

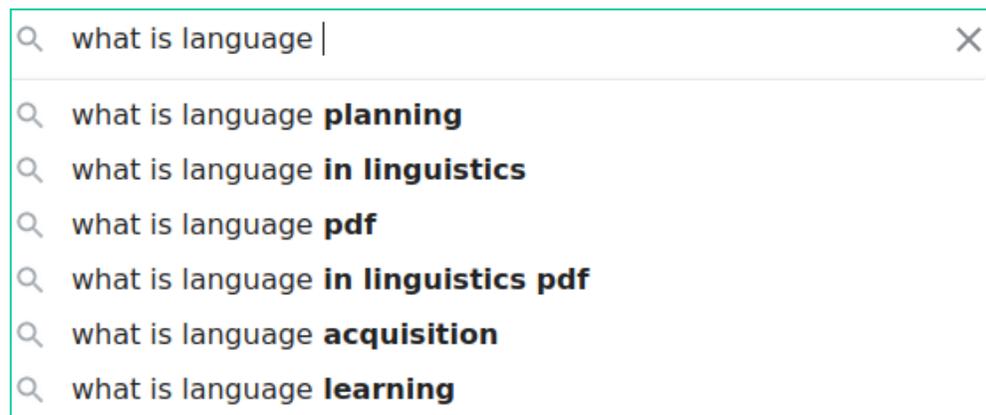
# ***Cours de Traitement Automatique du Langage***

## ***Chapitre 03 : Modèles de Langues*** ***DI5 – 2021-22***

# Introduction

## Modèle de langages

- La tâche centrale du TAL
- Objectif
  - Estimer les probabilités de séquences de mots



- Probabilité d'une phrase:  $P(S) = P(w_1, w_2, \dots, w_n)$ 
  - Auto-completion :  $P(\text{traitement automatique de l'information}) > P(\text{traitement automatique de l'eau})$
  - Traduction automatique :  $\text{My tall brother} \rightarrow P(\text{Mon grand frere}) > P(\text{Mon haut frere})$
  - Correction des fautes grammaticales :  $P(\text{Un objet qu'on puisse emporter}) > P(\text{Un objet qu'ont puisse emporter})$
  - Reconnaissance de paroles :  $P(\text{Jeudi matin}) > P(\text{Je dis matin})$
  - ...

# Introduction

## Approche probabiliste

- Probabilité d'occurrence d'un mot dans une phrase :  $P(w_n|w_1, \dots, w_{n-1})$
- Formule des probabilités composées (décomposition)
  - $P(w_1 \dots w_n) = P(w_1) \cdot P(w_2/w_1) \cdot P(w_3/w_1, w_2) \dots P(w_n/w_1, \dots, w_{n-1})$
  - $P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i|w_{<i})$
- Les probabilités sont estimées en comptant les mots dans des corpus (apprentissage)

### Exemple de probabilité d'une phrase

$$P(\text{je travaille à l'ESI}) = P(\text{je})P(\text{travaille|je})P(\text{à|je travaille}) \dots P(\text{ESI|je travaille à l'})$$

$$P(\text{ESI|je travaille à l'}) = \frac{C(\text{je travaille à l'ESI})}{C(\text{je travaille à l'})}$$

Où

$C$  est la fonction qui compte le nombre d'occurrences d'une expression dans un corpus d'entraînement

- Il faut un grand corpus pour estimer cette probabilité

# Introduction

---

## Approche probabiliste

- Il faut de grands corpus pour estimer ces probabilités !
- Mais :
  - Il y a trop de phrases possibles et pas assez de données !
  - Impossible d'avoir un corpus de texte contenant toutes les phrases.
  - Parmi les phrases qui n'apparaissent pas, certaines sont probables, d'autres non,
  - Cela s'aggrave lorsque les séquences deviennent plus longues.
  - Solution → Limiter leur longueur !

# Modèle n-gram

## Hypothèse de Markov (Markov assumption)

- Un état futur ne dépend que de l'état présent
  - $P(x_n | x_1, \dots, x_{n-1}) \approx P(x_n | x_{n-1})$
- Cas général avec k états passés
  - $P(x_n | x_1, \dots, x_{n-1}) \approx P(x_n | x_{n-k}, \dots, x_{n-1})$
- Estimation de probabilité en utilisant les N-gram ( $\rightarrow k = n - 1$ )
  - $P(x_1, \dots, x_n) \approx \prod P(x_i | x_{i-k+1}, \dots, x_{i-1})$
- Modèle uni-gram :  $P(w_n | w_1, \dots, w_{n-1}) \approx P(w_n) \rightarrow$  proba des mots dans un/des corpus
- Modèle bi-gram :  $P(w_n | w_1, \dots, w_{n-1}) \approx P(w_n | w_{n-1}) \rightarrow$  proba du mot suivant sachant un mot
- Modèle tri-gram :  $P(w_n | w_1, \dots, w_{n-1}) \approx P(w_n | w_{n-2}, w_{n-1}) \rightarrow \dots$
- Difficile d'aller bcp plus loin que le modèle Trigram.... Pourquoi ?
- Exemple Tri-grams

$$P(\text{What is language modeling}) \approx P(\text{What}) \times P(\text{is} | \text{What}) \\ \times P(\text{language} | \text{What, is}) \times P(\text{modeling} | \text{is, language})$$

# Modèle n-gram

## Estimateur de maximum de vraisemblance sur bi-grammes

$C$  est le nombre d'occurrences des N-grammes dans le corpus

$$\text{Bi-grammes : } P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

### Exemple d'un corpus d'entraînement

- $\langle s \rangle$ un ordianteur peut vous aider  $\langle /s \rangle$
- $\langle s \rangle$ il veut vous aider  $\langle /s \rangle$
- $\langle s \rangle$ il veut un ordinateur  $\langle /s \rangle$
- $\langle s \rangle$ il peut nager  $\langle /s \rangle$

