

# Data-Intensive Distributed Computing

CS 431/631 451/651 (Fall 2021)

## Part 4: Analyzing Text (1/2)

Ali Abedi

These slides are available at <https://www.student.cs.uwaterloo.ca/~cs451>

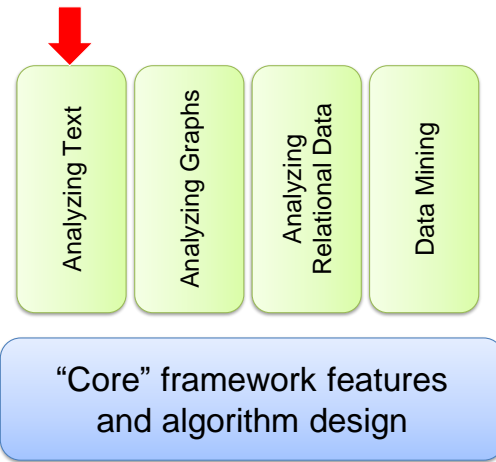


This work is licensed under a Creative Commons Attribution-NonCommercial-Share Alike 3.0 United States  
See <http://creativecommons.org/licenses/by-nc-sa/3.0/us/> for details

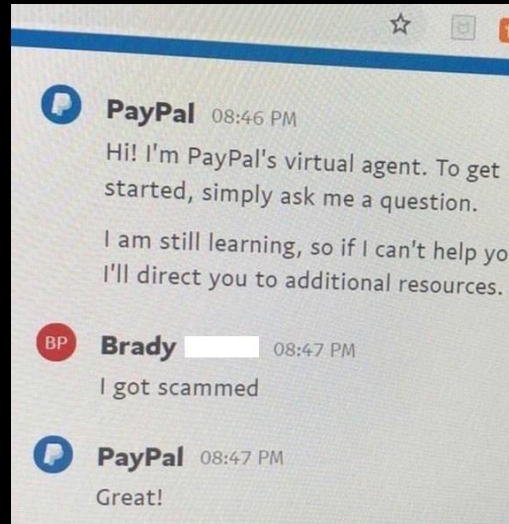
## Structure of the Course

“Core” framework features  
and algorithm design

## Structure of the Course



# Natural Language Processing



**Pairs. Stripes.**  
Seems pretty trivial...

More than a “toy problem”?  
Answer: language models

## Language Models

$P(w_1, w_2, \dots, w_T)$     Assigning a probability to a sentence

Why?

- **Machine translation**
  - $P(\text{High winds tonight}) > P(\text{Large winds tonight})$
- **Spell Correction**
  - $P(\text{Waterloo is a great city}) > P(\text{Waterloo is a grate city})$
- **Speech recognition**
  - $P(\text{I saw a van}) > P(\text{eyes awe of an})$

Slide: from Dan Jurafsky

Sentence with T words - assign a probability to it

## Language Models

$$P(w_1, w_2, \dots, w_T)$$
$$= P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \dots P(w_T|w_1, \dots, w_{T-1})$$

[chain rule]

$$P(\text{"Waterloo is a great city"}) =$$
$$P(\text{Waterloo}) \times P(\text{is} | \text{Waterloo}) \times P(\text{a} | \text{Waterloo is})$$
$$\times P(\text{great} | \text{Waterloo is a})$$
$$\times P(\text{city} | \text{Waterloo is a great})$$

Is this tractable?

