

Offline Evaluation of Term Utility Functions

Alexander K. Seewald, Christian Holzbaur and Gerhard Widmer

Austrian Research Institute for Artificial Intelligence

Schottengasse 3, A-1010 Wien, Austria

[alexsee,christian,gerhard]@ai.univie.ac.at

Abstract

In this paper we investigate characterizing ontology nodes corresponding to human-comprehensible concepts from the tool Melvil by a set of terms. We choose a variety of term utility functions, commonly use in text mining, to determine relative importance of terms for the task of deciding if a given document is part of a certain concept or not. We evaluated each utility function both quantitatively by considering precision and recall of the top ten terms returned and qualitatively by analyzing which of the original patterns and obviously related terms were recovered. This approach could be used to suggest promising terms to a human ontology editor during creation of a new node. Our results look somewhat promising but still needful of improvement – so we also report on probable causes of unsatisfactory results.

1 Introduction

The tool melvil allows to create ontology nodes, each one based on a human-comprehensible concept, and organize them in an arbitrary hierarchy. Each ontology node has an associated regular expression pattern which is used to retrieve all relevant documents of the corresponding concept. A concept such as *Internet* may have a pattern such as `\binternet1\b|\bweb\b|\bwww\b`, which consists of several subpatterns separated by `|`. Not all of these patterns must be single words as in our example – multiword patterns also appear frequently, at least in the research ontology which was provided by uma. Notice also that `\b` is used to denote word boundaries – in our example, the subpatterns match whole words only.

In this paper, we investigate characterizing concepts by a short list of appropriate terms, i.e. single words. We consider a variety of term utility functions, each of which map every tuple (term, concept) to a numeric value which signifies the usefulness of the respective term to decide if documents are part of the respective concept or not.

Such a characterization by terms could also suggest new single word patterns to the user during creation of a new ontology node, if real-time performance is achievable.

We will first give an overview about the research ontology, March15, which was provided by uma, focussing on common data statistics. Afterwards, we will introduce the term utility functions used throughout our experiments. These are more concisely referred to as *measures*. Then we will shortly describe the experimental setup, discuss major experimental results in the Results section and discuss minor experimental results and other issues in the Discussion section. At last, we will conclude this paper.

¹i.e. the term *internet*.

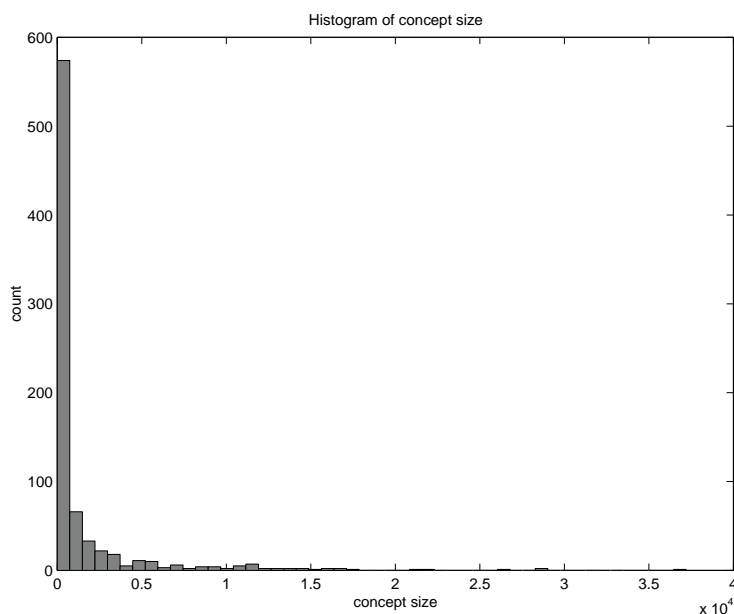


Figure 1: This is a histogram of concept size. Notice that the concept size is scaled by 10^4 , so a value of 1 stands for 10,000 documents on the x-axis.

2 Data Analysis

The March15 ontology contains 793 nodes resp. concepts. On average, these concepts contain $1,395.5 \pm 3,583.9$ documents – Figure 2 shows a histogram of concept size. The top ten concepts account for 230,161 documents². Since there are only 55,047 unique documents, this means that most documents are assigned to many concepts. The overlap³ for these ten concepts alone is already 4.18. About three quarters (76.0%) of all concepts contain less than a thousand documents; there are even 61 concepts which do not contain any documents! This exponential falloff in concept size is responsible for the mentioned abnormally high standard deviation.

Hierarchical relations between concepts can be modelled by the user – however, local regular expression at every node are used to determine the documents which are part of each node, independent of all other nodes; so this hierarchy is not relevant for the purposes of our report.⁴ Nevertheless, we have noted that the given hierarchy is not a forest⁵ as usually expected from an ontology but a acyclic⁶ graph, e.g. node *sms* has four parent nodes. Many nodes are thus reachable on multiple paths. Some top-level nodes such as *altavista*, *region* and *technology* do not contain any documents and seem to have been created simply as container for children nodes.

There are 55,047 unique documents altogether which usually belong to multiple concepts – the average overlap over all concepts is 20.08. The documents are partitioned into 187

²20.8% of the combined size of all concepts, which is 1,105,270.

³ $overlap = \frac{\sum conceptSize(c)}{number\ of\ unique\ documents}$

⁴The term *ontology* may thus be considered slightly inappropriate in the current implementation of melvil.

⁵Forest = a set of disconnected trees. A tree is a acyclic graph where every node has at most one parent node.

⁶Cycles seem possible in practice, if the ontology editor does not prevent this.

<i>conceptId</i>	<i>numDocs</i>	<i>numPatts</i>	<i>numKids</i>	<i>numParents</i>	<i>level(s)</i>
internet	37211	3	6	2	1,4
e_mail	28372	2	0	2	4,7
business unit	28360	4	2	1	1
search	26270	2	5	2	3,4,5,6,8
mobile	22035	4	9	1	2
media	21551	3	8	1	4,6,7
software	17119	1	6	1	2
mobile phone	16654	4	0	1	4
investment/investor	16462	2	2	1	4
network	16127	2	5	1	3,6
artificial intelligence	239	4	5	2	3,6
online community	350	3	0	1	5
manufacturer	379	4	4	1	3
wlan	974	3	0	1	4
market capitalisation	504	2	0	1	4
cryptography	831	3	0	1	4
mobile portal	474	4	13	2	3,5,6
online gaming	534	2	0	3	4,5,6,7
knowledge management	758	3	0	2	4,5,6,7,9
price earning ratio	392	4	0	1	4

Table 2: This table shows the ontology nodes resp. concepts chosen for our experiments. The top half shows the large concepts, the bottom half the small ones. *numPatts* is the number of patterns used for indexing. *numKids* and *numParents* shows the number of children resp. parent nodes. Level(s) shows the zero-based level index for each concept – if a node is accessible via multiple paths and/or multiple root nodes, more than one level has to be shown.

- *oddsRatio*, which is a commonly used feature in information retrieval (van Rijsbergen, Harper & Porter, 1981). In our case, when removing the logarithm which is irrelevant for relative ranking of terms, this simplifies to $\frac{a \cdot d}{b \cdot c}$.
- *odds2* is one of the many measures inspired by the original Odds Ratio formula, i.e. $\frac{a+c}{N} \log_2 \frac{ac+ad}{ac+bc}$ where $N = a + b + c + d$ is the total number of documents. It is equivalent to *FreqLogP* in (Mladenic, 1998). *odds2* and *oddsRatio* were two of three measures found to be superior in the mentioned paper. The third measure was based on exponentiation and deemed to be too costly for implementation.
- Precision (*prec*) – the ratio of documents belonging to the concept among all documents which include the term, i.e. $\frac{a}{a+c}$.
- Recall (*recall*) – the ratio of documents which include the term, among all documents belonging to the concept, i.e. $\frac{a}{a+b}$.
- *prec * recall (PR)* which trades off recall and precision. Usually, we want both high precision and high recall – our formula is a simple way to capture this relation.
- SimpleRatio1 (*sR1*) is $\frac{a}{N}$ where $N = a + b + c + d$ is the overall number of documents.
- SimpleRatio2 (*sR2*) is $\frac{s \cdot f \cdot C}{N}$, which prefers those terms which appear very frequently.
- SimpleRatio3 (*sR3*) is $\frac{s \cdot f \cdot C}{\neg s \cdot f \cdot C + 1}$, which prefers those terms appearing frequently within the concept, but seldom without.

In initial experiments we found that *prec* is unusable since about 20% of the terms we looked at have the maximum precision of $prec = 1.0$ ($a > 0$ and $c = 0$), which makes it impossible to determine a stable relative ranking!¹³ We also found that *sR1* and *sr2* perform very poorly, mostly selecting stop words. So we removed the three mentioned measures which still leaves us with seven measures for offline our evaluation. A simplistic way to evaluate these term utility measures is to look at which terms correspond to the indexing patterns. But since using single terms instead of regular expressions is a crude approximation at best, some precision is inevitably lost. So we used precision and recall, averaged over the top ten terms, as fairer evaluation.

4 Experimental Setup

As we already mentioned, due to resource constraints we used only the largest ten jobs¹⁴ in our experiments. These account for 69% of all unique documents. For the same reason we were unable to evaluate all 793 concepts, so we chose twenty concepts for further investigation: the ten largest concepts by size and also ten smaller concepts with 200-1000 indexed documents which were arbitrarily chosen. Details can be found in Table 3.

For these twenty concepts, we computed the contingency tables plus *sfC* and $\neg sfC$ for each term, by summing over all chosen jobs.¹⁵ This data was extracted from Melvil and then imported into MATLAB for further offline analysis, including computing all our term utility functions, statistical analysis and visualization. In MATLAB, computing all seven measures for ten concepts and all 428,173 terms took about thirty seconds.

5 Results

We first determined the top ten terms for each measure – these are shown in Tables 8 and 8 in the Appendix. Terms which are matched by the original indexing patterns are shown in **bold**. Multiword patterns were broken up into single word patterns at every place where a word boundary may appear, e.g. `e.?mail` maps to `e|mail|email`. For all but one concepts, at least one indexing pattern is recovered. Detailed qualitative results are discussed in the next section.

We then computed precision and recall for each concept and measure, averaged over the top ten terms. The complete results can be found in Table 5. For comparing the different measures at one glance, we have also computed the grand mean and standard deviation of recall and precision over large and small concepts separately, these are visualized in Figure 4. High precision and high recall – which is what we usually want – is found towards the top right of each graph. Detailed plots which show the top ten terms for each concept and measure can be found in the Appendix.

All measures seem to work quite well on our ten large concepts, but perform more poorly on the ten small ones. However, the relative order between measures is still quite similar – two adjacent measures ($IG=\circ$, $PR=\square$) switch places and *odds2* (*) moves from top left (high precision, low recall) to bottom right (low precision, high recall).

Using *PR*, i.e. *prec* multiplied by *recall* seems to work quite well, so we have extended the formula to $prec^x * recall^{2-x}$ and plotted grand mean precision and recall for ten different

¹³I.e. we would have to randomly choose ten best terms out of about 80,000 equally good ones.

¹⁴jobIds = 0,14,100,91,97,16,13,94,92,2

¹⁵We also investigated normalizing the job size prior to combining the contingency tables – the results were mixed, see Discussion.

