

Document structure:

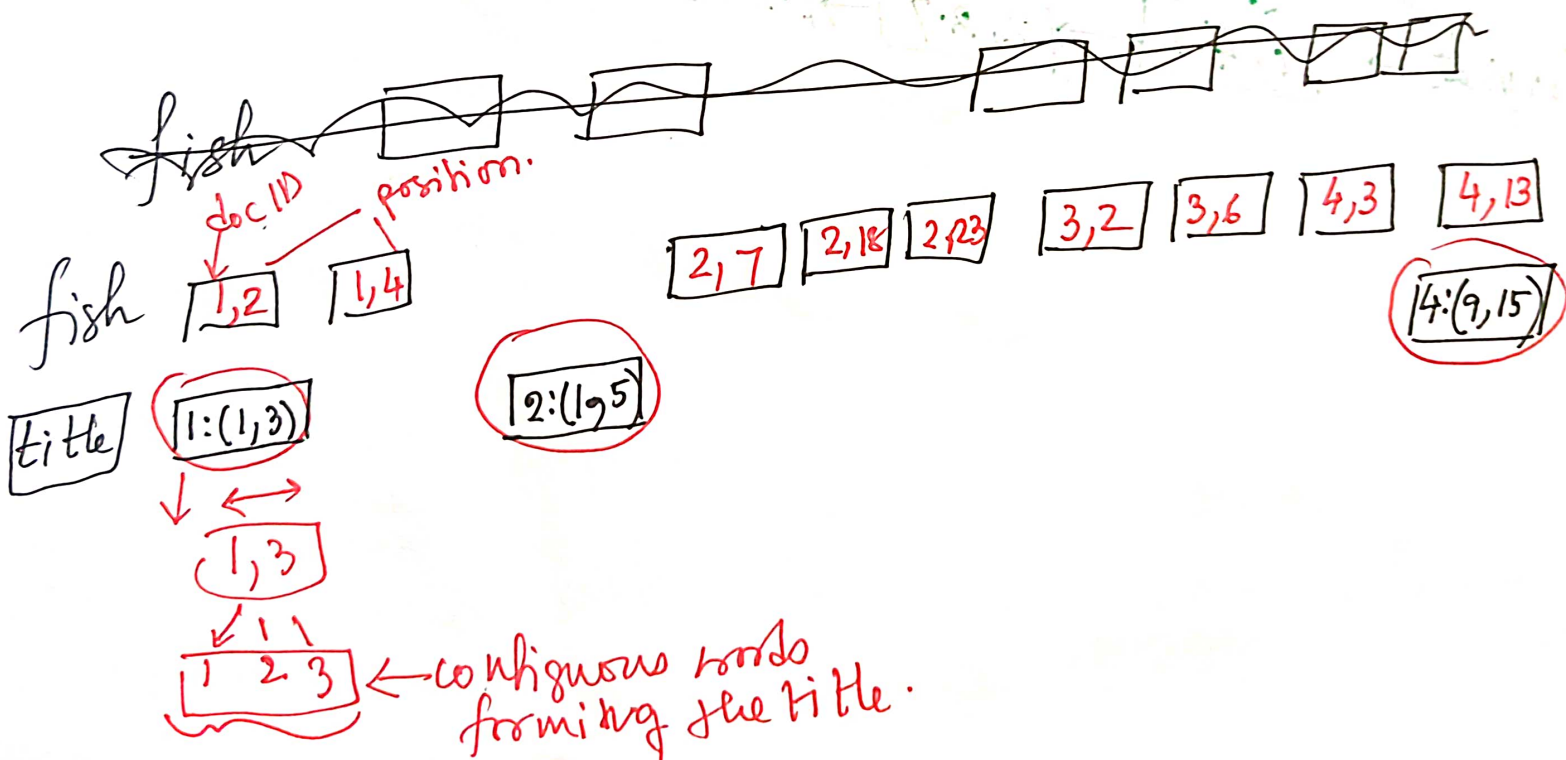
↳ fields (special):

↳ date.

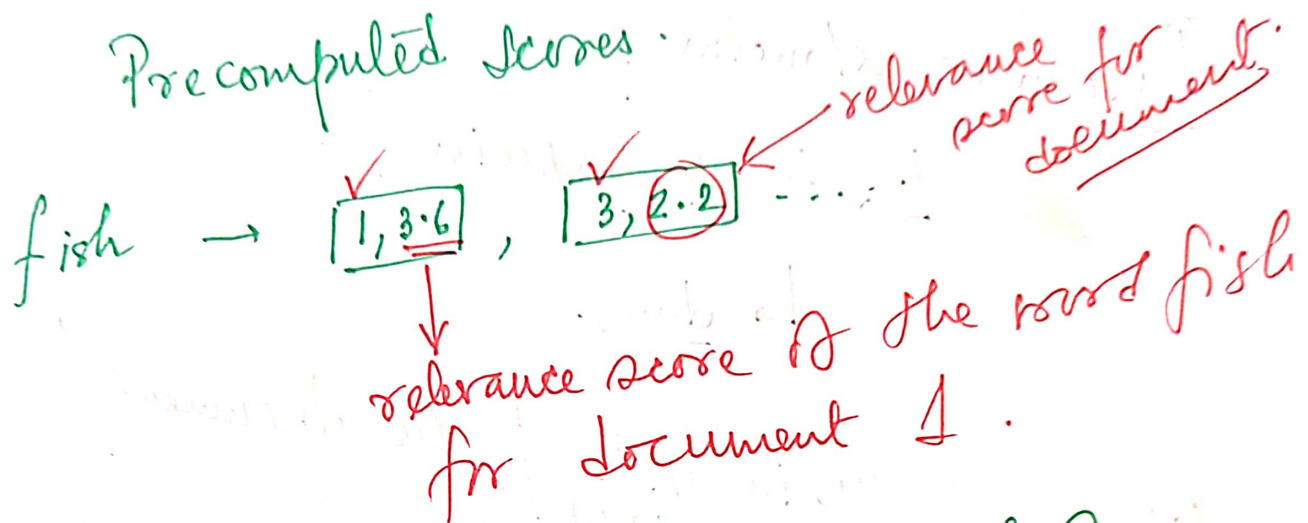
titles. (main title of the document/
section titles).

