

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
 - Rec Engine
 - App Workflow (Flask - Kafka - Rec Engine - Kafka - Flask)
- Live Demo
- Lessons Learned and Future Work
- Acknowledgement

Background



EDA



Recommendation
Algorithms

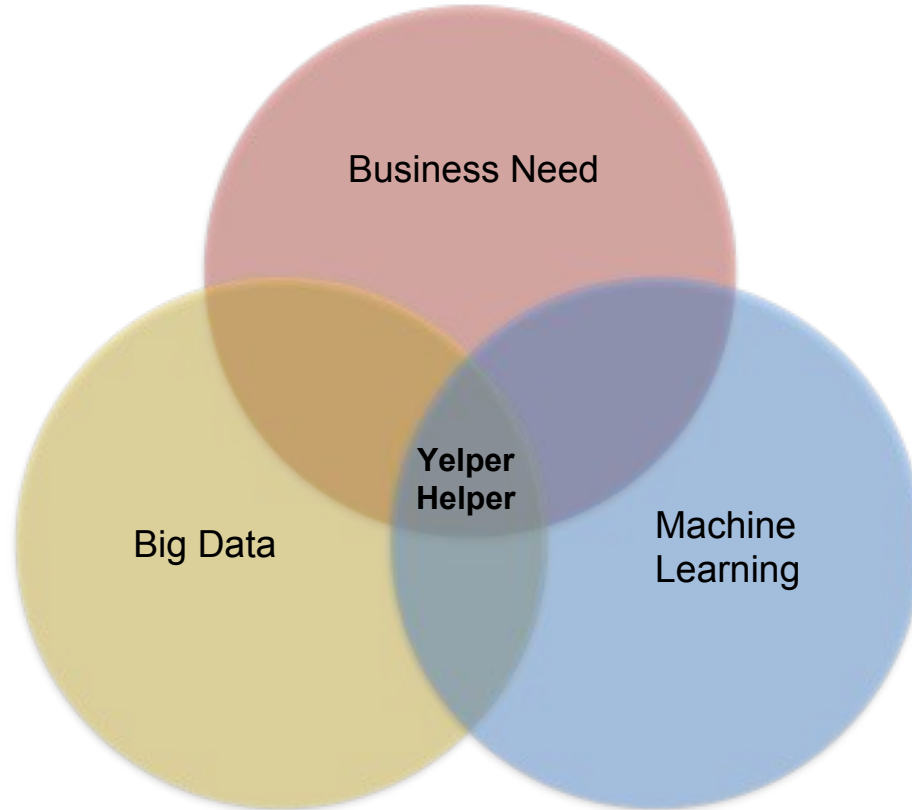


App Pipeline

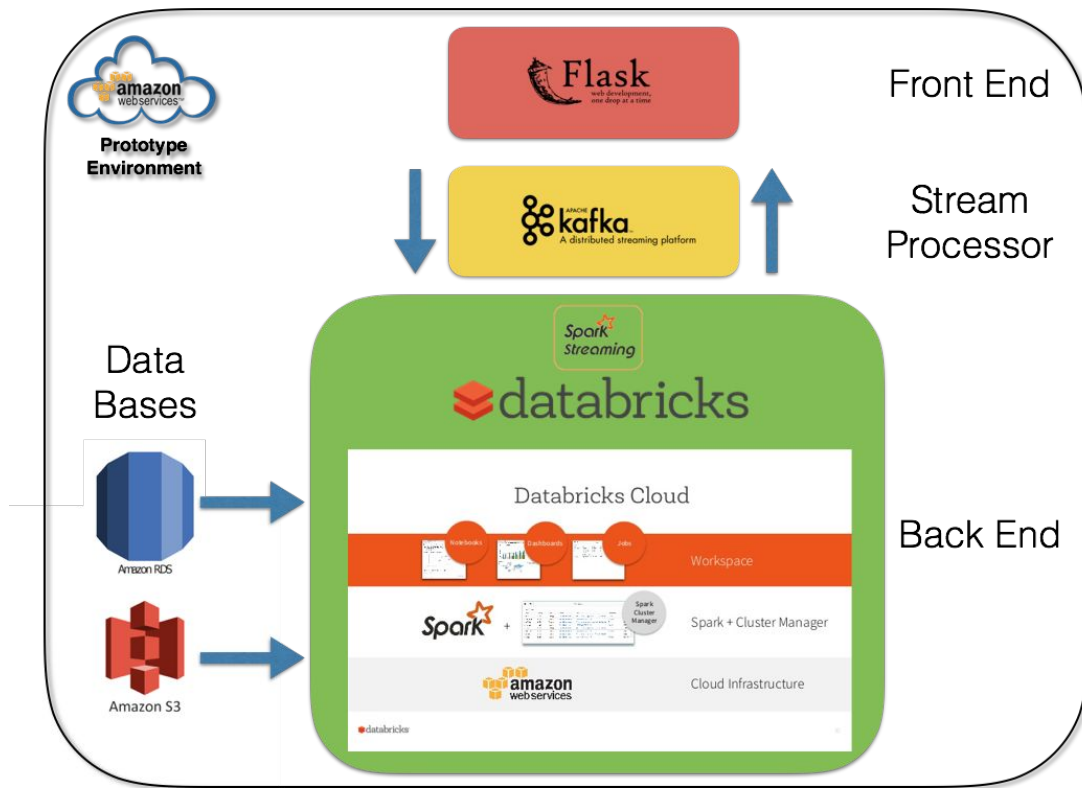


Demo

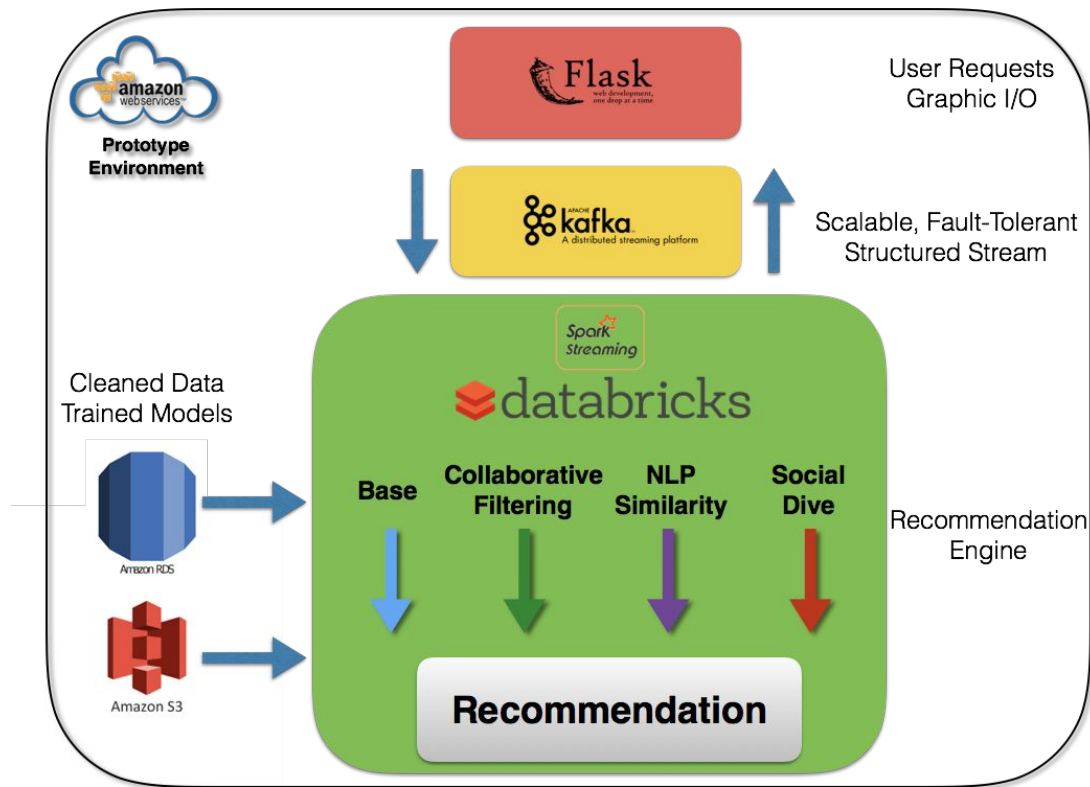
Motivation



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

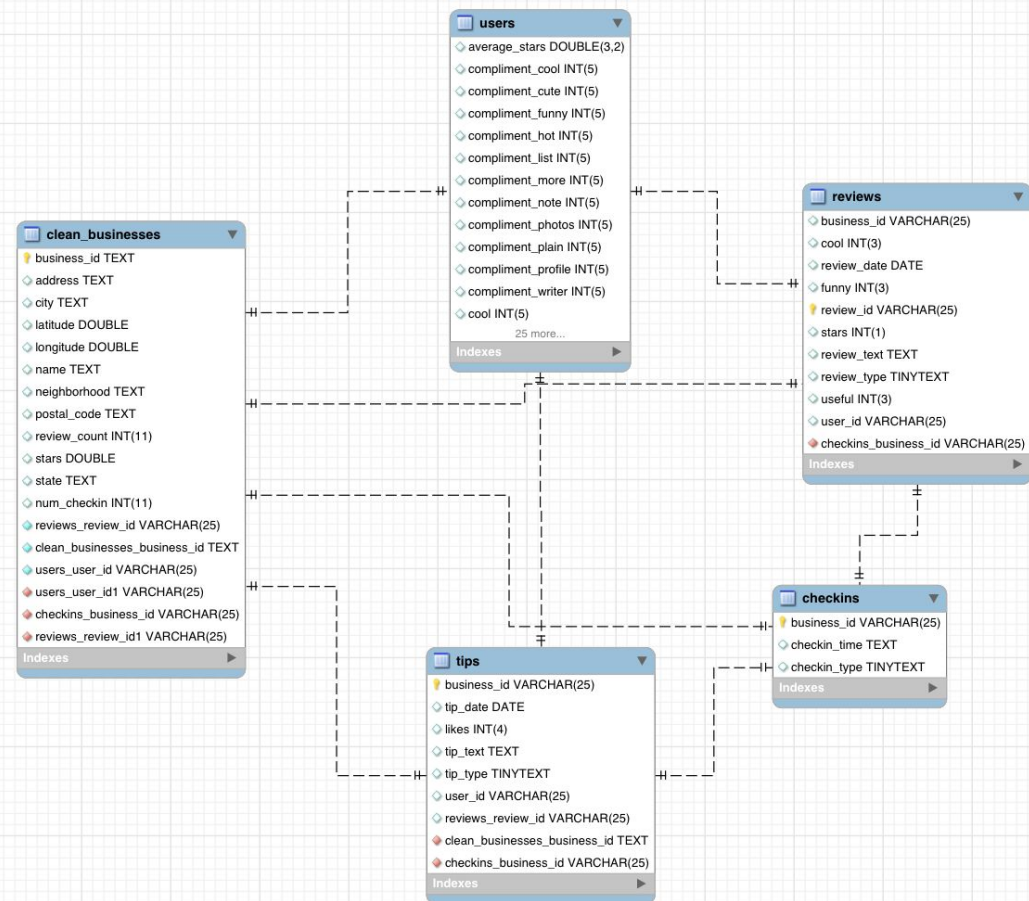
5 Gb text data
json format

Amazon RDS - MySQL Database

- RDS = Relational Database Service
- Main tables include:
 - **Reviews**
 - **Tips**
 - **Businesses**
 - **Check-ins**
 - **Users**
- 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.



