

The Yelp Machine

NYC Data Science Academy | Capstone Project | Team PC1

Aiko Liu | Amy Chen | David Steinmetz | Greg Domingo

The team



Aiko Liu

With quantitative training in math/physics, focus on the application of machine learning techniques to finance, big data and beyond



Amy Chen

Devoted to using data visualization and machine learning techniques for social and business innovations



David Steinmetz

Passionate about creating value by distilling data into actionable information, particularly through visualization



Greg Domingo

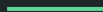
Keen interest in innovation with Data Science as one of the leading edges of innovation space

Agenda

Overview and Context

Explanation of the App

Next Steps

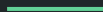


Agenda

Overview and Context

Explanation of the App

Next Steps



You want to go out to eat with a friend
but sifting through restaurant
listings is frustrating
and difficult



Lugo Cucina

4.0 ★★★★★ (45)

\$\$ · Italian · Pennsylvania Plaza

Madison Square Garden-area Italian cafe

Opens at 8:00 AM



Casa Nonna

4.4 ★★★★★ (90)

\$\$ · Italian · W 38th St

Italian dining in a spacious venue

Open until 9:30 PM



Uncle Jack's Steakhouse - Westside

4.0 ★★★★★ (83)

\$\$\$ · Steak · 9th Ave

Big steaks & traditional chophouse fare

Open until 11:00 PM



Club Bar & Grill

\$\$\$ · Grill · Pennsylvania Plaza

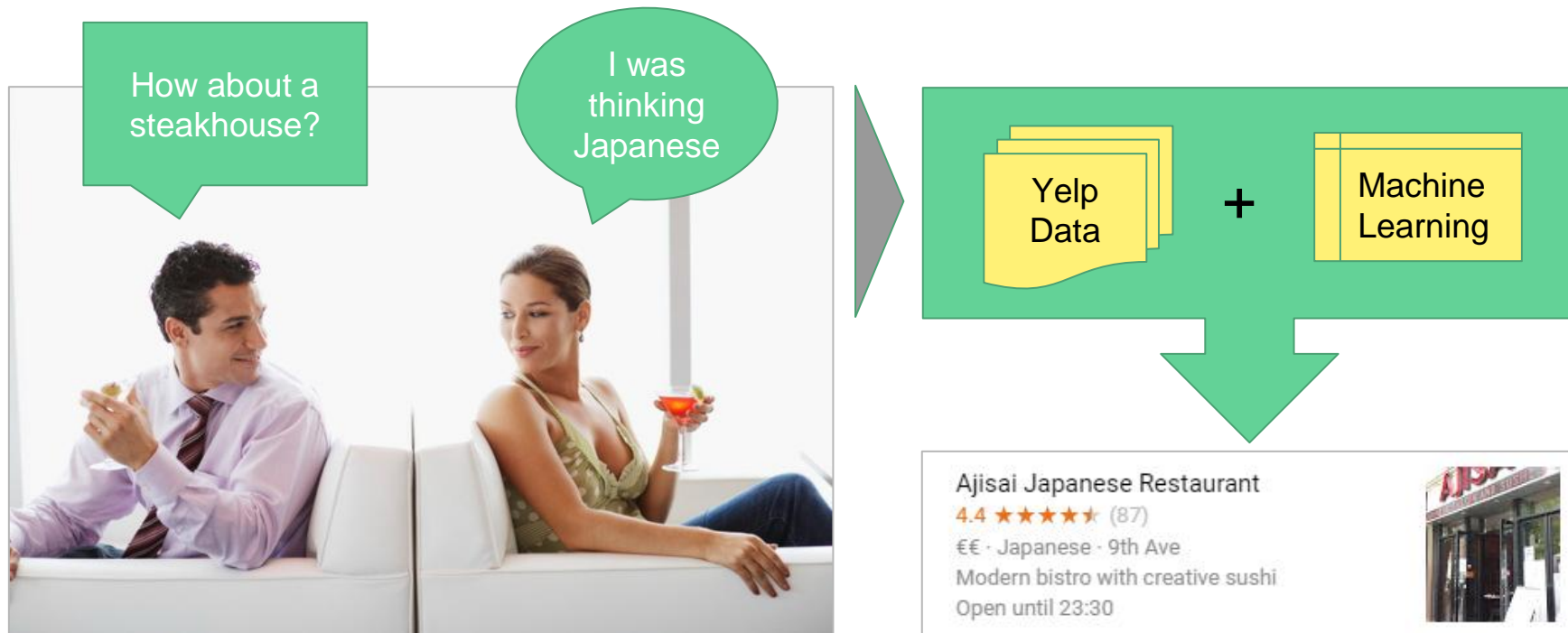
Elegant spot for drinks & American eats



Showing results 1 - 20



Could Yelp data be used in conjunction with machine learning to find a restaurant which will suit the tastes of two people?



**Amy** amy17519

A very picky Asian girl

**David** davidsteinmetz

He just wants protein for every meal

Cuisine

YELP NEARBY!

- | | | |
|---|---|---------------------------------|
| <input type="checkbox"/> Vietnamese | <input checked="" type="checkbox"/> Mexican | <input type="checkbox"/> Thai |
| <input type="checkbox"/> Japanese | <input type="checkbox"/> Italian | <input type="checkbox"/> Indian |
| <input checked="" type="checkbox"/> Chinese | <input type="checkbox"/> Seafood | <input type="checkbox"/> Pizza |
| <input type="checkbox"/> Steakhouses | <input checked="" type="checkbox"/> Burgers | <input type="checkbox"/> French |
| <input type="checkbox"/> American (New) | <input type="checkbox"/> American (Traditional) | <input type="checkbox"/> Greek |

Recommendations

Xi'an Famous Foods | Yelp It!

81 St Marks Pl | +1-212-786-2068

Rating: 4 | Price: 2

Taco Bandito | Yelp It!

