

CS 589 Lecture 2: Probabilistic Retrieval Models - Complete Study Guide

Course: CS 589 - Information Retrieval

Lecture: 2

Topics: Probability Ranking Principle, RSJ Model, BM25, Language Model-Based Retrieval

Table of Contents

1. Fundamentals: Random Variables & Probability
 2. Probability Ranking Principle (PRP)
 3. Model 1: Robertson & Sparck Jones (RSJ)
 4. Model 2: BM25
 5. Model 3: Language Model-Based Retrieval
 6. Practice Problems with Solutions
 7. Quiz Questions & Detailed Solutions
 8. Video Resources
 9. Quick Reference Guide
 10. Study Tips for Midterm
-

1. Fundamentals: Random Variables & Probability

1.1 Random Variables: The Biased Coin Example

Scenario: You have a biased coin that comes up heads with some unknown probability.

Observation sequence: H, T, H, H, T, T, H, T, H, T, H, H

- Heads: 7 times
- Tails: 5 times

Goal: Estimate the probability of heads

Model: Bernoulli Distribution

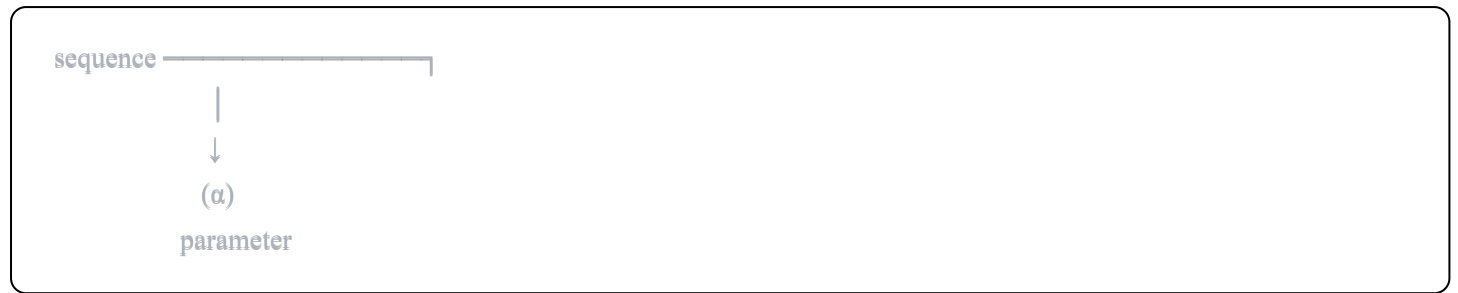
$$P(\text{Heads}) = \alpha$$

$$P(\text{Tails}) = 1 - \alpha$$

Parameter: α (what we want to estimate)

Observation: The sequence we observed

Graphical representation:



1.2 Maximum Likelihood Estimation (MLE)

Question: Given observations, what's the best estimate for α ?

Answer: The value that maximizes the probability of observing our data

Probability of our sequence:

$$\begin{aligned} P(\text{sequence}) &= \alpha \times (1-\alpha) \times \alpha \times \alpha \times (1-\alpha) \times (1-\alpha) \times \alpha \times (1-\alpha) \times \alpha \times (1-\alpha) \times \alpha \times \alpha \\ &= \alpha^{\text{count}(\text{heads})} \times (1-\alpha)^{\text{count}(\text{tails})} \\ &= \alpha^7 \times (1-\alpha)^5 \end{aligned}$$

To maximize, take derivative and set to zero:

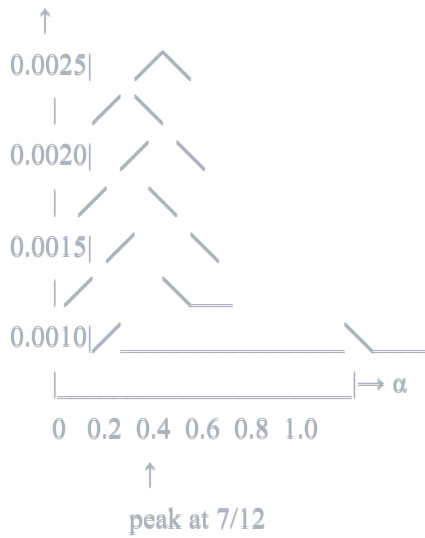
$$d/d\alpha [\alpha^7 \times (1-\alpha)^5] = 0$$

Result (by calculus):

$$\begin{aligned} \hat{\alpha} &= \text{count}(\text{heads}) / [\text{count}(\text{heads}) + \text{count}(\text{tails})] \\ &= 7 / (7 + 5) \\ &= 7/12 \\ &\approx 0.583 \end{aligned}$$

Visualization of likelihood:

Probability



Key insight: The MLE is just the empirical frequency!

1.3 Bayes' Rule

Chain Rule (Joint Distribution)

$$P(A, B) = P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

Bayes' Rule

From the chain rule, we can derive:

$$P(A|B) = P(B|A)P(A) / P(B)$$

Or equivalently:

$$P(A|B) = \frac{P(B|A) \times P(A)}{\sum_{\{X \in \{A, \bar{A}\}\}} P(B|X)P(X)}$$

Terminology

- **$P(A|B)$** : Posterior (what we want to know)
- **$P(B|A)$** : Likelihood (probability of data given hypothesis)
- **$P(A)$** : Prior (what we knew before seeing data)
- **$P(B)$** : Evidence (normalizing constant)

Proportionality

Since $P(B)$ doesn't depend on A :

$$P(A|B) \propto P(B|A)P(A)$$

Example: Medical diagnosis

$$P(\text{disease}|\text{symptoms}) = P(\text{symptoms}|\text{disease}) \times P(\text{disease}) / P(\text{symptoms})$$

1.4 Random Variables in Information Retrieval

Mapping to IR

Coin Toss	Information Retrieval
Outcome: H or T	Relevance: 0 or 1
$\alpha = P(\text{Heads})$	$P(\text{rel}=1 q,d)$ = probability document is relevant
Sequence of tosses	Collection of documents

Key notation

- **q**: query
- **d**: document
- **rel**: relevance (0 = not relevant, 1 = relevant)

Example

Query: "artificial intelligence"

Documents in collection:

```
d1 = [artificial, intelligence, machine, intelligence, information, retrieval]
d2 = [covid, patient, virus, ...]
d3 = [machine, learning, deep, neural, ...]
...
```

Observations:

```
rel = [0, 1, 0, 0, 0, 0, 0, 0]
      | |
      | |
      | └─ d2 is relevant
      └─ d1 is not relevant (for this query)
```

Our goal: Estimate $P(\text{rel}=1|q, d)$ for each document

2. Probability Ranking Principle (PRP)

2.1 The Principle

Statement:

Documents should be ranked by their probability of relevance $P(\text{rel}=1|q,d)$ in decreasing order.

Theorem (Ripley 1996):

The PRP is optimal in the sense that it minimizes expected loss.

Why is this important?