Read/Write Speed

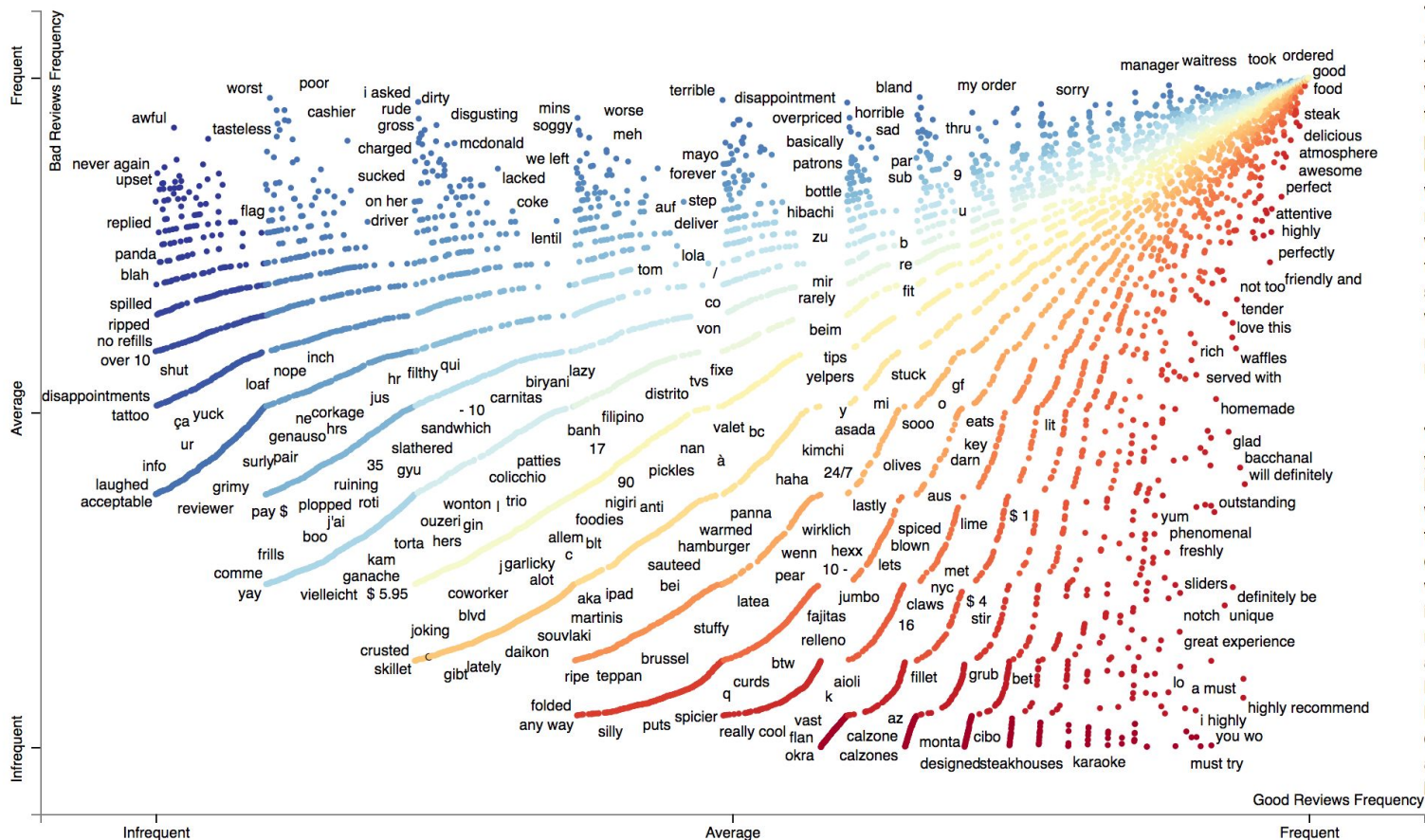
Monitoring



	CURRENT VALUE	THRESHOLD	LAST HOUR
CPU	1%	<div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>
Memory	35.2 MB	<div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>
Storage	11,500 MB	<div><div></div><div></div></div>	<div><div></div><div></div><div></div></div>

	CURRENT VALUE	LAST HOUR
Read IOPS	0/sec	<div><div></div><div></div><div></div></div>
Write IOPS	0.367/sec	<div><div></div><div></div><div></div></div>
Swap Usage	30 MB	<div><div></div><div></div><div></div></div>

Word Frequency per Review Type



Top Bad Reviews

awful
tasteless
worst
the worst
poor
burnt
we asked
not even
will never
frozen
sucks
very disap
managem
never aga

Top Good Reviews

you wo
highly recommend
would definitely
to die
die for
well worth
i highly
must try
a must
loved it
hands down
definitely recommend
are amazing
boba

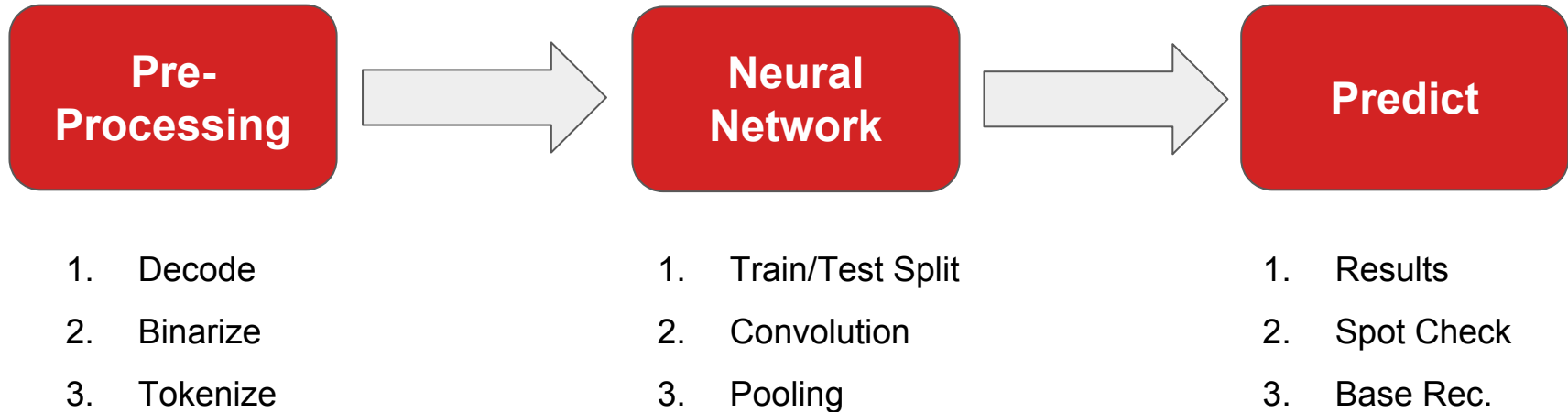
Characteristic

ramen
bacchanal
waitress
tacos
yelp
tasted
waiter
fries
wagyu
frites
flavorful
sashimi
overcooked
ribeye
groupon
ayce
craftsteak
filet
flavorless
guac
soggy
buffets
poutine
brisket
crispy
sushi
calamari
hibachi
appetizer
overpriced

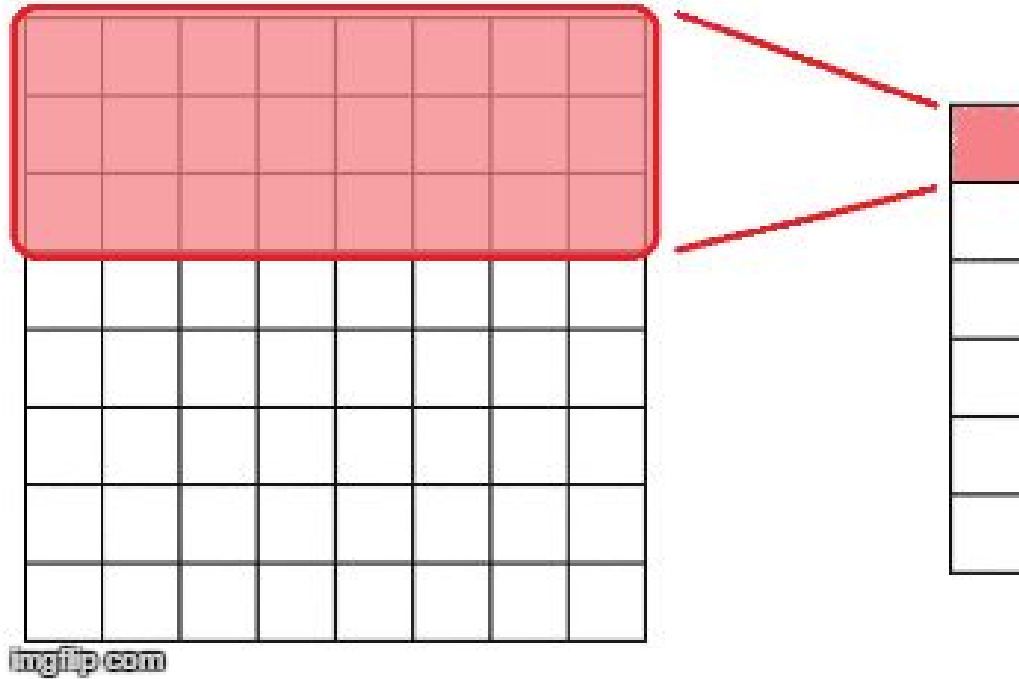
Bad Reviews document count: 1,000; word count: 143,367

