# Evaluation of IR Systems

Günther Specht
Eva Zangerle

Summer Term 2019

# Evaluation of IR Systems

- How to measure how 'good' a retrieval system performs?
- Without any evaluation hardly any improvements are possible
- Need to justify changes


- Today:
  - Evaluation overview
  - Evaluation metrics for (automatic) evaluations
  - Kappa
  - Accuracy
  - Precision@k, Recall@k
  - Mean Average Precision
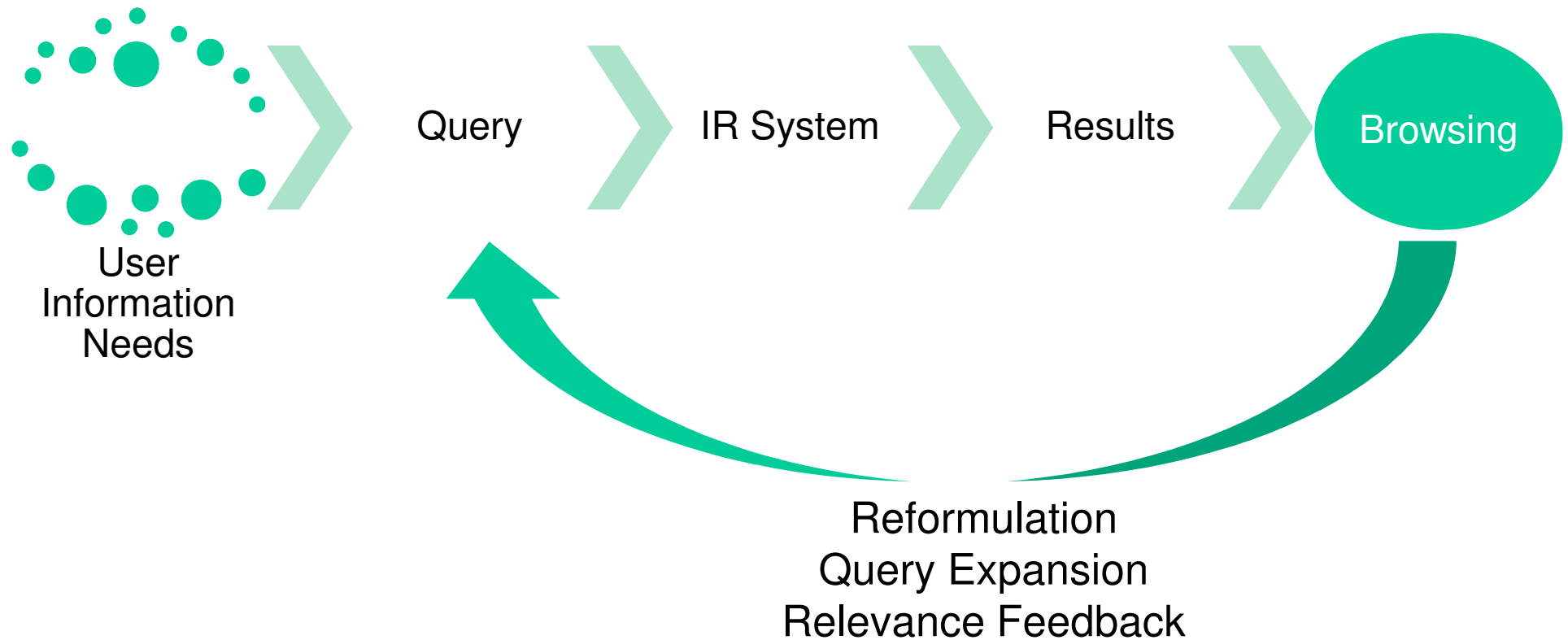  - User-based Evaluation
  - Crowdsourcing

# Evaluation of IR Systems

- First attempt
  - Speed of query evaluation
    - Queries per second evaluated
  - Speed of indexing
    - Number of tweets added to the index per second
  - Footprint of index

  - Query expressivity
    - No wildcards
    - Simple keyword search

- Factors above can be measured
- What about the user?

- Most important measure: user happiness ☺
  - **Relevance of retrieved documents**
  - Speed of response, latency
  - Size of indexed documents

  - Also: user interface design
    - Responsiveness
    - Clarity
    - Layout

# Information Needs



User
Information
Needs

Query

IR System

Results

Browsing

Reformulation
Query Expansion
Relevance Feedback

# User Happiness

- What are criteria for user happiness?

  in
    - (Web) search engine
    - Online shop
    - Social network

# Evaluation of IR Systems

- **Relevance**
  - Relevance related to information need (not the query itself)
    - Document containing all search keywords may not be relevant at all in regards to information need
    - Cf. apple vs. Apple

  - Binary (relevant or non-relevant)

- Translation of information need to query:
  - Sometimes difficult. Example:
    Information need: Information on whether drinking red wine is more effective at reducing your risk of heart attacks than white wine.

