



BM1: Advanced Natural Language Processing

University of Potsdam

Tatjana Scheffler

tatjana.scheffler@uni-potsdam.de

October 28, 2016



Today

- n-grams
- Zipf's law
- Ianguage models



Maximum Likelihood Estimation

- We want to estimate the parameters of our model from frequency observations. There are many ways to do this. For now, we focus on maximum likelihood estimation, MLE.
- Likelihood L(O ; p) is the probability of our model generating the observations O, given parameter values p.
- Goal: Find value for parameters that maximizes the likelihood.



Bernoulli model

Let's say we had training data C of size N, and we had N_H observations of H and N_T observations of T.

likelihood
$$L(C) = \prod_{i=1}^{N} P(w_i \mid p) = \prod_{i=1}^{N} p^{N_H} (1-p)^{N_T}$$

log-likelihood
 $\ell(C) = \log L(C) = \sum_{i=1}^{N} \log P(w_i \mid p) = N_H \log p + N_T \log(1-p)$



Likelihood functions





Logarithm is monotonic



р

Observation: If $x_1 > x_2$, then $\ln(x_1) > \ln(x_2)$.

Therefore, argmax L(C) = argmax I(C)

р



Maximizing the log-likelihood

□ Find maximum of function by setting derivative to zero:

$$\ell(C) = N_H \log p + N_T \log(1-p)$$

 $rac{d\ell(C)}{dp} = rac{N_H}{p} - rac{N_T}{1-p}$

Solution is $p = N_H / N = f(H)$.



Language Modelling



Let's play a game

- I will write a sentence on the board.
- Each of you, in turn, gives me a word to continue that sentence, and I will write it down.



Let's play another game

- You write a word on a piece of paper.
- You get to see the piece of paper of your neighbor, but none of the earlier words.
- In the end, I will read the sentence you wrote.



Statistical models for NLP

Generative statistical model of language:

prob. dist. P(w) over NL expressions that we can observe.

- w may be complete sentences or smaller units
- will later extend this to pd P(w, t) with hidden random variables t
- Assumption: A corpus of observed sentences w is generated by repeatedly sampling from P(w).
- We try to estimate the parameters of the prob dist from the corpus, so we can make predictions about unseen data.



Example





- A language model LM is a probability distribution P(w) over words.
- Think of it as a random process that generates sentences word by word:

$$X_1 \qquad X_2 \qquad X_3 \qquad X_4 \qquad \dots$$



- A language model LM is a probability distribution P(w) over words.
- Think of it as a random process that generates sentences word by word:





- A language model LM is a probability distribution P(w) over words.
- Think of it as a random process that generates sentences word by word:



- A language model LM is a probability distribution P(w) over words.
- Think of it as a random process that generates sentences word by word:





- A language model LM is a probability distribution P(w) over words.
- Think of it as a random process that generates sentences word by word:





Our game as a process

Each of you = a random variable X_t ;

event " $X_t = w_t$ " means word at position t is w_t .

- □ When you chose w_t , you could see the outcomes of the previous variables: $X_1 = w_1$, ..., $X_{t-1} = w_{t-1}$.
- Thus, each X_t followed a pd

$$P(X_{t} = w_{t} | X_{1} = w_{1}, \dots, X_{t-1} = w_{t-1})$$



Our game as a process

Assume that X_t follows some given pd

$$P(X_{t} = w_{t} | X_{1} = w_{1}, \dots, X_{t-1} = w_{t-1})$$

Then probability of the entire sentence (or corpus) w = w₁ ... w_n is

$$P(w_1 \dots w_n) = P(w_1)P(w_2 | w_1)P(w_3 | w_1, w_2) \dots P(w_n | w_1, \dots, w_{n-1})$$



Parameters of the model

- Our model has one parameter for
 - $P(X_{t} = w_{t} | w_{1}, ..., w_{t-1})$ for all t and $w_{1}, ..., w_{t}$.
- Can use maximum likelihood estimation:

$$P(w_t \mid w_1, \dots, w_{t-1}) = rac{C(w_1 \dots w_{t-1} w_t)}{C(w_1 \dots w_{t-1})}$$

- Let's say a natural language has 10⁵ different words. How many tuples w₁, ... w_t of length t?
 t = 1: 10⁵
 t = 2: 10¹⁰ different contexts
 - □ t = 3: 10¹⁵; etc.



