# LIVE SESSIONS

## ABSTRACT WORD AND DOCUMENT MODEL

word → sequence of characters
separated by **delimiters** → some languages don't have easy delimiters

↳ Sometimes case sensitive    AT&T → at

WORD ⟷ INTEGER WORD IDs

## WORD AND DOCUMENT IDs

DOCUMENT IDs → Randomly assigned
⇓
Better assignment for Index Compression

MAP: WORD → WORD IDs
↳ Can be compressed

INDEXING
Transposing + Compression ⟹  BINARY MATRIX
Doc IDs × Word IDs

### DOCUMENT REPRESENTATION

SEQUENCE          SET
↓                 ↓
Complete Representation    Loses Information

## INCIDENT VECTORS AND BOOLEAN QUERIES :

Phrase Queries
↓
Need to keep track of positions of words

↳ Document is either relevant or irrelevant ⟹ No notion of relative relevance
→ Emails
→ Library Catalogue

## INVERTED INDEX : → Variable-size posting list ⟍

*INDEXER STEPS  ① Tokenisation
② Sort → a. Terms
          b. Document ID
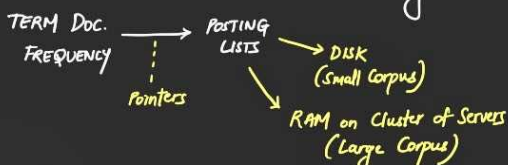③ Dictionary and Postings → Also calculate document frequency

Can be Compressed ⟸ [

### DATA STRUCTURE
- Variable Arrays
- Linked lists

STORAGE

TERM DOC.    →  POSTING
FREQUENCY       LISTS → DISK (Small Corpus)
   ┊ Pointers        ↘ RAM on Cluster of Servers (Large Corpus)

USE OF DOCUMENT FREQUENCY
↓
INTERSECTION always reduces set size.
* If you start with a smaller set, overall process will consume less time.

## QUERY OPTIMISATION :
- Merging / Process in increasing order of frequency

### STREAMS OR CHANNELS
- Raw text streams
- Lowercase token stream
- POS stream
  ↳ Compound Token
* Named Entity Tag. Ex. PERSON GEOLOCATION

### POSITIONAL QUERIES :
- "A" within x words of "B"
- Phrase "A B"
* SOLUTION: Retain relative offset of position in indexes

⇒ Why Compression? Disk transfer is slow but decompression algorithms are fast
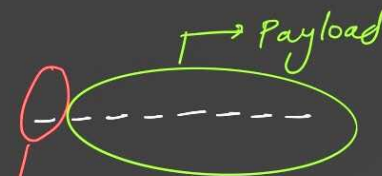
# COMPRESSION TECHNIQUES

→ Higher bit requirement

① Store GAPS vs DOC IDS directly

                                                 ⇒ VARIABLE LENGTH CODES
to reduce size

CASE ① ⇒ High frequency word → Gaps will be low.

CASE ② ⇒ Low frequency word → Few large gaps ⇒ lower bit consumption.

# VARIABLE LENGTH ENCODING :

→ Payload

① **VARIABLE BYTE CODE** ⟶ For every byte

                                           Continuation bit ←

      * Not very efficient for
small gaps → Gap ~ 5 (3-bits)

② **UNARY CODE** ⟶ n 1s with a final 0. ⎤→ Useful for cases where
    ⇒ Not very efficient in general.             number is very small

③ **GAMMA CODE**   - Cutoff leading bit               $Pr(N=n) \propto 2^{-n}$

        - Encode length of offset and append in front.

Ex. 13 → 1110 101          ↳ In Unary     BIT REQUIREMENT → $2\lfloor \log G \rfloor + 1$

**Q. Find the probability for which gamma code is optimal?**

# COMPRESSING THE TERM LIST :

→ DICTIONARY-AS-A-STRING : → Pointer to the next word shows end of current word
→ Can do binary search over the string

SPACE: 4 bytes for frequency
4 bytes for pointer to Posting
3 bytes per term pointer
Avg 8 bytes in term string

400k terms × 19 ⇒ 7.6 MB    vs.   11.2MB for fixed width

→ BLOCKING : Adjacent words share prefix

**EXAMPLE**
ADAM
ADAMANT ⟩ CANCEL 0 APPEND ANT
ADVERSE ⟩ CANCEL 5 APPEND VERSE
vs 7.1 MB

What is k = 4 ?

