

Extended Probabilistic Models

Kuan-Yu Chen (陳冠宇)

2020/10/23 @ TR-313, NTUST

Homework 1 – VSM

- In this project, we have

- 50 Queries

- 4191 Documents

- Our goal is to implement a vector space model

$$\text{sim}(q, d_j) = \cos(\theta) = \frac{\vec{q} \cdot \vec{d}_j}{|\vec{q}| |\vec{d}_j|} = \frac{\sum_{w_i \in V} w_{i,q} \times w_{i,j}}{\sqrt{\sum_{w_i \in V} w_{i,q}^2} \times \sqrt{\sum_{w_i \in V} w_{i,j}^2}}$$

47 teams · 7 days to go

Overview Data Notebooks Discussion **Leaderboard** Rules Team Host My Submissions [Submit Predictions](#)

Public Leaderboard Private Leaderboard

This leaderboard is calculated with all of the test data. [Raw Data](#) [Refresh](#)

#	Team Name	Notebook	Team Members	Score ?	Entries	Last
1	__incognito__			0.77476	24	10h
2	M10915201_陳牧凡			0.75664	11	4d
3	M10815048_張晏銘			0.74631	14	3d
4	M10915045_施信宏			0.74525	4	1d
5	Happy Baseline			0.74386	39	3d
6	M10915036_王繹威			0.73619	5	1d
7	M10915012_黃偉愷			0.72341	47	8h
8	M10909118_蔡旭真			0.72026	2	8h
9	B10615034_黃柏翰			0.71634	27	10h
10	80847002S_羅天宏			0.71302	31	4d

Review

- Boolean Model

- Probabilistic Model

- Binary Independence Model

$$sim(d_j, q) \propto \sum_{w_i \in d_j \& w_i \in q} \log \frac{P(w_i | R_q)}{1 - P(w_i | R_q)} + \log \frac{1 - P(w_i | \bar{R}_q)}{P(w_i | \bar{R}_q)}$$

- Robertson-Sparck Jones Equation

$$sim(d_j, q) \equiv \sum_{w_i \in d_j \& w_i \in q} \log \left(\frac{r_i + 0.5}{R_q - r_i + 0.5} \cdot \frac{N - R_q - n_i + r_i + 0.5}{n_i - r_i + 0.5} \right)$$

- TF-IDF

- Term Frequency

- Inverse Document Frequency

- Overlap Score Model $sim(q, d_j) = \sum_{w_i \in q} TF - IDF_{i,j}$

- Vector Space Model $sim(q, d_j) = \cos(\theta) = \frac{\vec{q} \cdot \vec{d}_j}{|\vec{q}| |\vec{d}_j|} = \frac{\sum_{w_i \in V} k_{i,q} \times k_{i,j}}{\sqrt{\sum_{w_i \in V} k_{i,q}^2} \times \sqrt{\sum_{w_i \in V} k_{i,j}^2}}$

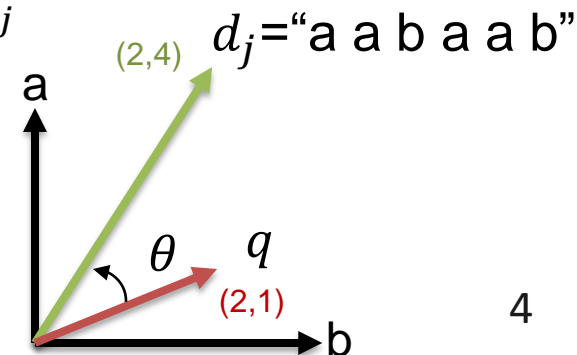
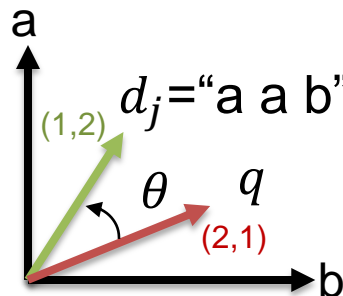
Document Length Normalization

- Longer documents can broadly be lumped into two categories
 - Verbose documents that essentially **repeat the same content**
 - The longer the document, the more the information?
 - Documents **covering multiple different topics**
 - The term frequency cannot really reveal the document

$$sim_{Probabilistic}(d_j, q) \equiv \sum_{w_i \in d_j \& w_i \in q} \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

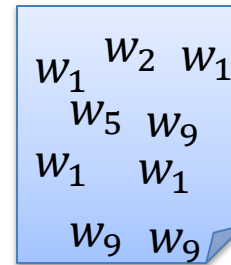
$$sim_{VSM}(q, d_j) \equiv \cos(\theta) = \frac{\vec{q} \cdot \vec{d_j}}{|\vec{q}| |\vec{d_j}|} = \frac{\sum_{w_i \in V} w_{i,q} \times w_{i,j}}{\sqrt{\sum_{w_i \in V} w_{i,q}^2} \times \sqrt{\sum_{w_i \in V} w_{i,j}^2}}$$

Scheme	Document Term Weight	Query Term Weight
1	$tf_{i,j} \times \log \frac{N}{n_i}$	$\left(0.5 + 0.5 \frac{tf_{i,q}}{\max_i (tf_{i,q})} \right) \times \log \frac{N}{n_i}$
2	$1 + tf_{i,j}$	$\log \left(1 + \frac{N}{n_i} \right)$
3	$(1 + tf_{i,j}) \times \log \frac{N}{n_i}$	$(1 + tf_{i,q}) \times \log \frac{N}{n_i}$

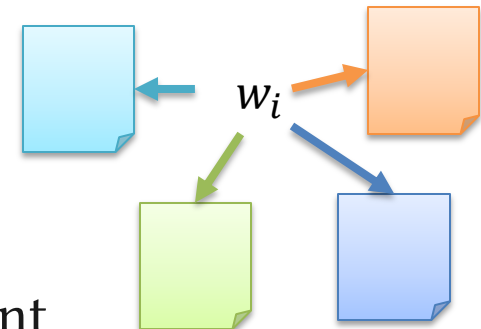


TF, IDF, and Document Length

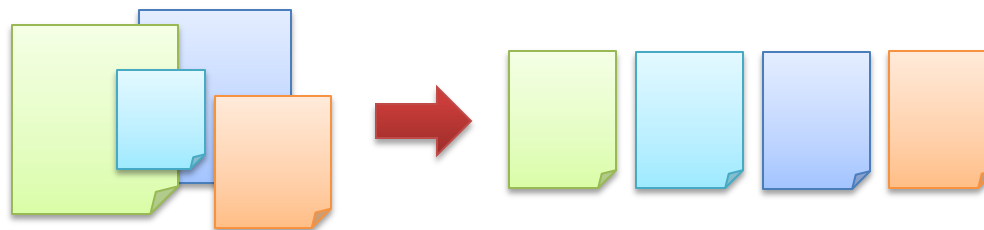
- Term Frequency
 - Based on the observations, high frequency terms are important for describing documents



