

Authorship Identification with Amazon Review Data

Xiaohui Chen

I love this cellphone

This cellphone worked
like a bomb and it
totally ruined my life!!!

Whatever, I don't care



Buyer 1



Buyer 2



Buyer 3

Steps

- Feature Extraction
- Build ML Model
- Tune Model
- Prediction

An Anonymous Text comes in for another product

- “This laptop looks as if it is going to explode because it is always overheating. My life is destroyed by that thing!!!”
- Most likely from buyer 2

Dataset

- data in which all users and items have at least 5 reviews (41.13 million reviews)
- 30GB file size

Sample Review Text

- { "reviewerID": "A2SUAM1J3GNN3B", "asin": "0000013714", "reviewerName": "J. McDonald", "helpful": [2, 3], "reviewText": "I bought this for my husband who plays the piano. He is having a wonderful time playing these old hymns. The music is at times hard to read because we think the book was published for singing from more than playing from. Great purchase though!", "overall": 5.0, "summary": "Heavenly Highway Hymns", "unixReviewTime": 1252800000, "reviewTime": "09 13, 2009" }

Steps for Feature Extraction

- Add review ID, remove unnecessary features, calculate helpful rate
- Sentiment Analysis (polarity, subjectivity)
- Use map reduce to calculate tfidf scores

