# Yelp Business Data Analysis: Feature Selection, Ratings Prediction, Sentiment Analysis

Nikhil Patil

# Problem Description

❖ Feature Selection

❖ Ratings Prediction

❖ Sentiment Analysis

❖ Geospatial Visualization of Ratings

# Dataset Description:

- ❖ Sub-datasets: business, user review, check-in. Total size: 5 GB.
- ❖ Business: 144K businesses with stars, attributes, categories
- ❖ E.g. parking availability, happy hour, drive thru, restaurants table service etc.
- ❖ Review: 4 M user comments, useful votes
- ❖ Check-in: check-in counts
- ❖ Challenges:
  - ➢ Data extraction
  - ➢ Integration of sub-datasets
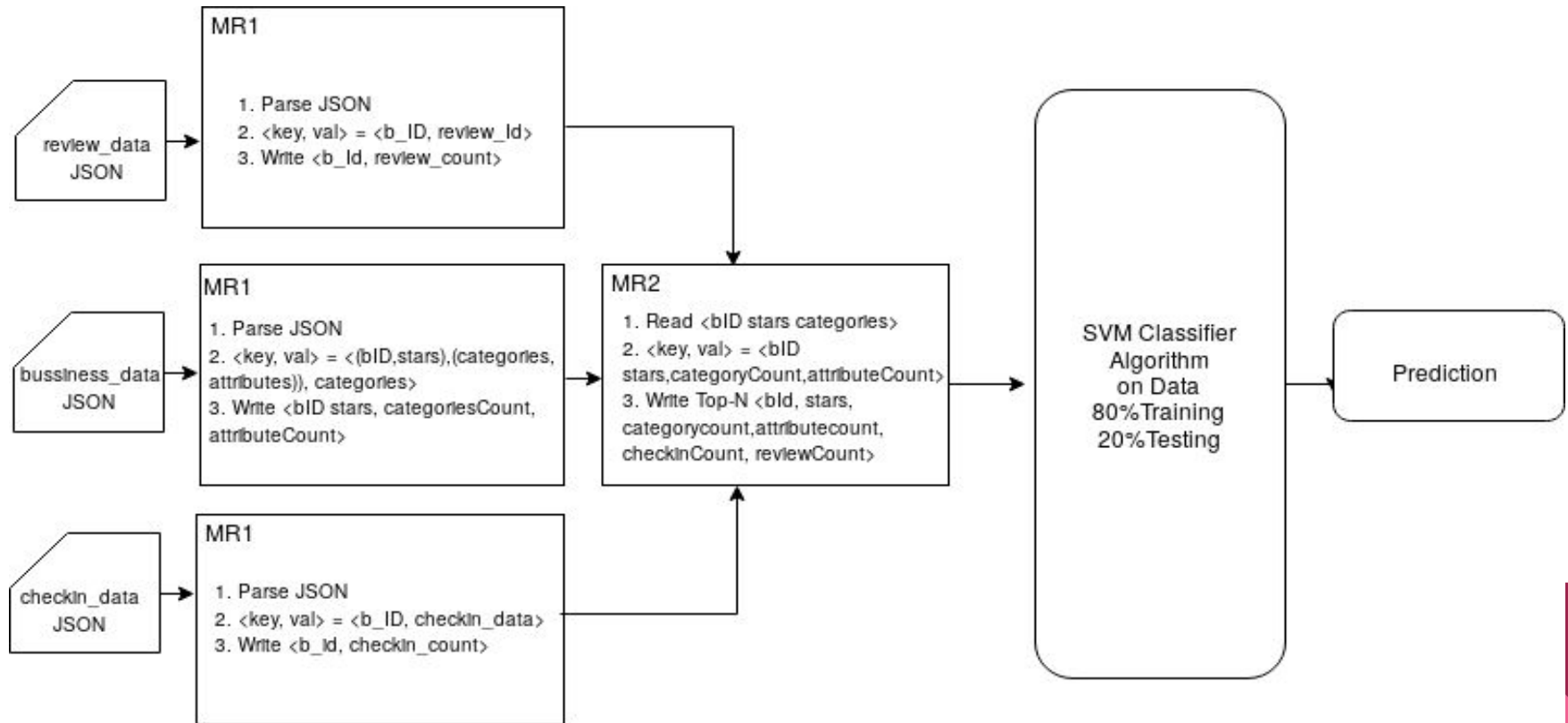  - ➢ Categorical features: quantification of True or False

# Feature Selection

❖ Map Reduce Jobs.