Index these special fields separately.

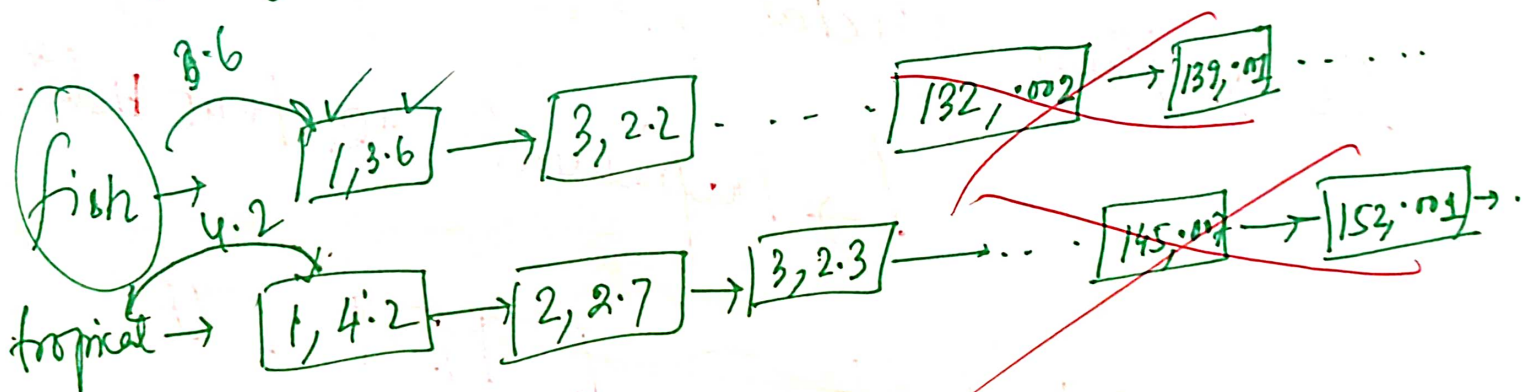
↳ implement "extent & info".



Precomputed scores.



Merge \rightarrow how does this help?



Abstract Model of Ranking for search / retrieval.

→ key words.

Fred's Tropical Fish Shop is the best place to find tropical fish at low, low prices. Whether you're looking for a little fish or a big fish we've got what you need. We even have fake seaweed for your fish tank (and little surf boards too)

9.7 fish

4.2 tropical

22.1 tropical fish

8.2 seaweed

4.2 surf boards

tropical / keep good features.

14 Incoming links

3 days since last updated

Quality features.

tropical fish Query

5.2 fish

3.4 tropical

9.9 tropical fish

1.2 chichlids

0.7 baskets.

1.2 incoming links.

0.2 update content

(R) Ranking function

Doc → [9.7 4.2 22.1]

Q → [5.2 3.4, 9.9]

document feature (fi)

Query feature (qi)

$$R(Q, D) = \sum_i q_i(Q) f_i(D)$$

tropical fish aquarium. | a |

Question(s): What is size of the result set for this query?

Query: "a b c" \rightarrow how many pages will have all three terms.

f_a = frequency of the word 'a' in the corpus.

f_b = frequency of the word 'b' in the corpus.

f_c = frequency of the word 'c' in the corpus.

occurrence of a, b, c are independent & N is the number of tokens in corpus. What is the expected frequency of "a b c"?

prob of a's occurrence = f_a/N

prob of b's " = f_b/N

prob of c's " = f_c/N

prob of occurrence of abc $\rightarrow ?$

$$f_a/N \times f_b/N \times f_c/N$$

Expected frequency

$$f_{abc} = N \times (f_a/N \times f_b/N \times f_c/N)$$
$$= (f_a \cdot f_b \cdot f_c) / N^2$$

Independence assumption is bad.

<u>Nodo</u>	<u>Doc freq.</u>
tropical	120,900
fish	1,131,855
aquarium	26,480
breeding	81,885
tropical fish	18,472
tropical aquarium	1921
fish breeding	5510 36,427
aquarium breeding	9722 18,48
tropical fish aquarium	36427 1529
tropical fish breeding	3629

$$P(a \cap b \cap c) = P(a \cap b) \cdot \underbrace{P(c|a \cap b)}_{\rightarrow \text{chain rule.}}$$

$$\min \left[P(c|a), P(c|b) \right]$$

$$P(a \cap b \cap c) = P(a \cap b) \cdot \min [P(c|a), P(c|b)]$$

$$\frac{f_a}{f_{aq}} \cap \frac{f_b}{f_{fish}} \cap \frac{f_c}{f_{tropical}}$$

$$= \frac{f_{aq \cap fish}}{f_{aq}} \times \min \left[\frac{f_a}{f_{aq \cap tropical}}, \frac{f_b}{f_{fish \cap tropical}} \right]$$

$$= \boxed{9722 \times 1921 / 26840}$$

→ 8 documents
find out some C triplets. ← containing aquarium

Instead of processing 26480 docs

→ 3000 docs.

→ 258 cases where
(tropical fish aquarium)

$S = 3000$, $C = 258$ →

3000 → 25
26480 → $\frac{25}{3000} \times 26480$

25%