- $P(\text{peut}|il) =$  

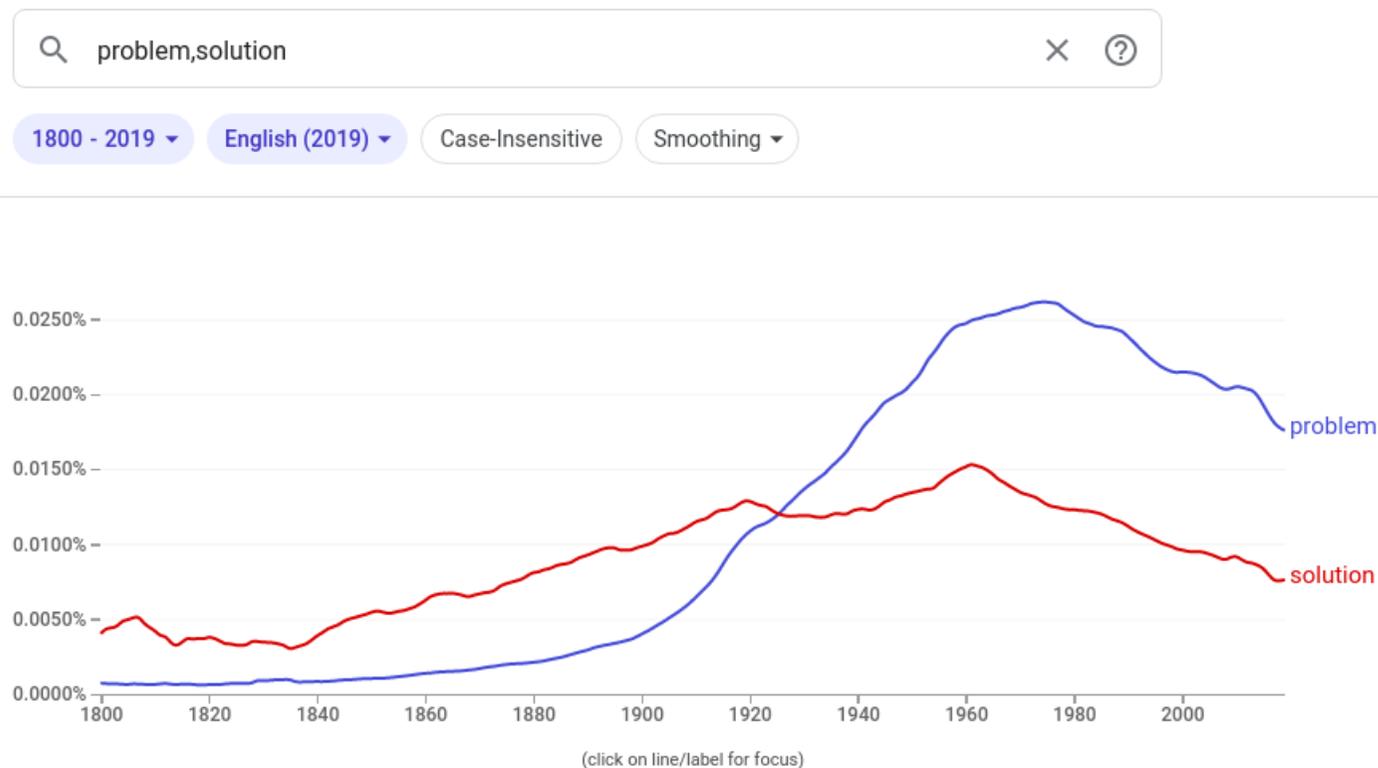
- $P(\langle s \rangle il \text{ peut vous aider } \langle /s \rangle) =$



# Modèle n-gram

- Google Books N-gram Viewer
  - <https://books.google.com/ngrams>
- Modèles prétraités a partir des livres :
  - <https://storage.googleapis.com/books/ngrams/books/datasetsv3.html>

## Google Books Ngram Viewer



# Modèle n-gram

---

Prenons un exemple : Berkeley Restaurant Project sentences

- can you tell me about any good cantonese restaurants close by
- mid priced thai food is what i'm looking for
- tell me about chez panisse
- can you give me a listing of the kinds of food that are available
- ....
- i'm looking for a good place to eat breakfast
- when is caffe venezia open during the day

Au total 9222 phrases dans ce corpus

# Modèle n-gram

---

## Enumération des bigrams

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

# Modèle n-gram

---

Normalisation selon les unigrams:

- Result: 

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

# Modèle n-gram

---

Estimation de la probabilité d'une phrase ?

$P(\langle s \rangle \text{ I want english food } \langle /s \rangle) =$

$P(\text{I} | \langle s \rangle)$

$\times P(\text{want} | \text{I})$

$\times P(\text{english} | \text{want})$

$\times P(\text{food} | \text{english})$

$\times P(\langle /s \rangle | \text{food})$

$= .000031$

# Modèle n-gram

---

## Quels intérêts ?

- $P(\text{english} | \text{want}) = .0011$
- $P(\text{chinese} | \text{want}) = .0065$
- $P(\text{to} | \text{want}) = .66$
- $P(\text{eat} | \text{to}) = .28$
- $P(\text{food} | \text{to}) = 0$
- $P(\text{want} | \text{spend}) = 0$
- $P(i | \langle s \rangle) = .25$

## Remarque pratique au passage

- On fait tous les calculs en  $\log()$  pour éviter le dépassement de capacité de précision et addition plus rapide que multiplication

$$\log(p_1 \cdot p_2 \cdot p_3 \cdot p_4) = \log p_1 + \log p_2 + \log p_3 + \log p_4$$

# Modèle n-gram - Lissage (smoothing)

---

- **Problème des N-gram absents dans le corpus d'entraînement**
  - $P(\langle s \rangle \text{il veut nager} \langle /s \rangle) = P(\text{il} | \langle s \rangle) P(\text{veut} | \text{il}) P(\text{nager} | \text{veut}) P(\langle /s \rangle | \text{nager}) = 3/4 \cdot 2/3 \cdot 0/1 \cdot 1/1 = 0$
- Plus N est grand plus le risque augmente, MAIS
- Si on utilise des petits N-gram  $\Rightarrow$  Perte de l'information
  - Les langues permettent des dépendances long terme
  - L'ordinateur que j'ai utilisé hier a Tours pendant la séance du cours a planté
- Si on utilise des grands N-gram  $\Rightarrow$  Complexité élevée du modèle
  - Il faut un corpus plus grand
  - Représentation des N-gram :  $V^N$  ou V est la taille du vocabulaire et N est le nombre de grams
- **Techniques de lissage, 2 intuitions :**
  - emprunter une petite portion des probabilités des N-grams existants pour former une probabilité aux N-gram absents (Laplace)
  - Réduire la taille des N-grams pour estimer les proba des N-grams qui sont absents (back-off)

# Modèle n-gram

## Exemple de techniques de lissage

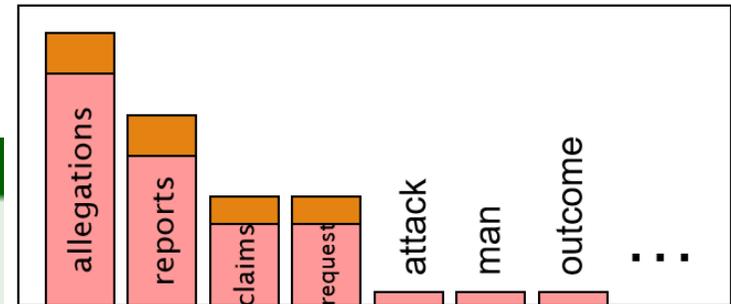
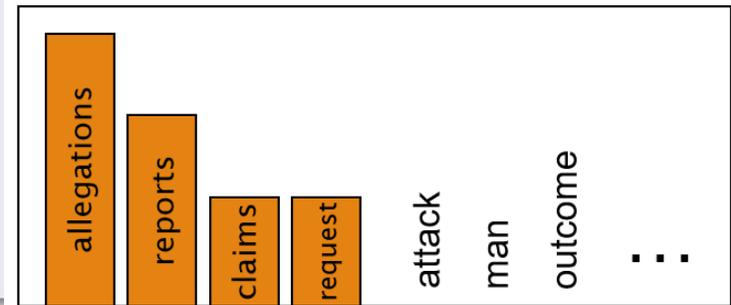
### Lissage de Lidstone (Bi-grammes comme exemple)

$$P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n) + \alpha}{C(w_{n-1}) + \alpha V}$$

