

# Should I invest it?

## Predicting future success of restaurants using dataset

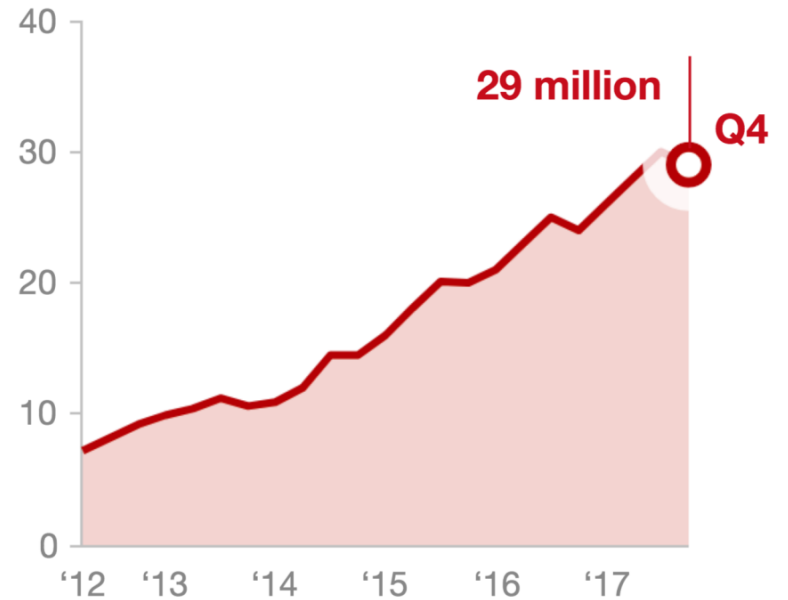
Xiaopeng Lu, Jiaming Qu

PEARC' 18

# INTRODUCTION

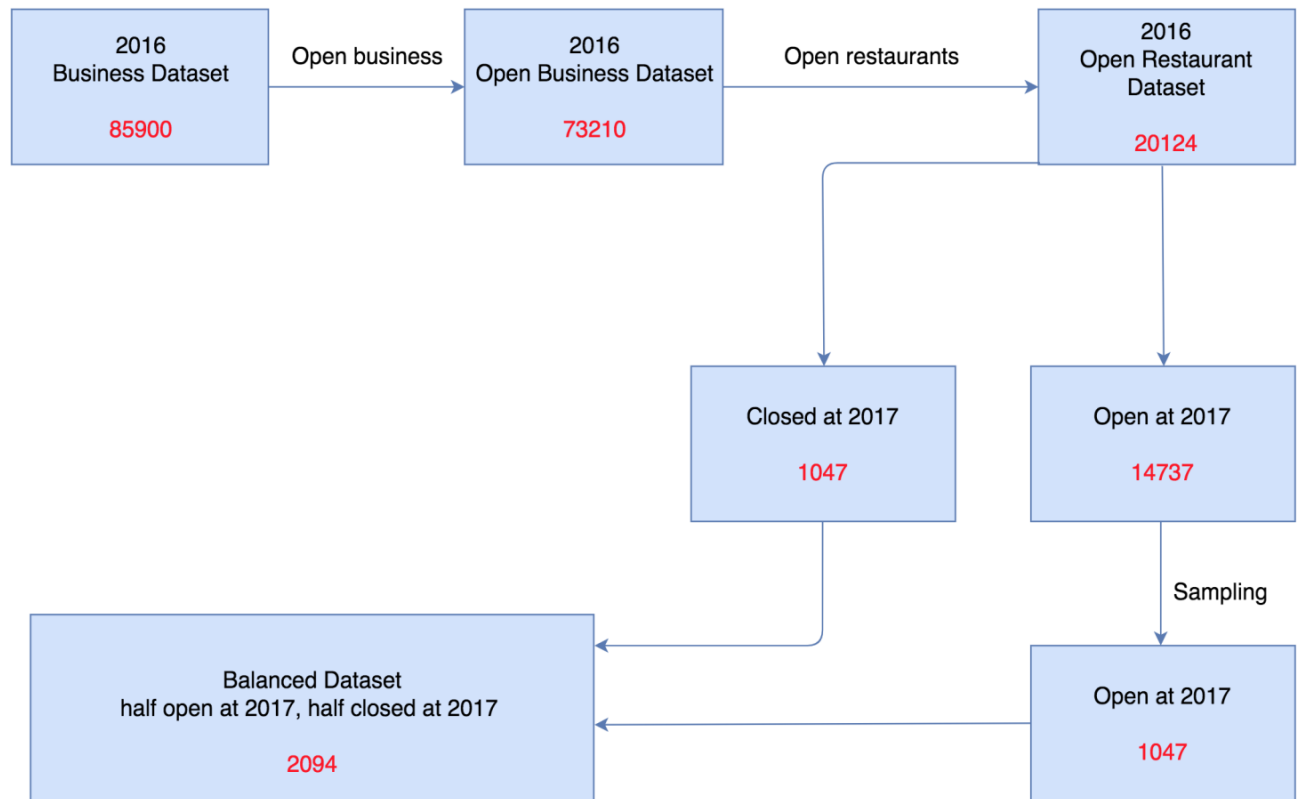
- More and more people choose Yelp to help making daily decisions
- It would be fun to see if the future development of certain restaurants can be predicted through current data
- Might help investors make better decisions