- Inverse Document Frequency
 - IDF is used to demonstrate the specification of the term



- Document Length Normalization
 - Unify the information quantity of each document



Best Match Models

Best Match Models

- Best Match models were created as the results of a series of experiments on variations of the probabilistic model
- A good term weighting is based on three principles
 - inverse document frequency
 - term frequency
 - document length normalization
- The classic probabilistic model covers only the first of these principles

$$sim_{Probabilistic}(d_j, q) \propto \sum_{w_i \in d_j \& w_i \in q} \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

- This reasoning led to a series of experiments with the Okapi system, which led to the “BM25” ranking formula

Best Match 1 – BM1

- At first, the Okapi system used the Equation below as ranking formula

$$sim_{BM1}(d_j, q) \equiv \sum_{w_i \in \{d_j \cap q\}} \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

- The equation is used in the classic probabilistic model
 - No relevance information can be leverage to estimate a fully probabilistic estimation
 - Consequently, a simple variant is derived
- It was referred to as the BM1 formula

$$\begin{aligned} sim(d_j, q) &\propto \sum_{w_i \in d_j \& w_i \in q} \log \left(\frac{r_i + 0.5}{R_q - r_i + 0.5} \cdot \frac{N - R_q - n_i + r_i + 0.5}{n_i - r_i + 0.5} \right) \\ &\approx \sum_{w_i \in d_j \& w_i \in q} \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right) \end{aligned}$$

Best Match 15 – 1

- The first idea for improving the ranking was to introduce a **term-frequency** factor $\mathcal{F}_{i,j}$ in the BM1 formula

- For document d_j

$$\mathcal{F}_{i,j} = S_1 \times \frac{tf_{i,j}}{K_1 + tf_{i,j}}$$

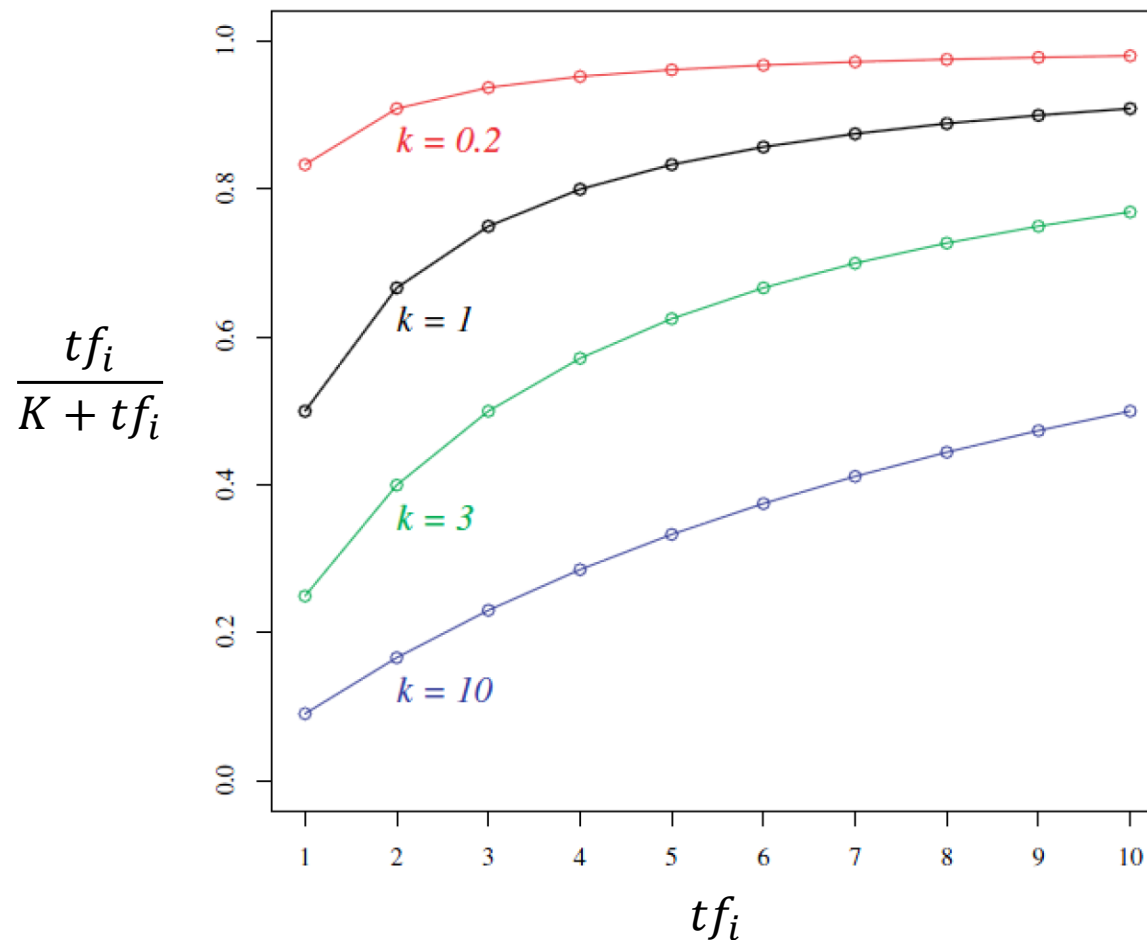
- For query q

$$\mathcal{F}_{i,q} = S_3 \times \frac{tf_{i,q}}{K_3 + tf_{i,q}}$$

- $tf_{i,j}$ (and $tf_{i,q}$) is the frequency of term w_i within d_j (and q)
 - K_1 and K_3 are constants setup experimentally for each collection
 - S_1 and S_3 are scaling constants, normally set to $S_1 = K_1 + 1$ and $S_3 = K_3 + 1$

Best Match 15 – 2

- $\frac{tf_i}{K+tf_i}$ is a saturation function
 - The resulting score is between 0 and 1



Best Match 15 – 3

- Next, a correction factor $G_{j,q}$ dependent on the document and query lengths was introduced

$$G_{j,q} = K_2 \times \text{len}(q) \times \frac{\text{avg}_{doclen} - \text{len}(d_j)}{\text{avg}_{doclen} + \text{len}(d_j)}$$

- $\text{len}(q)$ is the query length (number of terms in the query)
- $\text{len}(d_j)$ is the document length
- avg_{doclen} is the average length of documents in the collection
- K_2 is a constant