- Provides theoretical foundation for retrieval
- Tells us what to estimate (relevance probability)
- Doesn't tell us HOW to estimate it (that's where models come in)

2.2 From Counting to Probability

Naive approach

$$P(\text{rel}=1|q,d) = \text{count}(\text{rel}=1, q, d) / \text{count}(q, d)$$

Problems

1. **Not enough data:** Most (q,d) pairs never observed
2. **Cannot adapt to new queries:** Need relevance judgments for every new query

Example: If we have 1 million documents and 1000 queries:

- Possible (q,d) pairs: 1 billion
- Actually judged: maybe 50,000
- 99.995% of pairs have no data!

2.3 Solution: Use Bayes' Rule

Transform the problem:

$$P(\text{rel}=1|q,d) = P(d|\text{rel}=1,q) \times P(\text{rel}=1) / P(d)$$

Using odds ratio (helps cancel terms):

$$O(\text{rel}=1|q,d) = P(\text{rel}=1|q,d) / P(\text{rel}=0|q,d)$$

Property of odds:

$$a/(1-a) = 1/(1-a) - 1$$

So if we can rank by odds, we can rank by probability!

Why odds? Because:

$$O(\text{rel}=1|q,d) = [P(d|\text{rel}=1,q) \times P(\text{rel}=1)] / [P(d|\text{rel}=0,q) \times P(\text{rel}=0)]$$

The $P(d)$ cancels out! Now we have:

$$O(\text{rel}=1|q,d) \propto P(d|\text{rel}=1,q) / P(d|\text{rel}=0,q)$$

3. Model 1: Robertson & Sparck Jones (RSJ)

3.1 Key Assumptions

Independence assumption: Words in documents are independent given relevance

$$P(d|\text{rel},q) = \prod_{i \in V} P(w_i|\text{rel},q)$$

Binary occurrence: Each word either appears ($w_i=1$) or doesn't ($w_i=0$)

3.2 Complete Derivation

Step 1: Start with odds

$$O(\text{rel}=1|q,d) = P(d|\text{rel}=1,q)P(\text{rel}=1) / [P(d|\text{rel}=0,q)P(\text{rel}=0)]$$

Step 2: Apply independence assumption

$$P(d|\text{rel}=1,q) = \prod_{i \in V} P(w_i|\text{rel}=1,q)$$

Step 3: Split into words in document vs. not in document

For each word, either $w_i=1$ (in doc) or $w_i=0$ (not in doc):

$$P(d|\text{rel}=1,q) = \prod_{w_i=1} P(w_i=1|\text{rel}=1,q) \times \prod_{w_i=0} P(w_i=0|\text{rel}=1,q)$$

Similarly:

$$P(d|\text{rel}=0,q) = \prod_{w_i=1} P(w_i=1|\text{rel}=0,q) \times \prod_{w_i=0} P(w_i=0|\text{rel}=0,q)$$

Step 4: Form the odds ratio

$$O(\text{rel}=1|q,d) = [\prod_{i=1}^n \{w_i=1\} P(w_i=1|\text{rel}=1,q) / P(w_i=1|\text{rel}=0,q)] \\ \times [\prod_{i=1}^n \{w_i=1\} P(w_i=0|\text{rel}=1,q) / P(w_i=0|\text{rel}=0,q)] \\ \times [\prod_{i=1}^n \{w_i=0\} P(w_i=1|\text{rel}=1,q) / P(w_i=1|\text{rel}=0,q)] \\ \times [\prod_{i=1}^n \{w_i=0\} P(w_i=0|\text{rel}=1,q) / P(w_i=0|\text{rel}=0,q)]$$

Step 5: Key simplification

The last two products (over $w_i=0$) don't depend on the document! They're the same for all documents, so for ranking purposes, we can ignore them.

$$O(\text{rel}=1|q,d) \propto \prod_{i=1}^n \{w_i=1\} [P(w_i=1|\text{rel}=1,q) / P(w_i=1|\text{rel}=0,q)] \\ \times [P(w_i=0|\text{rel}=1,q) / P(w_i=0|\text{rel}=0,q)]$$

Using $P(w_i=0|\text{rel},q) = 1 - P(w_i=1|\text{rel},q)$:

$$O(\text{rel}=1|q,d) \propto \prod_{i=1}^n \{w_i=1\} [P(w_i=1|\text{rel}=1,q) / P(w_i=1|\text{rel}=0,q)] \\ \times [(1 - P(w_i=1|\text{rel}=1,q)) / (1 - P(w_i=1|\text{rel}=0,q))]$$

Step 6: Take logarithm (preserves ranking)

$$\log O(\text{rel}=1|q,d) = \sum_{i=1}^n \{w_i=1\} \log[\alpha_i(1-\beta_i) / (\beta_i(1-\alpha_i))]$$

where:

- $\alpha_i = P(w_i=1|\text{rel}=1,q)$: probability word i appears in relevant documents
- $\beta_i = P(w_i=1|\text{rel}=0,q)$: probability word i appears in non-relevant documents

3.3 Final RSJ Formula

Ranking function

$$\text{score}^{\text{RSJ}}(q,d) = \sum_{i=1}^n \{w_i=1\} \log[\alpha_i(1-\beta_i) / (\beta_i(1-\alpha_i))]$$

Estimation with smoothing (avoid zeros)

$$\alpha_i = [\text{count}(w_i=1, q, \text{rel}=1) + 0.5] / [\text{count}(q, \text{rel}=1) + 1]$$

$$\beta_i = [\text{count}(w_i=1, q, \text{rel}=0) + 0.5] / [\text{count}(q, \text{rel}=0) + 1]$$

Intuition

- If $\alpha_i > \beta_i$: word appears more in relevant docs \rightarrow positive contribution
- If $\alpha_i < \beta_i$: word appears more in non-relevant docs \rightarrow negative contribution

- The log ratio amplifies these differences

3.4 Example Calculation

Scenario:

- Query: "machine learning"
- We have 100 relevant documents, 900 non-relevant documents
- "machine": appears in 80 relevant, 200 non-relevant
- "learning": appears in 70 relevant, 150 non-relevant

Calculate α_i and β_i

For "machine":

$$\alpha_{\text{machine}} = (80 + 0.5) / (100 + 1) = 80.5/101 \approx 0.797$$
$$\beta_{\text{machine}} = (200 + 0.5) / (900 + 1) = 200.5/901 \approx 0.222$$

For "learning":

$$\alpha_{\text{learning}} = (70 + 0.5) / (100 + 1) = 70.5/101 \approx 0.698$$
$$\beta_{\text{learning}} = (150 + 0.5) / (900 + 1) = 150.5/901 \approx 0.167$$

Score for document containing both words

$$\begin{aligned} \text{score} &= \log[0.797 \times (1 - 0.222) / (0.222 \times (1 - 0.797))] \\ &\quad + \log[0.698 \times (1 - 0.167) / (0.167 \times (1 - 0.698))] \\ &= \log[0.797 \times 0.778 / (0.222 \times 0.203)] \\ &\quad + \log[0.698 \times 0.833 / (0.167 \times 0.302)] \\ &= \log[0.620 / 0.045] + \log[0.581 / 0.050] \\ &= \log[13.78] + \log[11.62] \\ &= 2.62 + 2.45 \\ &= 5.07 \end{aligned}$$

3.5 Advantages and Limitations

Advantages

- Probabilistic foundation
- Can incorporate relevance feedback

- Theoretically justified

Limitations

- Requires relevance judgments
 - Binary occurrence only (doesn't use term frequency)
 - Independence assumption may not hold
-

4. Model 2: BM25 (Okapi BM25)

4.1 Motivation

Problem with RSJ: Only uses binary occurrence (word present/absent)

Example:

```
d1: "machine learning machine learning machine"  
d2: "machine"
```

RSJ treats both the same! But d1 clearly emphasizes "machine" more.