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

Sentiment Analysis

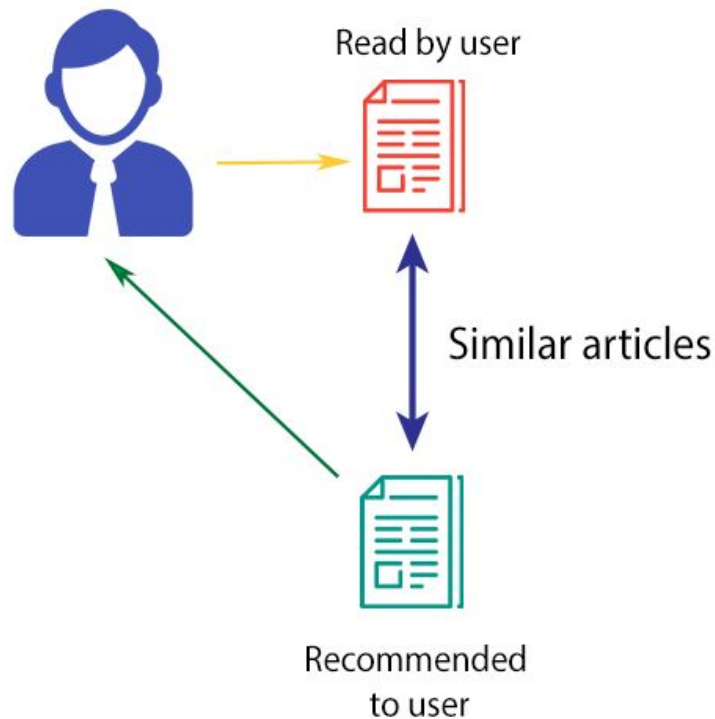


Convolution and Pooling

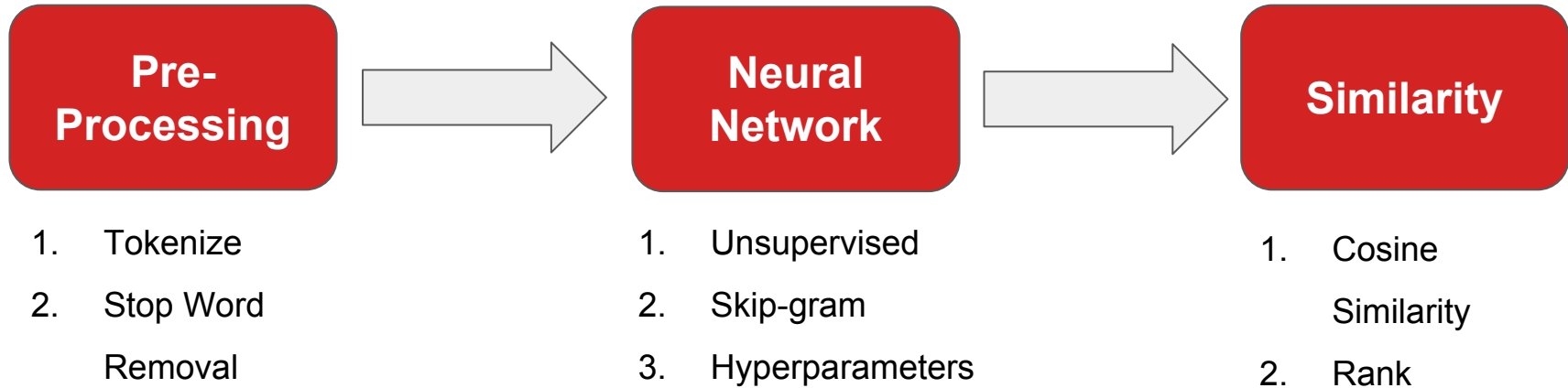


Content-Based Recommendations

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



Cosine Similarity using Word2Vec



Word Similarity

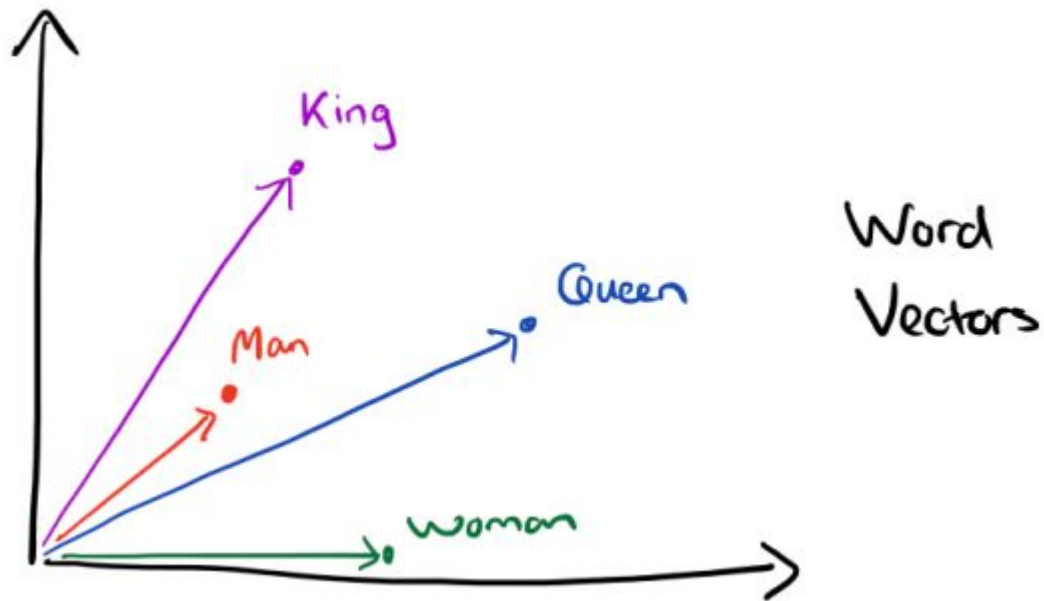
Cmd 28

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

word	similarity
insane	0.585583245392664
mean	0.570725045955314
weekends	0.5407901942365788
cause	0.5246057760341053
ridiculously	0.5229965692070889

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

Word Algebra



Word Algebra



Word Algebra



Word Algebra

Word Algebra (visual?) (latex equation)

```
word_algebra(add=[u'filet_mignon', u'seafood'], subtract=[u'beef'])
```