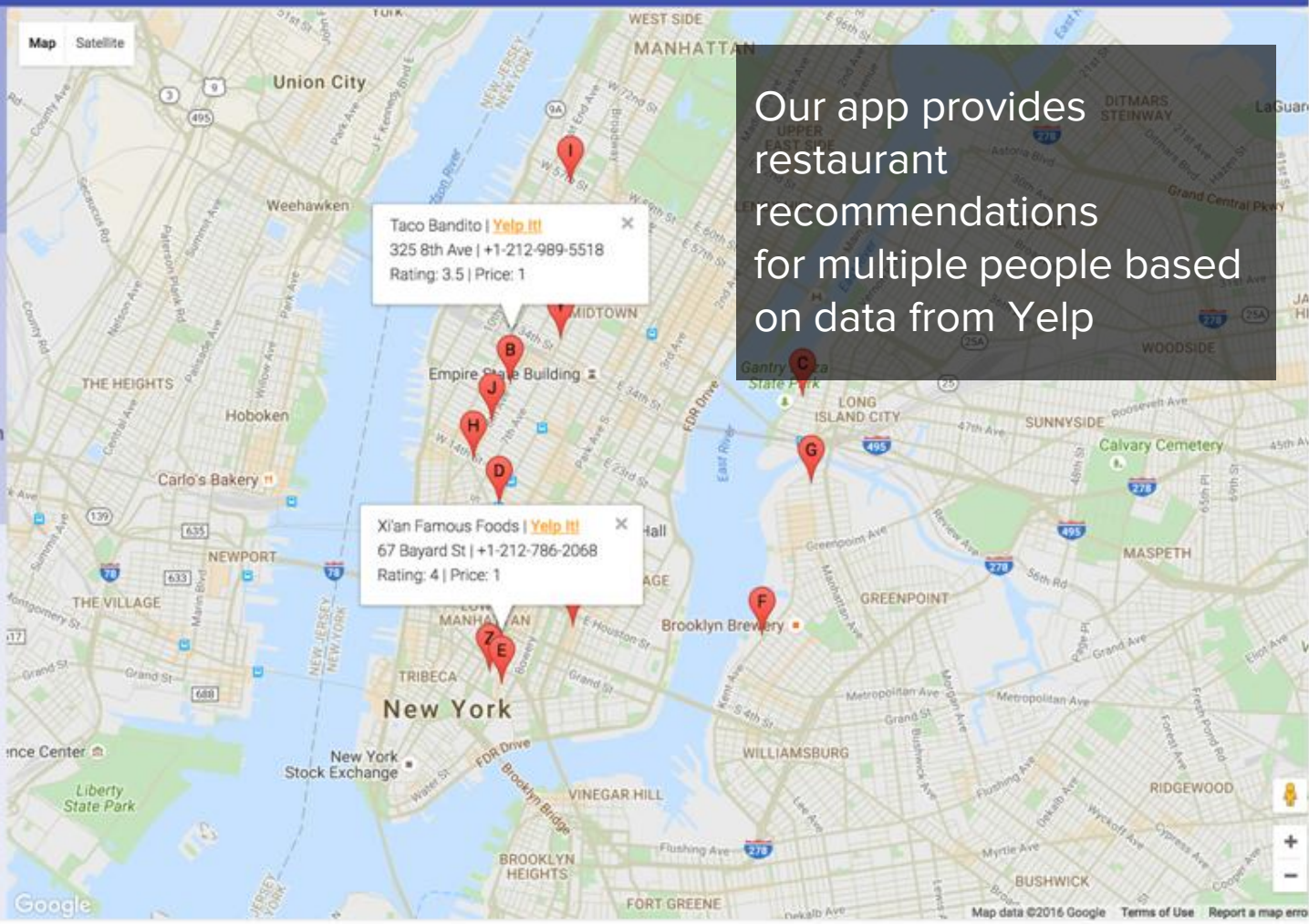
325 8th Ave | +1-212-989-5518

Rating: 3.5 | Price: 1

Skinny's Cantina | Yelp It!

4705 Center Blvd | +1-718-729-8300

Rating: 3 | Price: 3

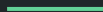


Agenda

Overview and Context

Explanation of the App

Next Steps



Our App marries a Flask front end with a Python back end to provide recommendations

Front end

