

## Eating Mate Recommendation for Yelp Users

Hao Xu, Zehao Song, Yutong Tie

Department of Computer Science

Columbia University

{hx2208, zs2324, yt2549}@columbia.edu

**Abstract**—Yelp has been more and more prevalent in making recommendations for people’s daily life, like where to eat or where to shop, however, little has been done regards making recommendations for who they could go with. We proposed that eating mate could be recommended based on the review history of each user. Therefore, we collected review data from Yelp to make eating mate recommendations, which included a collaborative filtering of the review data, linear regression on a single user and business data, and finding K nearest neighbors of each user. The experiment on data points confirms that it is a good way to make recommendations for eating mate. Since the current experiment is based on review data on partial cities in the United States, more tests are needed for other cities, and in order to have a better representation for recommendation, a User Interface is needed in the future.

**Keywords**-Yelp; eating mate; recommendation; collaborative filtering

### I. INTRODUCTION

Yelp is a platform about sharing reviews and making recommendations. On the one hand, users of Yelp could submit reviews on different shops and restaurants using a one to five-star rating. On the other hand, Yelp would make recommendations for users based on their past review history. McNichol states Yelp has become “one of the most import sites on the Internet” [1]. According to the statistics provided by Yelp on 2016, it had about 102 million total reviews and a monthly average of 92 million unique visitors in the second quarter of 2016 [2].

Though Yelp has strong functionalities in helping people find a place to eat, it has less tricks in finding people an eating mate. It is a great business loss since under most circumstance, people would not like to go eating alone, no matter how good the business is.

Based on this situation, this project is intended to make recommendations for people who share similar preferences to have meal together based on each individual’s review history. We are dedicated to dig into each user’s tastes which hides behind their

reviews, and hopefully, to analysis taste distributions of each city.

Finally, the application should work like: 1). For a new user, we would like this user answer some question so we can know about his preference and make eating mate recommendation based on his answers. 2). For an existed user, we would make recommendation directly based on his rating history. The final out would be several eating mates and corresponding restaurant.

### II. RELATED WORKS

I do not think there is any previous work that makes recommendation for eating mates, however, “people you may know” is a section that has been widely implemented in social websites. Besides, most E-commerce websites, like Amazon and eBay, make recommendations for goods that people would like to shop. These kinds of recommendation systems are in different tunes rendered with equal skill as eating mate recommendation. I would list a few websites and how they make recommendations in the following paragraphs.

Facebook [3] is dedicated to connect everyone who is using their website. According to Vignesh’s answer on Quora [4], Facebook makes friend recommendations on following 5 aspects: 1). Right after registration, every user would be asked to fill in a questionnaire about his background, like education, region, interest, etc. Based on that information, Facebook would recommend people who went to same school or work at the same place to that user. 2). Once has some friends, Facebook would recommend one’s friends to him. Say if A and B are friends, B and C are friends, then Facebook would recommend A to C and C to A at the same time. 3). Facebook would recommend two people in the same group to be friend. For example, if user A is in Columbia Engineering group and user B is also in that group, then Facebook would recommend A to B and B to A at the same time. 4). Facebook would also recommend two people like

same public page to be friend. For example, if user A likes Bayern Munich and user B likes Bayern Munich as well, then Facebook would recommend them to be friend, however, if A likes Bayern Munich and B likes Borussia Dortmund, it might not recommend them to be friend due to the counter relationship between these two teams.

LinkedIn [5] is designed for people to build up their career circles and trying to help people in job-hunting. LinkedIn makes recommendations for their user based on following 3 aspects: 1). Like Facebook, LinkedIn makes friend recommendation based on users' basic information, like education, region, etc. It would try to connect users from same school, same area or same company. 2). Unlike Facebook, LinkedIn is intended to build up a more real network among people than Facebook, which guarantees the quality of connections in LinkedIn. LinkedIn would import contacts via each users' email address and recommend friends' connections to each user.

Amazon's [6] recommendation starts from finding a set of customers who have items overlap with the user's purchase and rate history. Then it uses different algorithms aggregates items from those customers, remove the items that has already been purchased or rated, and recommends items that are left to the user. There are basically 2 algorithms being applied in Amazon's recommendation. 1). Traditional collaborative filtering. Each customer is represented as an N-dimensional vector. For example, there are two users A and B, the similarity of A and B is calculated by  $\cos(A, B)$ . 2). There are many attributes that described a specific user and the recommendation is based on different cluster that the user belongs to. Therefore, the recommendation problem is simplified to a classification problem. Amazon uses those attributes to divide customers into different clusters.

In our project, we choose to do similar as what Amazon does. We represent each user in Yelp as feature vectors and make recommendations based on the nearest neighbor of that user.

### III. SYSTEM OVERVIEW

In this section, we would ramp up on how we scarp data from Yelp, how our recommendation system works and how the dataset looks like.

### Data Collection.

In order to run scrapy, we would need to generate the URL that contains the data we want. In this case, we need to get business\_id for each city. Yelp provides its Search API for developers to search business\_ids that they are interested.

#### Sample Request:

```
https://api.yelp.com/v2/search?term=food&location=San+Francisco
```

This is how the sample request looks like. In this request, we are fetching business ids in San Francisco that sell food. We fetch data by changing corresponding attributes in the URL. The response from Yelp server is a JSON object. Below is a sample response:

```
"businesses": [
  {
    "categories": [
      [
        "Local Flavor",
        "localflavor"
      ],
      [
        "Mass Media",
        "massmedia"
      ]
    ],
    "display_phone": "+1-415-908-3801",
    "id": "yelp-san-francisco",
    "image_url": "http://s3-media3.fl.yelpcdn.com/bphoto/nQK-6_vZMt5n8zsAS94ew/ms.",
    "is_claimed": true,
    "is_closed": false,
```

This JSON object contains the import information we want, like business id, name, geo location, etc.

After getting enough number of businesses, we start to gather the users' reviews of each business by a python crawler implemented via Scrapy Framework.

Below is how we generate the target URL and fire the request:

```
urls = []
with open('yelp_scrapy/data/%s.jl'%input_file) as f:
    for line in f:
        r = json.loads(line)
        urls.append("https://www.yelp.com/biz/%s"%r[0])
f.close()
for url in urls:
    yield scrapy.Request(url=url, callback=self.parse_review)
```

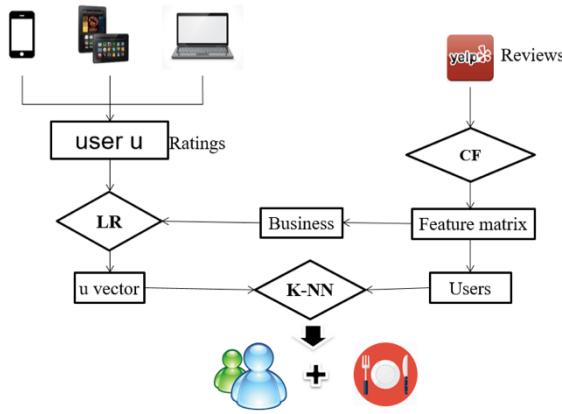
The code above reads lines in generated business id lists in previous step and append it to regular Yelp URL and fire the URL. Then the code parses the HTML code and find the reviews we want in each single page to collect review data. For more detailed

code, take a look at our Github: [http://www.github.com/tsszh/yelp\\_recommendation/blob/master/scrapy/yelp\\_scrapy/spiders/review\\_spider.py](http://www.github.com/tsszh/yelp_recommendation/blob/master/scrapy/yelp_scrapy/spiders/review_spider.py)

### Workflow.

After collecting large enough data, we start to analysis the data and make recommendation for each user. The whole process is as follow:

Figure 1 Eating Mate Recommendation System Overview



The analysis starts from the review data we collected from Yelp.

Firstly, we do a collaborative filtering(CF) over the review data (user\_id, business\_id and rate) to generate the feature matrixs. Save both business feature & user feature matrix for future usage.

In principle, to make recommendation for eating mates, we just need to select users with the minimum distance in terms of user features, which is a typical k-nearest neighbor (K-NN) algorithm.

However, for new users, which is not included in the training sets, the user feature vector is not available. To get the feature vecotr of new user, a typical linear regression algoritm is applied to the users' rating history and corresponding business feature matrix. After getting the best estimation of the user feature vector, apply K-NN algoritm to find similar users from previous User Feature Matrix. Then, we got the eating mates.