Average monthly mobile app unique users



# DATASET DESCRIPTION

- Two databases with identical fields but different release time (2016,2017)
- **Aim to get restaurants closed in this one year period**



# FEATURE ENGINEERING

Text Features	Unigram	Good
		Bad
	Bigram	Sanitation
		Location
		Service
		Taste
	Non-text Features	Trend
Business		Review count
		Chain restaurant
		Return guest count
		Restaurant type
Location		Nearby restaurant comparison
		(City economic status)

## TEXT FEATURES - Unigram (2)

- Using a sentiment dictionary to **catch** certain sentiment words
  - eg. **“unigram\_good”**: 'love', 'nice', 'delicious', 'amazing', 'top', 'favorite', etc.

**“unigram\_bad”**: 'nasty', 'noisy', 'disappoint', 'cockroach', 'fly', 'mosquito',

etc.

- Count number of word occurrence for all reviews with same business
- NOTICE: only TWO features generated finally

## A simple example...

restaurant name	reviewer	restaurant name	uni_good	uni_bad
Outback Steakhouse	Jack	The food here is <b>amazing!</b> One of my <b>favorite</b> restaurant in Chapel Hill. The environment is a little bit <b>noisy</b> however...	2	1
Burger King	Andrew	The food here is <b>amazing!</b> One of my <b>favorite</b> restaurant in Chapel Hill. The environment is a little bit <b>noisy</b> however...	4	2
	Sam	The food here is <b>amazing!</b> One of my <b>favorite</b> restaurant in Chapel Hill. The environment is a little bit <b>noisy</b> however...		

## TEXT FEATURES - Bigram (8)

- Want to discover which parts are critical for business success
- Construct Bigram features by different categories
  - Sanitation (2)
  - Location (2)
  - Service (2)
  - Taste (2)
- Find co-occurrence of pair of words in each sentence

## Bigram - Sanitation (2)

- “sanitation\_good”
  - eg. environment...clean, atmosphere...quiet, etc.
- “sanitation\_bad”
  - eg. environment...nasty, table...dirty, etc.



## Another example :)

restaurant name	reviewer	restaurant name	sani_good	sani_bad
Outback Steakhouse	Jack	I love the <b>atmosphere</b> here, its just so <b>quiet</b> ... The overall <b>environment</b> is <b>clean</b> except the table is a bit <b>dirty</b> .	2	1
Burger King	Andrew	Won't come here again! <b>Bad environment</b> and <b>nasty</b> <b>floor</b> . 1/5	0	3
	Sam	Love the burger. The <b>environment</b> is <b>bad</b> though.		

## Bigram - Service (2)

- “Service\_good”
  - eg. waiter...helpful,service...fantastic, etc.
- “Service\_bad”
  - eg. waitress...worst, staff...disrespect, etc.

