

Information Retrieval

Statistics of Text

James Allan

University of Massachusetts, Amherst

CMPSCI 646 (NTU ST770-A)

Fall 2002

All slides copyright © Bruce Croft and/or James Allan

Outline

- Zipf distribution
- Vocabulary growth
- Communication Theory
- Collocation as a basis for lexicography
- Markov models for part-of-speech tagging
- Language models for speech recognition

Zipf's Law

- A few words occur very often
 - 2 most frequent words can account for 10% of occurrences
 - top 6 words are 20%, top 50 words are 50%
- Many words are infrequent
- “Principle of Least Effort”
 - easier to repeat words rather than coining new ones
- Rank · Frequency \approx Constant
 - $p_r = (\text{Number of occurrences of word of rank } r)/N$
 - N total word occurrences
 - probability that a word chosen randomly from the text will be the word of rank r
 - for D unique words $\sum p_r = 1$
 - $r \cdot p_r = A$
 - $A \approx 0.1$

George Kingsley Zipf, 1902-1950
Linguistic professor at Harvard

Example of Frequent Words

	Frequent Word	Number of Occurrences	Percentage of Total
Artifact of InQuery's stemming technique	the	7,398,934	5.9
	of	3,893,790	3.1
	to	3,364,653	2.7
	and	3,320,687	2.6
	in	2,311,785	1.8
	is	1,559,147	1.2
	for	1,313,561	1.0
	The	1,144,860	0.9
	that	1,066,503	0.8
	said	1,027,713	0.8

Frequencies from 336,310 documents in the 1GB TREC Volume 3 Corpus
125,720,891 total word occurrences; 508,209 unique words

Zipf's Law and H.P.Luhn

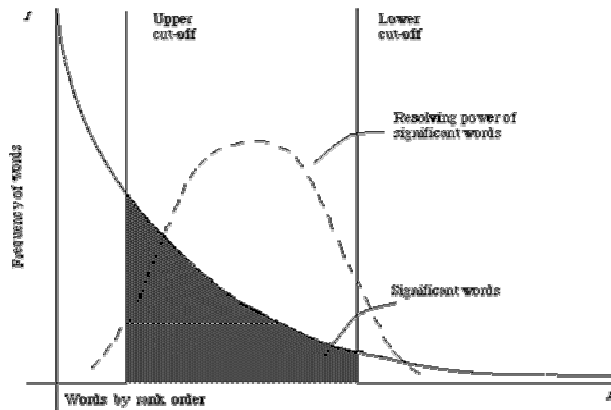


Figure 2.1. A plot of the Zipf-like curve relating f , the frequency of occurrence and r , the rank order (adapted from Schmidt, page 120)

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Examples of Zipf

Word	Freq	r	Pr	r*Pr
the	15659	1	6.422	0.0642
of	7179	2	2.944	0.0589
to	6287	3	2.578	0.0774
a	5830	4	2.391	0.0956
and	5580	5	2.288	0.1144
in	5245	6	2.151	0.1291
that	2494	7	1.023	0.0716
for	2197	8	0.901	0.0721
was	2147	9	0.881	0.0792
with	1824	10	0.748	0.0748
his	1813	11	0.744	0.0818
is	1800	12	0.738	0.0886
he	1687	13	0.692	0.0899
as	1576	14	0.646	0.0905
on	1523	15	0.625	0.0937
by	1443	16	0.592	0.0947
at	1318	17	0.541	0.0919
it	1232	18	0.505	0.0909
from	1217	19	0.499	0.0948
but	1136	20	0.466	0.0932
u	949	21	0.389	0.0817
had	937	22	0.384	0.0845
last	909	23	0.373	0.0857
be	906	24	0.372	0.0892
who	883	25	0.362	0.0905

