

# Introduction to **Information Retrieval**

CS276

Information Retrieval and Web Search

Pandu Nayak and Prabhakar Raghavan

Lecture 8: Evaluation

# This lecture

---

- How do we know if our results are any good?
  - Evaluating a search engine
    - Benchmarks
    - Precision and recall
- Results summaries:
  - Making our good results usable to a user

# **EVALUATING SEARCH ENGINES**

# 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
- **Expressiveness of query language**
  - Ability to express complex information needs
  - Speed on complex queries
- **Uncluttered UI**(Minimizing Complexity In User Interfaces)
- **Is it free?**

# Measures for a search engine

---

- All of the preceding criteria are *measurable*: we can quantify speed/size
  - we can make expressiveness precise
- **The key measure: user happiness**
  - What is this?
  - Speed of response/size of index are factors
  - But blindingly fast, useless answers won't make a user happy
- Need a way of quantifying user happiness

# Measuring user happiness

---

- Issue: **who is the user we are trying to make happy?**
  - Depends on the setting
- Web engine:
  - User finds what s/he wants and returns to the engine
    - Can measure rate of return users
  - User completes task – search as a means, not end
- eCommerce site: user finds what s/he wants and buys
  - Is it the end-user, or the eCommerce site, whose happiness we measure?
  - Measure time to purchase, or fraction of searchers who become buyers?

# Measuring user happiness

---

- Enterprise (company/govt/academic): Care about “user productivity”
  - How much time do my users save when looking for information?
  - Many other criteria .....

# Happiness: elusive to measure

---

- Most common proxy: **relevance of search results**
- But how do you measure relevance?
- We will detail a methodology here, then examine its issues
- **Relevance measurement requires 3 elements:**
  1. A benchmark document collection
  2. A benchmark suite of queries
  3. A usually binary assessment of either Relevant or Nonrelevant for each query and each document
    - Some work on more-than-binary, but not the standard



# 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., Information need: *I'm looking for information on whether drinking green tea is more effective at reducing your risk of heart attacks than black tea .*
- Query: ***tea green black heart attack 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, Relevant or Nonrelevant
  - 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  
=  $P(\text{relevant} | \text{retrieved})$
- **Recall:** fraction of relevant docs that are retrieved  
=  $P(\text{retrieved} | \text{relevant})$

	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
- Why is this not a very useful evaluation measure in IR?

# Why not just use accuracy?

- How to build a 99.9999% accurate search engine on a low budget....

A screenshot of a web browser showing the search engine 'snoogle.com'. The logo is in a colorful, playful font. Below the logo is a search bar with the text 'Search for:' and an empty input field. Below the input field, it says '0 matching results found.' in a blue, italicized font.

snoogle.com

Search for:

*0 matching results found.*

- People doing information retrieval *want to find something* and have a certain tolerance for junk.

# 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

# Difficulties in using precision/recall

---

- Should average over large document collection/query ensembles
- Need human relevance assessments
  - People aren't reliable assessors
- Assessments have to be binary
  - Nuanced assessments?
- Heavily skewed by collection/authorship
  - Results may not translate from one domain to another

# A combined measure: $F$

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

$$F = \frac{1}{\alpha \frac{1}{P} + (1-\alpha) \frac{1}{R}} = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

- People usually use balanced  $F_1$  measure
  - i.e., with  $\beta = 1$  or  $\alpha = \frac{1}{2}$
  - with  $\beta = 1$  or  $\alpha = 1$  =====>  $F = 2RP/(R+P)$



# Example

---

Assume that a collection contains 100 relevant documents, and a retrieval system returns 20 relevant documents and 10 non-relevant documents. Calculate:

- i) Precision
- ii) Recall
- iii)  $F_1$  score

## ■ Given Data:

- Total number of relevant documents = 100
- Number of relevant documents retrieved = 20
- Number of non-relevant documents retrieved = 10
- Total number of documents retrieved = 20+10=30

To calculate the **precision**, **recall**, and **F1 score**, let us define the relevant terms:

1. **Precision (P)**: The fraction of retrieved documents that are relevant.

$$\text{Precision} = \frac{\text{Number of relevant documents retrieved}}{\text{Total number of documents retrieved}}$$

2. **Recall (R)**: The fraction of relevant documents that are retrieved.

$$\text{Recall} = \frac{\text{Number of relevant documents retrieved}}{\text{Total number of relevant documents}}$$

3. **F1 Score**: The harmonic mean of precision and recall.

