# Event Detection in Newspaper Texts

Nikola Ljubešić, senior research assistant
Department of Information Sciences
University of Zagreb

JOTA, 28 October 2010
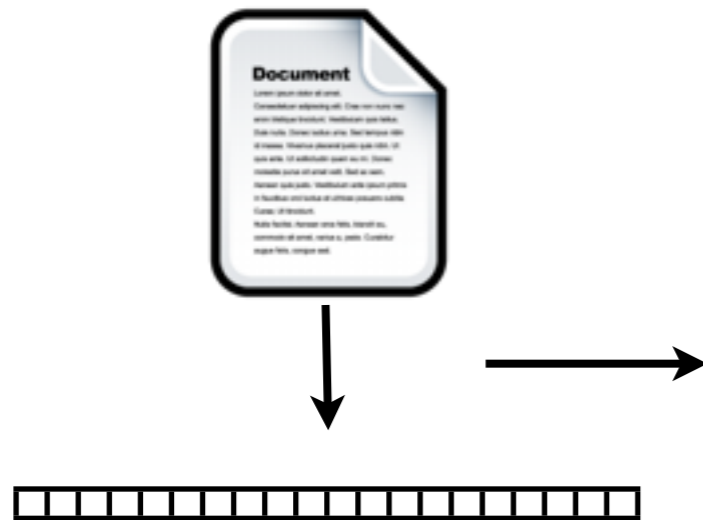
# Overview

# Event detection

- event - a particular thing that happens at a specific time and place (TDT, 2004)

- event detection - process of detecting an event description in a piece of information

- part of the topic detection and tracking problem set

- document : event == 1 : 1?

# Classification problem

- events are categories - classification task

1. unknown classification schema - solvable only by unsupervised classification - clustering

2. unknown number of events - unknown number of classes - hierarchical clustering
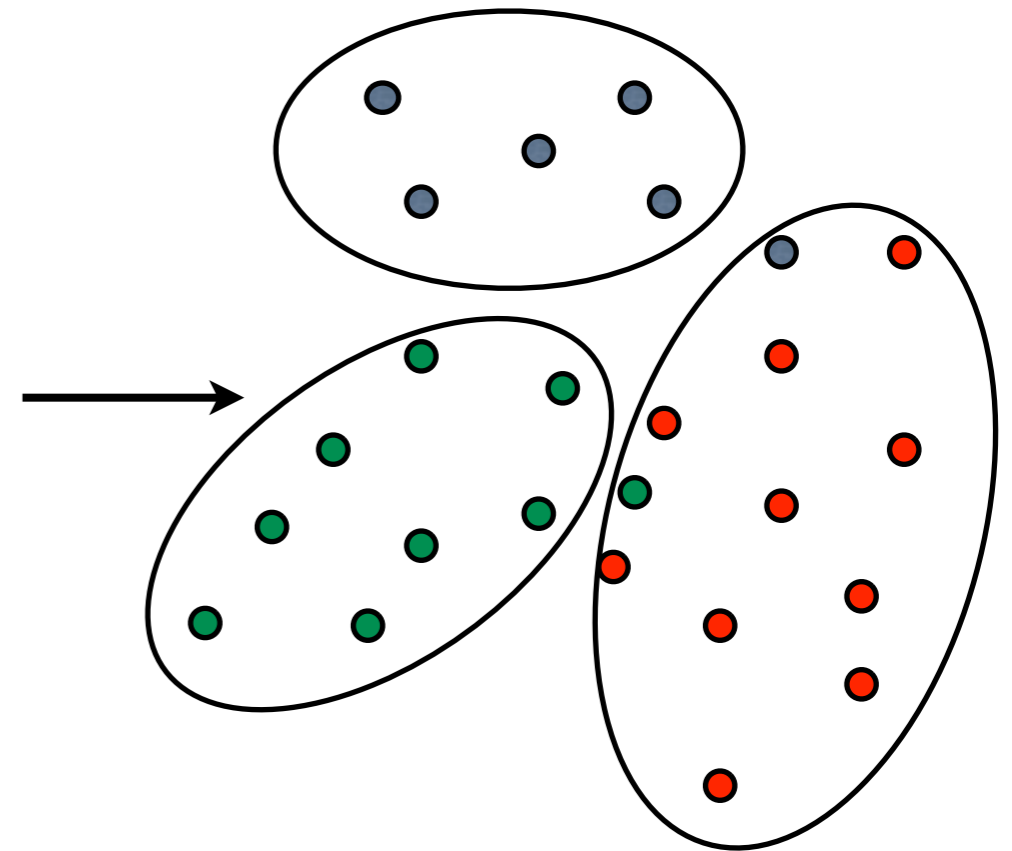
# Document clustering

