


# Information Retrieval

1

## Vannevar Bush

- Director of the Office of Scientific Research and Development (1941-1947)



Vannevar Bush, 1890-1974

- End of WW2 - what next big challenge for scientists?

## Historic Vision

“A memex is a device in which an individual stores all his books, records, and communications, and which is mechanized so that it may be consulted with exceeding speed and flexibility. It is an enlarged intimate supplement to his memory.”

Vannevar Bush, *As we may think*, *Atlantic Monthly*, July 1945.

## Prophetic: Hypertext

- "associative indexing, the basic idea of which is a provision whereby any item may be caused at will to select immediately and automatically another. This is the essential feature of the memex. The process of tying two items together is the important thing."

Vannevar Bush, *As we may think*, *Atlantic Monthly*, July 1945

## Prophetic: Wikipedia et al

- "Wholly new forms of encyclopedias will appear, ready made with a **mesh of associative trails** running through them, ready to be dropped into the memex and there amplified."

Vannevar Bush, *As we may think*, *Atlantic Monthly*, July 1945

## Vision

“ This is a much larger matter than merely the extraction of data for the purposes of scientific research; it involves the entire process by which man profits by his inheritance of acquired knowledge”

Vannevar Bush, *As we may think*, *Atlantic Monthly*, July 1945

## Historic Goals

“Google's mission is to organize the world's information and make it universally accessible and useful” [Larry Page, Sergey Brin, Google's mission statement, ~ 1998.](#)

“A memex is a device in which an individual stores all his books, records, and communications, and which is mechanized so that it may be consulted with exceeding speed and flexibility. It is an enlarged intimate supplement to his memory.” [Vannevar Bush, As we may think, Atlantic Monthly, July 1945.](#)

## Information Retrieval Problem

- User wants information from a collection of “objects”: **information need**
  - User formulates need as a “query”
    - Language of information retrieval system
  - System finds objects that “satisfy” query
  - System presents objects to user in “**useful form**”
  - User determines which objects from among those presented are **relevant** **CRITICAL NOTION**
- Define each of the words in quotes  
➤ Develop algorithms

8

## Think first about text documents

Although search has changed, **classic techniques** still provide **foundations**

– our starting point

- Early digital searches – digital card catalog:
  - subject classifications, keywords
- “Full text” : words + natural language syntax
  - No “meta-structure”
- Classic study
  - Gerald Salton SMART project 1960’ s

9

## Scaling

- What are attributes changing from 1960’ s to online searches of today?
  
- How do they change problem?

10

## Develop models

Begin with document:

How do we view document contents?

11

## Modeling: “*query*”

How do we want to express a query?

What does it mean?

12

## Classic Info Retrieval continued

foundational techniques for text documents

### Last time:

the information retrieval process

information need => query => results => relevance

brainstormed models for documents and queries

### Today:

Precise models documents and queries

the Boolean model

ranking and the vector model

(extending models, techniques for modern search)

13

## Modeling “document”

- Basic element: **term**  
single word or other character string
- Document is a \_\_\_\_\_ of terms  
Choose one
  - Set
    - Bag
    - Sequence

14

## Modeling: “query”

- Basic query is **one term**
- Multi-term query is (choose one):
  - Set of terms
  - Sequence of terms
    - multiplicity?
    - Other constraints?
  - Boolean combination of terms

15

## Modeling: “satisfying”

- What determines if document satisfies query?
- That depends ....
  - Document model
  - Query model
  - definition of “satisfying” can still vary
- **START SIMPLE**
  - better understanding
  - Use components of simple model later

16

### Present results in “useful form”

- most basic: give **list of results**
- **meaning** of order of list? => **RANKING**
- **Goals of ranking**
  - Order documents that **satisfy a query** by **how well match the query**
  - **Capture relevance** to **user** by algorithmic method of ordering