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

## i) Precision:

$$\text{Precision} = \frac{\text{Number of relevant documents retrieved}}{\text{Total number of documents retrieved}}$$

$$\text{Precision} = \frac{20}{30} = 0.6667 \text{ (or 66.67\%)}$$

## ii) Recall:

$$\text{Recall} = \frac{\text{Number of relevant documents retrieved}}{\text{Total number of relevant documents}}$$

$$\text{Recall} = \frac{20}{100} = 0.2 \text{ (or 20\%)}$$

## iii) F1 Score:

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{F1 Score} = 2 \times \frac{0.6667 \times 0.2}{0.6667 + 0.2} = 2 \times \frac{0.1333}{0.8667} \approx 0.3077 \text{ (or 30.77\%)}$$

Self Study (Slides 21-27)

# **CREATING TEST COLLECTIONS FOR IR EVALUATION**

# Test Collections

**TABLE 4.3 Common Test Corpora**

<i>Collection</i>	<i>NDocs</i>	<i>NQrys</i>	<i>Size (MB)</i>	<i>Term/Doc</i>	<i>Q-D RelAss</i>
ADI	82	35			
AIT	2109	14	2	400	>10,000
CACM	3204	64	2	24.5	
CISI	1460	112	2	46.5	
Cranfield	1400	225	2	53.1	
LISA	5872	35	3		
Medline	1033	30	1		
NPL	11,429	93	3		
OSHMED	34,8566	106	400	250	16,140
Reuters	21,578	672	28	131	
TREC	740,000	200	2000	89-3543	» 100,000

# From document collections to test collections

---

- Still need
  - Test queries
  - Relevance assessments
- Test queries
  - Must be germane to docs available
  - Best designed by domain experts
  - Random query terms generally not a good idea
- Relevance assessments
  - Human judges, time-consuming
  - Are human panels perfect?

# Kappa measure for inter-judge (dis)agreement (Self Study)

---

- Kappa measure
  - Agreement measure among judges
  - Designed for categorical judgments
  - Corrects for chance agreement
- $\text{Kappa} = [ P(A) - P(E) ] / [ 1 - P(E) ]$
- $P(A)$  – proportion of time judges agree
- $P(E)$  – what agreement would be by chance
- Kappa = 0 for chance agreement, 1 for total agreement.
- $P(E) = P(\text{nonrelevant})^2 + P(\text{relevant})^2$

$P(A)? P(E)?$ 

# Kappa Measure: Example (Self Study)

Number of docs	Judge 1	Judge 2
300	Relevant	Relevant
70	Nonrelevant	Nonrelevant
20	Relevant	Nonrelevant
10	Nonrelevant	Relevant

$$P(E) = P(\text{nonrelevant})^2 + P(\text{relevant})^2$$

$$\text{Kappa} = [P(A) - P(E)] / [1 - P(E)]$$

judges agree

- $P(A) = 370/400 = 0.925$
- $P(\text{nonrelevant}) = (10+20+70+70)/800 = 0.2125$
- $P(\text{relevant}) = (10+20+300+300)/800 = 0.7878$
- $P(E) = 0.2125^2 + 0.7878^2 = 0.665$
- $\text{Kappa} = (0.925 - 0.665)/(1 - 0.665) = 0.776$



# Kappa Example (Self Study)

---

- $\text{Kappa} > 0.8$  = good agreement
- $0.67 < \text{Kappa} < 0.8 \rightarrow$  “tentative conclusions”
- Depends on purpose of study
- For  $>2$  judges: average pairwise kappas

# TREC

---

- TREC Ad Hoc task from first 8 TRECs is standard IR task
  - 50 detailed information needs a year
  - Human evaluation of pooled results returned
  - More recently other related things: Web track, HARD
- A TREC query (TREC 5)
  - <top>
  - <num> Number: 225
  - <desc> Description:  
What is the main function of the Federal Emergency Management Agency (FEMA) and the funding level provided to meet emergencies?  
Also, what resources are available to FEMA such as people, equipment, facilities?
  - </top>

# Standard relevance benchmarks:

## Others

---

- GOV2
  - Another TREC/NIST collection
  - 25 million web pages
  - Largest collection that is easily available
  - But still 3 orders of magnitude smaller than what Google/Yahoo/MSN index
- NTCIR
  - East Asian language and cross-language information retrieval
- Cross Language Evaluation Forum (CLEF)
  - This evaluation series has concentrated on European languages and cross-language information retrieval.
- Many others

