

The most important concepts

1. Search: BM25 (takes into account both tf and doc length normalization to determine the relevance)

Components:

1. TF (How often term occurs)

2. Length of doc

3. IDF (measure the significance of a term across the entire doc collection.)

2. Embeddings: Word2vec, fastText, GLOVE, BERT

3. Precision/Recall

① Introduction to Multimedia Retrieval

- | | |
|--|--|
| <ul style="list-style-type: none">• What is Retrieval?• What do we want to achieve?- | <ol style="list-style-type: none">1. Task2. General Retrieval Process3. Retrieval models (7)4. Metadata (Two ways for extraction) |
|--|--|

(2)

Evaluation

1. How do validate a search algorithm?
2. How do compare boolean retrieval model?
3. How to compare model with order?
4. How to compare retrieval graded?
5. Confusion Matrix: def?
6. How to find threshold for model?

1. Benchmarks (Sparse and dense)

2. Boolean Retrieval

a) What

b) Precision and recall

c) F-measure

d) Micro- and macro evalution

$$\frac{F_1}{P} \cdot \frac{R}{P}$$

3. Retrieval with order

a) Def

b) MRR

c) P-R curve

d) "Analysing P-R":

- 1) System efficiency
- 2) P@k
- 3) R-precision
- 4) AP
- 5) MAP

4. Retrieval with Graded Relevance

1) def

2) CG-k = Eval:

3) DCG-k = \sum

4) Normalized DCG

5. Confusion Matrix (TP, TN, FP, FN, ACC, Sensitivity, Specificity)

6. Optimizing Hyperparameters

1) Threshold

2) Pdt (sensitivity + specificity)

3) ROC-curve, AUC

③ Classic Text Retrieval



1. Fundamentals of Retrieval

I. Offline phase: (Aim: create vec. with features)

1. Extract:

- (metadata, url, text, html tags) • Clean document

2. Split (Fixed sizes, min. tokens, metadata or structure, semantic splitting)

3. Tokenize (Seq of characters)

(Chars, words, n-grams and phrases)

* Stemming → Porter (5 steps with ~20 rules)

$$[C](VC)^m[V]$$

* Tagging of tokens

d) Summarizing → creating vector with features.

1. Create vector

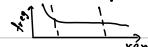
- set-of-words
- bag-of-words

* inverted file method to store.

2. Eliminate unimportant terms

1) Stop words

2) Zipf's Law $p(r) = \frac{c}{r}$, $c = \frac{1}{\sum r}$, M = # distinct terms in vocabulary



3) "Similarities". C-centroid doc all N terms with average freq

$$dp(t) = Q_t - Q_{\text{without } t} ; dp(t) \text{ large } \Rightarrow \text{import} ; Q = \sum_i^n \sin(D_i, t)$$

4) Inverse document frequency (idf)

$$idf(t) = \log \frac{N+1}{df(t)+1} ; \Rightarrow \begin{array}{l} \text{Weighted doc} \\ \text{document freq.} \end{array}$$

$t_f(D_i, t)$
number of
occurrences of
term t in D_i

③ Classical Text Retrieval

2. Text Retrieval Models

1. Standard Boolean Model

+ ranking
+ partial matches

- 1)
- 2)

2. Extended Boolean Model (Similarity)

- no theoretic background
- hard to express complex info

- 1) use $t\cdot f(D_i, t_j)$ and $\text{idf}(t_j) + \text{sim}(q, D_i) = 1 - d_{i,j}(e, k)$
- Fuzzy Algebraic, Fuzzy set, p-norm-model, soft boolean

3. Vector Space Retrieval

- no theoretic background
- independence of terms, but it is not always true.

- $d_{i,j} = t\cdot f \cdot \text{idf}$ - term-doc-matrix A

$$q_j = t\cdot f \cdot \text{idf}$$

- Search: $1. \text{sim}(Q, D_i) = q \cdot d_i = \sum q_j \cdot d_{ij}$ - inner vector product

$$2. \text{Sim} = \frac{q \cdot d_i}{\|q\| \cdot \|d_i\|} - \text{cosine}$$

• Actually we use only direction of vectors!