Où  $V$  est la taille du vocabulaire du modèle

$\alpha = 1$  : **Lissage de Laplace**

$\alpha = 0.5$  : **Loi de Jeffreys-Perks**



### Exemple : lissage de Laplace

- Le corpus contient 8 mots différents
- Il y a  $8^2 = 64$  bi-grammes possibles  $\rightarrow$  Prenons  $V \approx 64$

$P(\langle s \rangle \text{il veut nager} \langle /s \rangle) =$

$$P(\text{il}|\langle s \rangle)P(\text{veut}|\text{il})P(\text{nager}|\text{veut})P(\langle /s \rangle|\text{nager}) = \frac{3+1}{4+64} \frac{2+1}{3+64} \frac{0+1}{1+64} \frac{1+1}{1+64}$$

# Modèle n-gram

## Exemple de techniques de lissage

$$P(\text{is language}) = 0 \Rightarrow P(\text{is language}) \approx P(\text{language})$$

### Interpolation (Tri-grammes comme exemple)

$$P_I(w_n|w_{n-2}w_{n-1}) = \lambda_3 P(w_n|w_{n-2}w_{n-1}) + \lambda_2 P(w_n|w_{n-1}) + \lambda_1 P(w_n)$$

Où  $\sum_i \lambda_i = 1$

$\lambda_3$ ,  $\lambda_2$  et  $\lambda_1$  sont estimés en utilisant un autre corpus de réglage

### Back-off de Katz (Tri-grammes comme exemple)

$$P_{BO}(w_n|w_{n-2}w_{n-1}) = \begin{cases} P^*(w_n|w_{n-2}w_{n-1}) & \text{si } C(w_{n-2}w_{n-1}w_n) > 0 \\ \alpha(w_{n-2}w_{n-1})P_{BO}(w_n|w_{n-1}) & \text{sinon} \end{cases}$$

Où :

$P^*$  est la probabilité réduite (la réduction sera distribuée sur les probabilités des N-grammes de l'ordre inférieur)

$\alpha$  est une fonction qui distribue la réduction selon le contexte

# Modèle n-gram

## Reprenons notre exemple

- Lissage de Laplace

$$P^*(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n) + 1}{C(w_{n-1}) + V}$$

Avec  $V = 1467$

	i	want	to	eat	chinese	food	lunch	spend
i	6	828	1	10	1	1	1	3
want	3	1	609	2	7	7	6	2
to	3	1	5	687	3	1	7	212
eat	1	1	3	1	17	3	43	1
chinese	2	1	1	1	1	83	2	1
food	16	1	16	1	2	5	1	1
lunch	3	1	1	1	1	2	1	1
spend	2	1	2	1	1	1	1	1

	i	want	to	eat	chinese	food	lunch	spend
i	0.0015	0.21	0.00025	0.0025	0.00025	0.00025	0.00025	0.00075
want	0.0013	0.00042	0.26	0.00084	0.0029	0.0029	0.0025	0.00084
to	0.00078	0.00026	0.0013	0.18	0.00078	0.00026	0.0018	0.055
eat	0.00046	0.00046	0.0014	0.00046	0.0078	0.0014	0.02	0.00046
chinese	0.0012	0.00062	0.00062	0.00062	0.00062	0.052	0.0012	0.00062
food	0.0063	0.00039	0.0063	0.00039	0.00079	0.002	0.00039	0.00039
lunch	0.0017	0.00056	0.00056	0.00056	0.00056	0.0011	0.00056	0.00056
spend	0.0012	0.00058	0.0012	0.00058	0.00058	0.00058	0.00058	0.00058

# Modèle n-gram

## Reprenons notre exemple

- Reconstruction du tableau à partir des nouvelles proba

$$c^*(w_{n-1}w_n) = \frac{[C(w_{n-1}w_n) + 1] \times C(w_{n-1})}{C(w_{n-1}) + V}$$

	i	want	to	eat	chinese	food	lunch	spend
i	3.8	527	0.64	6.4	0.64	0.64	0.64	1.9
want	1.2	0.39	238	0.78	2.7	2.7	2.3	0.78
to	1.9	0.63	3.1	430	1.9	0.63	4.4	133
eat	0.34	0.34	1	0.34	5.8	1	15	0.34
chinese	0.2	0.098	0.098	0.098	0.098	8.2	0.2	0.098
food	6.9	0.43	6.9	0.43	0.86	2.2	0.43	0.43
lunch	0.57	0.19	0.19	0.19	0.19	0.38	0.19	0.19
spend	0.32	0.16	0.32	0.16	0.16	0.16	0.16	0.16

- Comparaison avec les données initiales

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

# Modèles neuronaux

---

## Idées principales → CF Word Embeddings

- Vecteurs de caractéristiques continus : Chaque mot est représenté par un vecteur caractéristique de dimension  $d \ll |V|$
- Fonction de probabilité : La probabilité du mot suivant est exprimée en tant que fonction continue des caractéristiques du mot dans le contexte actuel
- Apprentissage conjoint : Les paramètres de la représentation du mot et de la fonction de probabilité sont appris conjointement.

## Continuité + hypothèse distributionnelle implique :

- qu'un petit changement dans le vecteur de caractéristiques induira un petit changement dans la probabilité.
- Par conséquent, les mises à jour causées par la présence de la phrase suivante dans les données d'apprentissage augmenteront la probabilité de toutes les phrases 'voisines'.

*A dog was running in a room*

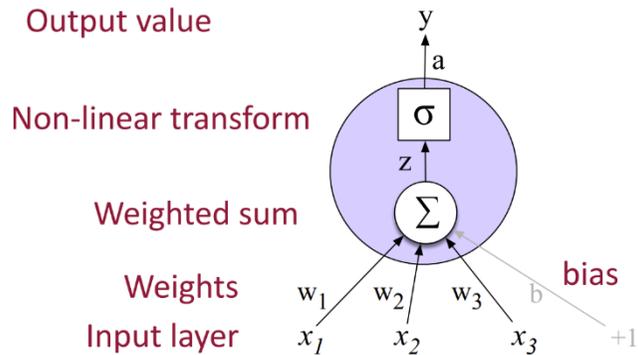
- Le fait d'avoir également la phrase suivante fera en sorte que les caractéristiques des mots (*cat*, *dog*) se rapprochent les unes des autres.

*A cat was running in a room*

# Modèles neuronaux

## Feedforward neural networks (RN à propagation avant / MLP - CF cours RN)

- Réseaux entièrement connectés → Fully-connected networks ( $\underline{W}$ ,  $\underline{U}$ )
- Fenêtre glissante de taille fixe (taille = N-Gram)
- Remplacement des “hand-crafted features” (inputs) par du “word embedding” (via  $\underline{E}$ )



