# Capstone project task 1

## Abstract

In order to finish this task, I used python and jupyter notebook to develop and run my code.
I used graphlab to run text alalytic algorithms and matplot to visualize.
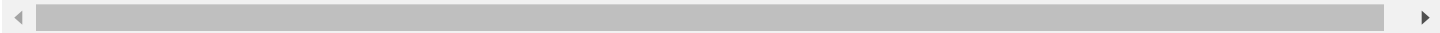
## Implementation

At first I imported all dataset. Because of the task requirement, I think that only the text reviews themselves are important, so I remove columns like: 'votes', 'data' and id columns.

```
reviews = gl.SFrame.read_json('yelp_dataset_challenge_academic_dataset/yelp_academic_dataset_
review.json',orient='lines')
reviews.remove_columns(['business_id','date','review_id','user_id','votes','type'])
```

After that, I created a dictionary map each word appeared in the review with its tf_idf score. I also removed stop words from the dictionary as they are mostly invaluable.

```
reviews['word_count'] = gl.text_analytics.count_words(reviews['text'],delimiters = delimiters
dict_trim_by_keys(stopwords, exclude=True)
reviews['tf_idf'] = gl.text_analytics.tf_idf(reviews['word_count'])
```

Because the data is quite big, I only use 10% of all reviews (about 110,000 reviews) to analyse. The data after preprocessing look like this:

```
sample = reviews.sample(.1,seed = 317)
sample.head()
```

| | stars | text | word_count | tf_idf |
|---|---|---|---|---|
| | 5 | Dr. Eric Goldberg is a fantastic doctor who has ... | {'accessible': 1L, 'fantastic': 1L, ... | {'accessible': 6.440940320962146, ... |
| | 3 | Ate a Saturday morning breakfast at the Pine ... | {'saturday': 1L, 'madison': 1L, 'qualms': ... | {'saturday': 3.7162420980736974, ... |
| | 5 | i rarely give five star reviews but for what ... | {'breakfasts': 1L, 'give': 1L, 'taste': 1L, ... | {'breakfasts': 6.55469249425622, 'gi ... |
| | 3 | This is definitely not your usual truck stop. ... | {'good': 1L, 'rude': 1L, 'restaurant': 1L, ... | {'good': 0.9400482737462729, ... |
| | 4 | Unlimited hot coffee. I don't have any ... | {'generous': 1L, 'good': 1L, 'restaurant': 1L, ... | {'generous': 4.523116918480465, ... |
| | 5 | Californians are all about the In-N-Out, w ... | {'love': 1L, 'family': 1L, 'feel': 1L, 'almo ... | {'love': 1.8910237607652902, ... |
| | 4 | My meaty goodness...why isn't Culvers a ... | {'cheese': 1L, 'picnic': 1L, 'chain': 1L, 'bac ... | {'cheese': 2.5850249953356883, ... |
| | 5 | Love it!!!!! Love it!!!!!! love it!!!!!!! ... | {'love': 4L, "culver's": 1L} ... | {'love': 7.564095043061161, ... |
| | 3 | If you live in Madison, you too probably have a ... | {'charter': 1L, 'love': 1L, 'inadequate': 1L, ... | {'charter': 9.46779250322994, 'lo ... |
| | 2 | My arch-enemy.\n\nI've never encountered a ... | {'contacted': 1L, 'extra': 1L, 'lack': 1L, ... | {'contacted': 6.237942631329764, ... |

[10 rows x 4 columns]

## Task 1.1:

In this task, I will use 2 different methods of lda topic modeling named Collapsed Gibbs sampling (cgs) and AliasLDA method (alias) to create 10 topics each model. After that, I will creat visualization base on 2 model topics and compare them.

```python
def model_cgs(data):
    return gl.topic_model.create(data,num_topics=topic_size, num_iterations=200,print_interval=50,method='cgs')
def model_alias(data):
    return gl.topic_model.create(data, num_topics=topic_size, num_iterations=50,print_interval=50,method='alias')

topic_model_cgs = model_cgs(sample['tf_idf'])

topic_model_alias = model_alias(sample['tf_idf'])
```

| Topic | Sample data using alias method | Sample data using cgs method |
|---|---|---|
| 0 |  |  |