4. Probabilistic Retrieval (Formal approach)

- assumptions
- feedback

- 1) BIR: 3 assumptions $\Rightarrow \text{sim}(Q, D_i) \sim \sum c_j, c_j = \log \frac{r_j \cdot (1 - r_j)}{n_j \cdot (1 - r_j)}$

1. Initial step (frequency) a unique case because there are no annotations

2. Feedback step (user annotation)

5. BM25 (Vector space model + probabilistic approach)

- Take into account length of doc + frequency

3. Index Structure

Inverted index

- 1. Boolean model

2. BIR: a) doc-at-a-time b) term-at-a-time

+ we can add predicates

- 3. For Vector Space model

same, but $+(freq, sim-f, \text{idf})^{\text{add}}$

6. Lucene scalability

- 1) segment level
- 2) sharding
- 3) replicate shards
- 4) multiple regions

④ Advanced Text Retrieval

① Chunking Text

- 1) Split into fixed-sized chunks (langchain library)
- semantic 2) Splitting at sentence boundaries (NLTK: Punkt; spaCy)
size control 3) Splitting on structure (Merge small parts + overlapping for semantic)
- 4) Semantic splitting (1. Def similarity (Vector space model or embeddings), 2. Set max/min chunk size, split into sentences NLTK or spaCy; merge chunks based on similarity)

② Tokenization Revisited

- 1) Divide by words: naive, nltk, spaCy, for script continua Amazon Codewhisperer
- 2) Continuous stream: HMM, overlapping seq + sub-seq search $t_1(t_1, t_2)$
- 3) n-grams: Pointwise Mutual Information (PMI); $\log \frac{P(t_1, t_2)}{P(t_1) \cdot P(t_2)}$
4. Likelihood Ratios (LHR):
 - $H_1: P(t_2 | t_1) = P(t_2 | \neg t_1) = p$
 - $H_2: P(t_1, t_2) = p_1 \cdot p_2 = P(t_2 | t_1)$

⚠ One-hot vectors → Embedding + positional encoding

- (transformers library) 4) BPE tokenization: 1) initialization (chars or subword units)

split → merge

- 2) freq count (bag-of-words)
- 3) merge most freq pair
- 4) repeat

5) tokenize text

- 5) Word Piece Tokenization (Extended BPE in 2 steps:
1. position beginning + middle
2. pairs frequency)

③ Lemmatization and linguistic Transformation

1. Rule based stemmers (Porter, Lancaster, Snowball)

2. Dict-based stemmers (WordNet, spaCy) ← include POS

* Words combination: split ← rule based + freq. of the components

* Synonyms → extended dict with predet synonyms

* Homonyms → POS, Mh., most common meaning

* Hyponyms (broad) → weights; Faceted search, expand queries, relevance ranking

* Hypnyms (narrow)

* Spelling mistakes: spellchecker

④ Part of Speech (POS and NER (name entity recognition))

1. Rule-based POS tagging (like -ing → verb)
2. Stochastic POS tagging (HMM)
3. Transformation-based POS (rule-based + correcting the errors iteratively)
4. DH POS tagging

⑤ Latent Semantic Analysis (LSA)

1. Latent Semantic Indexing (LSI)

$$1. \text{ SVD: } A = U S V^T \quad \begin{matrix} \text{doc} \\ \text{term} \end{matrix} = \begin{matrix} \text{U} \\ \vdots \end{matrix} \begin{matrix} S \\ \text{X}_k X_k^T \end{matrix} \begin{matrix} V^T \\ \vdots \end{matrix}$$

$$A = S_1(U_1 V_1^T) + S_2(U_2 V_2^T) + \dots + S_r(U_r V_r^T)$$

• LSI treats queries as mini-docs: $q^T = q^T U_k S_k^T$

• Sim. fun (dot-product or cosine measure) to compare the query

2. Embeddings (gensim library, spacy)

1. Word2vec

a) Skip-Gram Model ($\text{word} \rightarrow \text{word}$)

b) Continuous Bag of Words (CBOW) Model ($\text{word} \rightarrow \text{word}$)