Sparse data problem

typical corpus size:

- Brown corpus: 10⁶ tokens
- Gigaword corpus: 10⁹ tokens
- Problem exacerbated by Zipf 's Law:
 - Order all words by their absolute frequency in corpus (rank 1 = most frequent word).
 - Then rank is inversely proportional to absolute frequency; i.e., most words are really rare.
 - Zipf's Law is very robust across languages and corpora.



Interlude: Corpora



Terminology

- N = corpus size; number of (word) tokens
- V = vocabulary; number of (word) types
- hapax legomenon = a word that appears exactly once in the corpus



An example corpus

Es war einmal ein Müller, der hatte drei Söhne, seine Mühle, einen Esel und einen Kater; die Söhne mußten mahlen, der Esel Getreide holen und Mehl forttragen, die Katze dagegen die Mäuse wegfangen. Als der Müller starb, teilten sich die drei Söhne in die Erbschaft: der älteste bekam die Mühle, der zweite den Esel, der dritte den Kater; weiter blieb nichts für ihn übrig. Da war er traurig und sprach zu sich selbst:

Tokens: 86

Types: 53



Frequency list

Es	1	die	5
war	2	mußten	1
einmal	1	mahlen	1
ein	1	Getreide	1
Müller	2	holen	1
,	8	Mehl	1
der	5	forttragen	1
hatte	1	Katze	1
drei	2	dagegen	1
Söhne	3	Mäuse	1
seine	1	wegfangen	1
Mühle	2		3
einen	2	Als	1
Esel	3	starb	1
und	3	teilten	1
Kater	2	sich	
:	2		



Frequency list

Туре	abs.	relativ	Туре	abs.	relativ	Туре	abs.	relativ
,	8	0,083	Als	1	0,0104	in	1	0,01
der	5	0,052	älteste	1	0,0104	Katze	1	0,01
die	5	0,052	bekam	1	0,0104	Mäuse	1	0,01
• 4	3	0,031	blieb	1	0,0104	Mehl	1	0,01
Esel	3	0,031	Da	1	0,0104	mußten	1	0,01
Söhne	3	0,031	dagegen	1	0,0104	mahlen	1	0,01
und	3	0,031	dritte	1	0,0104	nichts	1	0,01
:	2	0,021	ein	1	0,0104	seine	1	0,01
;	2	0,021	einmal	1	0,0104	selbst	1	0,01
den	2	0,021	er	1	0,0104	sprach	1	0,01
drei	2	0,021	Erbschaft	1	0,0104	starb	1	0,01
einen	2	0,021	Es	1	0,0104	teilten	1	0,01
Kater	2	0,021	forttragen	1	0,0104	traurig	1	0,01
Mühle	2	0,021	für	1	0,0104	übrig	1	0,01
Müller	2	0,021	Getreide	1	0,0104	wegfangen	1	0,01
sich	2	0,021	hatte	1	0,0104	weiter	1	0,01
war	2	0,021	holen	1	0,0104	zu	1	0,01
			ihn	1	0,0104	zweite	1	0,01



Frequency profile

Туре	Rang	abs.	relativ	Туре	Rang	abs.	relativ	Туре	Rang	abs.	relativ
,	1	8	0,083	Als	53	1	0,0104	in	53	1	0,01
der	3	5	0,052	älteste	53	1	0,0104	Katze	53	1	0,01
die	3	5	0,052	bekam	53	1	0,0104	Mäuse	53	1	0,01
	7	3	0,031	blieb	53	1	0,0104	Mehl	53	1	0,01
Esel	7	3	0,031	Da	53	1	0,0104	mußten	53	1	0,01
Söhne	7	3	0,031	dagegen	53	1	0,0104	mahlen	53	1	0,01
und	7	3	0,031	dritte	53	1	0,0104	nichts	53	1	0,01
:	17	2	0,021	ein	53	1	0,0104	seine	53	1	0,01
;	17	2	0,021	einmal	53	1	0,0104	selbst	53	1	0,01
den	17	2	0,021	er	53	1	0,0104	sprach	53	1	0,01
drei	17	2	0,021	Erbschaft	53	1	0,0104	starb	53	1	0,01
einen	17	2	0,021	Es	53	1	0,0104	teilten	53	1	0,01
Kater	17	2	0,021	forttragen	53	1	0,0104	traurig	53	1	0,01
Mühle	17	2	0,021	für	53	1	0,0104	übrig	53	1	0,01
Müller	17	2	0,021	Getreide	53	1	0,0104	wegfangen	53	1	0,01
sich	17	2	0,021	hatte	53	1	0,0104	weiter	53	1	0,01
war	17	2	0,021	holen	53	1	0,0104	zu	53	1	0,01
				ihn	53	1	0,0104	zweite	53	1	0,01