After getting the eating mates, we still need to recommend the restaurants for them. The method is

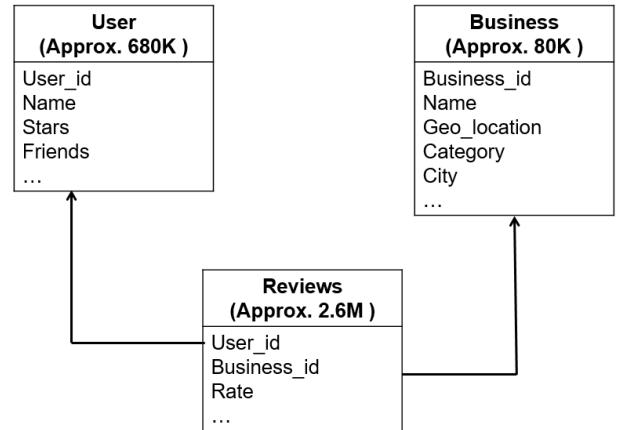
rather simple, taking advantage of the user vectors and business matrix to predict the rating for all business and select the ones with highest ratings. Then, find the shared restaurant between eating mates.

Besides, we could predict the taste or favor of a perticular user. The basic idea is to calculate the categories with the highest term freuecy (Chinese, Pizza, etc...) asscoiated with the user.

### Dataset Overview

There are mainly 3 tables in our database system. They are User data, Business data, and Reivew data.

Figure 2 Dataset Relation Schema Overview



There are about 680K user tuples, 80K business tuples and 2.6M review tuple in total has been collected so far in order to make recommendations.

In User table, each user is assigned a unique key, which is user\_id (primary key). Other information about each user, like name, his friend, his location, etc. has also been included in the table defined above.

In Business table, each business is assigned a unique key, which is business\_id (primary key). Other information about each business, such as name, geolocation, category, etc. is also included in the table defined above.

In review table, each review record is defined by a combination of user\_id and business\_id (primary key). For instance, tuple (1, 2, 3, ...) means that user whose id is 1 gives business whose id is 2 a rate 3.

Other information like rate, comment, timestamp, etc. is also included in the table.

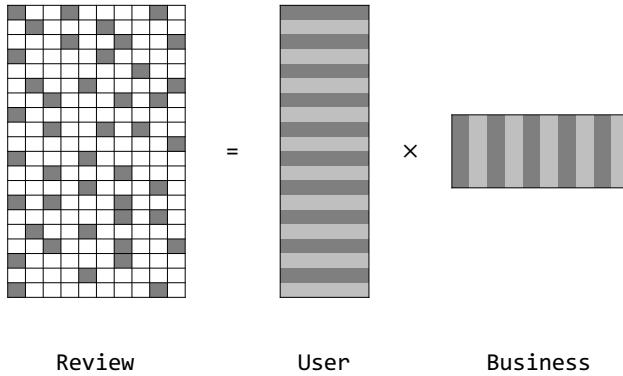
#### IV. ALGORITHM

Show algorithm and tools you have used to solve your problem.

##### 1. Collaborative Filtering (CF)

As shown in Figure 3, the input is the Review Sparse Matrix, each row stands for a specific user and each column represents a specific business. For each cell, the value is the rating in float type if there is a review for corresponding user\_id & business\_id, or zero otherwise.

Figure 3 Illustration of Collaborative Filtering (CF)



The goal of the algorithm is to decompose the Review Matrix ( $R^{n \times d}$ ) into two Rank-k matrix, User Matrix ( $R^{n \times k}$ ) and Business Matrix ( $R^{k \times d}$ ), with minimum square error,

$$\min_{R_{ij} \in \text{Review}} \sum_{R_{ij} \in \text{Review}} (R_{ij} - \langle U_i, B_j \rangle)^2$$

in which  $UB$  is the most similar matrix compared to Review Matrix.

Each row of User Matrix stands for a particular user vector ( $U_i$ ). Each column of Business Matrix stands for a particular business vector ( $B_j$ ). After getting these two feature matrix, it would be simple to calculate the predicted rating of User  $i$  on Business  $j$ .

$$R_{ij} = \langle U_i, B_j \rangle$$

To get User & Business Matrix, we use **Alternating Least Squares (ALS)** algorithm. The basic idea is:

- (1) Suppose Business Matrix is given. Each row of User Matrix can be calculated via simple Linear Regression.
- (2) Suppose User Matrix is given. Each row of Business Matrix can be calculated via simple Linear Regression.
- (3) Alternate between these two simple problem.

```
Initialize  $u_i \in R^k$  and  $v_j \in R^k$ 
for  $t = 1, 2, \dots, \text{NumIter}$ :
    For each user  $i = 1, 2, 3, \dots, n$ :
         $u_i = \text{argmin} \sum_{(i,j) \in \Omega} (\langle u, v_j \rangle - R_{ij})^2$ 
    For each business  $j = 1, 2, 3, \dots, d$ :
         $v = \text{argmin} \sum_{(i,j) \in \Omega} (\langle u_i, v \rangle - R_{ij})^2$ 
```

Actually, Spark provide this model-based recommendation algorithm.

```
from pyspark.mllib.recommendation import ALS,
MatrixFactorizationModel, Rating

ALS_setting = {
    'rank': 20,
    'numIterations': 20
}

model = ALS.train(
    ratings,
    ALS_setting['rank'],
    ALS_setting['numIterations']
)

vBusiness = model.productFeatures()
vUser = model.userFeatures()
```

By using Spark Library, CF-ALS becomes straightforward. There are two parameters in the model to be tuned, rank & numIterations. Rank is the dimension of User & Business Vector. NumIterations is the number of iterations in ALS. The choice of these two parameters will be discussed in Experiment Results Section.

## 2. Linear Regression (LR)

**Input:** Review history of unknown user ( $b \in R^m$ ) and Corresponding Business Feature Matrix ( $A \in R^{m \times k}$ ). Note,  $m$  is the number of reviews given by user, which might be smaller than  $k$ , rank of the user vector. If  $m < k$ , it is impossible to get reasonable LR result. Thus, we only provide recommendation for users who rates at least  $k$  restaurants.

**Output:** User Feature Vector of unknown user ( $u \in R^k$ ).

**Objective:**  $\min \|Au - b\|^2$

### Algorithm I: Ordinary Least Squares

Direct Mathematic Calculation via:

$$u = (A^T A)^{-1} A^T b$$

Pros: Accurate, Simple, Fast for small matrix

Cons: Slow for large matrix

### Algorithm II: Gradient Decent

Iteration method. Provided by Spark Library.

Pros: Fast for large matrix

Cons: Less accurate for small matrix

The number of reviews given by user is limited. Normally,  $m < 100$ . Obviously, ordinary least squares is much suitable for this circumstance.

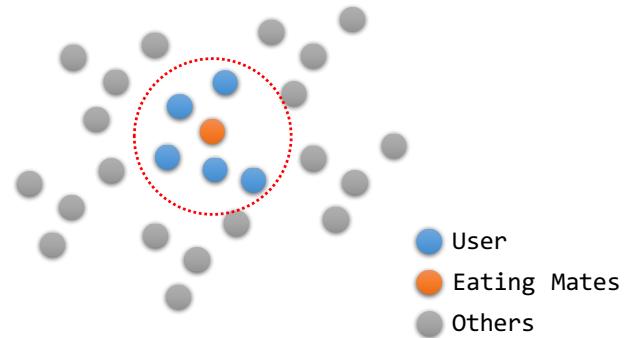
```
import numpy as np

vb, mA = getMatrix(compactM)
vu = np.linalg.inv(mA.T.dot(mA)) \
    .dot(mA.T) \
    .dot(vb)
```

The calculation is direct and simple with the help of numpy library.