Goal: Incorporate term frequency (TF) into probabilistic model

4.2 The Eliteness Hypothesis

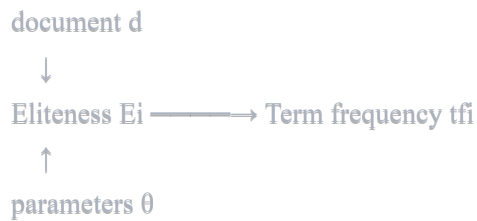
New hidden variable: Eliteness (E_i)

Definition: A document-term pair (d, w_i) is "elite" if the document is substantially about the concept denoted by the term.

Example:

- Query: "artificial intelligence"
- d1: "...artificial intelligence systems...learning...neural..."
 - "artificial": ELITE (document is about AI)
 - "intelligence": ELITE
 - "learning": ELITE
 - "the": NOT ELITE
- d2: "...studying the artificial flowers..."
 - "artificial": NOT ELITE (different meaning)

Graphical model:



Key insight: Term occurrence depends on eliteness

4.3 Two-Poisson Mixture Model

Model term frequency as mixture of two Poisson distributions:

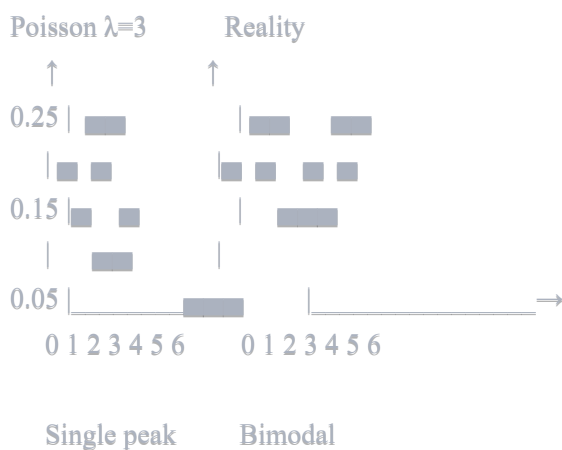
$$P(w_i=t_{fi}|q, rel=1) = \underbrace{\pi \times [\lambda^{t_{fi}} \times e^{(-\lambda)} / t_{fi}!]}_{\text{ELITE (high frequency)}} + \underbrace{(1-\pi) \times [\mu^{t_{fi}} \times e^{(-\mu)} / t_{fi}!]}_{\text{NON-ELITE (low frequency)}}$$

Parameters:

- π : probability document is elite for this term
- λ : rate parameter for elite distribution (large, e.g., 3-5)
- μ : rate parameter for non-elite distribution (small, e.g., 0.5)

Why two Poisson?

One Poisson can't capture "burstiness":



Intuition: Words either appear rarely (background) or appear frequently (topical)

4.4 Approximation: Saturation Function

Problem: Three unknown parameters (π , λ , μ) are hard to estimate

Solution: Design a parameter-free approximation!

The 2-Poisson mixture behaves like a saturation function:

$$c_i^{\text{elite}}(\text{tfi}) \approx \text{tfi} / (k1 + \text{tfi})$$

Why this works

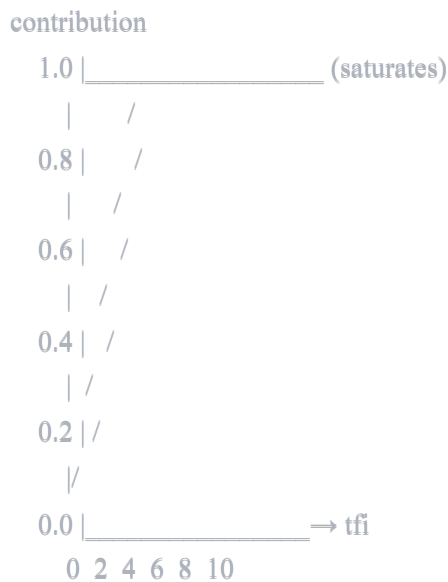
As $\text{tfi} \rightarrow 0$:

$$\text{tfi} / (k1 + \text{tfi}) \rightarrow 0$$

As $\text{tfi} \rightarrow \infty$:

$$\text{tfi} / (k1 + \text{tfi}) \rightarrow 1$$

Visualization:



Different $k1$ values:

$x/(0.2+x)$ — (fast saturation)

$x/(1+x)$ — (medium)

$x/(3+x)$ — (slow)

$x/(10+x)$ — (very slow)

Different curves for different $k1$:

$k1 = 0.5$: Very fast saturation (aggressive)

$k1 = 1.2$: Fast saturation (typical lower bound)

$k1 = 2.0$: Medium saturation (typical upper bound)

$k1 = 3.0$: Slow saturation

4.5 Complete BM25 Formula

Final scoring function:

$$\text{score}^{\text{BM25}}(q,d) = \sum_{w_i \in q} \log(N/\text{df}_i) \times [\text{tf}_i(k_1+1)] / [k_1((1-b)+b|d|/\text{avgdl}) + \text{tf}_i]$$

Component analysis

1. IDF term:

$$\log(N/\text{df}_i)$$

- N = total number of documents
- df_i = number of documents containing term i
- Exactly like TF-IDF!

2. Normalized TF term:

$$[\text{tf}_i(k_1+1)] / [k_1 \times \text{normalization} + \text{tf}_i]$$

3. Document length normalization (pivoting):

$$\text{normalization} = (1-b) + b \times (|d|/\text{avgdl})$$

- $|d|$ = document length
- avgdl = average document length in collection
- $b \in [0,1]$ controls strength (typically $b=0.75$)

4.6 Breaking Down the Formula

Let's see each component's effect

Scenario:

- $N = 1,000,000$ documents
- $\text{avgdl} = 1,000$ words
- $k_1 = 1.5$, $b = 0.75$

Term: "neural"

- $\text{df} = 10,000$ documents contain it
- $\text{IDF} = \log(1,000,000/10,000) = \log(100) \approx 4.6$

Document d_1 :

- Length: 500 words (short)
- "neural" appears 3 times

$$\text{normalization} = (1-0.75) + 0.75 \times (500/1000) = 0.25 + 0.375 = 0.625$$

$$\begin{aligned}\text{TF component} &= [3 \times (1.5 + 1)] / [1.5 \times 0.625 + 3] \\ &= [3 \times 2.5] / [0.9375 + 3] \\ &= 7.5 / 3.9375 \\ &= 1.905\end{aligned}$$

$$\text{score contribution} = 4.6 \times 1.905 = 8.76$$

Document d2:

- Length: 3000 words (long)
- "neural" appears 5 times

$$\text{normalization} = (1-0.75) + 0.75 \times (3000/1000) = 0.25 + 2.25 = 2.5$$

$$\begin{aligned}\text{TF component} &= [5 \times (1.5 + 1)] / [1.5 \times 2.5 + 5] \\ &= [5 \times 2.5] / [3.75 + 5] \\ &= 12.5 / 8.75 \\ &= 1.429\end{aligned}$$

$$\text{score contribution} = 4.6 \times 1.429 = 6.57$$

Notice: Despite d2 having higher TF (5 vs 3), d1 scores higher because it's shorter!