## Bigram - Location (2)

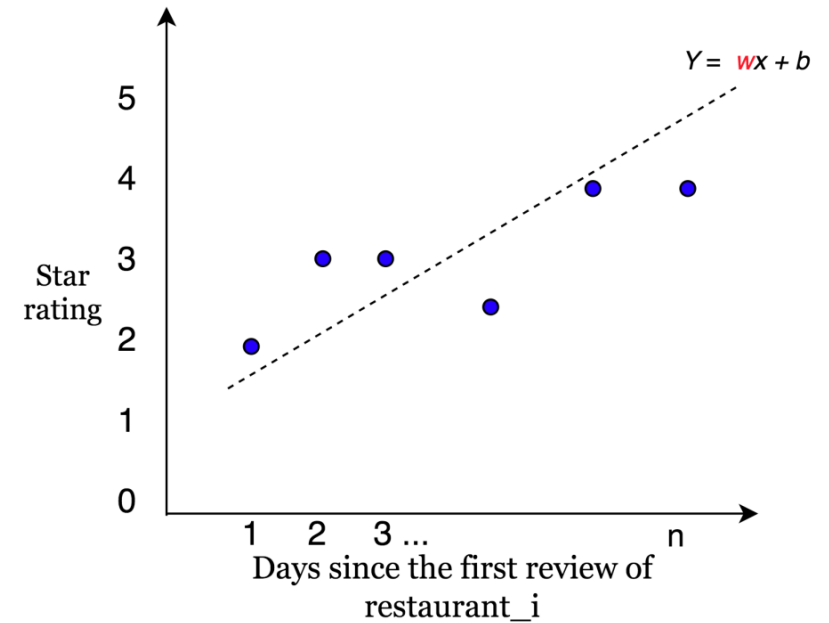
- “location\_good”
  - eg. place...cool, parking...easy, etc.
- “location\_bad”
  - eg. place...crowded, bar...boring, etc.

## Bigram - Taste (2)

- “Taste\_good”
  - eg. drink...best, dessert...wonderful, etc.
- “Taste\_bad”
  - eg. food...nasty, appetizer...disgusting, etc.

# NON-TEXT FEATURES (5)

- Trend
  - Star gain/loss coefficients
- Business
  - Review count
  - Chain restaurant
  - Return guest count
  - Restaurant type
- Location feature
  - Nearby restaurants comparison (not finished)
  - City economic status (failed)



Final Feature table looks like...

restaurant_id	uni_good	...	star_coeff	chain	...	Open_2017
0001	...	...	...	...	...	True
0002	...	...	...	...	...	False
0003	...	...	...	...	...	True
0004	...	...	...	...	...	True
...	...	...	...	...	...	...
2094	...	...	...	...	...	True

# EXPERIMENT

- 10-fold Cross-Validation
- Logistic Regression
- Feature ablation study
- Accuracy, Precision, Recall, Precision-Recall curve

RESULT...





