

Capstone Presentation

A Hybrid Recommendation System with Yelp Challenge Data

Team Yelper_Helper

S. Kickham, S. O'Mullane, R. Rad, A. Rubino, C. Shi

NYC Data Science Academy Cohort 9

06/20/2017

Outline

- Motivation
- Project Summary
- Data Source
- Data Preparation and Exploratory Analysis
- Hybrid Recommendation System
 - Natural Language Processing (NLP)
 - Collaborative Filtering (CF)
 - Social Network
 - Location-based
- Data Pipeline
 - Recommendation Engine
 - App Workflow (Flask - Kafka - Rec Engine - Kafka - Flask)
- Live Demo
- Lessons Learned and Future Work
- Acknowledgement

Background



EDA



Recommendation
Algorithms

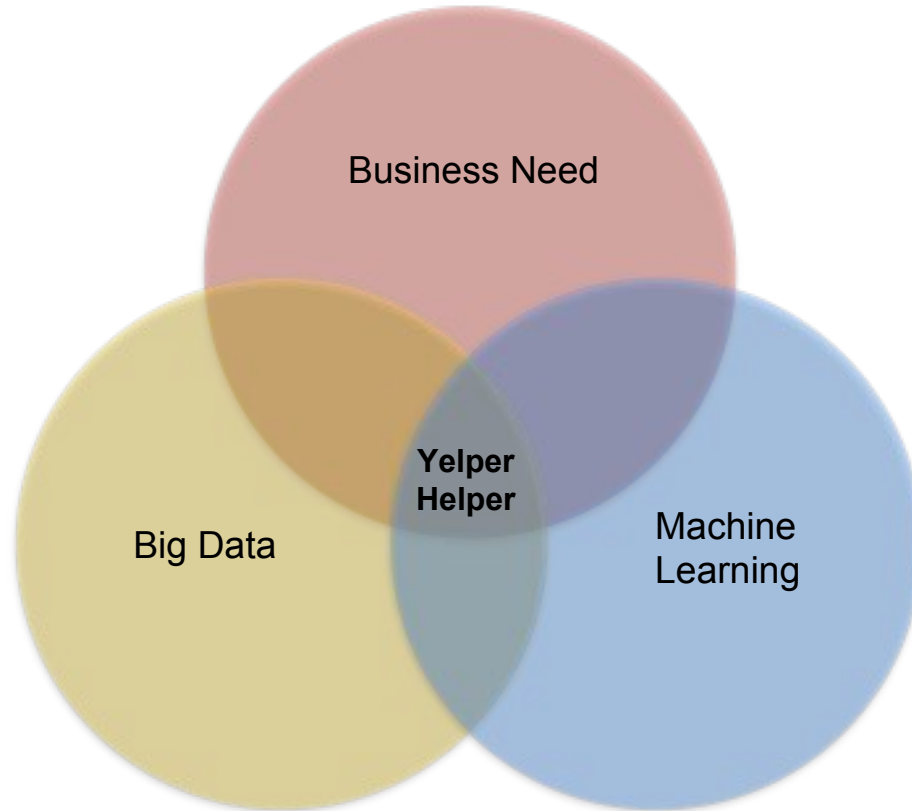


App Pipeline

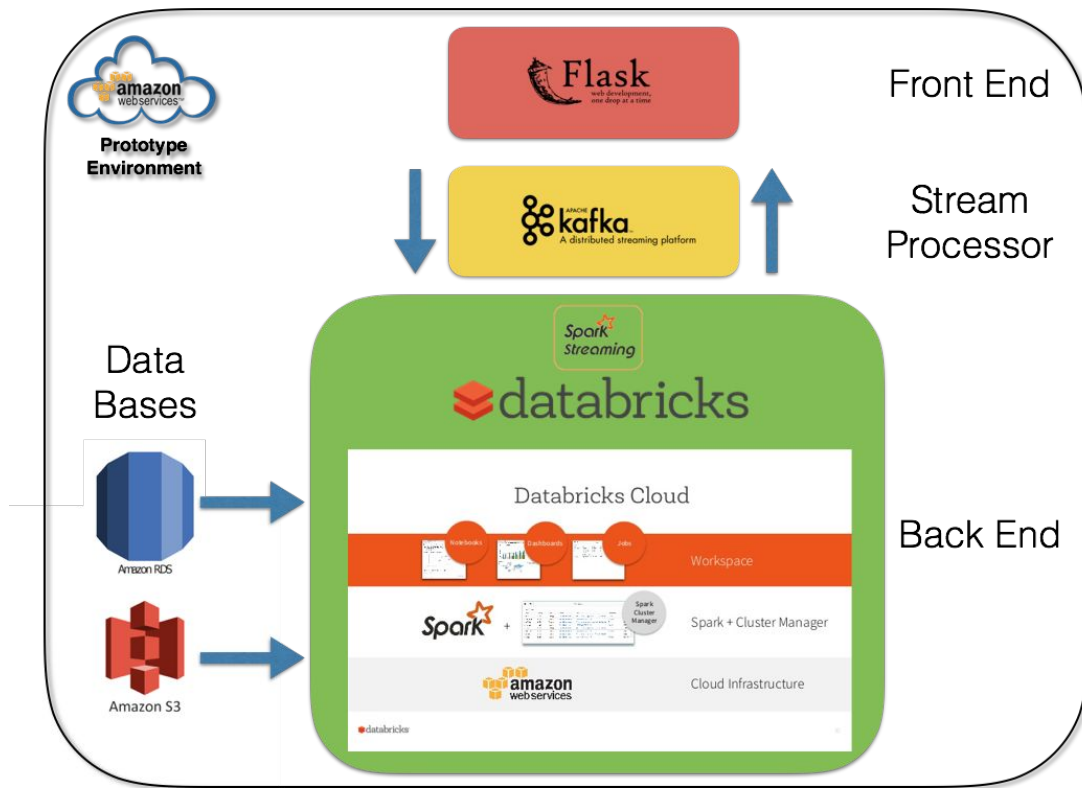


Demo

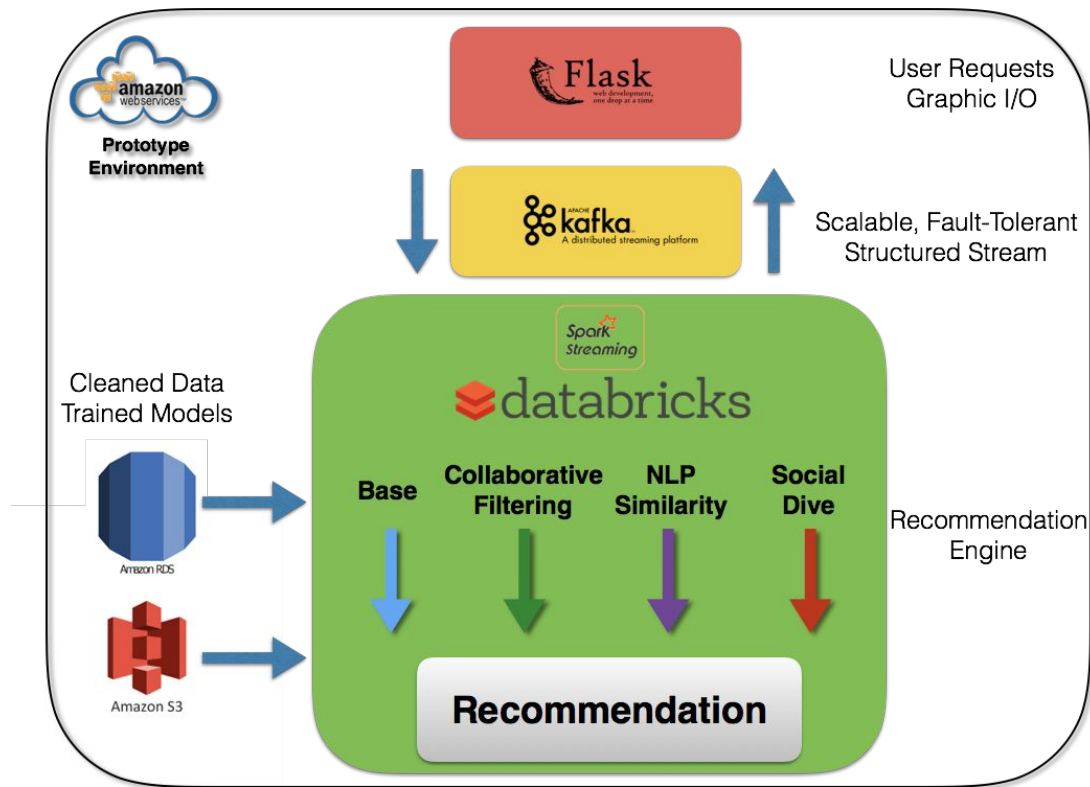
Motivation



Project Architecture



Project Architecture



Data Source

Yelp Dataset Challenge

Round 9 Of The Yelp Dataset Challenge: Our Largest Yet!