Best Match 15 – 4

- To put everything together, we can obtain the BM15 model

$$sim_{BM15}(d_j, q) \equiv G_{j,q} + \sum_{w_i \in \{d_j \cap q\}} \mathcal{F}_{i,j} \times \mathcal{F}_{i,q} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

$$sim_{BM15}(d_j, q) \equiv K_2 \times \frac{\text{len}(q) \times (\text{avg}_{doclen} - \text{len}(d_j))}{\text{avg}_{doclen} + \text{len}(d_j)} + \sum_{w_i \in \{d_j \cap q\}} \frac{S_1 \times tf_{i,j}}{K_1 + tf_{i,j}} \times \frac{S_3 \times tf_{i,q}}{K_3 + tf_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

Tunable Parameters

Document Length Normalization (Correction Factor)

Term Frequency

Inverse Document Frequency

Best Match 11

- A variant was to modify the $\mathcal{F}_{i,j}$ factor by adding **document length normalization** to it

$$\mathcal{F}'_{i,j} = S_1 \times \frac{tf_{i,j}}{\frac{K_1 \times \text{len}(d_j)}{\text{avg}_{doclen}} + tf_{i,j}}$$

- The model is named BM11

$$\text{sim}_{BM11}(d_j, q) \equiv G_{j,q} + \sum_{w_i \in \{d_j \cap q\}} \mathcal{F}'_{i,j} \times \mathcal{F}_{i,q} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

$$\begin{aligned} \text{sim}_{BM11}(d_j, q) \\ \equiv K_2 \times \frac{\text{len}(q) \times (\text{avg}_{doclen} - \text{len}(d_j))}{\text{avg}_{doclen} + \text{len}(d_j)} + \sum_{w_i \in \{d_j \cap q\}} \frac{S_1 \times tf_{i,j}}{K_1 \times \frac{\text{len}(d_j)}{\text{avg}_{doclen}} + tf_{i,j}} \times \frac{S_3 \times tf_{i,q}}{K_3 + tf_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right) \end{aligned}$$

Tunable Parameters

Document Length Normalization (Correction Factor)

Term Frequency

Inverse Document Frequency

BM1, BM15, and BM11

- Introduction of these three factors led to various BM (Best Matching) formulas, as follows:

$$sim_{BM1}(d_j, q) \equiv \sum_{w_i \in \{d_j \cap q\}} \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

$$sim_{BM15}(d_j, q) \equiv G_{j,q} + \sum_{w_i \in \{d_j \cap q\}} \mathcal{F}_{i,j} \times \mathcal{F}_{i,q} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

$$sim_{BM11}(d_j, q) \equiv G_{j,q} + \sum_{w_i \in \{d_j \cap q\}} \mathcal{F}'_{i,j} \times \mathcal{F}_{i,q} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

- Experiments using TREC data have shown that BM11 outperforms BM15

Empirical Considerations

- Several empirical considerations have been discussed
 - Empirical evidence suggests that a best value of K_2 is 0, which eliminates the $G_{j,q}$ factor from these equations

$$\begin{aligned} & \text{sim}_{BM11}(d_j, q) \\ & \equiv K_2 \times \frac{\text{len}(q) \times (\text{avg}_{doclen} - \text{len}(d_j))}{\text{avg}_{doclen} + \text{len}(d_j)} + \sum_{w_i \in \{d_j \cap q\}} \frac{S_1 \times \text{tf}_{i,j}}{K_1 \times \frac{\text{len}(d_j)}{\text{avg}_{doclen}} + \text{tf}_{i,j}} \times \frac{S_3 \times \text{tf}_{i,q}}{K_3 + \text{tf}_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right) \end{aligned}$$

- Further, good estimates for the scaling constants S_1 and S_3 are $K_1 + 1$ and $K_3 + 1$, respectively
- Empirical evidence also suggests that making K_3 very large is better
 - When K_3 is very large, $\mathcal{F}_{i,q}$ factor is reduced simply to $\text{tf}_{i,q}$

$$\mathcal{F}_{i,q} = S_3 \times \frac{\text{tf}_{i,q}}{K_3 + \text{tf}_{i,q}} = \frac{(K_3 + 1)}{K_3 + \text{tf}_{i,q}} \times \text{tf}_{i,q} \approx \text{tf}_{i,q}$$

BM1, BM15, and BM11 Formulas

- By following the considerations, BM models lead to simpler equations

$$sim_{BM1}(d_j, q) \equiv \sum_{w_i \in \{d_j \cap q\}} \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

$$sim_{BM15}(d_j, q) \equiv \sum_{w_i \in \{d_j \cap q\}} \frac{(K_1 + 1) \times tf_{i,j}}{K_1 + tf_{i,j}} \times \frac{(K_3 + 1) \times tf_{i,q}}{K_3 + tf_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

$$sim_{BM11}(d_j, q) \equiv \sum_{w_i \in \{d_j \cap q\}} \frac{(K_1 + 1) \times tf_{i,j}}{\frac{K_1 \times len(d_j)}{avg_{doclen}} + tf_{i,j}} \times \frac{(K_3 + 1) \times tf_{i,q}}{K_3 + tf_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

Best Match 25 – BM25

- The only difference between BM15 and BM 11 is the estimation of the term frequency

$$\mathcal{F}_{i,j}^{BM15} = \frac{(K_1 + 1) \times tf_{i,j}}{K_1 + tf_{i,j}} \quad \mathcal{F}_{i,j}^{BM11} = \frac{(K_1 + 1) \times tf_{i,j}}{\frac{K_1 \times \text{len}(d_j)}{\text{avg}_{doclen}} + tf_{i,j}}$$

- BM25 is proposed to combine BM15 and BM11

$$\mathcal{F}_{i,j}^{BM25} = \frac{(K_1 + 1) \times tf_{i,j}}{K_1 \left[(1 - b) + b \times \frac{\text{len}(d_j)}{\text{avg}_{doclen}} \right] + tf_{i,j}}$$

- b is a constant with values in the interval $[0,1]$
 - If $b = 0$, it reduces to the BM15 term frequency factor
 - If $b = 1$, it reduces to the BM11 term frequency factor
 - For values of b between 0 and 1, the equation provides a combination of BM11 with BM15

BM25 Formula

- To sum up, the BM25 model can be written as:

$$\text{sim}_{BM25}(d_j, q) \equiv \sum_{w_i \in \{d_j \cap q\}} \frac{(K_1 + 1) \times tf_{i,j}}{K_1 \left[(1 - b) + b \times \frac{\text{len}(d_j)}{\text{avg}_{doclen}} \right] + tf_{i,j}} \times \frac{(K_3 + 1) \times tf_{i,q}}{K_3 + tf_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

Tunable Parameters

Document Length Normalization

Term Frequency

Inverse Document Frequency

- b should be kept closer to 1 to emphasize the document length normalization effect present in the BM11 formula
 - $b = 0.75$ is a reasonable assumption
- Constants values (i.e., K_1 , K_3 , and b) can be fine tuned for particular collections through proper experimentation

