

-10-2005

M-60's

M-90's

Retrieval Models: Language Model

- ① Create a probability model for each document
- ② Often create a probability model for query
- ③ Compare ① & ② - docs which are "closer" are ranked higher

Simple view

$$\vec{P}_{t;d} = \frac{tf_{t;d}}{\sum_i tf_{i;d}}$$

ex. The quick brown dog dog
 t: the, quick, brown, dog
 $(\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{2}{5})$

q = "brown dogs"

$$? Pr[(brown, dog) | d] = \prod_{t \in q} Pr[t | d] = \frac{1}{5} \cdot \frac{2}{5} = \frac{2}{25}$$

Mathematical Aside

$Pr(s)$ - probability of ("event") s

$Pr(s|M)$ - probability of s given M

Independence: $Pr(s_1, s_2) = Pr(s_1) \cdot Pr(s_2)$

$Pr(q_1, q_2, q_3 \dots | M) = Pr(q_1 | M) \cdot Pr(q_2 | M) \dots = \prod_i Pr(q_i | M)$

Bayes Rule: $Pr(A|B) = \frac{Pr(B|A) \cdot Pr(A)}{Pr(B)} = \frac{Pr(B|A) \cdot Pr(A)}{\sum_j Pr(B|A_j) \cdot Pr(A_j)}$

Coding / Information Theory

① Coding & Entropy:

- suppose events generated by dist (P_1, P_2, \dots, P_n)

- optimum code length for shortest avg code length: $\log \frac{1}{P_i}$

- avg code: $= P_1 \log \frac{1}{P_1} + P_2 \log \frac{1}{P_2} \dots = \sum_i P_i \log \frac{1}{P_i} = H(\vec{p}) = \text{entropy of dist}$

Relative Entropy (cross-entropy, KL-distance)

$$KL(\vec{p} \parallel \vec{q}) = \sum_i P_i \log \frac{1}{P_i} - \sum_i P_i \log \frac{1}{Q_i} = \sum_i P_i \log \frac{P_i}{Q_i}$$

- how many bits are wasted by encoding p with q 's encoding

- not symmetric

all in $\log_2 = \log$

$$\log_a b = \frac{\log b}{\log a}$$

$$\log_2 = \log$$

$$\log$$

$$\ln = \log_e$$

And M.J

Language Models for Documents

Unigram model: $Pr(t|d) = \text{relative frequency of term in doc.}$
$$= \frac{tf_{t,d}}{\sum_{t \in d} tf_{t,d}}$$

Higher order models

- n-gram: 2-gram - look at word pairs...
- caching: $w_1 w_2$ based on a window
- grammar: based on a parse tree

don't help significantly
We will consider unigram models & "smooth" them later

Now, assume we don't model queries, just compute prob of generating query from model (S).

Two-ways

- (1) Multi-nomial: $Pr[q_1, q_2, \dots, q_n | M] = \prod Pr(q_i | M)$
- (2) Multiple Bernoulli: $Pr[q_1, q_2, \dots, q_n | M] = \prod_{w \in q_i} P(w|M) \prod_{w \notin q_i} (1 - P(w|M))$

Generate doc from query

$$P(D|q) = \frac{P(q|D) \cdot P(D)}{P(q)}$$

$P(D)$ - belief in how similar to D is

Doesn't affect rank of docs (corpus)

- $id = 0 = 0 = 0$
- $id = 1 = 1 = 1$
- $id = 2 = 2 = 2$
- $id = 3 = 3 = 3$
- $id = 4 = 4 = 4$
- $id = 5 = 5 = 5$
- $id = 6 = 6 = 6$
- $id = 7 = 7 = 7$
- $id = 8 = 8 = 8$
- $id = 9 = 9 = 9$
- $id = 10 = 10 = 10$

L.M. cont.

5/11/2005

Issues w/ L.M.