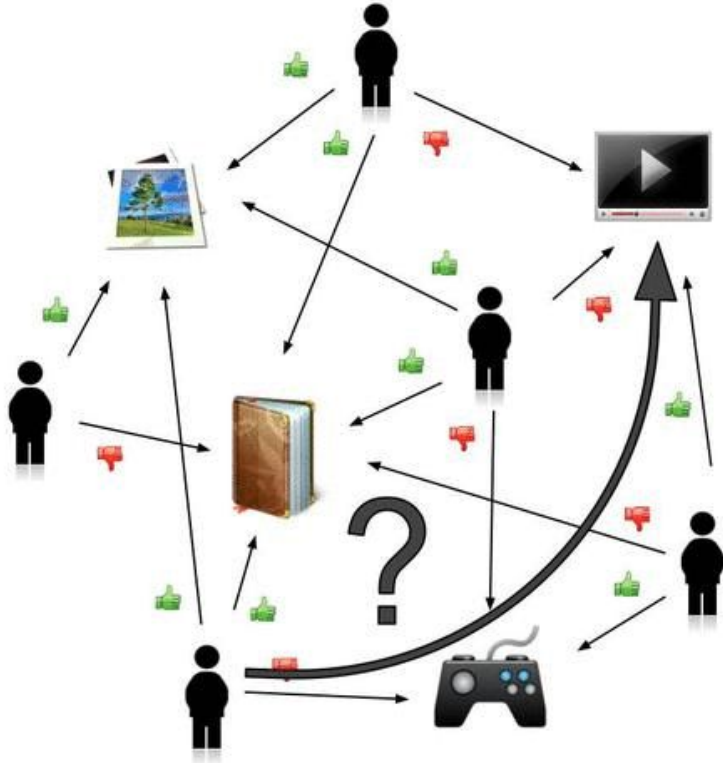
```
word_algebra(add=[u'fork', u'soup'])
```


























```
word_algebra(add=['drink', 'barley'])
```

NLP Summary (easy visual/flow)

- Sentiment Analysis as feature in Base Recommendation
- Cosine Similarity for Content-Based Recommendation