The Challenge Dataset:

- **4.1M** reviews and **947K** tips by **1M** users for **144K** businesses
- **1.1M** business attributes, e.g., hours, parking availability, ambience.
- Aggregated check-ins over time for each of the **125K** businesses
- **200,000** pictures from the included businesses

Get the Data

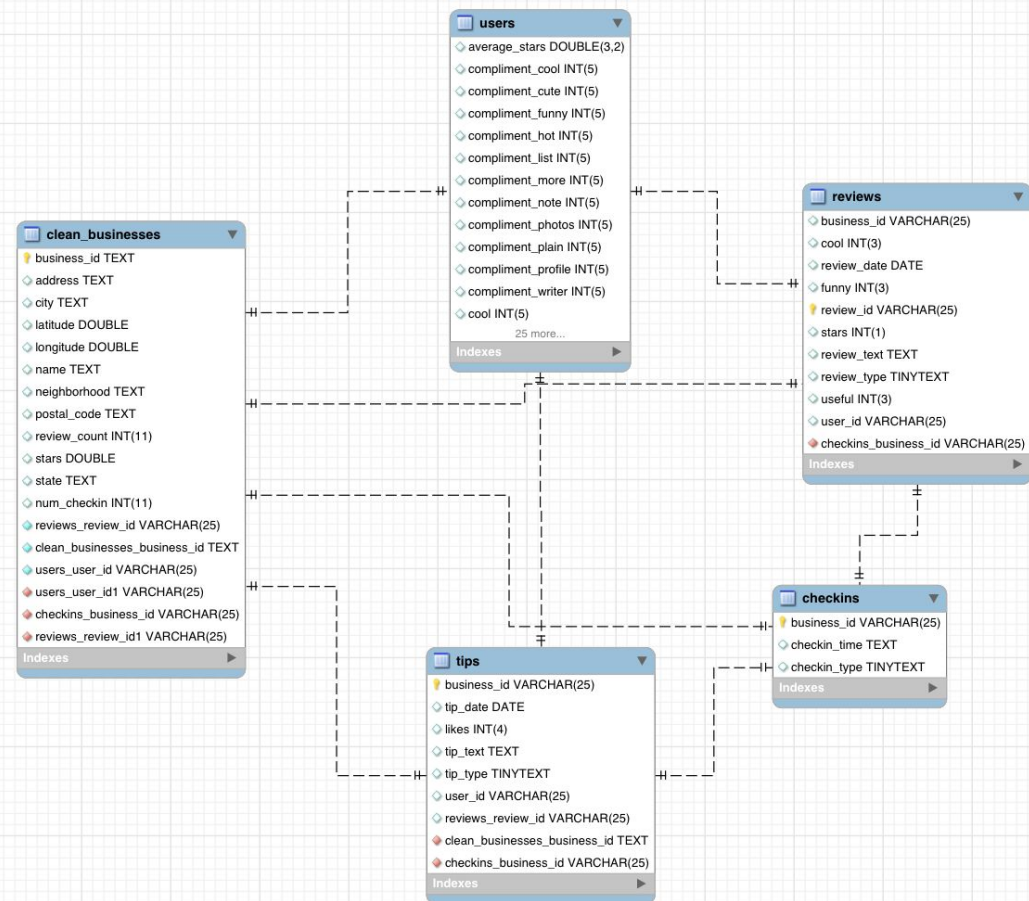
5 Gb text data
json format

Amazon RDS - MySQL Database

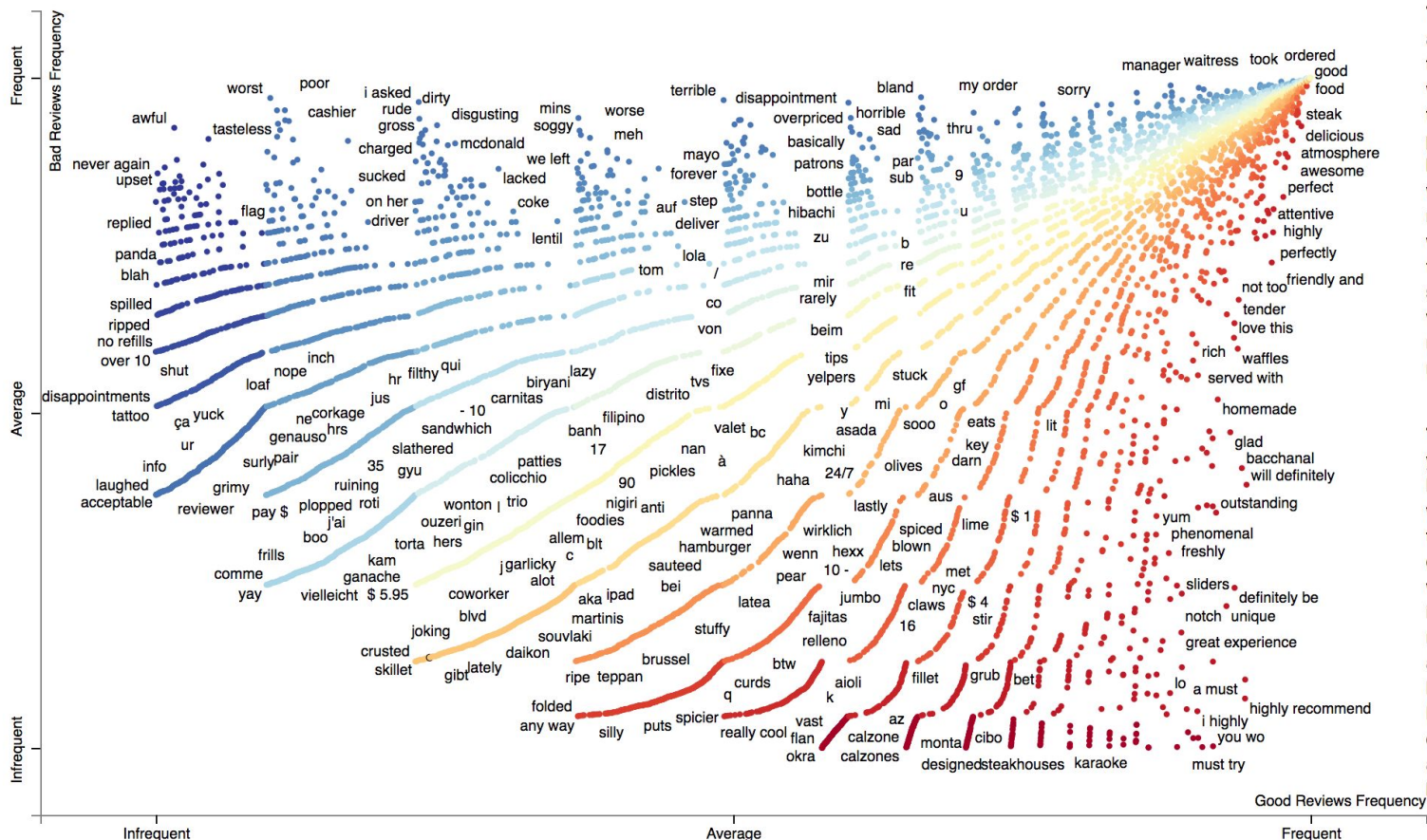
- RDS = Relational Database Service
- Main tables include:
 - **Reviews - 4.1M**
 - **Tips - 947K**
 - **Businesses - 144K**
 - **Checkins - 125K**
 - **Users - 1M**
- The main tables allow for simple subtable creation, which allows for fast requests and data loads.

