

Eat, Rate and Love — An Exploration of Yelp

Gurui Huang, Boyan Wu, Maoyi Song, Hanyi Chen

MSBA324 Web and Social Analytics | Instructor: Nabanita Talukdar



INTRODUCTION

Yelp is a company that published crowd-sourced reviews about local business and it also provides online reservation service (Yelp about us, 2018). User may use Yelp application to search local business, such as restaurants or schools and use 5 stars rating system to submit their reviews.

Data Clean- Process missing data and use SQL to build new table to analyze.

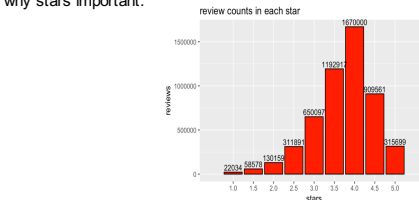
Overall Analysis – Descriptive Analysis about “Stars” and “Reviews”.

Regression Analysis – Using ANOVA and T-test to find out if independent variable “check in numbers” has difference when meeting different dependent variables.

OVERALL ANALYSIS

1. High star rating business exploration & suggestion.

why stars important:



Stars	Review Rate
1.0	0.4%
1.5	1.1%
2.0	2.5%
2.5	5.9%
3.0	12.4%
3.5	22.7%
4.0	31.7%
4.5	17.3%
5.0	6%

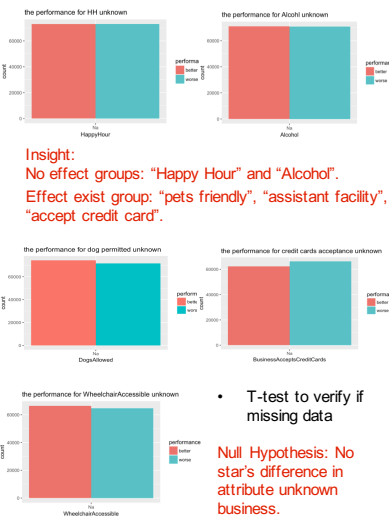
Insight: business with higher star attracts more reviewers, and may has more clients in general;

2. Exploration on business attributes to find out possible improvement direction.

- Attribute selection from data exploration: attributes selected as number of validate data greater than 5000 to be meaningful;

Variable Name	Available Data	Remark
“dog allowed”	6,005	Pets friendly
“Happy Hour”	6,182	
“Credit Card Acceptance”	23,581	Payment method
“wheelchair accessible”	20,947	Assistance facility
“alcohol”	10,412	

- Metric determination: assign “better” (star greater than 3) and “worse” (star 3 or less) performance to business.
- NA’s affection: difference in performance (stars) will affect the interpretation’s power.



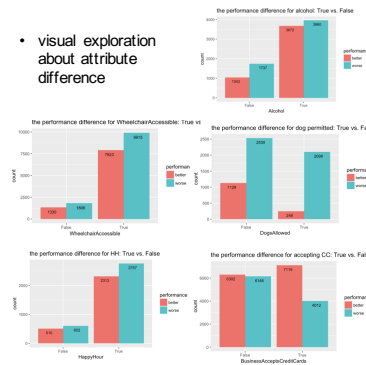
Insight:
No effect groups: “Happy Hour” and “Alcohol”.
Effect exist group: “pets friendly”, “assistant facility”, “accept credit card”.

- T-test to verify if missing data

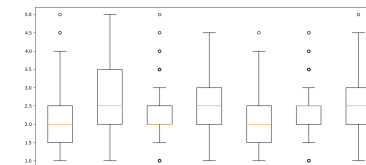
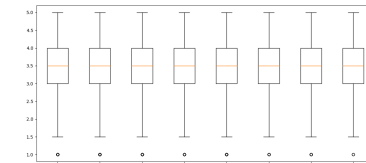
Null Hypothesis: No star’s difference in attribute unknown business.

Attributes:	P-value	Conclusion
HH	0.97	Fail to reject null hypothesis; no difference
Alcohol	0.25	
Wheelchair Accessibility	1.18e-7	
Dog Allowed	1.59e-12	
Credit Card Accepts	2.2e-16	Reject null hypothesis; difference exists

- visual exploration about attribute difference



SPECIFIC CASE ANALYSIS



- T-test (star & check-in database)

Attributes	Conclusion
MC VS Subway	
T-statistic: -13.856	H0: McDonald's star is lower than Subway. Accept H0. The result is significant, so McDonald's stars have significant difference with Subway.
P-value: 7.38	
MC VS Burger King	
T-statistic: -1.57	H0: MC star is lower than Burger King. Reject H0. The result is not significant; different exist
P-value: 0.12	
Weekdays VS Weekend	
p-value < 2.2e-16	H0: the check-in numbers between weekdays and weekends are the same. Reject H0; different exist

