

Economical Evaluation of Systems: Users as a Signal

Kevin Roitero
Research Talk
New York, 7 June 2019



About Myself

About Myself



November 2016 -Today (1 November 2019)

Ph.D. Candidate at University of Udine, Italy.

Work on Information Retrieval, Crowdsourcing, Artificial Intelligence, Statistical modeling, Transfer Learning, Data mining and Analysis. Supervisor: Stefano Mizzaro.



April 2017 - June 2017

Visiting period at University Of Sheffield, UK.

Work on Crowdsourcing and statistical modeling, with a focus on user interaction and engagement, and agreement metrics.



March 2018 - July 2018

Visiting period at RMIT University, Australia.

Work on Agreement metrics, relevance scales, and Crowdsourcing Evaluation.



September 2018 - December 2018

Research Intern at Spotify London.

Work on Statistical Modeling, Machine and Transfer Learning for BaRT.

Outline

Outline

- Measure Agreement
- Incorporate Agreement in Evaluation
- Effect of Judgment Scales on Agreement
- Bias
 - Bias in (IR) Evaluation
 - Bias in Scholarly Publishing
- Economical Evaluation of Systems

Measure Agreement

Setting

- micro-task crowdsourcing
- many workers do the same task
- agreement among workers can / should be leveraged
- leveraging agreement can be useful for:
 - estimating the reliability of collected data
 - understanding behavior of the workers

Agreement Formalization

	Item ₁	Item ₂	...	Item _j	...	Item _N
Assessor ₁	$r_{1,1}$	$r_{1,2}$...	$r_{1,j}$...	$r_{1,N}$
Assessor ₂	$r_{2,1}$	$r_{2,2}$...	$r_{2,j}$...	$r_{2,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor _i	$r_{i,1}$	$r_{i,2}$...	$r_{i,j}$...	$r_{i,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor _M	$r_{M,1}$	$r_{M,2}$...	$r_{M,j}$...	$r_{M,N}$

Agreement Formalization

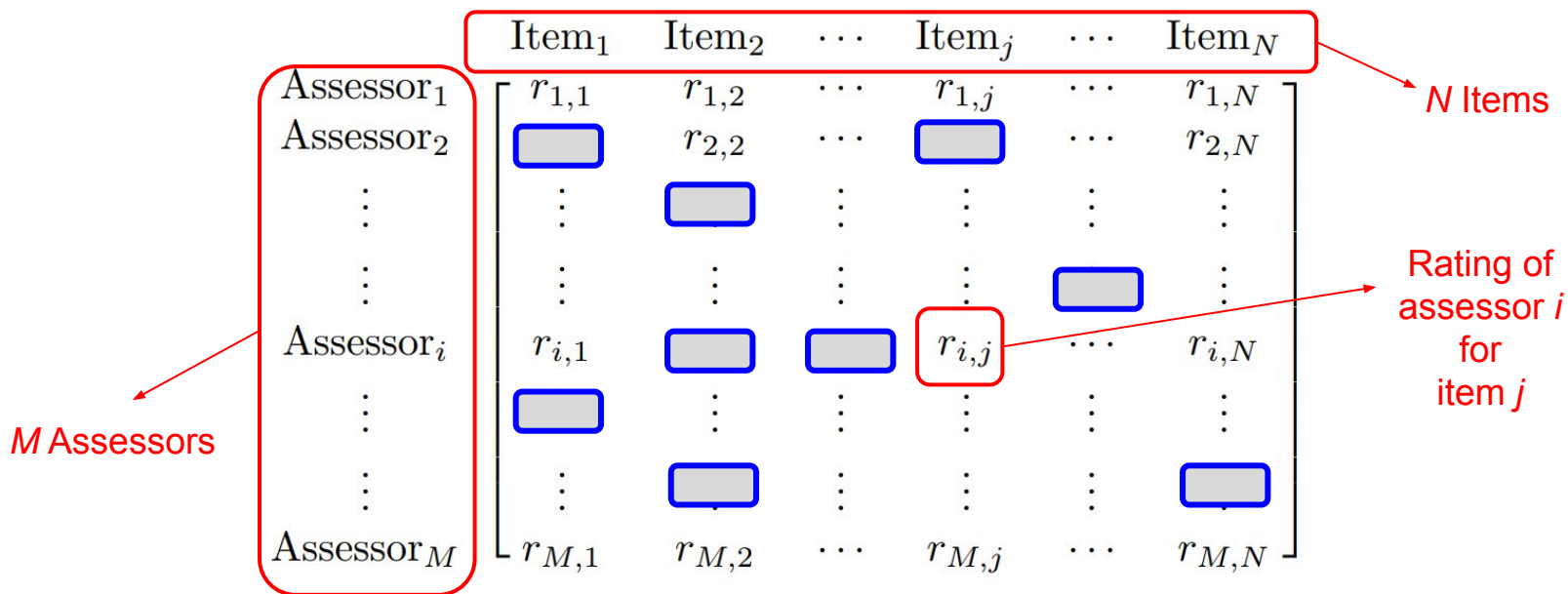
	Item ₁	Item ₂	...	Item _j	...	Item _N
Assessor ₁	$r_{1,1}$	$r_{1,2}$...	$r_{1,j}$...	$r_{1,N}$
Assessor ₂	$r_{2,1}$	$r_{2,2}$...	$r_{2,j}$...	$r_{2,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor _i	$r_{i,1}$	$r_{i,2}$...	$r_{i,j}$...	$r_{i,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
Assessor _M	$r_{M,1}$	$r_{M,2}$...	$r_{M,j}$...	$r_{M,N}$

M Assessors

N Items

Rating of assessor i for item j

Agreement Formalization



This matrix is often **very** sparse in crowdsourcing

There are Several Agreement Measures

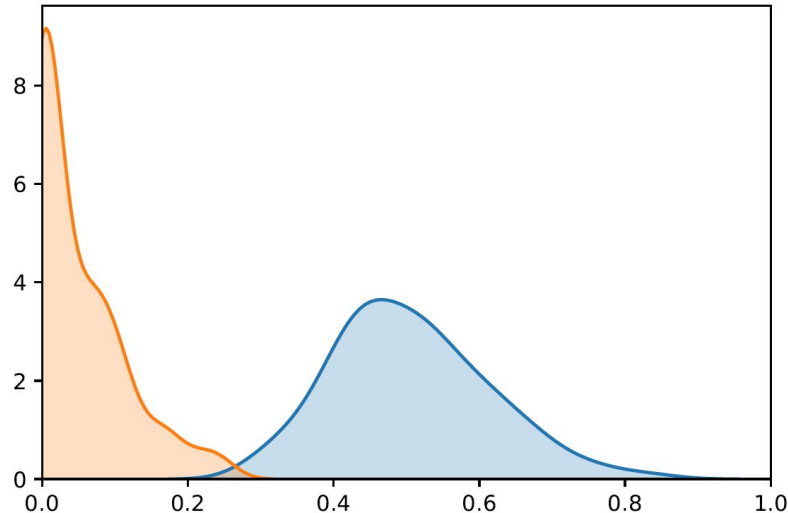
- Percentage Agreement (PA)
- Scott's π
- Cohen's κ
- Intraclass Correlation Coefficient (ICC)
- Fleiss κ
- Krippendorff's Alpha

Current Agreement Measures Are Inadequate

- measures often borrowed from other scenarios with **different assumptions** (which usually do not hold for crowdsourcing):
 - one assessor rates all items
 - all assessors rate all items
 - limited and fixed (= known) number of assessors
- measures are often designed for estimating **data reliability**, not **agreement**
 - **reliability**: the capacity of any measurement tool to differentiate between respondents when measured twice under the same conditions. [Berchtold]
 - **agreement**: the capacity of any other measurement tool applied twice on the same respondents under the same conditions to provide strictly identical results. [Berchtold]
 - reliability can be considered as a necessary but not sufficient condition to demonstrate agreement. [Berchtold]

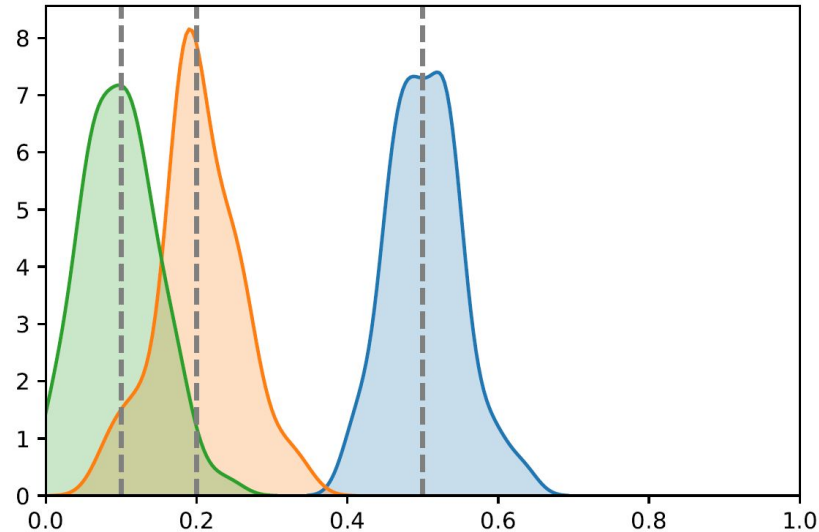
Problems

- there is more variability of judgments in the centre of the scale w.r.t. scale boundaries.
→ can lead to over-estimate agreement close to scale boundaries.



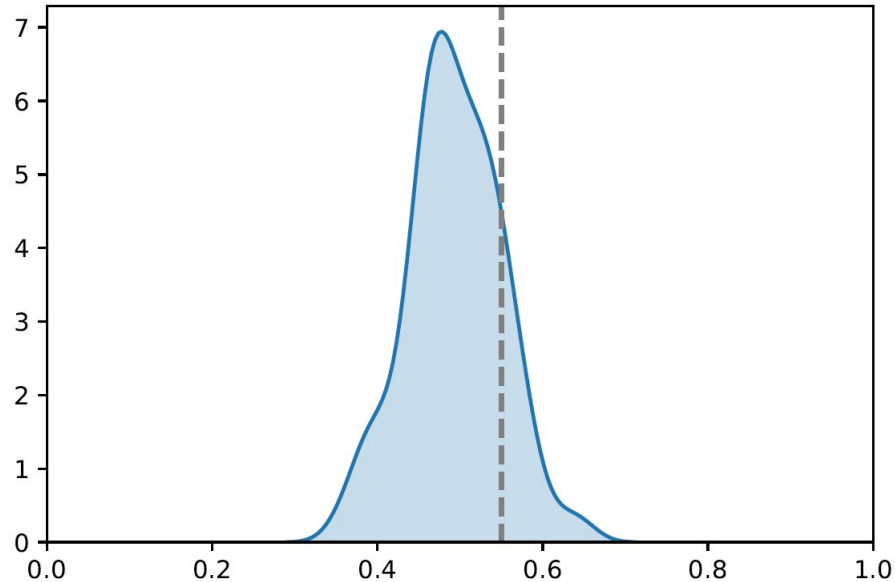
Problems

- the concentration point can be different for different items
→ can lead to over/under-estimate agreement



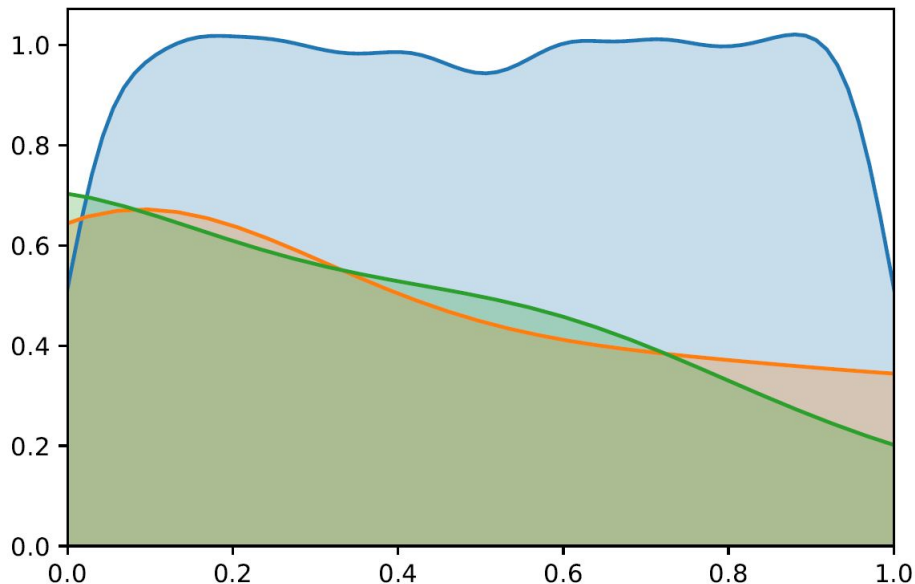
Problems

- additional information is often not considered (e.g., gold questions)