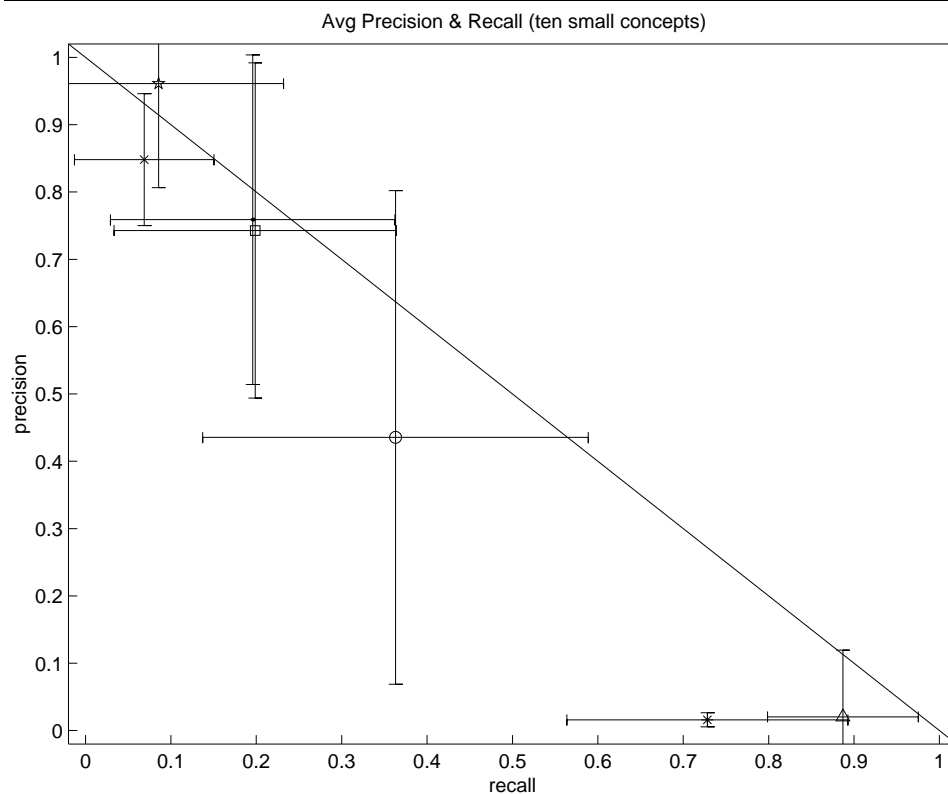
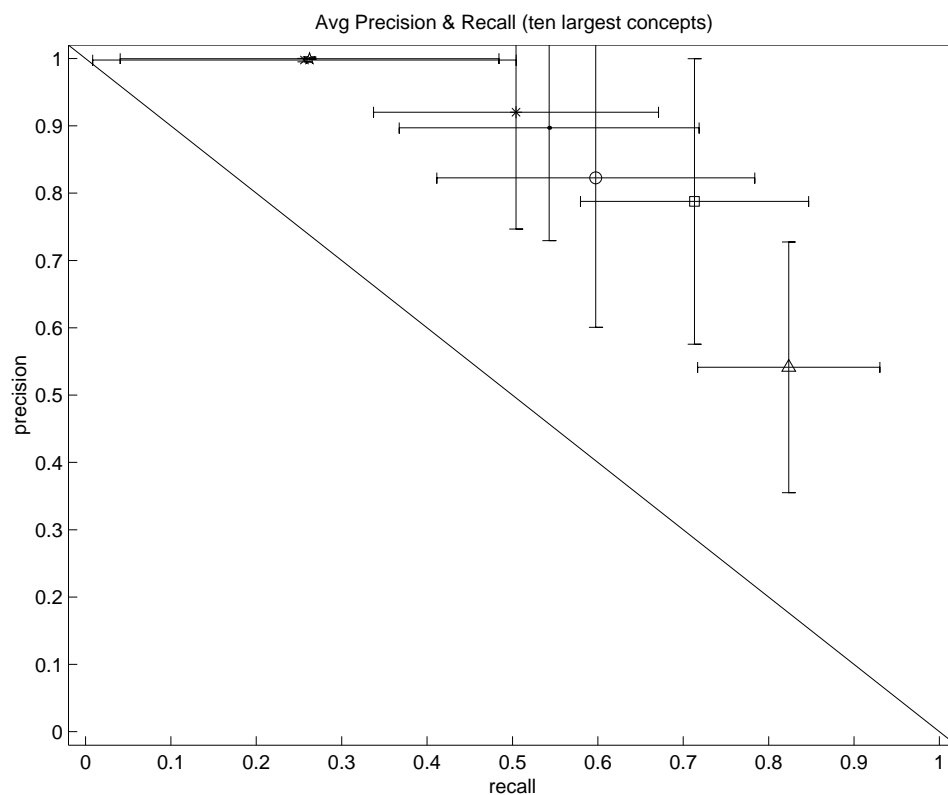


Figure 2: These plots show the average precision and recall \pm standard deviation for the ten largest concepts (top) and our ten small concepts (bottom) – one point from each measure, averaged over the top ten terms. • = χ^2 , ◦=IG, ×=oddsRatio, *=odds2, ◻=PR, Δ=recall and *sR3.

values of x in Figure 5. This gives competitive results and also allows for user-defined explicit trade-off between precision and recall. For $x = 0$ the formula is equivalent to *recall* and for $x = 1$ to *PR*. x should always be smaller than 2 because then the formula would be equivalent to *prec* and $prec = 1.0$ for about 20% of all terms as we mentioned earlier, so the obtained ranking would be quite random.

6 Discussion

In Table 8 concerned with the large concepts, some indexing patterns are recovered and usually appear near the top. Also, some obviously related patterns are found, e.g. *sms* and *wap* for concept handy, *kapital* for business unit and *microsoft* for software which is even found by more than half of our measures, always on second rank after *software* itself. For concept search, no indexing pattern could be recovered, but at least some related terms such as *recherche* and *stichwortsuche* are found.

In Table 8 concerned with the small concepts, we see a similar picture: some indexing patterns are recovered – sometimes even quite many e.g. for concepts manufacturer and cryptography. Related terms are also quite apparant: e.g. *blackbox* for concept community, *802 11a*¹⁶ for concept WLAN, many names of popular online games for concept online gaming and *kampfroboter*, *wunderwaffe* and *arbeitsklaven* for concept artificial intelligence. We were also quite surprised to notice *trappl* as third-most relevant term for concept artificial intelligence by measure *sR3* which otherwise performs poorly. Another less pleasant surprise was that the term *fuck* appears on second place in three measures for concept online community. When we look at precision and recall, we see that other top terms such as *blackbox* ($p=0.82, r=0.18$), *community* ($p=0.71, r=0.14$) and *fuck* ($p=0.72, r=0.19$) seem quite similar. Although community was used as indexing pattern for this concept, its precision is less than 1.0 – another indication that the full text index does not index whole words as we mentioned earlier.

When we normalized the job size before combining the respective contingency tables, we encountered mixed results: On the one hand, the terms corresponding to indexing patterns sometimes move more to the top¹⁷, and sometimes one or more additional indexing pattern are found¹⁸ On the other hand, up to four indexing patterns¹⁹ are lost and many of the related words disappear. It seems that this variation works better on the large concepts than on the small ones. Details can be found in Tables 8 and 8. The average precision and recall is very similar between normalized and non-normalized case – in fact, the grand average precision and recall plots were so similar to Figure 5 as to be almost indistinguishable. So we are inclined to prefer the original non-normalized version by Ockham’s Razor.

For our experiments, computing the contingency table, sfC and $\neg sfC$ took an average of 6.75 minutes per concept. To achieve realtime performance, a computational speedup of at least two orders of magnitude is necessary. Notice also that we only used the top ten jobs in our experiments instead of all 187 which may be problematic when very heterogenous document sources are employed.

A possible way to achieve this would be to precompute these values efficiently. This would also results in an additional advantage would be that from our six values, many basic statistics such as precision, recall, overall term frequency and term frequency within con-

¹⁶IEEE 802.11a WLAN standard

¹⁷mobile, media, network, price earning ratio

¹⁸e_mail, investor, artificial intelligence, search

¹⁹for concept manufacturer and measure χ^2

Measure	conceptId	avg prec.	avg recall
χ^2	internet	0.947±0.041	0.457±0.160
	e_mail	0.765±0.296	0.550±0.237
	business unit	0.979±0.060	0.426±0.163
	search	0.900±0.079	0.642±0.129
	mobile	0.963±0.053	0.618±0.060
	media	0.958±0.060	0.581±0.123
	software	0.625±0.131	0.513±0.184
	mobile phone	0.985±0.021	0.650±0.091
	investment/inv.	0.855±0.209	0.374±0.238
	network	0.993±0.010	0.619±0.048
IG	internet	0.934±0.045	0.479±0.148
	e_mail	0.765±0.296	0.550±0.237
	business unit	0.687±0.215	0.719±0.261
	search	0.803±0.238	0.632±0.211
	mobile	0.929±0.074	0.644±0.050
	media	0.955±0.059	0.583±0.122
	software	0.602±0.142	0.551±0.171
	mobile phone	0.963±0.068	0.667±0.094
	investment/inv.	0.614±0.306	0.523±0.281
	network	0.975±0.035	0.627±0.048
oddsRatio	internet	0.999±0.000	0.126±0.100
	e_mail	0.998±0.000	0.074±0.075
	business unit	0.999±0.000	0.262±0.238
	search	0.999±0.000	0.267±0.134
	mobile	0.997±0.001	0.274±0.247
	media	0.999±0.001	0.329±0.201
	software	0.994±0.005	0.115±0.299
	mobile phone	0.999±0.002	0.568±0.189
	investment/inv.	0.995±0.003	0.121±0.200
	network	0.997±0.002	0.426±0.275
odds2	internet	0.995±0.005	0.330±0.170
	e_mail	0.976±0.026	0.443±0.171
	business unit	0.998±0.002	0.397±0.149
	search	0.994±0.011	0.439±0.140
	mobile	0.991±0.007	0.572±0.079
	media	0.997±0.002	0.500±0.051
	software	0.580±0.144	0.572±0.154
	mobile phone	0.998±0.003	0.617±0.093
	investment/inv.	0.683±0.268	0.553±0.256
	network	0.991±0.011	0.619±0.048
PR	internet	0.777±0.110	0.733±0.145
	e_mail	0.802±0.160	0.761±0.154
	business unit	0.702±0.104	0.857±0.022
	search	0.822±0.135	0.736±0.120
	mobile	0.946±0.067	0.633±0.054
	media	0.894±0.120	0.633±0.123
	software	0.493±0.195	0.728±0.178
	mobile phone	0.963±0.068	0.667±0.094
	investment/inv.	0.485±0.191	0.763±0.138
	network	0.993±0.010	0.619±0.048
recall	internet	0.745±0.120	0.745±0.133
	e_mail	0.722±0.165	0.789±0.116
	business unit	0.635±0.053	0.888±0.051
	search	0.709±0.177	0.807±0.113
	mobile	0.440±0.086	0.824±0.088
	media	0.536±0.154	0.825±0.095
	software	0.408±0.205	0.793±0.128
	mobile phone	0.392±0.039	0.919±0.035
	investment/inv.	0.405±0.077	0.797±0.110
	network	0.420±0.092	0.849±0.079
sR3	internet	1.000±0.000	0.073±0.033
	e_mail	1.000±0.000	0.042±0.028
	business unit	1.000±0.000	0.344±0.015
	search	1.000±0.000	0.364±0.016
	mobile	0.999±0.001	0.148±0.218
	media	1.000±0.000	0.477±0.021
	software	1.000±0.000	0.020±0.000
	mobile phone	1.000±0.000	0.582±0.026
	investment/inv.	0.999±0.001	0.099±0.092
	network	0.999±0.002	0.472±0.216

Measure	conceptId	avg prec.	avg recall
χ^2	art.int.	0.759±0.238	0.179±0.130
	online c.	0.615±0.196	0.217±0.174
	manufacturer	0.625±0.367	0.186±0.188
	wlan	0.836±0.202	0.212±0.153
	market c.	0.611±0.256	0.186±0.236
	cryptography	0.956±0.125	0.117±0.102
	mobile p.	0.892±0.142	0.133±0.061
	online g.	0.800±0.097	0.305±0.026
	know. man.	0.805±0.284	0.132±0.208
	price earn.r.	0.690±0.237	0.291±0.214
IG	art.int.	0.577±0.327	0.210±0.120
	online c.	0.548±0.217	0.226±0.170
	manufacturer	0.170±0.295	0.488±0.196
	wlan	0.533±0.350	0.320±0.178
	market c.	0.137±0.303	0.704±0.051
	cryptography	0.441±0.416	0.260±0.132
	mobile p.	0.487±0.404	0.225±0.126
	online g.	0.723±0.221	0.347±0.104
	know. man.	0.338±0.419	0.273±0.200
	price earn.r.	0.401±0.364	0.579±0.267
oddsRatio	art.int.	0.783±0.071	0.073±0.121
	online c.	0.829±0.044	0.064±0.053
	manufacturer	0.861±0.029	0.024±0.007
	wlan	0.933±0.016	0.104±0.137
	market c.	0.869±0.031	0.031±0.039
	cryptography	0.958±0.020	0.043±0.016
	mobile p.	0.652±0.093	0.121±0.031
	online g.	0.943±0.028	0.136±0.095
	know. man.	0.787±0.034	0.024±0.014
	price earn.r.	0.867±0.032	0.068±0.094
odds2	art.int.	0.005±0.001	0.510±0.149
	online c.	0.008±0.002	0.709±0.145
	manufacturer	0.016±0.011	0.840±0.103
	wlan	0.024±0.011	0.715±0.118
	market c.	0.026±0.004	0.777±0.024
	cryptography	0.028±0.005	0.830±0.117
	mobile p.	0.001±0.000	0.717±0.131
	online g.	0.018±0.003	0.796±0.107
	know. man.	0.011±0.003	0.512±0.145
	price earn.r.	0.021±0.005	0.879±0.041
PR	art.int.	0.759±0.238	0.179±0.130
	online c.	0.615±0.196	0.217±0.174
	manufacturer	0.542±0.366	0.204±0.181
	wlan	0.808±0.196	0.216±0.150
	market c.	0.611±0.256	0.186±0.236
	cryptography	0.956±0.125	0.117±0.102
	mobile p.	0.842±0.182	0.142±0.060
	online g.	0.800±0.097	0.305±0.026
	know. man.	0.805±0.284	0.132±0.208
	price earn.r.	0.690±0.237	0.291±0.214
recall	art.int.	0.003±0.000	0.815±0.089
	online c.	0.005±0.001	0.854±0.105
	manufacturer	0.009±0.001	0.949±0.026
	wlan	0.016±0.001	0.893±0.048
	market c.	0.112±0.312	0.887±0.064
	cryptography	0.023±0.003	0.954±0.017
	mobile p.	0.001±0.000	0.883±0.081
	online g.	0.015±0.005	0.887±0.076
	know. man.	0.006±0.001	0.776±0.096
	price earn.r.	0.012±0.002	0.972±0.016
sR3	art.int.	0.998±0.007	0.086±0.126
	online c.	0.943±0.181	0.019±0.010
	manufacturer	1.000±0.000	0.075±0.114
	wlan	1.000±0.000	0.081±0.022
	market c.	1.000±0.000	0.092±0.264
	cryptography	1.000±0.000	0.107±0.107
	mobile p.	0.875±0.270	0.087±0.080
	online g.	0.995±0.011	0.117±0.035
	know. man.	0.892±0.249	0.093±0.218
	price earn.r.	0.908±0.262	0.099±0.250

