

# BIZTRENDER

**Search & Visualization on  
Yelp Dataset**

**INSTRUCTOR : Prof James Abello**

**TA : Mr Fangda Han**

**GRADER : Mr Nishant Kumar**



**GROUP 9:**

**ABHISHEK CHAITANYA PRANEETH**  
**ab2083 cd817 pt357**





# THE DATASET

Public data from Yelp Dataset  
Challenge

<https://www.yelp.com/dataset>

Raw data uncompressed – 8.6 GB

Business.json – 192, 609 records

Checkin.json – 161, 950 records

Review.json – 6, 685, 900 records

Tip.json – 1, 223, 094 records

User.json – 1, 637, 138 records

Data at rest

---

# DATA SNAPSHOT

## BUSINESS

```
"business_id" : "6KgGE8B1RsR7jc9R5nuH0Q",
"name" : "Ruby Tuesday",
"address" : "4 E University Dr",
"city" : "Tempe",
"state" : "AZ",
"postal_code" : "85281",
"latitude" : 33.422192,
"longitude" : -111.939615,
"stars" : 2.5,
"review_count" : 9,
"is_open" : 0,
"attributes" : {
  "RestaurantsTakeOut" : "True",
  "RestaurantsPriceRange2" : "2",
  "OutdoorSeating" : "False",
  "BusinessAcceptsCreditCards" : "True",
  "Alcohol" : "'full_bar'"
},
"categories" : "American (Traditional), Restaurants",
"hours" : null
```

# DATA SNAPSHOT

CHECKIN

```
{ "_id" : ObjectId("5e7e48038a7a0519ca17e838"), "business_id" : "--KCl2FvVQpvjzm  
ZSPyviA", "date" : "2011-07-29 16:53:35, 2011-08-25 00:19:18, 2012-04-20 21:15:3  
9, 2012-06-07 23:27:31, 2012-07-01 03:42:35, 2012-07-07 17:51:12, 2012-08-08 21:  
01:48, 2012-09-01 05:46:24, 2012-09-08 16:13:40, 2012-09-19 15:25:58, 2012-09-19  
21:00:07, 2012-10-10 21:14:56, 2012-10-18 22:22:03, 2012-11-09 00:49:39, 2012-1  
1-09 01:01:09, 2012-12-09 20:02:49, 2013-01-25 20:37:21, 2013-02-02 23:30:26, 20  
13-02-06 18:49:12, 2013-02-08 00:30:04, 2013-05-08 22:25:18, 2013-05-26 18:03:13  
, 2013-07-23 21:55:47, 2013-09-30 01:09:59, 2013-10-07 01:58:06, 2013-11-03 21:2  
8:48, 2013-11-11 00:23:15, 2013-11-27 02:22:58, 2014-01-11 20:36:30, 2014-03-25  
21:13:10, 2014-07-13 19:22:07, 2014-07-25 02:39:57, 2014-08-01 18:18:10, 2014-09  
-26 19:26:08, 2014-10-14 02:10:07, 2014-10-29 22:42:55, 2014-11-03 22:28:15, 201  
4-11-03 22:38:18, 2014-11-12 23:18:24, 2014-11-23 01:34:28, 2014-12-28 03:25:39,  
2014-12-29 19:20:53, 2015-01-03 01:39:55, 2015-01-19 13:59:53, 2015-01-25 00:54  
:54, 2015-01-27 18:01:20, 2015-02-03 23:42:25, 2015-02-14 22:43:09, 2015-04-15 2  
1:39:15, 2015-04-17 23:08:28, 2015-05-03 00:56:38, 2015-06-09 21:53:07, 2015-06-  
20 18:48:15, 2015-06-23 21:20:17, 2015-08-13 15:25:38, 2015-08-22 23:49:29, 2015  
-08-24 00:26:18, 2015-08-30 02:41:54, 2015-10-10 19:24:41, 2015-10-14 01:35:36,  
2015-10-16 20:35:57, 2015-11-02 18:11:38, 2015-11-05 22:04:20, 2015-11-15 21:49:  
41, 2015-11-18 21:29:38, 2015-12-12 17:22:05, 2016-01-04 22:49:06, 2016-01-25 00  
:07:44, 2016-02-26 17:07:59, 2016-04-17 01:56:05, 2016-04-26 19:21:53, 2016-05-1  
4 02:27:43, 2016-05-28 17:36:58, 2016-07-09 19:49:57, 2016-07-23 22:52:28, 2016-  
07-25 00:08:46, 2016-08-04 16:10:09, 2016-08-21 17:32:39, 2016-10-31 16:54:14, 2  
017-01-14 04:50:47, 2017-03-12 19:55:15, 2017-04-29 00:08:44, 2017-05-25 18:23:3  
0, 2017-05-28 02:45:12, 2017-05-29 22:51:23, 2017-06-18 20:35:01, 2017-06-23 23:  
28:38, 2017-07-23 00:54:54, 2017-08-13 04:32:54, 2017-09-11 03:30:24, 2017-10-14  
21:24:05, 2017-11-11 19:40:30, 2017-11-25 00:21:50, 2018-02-13 22:17:27, 2018-0  
5-11 20:15:30, 2018-06-19 21:29:11, 2018-06-23 23:03:05, 2018-08-12 00:11:11, 20  
18-08-25 18:29:50, 2018-09-08 18:42:06" }
```



