**PURDUE** UNIVERSITY. | Department of Computer Science

# CS47300: Web Information Search and Management

*Web Search*
Prof. Chris Clifton
16 September 2020
*Some slides courtesy*
*Manning, Raghavan, and Schütze*

Indiana
Center for
Database
Systems

---

**PURDUE** UNIVERSITY.
Department of Computer Science

# Challenge:
## "The need behind the query"

- Semantic analysis
  - Query language determination
    - Auto filtering
    - Different ranking (if query in Japanese do not return English)
  - Hard & soft (partial) matches
    - Personalities (triggered on names)
    - Cities (travel info, maps)
    - Medical info (triggered on names and/or results)
    - Stock quotes, news (triggered on stock symbol)
    - Company info
    - Etc.
  - Natural Language reformulation
  - Integration of Search and Text Analysis

# Answering "the need behind the query": Context

- Context determination
  - spatial (user location/target location)
  - query stream (previous queries)
  - personal (user profile)
  - explicit (user choice of a vertical search, )
  - implicit (use Google from France, use google.fr)

- Context use
  - Result restriction
    - Kill inappropriate results
  - Ranking modulation
    - Use a "rough" generic ranking, but personalize later

# The spatial context: Geo-search

- Two aspects
  - Geo-coding -- encode geographic coordinates to make search effective
  - Geo-parsing -- the process of identifying geographic context.
- Geo-coding
  - Geometrical hierarchy (squares)
  - Natural hierarchy (country, state, county, city, zip-codes, etc)
- Geo-parsing
  - Pages (infer from phone nos, zip, etc). About 10% can be parsed.
  - Queries (use dictionary of place names)
  - Users
    - Explicit (tell me your location -- used by NL, registration, from ISP)
    - From IP data
  - Mobile phones
    - Many sources of highly accurate location

# Answering "the need behind the query": Context

- Context determination
  - spatial (user location/target location)
  - query stream (previous queries)
  - personal (user profile)
  - explicit (user choice of a vertical search, )
  - implicit (use Google from France, use google.fr)

- Context use
  - Result restriction
    - Kill inappropriate results
  - Ranking modulation
    - Use a "rough" generic ranking, but personalize later

# Context transfer

No transfer


Context transfer

# Transfer from search results

Web | Images | Video | Directory | Local | News

**YAHOO! SEARCH** naive bayes performance

**My Web** BETA    My Search History OFF | On                Subscriptions (New)    Shortc

**Search Results**                        Results **1 - 10** of about **75,200** for **naive bayes perfor**

1. An analysis of data characteristics that affect **naive Bayes performance** (PDF)
   ... An analysis of data characteristics that affect **naive Bayes performance** ... ally a better indicator of **naive Bayes performance** than the ...
   www.research.ibm.com/PM/icml01.pdf - 357k - View as html - More from this site - Save - Block

2. Improving the **Performance** of **Naive Bayes** for Text Classification (PDF)
   Improving the **Performance** of **Naive Bayes** for. Text Classification. Yirong Shen and Jing Jiang. CS224N Spring 2003. Abstract. We seek to improve the **performance** of the **naive Bayes** classifier
   www-nlp.stanford.edu/courses/cs224n/2003/fp/yirong99/report.pdf - 135k - View as html - More from this site - Save - Block

---

icml01.pdf (application/pdf Obje...

