

LEARNING A LANGUAGE MODEL

Assumiamo di avere un **training corpus**. Per tale corpus sia $P(s)$ il **LM perfetto**. C'è un modo per calcolare un LM approssimato $Q(s)$ più vicino possibile a $P(s)$, usando il training corpus?

Possiamo raggiungere il nostro obiettivo risolvendo il seguente problema di ottimizzazione:

$$\underset{Q}{\text{argmin}} L(P, Q)$$

Dove L è la **loss function** che specifica l'errore che abbiamo usando $Q(s)$ invece di $P(s)$.

Dobbiamo trovare L . Essa deve confrontare due distribuzioni di probabilità, per farlo usiamo l'**entropia**.

INFORMATION ENTROPY

Dato una variabile aleatoria discreta X , con probabilità $P(x)$, definiamo l'**information content** $I[x]$ di un outcome x come:

$$I[x] = -\log P[x]$$

La $I[x]$ può essere vista come una funzione $I(x)$, così che possiamo considerare la variabile aleatoria (trasformata) $I(x)$ e calcolarne l'**expectation** $E[I(x)]$, che è anche chiamata **entropia** $H(P)$ della variabile aleatoria X con prob. P :

$$H(P) = E[I(x)] = E[-\log P(x)] = -\sum_x P[x] \log P[x]$$

Nel caso dei LM, questo concetto di $H(P)$ va adattato alla sequenza di parole dipendenti dalla perplexity PP. Bisogna tradurre l'entropia dell'informazione $H(P)$ con un'altra detta **entropia del linguaggio**.

CROSS-ENTROPY di $Q(x)$

Vogliamo capire quanta informazione $Q(x)$ copia, rispetto al modello perfetto $P(x)$. Per farlo calcoliamo la Cross-entropy di $Q(x)$ rispetto a $P(x)$

$$H(P, Q) = - \sum_x P[x] \log Q[x]$$

$Q[x]$ è usata solo come sorgente dell'informazione. Dobbiamo calcolare la perdita di informazione nell'usare Q al posto di P . Calcoliamo la seguente differenza detta **Loss-Information Content**:

$$\begin{aligned} H(P, Q) - H(P) &= \\ - \sum_x P[x] \log Q[x] + \sum_x P[x] \log P[x] &= \\ = \sum_x P[x] (\log P[x] - \log Q[x]) &= \\ = \sum_x P[x] \log \frac{P[x]}{Q[x]} \equiv D_{KL}(P \parallel Q) \end{aligned}$$

Questa $D_{KL}(P \parallel Q)$ è detta **relative entropy** o anche **KULLBACK-LIBLER (KL) DIVERGENCE**.

Non è una metrica ma una **semi-metrica**.

Disuguaglianza di Gibb

$$D_{KL}(P \parallel Q) \geq 0$$

Dove l'uguaglianza si verifica solo se $P \equiv Q$.

In generale D_{KL} non è simmetrica:

$$D_{KL}(P \parallel Q) \neq D_{KL}(Q \parallel P)$$

Dunque NON È UNA DISTANZA

ENTROPIA DI UN LM

Abbiamo visto come il nostro obiettivo sia quello di confrontare due distribuzioni di probabilità P e Q .

Assumiamo di avere n dati, e consideriamo la sequenza casuale di parole $S_n = w_1 \dots w_n$.

Assumiamo poi $P(S_n)$ come LM di S_n . Invece di calcolare l'entropia di S_n , che dipende da n , calcoliamo la **PRE-WORD-ENTROPY** (o **entropy rate**):

$$H'(S_n) = -\frac{1}{n} \sum_{S_n} P[S_n] \log P[S_n]$$

Che è una somma di tutte le frasi di lunghezza n .
Definiamo ora **l'entropia di un LM** $P(s)$:

$$H(P) \triangleq \lim_{n \rightarrow \infty} H'(S_n) = -\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{S_n} P[S_n] \log P[S_n]$$

Perciò noi non conosciamo $P(s)$ ma solo la sua approssimazione $Q(s)$. Per cui definiamo la **cross entropy** $H(P, Q)$:

$$H(P, Q) = - \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{S_n} P[S_n] \log Q[S_n]$$

Ancora abbiamo nella formula $P(s)$ che non siamo in grado di calcolare.

TEOREMA SHANNON-McMILLAN-BREIMAN

Se un LM è **stazionario** e **ergodico** allora si può calcolare la cross-entropy:

$$H(P, Q) = - \lim_{n \rightarrow \infty} \sum_{S_n} \frac{1}{n} \log Q[S_n]$$

Con S_n **singola** sequenza randomica di parole.

Stationary LM

Un LM si dice stazionario sse:

$$P[w_{t+1} \dots w_{t+n}] = P[w_1 \dots w_n]$$

Per ogni sequenza di parole e per ogni parametro di shift t .
Il che vuol dire che la distribuzione di probabilità per le parole a tempo t non cambia a tempo $t+1$. Modelli come il bigram sono stazionari. I linguaggi naturali non lo sono perché la probabilità della prossima parola può dipendere da eventi (o parole) arbitrariamente distanti nel tempo. Noi, infatti, stiamo **approssimando**.

Ergodic LM

Un LM è definito ergodico se l'expectation di una parola $\mathbb{E}[w_i]$ su $P(S_n)$ può essere calcolata come la media temporale della singola lunga sequenza casuale

$w_1 w_2 \dots$

In realtà la parola w_i potrebbe non apparire mai in una lunga sequenza, portando $\mathbb{E}[w_i] = 0$. Questa proprietà vuole che $\mathbb{E}[w_i] > 0$.

NEGATIVE LOG LIKELIHOOD

Quindi, senza conoscere P , e data una singola lunga sequenza $w_1 \dots w_T$ generata tramite $P(s)$, si approssima:

$$H(P, Q) \approx -\frac{1}{T} \log Q[w_1 \dots w_T]$$

Detta **negative log likelihood (NLL)** di $Q(S_n)$

PERPLEXITY e CROSS ENTROPY

La lunga sequenza $w_1 \dots w_T$ è detta **evaluation corpus** e ne possiamo calcolare la perplexity $PP(w_1 \dots w_T)$.

Usando un trucco matematico possiamo dire:

$$\begin{aligned} PP(w_1 \dots w_T) &= \exp(\log PP(w_1 \dots w_T)) \\ &= \exp(\log Q[w_1 \dots w_T]^{-1/T}) = \exp(-\frac{1}{T} \log Q[w_1 \dots w_T]) \\ &= \exp(H(P, Q)) \end{aligned}$$

Quindi perplexity e cross-entropy sono strettamente correlate, ma non uguali! Da notare che crescono e decrescono insieme.

MINIMIZZIAMO LA LOSS FUNCTION

Il problema si riduce nel risolvere:

$$\begin{aligned} & \underset{Q}{\operatorname{argmin}} (L(P, Q)) \\ & \downarrow \\ & \underset{Q}{\operatorname{argmin}} (D_{KL}(P \parallel Q)) \\ & \downarrow \\ & \underset{Q}{\operatorname{argmin}} (H(P, Q) - H(P)) \\ & \downarrow \text{perché } H(P) \text{ non dipende da } Q \\ & \underset{Q}{\operatorname{argmin}} (H(P, Q)) \\ & \downarrow \\ & \underset{Q}{\operatorname{argmin}} \left(-\frac{1}{T} \log Q[w_1 \dots w_T] \right) \\ & \downarrow \\ & \underset{Q}{\operatorname{argmax}} \left(+\frac{1}{T} \log Q[w_1 \dots w_T] \right) \end{aligned}$$