# DATA SNAPSHOT

## REVIEW

```
{
  "_id" : ObjectId("5e7e48408a7a0519ca1a60c9"),
  "review_id" : "2TzJjDVDEuAW6MR5Vuc1ug",
  "user_id" : "n6-Gk65cPZL6Uz8qRm3NYw",
  "business_id" : "WTqjgwHLXbSFevF32_DJvW",
  "stars" : 5,
  "useful" : 3,
  "funny" : 0,
  "cool" : 0,
  "text" : "I have to say that this office really has it together, they are so organized and friendly! Dr. J. Phillipp is a great dentist, very friendly and professional. The dental assistants that helped in my procedure were amazing, Jewel and Bailey helped me to feel comfortable! I don't have dental insurance, but they have this insurance through their office you can purchase for $80 something a year and this gave me 25% off all of my dental work, plus they helped me get signed up for care credit which I knew nothing about before this visit! I highly recommend this office for the nice synergy the whole office has!",
  "date" : "2016-11-09 20:09:03"
}
```

```
{
  "_id" : ObjectId("5e7e4ab88a7a0519ca806599"),
  "user_id" : "8zXgNCKkusOHMwh5Tj_8yCQ",
  "business_id" : "qrSsS0pk7SL67MP5nN8tlg",
  "text" : "More parking in the back of the restaurant",
  "date" : "2013-12-22 05:10:12",
  "compliment_count" : 0
}
```

