# Introduction to
# **Information Retrieval**

## **Evaluation IR systems**

# Measures for a search engine

- How fast does it index
  - Number of documents/hour
  - (Average document size)
- How fast does it search
  - Latency as a function of index size

Efficiency

- Quality of results
  - Precision
  - Recall
  - F-measure

Effectiveness

- Expressiveness of query language
  - Ability to express complex information needs

Usability

# Evaluating an IR system

- Note: the **information need** is translated into a **query**
- Relevance is assessed relative to the **information need** *not* the **query**
- E.g., <u>Information need</u>: *I'm looking for information on whether travelling by train from Cairo to Assuit is more effective than flying.*
- <u>Query</u>: ***travelling by train from Cairo to Assuit effective***
- Evaluate whether the doc addresses the information need, not whether it has these words

# Standard relevance benchmarks

- TREC - National Institute of Standards and Technology (NIST) has run a large IR test bed for many years
- Reuters and other benchmark doc collections used
- "Retrieval tasks" specified
  - sometimes as queries
- Human experts mark, for each query and for each doc, <u>Relevant</u> or <u>Nonrelevant</u>
  - or at least for subset of docs that some system returned for that query

# Unranked retrieval evaluation: Precision and Recall

- **Precision**: fraction of retrieved docs that are relevant = (relevant retrieved / retrieved)
- **Recall**: fraction of relevant docs that are retrieved = (relevant retrieved/relevant)

|  | Relevant | Nonrelevant |
|---|---|---|
| Retrieved | tp | fp |
| Not Retrieved | fn | tn |

- Precision P = tp/(tp + fp)
- Recall     R = tp/(tp + fn)

# Should we instead use the accuracy measure for evaluation?

- Given a query, an engine classifies each doc as "Relevant" or "Nonrelevant"
- The **accuracy** of an engine: the fraction of these classifications that are correct
  - (tp + tn) / ( tp + fp + fn + tn)
- **Accuracy** is a commonly used evaluation measure in machine learning classification work

# Precision/Recall

- You can get **high recall** (but low precision) by **retrieving all docs for all queries**!
- Recall is a non-decreasing function of the number of docs retrieved

- In a good system, precision decreases as either the number of docs retrieved or recall increases
  - This is not a theorem, but a result with strong empirical confirmation

# A combined measure: *F*

- Combined measure that assesses precision/recall tradeoff is **F measure** (weighted harmonic mean):

$$F = \frac{1}{a\,\dfrac{1}{P} + (1 - a)\,\dfrac{1}{R}}$$

- People usually use balanced $F_1$ measure
    - i.e., with ✔ = ½

$$F_1 = \frac{2\,PR}{P + R}$$

# Evaluating ranked results

- Evaluation of ranked results:
  - The system can return any number of results
  - By taking various numbers of the top returned documents (levels of recall), the evaluator can produce a *precision-recall curve*

# Averaging over queries

- A precision-recall graph for one query isn't a very sensible thing to look at
- You need to average performance over a whole bunch of queries.
- But there's a technical issue:
  - Precision-recall calculations place some points on the graph

# Typical (good) 11 point precisions