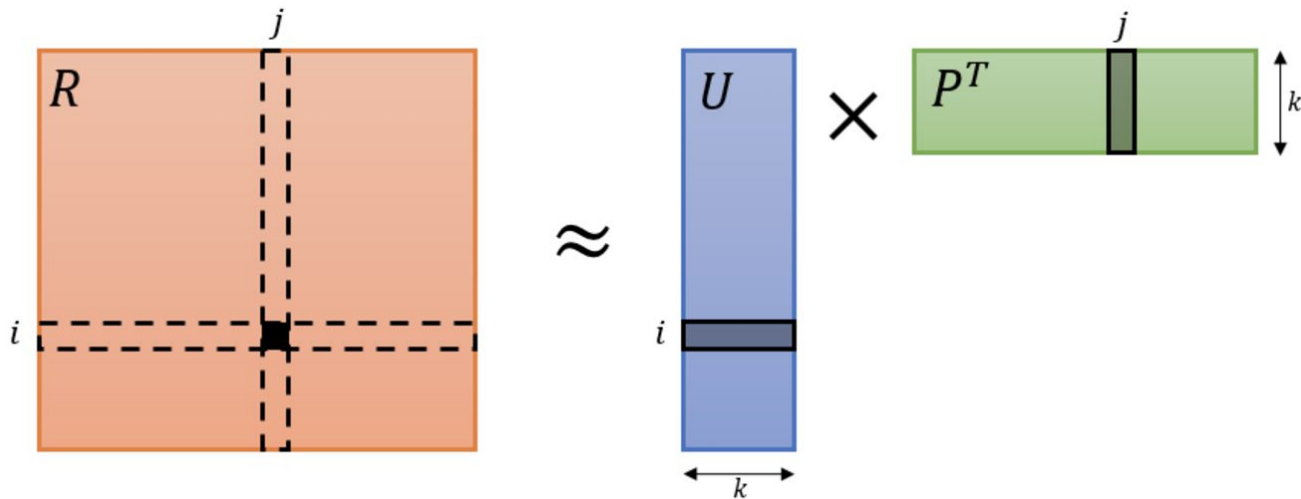
Collaborative Filtering 1



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

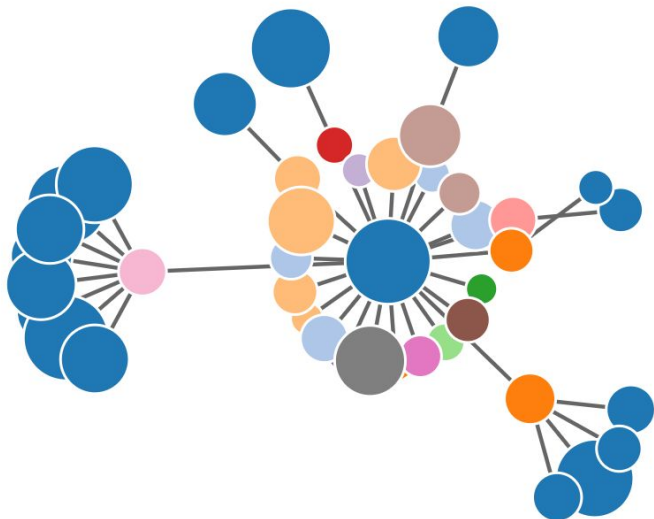
```
{
  "user_id": "encrypted user id",
  "name": "first name",
  "review_count": "number of reviews",
  "yelping_since": "date formatted like \"2009-12-19\"",
  "friends": ["an array of encrypted ids of friends"],
  "useful": "number of useful votes sent by the user",
  "funny": "number of funny votes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



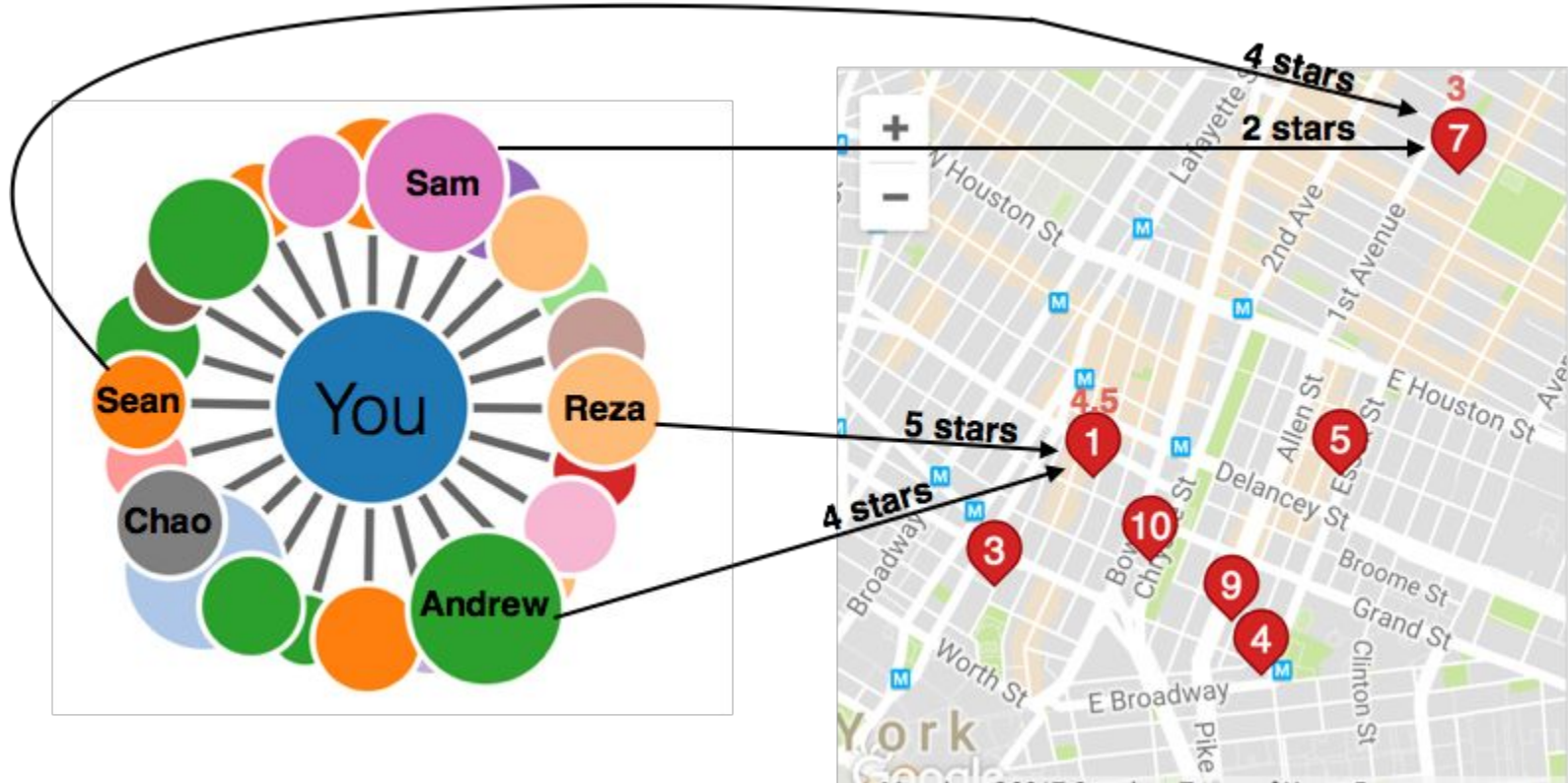
friends.show(50)

edgeDF

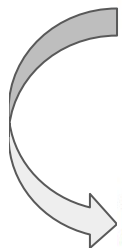
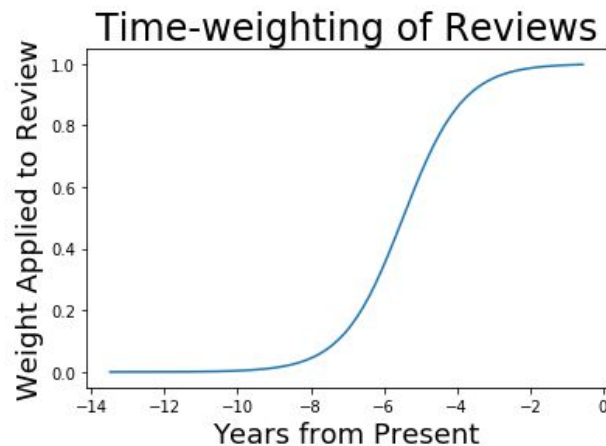
u0	u1
1003979	1024458
297574	1003979
970343	1003979
320810	1003979
188106	1003979
492231	1003979
352654	1003979



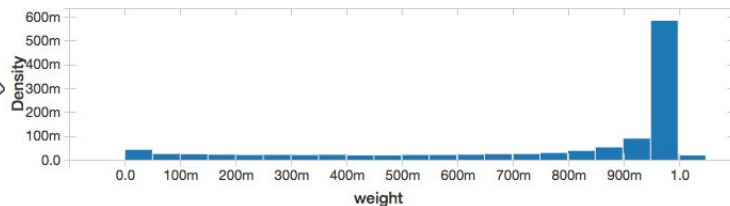