4.7 Parameter Effects

Effect of k1

k1 = 1.2: More aggressive saturation

$$\text{tf}=2 \rightarrow 2(2.2)/(1.2+2) = 1.375$$

$$\text{tf}=5 \rightarrow 5(2.2)/(1.2+5) = 1.774$$

k1 = 2.0: Less aggressive saturation

$$\text{tf}=2 \rightarrow 2(3)/(2+2) = 1.5$$

$$\text{tf}=5 \rightarrow 5(3)/(2+5) = 2.143$$

Effect of b

b = 0: No length normalization

Short doc (500 words): norm = 1.0

Long doc (3000 words): norm = 1.0

b = 1: Full length normalization

Short doc: norm = $500/1000 = 0.5$

Long doc: norm = $3000/1000 = 3.0$

b = 0.75: Typical setting (partial normalization)

Short doc: norm = $0.25 + 0.75 \times 0.5 = 0.625$

Long doc: norm = $0.25 + 0.75 \times 3.0 = 2.5$

4.8 BM25F: Multi-Field Extension

Motivation: Documents have structure (title, body, metadata)

Example: StackOverflow question

Title: "How to cite presentation slides?"

Body: "A friend has made some nice slides..."

Tags: ["citations", "academic"]

Different fields have different importance!

BM25F formula

$$\text{score}^{\text{BM25F}}(q,d) = \sum_{w_i \in q} \log(N/\text{df}_i) \times [\text{tf}_i^F(k_1+1)] / [k_1((1-b)+b|\text{d}|^F/\text{avgdl}^F) + \text{tf}_i^F]$$

Weighted combination:

$$\text{tf}_i^F = \sum_f \alpha_f \times \text{tf}_{i,f}$$

$$|\text{d}|^F = \sum_f \alpha_f \times |\text{d}|_f$$

$$\text{avgdl}^F = \sum_f \alpha_f \times \text{avgdl}_f$$

Example weights:

$\alpha_{\text{title}} = 3.0$ (title is very important)

$\alpha_{\text{body}} = 1.0$ (body is baseline)

$\alpha_{\text{tags}} = 2.0$ (tags are important)

Calculation

Document:

- Title (5 words): "cite presentation slides"
- Body (100 words): "...slides...cite..."

- Tags: ["citations"]

Query: "cite presentation"

$$\begin{aligned} \text{tf_cite}^F &= 3.0 \times 1 \text{ (title)} + 1.0 \times 1 \text{ (body)} + 2.0 \times 0 \text{ (tags)} = 4.0 \\ \text{tf_presentation}^F &= 3.0 \times 1 \text{ (title)} + 1.0 \times 0 \text{ (body)} + 2.0 \times 0 \text{ (tags)} = 3.0 \\ |d|^F &= 3.0 \times 5 + 1.0 \times 100 + 2.0 \times 2 = 119 \end{aligned}$$

4.9 Advantages and Limitations

Advantages

- State-of-the-art non-learning model
- Still widely used (Elasticsearch, Solr)
- Incorporates TF with principled normalization
- Few parameters to tune

Limitations

- Parameters (k1, b) hard to estimate theoretically
- No explicit relevance feedback mechanism (unlike RSJ)
- Still bag-of-words (no phrase information)

5. Model 3: Language Model-Based Retrieval

5.1 Key Paradigm Shift

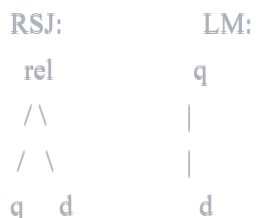
RSJ/BM25 approach:

$$P(\text{rel}=1|q,d) \text{ "Is document relevant given query?"}$$

Language model approach:

$$P(q|d) \text{ "Could this document generate this query?"}$$

Graphical models:



Intuition:

- Imagine each document is a "language model" (probability distribution over words)
- Rank documents by likelihood they would generate the query

5.2 Statistical Language Model

Definition: A probability distribution over word sequences

Examples:

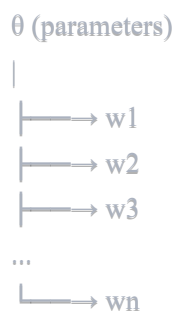
$P(\text{"Today is Wednesday"}) \approx 0.001$ (reasonable sentence)
 $P(\text{"Today Wednesday is"}) \approx 0.000000001$ (ungrammatical)
 $P(\text{"Eigenvalue positive is"}) \approx 0.0000001$ (weird)

Unigram language model:

- Generate each word independently
- $P(w_1 w_2 \dots w_n) = P(w_1) \times P(w_2) \times \dots \times P(w_n)$

Parameters: $\{P(w_1), P(w_2), \dots, P(w_N)\}$ where $\sum P(w_i) = 1$

Graphical model:



Example

Language model:

$P(\text{"the"}) = 0.05$
 $P(\text{"is"}) = 0.02$
 $P(\text{"cat"}) = 0.001$
 $P(\text{"eigenvalue"}) = 0.00001$
...

Generate sentence:

$P(\text{"the cat is"}) = P(\text{"the"}) \times P(\text{"cat"}) \times P(\text{"is"})$
 $= 0.05 \times 0.001 \times 0.02$
 $= 0.000001$

5.3 Query Likelihood Model

Ranking function:

$$\text{score}(q,d) = P(q|d) = \prod_{w_i \in q} P(w_i|d)$$

Taking log (preserves ranking, easier to compute):

$$\text{score}(q,d) = \log P(q|d) = \sum_{w_i \in q} \log P(w_i|d)$$

Example

Query: "presidential campaign news"

Document d: "...news of presidential campaign..... ...presidential candidate..."

Estimate $P(w_i|d)$ by MLE:

$$P_{\text{MLE}}(w_i|d) = \text{count}(w_i, d) / |d|$$

$$P(\text{"presidential"}|d) = 2/100 = 0.02$$

$$P(\text{"campaign"}|d) = 1/100 = 0.01$$

$$P(\text{"news"}|d) = 1/100 = 0.01$$

$$\text{score}(q,d) = \log(0.02) + \log(0.01) + \log(0.01)$$

$$= -3.91 + (-4.61) + (-4.61)$$

$$= -13.13$$

Problem: What if query word doesn't appear in document?

$$P(\text{"election"}|d) = 0/100 = 0$$

$$\log(0) = -\infty$$

$$\text{Entire score} = -\infty$$

This is catastrophic! We need smoothing.