## POSITIONAL INFORMATION

d gaps ← [DOC ID ; POSITIONS] → P gaps  * Occurence of words in a document are "BURSTY"

# LIVE SESSION — 4

- GFS /HDFS
- MAP REDUCE → Bulk synchronous parallel computation paradigm
- SSTABLE → Sorted String Table : Immutable Key-Value Pair
- BIG TABLE (Hbase)
- PERCOLATOR

→ INDEX CONSTRUCTION
MAP: collection → list(termID, docID)
REDUCE: reduce (<term IDI, list(docID) >...
→ (postings list1, postings list 3 ...)

PARSERS : → Reads document (1 at a time) and emits (term, doc) pairs
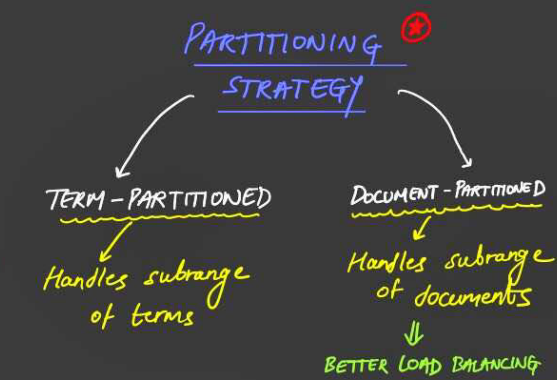
MAP REDUCE FOR WORD COUNT :

| MAP | REDUCE (<t, list(counts)>) |
|---|---|
| for each doc d | total ← 0 |
| for each word t | for each count in |
| output (t,1) | total += count |
| | output < t, (total)> |

**PARTITIONING** ⊛ **STRATEGY**

TERM-PARTITIONED
Handles subrange of terms

DOCUMENT-PARTITIONED
Handles subrange of documents
⇓
BETTER LOAD BALANCING

Network Overhead ← But disadvantage is broadcasting query to all partitions

* Query log is so skewed that great amount of computation is required for most popular words

## OTHER REAL WORD ISSUES :

→ Source format and language detection
→ Sentence and word delimiter → Abreviation vs Full stop
→ Case normalisation (MIT vs mit)
→ Morphological normalization ("stemming")
→ Compound word detection
→ Multilingual dictionary.

* Partitioning for fault tolerance and memory constraint

# RELEVANCE RANKING:
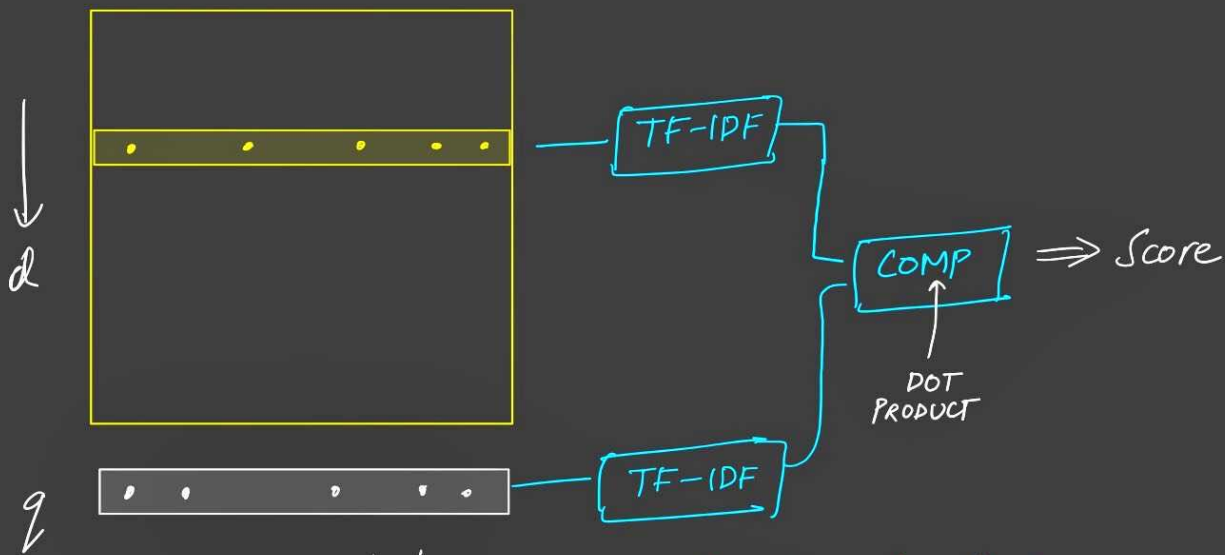
* Writing precise boolean queries not good for most users

**NEED:**
Average web query $\Rightarrow$ 2-3 words
Bits of entropy in query = 10 bits  (5 bits per word)
Corpus > 10 G Docs    ~ 1024 items
Query Result > 10 M hits

