | 70.64                                  | 83.11                | 74.80            | 74.80                                      |
| <b>Cipher 3</b> | <b>21.30</b>                       | <b>21.30</b>         | <b>21.30</b>     | <b>21.30</b>                               | 21.29                                  | 21.29                | 21.29            | 21.29                                      |

Reference: [1] ["Cryptanalysis of Classic Ciphers Using Hidden Markov Models" by Rohit Vobbilisetty](#)