

In Search of Disruptive Ideas: A Survey for Outlier Detection Techniques in Crowdsourcing Innovation Platforms

Overview

[research context / application area]

Idea Management Systems = online collaborative tool to **collect ideas from many people** (e.g. clients of a company)

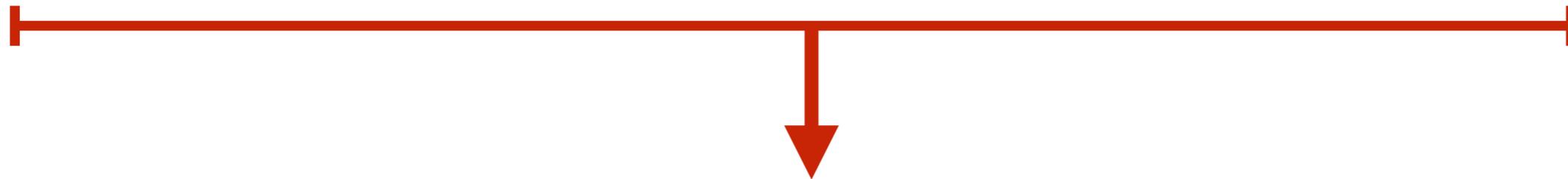
The screenshot displays the IdeaStorm website interface. At the top, the browser address bar shows 'www.ideastorm.com'. The main header features the IdeaStorm logo with the tagline 'Where Your Ideas Reign' and a notification stating 'There are currently no Storm Sessions active. Stay tuned!'. Below the header, a navigation bar includes 'IdeaStorm' and 'Storm Sessions'. A central row of four buttons offers actions: 'View' (All posted ideas by the community), 'Post' (Your idea for Dell products or services), 'Vote' (Promote or demote ideas), and 'See' (Your ideas in action). Below this, a 'Sort Ideas By' section shows 'Popular Ideas', 'Recent Ideas', and 'Top Ideas' as selected options, along with a dropdown menu set to 'All'. The main content area lists ideas with their IDs and titles. The first idea, ID 118030, is titled 'Pre-Installed OpenOffice | alternative to MS Works & MS Office' by user 'dhart' from Feb 17, 2007. It includes a description: 'Provide OpenOffice.org for free pre-installation alongside Microsoft Works and Microsoft Office. OpenOffice.org is more capable than Microsoft Works, and a serious competitor to Microsoft Office, at a fraction of the cost (it's free!)' and a note: 'OpenOffice.org can open, create, edit and save Microsoft Word, Excel and PowerPoint files.' The second idea, ID 106450, is titled 'Pre-Installed Linux | Ubuntu | Fedora | OpenSUSE | Multi-Boot' by 'dhart' from Feb 16, 2007, with a description: 'Offer the 3 top free Linux versions for free pre-installation on all Dell PCs. Quality free and open source software drastically lowers the cost of new PCs, and helps prevent...'. On the right side, there is a 'Login to IdeaStorm' form with fields for 'Username:' and 'Password:', a 'Login' button, and a 'Forgot Password?' link. Below the login form is a section titled 'Your Ideas in Action' featuring a 'IdeaStorm Recap - 1/22/2010' post with the text: 'Happy 2010 everyone! I know I'm a little late with the holiday greetings, but there is a lot to share on...'

Overview

[research context / application area]

Idea Management Systems **problems**

- Lots of contributions
- Lots of duplicates, similar ideas etc.
- Lots of simple or obvious input



difficult to choose the best innovations

Overview

[problem - hypothesis - approach]

- **problem:** pick best candidates for interesting/disruptive innovations
- **hypothesis:** good ideas are rare outliers that stand out from the majority of other proposals
- **approach:** use outlier detection algorithms on idea text to detect the most anomalous ideas

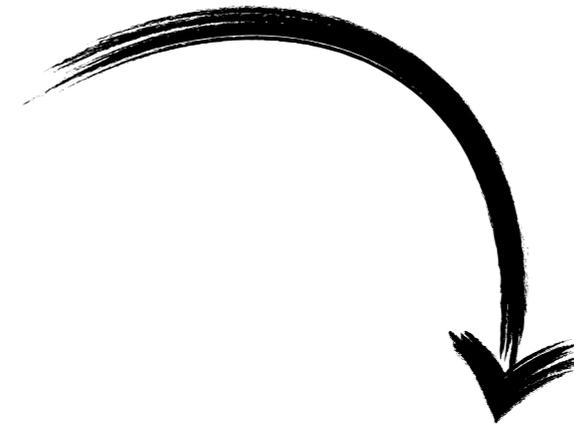
what is disruptive innovation / how to find it ?