document
formalization

distance
matrix

clustering

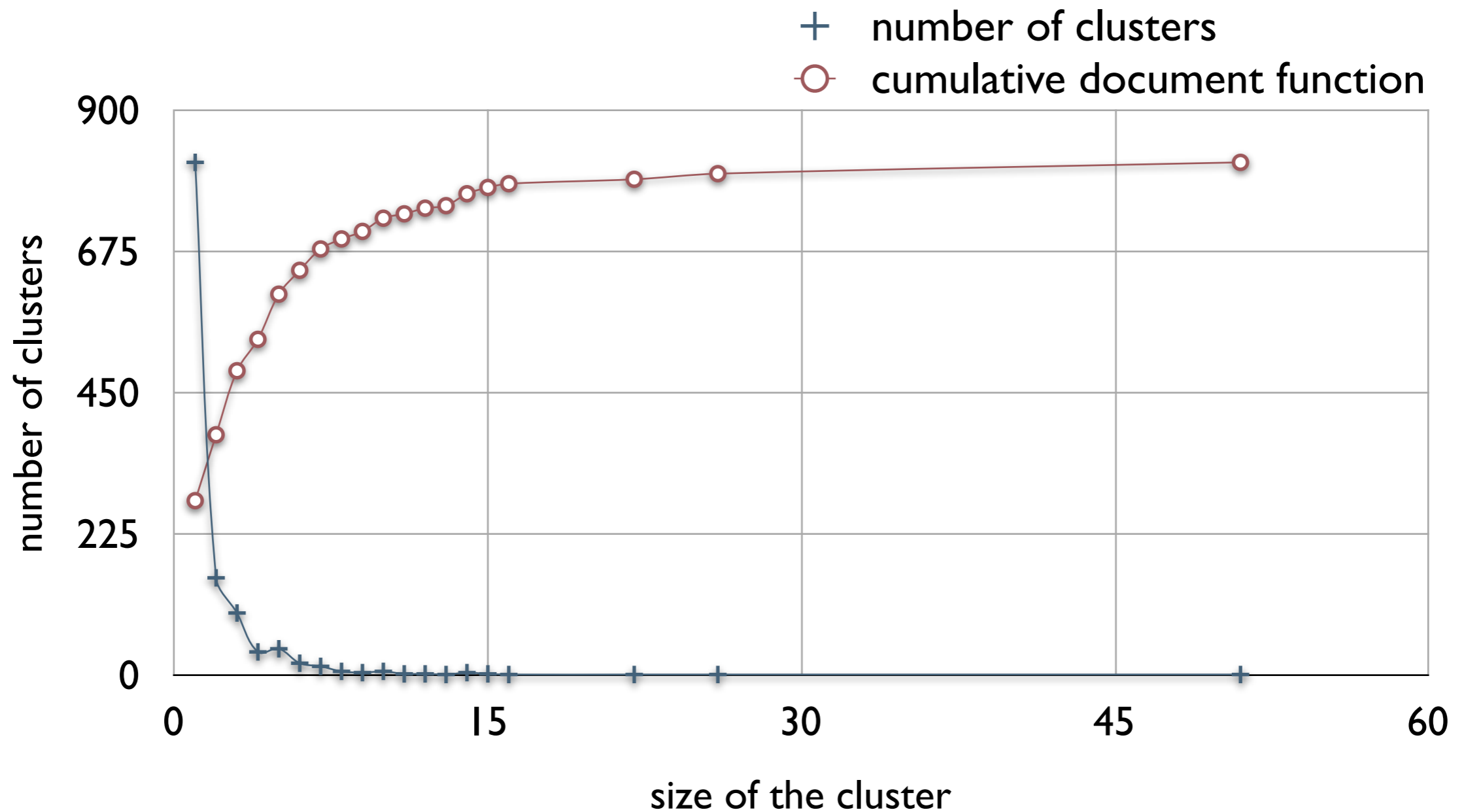|    | d1 | d2 | d3 | d4 | ... |
|----|----|----|----|----|-----|
| d1 |    |    |    |    |     |
| d2 |    |    |    |    |     |
| d3 |    |    |    |    |     |
| d4 |    |    |    |    |     |
| ...|    |    |    |    |     |