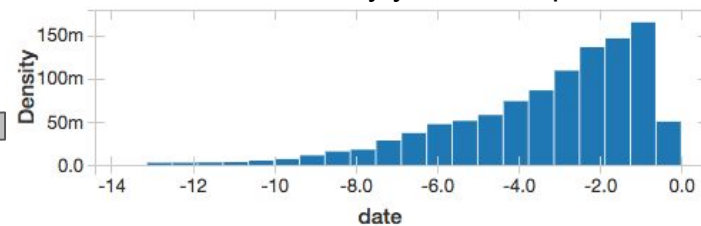
Yelp Social Network 2



Review-level adjustments



Reviews binned by years from present



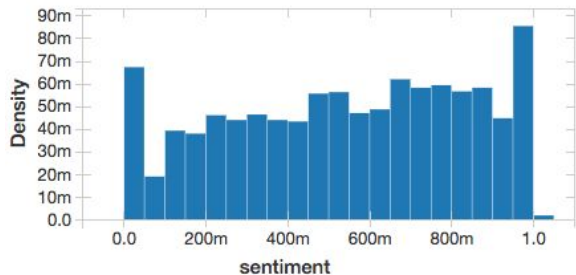
Newest 60% of reviews unaffected by time-weighting

How relevant is a 10 year old review?

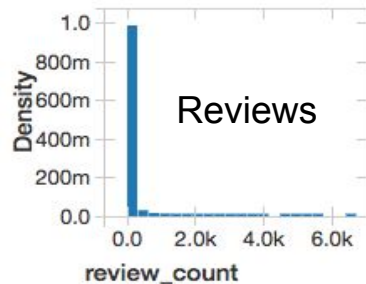
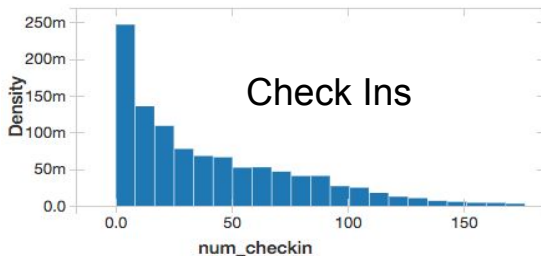
Business-level adjustments

Log transformations for popularity measures

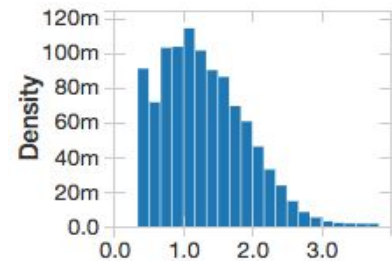
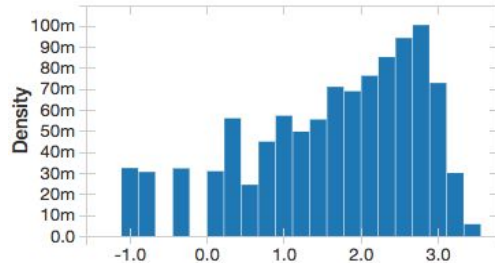
Multiplicative factor based on NLP sentiment



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



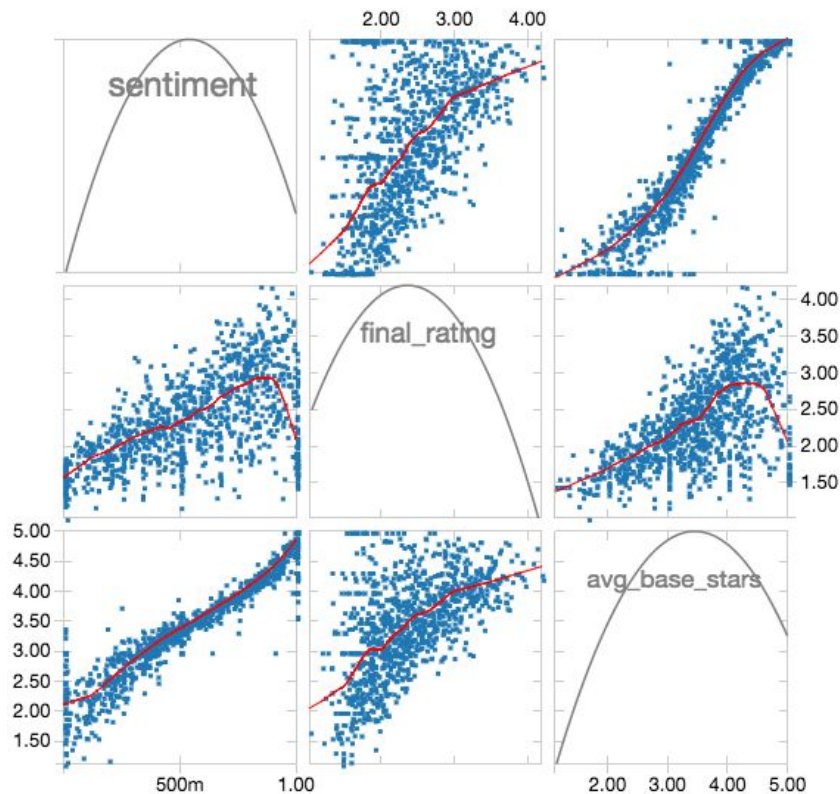
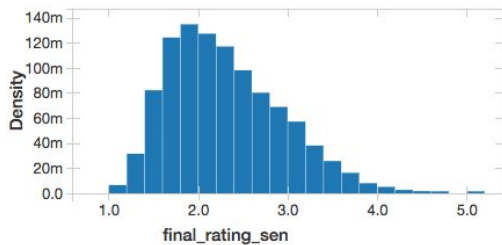
Log



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

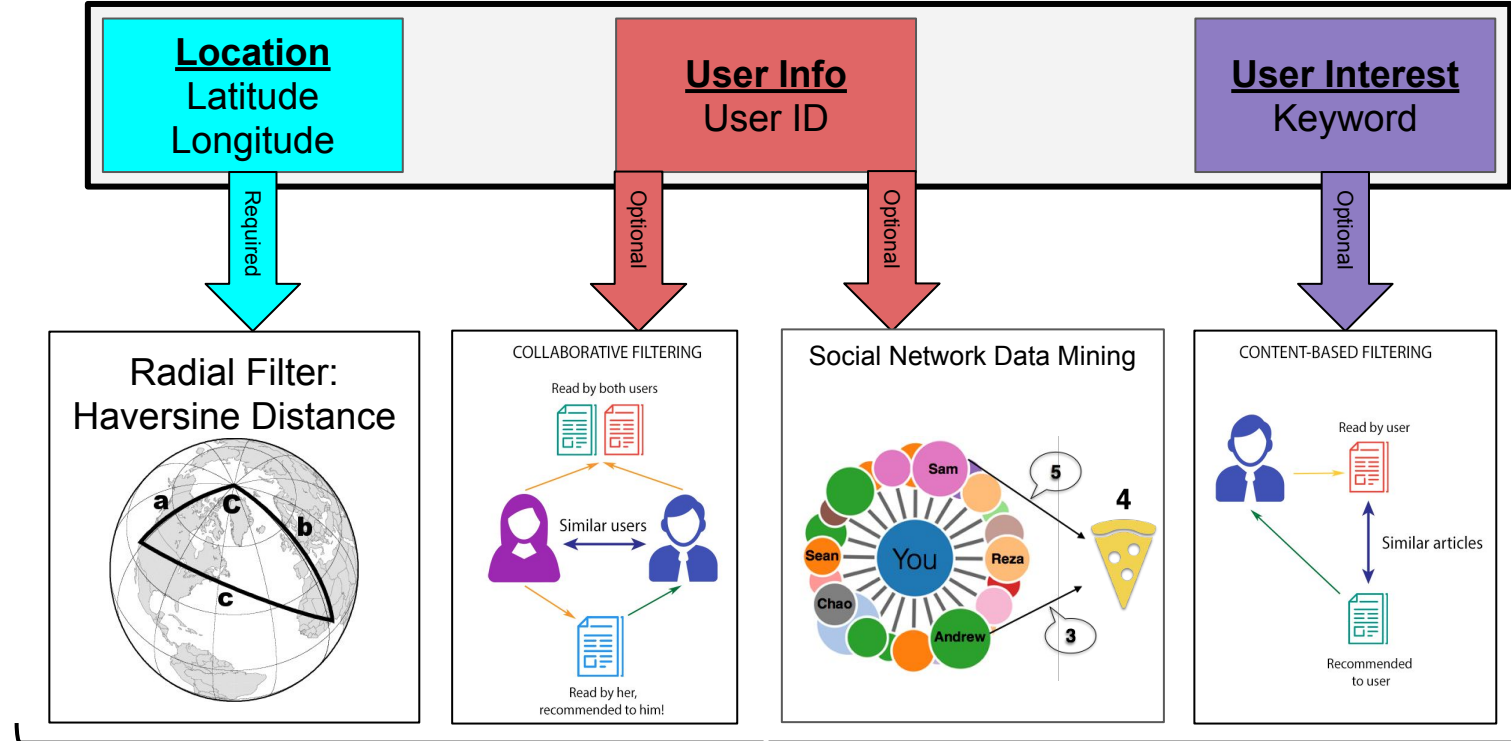
