



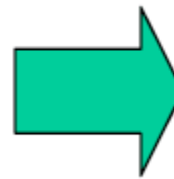
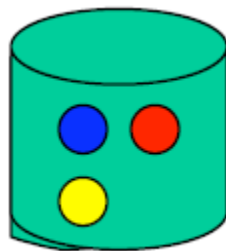
estimation

- Want to estimate MQ and/or MD from Q and/or D
- General problem:
 - given a string of text S ($= Q$ or D), estimate its language model MS
 - S is commonly assumed to be an i.i.d. random sample from MS
 - Independent and identically distributed
- Basic Language Models
 - maximum-likelihood estimator and the zero frequency problem
 - discounting, interpolation techniques
 - Bayesian estimation

maximum likelihood

- count relative frequencies of words in S
- maximum-likelihood property:
 - assigns highest possible likelihood to the observation
- unbiased estimator:
 - if we repeat estimation an infinite number of times with different starting points S , we will get correct probabilities (on average)
 - this is not very useful...

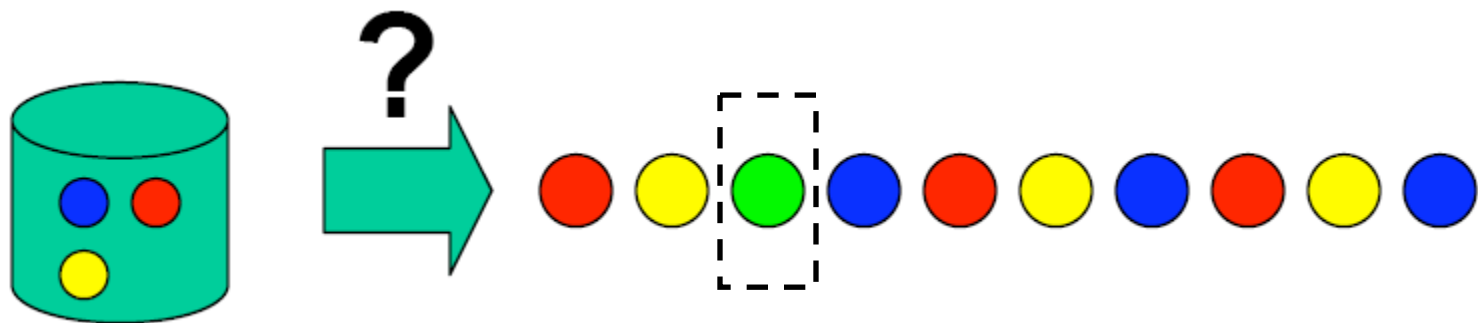
$$P_{ml}(w|M_S) = \#(w, S) / |S|$$



$$\begin{aligned} P(\bullet) &= 1/3 \\ P(\bullet) &= 1/3 \\ P(\bullet) &= 1/3 \\ P(\bullet) &= 0 \\ P(\bullet) &= 0 \end{aligned}$$

zero-frequency problem

- Suppose some event not in our observation S
 - Model will assign zero probability to that event
 - And to any set of events involving the unseen event
- Happens very frequently with language
- It is incorrect to infer zero probabilities
 - especially when creating a model from short samples





Laplace smoothing

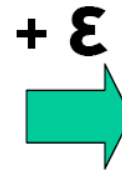
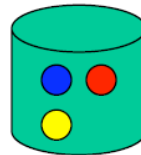
- count events in observed data
 - add 1 to every count
 - renormalize to obtain probabilities
 - it corresponds to uniform priors
-
- if event counts are (m_1, m_2, \dots, m_k) with $\sum_i m_i = N$ then
max likelihood estimates are $(\frac{m_1}{N}, \frac{m_2}{N}, \dots, \frac{m_k}{N})$
laplace estimates are $(\frac{m_1+1}{N+k}, \frac{m_2+1}{N+k}, \dots, \frac{m_k+1}{N+k})$