Further Analysis – 1

- The Okapi BM25 is a state-of-the-art retrieval function for nearly two decades
- The formula can be presented:
 - BM25 weighting

$$\sum_{w_i \in \{d_j \cap q\}} \frac{(K_1 + 1) \times tf_{i,j}}{K_1 \left[(1 - b) + b \times \frac{\text{len}(d_j)}{\text{avg}_{doclen}} \right] + tf_{i,j}} \times \frac{(K_3 + 1) \times tf_{i,q}}{K_3 + tf_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

Document term weighting Query term weighting Discriminative power

Further Analysis – 2

- The key component of BM25 contributing to its success is its term frequency normalization formula:

$$\frac{(K_1 + 1) \times tf_{i,j}}{K_1 \left[(1 - b) + b \times \frac{\text{len}(d_j)}{\text{avg}_{doclen}} \right] + tf_{i,j}} = \frac{(K_1 + 1) \times tf'_{i,j}}{K_1 + tf'_{i,j}}$$

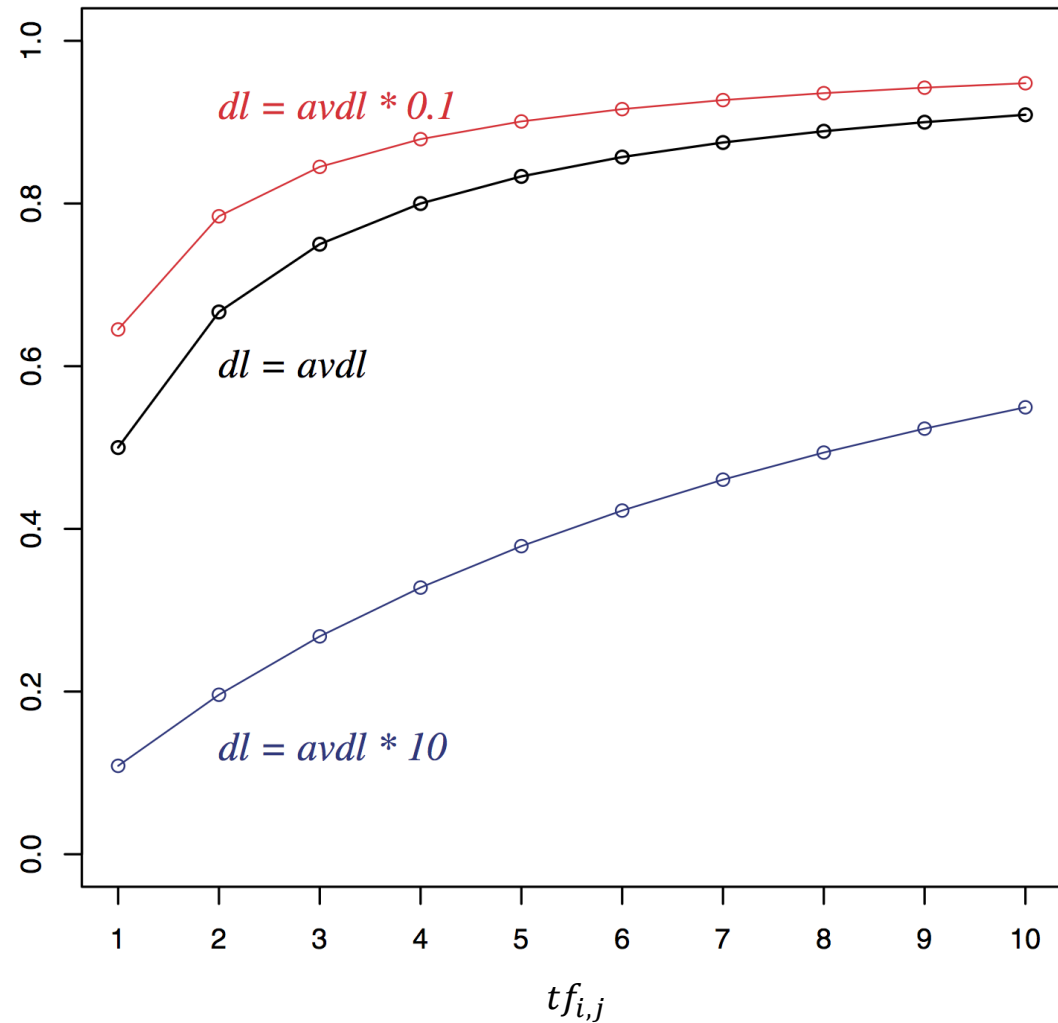
- $tf'_{i,j}$ is the normalized term frequency by document length using **pivoted length normalization**

$$tf'_{i,j} = \frac{tf_{i,j}}{1 - b + b \times \frac{\text{len}(d_j)}{\text{avg}_{doclen}}}$$

$$\text{sim}_{BM25}(d_j, q) \equiv \sum_{w_i \in \{d_j \cap q\}} \frac{(K_1 + 1) \times tf'_{i,j}}{K_1 + tf'_{i,j}} \times \frac{(K_3 + 1) \times tf_{i,q}}{K_3 + tf_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

Further Analysis – 3

$$\frac{tf_{i,j}}{K_1 \left[(1 - b) + b \times \frac{\text{len}(d_j)}{\text{avg}_{\text{doclen}}} \right] + tf_{i,j}}$$



Pros and Cons

- Advantages

- Unlike the probabilistic model, the BM25 formula can be computed without relevance information
- There is consensus that BM25 outperforms the classic vector model for general collections

$$sim_{BM25}(d_j, q) \equiv \sum_{w_i \in \{d_j \cap q\}} \frac{(K_1 + 1) \times tf'_{i,j}}{K_1 + tf_{i,j}} \times \frac{(K_3 + 1) \times tf_{i,q}}{K_3 + tf_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

- Disadvantages

- When a document is very long, we can see that $tf'_{i,j}$ could be very small and approach zero!
 - The presence of w_i in a very long document fails to differentiate clearly from other documents where w_i is absent
 - This suggests that **those very long documents can be overly penalized**

$$tf'_{i,j} = \frac{tf_{i,j}}{1 - b + b \times \frac{len(d_j)}{avg_{doclen}}}$$

Boosting Very Long Documents

- In order to avoid overly-penalizing very long documents, one heuristic way to achieve this goal is to define:

$$\begin{cases} \frac{(K_1 + 1) \times [tf'_{i,j} + \delta]}{K_1 + [tf'_{i,j} + \delta]} & , if \ tf'_{i,j} > 0 \\ 0 & , otherwise \end{cases}$$

$$sim_{BM25L}(d_j, q) \equiv \sum_{w_i \in \{d_j \cap q\}} \frac{(K_1 + 1) \times [tf'_{i,j} + \delta]}{K_1 + [tf'_{i,j} + \delta]} \times \frac{(K_3 + 1) \times tf_{i,q}}{K_3 + tf_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