Ranked Retrieval $\rightarrow$ Between $\wedge$ and $\vee$
* Returns results according to relevance

**BOOLEAN SEARCH:** Feast or Famine
* Intersection Caching : Performance boost



$d$

$q$

1-Bit Similarity
Jaccard Coefficient    $\Big\}\rightarrow$  Different words with diverse information content are all treated at par.
Cosine similarity

$\rightarrow$ Can store in dictionary

**RARITY:** $\longrightarrow$ Encoded in Inverse Document Frequency (IDF)

**TF × IDF:** (A) Rescaling word dimension according to relevance.

$\ell_w(x) = TF(x,w) \cdot IDF(w)$

Normalise to final vector : $x_w = \ell_w(x) / L(x)$

$L(x) = \sqrt{\sum_w \ell_w(x)}$

## BASIC TFIDF VECTOR SPACE SCORING:

$\rightarrow$ Init accumulator map : score [docid]
$\rightarrow$ In decreasing order of IDF order of query words  get $(x, TF(x,w))$
    $\quad$ score[x] += TF(x,w) * IDF(w)
                                        $\Big\}\rightarrow$ Report top k
$\rightarrow$ Divide each score by length of score $\Rightarrow$ COSINE

# COMPUTATION TIME

QUERIES
- → Rare terms
  - * Accumulator management
  - * Sparse of dense map?
- → Frequent terms
  - * Bit processing for decompressing postings
  - * Arithmetic to update accumulators
  - * Wasted effort in computing score to throw away

## SCORE/IMPACT ORDERING

→ Order postings by decreasing impact

⊗ PROBLEM: Each word will have a different list

## TERM IMPACT :
Term impact $= \dfrac{TF(x,w) * IDF(w)}{L(x)}$

* Documents sorted with Impact IDs

Q. When people do impact ordered posting list, do they do something clever to encode document ID?

INVERTED INDEX
- → Doc ID Ordered Posting
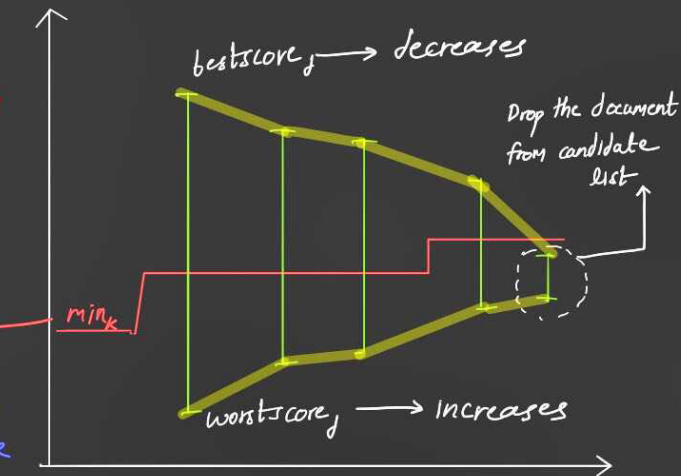  - * DAAT
- → Impact Ordered
  - * TAAT
    - → Threshold Algorithm
    - → quit/continue

→ Completely score a document before moving to next

*** Introduce randomness to counter PRESENTATION BIAS

## QUIT / CONTINUE HEURISTICS

QUIT : Once |score| exceeds some value, quit and report answers

CONTINUE : Stop creating new accumulators but continue processing for remaining words



bestscore$_j$ → decreases

Drop the document from candidate list

increases ← min$_k$

worstscore$_j$ → increases

## RELEVANCE SUPERVISION
- → Regression
- → Ordinal Regression
- → Complete Rank Order
- → Prefix of rank order
- → Pairwise preferences

## LOSS AND REWARD :

① L1 or L2   $|\hat{y} - y|^2$  or  $|\hat{y} - y|$

$y$ → Gold Score
$\hat{y}$ → Predicted Score

② RELATIVE AGGREGATED GOODNESS

$$\dfrac{\sum_{v \in \hat{T}_k} y(v)}{\sum_{v \in T_k} y(v)} \in [0,1]$$

③ PAIR PREFERENCE VIOLATION

If $u < v$ but $y(u) > y(v)$
⇑
Count number of violations

④ RANK CORRELATION

* Compute rank correlation with unrealistic ground truth ranking.