discounting methods

- Laplace smoothing

- Lindstone correction

- add ϵ to all count,
renormalize



- absolute discounting
 - subtract ϵ , redistribute probab mass

$$P(\bullet) = (1 + \epsilon) / (3 + 5\epsilon)$$

$$P(\bullet) = (1 + \epsilon) / (3 + 5\epsilon)$$

$$P(\bullet) = (1 + \epsilon) / (3 + 5\epsilon)$$

$$P(\bullet) = (0 + \epsilon) / (3 + 5\epsilon)$$

$$P(\bullet) = (0 + \epsilon) / (3 + 5\epsilon)$$



discounting methods

- Held-out estimation
 - Divide data into training and held-out sections
 - In training data, count N_r , the number of words occurring r times
 - In held-out data, count T_r , the number of times those words occur
 - $r^* = T_r/N_r$ is adjusted count (equals r if training matches held-out)
 - Use r^*/N as estimate for words that occur r times
- Deleted estimation (cross-validation)
 - Same idea, but break data into K sections
 - Use each in turn as held-out data, to calculate $T_r(k)$ and $N_r(k)$
 - Estimate for words that occur r times is average of each
- Good-Turing estimation
 - From previous, $P(w|M) = r^* / N$ if word w occurs r times in sample
 - In Good-Turing, steal total probability mass from next most frequent word
 - Provides probability mass for words that occur $r=0$ times
 - Take what's leftover from $r>0$ to ensure adds to one

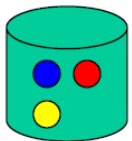
$$TPM(r + 1) = N_{r+1} \cdot \frac{r + 1}{N}$$

$$\begin{aligned} P(w_r|M) &= TPM(r + 1)/N_r \\ &= \frac{N_{r+1}}{N_r} \cdot \frac{r + 1}{N} \end{aligned}$$

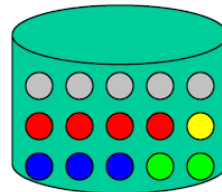
interpolation methods

- Problem with all discounting methods:
 - discounting treats unseen words equally (add or subtract ϵ)
 - some words are more frequent than others
- Idea: use background probabilities
 - “interpolate” ML estimates with General English expectations (computed as relative frequency of a word in a large collection)
 - reflects expected frequency of events

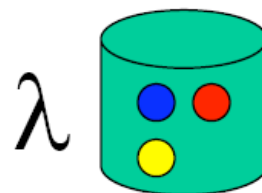
ML estimate



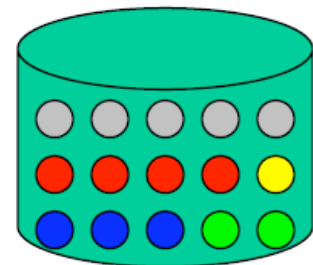
background probability



final estimate =



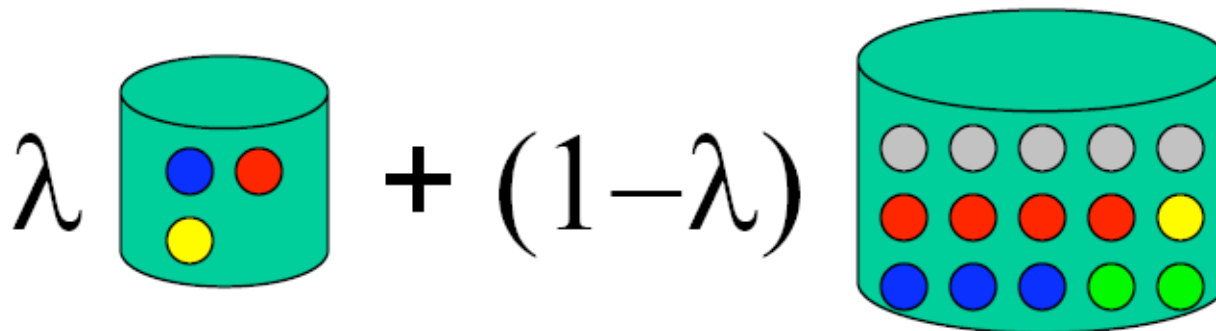
+ (1- λ)





Jelinek Mercer smoothing

- Correctly setting λ is very important
- Start simple
 - set λ to be a constant, independent of document, query
- Tune to optimize retrieval performance
 - optimal value of λ varies with different databases, query sets, etc.