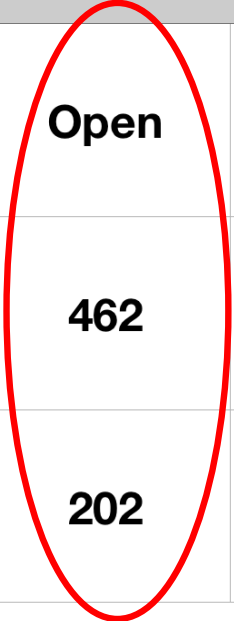
# RESULTS

Accuracy: **62.34%**

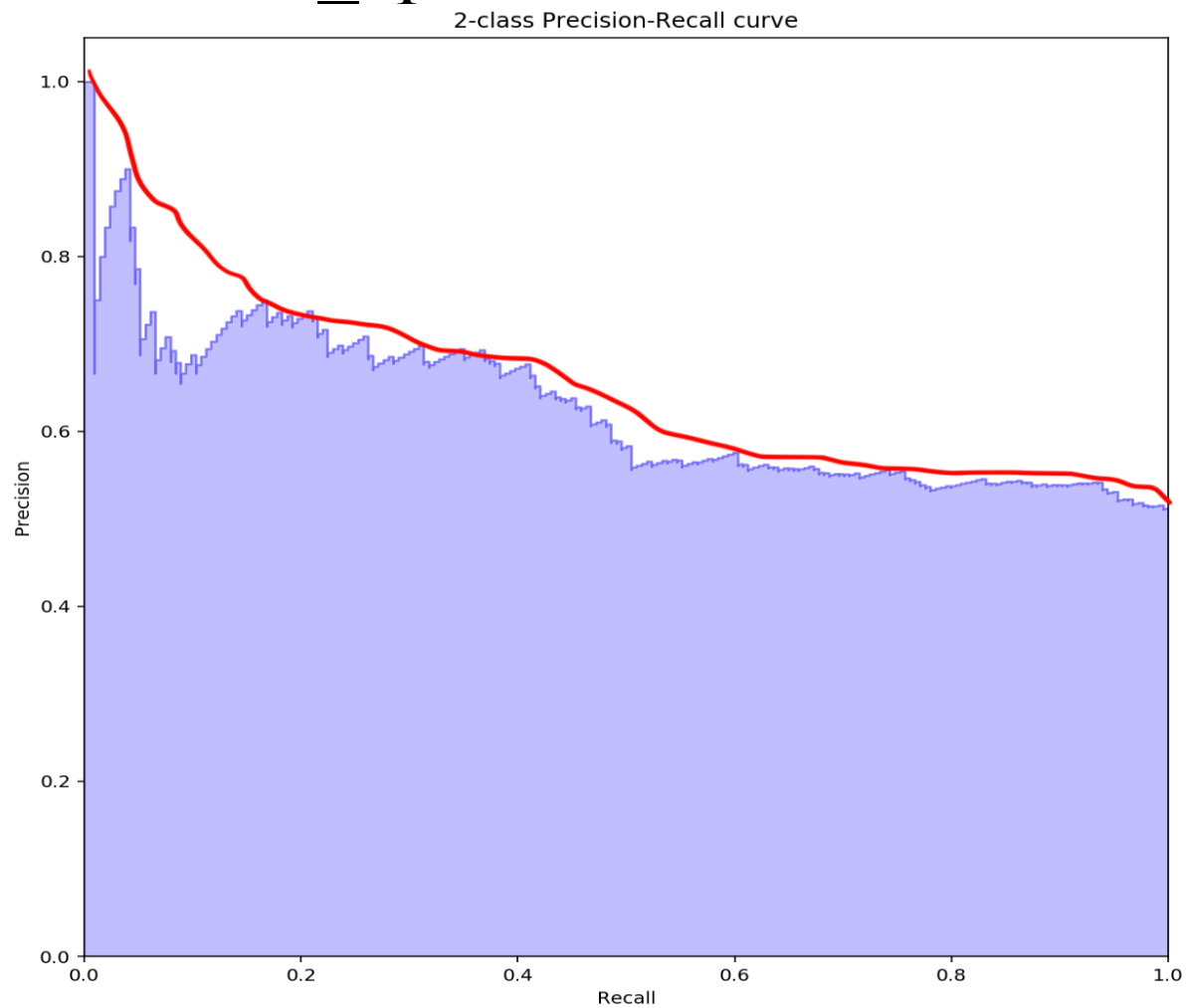
Precision (for open): **0.696**

Recall: 0.442

		Predicted	
		Open	Closed
True	Open	462	584
