

N-GRAM MODELS

Supponiamo di voler allenare un LM in grado di generare frasi in inglese composte da 4 lettere. Per esempio:

$$P[\text{We will rock you}] = P[\text{we}] \cdot P[\text{will} | \text{we}] \cdot P[\text{rock} | \text{we, will}] \cdot P[\text{you} | \text{we, will, rock}]$$

Per semplicità, assumiamo $n=2$. Da cui:

- 1) La probabilità di una parola dipende solo dalla parola che la precede, se presente
- 2) La probabilità di una parola, non dipende dalla posizione della parola nella frase.

Esempio:

$$1) P[\text{rock} | \text{we, will}] = P[\text{rock} | \text{will}]$$

$$2) P[w_3 = \text{rock} | w_2 = \text{will}] = P[w_{i+1} = \text{rock} | w_i = \text{will}] \quad 1 < i < 3$$

In generale la probabilità frase diventa:

$$P[\text{we}] \cdot P[\text{will} | \text{we}] \cdot P[\text{rock} | \text{will}] \cdot P[\text{you} | \text{rock}]$$

Estendiamo per $n > 2$:

- 1) La probabilità di una parola dipende solo dalle $n-1$ parole che la precedono, se presenti
- 2) La probabilità di una parola, non dipende dalla posizione della parola nella frase.

FORMALMENTS

Un n -gram model basato su V è un LM che assegna a $w_1 \dots w_k$ le seguenti probabilità:

1) se $n=1$:

$$P[w_1 \dots w_k] = \prod_{i=1}^k P[w_i] : \text{UNIGRAM}$$

2) se $n > 1$:

$$P[w_1 \dots w_k] = P[w_1] \prod_{i=2}^k P[w_i | w_{i-1} \dots w_1]$$

Che è circa uguale:

$$\approx P[w_1] \prod_{i=2}^k P[w_i | \underbrace{w_{i-1} \dots w_{i-n+1}}_{\text{le } n \text{ parole precedenti}}]$$

Stimiamo le probabilità dell' n -gram:

$$P[w_i] \approx \frac{\text{count}(w_i)}{\sum_{w \in V} \text{count}(w)} \rightsquigarrow \text{tf}$$

$$P[w_j | w_i \dots w_{j-1}] = \frac{\text{count}(w_i \dots w_{j-1} w_j)}{\sum_{w \in V} \text{count}(w_i \dots w_{j-1} w)}$$

Ovvero il numero di volte in cui è presente nel corpus la sequenza analizzata seguita da w_j diviso il numero di volte quella stessa sequenza segue qualsiasi parola (Inclusa w_j)

PAROLE NON PRESENTI NEL CORPUS

Per quanto riguarda parole usate a tempo di valutazione che non sono state viste a tempo di training?

The out of vocabulary words?

Dobbiamo introdurre un nuovo token speciale **<UNK>** per le parole sconosciute ma presenti nel evaluation corpus. Bisogna convertire ogni parola non vista a tempo di training ma presente nell'evaluation corpus con **<UNK>**, e riaviamo l'evaluation come training per per stimare la probabilità del token **<UNK>**.

PROBLEMA DELLA SPARSITA'

Il problema più grande degli n-gram model è la **sparsità**. Per evitarla si usa uno **smoothing filter**. Il filtro più comune è il **LAPLACE SMOOTHING** (o anche detto **add-one smoothing**).

Questo filtro per il bi-gram si presenta così:

$$P[w_i | w_{i-1}] = \frac{\text{count}(w_i, w_{i-1}) + 1}{\sum_{w \in V} \text{count}(w_{i-1}, w) + |V|}$$

È il filtro più facile.

LEARNING A BIGRAM

L'obiettivo, nel bigram, è calcolare:

$$P[w_i | w_{i-1}] = \frac{\text{count}(w_{i-1}, w_i)}{\sum_{w \in V} \text{count}(w_{i-1}, w)}$$

Per allenare il modello possiamo organizzare il tutto nella matrice dei pesi $W \in \mathbb{Z}^{V \times V}$ dove:

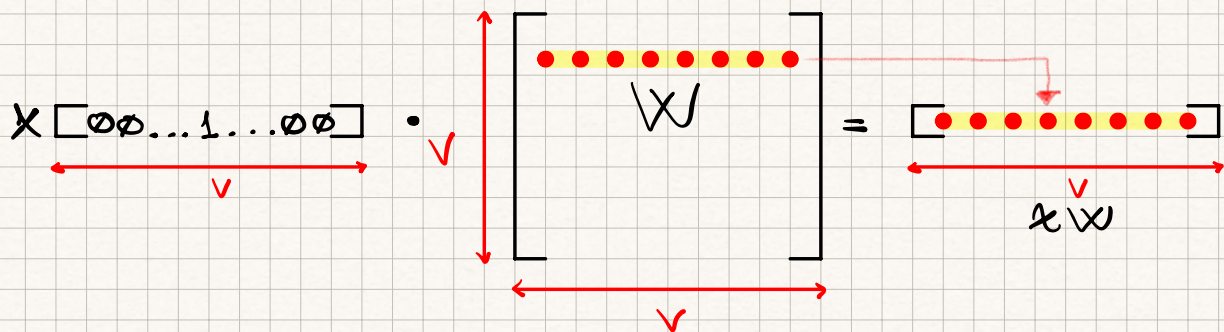
$$w_{ij} = \text{count}(w_i, w_j)$$

ovvero il numero di volte in cui w_i segue w_j nel corpus.
 W è la **matrice delle co-occorrenze**. Dove è facile intuire che:

$$\sum_{j=1}^V w_{ij} = \text{count}(w_i)$$

ONE-HOT VECTOR

Il modo più semplice per rappresentare matematicamente una parola è il one-hot vector $x \in \mathbb{Z}^{1 \times V}$ (Vettore riga). Che è un vettore pieno di zeri con un 1 a indicare la posizione in V della parola w_i : $x_i = 1$. Che poi indica la parola stessa in V .



X di fatto seleziona l' i -esima riga di W , corrispondente alla parola i : $x_i = 1$.

$$P[w_j | w_i] = \frac{\text{Count}(w_i w_j)}{\text{Count}(w_i)} = \frac{W_{ij}}{\sum_{j=1}^V W_{ij}} = \frac{W_{ij}}{w_j}$$

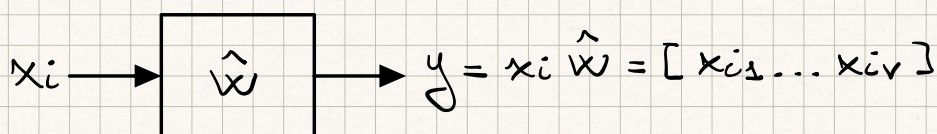
Parola in V
elemento matrice W

LEARNING THE MATRIX W

Possiamo imparare la matrice W , inizializzandola, randomicamente.

La nostra rete neurale, prenderà una parola in ingresso w_i e produrrà in output la parola w_j (che capita la più probabile per seguire w_i). Rappresenteremo w_i e w_j con i corrispondenti one-hot vector x_i e x_j .

La nostra NN ha un singolo livello lineare $\hat{W} \in \mathbb{R}^{V \times V}$ (Per semplicità incorporiamo il bias a \hat{W}):



Nota: dato che \hat{W} contiene inizialmente cose a caso (positive ma anche negative) non possiamo direttamente associare y ai conteggi delle parole come nella W ideale. Per tanto ne facciamo l'esponente per avere un qualcosa di simile che possiamo interpretare come conteggio:

$$y = x_i \hat{W} = [\hat{W}_{i1} \dots \hat{W}_{iV}] \rightarrow \text{log-count} \circ \text{logit}$$

$$z = [e^{w_{i1}}, \dots, e^{w_{iV}}] \rightarrow \text{"Counts"}$$

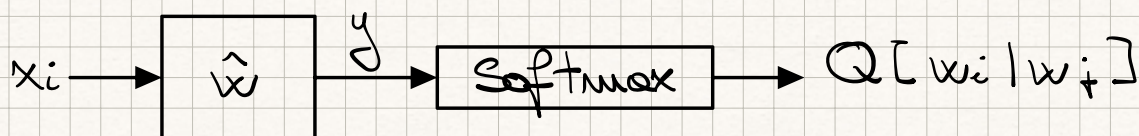
Ora possiamo **normalizzare** le componenti di z per poter avere una distribuzione di probabilità di V . Possiamo assumere che questa distribuzione di probabilità sia il nostro LM approssimato Q :

$$Q[w_j | w_i] = \frac{\hat{e}^{w_{ij}}}{\sum_{j=1}^V \hat{e}^{w_{ij}}}$$

Abbiamo fatto tutta sta menata per riuscire ad avere pesi la cui somma è 1. Queste due operazioni di **"esponenziale"** e **normalizzare** vengono sempre insieme e le riassumiamo nella singola operazione di **SOFTMAX**:

$$y = x_i \hat{w}$$

$$Q[w_i | w_j] = \text{Softmax}(y)$$



La loss da usare è **NLL**:

$$\mathcal{L}^{bi} = -\frac{1}{T} \sum_{j=1}^T \log Q[w_i | w_j]$$