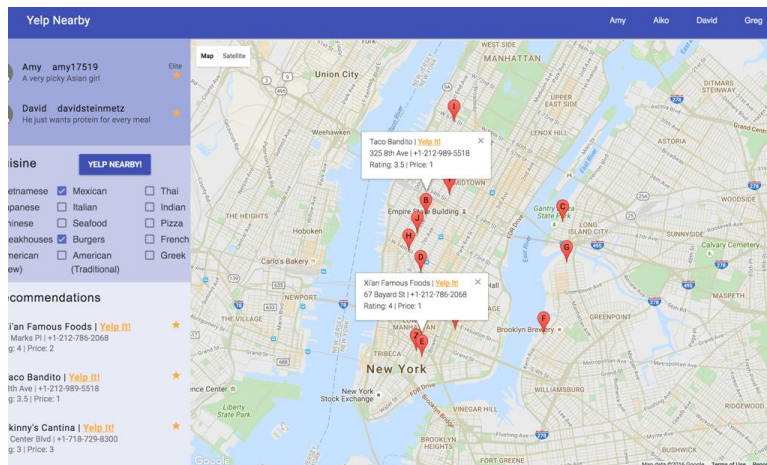
Flask
Python
Microframework

Yelp Nearby
A multiuser restaurant
recommendation engine

Jinja2
Templating

+

**HTML /
JavaScript**
Programming
Languages



Back end

Python
Classes,
Functions,
Modules

**Yelp
API**

Ext

**GraphLab
Machine Learning**

Internal

**Homemade
Class/Func/Mod**

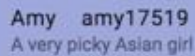
Two users log into our
app on the login page