[hypothesis theoretical grounding]

innovation literature

“ **disruptors** deliver innovations for overlooked market segments , while **market leaders** address their most demanding customers via incremental innovation “

["What is disruptive innovation?" Christensen, Raynor, Harvard Business Review, 2015]



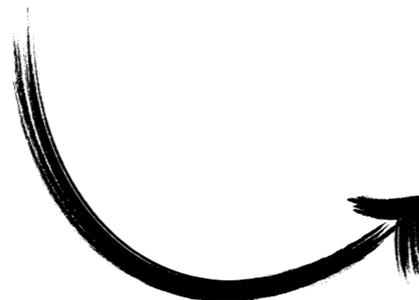
IMS hypothesis

standard criteria/ metrics of IMS:

- **favour success as perceived by the entrenched market leader** point of view
- **could overlook disruptive ideas** that get **no popular support**
- metrics bring up ideas of **most vocal customers**
- less-demanding customers are less vocal and not equally participatory in IMS, (*ie. low-end foothold; or non existent customers ie. new market-footholds [Harvard Review]*)

“ **not all disruptive ideas have to lead to success** “

["What is disruptive innovation?" Christensen, Raynor, Harvard Business Review, 2015]



disruptive
origin

disruptive
vs. success

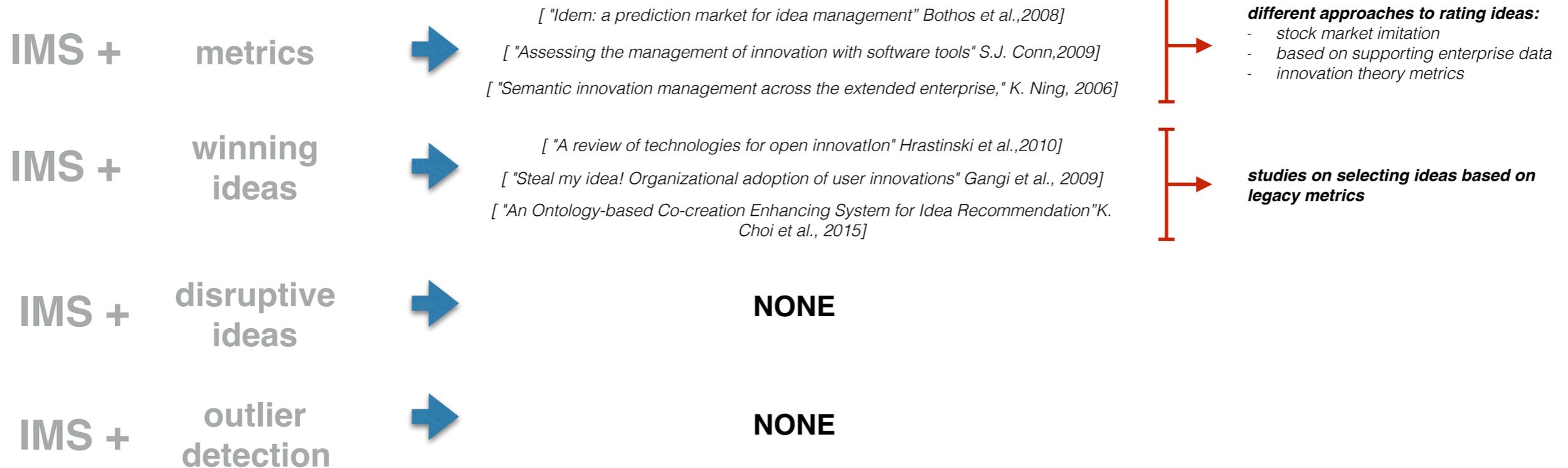
Approach: creating metric for disruptiveness of ideas
[how to find and evaluate the best outlier detection]

- 1. survey** available outlier detection algorithms
- 2. pick** the most representative candidates based on previous applications
- 3. apply + eval** algorithms for idea management problem using two different public datasets (alg picks vs. manual annotation)
- 4. compare** results of different algorithms
- 5. recommend** the best approach