Word	Freq	r	Pr	r*Pr
has	880	26	0.361	0.0938
not	875	27	0.359	0.0969
an	863	28	0.354	0.0991
s	862	29	0.354	0.1025
have	860	30	0.353	0.1058
were	858	31	0.352	0.1091
their	812	32	0.333	0.1066
are	807	33	0.331	0.1092
one	742	34	0.304	0.1035
they	679	35	0.278	0.0975
its	668	36	0.274	0.0986
all	646	37	0.265	0.098
week	626	38	0.257	0.0976
government	582	39	0.239	0.0931
when	577	40	0.237	0.0947
would	572	41	0.235	0.0962
been	554	42	0.227	0.0954
out	553	43	0.227	0.0975
new	544	44	0.223	0.0982
which	539	45	0.221	0.0995
up	539	45	0.221	0.0995
more	535	47	0.219	0.1031
into	516	48	0.212	0.1016
only	504	49	0.207	0.1013
will	488	50	0.2	0.1001

Top 50 words from 423 short TIME magazine articles
(243,836 word occurrences, lowercased, punctuation removed, 1.6 MB)

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Examples of Zipf

Word	Freq	r	Pr(%)	r*Pr
the	2,420,778	1	6.488	0.0649
of	1,045,733	2	2.803	0.0561
to	968,882	3	2.597	0.0779
a	892,429	4	2.392	0.0957
and	865,644	5	2.32	0.116
in	847,825	6	2.272	0.1363
said	504,593	7	1.352	0.0947
for	363,865	8	0.975	0.078
that	347,072	9	0.93	0.0837
was	293,027	10	0.785	0.0785
on	291,947	11	0.783	0.0861
he	250,919	12	0.673	0.0807
is	245,843	13	0.659	0.0857
with	223,846	14	0.6	0.084
at	210,064	15	0.563	0.0845
by	209,586	16	0.562	0.0899
it	195,621	17	0.524	0.0891
from	189,451	18	0.508	0.0914
as	181,714	19	0.487	0.0925
be	157,300	20	0.422	0.0843
were	153,913	21	0.413	0.0866
an	152,576	22	0.409	0.09
have	149,749	23	0.401	0.0923
his	142,285	24	0.381	0.0915
but	140,880	25	0.378	0.0944

Word	Freq	r	Pr(%)	r*Pr
has	136,007	26	0.365	0.0948
are	130,322	27	0.349	0.0943
not	127,493	28	0.342	0.0957
who	116,364	29	0.312	0.0904
they	111,024	30	0.298	0.0893
its	111,021	31	0.298	0.0922
had	103,943	32	0.279	0.0892
will	102,949	33	0.276	0.0911
would	99,503	34	0.267	0.0907
about	92,983	35	0.249	0.0872
i	92,005	36	0.247	0.0888
been	88,786	37	0.238	0.0881
this	87,286	38	0.234	0.0889
their	84,638	39	0.227	0.0885
new	83,449	40	0.224	0.0895
or	81,796	41	0.219	0.0899
which	80,385	42	0.215	0.0905
we	80,245	43	0.215	0.0925
more	76,388	44	0.205	0.0901
after	75,165	45	0.201	0.0907
us	72,045	46	0.193	0.0888
percent	71,956	47	0.193	0.0906
up	71,082	48	0.191	0.0915
one	70,266	49	0.188	0.0923
people	68,988	50	0.185	0.0925

Top 50 words from 84,678 Associated Press 1989 articles
(37,309,114 word occurrences, lowercased, punctuation removed, 266MB)

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Predicting Occurrence Frequencies

- A word that occurs n times has rank $r_n = AN/n$
- Several words may occur n times
- Assume rank given by r_n applies to *last* of the words that occur n times
- r_n words occur n times or more
- r_{n+1} words occur $n+1$ times or more
 - Note: $r_n < r_{n+1}$ since words that occur frequently are at the start of list (lower rank)
- The number of words that occur exactly n times is

$$I_n = r_n - r_{n+1} = AN/n - AN/(n+1) = AN / (n(n+1))$$
- Highest ranking term occurs once and has rank $D =$

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Example of Occurrence Frequencies

Number of Occurrences (n)	Predicted Proportion of Occurrences $1/n(n+1)$	Actual Proportion occurring n times I_n/D	Actual Number of Words occurring n times
1	.500	.402	204,357
2	.167	.132	67,082
3	.083	.069	35,083
4	.050	.046	23,271
5	.033	.032	16,332
6	.024	.024	12,421
7	.018	.019	9,766
8	.014	.016	8,200
9	.011	.014	6,907
10	.009	.012	5,893

Frequencies from 336,310 documents in the 1GB TREC Volume 3 Corpus
125,720,891 total word occurrences; 508,209 unique words

Does Real Data Fit Zipf's Law?