Problems

- different ideas of “**agreement by chance**” definition
- correction by chance assumptions are often violated in crowdsourcing setting



Real Problems with State-of-the-Art Measures

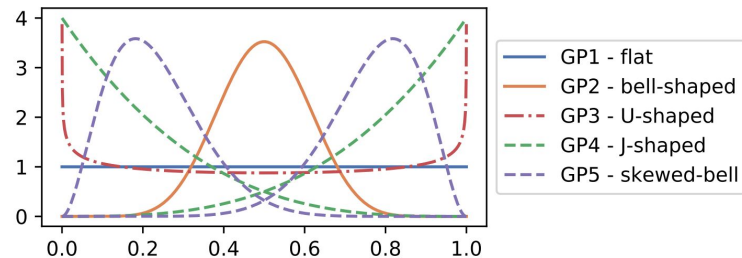
- Percentage Agreement (PA)
 - does not consider agreement by chance
 - works only with nominal data
 - depends on the scale granularity (can not compare different scales)
- Scott's π and Cohen's κ
 - work only with two assessors
 - work only with nominal data
- Intraclass Correlation Coefficient (ICC)
 - assessor have same marginal probability of an answer (not true in crowdsourcing)
 - equivalent to weighted Cohen's κ
- Fleiss κ
 - Generalizes κ to multiple assessors (i.e., shares the same issues)

Real Problems with State-of-the-Art Measures

- Krippendorff's Alpha: an attempt to generalize previous metrics
 - **Random guessing can have high agreement**
 - **Random guessing may have more agreement than honest coding**
 - High agreement, low reliability
 - Zero change in percentage agreement causing radical drop in reliability.
 - **Eliminating disagreements does not improve agreement**
 - Honest work as bad as coin flipping.
 - Two datasets: same quality, same agreement; but higher reliability in one.
 - punishing larger sample and replicability (i.e., data quantity dependent)
 - **“reverse answer” problem** $([1, 0, 0, 0, 1] \neq [0, 1, 1, 1, 0])$

Our Measure: Φ

- “*agreement is the amount of concentration around a data value*”
- if not agreement, we have disagreement \rightarrow negative agreement
- in practice:
 - first, we fit a distribution over the histogram of the ratings
 - then, we measure the dispersion of such distribution
- the fitting distribution has to be general enough to capture:
 - random judgments \rightarrow flat distribution
 - agreement \rightarrow bell-shaped distribution
 - agreement around scale boundaries \rightarrow J-distribution
 - disagreement \rightarrow U shaped distribution
 - \rightarrow use a Beta function
- we should have a minimal number of parameters, to avoid overfitting



Our Measure: Φ

- we use a Beta distribution to model our scenario: $B(a, b)$
- we re-parametrize the distribution in terms of the mean value μ and the precision p as $\mu = \frac{a}{a+b}$; $p = a + b$
- now, we can treat separately mean and dispersion
- we can have a metric that is agnostic of the mean value
- then, we transform to have values in the $[-1, +1]$ range:


$$\Phi = 1 - 2^{\frac{-p \log 2}{2}}$$

Our Measure: Φ

- we use Bayesian inference to compute Φ :

$$P(\vec{\mu}, \Phi | X) = \prod_{i=1}^N \prod_{j=1}^M B(X_{i,j} | \mu_i, \Phi)^{O_{ij}} \prod_{i=1}^N \mathcal{N}(1/2, \sigma_{\mu}^2 \mathbf{I}) \mathcal{N}(0, \sigma_{\Phi}^2) C,$$

probability of observing the mean values, with a common dispersion, given the observed data



- Then, we estimate Φ using


$$\hat{\Phi} = \arg \max_{\Phi} P(\vec{\mu}, \Phi | X).$$

Our Measure: Φ

- we use Bayesian inference to compute Φ :

$$P(\vec{\mu}, \Phi | X) = \prod_{i=1}^N \prod_{j=1}^M B(X_{i,j} | \mu_i, \Phi)^{O_{ij}} \prod_{i=1}^N \mathcal{N}(1/2, \sigma_{\mu}^2 \mathbf{I}) \mathcal{N}(0, \sigma_{\Phi}^2) C,$$

probability of observing the mean values, with a common dispersion, given the observed data



- Then, we estimate Φ using

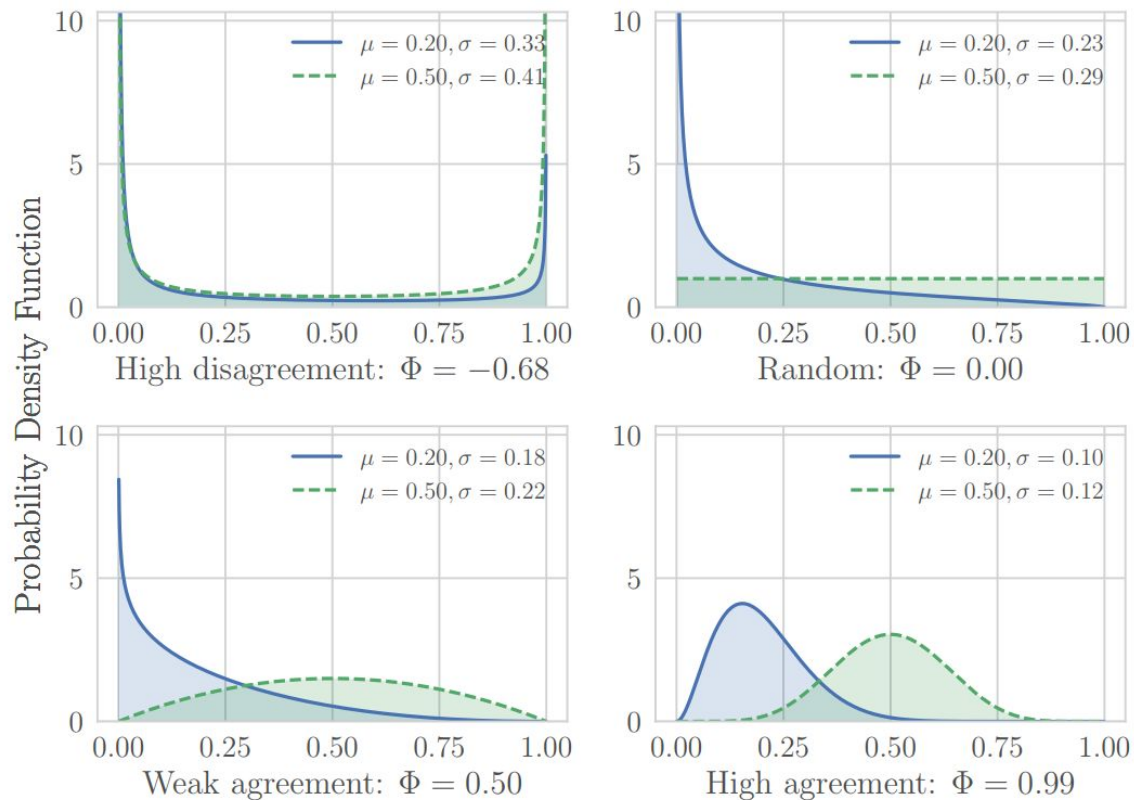
$$\hat{\Phi} = \arg \max_{\Phi} P(\vec{\mu}, \Phi | X).$$

the formula can change to incorporate custom ground truth

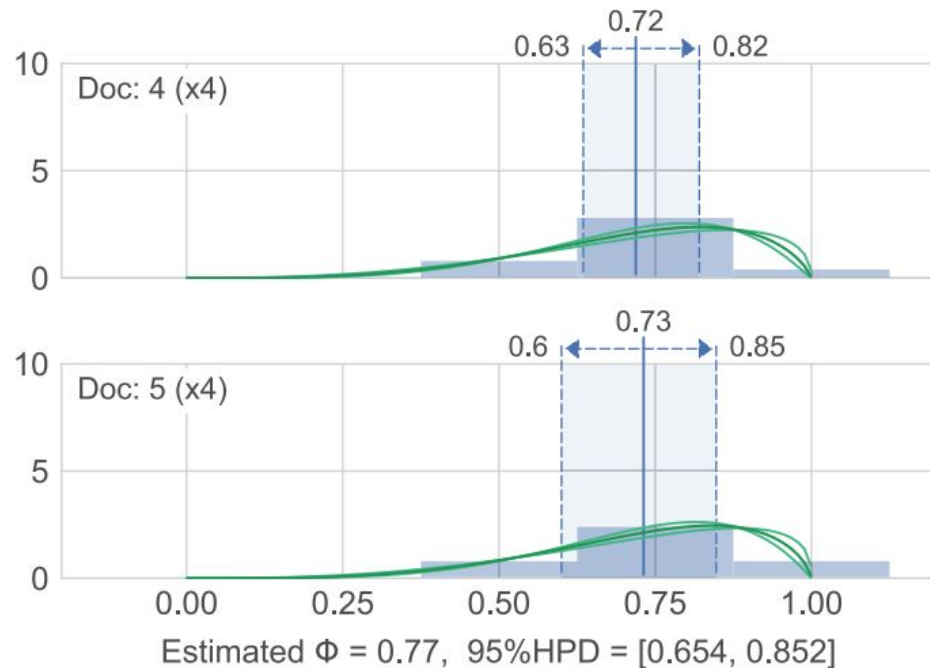
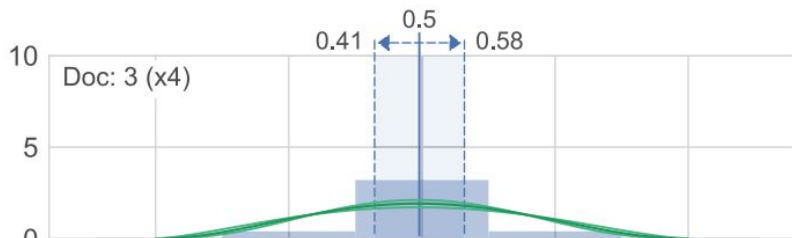
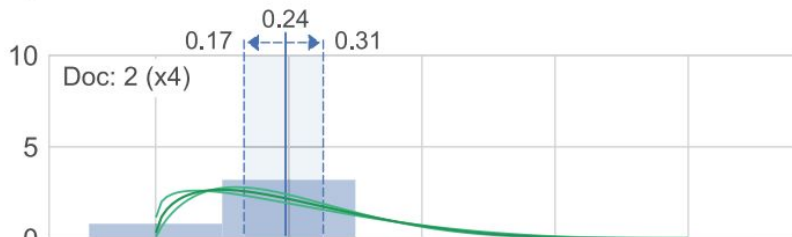
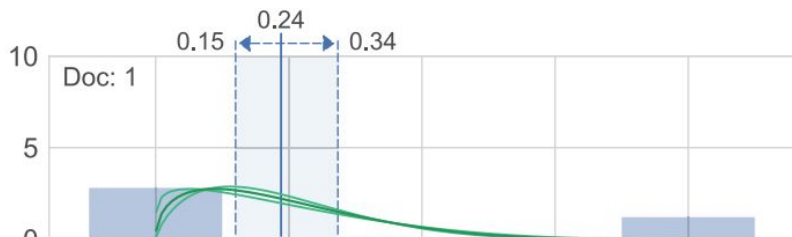
Φ Interpretation

- **High Disagreement.** When $\Phi < 0$, there is no central tendency value but rather a tendency to exclude a central area (polarized behavior)
- **Random.** When $\Phi=0$, the behavior is equivalent with a unbounded uniform process censored on the scale
- **Weak Agreement.** When $0 < \Phi \leq 0.5$, the distribution has no inflection point, but there is a unique central tendency or a dispersion that is smaller than a uniform process
- **High Agreement.** When $\Phi > 0.5$, the distribution is bell shaped with two inflection points, more narrow around the mean as Φ grows

