

Document Classification

- Document -1 vector: $d1=(x_1, \dots x_N)$, $x_i \in \mathbb{R}$
- Document -2 vector: $d2=(y_1, \dots y_N)$, $y_i \in \mathbb{R}$
- Design a Similarity function: $f(d1, d2)$
 - A good document similar function should identify relevant documents
- Similarity function based on
 - Bag of Words
 - Term Frequency (TF)
 - Document Length

Many Similarity Function

- Distance Function (q, d)
 - **Vector space model**
 - Probabilistic models $P(s=1 \mid q, d)$
 - Probabilistic inference model $f(q, d) = p(d \rightarrow q)$
 - Axiomatic model
 - Deep Learning

Vector Space Model (VSM) Assumptions

- Terms are assumed to be orthogonal (independent from each other)
- Term: basic concepts such as word or phrase (n-gram)
- N terms define an N-dimensional space
- Vector Placements base on Terms (w_1, \dots, w_N)
- Term weight in query and document indicates how well the term characterizes the document or query.

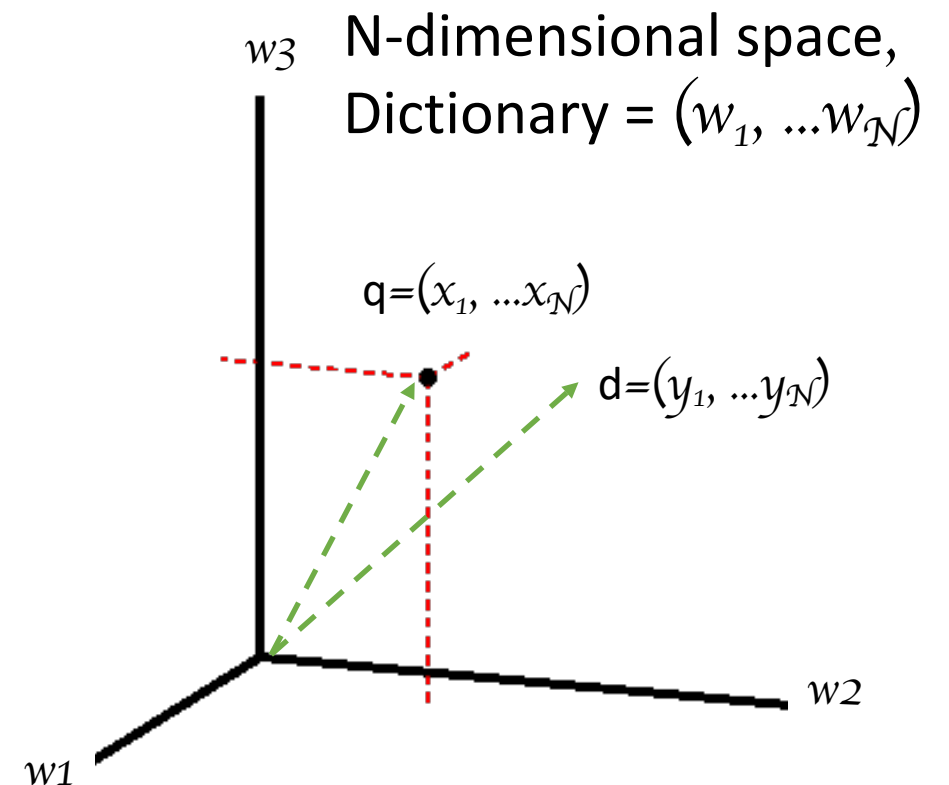
We need to define a similarity function Distance Function (q, d)
for measuring the relevance.

Bit Vector Presentation

- Represent a document or query by a **term vector**
- Query vector: $q=(x_1, \dots x_N)$, $x_i \in \{0,1\}$ is *query term weight*
- Document vector: $d=(y_1, \dots y_N)$, $y_i \in \{0,1\}$ is *document term weight*

1: if word w_i is present

0: if word w_i is not present



Similarity \rightarrow Distance Function: DoT Product

Dictionary = (w_1, \dots, w_N)

$q = (x_1, \dots, x_N), x_i \in \{0, 1\}$

$d = (y_1, \dots, y_N), y_i \in \{0, 1\}$

IF word w_i is present in q or d THEN $x_i = 1$

IF word w_i is absent in q or d THEN $x_i = 0$

$$D(q, d) = \vec{q} \cdot \vec{d} = x_1 y_1 + x_2 y_2 + \dots + x_n y_n = \sum_{i=1}^N x_i y_i$$

Example 1

- I hate running
- I like NLP
- I like deep learning

Terms = {I, like, hate, Deep, Learning, NLP, flying}

Example 2

Query = news about sport campaign

Document 1 = .. news about food campaign

Document 2 = .. news of sport campaign

Document 3 = .. news of sport campaign ... sport activities ..

