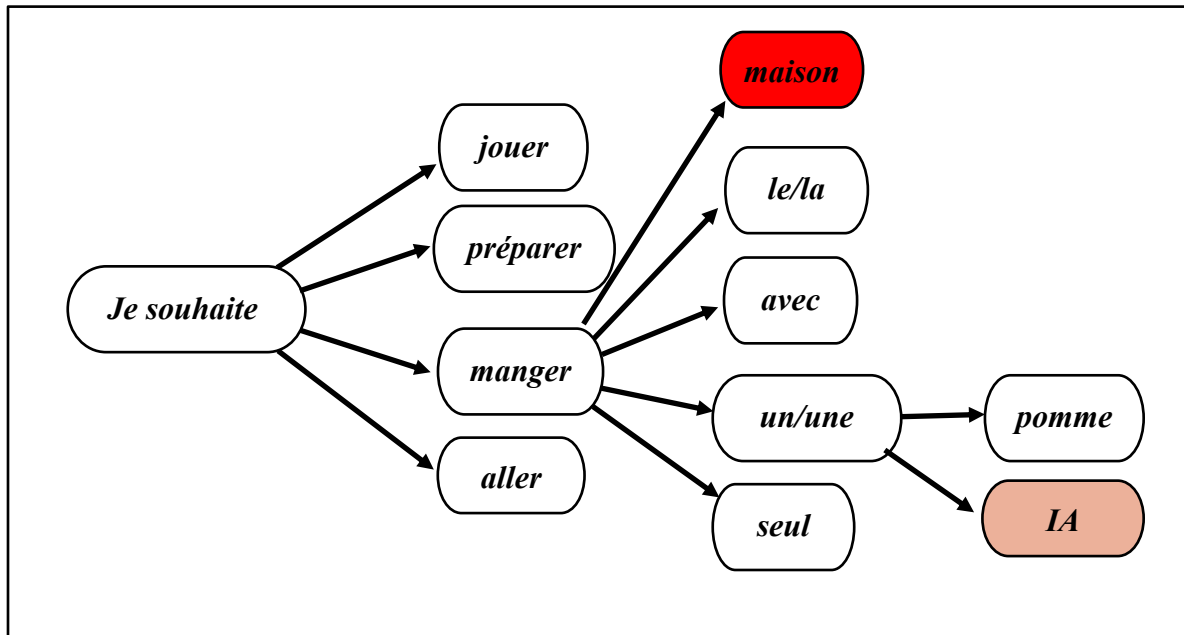


## Chapitre 2 : Modèles statistiques de langage (Langue Model)

# Ces systèmes de complétion ont un nom

- Modèle de langage - Modèle de langue – Language Model :  
vise à prédire le mot suivant dans une séquence de mots avec une certaine probabilité.



Sachant un texte (une séquence de mots – contexte-) prédire le mot qui suit

LM : (Language Model) - Modèle de langage

LLM : (Large Language Model) - Grands Modèles de langage

# Prise en compte de l'ordre des mots

- Un modèle de langage est un modèle probabiliste
  - Capable de prédire le mot manquant (suivant) dans une séquence de mots.
    - $P(w_{t+1} | w_1, w_2, \dots, w_t)$
    - Ex:
      - $P(x | \text{en, sibérie, il, fait}) \rightarrow$  trouver x ?
  - Assigne la probabilité d'observer une séquence de mots dans une langue
    - $P(w_1, w_2, \dots, w_n) \rightarrow$   
probabilité d'observer la séquence  $w_1, w_2, \dots, w_n$
    - Ex:
      - P (un garçon mange une pomme)
      - P (une pomme mange un garçon)

Comment estimer (calculer) ces probabilités ?

# Prise en compte de l'ordre des mots

- Comment estimer (prédire) ces probabilités ?