5.4 Smoothing Methods

Core idea: Mix document language model with collection language model

Collection language model:

$$P(w_i|C) = \sum_{d \in C} \text{count}(w_i, d) / \sum_{d \in C} |d|$$

Method 1: Dirichlet Smoothing

Formula:

$$P_s(w_i|d) = [\text{count}(w_i, d) + \mu \times P(w_i|C)] / [|d| + \mu]$$

Rewrite as weighted combination:

$$P_s(w_i|d) = \underbrace{[|d|/(|d|+\mu)]}_{\text{weight on document}} \times P_{\text{MLE}}(w_i|d) + \underbrace{[\mu/(|d|+\mu)]}_{\text{weight on collection}} \times P(w_i|C)$$

Parameter μ : Controls smoothing strength

- μ small (e.g., 100): Trust document more
- μ large (e.g., 5000): Trust collection more
- Typical value: $\mu \approx 2000$

Intuition:

Document only
↓
Long doc ($|d|=5000$): $5000/7000 = 71\%$ document, 29% collection
Short doc ($|d|=100$): $100/2100 = 5\%$ document, 95% collection
↑
More smoothing needed

Example calculation:

Document d (100 words):

- "neural" appears 3 times
- $P(\text{"neural"}|C) = 0.001$ (collection)
- $\mu = 2000$

$$\begin{aligned} P_s(\text{"neural"}|d) &= [3 + 2000 \times 0.001] / [100 + 2000] \\ &= [3 + 2] / 2100 \\ &= 5/2100 \\ &= 0.00238 \end{aligned}$$

For unseen word "quantum":

$$\begin{aligned} P_s(\text{"quantum"}|d) &= [0 + 2000 \times 0.0001] / [100 + 2000] \\ &= 0.2 / 2100 \\ &= 0.000095 \end{aligned}$$

Still non-zero!

Method 2: Jelinek-Mercer Smoothing

Formula:

$$P_s(w_i|d) = \lambda \times P_{MLE}(w_i|d) + (1-\lambda) \times P(w_i|C)$$

Fixed interpolation weight:

- $\lambda \in [0, 1]$
- Typical: $\lambda \approx 0.7$

Difference from Dirichlet:

- Dirichlet: Adapts to document length
- JM: Same λ for all documents

Example:

$\lambda = 0.7$, same document as before:

$$\begin{aligned} P_s(\text{"neural"}|d) &= 0.7 \times (3/100) + 0.3 \times 0.001 \\ &= 0.7 \times 0.03 + 0.0003 \\ &= 0.021 + 0.0003 \\ &= 0.0213 \end{aligned}$$

For unseen "quantum":

$$\begin{aligned} P_s(\text{"quantum"}|d) &= 0.7 \times 0 + 0.3 \times 0.0001 \\ &= 0.00003 \end{aligned}$$

5.5 Ranking Formula (Dirichlet)

Starting from:

$$\log P(q|d) = \sum_{w_i \in q} \log P_s(w_i|d)$$

Substitute Dirichlet smoothing:

$$= \sum_{w_i \in q} \log \left(\frac{\text{count}(w_i, d) + \mu \times P(w_i|C)}{|d| + \mu} \right)$$

Simplify (see derivation link: <https://www.overleaf.com/read/jbztxtnbzwcx>):

$$= \underbrace{\sum_{w_i \in q, w_i \in d} \log(\text{count}(w_i, d) / [\mu \times P(w_i | C)])}_{\text{depends on matching}} + \underbrace{|q| \times \log[\mu / (\mu + |d|)]}_{\text{document length penalty}}$$

For ranking (ignore constant $|q|$):

$$\text{score}^{\text{Dir}}(q, d) = \sum_{w_i \in q, w_i \in d} \log(1 + \text{count}(w_i, d) / (\mu \times P(w_i | C))) + \log(\mu / (\mu + |d|))$$

Intuition:

1. For each query term in document, get score based on its frequency
2. Longer documents get penalized by log term

5.6 Feedback Language Model

Problem: Query is often too short and ambiguous

Example: Query: "apple"

- Could mean: fruit, company, music, records, pie, tree...

Solution: Use pseudo-relevance feedback

Pseudo-relevance feedback

1. Retrieve top-k documents with initial query
2. Assume these are relevant (even without judgments)
3. Expand query using these documents

Model: Feedback creates improved query model

$$\theta_q^{\text{FB}} = \lambda \times \theta_q + (1 - \lambda) \times \theta_q^{\text{F}}$$

where:

- θ_q = original query model
- θ_q^{F} = model learned from feedback documents
- $\lambda \in [0, 1]$ = interpolation weight

How to estimate θ_q^{F} ?

EM Algorithm:

Assume feedback documents are mixture:

$$P(d | \theta_q^{\text{F}}) = \prod_{w \in d} [\theta_q^{\text{F}}(w)]^{\text{count}(w, d)}$$

E-step: Estimate which words came from topic vs. background

$$P(z=\text{topic}|w,d) = [\theta q^F(w) \times \lambda] / [\theta q^F(w) \times \lambda + P(w|C) \times (1-\lambda)]$$

M-step: Update θq^F

$$\theta q^F(w) \propto \sum_{d \in \text{feedback}} \text{count}(w,d) \times P(z=\text{topic}|w,d)$$

Iterate until convergence

5.7 Complete Example

Query: "airport security"

Step 1: Initial retrieval using query likelihood

Retrieve top 5 documents

Step 2: Extract top terms from feedback documents

Document 1: "...airport security measures... passengers... screening..."

Document 2: "...TSA procedures...airport... safety..."

Document 3: "...security checkpoints... baggage... inspection..."

Step 3: Learn θq^F using EM

Discovered related terms:

$$P(\text{"passengers"}|\theta q^F) = 0.08$$

$$P(\text{"screening"}|\theta q^F) = 0.06$$

$$P(\text{"TSA"}|\theta q^F) = 0.05$$

$$P(\text{"checkpoints"}|\theta q^F) = 0.04$$

...

Step 4: Combine with original query

$$\theta q^{FB}(\text{"airport"}) = 0.7 \times 0.5 + 0.3 \times 0.03 = 0.359$$

$$\theta q^{FB}(\text{"security"}) = 0.7 \times 0.5 + 0.3 \times 0.02 = 0.356$$

$$\theta q^{FB}(\text{"passengers"}) = 0.7 \times 0 + 0.3 \times 0.08 = 0.024 \leftarrow \text{NEW}$$

$$\theta q^{FB}(\text{"screening"}) = 0.7 \times 0 + 0.3 \times 0.06 = 0.018 \leftarrow \text{NEW}$$

Step 5: Rerank using expanded query

5.8 Advantages and Limitations

Advantages

- Flexible framework

- Can easily incorporate feedback
- Probabilistic foundation
- State-of-the-art performance with good smoothing

Limitations

- Query-document equivalence assumption unrealistic
- Unigram model (no phrases)
- Smoothing parameters must be tuned

5.9 Comparison: BM25 vs Language Models

Empirical results (Bennett et al. 2008):

On some datasets:

- BM25 outperforms LM (e.g., TREC TD)
- LM with feedback outperforms baseline LM

Typical performance ranking:

BM25 \approx LM-Dirichlet > LM-JM > RSJ

When to use which:

- **BM25:** Simple, robust, well-tested
- **LM:** Want feedback, need theoretical framework
- **BM25F:** Multi-field documents
- **Feedback LM:** Query expansion is critical

6. Practice Problems with Solutions

Problem 1: MLE for Biased Coin

Question: You flip a coin 20 times:

- Heads: 13 times
- Tails: 7 times

What's the MLE for $P(\text{Heads})$?