- ① How to model docs? - prob dist over words
- ② How to model queries? - prob dist over words
- ③ How to "assess" similarity of doc to query (or vice versa)

2 methods: $P(Q|D)$ * $P(D|Q)$

third: $\text{sim}(d, q)$ - use KL-distance

Let D be a dist corresponding to a doc * - likely to perform

Let Q be a dist corresponding to a query * - likely to perform

$\Rightarrow \text{KL}(Q||D)$ - "distance" from query to doc *
 $= \sum_i q_i \log \frac{q_i}{d_i}$

Smoothing

Issue: how to gen dist over document

- probably not enough data to build good model just on doc itself
- should incorporate some "prior belief" based on dist in entire corpus
- have to move smoothly from prior belief to totally evidence based

Generalize using a coin

Laplace Smoothing - determining estimates for probabilities from data

- Problem: Given lots of data, the maximum likelihood estimate is good. For limited data, this is often poor

Slowly...

Consider 5 possibilities:

- Hypothesis $\left\{ \begin{array}{l} H_1 = Pr = 0 = h_1 \\ H_2 = Pr = \frac{1}{4} = h_2 \\ H_3 = Pr = \frac{1}{2} = h_3 \\ H_4 = Pr = \frac{3}{4} = h_4 \\ H_5 = Pr = 1 = h_5 \end{array} \right.$

Evidence: E : the sun has risen each of the last N days

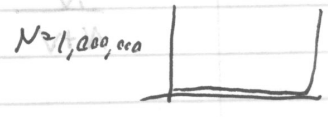
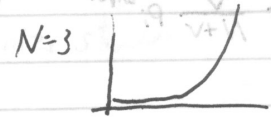
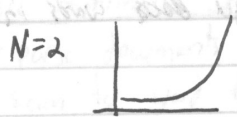
$q = \frac{1-M}{N} = \text{Good News, bad news}$

$P[E|H_i] = h_i^N$

$P[H_i] = \frac{1}{5}$
uniform, no assumptions

$$Pr[H_i|E] = \frac{Pr[E|H_i] \cdot Pr[H_i]}{\sum_j Pr[E|H_j] \cdot Pr[H_j]} = \frac{Pr[E|H_i] \cdot Pr[H_i]}{\sum_j h_j^N \cdot \frac{1}{5}} = \frac{h_i^N \cdot \frac{1}{5}}{\sum_j h_j^N \cdot \frac{1}{5}} = \frac{h_i^N}{\sum_j h_j^N}$$

eg: $N=0$ $P[H_i|E] = \frac{1}{5}$
 $N=1$ $P[H_i|E] = \frac{h_i}{\sum_j h_j} = \frac{h_i}{0 + \frac{1}{5} + \frac{1}{5} + \frac{1}{5} + 1} = \frac{h_i}{\frac{3}{5} + 1} = \frac{h_i}{\frac{8}{5}} = \frac{5}{8} h_i$



What's the answer: $1 \cdot (.4) + (\frac{3}{4}) \cdot (.3) + (\frac{1}{2}) \cdot (.2) + (\frac{1}{4}) \cdot (.1) + 0(0) = \frac{3}{4}$
 -break even point

$p(x) = \frac{1}{N+1}$

Now consider all possible hypotheses $x \in [0, 1]$

$$Pr[x|E] = \frac{Pr[E|x] \cdot 1}{\int_{x=0}^1 Pr[E|x] \cdot 1 dx} = \frac{x^N \cdot 1}{\int_{x=0}^1 x^N dx} = \frac{x^N}{\frac{x^{N+1}}{N+1} \Big|_0^1} = \frac{x^N}{\frac{1}{N+1}} = (N+1)x^N$$

Expected Answer:

$$\int_0^1 (N+1)x^N \cdot x dx = (N+1) \int_0^1 x^{N+1} dx = (N+1) \frac{x^{N+2}}{N+2} \Big|_0^1 = \frac{N+1}{N+2} = p_i$$