# Critique of pure relevance

---

- **Relevance vs Marginal Relevance**
  - A document can be redundant even if it is highly relevant
  - Duplicates
  - The same information from different sources
  - **Marginal relevance is a better measure of utility for the user.**
- But harder to create evaluation set

# Can we avoid human judgment?

---

- **No**
- Makes experimental work hard
  - Especially on a large scale
- **In some very specific settings, you can use proxies**
  - **E.g.: for approximate vector space retrieval, we can compare the cosine distance closeness of the closest docs to those found by an approximate retrieval algorithm**
- But once we have test collections, we can reuse them (so long as we don't overtrain too badly)

**Marginal Relevance**

# Evaluation at large search engines

---

- Search engines have test collections of queries and hand-ranked results
- Recall is difficult to measure on the web
- Search engines often use **precision** at top  $k$ , e.g.,  $k = 10$

**Search engines also use non-relevance-based measures.**

- Clickthrough on first result
  - Not very reliable if you look at a single clickthrough (you may realize after clicking that the summary was misleading and the document is nonrelevant)
- Studies of user behavior in the lab
- A/B testing

**Click-through** rate is the ratio of users who click on a specific link to the number of total users who view a page, email, or advertisement.

# A/B testing

---

- Purpose: Test a single innovation
- Prerequisite: You have a large search engine up and running.
- Have most users use old system
- Divert a small proportion of traffic (e.g., 1%) to the new system that includes the innovation
- Evaluate with an “automatic” measure like clickthrough on first result
- Now we can directly see if the innovation does improve user happiness.
- Probably the evaluation methodology that large search engines trust most
- In principle less powerful than doing a multivariate regression analysis, but easier to understand

# RESULTS PRESENTATION



# Result Summaries

---

- Having ranked the documents matching a query, we wish to present a results list
- Most commonly, a list of the document titles plus a short summary, aka “10 blue links”

## [John McCain](#)

**John McCain** 2008 - The Official Website of **John McCain's** 2008 Campaign for President ... African American Coalition; Americans of Faith; American Indians for **McCain**; Americans with ...  
[www.johnmccain.com](http://www.johnmccain.com) · [Cached page](#)

## [JohnMcCain.com - McCain-Palin 2008](#)

**John McCain** 2008 - The Official Website of **John McCain's** 2008 Campaign for President ... African American Coalition; Americans of Faith; American Indians for **McCain**; Americans with ...  
[www.johnmccain.com/Informing/Issues](http://www.johnmccain.com/Informing/Issues) · [Cached page](#)

## [John McCain News- msnbc.com](#)

Complete political coverage of **John McCain**. ... Republican leaders said Saturday that they were worried that Sen. **John McCain** was heading for defeat unless he brought stability to ...  
[www.msnbc.msn.com/id/16438320](http://www.msnbc.msn.com/id/16438320) · [Cached page](#)

## [John McCain | Facebook](#)

Welcome to the official Facebook Page of **John McCain**. Get exclusive content and interact with **John McCain** right from Facebook. Join Facebook to create your own Page or to start ...  
[www.facebook.com/johnmccain](http://www.facebook.com/johnmccain) · [Cached page](#)

# Summaries

---

- The title is often automatically extracted from document metadata. What about the summaries?
  - This description is crucial.
  - User can identify good/relevant hits based on the description.
- Two basic kinds:
  - Static
  - Dynamic
- A **static summary** of a document is always the same, regardless of the query that hit the doc
- A **dynamic summary** is a *query-dependent* attempt to explain why the document was retrieved for the query at hand