    Query: wine and red and white and heart and attack and effective

# Measuring Ad-hoc Effectiveness

- Test collection required
  - Document collection
  - Queries:  Test suite of information needs (queries expressing these information needs)
  - Relevancy judgements for query-document pairs
    - Team of human judges
    - At least one judge per document

  - Ground truth

# Test Collections, e.g

- Cranfield Collection
  - First evaluation collection
  - Collection is fully evaluated
  - 1.400 documents (abstracts of aerodynamics journals)
  - 225 queries
  - Imporant, since first one; today: too small.

- TREC
  - Text Retrieval Conference
  - 1.89 mio documents (newswire articles) in total
  - 450 information needs (topics) in total
  - TREC 6-8: 150 information needs for 528,000 documents
  - trec.nist.gov

- PAN Workshops
  - For author ship attribution, plagiarism detection etc.

# TREC Example

```
<top>
<num> 52 </num>
<title> Topic:  Accusations of Cheating by Contractors on U.S. Defense
Projects</title>
<desc> Description:
Document will refer to an alleged illegality committed by any entity
seeking a contract on behalf of the U.S. Military Forces. </desc>
<narr> Narrative:
A relevant document will mention an alleged impropriety or improprieties
by individuals or entities (companies, corporations), either domestic or
foreign, attempting to provide services or products related to the U.S.
military effort.  To be relevant, document will specifically:
(1) identify the wrongdoer, and
(2) describe in at least general terms (e.g., overcharging, bribing,
swindling, illegal gifts, bid-riggings found, insider information) the
nature of the wrongdoing. </narr>
</top>
```

# Assessing Relevance

- Manual task
- User information needs have to be realistic (not randomly generated keywords) → realistic distribution
- Relevance assessment tedious, costly, time-consuming

- Humans for gold standard judgement?
  - Different opinions, backgrounds, expections
  - But: systems are used by the same users
- Multiple judges for relevance assessment

- Measurement of agreement of judges on relevance judgements?

# Kappa Statistic

- Measure for agreement of judges: Kappa statistic

$$kappa = \frac{P(A) - P(E)}{1 - P(E)}$$

- P(A) proportion of the times the judges agreed
- P(E) proportion of the times they would be expected to agree by chances (next slide)
  - Two-class decision: chance of agreement is 0.5

- kappa = 1 for total agreement
- kappa = 0 agreement per chance (dt. Zufall)

# Kappa Statistic

| | | Judge 2 | | |
|---|---|---|---|---|
| | | yes | no | total |
| Judge 1 | yes | 300 | 20 | 320 |
| | no | 10 | 70 | 80 |
| | total | 310 | 90 | 400 |

- P(A) = (300 + 70) / 400 = 0.925
- P(relevant) = (320 + 310) / (400 + 400) = 0.7878
- P(nonrelevant) = (80 + 90) / (400 + 400) = 0.2125
- P(E) = P(relevant)$^2$ + P(nonrelevant)$^2$ = 0.665
- Kappa = 0.776

# Kappa Statistic

Manning:
- 0    -0.66 = seen as dubious base for evaluation
- 0.67-0.80 = fair agreement
- 0.81-1.00 = good agreement
- 

Landis and Koch:
- 0.41-0.60 = „mittelmäßige (moderate) Übereinstimmung",
  0.61-0.80 = „beachtliche (substantial) Übereinstimmung",
  0.81-1.00 = „(fast) vollkommene ((almost) perfect)
  Übereinstimmung

- Within TREC dataset, normally kappa > 0.67, < 0.8

# Unranked Retrieval Evaluation

User

Information System

| | Relevant | Not Relevant |
|---|---|---|
| Retrieved | True Positive | False Positive |
| Not Retrieved | False Negative | True Negative |

- ## Precision
  - Fraction of relevant documents which are retrieved
  - $precision = \dfrac{\#\ relevant\ retrieved\ docs}{\#retrieved\ docs} = \dfrac{\#tp}{\#tp+\#fp}$

- ## Recall
  - Fraction of documents which are relevant
  - $recall = \dfrac{\#\ relevant\ retrieved\ docs}{\#relevant\ docs} = \dfrac{\#tp}{\#tp+\#fn}$

|  | relevant | not relevant |
|---|---|---|
| retrieved | true positive | false positive |
| not retrieved | false negative | true negative |

# Accuracy

- Accuracy
  - Number of correct classifications

  - $accuracy = \dfrac{\#tp + \#tn}{\#tp + \#tn + fn + \#fp}$

  - For most queries, 99% of all documents are irrelevant
  - Optimize accuracy: label all documents as non-relevant → accuracy = 99%
    But: empty resultset not satisfying for user

- Recall vs. Precision depends on aim of information system
  - Consider web search: better optimize for recall or for precision?