2. GloVe (self-supervised, considers word co-occurrences within context windows.)

3. fastText (sub-word representations for center word) $\xrightarrow{\text{skip-gram}} \text{word} \rightarrow \text{word}$

4. BERT (NN based embedding; transformer based)

3. Embeddings in Text Retrieval

1. Retriever-Ranker ($\text{BM25} + \text{embedding}$) $\xrightarrow{\text{ranker}}$

2. Semantic Search with Reader

3. Retrieval-Augmented Generation (RAG)

$(q \rightarrow \text{Retriever} \rightarrow \text{Generator})$

4. Language Detection

• Naive Bayes (posterior = $\frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}}$) multinomial

5. Sentiment Analysis

1. Naive Bayes (Bernoulli distr.)

2. Text-CNN

3. Transformer-based classification

5 Web Social

• Types of search queries:

1. Navigational search queries (reach specific website)
2. Transactional search q. (ended purchasing smth)
3. Informational search q. (info)

① Additional Features (compare to classical text retrieval)

- 1) Proximity of query terms in docs
- 2) Penalties for low-quality pages
- 3) Language and location awareness
- 4) Term boosting if it occurs in link
- 5) Term occurrence weighting based on doc
- 6) Proximity of query terms in doc

② Page Rank ()

- rank represents the likelihood of a random surfer visiting that page. $PR(p) = \frac{1-d}{N} + d \sum_{q \in M_p} \frac{PR(q)}{L(q)}$
- solving iteratively \rightarrow fast due to sparse matrix \rightarrow linear

③ Hypertlinked-Induced Topic Search (HITS)

- Authorities $a(p)$ and Hubs $h(p)$; Problems:
 - 1. Overpowered Authorities and Hubs
 - 2. Automated, non-related Links
 - 3. Query Term Frequency
- ↓ modification
 - 1. One domain - one voice
 - 2. Filter unrelated pages, using reference doc, cosine and thresholds.

④ Other relationships

1. Co-citation analysis
2. Stochastic Approach for Link Structure Analysis (SA.LSA) \leftarrow HITS extension
 - Bipartite graph \rightarrow two-step markov chain \rightarrow probability
3. Sim Rank (Similarity between nodes, based on their connections)



⑥ Vector Search

Motivation:

Idea: each doc D_i yields an N -dim embedding vector d_i . we have query vec. q . We want to compare.

1. Best Match Problem:

Given d_i and q , find \hat{d}^* that optimizes the similarity fun.

$$\hat{d}_{\text{dot}} = \underset{d}{\operatorname{argmax}} q \cdot d_i ; \quad d_{\text{cos}}^* = \underset{d}{\operatorname{argmax}} \frac{q \cdot d_i}{\|q\| \|d_i\|}$$

2. Normalization of vectors: (Gaussian normalization)

$$\hat{v}_j = \frac{v_j - \mu_j}{\sigma_j} \rightarrow (0.13, [0, 1000])$$

- Using eigenvalues and eigenvalues we can rotate and scale the original space to get normalized space

$$3. \text{MS}(q) = \underset{p}{\operatorname{argmax}} \text{Sim}(p, q), \quad \text{NN}(q) = \underset{p}{\operatorname{argmin}} \text{dist}(p, q)$$

3. Low-dim Search Structure

Voronoi Cell

Gridfile

Nearest Neighbor Problem

R-Tree

- Optimal Search Method by Somp (using priority queue)

- We can utilize lower-upper bounds to "navigate" a search

4. Curse of Dimensionality

As the number of dims increases, the data becomes increasingly sparse. In high-dim spaces, most data points concentrate near the corners of the hypercube representing the feature space.