YelpNearby

amy17519

davidsteinmetz

Login

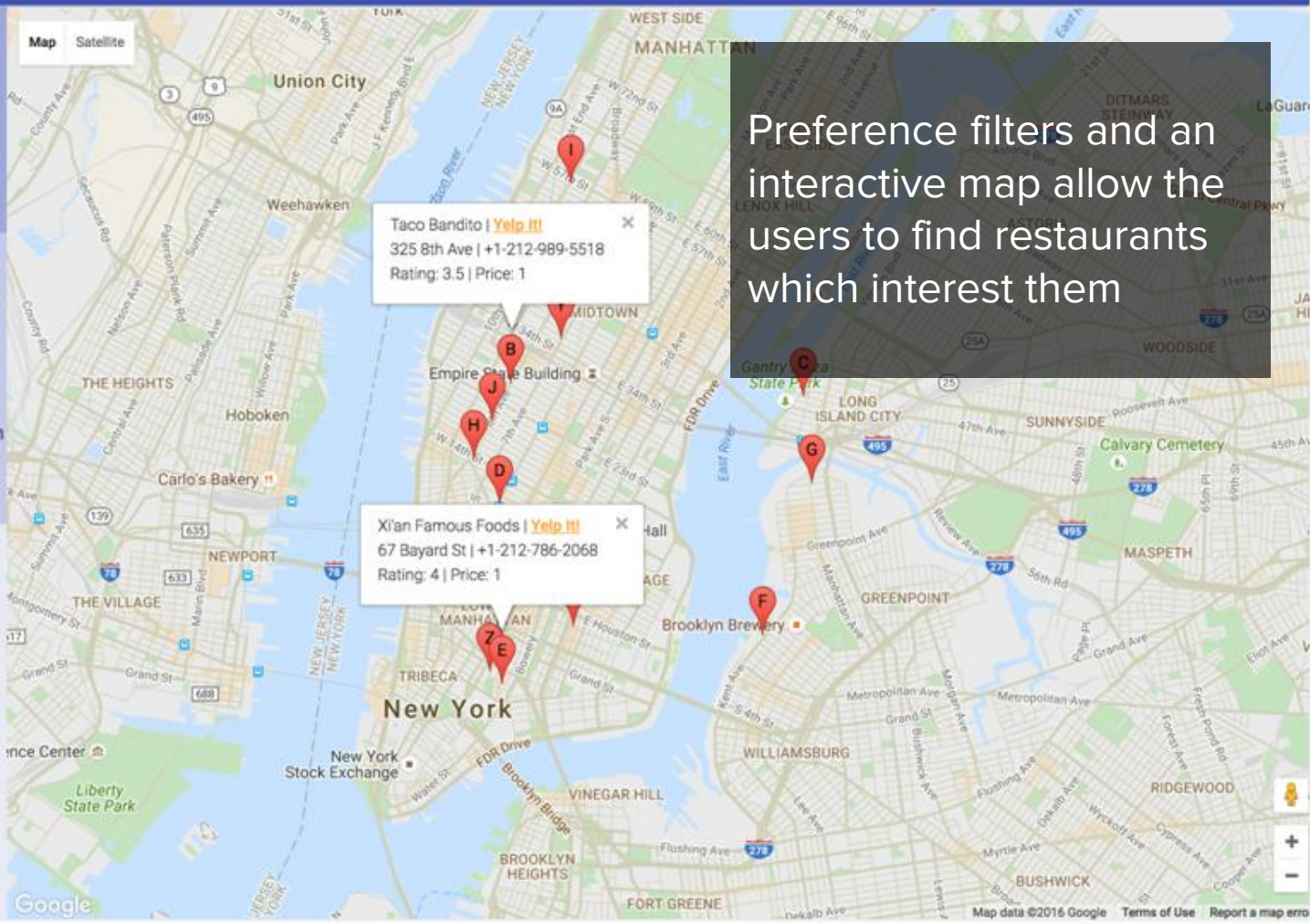


YELP NEARBY!

- | | | |
|---|---|---------------------------------|
| <input type="checkbox"/> Vietnamese | <input checked="" type="checkbox"/> Mexican | <input type="checkbox"/> Thai |
| <input type="checkbox"/> Japanese | <input type="checkbox"/> Italian | <input type="checkbox"/> Indian |
| <input checked="" type="checkbox"/> Chinese | <input type="checkbox"/> Seafood | <input type="checkbox"/> Pizza |
| <input type="checkbox"/> Steakhouses | <input checked="" type="checkbox"/> Burgers | <input type="checkbox"/> French |
| <input type="checkbox"/> American | <input type="checkbox"/> American | <input type="checkbox"/> Greek |
| (New) | (Traditional) | |