MySQL Schema Viz

- Primary Keys for all tables:
 - Businesses: business_id
 - Users: user_id
 - Reviews: review_id
 - Tips: user_id
 - Checkins: business_id
- RDS allowed us to create relationships between tables.



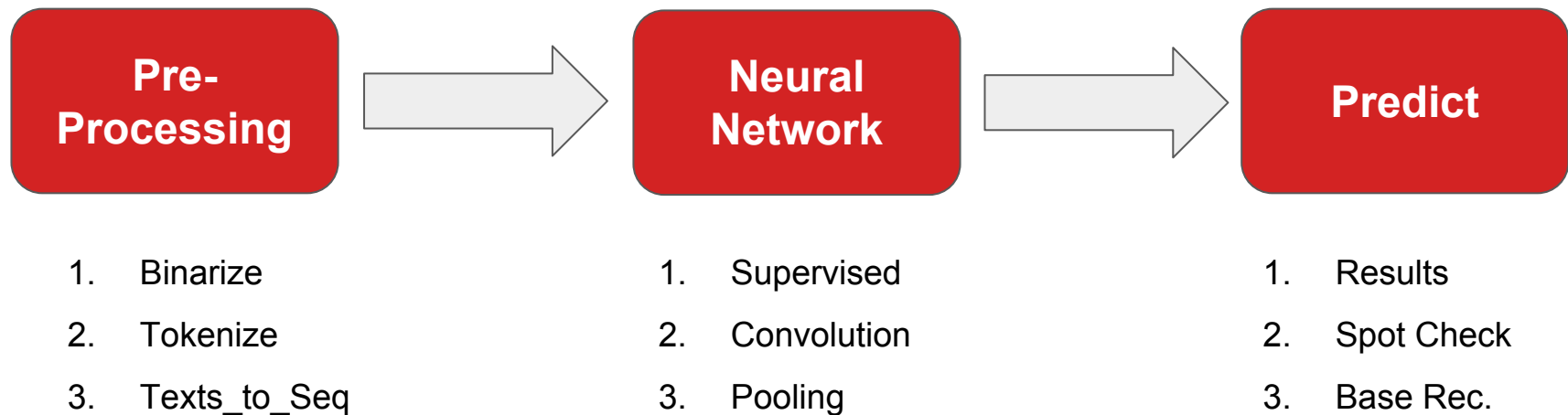
Word Frequency per Review Type



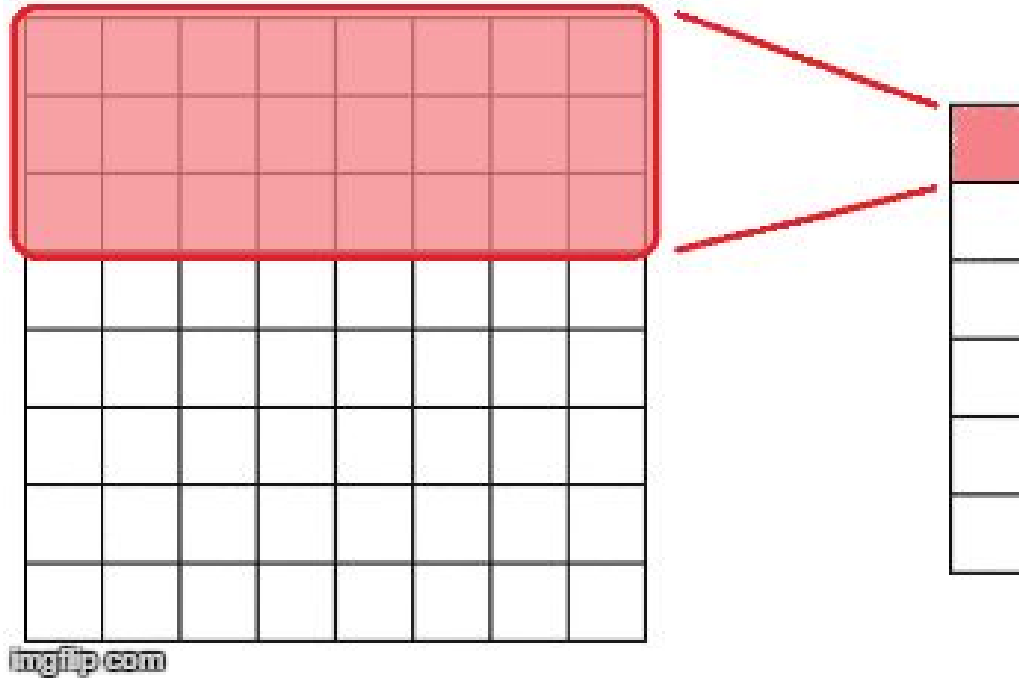
Bad Reviews document count: 1,000; word count: 143,367