Sentence with T words - assign a probability to it

$$P(A,B) = P(B) P(A|B)$$

## Approximating Probabilities: $N$ -Grams

Basic idea: limit history to fixed number of  $(N - 1)$  words  
(Markov Assumption)

$$P(w_k | w_1, \dots, w_{k-1}) \approx P(w_k | w_{k-N+1}, \dots, w_{k-1})$$

$N=1$ : Unigram Language Model

$$P(w_k | w_1, \dots, w_{k-1}) \approx P(w_k)$$

$$\Rightarrow P(w_1, w_2, \dots, w_T) \approx P(w_1)P(w_2) \dots P(w_T)$$



## Approximating Probabilities: $N$ -Grams

Basic idea: limit history to fixed number of  $(N - 1)$  words  
(Markov Assumption)

$$P(w_k | w_1, \dots, w_{k-1}) \approx P(w_k | w_{k-N+1}, \dots, w_{k-1})$$

$N=2$ : Bigram Language Model

$$P(w_k | w_1, \dots, w_{k-1}) \approx P(w_k | w_{k-1})$$

$$\Rightarrow P(w_1, w_2, \dots, w_T) \approx P(w_1 | \langle S \rangle) P(w_2 | w_1) \dots P(w_T | w_{T-1})$$

Since we also want to include the first word in the bigram model, we need a dummy beginning of sentence marker  $\langle s \rangle$ . We usually also have an end of sentence marker but for the sake of brevity, I don't show that here.

## Approximating Probabilities: $N$ -Grams

Basic idea: limit history to fixed number of  $(N - 1)$  words  
(Markov Assumption)

$$P(w_k | w_1, \dots, w_{k-1}) \approx P(w_k | w_{k-N+1}, \dots, w_{k-1})$$

$N=3$ : Trigram Language Model

$$P(w_k | w_1, \dots, w_{k-1}) \approx P(w_k | w_{k-2}, w_{k-1})$$

$$\Rightarrow P(w_1, w_2, \dots, w_T) \approx P(w_1 | \langle S \rangle) \dots P(w_T | w_{T-2} w_{T-1})$$

## Building $N$ -Gram Language Models

Compute maximum likelihood estimates (MLE) for  
Individual  $n$ -gram probabilities

Unigram  $P(w_i) = \frac{C(w_i)}{N}$

Bigram  $P(w_i, w_j) = \frac{C(w_i, w_j)}{N}$

$$P(w_j|w_i) = \frac{P(w_i, w_j)}{P(w_i)} = \frac{C(w_i, w_j)}{\sum_w C(w_i, w)} = \frac{C(w_i, w_j)}{C(w_i)}$$

Generalizes to higher-order  $n$ -grams

State of the art models use ~5-grams

We already know how to do this in MapReduce!



**SOMETIMES I'LL START A  
SENTENCE AND I DON'T EVEN  
KNOW WHERE IT'S GOING.**

Estimating Probability Distribution  
Sparsity problem

## Example: Bigram Language Model

```
<s> I am Sam </s>  
<s> Sam I am </s>  
<s> I do not like green eggs and ham </s>
```

Training Corpus

$$P(I | \langle s \rangle) = 2/3 = 0.67$$

$$P(\text{am} | I) = 2/3 = 0.67$$

$$P(\langle s \rangle | \text{Sam}) = 1/2 = 0.50$$

...

$$P(\text{Sam} | \langle s \rangle) = 1/3 = 0.33$$

$$P(\text{do} | I) = 1/3 = 0.33$$

$$P(\text{Sam} | \text{am}) = 1/2 = 0.50$$

### Bigram Probability Estimates

Note: We don't ever cross sentence boundaries

## Data Sparsity

$$P(I | \langle s \rangle) = 2/3 = 0.67$$

$$P(\text{Sam} | \langle s \rangle) = 1/3 = 0.33$$

$$P(\text{am} | I) = 2/3 = 0.67$$

$$P(\text{do} | I) = 1/3 = 0.33$$

$$P(\langle s \rangle | \text{Sam}) = 1/2 = 0.50$$

$$P(\text{Sam} | \text{am}) = 1/2 = 0.50$$

...