# Gold standard

- 2,398 documents published on 17 Croatian news portals in three days

- two annotators, application developed for that purpose

- pooling - using a combination of all similarity metrics to obtain a candidate list

- built 1,214 and 955 clusters

# Inter-annotator agreement

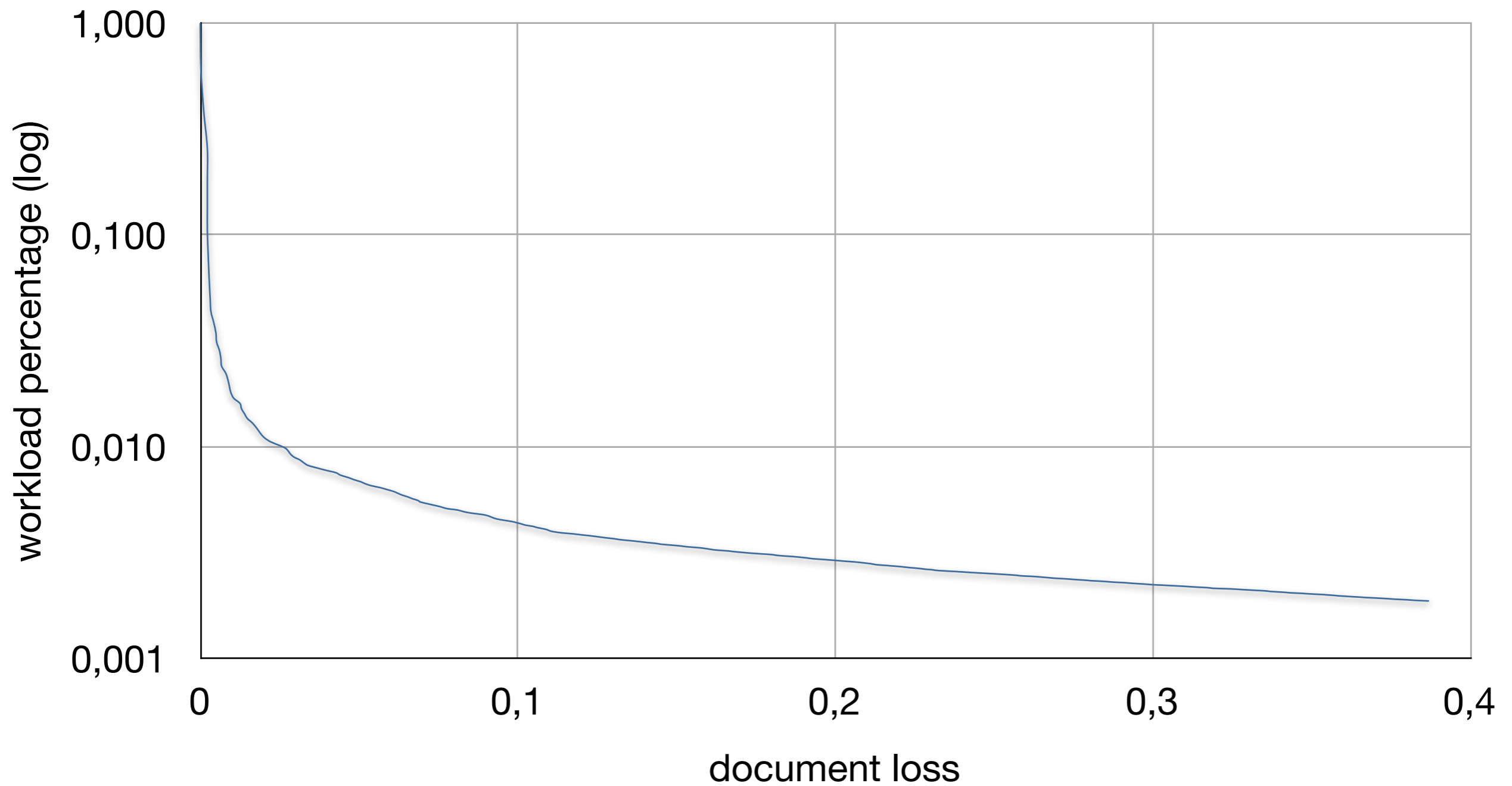| kappa | $\varkappa = \dfrac{2 \cdot |A_1 \cap A_2|}{|A_1| + |A_2|}$ | 0,684 |
|-------|------|-------|
| modified kappa | $\varkappa_{\mathrm{mod}} = \dfrac{|A_1 \cap A_2|}{\min(|A_1|, |A_2|)}$ | 0,91 |