42%    Search PDF    Hide

Finished searching for:
**naive bayes performance**

Total instances found:
**126**

New Search

An analysis of data characteristics that affect **naive Bayes performance**

Irina Rish                                    RISH@US.IBM.COM
Joseph Hellerstein                            HELLERS@US.IBM.COM
Jayram Thathachar                             JAYRAM@US.IBM.COM
IBM T.J. Watson Research Center 30 Saw Mill River Road, Hawthorne, NY 10532

**Abstract**

Results:

affect **naive Bayes performance** I
the **naive Bayes** classifier is remark

Done

Save and View this PDF in Reader

**Find a word in the current PDF document**

1 of 8

## Answering "the need behind the query": Context

- Context determination
  - spatial (user location/target location)
  - query stream (previous queries)
  - personal (user profile)
  - explicit (user choice of a vertical search, e.g., Amazon or eBay)
  - implicit (use Google from France, use google.fr)

- Context use
  - Result restriction
    - Kill inappropriate results
  - Ranking modulation
    - Use a "rough" generic ranking, but personalize later

## Result Restriction

- Geographic restrictions
  - Holocaust denial in Germany
  - Imagery that may be illegal in some jurisdictions, accepted in others
- Age restrictions
  - COPPA



"On the Internet, nobody knows you're a dog."

51

# Web Crawler

- Finds and downloads web pages automatically
  - provides the collection for searching
- Web is huge and constantly growing
- Web is not under the control of search engine providers
- Web pages are constantly changing
- Crawlers also used for other types of data

# Retrieving Web Pages

- Web crawler client program connects to a *domain name system* (DNS) server
- DNS server translates the hostname into an *internet protocol* (IP) address
- Crawler then attempts to connect to server host using specific *port*
- After connection, crawler sends an HTTP request to the web server to request a page
  - usually a GET request

# Crawling the Web



# Web Crawler

- Starts with a set of *seeds*, which are a set of URLs given to it as parameters
- Seeds are added to a URL request queue
- Crawler starts fetching pages from the request queue
- Downloaded pages are parsed to find link tags that might contain other useful URLs to fetch
- New URLs added to the crawler's request queue, or *frontier*
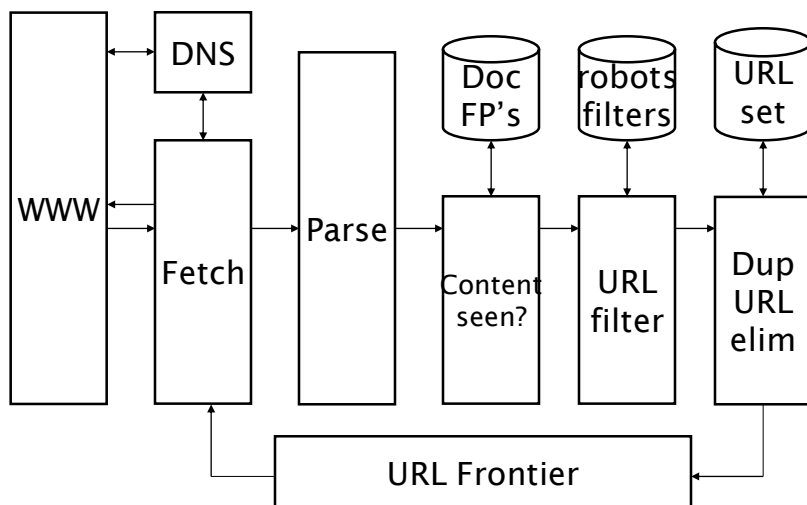
## Processing steps in crawling

- Pick a URL from the frontier  ← Which one?
- Fetch the document at the URL
- Parse the URL
  - Extract links from it to other docs (URLs)
- Check if URL has content already seen
  - If not, add to indexes

  E.g., only crawl .edu, obey robots.txt, etc.
- For each extracted URL
  - Ensure it passes certain URL filter tests
  - Check if it is already in the frontier (duplicate URL elimination)

75

## Basic crawl architecture

76

## What any crawler *must* do

- Be <u>Robust</u>: Be immune to spider traps and other malicious behavior from web servers

- Be <u>Polite</u>: Respect implicit and explicit politeness considerations

## What any crawler *should* do

- Be capable of <u>distributed</u> operation: designed to run on multiple distributed machines
- Be <u>scalable</u>: designed to increase the crawl rate by adding more machines
- <u>Performance/efficiency</u>: permit full use of available processing and network resources