### Bigram Probability Estimates

$P(I \text{ like ham})$

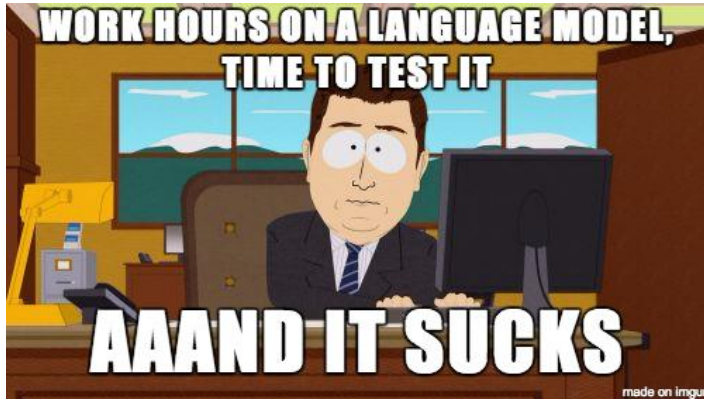
$$= P(I | \langle s \rangle) P(\text{like} | I) P(\text{ham} | \text{like}) P(\langle s \rangle | \text{ham})$$

$$= 0$$

*Why is this bad?*

**Issue: Sparsity!**

Why is the 0 bad ?





## Solution: Smoothing



Zeros are bad for any statistical estimator

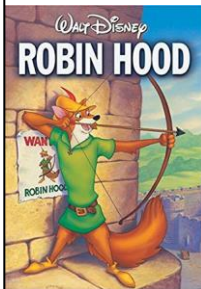
Need better estimators because MLEs give us a lot of zeros

A distribution without zeros is “smoother”

The Robin Hood Philosophy: Take from the rich (seen  $n$ -grams)  
and give to the poor (unseen  $n$ -grams)

Need better estimators because MLEs give us a lot of zeros

A distribution without zeros is “smoother”



Lots of techniques:

Laplace, Good-Turing, Katz backoff, Jelinek-Mercer

Kneser-Ney represents best practice

# Laplace Smoothing

Learn fancy words for  
simple ideas!

Simplest and oldest smoothing technique  
Just add 1 to all  $n$ -gram counts including the unseen ones  
So, what do the revised estimates look like?