Good Reviews document count: 1,000; word count: 112,809

Sentiment Analysis

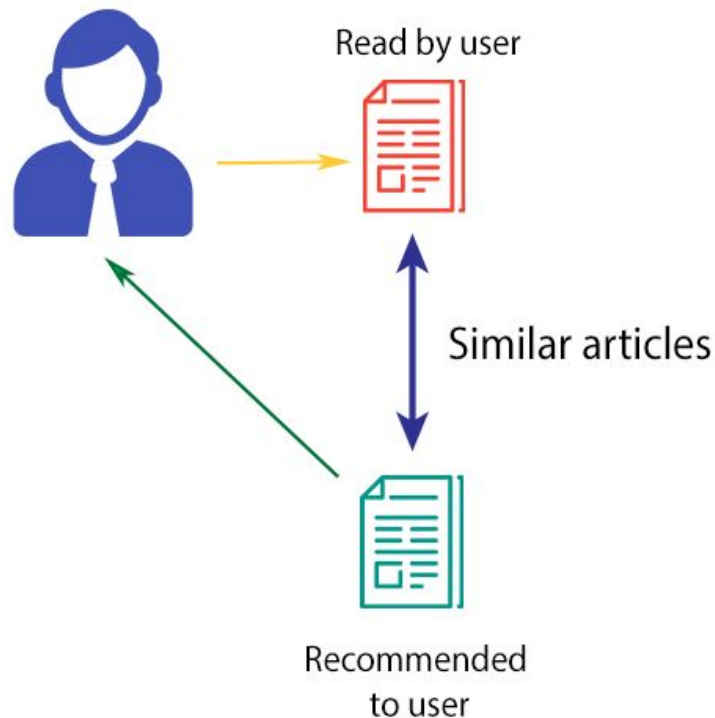


Convolution and Pooling

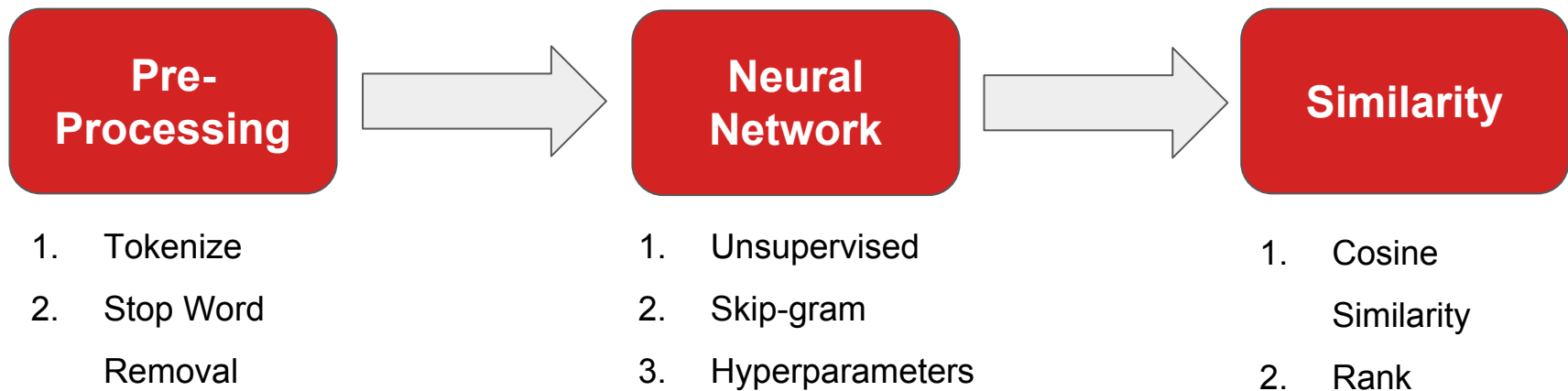


Content-Based Recommendations

- Item-to-item based on user profile
- Keyword soft filtering
- Average Word2Vec



Cosine Similarity using Word2Vec



Word Similarity

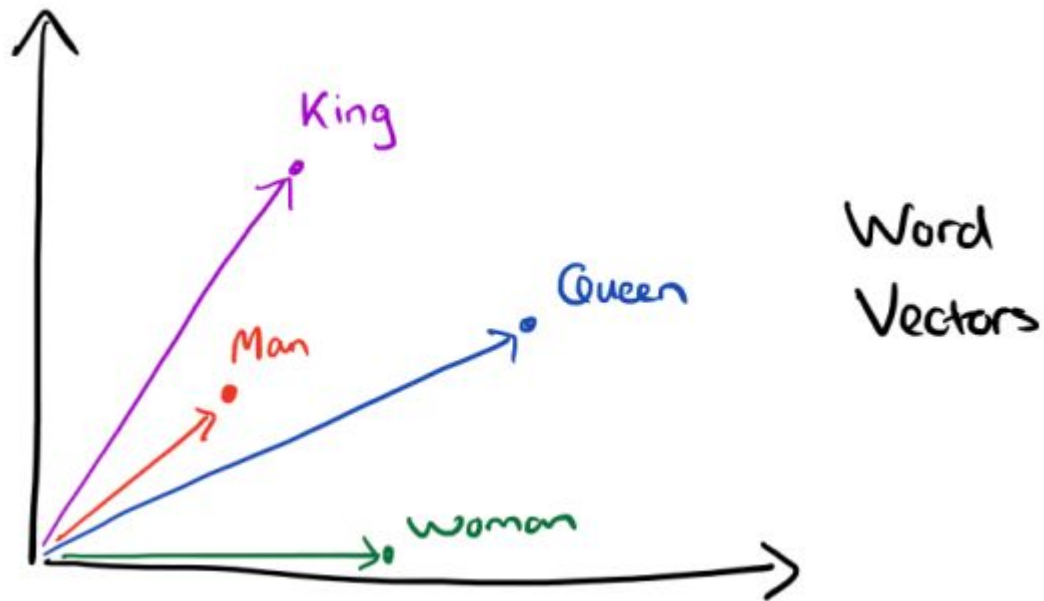
Cmd 28

```
1 #test similarity between words
2 synonyms = model.findSynonyms("chinese", 5)
3 synonyms.show(5)
```

word	similarity
asian	0.7530649328506537
cantonese	0.7275205535026946
filipino	0.7259573442238313
shanghainese	0.6988964525834206
korean	0.6691107471966923

Command took 0.48 seconds -- by skickham@gmail.com at 6/15/2017, 2:59:40 PM on My Cluster

Word Similarity



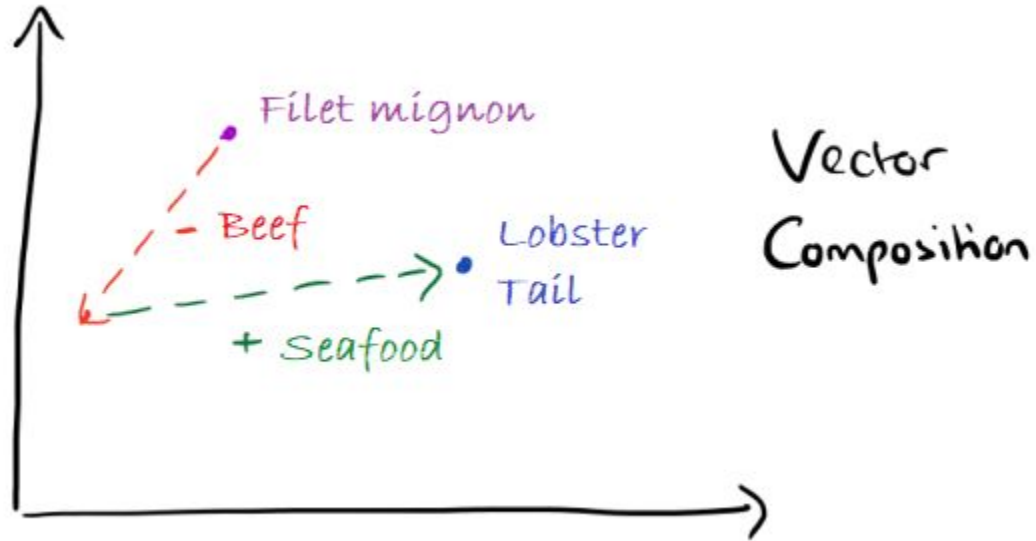
Word Algebra



Word Algebra: Food



Word Algebra: Food



Word Algebra: Examples

'filet_mignon' - 'beef' + 'seafood' =

'fork' + 'soup' =

'drink' + 'barley' =

```
In [*]: word_algebra(add=[u'filet_mignon', u'seafood'], subtract=[u'beef'])  
word_algebra(add=[u'fork', u'soup'])  
word_algebra(add=['drink', 'barley'])
```

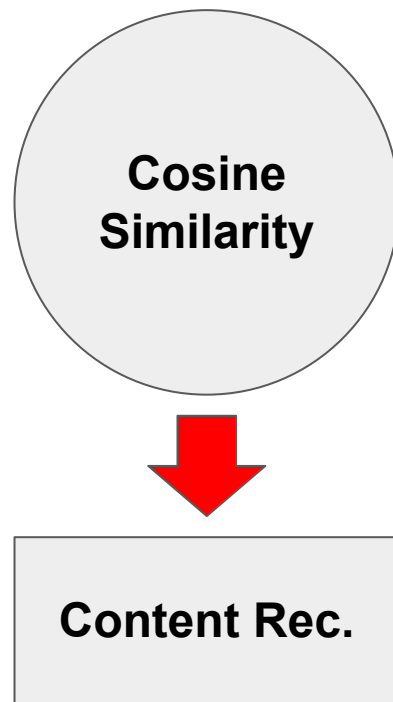
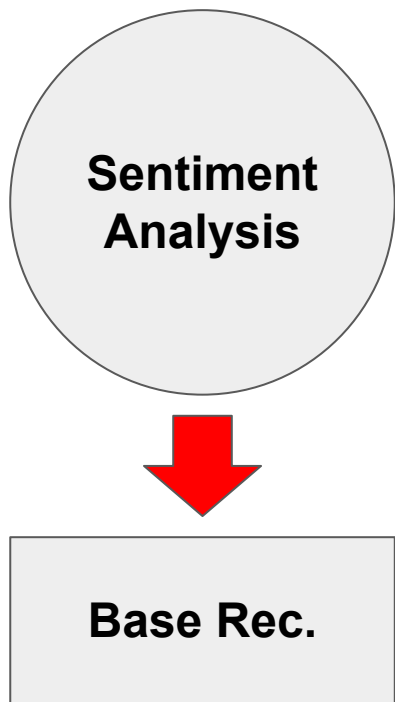
Word Algebra: Examples