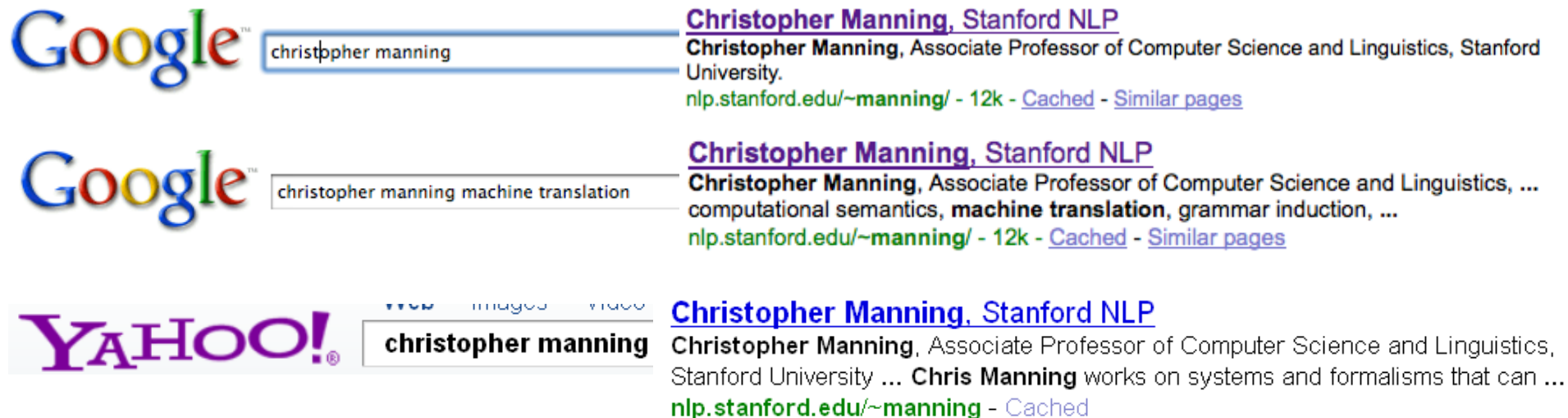
# Static summaries

---

- In typical systems, the static summary is a subset of the document
- Simplest heuristic: the first 50 (or so – this can be varied) words of the document
  - Summary cached at indexing time
- More sophisticated: extract from each document a set of “key” sentences
  - Simple NLP heuristics to score each sentence
  - Summary is made up of top-scoring sentences.
- Most sophisticated: NLP used to synthesize a summary
  - Seldom used in IR; cf. text summarization work

# Dynamic summaries

- Present one or more “windows” within the document that contain several of the query terms
  - “KWIC” snippets: Keyword in Context presentation



The image displays three search engine results for the query 'christopher manning' and 'christopher manning machine translation'. The first two results are from Google, and the third is from Yahoo!. Each result shows the search engine logo, the search query in a text box, and the search results. The Google results show a snippet of text from a document, while the Yahoo! result shows a snippet of text from a document. The search results are as follows:

**Google**  [Christopher Manning, Stanford NLP](#)  
**Christopher Manning**, Associate Professor of Computer Science and Linguistics, Stanford University.  
[nlp.stanford.edu/~manning/](http://nlp.stanford.edu/~manning/) - 12k - [Cached](#) - [Similar pages](#)

**Google**  [Christopher Manning, Stanford NLP](#)  
**Christopher Manning**, Associate Professor of Computer Science and Linguistics, ...  
computational semantics, **machine translation**, grammar induction, ...  
[nlp.stanford.edu/~manning/](http://nlp.stanford.edu/~manning/) - 12k - [Cached](#) - [Similar pages](#)

**YAHOO!**  [Christopher Manning, Stanford NLP](#)  
**Christopher Manning**, Associate Professor of Computer Science and Linguistics, Stanford University ... **Chris Manning** works on systems and formalisms that can ...  
[nlp.stanford.edu/~manning/](http://nlp.stanford.edu/~manning/) - [Cached](#)