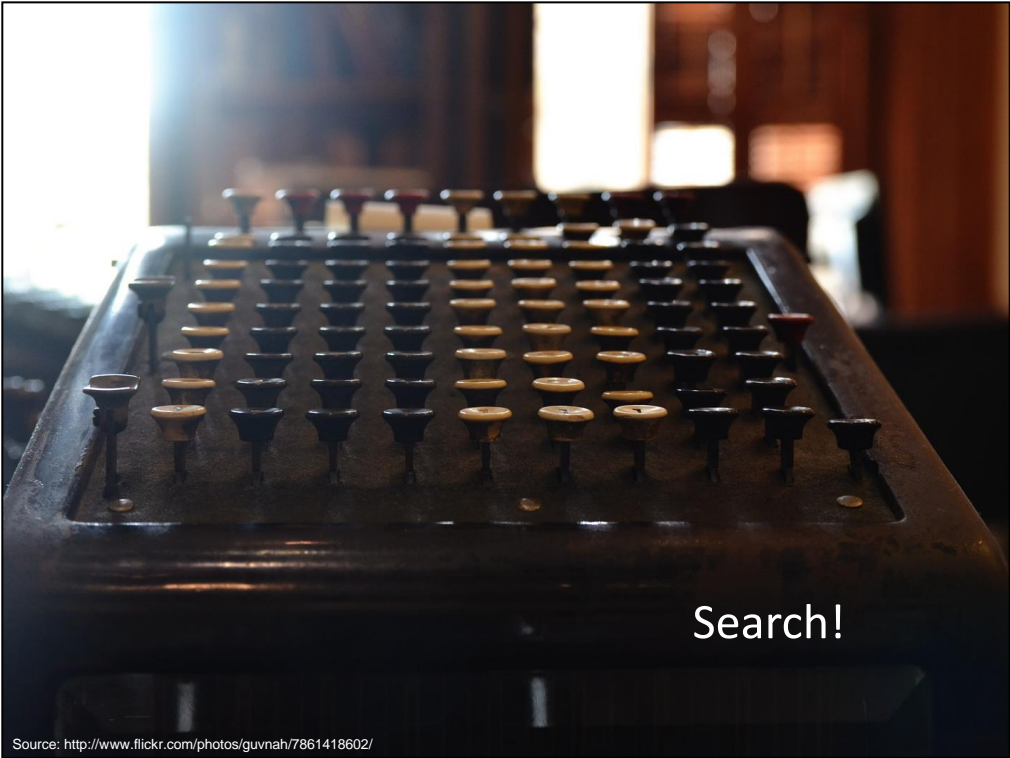
## Laplace Smoothing

$$P_{MLE}(w_i) = \frac{C(w_i)}{N} \xrightarrow{\text{Unigrams}} P_{LAP}(w_i) = \frac{C(w_i) + 1}{N + V}$$

$$P_{MLE}(w_i, w_j) = \frac{C(w_i, w_j)}{N} \xrightarrow{\text{Bigrams}} P_{LAP}(w_i, w_j) = \frac{C(w_i, w_j) + 1}{N + V^2}$$

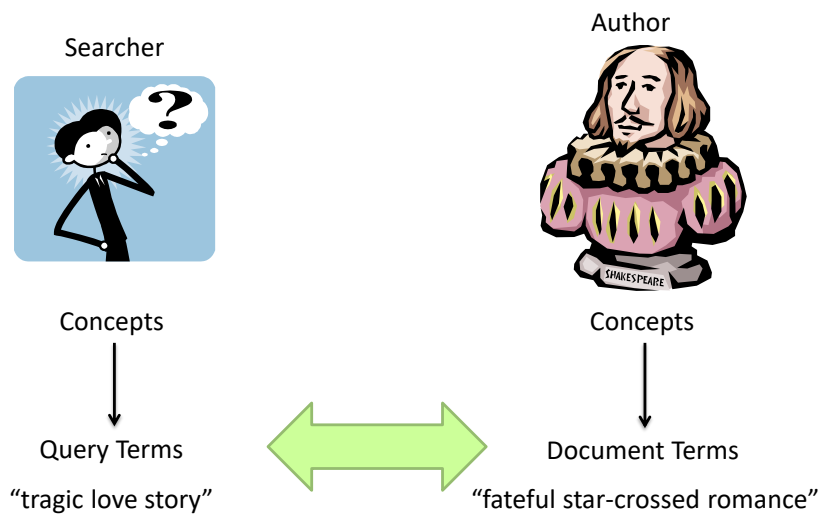
What if we don't know  $V$ ?

You have to make sure that the joint is well-formed and understand how the conditional probability formula is derived.



Source: <http://www.flickr.com/photos/guvnah/7861418602/>

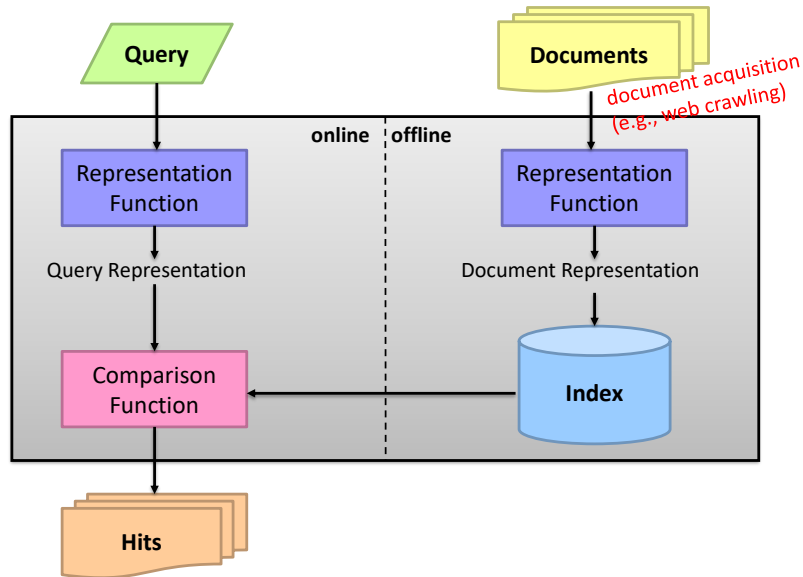
## The Central Problem in Search



Do these represent the same concepts?

