

Vector Space Model, Probabilistic Model

국민대학교 컴퓨터공학부

강 승 식

Vector Space Model

- Motivation
 - Assign *non-binary* weights to index terms
 - A framework in which *partial matching is possible*
 - Instead of attempting to predict whether a document is relevant or not
 - *Rank the documents* according to their degree of similarity to the query

Similarity: Query and Document

$$\vec{q} = (w_{1q}, w_{2q}, \dots, w_{tq}) \quad w_{iq} \geq 0$$

$$\vec{d}_j = (w_{1j}, w_{2j}, \dots, w_{tj}) \quad w_{ij} \geq 0$$

$$\text{sim}(d_j, q) = \frac{\vec{d}_j \cdot \vec{q}}{|\vec{d}_j| \times |\vec{q}|} = \frac{\sum_{i=1}^t w_{ij} \times w_{iq}}{\sqrt{\sum_{i=1}^t w_{ij}^2} \times \sqrt{\sum_{i=1}^t w_{iq}^2}}$$

$$0 \leq \text{sim}(d_j, q) \leq 1 \quad (\text{cosine similarity})$$

$|\vec{q}|$: Does not affect the ranking

$|\vec{d}_j|$: Normalization in the space of the documents

Tf-idf

- Clustering Problem

- Intra-cluster similarity
 - What are the features which better *describe* the objects
- Inter-cluster similarity
 - What are the features which better *distinguish* the objects

- IR Problem

- Intra-cluster similarity (*tf* factor)
 - Raw frequency of a term k_i inside a document d_j
- Inter-cluster similarity (*idf* factor)
 - Inverse of the frequency of a term k_i among the documents

Weighting Scheme

- Term Frequency (*tf*)

- Measure of *how well that term describes the document* contents

$$f_{ij} = \frac{freq_{ij}}{\max_l freq_{lj}} \quad (freq_{ij} : \text{Raw frequency of term } k_i \text{ in the document } d_j)$$

- Inverse Document Frequency (*idf*)

- *Terms which appear in many documents are not very useful for distinguishing a relevant document from a non-relevant one*

$$idf_i = \log \frac{N}{n_i}$$

n_i : Number of documents in which the index term k_i appears

N : Total number of documents

- Best known index term weighting scheme
 - Balance *tf* and *idf* (*tf-idf* scheme)

$$w_{ij} = f_{ij} \times idf_i$$

- Query term weighting scheme

$$w_{iq} = (0.5 + 0.5 f_{iq}) \times idf_i$$

Example

Q : "gold silver truck"

D_1 : "Shipment of gold damaged in a fire"

D_2 : "Delivery of silver arrived in a silver truck"

D_3 : "Shipment of gold arrived in a truck"

$$idf_i = \log \frac{N}{n_i}$$

$$w_{ij} = f_{ij} \times idf_i$$

$$w_{iq} = f_{iq} \times idf_i$$

Term	<i>a</i>	<i>arrived</i>	<i>damaged</i>	<i>delivery</i>	<i>fire</i>	<i>gold</i>	<i>in</i>	<i>of</i>	<i>silver</i>	<i>shipment</i>	<i>truck</i>
<i>idf</i>	0	.176	.477	.477	.477	.176	0	0	.477	.176	.176

	t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	t_9	t_{10}	t_{11}
D_1	0	0	.477	0	.477	.176	0	0	0	.176	0
D_2	0	.176	0	.477	0	0	0	0	.954	0	.176
D_3	0	.176	0	0	0	.176	0	0	0	.176	.176
Q	0	0	0	0	0	.176	0	0	.477	0	.176

	t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	t_9	t_{10}	t_{11}
D_1	0	0	.477	0	.477	.176	0	0	0	.176	0
D_2	0	.176	0	.477	0	0	0	0	.954	0	.176
D_3	0	.176	0	0	0	.176	0	0	0	.176	.176
Q	0	0	0	0	0	.176	0	0	.477	0	.176

$$Sim(Q, D_j) = \sum_{i=1}^n w_{iq} \times w_{ij}$$

$$\begin{aligned}
Sim(Q, D_1) &= (0)(0) + (0)(0) + (0)(0.477) + (0)(0) + (0)(0.477) \\
&+ (0.176)(0.176) + (0)(0) + (0)(0) + (0.477)(0) + (0)(0.176) + (0.176)(0) \\
&= (0.176)^2 \approx 0.031
\end{aligned}$$

$$Sim(Q, D_2) = (0.954)(0.477) + (0.176)^2 \approx 0.486$$

$$Sim(Q, D_3) = (0.176)^2 + (0.176)^2 \approx 0.062$$

Hence, the ranking would be D_2, D_3, D_1

Probabilistic Model

Probabilistic Model

- Motivation
 - Attempt to capture the IR problem within a probabilistic framework
 - Assumption (Probabilistic Principle)
 - Probability of relevance depends on the query and the document representations only
 - There is a subset of all documents which the user prefers as the answer set for the query q
 - Such an *ideal answer set* is labeled R and should maximize the overall probability of relevance to the user
 - Documents in the set R are predicted to be *relevant* to the query
 - Documents not in this set are predicted to be *non-relevant*
- (\overline{R})