Solution:

$$\begin{aligned}\alpha &= \text{count}(\text{heads}) / [\text{count}(\text{heads}) + \text{count}(\text{tails})] \\ &= 13 / (13 + 7) \\ &= 13/20 \\ &= 0.65\end{aligned}$$

Problem 2: Bayes Rule Application

Question:

- $P(\text{disease}) = 0.01$ (1% of population has disease)
- $P(\text{positive test}|\text{disease}) = 0.95$ (95% sensitivity)
- $P(\text{positive test}|\text{no disease}) = 0.05$ (5% false positive)

If someone tests positive, what's $P(\text{disease}|\text{positive test})$?

Solution:

Use Bayes:

$$P(\text{disease}|\text{pos}) = P(\text{pos}|\text{disease}) \times P(\text{disease}) / P(\text{pos})$$

Calculate $P(\text{pos})$:

$$\begin{aligned}P(\text{pos}) &= P(\text{pos}|\text{disease}) \times P(\text{disease}) + P(\text{pos}|\text{no disease}) \times P(\text{no disease}) \\ &= 0.95 \times 0.01 + 0.05 \times 0.99 \\ &= 0.0095 + 0.0495 \\ &= 0.059\end{aligned}$$

Therefore:

$$\begin{aligned}P(\text{disease}|\text{pos}) &= (0.95 \times 0.01) / 0.059 \\ &= 0.0095 / 0.059 \\ &\approx 0.161 \text{ or } 16.1\%\end{aligned}$$

Surprising! Despite positive test, only 16% chance of disease (because disease is rare).

Problem 3: RSJ Scoring

Given:

- Query: "machine learning"
- 50 relevant documents, 950 non-relevant
- "machine": 40 relevant, 100 non-relevant

- "learning": 35 relevant, 80 non-relevant

Calculate RSJ score for document containing both terms.

Solution:

Calculate α and β with smoothing:

$$\alpha_{\text{machine}} = (40 + 0.5) / (50 + 1) = 40.5/51 \approx 0.794$$

$$\beta_{\text{machine}} = (100 + 0.5) / (950 + 1) = 100.5/951 \approx 0.106$$

$$\alpha_{\text{learning}} = (35 + 0.5) / (50 + 1) = 35.5/51 \approx 0.696$$

$$\beta_{\text{learning}} = (80 + 0.5) / (950 + 1) = 80.5/951 \approx 0.085$$

Calculate score:

$$\text{score} = \log[0.794 \times (1 - 0.106) / (0.106 \times (1 - 0.794))] + \log[0.696 \times (1 - 0.085) / (0.085 \times (1 - 0.696))]$$

$$= \log[0.794 \times 0.894 / (0.106 \times 0.206)] + \log[0.696 \times 0.915 / (0.085 \times 0.304)]$$

$$= \log[0.710 / 0.022] + \log[0.637 / 0.026]$$

$$= \log[32.27] + \log[24.50]$$

$$= \log[32.27] + \log[24.50]$$

$$= 3.47 + 3.20$$

$$= 6.67$$

$$= 6.67$$

Problem 4: BM25 Calculation

Given:

- $N = 10,000$ documents
- $\text{avgdl} = 500$ words
- $k1 = 1.5$, $b = 0.75$
- Term "neural": $\text{df} = 100$
- Document: $\text{length} = 300$, "neural" appears 4 times

Calculate BM25 score contribution.

Solution:

IDF:

$$\text{IDF} = \log(N/\text{df}) = \log(10,000/100) = \log(100) \approx 4.605$$

Document length normalization:

$$\begin{aligned}\text{norm} &= (1-b) + b \times (|d|/\text{avgdl}) \\ &= (1-0.75) + 0.75 \times (300/500) \\ &= 0.25 + 0.75 \times 0.6 \\ &= 0.25 + 0.45 \\ &= 0.70\end{aligned}$$

TF component:

$$\begin{aligned}\text{TF_score} &= [4 \times (1.5 + 1)] / [1.5 \times 0.70 + 4] \\ &= [4 \times 2.5] / [1.05 + 4] \\ &= 10 / 5.05 \\ &= 1.980\end{aligned}$$

Total:

$$\begin{aligned}\text{score} &= \text{IDF} \times \text{TF_score} \\ &= 4.605 \times 1.980 \\ &= 9.118\end{aligned}$$

Problem 5: Language Model Smoothing

Given:

- Document: "neural networks deep learning neural" (5 words)
- Query: "neural quantum"
- $P(\text{"neural"}|C) = 0.002$
- $P(\text{"quantum"}|C) = 0.0001$
- $\mu = 1000$

Calculate score using Dirichlet smoothing.

Solution:

For "neural" (appears 2 times):

$$\begin{aligned}
 P_s(\text{"neural"}|d) &= [2 + 1000 \times 0.002] / [5 + 1000] \\
 &= [2 + 2] / 1005 \\
 &= 4/1005 \\
 &= 0.00398
 \end{aligned}$$

For "quantum" (appears 0 times):

$$\begin{aligned}
 P_s(\text{"quantum"}|d) &= [0 + 1000 \times 0.0001] / [5 + 1000] \\
 &= 0.1 / 1005 \\
 &= 0.0000995
 \end{aligned}$$

Score:

$$\begin{aligned}
 \text{score} &= \log(0.00398) + \log(0.0000995) \\
 &= -5.53 + (-9.22) \\
 &= -14.75
 \end{aligned}$$

Problem 6: Feedback Language Model

Given:

- Original query model: $P(\text{"apple"}|\theta_q) = 0.5$, $P(\text{"pie"}|\theta_q) = 0.5$
- Feedback documents yield: $P(\text{"apple"}|\theta_q^F) = 0.3$, $P(\text{"pie"}|\theta_q^F) = 0.2$, $P(\text{"recipe"}|\theta_q^F) = 0.15$
- $\lambda = 0.7$

Calculate θ_q^{FB} .

Solution:

$$\begin{aligned}
 P(\text{"apple"}|\theta_q^{FB}) &= 0.7 \times 0.5 + 0.3 \times 0.3 = 0.35 + 0.09 = 0.44 \\
 P(\text{"pie"}|\theta_q^{FB}) &= 0.7 \times 0.5 + 0.3 \times 0.2 = 0.35 + 0.06 = 0.41 \\
 P(\text{"recipe"}|\theta_q^{FB}) &= 0.7 \times 0 + 0.3 \times 0.15 = 0 + 0.045 = 0.045 \leftarrow \text{NEW TERM!}
 \end{aligned}$$

Notice: "recipe" wasn't in original query but now has probability 0.045

7. Quiz Questions & Detailed Solutions

Quiz Question 1 (From Lecture 2 Quiz)

Question: Suppose we have one query and two documents:

- $q = \text{"covid 19"}$

- d1 = "covid patient"
- d2 = "19 99 car wash"
- d3 = "19 street covid testing facility is reopened next week"

If we use the vector space model **without IDF**, which document has the highest and the lowest score?