Similarity \rightarrow Distance Function: DoT Product

Improved SVM with Term Frequency Weighting

Dictionary = (w_1, \dots, w_N)

$q = (x_1, \dots, x_N)$, $x_i = \text{count of word } w_i \text{ in } q = c(w_i, q)$

$d = (y_1, \dots, y_N)$, $y_i = \text{count of word } w_i \text{ in } d = c(w_i, d)$

$$D(q, d) = \vec{q} \cdot \vec{d} = x_1 y_1 + x_2 y_2 + \dots + x_n y_n = \sum_{i=1}^N c(w, q) c(w, d)$$

IDF Weighting

Penalizing popular terms

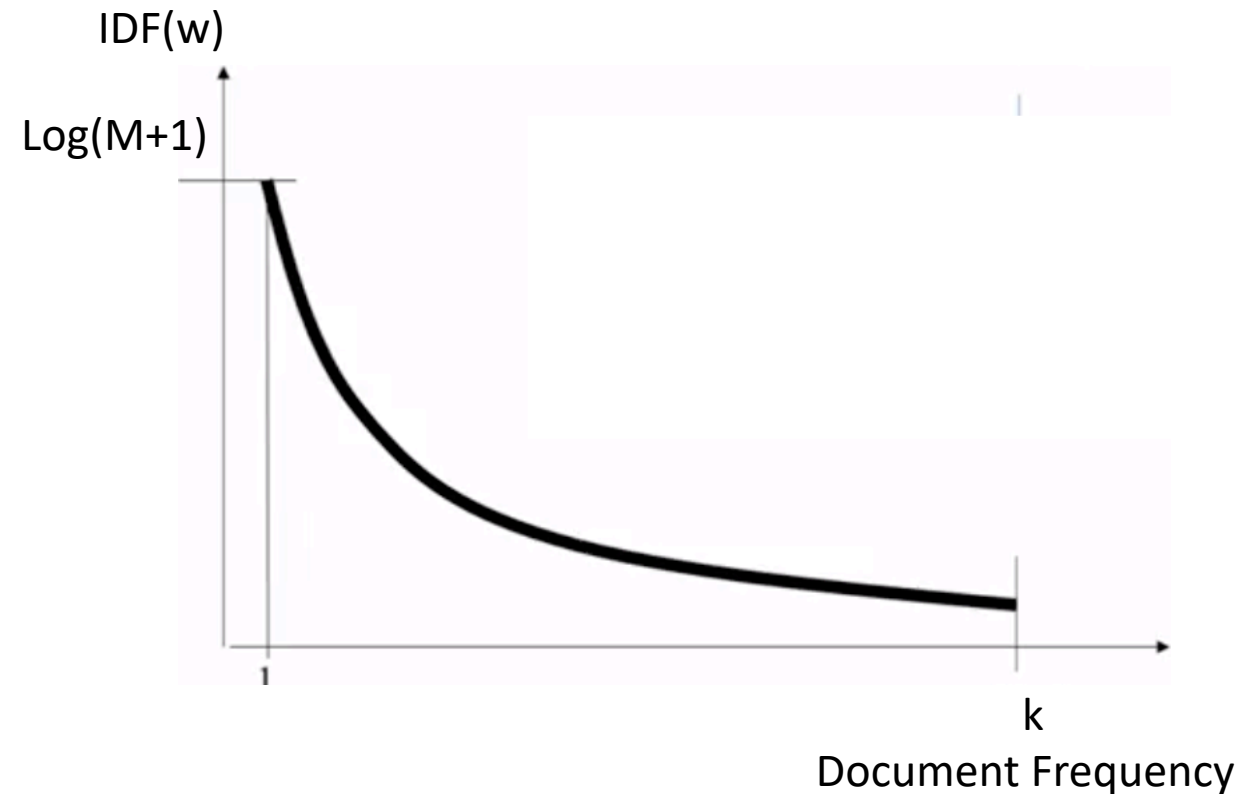
$$\text{IDF}(\mathbf{w}) = \log \frac{M+1}{k(\mathbf{w})}$$

M = total number of docs in collection

$df(\mathbf{w})$ = total number of docs containing \mathbf{w}
(Document frequency)

$k(\mathbf{w}) < M$

$k(\mathbf{w}) \rightarrow M$ and M is large $\text{IDF}(\mathbf{w}) \rightarrow 0$



Similarity \rightarrow Distance Function: DoT Product

Further Improvement of Vector Placement

Inverse Document Frequency (IDF)

Dictionary = (w_1, \dots, w_N)

$q = (x_1, \dots, x_N)$, $x_i = \text{count of word } w_i \text{ in } q = c(w, q)$

