

# Information Retrieval

Andreas Rauber

<http://www.ifs.tuwien.ac.at/~andi>

## Outline

- retrieval
- relevance feedback
- automatic text classification (ATC)
- evaluation
- clustering

## Retrieval

- documents and queries
- retrieval: find documents satisfying query
- different retrieval models
  - boolean (exact match)
  - vector space model
  - probabilistic models

## Retrieval: exact vs. best match

- exact match:
  - precise query
  - retrieved documents satisfy query criteria
  - result: unordered set of documents
  - efficient
  - predictable, clear evaluation criteria
  - works with clearly identifiable goals
  - good performance in specific domains
  - queries difficult to formulate
  - does not work too well in many general-purpose applications

## Retrieval: exact vs. best match

- best-match:
  - query describes optimal document
  - retrieved documents satisfy criterium as far as possible
  - result: ranked set of documents
  - works well with unclear criteria
  - may return irrelevant documents
  - harder to evaluate wha a certain document was returned
  - seems to outperform exact match in many application scenarios

## Retrieval: boolean model

- most common exact match model
- still widely used
- supports a range of boolean operators:
  - and, or, not
  - proximity
  - position
  - regular expressions

## Retrieval: WESTLAW

- (slides taken from lecture "Retrieval Models" Marten de Rijke et al., Univ. Amsterdam)
- large commercial system
- legal material, news, stock exchange data
- operational since 1974
- approx. 700.000 users, 5-7 TB data
- supports exact-match
- best-match added in 1992

## Retrieval: WESTLAW

- supports range of operators:
  - phrases: "West Publishing"
  - word proximity: West /5 Publishing
  - same sentence: Massachusetts /s technology
  - same paragraph: "information retrieval" /p "exact match"
  - restrictions: (DATE(AFTER 1992 & BEFORE 1995))
- term expansion:
  - wild card: THOM\*ON
  - truncation: THOM!
- queries according to document structure (fields)

(© lecture "Retrieval Models" Marten de Rijke et al., Univ. Amsterdam)

## Retrieval: WESTLAW

- query examples:
  - What is the statute of limitations in cases involving the federal tort claims act?
    - LIMIT! /3 STATUTE ACTION /S FEDERAL /2 TORT /3 CLAIM
  - What factors are important in determining what constitutes a vessel for purposes of determining liability of a vessel owner for injuries to a seaman under the "Jones Act" (46 USC 688)?
    - (741 +3 824) FACTOR ELEMENT STATUS FACT /P VESSEL SHIP BOAT /P (46 +3 688) "JONES ACT" /P INJUR! /S SEAMAN CREWMAN WORKER

(© lecture "Retrieval Models" Marten de Rijke et al., Univ. Amsterdam)

## Retrieval: WESTLAW

- query examples:
  - Are there any cases which discuss negligent maintenance or failure to maintain aids to navigation such as lights, buoys, or channel markers?
    - NOT NEGET! FAIL! NEGLIG! /5 MAINT! REPAIR! /P NAVIGAT! /5 AID EQUIP! LIGHT BUOY "CHANNEL MARKER"
  - What cases have discussed the concept of excusable delay in the application of statutes of limitations or the doctrine of laches involving actions in admiralty or under the "Jones Act" or the "Death on the High Seas Act"?
    - EXCUS! /3 DELAY /P (LIMIT! /3 STATUTE ACTION) LACHES / P "JONES ACT" "DEATH ON THE HIGH SEAS ACT" (46 +3 761)

(© lecture "Retrieval Models" Marten de Rijke et al., Univ. Amsterdam)

## Retrieval: boolean model

- develop query incrementally
- add new terms until results satisfy information need
- relatively long queries
- rather complex queries
- require detailed knowledge on query and document domain
- mostly for trained experts
- easier to control

## Retrieval: vector space model

- best-match
- indexing: high-dimensional feature space
- query: describes "optimal" document, information need
- documents are vectors
- queries are vectors
- return documents that are closest to query in high-dimensional feature space
- measure similarity in two ways:
  - distance between points
  - angle between vectors

## Retrieval: vector space model

- similarity: distance between end points
- range of metrics

- L1 (city block):

- L2 (euclidean distance):

$$m = \sqrt{\sum_{i=1}^n (d_i - q_i)^2}$$

- minkovsky metric (general):

$$m = \left( \sum_{i=1}^n (d_i - q_i)^r \right)^{1/r}$$

## Retrieval: vector space model

- vector normalization
- normalize to unit length

- vector length

$$\|d\| = \sqrt{\sum_{i=1}^n d_i^2}$$

- normalize:

$$\frac{d_i}{\|d\|}$$

## Retrieval: vector space model

- given unit vectors  $d, q$ :

