

# 3: Statistical Properties of Language

## Machine Learning and Real-world Data

Simone Teufel and Ann Copestake

Computer Laboratory  
University of Cambridge

Lent 2017

# Last session: Naive Bayes Classifier

- You built a smoothed and an unsmoothed NB classifier.
- You evaluated them in terms of accuracy.
- The unsmoothed classifier mostly produced equal probabilities = 0.
- In the smoothed version, this problem has been alleviated.
- Why are there so many zero frequencies, and why does smoothing work?

# Statistical Properties of Language I

- How many frequent vs. infrequent terms should we expect in a collection?
- **Zipf's law** states that there is a direct inverse relationship between a word's frequency rank and the absolute value of that frequency.
- This is an instance of a Power Law.
- The law is astonishingly simple ...
- ... given how complex the sentence-internal constraints between some of the words concerned are.

# Statistical Properties of Language II

- **Heaps' law** concerns the relationship between all items of a language and unique items of a language.
- There is an exponential relationship between the two.
- This is also surprising because one might expect saturation.
- Surely at some point all words of a language have been “used up”?

# Frequencies of words

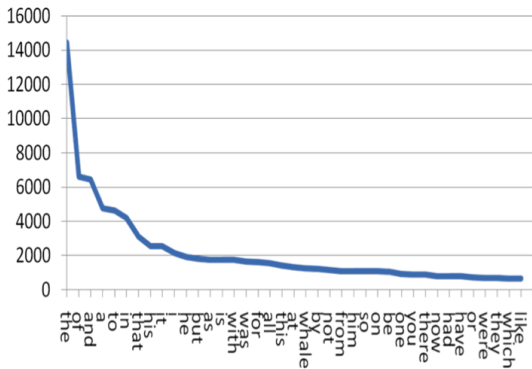
## Zipf's law:

There is a direct inverse relationship between a word's frequency rank and the absolute value of that frequency.

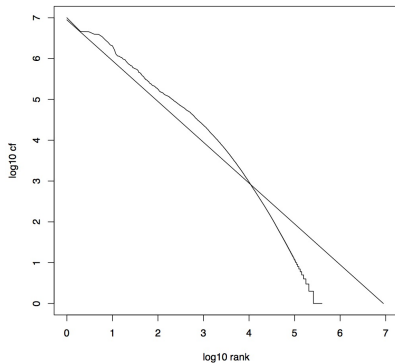
$$f_w \approx \frac{k}{r_w^\alpha}$$

- $f_w$ : frequency of word  $w$
- $r_w$ : frequency rank of word  $w$
- $\alpha, k$ : constants (language-dependent)
  - $\alpha$  around 1 for English, 1.3 for German
- Zipf's Law means that in language, there are **a few very frequent** terms and **very many very rare** terms.

# Zipf's Law



# Zipf's Law in log-log space



(Reuters dataset)

# Zipf's Law: Examples from 5 Languages

Top 10 most frequent words in some large language samples:



# Zipf's Law: Examples from 5 Languages

Top 10 most frequent words in some large language samples:

## English

1	the	61,847
2	of	29,391
3	and	26,817
4	a	21,626
5	in	18,214
6	to	16,284
7	it	10,875
8	is	9,982
9	to	9,343
10	was	9,236

BNC,  
100Mw

# Zipf's Law: Examples from 5 Languages

Top 10 most frequent words in some large language samples:

## English

1	the	61,847
2	of	29,391
3	and	26,817
4	a	21,626
5	in	18,214
6	to	16,284
7	it	10,875
8	is	9,982
9	to	9,343
10	was	9,236

## German

1	der	7,377,879
2	die	7,036,092
3	und	4,813,169
4	in	3,768,565
5	den	2,717,150
6	von	2,250,642
7	zu	1,992,268
8	das	1,983,589
9	mit	1,878,243
10	sich	1,680,106

BNC,  
100Mw

“Deutscher  
Wortschatz”,  
500Mw

# Zipf's Law: Examples from 5 Languages

Top 10 most frequent words in some large language samples:

English		German		Spanish				
1	the	61,847	1	der	7,377,879	1	que	32,894
2	of	29,391	2	die	7,036,092	2	de	32,116
3	and	26,817	3	und	4,813,169	3	no	29,897
4	a	21,626	4	in	3,768,565	4	a	22,313
5	in	18,214	5	den	2,717,150	5	la	21,127
6	to	16,284	6	von	2,250,642	6	el	18,112
7	it	10,875	7	zu	1,992,268	7	es	16,620
8	is	9,982	8	das	1,983,589	8	y	15,743
9	to	9,343	9	mit	1,878,243	9	en	15,303
10	was	9,236	10	sich	1,680,106	10	lo	14,010

BNC,  
100Mw

“Deutscher  
Wortschatz”,  
500Mw

subtitles,  
27.4Mw

# Zipf's Law: Examples from 5 Languages

Top 10 most frequent words in some large language samples:

English	German	Spanish	Italian
1 the 61,847	1 der 7,377,879	1 que 32,894	1 non 25,757
2 of 29,391	2 die 7,036,092	2 de 32,116	2 di 22,868
3 and 26,817	3 und 4,813,169	3 no 29,897	3 che 22,738
4 a 21,626	4 in 3,768,565	4 a 22,313	4 è 18,624
5 in 18,214	5 den 2,717,150	5 la 21,127	5 e 17,600
6 to 16,284	6 von 2,250,642	6 el 18,112	6 la 16,404
7 it 10,875	7 zu 1,992,268	7 es 16,620	7 il 14,765
8 is 9,982	8 das 1,983,589	8 y 15,743	8 un 14,460
9 to 9,343	9 mit 1,878,243	9 en 15,303	9 a 13,915
10 was 9,236	10 sich 1,680,106	10 lo 14,010	10 per 10,501
BNC, 100Mw	“Deutscher Wortschatz”, 500Mw	subtitles, 27.4Mw	subtitles, 5.6Mw

# Zipf's Law: Examples from 5 Languages

Top 10 most frequent words in some large language samples:

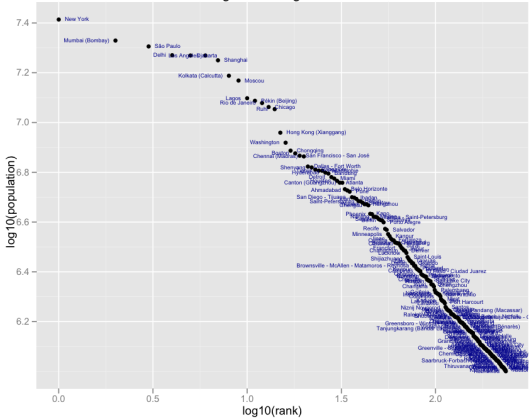
English	German	Spanish	Italian	Dutch
1 the 61,847	1 der 7,377,879	1 que 32,894	1 non 25,757	1 de 4,770
2 of 29,391	2 die 7,036,092	2 de 32,116	2 di 22,868	2 en 2,709
3 and 26,817	3 und 4,813,169	3 no 29,897	3 che 22,738	3 het/'t 2,469
4 a 21,626	4 in 3,768,565	4 a 22,313	4 è 18,624	4 van 2,259
5 in 18,214	5 den 2,717,150	5 la 21,127	5 e 17,600	5 ik 1,999
6 to 16,284	6 von 2,250,642	6 el 18,112	6 la 16,404	6 te 1,935
7 it 10,875	7 zu 1,992,268	7 es 16,620	7 il 14,765	7 dat 1,875
8 is 9,982	8 das 1,983,589	8 y 15,743	8 un 14,460	8 die 1,807
9 to 9,343	9 mit 1,878,243	9 en 15,303	9 a 13,915	9 in 1,639
10 was 9,236	10 sich 1,680,106	10 lo 14,010	10 per 10,501	10 een 1,637
BNC, 100Mw	"Deutscher Wortschatz", 500Mw	subtitles, 27.4Mw	subtitles, 5.6Mw	subtitles, 800Kw

# Other collections (allegedly) obeying power laws

- Sizes of settlements
- Frequency of access to web pages
- Income distributions amongst top earning 3% individuals
- Korean family names
- Size of earth quakes
- Word senses per word
- Notes in musical performances
- ...