17

### (pure) Boolean Model of IR

- Document: *set* of terms
- Query: Boolean expression over terms
- Satisfying:
  - Doc. **evaluates** to “true” on single-term query if contains term
  - Evaluate doc. on expression query as you would any Boolean expression
  - **doc satisfies query if evals to true on query**

18

### Boolean Model example

**Doc 1:** “Computers have brought the world to our fingertips. We will try to understand at a basic level the **science** -- old and new -- underlying this new Computational Universe. Our quest takes us on a broad sweep of scientific **knowledge** and related technologies... Ultimately, this study makes us look anew at ourselves -- our genome; language; music; “**knowledge**”; and, above all, the mystery of our intelligence. (cos 116 description)

**Doc 2:** “An introduction to computer science in the context of scientific, engineering, and commercial applications. The goal of the course is to teach basic principles and practical issues, while at the same time preparing students to use computers effectively for applications in computer science ...” (cos 126 description)

**Query:**  
(principles AND knowledge) OR (science AND engineering)

0	1	1	0
---	---	---	---

**Doc 1: FALSE**

19

### Boolean Model example

**Doc 1:** “Computers have brought the world to our fingertips. We will try to understand at a basic level the science -- old and new -- underlying this new Computational Universe. Our quest takes us on a broad sweep of scientific knowledge and related technologies... Ultimately, this study makes us look anew at ourselves -- our genome; language; music; “knowledge”; and, above all, the mystery of our intelligence. (cos 116 description)

**Doc 2:** “An introduction to computer **science** in the context of scientific, **engineering**, and commercial applications. The goal of the course is to teach basic **principles** and practical issues, while at the same time preparing students to use computers effectively for applications in computer **science** ...” (cos 126 description)

**Query:**  
(principles AND knowledge) OR (science AND engineering)

1	0	1	1
---	---	---	---

**Doc 2: TRUE**

20

### Boolean Model example

**Doc 1:** “Computers have brought the world to our fingertips. We will try to understand at a basic level the **science** -- old and new -- underlying this new Computational Universe. Our quest takes us on a broad sweep of scientific **knowledge** and related technologies... Ultimately, this study makes us look anew at ourselves -- our genome; language; music; **knowledge**”; and, above all, the mystery of our intelligence. (cos 116 description)

**Doc 2:** “An introduction to computer science in the context of scientific, engineering, and commercial applications. The goal of the course is to teach basic principles and practical issues, while at the same time preparing students to use computers effectively for applications in computer science ...” (cos 126 description)

**Query:** (principles OR knowledge) AND (science AND NOT(engineering))

**Doc 1:** (0 OR 1 ) AND (1 AND NOT(0 )) **TRUE** 21

### Boolean Model example

**Doc 1:** “Computers have brought the world to our fingertips. We will try to understand at a basic level the science -- old and new -- underlying this new Computational Universe. Our quest takes us on a broad sweep of scientific knowledge and related technologies... Ultimately, this study makes us look anew at ourselves -- our genome; language; music; "knowledge"; and, above all, the mystery of our intelligence. (cos 116 description)

**Doc 2:** “An introduction to computer **science** in the context of scientific, **engineering**, and commercial applications. The goal of the course is to teach basic **principles** and practical issues, while at the same time preparing students to use computers effectively for applications in computer **science** ...” (cos 126 description)

**Query:** (principles OR knowledge) AND (science AND NOT(engineering))

**Doc 2:** (1 OR 0 ) AND (1 AND NOT(1 )) **FALSE** 22

### (pure) Boolean Model of IR: how “present results in useful form”

- most basic: give **list of results**
- **meaning** of order of list? => **RANKING?**
- There is **no sense of ranking** in pure Boolean model
  - need idea in addition to “satisfying documents”: **generalize model**

23

**Doc 1:** “Computers have brought the world to our fingertips. We will try to understand at a basic level the **science** -- old and new -- underlying this new Computational Universe. Our quest takes us on a broad sweep of scientific **knowledge** and related technologies... Ultimately, this study makes us look anew at ourselves -- our genome; language; music; **knowledge**”; and, above all, the mystery of our intelligence. (cos 116 description)