So, if N successes in N trials, belief in success is $\frac{N+1}{N+2}$
 Generalization: ① If M successes in N trials, belief in success is $\frac{M+1}{N+2}$ for each

② if M_i occurrences of outcome i in N trials for K total possible outcomes:

$$\frac{M_i + 1}{N + K} = p_i$$

$$\text{Max Likelihood} = \frac{M_i}{N} = p_i$$

Generalizations of Laplace Smoothing

① Jelinek-Mercer - take into account non-uniform prior (eg. corpus prior)

$$p_i = \lambda p_i^{ML} + (1-\lambda) p_i^{corpus} \quad \lambda = 1 - \frac{1}{N+1}$$

② Dirichlet Smoothing - μ -arbitrary - generalization of K

$$\frac{N}{N+\mu} p_i^{ML} + \frac{\mu}{N+\mu} p_i^{corpus}$$

③ Witten-Bell

$$\frac{N}{N+V} p_i^{ML} + \frac{V}{N+V} p_i^{corpus}$$

N = length of doc

V = # unique words in doc

$$\frac{1}{N} = \left(\frac{1}{N}\right) \left(\frac{1}{N}\right) + \left(\frac{1}{N}\right) \left(\frac{1}{N}\right) + \dots + (1) \cdot \left(\frac{1}{N}\right)$$

$$P(x) = \prod_{i=1}^N p_{x_i} = \prod_{i=1}^N \left(\frac{N}{N+V} p_{x_i}^{ML} + \frac{V}{N+V} p_{x_i}^{corpus} \right)$$

$$\frac{N!}{(N+V)^N} \prod_{i=1}^N p_{x_i}^{ML} \left(\frac{V}{N+V} \right)^V$$

① If M occurs in N trials, point in success is M .
 ② If M occurs i in N trials per N trials.

$$p_i = \frac{M+1}{N+K}$$

S-24-2004

Language Model Example

D1 The cat in the hat

D2 One fish two fish red fish blue fish

D3 The fish in the red hat

3 documents in corpus

Corpus Vocab - 9 words (unique)

	The	cat	in	hat	one	fish	two	red	blue	Doc Length
D1	2	1	1	1	0	0	0	0	0	5
D2	0	0	0	0	1	4	1	1	1	8
D3	2	0	1	1	0	1	0	1	0	6

Max Likelihood Model	D1	$\frac{2}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	0	0	0	0	0
	D2	0	0	0	0	$\frac{1}{8}$	$\frac{4}{8}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{8}$
	D3	$\frac{2}{6}$	0	$\frac{1}{6}$	$\frac{1}{6}$	0	$\frac{1}{6}$	0	$\frac{1}{6}$	0

Query: "red fish"

- Score associated w/ doc is prob. query words would be generated when drawing words at random according to LM.

$$\text{Score}_{ML}(D1) = 0 \cdot 0 = 0$$

$$\text{Score}_{ML}(D2) = \frac{1}{8} \cdot \frac{4}{8} = \frac{4}{64} = \frac{1}{16}$$

$$\text{Score}_{ML}(D3) = \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{36}$$

Smoothing

Laplace: $\frac{\# \text{ Successes} + 1}{\# \text{ trials} + \# \text{ outcomes (possible)}}$

Laplace Smoothing

	the	cat	in	hat	one	fish	two	red	blue
D1	$\frac{2+1}{5+9}$	$\frac{1+1}{5+9}$	$\frac{1+1}{5+9}$	$\frac{1+1}{5+9}$	$\frac{0+1}{5+9}$	$\frac{2+1}{5+9}$	$\frac{2+1}{5+9}$	$\frac{0+1}{5+9}$	$\frac{0+1}{5+9}$
D2
D3

Score Laplace (D1) = $\frac{1}{14} \cdot \frac{1}{14} =$ small but not 0

Score Laplace (D2) = $\frac{1+1}{8+9} \cdot \frac{4+1}{8+9} = \frac{2}{17} \cdot \frac{5}{17}$