Plotting corpus frequencies

Number of types	rank	frequency
1	1	8
2	3	5
4	7	3
10	17	2
36	53	1

How many different words in the corpus are there with each frequency?



Plotting corpus frequencies

x-axis: rank

□ y-axis: frequency





Some other corpora





Zipf's Law

Zipf's Law characterizes the relation between frequent and rare words:

f(w) = C / r(w)or equivalently: f(w) * r(w) = C

- Frequency of lexical items (words types) in a large corpus is inversely proportional to their rank.
- Empirical observation in many different corpora
- Brown corpus:
 - half of all types are hapax legomena



Effects of Zipf's Law

Lexicography:

- □ Sinclair (2005): need at least 20 instances
- BNC (10⁸ Tokens): <14% of words appear 20 times or more

Speech synthesis:

- may accept bad output for rare words
- but most words are rare! (at least 1 per sentence)
- Vocabulary growth:
 - vocabulary growth of corpora is not constant
 - G = #hapaxes / #tokens



Back to Language Models



Independence assumptions

- Let's pretend that word at position t depends only on the words at positions t-1, t-2, ..., t-k for some fixed k (Markov assumption of degree k).
- Then we get an n-gram model, with n = k+1:
 - $P(X_{t} | X_{1},...,X_{t-1}) = P(X_{t} | X_{t-k},...,X_{t-1})$ for all t.
- Special names for unigram models (n = 1), bigram models (n = 2), trigram models (n = 3).



Independence assumption

- We assume independence of X_t from events that are too far in the past, although we know that this assumption is incorrect.
- Typical tradeoff in statistical NLP:
 - if model is too shallow, it won't represent important linguistic dependencies
 - if model is too complex, its parameters can't be estimated accurately from the available data





Tradeoff in practice

In person	she		was		inferior		to		both		sisters	
2-gram	$P(\cdot pe$	erson)	$P(\cdot sh$	ie)	$P(\cdot was)$		P(· i≀	nferior)	$P(\cdot to)$		$P(\cdot both)$	
$ \begin{array}{c} 1 \\ 2 \\ 3 \\ 4 \\ \\ 23 \\ \\ 41 \\ \\ 293 \\ \end{array} $	and who to in she	0.099 0.099 0.076 0.045 0.009	had was	0.141 0.122	not a the to	0.065 0.052 0.033 0.031	to	0.212	be the her have Mrs what	0.111 0.057 0.048 0.027 0.006 0.004 0.004	of to in and she sisters	0.066 0.041 0.038 0.025 0.009 0.006
					inferior	0						



Tradeoff in practice

In person	she	was	inferior	to	both	sisters
3-gram 1 2 3 4 ∞	P(• In,person) Unseen	<i>P</i> (· <i>person,she</i>) did 0.5 was 0.5	P(· she,was) not 0.057 very 0.038 in 0.030 to 0.026	P(· was,inf.) Unseen	P(· inferior,to) the 0.286 Maria 0.143 cherries 0.143 her 0.143 both 0	P(· to,both) to 0.222 Chapter 0.111 Hour 0.111 Twice 0.111 sisters 0



Tradeoff in practice

In person	she	was	inferior	to	both	sisters
4-gram 1 ∞	P(· <i>u,I,p</i>) Unseen	P(+ I,p,s) Unseen	<i>P</i> (· <i>p,s,w</i>) in 1.0 inferior 0	P(· s,w,i) Unseen	P(· <i>w,i,t</i>) Unseen	P(+ <i>i,t,b</i>) Unseen



Conclusion

- Statistical models of natural language
- Language models using n-grams
- Data sparseness is a problem.



next Tuesday

smoothing language models