**Doc 2:** “An introduction to computer **science** in the context of scientific, **engineering**, and commercial applications. The goal of the course is to teach basic **principles** and practical issues, while at the same time preparing students to use computers effectively for applications in computer **science** ...” (cos 126 description)

**Query:** (principles OR knowledge) AND (science OR engineering)

Doc 1:	0	1	1	0	<b>TRUE</b>
Doc 2:	1	0	1	1	<b>TRUE</b>
			<b>RANK?</b>		

24

## Restrict Boolean Model

- **AND model:** query is the AND of a set of query terms: term\_1 AND term\_2 AND...
  - just need specify set of terms
  - This model used by **current search engines**
- **OR model:** query is the OR of a set of query terms: term\_1 OR term\_2 OR ...
  - just need specify set of terms
  - This **original model** for IR development
    - why?

25

## Simple Model with Ranking

- Document: **bag** of terms - count occurrences
- Query: **set** of terms
- Satisfying: **OR** model
- Ranking: **numerical score** measuring degree to which document satisfies query
  - some choices:
    - one point for each query term in document
    - one point for **each occurrence** of a query term in document
- Documents returned in **sorted list** by decreasing score

26

## Simple Model: example

**Doc 1:** "Computers have brought the world to our fingertips. We will try to understand at a basic level the **science** -- old and new -- underlying this new Computational Universe. Our quest takes us on a broad sweep of scientific **knowledge** and related technologies... Ultimately, this study makes us look anew at ourselves -- our genome; language; music; "**knowledge**"; and, above all, the mystery of our intelligence. (cos 116 description)

Frequencies:

**science** 1; **knowledge** 2; **principles** 0; **engineering** 0

**Doc 2:** "An introduction to computer **science** in the context of scientific, **engineering**, and commercial applications. The goal of the course is to teach basic **principles** and practical issues, while at the same time preparing students to use computers effectively for applications in computer **science** ..." (cos 126 description)

Frequencies:

**science** 2; **knowledge** 0; **principles** 1; **engineering** 1

27

## Generalize Simple Model: The Vector Model

- Have a **lexicon** (aka *dictionary*) of all terms appearing in the collection of documents
  - $m$  terms in all, number 1, ...,  $m$
- Document: an  $m$ -dimensional **vector**
  - $i^{\text{th}}$  entry of the vector is a real-valued **weight** (importance of ) term  $i$  in the document
- Query: an  $m$ -dimensional **vector**
  - The  $i^{\text{th}}$  entry of the vector is a real-valued **weight** (importance of ) term  $i$  in the query

28

## Vector Model: Satisfying & Ranking

- Satisfying:
  - Each document is scored as to the degree it satisfies query (higher better)
  - there is **no inherent notion of satisfying**
  - typically doc satisfies query if score is > threshold
- Ranking:
  - Documents are returned in **sorted list** decreasing by score:
    - Include only highest  $n$  documents, some  $n$ ?

29

## Where get dictionary of $t$ terms?

- Pre-determined dictionary.
  - How sure get all terms?
- Build lexicon when collect documents
  - What if collection dynamic: add terms?

30

## How compute score

Calculate a **vector function** of the **document vector** and the **query vector**

Choices:

### 1. distance between the vectors:

$$\text{Dist}(\mathbf{d}, \mathbf{q}) = \sqrt{(\sum_{i=1}^t (\mathbf{d}_i - \mathbf{q}_i)^2)}$$

- Is *dissimilarity* measure
- Not normalized: Dist ranges [0, inf.)
- Fix: use  $e^{-\text{Dist}}$  with range (0,1]
- Is it the right sense of difference?