Work so far

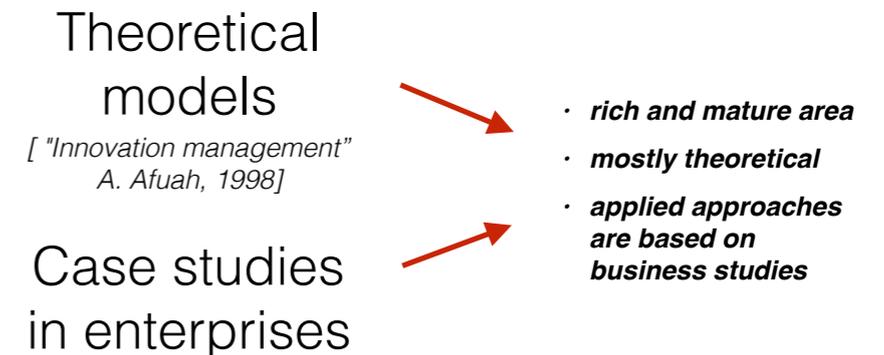
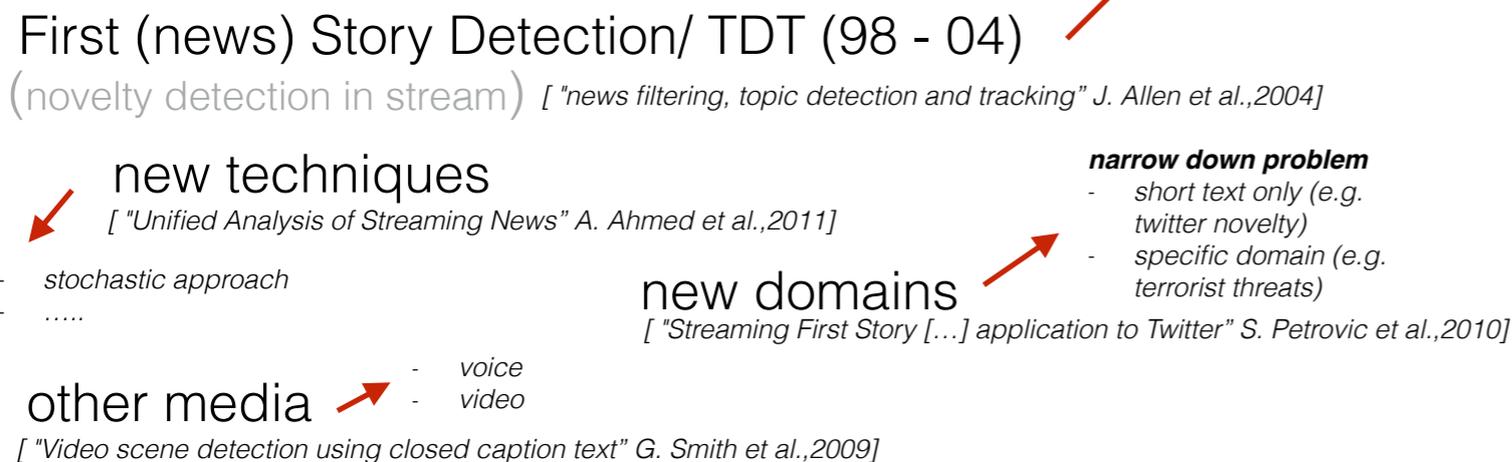
[related / past work]

Idea Management Systems (key area)



Outlier Detection + text (subarea)

Innovation Management (subarea)



Survey + alg picks

[different taxonomies / evaluations]

multitude of algorithms across years



many surveys / taxonomies to classify SoA

CLASSIFICATION [Ji Zhang, 2013]

1. Statistical (probabilistic)
2. Distance (proximity) based
3. Density based
4. Clustering
5. High dimensional

CLASSIFICATION [Aggarwal, 2013]

1. Extreme Value Analysis
2. Probabilistic and statistical
3. Linear
4. Proximity based
 - 4.1. Clustering
 - 4.2. Density
 - 4.3. Nearest neighbour
5. Information Theory based
6. High dimensional

CLASSIFICATION [Chandola, 2008]

1. Classification Based
2. Clustering Based
3. Nearest Neighbour Based
 - 3.1. kNN
 - 3.2. Density
4. Statistical
5. Information Theoretic
6. Spectral

choose the categories that repeat across surveys
pick one algorithm per each category to evaluate

Evaluated Algorithms

[evaluation outline]

1.Distance Based: kNN (k Nearest Neighbours)

1. Feature vector generation:

1. TF-IDF
2. WORD2VEC
3. LDA / VEM
4. LDA / Gibbs

2. Distance measures:

1. Cosine
2. Manhattan
3. Euclidean

2.Probabilistic / statistical: LDA (Latent Dirichlet Allocation)

3.Density Based: LOF (Local Outlier Factor)

1.Distance measures: Cosine, Manhattan, Euclidean

4.Clustering: kMeans / kMedoids

1.Distance Measures: Cosine, Manhattan, Euclidean

Evaluation datasets

[two different scenarios]



1. Dell IdeaStorm:

1. Ideas: new equipment, software for PC manufacturer business

2. Innovators: customers

3. Stats:

- 15,000 ideas (207 implemented)
- 2,000 users



2. Starbucks Ideas:

1. Ideas: new drinks, food, changes in offering for coffee chain

2. Innovators: customers / store owners

3. Stats:

- 10,000 ideas (1069 implemented)
- 3,000 users

Evaluation - dataset labels

[manual annotation]

Idea title

Idea textual desc

manual annotation

3 innovation metrics:

- Implementation cost
- Potential profit
- Market size

1 overall rating

- Breakthrough

1 - 10
Likert
scale

based on **innovation management** theory

Which ideas to annotate ?

legacy metrics ranking

- Vote count
 - 10 top
 - 10 middle
 - 10 bottom
- Comment count
 - 10 top
 - 10 middle
 - 10 bottom
- 10 implemented (random pick)
- 10 unimplemented (random pick)

outlier metrics ranking

- 10 top for every algorithm / configuration tested

~1000 ideas annotated
per dataset

Evaluation metrics

[assessment of results quality]

What makes a good ranking ?

correlation with manual eval results

shows if the overall ordering reflects the expected one (ie 1,2,3,4... 5000 etc. if idea count = 5000 , as ranked by breakthrough rating)

precision@10 vs. manual ranking

shows how well the outlier ranking works for the top outliers (most important ones for organization stand point)

