

N-gram language model

→ $P(w|h)$

exp: $P(\text{the | its water is so transparent that})$

$P(\text{+ its water is so transparent that the})$

$C(\text{Its water is so transparent that})$
→ 3.2

1

Avec un grand corpus on peut estimer la proba apres de l'Equation 3.2.

même si cette methode marche bien, mais le web n'est pas assez grand pour nous donner une bonne estimation.

parce que (this is because language is creative, new sentences are created all the time and we won't always be able to count entire sentences.)

Etats la plus petite modification de l'historique de

$P(w|h)$ fait que le count va valoir zero (ex: Walden Pond's water is so transparent that the " instead of "Its water is so transparent that the")

si on cherche la distribution jointe d'un groupe de mots il faut alors compter alors toute les possibilités pour la longueur de notre sequence et voir tous les cas favorables. Ce qui fait beaucoup de chose a estimer.

→ une façon plus intelligente d'estimer la proba $P(w|h)$

$$P(x_1, \dots, x_n) = \prod_{k=1}^n P(x_k | x_1, \dots, x_{k-1})$$

Changement de proba

le langage est creatif, c'est un petit changement introduit une chose qu'on ne connaît pas. Au lieu de compter, on se base sur tout l'historique avant le mot w on approxime l'historique par quelques mots

bigram model

$$IP(W_n | W_1; n-1) \approx IP(W_n | W_{n-1})$$

↗
Markov assumption

$$P(W_n | W_{1:n-1}) \approx IP(W_n | W_{n-n+1:n-1})$$

→ How to estimate bigram or ngram probabilities.

- MLE (Maximum likelihood estimation)

$$\begin{aligned} P(W_n | W_{n-1}) &= \frac{C(W_{n-1} W_n)}{\sum_w C(W_{n-1} w)} \\ &= \frac{C(W_{n-1} W_n)}{C(W_{n-1})} \end{aligned}$$

$$IP(W_n | W_{n-n+1:n-1}) = \frac{C(W_{n-n+1:n-1} W_n)}{C(W_{n-n+1:n-1})}$$

↗
relative frequency

MLE

→ end to end evaluation (performance evaluation)

Fig such as build, robustness

We need the end symbol to make a bigram a true probability distribution. Without an end symbol the sentence probabilities for all sentences of a given length would sum to one. This model would define an infinite set of probability distributions, with one distribution per sentence length. Exo 3.5

2

$P(T/M)$

↑
maximizes (parameter set to proba)

bigram captures syntactic phenomena

like
what comes after cat is noun or an adjective
or what comes after to is

→ Some may be cultural than linguistic

→ log-probabilities
we need to pad sentence

adding in log-space = multiplying in linear space

$$P_1 \times P_2 \times \dots \times P_n = \exp(\log P_1 + \dots + \log P_n)$$

- avg, perplexity, entropy

avg

3

$$\sum \frac{h_{\text{word}}(\text{word})}{N_{\text{word}}} \times \log(\text{prob})$$

Il faut que la distribution
maximiser il faut que la distribution
l'ensemble dans le train et test soient ressemblant.

les modèles statistiques sont utiles
si la base d'entraînement et de test sont
très différents des points de vue de leur distribution.

$$P_{\text{Laplace}}(\text{unigram}) = \frac{N_{\text{train}}(\text{unigram}) + k}{N_{\text{train}} + kV}$$

Laplace

$$= \underbrace{\frac{N_{\text{train}}}{N_{\text{train}} + kV}}_{\text{poids}} \underbrace{\frac{h_{\text{train}}(\text{unigram})}{N_{\text{train}}}}_{\text{unigram distribut}} + \underbrace{\frac{kV}{N_{\text{train}} + kV}}_{\text{poids}} \times \underbrace{\frac{1}{V}}_{\text{uniform probability}}$$

une bonne interprétation de modèle
on a combiné deux modèles avec un certain ensemble
de poids pour trouver notre modèle
final.

$K=0$ un smooth model

$K \nearrow$ diminue la proba des Unigram les plus communs jusqu'au niveau où on obtient un model uniform.

Effet de l'Interpolat

Plus de test-thain plus on a besoin d'interpoler (combinaison avec unigram et uniform distribution)

$$\sum \frac{n_{\text{eval}} \times \log(IP)}{n}$$

$$n \log \frac{IP(U_n)}{IP(U_n)} + n \log IP(U_n)$$

N gram length

| | 0 | 1 | 2 | 3 | 4 | 5 |
|-------|---------------------|--|---|---|---|---|
| I | \downarrow | $\frac{n(\text{BJI})}{n(\text{I})}$ | | | | |
| have | \downarrow | $\frac{n(\text{I have})}{n(\text{I})}$ | | | | |
| a | $\frac{1}{\sqrt{}}$ | $\frac{n(\text{the king})}{n(a)}$ | | | | |
| them | \downarrow | $\frac{n(\text{the king})}{n(a)}$ | | | | |
| [END] | \downarrow | $\frac{n(\text{the king})}{n(a)}$ | | | | |
| We | \downarrow | $\frac{n(\text{the king})}{n(a)}$ | | | | |

$$\text{ngram_encl} = \text{token_position} + 1$$

$\frac{0+1}{1+1}$

Je suis ton par

'Je suis' 'km' 'par'

Je point 0

ngn 1 P(50 | 55)

2
 3
 4
 5

Washburn

Not a student

Run, Run, Run
and don't stop

(1) अथर्ववेद

1881. The morning of the 1st of July. 1st of July.