- A shifted version by adding a shift parameter $\delta > 0$
- We can notice that the modification has a **positive lower bound** for $tf'_{i,j} > 0$

$$\frac{(K_1 + 1) \times \delta}{K_1 + \delta}$$

$$tf'_{i,j} = \frac{tf_{i,j}}{1 - b + b \times \frac{\text{len}(d_j)}{\text{avg}_{doclen}}}$$

Language Models

Language Modeling

- A goal of statistical language modeling is to learn the joint probability function of sequences of words in a language

$$P(w_1, w_2, \dots, w_T)$$

- A statistical model of language can be represented by the conditional probability of the next word given all the previous ones (**chain rule**)

$$\begin{aligned} P(w_1, w_2, \dots, w_T) &= \prod_{t=1}^T P(w_t | w_1, w_2, \dots, w_{t-1}) \\ &\approx \prod_{t=1}^T P(w_t | w_{t-n+1}, \dots, w_{t-1}) \end{aligned}$$

- Such statistical language models have already been found useful in many technological applications involving natural language

N-gram

- Assume words (terms)
 - are independent of each other
 - follow a multinomial distribution
- Unigram
 - Each word occurs independently of the other words
 - The so-called “bag-of-words” model

$$P(w_1, w_2, \dots, w_T) = P(w_1) \cdot P(w_2) \cdots P(w_T) = \prod_{t=1}^T P(w_t)$$

- Bigram

$$P(w_1, w_2, \dots, w_T) = P(w_1)P(w_2|w_1) \cdots P(w_T|w_{T-1}) = P(w_1) \prod_{t=2}^T P(w_t|w_{t-1})$$

- Unigram model is the most popular choice in IR
 - IR does not directly depend on the structure of sentences

NN-based Language Models

- Pre-trained Language Representations (2018)
- Paragraph Representation Models(2014)
- Word Representations(2013, 2014)
- Long Short-Term Memory Language Model(2012)
- Recurrent Neural Network Language Model(2010)
- C&W Neural Network Language Model(2008)
- Log-bilinear Language Model(2007)
- Neural Probabilistic Language Model(2001)

Continuous Language Models

- Tied-Mixture Language Model(2009)
- Continuous Topic Language Model(2008)
- Gaussian Mixture Language Model(2007)

Topic Models

- Latent Dirichlet Allocation(2003)
- Probability Latent Semantic Analysis(1999)
- Latent Semantic Analysis(1997)

Word-Regularity Models

- Regularized Mixture Model(2006)
- Three Mixture Model(2002)
- Simple Mixture Model(2001)
- Relevance-based Language Model(2001)
- Latent Maximum Entropy Model(2001)
- Discriminative Training Language Model(2000)
- Mixture-Based Language Model(1997)
- Aggregate Language Model(1997)
- Mixed-Order Markov Model(1997)
- Structured Model(1997)
- Minimum Word Error Training Language Model(2005)
- Global Conditional Log-linear Model(2007)
- Pseudo-conventional N -gram Model(2008)
- Round-robin Discriminative Language Model(2011)
- Maximum Entropy Model(1994)
- Skipping Model(1993)
- Trigger-based Model(1993)
- Class-based Model(1992)
- Cache-based Model(1988)
- N -gram Model

2000 2002 2004 2006 2008 2010 2012 2014 2016 2018 2020

NN-based Language Models



Word/ Paragraph/
Language
Representations (2013~)

Neural Network Language Models (2001~)

Continuous Language Models



Continuous
Language Models
(2007~2009)

Topic Models (1997~2003)



Topic Models

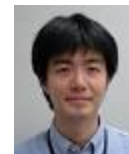
Query Language
Models (2001~2006)



Word-Regularity Models

Discriminative Language Models (2000~2011)

Word-Regularity
Models (~1997)



2000

2002

2004

2006

2008

2010

2012

2014

2016

2018

2020

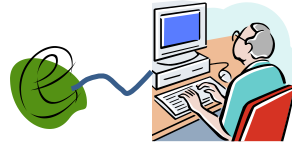
SESAME STREET

- The Sesame Street almost dominates the NLP community now!

Deep Neural Networks Are Our Friends



Language Modeling for Retrieval



- (Statistical) language models (LM) have been widely used for **speech recognition** and **language (machine) translation** for more than thirty years
- However, their use for information retrieval started only in 1998
 - Basically, a query is considered generated from an “ideal” document that satisfies the information need
 - Query-Likelihood Measure
 - KL-Divergence Measure

Query Likelihood Measure

- In the query likelihood retrieval model, we rank documents by the probability that the query could be generated by the document language model

$$P(d_j|q) = \frac{P(q|d_j)P(d_j)}{P(q)} \propto P(q|d_j)P(d_j)$$

$$\approx P(q|d_j) \approx \prod_{i=1}^{|q|} P(w_i|d_j)$$

Document Model

- The user has a prototype (ideal) document in mind, and generates a query based on words that appear in this document
- A document is treated as a model to predict (generate) the query

Document Model – Unigram.

- Use each document itself a sample for estimating its corresponding unigram model
 - The unigram model is a multinomial distribution
 - If Maximum Likelihood Estimation (MLE) is adopted

$$P(w_i|d_j) = \frac{c(w_i, d_j)}{|d_j|}$$

d_j		
w_2	w_1	w_3
w_1	w_1	w_2
w_3	w_3	w_3
w_1	w_2	w_3

$$P(w_1|d_j) = \frac{4}{12}$$

$$P(w_2|d_j) = \frac{3}{12}$$

$$P(w_3|d_j) = \frac{5}{12}$$

$$\begin{aligned} P(q|d_j) &= P(w_1, w_2, w_3|d_j) \\ &= P(w_1|d_j)P(w_2|d_j)P(w_3|d_j) \\ &= \frac{4}{12} \times \frac{3}{12} \times \frac{5}{12} \end{aligned}$$

Document Model – Unigram..

- Use each document itself a sample for estimating its corresponding unigram model
 - The unigram model is a multinomial distribution
 - If Maximum Likelihood Estimation (MLE) is adopted

$$P(w_i|d_j) = \frac{c(w_i, d_j)}{|d_j|}$$

d_j		
w_2	w_1	w_3
w_1	w_1	w_2
w_3	w_3	w_3
w_1	w_2	w_3

$$P(w_1|d_j) = \frac{4}{12}$$

$$P(w_2|d_j) = \frac{3}{12}$$

$$P(w_3|d_j) = \frac{5}{12}$$

$$P(w_4|d_j) = \frac{0}{12}$$

$$\begin{aligned} P(q|d_j) &= P(w_1, w_1, w_4|d_j) \\ &= P(w_1|d_j)P(w_1|d_j)P(w_4|d_j) \\ &= \frac{4}{12} \times \frac{4}{12} \times 0 = 0 \end{aligned}$$

Zero-probability Problem!
Data Sparseness

Smoothing by Background Model.