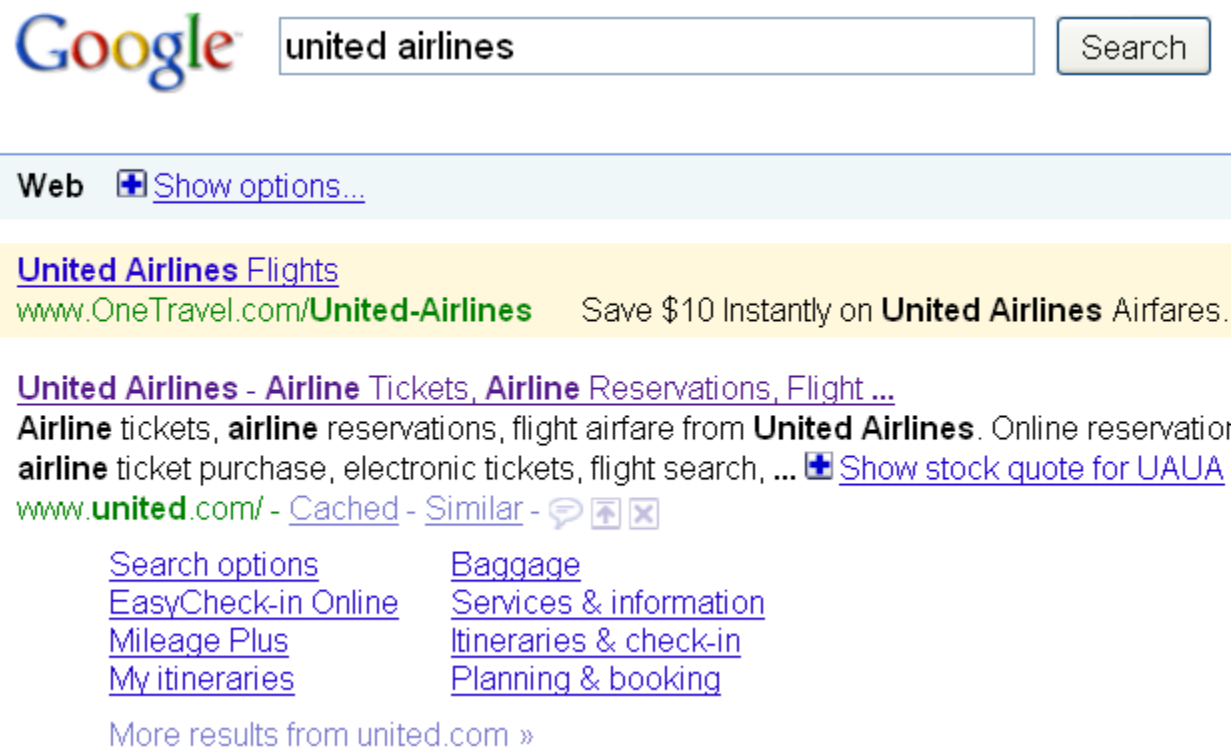
# Techniques for dynamic summaries

---

- Find small windows in doc that contain query terms
  - Requires fast window lookup in a document cache
- Score each window wrt query
  - Use various features such as window width, position in document, etc.
  - Combine features through a scoring function
  - Challenges in evaluation: judging summaries
  - Easier to do pairwise comparisons rather than binary relevance assessments

# Quicklinks

- For a *navigational query* such as ***united airlines*** user's need likely satisfied on [www.united.com](http://www.united.com)
- Quicklinks provide navigational cues on that home page



The screenshot shows a Google search interface with the query 'united airlines' entered in the search bar. Below the search bar, there is a 'Web' tab and a link to 'Show options...'. The search results are displayed on a yellow background. The first result is 'United Airlines Flights' with the URL 'www.OneTravel.com/United-Airlines' and a snippet 'Save \$10 Instantly on United Airlines Airfares.'. Below this, there is a link to 'United Airlines - Airline Tickets, Airline Reservations, Flight ...' with a snippet 'Airline tickets, airline reservations, flight airfare from United Airlines. Online reservation airline ticket purchase, electronic tickets, flight search, ...' and a link to 'Show stock quote for UAUA'. At the bottom, there are several quicklinks: 'Search options', 'EasyCheck-in Online', 'Mileage Plus', 'My itineraries', 'Baggage', 'Services & information', 'Itineraries & check-in', 'Planning & booking', and a link to 'More results from united.com »'.

Web [Show options...](#)

**United Airlines Flights**  
[www.OneTravel.com/United-Airlines](http://www.OneTravel.com/United-Airlines) Save \$10 Instantly on **United Airlines** Airfares.

**United Airlines - Airline Tickets, Airline Reservations, Flight ...**  
 Airline tickets, airline reservations, flight airfare from **United Airlines**. Online reservation  
 airline ticket purchase, electronic tickets, flight search, ... [Show stock quote for UAUA](#)  
[www.united.com/](http://www.united.com/) - [Cached](#) - [Similar](#) - [Show stock quote for UAUA](#)

[Search options](#) [Baggage](#)  
[EasyCheck-in Online](#) [Services & information](#)  
[Mileage Plus](#) [Itineraries & check-in](#)  
[My itineraries](#) [Planning & booking](#)

[More results from united.com »](#)

united airlines

Search Pad

SearchScan - On

102,000,000 results for united airlines:

Show All

United Air Lines

Wikipedia

Also try: [united airlines reservations](#), [united airlines flight](#), [More...](#)

