U F *m* G  UNIVERSIDADE FEDERAL
DE MINAS GERAIS

Ranking Models

# Introduction

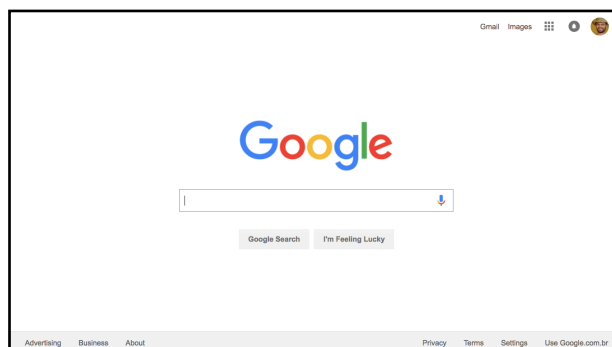Rodrygo L. T. Santos
rodrygo@dcc.ufmg.br
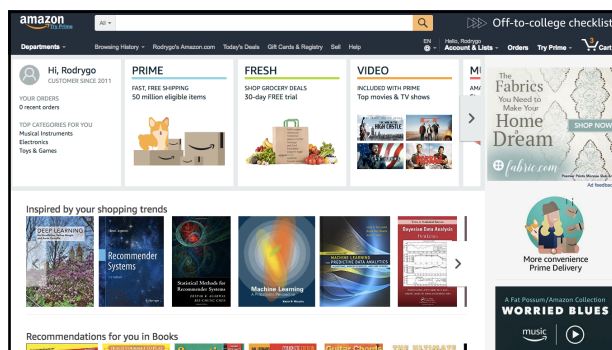
---

## Information retrieval

❝ *Information retrieval is a field concerned with the structure, analysis, organization, storage, searching, and retrieval of information.*
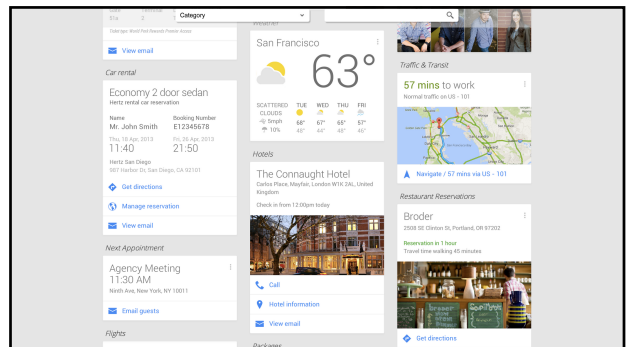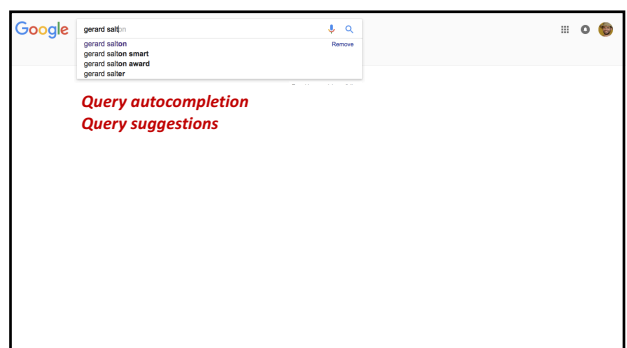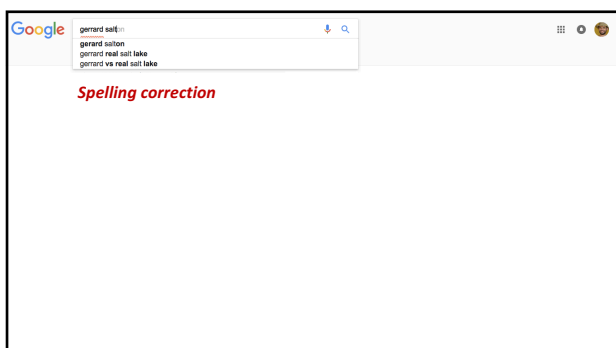
Gerard Salton, 1968