- $\cos = 0$ : no similarity
- $\cos = 1$ : identical docs

- cosine similarity

$$\text{sim}(d, q) = \frac{\sum_{i=1}^n q_i \cdot d_i}{\sqrt{\sum_{i=1}^n q_i^2} \cdot \sqrt{\sum_{i=1}^n d_i^2}} = (\|q\|^{-1} \cdot q) \cdot (\|d\|^{-1} \cdot d)$$

## Retrieval: vector space model

- other similarity measures

- Jacquard measure
- Dice's coefficient

## Retrieval: probabilistic models

- goal: identifying documents that are relevant
- 2-class classification problem
- model probability that document belongs to relevant class
- different ways of modeling probability
- basic assumption: relevance of document is independent of other docs in collection
- use Bayes rule:  $P(R|d) = P(d|R) \cdot P(R) / P(d)$

## Outline

- retrieval
- relevance feedback
- automatic text classification (ATC)
- evaluation
- clustering

## Relevance Feedback

- retrieval: one-stop shop  
if results not satisfactory: re-try  
re-fine: add/remove terms
- interactive retrieval
- relevance feedback: have computer add/remove terms automatically
- from result set: show which ones are relevant/not relevant
- algorithm infers how to weight terms
- active learning

## Relevance Feedback

- problem: not possible to show many
  - select few documents for user to rate
- problem: which examples to choose?
  - most positive:  
probably little additional information  
not useful to capture range of positives
  - most borderline:  
most difficult to decide  
may not lead to most relevant returned

## Relevance Feedback: rocchio

- re-write query by adding terms from relevant documents and subtracting terms from non-relevant documents

$$w_{ki} = \beta \cdot \sum_{\{d_j \in POS_i\}} w_{kj} - \gamma \cdot \sum_{\{d_j \in NEG_i\}} w_{kj}$$

- typical values:  $\beta = 4$ ,  $\gamma = 1$

## Relevance Feedback: blind

- normally, user defines relevant docs
- blind relevance feedback: top-n documents are considered relevant
- Rocchio:  $\beta = 1$ ,  $\gamma = 0$ ,  
only use pos. examples

## Outline

- retrieval
- relevance feedback
- automatic text classification (ATC)
- evaluation
- clustering

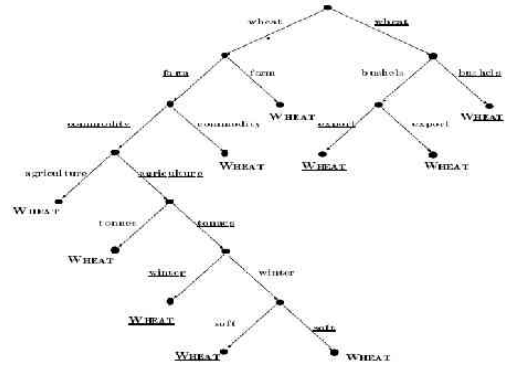
## ATC

- assigning documents to pre-defined classes
- relevant vs. not-relevant
- target function  $f: D \times C \rightarrow \{t, f\}$
- approximated by classifier:
- categories are just symbolic labels
- no exogenous knowledge
- binary, single label, multi-label
- document pivoted vs. category pivoted

# ATC

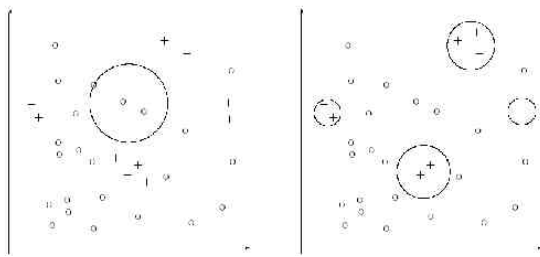
- range of machine learning algorithms
  - knn-classifiers
  - decision trees
  - rocchio
  - naive bayes
  - support vector machines
  - ...
- multi-label: set of individual binary classifiers
- classifier committees