## 3. K-Nearest Neighbor (kNN)

Figure 4 Illustration of K-Nearest Neighbor Algorithm



After getting the user's feature vector, we could give the recommendation by finding the most similar users (minimum norm-2 distance) in User Feature Matrix. Spark doesn't have default library for K-NN algorithm. So, we implemented the naïve K-NN by ourselves in Spark.

```
def knn(vUser, vu):
    return vUser.map(lambda r: (
        r[0],
        np.linalg.norm(vu-r[1])
    ))
    .top(k, key=lambda r: -r[1])
```

## 4. Restaurant Recommendation

**Input:** User Vector & Business Matrix

**Output:** List of recommended restaurants

(1) Predict the rating for all restaurants by

$$R_j = u^T \cdot U_j$$

(2) Select the top 10 restaurants with highest ratings

(3) Select the top shared restaurants for eating mates as the recommendation

```
def recomd(vBusiness, vu, num=10):
    return map(
        lambda r: r[0],
        vBusiness.map(
```

```

lambda r: (
    r[0],
    np.array(r[1]).dot(vu))
).top(num, key=lambda r: r[1])
)

def sharedRecom(vu1, vu2):
    pool = set( recomd(vBusiness, vu1) )
    return filter(
        lambda r: r in pool,
        recomd(vBusiness, vu2)[:1]
)

```

## 5. Taste Analysis

**Input:** User Vector & Business Matrix

**Output:** List of (Category, Weight)

```

from operator import add

def getTaste(vBusiness, vu):
    rest = recomd(vBusiness, vu, num=1000)
    taste = sc.parallelize(rest)\.
        flatMap(
            lambda r: business_dict[r][1])
        .map(
            lambda r: (r, 1))
        .reduceByKey(add)\.
        takeOrdered(
            50, key=lambda r: -r[1])
)

```

- (1) Get top 1000 recommended restaurants
- (2) Count the term frequency of each category associated with these restaurants
- (3) Select the Top 50 categories.

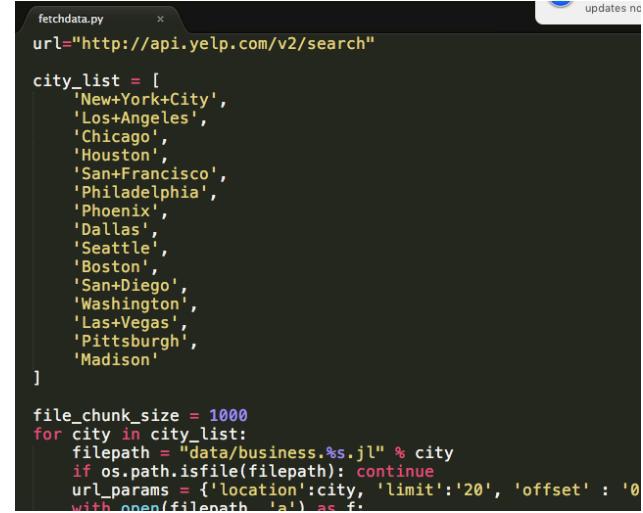
## V. SOFTWARE PACKAGE DESCRIPTION

Describe the software package that is going to be open sourced. Show some screen shots of how user use it and or UI.

### 1. Yelp API

Path: scrapy/yelp\_fetcher/fetchdata.py

First, add the city name to the `city_list`, eg, "Washington".



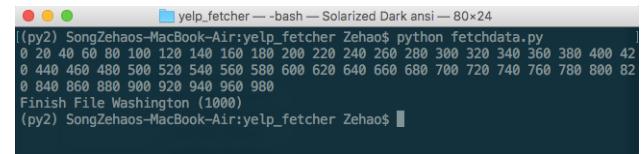
```

fetchdata.py
url="http://api.yelp.com/v2/search"
city_list = [
    'New+York+City',
    'Los+Angeles',
    'Chicago',
    'Houston',
    'San+Francisco',
    'Philadelphia',
    'Phoenix',
    'Dallas',
    'Seattle',
    'Boston',
    'San+Diego',
    'Washington',
    'Las+Vegas',
    'Pittsburgh',
    'Madison'
]

file_chunk_size = 1000
for city in city_list:
    filepath = "data/business.%s.jl" % city
    if os.path.isfile(filepath): continue
    url_params = {'location':city, 'limit':'20', 'offset' : '0
    with open(filepath, 'a') as f:

```

Second, run the scripts with python 2.

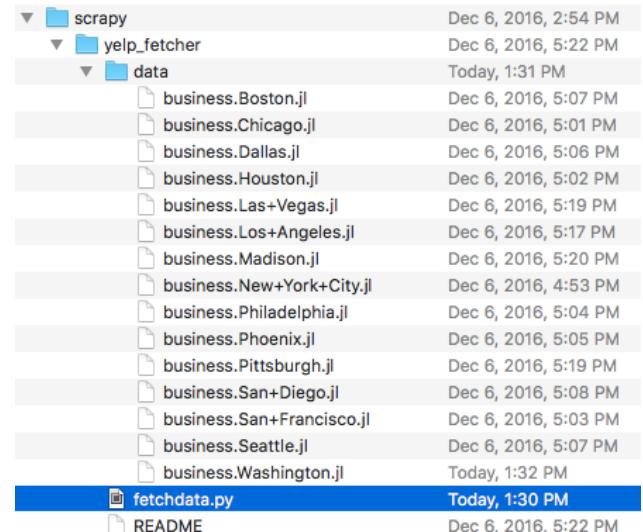


```

SongZehao-MacBook-Air:yelp_fetcher Zehao$ python fetchdata.py
0 20 40 60 80 100 120 140 160 180 200 220 240 260 280 300 320 340 360 380 400 42
0 440 460 480 500 520 540 560 580 600 620 640 660 680 700 720 740 760 780 800 82
0 840 860 880 900 920 940 960 980
Finish File Washington (1000)
SongZehao-MacBook-Air:yelp_fetcher Zehao$ 

```

Now, the results will be stored in "data/business.Washington.jl"



File/Folder	Last Modified
scrapy	Dec 6, 2016, 2:54 PM
yelp_fetcher	Dec 6, 2016, 5:22 PM
data	Today, 1:31 PM
business.Boston.jl	Dec 6, 2016, 5:07 PM
business.Chicago.jl	Dec 6, 2016, 5:01 PM
business.Dallas.jl	Dec 6, 2016, 5:06 PM
business.Houston.jl	Dec 6, 2016, 5:02 PM
business.Las+Vegas.jl	Dec 6, 2016, 5:19 PM
business.Los+Angeles.jl	Dec 6, 2016, 5:17 PM
business.Madison.jl	Dec 6, 2016, 5:20 PM
business.New+York+City.jl	Dec 6, 2016, 4:53 PM
business.Philadelphia.jl	Dec 6, 2016, 5:04 PM
business.Phoenix.jl	Dec 6, 2016, 5:05 PM
business.Pittsburgh.jl	Dec 6, 2016, 5:19 PM
business.San+Diego.jl	Dec 6, 2016, 5:08 PM
business.San+Francisco.jl	Dec 6, 2016, 5:03 PM
business.Seattle.jl	Dec 6, 2016, 5:07 PM
business.Washington.jl	Today, 1:32 PM
fetchdata.py	Today, 1:30 PM
README	Dec 6, 2016, 5:22 PM

Here is the results:

```
▶ business.Washington.jl ×
1 ["lincoln-memorial-washington-2", "Lincoln Memorial", 5.0,
  [{"Landmarks & Historical Buildings", "landmarks"}], {
  {"latitude": 38.8993459813715, "longitude": -77.052001422584}, 716]
2 ["samis-coffee-kiosk-washington", "Sami's Coffee Kiosk", 5.0,
  [{"Coffee & Tea"}, {"coffee"}], {"latitude": 38.89647,
  "longitude": -77.04813}, 38]
3 ["vietnam-veterans-memorial-washington", "Vietnam Veterans
  Memorial", 5.0, [{"Landmarks & Historical Buildings",
  "landmarks"}], {"latitude": 38.8910251787901, "longitude": -77.0476747857971}, 236]
4 ["un-ne-sais-quoi-washington", "Un Je Ne Sais Quoi", 5.0,
  [{"Desserts", "desserts"}, {"Cafes", "cafes"}, {"Bakeries",
  "bakeries"}], {"latitude": 38.9087384, "longitude": -77.0425771, 116]
```

Each line is a record of business, including business\_id, business\_name, overall rating, categories, location and review count.