[From R.Mooney, UT.Austin]

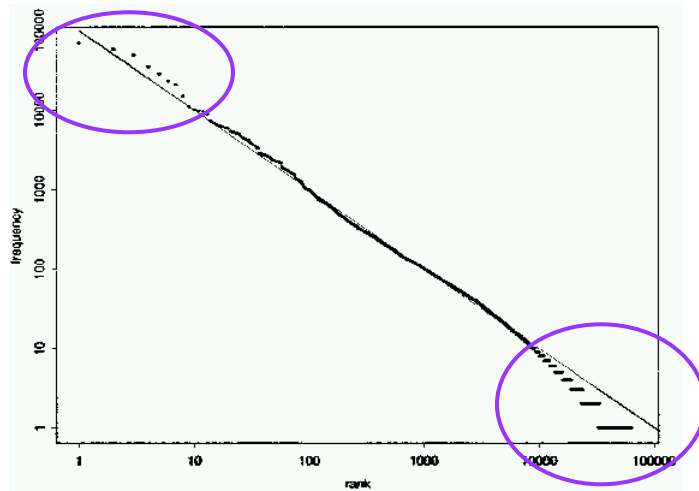
- A law of the form $y = kx^c$ is called a power law.
- Zipf's law is a power law with $c = -1$
 - $r = (AN) \cdot n^{-1}$
 - AN is a constant for a fixed collection
- On a log-log plot, power laws give a straight line with slope c .

$$\log(y) = \log(kx^c) = \log k + c \log(x)$$

- Zipf is quite accurate except for very high and low rank.

Fit to Zipf for Brown Corpus

[From R.Mooney, UT.Austin]



$$k = 100,000$$

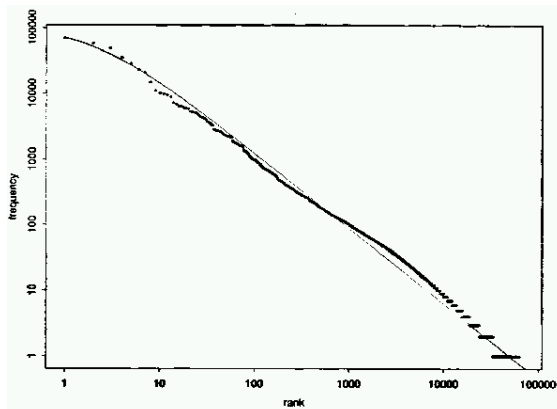
CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Mandelbrot (1954) Correction

[From R.Mooney, UT.Austin]

- The following more general form gives bit better fit
 - Adds a constant to the denominator
 - $y = k(x+t)^c$
- Here,
 $r = (AN) \cdot (n+t)^{-1}$



Mandelbrot's function on Brown corpus

$$k = 10^{5.4}, C = -1.15, t = 100$$

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Explanations for Zipf's Law

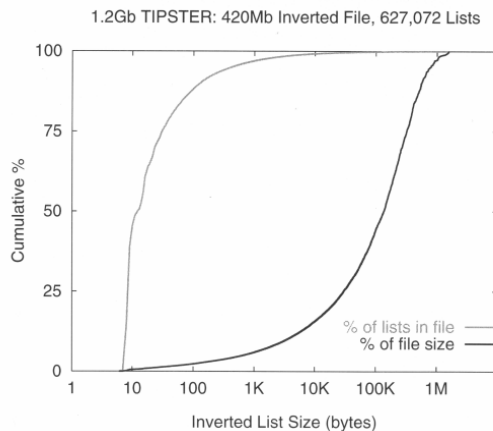
[From R.Mooney, UT.Austin]