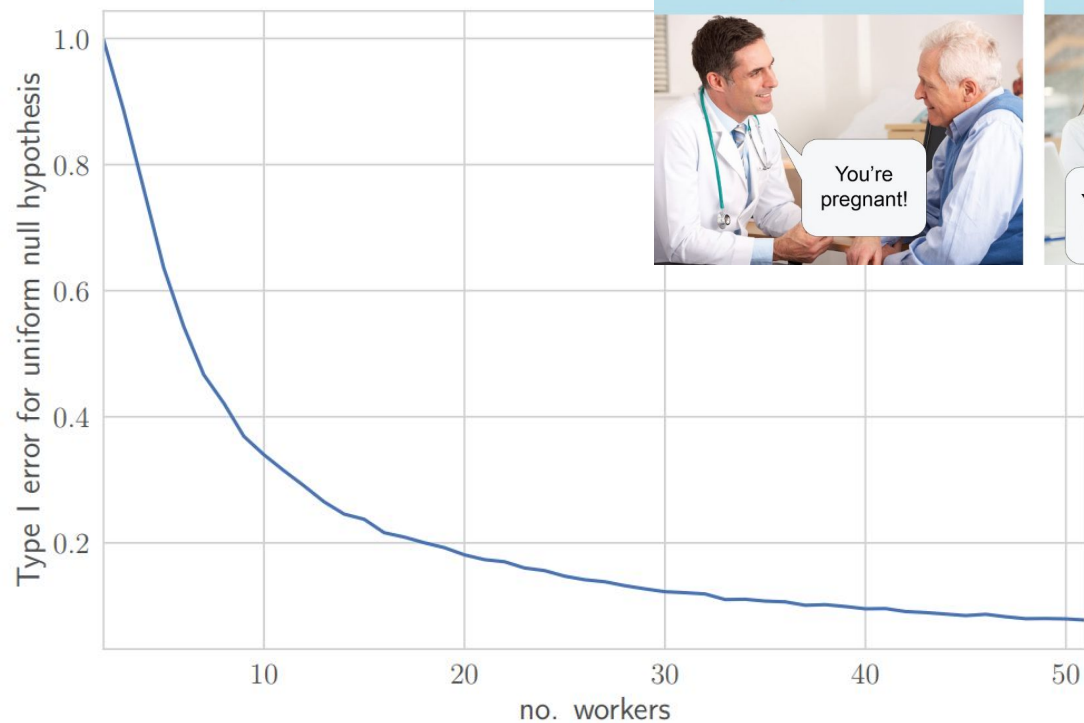
Examples of Φ Shapes



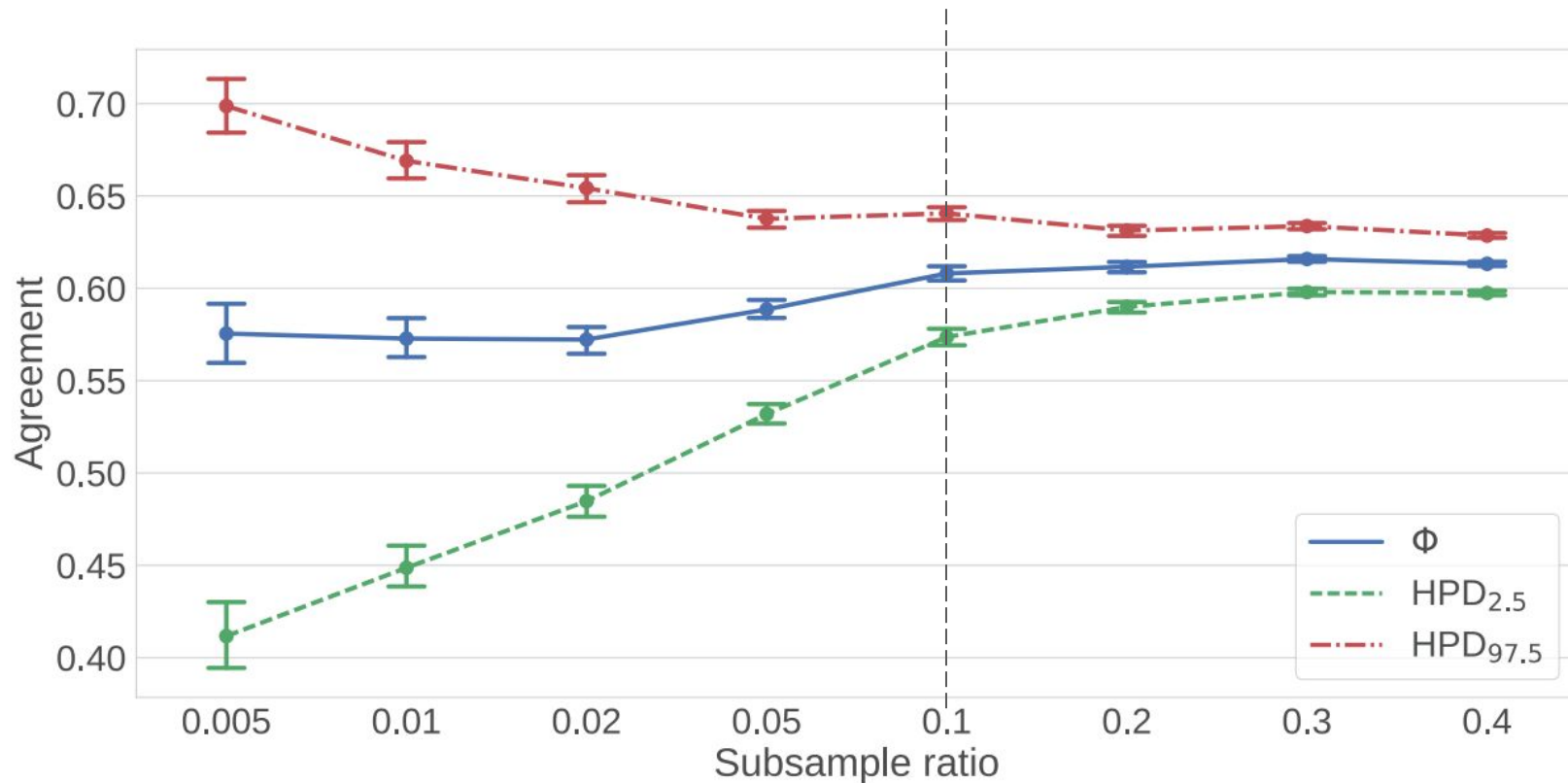
Φ in Action on Real Data



Robustness of Φ



Confidence Interval, Robustness of Φ



Incorporating Agreement in Evaluation

Setting

- test collection evaluation scenario:
 - Document collection
 - Information Needs \approx Queries (called topics)
 - A set of Information Retrieval systems (called runs)
- each system retrieves a ranked list of documents for each topic
- human made (= from experts) relevance judgments for a subset of documents of each topic
- metrics (such as Precision, Recall, NDCG, etc.) are computed
- systems are then ranked

Setting

- test collection evaluation scenario:
 - Document collection
 - Information Needs \approx Queries (called topics)
 - A set of Information Retrieval systems (called runs)
- each system retrieves a ranked list of documents for each topic
- human made (= from experts) relevance judgments for a subset of documents of each topic
- metrics (such as Precision, Recall, NDCG, etc.) are computed
- systems are then ranked



Very Expensive (Money + Time). A more economical evaluation process is needed.

Evaluation of IR systems

N Topics / Queries

M Systems

	Topic ₁	Topic ₂	...	Topic _j	...	Topic _N
system ₁	$AP_{1,1}$	$AP_{1,2}$...	$AP_{1,j}$...	$AP_{1,N}$
system ₂	$AP_{2,1}$	$AP_{2,2}$...	$AP_{2,j}$...	$AP_{2,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
system _i	$AP_{i,1}$	$AP_{i,2}$...	$AP_{i,j}$...	$AP_{i,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
system _M	$AP_{M,1}$	$AP_{M,2}$...	$AP_{M,j}$...	$AP_{M,N}$

Average
Precision
of
system i
for
topic j

Test Collection Evaluation Using Crowdsourcing

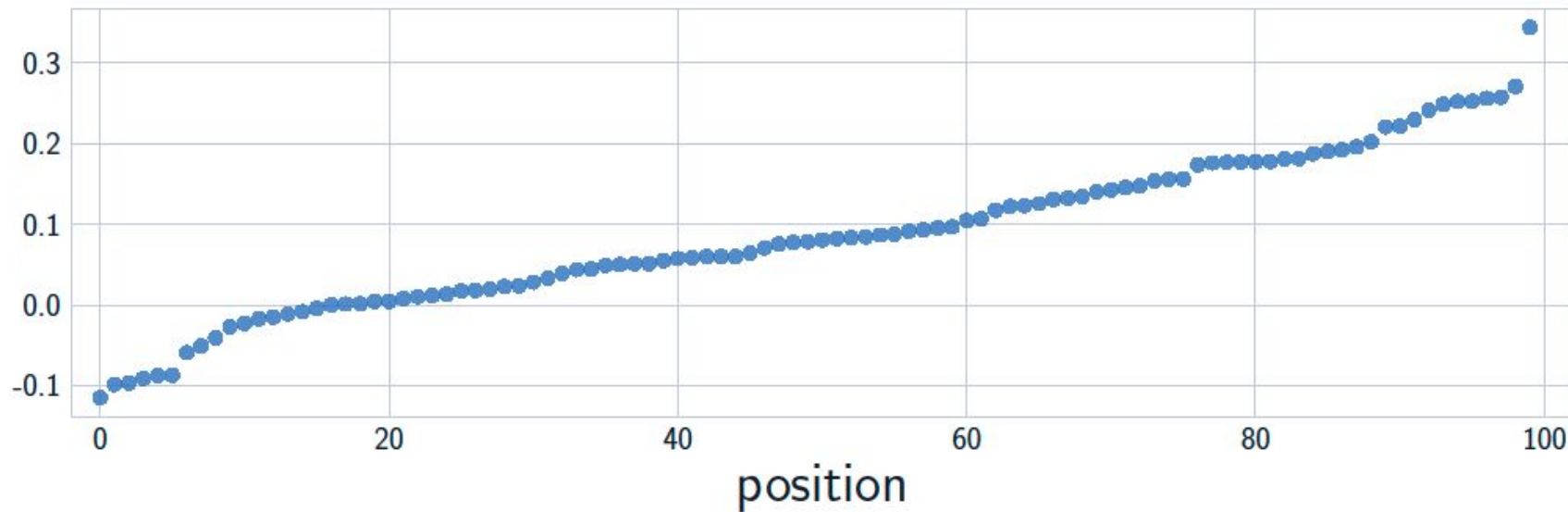
- relevance assessment is a crucial activity in IR effectiveness evaluation.
- relevance is often uncertain:
 - ambiguous query
 - terms with multiple meanings
 - unclear or ambiguous information needs
 - non-ideal work conditions of relevance assessors
 - inadequate or even malicious assessors
 - ...
- even worse when crowdsourcing is used!

Test Collection Evaluation Using Crowdsourcing

- when crowdsourcing relevance assessments:
 - each document is redundantly assessed by several crowd workers
 - workers judge the relevance of the document to a specific topic
 - all the relevance judgments for a document by different workers are aggregated to compute a final relevance score
- → only the final aggregated scores are used to compute IR evaluation metrics.
- → agreement information is lost!