## 2. Review Spider

Path: scrapy/yelp\_scrapy

First, copy the business list file collected in the previous step into “scrapy/yelp\_scrapy/data” folder and change the input\_file variable to proper path.

```
OPEN FILES
FOLDERS
  ▶ yelp_scrapy
    ▶ yelp_scrapy
      ▶ __pycache__
        ▶ __init__.py
      ▶ data
        ▶ business.Boston.jl
    ▶ spiders
      ▶ __pycache__
        ▶ __init__.py
      ▶ review_spider.py
    ▶ items.py
    ▶ log.py
    ▶ pipelines.py
    ▶ settings.py
    ▶ error.log.txt
    ▶ info.log.txt
    ▶ scrapy.cfg
    ▶ start.sh
  ▶ review_spider.py
    ▶ __init__.py
    ▶ items.py
    ▶ log.py
    ▶ pipelines.py
    ▶ settings.py
    ▶ error.log.txt
    ▶ info.log.txt
    ▶ scrapy.cfg
    ▶ start.sh
```

```
review_spider.py ×
1 import scrapy, re, json, logging
2 from yelp_scrapy.items import Review
3 import yelp_scrapy.log
4
5 logger = logging.getLogger('my-logger')
6
7 input_file = 'business.Boston'
8
9 class ReviewSpider(scrapy.Spider):
10     name = "reviews"
11
12     def __init__(self, category=None, *args, **kwargs):
13         super(ReviewSpider, self).__init__(*args, **kwargs)
14         self.REVIEW_LIMIT = 1000
15         self.countDict = {}
16
17     def start_requests(self):
18         urls = []
19         with open('yelp_scrapy/data/%s.jl'%input_file) as input_file:
20             for line in f:
21                 r = json.loads(line)
22                 urls.append("https://www.yelp.com/biz/%s" % r['id'])
```

Then, in terminal, cd to the root folder of yelp\_scrapy and run “sh start.sh”.

```
(myz) SongZhehao-MacBook-Air:yelp_scrapy.Zhehao$ ls
error_log.txt info.log.txt scrapy.cfg start.sh    yelp_scrapy
(myz) SongZhehao-MacBook-Air:yelp_scrapy.Zhehao$ sh start.sh
2016-12-22 13:43:22 [scrapy] INFO: Scrapy 1.1.1 started (bot: yelp_scrapy)
2016-12-22 13:43:22 [scrapy] INFO: Overridden settings: {'DOWNLOAD_DELAY': 0.5, 'BOT_NAME': 'yelp_scrapy', 'NEMSPIDER_MODULE': 'yelp_scrapy.spiders', 'LOG_LEVEL': 'INFO', 'SPIDER_MODULES': ['yelp_scrapy.spiders']}
2016-12-22 13:43:22 [scrapy] INFO: Enabled extensions:
['scrapy.extensions.logstats.LogStats']
['scrapy.extensions.telnet.TelnetConsole']
2016-12-22 13:43:22 [scrapy] INFO: Enabled downloader middlewares:
['scrapy.downloadermiddlewares.httpauth.HttpAuthMiddleware',
 'scrapy.downloadermiddlewares.downloadtimeout.DownloadTimeoutMiddleware',
 'scrapy.downloadermiddlewares.useragent.UserAgentMiddleware',
 'scrapy.downloadermiddlewares.retry.RetryMiddleware',
 'scrapy.downloadermiddlewares.defaultheaders.DefaultHeadersMiddleware',
 'scrapy.downloadermiddlewares.redirect.MetaRefreshMiddleware',
 'scrapy.downloadermiddlewares.httpcompression.HttpCompressionMiddleware',
 'scrapy.downloadermiddlewares.redirect.RedirectMiddleware',
 'scrapy.downloadermiddlewares.cookies.CookiesMiddleware',
 'scrapy.downloadermiddlewares.chunked.ChunkedTransferMiddleware',
 'scrapy.downloadermiddlewares.stats.DownloaderStats']
2016-12-22 13:43:22 [scrapy] INFO: Enabled spider middlewares:
['scrapy.spidermiddlewares.httperror.HttpErrorMiddleware',
 'scrapy.spidermiddlewares.offsite.OffsiteMiddleware',
 'scrapy.spidermiddlewares.referer.RefererMiddleware',
 'scrapy.spidermiddlewares.urllength.UrlLengthMiddleware',
 'scrapy.spidermiddlewares.depth.DepthMiddleware']
2016-12-22 13:43:22 [scrapy] INFO: Enabled item pipelines:
['yelp_scrapy.pipelines.YelpPipeline']
2016-12-22 13:43:22 [scrapy] INFO: Spider opened
2016-12-22 13:43:22 [scrapy] INFO: Crawled 0 pages (at 0 pages/min), scraped 0 items (at 0 items/min)
2016-12-22 13:43:24 [my-logger] INFO: Success - https://www.yelp.com/biz/polaris-coffee-boston
2016-12-22 13:43:25 [my-logger] INFO: Success - https://www.yelp.com/biz/grandens-distilling-boston
2016-12-22 13:43:25 [my-logger] INFO: Success - https://www.yelp.com/biz/piperi-mediterranean-grill-boston
2016-12-22 13:43:27 [my-logger] INFO: Success - https://www.yelp.com/biz/norman-b-leventhal-park-boston
```

**Note:** Yelp might ban your IP.

The results are saved in “yelp\_scrapy/data/business.Boston.reviews.jl”.

```
OPEN FILES
FOLDERS
  ▶ yelp_scrapy
    ▶ __pycache__
      ▶ __init__.py
      ▶ review_spider.py
    ▶ data
      ▶ business.Boston.jl
      ▶ business.Boston.reviews.jl
    ▶ spiders
      ▶ __pycache__
        ▶ __init__.py
      ▶ items.py
      ▶ log.py
      ▶ pipelines.py
      ▶ settings.py
      ▶ error.log.txt
      ▶ info.log.txt
      ▶ scrapy.cfg
      ▶ start.sh
```

```
business.Boston.reviews.jl ×
1 [{"FZBTKAZExOpTCYRV22W0", "William G.", "5.0",
  "polaris-coffee-boston", "I walked in here with the
  impression that it was a coffee shop. Indeed, they sell
  coffee and have a brew bar. However, during my brief
  chat with the owner, we chatted a bit about the borough and where to
  go. It turns out I was in Little Italy, and it has a ton of
  restaurants and coffee shops. After a few minutes of
  conversation, I went to look for a coffee shop for an iced
  latte. I found one and asked about the ratings and cleanliness
  of this place, and their good variety. If you're looking to
  buy coffee, but not necessarily having it brewed, then check
  this place out.\n\nS. They DO take credit cards."}]
2 [{"ZGQDg_smkVjIZ2Vg2K8g", "Jae T.", "4.0",
  "polaris-coffee-boston", "I would translate this place
  as 'Ice coffee cafe'. I would say that ice coffee is
  excellent !!!\nMust have !!! To know that living in N.P.N.Hs
  is possible after having a cup of this !!!\nIce Superman stuff !!! I
  can live there after I drink this !!!\nImpossible thinking
  goes away !!!\nIce coffee fan only !!!"]
3 [{"s3jaex-0pUfU0E4PCJy0g", "Dorian G.", "5.0",
  "polaris-coffee-boston", "I just had a latte here and it was delicious! The coffee was strong and well-made. The service was friendly and efficient. I will definitely be back! Thanks for the great coffee!"}]
```

Each line is a record of review including user\_id, user\_name, rating and review content.

## 3. Train Model

Path: spark/recommendation.ipynb

Open the file in Jupyter Notebook and run the scripts in order.

The configure the ALS\_setting to change the training parameters.

**Import Library**

```
In [5]: from pyspark.mllib.recommendation import ALS, MatrixFactorizationModel, Rating
from pyspark.mllib.regression import LabeledPoint, LinearRegressionWithSGD, LinearRegressionModel
import numpy as np
import json, time
from operator import add
```

**Define File Path & ALS Setting**

```
In [6]: filepat = "data/reviews.jl"
save_dir = "target/"
ALS_setting = {
  'rank': 10,
  'numIterations': 20
}
```

**Load & Process Data**

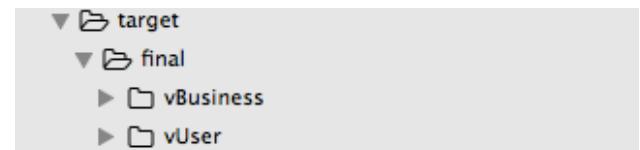
```
In [7]: # Load and parse the data.
data = sc.textFile(filepat)
ratings = data.map(lambda l: json.loads(l))\
.map(lambda l: Rating(int(l[0]), int(l[1]), float(l[2])).cache())
```

Train & save the model.

**Build the complete CF Model**

```
In [6]: model = cf_als(ratings, ALS_setting)
Training CF Model
Time Cost: 177.23s
In [10]: saveModel(model, "final")
Number of Business: 85539
Number of User: 686556
Time Cost: 14.81s
```

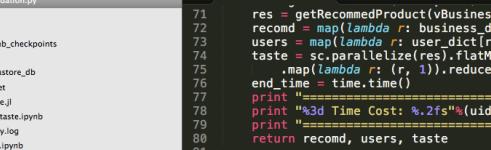
The saved User & Business Feature Matrix.