Table 3: This shows the average precision and recall over the top ten terms for the ten largest concepts (left) and our arbitrarily chosen ten small concepts (right)

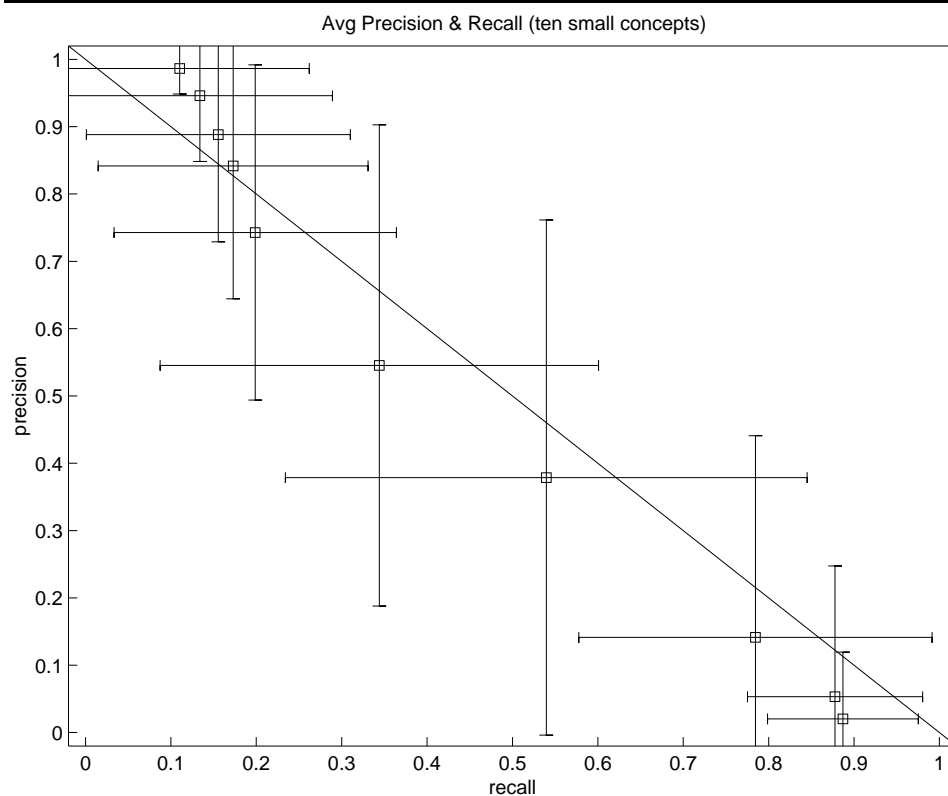
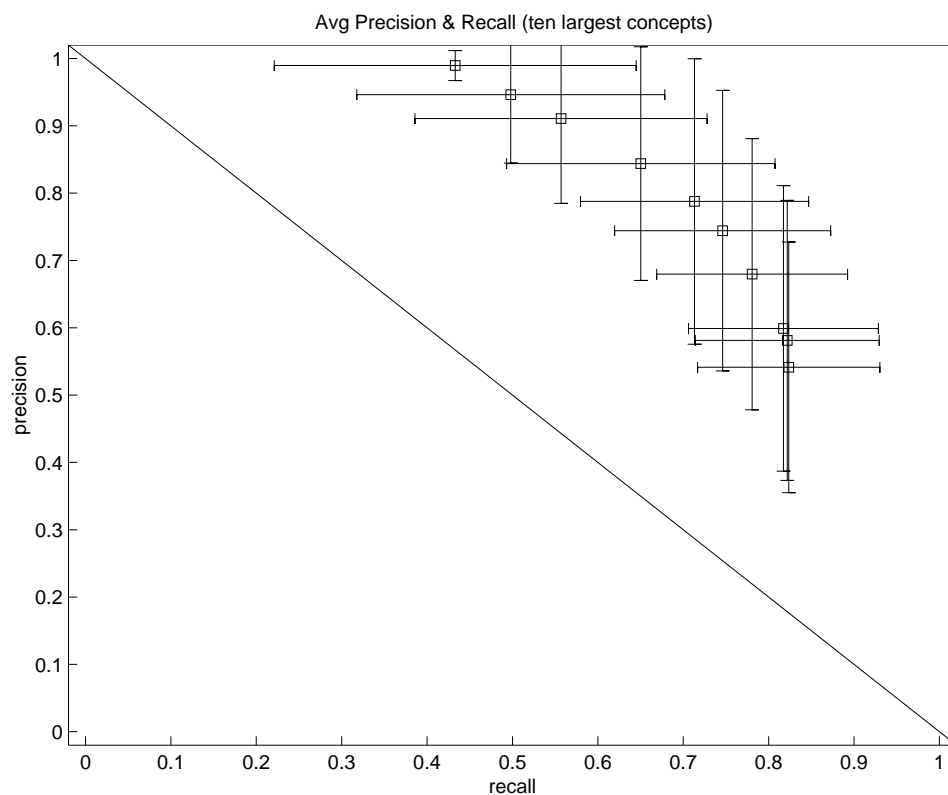


Figure 3: These plots show the average precision and recall \pm standard deviation for the ten largest concepts (top) and our ten small concepts (bottom). The measures used here are $prec^x * recall^{2-x}$ where $0 \leq x < 2$ and steps of 0.2 were used, yielding ten points. $x = 0 \equiv recall$, $x = 1 \equiv PR$.

cept etc. can be instantly computed, thus allowing to give the user instant feedback during ontology editing.

Another idea would be to reduce the vocabulary to more manageable level – e.g. removing those terms which appear in less than ten documents would reduce the number of terms by an order of magnitude.

7 Conclusion

As we have seen, it is in principle possible to characterize ontology nodes by single words even in the presence of multi-word patterns. The results are somewhat promising, but in need of improvement. Our research would benefit from the following, roughly in order of importance.

- Reducing the terms by an order of magnitude. E.g. removing those terms which appear in less than 0.02% of all documents would be a feasible option.
- A Full-Text Index which scans for words and not for substrings would enable a more precise estimation of recall and precision.
- Precomputed contingency tables for each term and concept definition – or a way to compute these tables fast enough for realtime feedback.
- For people starting new ontologies with single word patterns, precomputed contingency tables for each combination of two terms would also be useful to approximate both multi-word patterns and multiple single word patterns. However, the memory requirements make this quite unfeasible – e.g. even for just 40,000 terms this would mean $6 * (40,000)^2 = 9.6 * 10^9$ values 40 Gigabytes memory consumption.

Acknowledgements

We want to thank Reinhard Schwab for optimizing our experimental java code and for giving valuable hints concerning runtime improvements. We also want to thank an anonymous colleague for being on vacation, so we could use his much faster computer for our evaluation of results.

References

- Mladenic, D., (1998) Feature subset selection in text-learning, Proceedings of 10th European Conference on Machine Learning, 1998.
- van Rijsbergen, C.J., Harper, D.J., Porter, M.F., The selection of good search terms, *Information Processing and Management*, 17, pp.77–91, 1981.
- Yang, Y., Pedersen J.P. (1997). A Comparative Study on Feature Selection in Text Categorization Proceedings of the Fourteenth International Conference on Machine Learning (ICML'97), pp.412-420.

8 Appendix

conceptId	χ^2	IG	oddsRatio	odds2	PR	recall	sR3
internet	internet com web it contact from copyright privacy about sites	internet com web it contact from copyright privacy for about	lebensversicherers mamax mysimon techrepublic gamespot shopper corrections aktivierungsschlüssel mediadaten bandwidth	internet web mysimon techrepublic gamespot corrections comparisons contact sites reviews	internet in 2001 com an 1 online jobs at it	in 2001 an internet 1 com at jobs online die	dispatch newshilght catchup newslinksmall boldlink storyurl openmpwin bigopen scrollable gatrixx
e_mail	mail version e jobs td kontakt archiv latest tools color	mail version td e jobs kontakt archiv color latest tools	9pt fußballverein lebensversicherers mamax hoppenstedt marketperformer javaenabled javaon gflshplugintargetversion gflshpluginname	mail kontakt version hoppenstedt email leserbrief drucken verschicken archiv kurssuche	mail version jobs e in 2001 news kontakt an archiv	in 2001 an jobs e version news mail 1 die	boldlink companyname smalltext mediumbold 7pt smalltext2 newmediasales stammzellen ft26xx3044x11 ft29xx3144xx1
business unit	unternehmen weitere kapital fd partnersites writelayersn writerecherchen writzeitungn writemultichanneln writeregistriern	unternehmen weitere contact das für von das mit kontakt der und	leserbriefe unternehmen creditreform stichwortsuche icra graumarkt marketperformer mdax lbbw unternehmen	unternehmen leserbriefe creditreform stichwortsuche graumarkt fd guided persönliches zeitungs neuemissionen	unternehmen der die und von das für mit den im	2001 in die der und von mit den an	börsenticker hoppenstedt analystenstimmen munzinger partnersites leserbriefe sitemap logout menusactive writelayersn
search	tools latest news pda times mail home edition jobs fi nancial	news tools latest pda mail home times e td	leserbriefe partnersites research dispatch analystenstimmen icra guided graumarkt marketperformer creditreform	leserbriefe partnersites latest analystenstimmen guided graumarkt recherche tools creditreform topics	news tools latest jobs e mail pda home in service	in 2001 an news jobs e us tools latest 1	börsenticker hoppenstedt munzinger stichwortsuche sitemap logout menusactive writelayersn writerecherchen writzeitungn
mobile	mobile policy headlines pt featured 1995 topics search wireless services	mobile policy headlines pt search featured services 1995 privacy about	topics wireless 5s mobile catchup openwindow unescape getsurveyfi le surveyfi le pdt	mobile topics wireless headlines pt featured policy corrections mysimon gamespot	mobile policy headlines pt featured search 1995 services topics privacy	in 2001 a us com top at an s new	dispatch emazing writeexchangelink metricom convertedsymbol nextcard wireless xis topics aprs
media	media medien edition times true partnersites writelayersn writerecherchen writzeitungn writemultichanneln	media medien edition times true überblick partnersites writelayersn writerecherchen writzeitungn	leserbriefe analystenstimmen stichwortsuche graumarkt marketperformer creditreform realmedia redir guided useragent	leserbriefe medien analystenstimmen stichwortsuche graumarkt creditreform guided zeitungs fd persönliches	media medien edition times true version home überblick partnersites writelayersn	in 2001 an media news version jobs document write die	börsenticker hoppenstedt munzinger partnersites sitemap logout menusactive writelayersn writerecherchen writzeitungn
software	software microsoft pc web desktop product this sites computer reserved	software microsoft pc web desktop product it this computer services	software untermieter klingeling fi bel bauzulieferer malade papierchen glitzernde lpez schlafmützen	software microsoft web pc it services this to tech computer	software microsoft web pc it in 2001 com a for	in 2001 software an at new us com 1 a	geldverwalter witwenrente sponsort hilf freudentänze papiertiger fondspolice goldies erdrückende konjunkturgewitter
mobile phone	handy stellenmarkt branchen sms wap writelayersn writerecherchen writzeitungn writemultichanneln writeregistriern	handy stellenmarkt branchen sms abonnieren wap persönliches writelayersn writerecherchen writzeitungn	kurssuche börsenticker leserbriefe partnersites analystenstimmen stichwortsuche graumarkt creditreform marketperformer handy	kurssuche börsenticker leserbriefe partnersites analystenstimmen handy stichwortsuche graumarkt creditreform guided	handy stellenmarkt branchen sms wap abonnieren writelayersn writerecherchen writzeitungn writemultichanneln	in 2001 die der und das an für mit den	hoppenstedt munzinger sitemap logout menusactive writelayersn writerecherchen writzeitungn writemultichanneln writeregistriern
investment/investor	investor cfd6b2 9ca380 rundschau dispatch a line context weekly starting	investor a line top new s context cfd6b2 9ca380 http	dispatch investment todays personalia investor lebhaft drahtseilakt reichtmacher sponsort untermieter	investor dispatch a rundschau line top investment s new time	investor a top line new s in 2001 us time	in 2001 an a new top at us s line	9ca380 cfd6b2 wirtschaftit dispatch tandard investment maxi geldverwalter schwalben witwenrente
network	networks hot mysimon corrections comparisons techrepublic gamespot reserved reviews send	networks hot featured reserved topics send comparisons mysimon reviews corrections	corrections techrepublic gamespot comparisons networks reviews network kary heavenly netzwerk	corrections techrepublic gamespot networks comparisons reviews reserved send zdnet hot	networks hot reserved mysimon comparisons corrections techrepublic gamespot reviews send	2001 in us news com a s all for new	mysimon dispatch catchup techrepublic gamespot corrections networks comparisons 5s network

