

Elements of elements of information retrieval

Czyli jak trudny był doktorski egzamin przedmiotowy?

Dominika Tkaczyk

ICM, Uniwersytet Warszawski

21 June 2016



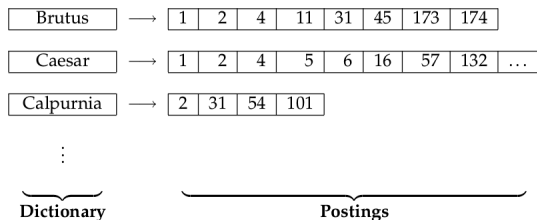
“**Information retrieval (IR)** is finding material (usually documents) of an unstructured nature (usually text) that satisfies an information need from within large collections (usually stored on computers).”

Basic terms:

- collection, document, term
- information need, query, relevance
- effectiveness, precision, recall

The whole presentation is based on: **Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Introduction to Information Retrieval, Cambridge University Press. 2008**, available at <http://nlp.stanford.edu/IR-book/>

Boolean retrieval model



- documents are **sets of words**
- query is in the form of a **Boolean expression of terms**, eg. "Brutus AND Caesar AND NOT Calpurnia"
- term-document incidence matrix vs. **inverted index**
- what if the query retrieves a lot of documents?

Advanced Search

Select items from [?](#)

Where	<input type="text" value="Title"/>	<input type="text" value="matches any"/>	of the following words or phrases:	<input type="text"/>	<input type="button" value="-"/>	<input type="button" value="+"/>
Where	<input type="text" value="Author"/>	<input type="text" value="matches all"/>	of the following words or phrases:	<input type="text"/>	<input type="button" value="-"/>	<input type="button" value="+"/>
Where	<input type="text" value="Abstract"/>	<input type="text" value="matches none"/>	of the following words or phrases:	<input type="text"/>	<input type="button" value="-"/>	<input type="button" value="+"/>
Where	<input type="text" value="Publication Year"/>	<input type="text" value="is in the range"/>	<input type="text" value="1947"/>	to	<input type="text" value="2016"/>	<input type="button" value="-"/> <input type="button" value="+"/>

[\[clear\]](#)

[\[sign in required to save query\]](#) [\[show query syntax\]](#)

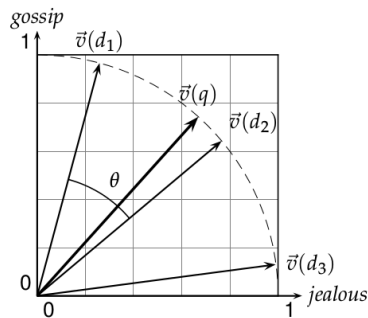
- $\text{score} = \sum_{i=1}^n g_i s_i$, where $g_i \in [0, 1]$ and $\sum_{i=1}^n g_i = 1$
- weights g_i are provided by the user or learned from the training set

Term weighting

- **term frequency:** $tf_{t,d}$ = how many times the term appears in the document
- **document frequency:** df_t = how many documents in the collections contain the term
- **inverse document frequency:** $idf_t = \log \frac{N}{df_t}$
- **term weight:** $tf-idf_{t,d} = tf_{t,d} \times idf_t$

Vector space model

- dimensions \sim (chosen) terms
- documents \sim vectors of tf-idf
- query \sim vector in the same space
- similarity between a document and a query \sim vectors' cosine similarity
- similarity between documents \sim vectors' cosine similarity



Tf-idf variants

Term frequency		Document frequency		Normalization	
n (natural)	$tf_{t,d}$	n (no)	1	n (none)	1
l (logarithm)	$1 + \log(tf_{t,d})$	t (idf)	$\log \frac{N}{df_t}$	c (cosine)	$\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$
a (augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{\max_i(tf_{i,d})}$	p (prob idf)	$\max\{0, \log \frac{N - df_t}{df_t}\}$	u (pivoted unique)	$1/u$
b (boolean)	$\begin{cases} 1 & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$			b (byte size)	$1/CharLength^\alpha, \alpha < 1$
L (log ave)	$\frac{1 + \log(tf_{t,d})}{1 + \log(\text{ave}_{t \in d}(tf_{t,d}))}$				

We need:

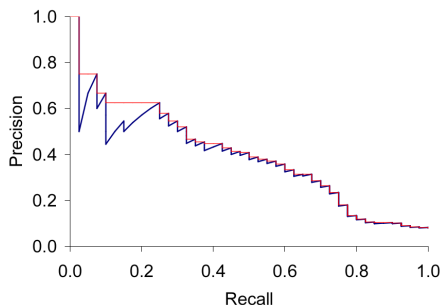
- a document collection
- a test suite of information needs, expressible as queries
- a set of relevance judgments, usually a binary assessment of either *relevant* or *nonrelevant* for each query-document pair

Classification-related measures

- precision = $\frac{\#(\text{relevant items retrieved})}{\#(\text{retrieved items})} = P(\text{relevant} \mid \text{retrieved})$
- recall = $\frac{\#(\text{relevant items retrieved})}{\#(\text{relevant items})} = P(\text{retrieved} \mid \text{relevant})$
- $F_1 = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$
- accuracy = $\frac{\#(\text{relevant items retrieved}) + \#(\text{nonrelevant items not retrieved})}{\#(\text{all items})}$

Ranked retrieval measures

- 11-point interpolated average precision
- Mean Average Precision (MAP)
- precision at k
- R-precision

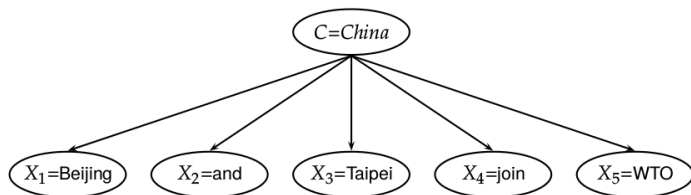


- How reliable are judges' assessments?
- How to know the set of relevant documents in a large collection?
- Marginal relevance: is a document still useful after the user has looked at certain other documents?
- A/B testing, clickthrough log analysis

Text classification

- assigning topics/categories to documents
- vertical search engines restrict searches to a particular topic
- rule-based vs ML-based
- the goal in text classification is high accuracy on new data

Multinomial Naïve Bayes



- $c_d = \arg \max_{c \in C} P(c|d)$
- $P(c|d) \propto P(c) \prod_{1 \leq k \leq n_d} P(t_k|c)$
- parameters estimation: MLE + add-one smoothing

Why is Bayes naïve?