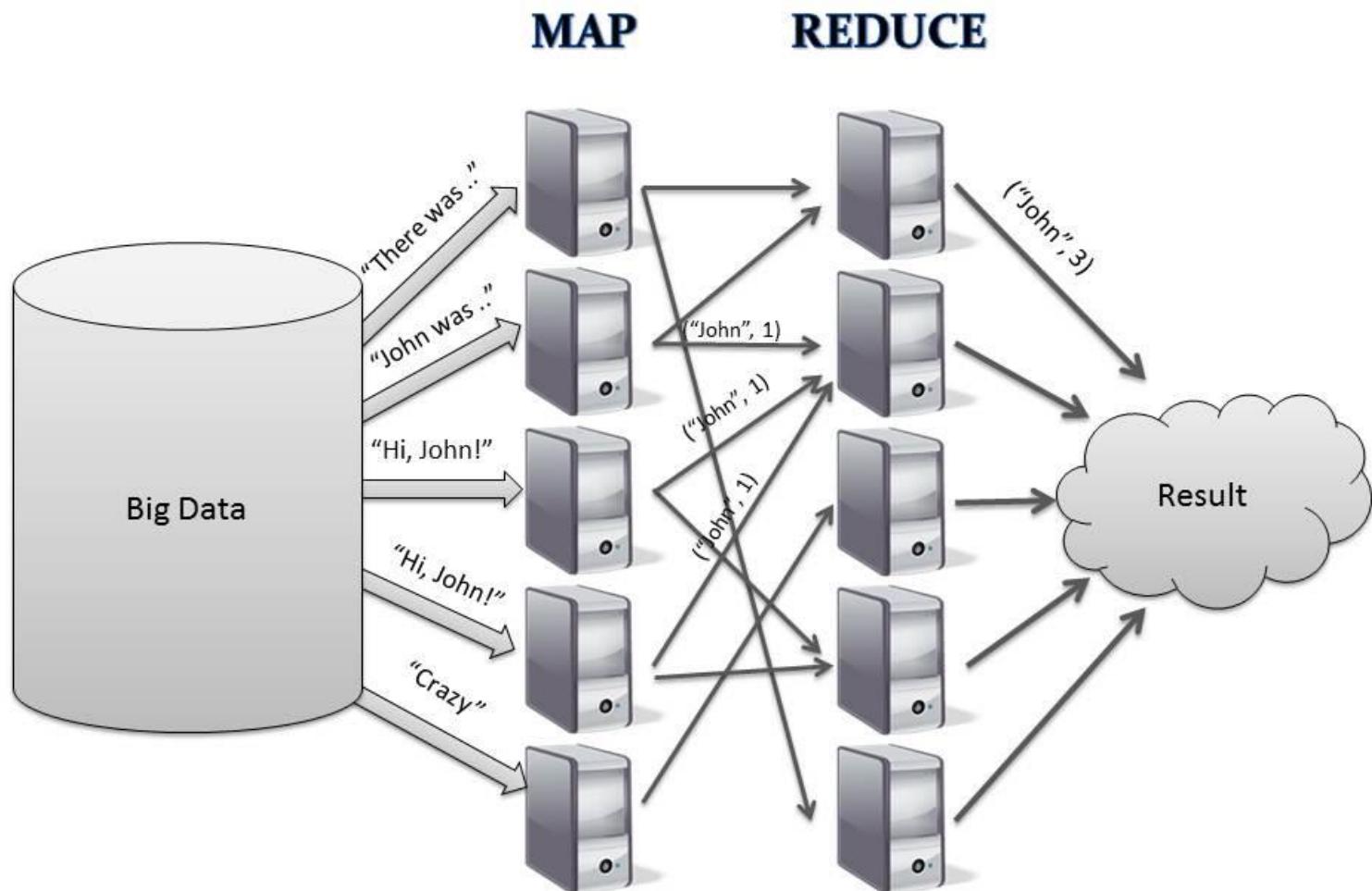
TF-IDF

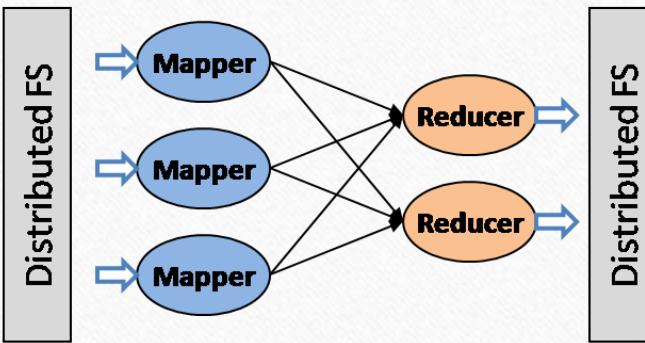
- Wikipedia: In information retrieval, tf–idf, short for term frequency–inverse document frequency, is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus. It is often used as a weighting factor in information retrieval and text mining.

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}}$$

$$idf_i = \log \frac{|D|}{|\{j : t_i \in d_j\}|}$$

$$tfidf_{i,j} = tf_{i,j} \times idf_i$$





$$tf_{wd} = \frac{wordCount_{wd}}{wordsPerDoc_d}$$

$$idf_{wd} = \log(\frac{totalDocs}{docsPerWord_w})$$

$$tf-idf_{wd} = tf_{wd} * idf_{wd}$$

Round 1
wordCount

$$\{docId \Rightarrow [words]\}$$

$$\{[word, docId] \Rightarrow 1\}$$

$$\{[word, docId] \Rightarrow [1, 1 \dots]\}$$

$$\{[word, docId] \Rightarrow wordCount\}$$

Round 2
wordCount per doc

$$\{[word, docId] \Rightarrow wordCount\}$$

$$\{docId \Rightarrow [word, wordCount]\}$$

$$\{docId \Rightarrow [[word1, wordCount1], [word2, wordCount2] \dots]\}$$

$$\{[word, docId] \Rightarrow [wordCount, wordsPerDoc]\}$$

Round 3
docCount per word
 $tfidf_{wd} =$

$$(wordCount_{wd} / wordsPerDoc_d) * \log(\frac{totalDocs}{docsPerWord_w})$$

$$\{word \Rightarrow [[doc1, wc1, wpd1], [doc2, wc2, wpd2] \dots]\}$$

$$\{[word, docId] \Rightarrow [wordCount, wordsPerDoc, docsPerWord]\}$$

$$\{[word, docId] \Rightarrow tfIdf\}$$

Map Reduce Using Amazon EMR

- EMR: Elastic Map Reduce
- Configure to support Map Reduce
- Currently stuck on bootstrapping the EC2 instances

After Feature Extraction

- Put Data into Amazon Redshift (simple, train-test split)
- Use AWS Machine Learning to build the model (multi-class classification, can train dataset up to 100GB)
- For multiclass classification, Amazon ML uses multinomial logistic regression (multinomial logistic loss + SGD).

Questions?