Why is IR hard? Because language is hard!

# Abstract IR Architecture



# How do we represent text?

Remember: computers don't "understand" anything!

## "Bag of words"

Treat all the words in a document as index terms  
Assign a "weight" to each term based on "importance"  
(or, in simplest case, presence/absence of word)  
Disregard order, structure, meaning, etc. of the words  
Simple, yet effective!

## Assumptions

Term occurrence is independent  
Document relevance is independent  
"Words" are well-defined

## What's a word?

天主教教宗若望保祿二世因感冒再度住進醫院。  
這是他今年第二度因同樣的病因住院。

الناطق باسم -وقال مارك ريجيف  
إن شارون قبل -الخارجية الإسرائيلية  
الدعوة وسيقوم للمرة الأولى بزيارة  
تونس، التي كانت لفترة طويلة المقر  
1982.الرسمي لمنظمة التحرير الفلسطينية بعد خروجها من لبنان عام

Выступая в Мещанском суде Москвы экс-глава ЮКОСа  
заявил не совершал ничего противозаконного, в чем  
обвиняет его генпрокуратура России.

भारत सरकार ने आर्थिक सर्वेक्षण में वित्तीय वर्ष 2005-06 में सात फ़ीसदी  
विकास दर हासिल करने का आकलन किया है और कर सुधार पर ज़ोर  
दिया है

日米連合で台頭中国に対処...アーミテージ前副長官提言

조재영 기자= 서울시는 25일 이명박 시장이 `행정중심복합도시' 건설안  
에 대해 `군대라도 동원해 막고싶은 심정'이라고 말했다는 일부 언론의  
보도를 부인했다.



# Sample Document

## McDonald's slims down spuds

Fast-food chain to reduce certain types of fat in its french fries with new cooking oil.

NEW YORK (CNN/Money) - McDonald's Corp. is cutting the amount of "bad" fat in its french fries nearly in half, the fast-food chain said Tuesday as it moves to make all its fried menu items healthier.

But does that mean the popular shoestring fries won't taste the same? The company says no. "It's a win-win for our customers because they are getting the same great french-fry taste along with an even healthier nutrition profile," said Mike Roberts, president of McDonald's USA.

But others are not so sure. McDonald's will not specifically discuss the kind of oil it plans to use, but at least one nutrition expert says playing with the formula could mean a different taste.

Shares of Oak Brook, Ill.-based McDonald's (MCD: down \$0.54 to \$23.22, Research, Estimates) were lower Tuesday afternoon. It was unclear Tuesday whether competitors Burger King and Wendy's International (WEN: down \$0.80 to \$34.91, Research, Estimates) would follow suit. Neither company could immediately be reached for comment.

...

## "Bag of Words"

14 × McDonalds

12 × fat

11 × fries

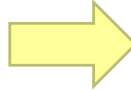
8 × new

7 × french

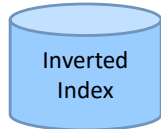
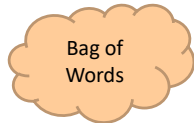
6 × company, said, nutrition

5 × food, oil, percent, reduce,  
taste, Tuesday

...



## Counting Words...



case folding, tokenization, stopword removal, stemming

~~syntax~~, ~~semantics~~, ~~word knowledge~~, etc.