**TIP**

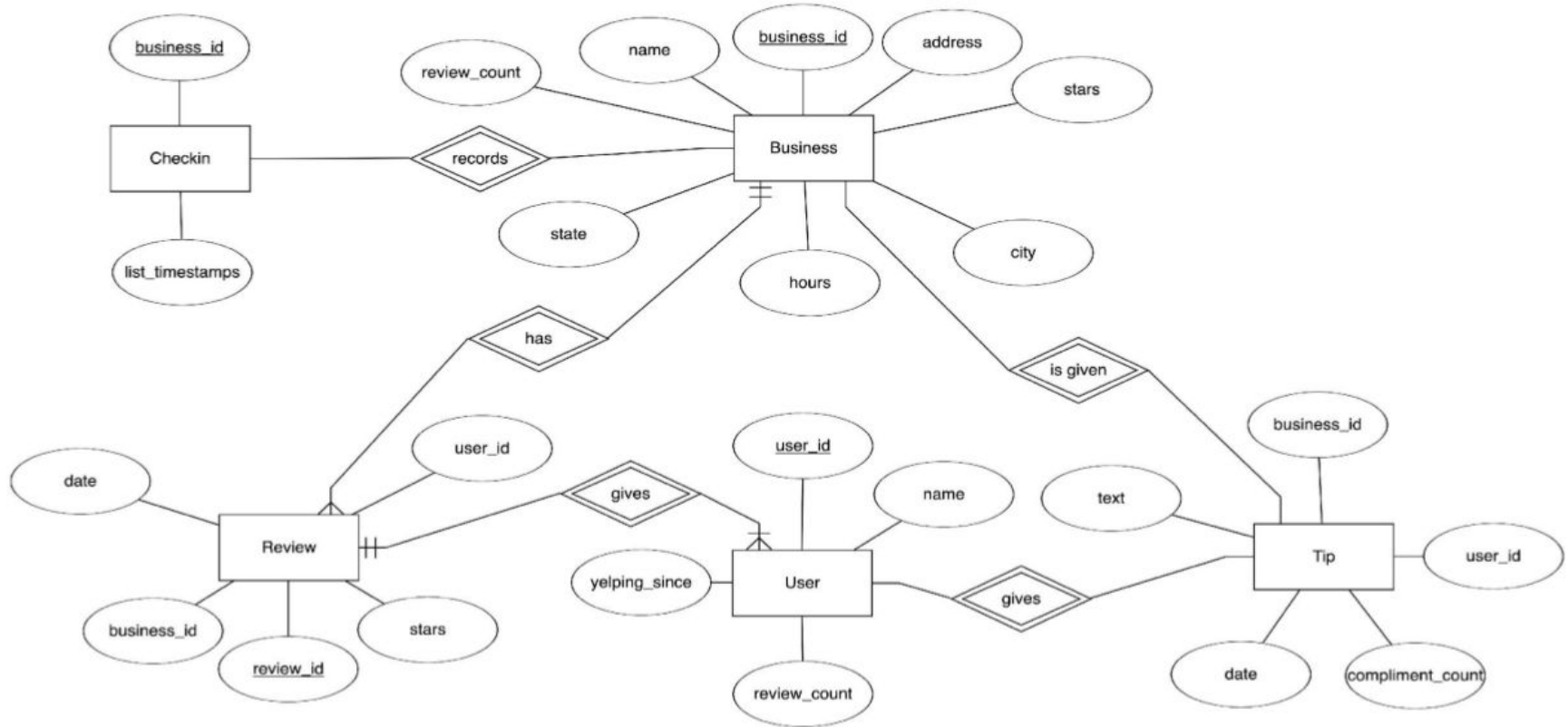
# DATA SNAPSHOT

## USER

```
{
  "id": "ObjectID(\"5e7e4b038a7a0519ca930f4a\")",
  "user_id": "l6BmjZMeQD3rDxWUbiAiw",
  "name": "Rashni",
  "review_count": 95,
  "yelping_since": "2013-10-08 23:11:33",
  "useful": 84,
  "funny": 17,
  "cool": 25,
  "elite": "2015,2016,2017",
  "friends": "c78V-rj8NQcQj0I8KP3UEA, aLRMgPcngYSCJ5naFRBz5g, ajcnq75Z5xxkvUSmmJ1bCg, BSMAMP2-wMzCkhtFq9ToNg, jka10dk9ygX76hJG0gfpZQ, du
t0e4xvme7Q5les0yCHQA, l4l5lBnK356zBua7B-UJ6Q, 0HicM00s-M_gl2e0-zES4Q, _uI57wL2fLyftrcSFpfSGQ, T4_Qd0YwbC3co6WSMw4vxg, iBRoLWptWmsI1kdbE9ORSA, x
jrUcid6Ymq0DoTJELkYyw, GqadWVzJ6At-vgLzK_SkGA, DvB13VjBmSnbFXBVBsKmdA, vRP9nQkYTeNi0dJtxZLVhg, gT0A1iN3eeQ08EMAjJhwQtw, 6yCWjFtp_AD4x93Wawmnw,
IdKzpNnib-JlViKv8_Gt5g, 3Bv4_JxHXq-gVLOxYMQX0Q, ikQyfuiviYh8T0us7wiFQ, fiGGltNaB7K5DR1jF3dOmG, tgeFUCHlh7v8bZfVl2-hjQ, -9-9oyXlqsMG2he5xIwDLQ,
Adj9fBPVJad8vSs-mIP7gw, Ce49RY8CKXvSfifxRYFTsw, M1_7TLi8Cbda89nFLH4iW, wFsnv-hqbW_F5-IrqfBN6g, 0Q1L7zXHocaUZ2gsG2XJeg, cBFgmOCBdhYa0xoFEAzp_g
, VrD_AgiFvzqtLR15vIr3SQ, cpE-7HK514Sr5vpSen9CEQ, FIUYelhPFB-zIKlt0ygIZg, CQAL1hvsLMCzuJf9AglsXw, 1KnY1wr15wFEWIRLB9IS6g, QWFQ-kXBilbid-lm5Jr3d
Q, nymT8liFugCrM16lTy0ZfQ, qj69bdd885heDvUPCyHd2Q, DysCZZcgbdrLHgEovk5y9w, LZMJIDuvht9Dy4KyquLXYa, b_9Gn7ws93AoPZPR0dIJQq, N07g1IaLh0_6suJtiSRe
4w, YdfPX_7DxSnKvvdCJ57i0w, 8GYryZPD22W7WgQ8kvMKEQ, cpQmAg0Watghp14h1pn1dQ, EnchhymLYMqftCRjqvVmw, -JdfKhFktE7Zs9BMDfCPeQ, uWhC9eof98zPkvSalga