- conditional independence assumption

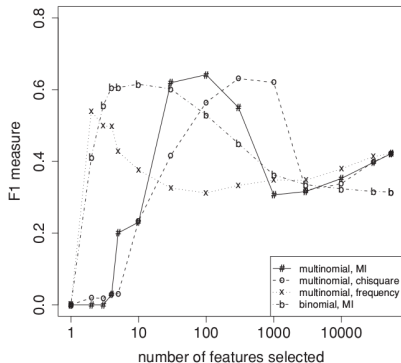
$$P(t_1, \dots, t_{n_d} | c) = \prod_{1 \leq k \leq n_d} P(X_k = t_k | c)$$

- positional independence

$$P(X_{k_1} = t | c) = P(X_{k_2} = t | c)$$

Feature selection

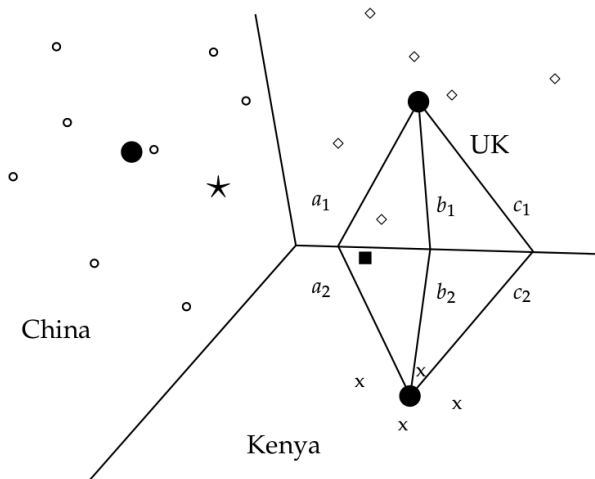
- decreases the vocabulary size
- eliminates noise features
- select k terms with the highest **utility measure**
 - mutual information
 - χ^2 test
 - frequency-based



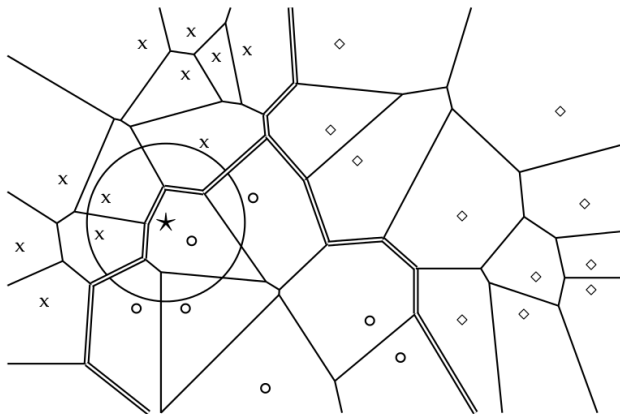
Vector space classification

- represent documents in **vector space model** as before
- now we can apply **various classifiers** to **numerical vectors of features**

Rocchio classification



kNN classification



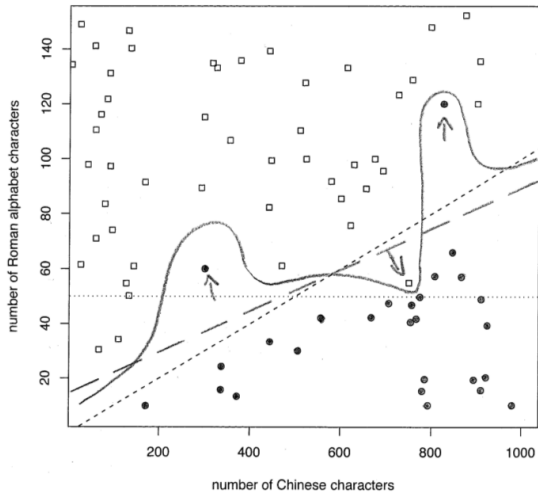
The usual stuff:

- precision, recall, F1
- microaveraging vs. macroaveraging

$$\begin{aligned}\text{learning-error}(\Gamma) &= E_{\mathbf{D}}[\text{MSE}(\Gamma_{\mathbf{D}})] \\ &= E_{\mathbf{D}}E_d[\Gamma_{\mathbf{D}}(d) - P(c|d)]^2 \\ &= E_d[\text{bias}(\Gamma, d) + \text{variance}(\Gamma, d)]\end{aligned}$$

$$\begin{aligned}\text{bias}(\Gamma, d) &= [P(c|d) - E_{\mathbf{D}}\Gamma_{\mathbf{D}}(d)]^2 \\ \text{variance}(\Gamma, d) &= E_{\mathbf{D}}[\Gamma_{\mathbf{D}}(d) - E_{\mathbf{D}}\Gamma_{\mathbf{D}}(d)]^2\end{aligned}$$

Bias-variance tradeoff



Thank you

Thank you!

Dominika Tkaczyk
d.tkaczyk@icm.edu.pl