#### 4. Recommendation

Path: spark/recommendation.py

First, modify the range of users to be recommended and the save file path.



The screenshot shows the PyCharm IDE interface. The top bar displays the title "recommendation.py - spark". The left sidebar lists the project structure with files like "recommendation.py", ".ipynb\_checkpoints", "data", "metastore\_db", "target", "cache.jl", "city\_taste.ipynb", "derby.log", "final.ipynb", "parse\_business.py", "parse\_user.py", "recommendation.ipynb", and "user\_recomm.ipynb". The main code editor window contains the Scala code for the "recommendation.py" script. The code uses parallel processing and reduceByKey operations to calculate recommended products based on user-business interactions.

```
res = getRecommendedProduct(vBusiness, vu, 100)
recomd = map(lambda r: business_dict[r[0]], users = map(lambda r: user_dict[r[0]], getTaste = sc.parallelize(res).flatMap(lambda r: map(lambda t: (r, t), 1).reduceByKey(add).end_time = time.time()
print "-----"
print "Total Time Cost: %.2fs" % (uid, end_time)
print "-----"
return recomd, users, taste

with open('cache1jl', "a") as f:
    for i in range(1,10):
        print "-----"
        print "Start Calculation On %d" % i
        print "-----"
        r,u,t = getTaste(i)
        f.write(json.dumps([i, r, u, t]))
        f.write("\n")
        f.flush()
f.close()
```

Second, run the scripts “spark-submit recommendation.py”.

```
16/12/22 14:11:32 INFO DAGScheduler: Job 63 finished: takeOrdered at /Users/Zehao/Documents/workspace/yelp_recommenderation/spark/recommendation.py:75, took 0.108457 s
=====
9 Time Cost: 22.3s

16/12/22 14:11:33 INFO SparkContext: Invoking stop() from shutdown hook
16/12/22 14:11:33 INFO SparkUI: Stopped Spark Web UI at http://192.168.0.7:4040
16/12/22 14:11:33 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
16/12/22 14:11:33 INFO MemoryStore: MemoryStore cleared
16/12/22 14:11:33 INFO BlockManager: BlockManager stopped
16/12/22 14:11:33 INFO BlockManagerMaster: BlockManagerMaster stopped
16/12/22 14:11:33 INFO OutputCommitCoordinator$OutputPortCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
16/12/22 14:11:33 INFO SparkContext: Successfully stopped SparkContext
16/12/22 14:11:33 INFO ShutdownHookManager: Shutdown hook called
16/12/22 14:11:33 INFO ShutdownHookManager: Deleting directory /private/var/folders/43/11r6nts1g17gkcy5fkfz2w000gn/T/spark-1773d74-87bd-44c6-b9ff-b3e37a-f1f140c3-ba18-0758d060e9fe
16/12/22 14:11:33 INFO ShutdownHookManager: Deleting directory /private/var/folders/43/11r6nts1g17gkcy5fkfz2w000gn/T/spark-1773d74-87bd-44c6-b9ff-b3e37a-f1f140c3-ba18-0758d060e9fe
(pyz) Zehao-Zebras-MacBook-Pro:~ spark Zehao
```