qJw, eyTLNDdaiPatfe6mheIZ0g, VfHq0o73aKs0DvfAhWAQtg, kvD5tICngLaaQDujiSFwupA, dXacwEhqI9-3_XT6JeH00g, NfU0zDaTMEQ4-X9dbQWd9A, cTHWBdJ5KbctSUIvWs
gFwx, 3IEtCbSDF5t7RKZ20T6s9A, HJJXTrp6UybyYpQd9DA0JA, JaXogQFVjzGRAeBvzamBHg, NUonfKkj51iVqnNITtgXZg, D5vaJAYp0s0rGfsj9qvsMA, H27Ecbwuu4FGAllgI
Courw, 5SHrLmMiE4u8FWYwKNEoTw, Io36Y3xwQcIX9rYvPcyFXQ, J5mcqh8KxYpqjaLBNlwcig, -nTB3_08g06fD0GT8AtDBQ, wMpFA461lKh8oFns_5p65A, RZGFJHeomGJCWp3X
cL3eJa, ZoQ5zzXoSP1Rx0D4Amv9Bg, qzM0EB0SkuuGIFv0adJQAQ, HuM6vvuvuken-fPZ7d4oLA, H3oukHpGpn9n_mJwSDSQyQ, PkmsJsQ8FIZe8eh8c_u96g, wSByvbwME4MzgKJ
aFyfvNg, YEVqknoDmrHAoUbHX0nPnA, li3vsK1XAPmeJYAUTYfLHQ, MKC8yX10glbPYt00b4PECw, fQPH6W9fKsi27gkuUPnFaA, amrCMrDsoRetYFg2kwwdFA, 84dVQ6n6r2ezNa
Tuc7RkKA, yw9QjWY0olv5-urKv3t_Kw, 5XJDj7c3eoidfQ3jW18Zgw, tXSc6a6pIDctvwyBeu7Aqq, HFbbdCyyqP9xPKULxcIdg, hTUV5oh2do6Z30ppPuuiJA, gSqonG9J4fNM-
fl_FE71AA, pd9mgTFpBTg5F9x-MsczNg, j3VE22V2GcHi8UzxfLfw, NYXLMW-T-3V4Jqr4r-i0Wg, btXgAZedX8IWhMfA7Xkg, -Hp5mPLiRJNFnyeXSYgzag, P6-DwVg6-t2J
uQWtUEK0iQ, OI2TvxYvZrAodBG_RF53Xw, bHxf_VPKmZur1Bier-6A2A, Et_Sb39cVm81_Xe9HDM8ZQ, 5HwGL2UyYbaRq8aD6YC-fA, ZK228WMcCKLo5thcjd7rdw, lTf8wojwfm0
NOi7d0iz3Nw, bTYRXQYNJjpecfLNHtFH0A, Kgo42Fzpw_dXFgDKoewbtg, MNk_1Q_dq0Y3xxHZKe08VQ, AlwD504T9k0m5lkg3k5g6Q",
  "fans": 5,
  "average_stars": 4.03,
  "compliment_hot": 2,
  "compliment_more": 0,
  "compliment_profile": 0,
  "compliment_cute": 0,
  "compliment_list": 0,
  "compliment_note": 1,
  "compliment_plain": 1,
  "compliment_cool": 1,
  "compliment_funny": 1,
  "compliment_writer": 2,
  "compliment_photos": 0
}
```