Score Laplace (D3) =

Witten-Bell smoothing

$LM(D) = \frac{\# \text{ words in doc}}{\# \text{ words in doc} + \# \text{ unique words in doc}} \cdot LM_{ML} + \frac{\# \text{ unique words in doc}}{\# \text{ unique words in doc} + \# \text{ words in doc}} \cdot LM_{corpus}$

need corpus language model...

often corpus LM is average of ML LMs

corpus	the	cat	in	hat	one	fish	two	red	blue
	$\frac{\frac{2}{5} + 1 + \frac{1}{5}}{3}$	$\frac{\frac{1}{5} + 1 + 0}{3}$							
$LM_{WB}(D1)$	$\frac{5}{5+4} \cdot \frac{2}{5} + \frac{4}{5+4} \cdot \frac{2+1}{3}$								

optional for Witten-Bell smoothing

$V = \# \text{ unique words in doc}$ Heaps Law: $V = KN^\beta$ $\beta \approx \frac{1}{2}$

$N = \# \text{ words in doc}$

$\frac{N}{N+V} \approx \frac{N}{N+KN^\beta}$ $\frac{V}{N+V} = \frac{KN^\beta}{N+KN^\beta}$ - shrinks w/ N

grows with N

S-24-2004

Language Model Example

- D1 The cat in the hat
 D2 One fish two fish red fish blue fish
 D3 The fish in the red hat

3 documents in corpus

Corpus Vocab - 9 words (unique)

	The	cat	in	hat	one	fish	two	red	blue	Doc Length
D1	2	1	1	1	0	0	0	0	0	5
D2	0	0	0	0	1	4	1	1	1	8
D3	2	0	1	1	0	1	0	1	0	6

Max Likelihood Model	D1	$\frac{2}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	0	0	0	0	0
	D2	0	0	0	0	$\frac{1}{8}$	$\frac{4}{8}$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{8}$
	D3	$\frac{2}{6}$	0	$\frac{1}{6}$	$\frac{1}{6}$	0	$\frac{1}{6}$	0	$\frac{1}{6}$	0

Query: "red fish"

- Score associated w/ doc is prob. query words would be generated when drawing words at random according to LM.

$$\text{Score}_{ML}(D1) = 0 \cdot 0 = 0$$

$$\text{Score}_{ML}(D2) = \frac{1}{8} \cdot \frac{4}{8} = \frac{4}{64} = \frac{1}{16}$$

$$\text{Score}_{ML}(D3) = \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{36}$$

AND OR

- register smoothing

Smoothing

Laplace: $\frac{\# \text{ Successes} + 1}{\# \text{ trials} + \# \text{ Outcomes (possible)}}$

Laplace Smoothing

D1:	the	cat	in	hat	one	fish	two	red	blue
D1	$\frac{2+1}{5+9}$	$\frac{1+1}{5+9}$	$\frac{1+1}{5+9}$	$\frac{1+1}{5+9}$	$\frac{0+1}{5+9}$	$\frac{2+1}{5+9}$	$\frac{2+1}{5+9}$	$\frac{0+1}{5+9}$	$\frac{0+1}{5+9}$
D2									
D3									

Score Laplace (D1) = $\frac{1}{14} \cdot \frac{1}{14}$ = small but not 0

Score Laplace (D2) = $\frac{1+1}{8+9} \cdot \frac{4+1}{8+9} = \frac{2}{17} \cdot \frac{5}{17}$

Score Laplace (D3) =

Witten-Bell Smoothing

$LM(D) = \frac{\# \text{ words in doc}}{\# \text{ words in doc} + \# \text{ unique words in doc}} \cdot LM_{\text{doc}} + \frac{\# \text{ unique words in doc}}{\# \text{ words in doc} + \# \text{ unique words in doc}} \cdot LM_{\text{corpus}}$

need corpus language model...