- Smooth the document-specific unigram model with a collection model
 - Usually named the “background model”
 - The background model can be estimated in a similar way as what we do for the document unigram model

$$P(w_i|BG) = \frac{c(w_i, collection)}{|collection|}$$

- Two representative language model smoothing methods
 - Linear Interpolation (Jelinek-Mercer smoothing)

$$P'(w_i|d_j) = \lambda \cdot P(w_i|d_j) + (1 - \lambda) \cdot P(w_i|BG)$$

- Bayesian Smoothing with Dirichlet Prior

$$P'(w_i|d_j) = \frac{c(w_i, d_j) + \mu \cdot P(w_i|BG)}{|d_j| + \mu}$$

Smoothing by Background Model..

- Linear Interpolation (Jelinek-Mercer smoothing) is the popular one

$$P'(w_i|d_j) = \lambda \cdot P(w_i|d_j) + (1 - \lambda) \cdot P(w_i|BG)$$

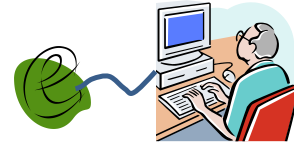
- The role of the background model
 - Help to solve zero-probability problem
 - Play a role as IDF?
 - Not clear!? But it is essential to the good performance

Smoothing and IDF

- Smoothing acts as the IDF factor

$$\begin{aligned} \log P'(q|d_j) &\approx \log \prod_{i=1}^{|q|} P'(w_i|d_j) = \sum_{w_i \in V} \log P'(w_i|d_j)^{c(w_i,q)} \\ &= \sum_{w_i \in d_j} \log P'(w_i|d_j)^{c(w_i,q)} + \sum_{w_i \notin d_j} \log P'(w_i|d_j)^{c(w_i,q)} \\ &= \sum_{w_i \in d_j} \log [\lambda \cdot P(w_i|d_j) + (1 - \lambda) \cdot P(w_i|BG)]^{c(w_i,q)} + \sum_{w_i \notin d_j} \log [(1 - \lambda) \cdot P(w_i|BG)]^{c(w_i,q)} \\ &= \sum_{w_i \in d_j} \log [\lambda \cdot P(w_i|d_j) + (1 - \lambda) \cdot P(w_i|BG)]^{c(w_i,q)} \\ &\quad + \sum_{w_i \in V} \log [(1 - \lambda) \cdot P(w_i|BG)]^{c(w_i,q)} - \sum_{w_i \in d_j} \log [(1 - \lambda) \cdot P(w_i|BG)]^{c(w_i,q)} \\ &= \sum_{w_i \in d_j} \log \left(\frac{\lambda \cdot P(w_i|d_j) + (1 - \lambda) \cdot P(w_i|BG)}{(1 - \lambda) \cdot P(w_i|BG)} \right)^{c(w_i,q)} \\ &\quad + \sum_{w_i \in V} \log [(1 - \lambda) \cdot P(w_i|BG)]^{c(w_i,q)} \end{aligned}$$

KL-Divergence Measure



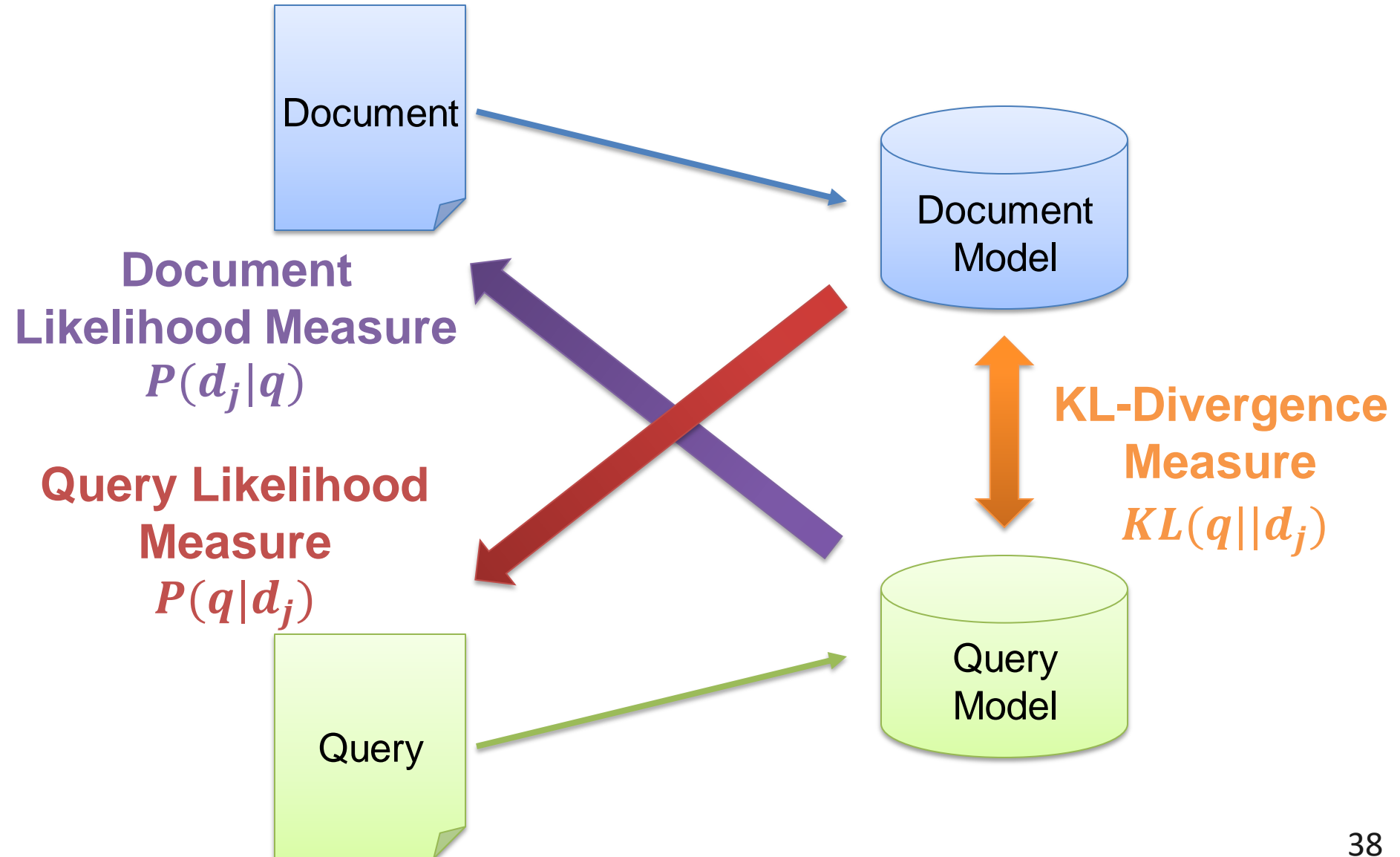
- Another basic formulation of LM for IR is the Kullback-Leibler (KL)-Divergence measure

$$KL(q||d_j) = \sum_{w \in V} P(w|q) \log \frac{P(w|q)}{P(w|d_j)} \propto - \sum_{w \in V} P(w|q) \log P(w|d_j)$$