output distribution

$$\hat{y} = \text{softmax}(Uh + b_2) \in \mathbb{R}^{|V|}$$

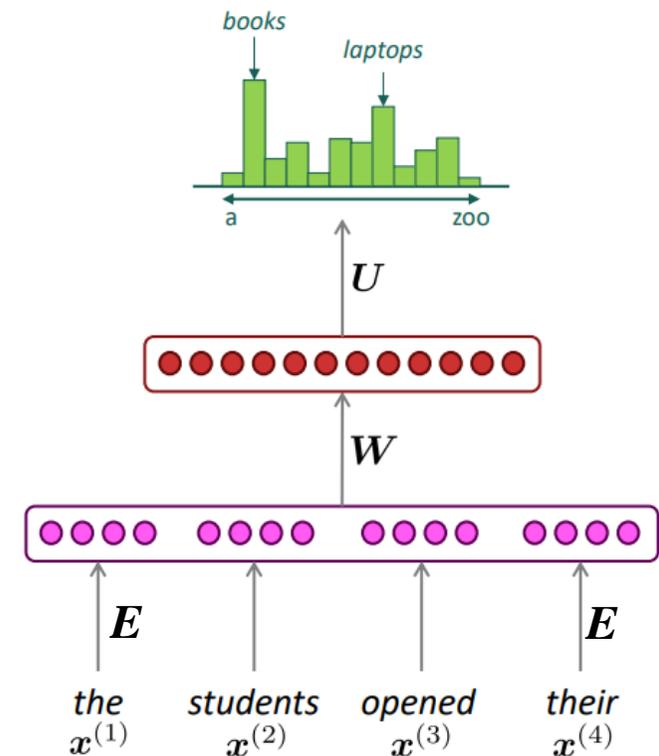
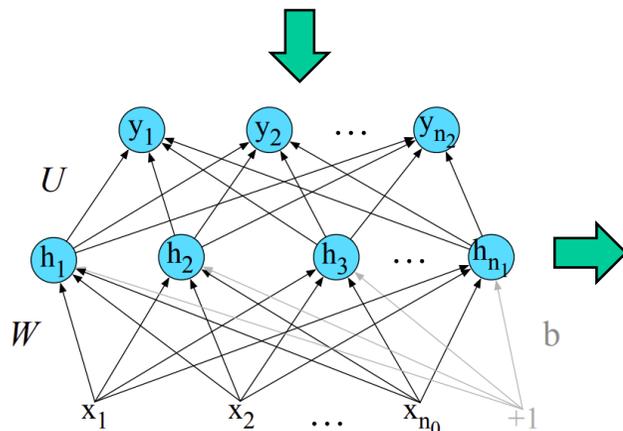
hidden layer

$$h = f(We + b_1)$$

concatenated word embeddings

$$e = [e^{(1)}; e^{(2)}; e^{(3)}; e^{(4)}]$$

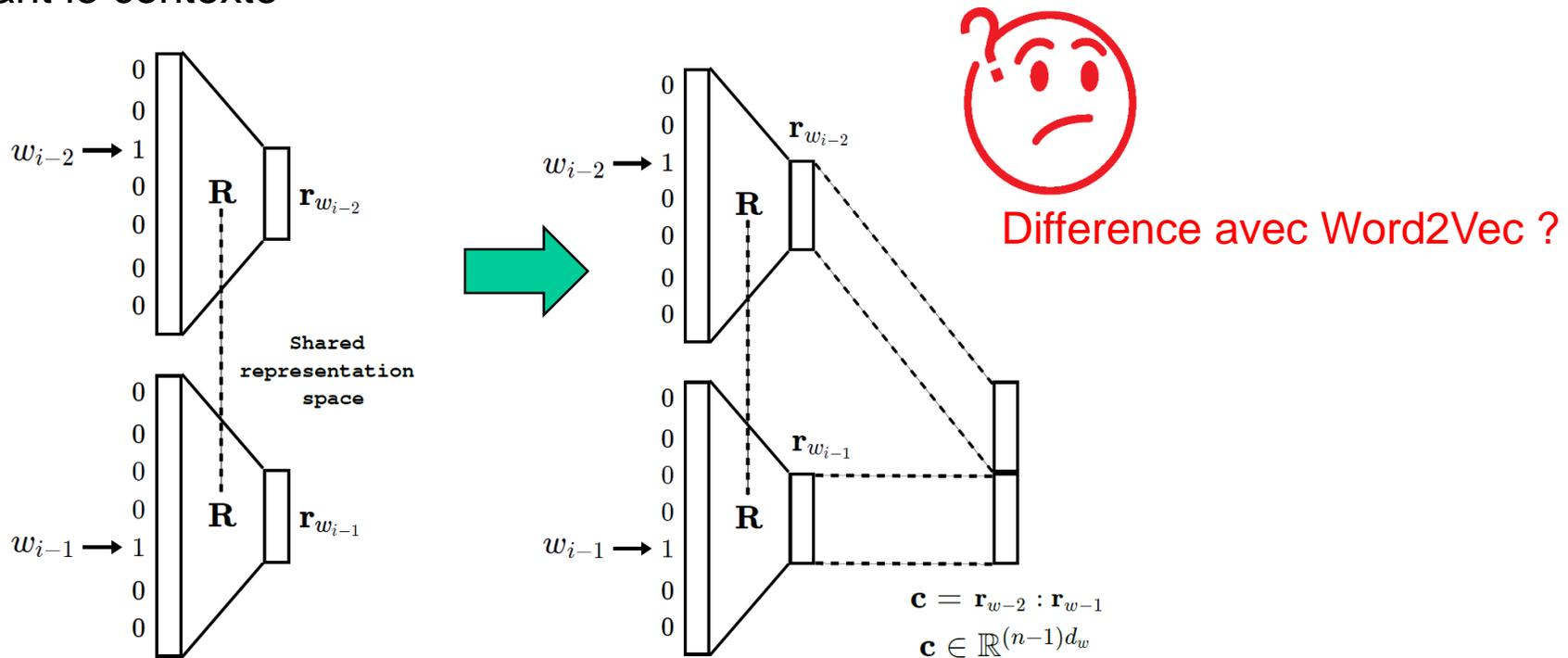
words / one-hot vectors  
 $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$



# Modèles neuronaux

## 1/ Word Embeddings de N-grams

- L'entrée est un ensemble de  $N-1$  vecteurs one-shot du vocabulaire  $V$
- Couche dense connectée à une couche plus petite, de dimension  $d_w$ .
- $R$  contient les paramètres du embedding de mots.
- Les embedding de mots **sont fusionnés (concaténés)** en un unique vecteur  $c$  représentant le contexte



2 embedding pour des Tri-grams (N=3)