⑥ PREFIX RANK CORRELATION

$m$ pairs $v, w$ from $T_k \cup \hat{T}_k$

Concordant pairs ⇒ $(y(v) - y(w))(\hat{y}(v) - \hat{y}(w)) > 0$

Discordant pairs ⇒ $(y(v) - y(w))(\hat{y}(v) - \hat{y}(w)) < 0$

CONCEPT
- → Precision @ k
- → Recall @ k

# MEAN RECIPROCAL RANK

$$MRR = \frac{1}{|Q|} \sum_{q \in Q} \frac{1}{r_q}$$

*Dropping from $1 \to 2$ as bad as $2 \to \infty$

⭐⭐ Good for navigation queries

Truncated at rank $k \Rightarrow MRR = \frac{1}{|Q|} \sum_{\substack{q \in Q \\ \xi \leq k}} (1/r_q)$

*Mean rank is not a good metric.
1 "hard" query can mess up entire score

## ROC $\to$ Receiver Operating Characteristic Curve

Plotting true positive @ $i = (n_i^+/n^+)$
vs false positive @ $i = (n_i^-/n^-)$

Good Ranking functions $\Rightarrow$ Area close to 1
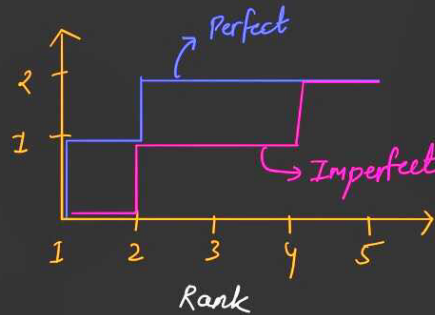
Effectively random $\Rightarrow \boxed{1/2}$



True positive rate (y-axis) vs False positive rate (x-axis)

### AUC
⭐ Area under the curve.

Another AUC Example $\Rightarrow$

\# Good vs Rank



\#Good vs Rank — Perfect, Imperfect

⭐ Area between 2 plots related to number of Discordant Pairs Q

## CONCORDANT AND DISCORDANT PAIRS :

Following is the ranking assigned by engine :
$$\Rightarrow 1 \leq P_1 < P_2 \cdots P_R$$

$$\sum_{i=1}^{R} (P_i - 1) - (i - 1) = Q$$

$$r_{engine} = d_1^+, d_2^-, d_3^+, d_4^+, d_5^-, d_6^-, d_7^+, d_8^-$$

$$r_{ideal} = d_1^+, d_3^+, d_4^+, d_7^+ ; d_2^-, d_5^-, d_6^-, d_8^-$$

$$\Rightarrow \sum_{i=1}^{R} P_i = Q + \frac{R(R+1)}{2}$$

$$= Q + \binom{R+1}{2}$$

## AVERAGE PRECISION :

*Bounding average precision given Q
$\to$ Minimise average precision subject to Q

$$Avg \, Prec = \frac{1}{R} \sum_{i=1}^{R} \left(\frac{i}{P_i}\right)$$

$$\mathcal{L}(P_1, \cdots P_R ; \lambda) = \frac{1}{R} \sum_{i=1}^{R} \frac{i}{P_i} + \lambda \left(\sum_{i=1}^{R} P_i - Q - \binom{R+1}{2}\right)$$

$$\frac{\partial \mathcal{L}}{\partial P_i} = -\frac{i}{R P_i^2} + \lambda = 0 \Rightarrow P_i^* = \sqrt{\frac{i}{R\lambda}}$$

$$\Rightarrow Avg \, Prec \, (r_{engine} ; r_{ideal}) \geq \left(\sum_{i=1}^{R} \sqrt{i}\right)^2 \Big/ R \left(Q + \binom{R+1}{2}\right)$$

False
Negative

False
Positive

$Recall = \dfrac{TP}{FN+TP}$

$Precision = \dfrac{TP}{FP+TP}$

→ HARMONIC MEAN

$F_1 = \dfrac{2RP}{R+P}$

Maximise $F_1$ ⟶ $F_1$ Score

True Positive

## NORMALISED DISCOUNTED CUMULATIVE GAIN:

COMULATIVE

$$NDCG_q = Z_q \sum_{j=1}^{L} \dfrac{2^{r_q(j)}-1}{\log(1+j)}$$

→ GAIN

$r_q(j) \in \{0,1\}$

IRRELEVANT ⟵⟶ RELEVANT

→ RANK DISCOUNT

Normalization
so that perfect ordering
has value 1

## WORD EMBEDDINGS:

word → real vector $R^{300}$

⊛ Similar meaning vectors have higher cosine
and low distances

Document ⟹ 300 dimensional vector
(DENSE)

Clustering requires LOCALLY SENSITIVE HASH FUNCTION
⇓
** Similarity Search

Q. Why hash functions?

⊛ Tree-based indices do not do well in
higher dimensions