Table 4: This table shows the top ten terms according to each measure for the largest ten concepts. The first (uppermost) term has rank one, the next lower one has rank two and so on.

conceptId	X ²	IG	oddsRatio	odds2	PR	recall	sR3
artificial intelligence	ai intelligenz ki artificial künstliche seminarvorträge kampfroboter wunderwaffe arbeitsklaven privatstiftungen	ai intelligenz artificial künstliche ki intelligence seminarvorträge privatstiftungen künstlicher wölf	ai ziegelindustrie wunderwaffe arbeitsklaven crossbar turing milky descendant fluents blinding	man height ist eine kann können bin online welt werden	ai intelligenz ki artificial künstliche seminarvorträge kampfroboter wunderwaffe arbeitsklaven privatstiftungen	2001 in der die und von den das für mit	ki ai trapp dunietz hutchens vermobil seminarvorträge goren treister dfki
online community	blackbox fuck zwischenablagen playstations handtasche elektron zumindestens community mond minigehäuse	community fuck blackbox playstations zumindestens mond schett zwischenablagen schrott handtasche	oberschalen kemco pagern dreiste tabellenkalkulationen blackbox heimkehrende spatnik7 zwischenablagen faxnachrichten	internet online web software microsoft pc net at on 5	blackbox fuck zwischenablagen playstations handtasche elektron community zumindestens mond minigehäuse	2001 in internet online an at l us com e	musicity konop pocketmail vitos mnet phonoverband markle wissenszentren thielen mcowen
manufacturer	handyhersteller mobiltelefonhersteller sendo handyherstellern mangment mobiltelefone zusammenschaltung aktienverkauf duftchip handys	handyhersteller handys nokia hersteller ericsson mobiltelefone handy mobilfunk gprs fi nnsiche	z100 bistum netzbau farbiges fi m jussi mobiltelefongeschäft telefonino abbauten typografi e	handy als für markt document das nokia bei mit den	handyhersteller mobiltelefonhersteller sendo handyherstellern mobiltelefone handys mangment zusammenschaltung aktienverkauf fi nnsiche	2001 in der die für und das mit von den	handyhersteller mobiltelefonhersteller handyherstellern deighton pretec duftchip adventskalender würzburg besinnlichen wohlgeruch
wlan	wlan lan 802 11a bankcomputer grafikchips subnotebook kummernummer kummermail radau	wlan lan 802 11a wireless grafikchips lans bluetooth notebooks bankcomputer	wlan computersäulen döw überwachungsservice methangas funklösung internetexplorer computerschrott mönche mov	wireless parent neue microsoft wird für den var mit ist	wlan lan 802 11a bankcomputer grafikchips subnotebook kummernummer kummermail komma	2001 in die der und den mit für von das	bankcomputer kummernummer kummermail netstumbler radau wlans prozessortechnologie informatikbegriffe betastadium kuten
market capitalisation	marktkapitalisierung capitalization rangliste zwischenbericht autokonzerne oica vda monatszahlen reinking börsenumsatz	marktkapitalisierung börsen neuemissionen persönliches sitemap logout munzinger writelayersn writerecherchen writezeitungen	capitalization angeschwollen umschreibung oica wertpapiermärkten schwieg zeichnungsgewinn auswahlkriterium ibiza bewertungs	fi nancial latest tools aktien times euro 7 unternehmen pda home	marktkapitalisierung capitalization rangliste zwischenbericht autokonzerne oica vda reinking monatszahlen börsenumsatz	2001 in an die der und mit den für markt	marktkapitalisierung kursänd behrnt die furse kleiman goy skyy thelemann indexrevision
cryptography	verschlüsselung verschlüsseln verschlüsselte verschlüsselt verschlüsselungs verschlüsseltem	verschlüsselung verschlüsseln netzpolitik webstandard hacker mails fi rewall verschlüsselte nsa umfrage	animator freundet maßvoll virenliste winlinux internetüberwachung maschen inkompatiblen musikfan mobilfunkvertrag	internet ist werden auf date für oder das den vor	verschlüsselung verschlüsseln verschlüsselte verschlüsselt verschlüsselungs verschlüsseltem animator freundet kompostierbaren	die in und der 2001 für von den das auf	verschlüsselung verschlüsseln verschlüsselt verschlüsselte verschlüsselungs verschlüsseltem kryptographie kompostierbaren verschlüsselungssoftware bildschirmchoners
mobile portal	mobifunkportal contentbereich betriebsystemunterstützung travelchannel unterhaltsames gevey internetfernsehen datenfähigkeit communitysektor umsatzsteigerung	mobifunkportal jamba travelchannel kundrun gerätes wirtschaftspresse stadtplan contentbereich betriebsystemunterstützung portal	unterhaltsames gevey internetfernsehen travelchannel skywire handyabsatzes playerlösung sugarman kundrun faircar	handy com de mobile soll einer internet i an 2002	mobifunkportal contentbereich betriebsystemunterstützung travelchannel unterhaltsames gevey internetfernsehen datenfähigkeit kundrun communitysektor	in an 2001 com der den und das die von	mobifunkportal microphones contentbereich netsize betriebsystemunterstützung targus datenfähigkeit uninstalled chesnais wichmann
online gaming	startopia unreal xtreme racer dungeon eidos siege quake commandos desperados	games unreal startopia xtreme dungeon racer siege eidos desperados kabuto	spielehits undying barkers rennbahn einköpfige cybersex interplay gehirne startopia xtreme	online the of top line a microsoft 4 to was	unreal startopia xtreme racer dungeon eidos siege quake commandos desperados	in online 2001 the top a at of an new	onlinespiele teken elektroschocks flirtspielchen biershooter nashorn fallschirmspringen spielehits multimedibrille barkers
knowledge management	km wissensmanagement karraker kmart richtfunkantenne dissuaded wmt sensational cmb begriffswelt	km wissensmanagement kmart bluelight mart email karraker knowledge ansprechpartner druckerfreundliche	sensational cmb karraker ambulanten jcp troger sprachsoftware kommunikationsebene sprachlösungen glasfaserleitung	email online com dieser version meldungen seite wien pte produkte	km wissensmanagement karraker kmart richtfunkantenne dissuaded wmt sensational cmb begriffswelt	in 2001 l at online com an die und von	km wissensmanagement artifact x83 langlauf nbsp artibrain airlancer audigy siriri
price earning ratio	kgv aktienrückkäufe rechtfertigen sommerrally lignum brainlab winkt frauenförderung gewinnhalbierung deepwater	kgv rechtfertigen aktienrückkäufe lignum sommerrally brainlab fonds kurse aktien börsen	gewinnwarnungsreigen emagine aktienrückkäufe irrglaube stammzelltherapien profi tablere gewinnhalbierung bewertungs hochkapitalisierten rückschlagpotenzial	wir aktien unternehmen uns mehr zum aus über heute bei	kgv aktienrückkäufe rechtfertigen sommerrally lignum brainlab winkt frauenförderung gewinnhalbierung deepwater	2001 in die der und für von auf bei im	kgv frauenförderung kursänd punktabzug skyy steinemann prozeß gejamm emagine jurka

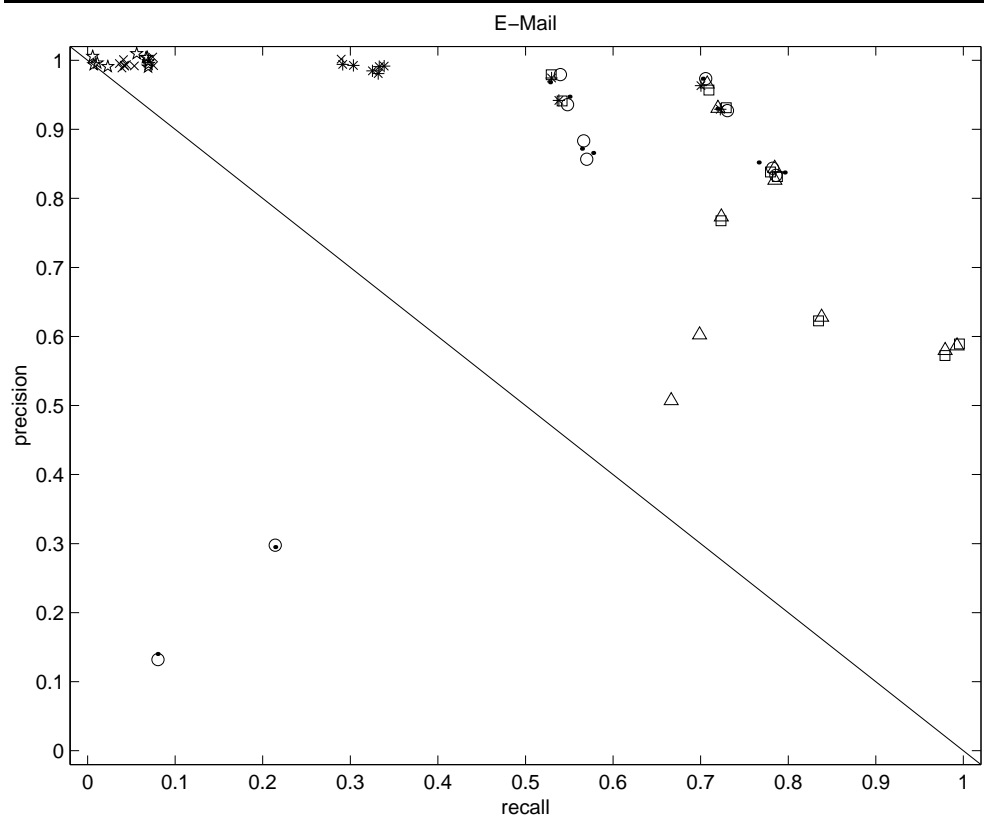
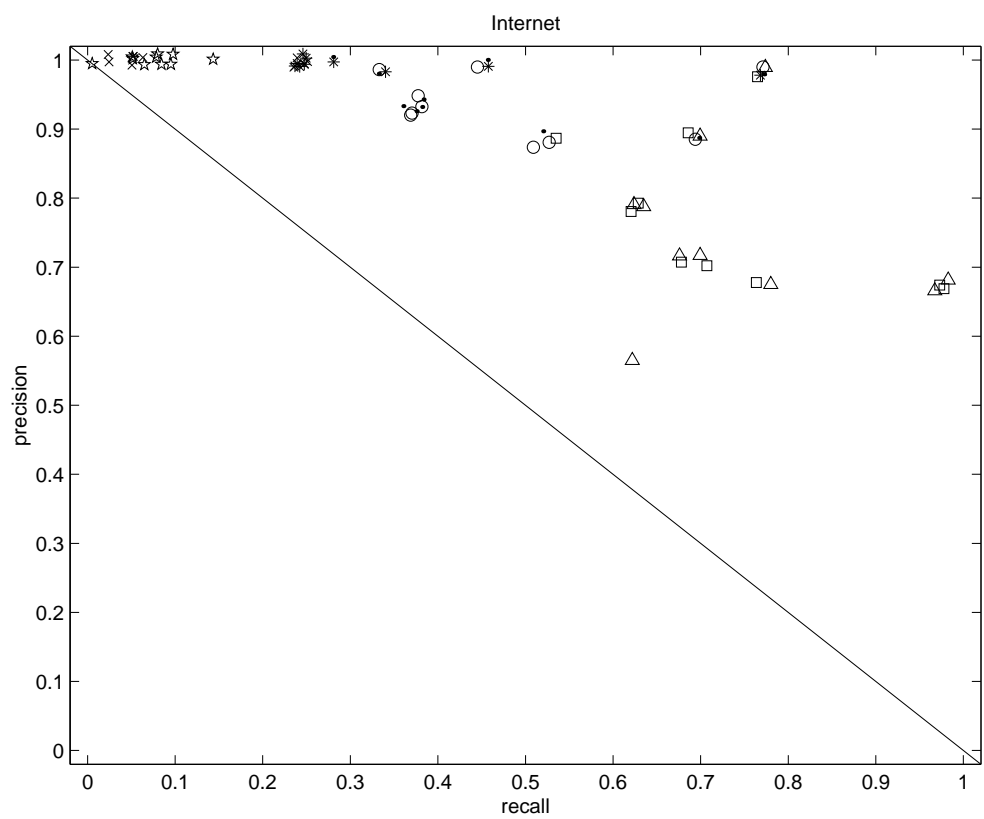
Table 5: This table shows the top ten terms according to each measure for our ten small concepts. The first (uppermost) term has rank one, the next lower one has rank two and so on.