- biggest story of May 3 2009 - the Myanmar cyclone

- annotator 1 - one cluster with 52 documents

- annotator 2 - three clusters - the catastrophe, first rescue operations, Croatian Red cross reaction

# Workload-recall trade-off
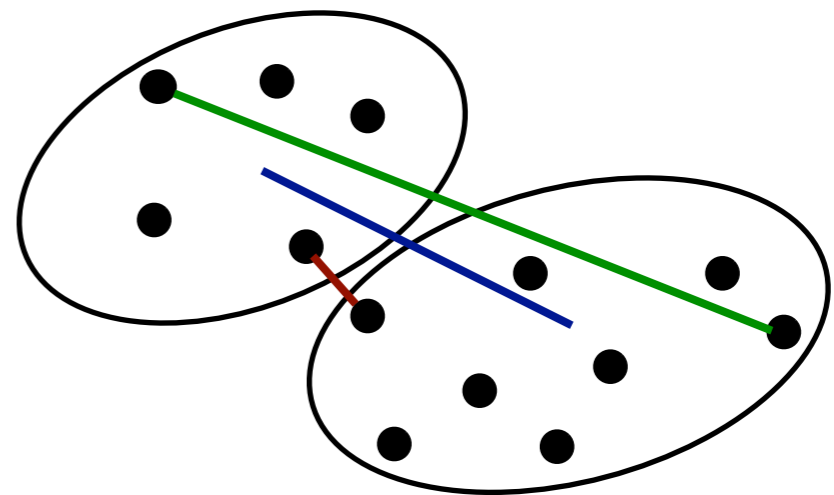
# Experimental setup

- 14 categorical variables with 2-6 levels - 2,073,600 experiments

- huge search space - independence assumption

- variable categories:

  - clustering algorithm

  - distance metrics

  - feature weight measures

  - feature selection and extraction methods

  - reference corpus significance

# Evaluation measures

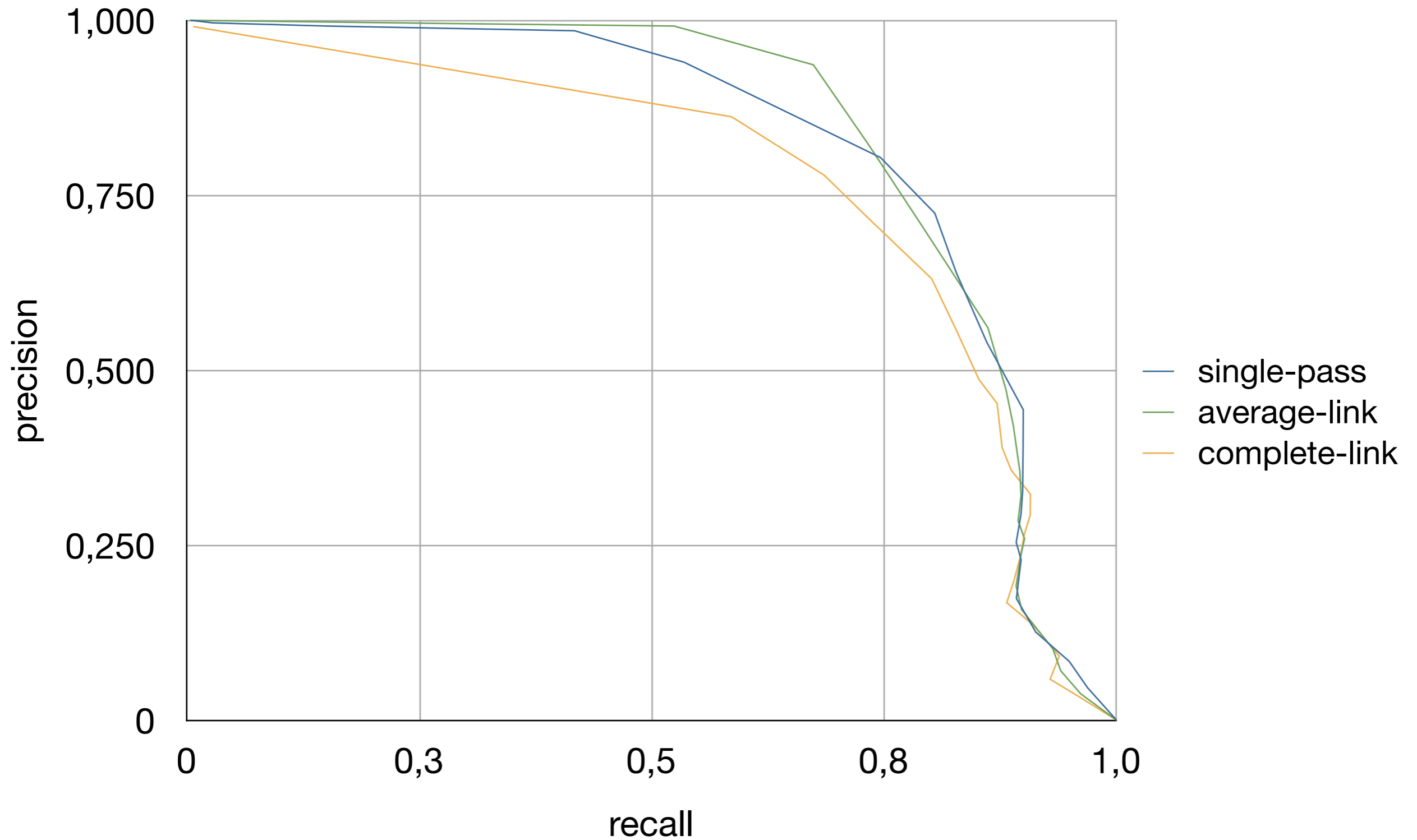| | |
|---|---|
| purity | $purity(\Omega,C) = \frac{1}{N}\sum_k \max |\omega_k \cap c_j|$ |
| normalized mutual information | $NMI(\Omega,C) = \frac{I(\Omega;C)}{[H(\Omega)+H(C)]\cdot 0.5}$ |
| rand index (accuracy) | $RI = \frac{TP+TN}{TP+FP+TN+FN}$ |
| precision, recall | $P = \frac{TP}{TP+FP} \quad R = \frac{TP}{TP+FN}$ |
| F$\beta$ | $F_\beta = \frac{(\beta^2+1)PR}{\beta^2 P+R}$ |