Dirichlet smoothing

- Problem with Jelinek-Mercer:
 - longer documents provide better estimates
 - could get by with less smoothing
- Make smoothing depend on sample size
- N is length of sample = document length
- μ is a constant

$$\underbrace{N / (N + \mu)}_{\lambda} \text{ } \text{cylinder with 3 dots} + \underbrace{\mu / (N + \mu)}_{(1-\lambda)} \text{ } \text{cylinder with 15 dots}$$

The diagram illustrates the Dirichlet smoothing formula. It shows two terms being added together. The first term is $N / (N + \mu)$, which is labeled as λ . This term is represented by a small teal cylinder containing three colored dots (blue, red, and yellow). The second term is $\mu / (N + \mu)$, which is labeled as $(1-\lambda)$. This term is represented by a larger teal cylinder containing fifteen colored dots (red, yellow, blue, and green).

Witten-Bell smoothing

- A step further:
 - condition smoothing on “redundancy” of the example
 - long, redundant example requires little smoothing
 - short, sparse example requires a lot of smoothing
- Derived by considering the proportion of new events as we walk through example
 - N is total number of events = document length
 - V is number of unique events = number of unique terms in doc

$$\underbrace{N / (N + V)}_{\lambda} \text{ } \text{cylinder with 3 colored balls} + \underbrace{V / (N + V)}_{(1-\lambda)} \text{ } \text{cylinder with 15 colored balls}$$



interpolation vs back-off

- Two possible approaches to smoothing
- Interpolation:
 - Adjust probabilities for all events, both seen and unseen
- Back-off:
 - Adjust probabilities only for unseen events
 - Leave non-zero probabilities as they are
 - Rescale everything to sum to one: rescales “seen” probabilities by a constant
- Interpolation tends to work better
 - And has a cleaner probabilistic interpretation

Two-stage smoothing



Query = "the algorithms for data mining"

d1:	0.04	0.001	0.02	0.002	0.003
d2:	0.02	0.001	0.01	0.003	0.004

Two-stage smoothing



Query = "the algorithms for data mining"

d1:	0.04	0.001	0.02	0.002	0.003
d2:	0.02	0.001	0.01	0.003	0.004

$p(\text{"algorithms"}|d1) = p(\text{"algorithm"}|d2)$

$p(\text{"data"}|d1) < p(\text{"data"}|d2)$

$p(\text{"mining"}|d1) < p(\text{"mining"}|d2)$

But $p(q|d1) > p(q|d2)$!

Two-stage smoothing



Query = "the algorithms for data mining"

d1:	0.04	0.001	0.02	0.002	0.003
d2:	0.02	0.001	0.01	0.003	0.004

$p(\text{"algorithms"}|d1) = p(\text{"algorithm"}|d2)$

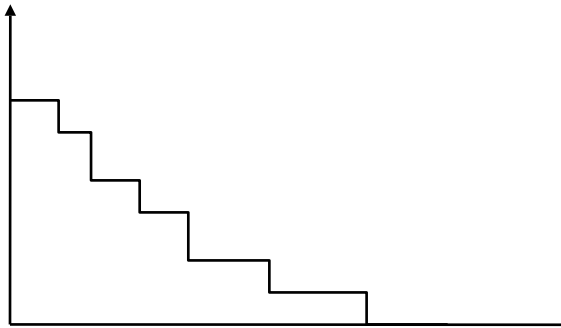
$p(\text{"data"}|d1) < p(\text{"data"}|d2)$

$p(\text{"mining"}|d1) < p(\text{"mining"}|d2)$

But $p(q|d1) > p(q|d2)!$

We should make $p(\text{"the"})$ and $p(\text{"for"})$ **less different** for all docs.