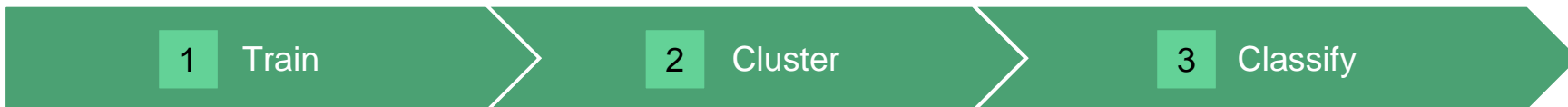
Xi'an Famous Foods | [Yelp](#) [It!](#)
81 St Marks Pl | +1-212-786-2068
Rating: 4 | Price: 2

- Taco Bandito | [Yelp It!](#)**
325 8th Ave | +1-212-989-5518
Rating: 3.5 | Price: 1
- Skinny's Cantina | [Yelp It!](#)**
4705 Center Blvd | +1-718-729-8300
Rating: 3 | Price: 3



Preference filters and an interactive map allow the users to find restaurants which interest them

The recommendation system works in a pipeline of three processes



Data about users,
business, reviews
from online Yelp
Challenge

Collaborative
Filtering
recommends
restaurants for
specific users

Clusters needed to
extend model to
new locales

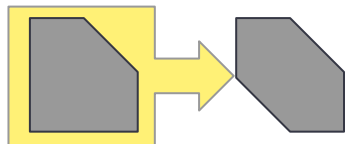
Density-based
scanning used to
create clusters

Restaurants are
clustered based on

Locally available
restaurants are
classified into the
clusters of the
predicted
recommendations

A collaborative filtering model was chosen because it incorporates information from users who make similar reviews

Content-based systems



Predicts similar items

Advantages

1. Uses the items' content to predict the user's interest
2. Recommendation quality improves as the review/item content data cumulates

Disadvantages

1. Impossible to predict the totally distinct types of items the particular user has never expressed interest in
2. Limited by the collected items' info in making recommendation (New Item?)

Collaborative filtering

User A

User B

Rest. 1



Rest. 2

Rest. 2

Rest. 3

Rest. 3



Rest. 4

Predicts items from user preferences and from similar users

Advantages

1. Predict items through similar user patterns, even if the particular user has a short review history
2. Works without item attributes
3. 'Outside the Box' recommendation

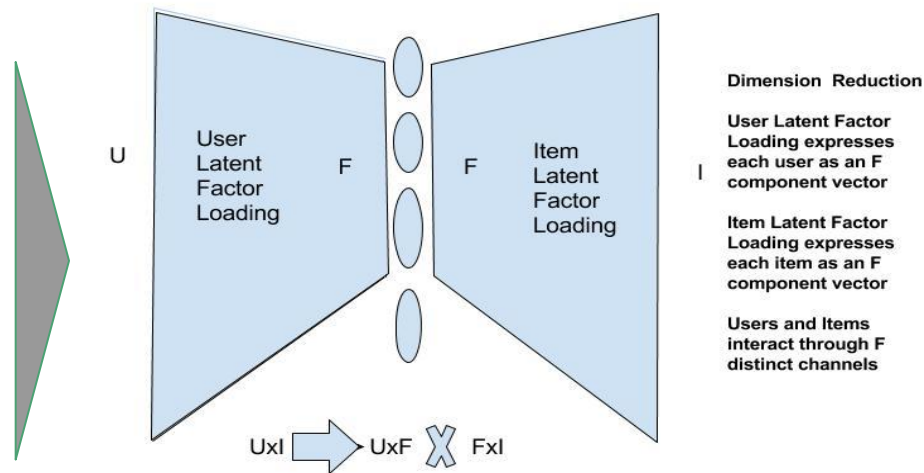
Disadvantages

1. Cold Start for the new users
2. Sparse Ratings on the same item
3. Recommendations are difficult for users with distinct tastes; these users are called black sheep or gray sheep.