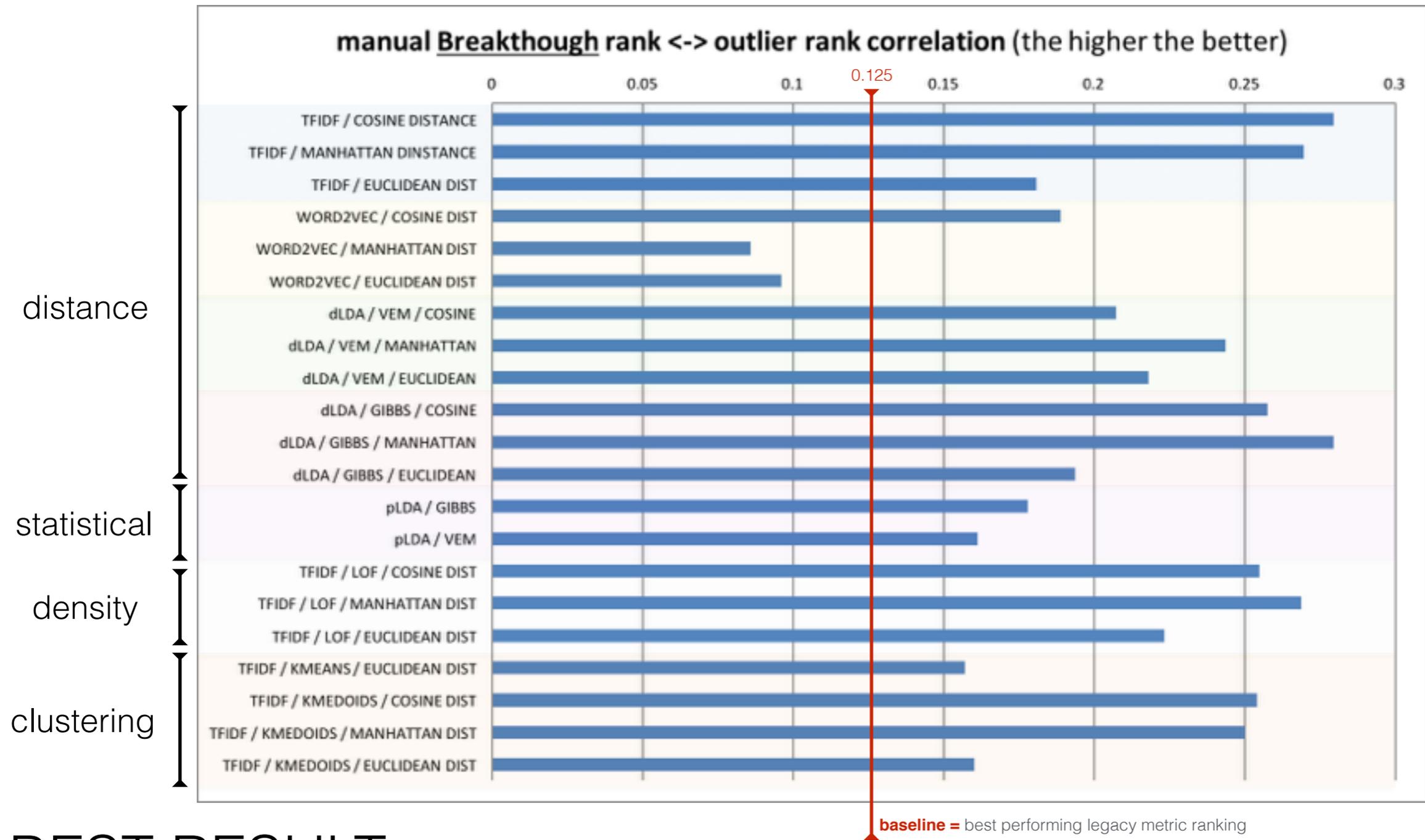
[Evaluating Recommendation Systems, Guy Shani and Asela Gunawardana, Microsoft Research, 2009]

additional extended analysis

- [distance/density] comparison of effectiveness for different neighbourhood settings
- [probabilistic] comparison for different topic optimisation settings
- [clustering] comparison for different cluster sizes / iterations / feature vectors

Results (IdeaStorm)

[correlation of algorithm rankings vs. manual picks]