The results are saved into “cache1.jl”.

```
recommendation.py      cache1.json
1 [1, ["Sol Caribe Puerto Rican Cuisine", "Scottsdale", "Island Fusion, Las Vegas"], ["Papa Del's Pizza", "Champaign", "Black Rock Pizza, Urbana"], ["Lupitas Carniceria & Tortilleria, Charlotte"], [{"Jeremy", "Henri", "Jake", "Cathy", "Russel"}, [{"American (Traditional)", 108}, {"Pizza", 106}, {"Mexican", 92}, {"Sandwiches", 90}, {"Italian", 68}, {"American (New)", 66}, {"Chinese", 65}, {"Fast Food", 64}, {"Burgers", 45}, {"Cafes", 42}, {"Coffee & Tea", 40}, {"Deli's", 37}, {"Seafood", 34}, {"Indian", 32}, {"Sushi Bars", 31}, {"Thai", 29}, {"Japanese", 29}, {"Mediterranean", 28}, {"Chicken Wings", 27}, {"Pubs", 27}, {"British", 27}, {"Barbeque", 26}, {"Salad", 22}, {"Greek", 19}, {"Asian Fusion", 19}, {"Latin American", 19}, {"Sports Bars", 19}, {"Steakhouses", 18}, {"Vietnamese", 16}, {"French", 14}, {"Middle Eastern", 13}, {"Fish & Chips", 12}, {"Gastropubs", 12}, {"Bakeries", 12}, {"Wine Bars", 12}, {"Specialty Food", 12}, {"Buffets", 11}, {"Lounges", 11}, {"Caribbean", 11}, {"Vegetarian", 11}, {"Event Planning & Services", 10}, {"Southern", 10}, {"Korean", 9}, {"Juice Bars & Smoothies", 8}, {"Food Trucks", 8}, {"Tex-Mex", 7}, {"Soup", 7}, {"Peruvian", 7}, {"Soul Food", 7}, {"Gluten-Free", 6}]]  
2 [2, {"Yak & Yeti", "Edinburgh", "99 Pub and Grill, Avondale", "Cabo's Bar and Grill, Casa Grande", "Field, Edinburgh", "Corral West Arena, Goodyear"}, {"Brian", "Melanie", "daniel", "A", "Mike"}, [{"Mexican", 151}, {"Pizza", 116}, {"Chinese", 114}, {"Italian", 113}, {"American (Traditional)", 108}, {"Seafood", 94}, {"Fast Food", 88}, {"Buffets", 85}, {"Deli's", 81}, {"Gastropubs", 79}, {"Wine Bars", 76}, {"Lounges", 74}, {"Middle Eastern", 73}, {"Sports Bars", 72}, {"Asian Fusion", 71}, {"Caribbean", 70}, {"Lounging", 69}, {"Peruvian", 68}, {"Mexican Fusion", 67}, {"Italian Fusion", 66}, {"American (New)", 66}, {"Chinese Fusion", 65}, {"French Fusion", 64}, {"Middle Eastern Fusion", 63}, {"Sports Bar", 62}, {"Asian Bar", 61}, {"Mexican Bar", 60}, {"Italian Bar", 59}, {"American Bar", 58}, {"Chinese Bar", 57}, {"French Bar", 56}, {"Middle Eastern Bar", 55}, {"Sports Bar", 54}, {"Asian Bar", 53}, {"Mexican Bar", 52}, {"Italian Bar", 51}, {"American Bar", 50}, {"Chinese Bar", 49}, {"French Bar", 48}, {"Middle Eastern Bar", 47}, {"Sports Bar", 46}, {"Asian Bar", 45}, {"Mexican Bar", 44}, {"Italian Bar", 43}, {"American Bar", 42}, {"Chinese Bar", 41}, {"French Bar", 40}, {"Middle Eastern Bar", 39}, {"Sports Bar", 38}, {"Asian Bar", 37}, {"Mexican Bar", 36}, {"Italian Bar", 35}, {"American Bar", 34}, {"Chinese Bar", 33}, {"French Bar", 32}, {"Middle Eastern Bar", 31}, {"Sports Bar", 30}, {"Asian Bar", 29}, {"Mexican Bar", 28}, {"Italian Bar", 27}, {"American Bar", 26}, {"Chinese Bar", 25}, {"French Bar", 24}, {"Middle Eastern Bar", 23}, {"Sports Bar", 22}, {"Asian Bar", 21}, {"Mexican Bar", 20}, {"Italian Bar", 19}, {"American Bar", 18}, {"Chinese Bar", 17}, {"French Bar", 16}, {"Middle Eastern Bar", 15}, {"Sports Bar", 14}, {"Asian Bar", 13}, {"Mexican Bar", 12}, {"Italian Bar", 11}, {"American Bar", 10}, {"Chinese Bar", 9}, {"French Bar", 8}, {"Middle Eastern Bar", 7}, {"Sports Bar", 6}, {"Asian Bar", 5}, {"Mexican Bar", 4}, {"Italian Bar", 3}, {"American Bar", 2}, {"Chinese Bar", 1}], [{"name": "Sol Caribe Puerto Rican Cuisine", "city": "Scottsdale", "cuisine": "Island Fusion, Las Vegas"}, {"name": "Papa Del's Pizza", "city": "Champaign", "cuisine": "Black Rock Pizza, Urbana"}, {"name": "Lupitas Carniceria & Tortilleria", "city": "Charlotte", "cuisine": "Charlotte"}, {"name": "Jeremy", "city": null, "cuisine": "Henri, Jake, Cathy, Russel"}, {"name": "American (Traditional)", "city": null, "cuisine": "108"}, {"name": "Pizza", "city": null, "cuisine": "106"}, {"name": "Mexican", "city": null, "cuisine": "92"}, {"name": "Sandwiches", "city": null, "cuisine": "90"}, {"name": "Italian", "city": null, "cuisine": "68"}, {"name": "American (New)", "city": null, "cuisine": "66"}, {"name": "Chinese", "city": null, "cuisine": "65"}, {"name": "Fast Food", "city": null, "cuisine": "64"}, {"name": "Burgers", "city": null, "cuisine": "45"}, {"name": "Cafes", "city": null, "cuisine": "42"}, {"name": "Coffee & Tea", "city": null, "cuisine": "40"}, {"name": "Deli's", "city": null, "cuisine": "37"}, {"name": "Seafood", "city": null, "cuisine": "34"}, {"name": "Indian", "city": null, "cuisine": "32"}, {"name": "Sushi Bars", "city": null, "cuisine": "31"}, {"name": "Thai", "city": null, "cuisine": "29"}, {"name": "Japanese", "city": null, "cuisine": "29"}, {"name": "Mediterranean", "city": null, "cuisine": "28"}, {"name": "Chicken Wings", "city": null, "cuisine": "27"}, {"name": "Pubs", "city": null, "cuisine": "27"}, {"name": "British", "city": null, "cuisine": "27"}, {"name": "Barbeque", "city": null, "cuisine": "26"}, {"name": "Salad", "city": null, "cuisine": "22"}, {"name": "Greek", "city": null, "cuisine": "19"}, {"name": "Asian Fusion", "city": null, "cuisine": "19"}, {"name": "Latin American", "city": null, "cuisine": "19"}, {"name": "Sports Bars", "city": null, "cuisine": "19"}, {"name": "Steakhouses", "city": null, "cuisine": "18"}, {"name": "Vietnamese", "city": null, "cuisine": "16"}, {"name": "French", "city": null, "cuisine": "14"}, {"name": "Middle Eastern", "city": null, "cuisine": "13"}, {"name": "Fish & Chips", "city": null, "cuisine": "12"}, {"name": "Gastropubs", "city": null, "cuisine": "12"}, {"name": "Bakeries", "city": null, "cuisine": "12"}, {"name": "Wine Bars", "city": null, "cuisine": "12"}, {"name": "Specialty Food", "city": null, "cuisine": "12"}, {"name": "Buffets", "city": null, "cuisine": "11"}, {"name": "Lounges", "city": null, "cuisine": "11"}, {"name": "Caribbean", "city": null, "cuisine": "11"}, {"name": "Vegetarian", "city": null, "cuisine": "11"}, {"name": "Event Planning & Services", "city": null, "cuisine": "10"}, {"name": "Southern", "city": null, "cuisine": "10"}, {"name": "Korean", "city": null, "cuisine": "9"}, {"name": "Juice Bars & Smoothies", "city": null, "cuisine": "8"}, {"name": "Food Trucks", "city": null, "cuisine": "8"}, {"name": "Tex-Mex", "city": null, "cuisine": "7"}, {"name": "Soup", "city": null, "cuisine": "7"}, {"name": "Peruvian", "city": null, "cuisine": "7"}, {"name": "Soul Food", "city": null, "cuisine": "7"}, {"name": "Gluten-Free", "city": null, "cuisine": "6"}], [{"name": "Yak & Yeti", "city": "Edinburgh", "cuisine": "99 Pub and Grill, Avondale"}, {"name": "Cabo's Bar and Grill", "city": "Casa Grande", "cuisine": "Field, Edinburgh"}, {"name": "Corral West Arena", "city": "Goodyear", "cuisine": "Brian, Melanie, daniel"}, {"name": "Sports Bar", "city": null, "cuisine": "A, Mike"}, {"name": "Mexican", "city": null, "cuisine": "151"}, {"name": "Italian", "city": null, "cuisine": "116"}, {"name": "American (Traditional)", "city": null, "cuisine": "114"}, {"name": "Chinese", "city": null, "cuisine": "113"}, {"name": "Mexican Fusion", "city": null, "cuisine": "108"}, {"name": "Italian Fusion", "city": null, "cuisine": "94"}, {"name": "Fast Food", "city": null, "cuisine": "88"}, {"name": "Buffets", "city": null, "cuisine": "85"}, {"name": "Deli's", "city": null, "cuisine": "81"}, {"name": "Gastropubs", "city": null, "cuisine": "79"}, {"name": "Wine Bars", "city": null, "cuisine": "76"}, {"name": "Lounges", "city": null, "cuisine": "74"}, {"name": "Middle Eastern", "city": null, "cuisine": "73"}, {"name": "Sports Bar", "city": null, "cuisine": "72"}, {"name": "Asian Fusion", "city": null, "cuisine": "71"}, {"name": "Caribbean", "city": null, "cuisine": "70"}, {"name": "Lounging", "city": null, "cuisine": "69"}, {"name": "Peruvian", "city": null, "cuisine": "68"}, {"name": "Mexican Fusion", "city": null, "cuisine": "67"}, {"name": "Italian Fusion", "city": null, "cuisine": "66"}, {"name": "American (New)", "city": null, "cuisine": "65"}, {"name": "Chinese Fusion", "city": null, "cuisine": "64"}, {"name": "French Fusion", "city": null, "cuisine": "63"}, {"name": "Middle Eastern Fusion", "city": null, "cuisine": "62"}, {"name": "Sports Bar", "city": null, "cuisine": "61"}, {"name": "Asian Bar", "city": null, "cuisine": "60"}, {"name": "Mexican Bar", "city": null, "cuisine": "59"}, {"name": "Italian Bar", "city": null, "cuisine": "58"}, {"name": "American Bar", "city": null, "cuisine": "57"}, {"name": "Chinese Bar", "city": null, "cuisine": "56"}, {"name": "French Bar", "city": null, "cuisine": "55"}, {"name": "Middle Eastern Bar", "city": null, "cuisine": "54"}, {"name": "Sports Bar", "city": null, "cuisine": "53"}, {"name": "Asian Bar", "city": null, "cuisine": "52"}, {"name": "Mexican Bar", "city": null, "cuisine": "51"}, {"name": "Italian Bar", "city": null, "cuisine": "50"}, {"name": "American Bar", "city": null, "cuisine": "49"}, {"name": "Chinese Bar", "city": null, "cuisine": "48"}, {"name": "French Bar", "city": null, "cuisine": "47"}, {"name": "Middle Eastern Bar", "city": null, "cuisine": "46"}, {"name": "Sports Bar", "city": null, "cuisine": "45"}, {"name": "Asian Bar", "city": null, "cuisine": "44"}, {"name": "Mexican Bar", "city": null, "cuisine": "43"}, {"name": "Italian Bar", "city": null, "cuisine": "42"}, {"name": "American Bar", "city": null, "cuisine": "41"}, {"name": "Chinese Bar", "city": null, "cuisine": "40"}, {"name": "French Bar", "city": null, "cuisine": "39"}, {"name": "Middle Eastern Bar", "city": null, "cuisine": "38"}, {"name": "Sports Bar", "city": null, "cuisine": "37"}, {"name": "Asian Bar", "city": null, "cuisine": "36"}, {"name": "Mexican Bar", "city": null, "cuisine": "35"}, {"name": "Italian Bar", "city": null, "cuisine": "34"}, {"name": "American Bar", "city": null, "cuisine": "33"}, {"name": "Chinese Bar", "city": null, "cuisine": "32"}, {"name": "French Bar", "city": null, "cuisine": "31"}, {"name": "Middle Eastern Bar", "city": null, "cuisine": "30"}, {"name": "Sports Bar", "city": null, "cuisine": "29"}, {"name": "Asian Bar", "city": null, "cuisine": "28"}, {"name": "Mexican Bar", "city": null, "cuisine": "27"}, {"name": "Italian Bar", "city": null, "cuisine": "26"}, {"name": "American Bar", "city": null, "cuisine": "25"}, {"name": "Chinese Bar", "city": null, "cuisine": "24"}, {"name": "French Bar", "city": null, "cuisine": "23"}, {"name": "Middle Eastern Bar", "city": null, "cuisine": "22"}, {"name": "Sports Bar", "city": null, "cuisine": "21"}, {"name": "Asian Bar", "city": null, "cuisine": "20"}, {"name": "Mexican Bar", "city": null, "cuisine": "19"}, {"name": "Italian Bar", "city": null, "cuisine": "18"}, {"name": "American Bar", "city": null, "cuisine": "17"}, {"name": "Chinese Bar", "city": null, "cuisine": "16"}, {"name": "French Bar", "city": null, "cuisine": "15"}, {"name": "Middle Eastern Bar", "city": null, "cuisine": "14"}, {"name": "Sports Bar", "city": null, "cuisine": "13"}, {"name": "Asian Bar", "city": null, "cuisine": "12"}, {"name": "Mexican Bar", "city": null, "cuisine": "11"}, {"name": "Italian Bar", "city": null, "cuisine": "10"}, {"name": "American Bar", "city": null, "cuisine": "9"}, {"name": "Chinese Bar", "city": null, "cuisine": "8"}, {"name": "French Bar", "city": null, "cuisine": "7"}, {"name": "Middle Eastern Bar", "city": null, "cuisine": "6"}]
```