Latent Matrix Factorization is the key component of collaborative filtering

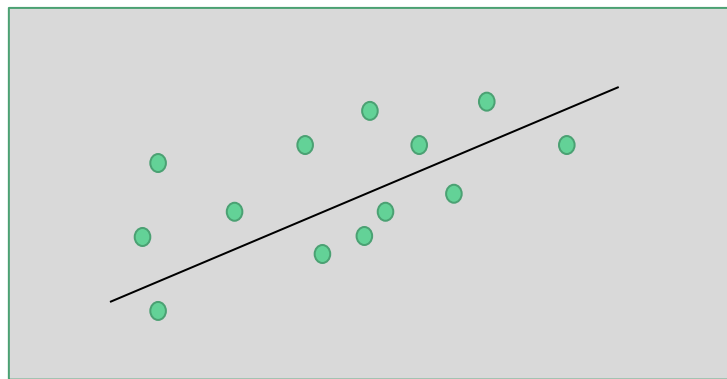
	Rest. 1	Rest. 2	Rest. 3	Rest. 4
User 1	5		2	
User 2		3	4	1
User 3	1			4

Numbers in the table are the rating the user gave the restaurant on a scale of 1-5

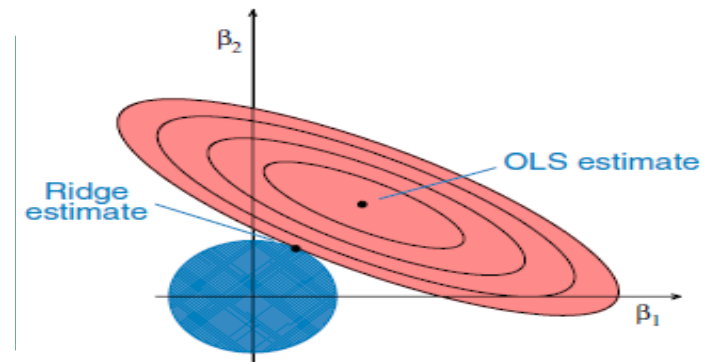


The matrix to the left is factorized

Latent Matrix Factorization adds to two well-known machine learning techniques: linear regression & L2 regularization



Linear regression



Ridge regression

Latent Matrix Factorization

1. Matrix Factorization, the core of CF, can work without side inputs from the users, items, capturing the user-item interaction through factorizing the sparse user-item matrix
2. The linear regression upon the side information reduces the model estimation residuals
3. The L2 regularization, known as Ridge regression in the context of MLR, controls the stability of the model fit and prevents over-fit
4. Three parts are combined into the single equation system (graphlab)

Additional features can be included as side information in Latent Factorization to train the model

Feature Name	Feature Equation	Why it's included
User_EliteYears	$1 * \text{years_elite}$	Elite users have outsized influence on ratings
User_AvgRating	$\text{mean}(\text{rating})$	Different users have different rating standards
User_Num_Review	$\log(\text{u_num_reviews}+1)$	The indicator of the user's engagement on yelp
User_Location	city/state of the reviews	The reviews from the same location may be similar
Rest_AvgRating	$\text{mean}(\text{user rating})$	The reviews' consensus on the restaurant quality
Rest_Num_Review	$\log(\text{r_num_reviews}+1)$	The attention the business get from the reviewers
Rest_Aggr_EliteYear	Sum of the	The attention the business receives from the leaders among the reviewers
Rest_Location	city/state of the Rest.	The location of the restaurants are highly correlated with the residence of the reviewers

Making restaurant recommendations outside the limited region in the dataset posed a problem

Recommender system only maps to restaurants in the original dataset

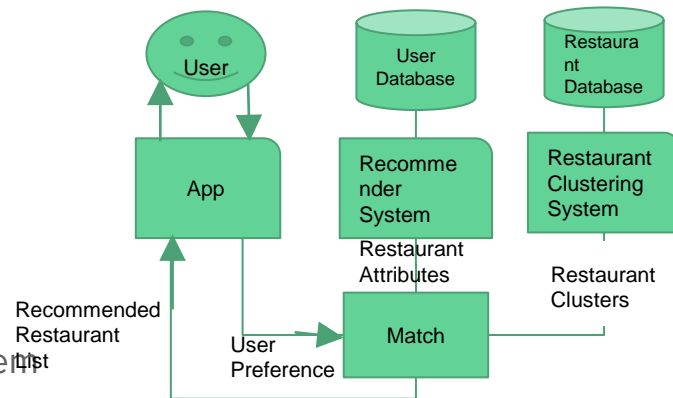
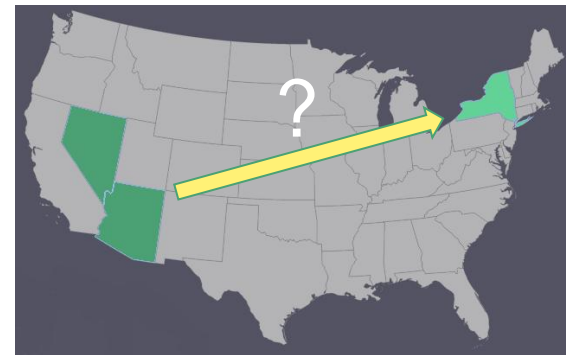
Original dataset does not include major cities like New York and San Francisco