BEST RESULT

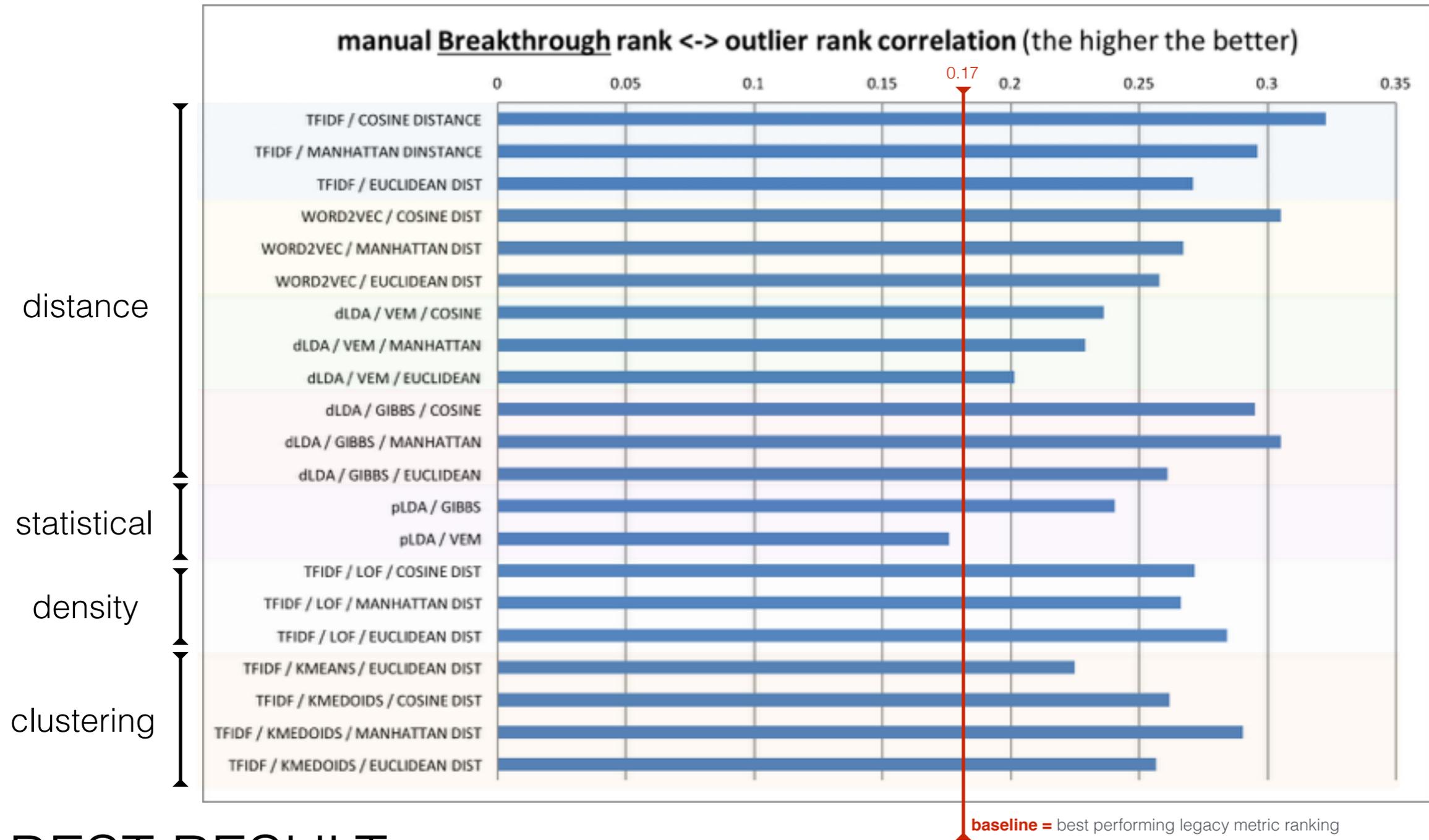
DISTANCE BASED OUTLIER DETECTION | **TF-IDF + COSINE = 0.28** -> **MEDIUM*** correlation with manual scoring

VS. 0.12 legacy score -> **WEAK*** correlation with manual scoring

*Cohen correlation scale for social sciences (Cohen,xxxx)

Results (Starbucks)

[correlation of algorithm rankings vs. manual picks]