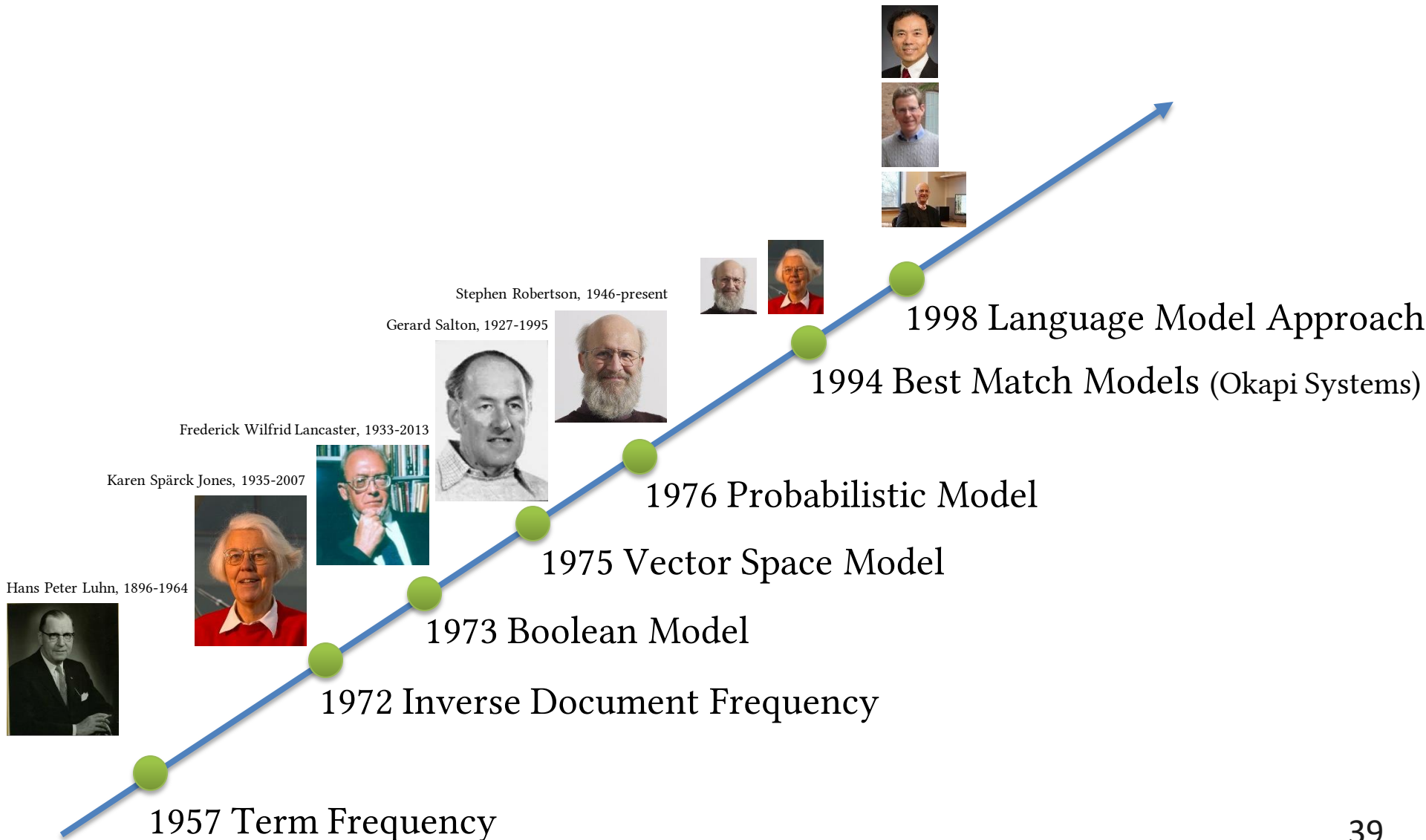
- A query is treated as a **probabilistic model** rather than simply an **observation**
- KL-divergence supports us to achieve a better result by considering **both** query and document models
- KLM can be degenerated to QLM

$$\begin{aligned} KL(q||d_j) &\propto - \sum_{w \in V} P(w|q) \log P(w|d_j) = - \sum_{w \in V} \log P(w|d_j)^{\frac{c(w,q)}{|q|}} \\ &\propto - \sum_{w \in V} \log P(w|d_j)^{c(w,q)} = - \log P(q|d_j) \end{aligned}$$

Three Ways of LM Approaches for IR



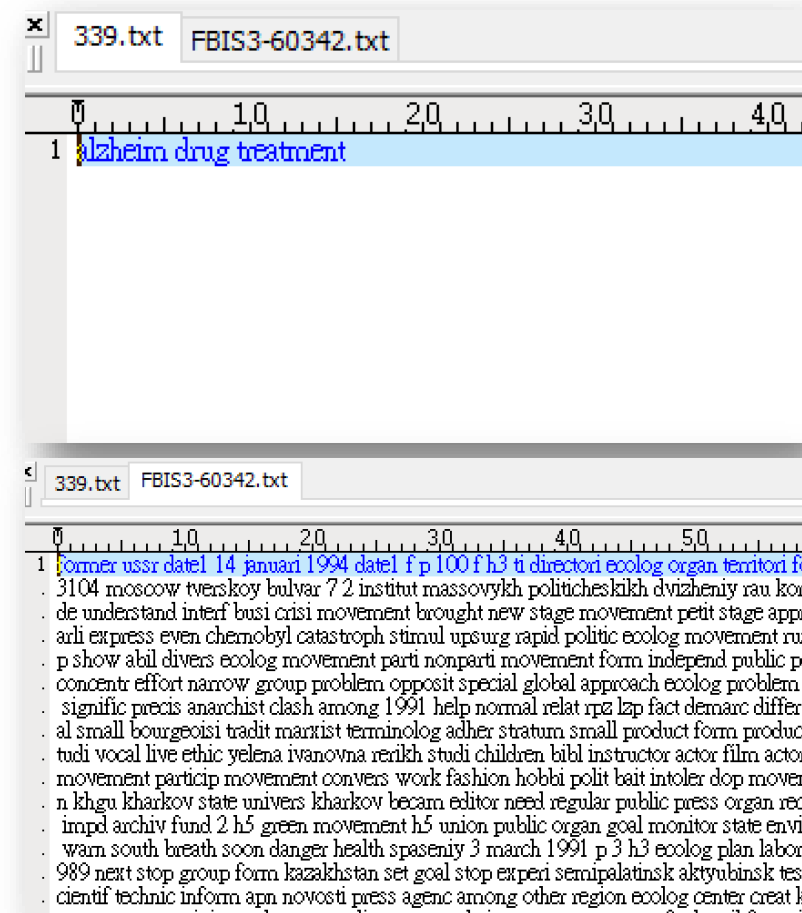
The Evolution



Homework 2 – Best Match Models

Homework 2 - Description.