Final Location-only Ratings

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

User ID

Keyword

Longitude

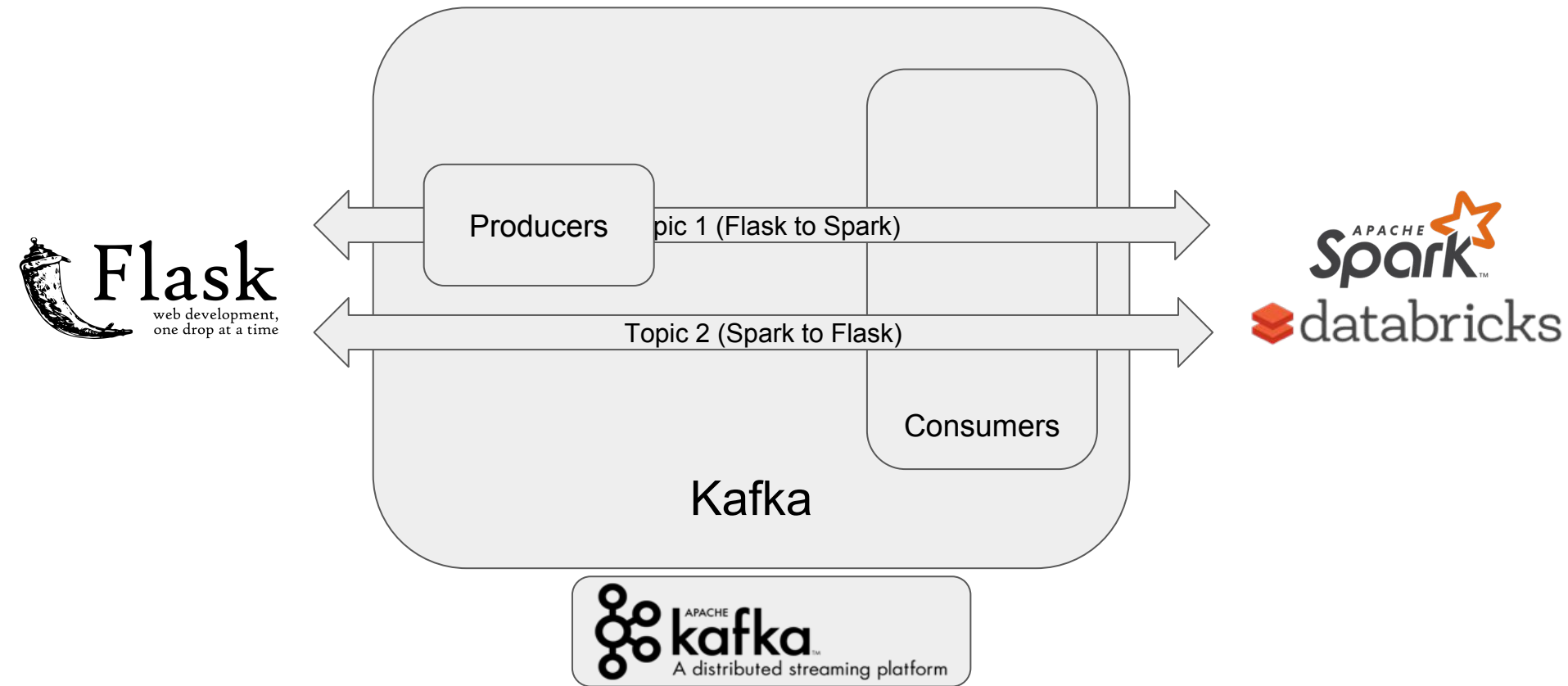
Latitude

Search Radius

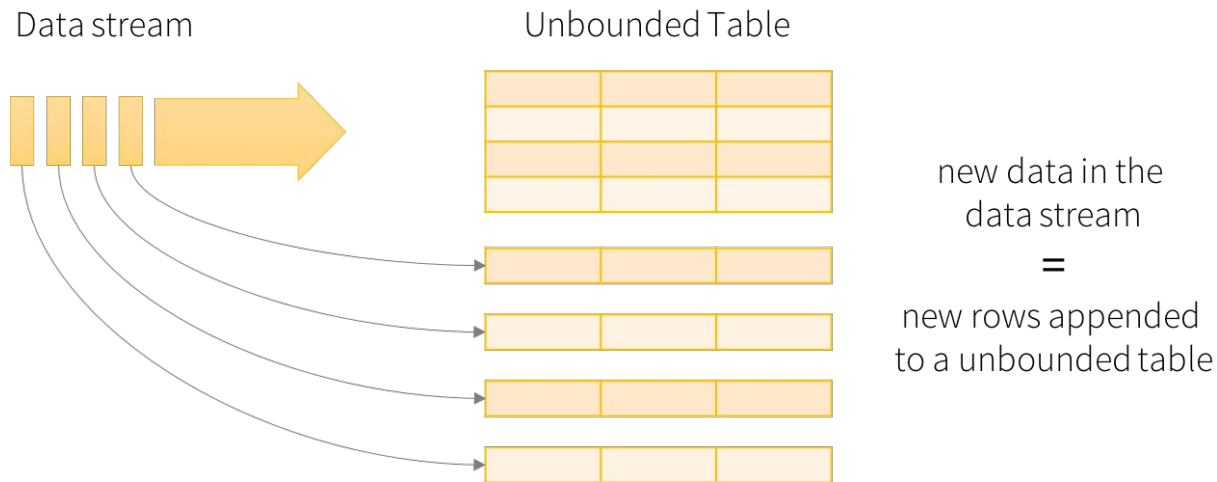


So What's Kafka?

Data Pipeline: 2. Kafka



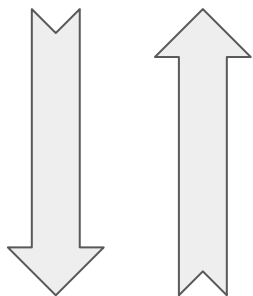
Data Pipeline: 3. Spark



Data stream as an unbounded table

Spark's Structured Streaming

Data Pipeline: Overall



The screenshot shows a Jupyter Notebook with three cells. The first cell contains a Kafka producer code snippet that sends four messages to a topic named 'location'. The second cell contains a Kafka consumer code snippet that reads messages from the 'location' topic. The third cell shows the output of the consumer, displaying four messages with their keys, values, and timestamps.