Options:

- d1, d2
- d1, d3
- d3, d1
- d3, d2

Solution:

Without IDF, we're just counting term matches (TF only).

Query terms: {"covid", "19"}

d1: "covid patient"

- Contains: "covid" ✓
- Contains: "19" ✗
- **Matches: 1/2 query terms**

d2: "19 99 car wash"

- Contains: "covid" ✗
- Contains: "19" ✓
- **Matches: 1/2 query terms**

d3: "19 street covid testing facility is reopened next week"

- Contains: "covid" ✓
- Contains: "19" ✓
- **Matches: 2/2 query terms**

Scores (proportional to matches):

score(q, d1) \propto 1
 score(q, d2) \propto 1
 score(q, d3) \propto 2

Ranking:

$d3 \text{ (highest)} > d1 = d2 \text{ (tied for lowest)}$

Answer: d3 has highest, d1 and d2 tied for lowest

From quiz solutions: d1, d2

Quiz Question 2 (From Lecture 2 Quiz)

Question: What about using vector space model with IDF?

Solution:

Now we include **IDF:** $IDF(w) = \log(N/df(w))$

Document frequencies:

- "covid": appears in d1, d3 $\rightarrow df = 2$
- "19": appears in d2, d3 $\rightarrow df = 2$
- $N = 3$ (total documents)

IDF values:

$$IDF(\text{"covid"}) = \log(3/2) = \log(1.5) \approx 0.405$$

$$IDF(\text{"19"}) = \log(3/2) = \log(1.5) \approx 0.405$$

Scores:

d1: "covid patient"

$$\begin{aligned} \text{score} &= TF(\text{"covid"}) \times IDF(\text{"covid"}) \\ &= 1 \times 0.405 \\ &= 0.405 \end{aligned}$$

d2: "19 99 car wash"

$$\begin{aligned} \text{score} &= TF(\text{"19"}) \times IDF(\text{"19"}) \\ &= 1 \times 0.405 \\ &= 0.405 \end{aligned}$$

d3: "19 street covid testing facility is reopened next week"

$$\begin{aligned} \text{score} &= TF(\text{"covid"}) \times IDF(\text{"covid"}) + TF(\text{"19"}) \times IDF(\text{"19"}) \\ &= 1 \times 0.405 + 1 \times 0.405 \\ &= 0.810 \end{aligned}$$

Results:

- d3: 0.810 (highest)
- d1: 0.405 (tied lowest)
- d2: 0.405 (tied lowest)

Ranking:

$$d3 > d1 = d2$$

From quiz solutions: d1, d2

Quiz Question 3 (Feedback Language Model)

Question: Go to this notebook:

<https://drive.google.com/file/d/1APsRt8T41kLxNchQv2uJ3Ws5aZN0aI4Hk/view?usp=sharing>

Suppose query is "airport security" and $\lambda = 0.7$.

Use this notebook to answer the following question: what is the probability of "airport" using the language model-based retrieval model? (4 decimal places)

Given:

$$\theta^{FB} = \lambda \theta_q + (1-\lambda) \theta_d$$

Solution:

Using the feedback documents d1, d2, ..., calculate:

$$P(\text{"airport"}|\theta^{FB}) = 0.7 \times P(\text{"airport"}|\text{query}) + 0.3 \times P(\text{"airport"}|\text{feedback_docs})$$

Answer from quiz solutions: 0.38

Quiz Question 4 (Lambda Effect)

Question: How does increasing the value of λ change the probability of "all" in θ^{FB} ?

How does adding stop words to the query, e.g., "airport security" \rightarrow "the airport security" change the probability of "all" in θ^{FB} ?

Solution:

Part 1: Increasing λ

Formula:

$$\theta^{FB} = \lambda \theta_q + (1-\lambda) \theta_d$$

- Increasing $\lambda \rightarrow$ more weight on original query
- "all" is likely a stop word (appears in corpus but not query)
- Higher $\lambda \rightarrow$ **decreases** contribution from documents
- Higher $\lambda \rightarrow$ **decreases** "all" probability

Part 2: Adding stop words to query

If we add "the" to query:

- "the" will have higher probability in θ_q
- But "all" is not in the query
- The probability of "all" remains **the same**

Answer from quiz solutions: decrease, same

8. Video Resources

Probability & Statistics Foundations

1. Bayes' Theorem (3Blue1Brown)

- Title: "Bayes theorem, the geometry of changing beliefs"
- URL: <https://www.youtube.com/watch?v=HZGCoVF3YvM>
- **Why:** Best visual intuition for Bayes' rule
- **Duration:** 15 minutes

2. Maximum Likelihood Estimation (StatQuest)

- Title: "Maximum Likelihood, clearly explained!!!"
- URL: <https://www.youtube.com/watch?v=XepXtl9YKwc>
- **Why:** Simple explanation with examples
- **Duration:** 6 minutes

3. Probability Distributions (Khan Academy)

- URL: <https://www.khanacademy.org/math/statistics-probability>
- **Why:** Covers Bernoulli, Binomial, Poisson
- **Duration:** Series of short videos

Information Retrieval Specific

4. Stanford CS276 Lectures

- **Probabilistic IR:** Lecture 11
- **Language Models:** Lecture 12
- URL: <http://web.stanford.edu/class/cs276/>
- **Why:** Detailed derivations from experts
- **Duration:** 75 minutes each

5. BM25 Explanation (Elastic)

- Title: "BM25 The Next Generation of Lucene Relevance"
- URL: <https://www.elastic.co/blog/practical-bm25-part-2-the-bm25-algorithm-and-its-variables>
- **Why:** Practical implementation perspective
- **Duration:** Reading (15 min)

Mathematical Concepts

6. Logarithms Review (Khan Academy)

- **Why:** Essential for understanding log-odds, IDF
- URL: <https://www.khanacademy.org/math/algebra2/x2ec2f6f830c9fb89:logs>
- **Duration:** 30 minutes

7. Poisson Distribution (StatQuest)

- Title: "The Poisson Distribution, Clearly Explained!!!"
- URL: https://www.youtube.com/watch?v=cPOChr_kuQs
- **Why:** Understand the foundation of 2-Poisson model
- **Duration:** 11 minutes

8. EM Algorithm (Normalized Nerd)

- Title: "Expectation Maximization: how it works"
 - URL: <https://www.youtube.com/watch?v=iQoXFmbXRJA>
 - **Why:** Understand feedback language model estimation
 - **Duration:** 15 minutes
-

Quick Reference Guide

Key Formulas

RSJ:

$$\text{score} = \sum_{_} \{w_i=1\} \log[\alpha_i(1-\beta_i) / (\beta_i(1-\alpha_i))]$$

where:

$$\alpha_i = [\text{count}(w_i=1, q, \text{rel}=1) + 0.5] / [\text{count}(q, \text{rel}=1) + 1]$$

$$\beta_i = [\text{count}(w_i=1, q, \text{rel}=0) + 0.5] / [\text{count}(q, \text{rel}=0) + 1]$$

BM25:

$$\text{score} = \sum_{_} \{w_i \in q\} \log(N/\text{df}_i) \times [\text{tf}_i(k_1+1)] / [k_1((1-b)+b|d|/\text{avgdl}) + \text{tf}_i]$$