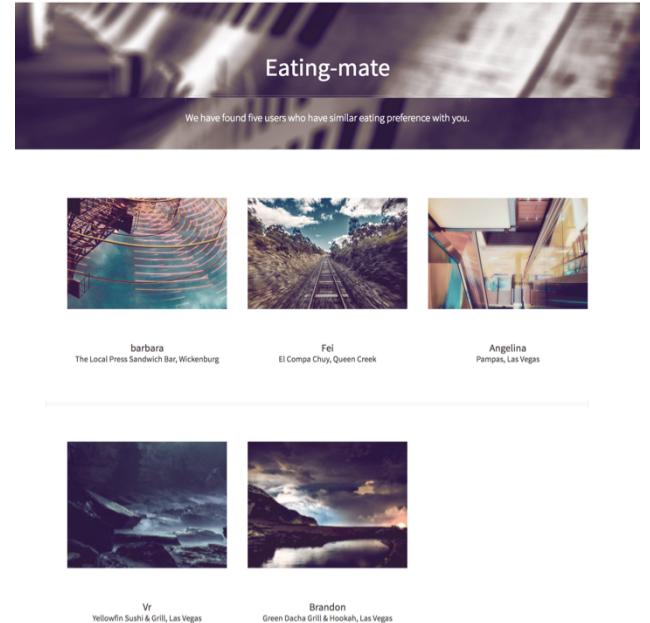
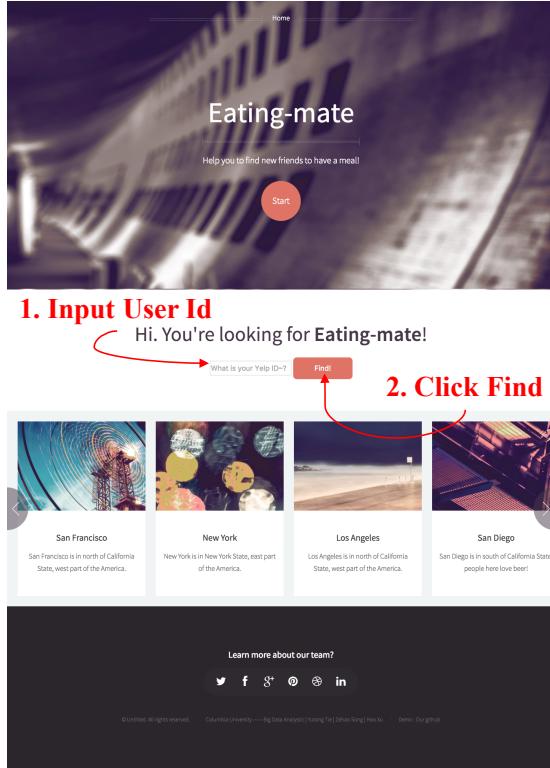
Each line is a recommendation for a particular user, including user\_id, list of recommended restaurants, list of eating mates and taste analysis.

## 5. Web App

Path: web/web.py

Start the Flask Server by “python web.py”. Note, python 2 is required.

```
(py2) SongZehao-MacBook-Air:web Zehao$ ls
cache.jl          data      recommendation.pyc static    web.py
(py2) SongZehao-MacBook-Air:web Zehao$ python web.py
running on 127.0.0.1:8000
Threaded True
 * Running on http://127.0.0.1:8000/ (Press CTRL+C to quit)
|
```



The analysis of the user's taste. The result is shown in word cloud.



Recommended eating mates and restaurants.

## VI. EXPERIMENT RESULTS

### 1. Collaborative Filtering (CF)



As the rank grows, the train error decreases and the time cost increases. Although the long training time is acceptable since we only need to train one time, we should choose a relatively small rank. Because the number of reviews provided by new users must larger than the rank, otherwise the Linear Regression will fail. So. **Rank 20** is chosen.



Number of iteration of ALS is the second parameter to be tuned. The larger the iterNum, the lower the train error & the higher the time cost. Since we only need to train once, we prefer to chose a relatively high numIter. Thus, *numIteration* = 20 is chosen.

### Builde CF Model based on training data

```
In [8]: trainModel = cf_als(train_ratings, ALS_setting)
Training CF Model
Time Cost: 67.94s
```

```
In [9]: cf_als_eval(trainModel, train_ratings)
Evaluating CF Model
Mean Squared Error = 0.202851587955
Time Cost: 117.99s
```

```
In [35]: saveModel(trainModel, "full")
Number of Business: 85538
Number of User: 685557
Time Cost: 14.18s
```

The user and business feature matrixes are saved as pickle file for future usage.

```
FOLDERS
▼ target
  ▷ final
    ▷ vBusiness
    ▷ vUser
  ▷ full
    ▷ vBusiness
    ▷ vUser
  ▷ train
    ▷ vBusiness
    ▷ vUser
```

## 2. Linear Regression (LR)

```
In [20]: vu, err = getPredictedUserFeature(9, test_ratings, vBusiness)
vu
```

```
Out[20]: array([ 0.85343739, -1.01137154, -1.61861862, -0.8111291, -1.14654225,
   -0.84237122, -0.03597422,  0.3464966, -0.56901162, -0.25857191])
```

### 3. K-NN & Recommend

Get the eating mates of a particular user (*vu*)

```
In [22]: getKNN(vUser, vu, 1000)
Out[22]: [(111900, 0.5708298747401556),
(470260, 0.5805114246678903),
(95447, 0.6244813045550115),
(299306, 0.62680960603461477),
(169270, 0.63099199680208717),
(358501, 0.63812444296578785),
(172936, 0.64374111614503959),
(546037, 0.64735737398127591),
(1070, 0.65324326022194978),
(475477, 0.65700562683601915),
```

Eating mate user id, similarity & recommended shared reataurants.

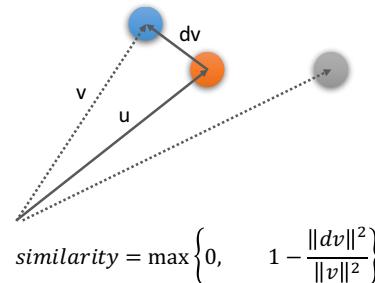
```
In [17]: evalEatingMateRecomd(9, test_ratings, trainModel)
111900 79% [70360, 76034, 54231, 82374, 75712, 57194]
470260 78% [70360, 78534, 83881, 76034, 83340]
95447 77% [70360, 82374, 54231, 76034, 26317]
299306 77% [70360, 83340, 82374, 75712, 54231]
169270 77% [70360, 76034, 54231, 83340]
358501 76% [70360, 82374, 54231, 83340, 57194, 78534, 26317]
172936 76% [70360, 26317, 78534, 57194, 83340, 82374, 54231]
546037 76% [70360, 57194, 83340, 78534, 76034, 26317, 83881]
1070 76% [54231, 78534, 70360, 76034, 82374]
475477 76% [70360, 57194, 76034, 26317]
Time Cost: 35.28s
MSE of User Vector: 1.11
```

### Evaluation:

The evaluation of recommendation system is hard. In real world, companies could use A/B test to compare different recommendation algorithm. Obviously, we don't have such circumstance.