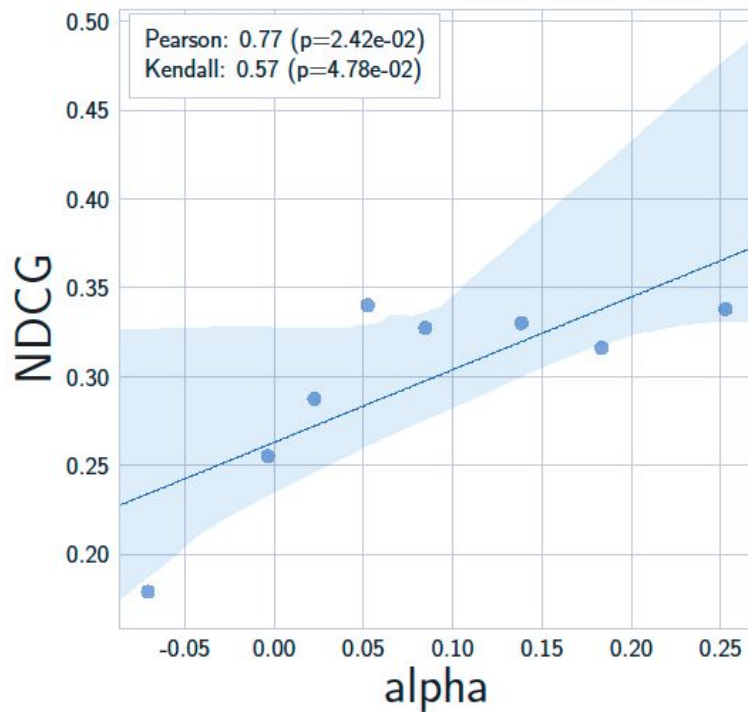
- (log all user behavior, because why not)

Different Topics, Different Agreement



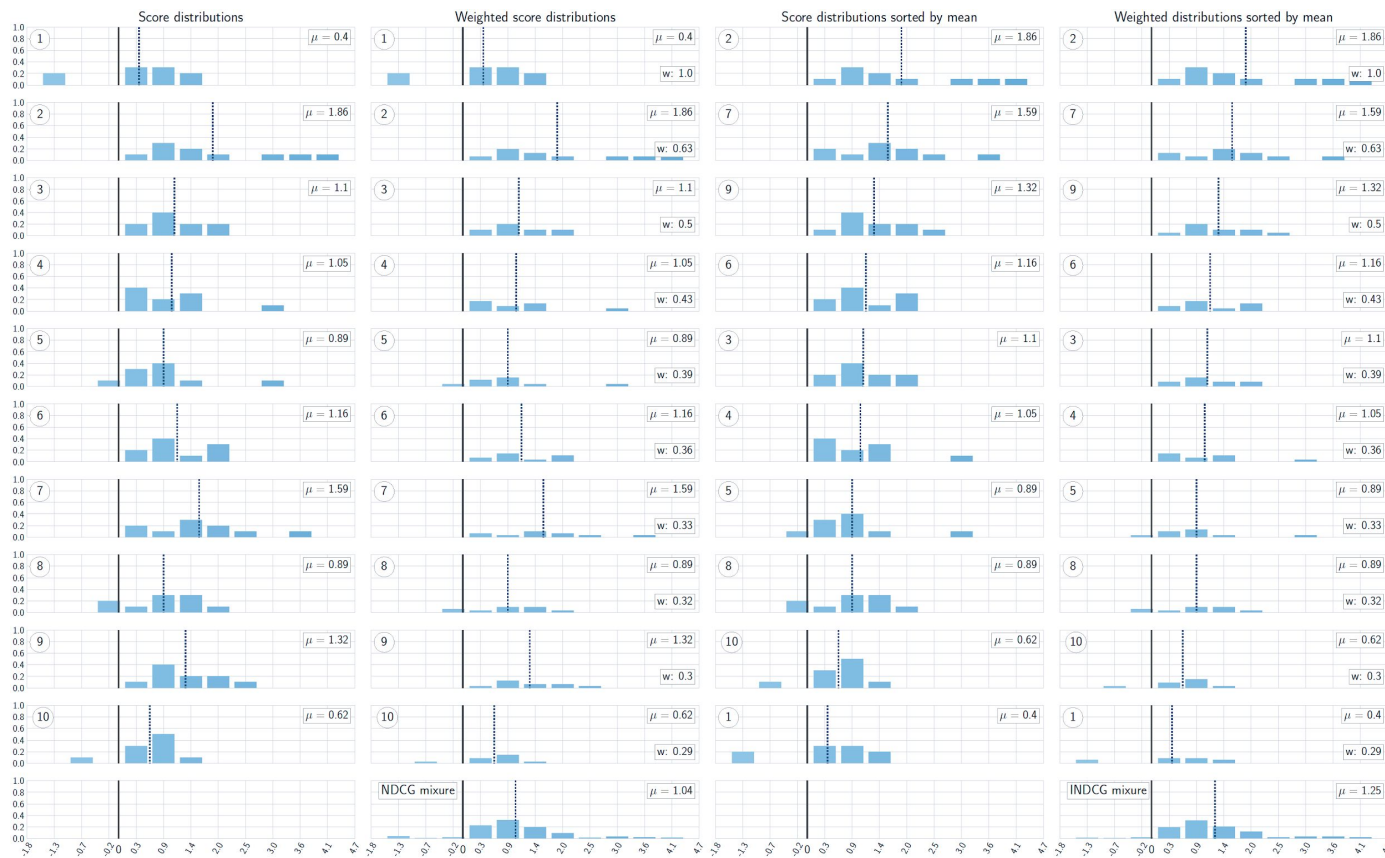
Different topics have a very different agreement

Agreement over Relevance



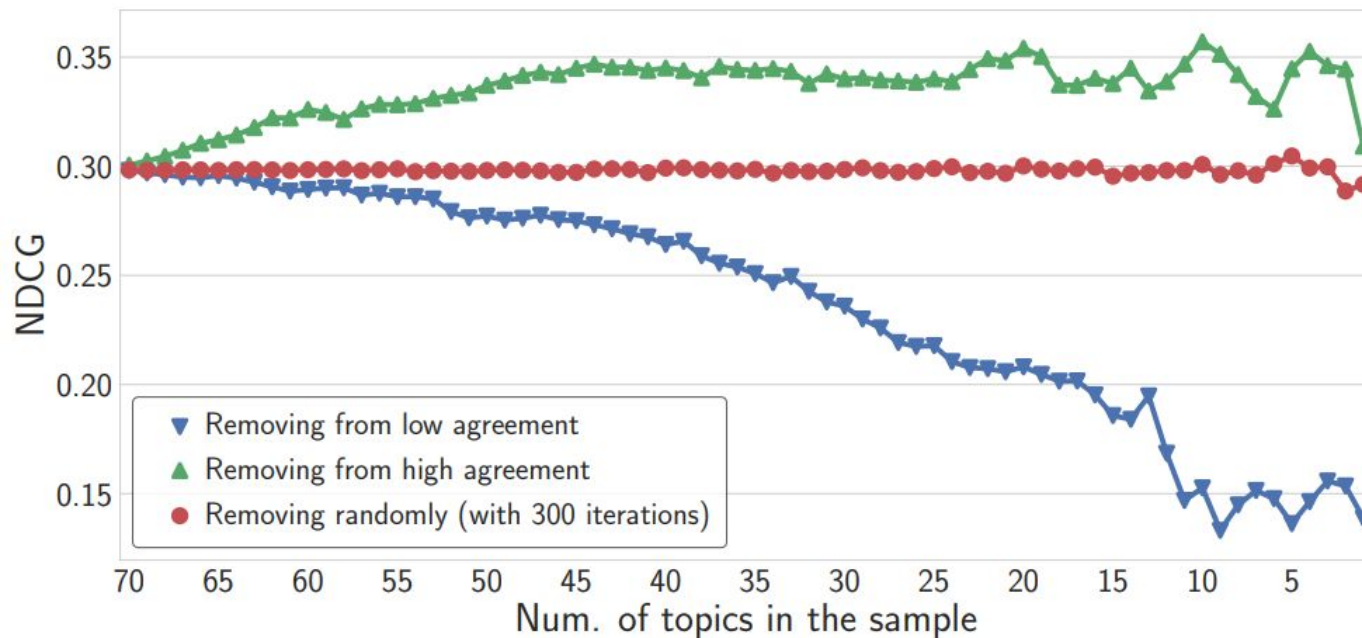
there is a relationship between agreement and effectiveness

Incorporate Agreement in the Evaluation



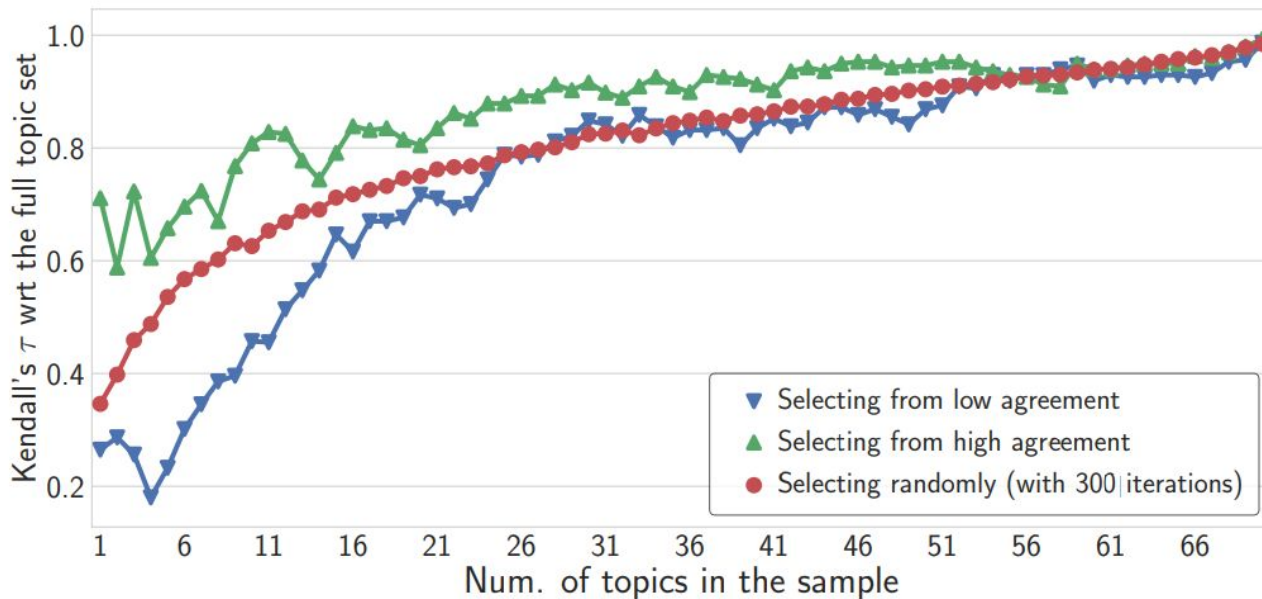
ad-hoc metrics (such as modified NDCG) can naturally incorporate agreement

Incorporate Document Agreement in the Evaluation



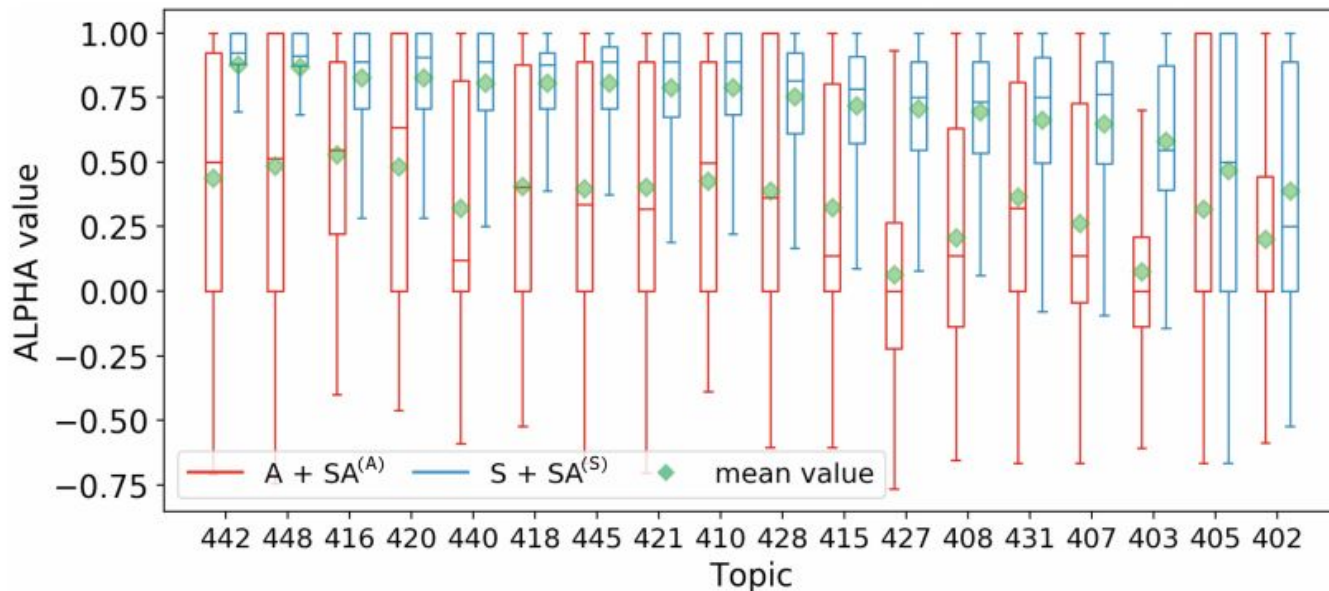
Documents with less/more agreement impact differently in the evaluation