31

## How compute score, continued

### 2. angle between the vectors:

$$\text{Dot product: } \mathbf{d} \cdot \mathbf{q} = \sum_{i=1}^t (\mathbf{d}_i * \mathbf{q}_i)$$

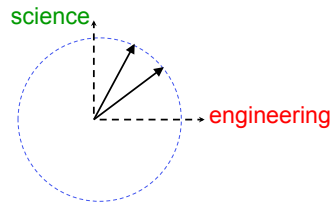
- Is *similarity* measure
- Not normalized: dot product ranges (-inf., inf.)
- Fix: use normalized dot product, range [-1,1]
 
$$(\mathbf{d} \cdot \mathbf{q}) / (|\mathbf{d}| * |\mathbf{q}|) \quad \text{where } |\mathbf{v}| = \sqrt{\sum_{i=1}^t (\mathbf{v}_i^2)}$$
 the length of  $\mathbf{v}$
- In practice vector components are non-negative so range is [0,1]
- This **most commonly used function for score**

32



## Normalizing vectors

- If use unit vectors,  $\mathbf{d} / |\mathbf{d}|$  and  $\mathbf{v} / |\mathbf{v}|$  some of issues go away



33

## The Simple Model as a Vector Model

- **Document:** an  $m$ -dimensional vector
  - $i^{\text{th}}$  entry of the vector is the number of times term  $i$  appears in the document
- **Query:** an  $m$ -dimensional vector
  - The  $i^{\text{th}}$  entry of the vector is 1 if term  $i$  in the query, 0 otherwise
- **Vector function:** dot product

34

How compute weights  $d_i$  and  $q_i$ ?

First:  
observations about this  
model?

35

## Vector model: Observations

- Have matrix of terms by documents
  - ⇒ Can use linear algebra
- Queries and documents are the same
  - ⇒ Can compare documents same way
    - Clustering documents
- Document with only some of query terms can score higher than document with all query terms

36

## How compute weights

- Vector model *could* have weights assigned by **human intervention**
  - may add **meta-information**
  - User setting **query weights** might make sense
    - User decides **importance** of terms in own search
  - Humans setting **document weights?**
    - Who? Billions+ of documents
- Return to model of documents as **bag of words** – calculate weights
  - Function mapping bag of words to vector

37

## Calculations on board:

- General notation:
  - $w_{jd}$  is the weight of term  $j$  in document  $d$
  - $freq_{jd}$  is the # of times term  $j$  appears in doc  $d$
  - $freq_{jC}$  is the # of times term  $j$  appears in the collection  $C$  of documents (collection frequency)
  - $length_d$  is the total number of occurrences of terms in document  $d$  (word length)
  - $n_j$  = # docs containing term  $j$
  - $N$  = number of docs in collection

38

## Summary weight calculation

- General notation:
  - $w_{jd}$  is the weight of term  $j$  in document  $d$
  - $freq_{jd}$  is the # of times term  $j$  appears in doc  $d$
  - $n_j$  = # docs containing term  $j$
  - $N$  = number of docs in collection
- Classic “**term frequency inverse document frequency (tf-idf)**” definition of weight, normalized:

$$u_{jd} = freq_{jd} * \log(N/n_j)$$

$$w_{jd} = \frac{u_{jd}}{(\sum_{i=1}^t (u_{id}^2))^{1/2}}$$

39

## Weight of query components?

- **Set** of terms, **some choices**:
  1.  $w_{jq} = 0$  or 1
  2.  $w_{jq} = freq_{jq} * \log(N/n_j)$   
= 0 or  $\log(N/n_j)$
- **Bag** of terms
  - Analyze like document
  - Some queries are prose expressions of **information need**

*Do we want idf term in both document weights and query weights?*

40

### Vector Model example

**Doc 1:** "Computers have brought the world to our fingertips. We will try to understand at a basic level the **science** -- old and new -- underlying this new Computational Universe. Our quest takes us on a broad sweep of scientific **knowledge** and related technologies... Ultimately, this study makes us look anew at ourselves -- our genome; language; music; "**knowledge**"; and, above all, the mystery of our intelligence. (cos 116 description)