So how do we recommend restaurants outside the areas in the dataset or in areas with very few reviews

Solution: Cluster Analysis

Solution Concept:

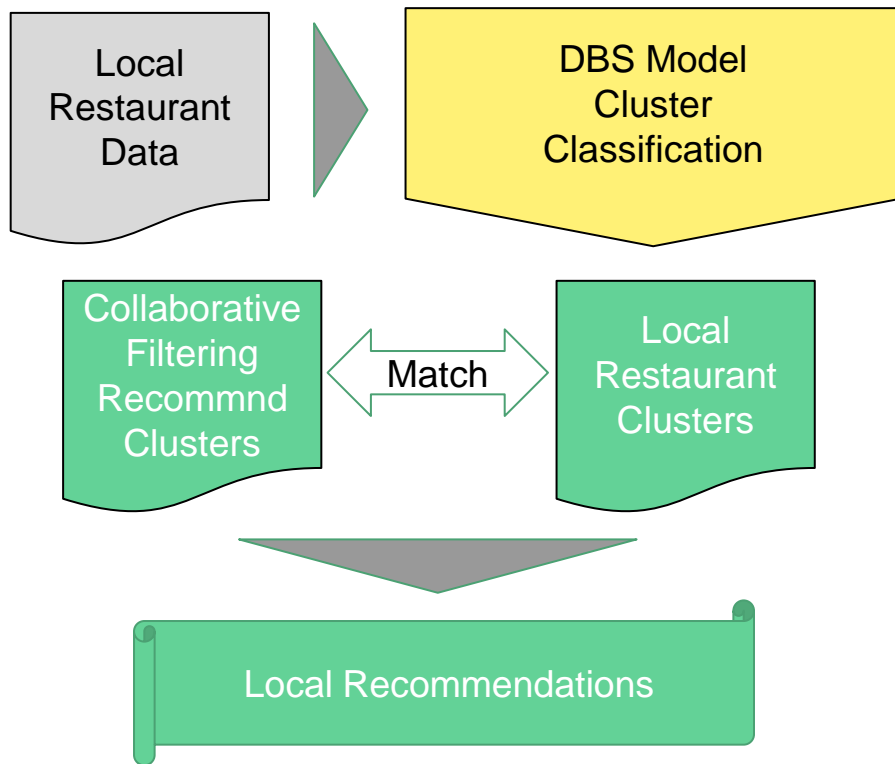
Get attributes of restaurants selected from recommender system



Density-based scanning was chosen to cluster all restaurants in the data set

Algorithm	Advantages	Disadvantages
K-Means	<ul style="list-style-type: none">• K-means works well when the shape of clusters are hyper-spherical• Computationally efficient	<ul style="list-style-type: none">• May give different results every time it is run• Requires prior knowledge of number of clusters
Hierarchical Clustering	<ul style="list-style-type: none">• Gives recommended clusters• Repeatable results	<ul style="list-style-type: none">• Time complexity is quadratic
Density-based scanning	<ul style="list-style-type: none">• Can handle clusters of different shapes and sizes• Gives recommended clusters• Computationally more efficient than hierarchical cluster method	<ul style="list-style-type: none">• May have problem handling high dimensional data• May have problem dealing with data that has widely varying densities

Classifying locally available restaurants based on the DBS model solved the problem of a limited data set