❖ Frequently occurred attributes, categories in high star rated businesses: Suggestion to new businesses.

❖ User review count, available attributes count, available categories count, total check in counts and stars for each business.

❖ Machine learning approach to predict stars, based on these features.

```
BusinessAcceptsCreditCards    103200;   BikeParking        51914
BusinessParking        49767;            RestaurantsTakeOut        45544
GoodForKids    44115;                    RestaurantsGoodForGroups  40887
WheelchairAccessible  34799;             GoodForMeal        25093
RestaurantsTableService       24423;     Ambience           21965
HasTV          21233;                    OutdoorSeating     19510
RestaurantsReservations       17378;     Caters   14938
ByAppointmentOnly      14201;            RestaurantsDelivery       10162
BestNights     5632;                     HappyHour          5473
AcceptsInsurance       5255;             Music    4850
DogsAllowed    3089;                     DriveThru          2154
GoodForDancing         1940;             HairSpecializesIn         1079
CoatCheck      1011;                     RestaurantsCounterService 246
BusinessAcceptsBitcoin        178;       DietaryRestrictions       155
Corkage        140;                        BYOB    47
Open24Hours    29}
```

# Architecture: Feature Selection & Ratings Prediction

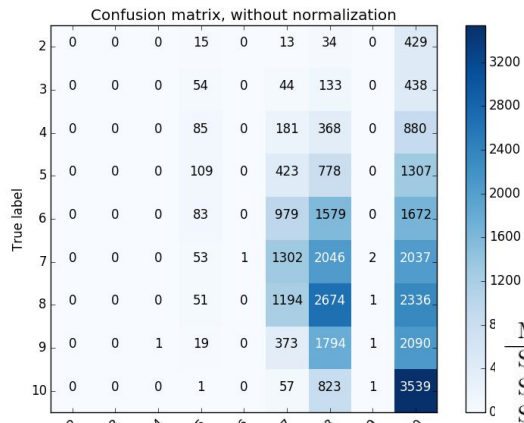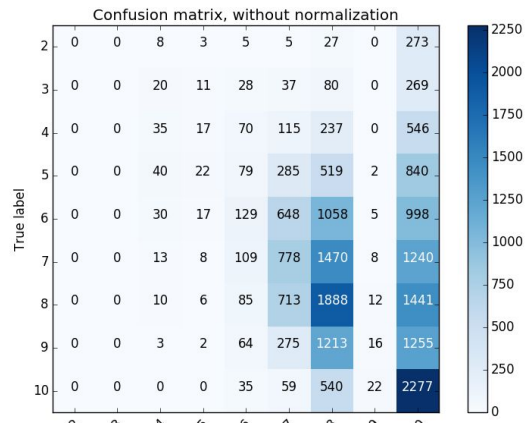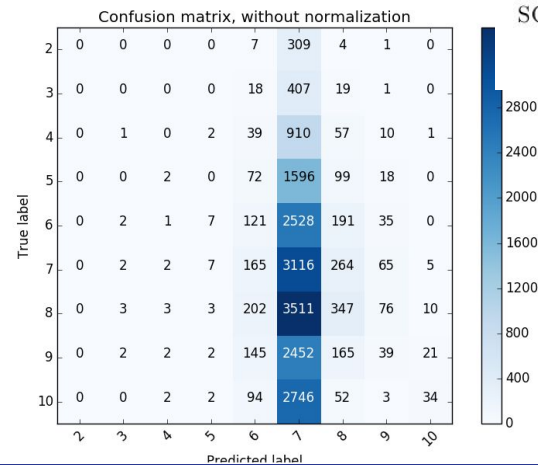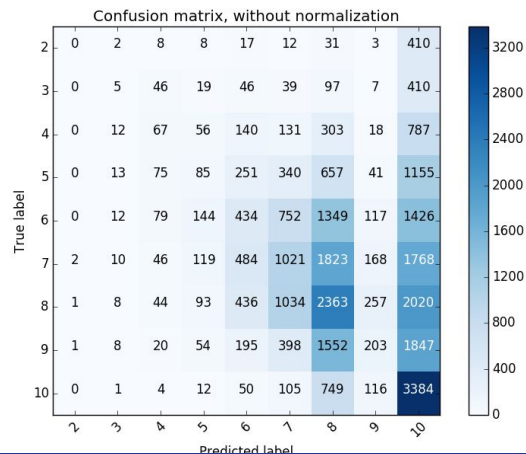# Ratings Prediction using Classification Algorithms