Incorporate Topic Agreement in the Evaluation



leveraging agreement can lead to save resources

Incorporate Crowd Agreement in the Evaluation



different users have a different level of agreement w.r.t. ground truth

Control Crowd Agreement in the Evaluation

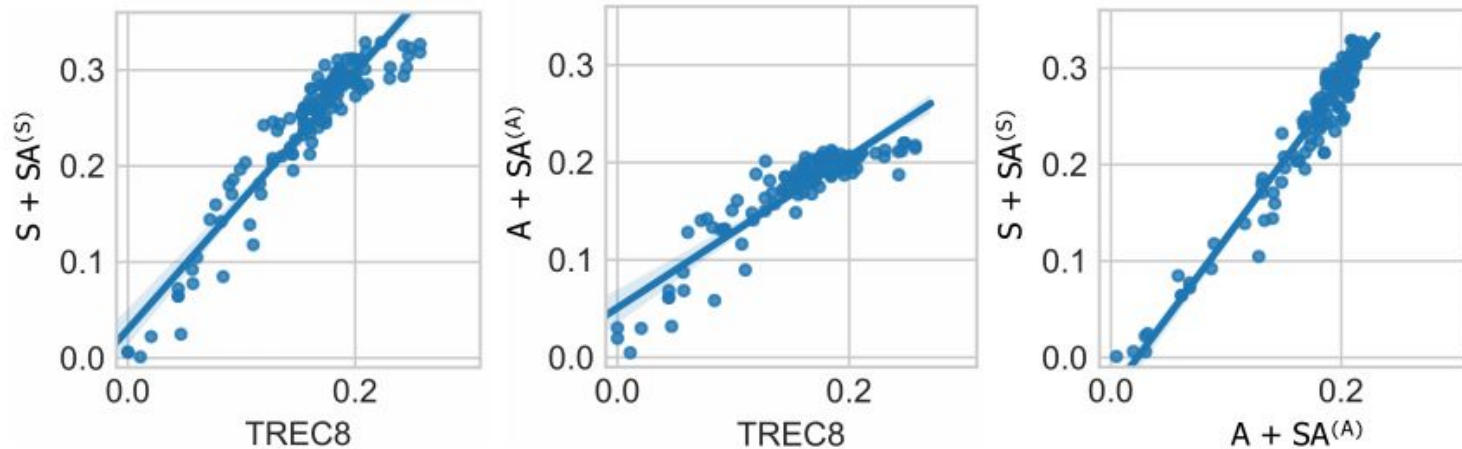
	F	Adj. <i>p</i> -value	ω^2
Two-way ANCOVA			
Number of Sessions			
Reward	76.07	$p < .001$	0.11
Task Length	113.01	$p < .001$	0.18
Reward:Task Length	43.35	$p < .001$	0.07
Number of Steps			
Reward	48.05	$p < .001$	0.10
Task Length	22.96	$p < .001$	0.05
Reward:Task Length	4.04	$p = .27$	0.01
AVG Time per Session			
Reward	0.08	$p = 1$	-0.01
Task Length	1.38	$p = 1$	0.01
Reward:Task Length	1.01	$p = 1$	0.01
One-way ANOVA			
Number of Sessions			
Quality Control	0.76	$p = 1$	-0.01
Number of Steps			
Quality Control	47.31	$p < .001$	0.09
AVG Time per Sessions			
Quality Control	65.47	$p < .001$	0.12

Short and highly paid tasks tend to maximize agreement.

Length is more important than reward.

External factors can influence User agreement

Incorporate Agreement in the System Evaluation



Different agreement levels lead to obtain less / more quality

Effect of Judgment Scale on Agreement

How many Relevance Judgments?

- Historically (until 90s) → two
- NDCG metric (SIGIR 2000) → multi-level judgements
- Sormunen (2002) → four
- Terabyte Track (2004) → three
- Tang et al. (2007) → seven (from 2 to 11)
- Web Track (2014) → six
- Magnitude Estimation (SIGIR 2016) → infinite ($0, +\infty$)
- S100 (SIGIR 2018) → $[0,100]$

	RELEVANT	NON-RELEVANT	
RETRIEVED	a	b	a + b
NOT RETRIEVED	c	d	c + d
	a + c	b + d	a + b + c + d = N (Total Collection)



A Jug
1140ml (40 fl oz)



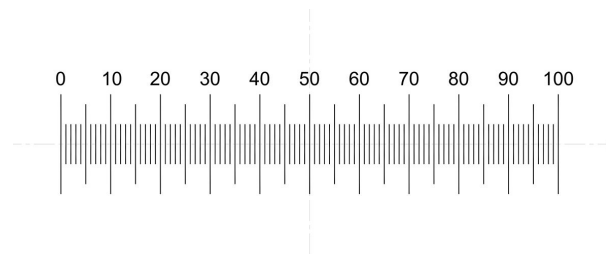
A Pint
570ml (20 fl oz)



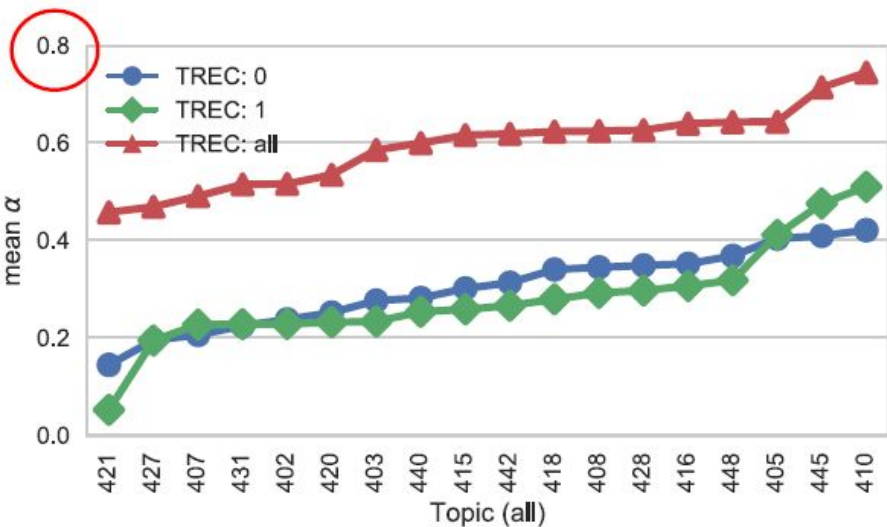
A Schooner
450ml (15 fl oz)



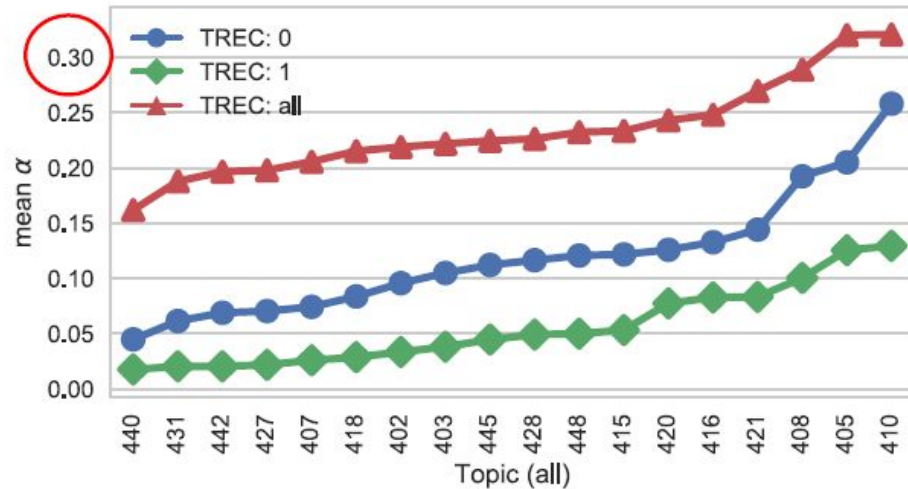
A Pot
285ml (10 fl oz)



Different Scales, Different Agreement With Experts



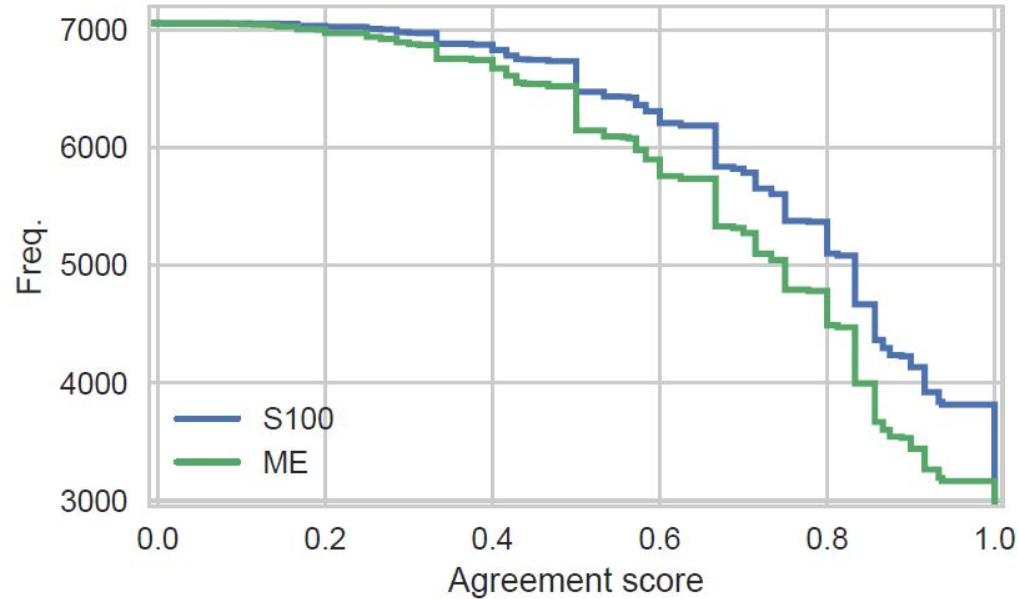
S100



ME

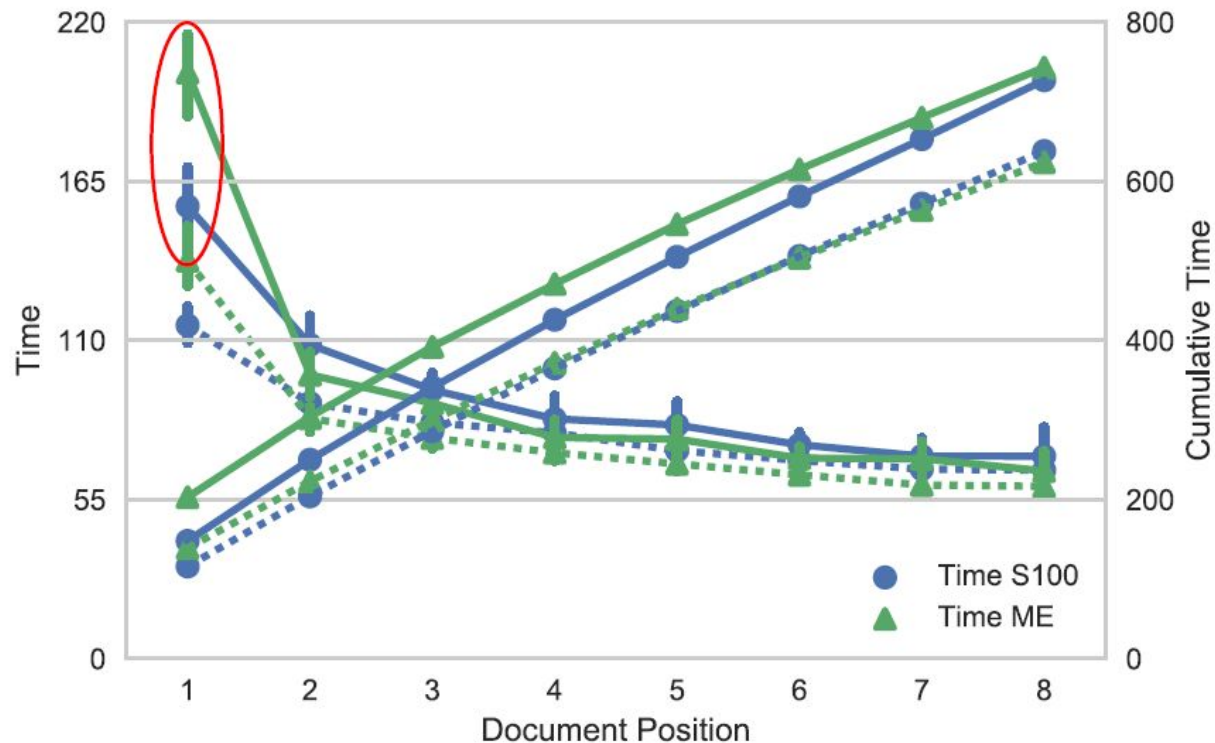
behavior is similar, but agreement is not!