REGRESSION ANALYSIS

- Model 1: checkins ~ hour + weekday (original data, numeric)
- Model 2: checkins ~ hour + weekday (cleaned data, string)

```
Call:
lm(formula = checkins ~ hour + weekday, data = df1)
lm(formula = checkins ~ hour + weekday, data = df2)

Residuals:
Min      1Q  Median      3Q      Max
-4.14   -3.22   -2.23   -0.82  1477.38

Coefficients:
(Intercept) 4.54404 0.03342 144.444    2e-16 ***
hour1000    0.00440 0.00079  22.254    2e-16 ***
hour100     -1.03436  0.00068  -21.434    2e-16 ***
hour10      -0.00451  0.00142  -30.234    2e-16 ***
hour10000   -1.94897  0.00248  -81.190    2e-16 ***
hour100000  -1.92395  0.00452  -42.503    2e-16 ***
hour1000000 -1.61203  0.00804  -22.394    2e-16 ***
hour10000000 -1.45032  0.01311  -13.832    2e-16 ***
hour100000000 -1.32082  0.02080  -10.124    2e-16 ***
hour1000000000 -1.23475  0.03040  -10.136    2e-16 ***
hour10000000000 -0.70033  0.03719  -16.473    2e-16 ***
hour100000000000 0.01713  0.03745  0.438    0.6497
hour1000000000000 1.34032  0.04074  33.023    2e-16 ***
hour10000000000000 -0.38030  0.03764  -10.231    2e-16 ***
hour100000000000000 -0.81277  0.03767  -21.703    2e-16 ***
hour1000000000000000 -0.86436  0.03733  -23.220    2e-16 ***
hour10000000000000000 -0.60084  0.03701  -16.141    2e-16 ***
hour100000000000000000 -1.00328  0.04030  -24.380    2e-16 ***
hour1000000000000000000 0.55441  0.04708  11.493    2e-16 ***
hour10000000000000000000 0.13416  0.03388  2.440    0.01277
hour100000000000000000000 -0.09071  0.06193  -1.443    0.14303
hour1000000000000000000000 -0.44259  0.06055  -44.310    2e-16 ***
hour10000000000000000000000 -0.97171  0.07494  -12.424    2e-16 ***
hour100000000000000000000000 -1.45060  0.08203  -17.474    2e-16 ***
hour1000000000000000000000000 -0.00021  0.02679  -0.046    0.95217
weekday1    0.53042  0.02051  27.328    2e-16 ***
weekday2    0.14468  0.02020  44.271    2e-16 ***
weekday3    -0.32055  0.02615  -12.451    2e-16 ***
weekday4    -0.38869  0.02463  -14.497    2e-16 ***
weekday5    -0.37428  0.02424  -14.292    2e-16 ***
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 13.9 on 391214 degrees of freedom
Multiple R-squared: 0.00109, Adjusted R-squared: 0.00030
F-statistic: 3941 on 3 and 391214 DF, p-value < 2.2e-16
```

- In Model1 both main effects and the interaction are significant. However, most of the coefficients are smaller than 1. Then, we simplified the model (Model2) to weekdays, weekend morning, noon and night.
- Results show the both model are significant. ANOVA compared. P<0.05 Model 1 is the best fit.

CONCLUSION

Attributes	Conclusions	Suggestion
“Dog Allowed”	Large loss in better group when businesses are considered to be pets friendly.	Tips for business owner: it may not be a good idea to allow dogs in order to achieve higher star.
“Happy Hour”	More business tends to have HH; the difference for two levels increases when business have “Happy Hour”	Tips for business owners: if HH applicable: HH might not result in increasing the chance to gain higher star level; Test design needed by Yelp to design promotional/make suggestions relating to having HH for business.
“Credit Card Acceptance”	Large loss in better business group measured for business dose accept credit card.	Tips for business owners: it may be a good idea to accept credit card in order to achieve higher star.
“Wheelchair accessible”	More business tends to have wheelchair friendly environment; the difference for two levels increases when business do have wheelchair accessible environment.	Not too much business meaning.
“Alcohol”	Difference of performance decreasing significantly as having the alcohol.	It is a good idea to have alcohol served in business to increase the chance
“Check-in”	Check-in numbers are negative correlation with weekdays, but positive correlation with weekends.	May be good idea to build brand effect, word of mouth; increasing the number of promotion and customer benefit during weekend.

Reference
Yelp dataset Retrieved from: <https://www.kaggle.com/yelp-dataset/yelp-dataset>.
Yelp about us (2018). Retrieved from <https://www.yelp.com/about>