# Clustering

- partitional vs. hierarchical

- retrospective vs. on-line

- linkage criterion in hierarchical algorithms
  - maximum - complete-link
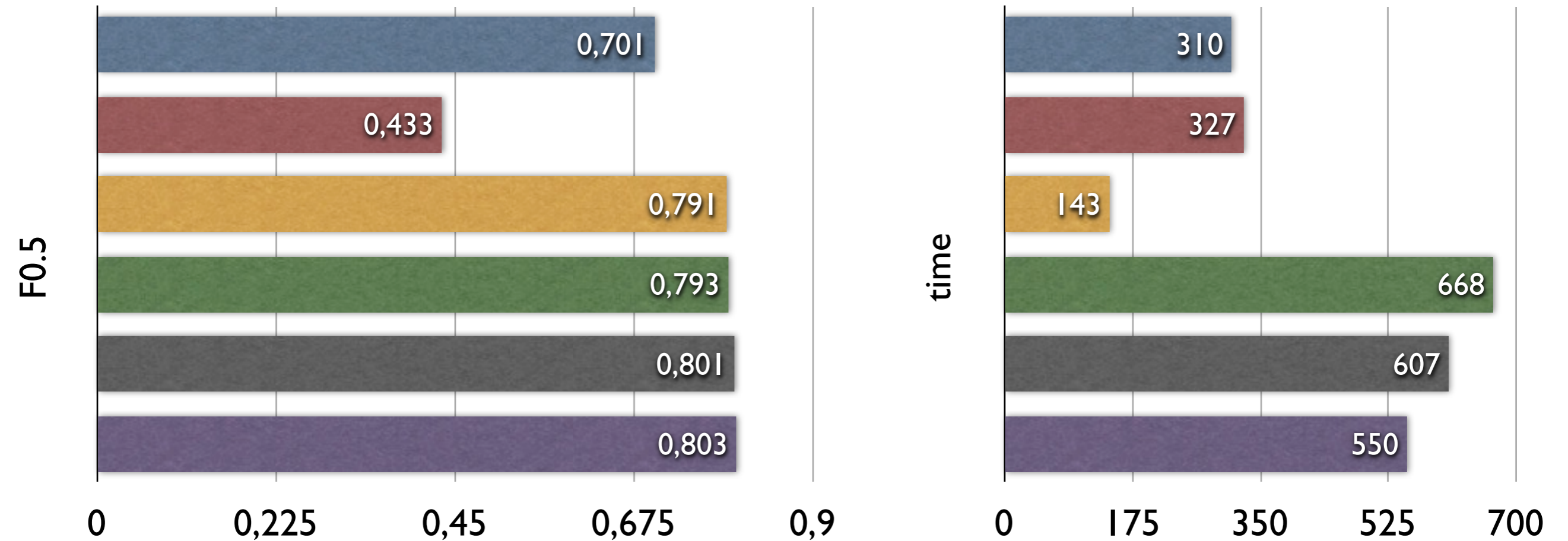  - minimum - single-link
  - mean - average-link