Different Scales, Different Agreement With Experts



even for individual workers different scales have different agreement levels

Different Scales, Different Human Effort



some scales are more intuitive / familiar / easy to use for workers → and produce also higher agreement levels

Conclusions / Ongoing

- different scales behave very differently for evaluation
- different scales are perceived very differently from users
- (ongoing) predict churning of crowd platforms users using their agreement level compute in real time using a (deep) neural network

Relevance to Spotify Research

- use (our) metric to compute agreement between users of spotify (on listening behavior, genres, etc.)
- embed agreement into Machine Learning / Retrieval / moderation / etc. evaluation measures
- predict user churning / unsubscribing from premium / etc. leveraging user agreement
- use a subset of data with high\less agreement to train ML systems with less data, and thus use more demanding system configurations.
- use the appropriate relevance scale to compute metrics in Spotify ecosystem
- leverage agreement for automatic content moderation for songs / playlist or when crowdsourcing tasks

Resources

- HCOMP: <https://aaai.org/ocs/index.php/HCOMP/HCOMP17/paper/viewFile/15927/15258>.
- ICTIR: <https://dl.acm.org/citation.cfm?id=3121060>
- SIGIR - S100: <https://dl.acm.org/citation.cfm?id=3210052>
- SIGIR - DC: <https://dl.acm.org/citation.cfm?id=3210229>
- WWW: <https://dl.acm.org/citation.cfm?id=3291035>

- (submitted) - CIKM - transformation of relevance scales
- (submitted) - TKDE - role of abandonment (and agreement) in Crowdsourcing
- (finalizing) - TOIS - the effect of relevance scales for (crowdsourcing) IR evaluation

Bias

Bias in (IR) Evaluation

Title

N Topics / Queries

M Systems

	Topic ₁	Topic ₂	...	Topic _j	...	Topic _N
system ₁	$AP_{1,1}$	$AP_{1,2}$...	$AP_{1,j}$...	$AP_{1,N}$
system ₂	$AP_{2,1}$	$AP_{2,2}$...	$AP_{2,j}$...	$AP_{2,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
system _i	$AP_{i,1}$	$AP_{i,2}$...	$AP_{i,j}$...	$AP_{i,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
system _M	$AP_{M,1}$	$AP_{M,2}$...	$AP_{M,j}$...	$AP_{M,N}$

Average
Precision
of
system i
for
topic j

Title

N Topics / Queries

M Systems

System Effectiveness

Topic Difficulty

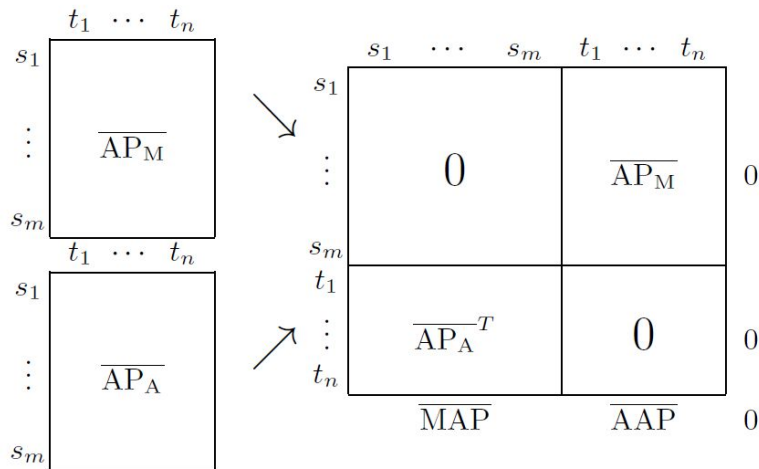
	Topic ₁	Topic ₂	...	Topic _j	...	Topic _N	MAP
system ₁	$AP_{1,1}$	$AP_{1,2}$...	$AP_{1,j}$...	$AP_{1,N}$	MAP ₁
system ₂	$AP_{2,1}$	$AP_{2,2}$...	$AP_{2,j}$...	$AP_{2,N}$	MAP ₂
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
system _i	$AP_{i,1}$	$AP_{i,2}$...	$AP_{i,j}$...	$AP_{i,N}$	MAP _i
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
system _M	$AP_{M,1}$	$AP_{M,2}$...	$AP_{M,j}$...	$AP_{M,N}$	MAP _M
AAP	$[AP_{1,1}$	$AP_{1,2}$...	$AP_{1,j}$...	$AP_{1,N}]$	

Building a Graph



how much a
system thinks a
topic is difficult

how much a topic
thinks a system is
effective

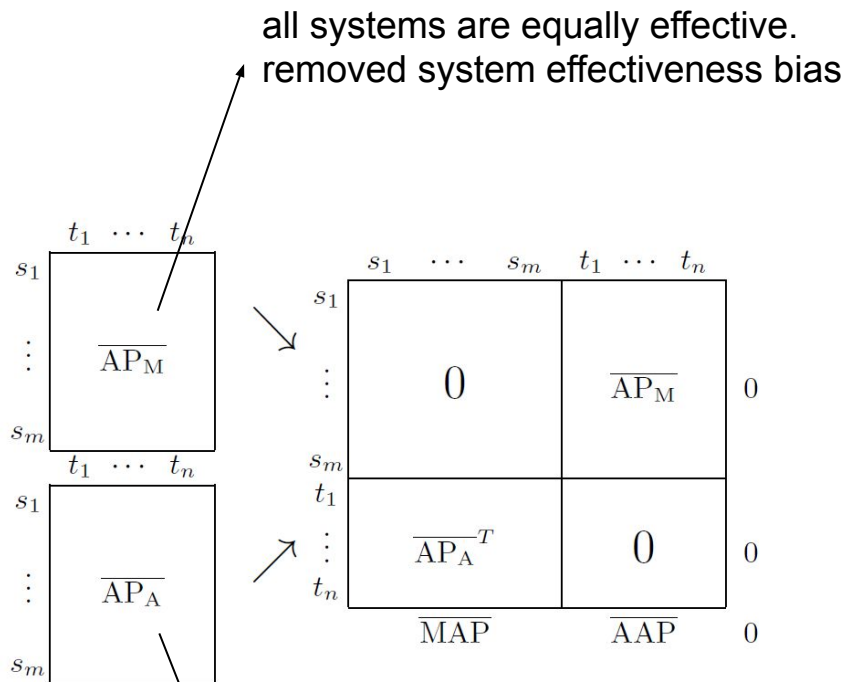


Building a Graph



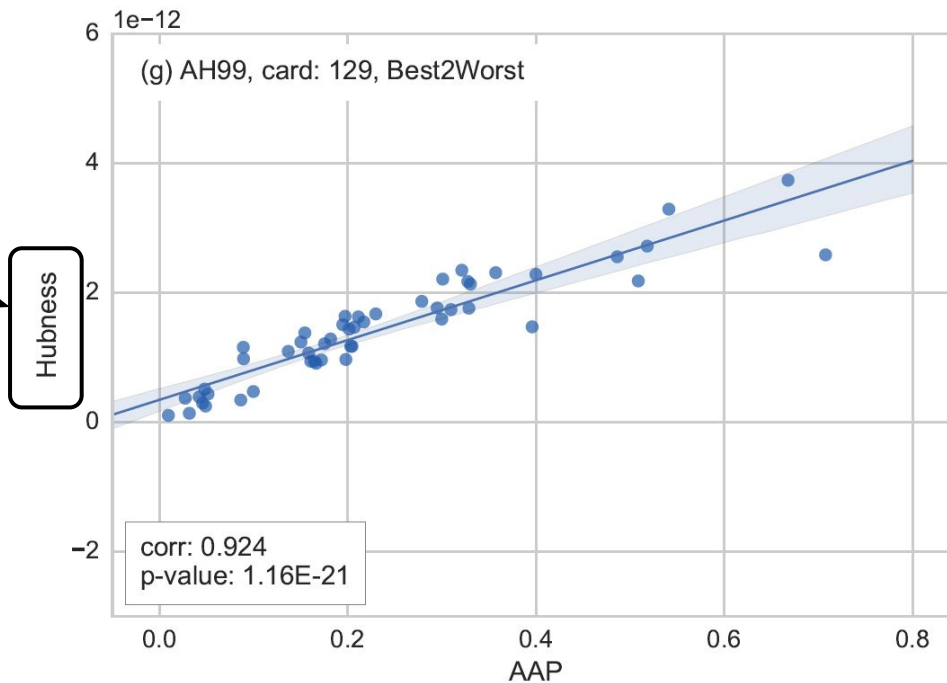
how much a
system thinks a
topic is difficult