# World city populations

world city populations for 8 countries  
log-size vs log-rank



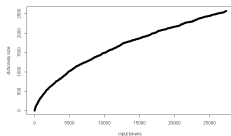
# Vocabulary size

## Heaps' Law:

The following relationship exists between the size of a vocabulary and the size of text that gave rise to it:

$$u_n = kn^\beta$$

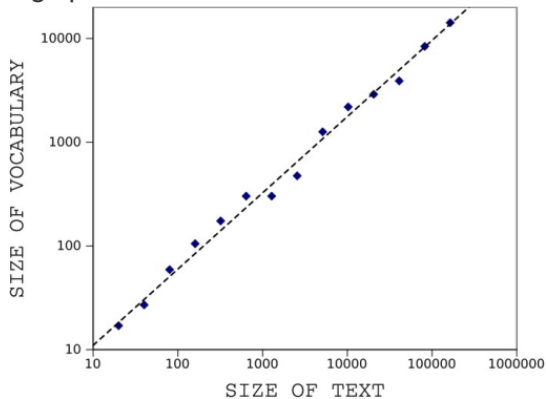
- $u_n$ : number of types (unique items); vocabulary size
- $n$ : number of tokens; text size
- $\beta, k$ : constants (language-dependent)
  - $\beta$  normally around  $\frac{1}{2}$
  - $30 \leq k \leq 100$





# Heaps' Law

- In log-log space:

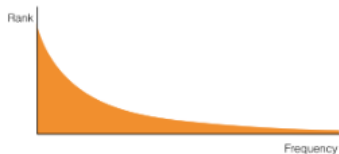


- Reasons for infinite vocabulary growth?

# Consequences for our experiment

- Zipf's law and Heaps' law taken together explain why smoothing is necessary and effective:
  - MLE overestimates the likelihood for seen words.
  - Smoothing redistributes some of this probability mass.

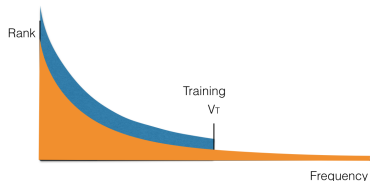
# The real situation



- Most of the probability mass is in the **long tail**.

# The situation according to MLE

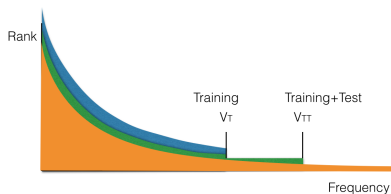
$$\hat{P}_{MLE}(w_i|c) = \frac{\text{count}(w_i, c)}{\sum_{w \in V_T} \text{count}(w, c)}$$



- With MLE, only seen words can get a frequency estimate.
- Probability mass is still 1.
- Therefore, the probability of seen words is a (big) overestimate.

# What smoothing does

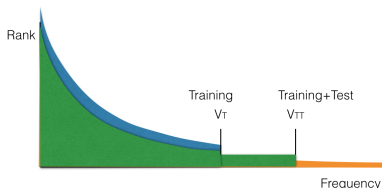
$$\hat{P}_S(w_i|c) = \frac{\text{count}(w_i, c) + 1}{(\sum_{w \in V_{TT}} \text{count}(w, c)) + |V_{TT}|}$$



- Smoothing redistributes the probability mass towards the real situation.
- It takes some portion away from the MLE overestimate for seen words.
- It redistributes this portion to a certain, finite number of unseen words (in our case, as a uniform distribution).

# What smoothing does

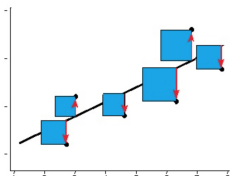
$$\hat{P}_S(w_i|c) = \frac{\text{count}(w_i, c) + 1}{(\sum_{w \in V_{TT}} \text{count}(w, c)) + |V_{TT}|}$$



- Smoothing takes some portion away from the MLE overestimate for seen words.
- It redistributes this portion to a certain, finite number of unseen words (in our case, as a uniform distribution).
- As a result, the real situation is approximated more closely.

# Your first task today

- Plot frequency vs frequency rank for larger dataset (i.e., visually verify Zipf's Law)
- Estimate parameters  $k, \alpha$  for Zipf's Law
- Use least-squares algorithm for doing so.
- Is  $\alpha$  really 1 for English?
  - There is much scientific discussion of this question.



# Your second task today

- Plot type/token ratio for IMDB dataset (verify Heaps' Law)



# Ticking today

- Task 2 – NB Classifier

- Introduction to Information Retrieval, Christopher C. Manning, Prabhakar Raghavan, Hinrich Schütze, Cambridge University Press, 2008. Section 5.1, pages 79-82.
- (Please note that  $\alpha = 1$  is assumed in the Zipf Law formula on page 82.)