BEST RESULT

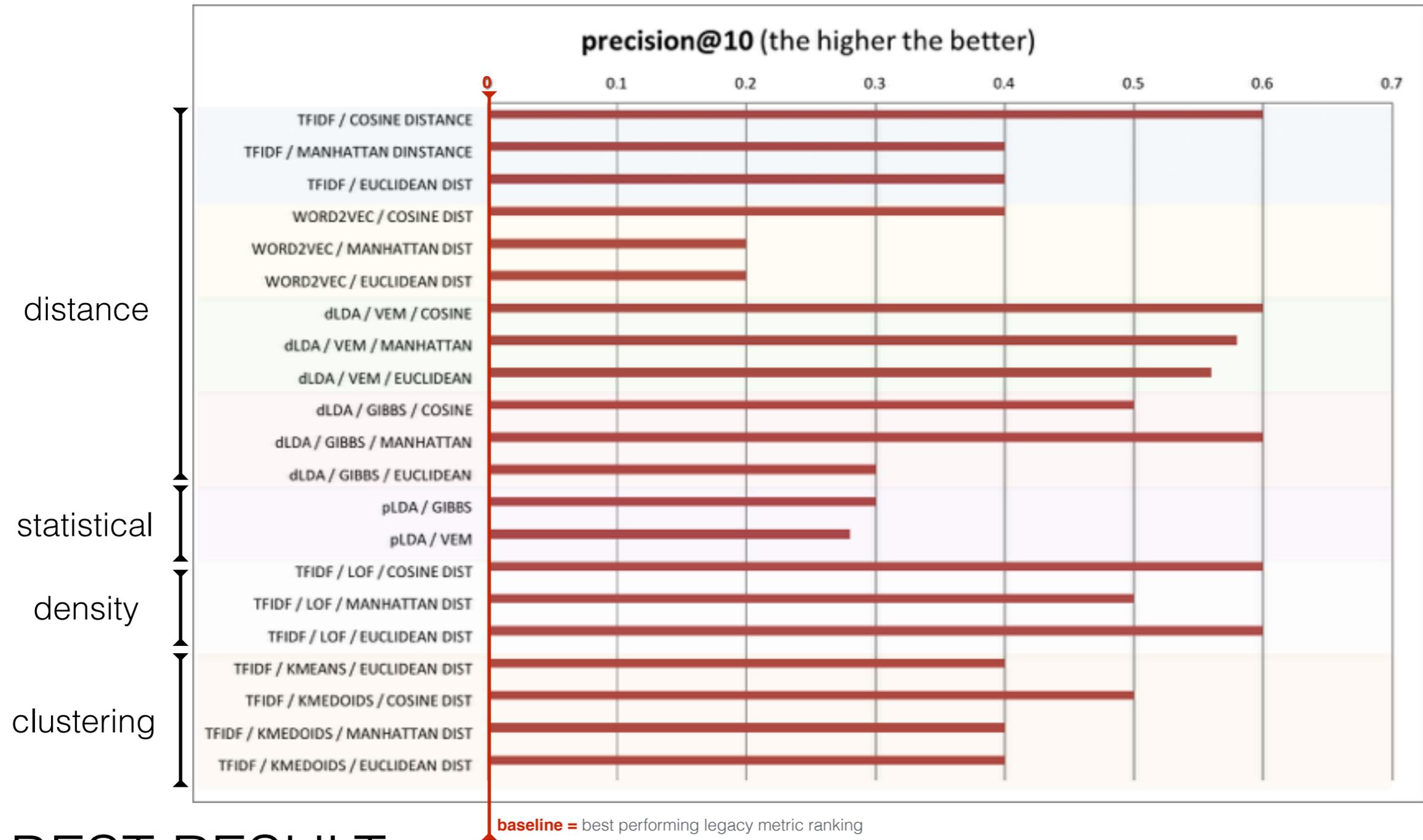
DISTANCE BASED OUTLIER DETECTION | **TF-IDF + COSINE = 0.32** -> **MEDIUM*** correlation with manual scoring

VS. 0.17 legacy score -> **WEAK*** correlation with manual scoring

*Cohen correlation scale for social sciences (Cohen,xxxx)

Results (IdeaStorm)

[precision@10 for algorithm rankings vs. manual picks]



BEST RESULT

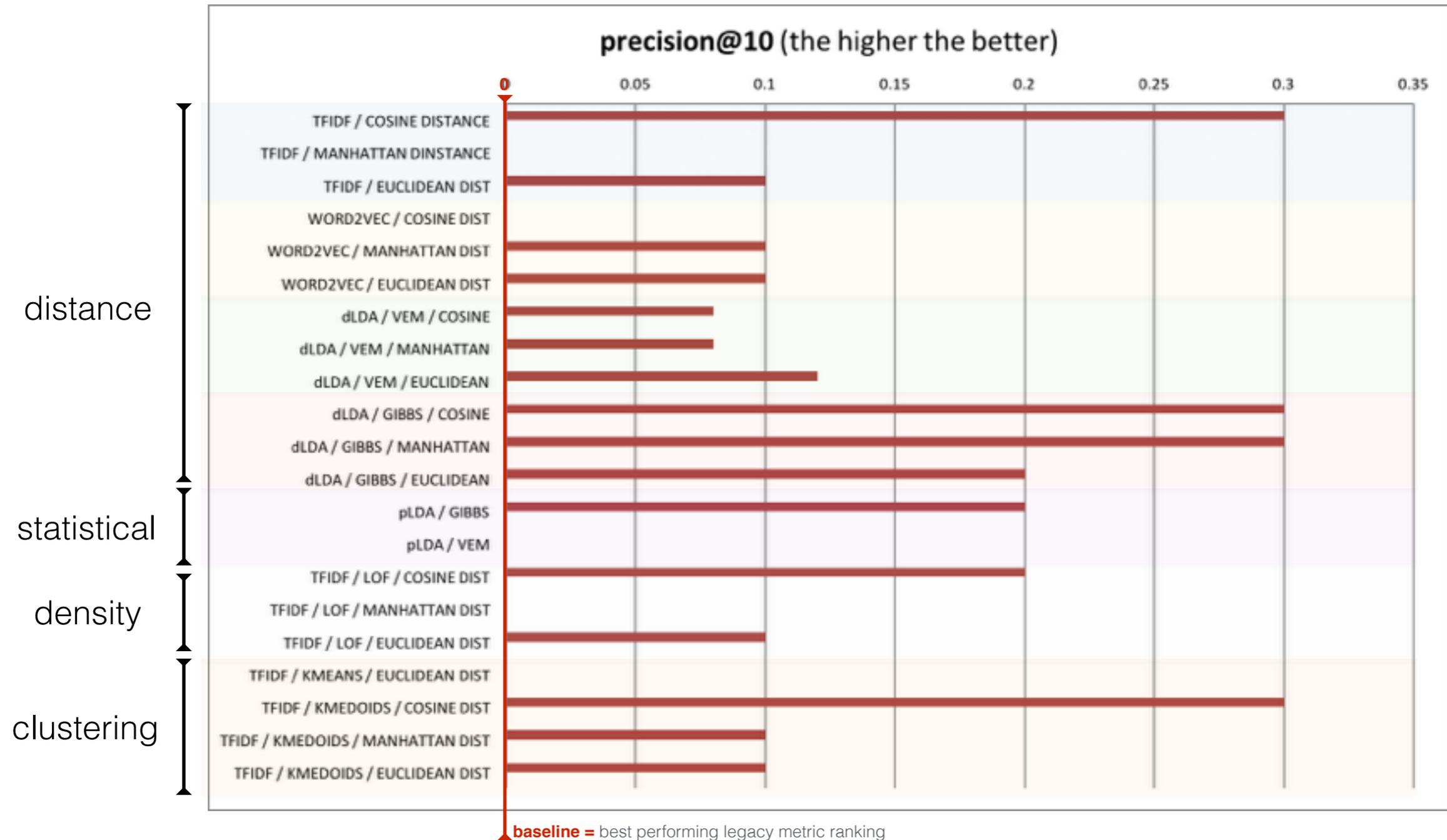
DISTANCE BASED OUTLIER DETECTION | TF-IDF + COSINE* = 0.6