| Topic | Sample data using alias method | Sample data using cgs method |
|---|---|---|
| 1 | pizza cheese chicken fries burger sauce sandwich bacon bread beef delicious chocolate cream salad | car show hair store customer service time told call work business |
| 2 | car dogs dog kids friendly work company call drive phone card days | great awesome love food service store beer friendly restaurant pizza |
| 3 | show staff seats nail movie dr love office amazing great fantastic | great food love awesome service pizza stay menu people friendly |
| 4 | room hotel stay vegas club strip casino pool rooms night | great food service awesome store love friendly restaurant room |
| 5 | bar breakfast drink beer drinks beers atmosphere pretty patio area | great food service love restaurant pizza friendly awesome store |
| 6 | minutes wait told bad time order line people service hour | great service food love pizza staff people friendly |
| 7 | store find stuff shop items mall price location selection | great service food awesome love pizza amazing friendly |
| 8 | sushi restaurant menu food dinner dishes steak shrimp buffet rolls | room hotel stay vegas casino people floor show |
| 9 | coffee hair job tea make wonderful cut owner | great food closed love awesome service nom amazing |

It is clear that alias method did better job than cgs method.
Model using alias method can differentiate topics very good and topics look very distinct: restaurants, coffee shops, stores, movie theatres, ...
By contrast, cgs method did not do well its job because more than half the topics scored the word *'great'* as the most important words. Therefore, in task 1.2 I will only use alias method.

## Task 1.2:

In this task, I will get about 100,000 positive reviews (5 stars) and about 100,000 negative reviews (1 or 2 stars) as sample and run topic models on them.

```
gl.canvas.set_target('ipynb')
reviews['stars'].show('Categorical')
```

### Most frequent items from <SArray>

| Value | Count | Percent |
|---|---|---|
| 5 | 406,045 | 36.078% |
| 4 | 342,143 | 30.4% |
| 3 | 163,761 | 14.551% |
| 1 | 110,772 | 9.842% |
| 2 | 102,737 | 9.128% |

It is clear that reviews with 4 or 5 stars account for 66% of original data and reviews with 3 or less stars share the remaining 33%. I will use 25% of positive reviews (5 stars) and 30% of negative reviews(less than 3 stars) as sample and run topic model on these samples.

```
positive_reviews = reviews[reviews['stars'] == 5]
sample_positive = positive_reviews.sample(.25 ,seed = 317)
sample_positive.shape

(101630, 4)


negative_reviews = reviews[reviews['stars'] < 3]
sample_negative = negative_reviews.sample(.5,seed = 317)
sample_negative.shape

(106853, 4)


positive_model_alias = model_alias(sample_positive['tf_idf'])

negative_model_alias = model_alias(sample_negative['tf_idf'])
```
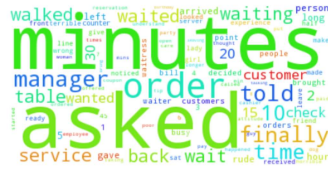
| Topic | Positive | Negative |
|---|---|---|

| 0 | show — care, office, dr, shows, studio, seats, doctor, feel, classes, tickets, dog, class, staff, yoga | bar — water, drinks, waitress, server, drink, tables, food, bartender |
| 1 | lunch — food, fries, breakfast, chicken, delicious, sandwich, burger, menu, order | room — rooms, bathroom, front, dirty, stayed, stay, hotel, desk, pool, floor |
| 2 | pizza — steak, cake, bread, dessert, chocolate, butter, wine, service | sushi — food, ordered, buffet, dish, steak, restaurant, dinner, chicken |
| 3 | car — customer, job, work, house, dog, professional, nails, care, call, highly, recommend | people — girls, group, guys, kids, line, club, guy, crowd, sucks |
| 4 | room — hotel, vegas, pool, strip, stay, casino, clean, view | rude — reviews, nails, money, waste, hair, bad, massage, horrible, service |
| 5 | store — coffee, shop, service, friendly, find, buy, items, employees, local | show — vegas, music, floor, night, casino, free, strip, loud |
| 6 | cool — awesome, big, stop, line, back, park, time, bring, people | store — location, shop, buy, items, prices, quality, service, terrible |
| 7 | sushi — food, beef, dishes, menu, rice, roll, spicy, chicken, fish | car — company, customer, call, called, told, phone, charge, office |
| 8 | bar — great, night, wine, beer, drinks, music, amazing, group, friends, atmosphere | pizza — ordered, burger, fries, taco, sandwich, cheese, chicken, salad |

| | | |
|---|---|---|
| 9 |  |  |

Even though the most scored words in models are words that related to business categories of the places being reviewed. We can notice that there are words like *love*, *awsome* or *amazing* in positive feedbacks while in negative reviews there are many negative words like *rude*, *waste*, or *horrible*.