# Modèles neuronaux

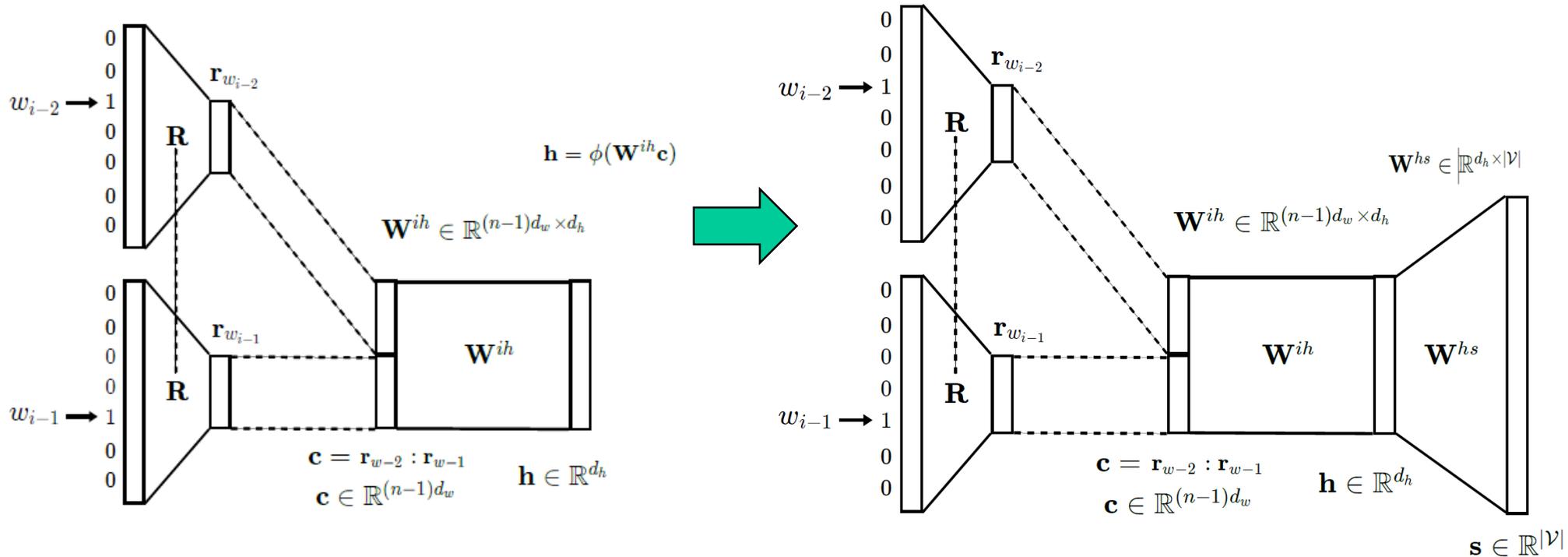
## 2/ Couches cachées → 3/ Output

- Étant donné la représentation du contexte, création d'1 représentation cachée

$$\mathbf{h} = \phi(\mathbf{W}^{ih} \cdot \mathbf{c})$$

- Étant donné  $\mathbf{h}$ , obtention de scores pour tous les mots de  $V$

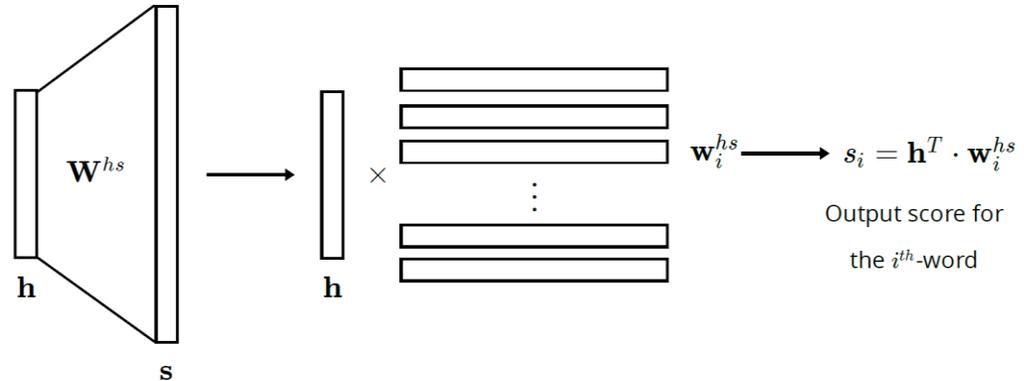
$$\mathbf{s} = \mathbf{W}^{hs} \cdot \mathbf{h}$$



# Modèles neuronaux

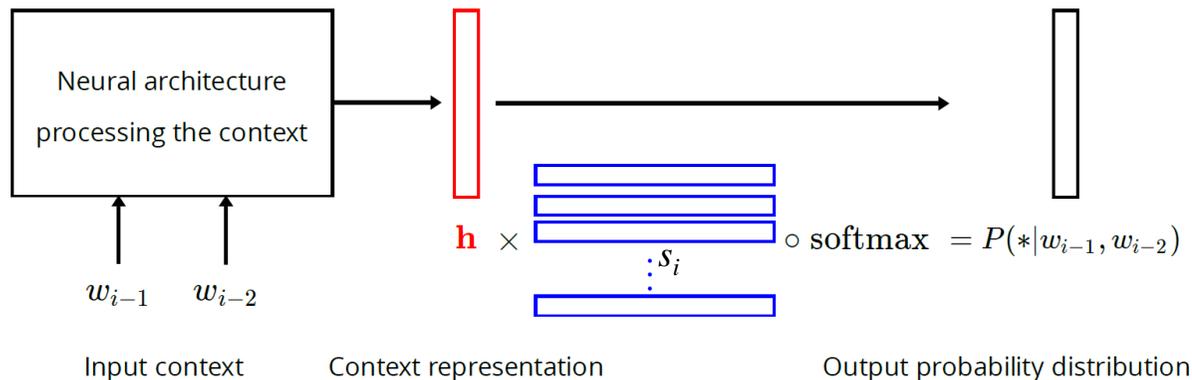
## 3/ Output

- La dernière couche peut être vue comme un produit scalaire entre  $h$  et un embedding des mots de sortie



- Ensuite, on estime les probabilités pour tous les mots d'indices  $k \in V$  étant donnés  $w_{i-1}, w_{i-2}$  avec la fonction softmax :

$$P(k|w_{i-1}, w_{i-2}) = \frac{\exp(\mathbf{s}_k)}{\sum_{l=1}^{|\mathcal{V}|} \exp(\mathbf{s}_l)} = \frac{\exp(\mathbf{h}_i^T \mathbf{w}_k^{hs})}{\sum_{l=1}^{|\mathcal{V}|} \exp(\mathbf{h}_i^T \mathbf{w}_l^{hs})}$$



# Modèles neuronaux

## 4/ Apprentissage

- $\theta = (R, W^{ih}, W^{hs})$  sont les paramètres du modèle **appris conjointement**
- Nous voulons que la distribution de probabilité de la sortie du modèle à chaque pas de temps (i) se rapproche de la vérité terrain
- En minimisant l'entropie croisée, à chaque étape, nous minimisons la log-vraisemblance négative de l'échantillon de données

$$d\left( \begin{array}{c} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{array}, \begin{array}{c} P_{\theta}^{(i)}(1|w_{<i}) \\ \vdots \\ P_{\theta}^{(i)}(w_i|w_{<i}) \\ \vdots \\ P_{\theta}^{(i)}(|\mathcal{V}||w_{<i}) \end{array} \right) ? \longrightarrow \begin{array}{l} \text{Cross-entropy}(P_*^{(i)}, P_{\theta}^{(i)}(*|w_{<i})) \\ = - \sum_{k=1}^{|\mathcal{V}|} P_*^{(i)}(k) \log(P_{\theta}^{(i)}(k|w_{<i})) \\ = -\log(P_{\theta}^{(i)}(w_i|w_{<i})) ! \end{array}$$