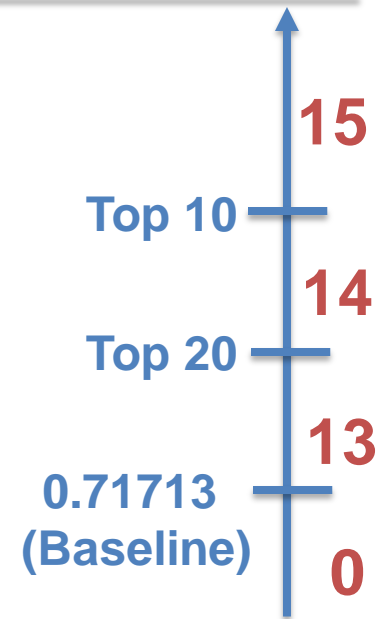
- In this project, we have
 - 50 Queries
 - 4191 Documents
 - Our goal is to implement a BM model
 - BM1
 - BM15
 - BM11
 - BM25
 - BM25L



$$\text{sim}_{\text{BM25}}(d_j, q) \equiv \sum_{w_i \in \{d_j \cap q\}} \frac{(K_1 + 1) \times \text{tf}_{i,j}}{K_1 \left[(1 - b) + b \times \frac{\text{len}(d_j)}{\text{avg}_{\text{doclen}}} \right] + \text{tf}_{i,j}} \times \frac{(K_3 + 1) \times \text{tf}_{i,q}}{K_3 + \text{tf}_{i,q}} \times \log \left(\frac{N - n_i + 0.5}{n_i + 0.5} \right)$$

Homework 2 – Description..

- The evaluation measure is MAP
 - The **hard** deadline is 11/5 23:59
 - Your point is depended on your performance!
 - Please submit a **report** and your **source codes** to the Moodle system, otherwise you will get 0 point
 - The report will be judged by TA, and the score is either 1 or 2
- You should
 - upload your answer file to kaggle
 - <https://www.kaggle.com/t/cd59d965ad1449159533515e1c4a239c>
 - The maximum number of daily submissions is 20
 - **Your team name is ID_Name**



M123456_陳冠宇

Homework 2 – Submission Format

339.txt FBIS3-60342.txt vsm_result.txt

0 10 20 30 40 50

1 Query, Retrieved Documents

2 301,FBIS3-23986 FBIS4-7811 FBIS3-21961 FBIS3-19646 FBIS4-68801 FT
589-0035 LA051689-0064 FT943-789 FBIS4-19393 FBIS4-67908 LA10218
941-2918 FT943-7361 FBIS4-42220 FBIS3-24411 FBIS3-10937 FBIS3-218
96 LA092689-0080 FBIS4-57737 LA121390-0041 FT934-14058 FT933-16
40689-0003 FBIS4-1162 FT921-1950 FT931-10293 FT942-6650 FT944-10
9-0052 FBIS4-26901 FBIS3-23332 LA030190-0006 FR941206-1-00140 FB
589-0098 FBIS3-34637 FBIS3-59697 LA100189-0090 FR940104-0-00032
LA072289-0132 FBIS3-21832 FBIS4-57626 LA120489-0113 FBIS4-5648
FBIS3-4085 FBIS4-3578 FBIS4-35797 FBIS3-4051 FBIS4-36330 LA08218
14179 FT942-1411 FT942-14018 LA051089-0101 FT942-13566 FT942-135
2389-0031 LA032290-0185 LA031190-0126 LA031090-0055 LA021590-0
3-13871 FT923-13808 FT922-9228 FT922-8871 FT922-15435 FT922-5464
T932-2564 FT932-2528 FT932-219 FT932-2181 FT933-10204 FT933-1024
3 302,LA043090-0036 FBIS3-60405 LA072890-0066 FR940126-2-00106 FB
S4-17546 FBIS3-24295 FBIS4-46780 FBIS4-38364 FT943-13796 FT932-5
4-5195 FBIS4-52033 FBIS4-52218 FR940112-2-00082 FR940111-2-00009
2760 FT911-2650 FT911-438 FT911-2421 FT911-2420 FT911-2300 FT911
FBIS3-33671 FBIS3-33570 FBIS3-33505 FBIS3-33438 FBIS3-33287 FB
S3-21817 FBIS3-21807 FBIS3-21790 FBIS3-21771 FBIS3-21769 FBIS3-2
S3-60576 FBIS3-59797 FBIS3-59792 FBIS4-21990 FBIS3-59784 FBIS3-5
90-0043 LA070589-0081 LA070490-0027 LA070490-0026 LA070489-0051
LA040690-0086 LA040689-0176 LA040689-0066 LA040689-0003 LA040
3 LA120190-0068 LA121190-0111 LA121190-0096 LA121090-0065 LA12
FT933-7170 FT933-6233 FT933-7165 FT933-6790 FT933-6767 FT933-673
FT933-12074 FT933-13176 FT933-13170 FT933-13059 FT933-13035 FT9
-17947 FT944-17896 FT944-17680 FT944-17268 FT944-17215 FT944-172
4 303,FT921-7107 LA122990-0029 FT941-15661 FT931-6554 FT944-128 LA
-715 FBIS4-68893 FBIS4-5195 FT942-11934 FBIS3-22049 FBIS4-1213 FB
-56195 FBIS4-67135 FBIS4-67140 FBIS4-67144 FR940114-2-00040 FR94
FT911-368 FT911-3436 FT911-3434 FT911-3409 FT911-336 FT911-2686
FBIS3-37248 FBIS3-36865 FBIS3-3612 FBIS3-3611 FBIS3-35836 FBIS3
FBIS3-21907 FBIS3-21906 FBIS3-21900 FBIS3-21886 FBIS3-21884 FB
9016 FBIS3-58943 FBIS3-58867 FBIS3-58831 FBIS3-58752 FBIS3-58751

Questions?



kychen@mail.ntust.edu.tw