# RAW DATA - ENTITIES

## ER DIAGRAM



# THE QUESTIONS

## USE CASE

### WHAT?

Discover trends in the dataset (business, ratings, reviews, visits, users). Support data-driven “human” decision making.

### WHY?

Customer vs Analyst. Platform for analysts.

### HOW?

Centralized datalake. Web application for search and visualization.



## Conceptual Flowchart

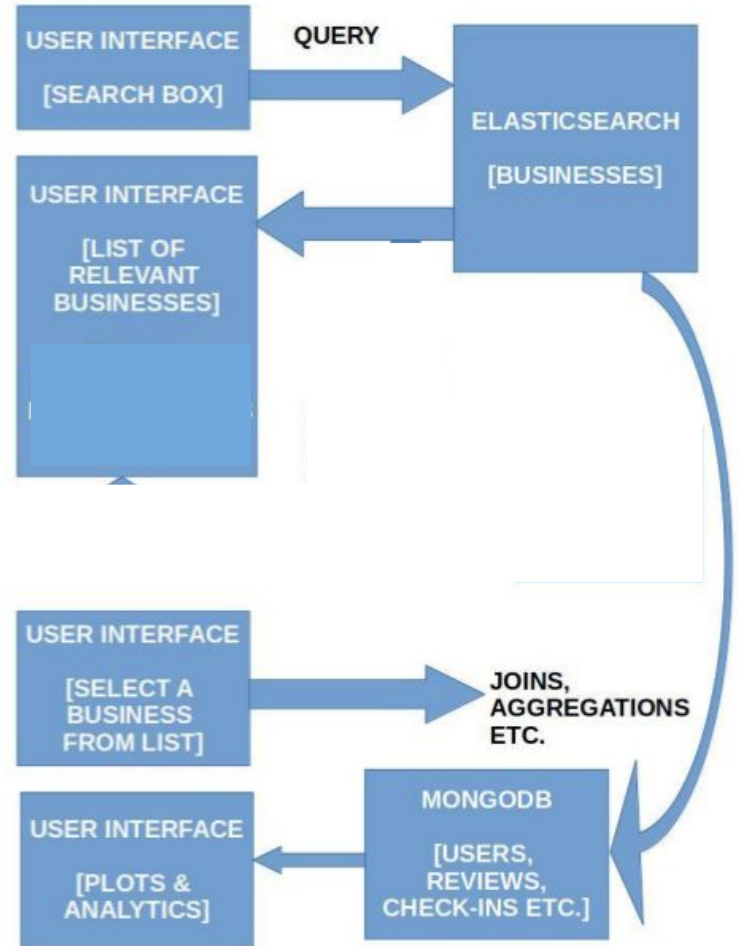
# PROPOSED SOLUTION

## Target Users

- BUSINESS ANALYSTS
- BUSINESS

## OWNERS

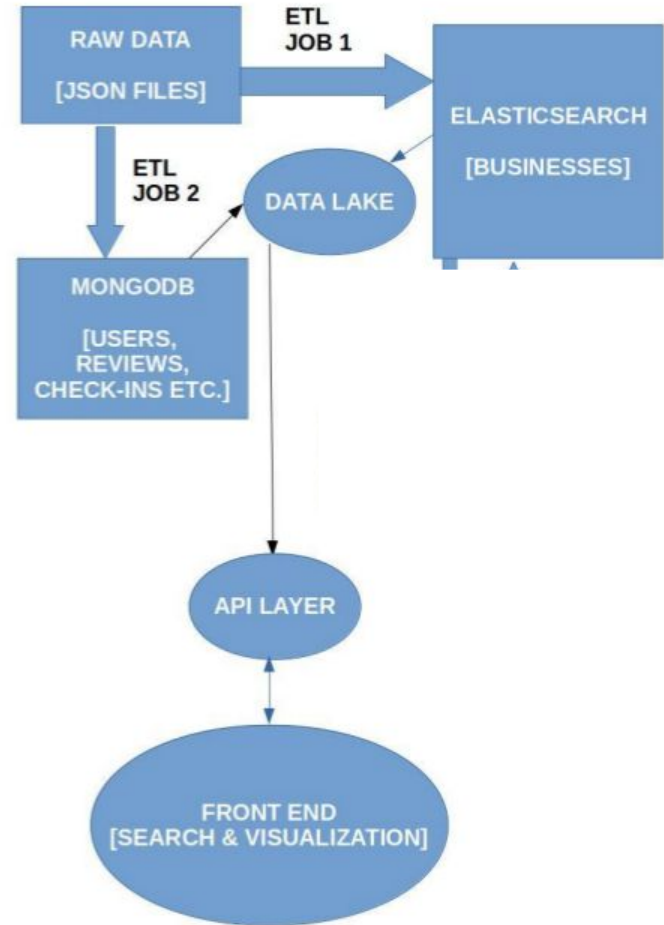
- REAL ESTATE  
CONSULTANTS



# MODE OF PROCESSING

- No pre-processing
- Data Lake : Elasticsearch + MongoDB
- DB Indexing
- Request Caching
- Streaming HTTP Response Body
- Filtering and Aggregation

## Data Lineage



# SOFTWARE STACK

## BACKEND

JAVA Spring Boot, MVC  
Architecture

## FRONTEND

Javascript, AngularJs,  
Highcharts

## SCRIPTING / COMMANDS

Bash, Mongo shell, Python

## DEPLOYMENT

—— Localhost, AWS EC2 & ES  
(experimental)

# VISUAL REPRESENTATIONS

- Use of standard UI fonts, sizes, color and texture palettes.
- Elements : Search box, click buttons, dropdown menus, unordered lists, navigation buttons, canvas divisions for plots.
- Interactive plots: 3D column chart with stacking and grouping, map (US and Canada) with latitudes/longitudes, packed bubble chart, dynamic column charts, dynamic time-series line chart, histogram, fixed placement column chart, and word cloud.
- Multi-colored labels and legends.



# INTERACTIVITY

- Bidirectional vertical sliders.
- Both textual and graphic representations for answered user queries.
- Mouse clicks, mouse hovering, mouse selection and their combinations.
- Menu driven plots and tabs.
- Mouse controlled zooming and panning.
- Hierarchical abstraction :  
Generic (trend summary) and specific (search) views.
- Built-in options for full-screen view and exporting plots.
- Different views linked through navigation using mouse clicks.

# KEY INSIGHTS

- Largest number of businesses in the dataset are either restaurants or food-based.
- Religious organizations is the only category that has no business with rating 1.
- For any category, most number of businesses are rated between 2 & 4.
- Smallest number of businesses in each category are the ones that are rated the worst at 1.
- All businesses lie in one of the clusters within these states/provinces - Alberta, Quebec or Ontario in Canada, & Nevada, Arizona, Illinois, Wisconsin, Ohio, Pennsylvania, North Carolina or South Carolina in the USA.

# KEY INSIGHTS

- Higher rating does not imply more visits for any business.
- Religious organizations receive the highest number of visits on Sunday, which is significantly higher than the visits received on any other days.
- Businesses receive more visits during late evening or early morning as compared to daylight hours.
- For poorly rated restaurants (1 to 2), the review sentiment is mostly negative with few positive scores in between.

# KEY INSIGHTS

- For middle-rated restaurants (2 to 3, the review sentiment is mostly centered (roughly symmetrical) about zero (neutral).
- For restaurants with good ratings (above 3), the review sentiment is mostly positive with few negative scores in between.
- For a specific business, the number negative reviews corresponds to the number of users rating the business below their average user rating.