Two-stage smoothing

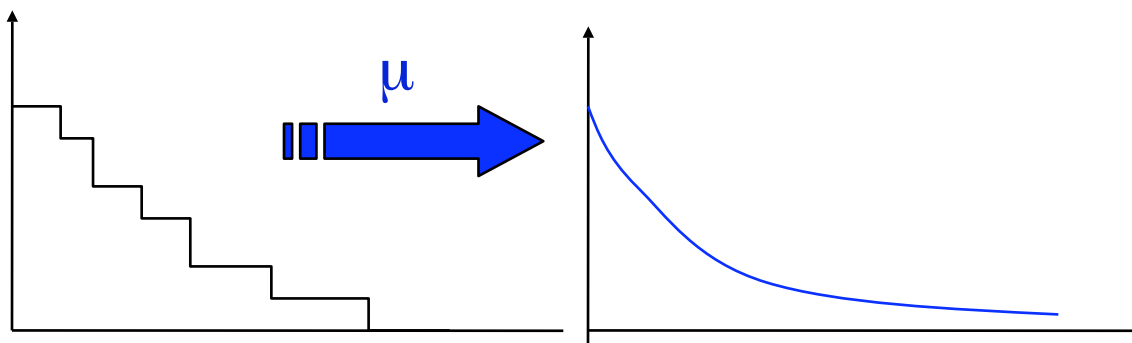


$$P(w|d) = \frac{c(w,d)}{|d|}$$

Two-stage smoothing

Stage-1

- Explain unseen words
- Dirichlet prior(Bayesian)



$$P(w|d) = \frac{c(w,d) + \mu p(w|C)}{|d| + \mu}$$

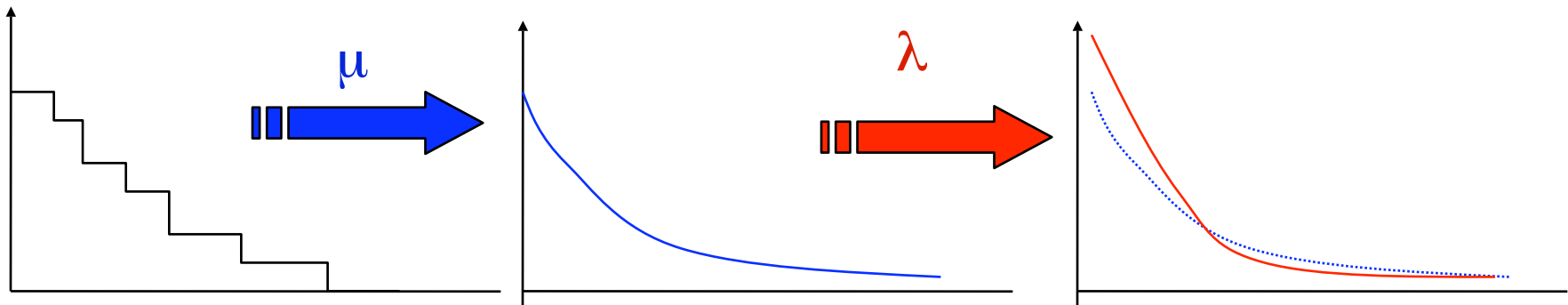
Two-stage smoothing

Stage-1

- Explain unseen words
- Dirichlet prior(Bayesian)

Stage-2

- Explain noise in query
- 2-component mixture



$$P(w|d) = (1-\lambda) \frac{c(w,d) + \mu p(w|C)}{|d| + \mu} + \lambda p(w|U)$$