Let's see whether the recommended eating mates is the real friend of the original user.

### Concept of Similarity



From the previous result, we could see that the recommended eating mates has a similarity larger than 70%.

If we check the real friends of the users, we could plot the histogram of their similarity between the base user.

**Count the similarity between user and his friends**

```

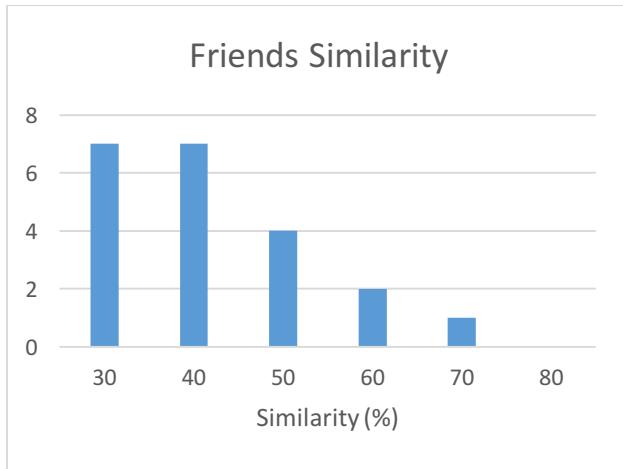
vu, err = getPredictedUserFeature(9, test_ratings, vBusiness)
sims = []
for fid in friends:
    if fid <= 1000: continue
    vf = getUserFeature(vUser, fid)
    dv = np.linalg.norm(vu-vf)
    sims.append( int(100 * getSimilarity(dv, vu)) )
print sims

[32, 39, 58, 11, 17, 13, 20, 0, 0, 5, 0, 28, 49, 19, 0, 47, 38, 0, 0, 29, 0, 0, 0, 11, 61, 0,
 4, 0, 3, 11, 24, 11, 0, 0, 0, 0, 8, 25, 5, 0, 0, 28, 20, 0, 18, 0, 0, 36, 0, 0, 36, 0, 0,
 0, 0, 16, 16, 17, 38, 0, 0, 48, 4, 0, 34, 0, 0, 0, 1, 0, 0, 0]

hist = np.histogram(sims, np.arange(0, 110, 10))
hist

(array([44, 11, 7, 7, 4, 2, 1, 0, 0, 0, 0]), array([ 0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100]))

```



It seems that the similarity between friends are not as high as the recommended eating mates. It's reasonable since in the model we only consider the ratings of businesses. But this result also indicates our naive CF recommendation system is not good enough. We should take into account more other features of users and businesses, like location and price. A simple solution is to recommend a potential restaurants pool of large size (eg. 1000) and then use other algorithm to filter the results. Constrained by the time, we didn't implement this idea.

#### 4. Taste Analysis

```
In [85]: recomd, users, taste = getTaste(9)
Time Cost: 13.86s

In [86]: json.dumps(taste)

Out[86]: "[{'Pizza": 117}, {"Mexican": 114}, {"Sandwiches": 108}, {"American (Traditional)": 78}, {"Italian": 77}, {"Fast Food": 74}, {"American (New)": 69}, {"Chinese": 65}, {"Japanese": 51}, {"Burgers": 45}, {"Cafe": 45}, {"Mediterranean": 38}, {"Sushi Bars": 35}, {"Delis": 35}, {"Thai": 31}, {"Seafood": 30}, {"American Fusion": 27}, {"Coffee & Tea": 27}, {"Steakhouse": 23}, {"Bedeque": 23}, {"Gyro": 23}, {"Hibachi": 22}, {"Bistro": 21}, {"Salad": 19}, {"Mexican": 17}, {"Indian": 18}, {"Greek": 17}, {"Pork EK": 17}, {"Vegan": 17}, {"French": 17}, {"Vegetarian": 17}, {"Latin American": 16}, {"Hot Dogs": 16}, {"Juice Bars & Smoothies": 15}, {"Sports Bars": 14}, {"Hawaiian": 14}, {"Vietnamese": 14}, {"Food Trucks": 14}, {"Gluten-Free": 13}, {"Desserts": 13}, {"Specialty Food": 13}, {"Pub": 13}, {"Ice Cream & Frozen Yogurt": 10}, {"Cocktail Bars": 10}, {"Korean": 9}, {"Wine Bars": 9}, {"Etnic food": 8}, {"Gourmet": 8}, {"Souf Food": 7}, {"Caribbean": 7}, {"Soup": 7}, {"Mexican": 7}]
```

After getting the list of categories frequency in the following format,

```
[  
    ["Category Name", 30]  
]
```

Generate the word cloud in front end via WordCloud2.



## VII. CONCLUSION

### Conclusion.

In this paper, we covered the entire workflow of our project, dive deep into the algorithms we chose to use, the dataset we collected in our own hand, the process of hyper parameter selection and evaluation of algorithms.

To sum up, we have built an eating mate recommendation system successfully. The system only needs some piece of information about user's reviews, then it would generate a feature vector of that user, find the nearest neighbor and finally output the eating mate that is recommended. We also deployed a site that supports eating mate recommendations. The result of recommendation is a combination of eating mate and a list of restaurant that could satisfy both users. According to several test results, the recommendation is reasonable and future customer feedback is to be collected.

### Contributions.

Major works for each team member are listed as follow:

Hao Xu (hx2208):

- 1). Implement a script to get business ids from Yelp with Yelp search API.

2). Preprocess the data collected by Scrapy and map each user and business into our preferred ids.

3). Final Presentation (Slides & Present).

4). Final documentation.

Zehao Song (zs2324):

1). Implement python crawler via Scrapy to collect data from Yelp by using business ids collected by Yelp API.

2). Project proposal presentation in front of the class.

3). Design and develop algorithms for recommendations and run hyper parameter experiments and evaluate the result.

4). Deploy a server that supports the recommendation.

5). Final documentation.

6). Video Demo

Yutong Tie (yt2549):

1). Design PowerPoint for project proposal presentation.

2). Develop a front end display that supports eating mate recommendation as well as word cloud display.

### Future works.

Due to the time limit, our project is not as perfect as we designed at the very beginning, therefore, several future works are needed to enhance the performance of our project:

1). More data

Currently, due to a limitation of the dataset we collected, we only support making recommendations for Las Vegas, San Francisco, Pittsburgh, Charlotte. In the future, we want to extend this project into the whole United States areas, which requires to collect review data all over the nation.

2). Scalability

While the size of dataset is increasing, we need to guarantee the performance of our algorithm, which should not be effected too much by the size of dataset and the running time of the whole process should not take too long, therefore, when the dataset is large, we should have method to deal with it, such as sampling, etc.

3). Dashboard

As we dive deep into the project, we found some interesting features other than eating mate could be derived from the dataset we collected. For example, city tastes. We could have a dashboard to display some interesting features, which could make our website more attractive.

4). Feedbacks

As mentioned in the final presentation, it is hard for us to evaluate the recommendation result. What we did in our experiment is to try to select the model that makes less training error. Therefore, we need some real data, which reflects how well our recommendation performs. We could collect feedbacks from users that whether they are satisfied about the recommendation and comments like how could we make improvements.

### **VIII. ACKNOWLEDGMENT**

We would like to thank Prof. Lin for this fantastic semester, the informative lectures and relax atmosphere.

We would like to thank each Teaching Assistant for their helps during office hours.

We would like to every team member for the productive ideas and supportive for each other.

### **IX. APPENDIX**

### **X. REFERENCES**

- [1] T. McNichol, "Word on the street," NYSE magazine, 2012. Retrieved from [http://www.nysemagazine.com/statics/NYSE\\_fall\\_2012.pdf](http://www.nysemagazine.com/statics/NYSE_fall_2012.pdf).
- [2] Yelp. (2016). *About us*. Retrieved from <https://www.yelp.com/about>.
- [3] Facebook. Available: <https://www.facebook.com>.

- [4] N. Vignesh, "How does Facebook's friend recommendation system work". Retrieved from <https://www.quora.com/How-does-Facebooks-friend-recommendation-system-work>
- [5] Liknedin. Available: <https://www.linkedin.com>
- [6] Amazon. Available: <https://www.amazon.com>
- [7] Spark Collaborative Filtering - RDD-based API. Available: <http://spark.apache.org/docs/latest/mllib-collaborative-filtering.html>