- Zipf's explanation was his "principle of least effort." Balance between speaker's desire for a small vocabulary and hearer's desire for a large one.
- Debate (1955-61) between Mandelbrot and H. Simon over explanation.
- Li (1992) shows that just random typing of letters including a space will generate "words" with a Zipfian distribution.
 - <http://linkage.rockefeller.edu/wli/zipf/>
 - Short words more likely to be generated

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

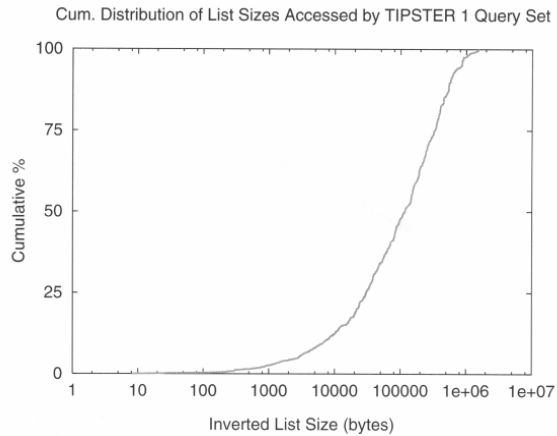
Size Distribution of Term Lists



CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Characteristics of Query Terms



CMPSCI 646

Copyright © Bruce Croft and/or James Allan

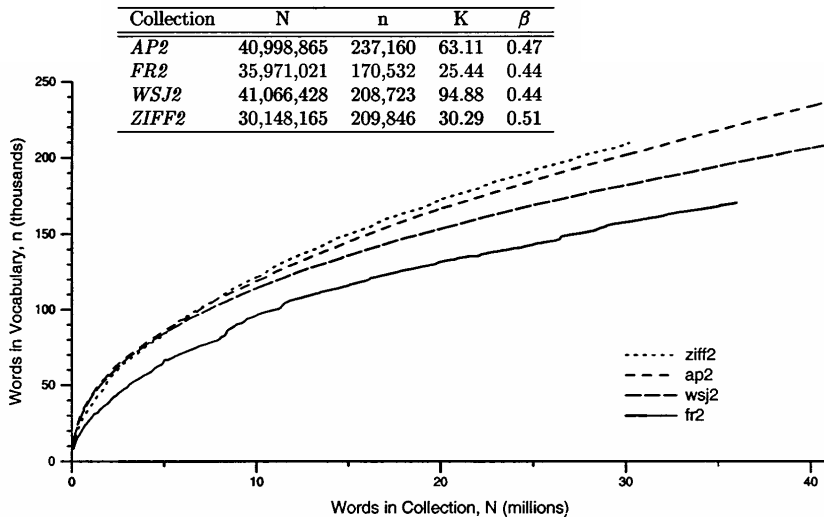
Vocabulary Growth

- How does the size of the overall vocabulary (number of unique words) grow with the size of the corpus?
 - Vocabulary has no upper bound due to proper names, typos, etc.
 - New words occur less frequently as vocabulary grows
- If V is the size of the vocabulary and the n is the length of the corpus in words:
 - $V = Kn^{\beta}$ ($0 < \beta < 1$)
- Typical constants:
 - $K \approx 10-100$
 - $\beta \approx 0.4-0.6$ (approx. square-root)
- Can be derived from Zipf's law by assuming documents are generated by randomly sampling words from a Zipfian distribution.