5. Quantization based NN-search

- VA-File \rightarrow grid with points

6. Approximate NN-Search (FAISS)

1. Vector Transformers (

2. Coarse quantizers

3. Encodings (Fine quantizers)

4. Refiners

]] VA-File - approximation using lower bit precisions, while maintaining

a small error.

$$\begin{array}{r} 1 \\ 2 \end{array} \begin{array}{r} 0.6 \\ 0.5 \end{array} \begin{array}{r} 0.8 \\ 0.0 \end{array} \rightarrow \begin{array}{r} 1011 \\ 0011 \end{array}$$

④ Basic Image Retrieval

def = illumination - measure energy per sec, accounting for the sensitivity of the human eye to different light wavelength.

- Image Normalization → Image Segmentation → Extraction → Aggregation
- Global feature extraction:
 1. Feature Set
 2. Feature Concatenation
 3. Feature Summary
- Segment images:
 1. Global feature
 2. Static Segmentation
 3. Object Segmentation

1. Image Normalization

1. Gray scale $i(x,y) : N^2 \rightarrow [0,1]$

2. Color image $i(x,y) : N^2 \rightarrow [0,1]^3 = \begin{bmatrix} R \\ G \\ B \end{bmatrix}$

3. Quantization to replace fixed-point numbers in image representation:

• True color (32 bit): $\hat{f}(x,y) : N^2 \rightarrow [0,255]^3$ $\hat{f}(x,y) = \frac{f(x,y)}{255}$

• Deep color (64): $\rightarrow [0,65535]$

4. Standardize the dataset

1. Rotate-invariant features

2. Histogram normalization

3. Grayscale transformation

4. Scaling

5. Affine Transformation

6. Convolution (Kernels, remember to normalize to get colors in range)

2. Color Information

1. Common color systems: 1. CIE 2. RGB 3. HSL/HSV 4. YUV 5. CMYK

2. Color: Brightness + Chromaticity (Hue, Saturation)

CIE - based on human perception, XYZ-color space

RGB - most common, large gamut range
SRGB

Conversion
RGB to CIE
gamma - compression
printing

device independent

CIE Lab color space

L_a (Green-Red), L_b (Blue-Yellow)

device independent

easy to perform color manipulations.

3. Color Features

1. Color Histogram

Compare: Manhattan, Euclidean

Quadratic ← to consider similarity

• Also consider luminance

• Consider only chromaticity

2. Color Moments (Lab)

as hist. resistant to noise, scaling, transform.
do not account for spatial relat. between colors
not independent.

3. Texture information:

1. Structural approach

1. Edge magnitude and directions (Sobel operator on smoothed Gaussian image).

2. Gradient Histogram (2 values per pixel - direction magnitude)

3. Gradient Moments

4. Law's Texture Energy (Combine 1D patterns to create 2D texture filters. This filters capture dif. aspects of texture, like edges, spots, waves).⁵

2. Statistical approach

3. Fourier approach

1. Gabor Moments (^{before we use rotation} Use to capture various frequency

ranges and directions. Reveals info about the presence of selected freq in frame) $G(x,y) = \frac{e^{-\frac{(x-x_0)^2}{2\sigma_x^2}}}{2\pi\sigma_x} \cdot e^{j2\pi f_x(x-x_0)}$
• We use it through convolution, each pixel mapped to energy value based on chosen G.filter's direction and scale.

4. Shape information

• Edge Detection

• Contour Detection

• Region-based Segmentation

Once shape is available we can calculate:

• Area • Centroid • Different Ratio⁶ length, shape resembles a circle etc.

1. HOG (Histogram of Oriented Gradients)

2. SIFT (Scale-Invariant Feature Transform)

⑧ Advanced Image Retrieval

1. Perceptron concept
2. SVM
3. Feed Forward
4. Error Function
5. SGD
6. Backpropagation (Vanishing gradient or exploding)
7. Overfitting and underfitting
8. Occam's razor
9. Activation Function (ReLU, leaky ReLU, parametric ReLU, exponential Linear Unit)
10. Generative Adversarial Network (GAN)
- L_1/L_2 regularization
11. Normalization
 1. Batch
 2. Layer
 3. Weight Standardization
12. Deep learning
 1. Convolution
 - 2.

--	--