Source: <http://www.flickr.com/photos/guvnah/7861418602/>

**Doc 1**  
one fish, two fish

**Doc 2**  
red fish, blue fish

**Doc 3**  
cat in the hat

**Doc 4**  
green eggs and ham

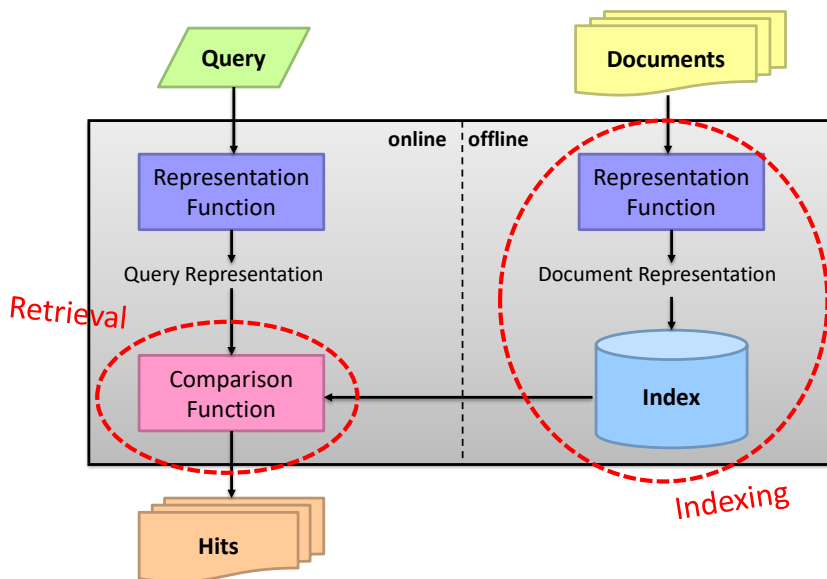
|       | 1 | 2 | 3 | 4 |
|-------|---|---|---|---|
| blue  |   | 1 |   |   |
| cat   |   |   | 1 |   |
| egg   |   |   |   | 1 |
| fish  | 1 | 1 |   |   |
| green |   |   |   | 1 |
| ham   |   |   |   | 1 |
| hat   |   |   | 1 |   |
| one   | 1 |   |   |   |
| red   |   | 1 |   |   |
| two   | 1 |   |   |   |

What goes in each cell?

boolean  
count  
positions

cs451

# Abstract IR Architecture



Doc 1

one fish, two fish

Doc 2

red fish, blue fish

Doc 3

cat in the hat

Doc 4

green eggs and ham

|       | 1 | 2 | 3 | 4 |
|-------|---|---|---|---|
| blue  |   | 1 |   |   |
| cat   |   |   | 1 |   |
| egg   |   |   |   | 1 |
| fish  | 1 | 1 |   |   |
| green |   |   |   | 1 |
| ham   |   |   |   | 1 |
| hat   |   |   | 1 |   |
| one   | 1 |   |   |   |
| red   |   | 1 |   |   |
| two   | 1 |   |   |   |

Indexing: building this structure

Retrieval: manipulating this structure

**Doc 1**      **Doc 2**      **Doc 3**      **Doc 4**  
**one fish, two fish**    **red fish, blue fish**    **cat in the hat**    **green eggs and ham**

|       | 1 | 2 | 3 | 4 |
|-------|---|---|---|---|
| blue  |   | 1 |   |   |
| cat   |   |   | 1 |   |
| egg   |   |   |   | 1 |
| fish  | 1 | 1 |   |   |
| green |   |   |   | 1 |
| ham   |   |   |   | 1 |
| hat   |   |   | 1 |   |
| one   | 1 |   |   |   |
| red   |   | 1 |   |   |
| two   | 1 |   |   |   |



|       |   |       |
|-------|---|-------|
| blue  | → | 2     |
| cat   | → | 3     |
| egg   | → | 4     |
| fish  | → | 1 → 2 |
| green | → | 4     |
| ham   | → | 4     |
| hat   | → | 3     |
| one   | → | 1     |
| red   | → | 2     |
| two   | → | 1     |

*postings lists*

# Indexing: Performance Analysis

Fundamentally, a large sorting problem

Terms usually fit in memory

Postings usually don't

How is it done on a single machine?

How can it be done with MapReduce?