'filet_mignon' - 'beef' + 'seafood' = lobster tail

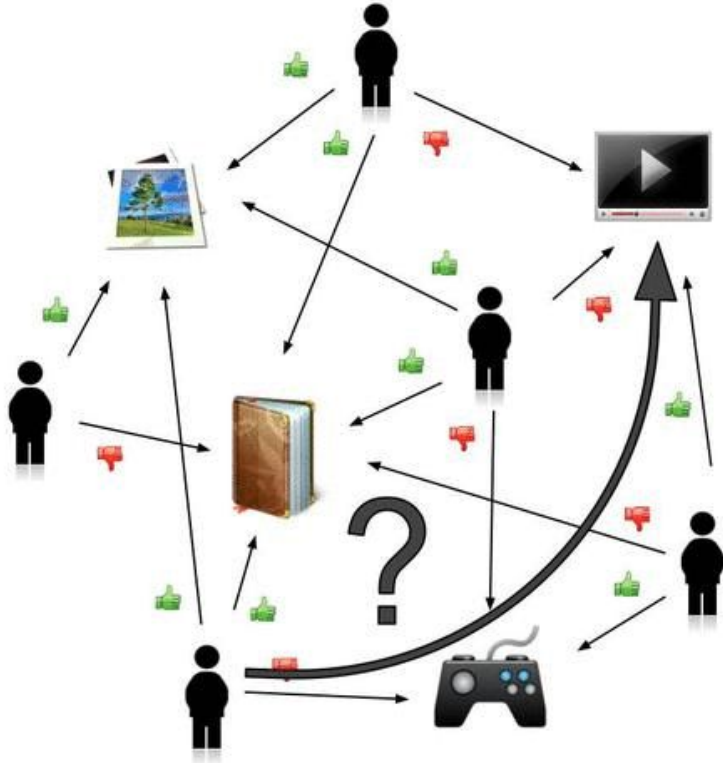
'fork' + 'soup' = spoon


























'drink' + 'barley' = beer

NLP Summary



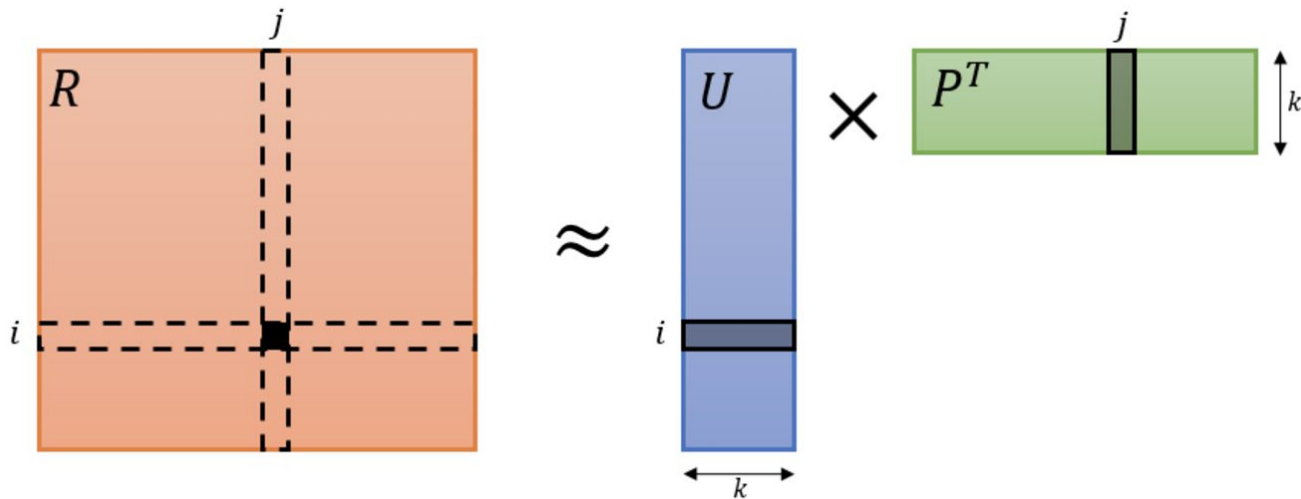
Collaborative Filtering 1



Collaborative Filtering 2

Matrix
Factorization



Cost Function

$$J = ||R - U \times P^T||_2 + \lambda (||U||_2 + ||P||_2)$$

Yelp Social Network 1

yelp_academic_dataset_user.json

nodeDF

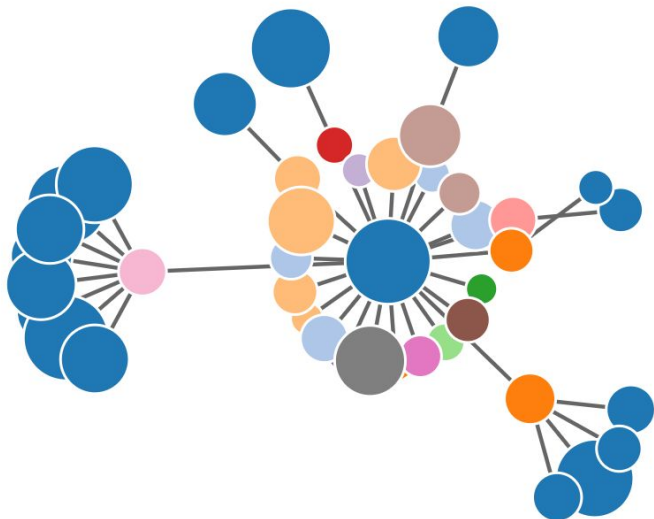
```
{
  "user_id": "encrypted user id",
  "name": "first name",
  "review_count": "number of reviews",
  "yelping_since": "date formatted like \"2009-12-19\"",
  "friends": ["an array of encrypted ids of friends"],
  "useful": "number of useful votes sent by the user",
  "funny": "number of funny votes sent by the user",
  "cool": "number of cool votes sent by the user",
  "fans": "number of fans the user has",
}
```