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Heaps' Law Data



CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Information Theory

- Shannon studied theoretical limits for data compression and transmission rate
- Compression limits given by Entropy (H)
- Transmission limits given by Channel Capacity (C)
- A number of language tasks have been formulated as a “noisy channel” problem
 - i.e., determine the most likely input given the noisy output
 - OCR
 - Speech recognition

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Shannon Game

- The President of the United States is George W. ...
- The winner of the \$10K prize is ...
- Mary had a little ...
- The horse raced past the barn ...
 - Period (end of sentence)
 - “whinnied” (garden path sentence)

Information Theory

- *Information content* of a message is dependent on the receiver’s prior knowledge as well as on the message itself
- How much of the receiver’s uncertainty (entropy) is reduced
- How predictable is the message

Information Theory

- Information content H is defined as a decreasing function $H(p)$ of the *a priori* probability p with which the message could be predicted
 - if receiver predicts message with probability 1, information content is zero
 - $H(1) = 0$
 - if prediction of message is probability 0, message would have infinite information content
 - $H(0)$ undefined
 - information content should be additive
 - $H(p_1 p_2) = H(p_1) + H(p_2)$
- $H(p) = -\log p$
- With logs base 2, unit of information content

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Information Theory

- Given n messages, the average or expected information content to be gained through receipt of one of the n possible messages is
$$\bar{H} = - \sum_{r=1}^n p_r \log p_r$$
- Average entropy is a maximum when messages are equally probable
 - e.g., average entropy associated with characters assuming equal probabilities
 - $\log 1/26 = 4.7$ bits
- Taking actual probabilities into account, entropy is 4.14 bits
- With bigram probabilities, reduces entropy to 3.56 bits
- Experiments with people give values around 1.3 bits
- Better models reduce the relative entropy or “perplexity”

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Information Theory

- For words $\bar{H} = - \sum_{r=1}^D (A/r) \log (A/r)$
- Approximations give $-(A \log_2 A) \log_e (2D+1)$
- For $D = 10,000$ $H = 9.5$
 50,000 10.9
 100,000 11.4 bits
- Equi-probable case gives
 $H = 13.3, 15.6$ and 16.6 bits

D is number of
unique words

Information Theory

- Consider word-probability distribution p_r which produces the smallest mean number of letters per word for a particular value of entropy H
- That is, minimize $\sum p_r m_r$ where
 - m_r is the length of the word with rank r
 - $\sum p_r = 1$ and
 - $H = - \sum p_r \log p_r = \text{constant}$
- Gives $p_r = A/(r + B)^\beta$ where A, B and β are fixed for a given subject vocabulary
 - Look familiar?
 - Mandelbrot's derivation
- Information theory has been used for compression, term weighting, and evaluation measures

Mutual Information

- Mutual information is a symmetric, non-negative measure of the common information in two random variables
- $I(X; Y) = H(X) - H(X|Y) = H(X) + H(Y) - H(X, Y)$

$$I(X; Y) = \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

- $I(X; Y) = D(p(x, y) || p(x)p(y))$ which is the relative entropy or Kullback-Leibler 'distance'

$$D(p | q) = \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)}$$

Collocation (Co-occurrence)

- Co-occurrence patterns of words and word classes reveal significant information about how a language is used
 - pragmatics
- Used in building dictionaries (lexicography) and for IR tasks such as phrase detection, query expansion, etc.
- Co-occurrence based on *text windows*
 - typical window may be 100 words
 - smaller windows used for lexicography, e.g. adjacent pairs or 5 words

Collocation and Linguistic Relations

Relation	Word x	Word y	Separation	
			mean	variance
fixed	<i>bread</i>	<i>butter</i>	2.00	0.00
	<i>drink</i>	<i>drive</i>	2.00	0.00
compound	<i>computer</i>	<i>scientist</i>	1.12	0.10
	<i>United</i>	<i>States</i>	0.98	0.14
semantic	<i>man</i>	<i>woman</i>	1.46	8.07
	<i>man</i>	<i>women</i>	-0.12	13.08
lexical	<i>refraining</i>	<i>from</i>	1.11	0.20
	<i>coming</i>	<i>from</i>	0.83	2.89
	<i>keeping</i>	<i>from</i>	2.14	5.53

Word Pair Statistics from 1988 AP Corpus (Church and Hanks)