First, let's characterize the problem size:

Size of vocabulary

Size of postings



## Vocabulary Size: Heaps' Law

$$M = kT^b$$

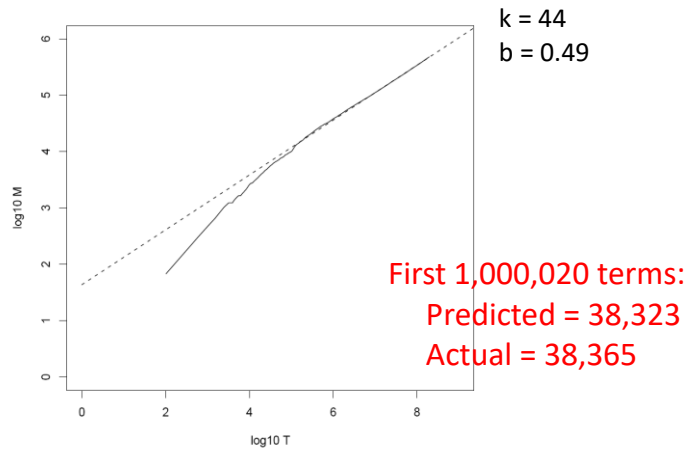
$M$  is vocabulary size  
 $T$  is collection size (number of documents)  
 $k$  and  $b$  are constants

Typically,  $k$  is between 30 and 100,  $b$  is between 0.4 and 0.6

Heaps' Law: linear in log-log space

**Surprise: Vocabulary size grows unbounded!**

## Heaps' Law for RCV1



Reuters-RCV1 collection: 806,791 newswire documents (Aug 20, 1996-August 19, 1997)

## Postings Size: Zipf's Law

$$f(k; s, N) = \frac{1/k^s}{\sum_{n=1}^N (1/n^s)}$$

$N$  number of elements  
 $k$  rank  
 $s$  characteristic exponent

Zipf's Law: (also) linear in log-log space

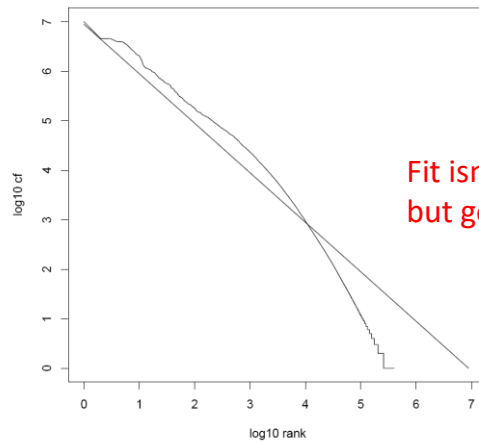
Specific case of Power Law distributions

In other words:

A few elements occur very frequently

Many elements occur very infrequently

## Zipf's Law for RCV1

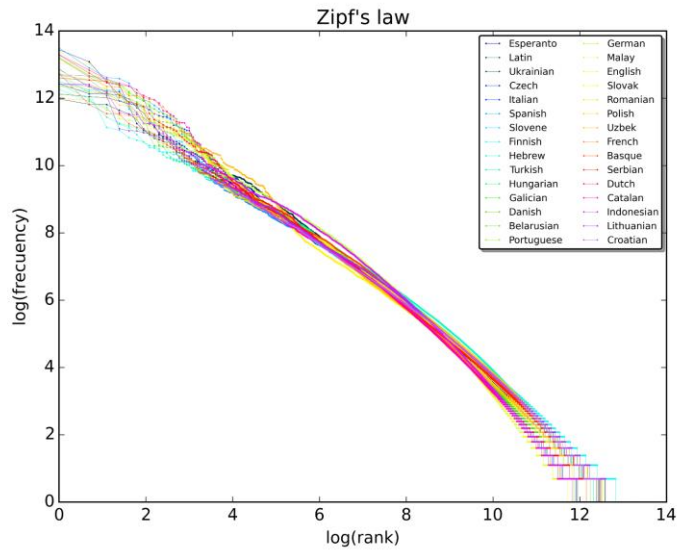


Fit isn't that good...  
but good enough!

