

Construisons un HMM depuis un corpus d'apprentissage<sup>3</sup> :

|   |   |     |   |       |   |       |
|---|---|-----|---|-------|---|-------|
| Y | = | DET | → | NC    | → | VERB  |
|   |   | ↓   |   | ↓     |   | ↓     |
| X | = | la  |   | bonne |   | soupe |
| Y | = | DET | → | ADJ   | → | NC    |
|   |   | ↓   |   | ↓     |   | ↓     |
| X | = | la  |   | bonne |   | soupe |
| Y | = | DET | → | NC    | → | VERB  |
|   |   | ↓   |   | ↓     |   | ↓     |
| X | = | la  |   | soupe |   | fume  |

<sup>3</sup>exemple repris de : <https://lattice.cnrs.fr/sites/itellier/cours-HMM-CRF.pdf>

Initiales :

- DET = 1
- Transitions :
- ADJ ⇒ NC = 1
  - DET ⇒
    - NC =  $\frac{2}{3}$
    - ADJ =  $\frac{1}{3}$
  - NC ⇒ VERB = 1

Émissions :

- ADJ ⇒ bonne = 1
- DET ⇒ la = 1
- NC ⇒
  - bonne =  $\frac{1}{3}$
  - soupe =  $\frac{2}{3}$
- VERB ⇒
  - soupe =  $\frac{1}{2}$
  - fume =  $\frac{1}{2}$

L'algorithme de Viterbi utilise la propriété de Markov pour donner la meilleure séquence. À chaque instant, on évalue quel est le meilleur "pas" à faire, la meilleure séquence étant simplement la suite des meilleurs pas (la séquence la plus probable est la suite des transitions les plus probables).

Algorithme de Viterbi

L'algorithme de Viterbi utilise la propriété de Markov pour donner la meilleure séquence. À chaque instant, on évalue quel est le meilleur "pas" à faire, la meilleure séquence étant simplement la suite des meilleurs pas (la séquence la plus probable est la suite des transitions les plus probables).

L'algorithme fonctionne de la façon suivante :

1. pour le 1er mot : on calcule, pour chaque étiquette, sa probabilité initiale étant donné le mot.
2. pour les autres : on calcule, pour chaque étiquette, les meilleures transitions depuis le mot précédent, on garde une trace de là où l'on vient (*back pointer*).
3. récupération la meilleure séquence en partant de la fin et en rebrousant chemin (*backtrack*).

HMM et morphosyntaxe : Algorithme de Viterbi (1.)

On commence avec une matrice "vide" (cellule vide = proba de 0).  
On calcule la probabilité du premier mot en multipliant la probabilité initiale et la probabilité d'émission...

|      |    |       |      |
|------|----|-------|------|
|      | la | bonne | fume |
| ADJ  | 0  |       |      |
| DET  | 1  |       |      |
| NC   | 0  |       |      |
| VERB | 0  |       |      |

- Initiales :
- DET = 1

HMM et morphosyntaxe : Algorithme de Viterbi (2.a)

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".  
Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible : on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      |    |       |      |
|------|----|-------|------|
|      | la | bonne | fume |
| ADJ  | 0  | →     |      |
| DET  | 1  | ↗     |      |
| NC   | 0  |       |      |
| VERB | 0  |       |      |

- Transitions :
- ADJ ⇒ NC = 1
  - DET ⇒
    - NC =  $\frac{2}{3}$
    - ADJ =  $\frac{1}{3}$
  - NC ⇒ VERB = 1

HMM et morphosyntaxe : Algorithme de Viterbi (2.b)

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".  
Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible : on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      |    |       |      |
|------|----|-------|------|
|      | la | bonne | fume |
| ADJ  | 0  | → 0   |      |
| DET  | 1  |       |      |
| NC   | 0  |       |      |
| VERB | 0  |       |      |

- Transitions :
- ADJ ⇒ NC = 1
  - DET ⇒
    - NC =  $\frac{2}{3}$
    - ADJ =  $\frac{1}{3}$
  - NC ⇒ VERB = 1

HMM et morphosyntaxe : Algorithme de Viterbi (2.c)

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".  
Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible : on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      |    |                     |      |
|------|----|---------------------|------|
|      | la | bonne               | fume |
| ADJ  | 0  | → 1 × $\frac{1}{3}$ |      |
| DET  | 1  |                     |      |
| NC   | 0  |                     |      |
| VERB | 0  |                     |      |

- Transitions :
- ADJ ⇒ NC = 1
  - DET ⇒
    - NC =  $\frac{2}{3}$
    - ADJ =  $\frac{1}{3}$
  - NC ⇒ VERB = 1

HMM et morphosyntaxe : Algorithme de Viterbi (2.c)

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".  
Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible : on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      |    |                 |      |
|------|----|-----------------|------|
|      | la | bonne           | fume |
| ADJ  | 0  | → $\frac{1}{3}$ |      |
| DET  | 1  |                 |      |
| NC   | 0  |                 |      |
| VERB | 0  |                 |      |

- Transitions :
- ADJ ⇒ NC = 1
  - DET ⇒
    - NC =  $\frac{2}{3}$
    - ADJ =  $\frac{1}{3}$
  - NC ⇒ VERB = 1

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".

Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible :  
on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      | la | bonne         | fume |
|------|----|---------------|------|
| ADJ  | 0  | $\frac{1}{3}$ |      |
| DET  | 1  |               |      |
| NC   | 0  |               |      |
| VERB | 0  |               |      |

Transitions :

- ADJ  $\Rightarrow$  NC = 1
- DET  $\Rightarrow$ 
  - NC =  $\frac{2}{3}$
  - ADJ =  $\frac{1}{3}$
- NC  $\Rightarrow$  VERB = 1

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".

Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible :  
on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      | la | bonne         | fume |
|------|----|---------------|------|
| ADJ  | 0  | $\frac{1}{3}$ |      |
| DET  | 1  |               |      |
| NC   | 0  |               |      |
| VERB | 0  |               |      |

Transitions :

- ADJ  $\Rightarrow$  NC = 1
- DET  $\Rightarrow$ 
  - NC =  $\frac{2}{3}$
  - ADJ =  $\frac{1}{3}$
- NC  $\Rightarrow$  VERB = 1

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".

Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible :  
on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      | la | bonne         | fume |
|------|----|---------------|------|
| ADJ  | 0  | $\frac{1}{3}$ |      |
| DET  | 1  |               |      |
| NC   | 0  |               |      |
| VERB | 0  |               |      |

Transitions :

- ADJ  $\Rightarrow$  NC = 1
- DET  $\Rightarrow$ 
  - NC =  $\frac{2}{3}$
  - ADJ =  $\frac{1}{3}$
- NC  $\Rightarrow$  VERB = 1

Une fois la meilleure transition trouvée, on la multiplie avec la probabilité d'émission.

|      | la | bonne                  | fume |
|------|----|------------------------|------|
| ADJ  | 0  | $\frac{1}{3} \times 1$ |      |
| DET  | 1  |                        |      |
| NC   | 0  |                        |      |
| VERB | 0  |                        |      |

Émissions :

- ADJ  $\Rightarrow$  **bonne = 1**
- DET  $\Rightarrow$  la = 1
- NC  $\Rightarrow$ 
  - bonne =  $\frac{1}{3}$
  - soupe =  $\frac{2}{3}$
- VERB  $\Rightarrow$ 
  - soupe =  $\frac{1}{2}$
  - fume =  $\frac{1}{2}$

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".

Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible :  
on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      | la | bonne         | fume |
|------|----|---------------|------|
| ADJ  | 0  | $\frac{1}{3}$ |      |
| DET  | 1  | 0             |      |
| NC   | 0  |               |      |
| VERB | 0  |               |      |

Transitions :

- ADJ  $\Rightarrow$  NC = 1
- DET  $\Rightarrow$ 
  - NC =  $\frac{2}{3}$
  - ADJ =  $\frac{1}{3}$
- NC  $\Rightarrow$  VERB = 1

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".

Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible :  
on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      | la       | bonne                  | fume |
|------|----------|------------------------|------|
| ADJ  | 0        | $\frac{1}{3}$          |      |
| DET  | <b>1</b> | 0                      |      |
| NC   | 0        | $1 \times \frac{2}{3}$ |      |
| VERB | 0        |                        |      |

Transitions :

- ADJ  $\Rightarrow$  NC = 1
- DET  $\Rightarrow$ 
  - **NC =  $\frac{2}{3}$**
  - ADJ =  $\frac{1}{3}$
- NC  $\Rightarrow$  VERB = 1

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".

Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible :  
on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      | la | bonne         | fume |
|------|----|---------------|------|
| ADJ  | 0  | $\frac{1}{3}$ |      |
| DET  | 1  | 0             |      |
| NC   | 0  | $\frac{2}{3}$ |      |
| VERB | 0  |               |      |

Transitions :

- ADJ  $\Rightarrow$  NC = 1
- DET  $\Rightarrow$ 
  - NC =  $\frac{2}{3}$
  - ADJ =  $\frac{1}{3}$
- NC  $\Rightarrow$  VERB = 1

Une fois la meilleure transition trouvée, on la multiplie avec la probabilité d'émission.

|      | la | bonne                            | fume |
|------|----|----------------------------------|------|
| ADJ  | 0  | $\frac{1}{3}$                    |      |
| DET  | 1  | 0                                |      |
| NC   | 0  | $\frac{2}{3} \times \frac{1}{3}$ |      |
| VERB | 0  |                                  |      |

Émissions :

- ADJ  $\Rightarrow$  bonne = 1
- DET  $\Rightarrow$  la = 1
- NC  $\Rightarrow$ 
  - **bonne =  $\frac{1}{3}$**
  - soupe =  $\frac{2}{3}$
- VERB  $\Rightarrow$ 
  - soupe =  $\frac{1}{2}$
  - fume =  $\frac{1}{2}$

Une fois la meilleure transition trouvée, on la multiplie avec la probabilité d'émission.

|      | la | bonne         | fume |
|------|----|---------------|------|
| ADJ  | 0  | $\frac{1}{3}$ |      |
| DET  | 1  | 0             |      |
| NC   | 0  | $\frac{2}{9}$ |      |
| VERB | 0  |               |      |

Émissions :

- ADJ  $\Rightarrow$  bonne = 1
- DET  $\Rightarrow$  la = 1
- NC  $\Rightarrow$ 
  - bonne =  $\frac{1}{3}$
  - soupe =  $\frac{2}{3}$
- VERB  $\Rightarrow$ 
  - soupe =  $\frac{1}{2}$
  - fume =  $\frac{1}{2}$

HMM et morphosyntaxe : Algorithme de Viterbi (2.h)

On calcule la meilleure transition pour chaque étiquette du mot suivant : "bonne".

Pour chaque étiquette à un instant  $t$  (étiquette <sub>$t$</sub> ) possible :  
on cherche la meilleure transition depuis l'instant  $t - 1$ .

|      | la | bonne         | fume |
|------|----|---------------|------|
| ADJ  | 0  | $\frac{1}{3}$ |      |
| DET  | 1  | 0             |      |
| NC   | 0  | $\frac{2}{9}$ |      |
| VERB | 0  | 0             |      |

Transitions :

- ADJ  $\Rightarrow$  NC = 1
- DET  $\Rightarrow$ 
  - NC =  $\frac{2}{3}$
  - ADJ =  $\frac{1}{3}$
- NC  $\Rightarrow$  VERB = 1

32

HMM et morphosyntaxe : Algorithme de Viterbi (2.i)

On continue jusqu'à remplir la matrice... On peut (enfin!) passer à la dernière étape.

|      | la | bonne         | fume          |
|------|----|---------------|---------------|
| ADJ  | 0  | $\frac{1}{3}$ | 0             |
| DET  | 1  | 0             | 0             |
| NC   | 0  | $\frac{2}{9}$ | 0             |
| VERB | 0  | 0             | $\frac{1}{9}$ |

33

HMM et morphosyntaxe : Algorithme de Viterbi (3.)

À présent, on parcourt la matrice en sens inverse pour reconstituer la meilleure séquence.

|      | la | bonne         | fume          |
|------|----|---------------|---------------|
| ADJ  | 0  | $\frac{1}{3}$ | 0             |
| DET  | 1  | 0             | 0             |
| NC   | 0  | $\frac{2}{9}$ | 0             |
| VERB | 0  | 0             | $\frac{1}{9}$ |

34

HMM et morphosyntaxe : Algorithme de Viterbi (3.)

À présent, on parcourt la matrice en sens inverse pour reconstituer la meilleure séquence.

|      | la | bonne         | fume          |
|------|----|---------------|---------------|
| ADJ  | 0  | $\frac{1}{3}$ | 0             |
| DET  | 1  | 0             | 0             |
| NC   | 0  | $\frac{2}{9}$ | 0             |
| VERB | 0  | 0             | $\frac{1}{9}$ |

34

HMM et morphosyntaxe : Algorithme de Viterbi (3.)

À présent, on parcourt la matrice en sens inverse pour reconstituer la meilleure séquence.

|      | la | bonne         | fume          |
|------|----|---------------|---------------|
| ADJ  | 0  | $\frac{1}{3}$ | 0             |
| DET  | 1  | 0             | 0             |
| NC   | 0  | $\frac{2}{9}$ | 0             |
| VERB | 0  | 0             | $\frac{1}{9}$ |

34

HMM et morphosyntaxe : Algorithme de Viterbi (3.)