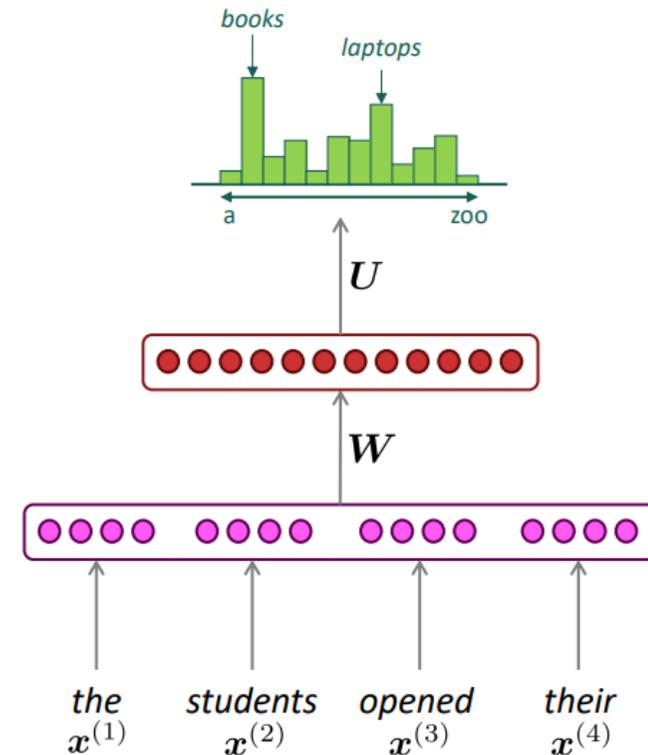
- Pour tous les échantillons de données  $D = (w^i)_{i=1..N}$ , l'objectif est l'estimation de la plus grande vraisemblance (MLE).

$$NLL(\theta) = - \sum_{i=1}^N \log(P_{\theta}^{(i)}(w_i|w_{<i}))$$

# Modèles neuronaux

## Bilan sur les Feedforward neural networks pour les ML

- Améliorations par rapport aux modèles n-gram :
  - Pas de problème de sparsité (lissage inutile)
  - Il n'est pas nécessaire de stocker tous les n-grams observés
- Problèmes restants :
  - La taille de fenêtre fixe est souvent trop petite
  - L'agrandissement de la fenêtre agrandit  $W$
  - La fenêtre n'est jamais assez grande !
  - $x_1$  et  $x_2$  sont multipliés par des poids complètement différents dans  $W$
  - Aucune symétrie dans la façon dont les entrées sont traitées
- Nous aurions besoin d'une architecture neuronale qui peut traiter n'importe quelle longueur de chaînes en entrée.



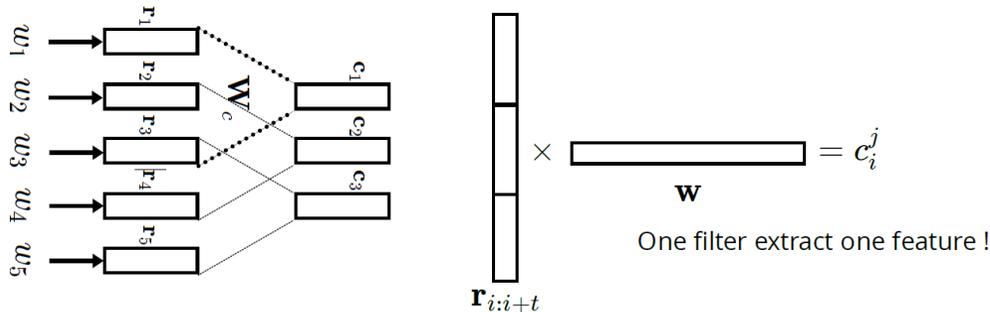


Explications ?

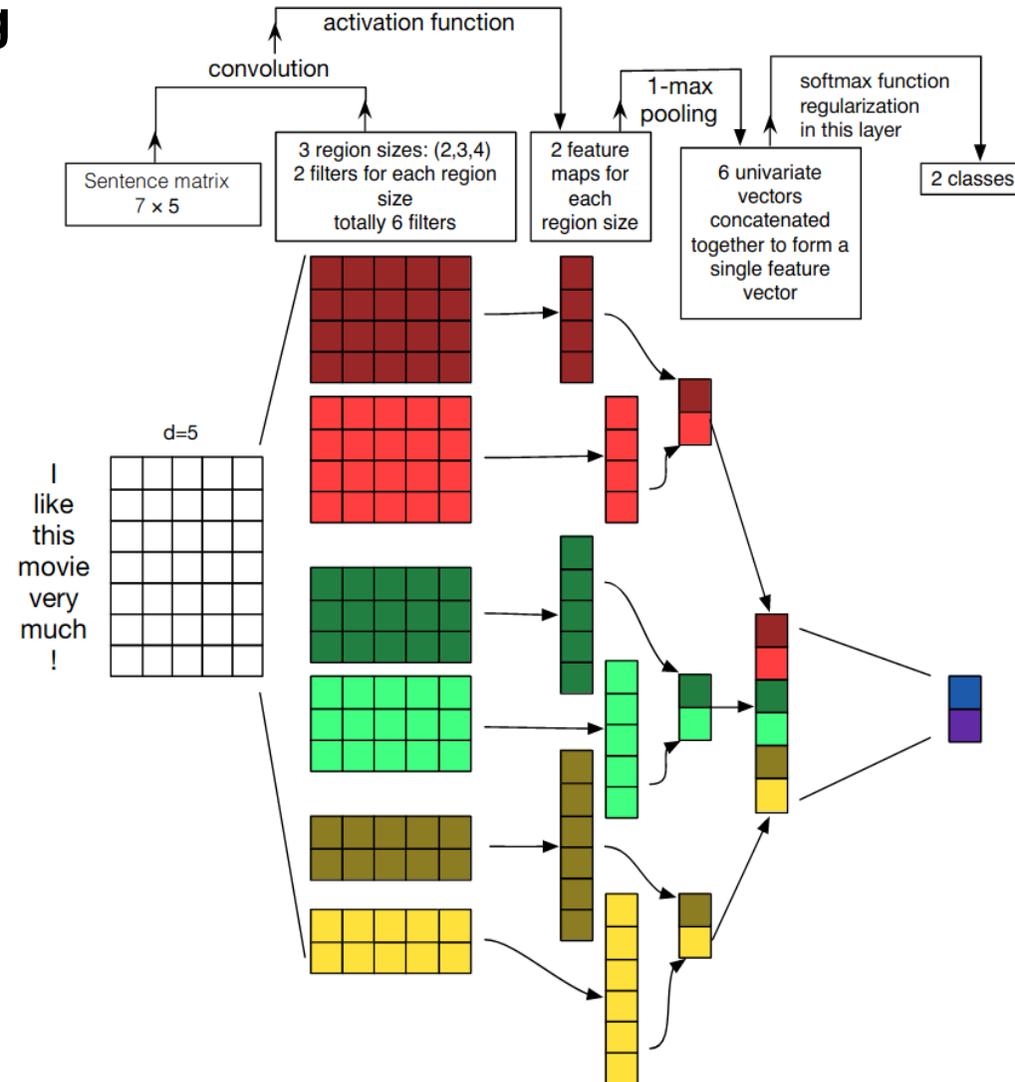
# Modèles neuronaux

## Approche avec Convolution et Pooling

- Les tâches de classification sur des séquences nécessitent une représentation globale
- L'idée est d'agglomérer consécutivement les caractéristiques obtenues sur des fenêtres de N-grams.
  - Convolution** : Calculer plusieurs représentations pour chaque fenêtre (sous-séquence possibles d'une certaine longueur)



- Pooling** : Regrouper des éléments dans une représentation plus synthétique...
  - max, mean, ...