Frequencies:

**science** 1; **knowledge** 2; **principles** 0; **engineering** 0

**Doc 2:** "An introduction to computer **science** in the context of scientific, **engineering**, and commercial applications. The goal of the course is to teach basic **principles** and practical issues, while at the same time preparing students to use computers effectively for applications in computer **science** ..." (cos 126 description)

Frequencies:

**science** 2; **knowledge** 0; **principles** 1; **engineering** 1

41

### Vector model example cont.

- Consider the 5 100-level and 200-level COS courses as the collection (109, 217, 226)
- Only other appearance of our 4 words is "**science**" once in 109 description.
- idf:
  - science**  $\ln(5/3) = .51$
  - engineering, principles, knowledge:**  $\ln(5/1) = 1.6$

42

Term by Doc. Table:  $freq_{jd} * \log(N/n_j)$

	Doc 1	Doc 2
<b>science</b>	.51	1.02
<b>engineering</b>		1.6
<b>principles</b>		1.6
<b>knowledge</b>	3.2	

43

Unnormalized dot product for query:  
**science, engineering, knowledge, principles**  
 using 0/1 query vector

- Doc 1: 3.71
- Doc 2: 4.22
- If documents have about same vector length, this right ratio for normalized (cosine) score

44

## Additional ways to calculate document weights

- Dampen frequency effect:  
 $w_{jd} = 1 + \log(\text{freq}_{jd})$  if  $\text{freq}_{jd} > 0$ ; 0 otherwise
- Use smoothing term to dampen effect:  
 $W_{jd} = a + (1-a) \text{freq}_{jd} / \max_p(\text{freq}_{pd})$ 
  - $a$  is typically .4 or .5
  - Can multiply second term by *idf*
- Effects for long documents (Section 6.4.4)

45

## Classic IR models - Taxonomy

### Well-specified models:

- ✓ Boolean
- ✓ Vector
- Probabilistic
  - based on probabilistic model of words in documents

46

## Probabilistic Model Brief Overview

- Binary Independence Model
  - Original model
  - Chapter 11
- Language Model
  - More commonly used
  - Chapter 12

47

## Key probability ideas

- $P(A)$  probability event A occurs
- $P(A,B)$  probability both A and B occur
- $P(A|B)$  probability A occurs given B occurs
- $P(A,B) = P(A|B)*P(B) = P(B|A)*P(A)$

48

## Language Model

- Is probability measure over strings over some “vocabulary”
- Want to estimate the model of a document
- Want to rank using this estimate
- Assume terms independent

$$P(t_1, t_2) = P(t_1) * P(t_2)$$

49

## Ranking function: main idea

- Rank using  $P(\text{doc}|\text{query}) = P(d|q)$
- $P(d|q) = P(q|d) * P(d)/P(q)$
- Estimate  $P(d|q)$ 
  - Don't care about numerical accuracy
  - **❖ Care that preserves ranking order!**
- Assume uniform distribution  $P(d) \Rightarrow$  ignore
- $P(q)$  constant for give query  $\Rightarrow$  ignore
- Left with  $\text{rank} \approx P(q|d)$
- Estimate  $P(q|d) = \prod_{t \text{ in } q} P(t|d)$  for terms  $t$   
 $= \prod_{t \text{ in } q} (\text{freq}_{td} / \text{length}_d)$

50

Extending  
classic information retrieval  
for today's possibilities

51

## Ranking

- What intuitive criteria?

52

## Enhanced document model

- First model: set of terms
  - term in/not in document
- Next: bag of terms
  - know **frequency** of terms in document
- Now: sequence of terms + **additional properties of terms**
  - sequence gives you **where term** in doc
    - derive **relative position** of multiple query terms
  - **Special use?** (e.g. in title, font, ... )
    - most **require “mark-up”**: tags, meta-data, etc.