Typical parameters:

$$k_1 = 1.5$$

$$b = 0.75$$

Language Model (Dirichlet):

$$\text{score} = \sum_{_} \{w_i \in q, w_i \in d\} \log(1 + \text{count}(w_i, d) / [\mu \times P(w_i|C)]) + \log[\mu / (\mu + |d|)]$$

Smoothing:

$$P_s(w_i|d) = [\text{count}(w_i, d) + \mu \times P(w_i|C)] / [|d| + \mu]$$

Typical parameter:

$$\mu = 2000$$

Feedback LM:

$$\theta_q^{\text{FB}} = \lambda \theta_q + (1-\lambda) \theta_q^{\text{F}}$$

Typical parameter:

$$\lambda = 0.7$$

Parameter Typical Values

BM25:

$$k_1 \in [1.2, 2.0] \quad (\text{typical: } 1.5)$$

$$b \in [0, 1] \quad (\text{typical: } 0.75)$$

Language Models:

$\mu \in [1000, 5000]$ (typical: 2000) - Dirichlet
 $\lambda \in [0.6, 0.9]$ (typical: 0.7) - JM & Feedback

Model Comparison

Feature	RSJ	BM25	LM
Term Frequency	Binary	Saturated	Linear with smoothing
Doc Length Norm	No	Yes (pivoting)	Yes (implicit in smoothing)
IDF	Implicit in α, β	Explicit $\log(N/df)$	Implicit in $P(w C)$
Feedback	Easy	Hard	Easy
Parameters	Hard to estimate	2 parameters	1-2 parameters
Performance	Baseline	SOTA	SOTA with feedback

When to Use Each Model

Use RSJ when:

- You have relevance judgments
- You want theoretical foundation
- Binary occurrence is sufficient

Use BM25 when:

- You need state-of-the-art non-learning model
- You want robust, well-tested performance
- You have multi-field documents (BM25F)

Use Language Models when:

- You need query expansion/feedback
- You want flexible probabilistic framework
- You're willing to tune smoothing parameters

Study Tips for Midterm

What to Memorize

1. Bayes' rule formula

$$P(A|B) = P(B|A)P(A) / P(B)$$

2. MLE formula

$$\alpha = \text{count}(\text{heads}) / \text{total}$$

3. Complete BM25 formula

$$\text{score} = \sum \log(N/df) \times [\text{tf}(k_1+1)] / [k_1((1-b)+b|d|/\text{avgdl}) + \text{tf}]$$

4. Dirichlet smoothing formula

$$P_s(w|d) = [\text{count}(w,d) + \mu P(w|C)] / [|d| + \mu]$$

5. What α_i and β_i represent in RSJ

- $\alpha_i = P(w_i=1|\text{rel}=1,q)$
- $\beta_i = P(w_i=1|\text{rel}=0,q)$

What to Understand Conceptually

1. Why we use odds ratio in RSJ derivation

- Cancels out $P(d)$ term
- Easier to work with products

2. Why 2-Poisson models term burstiness

- One Poisson can't capture bimodal distribution
- Words appear rarely (background) or frequently (topic)

3. Why document length normalization is needed

- Long documents accumulate higher scores
- Need to balance between no penalty and full penalty

4. Why language models need smoothing

- Zero probability for unseen words
- Would make entire query probability zero

5. How feedback improves retrieval

- Expands query with related terms
- Learns from pseudo-relevant documents

Practice Skills

1. Calculate RSJ score by hand (given α , β)
2. Calculate BM25 score by hand (given all parameters)
3. Apply Dirichlet smoothing (given document and collection)

4. **Explain when to use each model**

5. **Derive formulas from first principles**

Common Mistakes to Avoid

1. ❌ Confusing $P(q|d)$ vs $P(d|q)$
2. ❌ Forgetting to add smoothing constant (0.5 in RSJ)
3. ❌ Wrong order of operations in BM25
4. ❌ Forgetting log when working with language models
5. ❌ Mixing up λ (feedback weight) and μ (smoothing parameter)

Step-by-Step Problem Solving

For RSJ problems:

1. Calculate α for each term (with +0.5 smoothing)
2. Calculate β for each term (with +0.5 smoothing)
3. For each term in document: $\log[\alpha(1-\beta)/(\beta(1-\alpha))]$
4. Sum all terms

For BM25 problems:

1. Calculate IDF: $\log(N/df)$
2. Calculate normalization: $(1-b) + b(|d|/avgdl)$
3. Calculate TF component: $[tf(k+1)] / [k_1 \times norm + tf]$
4. Multiply: $IDF \times TF$ component
5. Sum for all query terms

For Language Model problems:

1. For each query term, check if in document
2. If yes: use $[count + \mu P(w|C)] / [|d| + \mu]$
3. If no: use $[\mu P(w|C)] / [|d| + \mu]$
4. Take log of each probability
5. Sum all logs

Exam Strategy

1. **Read the entire question first**
 - Identify what model is being used

- Note what parameters are given

2. Write down relevant formulas

- Don't try to remember while calculating
- Write them at the top of your work

3. Show your work

- Partial credit is real
- Label intermediate steps

4. Check units and ranges

- Probabilities must be $[0,1]$
- Log of probabilities is negative
- Scores can be any real number

5. Manage your time

- Don't get stuck on one problem
- Come back to difficult ones

Additional Resources

Official Course Materials

- **Stanford IR Book:** <https://nlp.stanford.edu/IR-book/>
 - Chapter 11: Probabilistic Information Retrieval
 - Chapter 12: Language Models
- **Overleaf Derivation:** <https://www.overleaf.com/read/jbztxtnbzwcx>
 - Complete mathematical derivations

Research Papers (Optional but Helpful)

1. RSJ Model:

- Robertson & Sparck Jones (1976)
- "Relevance weighting of search terms"

2. BM25:

- Robertson et al. (1994)
- "Okapi at TREC-3"

3. Language Models:

- Ponte & Croft (1998)

- "A language modeling approach to information retrieval"

4. Feedback Language Models:

- Zhai & Lafferty (2001)
- "Model-based feedback in the language modeling approach to IR"

Practice Problems Online

- Stanford CS276 Homeworks
 - TREC evaluation datasets
 - Kaggle information retrieval competitions
-

Summary Checklist

Before the midterm, make sure you can:

- ☐ Explain Bayes' rule and apply it
 - ☐ Calculate MLE from observations
 - ☐ State the Probability Ranking Principle
 - ☐ Derive RSJ model from first principles
 - ☐ Calculate RSJ score by hand
 - ☐ Explain eliteness hypothesis
 - ☐ Understand 2-Poisson mixture model
 - ☐ Calculate BM25 score by hand
 - ☐ Explain document length normalization
 - ☐ Understand saturation function
 - ☐ State the paradigm shift in LM approach
 - ☐ Explain why smoothing is necessary
 - ☐ Calculate Dirichlet smoothing by hand
 - ☐ Understand feedback language model
 - ☐ Compare RSJ, BM25, and LM models
 - ☐ Know when to use each model
-

Good luck on your midterm!

Last updated: [Current Date]

For questions or corrections, contact: [Your Email]