VS. 0 legacy score

*(similar dLDA + COSINE; dLDA + MANHATTAN; LOF + COSINE; LOF EUCLIDEAN)

Results (Starbucks)

[precision@10 of algorithm rankings vs. manual picks]



BEST RESULT

DISTANCE BASED OUTLIER DETECTION | TF-IDF + COSINE* = 0.3

*(similar dLDA + COSINE; dLDA + MANHATTAN; K-MEDOIDS + COSINE)

VS.

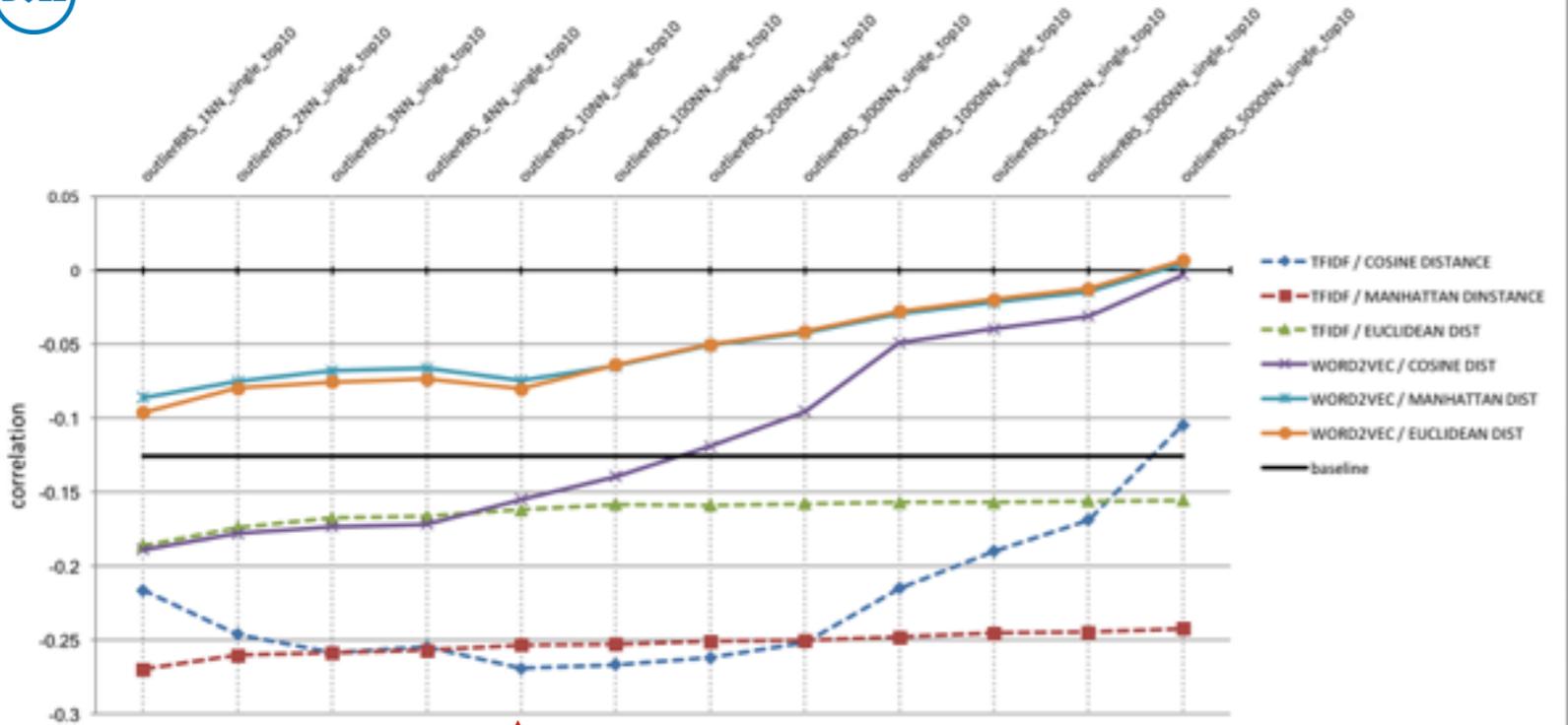
0 legacy score

Results / distance algorithms

[correlation of algorithm rankings vs. manual picks]