# F-Measure

- F-measure (or F1-measure)
  - Combination of recall and precision
  - Harmonic mean of recall and precision

  - $F1 = \dfrac{2 * P * R}{P + R}$

  - Why harmonic mean and not
    - Arithmetic mean (a+b)/2
    - Geometric mean (a*b)^1/2

    - Harmonic Mean is conservative, not that affected by outliers

# F-Measure



Precision (Recall fixed at 70%)

20

# Ranked Retrieval Evaluation

# Ranked Retrieval Evaluation

- Ranked Retrieval
- Precision, recall and F1 also applicable for ranked result sets

- But: ranking of documents not reflected
  → Precision, recall and F1 are set-based measures

# Single Value Summaries

- Single value summaries required to quickly grasp performance of information retrieval system

- <span style="color:red">Precision@k, recall@k</span>
  - Precision/recall of top-k results
  - Evaluate precision/recall performance of top-k ranked results
  - Mostly used for web search and recommender systems

Recall per Similarity Measure (scoreRank)

Precision per Similarity Measure (scoreRank)

# Precision-Recall Curve

- Precision-Recall curve for top-k results (blue line) (typically: with sawtooth shape)



What happened here?

Sawtooth shape: consider recall of documents k+1 and k if k+1 is nonrelevant

Figure taken from Manning, Raghavan, Schütze: Introduction to Information Retrieval, Cambridge University Press, 2008.

# Interpolated Precision

- **Interpolated precision**
  - No sawtooth shapes, smoothed curve

$$p_{inter} = \max_{r' \geq r} p(r')$$

  - Why?
  - You would look three documents further if you knew that there is a relevant document!

# Interpolated Precision

| recall | interpolated precision (red) |
|--------|------------------------------|
| 0.0 | 1.00 |
| 0.1 | 0.67 |
| 0.2 | 0.63 |
| 0.3 | 0.55 |
| 0.4 | 0.45 |
| 0.5 | 0.41 |
| 0.6 | 0.36 |
| 0.7 | 0.28 |
| 0.8 | 0.13 |
| 0.9 | 0.10 |
| 1.0 | 0.08 |

# Eleven Point Interpolated Avg. Precision

- **Eleven point interpolated average precision**
  - Computed for each recall level as in prev. slide
  - For all tests
  - Averaged over whole set of information needs at the eleven levels 0.0, 0.1, 0.2, ...0.9, 1.0

# Mean Average Precision

- Alternative:
- MAP – Mean Average Precision
  - Avg. precision = average precision value for top-k documents
  - Average of all information needs which are evaluated
  - Relevant documents for information need $q_j \in Q = \{d_1, \ldots d_{m\_j}\}$
  - R$_{jk}$ = set of ranked retrieval results (Q queries, $m\_j$ docs per Query)

$$MAP(Q) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} \frac{1}{m} \sum_{k=1}^{m_j} precision(R\_jk)$$

  - For single information need: avg. precision approx. area under uninterpolated precision-recall curve
  - Stable measure

# User-based Evaluation

# User-Based Evaluation

- Up until now: user happiness and performance assessed by automated tests and automatically computed metrics.

- What about not directly measurable influence factors?
  - Usability
  - User interface design
  - User interaction with system

- Changes on these factors have to be done carefully
- Pre-tests

# UI Changes



Figure taken from http://googlesystem.blogspot.co.at/

# Human Experimentation

- Small number of test users
- Make small changes to UI, observe human behaviour and assess preferences
  - Video
  - Eyetracking
  - Questionnaire
- Evaluations are run in a lab
- Human users have to be selected carefully, choice of users is crucial for evaluation outcome
- Costly
- Users know that they are observed and their behaviour is evaluated

# Eye Tracking



Figure taken from S. Castagnos, N. Jones, P. Pu: "Eye-tracking product recommenders' usage"

# Eye Tracking



Figure taken from S. Castagnos, N. Jones, P. Pu: "Eye-tracking product recommenders' usage"