Beaucoup de textes



Puis,

un simple comptage de séquences de mots  
(fréquences relatives de séquences)

- Retour sur le calcul de probabilités

- $P(w_1, w_2, \dots, w_m) = P(w_1) * P(w_2|w_1), \dots P(w_m|w_1, w_2, \dots, w_{m-1})$

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i | w_1, w_2, \dots, w_{i-1})$$

$$P(w_i | w_1, w_2, \dots, w_{i-1}) = \frac{\text{count}(w_1, w_2, \dots, w_{i-1}, w_i)}{\text{count}(w_1, w_2, \dots, w_{i-1}, u)}$$

# Estimation du modèle

- Modèles n-grammes → n étant la taille du modèle
  - soit  $s$  une observation (un texte) de  $m$  mots  $s = w_1 w_2 \dots w_m$
  - *Unigram* – (observe/génère des séquences de 1 mot)

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i) = \frac{\text{count}(w_i)}{(\sum_{u \in V} \text{count}(u))}$$

- *bigrams* – (observe/génère des séquences de deux mots)

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1} w_i)}{\sum_{u \in V} \text{count}(w_{i-1} u)}$$

- *Trigrams* – (génère/observe des séquences de 3 mots)

$$P(w_1, w_2, \dots, w_m) = \prod_{i=1}^m P(w_i | w_{i-2}, w_{i-1}) = \frac{\text{count}(w_{i-2} w_{i-1} w_i)}{\sum_{u \in V} \text{count}(w_{i-2} w_{i-1} u)}$$

# Estimation du modèle (exemple)

- D'où viennent ces probabilités :

- **Modèle unigram**



Text ... ..algorithm ...  
in a  
Text  
Ass  
...  
Text mining ... is  
effecient ...association  
... in database ... to  
query. Text mining , ...  
Text, .. Association  
mining .... ..

Extraire les  
unigrams

Total des occurrences =30

text 10  
mining 5  
association 3  
database 3  
algorithm 2  
query 1  
efficient 1  
.....

Estimation  
du Modèle

Calcul de la probabilité d'une séquence

$P(m)=$

text ?	→	10/30
mining ?	→	5/30
association ?	→	3/30
database ?	→	3/30
algorithm?	→	2/30
query?	→	1/30
efficient?	→	1/30

# Modèle de langage - Language Model

- D'où viennent ces probabilités :
  - Modèle bi gram ( $P(w_i | w_{i-1})$ )



Text ... ..algorithm ... in  
algorithm ... to query. Text

Text mining is effecient  
association algorithm in  
database ... to query. Text  
mining , ... Text, .. Association  
rule .... ..

$P(\text{mining}/\text{text}), P(\text{efficient}/\text{mining})$

text mining 3 →  $\text{count}(\text{text}, \text{mining}) / \text{count}(\text{text}, *)$   
mining is 2 →  $\text{count}(\text{mining}, \text{is}) / \text{count}(\text{mining}, *)$   
association rule 1  
is efficient 3  
algorithm data 2  
....  
query process 1  
.....

Total des occurrences =20



- Si un événement (un mot de la séquence) n'apparaît pas dans le modèle, le modèle lui assigne une probabilité 0

$$P(s | M) = \prod_{i=1}^l P(m_i | M) = 0, \quad \text{si} \quad \exists m_i / P(m_i | M) = 0$$

- Solution : assigner des probabilités différentes de zéro aux événements (mots) absents
  - → Lissage (Smoothing)

# Techniques de lissage

- Méthodes de « discounting »
  - Laplace correction, Lindstone correction, absolute discounting, leave one-out discounting, Good-Turing method
- Techniques d'Interpolation
  - Estimations de Jelinek-Mercer, Dirichlet

# Lissage par interpolation

- Interpolation (Jelinek-Mercer)
  - Combiner le modèle  $M$  avec un modèle plus général (Modèle de référence)

$$P_{JM}(t | M) = \lambda.P_{ML}(t | M) + (1 - \lambda)P_{ML}(t / REF)$$

- Pb. “Règlage” de  $\lambda$

## — Modèle lissé (JM)

$$P_{JM}(t | d) = \alpha \times P_{ML}(t | d) + (1 - \alpha) P_{ML}(t | C)$$

$$P(t | M_c) = p(t) = \frac{tf(t, C)}{\sum_i tf(t', C)}$$

### Exemple

$$P(\text{"text retrieval"}) = P(\text{text}) * P(\text{retrieval}) = (10/30) * (0)$$

Collection 100 documents (total tf ds C (2000)

$$tf(\text{retrieval}, C) = 6 \rightarrow P_{ml}(\text{retrieval} | C) = 6/2000$$

$$tf(\text{text}, C) = 25 \rightarrow P_{ml}(\text{text} | C) = 25/2000$$