Probabilistic Model

- Definition

Using Bayes' rule

$P(R), P(\bar{R})$ are the same
all the documents

$w_{ij} \in \{0,1\}, w_{iq} \in \{0,1\}$: index term weight variables are all binary

$$\text{sim}(d_j, q) = \frac{P(R | \vec{d}_j)}{P(\bar{R} | \vec{d}_j)} = \frac{P(\vec{d}_j | R) \times P(R)}{P(\vec{d}_j | \bar{R}) \times P(\bar{R})} \sim \frac{P(\vec{d}_j | R)}{P(\vec{d}_j | \bar{R})}$$

R : Set of documents known to be relevant

\bar{R} : Set of documents known to be non - relevant

$P(R | \vec{d}_j)$: Probability that the document d_j is relevant to the query q

$$\text{Bayes' rule: } P(a | b) = \frac{P(a \cap b)}{P(b)} = \frac{P(b | a)P(a)}{P(b)}$$

$\vec{P}(\vec{d}_j | R)$: Probability of randomly selecting the document d_j from the set R

$P(R)$: Probability that a document randomly selected from the entire collection is relevant.

Probabilistic Model

- Definition (Cont.)

$$sim(d_j, q) \sim \frac{\left(\prod_{g_i(\bar{d}_j)=1} P(k_i | R) \right) \times \left(\prod_{g_i(\bar{d}_j)=0} P(\bar{k}_i | R) \right)}{\left(\prod_{g_i(\bar{d}_j)=1} P(k_i | \bar{R}) \right) \times \left(\prod_{g_i(\bar{d}_j)=0} P(\bar{k}_i | \bar{R}) \right)}$$

Assuming
independence of
index terms

Taking logarithms,
ignoring constants,
 $P(k_i | R) + P(\bar{k}_i | R) = 1$

$$sim(d_j, q) \sim \sum_{i=1}^t w_{iq} \times w_{ij} \times \left(\log \frac{P(k_i | R)}{1 - P(k_i | R)} + \log \frac{1 - P(k_i | \bar{R})}{P(k_i | \bar{R})} \right)$$

$P(k_i | R)$: the probability that the index term k_i is present in a document randomly selected from the set R

$P(\bar{k}_i | R)$: the probability that the index term k_i is not present in a document randomly selected from the set R

Probabilistic Model

- Initial Probability

$$P(k_i | R) = 0.5$$

$$P(k_i | \bar{R}) = \frac{n_i}{N} \quad n_i : \text{number of documents which contain the index term } k_i$$

- Improving Probability

$$P(k_i | R) = \frac{V_i}{V} = \frac{V_i + 0.5}{V + 1} = \frac{V_i + \frac{n_i}{N}}{V + 1}$$
$$P(k_i | \bar{R}) = \frac{n_i - V_i}{N - V} = \frac{n_i - V_i + 0.5}{N - V + 1} = \frac{n_i - V_i + \frac{n_i}{N}}{N - V + 1}$$

Adjustment the
problems for small
values of V and V_i

V : subset of documents initially retrieved

V_i : subset of V which contain the index term k_i

Probabilistic Model

- Advantage
 - Documents are ranked in decreasing order of their *probability* of being relevant
- Disadvantage
 - *Guess the initial separation* of documents into relevant and non-relevant sets
 - Binary weights
 - Does not take into account the frequency with which an index term occurs inside a document
 - Independence assumption for index terms

Brief Comparison of Classic Models

- Boolean Model
 - Weakest classic method
 - Inability to recognize partial matches
- Vector Model
 - Popular retrieval model
- Vector Model vs. Probabilistic Model
 - Croft
 - Probabilistic model provides a better retrieval performance
 - Salton, Buckley
 - Vector model is expected to outperform the probabilistic model with general collections

Example of Probabilistic Model

- Query : $q(k_1, k_2)$
- 5 documents

	d_1	d_2	d_3	d_4	d_5
k_1	1	0	1	1	0
k_2	0	1	0	1	1

- Assume d_2 and d_4 are relevant

$$P(k_1 | R) = \frac{1}{2} \quad P(k_1 | \bar{R}) = \frac{2}{3} \quad P(k_2 | R) = 1 \quad P(k_2 | \bar{R}) = \frac{1}{3}$$

Example of Probabilistic Model

Q: “gold silver truck”

D1: “Shipment of gold damaged in a fire.”