# Modèles neuronaux

## RN Récurrents

- Idée : Les poids  $W$  appliqués aux entrées restent fixes au cours du temps
- **Présence de boucle dans les connexions entre neurones**
- **Les poids des connexions récurrentes ( $U$ ) sont fixes également**
- A un instant  $t$ , on calcule l'état  $y_t$  en se basant sur l'état précédent  $y_{t-1}$  et sur l'entrée actuelle  $x_t \rightarrow y_{t-1}$  **au lieu des mots contextes**
- Ce type de RNN peut se représenter sous la forme d'un réseau FW adapté  $\rightarrow$  utilisation des modes d'apprentissage classiques

$$e_t = E^T x_t$$

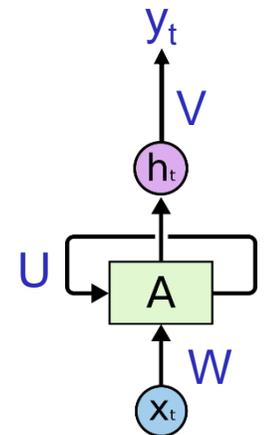
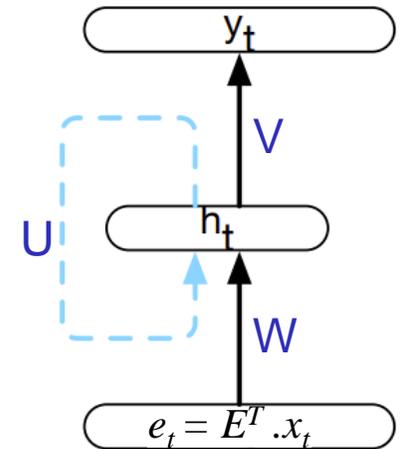
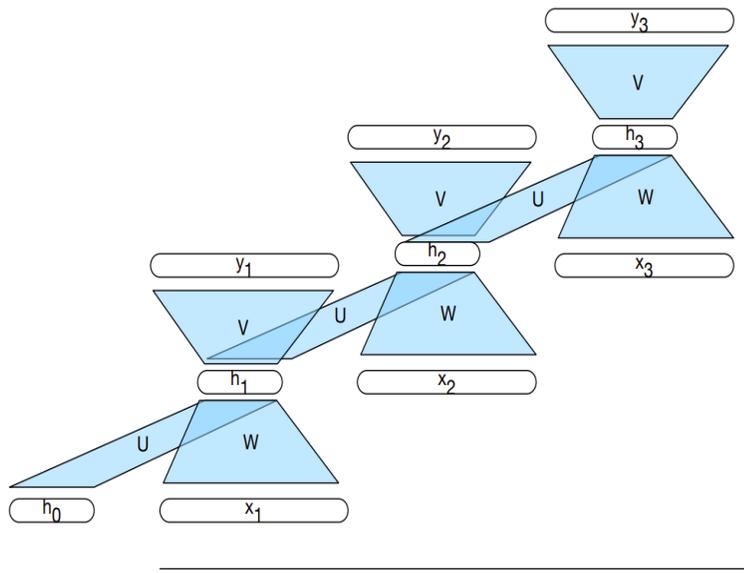
$$h_t = g(Uh_{t-1} + We_t)$$

$$y_t = \text{softmax}(Vh_t)$$

$$P(w_{1:n}) = \prod_{i=1}^n P(w_i | w_{1:i-1})$$

$$= \prod_{i=1}^n y_{w_i}^i$$

- $g() = \text{sigmoid}() \text{ ou } \text{tanh}()$



# Modèles neuronaux

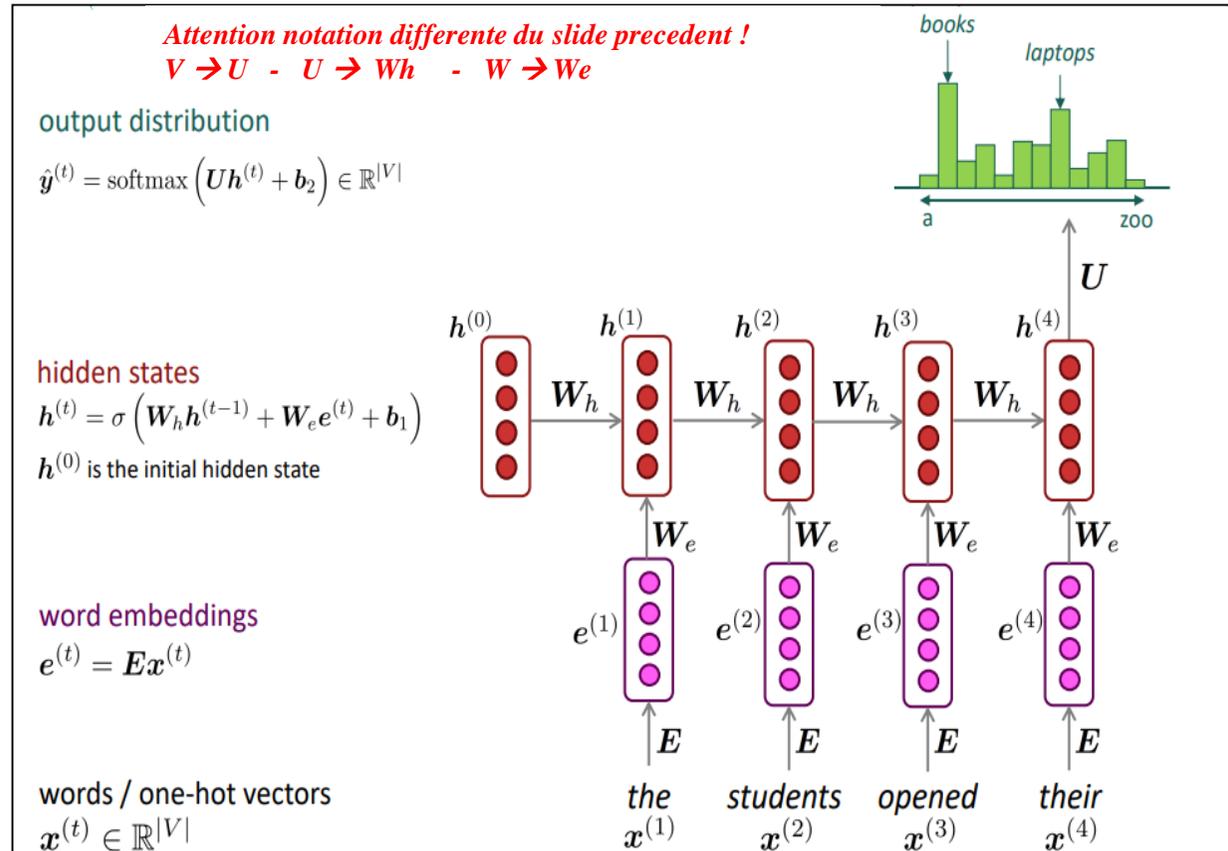
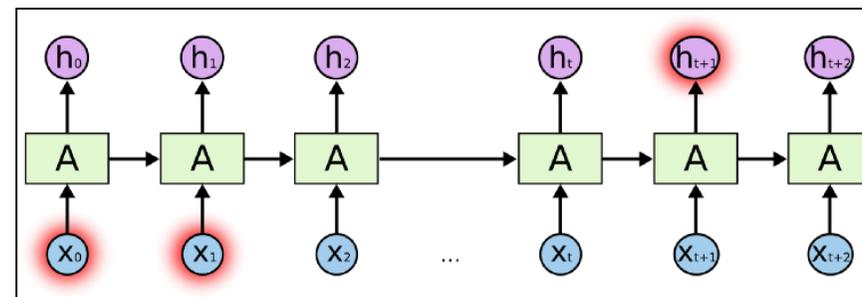
## RN Récurrents