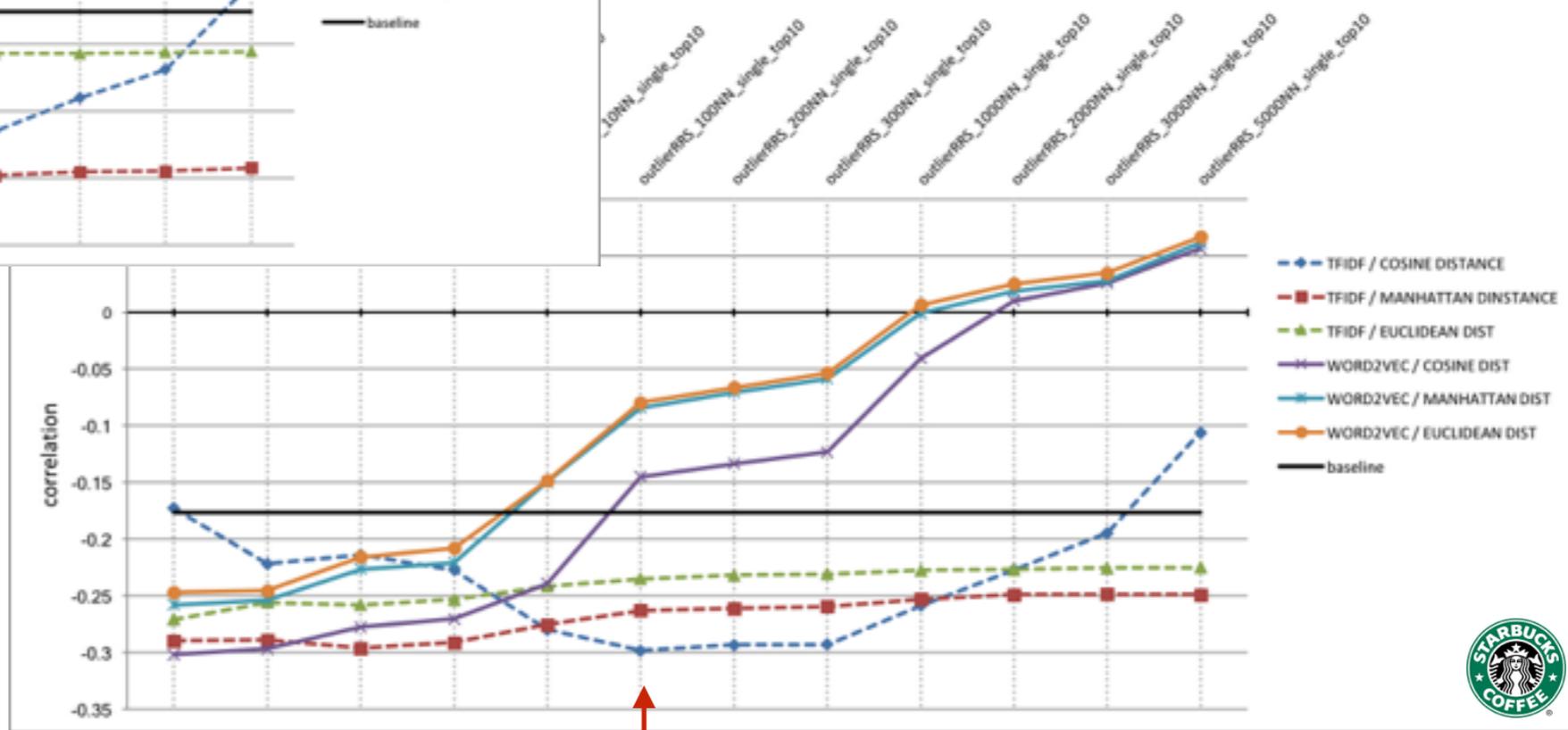
RRS-single - manual rank <-> outlier rank correlation (the lower the better)



best performance

closer look at best case
distance based algorithms

rank <-> outlier rank correlation (the lower the better)



best performance

Observations:

- Best performance | **k = 100** | dataset independent
- Overall behaviour | dataset independent



Conclusions

- **DISTANCE BASED ALGORITHMS**
perform best for particular problem discussed
- **STATISTICAL OUTLIER DETECTION**
performs worst and is also **hardest to tune**
- **CLUSTERING ALGORITHMS**, contrary to distance algorithms
significance of “k” outweighs any other parameter by big margin (in terms of accuracy impact)
- **ALL CASES** (almost) regardless of approach outlier detection brings **new metric quality** to Idea Management System