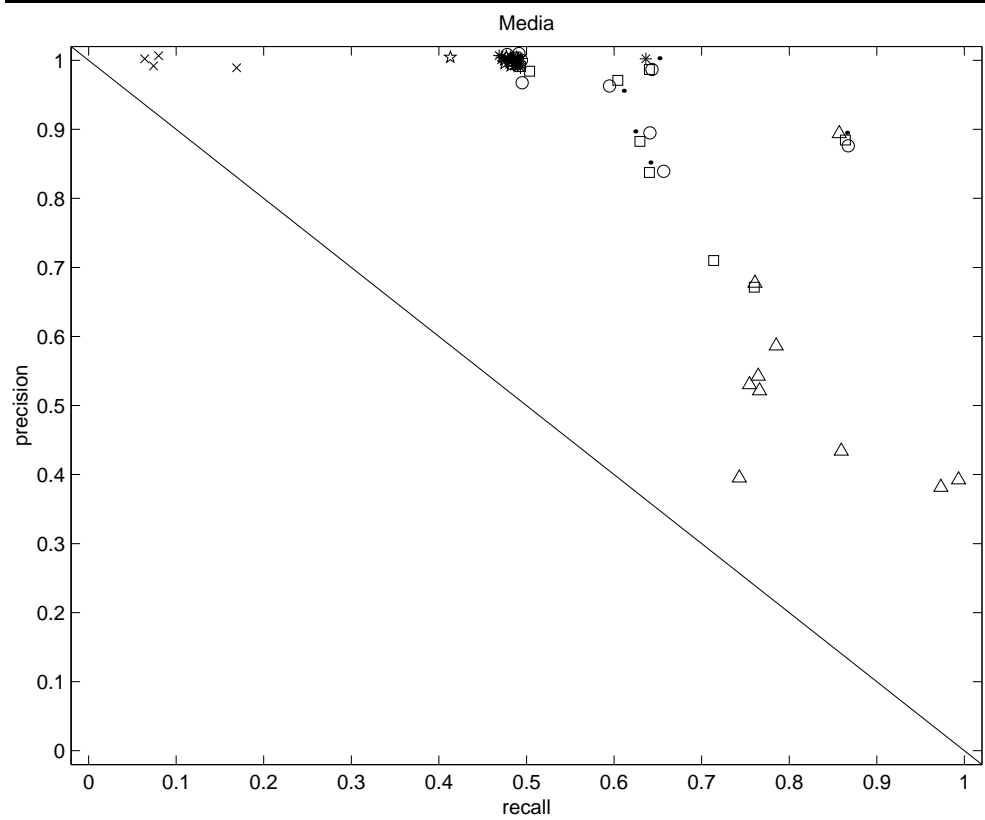
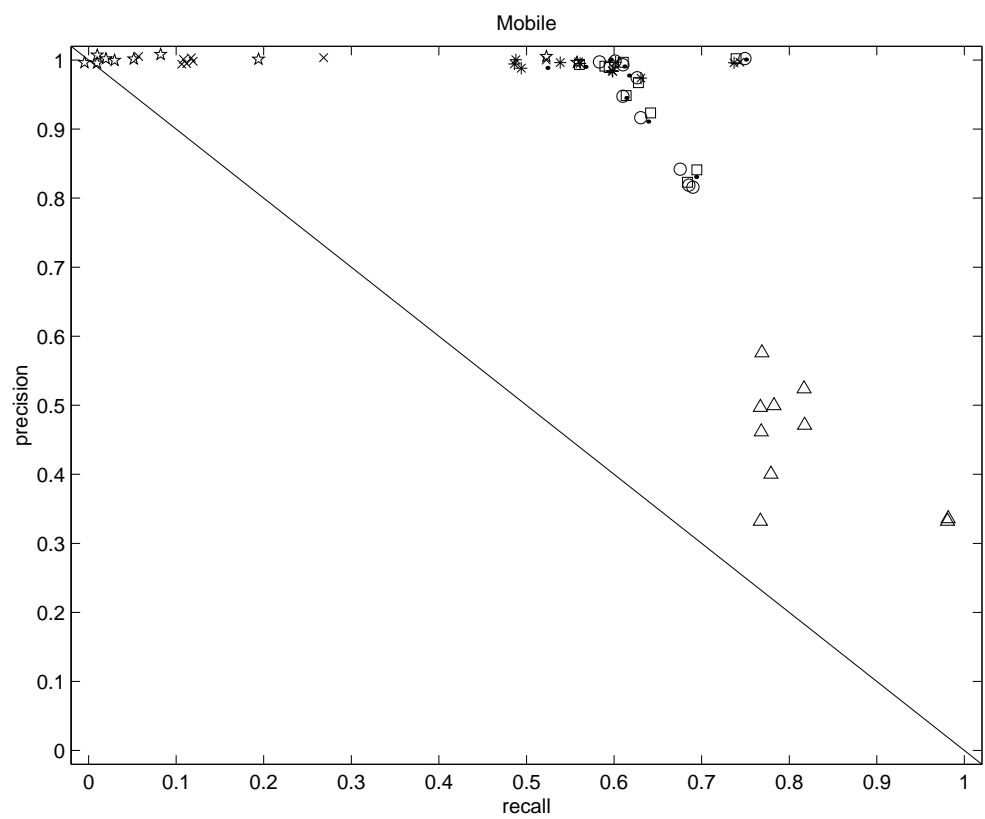
conceptId	χ^2	IG	oddsRatio	odds2	PR	recall	sR3
internet	internet com web wirtschaftblatt atx ifg easi vidx währungen kursliste	internet com wirtschaftblatt atx ifg easi vidx währungen smi kursliste	mediadaten menubar bookshop lebensversicherers mamax ck mysimon techrepublic susa gamespot	internet web contact com menubar mysimon techrepublic gamespot sites nachrichtendienst	internet in 2001 com 1 an for news new us	2001 in internet 1 an com for news document us	newshlight newslinksmall boldlink gatrixx storyurl argnr getday getmonth getdate ismonths
e_mail	version mail e email kontakt td jobs media home archiv	version mail e email td wirtschaftblatt ifg easi vidx atx	fi ndobj 9pt fußballverein lebensversicherers mamax hoppenstedt javaenbled javon gfshpluginname gfshpluginname	email mail version kontakt e printer hoppenstedt druckerfreundliche aussender friendly	version mail e in 2001 an jobs news us email	in 2001 an 1 version e com us news document	boldlink companyname smalltext mediumbold 7pt smalltext2 thiskind unescape newmediasales stammzellen
business unit	unternehmen weitere kontakt kapital ftd partnersites creditreform writelayersn writerecherchen writzeitungn	unternehmen contact weitere das für von im den mit quotes	leserbrieft creditreform stichwortsuche icra unternehmen börsengänge init 68k enterprise graumarkt	unternehmen leserbrieft creditreform stichwortsuche graumarkt ftd guided das zeitungs abonnenten einrichten	unternehmen der von 2001 und die in das für mit	2001 in die der an und von mit das für	börsenticker hoppenstedt analystenstimmen munzinger partnersites sitemap logout menusactive writelayersn writerecherchen
search	news tools latest times policy pt headlines 1995 edition featured	news tools latest times policy pt edition 1995 headlines search	leserbrieft partnersites besucher research realmedia icra redir analystenstimmen dispatch erotikangebote	topics featured headlines recherche leserbrieft partnersites pt 1995 policy latest	news tools in home us latest times service 2001 7	in 2001 news an 1 us document for com write	börsenticker hoppenstedt munzinger stichwortsuche sitemap logout menusactive writelayersn writerecherchen writzeitungn
mobile	mobile wireless policy pt headlines 1995 featured services search magazine	mobile wireless policy pt services headlines 1995 tech search featured	topics 5s wireless mobile subscribe advertisement maximize openwindow unescape getsurveyfile	mobile wireless topics pt headlines featured policy 1995 magazine related	mobile wireless policy pt services headlines tech 1995 search its	in 2001 for 1 us an news com a s	dispatch emazing writeexchanglink convertedsymbol xis nextcard aprs compq symbol1 symbol2
media	media edition times medien true navigator style home msie überblick	media edition times true medien home style navigator version 7	leserbrieft stichwortsuche realmedia redir analystenstimmen relevancy creditreform useragent illustration rm	medien leserbrieft media edition stichwortsuche msie analystenstimmen creditreform useragent appversion	media edition times true home version 7 service medien news	in 2001 media an 1 document news write us version	börsenticker hoppenstedt munzinger partnersites sitemap logout menusactive writelayersn writerecherchen writzeitungn
software	software microsoft computer web pc product it desktop this mobile	software microsoft web it computer pc product for this mobile	software papierchen klingeling bauzulieferer malade glitzernde lpez schlafmützen wahlzeit wurlitzer	software microsoft it web computer pc for more services this	software microsoft in it web 2001 for internet com web us	2001 in software an 1 com for internet us news	tagsüberblick webpromotion newscasts newslink international geldverwalter witwenrente sponsort hilf freundentänze
mobile phone	handy stellenmarkt sms branchen wap persönliches verschicken writelayersn writerecherchen writzeitungn	handy sms stellenmarkt branchen wap technik verschicken persönliches writelayersn writerecherchen	börsenticker kursuche leserbrieft partnersites stichwortsuche analystenstimmen handy creditreform icra schlussbericht	handy börsenticker kursuche leserbrieft partnersites stichwortsuche analystenstimmen creditreform graumarkt stellenmarkt	handy sms stellenmarkt branchen wap verschicken persönliches technik writelayersn writerecherchen	in 2001 die der und an das für mit den	hoppenstedt munzinger sitemap logout menusactive writelayersn writerecherchen writzeitungn writemultichanneln writeregistrierenn
investment/investor	investor investment a dispatch context line weekly investoren starting stay	investor a line top investment s time tech context new	today's dispatch investment sponsort drahtseilakt klingeling bauzulieferer malade lpez moravec	investor investment a top line s new time in us 2001	investor a top line s new time in us 2001	in 2001 an a new us top 1 s at	9ca380 cfd6b2 wirtschaftit tandard gilder maxi geldverwalter schwalben witwenrente glitzernde
network	networks network mysimon corrections comparisons techrepublic gamespot reserved reviews send	networks corrections inc services this today featured policy network reserved	techrepublic topics inc gamespot networks comparisons today's network eweek emachines	networks techrepublic corrections gamespot network comparisons reviews reserved zndnet send	networks network reserved topics mysimon comparisons today corrections techrepublic gamespot	in 2001 us for news com an s 1 a	mysimon dispatch catchup 5s chasm netzwerken newsmakers mariano gwendolyn bulletproofing