# ATC - decision trees



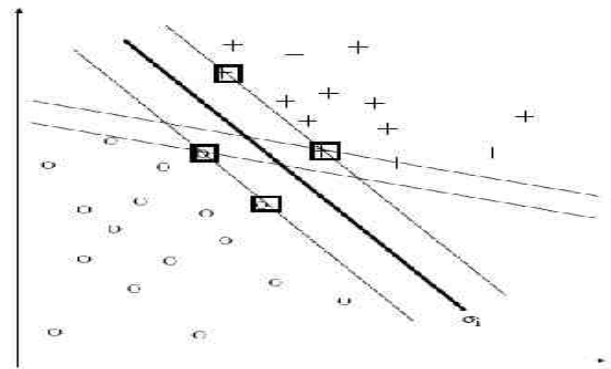
© F. Sebastiani: Machine Learning in Automated Text Categorization, ACM Computing Surveys, 34(1):1-47, 2001

# ATC: Rocchio vs. knn



© F. Sebastiani: Machine Learning in Automated Text Categorization, ACM Computing Surveys, 34(1):1-47, 2001

# ATC - SVM



© F. Sebastiani: Machine Learning in Automated Text Categorization, ACM Computing Surveys, 34(1):1-47, 2001

# ATC - dimensionality reduction

- dimensionality reduction
- global vs. local
- selection vs. extraction
  - information gain
  - chi-squared
  - mutual information
  - ...

# ATC: dimensionality reduction

Function	Denoted by	Mathematical form
Document frequency	$\#(t_k, c_i)$	$P(t_k c_i)$
DIA association factor	$\alpha(t_k, c_i)$	$P(c_i t_k)$
Information gain	$IG(t_k, c_i)$	$P(t_k, c_i) \log \frac{P(t_k, c_i)}{P(t_k) \cdot P(c_i)} + P(t_k, \bar{c}_i) \log \frac{P(t_k, \bar{c}_i)}{P(t_k) \cdot P(\bar{c}_i)}$
Mutual information	$MI(t_k, c_i)$	$\log \frac{P(t_k, c_i)}{P(t_k) \cdot P(c_i)}$
Chi-square	$\chi^2(t_k, c_i)$	$\frac{[P(t_k, c_i) \cdot P(\bar{t}_k, \bar{c}_i) - P(t_k, \bar{c}_i) \cdot P(\bar{t}_k, c_i)]^2}{P(t_k) \cdot P(\bar{t}_k) \cdot P(c_i) \cdot P(\bar{c}_i)}$
NGL coefficient	$NGL(t_k, c_i)$	$\sqrt{2} \cdot [P(t_k, c_i) \cdot P(\bar{t}_k, \bar{c}_i) - P(t_k, \bar{c}_i) \cdot P(\bar{t}_k, c_i)]$
Relevance score	$RS(t_k, c_i)$	$\frac{P(t_k c_i) - d}{\sqrt{P(t_k c_i) - d}}$
Odds Ratio	$OR(t_k, c_i)$	$\frac{P(t_k, \bar{c}_i) \cdot 11 - P(t_k, \bar{c}_i)}{(1 - P(t_k c_i)) \cdot P(t_k, c_i)}$
GSS coefficient	$GSS(t_k, c_i)$	$P(t_k, c_i) \cdot P(t_k, \bar{c}_i) - P(t_k, \bar{c}_i) \cdot P(t_k, c_i)$

© F. Sebastiani: Machine Learning in Automated Text Categorization, ACM Computing Surveys, 34(1):1-47, 2001

## ATC: dimensionality reduction

- term extraction:
  - term clustering
  - replace with cluster centroids
  - "synonyms"
- latent semantic indexing
- singular value decomposition
- creates new feature space
- orthogonal axes