53

## HTML mark-up example

`<h2> <font color="#A52A2A"> Communication </font></h2>`  
 This course will be essentially “paperless”. All assignments will be posted `<i>only</i>` on the course Web site. “Handouts” and copies of any transparencies used in class will be posted on the course Web site as well. Important announcements on all aspects of the course will be made on the `<a href="announce.html">` Announcements`</a>` page. `<b>`Students are responsible for monitoring the postings under “Announcements”. `</b>` Schedule changes will be made on the on-line `<a href="schedule.html">` schedule page`</a>`. and announced under “Announcements”. The only paper we will exchange is your solutions to the problem sets, which we will grade and hand back, the exam questions and your responses, and your project reports.

54

## yields

### Communication

This course will be essentially “paperless”. All assignments will be posted *only* on the course Web site (see Schedule and Readings). “Handouts” and copies of any transparencies used in class will be posted on the course Web site as well. Important announcements on all aspects of the course will be made on the [Announcements](#) page. **Students are responsible for monitoring the postings under “Announcements”**. Schedule changes will be made on the on-line [schedule page](#). and announced under “Announcements”. The only paper we will exchange is your solutions to the problem sets, which we will grade and hand back, the exam questions and your responses, and your project reports.

55

## Enhanced document model: restate

“sequence of terms + properties of terms”

⇓ WHY?

“set of (term, properties) pairs”

Properties:

- for each distinct term
  - Frequency of term in doc
    - Vector model of classic IR
- for individual occurrence of each term
  - **Where** in doc.
  - properties of use

56

## Model

- **Document**: set of (term,properties) pairs
- **Query**: **sequence** of terms
  - Can make more complicated
- **Satisfying**: AND model
  - relax if no document contains all?
- **Ranking**: wide open function
  - info beyond documents and query ?

57

## Data Structure for Collection

- for each document, keep list of:
  - **terms** appearing
    - **aggregate properties of term**  
e.g. frequency
  - **positions** at which each term occurs
    - **attributes** for each occurrence of term
- keep summary information for documents

58

## Data Structure for Collection: Invert

- for each term, keep list of:
  - **documents** in which it appears
    - **positions** at which it occurs in each doc.
    - **attributes** for each occurrence
- keep summary information for documents
- keep summary information for terms

59

## Inverted Index for Collection

- for each term, keep **POSTINGS LIST** of:
  - each **document** in which it appears
    - each **position** at which it occurs in doc. **POSTING**
    - **attributes** for each occurrence
- Core structure used by query evaluation and document ranking algorithms

60

### Index structure

$term_1:(doc\ ID\ (position,\ attributes)$   
 $(position,\ attributes),$   
 $\dots$   
 $(position,\ attributes)\ )$   
 $(doc\ ID\ (position,\ attributes)$   
 $(position,\ attributes),$   
 $\dots$   
 $(position,\ attributes)\ )$   
 $\dots$   
 $term_2:(doc\ ID\ (position,\ attributes)$   
 $(position,\ attributes),$   
 $\dots$   
 $(position,\ attributes)\ )$   
 $\dots$

61

### Models have seen

Model	Document	Query	Satisfy
Boolean	set of terms	Boolean expression over terms	evaluate boolean expression
Vector	t-dimensional vector	t-dimensional vector	vector measure of similarity Doc.s ranked by score
dictionary of t terms			
Extended	set of pairs (term, properties)	sequence of terms	Boolean AND Doc.s ranked; flexible scoring algorithm

62

### Query models advantages

- Boolean
  - No ranking in pure
  - + Get power of Boolean Algebra:
    - expressiveness
    - optimization of query forms
- Vector
  - + Query and document look the same
  - + Power of linear algebra
  - No requirement all terms present in pure

63

### Query models advantages

- Extended
  - + could use full Boolean Algebra to define satisfying documents
  - query and document not same
  - ranking arbitrary function of document and query
    - computational cost?

64