# Clustering algorithms

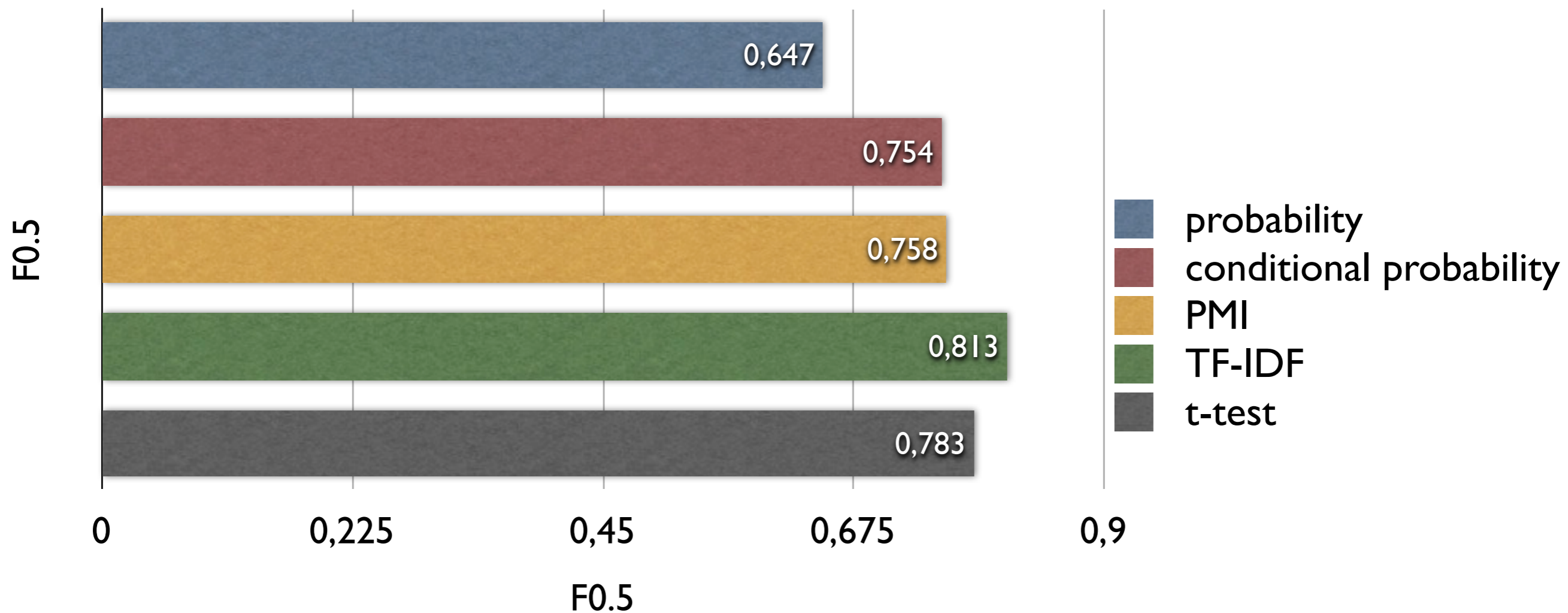| algorithm | linkage criterion | time complexity | on-line |
|---|---|---|---|
| hierarchical agglomerative | complete | $O(n^2 \log n)$ | no |
| hierarchical agglomerative | average | $O(n^2 \log n)$ | no |
| single-pass | single | $O(n)$ | yes |

# Precision-recall curve



precision (y-axis) vs recall (x-axis)