$d = (y_1, \dots, y_N)$, $y_i = \text{count of word } w_i \text{ in } d * \text{IDF}(w_i)$
 $= c(w, d) * \text{IDF}(w_i)$

$$\begin{aligned} D(q, d) &= \vec{q} \cdot \vec{d} = x_1 y_1 + x_2 y_2 + \dots + x_n y_n \\ &= \sum_{i=1}^N c(w, q) * c(w, d) * \text{IDF}(w) \end{aligned}$$

Similarity Function with TF-IDF Weighting

$$f(q, d) = \sum_{i=1}^N x_i y_i = \sum c(w, q) \underbrace{c(w, d)}_{\text{TF}(w)} \underbrace{\log \frac{M + 1}{df(w)}}_{\text{IDF}(w)}$$

Query: q

Document: d

All matches query words in d : $w \in q \cap d$

Document Frequency: $df(w)$

Total number of documents in collection: M

Count of word w in query q : $c(w, q)$

Count of word w in document d : $c(w, d)$

Example 2

Query = news about sport campaign

Document 1 = .. news about food campaign

Document 2 = .. news of sport campaign

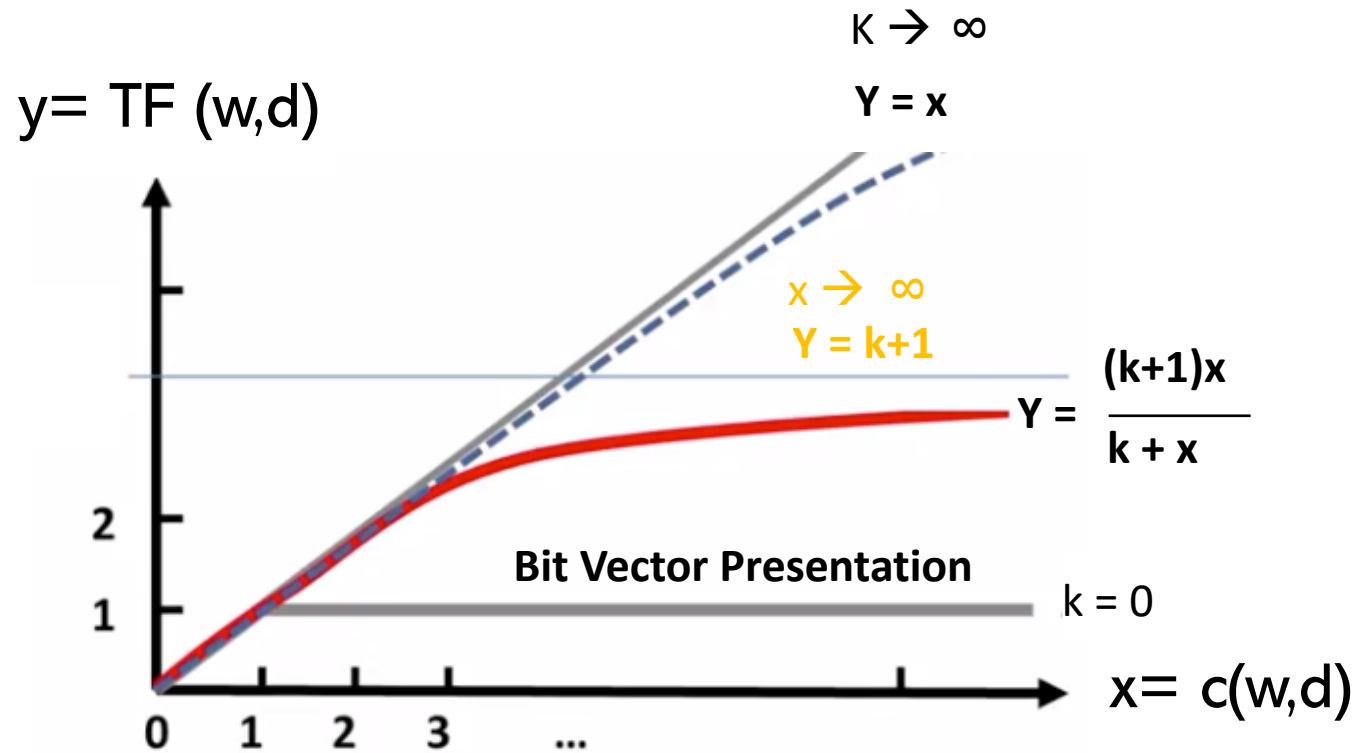
Document 3 = .. news of sport campaign ... sport activities ..