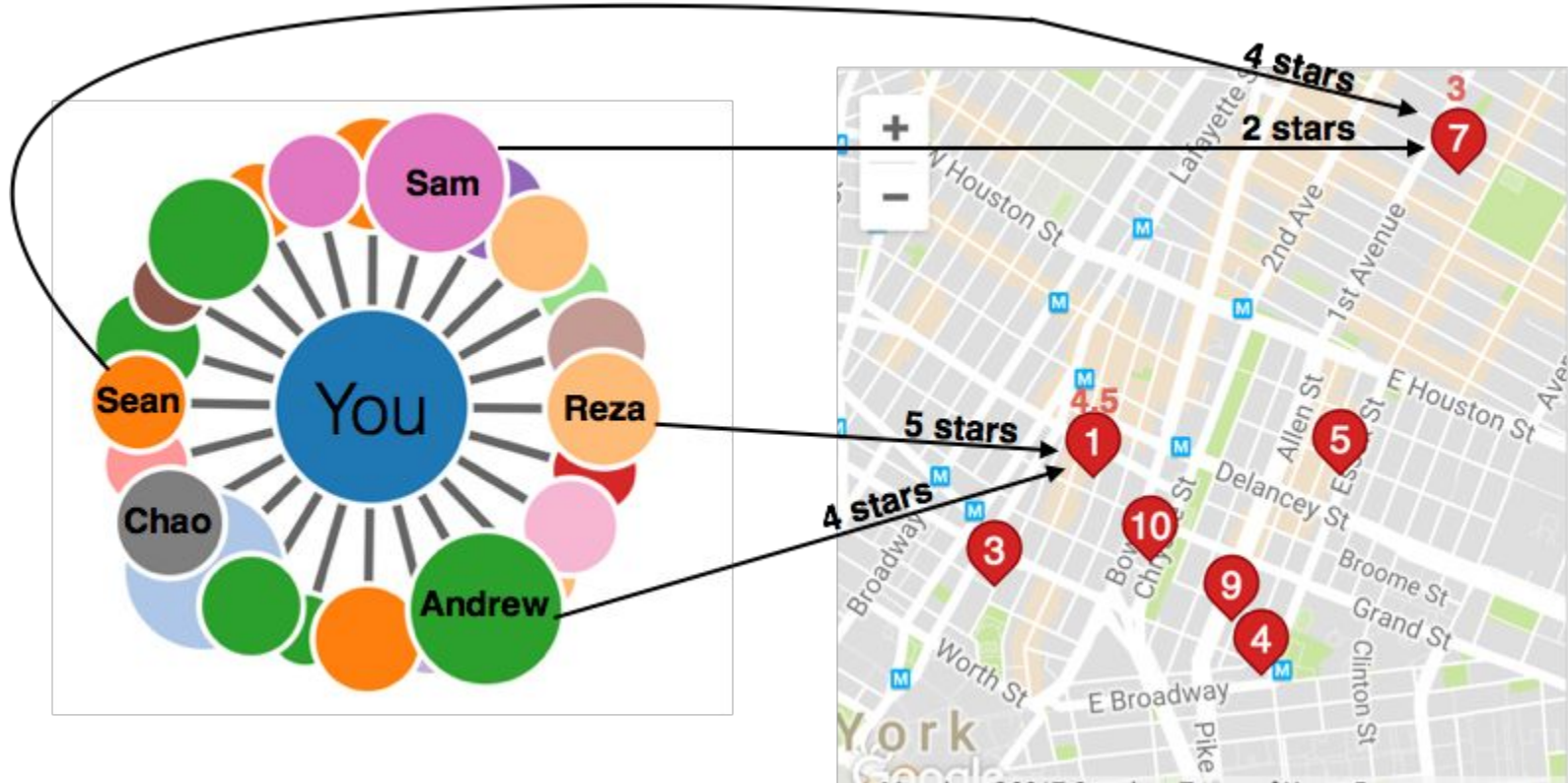
friends.show(50)

edgeDF

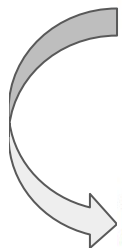
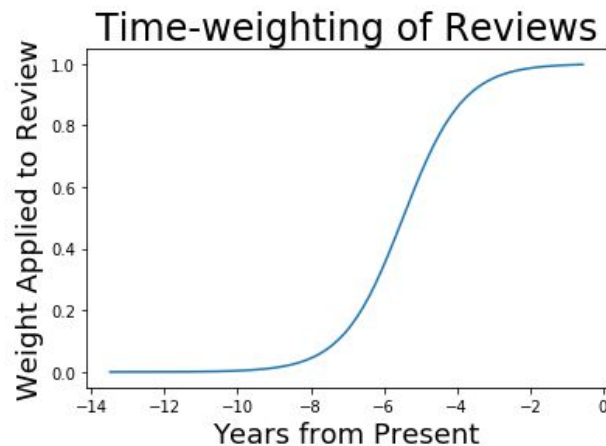
u0	u1
1003979	1024458
297574	1003979
970343	1003979
320810	1003979
188106	1003979
492231	1003979
352654	1003979



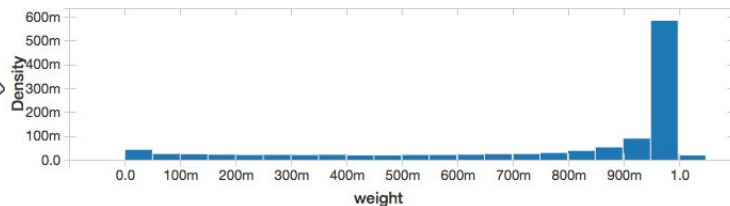
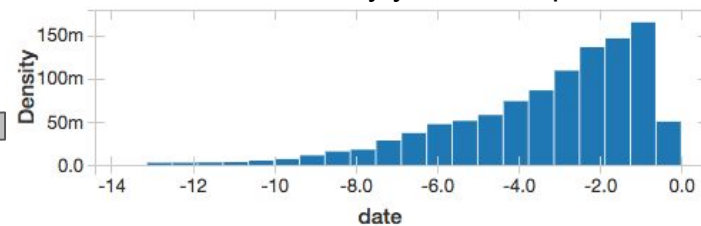
Yelp Social Network 2



Review-level adjustments



Reviews binned by years from present



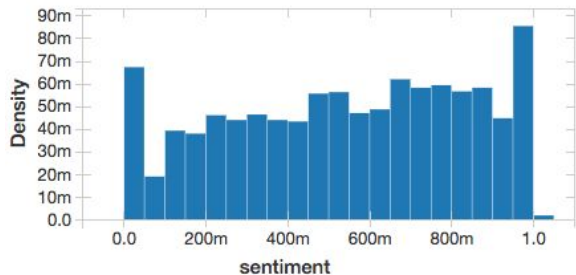
Newest 60% of reviews unaffected by time-weighting

How relevant is a 10 year old review?

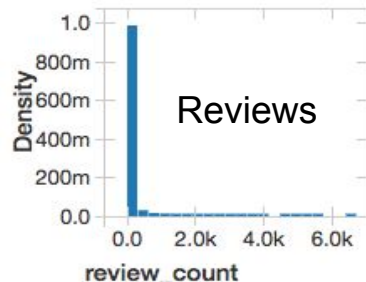
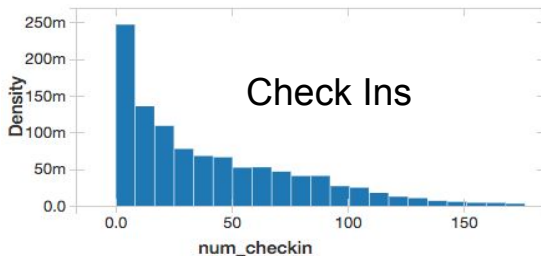
Business-level adjustments

Log transformations for popularity measures

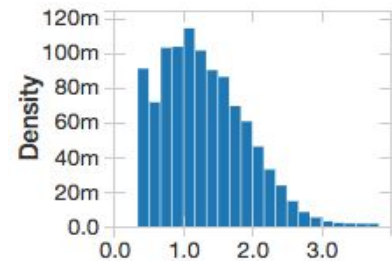
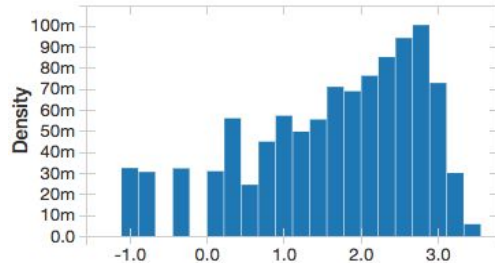
Multiplicative factor based
on NLP sentiment



$$\xi = \left(\frac{3}{4} + \frac{\text{sentiment}}{4} \right) \cdot \sum_{\text{all reviews}} f(t, \text{rating})$$