# Side-by-Side Panels

- Evaluate results produced by two systems side by side
- Perfect for search engine result lists
- Side by side: system A on left side of screen, system B on right side of screen (e.g. results for same query)

- Users are aware that their behaviour is evaluated
- Evaluations can be done anonymously

# Side-by-Side Panels



Figure taken from P. Bailey, P. Thomas, D. Hawking: Does brandname influence perceived search result quality? Yahoo!, Google, and WebKumara

# A/B and Multivariate Testing

Online tests
- Show modified version of (mostly) webpage to pre-selected amount of users and observe reaction
- A/B tests: compare two versions
- Multivariate tests: compare multiple versions
- Huge numbers of participating users possible, preselection based on e.g. demographics

- Users are not aware of being part of an evaluation
- Large number of users
- Many evaluations possible
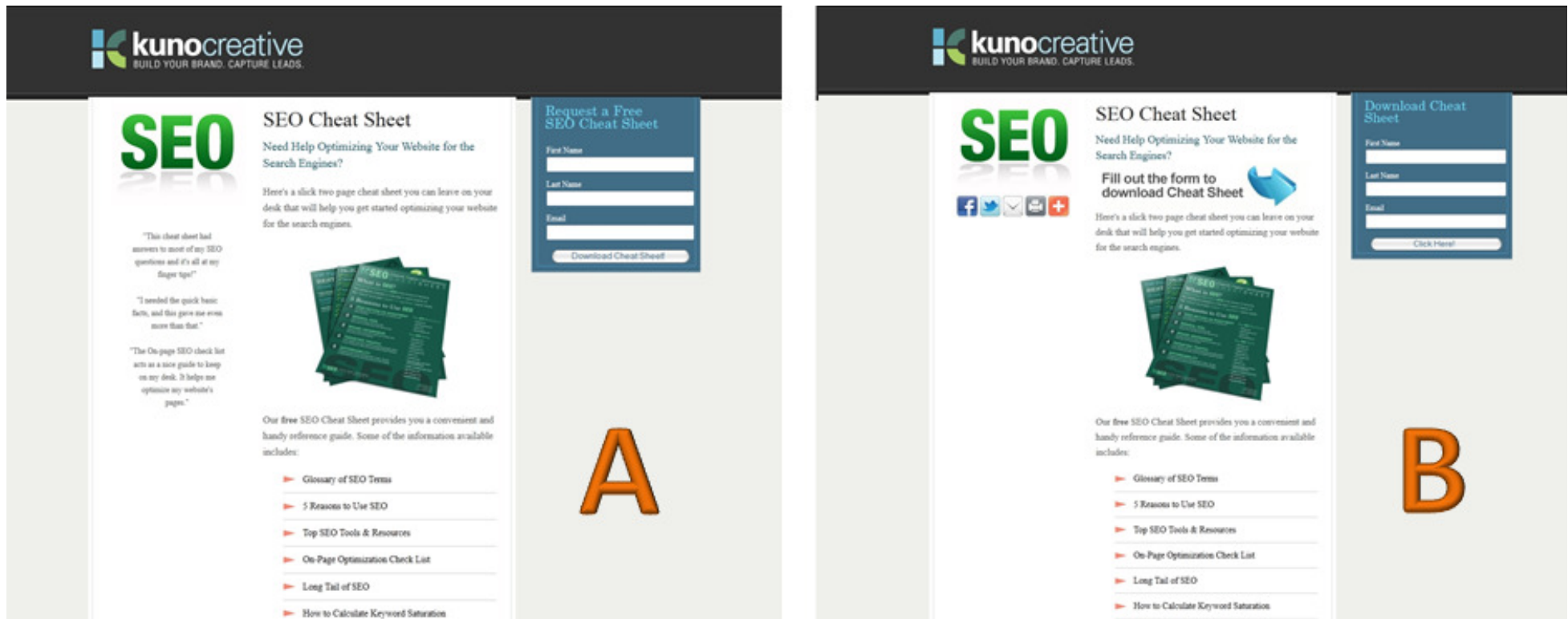
# A/B and Multivariate Testing

# A/B and Multivariate Testing

- Netflix A/B tests for PlayStation 3



Figure taken from http://gigaom.com/video/netflix-ui-innovation/

# Implicit Feedback / Evaluation

- Implicit metrics which can be collected without direct interaction with the user:
  - Time spent on page
  - Links clicked
  - Rank of item clicked
  - Maximum rank of clicked item
  - Link traversal
  - Items on wishlist
  - Items within basket
  - etc.

- Very useful for A/B testing