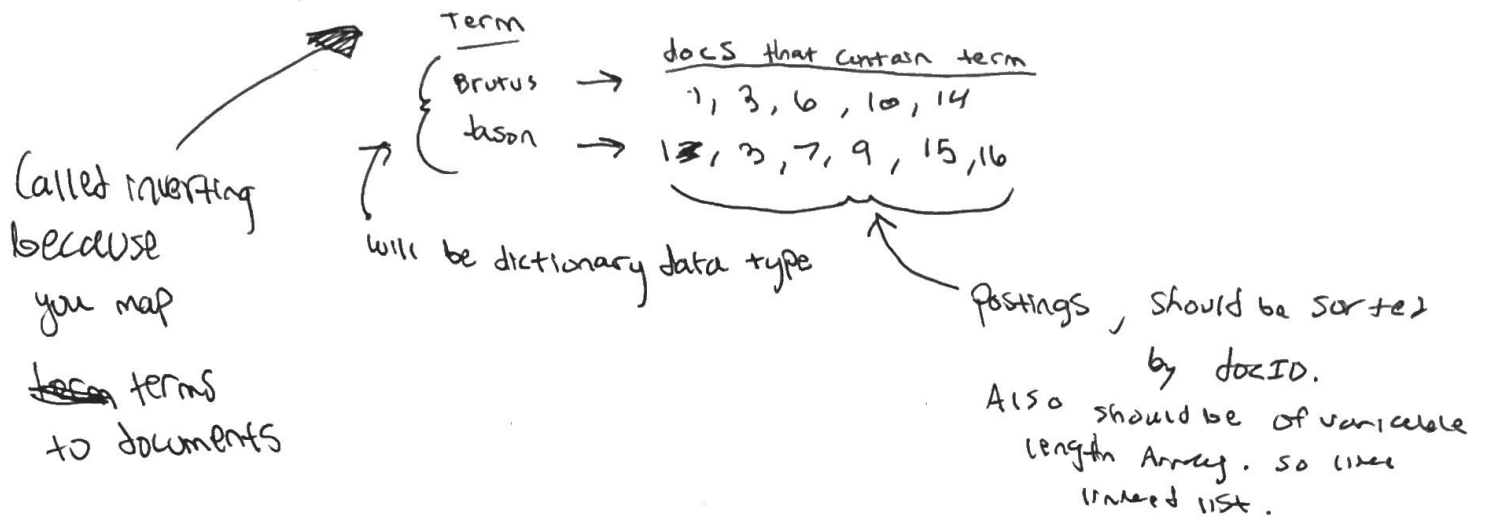# Indexing:

- Inverted Index: For each term t ~~t~~, we must store a list of all documents that contain t.

  - Identify each doc by docID



Called inverting because you map ~~term~~ terms to documents

Term
{ Brutus → docs that contain term
     1, 3, 6, 10, 14
  Jason → 1̶3̶, 3, 7, 9, 15, 16 }

will be dictionary data type

Postings, should be sorted by docID.
Also should be of variable length Array. So like linked list.

- Inverted index steps

  1) Collect documents to be indexed
  2) Tokenize documents.
     - go through each document and assign each term with the doc it's in.



doc 1
| Jason goes
| to store |

doc 2
| Kelsey
| goes here |

tokens
Jason — 1          Kelsey — 2
goes  — 1          goes — 2
to    — 1          here — 2
store — 1

  3) modify ~~tokens~~ to make things like "friends" to "frend"
  4) Sort Alphabetically by term.
  5) Index the term listings.
     - multiple term entries in a single document are merged
     - Split into dictionary and postings as shown at top of page.
     - Doc freq. information is added.

# INDEXING :

- Positional indexes :

  - In the postings of inverted index, store positions in which the tokens appear.
  
  $\langle$ term, # of docs containing term,
  
     doc#, position 1, position 2, etc ,
  
     doc #, position 1, position 2, etc,
  
     $\rangle$

- Zipf's law : frequency of any word is inversly proportional to its Rank in the frequency table.

  - Basically words like I, you, an appear a lot in documents, but words like "etymology" appear less so they are Ranked higher.
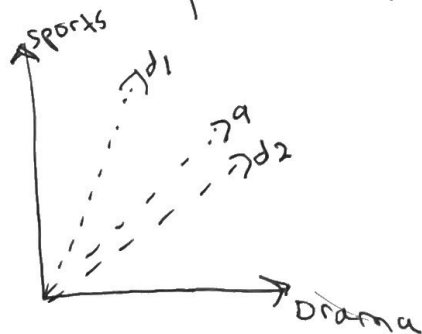
# RANKING:

- Vector space model:

### Boolean retrival vs Ranked Retrival

| Boolean retrival | Ranked Retrival |
|---|---|
| - good for expert users with persice understandy of collection. | - most users incapaloe of writing boolean quenes. |
| - easily comb through lot's of documents | - most urers dont want to wade through looos of results. |
| | - boolean retrival produoos too few or too many results |

- Represent document and query by concept vectors.
  - Each vector is one-dimension.

- measure similarity between query vector and document vector.



- x,y are concept vectors, which are independaut of the documents.

- Vectors have a weight of each concept and query has weight to defino importauce of concept.

- Two ways to define weghts
  - Tf: count of a term t in doc d.
    - wugh a term more if it occurs more in a document.
  - IDF: Assign a higher weight to rare terms.
  - Tf-idf vector: cross product of Tf and IDF
- Compute cosine similarity on Tf-Idf vectors.

# RANKING:

BM25 - Goal is to be sensitive to term frequency and document length while not adding too many parameters

- Normalize ~~document length~~ using ~~Norma~~ document length

- 2 parameters
  - term frequency scaling for document
  - ~~document length normalization~~
  - term freq. scaling for query

- BM25 is better than vsm because of the document length normalization.

# CLASSIFICATION:

3 ways to classify

- manually → consistent when problem size and team is small. Cannot be scaled to a large dataset
- Hand-coded rule-based classifiers
- supervised learning
  - └→ Classify document d to a class c. Classes are predefined.
    a training set of D documents is used with label in C.
    Naive bayes, k-nearest neighbors
    All need a hand-classified training data set.
- Evaluation
  Classification Accuracy - $r/n$    $n$ = # test docs
                           $r$ = # of test docs correctly classified
        recall,
         F1,

# CLASSIFICATION :

Rocchio - Create centroids which is an average vector of all documents that belong to a class.

~~Compare their~~ Attempt to classify a document to the most related centroid.

*worse than Naive Bayes*

KNN - define K the number of documents to consider when classifying a document,

*No training necisary.*
*scals well to large # of classes.*
*more accurate than Rocchio, Naive Bayes*

## Variance vs Bias —

bias - Simplifying assumptions made by the model

Variance — amount that the estimate of the target function will change given different brantraining data.

KNN high variance, low bias

Rochion/NB low variance, high bias

Naive Bayes — Assign probabilities to each class that the document inquestion belongs to that class.

*good dependable baseline for text classification*

$$\text{Probabilty for a class} = \frac{\left(\begin{array}{c}\text{# times term appears} \\ \text{in class's docs}\end{array}\right)}{\text{total # words in class's docs}}$$

*sometimes apply laplace smoothing "add 1"*