how much a topic
thinks a system is
effective



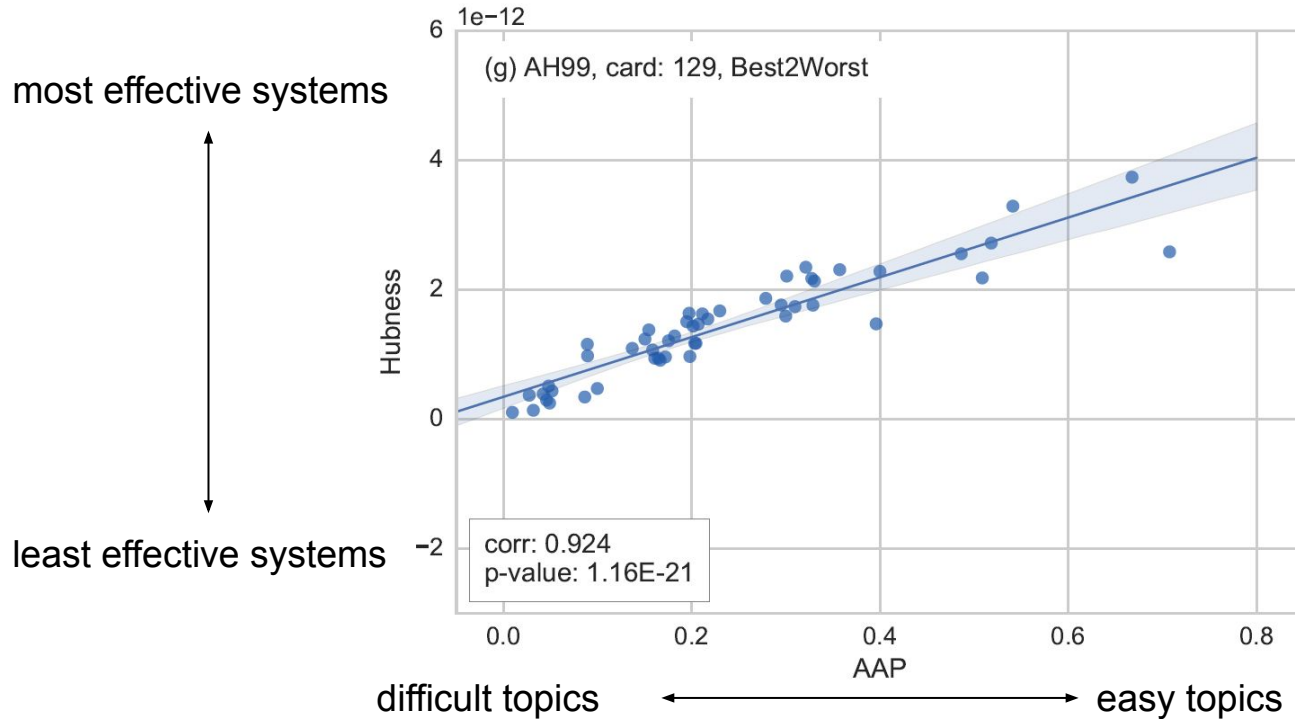
Results

ability of a topic to identify effective systems



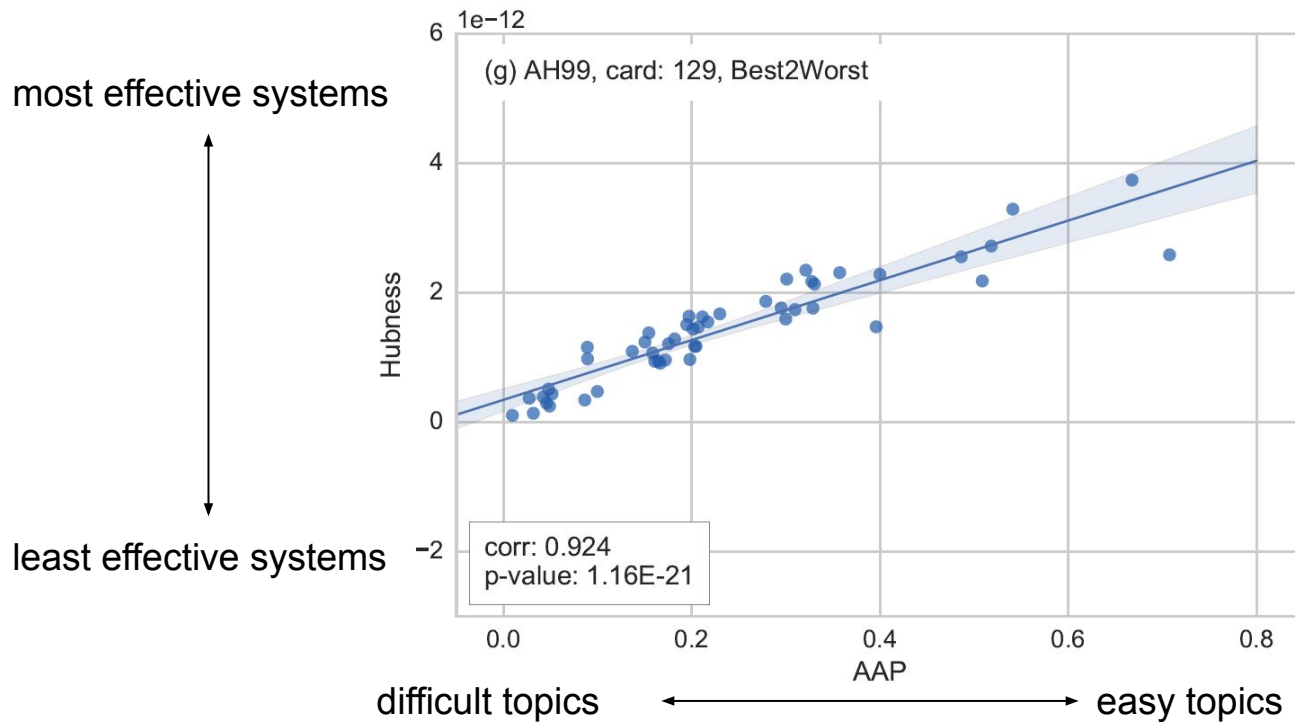
Results

•



Results

•



easy topics are better in identify the final rank of retrieval systems.

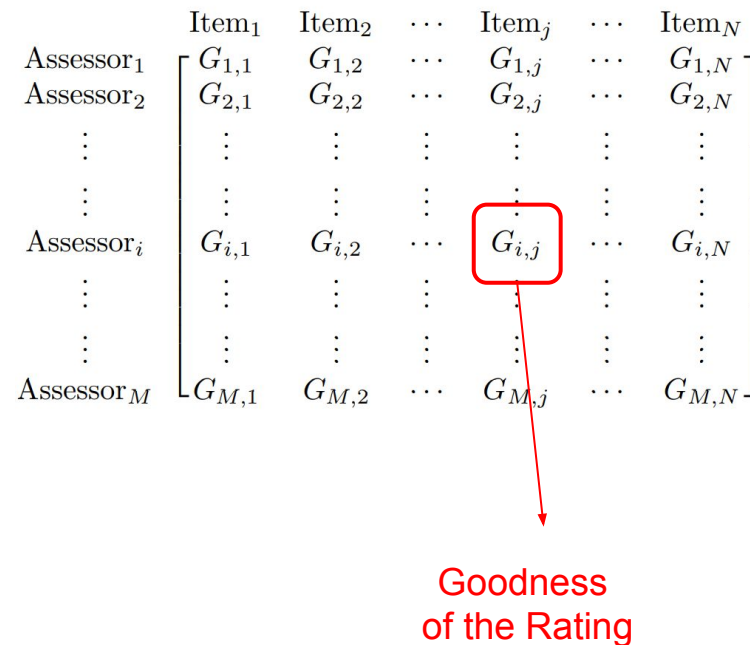
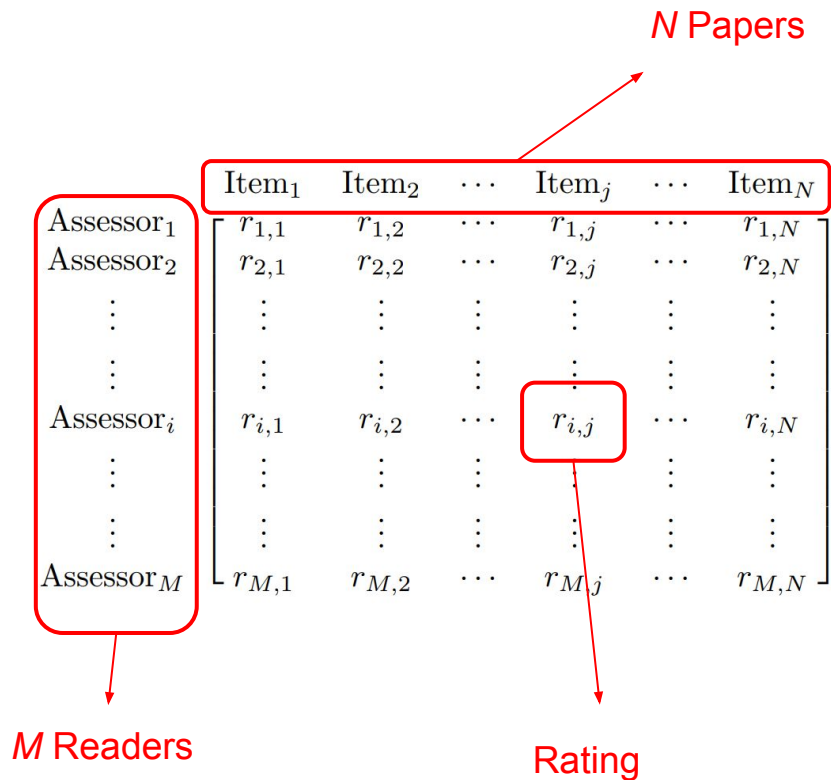
→ issue in the evaluation process, → it can be leveraged to find more informative topics, → bias identification

Bias in Scholarly publishing

Scholarly publishing / Readersourcing

- Scholarly publishing:
 - hard work
 - write + submit
 - Peer Review
 - accept, or goto(hard work)
- not the only possible model
 - open access
 - wikipedia like approaches
- peer review issues
 - slow, biased, subjective
 - unable to detect fraud
 - Page Rank being rejected from SIGIR1998...
 - ...
- readersourcing : Crowdsourcing peer review
 - readers read papers and express opinions
 - a formalized model

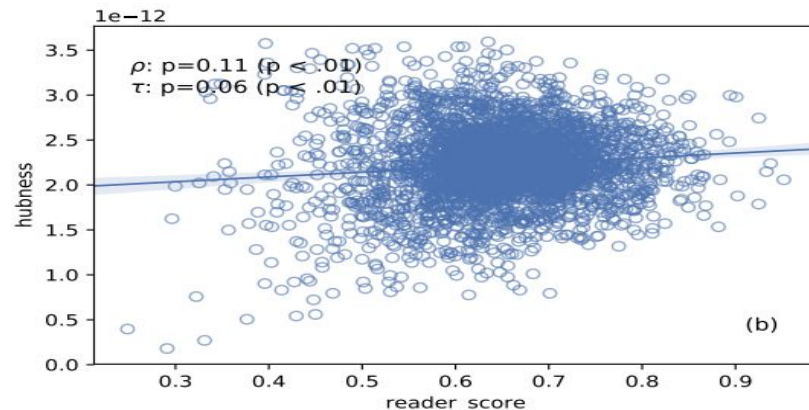
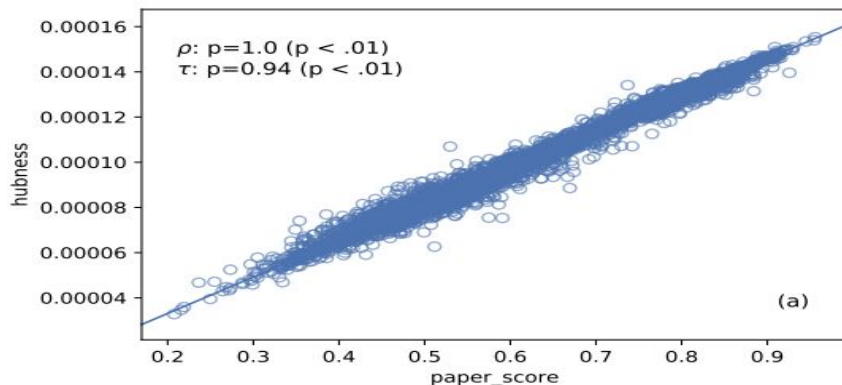
Scholarly publishing / Readersourcing



Removing Bias

- reader bias:
 - all readers vote with the same median
 - all reader vote with the same judgment scale
 - all readers vote with the same quality
 - (...)
- paper bias:
 - all papers have the same average judgment
 - all papers are judged using the same perceived scale
 - all papers have the same quality
 - (...)
- mixed approaches are possible

Results / Potential



- see if the model has some kind of bias
 - high quality papers are judged by high quality readers
 - readers that overestimate true judgment have an high quality
 - ...
- remove bias from the model
 - are the best papers still best papers after removing the bias?
 - ...

Relevance to Spotify

- remove / identify various user biases
 - listen lot of songs
 - free / premium
 - skips lots of songs
 - ...
- see if recommender system evaluation is biased for some users and correct the bias
- identify playlists that are biased in the recommendation phase:
 - playlists that are always recommended at lunch time
- force / embed bias in playlists commendations
 - If I like to have dinner listening to piano songs then ...

Resources

- ECIR: http://dx.doi.org/10.1007/978-3-319-56608-5_55.
- (Submitted) BIRNDL: bias in scholarly publishing models
- (finalizing) Experiments on bias in scholarly publishing for IIR
- (in progress) Experiments on bias in scholarly publishing for a Journal paper