# DEMO

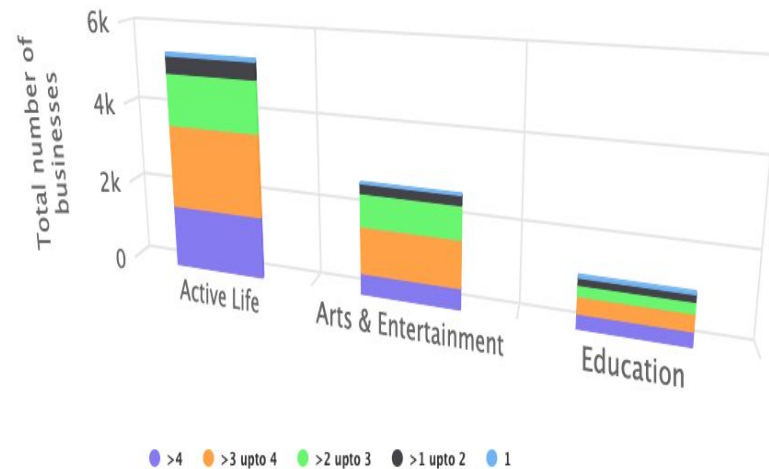
**VISUAL REPRESENTATIONS, INTERACTIVITY, INSIGHTS**

