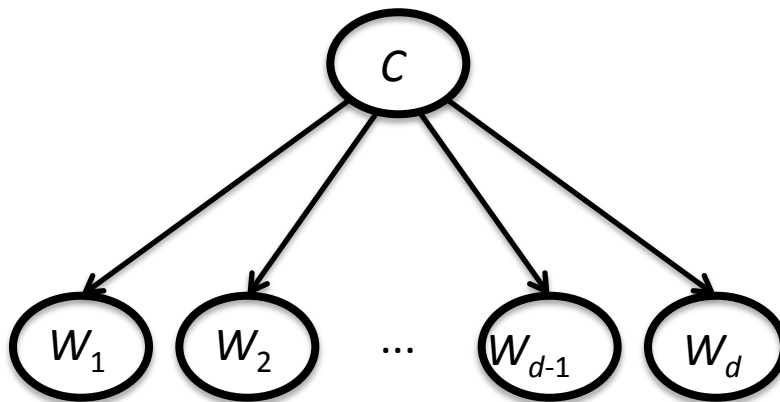


Modèle bayésien naïf multinomial

- Réseau bayésien: **modèle bayésien naïf multinomial**



- En général la **probabilité conjointe** d'un document $[W_1, \dots, W_d]$ ayant d mots et de sa catégorie C :

$$P([W_1, \dots, W_d], C) = P(C) \prod_i P(W_i \mid C)$$

Apprentissage du modèle

- Comment obtient-on les distributions $P(C)$ et $P(W_i | C)$?
 - ◆ on les obtient à partir de vraies données
 - ◆ on choisit $P(C)$ et $P(W_i | C)$ pour quelles reflètent les statistiques de ces données
- Soit un **corpus**, c.-à-d. un ensemble de T documents $\{ (D_t, c_t) \}$
 - ◆ chaque document D_t est une liste de mots $[w_1^t, \dots, w_d^t]$ de taille variable
 - ◆ c_t est la catégorie de D_t

$$\begin{aligned} P(C=c) &= (\text{nb. de documents de la catégorie } c) / (\text{nb. de documents total}) \\ &= |\{ t \mid c_t = c \}| / T \end{aligned}$$

$$\begin{aligned} P(W_i=w \mid C=c) &= \frac{\text{nb. de fois que } w \text{ apparaît dans les documents de la catégorie } c}{\text{nb. de mots total dans les documents de la catégorie } c} \\ &= \frac{\sum_{t \mid c_t=c} \text{freq}(w, D_t)}{\sum_{t \mid c_t=c} |D_t|} \end{aligned}$$

Lissage du modèle

- Selon la formule pour $P(W_i=w \mid C=c)$, un mot w aura une probabilité de 0 s'il n'apparaît jamais dans notre corpus
- Si un seul des $P(W_i=w \mid C=c) = 0$, alors tout $P(C=c, [w_1, \dots, w_d]) = 0$!
 - ◆ les mots rares vont beaucoup faire varier $P(C=c, [w_1, \dots, w_d])$ en général
- Pour éviter cette instabilité, deux trucs afin de **lisser la distribution $P(w \mid c)$**
 - ◆ on détermine un **vocabulaire** V de taille fixe, et on associe les mots qui ne sont pas dans ce vocabulaire au **symbole OOV** (*out of vocabulary*)
 - ◆ **lissage δ** : on ajoute une constante δ au numérateur, pour chaque mot

$$P(W_i=w \mid C=c) = \frac{\delta + \sum_{t \mid c_t=c} \text{freq}(w, D_t)}{\delta (|V|+1) + \sum_{t \mid c_t=c} |D_t|}$$

Lissage du modèle

- Exemple: soit le vocabulaire

$V = \{ \text{« Perceptron »}, \text{« , »}, \text{« un »}, \text{« apprentissage »} \}$

- La phrase

« Perceptron, un algorithme d'apprentissage. »

sera représentée par la liste de mots

[« Perceptron » , « , » , « un » , « OOV » , « OOV » , « apprentissage » , « OOV »]

w_1

w_2

w_3

w_4

w_5

w_6

w_7

- Les statistiques sont calculées à partir de cette représentation
 - ◆ on pourrait aussi enlever les mots « OOV » et les ignorer

Exemple

- Si, parmi tous les intra des années dernières (corpus de 426 mots)
 - ◆ « Perceptron » apparaîtrait 0 fois
 - ◆ « , » apparaîtrait 15 fois
 - ◆ « un » apparaîtrait 10 fois
 - ◆ « apprentissage » apparaîtrait 1 fois
 - ◆ « OOV » (tous les autres mots) apparaissent 400 fois
- Si on utilisait $\delta = 1$, alors
 - ◆ $P(\text{« Perceptron »} \mid C=\text{intra}) = (1 + 0) / (1 (4+1) + 426) = 1 / 431$
 - ◆ $P(\text{« , »} \mid C=\text{intra}) = (1 + 15) / (1 (4+1) + 426) = 16 / 431$
 - ◆ $P(\text{« un »} \mid C=\text{intra}) = (1 + 10) / (1 (4+1) + 426) = 11 / 431$
 - ◆ $P(\text{« apprentissage »} \mid C=\text{intra}) = (1 + 1) / (1 (4+1) + 426) = 2 / 431$
 - ◆ $P(\text{« OOV »} \mid C=\text{intra}) = (1 + 400) / (1 (4+1) + 426) = 401 / 431$

somme à 1

Prétraitement des données

- Comment choisir V
 - ◆ ne garder que **les mots les plus fréquents** (ex.: apparaissent au moins 10 fois)
 - ◆ **ne pas garder les mots trop communs**
 - » ne pas inclure la ponctuation
 - » ne pas inclure les déterminants (« un », « des », etc.)
 - » ne pas inclure les conjonction (« mais », « ou », etc.)
 - » ne pas inclure les pronoms (« je », « tu », etc.)
 - » ne pas inclure les verbes communs (« être », « avoir », « faire », etc.)
 - » etc.
 - ◆ utiliser une **forme normalisée des mots** (fusion de mots différents en un seul)
 - » **enlever les majuscules** (« Perceptron » → « perceptron »)
 - » **lemmatiser** les mots (« marchons » → « marcher »,
« suis » → « être », « est » → « être »)
- Il n'y a pas de recette universelle, le meilleur choix de V varie d'une application à l'autre