Reuters-RCV1 collection: 806,791 newswire documents (Aug 20, 1996-August 19, 1997)

Manning, Raghavan, Schütze, Introduction to Information Retrieval (2008)

# Zipf's Law for Wikipedia



Rank versus frequency for the first 10m words in 30 Wikipedias (dumps from October 2015)

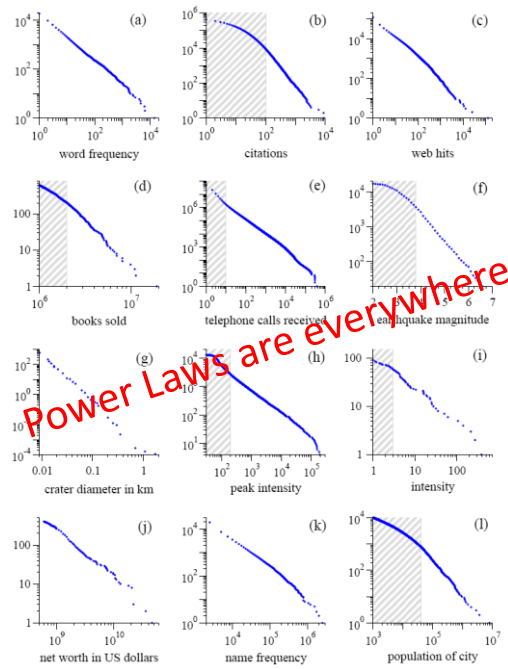


Figure from: Newman, M. E. J. (2005) "Power laws, Pareto distributions and Zipf's law." Contemporary Physics 46:323–351.

# MapReduce: Index Construction

Map over all documents

Emit *term* as key, (*docid*, *tf*) as value

Emit other information as necessary (e.g., term position)

Sort/shuffle: group postings by term

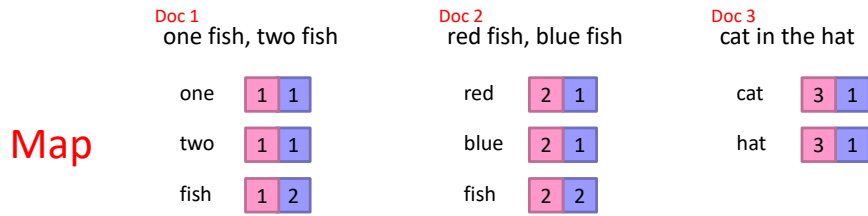
Reduce

Gather and sort the postings (typically by *docid*)

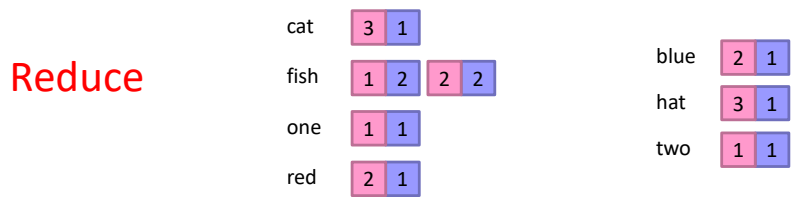
Write postings to disk

**MapReduce does all the heavy lifting!**

# Inverted Indexing with MapReduce



**Shuffle and Sort: aggregate values by keys**





## Inverted Indexing: Pseudo-Code

```
class Mapper {
  def map(docid: Long, doc: String) = {
    val counts = new Map()
    for (term <- tokenize(doc)) {
      counts(term) += 1
    }
    for ((term, tf) <- counts) {
      emit(term, (docid, tf))
    }
  }
}

class Reducer {
  def reduce(term: String, postings: Iterable[(docid, tf)]) = {
    val p = new List()
    for ((docid, tf) <- postings) {
      p.append((docid, tf))
    }
    p.sort()
    emit(term, p)
  }
}
```

What's the problem?

Stay tuned...