## APPLICATION VIEWS

### Wordcloud of Tips and Reviews



### Rating distribution by category





# APPLICATION VIEWS

Rating Distribution

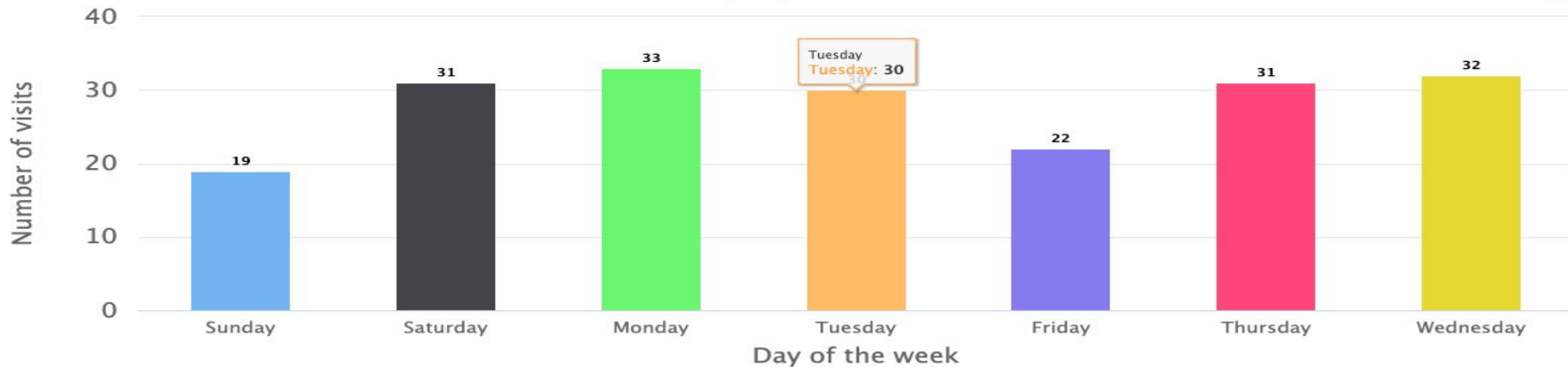
Geographic Distribution

Rating vs visits

Visits Distribution

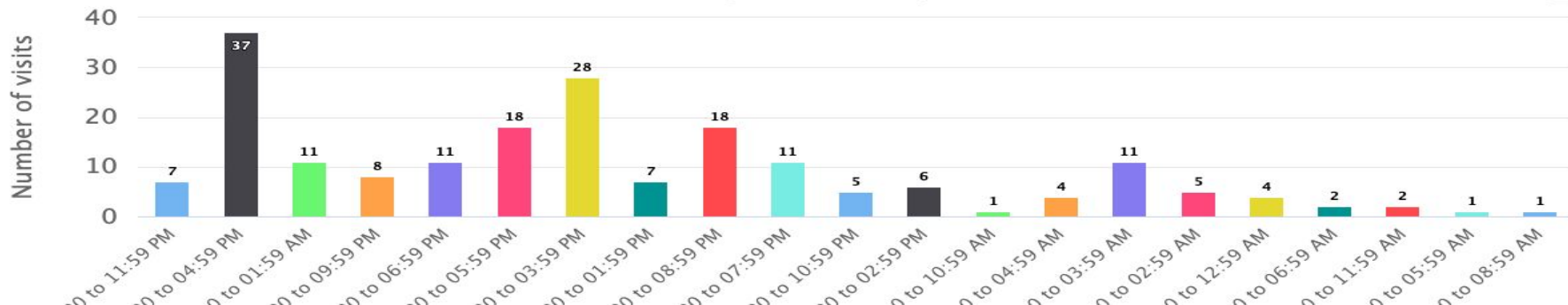
Review Sentiment vs Time

Visits by Day of the week



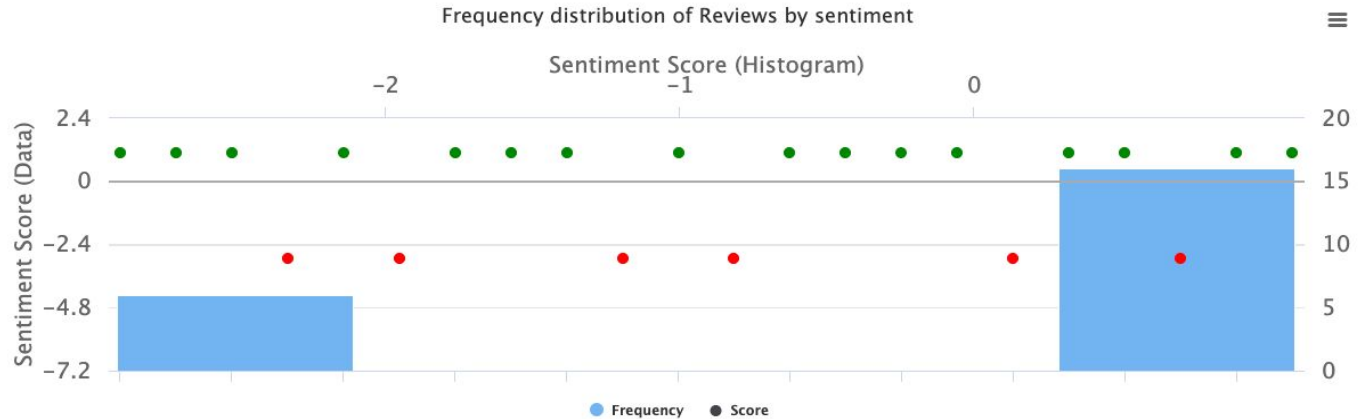
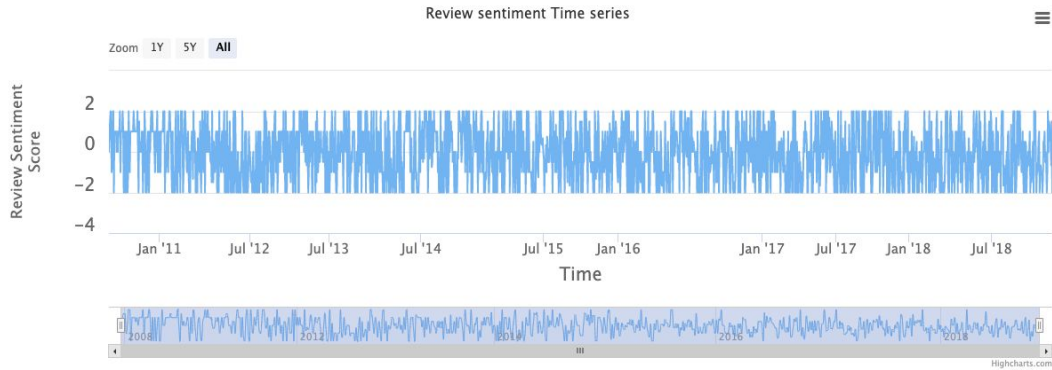
Highcharts.com