Collocation

- Typical measure used is the point version of the mutual information measure (compared to the expected value of I, sometimes called EMIM)

$$I(x, y) = \log \frac{p(x, y)}{p(x)p(y)}$$

- Paired t test also used to compare collocation probabilities

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

- Other tests such as Chi-square can also be used

A very small sample of the concordances to "strong" (from 1988 AP newswire)

f somebody catching it has become quite strong , ' ' the newspaper said . *E* *S* The Monitor said necessarily appear on the surface to be strong , ' ' said McGovern , who first drew attention in the the actress . *E* *S* Kristy is ' ' very strong , although she doesn't necessarily appear on the surf eratures . *E* *S* ' ' What we need is a strong , energetic , young , brilliant man , and that's what S* ' ' You know , the Soviet Union has a strong , energetic man , ' ' Cash told about 150 people who s rt showed . *E* *S* The impression of a strong , potentially inflationary economy was heightened by or the November election . *E* *S* ' ' A strong , well-financed Republican Party will be a benefit to mathematics is regarded in the West as strong , *E* *S* It is not known exactly what changes the Ce evious months . *E* *S* Sales were up a strong 1.2 percent in December and 0.3 percent in November , about Mr . Gorbachev and they welcomed strong American leadership of the NATO alliance . *E* *S* We et Ambassador Yuri Dubinin to receive a strong U.S . protest and that Defense Secretary Frank C . Ca uded Hughes ' direction . *E* *S* ' ' As strong and independent as I come off on the set , I need a d rtner . *E* *S* Our commercial ties are strong and of great benefit to people on both sides of the b analyst Linda Simard said crude opened strong at the start , picking up on moderate overnight gains f follow-through buying from Thursday's strong close . *E* *S* Early trading volume was light ahead is energetic person-to-person style and strong conservative message will make him the conservative a c population , always have maintained a strong cultural and ethnic identity . *E* *S* One of the Est themselves ... and we've got to have a strong defense . ' ' .End of Discourse *E* *S* .Story 88 /mur 0457 *E* *S* TAIPEI , Taiwan (AP) - A strong earthquake centered off Taiwan's eastern coast violen s , some analysts said the figures were strong enough to indicate consumers were not dragging the we s pointed toward the December report as strong evidence of the long-awaited reversal in the nation's 5.8 billion Canadian dollars largely on strong foreign sales of forest products . *E* *S* However , , and basically a black school that was strong in academics , ' ' Dade said . *E* *S* ' ' Before , we finishing third in Iowa , maintained a strong lead in New Hampshire - but he no longer had the huge

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Table 1: Some Interesting Associations with *strong* and *powerful* in the 1988 AP Corpus (N = 44.3 million)

I(x:y)	fx _y	fx	fy	x	y
10.47	7	7809	28	strong	northberly
9.76	23	7809	151	strong	showings
9.30	7	7809	63	strong	believer
9.22	14	7809	133	strong	second-place
9.17	6	7809	59	strong	runup
9.04	10	7809	108	strong	currents
8.85	62	7809	762	strong	supporter
8.84	8	7809	99	strong	proponent
8.68	15	7809	208	strong	thunderstorm
8.45	7	7809	114	strong	odor
8.66	7	1984	388	powerful	legacy
8.58	7	1984	410	powerful	tool
8.35	8	1984	548	powerful	storms
8.32	31	1984	2169	powerful	minority
8.14	9	1984	714	powerful	neighbor
7.98	9	1984	794	powerful	Tamil
7.93	8	1984	734	powerful	symbol
7.74	32	1984	3336	powerful	figure
7.54	10	1984	1204	powerful	weapon
7.47	24	1984	3029	powerful	post

CMPS

Table 2: An Example of the t-score

Strong w				Powerful w			
t	strong w	powerful w	w	t	strong w	powerful w	w
12.42	161	0	showing	-7.44	1	56	than
11.94	175	2	support	-5.60	1	32	figure
10.08	550	68	,	-5.37	3	31	minority
9.97	106	0	defense	-5.23	1	28	of
9.76	102	0	economy	-4.91	0	24	post
9.50	97	0	demand	-4.63	5	25	new
9.40	95	0	gains	-4.35	27	36	military
9.18	91	0	growth	-3.89	0	15	figures
8.84	137	5	winds	-3.59	6	17	presidency
8.02	83	1	opposition	-3.57	27	29	political
7.78	67	0	sales	-3.33	0	11	computers