often corpus LM is average of ML LMs

corpus	the	cat	in	hat	one	fish	two	red	blue
	$\frac{2+1}{5}$	$\frac{1+1}{3}$							
LM _{WB} (D1)	$\frac{5}{5+4} \cdot \frac{2}{5} + \frac{4}{5+4} \cdot \frac{2+1}{3}$								

variant for Witten-Bell smoothing

$V = \# \text{ unique words in doc}$ Heaps Law: $V = KN^\beta$ $\beta \approx \frac{1}{2}$

$N = \# \text{ words in doc}$

$\frac{N}{N+V} \approx \frac{N}{N+KN^\beta}$ $\frac{V}{N+V} = \frac{KN^\beta}{N+KN^\beta}$

grows with N -shrinks w/ N

Two ways:
 generate query \uparrow query \downarrow document model \downarrow k
 performs better \leftarrow (1) Multinomial : $\Pr[q_1, q_2, \dots, q_k | M] = \prod_{i=1}^k \Pr[q_i | M]$
 come first \leftarrow (2) Multiple Bernoulli
 generate these query words and no other word \leftarrow $\Pr[q_1, q_2, \dots, q_k | M] = \prod_{w_i \in q_1, q_2, \dots, q_k} \Pr(w_i | M) \cdot \prod_{w_i \notin q_1, q_2, \dots, q_k} (1 - \Pr(w_i | M))$
 generate query and nothing but query

Something in language modeling = parameter estimation in statistics.

\rightarrow Why calculate prob. of generating a query given a doc. but not prob. of a document given a query?

Generate doc from query instead:

$$\Pr(D | q) = \frac{\Pr(q | D) \Pr(D)}{\Pr(q)}$$

$\Pr(q)$ \rightarrow doesn't affect ranking of docs

or

\mathcal{L} can generate probability models for queries also.
 Then, we can use KL-distance to see how two distributions are different from each other.

Lecture #5

Issues w/ L.M.

- 1) How to model docs?
- 2) How to model queries?
- 3) How to "assess" sim. of doc. to query (or vice versa)?
- 4) \hookrightarrow two methods: $\Pr(q | D)$ hypotheses $x \in \{0, 1\}$
 $\Pr(D | q)$ $\rightarrow \Pr(x) = \text{uniform}$

third: sim(d, q) - use KL-distance

Let D be a distr. corresponding to a doc.

Let Q be a distr. corresponding to a query

$\Rightarrow KL(Q||D)$ = "distance" from query to doc.

$$= \sum_i q_i \log \frac{q_i}{d_i}$$

Laplace Smoothing

Determining estimates for probabilities from data.

- Problem: Given lots of data, the maximum likelihood estimate, e.g.

$$\Pr[\text{Heads}] = \frac{\# \text{heads}}{\# \text{trials}} \text{ is good.}$$

- for limited data, this is often poor.

Will sun rise tomorrow?

Consider 5 possibilities

$$\begin{cases} H_1 - \Pr = 0 = h_1 \\ H_2 - \Pr = 1/4 = h_2 \\ H_3 - \Pr = 1/2 = h_3 \\ H_4 - \Pr = 3/4 = h_4 \\ H_5 - \Pr = 1 = h_5 \end{cases}$$

Evidence E : The sun has risen each of the last N days.

$$\Pr[H_i|E] = \frac{\Pr[E|H_i] \cdot \Pr[H_i]}{\Pr[E]}$$

$$= \frac{\Pr[E|H_i] \cdot \Pr[H_i]}{\sum_j \Pr[E|H_j] \Pr[H_j]}$$

$$\begin{aligned} & \text{uniform} \\ & = \frac{h_i^N \cdot 1/5}{\sum_j h_j^N \cdot 1/5} = \frac{h_i^N}{\sum_j h_j^N} \end{aligned}$$