Economical Evaluation of (IR) Systems

Economical Evaluation of IR systems

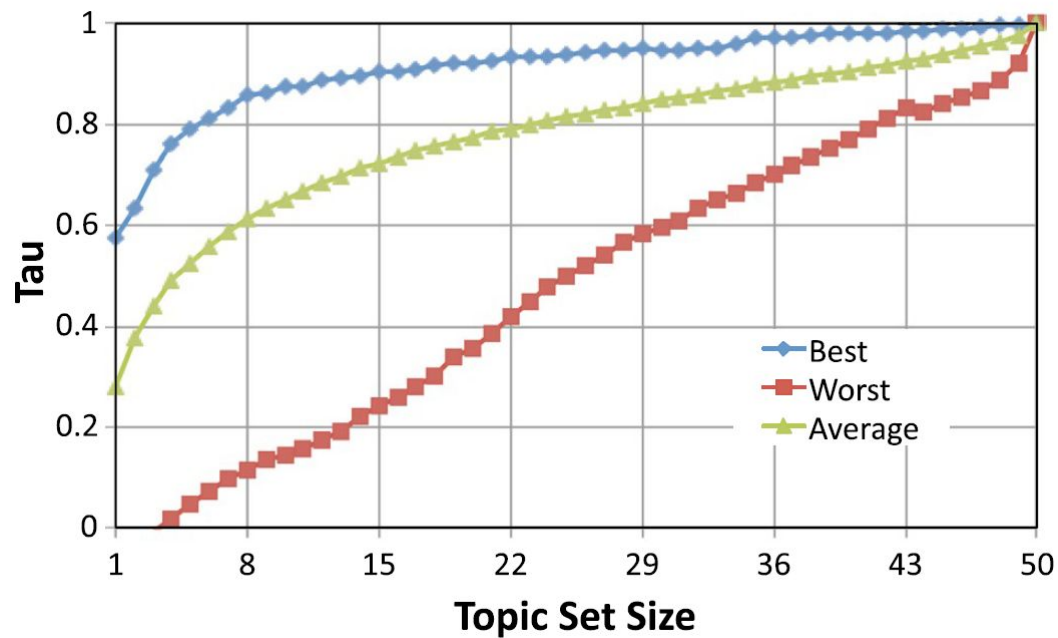
N Topics / Queries

M Systems

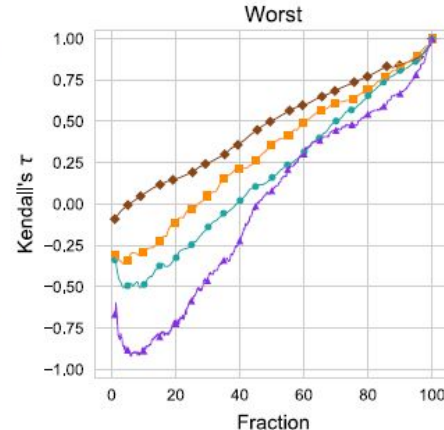
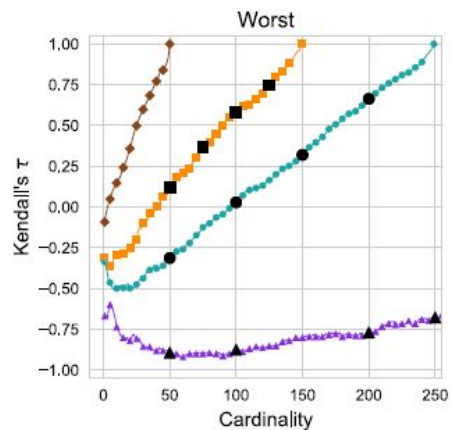
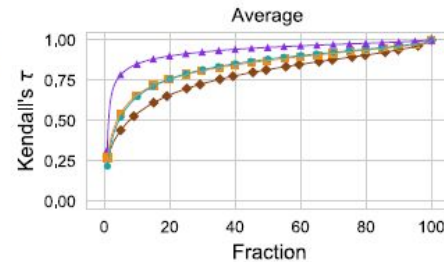
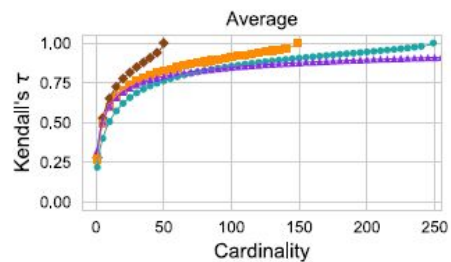
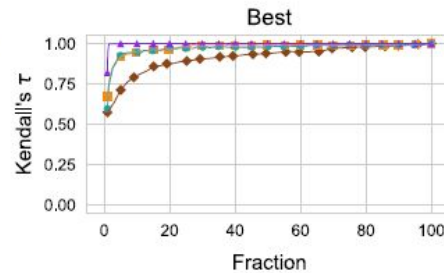
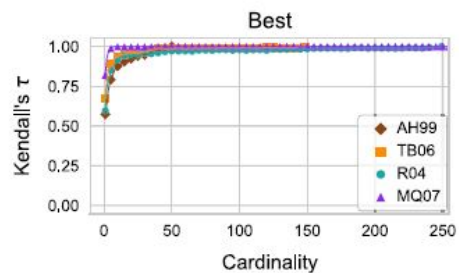
	Topic ₁	Topic ₂	...	Topic _j	...	Topic _N
system ₁	$AP_{1,1}$	$AP_{1,2}$...	$AP_{1,j}$...	$AP_{1,N}$
system ₂	$AP_{2,1}$	$AP_{2,2}$...	$AP_{2,j}$...	$AP_{2,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
system _i	$AP_{i,1}$	$AP_{i,2}$...	$AP_{i,j}$...	$AP_{i,N}$
⋮	⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮	⋮
system _M	$AP_{M,1}$	$AP_{M,2}$...	$AP_{M,j}$...	$AP_{M,N}$

Average
Precision
of
system i
for
topic j

What About using Few Topics?

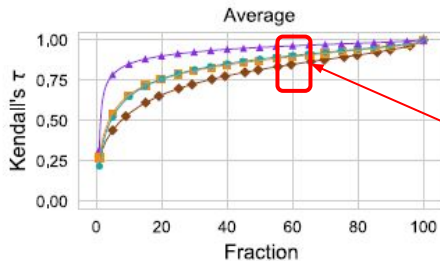
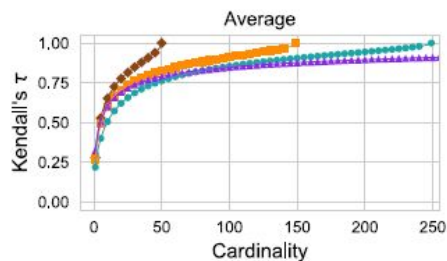
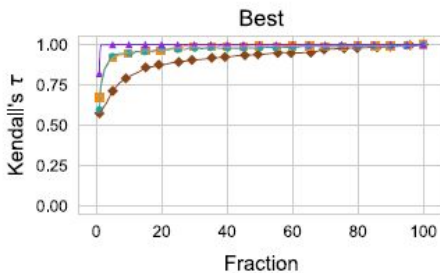
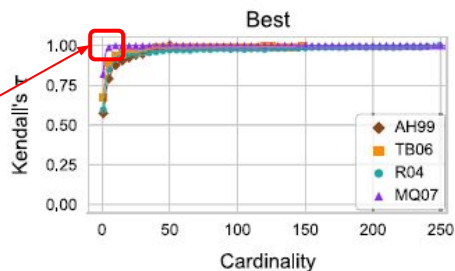


Few Topics

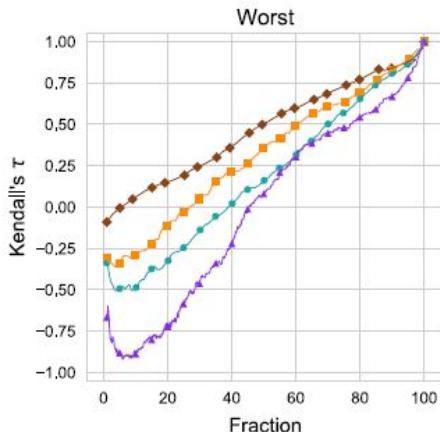
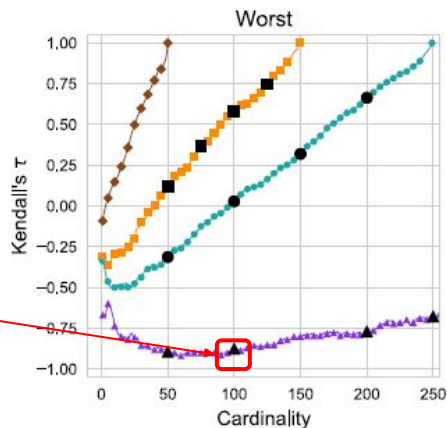


Few Topics

correlation ≈ 1 with 10% of the dataset



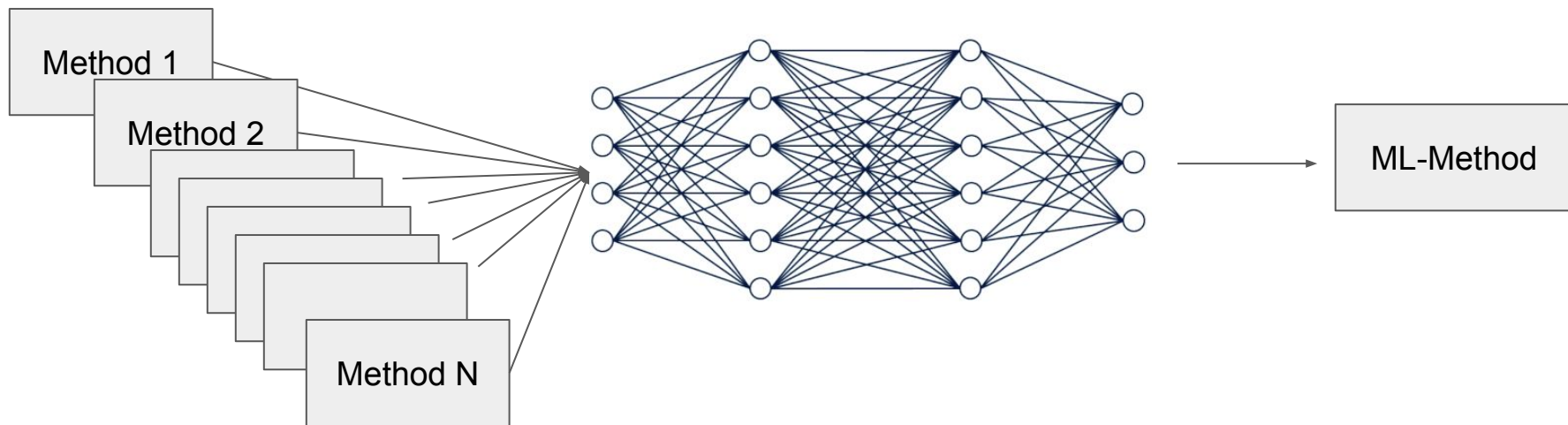
correlation ≈ -1 with 10-20% of the dataset



random guessing is not a good strategy

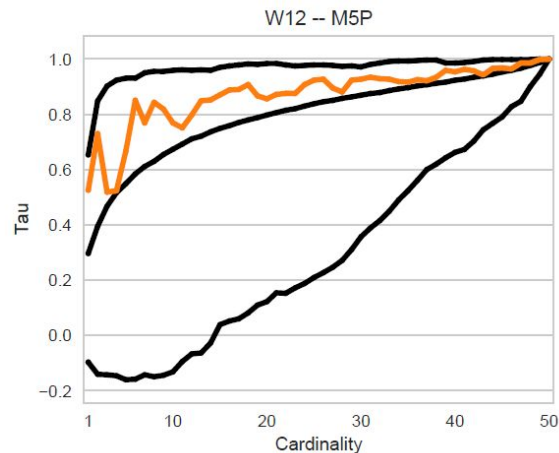
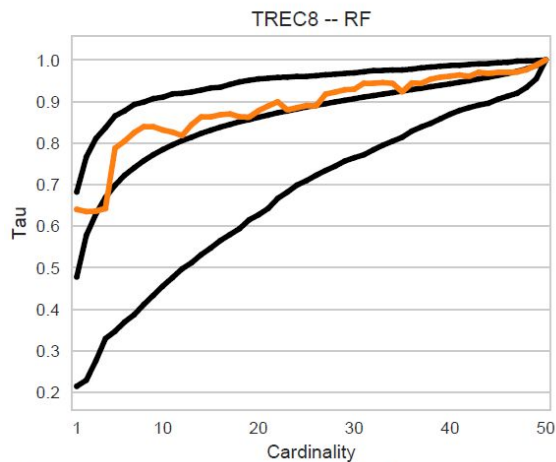
A practical approach to Few Topics

- Evaluation Without Relevance Judgments (EEWRJ)
 - evaluate retrieval systems based on their ranked lists of documents
 - no human effort \approx no cost!
- strategy based on individual EEWRJ methods + ML combinations of them



A practical approach to Few Topics

- Evaluation Without Relevance Judgments (EWRJ)
 - evaluate retrieval systems based on their ranked lists of documents
 - no human effort \approx no cost!
- strategy based on individual EWRJ methods + ML combinations of them



Relevance to Spotify

- train Spotify recommender system using (much) less data
- identify the subset of good (= more informative) users / playlists / etc.
- use fusion methods and ML combinations of weak EEWRJ-like signals
 - Spotify Search
 - ensembles / cascade of Recommender systems
 - ...

Resources

- IRJ: <https://rdcu.be/bFuZq>
- SIGIR2018: <https://dl.acm.org/citation.cfm?id=3210108>
- JDIQ : <https://dl.acm.org/citation.cfm?id=3239573>
- JDIQ : <https://dl.acm.org/citation.cfm?id=3241064>
- (rebuttal) - IPM - EEWJRJ overview and combination of approaches
- many ongoing experiments:
 - Neural Networks with custom loss to build a good model
 - GANs and Genetic Algorithms to generate fake data for data augmentation
 - Crowdsourcing query variations for retrieval fusion and for data augmentation
 - Multi Armed Bandits to select topics / Documents using “no relevance judgments” signals
 - ...

Thank you!