❖ Feature Spaces (X_train, test):
  ➢ Count-based (4): [business_id, review count, attribute count, category count, checkin count]
  ➢ Categorical (1 or 0) (31): [BusinessParking, HasTV, BYOB, etc.]
❖ Labels (y_train, test): Star ratings. [0,0.5,1,1.5,2,2.5,3,3.5,4,4.5,5] -> [0,1,2,3,4,5,6,7,8,9,10]
❖ Training Set: 75% (~110K); Test Set: 25% (~30K)
❖ Algorithms:
  ➢ Regression: linear regression
  ➢ Classification: SVM_linear, SVM_rbf, SGD

# Evaluation of Analysis:



❖ Precision = TP/(TP+FP)

❖ Recall = TP/(TP+FN)

❖ F1 = 2*P*R/(P+R)

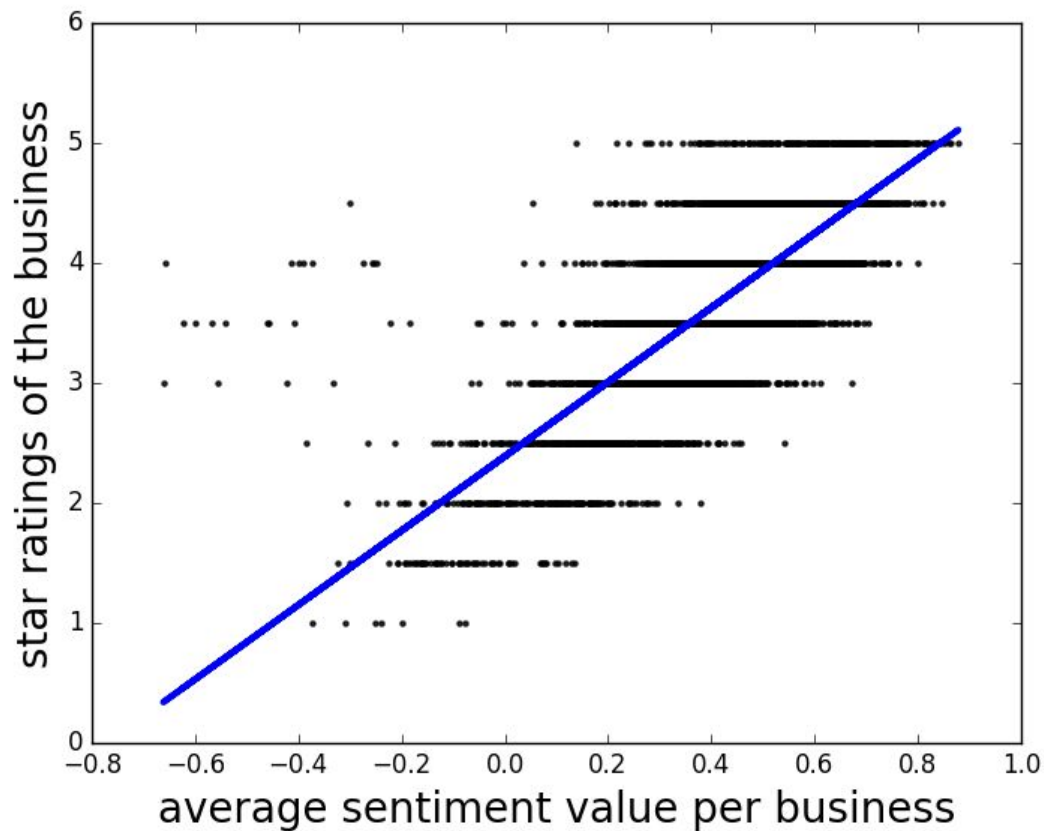| Model | Mean Precision | Mean Recall | Mean F1 score |
|---|---|---|---|
| SVM_linear (C=1) | 0.24 | 0.26 | 0.19 |
| SVM_rbf (C=1,g=1/31) | 0.19 | 0.25 | 0.18 |
| SVM_rbf (C=1000,g=0.3) | 0.23 | 0.25 | 0.20 |
| SGD | 0.21 | 0.20 | 0.12 |