Legend:
- single-pass
- average-link
- complete-link

# Feature selection

- character case and punctuation obsolete

- information in title more relevant, optimal repetition rate is four

- function words (IDF) - minor decrease in model and memory complexity

- hapax legomena - decreases number of dimensions drastically, memory 5-10%
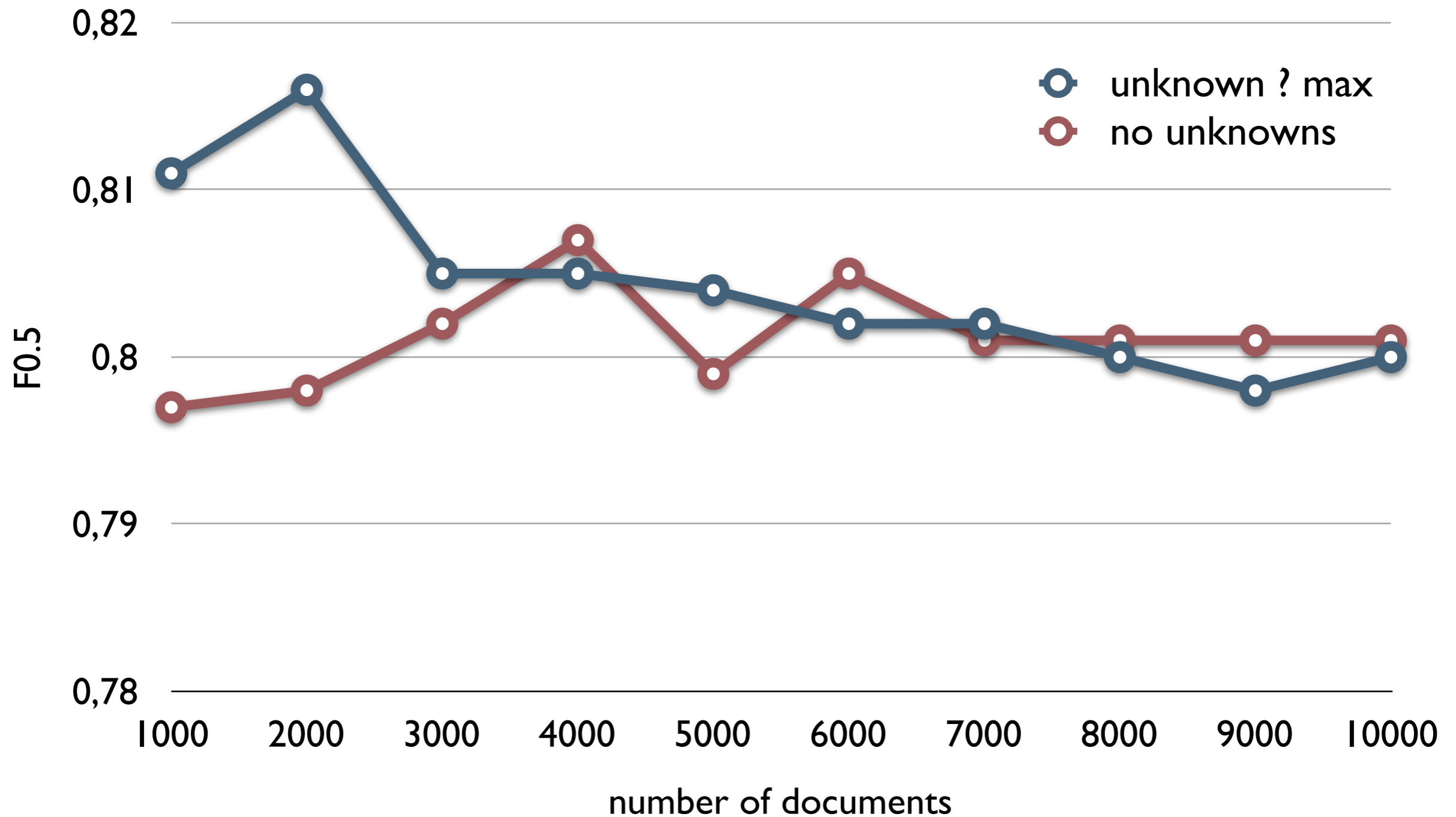
# Feature extraction

- stemming, POS tagging, lemmatization (two stemmers, TnT, HML)

- multi-word expressions (chi-square)

- named entity recognition (person and business entities)

- no significant improvement

# Heuristics

1. an event ranges on a one-day time span true in 83% of documents (non-singleton events)

2. one source reports only once about an event - true in 86% of documents (non-singleton events)

- implementing heuristics increases $F_{0.5}$, first heuristic simplifies calculation drastically

# Reference corpus



F0.5 versus number of documents, comparing "unknown ? max" and "no unknowns".

# Output example

Primorac saslušao studente (vijesti.hrt.hr)
Studenti nakon sastanka s Primorcem ipak ne odustaju od prosvjeda (index.hr)
Primorac pokušava izbjeći studentske prosvjede razgovorom s Rektorskim
zborom (business.hr)
Primorac primio organizatore studentskog štrajka (javno.com)
Studenti ne odustaju od najavljenog prosvjeda (dnevnik.hr)

VIDEO: Istukla i opljačkala susjedu zbog ljubomore (javno.com)
U stanu ju udarila palicom po glavi i opljačkala (index.hr)
Prijateljicu nevjenčanog supruga pretukla palicom i opljačkala (vecernji.hr)
Opalila je palicom u stanu i opljačkala (index.hr)

Sindikalna košarica u travnju 0,17 posto skuplja nego u ožujku (vecernji.hr)
Sindikalna košarica u travnju 0,17 posto skuplja (tportal.hr)
Sindikalna košarica u travnju 0,17 posto skuplja nego u ožujku (poslovni.hr)
Životni troškovi četveročlane obitelji 6206 kuna (business.hr)

# Further steps

- windowing technique $\Longleftrightarrow$ decay function

- feature position - features found at the beginning (in the first sentence?) should be given more weight

- write a Java API (Apache license)

- events $\Longrightarrow$ topics;
  event : document relationship

# Thank you!