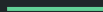
Feature used	Description
Ratings	Average Rating Of Restaurant Based on User Reviews
Price Range	Price Range For Restaurant
Review Counts	Number Of Reviews Of A Restaurant (transformed using log function)

Agenda

Overview and Context

Explanation of the App

Next Steps



The functionality of the app can be extended

- Extract information from the restaurant reviews using an NLP technique called Latent Dirichlet Allocation

- This data can be included in the clustering model to improve distinction between clusters

- Use app users' reviews to improve recommendations

- Include new users not existing in the data set

- Extend to larger groups of users

Thank you
for your attention

GraphLab Recommender Model

$$\text{score}(i, j) = \mu + w_i + w_j + \mathbf{a}^T \mathbf{x}_i + \mathbf{b}^T \mathbf{y}_j + \mathbf{U}_i^T \mathbf{V}_j,$$

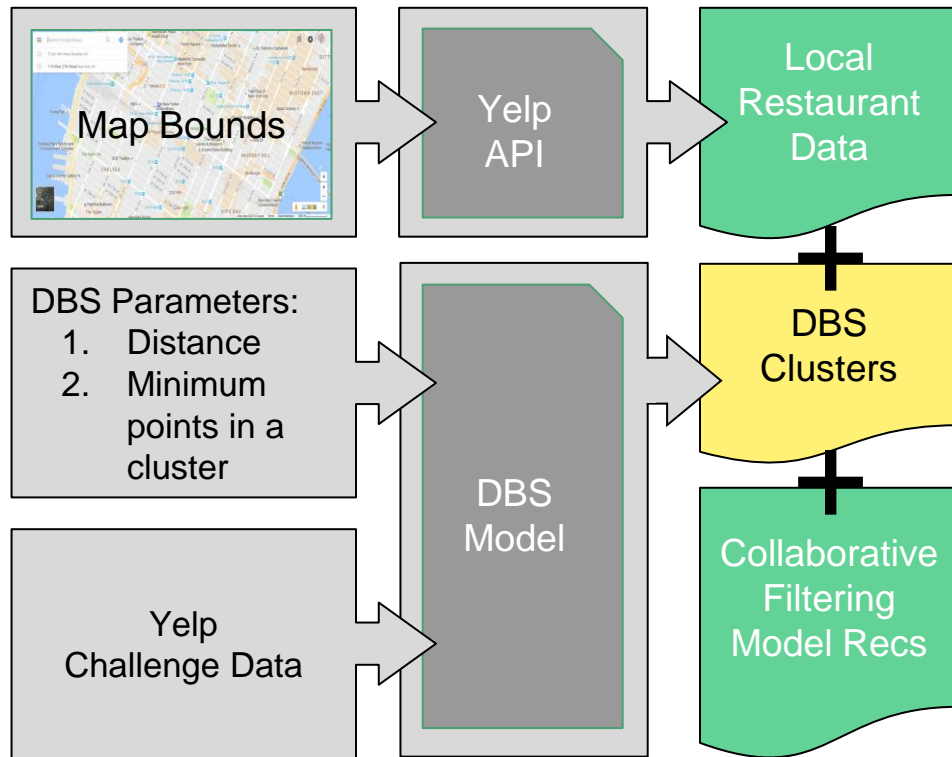
$$\text{Objective} = \min_{\mathbf{w}, \mathbf{a}, \mathbf{b}, \mathbf{V}, \mathbf{U}} \frac{1}{|\mathcal{D}|} \sum_{(i, j, r_{ij}) \in \mathcal{D}} \mathcal{L}(\text{score}(i, j), r_{ij})$$

$$+ \lambda_1 (\|\mathbf{w}\|_2^2 + \|\mathbf{a}\|_2^2 + \|\mathbf{b}\|_2^2) + \lambda_2 (\|\mathbf{U}\|_2^2 + \|\mathbf{V}\|_2^2)$$

\mathcal{L} = (Squared Error) Loss Function,

r_{ij} = rating of user i to item j .

Classifying locally available restaurants based on the DBS model solved the problem of a limited data set



Feature used	Description
Ratings	Average Rating Of Restaurant Based on User Reviews
Price Range	Price Range For Restaurant
Review Counts	Number Of Reviews Of A Restaurant (transformed using log function)

Local Recommendations