---

## Retrieval tasks

search

**1**
query

---

---

## Retrieval tasks

search        recommendation

**1**          **0**
query          query

---

**Retrieval tasks**



search → recommendation → anticipation

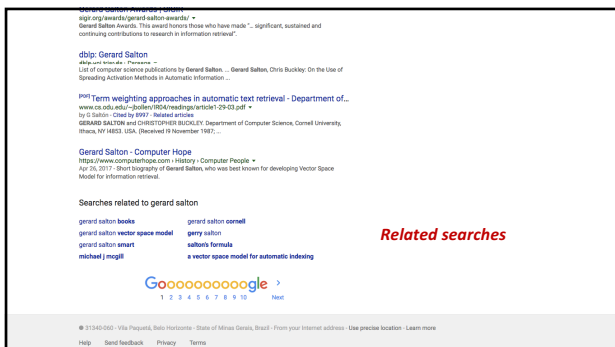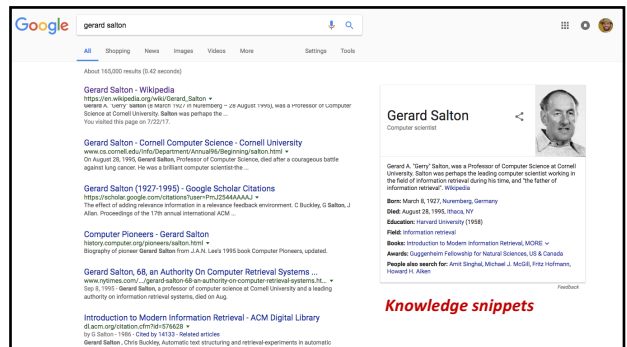| 1 query | 0 query | -1 query |



---

**Our focus: search**

A lot of people *………………………… 60k queries per second*
From a lot of places *………. whole planet (and beyond?)*
Using a lot of a devices *…………………. smart-you-name-it*
Looking for a lot of info *……………………… $10^{11}$ documents*

---

**What does a search engine do?**

---



*Spelling correction*

---



*Query autocompletion*
*Query suggestions*

**Vertical search results**

**Knowledge snippets**

**Related searches**

ten
blue
links

**The ranking problem**

Given
∘ Some evidence of the user's need
Produce
∘ A list of matching information items
∘ In decreasing order of relevance

**The ranking problem**

Given
∘ Some ~~evidence of the user's need~~ *query*
Produce
∘ A list of matching ~~information items~~ *documents*
∘ In decreasing order of relevance

**1) What documents do we show?**



**2) What order do we show them in?**



**2) What order do we show them in?**



$$f(q, d)$$

**Isn't it a solved problem?**



GOOGLE DOESN'T HAVE ALL THE ANSWERS

Credit: Anna Jumped

**Efficiency**

Efficiency is about doing something (good or bad) in an optimal way (i.e., faster or with fewer resources)

Key performance indicators

◦ *Query latency:* searching billions of documents
◦ *Query throughput:* serving thousands of users
◦ *Document latency:* serving freshly published content

**Effectiveness**

Effectiveness is about doing the right thing; it's about finding documents that are relevant to the user

Relevance is influenced by many factors
◦ Topical relevance vs. user relevance
◦ Task, context, novelty, style

Ranking models define *a view of* relevance

**Pursuing relevance**

Exact matching of words is not enough
◦ Queries are ambiguous; documents are ambiguous
◦ Many different ways to write the same thing
Even perfect matching is not enough
◦ Must counter adversarial content
◦ Must infer quality beyond content

**Assessing relevance**

Relevance is a user's prerogative
◦ We can observe changes in user behavior
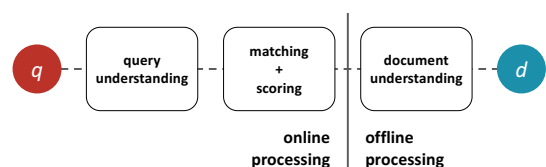◦ Or directly ask the user how we're doing
Evaluation is an empirical science
◦ It must be scientifically rigorous
◦ It must be economically viable

# What do search engineers do?

**The ranking problem**

q - - - - - - - d

$$f(q, d)$$

**Ranking pipeline**

q → query understanding → matching + scoring → document understanding → d

online processing | offline processing

## (Continuous) offline processing

Document acquisition

Document understanding

Document indexing

## Document acquisition

The Web is huge
◦ Trillions of known URLs, billions fetched

The Web is constantly evolving
◦ Updates, additions, deletions

Efficient crawling is key
◦ Must aim for coverage, but also freshness

## Document understanding

Documents come in many flavors
◦ Web pages, email, books, news stories, scholarly papers, text messages, Word™, Powerpoint™, PDF, forum postings, patents, IM sessions, etc.