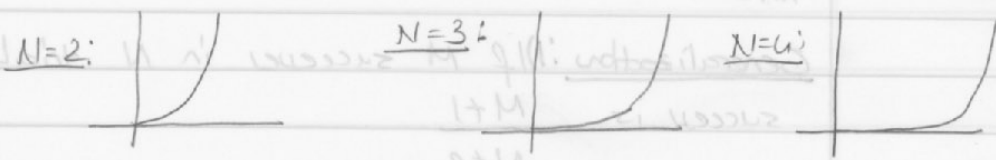
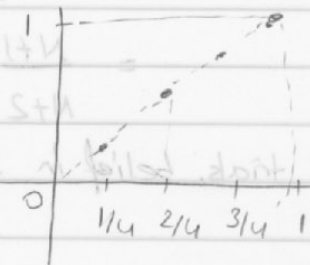
Emine Yilmaz



eg. $N=0$: $\Pr[H_i|E] = \frac{h_i^0}{\sum h_i^0} = \frac{1}{5}$

Diagram showing a horizontal axis with points at $1/4, 1/2, 3/4, 1$. A vertical line is drawn at $1/5$. There are four dots on the axis, one in each interval defined by the vertical line and the axis points.

$N=1$: $\Pr[H_i|E] = \frac{h_i^1}{\sum h_i^1} = \frac{h_i}{0 + 1/4 + 1/2 + 3/4 + 1} = \frac{h_i}{5/2} = \frac{2}{5} h_i$



$N=3$: $\Pr[H_i|E] = \frac{h_i^3}{\sum h_i^3} = \frac{h_i^3}{0 + (1/2)^3 + 1} = \frac{h_i^3}{5/4}$

Problem: Laplace (here) considers only a small set that has been seen. However, in reality, there can be much more different events happening.

Now consider all possible hypotheses $x \in [0, 1]$.

$$\Pr[X|E] = \frac{\Pr[E|X] \cdot \Pr[X]}{\int_{x=0}^1 \Pr[E|x] \cdot 1 dx}$$

$\Pr[X] = \text{uniform}$

Ernie's Diner



$$= \frac{x^N \cdot 1}{\int_0^1 x^N dx} = \frac{x^N}{\frac{x^{N+1}}{N+1} \Big|_0^1} = \frac{x^N}{\frac{1^{N+1}}{N+1}} = (N+1)x^N$$

My beliefs now are not just dots, but curves.

Expected answer:

$$\int_0^1 (N+1)x^N \cdot x dx = (N+1) \int_0^1 x^{N+1} dx = (N+1) \frac{x^{N+2}}{N+2} \Big|_0^1 = \frac{N+1}{N+2}$$

So, if N successes in N trials, belief in success is $\frac{N+1}{N+2}$.

Generalization 1: If M successes in N trials, belief in success is $\frac{M+1}{N+2}$.

2) If M_i occurrences of outcome i in N trials for K total possible outcomes

how many times I see that word $\leftarrow M_i + 1$
docs $\leftarrow N + k$

$$\text{Max. Likelihood} = \frac{M_i}{N} = p_i^{ML}$$

Generalizations of Laplace Smoothing

1) DeGroot-Mercer: take into account non-uniform prior (e.g. corpus prior)

$$p_i = \lambda \cdot p_i^{ML} + (1-\lambda) p_i^{corpus} \quad \lambda = 1 - \frac{1}{N+k}$$

$$\frac{\sum h_i \cdot \frac{1}{N+k}}{\sum h_i \cdot \frac{1}{N+k}}$$

Dirichlet Smoothing

$$\frac{N}{N+M} p_i^{ML} + \frac{M}{N+M} p_i^{\text{corpus}}$$

Witten-Bell:

$$\frac{N}{N+V} p_i^{ML} + \frac{V}{N+V} p_i^{\text{corpus}}$$

N = length of doc.

V = # unique words in doc.