Table 3: Answer Different Questions

Associated with strong					Associated with powerful				
I(strong; w)	t	strong w	powerful w	w	I(powerful; w)	t	strong w	powerful w	w
10.47	1.73	7	0	northerly	8.66	-2.53	1	7	legacy
9.76	3.12	23	1	showings	8.58	-2.67	0	7	tool
9.30	1.73	7	0	believer	8.35	-2.33	4	8	storms
9.22	2.98	14	0	second-place	8.32	-5.37	3	31	minority
9.17	1.51	6	0	runup	8.14	-3.02	0	9	neighbor
9.04	1.22	10	1	currents	7.98	-3.02	0	9	Tamil
8.85	7.45	62	0	supporter	7.93	-2.59	2	8	symbol
8.84	1.94	8	0	proportion	7.74	-3.89	0	15	figures
8.68	0.89	20	4	thunderstorms	7.54	-3.18	0	10	weapon
8.45	1.73	7	0	odor	7.47	-4.91	0	24	post

Table 8: What does a boat do?

(N = 24,677,658; f(x, y) ≥ 3).

I(x;y)	f(x,y)	f(x)	f(y)	x	y	I(x;y)	f(x,y)	f(x)	f(y)	x	y
11.01	16	984	194	boat/S	capsize/V	3.09	4	984	11768	boat/S	fail/V
9.30	51	984	2036	boat/S	sink/V	2.72	4	984	15244	boat/S	stop/V
8.17	3	984	262	boat/S	cruise/V	2.59	5	984	20894	boat/S	accord/V
7.40	6	984	890	boat/S	sail/V	2.54	4	984	17266	boat/S	reach/V
7.27	3	984	488	boat/S	tow/V	2.14	3	984	17074	boat/S	lose/V
7.18	3	984	518	boat/S	turn_in/V	2.09	6	984	35456	boat/S	leave/V
6.83	3	984	660	boat/S	collide/V	2.04	4	984	24410	boat/S	keep/V
6.61	3	984	772	boat/S	drown/V	2.04	6	984	36494	boat/S	kill/V
6.34	4	984	1238	boat/S	drag/V	1.69	6	984	46624	boat/S	be_in/V
6.28	3	984	968	boat/S	escort/V	1.61	3	984	24714	boat/S	put/V
6.04	4	984	1522	boat/S	overturn/V	1.38	8	984	77238	boat/S	take/V
5.90	5	984	2096	boat/S	rescue/V	1.36	3	984	29338	boat/S	hold/V
5.43	5	984	2902	boat/S	approach/V	1.28	4	984	41232	boat/S	use/V
4.64	16	984	16068	boat/S	carry/V	1.26	3	984	31506	boat/S	become/V
4.43	9	984	10470	boat/S	hit/V	0.94	19	984	247542	boat/S	have/V
4.18	4	984	5524	boat/S	travel/V	0.67	3	984	47214	boat/S	begin/V
3.86	6	984	10348	boat/S	pass/V	0.57	3	984	50766	boat/S	get/V
3.71	4	984	7656	boat/S	attack/V	0.17	4	984	89256	boat/S	do/V
3.48	3	984	6748	boat/S	injure/V	-0.35	26	984	830120	boat/S	be/V
3.38	4	984	9614	boat/S	fire/V	-0.35	3	984	95880	boat/S	make/V
3.30	3	984	7634	boat/S	operate/V	-3.38	4	984	1045494	boat/S	say/V

Copyright ©

Copyright © Bruce Croft and/or James Allan

Table 9: What do you typically do with *food* and *water*?