[United Airlines - Airline Tickets, Airline Reservations ...](#) (Nasdaq: [UAUA](#))  
Official site for **United Airlines**, commercial air carrier transporting people, property, and mail across the U.S. and worldwide.  
[www.united.com](#) - 65k - [Cached](#)

[Planning & Booking](#)  
[Itineraries & Check-in](#)  
[Mileage Plus](#)  
[Services & Information](#)

[Shop for Flights](#)  
[Special Deals](#)  
[Flight Status](#)  
[Customer Service](#)

[more results from united.com »](#)

united airlines

UNITED AIRLINES

[United Airline Fleet](#)

[United Airline Schedule](#)

[United Airlines Reservations](#)

[United Airline Jobs](#)

[Reference](#)

ALL RESULTS

[Cheap Flight Tickets](#) · [www.CheapOair.com](#)  
CheapOair - The Only Way to Go!! Find Over 18 Million Exclusive Fares.

[Fly United Airlines](#) · [www.OneTravel.com/United-Airline](#)  
Save \$10 Instantly on **United Airlines** Flights. Book Now, Hurry!

Best match

[United Airlines - Airline Tickets, Airline Reservations, Flight ...](#)  
[www.united.com](#) · Official site  
**Airline** tickets, **airline** reservations, flight airfare from **United Airlines**. Online reservations, **airline** ticket purchase, electronic tickets, flight search, fares and availability ...

[Flights](#)  
[Check In Online](#)  
[My itineraries](#)  
[Baggage](#)

[Redeem miles](#)  
[Children, pets, & assistance](#)  
[Change your travel plans](#)  
[Special deals](#)

Customer service 800-864-8331

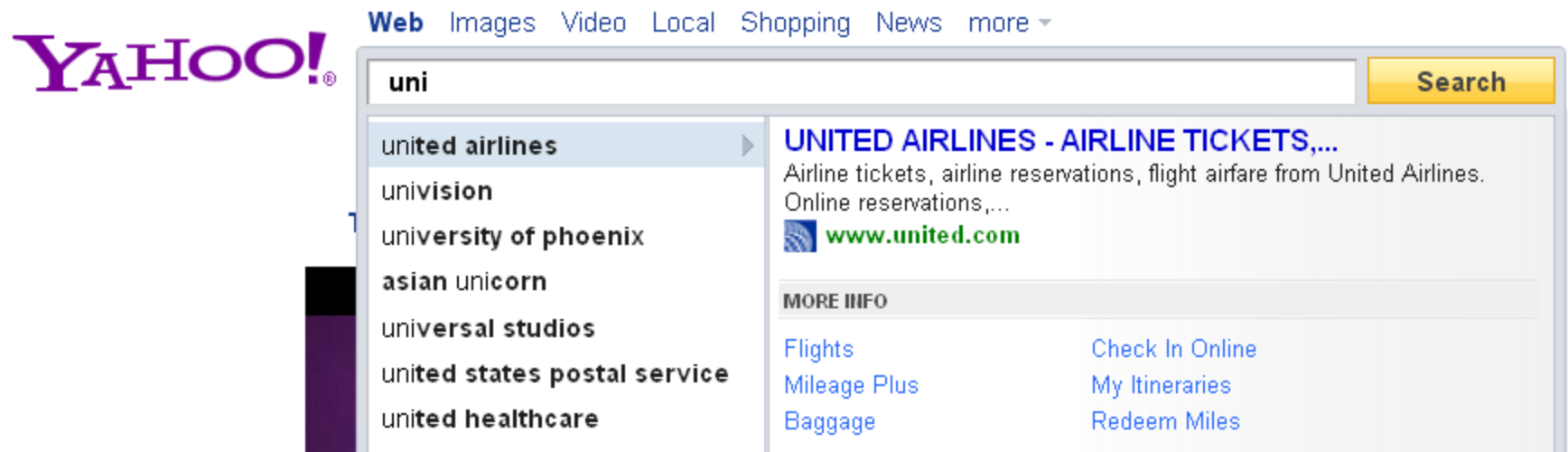
RELATED SEARCHES

[United Airlines Flight Status](#)

[US Airways](#)

[Continental Airlines](#)

# Alternative results presentations?





# Resources for this lecture

---

- IIR 8
- MIR Chapter 3
- MG 4.5
- Carbonell and Goldstein 1998. The use of MMR, diversity-based reranking for reordering documents and producing summaries. SIGIR 21.