	Closed	202	839



# Precision - Recall curve for label\_open

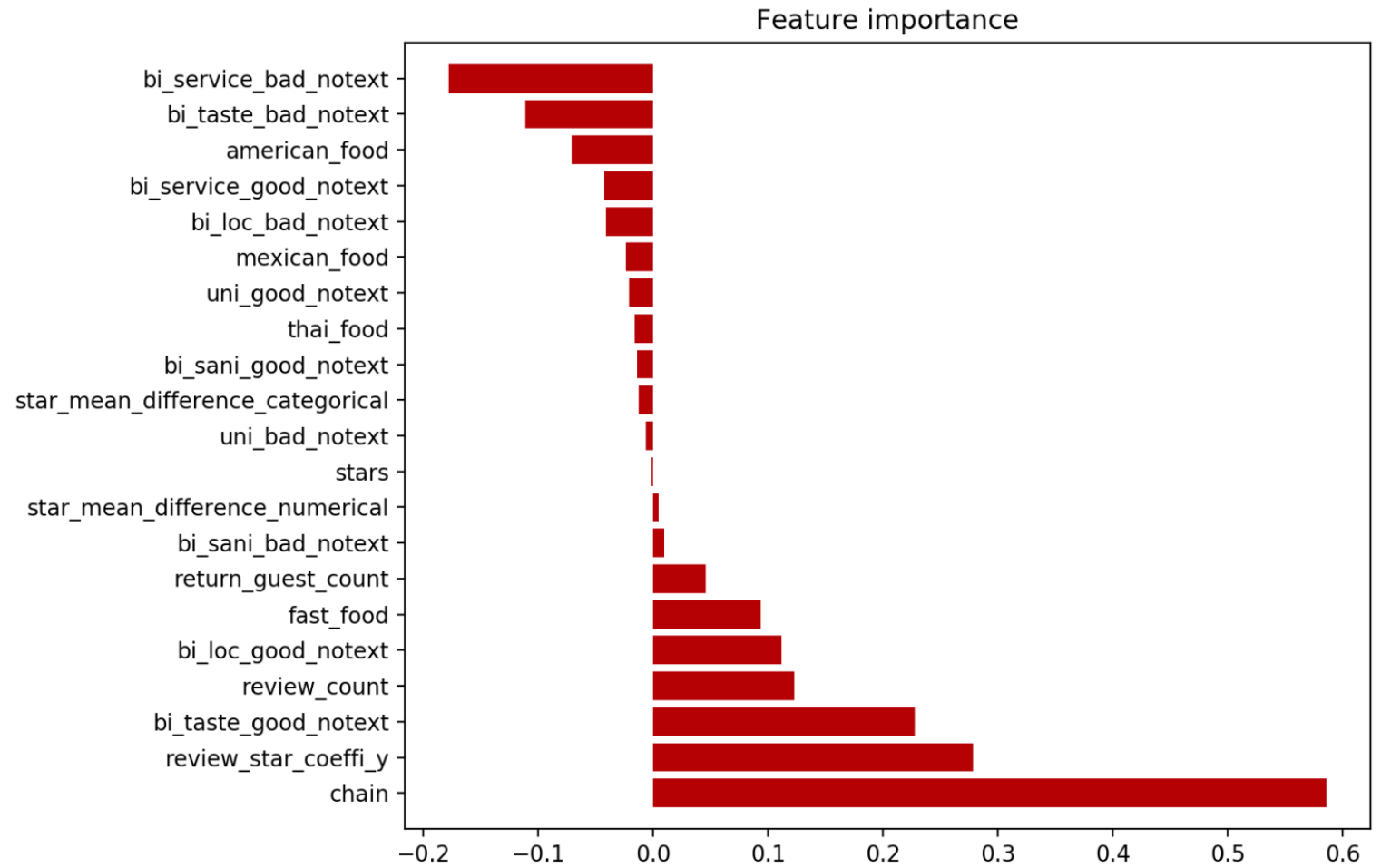


# Feature ablation study

- Business features are the most important
- Text features does not work as desired
  - Why?