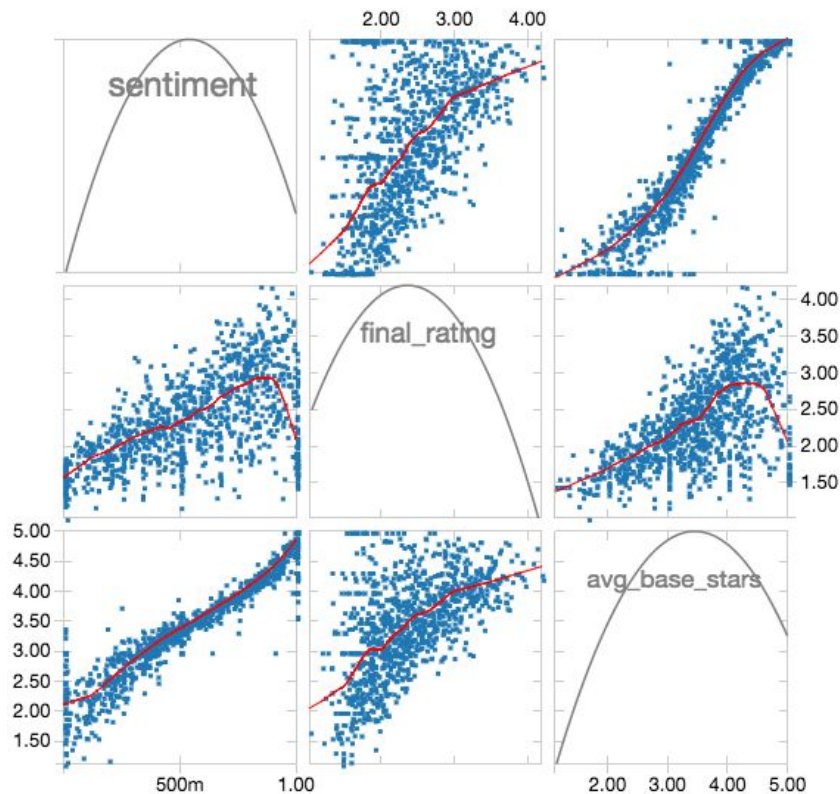
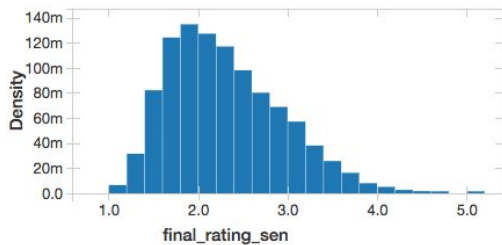
Log



final score = $\log_{10}(\text{review count}^{\xi} \cdot g(\text{check-in count}))$

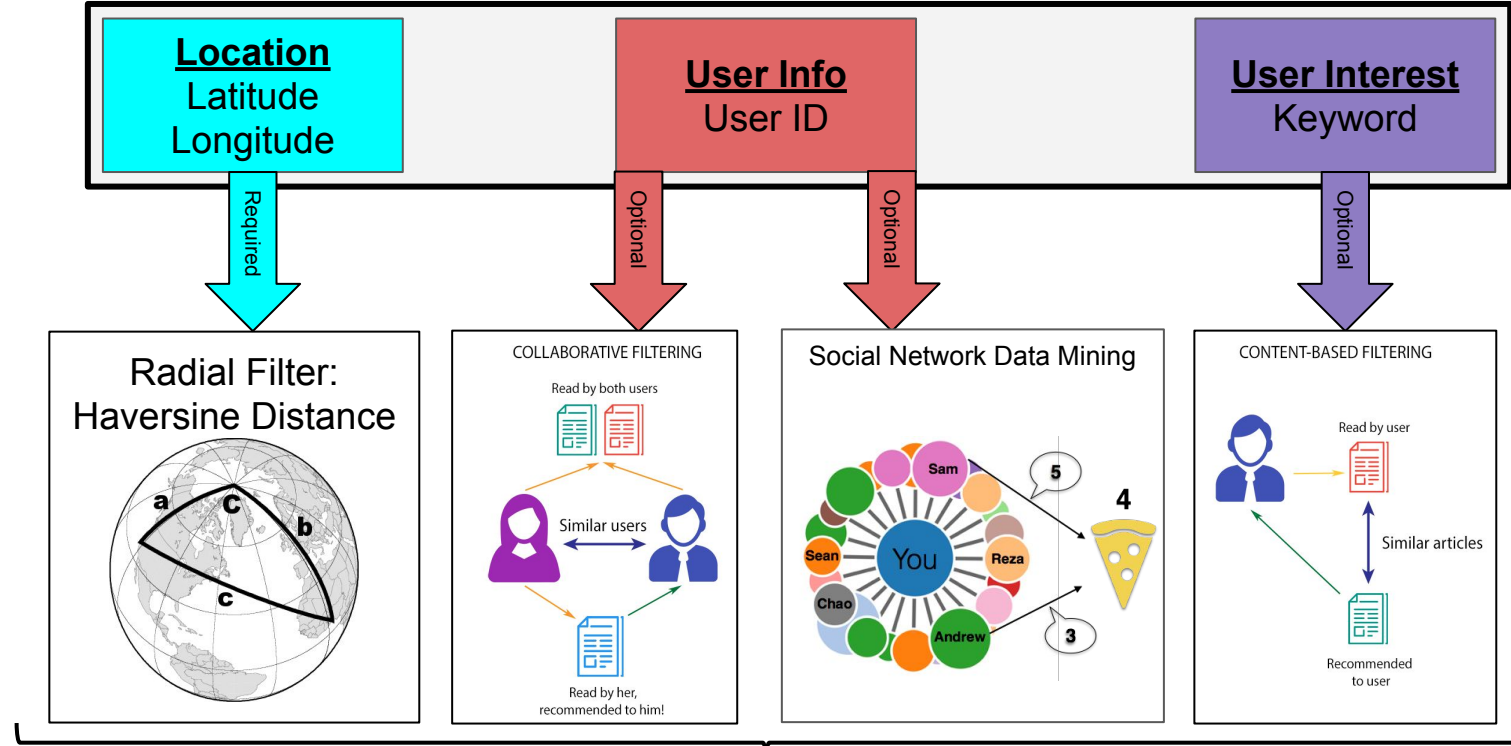
Final Location-only Ratings

Stars	Percentile
★★★★	99th (top 40)
★★★★☆	95th (top 300)
★★★	83rd
★★★☆☆	60th
★★	31st
★★☆	5th



Recommendation Engine

User Requests



Dynamic generation of up to 40 recommendations within 20s

Data Pipeline: 1. Flask

Map Satellite

Yelp Capstone Project

User ID
1003979

Keyword
tacos

Longitude
-115.150834

Latitude
36.1145

Radius (in meters)
2500

SUBMIT

Key:

- Location-Based
- Social
- Keyword
- Collaborative Filtering

Map data ©2017 Google 2 km Terms of Use Report a map error

User ID

Keyword

Longitude

Latitude

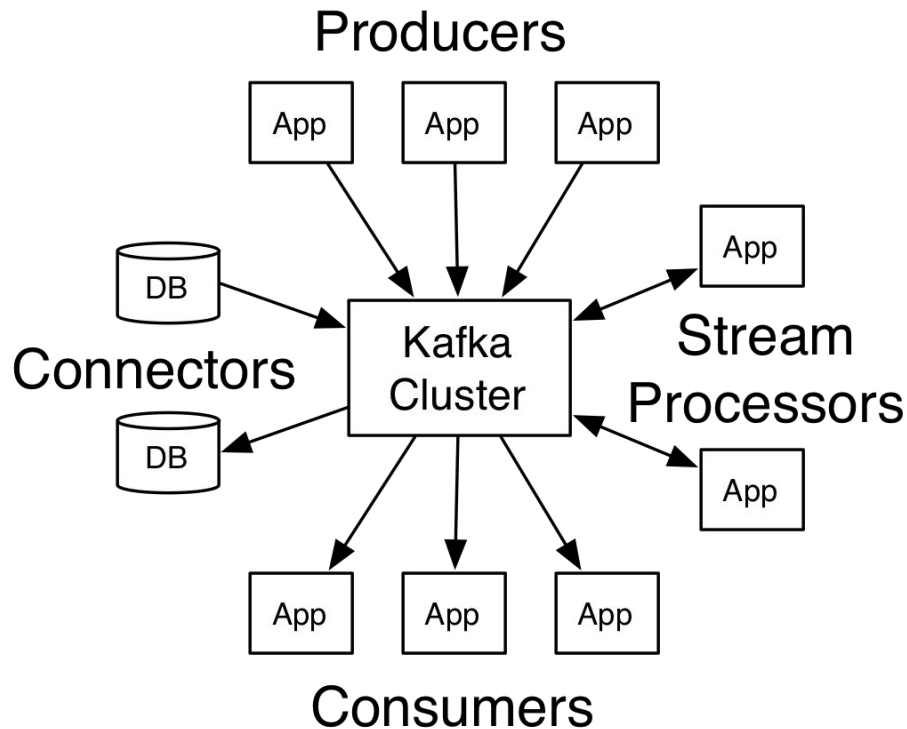
Search Radius



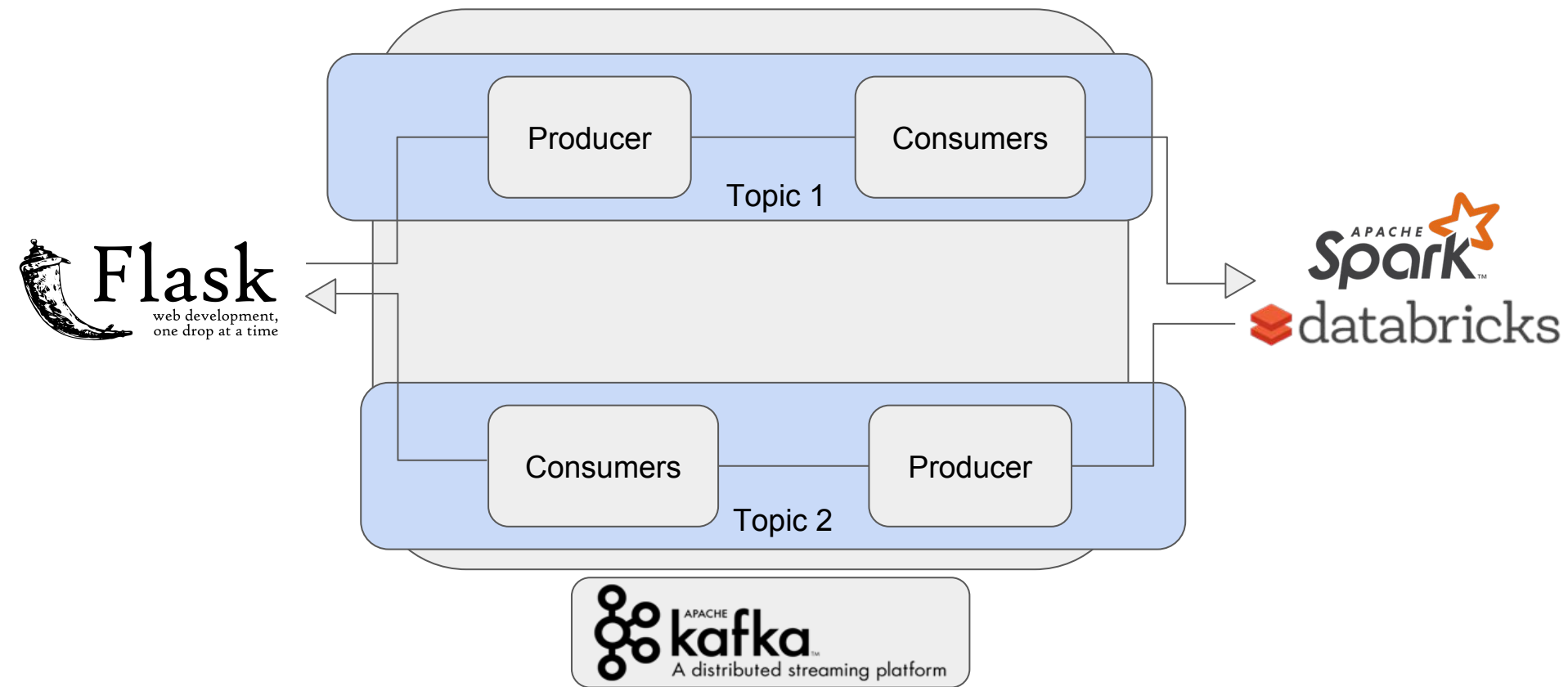
Data Pipeline: Kafka

Components:

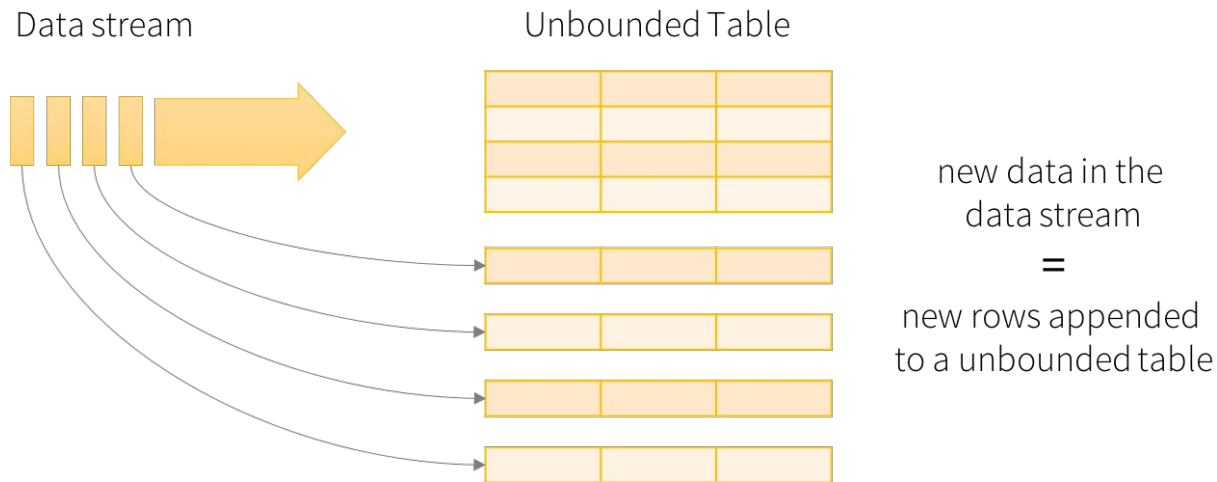
1. Brokers
2. Topics
3. Producers
4. Consumers



Data Pipeline: 2. Kafka



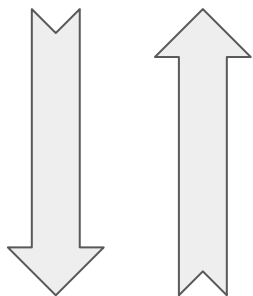
Data Pipeline: 3. Spark



Data stream as an unbounded table

Spark's Structured Streaming

Data Pipeline: Overall



```

# Read the locations.json file and print its contents
with open('locations.json') as f:
    locations = json.load(f)
    print(locations)

# Send the locations to Kafka
producer = kafka.KafkaProducer(bootstrap_servers='localhost:9092')
for location in locations:
    producer.send('locations', json.dumps(location).encode('utf-8'))
producer.flush()

# Read the locations from Kafka
consumer = kafka.KafkaConsumer('locations', bootstrap_servers='localhost:9092')
consumer.subscribe()
while True:
    message = consumer.poll(1)
    if message:
        print(message.value())

```

Demo

Lessons Learned and Future Work

Lessons

- Recommendation Algorithms
- Dealing with Big Data
- Working in a Startup-like Environment

Future Work

- Front End and Streaming
 - User experience enhancement
 - Restructure the data streams for better stability
- Back End Algorithms
 - Speed up current code
 - The very fun part has just begun

Acknowledgement

Guidance and Discussion from

- Shu Yan
- Yvonne Lau
- Zeyu Zhang

Inspiration from

Chuan Sun and Aiko Liu

Thanks!