Visits by Hour of the day





# APPLICATION VIEWS

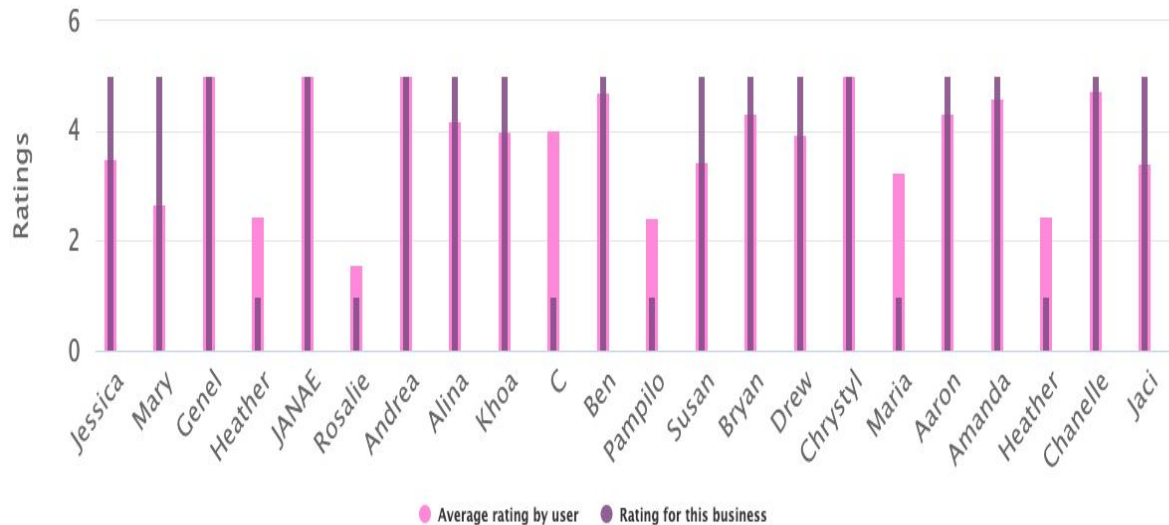


# APPLICATION VIEWS

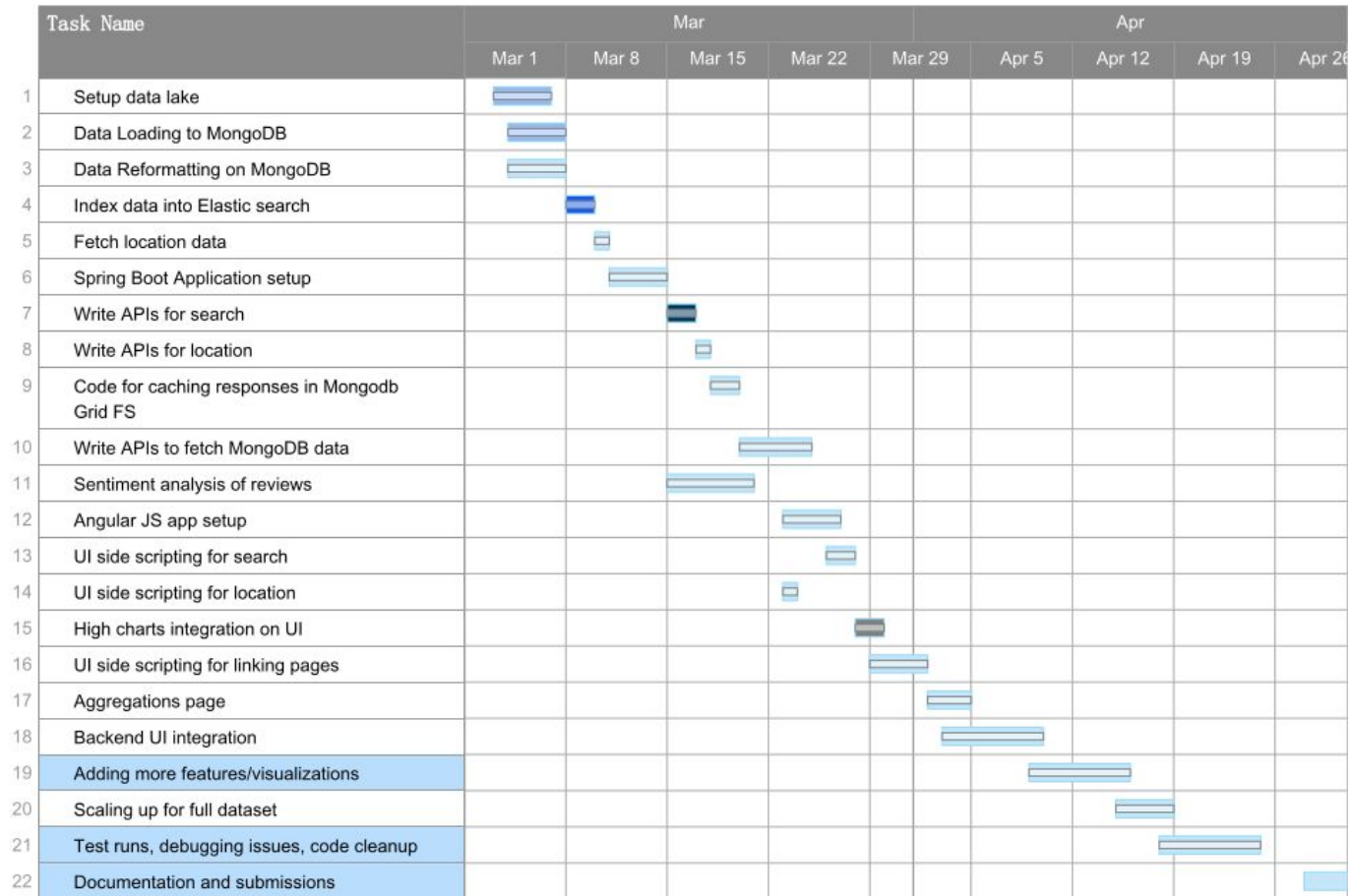
Business geo-distribution by category



Average User Ratings vs Actual Ratings



# PROGRESS - GANTT CHART



# FUTURE WORK

- Analyzing user-user relationships to discover insights into how one user's ratings affect other user's ratings.
- Adding user profiling and security to the application, with ability to assign different levels of visibility to each user.
- Adding location based data like population, demographics, income etc., and analyzing them against existing dimensions to discover more interesting trends.
- Make the solution generic by giving the user the ability to add data from multiple sources, so that the centralized datalake is further enriched and complex insights can be discovered.



# REFERENCES

- Inspiration : <http://people.ischool.berkeley.edu/~sayantan.satpati/yelp/>
- Yelp Dataset Challenge : <https://www.yelp.com/dataset/challenge>
- Yelp Data dictionary : <https://www.yelp.com/dataset/documentation/main>
- Highcharts : <https://www.highcharts.com/demo>
- Highmaps : <https://www.highcharts.com/maps/demo>

**THANK YOU**  
Q&A

---