## ATC results

System	Type	# of documents # of training documents # of test documents # of categories	Results reported by				
			P1	P2	P3	P4	P5
Walt	non-linear	Yana 1999	100	100	100		
Duchon	probabilistic	[Domini et al. 1999] [Duchon 1999]	418 (MAP)			192	418
	probabilistic	[Hara et al. 1999]					500
BN	probabilistic	[Hara 1999]	382			147	374
	probabilistic	[Li and Tommasi 1999]					174
Pa & BP	probabilistic	[Yang and Liu 1999]	390			155	390
	probabilistic	[Hachisu 1999]					174
RFB	decision rules	[Kopis et al. 1999]	372			161	374
	decision rules	[Kopis and James 1999]					161
Karna-Raven	decision rules	[Coker and Singer 1997]	360			120	327
	decision rules	[Li and Tommasi 1999]					120
CARR2	decision rules	[Niedtke and Gaus 1996]	339			150	339
	decision rules	[Mouliner et al. 1998]					150
Wang	decision rules	[Yana 1999]	347 (GD)			146	347
	decision rules	[Yang et al. 1999]					146
Wang	decision rules	[Duan et al. 1999]	347 (GD)			146	347
	decision rules	[Duan et al. 1998]					146
Rosen	local linear	[Yokoi and James 1999]	320			170	320
	local linear	[Yokoi and James 1999]					170
Rosen	local linear	[Duchon 1999]	320			170	320
	local linear	[Hara and Red 1999]					170
Rosen	local linear	[Li and Tommasi 1999]	320			170	320
	local linear	[Li and Tommasi 1999]					170
Rosen	local linear	[Yang et al. 1999]	320			170	320
	local linear	[Yang et al. 1999]					170
Rosen	local linear	[Werner et al. 1999]	320			170	320
	local linear	[Werner et al. 1999]					170
K-NN	nearest-neighbor	[Hara and Lee 1998]	320			170	320
	nearest-neighbor	[Duchon 1999]					170
K-NN	nearest-neighbor	[Hara and Red 1999]	320			170	320
	nearest-neighbor	[Yana 1999]					170
K-NN	nearest-neighbor	[Yang and Liu 1999]	320			170	320
	nearest-neighbor	[Duchon et al. 1999]					170
SVM	SVM	[Duchon 1999]	320			170	320
	SVM	[Li and Tommasi 1999]					170
SVM	SVM	[Yang and Liu 1999]	320			170	320
	SVM	[Yang and Liu 1999]					170
SVM	SVM	[Werner et al. 1999]	320			170	320
	SVM	[Werner et al. 1999]					170
SVM	SVM	[Hara et al. 1999]	320			170	320
	SVM	[Hara et al. 1999]					170

© F. Sebastiani: Machine Learning in Automated Text Categorization, ACM Computing Surveys, 34(1):1-47, 2001

## ATC collections

- Reuters (22173, 21578, ModLewis, ModApte)
- OHSUMED
- TREC
- 20-newsgroups

## Outline

- retrieval
- relevance feedback
- automatic text classification (ATC)
- evaluation
- clustering

## Evaluation

- training set - validation set - test set
- same split
- specific task or cross-validation
- same (or: specified) pre-processing
- same (or specified) indexing
- same evaluation measures

## Evaluation

- contingency table:

Classifier		expert judgments	
		YES	NO
classifier	YES	$TP = \sum_{i=1}^c TP_i$	$FP = \sum_{i=1}^c FP_i$
	NO	$FN = \sum_{i=1}^c FN_i$	$TN = \sum_{i=1}^c TN_i$

## Evaluation

- precision:  $\hat{\pi}^i = \frac{TP}{TP+FP} = \frac{\sum_{i=1}^{|C|} TP_i}{\sum_{i=1}^{|C|} (TP_i + FP_i)}$
- recall: (micro-averaged)  $\hat{\rho}^i = \frac{TP}{TP+FN} = \frac{\sum_{i=1}^{|C|} TP_i}{\sum_{i=1}^{|C|} (TP_i + FN_i)}$
- macro averaged: (first per category, then average)  $\hat{\pi}^M = \frac{\sum_{i=1}^{|C|} \hat{\pi}^i}{|C|}$   
 $\hat{\rho}^M = \frac{\sum_{i=1}^{|C|} \hat{\rho}^i}{|C|}$

## Evaluation

- accuracy:  $A = \frac{TP+TN}{TP+TN+FP+FN}$
- error:  $E = \frac{FP+FN}{TP+TN+FP+FN} = 1 - A$
- utility: 

Classifier	actual judgments	
	YES	NO
judgments: YES	TP	FP
judgments: NO	FN	TN

## Evaluation

- combined effectiveness measures
  - break-even: point at which  $F = P$
  - F :  $F_{\beta} = \frac{(\beta^2 + 1)\pi\rho}{\beta^2\pi + \rho}$

## Outline

- retrieval
- relevance feedback
- automatic text classification (ATC)
- evaluation
- clustering

## Clustering

- no class information available
- unsupervised learning
- identify groups of documents with similar characteristics
- used in a range of applications
  - pre-processing
  - novelty detection
  - document summarization
  - text analysis

## Clustering

- range of clustering techniques
- clustering vs. topology-preserving mapping
- visualization
- self-organizing maps
- more difficult to evaluate
- cluster validity measures and interpretations