## Document understanding

Documents have mostly textual content
◦ Some of it static, some dynamically rendered

And often some structure
◦ Title, body, url, anchor text for web pages
◦ Subject, sender, destination for email

## Document understanding

Documents carry meaning
◦ Term-based matching as a first approximation
◦ Several techniques to leverage semantics

Documents vary in quality
◦ *Genuinely:* accessibility, readability, authority, depth
◦ *Maliciously:* content farms, link farms

## Document indexing

Efficient retrieval through indexes
◦ Like the index of a book
  • For each word, a list of documents it appears on
◦ Broken up into shards of millions of documents
  • 1000s of shards for the web index
◦ Plus per-document metadata

### Online processing

Query understanding

Matching and scoring

Post-processing

### Query understanding

Keyword queries are often poor descriptions of the user's actual information need
∘ Interaction and context also matter

Query understanding techniques can help
∘ Query segmentation, query scoping
∘ Query relaxation, query expansion

### Query understanding

Query scoping through semantic annotation
∘ [**san jose** convention center]
∘ [**matt cutts**]

Query expansion through acronym expansion
∘ [**gm** trucks] → [**general motors** trucks]
∘ [**gm** corn] → [**genetically modified** corn]

### Matching and scoring

Send the query to all the shards

Each shard
∘ Finds matching documents
∘ Scores each query-document pair
∘ Sends back the top $n$ documents

Combine all the top documents and sort by score

### Post-retrieval adjustments

Host clustering, sitelinks

Near-duplicate removal, diversification

Spam demotions, copyright takedowns

### Course goals

Provide an in-depth account of ranking models and evaluation methods for information retrieval

Provide an exploration of recent advances and current research directions in the field

## Course scope

Focus on ranking
◦ Query-dependent ranking
◦ Query-independent ranking
◦ Machine-learned ranking
Focus on evaluation
◦ Offline, online, counterfactual

## Out-of-scope

We have dedicated courses for:
◦ Information retrieval
◦ Recommender systems
◦ Machine learning
◦ Data mining

## Course materials: textbooks

Search Engines: Information Retrieval in Practice
by B. Croft, D. Metzler, and T. Strohman

Introduction to Information Retrieval
by C. Manning, P. Raghavan, and H. Schütze

Modern Information Retrieval
by R. Baeza-Yates and B. Ribeiro-Neto

## Course materials: textbooks

Information Retrieval: Implementing and Evaluating
Search Engines
by S. Büttcher, C. Clarke, and G. Cormack

Text Data Management: A Practical Introduction to
Information Retrieval and Text Mining
by C. Zhai and S. Massung

## Course materials: surveys

Foundations and Trends in Information Retrieval
by several authors

Morgan & Claypool Synthesis Lectures on Information
Concepts, Retrieval, and Services
by several authors

## Other relevant material

General background
◦ Algorithms and data structures
◦ Basic statistics
◦ Basic linear algebra
Advanced readings
◦ Google Scholar is your friend

**Course grading (tentative)**

Exams: 40%

Project: 40%

Assignments: 20%

**Course website**

http://homepages.dcc.ufmg.br/~rodrygo

**References**

Google Search Statistics
Internet Live Stats, 2017

How Google Works: A Google Ranking Engineer's Story
Haahr, SMX West 2016

Ten blue links on Mars
Clarke et al., WWW 2017

**Writing assignment #0**

Fill in a short questionnaire describing your past experience and expectations related to the course
○ **Due Tue, Aug 15 @ 23:55 via Moodle**

U F *m* G   UNIVERSIDADE FEDERAL
DE MINAS GERAIS

Coming next…

# Search Architecture

Rodrygo L. T. Santos
rodrygo@dcc.ufmg.br