```

# Producer code
from kafka import KafkaProducer
producer = KafkaProducer(bootstrap_servers=['localhost:9092'])

# Send four messages
producer.send('location', {'latitude': 40.712776, 'longitude': -87.630188, 'name': 'Tokyo Rikiten Center'})
producer.send('location', {'latitude': 40.709156, 'longitude': -87.630188, 'name': 'Tokyo Rikiten Center'})
producer.send('location', {'latitude': 40.709156, 'longitude': -87.630188, 'name': 'Tokyo Rikiten Center'})
producer.send('location', {'latitude': 40.709156, 'longitude': -87.630188, 'name': 'Tokyo Rikiten Center'})

# Consumer code
from kafka import KafkaConsumer
consumer = KafkaConsumer('location', bootstrap_servers=['localhost:9092'], group_id='my-group')

# Read messages
for message in consumer:
    print(message)

```

The output of the consumer shows four messages with the following details:

key	value	timestamp
0	{'latitude': 40.712776, 'longitude': -87.630188, 'name': 'Tokyo Rikiten Center'}	2017-06-15 18:58:00
1	{'latitude': 40.709156, 'longitude': -87.630188, 'name': 'Tokyo Rikiten Center'}	2017-06-15 18:58:00
2	{'latitude': 40.709156, 'longitude': -87.630188, 'name': 'Tokyo Rikiten Center'}	2017-06-15 18:58:00
3	{'latitude': 40.709156, 'longitude': -87.630188, 'name': 'Tokyo Rikiten Center'}	2017-06-15 18:58:00

Demo

Lessons Learned and Future Work

Lessons

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

Future Work

- Front End:
- Back End:

Acknowledgement

Guidance and Discussion from

- Shu Yan
- Yvonne Lau
- Zeyu Zhang

Inspiration from

Chuan Sun and Aiko Liu

Thanks!