Table 3:  Precision and Recall score for the models

# Sentiment Analysis of Reviews on Business

❖ Sentiment score for each review

❖ Sentiment score = (n_Positive - n_Negative)/(n_Positive + n_Negative)

❖ Priority based sorting of users - stars, sentiment score, useful votes

❖ Top-N user suggestion.

❖ Average sentiment score for each business.

❖ Linear regression based star prediction using average sentiment score of business.
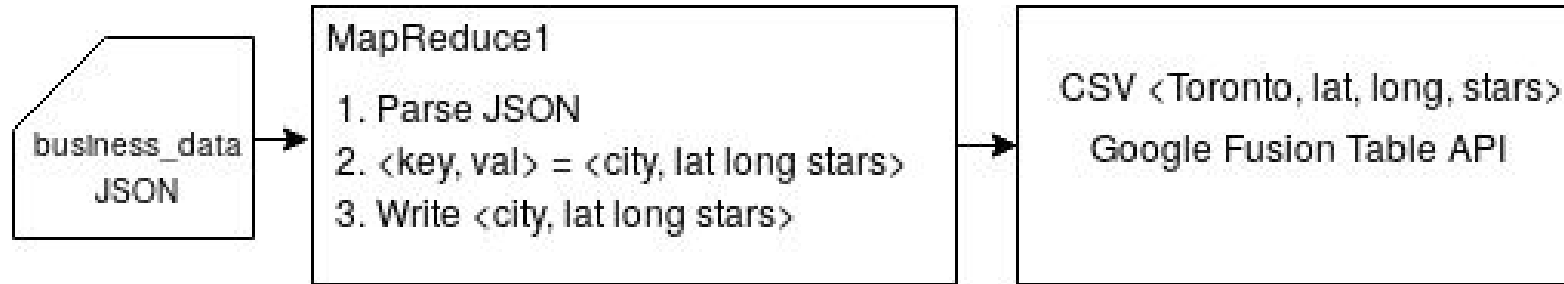
❖ Mean Square Error = 0.23

| BusinessId | Stars | Sentiment Value |
|---|---|---|
| -2ToCaDFpTNmmg3QFzxcWg | 1.5 | -0.047114883 |
| --DaPTJW3-tB1vP-PfdTEg | 3.5 | 0.55616015 |
| -AtzcXIwEP6yO7rM9CM9ww | 4.0 | 0.54652196 |
| -AxKgZHxyV-oBBHNyOESAg | 5.0 | 0.66816974 |
| HaRN97OjSnUJlHk21QAIXw | 1.5 | -0.0071489513 |
| -B8uga7IGEQijKERiwuz7A | 3.5 | 0.42080826 |
| 3quHwBY8XLalsx9m8e6Xuw | 5.0 | 0.60437346 |
| 3rJJeMOkVi_wEura5UCQqw | 3.0 | 0.38342175 |

# Geospatial Visualization of Ratings

- Visualization to reveal trends in star ratings within a city
- Parse business_data in MapReduce with 'city' as key
- Get all Business_id, Latitude, Longitude, and Stars
- Use Google fusion table API for heatmap

business_data
JSON

MapReduce1

1. Parse JSON
2. <key, val> = <city, lat long stars>
3. Write <city, lat long stars>

CSV <Toronto, lat, long, stars>

Google Fusion Table API

# Conclusion

- We did frequency based feature selection for attributes of business.
- Ratings prediction formulated as a 11-classification problem.
- To improve predictions:
  - Need better features
  - Need to optimize parameters for classification algorithms
- Sentiment Analysis: Unigrams-based sentiment results in positive correlation between sentiment value and the star ratings.
- Visualization helps to decide what areas, streets in a city have higher density of 4.5 or 5 stars.

# Questions?

Thank You.