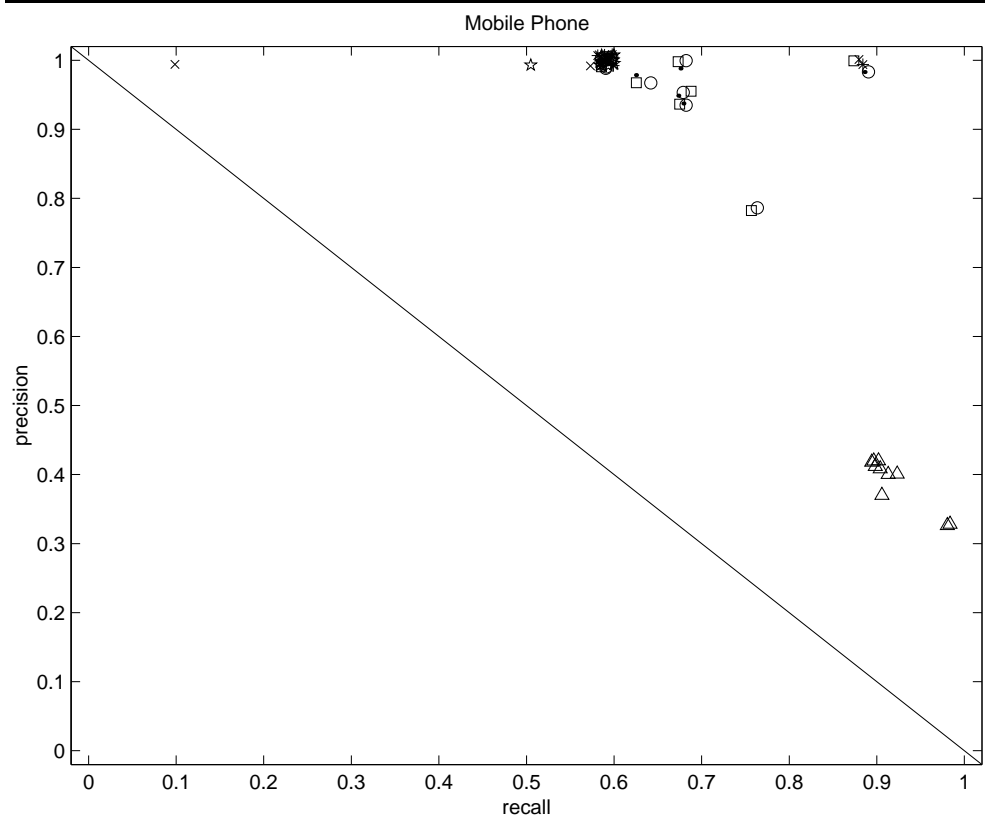
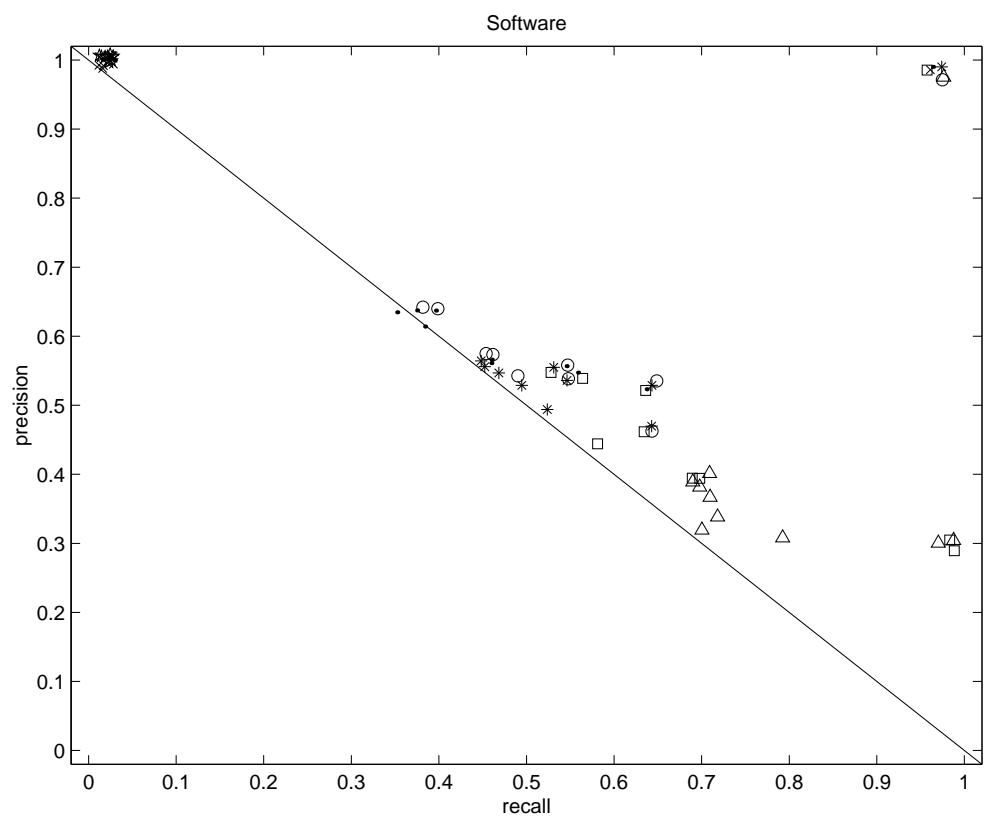
Table 6: This table shows the top ten terms according to each measure for the largest ten concepts when normalizing job size before combining the contingency tables. The first (up-permost) term has rank one, the next lower one has rank two and so on.

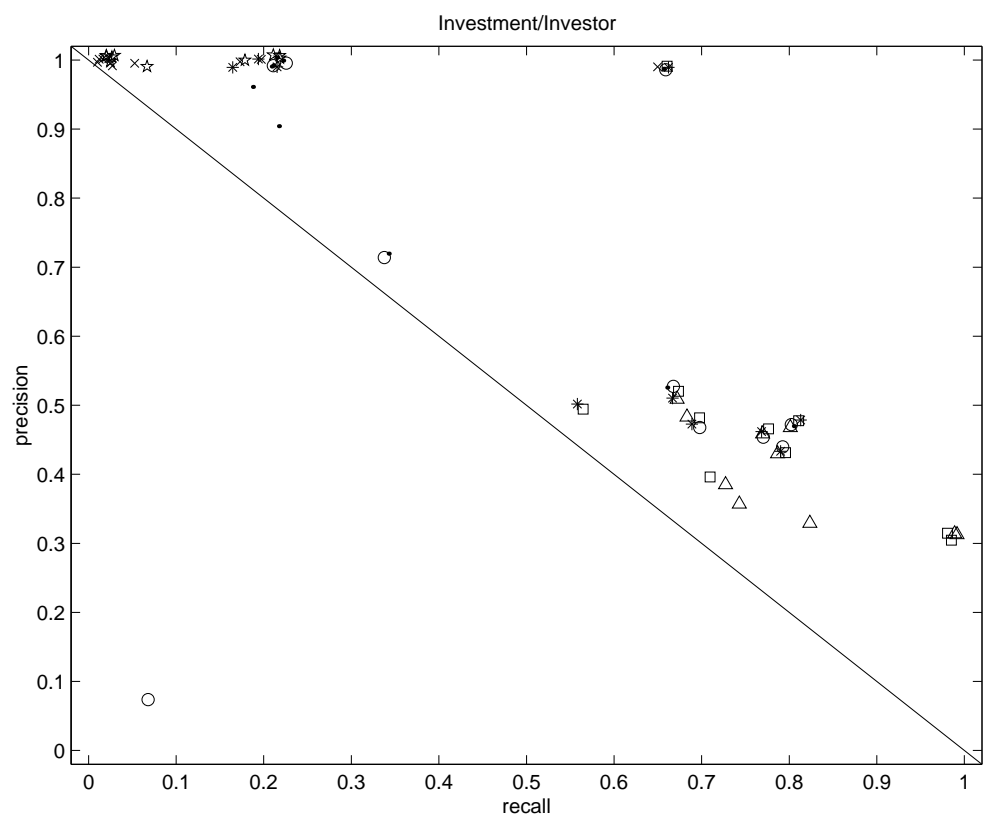
conceptId	X ²	IG	oddsRatio	odds2	PR	recall	sR3
artificial intelligence	ai intelligenz ki seminarvorträge künstliche artificial künstlicher privatstiftungen suchtechnologie dfki	ai intelligenz künstliche artificial intelligenz ki seminarvorträge privatstiftungen künstlicher zusammenfassung	ai zeitaufwendig roboterentwicklung monatsperformance ziegelindustrie luxinvest anlageinstrumenten nordinternet turing inspizieren	eine ist von bin der vor welt nur das kann	ai intelligenz ki seminarvorträge künstliche artificial künstlicher privatstiftungen suchtechnologie dfki	2001 in der 1 die und von das für den	ki dunietz seminarvorträge hutchens goren treister topicalnet dfki trapp vermbobil
online community	pocketmail kemco faxnachrichten oberschalen pagern handlicher nachstehen community auswechselbaren modisches	community pocketmail kemco faxnachrichten oberschalen pagern handlicher nachstehen auswechselbaren modisches	kemco faxnachrichten sputnik7 rocksänger oberschalen pagern spieleproduzent handlicher dreiste tabellenkalkulationen	online internet e software web pc com microsoft service new	pocketmail kemco faxnachrichten oberschalen pagern handlicher community nachstehen auswechselbaren modisches	2001 in online internet an 1 com community e at	pocketmail marke musicity stabilisierungsanleihe mindestlohnes konop targus thumbpad firmmaster vitos
manufacturer	handyhersteller duftchip pretec adventskalender würzbug besinnlichen wohlgeruch bistum betriebsystems sendo	handyhersteller handys nokia mobiltelefone handy mobiltelefon hersteller klingeltöne gprs mobilfunk	bistum betriebsystems textbasierten farbiges z100 potentiell fi nanstidningen energiehunger simo anzuhängen	handy für das dem ein markt hat der von die	handyhersteller duftchip pretec adventskalender würzbug besinnlichen wohlgeruch bistum betriebsystems sendo	2001 in die der und für das von mit an	handyhersteller pretec duftchip adventskalender würzbug besinnlichen wohlgeruch strahlungstests sommerprojekt nachwuchsforscher
wlan	wlan lan computersäulen auskunftssuchenden patrone gaspatrone überwachungsservice beobachtungs methangas gridpatrol	lan wlan wireless thiskind fi ndobj v4 mm versenden multimedia vereinigen	überwachungsservice computersäulen beobachtungs wistron mov stündiges wlan produktionsaktivitäten methangas lebenslänglich	wireless neue parent eine internet den wird für mit hat	wlan lan computersäulen auskunftssuchenden patrone gaspatrone überwachungsservice beobachtungs methangas gridpatrol	2001 in die 1 der und mit an den von	auskunftssuchenden patrone gaspatrone methangas busbestellung protege wlan federgewicht moorhen skycross
market capitalisation	marktkapitalisierung capitalization rangliste zwischenbericht börsentagen autokonzerne oica vda börsenumsatz reinking	marktkapitalisierung capitalization börsen schwieg sitemap out munzinger writelayersn writerecherchen writezeitungen	angeschwellen oica capitalization schwieg jama retailgeschäft sektorstrategie auswahlkriterium bewertungs mehrheitsentscheidungen	euro unternehmen times 7 am fi nancial tools vom home aktien	marktkapitalisierung capitalization rangliste zwischenbericht börsentagen autokonzerne oica vda reinking börsenumsatz	2001 in an die der und markt kapitalisierung von für mehr	marktkapitalisierung split kursänd behrnt punktabzug furse kleiman goy geldbrief skyy
cryptography	verschlüsselung verschlüsselte verschlüsseln verschlüsselt bildmitteilung thumbboard unterhaltungsmaschine inkompatiblen bedienbaren freihand	verschlüsselung verschlüsselte technologie netz stellt für und alle length auf	bedienbaren freihand mobilfunkvertrag inkompatiblen netzstandards animator freundet bereue freihändige archos	auf für werden ist den das mit eine und von	verschlüsselung verschlüsselte verschlüsseln verschlüsselt inkompatiblen bildmitteilung thumbboard unterhaltungsmaschine bedienbaren freihand	die und der in für den mit auf von das	verschlüsselung verschlüsselte verschlüsselt verschlüsseln kryptographie thumbboard bildmitteilung verschlüsselungssoftware unterhaltungsmaschine kryptografie
mobile portal	mobilmfunkportal contentbereich betriebsystemunterstützung zusatzservices pocketmail kemco gerätepreise unterhaltsames geyay internetfernsehen	kabellos pocketmail kemco adapter faxnachrichten oberschalen pagern handlicher nachstehen ideenaustausch	unterhaltsames geyay internetfernsehen playerlösung handyabsatzes tant zeichenerkennung notbook zusatzservices ideenaustausch	handy mobile i com einer de soll sich an var	mobilmfunkportal contentbereich betriebsystemunterstützung zusatzservices pocketmail kemco gerätepreise unterhaltsames geyay internetfernsehen	in an 2001 com 1 der den und das die	mobilmfunkportal targus vpen thumbpad microphones netsize contentbereich betriebsystemunterstützung uninstalled chesnais
online gaming	gaming unreal startopia xtreme racer dungeon vanuatu aufrütteln almanach desperados	gaming games playstation spiele game entertainment xbox unreal computerspiele online	vanuatu spielehits kommunikationsprogramms undying kampfpres virenprobleme drittanbieter aufrütteln rennbahn gehirne	online the microsoft of top line 4 a to entertainment	gaming unreal startopia xtreme racer dungeon vanuatu aufrütteln almanach desperados	online in 2001 the an if top internet for document	ussr onlinespiele tekkon elektroschocks fanaticism gighahertz twi flrtspielchen biershooter saboteure
knowledge management	km wissensmanagement kmart karraker ibi bluelight dissuaded lichtwellenleiter underestimating wmt	km kmart wissensmanagement bluelight mart retailing retailers email discount karraker	karraker ambulanten cmb glasfaserleitung underestimating lichtwellenleiter kernteam srdf asb google	email version online 11 com 13 at customer marketing one	km wissensmanagement kmart karraker ibi bluelight dissuaded lichtwellenleiter wmt underestimating	in 2001 1 com online km an at version 11	km wissensmanagement siriri langlauf sauter kph hevs b61 murzuk lc4
price earning ratio	kgv aktienrückkäufe rechtfertigen sommerrally lignum brainlab frauenförderung gewinnhalberung winkt trüber	kgv rechtfertigen aktienrückkäufe kurse lignum fonds brainlab graumarkt sommerrally aktien	gewinnwarnungsreigen emagine irrglaube stammzelltherapien profitablere schlussniveau aktienrückkäufe bookrunner gewinnhalberung bewertungs	bei unternehmen als sich aus einem zum über ein für	kgv aktienrückkäufe rechtfertigen sommerrally lignum brainlab frauenförderung winkt gewinnhalberung trüber	2001 die in punkt und für von mit den auf	kgv frauenförderung kursänd punktabzug erstnotizausblick geldbrief primärmärkten spwx skyy steinemann