### Avantages du RNN :

- Longueur d'entrée quelconque
- La taille du modèle n'augmente pas pour contexte d'entrée plus long
- Les mêmes poids sont appliqués à chaque pas de temps → symétrie dans la façon dont les entrées sont traitées
- Le calcul à l'étape t peut (en théorie) utiliser des informations de plusieurs étapes passées

### Inconvénients des RNN :

- Le calcul récurrent est lent
- Apprentissage plus difficile (descente de gradient)
- En pratique, il est difficile d'accéder à l'information de plusieurs étapes en arrière

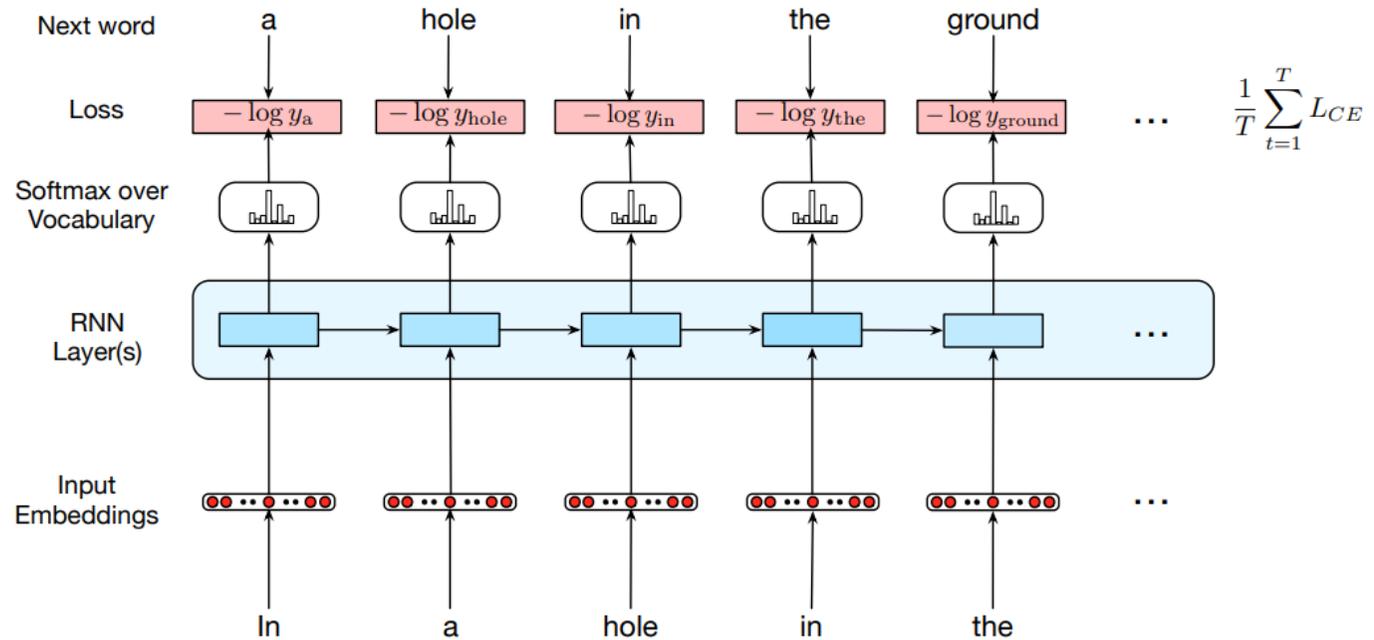


# Modèles neuronaux

## RN Récurrents

### Phase d'Apprentissage

- Loss = Moyenne des cross-entropies calculées sur une sequence



# Evaluation des modeles

## • Mesure de Perplexité

- Mesure l'incertitude : dans quelle mesure le modèle est-il surpris par de nouvelles données ?
- Mesure la qualité de prédiction d'un modèle sur un corpus de test
- Utilise la probabilité estimée sur chaque mot d'un corpus de test de taille T
- Maximisation des probabilités de l'ensemble des mots du corpus
- Le modèle avec une **perplexité minimale est le meilleur**

$$\text{perplexity} = \prod_{t=1}^T \left( \frac{1}{P_{\text{LM}}(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} \right)^{1/T}$$

← Normalized by number of words

⏟  
Inverse probability of corpus, according to Language Model

- Elle correspond à  $\exp(L(\theta))$  (Loss cross-entropique sur le corpus de test – cf slide 23)

$$\prod_{t=1}^T \left( \frac{1}{\hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)}} \right)^{1/T} = \exp \left( \frac{1}{T} \sum_{t=1}^T -\log \hat{\mathbf{y}}_{\mathbf{x}_{t+1}}^{(t)} \right) = \exp(L(\theta))$$

# Evaluation des modeles

---

- **Evaluation intrinsèque**
  - Evaluer le modèle par rapport à sa représentation du langage
  - Exemple, comparer deux modèles en se basant sur leur capacité à bien représenter un dataset de test
  - Ne garantit pas une bonne performance du modele pour une autre tache
- **Evaluation extrinsèque**
  - Evaluer le modèle par rapport à une autre tache
  - Les modèles de langages sont très souvent utiliser pour évaluer des modèles conçus pour d'autres taches
  - Exemple, la qualité de traduction automatique en utilisant un modèle de langage
  - Evaluation souvent très couteuse

# Bilan ML

---

## Modèles de Langues

- Tache centrale (prédiction mot suivant = classification multi-classes) pouvant servir de base pour de nombreuses autres tâches de classification : POS, NER, Sentiment, ...
- Modèles N-Grams
- Réseaux neuronaux à propagation avant, convolutionnels ou récurrents

## Limites des RN classiques pour la création de Modèles de Langues

- Taille limitée du contexte dans les réseaux à propagation avant
- Problème de disparition du gradient dans les réseaux récurrents
  - Limiter la taille du contexte
- Mots hors vocabulaire
  - Limiter le vocabulaire et marquer le reste des mots par ⟨UNK⟩

**Utilisation d'architectures / réseaux plus avancés : LSTM, GRU, Transformers, ...**

**→ CF suite du cours...**

# Bibliographie

---

- Eisenstein, J. (2018). Natural Language Processing.
- Jurafsky, D. (2020). Language modeling. Cours "CS 124: From Languages to Information". Stanford University.
- Mikolov, T., Karafiat, M., Burget, L., Cernocký, J., and Khudanpur, S. (2010). Recurrent neural network based language model. In Proc. of Inter speech.  
[http://www.fit.vutbr.cz/research/groups/speech/publi/2010/mikolov\\_interspeech2010\\_IS100722.pdf](http://www.fit.vutbr.cz/research/groups/speech/publi/2010/mikolov_interspeech2010_IS100722.pdf).
- Smith, N. (2020). Neural models. Cours "Natural Language Processing (CSE 517)". University of Washington.