D2: “Delivery of silver arrived in a silver truck.”

D3: “Shipment of gold arrived in a truck.”

- Assume that documents D2 and D3 are relevant to the query

N= number of documents in the collection

R= number of relevant documents for a given query Q

n= number of documents that contain term t

r= number of relevant documents that contain term t

variable	Gold	Silver	Truck
N	3	3	3
R	2	2	2
n	2	1	2
r	1	1	2

Example of Probabilistic Model (Cont.)

- Computing Term Relevance Weight

$$tr_j = \log\left(\frac{r + 0.5}{R - r + 0.5} \div \frac{n - r + 0.5}{(N - n) - (R - r) + 0.5}\right)$$

$$gold : \log\left(\frac{1 + 0.5}{2 - 1 + 0.5} \div \frac{2 - 1 + 0.5}{3 - 2 - 2 + 1 + 0.5}\right) = \log \frac{1}{3} = -0.477$$

$$silver : \log\left(\frac{1 + 0.5}{2 - 1 + 0.5} \div \frac{1 - 1 + 0.5}{3 - 1 - 2 + 1 + 0.5}\right) = \log \frac{1}{0.333} = 0.477$$

$$truck : \log\left(\frac{2 + 0.5}{2 - 2 + 0.5} \div \frac{2 - 2 + 0.5}{3 - 2 - 2 + 2 + 0.5}\right) = \log \frac{5}{0.333} = 1.176$$

- Similarity Coefficient for each document
 - D1: -0.477, D2: 1.653, D3: 0.699

Review of Probability Theory

Review of Probability Theory

- Try to predict whether or not a baseball team (e.g. LA Dodgers) will win one of its games
 - $P(\text{win}) = 0.5$, $P(\text{win} \mid \text{sunny}) = 0.75$, $P(\text{win} \mid \text{chanho}) = 0.6$
 - $P(\text{win} \mid \text{sunny}, \text{chanho}) = ?$
- Let's assume the independence of evidences
 - $P(w \mid s, c) = \alpha$, $P(w \mid s) = \beta$, $P(w \mid c) = \gamma$
- By Bayes' Theorem

$$\alpha = P(w \mid s, c) = \frac{P(w, s, c)}{P(s, c)} = \frac{P(s, c \mid w)P(w)}{P(s, c)}$$

We can't calculate it, so...

Probability
to lose

$$\frac{\alpha}{1 - \alpha} = \frac{P(s, c \mid w)P(w)}{P(s, c)} \div \frac{P(s, c \mid l)P(l)}{P(s, c)} = \frac{P(s, c \mid w)P(w)}{P(s, c \mid l)P(l)}$$

Review of Probability Theory (Cont.)

- Independence assumption

$$\begin{pmatrix} P(s, c | w) = P(s | w)P(c | w) \\ P(s, c | l) = P(s | l)P(c | l) \end{pmatrix}$$

- Making substitutions

$$\frac{\alpha}{1-\alpha} = \frac{P(s, c | w)P(w)}{P(s, c | l)P(l)} = \frac{P(s | w)P(c | w)P(w)}{P(s | l)P(c | l)P(l)}$$

$$\beta = P(w | s) = \frac{P(s | w)P(w)}{P(s)}, \quad \frac{\beta}{1-\beta} = \frac{P(s | w)P(w)}{P(s | l)P(l)}$$

$$\gamma = P(w | c) = \frac{P(c | w)P(w)}{P(c)}, \quad \frac{\gamma}{1-\gamma} = \frac{P(c | w)P(w)}{P(c | l)P(l)}$$

Review of Probability Theory (Cont.)

- Making substitutions (Cont.)

$$\begin{aligned}\frac{\alpha}{1-\alpha} &= \frac{P(s|w)P(c|w)P(w)}{P(s|l)P(c|l)P(l)} \\ &= \left(\frac{\beta}{1-\beta} \frac{P(l)}{P(w)} \right) \left(\frac{\gamma}{1-\gamma} \frac{P(l)}{P(w)} \right) \frac{P(w)}{P(l)} \\ &= \frac{\beta}{1-\beta} \times \frac{\gamma}{1-\gamma} \times \frac{P(l)}{P(w)} \\ &= \frac{0.75}{0.25} \times \frac{0.6}{0.4} \times \frac{0.5}{0.5} = 4.5\end{aligned}$$

- $\alpha = \frac{9}{11} = 0.818$