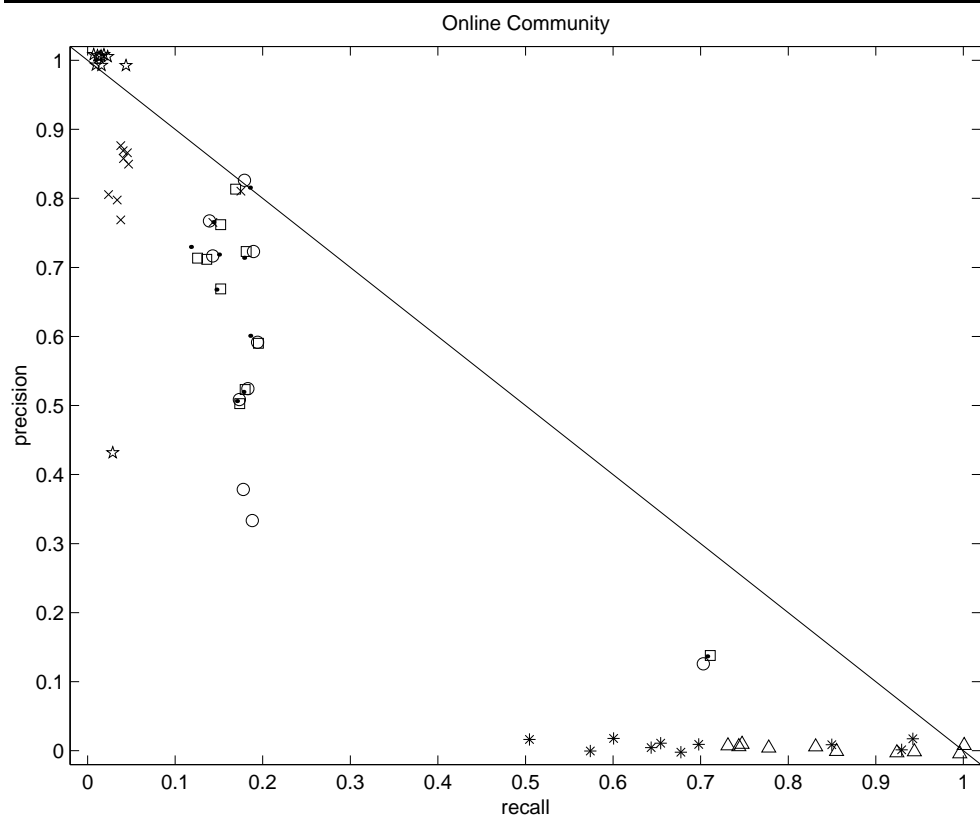
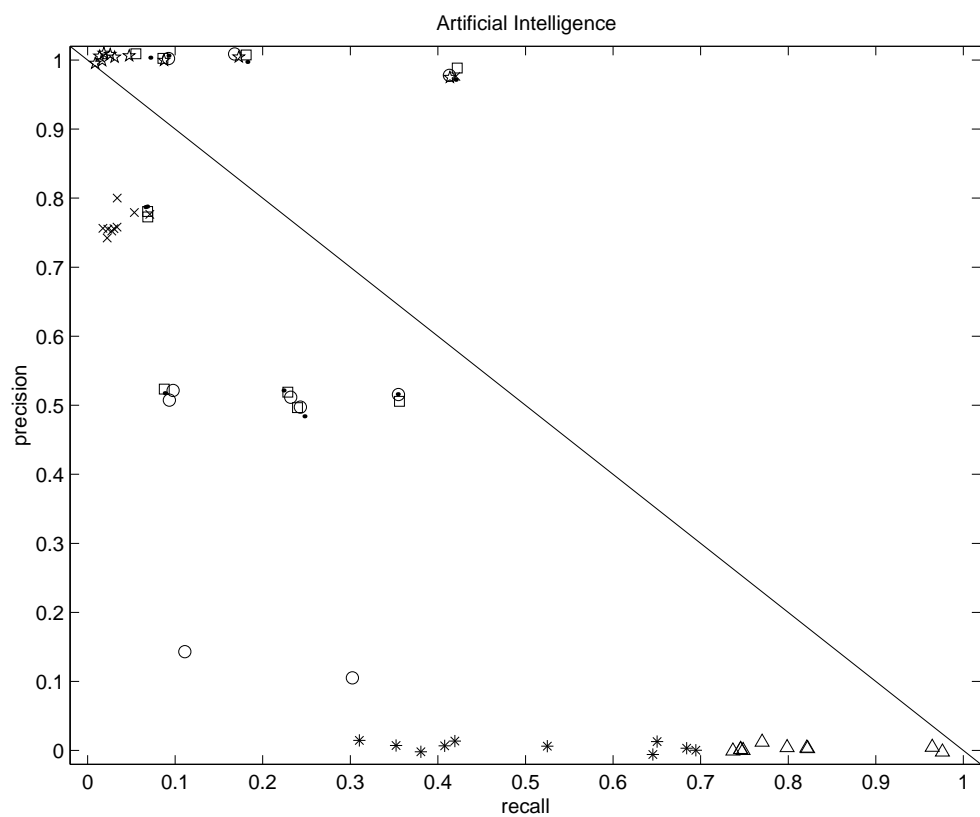
Table 7: This table shows our the top ten terms according to each measure for our ten small concepts when normalizing job size before combining the contingency tables. The first (up-permost) term has rank one, the next lower one has rank two and so on.

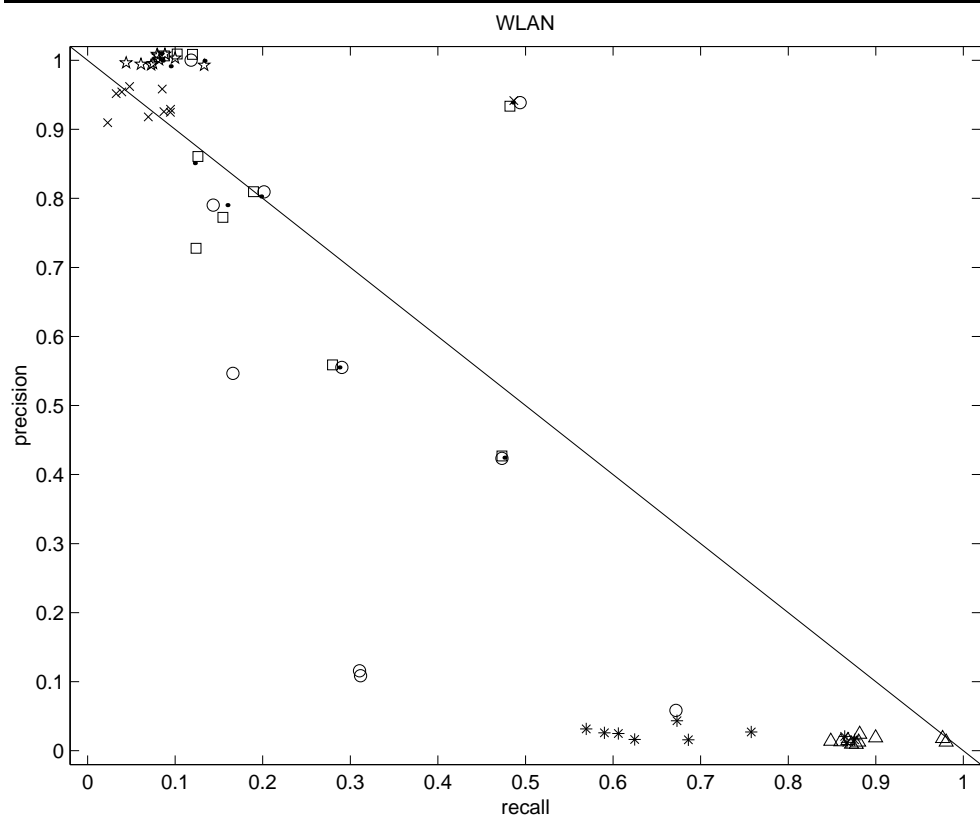
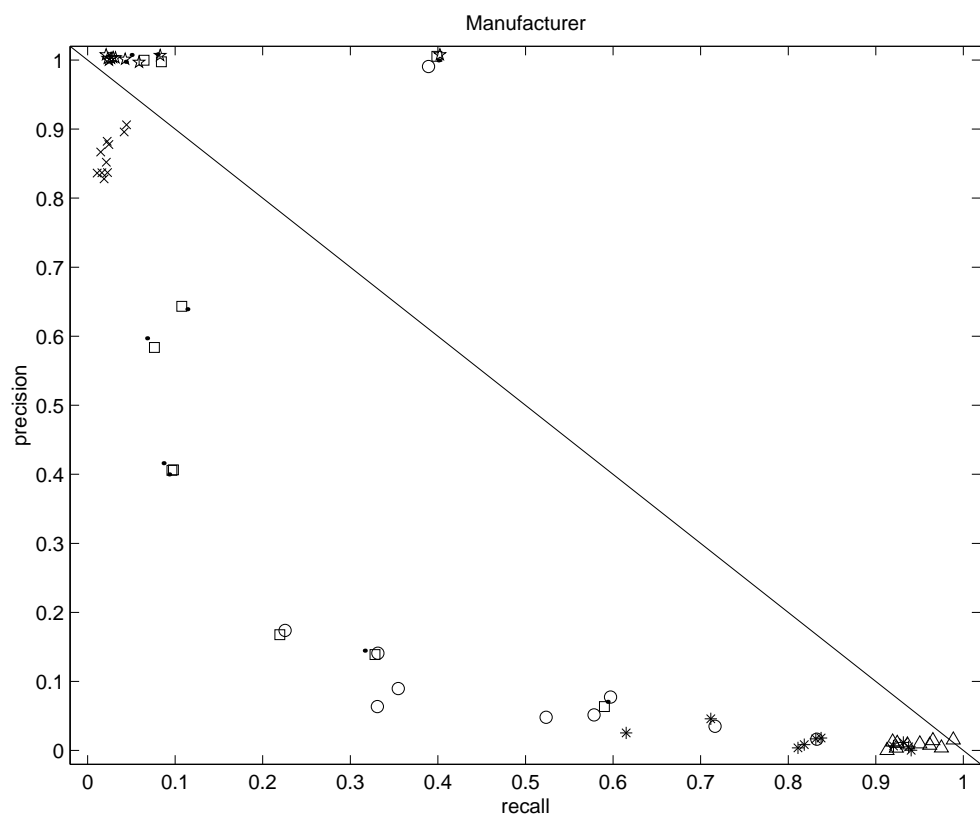