Document 4 = .. news of food campaign ... campaign ... campaign

$$f(q, doc\ 4) = \sum_{i=1}^N x_i y_i = \sum c(w, q) c(w, d) \log \frac{M+1}{df(w)}$$

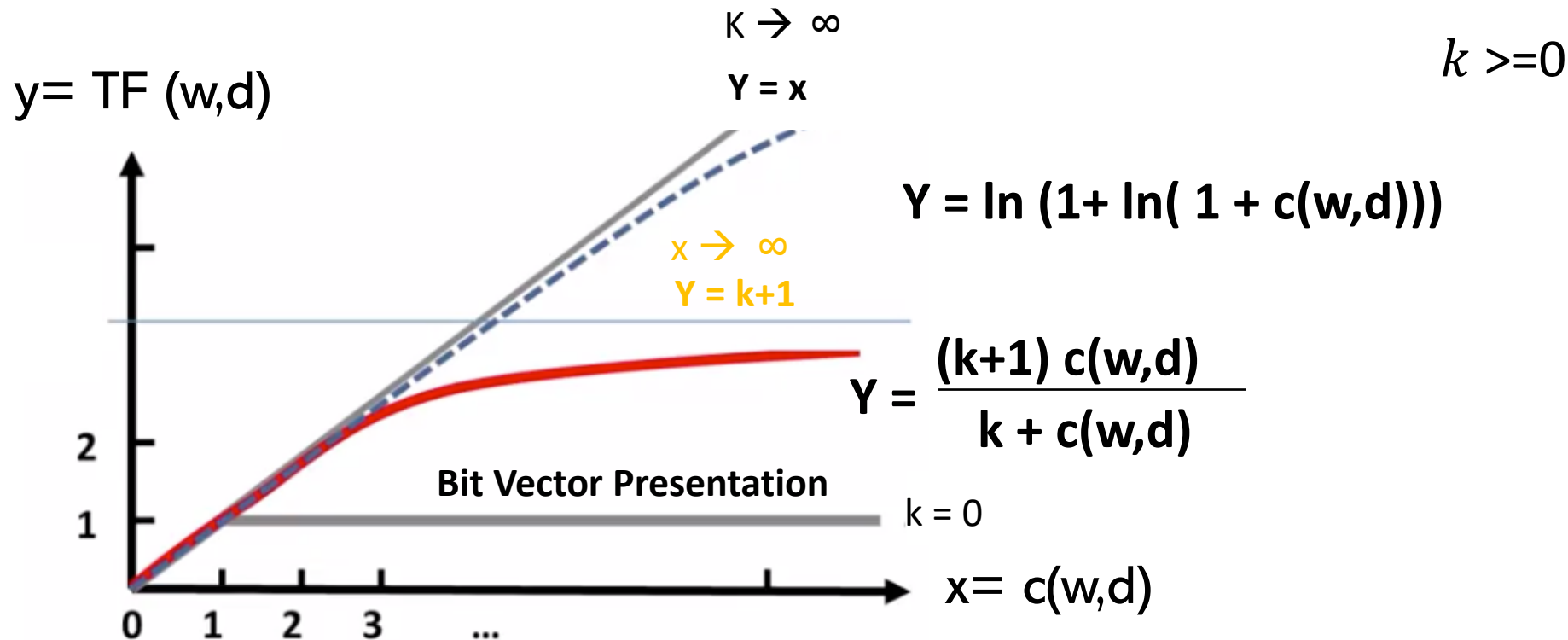
$$\underbrace{c("campaign", doc4)} * \log \frac{5}{4} =$$

BM25 Transformation



BM25 Transformation

$$f(q, d) = \sum_{i=1}^N x_i y_i = \sum c(w, q) c(w, d) \frac{(k + 1)}{k + c(w, d)} \log \frac{M + 1}{df(w)}$$



Document Length Normalization

A long document has a higher chance to match any query so we should penalize the document length

$$\text{Normalizer} = 1 - b + b \frac{d}{\text{Avg}(\text{all } d)}$$

b is [0, 1]

BM25/Okapi

A long document has a higher chance to match any query so we should penalize the document length

$$\text{Normalizer} = 1 - b + b \frac{d}{\text{Avg}(\text{all } d)}$$

b is [0, 1]

$$f(q, d) = \sum_{i=1}^N x_i y_i = \sum c(w, q) c(w, d) \frac{(k + 1)}{k + c(w, d)} \frac{1}{1 - b + b \frac{d}{\text{Avg}(\text{all } d)}} \log \frac{M + 1}{df(w)}$$

Pivoted Length Normalization VSM

$$f(q, d) = \sum_{i=1}^N x_i y_i = \sum c(w, q) \frac{\ln(1 + \ln(1 + c(w, d)))}{1 - b + b \frac{d}{\text{Avg}(\text{all } d)}} \log \frac{M + 1}{df(w)}$$