À présent, on parcourt la matrice en sens inverse pour reconstituer la meilleure séquence.

|      | la | bonne         | fume          |
|------|----|---------------|---------------|
| ADJ  | 0  | $\frac{1}{3}$ | 0             |
| DET  | 1  | 0             | 0             |
| NC   | 0  | $\frac{2}{9}$ | 0             |
| VERB | 0  | 0             | $\frac{1}{9}$ |

On a donc la meilleure annotation :

- la/DET bonne/NC fume/VERB
- probabilité =  $\frac{1}{9}$

Alors que cette séquence n'existe pas dans nos exemples!

34

Features

On peut bien sûr améliorer ce modèle en prenant en compte plus de features que simplement les formes.

En supposant encore que les features sont indépendantes, ça revient à considérer des probabilités d'émission composites :

$$P([0, 0]=\text{petit}, [-1, 0]=\text{Le}|\text{ADJ}) = P([0, 0]=\text{petit}|\text{ADJ}) \times P([-1, 0]=\text{Le}|\text{ADJ})$$

Cependant ce n'est pas le cadre le plus sympathique pour travailler.

- L'hypothèse d'indépendance est vite gênante
- On aimerait pouvoir utiliser des features composites

HMM

On peut penser aux modèles de Markov cachés comme un équivalent de Naïve Bayes pour l'étiquetage, avec les mêmes avantages

- Faciles et rapides pour l'entraînement et l'inférence
- Étonnamment efficace même pour peu de données

HMM

On peut penser aux modèles de Markov cachés comme un équivalent de Naïve Bayes pour l'étiquetage, avec les mêmes avantages

- Faciles et rapides pour l'entraînement et l'inférence
- Étonnamment efficace même pour peu de données

Et les mêmes inconvénients

- Interprétabilité limitée avec beaucoup de features
- Les hypothèses d'indépendance simplifient parfois trop le problème.

35

36

36

|   | MEMM  | MEMM   |
|---|---|--|
| <div>Modèles d'étiquetage</div> <div>MEMM et CRF</div>  | <p>Les <b>Maximum Entropy Markov Models</b> sont une reformulation des HMM où on modélise directement la probabilité <math>P(Y X)</math> en définissant manuellement ses composantes.</p>   | <p>Les <b>Maximum Entropy Markov Models</b> sont une reformulation des HMM où on modélise directement la probabilité <math>P(Y X)</math> en définissant manuellement ses composantes.</p> <p>Comme pour les HMM, on travaille sous l'hypothèse que les étiquettes sont émises suivant un processus Markovien, c'est à dire qu'une étiquette ne dépend directement que de l'étiquette précédente.</p> |
|   | 37  | 37   |
| Forme du modèle   | Forme du modèle   | Forme des probabilités   |
| <p>On écrit comme précédemment</p> $P(y_1, \dots, y_n   X) = \prod P(y_i   y_1, \dots, y_{i-1}, X)$   | <p>On écrit comme précédemment</p> $P(y_1, \dots, y_n   X) = \prod P(y_i   y_1, \dots, y_{i-1}, X)$ <p>Soit, avec l'hypothèse de <b>Markov</b></p> $P(y_1, \dots, y_n   X) = \prod P(y_i   y_{i-1}, X)$   | <p>On modélise les probabilités comme</p> $P(y_i   y_{i-1}, X) = \frac{e^{\phi(y_i, y_{i-1}, X)}}{\sum_c e^{\phi(c, y_{i-1}, X)}}$   |
|   | 38  | 38   |
| Forme des probabilités  | Forme des probabilités  | Inférence  |
| <p>On modélise les probabilités comme</p> $P(y_i   y_{i-1}, X) = \frac{e^{\phi(y_i, y_{i-1}, X)}}{\sum_c e^{\phi(c, y_{i-1}, X)}}$ <p>Avec</p> $\phi(y_i, y_{i-1}, X) = \sum_k w_k \phi_k(y_i, y_{i-1}, X)$ | <p>On modélise les probabilités comme</p> $P(y_i   y_{i-1}, X) = \frac{e^{\phi(y_i, y_{i-1}, X)}}{\sum_c e^{\phi(c, y_{i-1}, X)}}$ <p>Avec</p> $\phi(y_i, y_{i-1}, X) = \sum_k w_k \phi_k(y_i, y_{i-1}, X)$ <p>Où</p> <ul style="list-style-type: none"> <li>Les <math>\phi_k</math> sont des features définies manuellement (souvent binaires)</li> <li>Les <math>w_k</math> sont des poids appris</li> </ul> <p>On parle de modèle <b>log-linéaire</b>.</p> | <p>L'inférence avec un MEMM se fait exactement comme avec un HMM : on a des probas de transition, donc on peut soit appliquer une heuristique rapide soit Viterbi.</p> <p>L'apprentissage, c'est une autre paire de manches...</p>   |
|   | 39  | 40   |