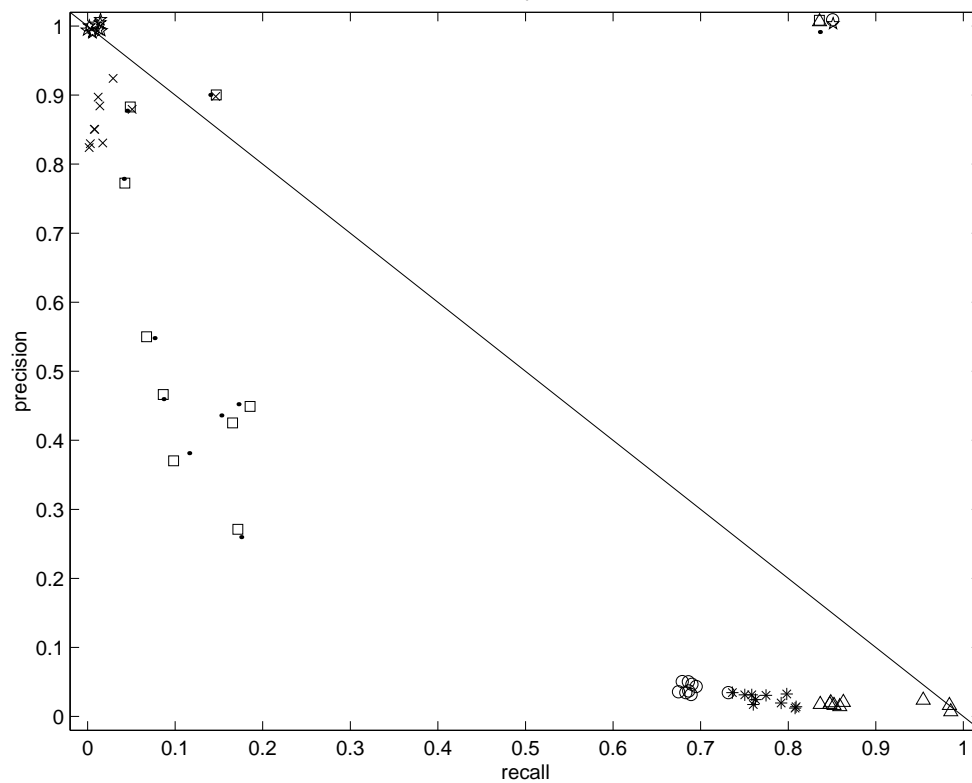




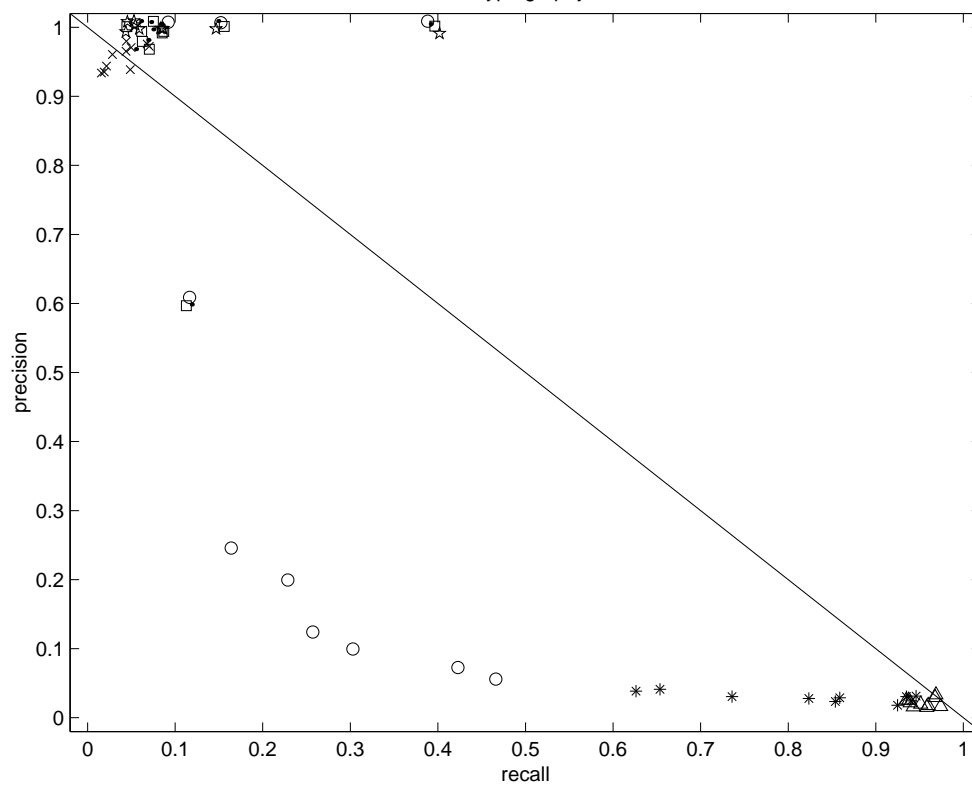


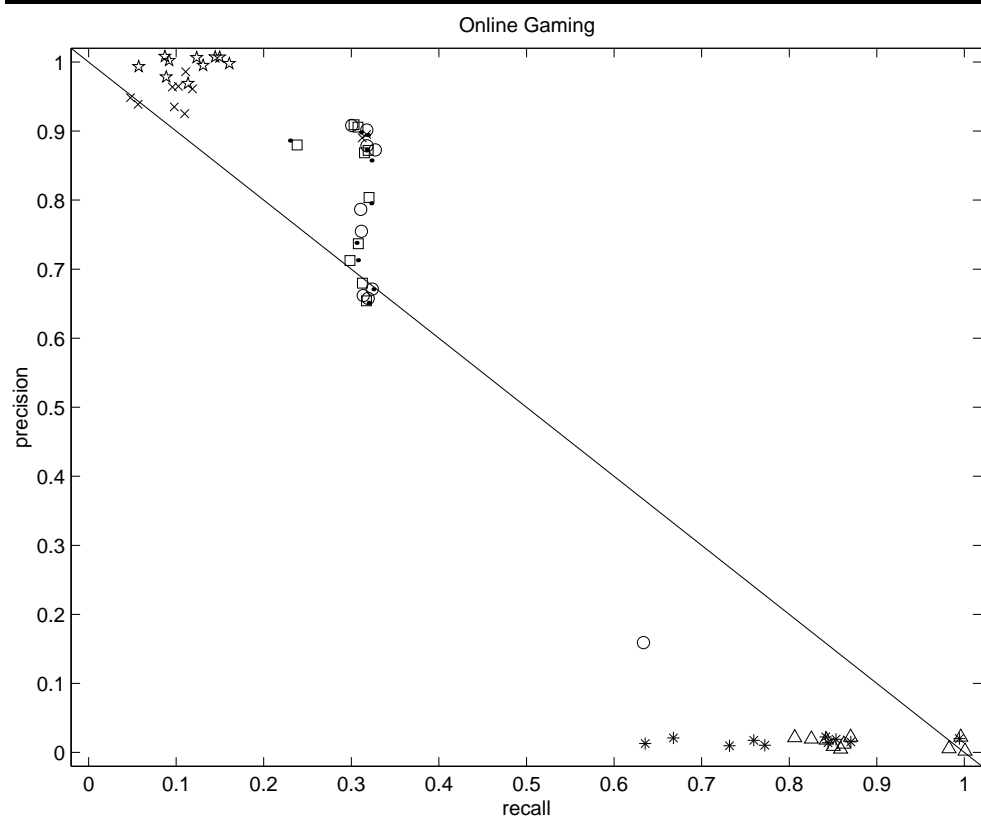
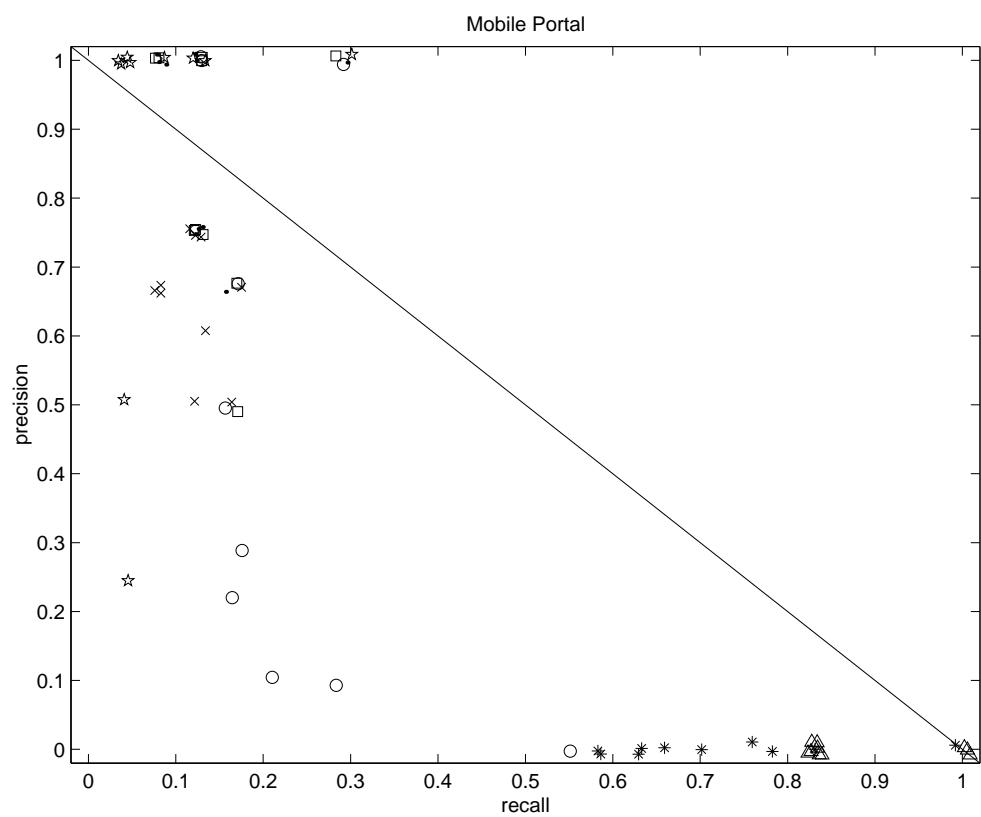


Market Capitalisation

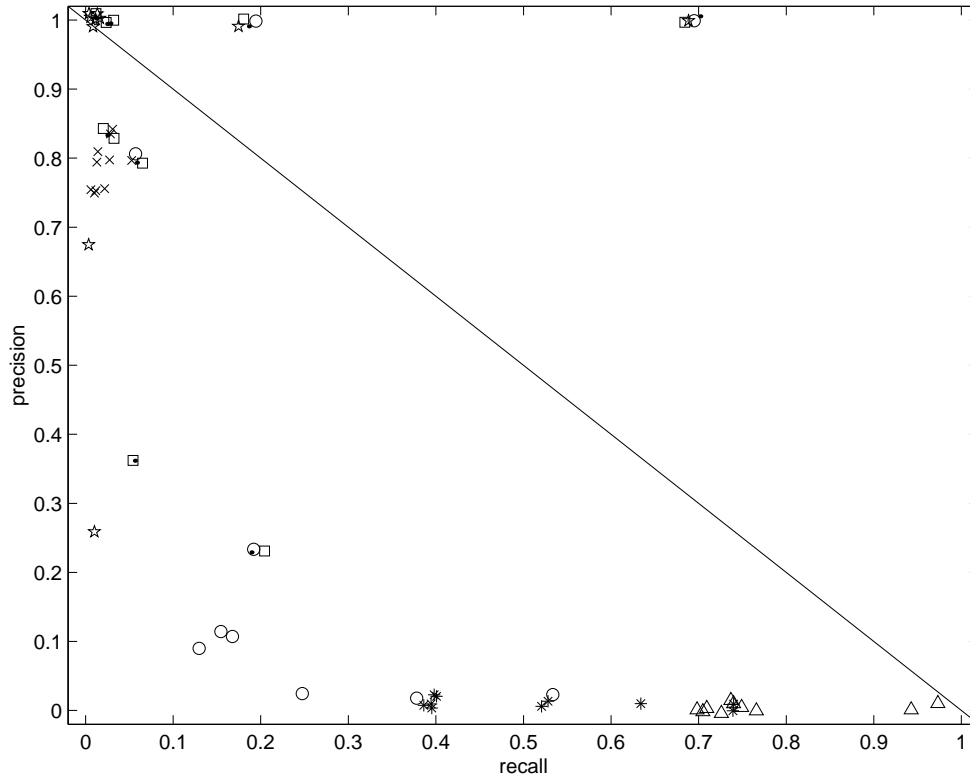


Cryptography





Knowledge Management



Price Earning Ratio