Computed over Parsed AP Corpus (N = 24.7 million SVO triples)

Associated with food						Associated with water					
I(x;y)	fx	fy	x	y		I(x;y)	fx	fy	x	y	
9.62	6	84	2240	hoard/V	food/O	9.05	16	208	3574	conserve/V	water/O
8.83	9	218	2240	go_without/V	food/O	8.98	18	246	3574	boil/V	water/O
7.68	58	3114	2240	eat/V	food/O	8.64	6	104	3574	ration/V	water/O
6.93	8	722	2240	consume/V	food/O	8.45	10	198	3574	pollute/V	water/O
6.42	6	772	2240	run_of/V	food/O	8.40	20	408	3574	contaminate/V	water/O
6.29	14	1972	2240	donate/V	food/O	8.37	38	794	3574	pump/V	water/O
6.08	17	2776	2240	distribute/V	food/O	7.86	6	178	3574	walk_on/V	water/O
5.14	51	15900	2240	buy/V	food/O	7.81	43	1320	3574	drink/V	water/O
4.80	53	21024	2240	provide/V	food/O	7.39	15	618	3574	spray/V	water/O
4.65	13	5690	2240	deliver/V	food/O	7.39	9	370	3574	poison/V	water/O

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Markov Models

- Modeling a sequence of events where probability depends on previous events
- Markov properties

- Limited Horizon

$$P(X_{t+1} = k | X_1, \dots, X_t) = P(X_{t+1} = k | X_t)$$

- Time invariant

$$= P(X_2 = k | X_1)$$

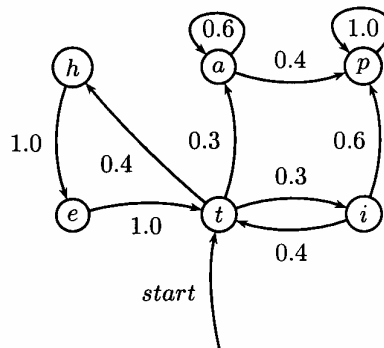
- A Markov chain is described by a transition probability matrix

$$a_{ij} = P(X_{t+1} = s_j | X_t = s_i)$$

CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Markov Models



(from Manning and Schutze)

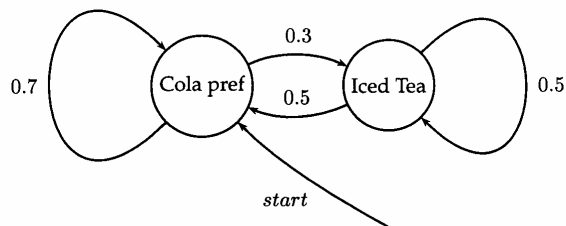
CMPSCI 646

Copyright © Bruce Croft and/or James Allan

Hidden Markov Models

- Don't know state sequence that the model passes through, only some probabilistic function of it
 - underlying events probabilistically generating surface events
- Both regular and hidden Markov models used for part of speech tagging
 - regular is trained using a tagged corpus
 - HMM approach assumes that an underlying Markov chain of parts of speech generates actual words in the text

HMM Example



Output probability given From state

	cola	iced tea	lemonade
CP	0.6	0.1	0.3
IP	0.1	0.7	0.2

(From Manning and Schutze)

Hidden Markov Models

- 3 basic questions
 - Given a model, how do we efficiently compute how likely a certain observation is?
 - Given an observation sequence and a model, how do we choose a state sequence that best explains the observations
 - Given an observation sequence and a space of possible models, how do we find the model that best explains the observations
- Viterbi algorithm commonly used for second problem
- Baum-Welch algorithm used for third problem

Language Models

- “Shannon game” - guess the next word in a text
- Particularly important for speech recognition, OCR
- n-gram models commonly used to estimate probabilities of words
 - unigram, bigram, trigram
 - n-gram model is equivalent to an (n-1)th order Markov model
- Estimates must be smoothed by, for example, interpolating combinations of n-gram estimates

$$P(w_n | w_{n-1}, w_{n-2}) = \lambda_1 P_1(w_n) + \lambda_2 P_2(w_n | w_{n-1}) + \lambda_3 P_3(w_n | w_{n-1}, w_{n-2})$$

- HMM algorithms can determine the optimal parameter settings