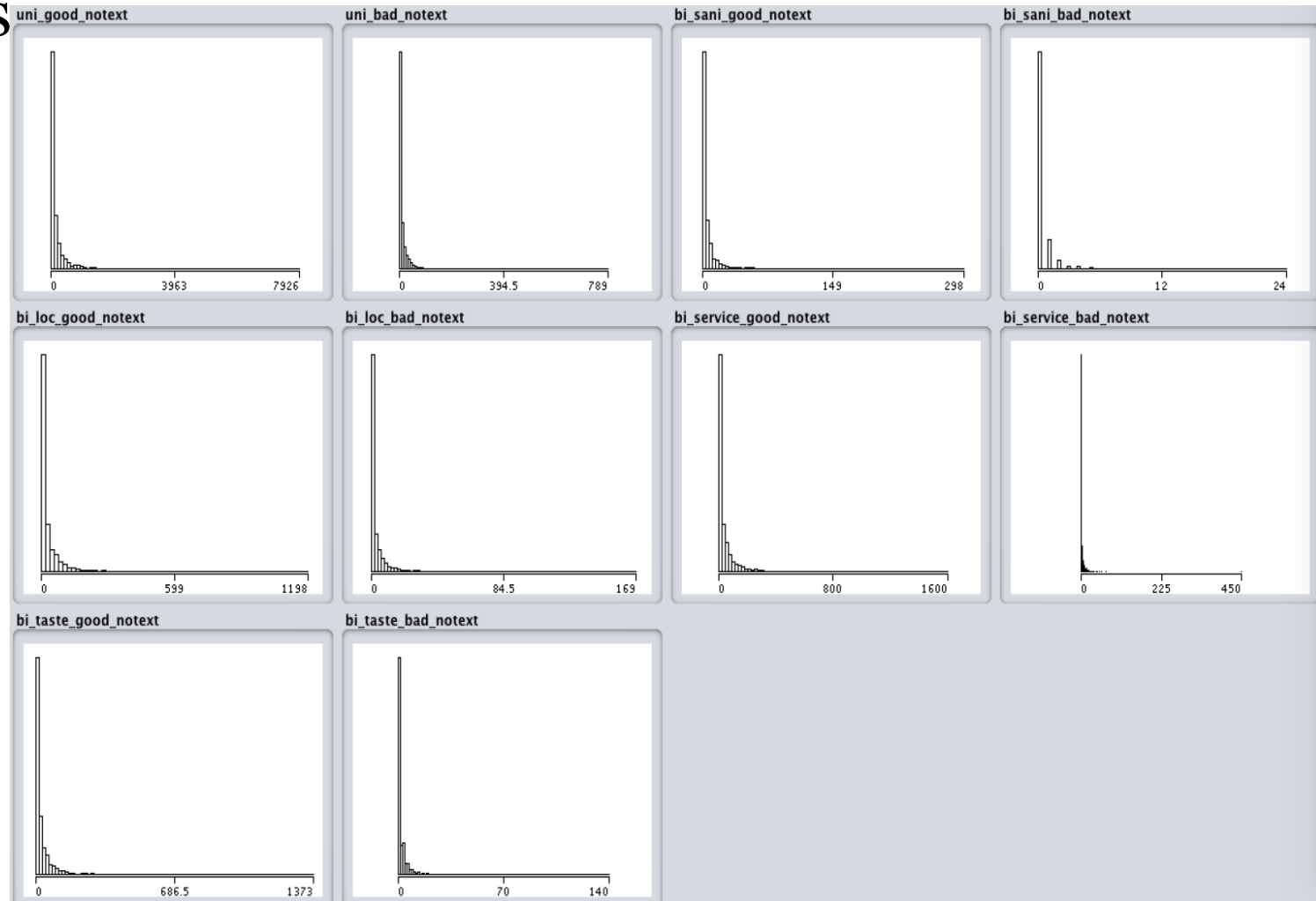
All features	Accuracy	Precision
All	0.6234	0.696
-Text features	0.6229 ▼ (-0.05%)	0.701
-Non-text features	0.5199 ▼ (-10.35%)	0.534
-Unigram	0.6243 ▲ (+0.09%)	0.696
-Bigram	0.6234	0.698
-Trend	0.6224 ▼ (-0.1%)	0.698
-Business	0.5141 ▼ (-10.93%)	0.520

# Error Analysis



# Error Analysis

- Too sparse
- Look back into dictionary



# Error Analysis

- potential solution: Add more words
- Look back into training set and do supervised feature selection

```
sani_noun_notext_list = ['sanitation', 'environment', 'health', 'hygiene', \
                        'surrounding', 'floor', 'table']
sani_good_notext_list = ['clean', 'quiet']
sani_good_notext_list.extend(uni_good_notext_list)
sani_bad_notext_list = uni_bad_notext_list

loc_noun_notext_list = ['location', 'place', 'bar', 'bartenders', 'bartender', \
                        'atmosphere', 'parking', 'beach']
loc_good_notext_list = ['easy', 'pleasant', 'pleased', 'fun']
loc_good_notext_list.extend(uni_good_notext_list)
loc_bad_notext_list = ['hard', 'busy', 'annoy', 'underground']
loc_bad_notext_list.extend(uni_bad_notext_list)

service_noun_notext_list = ['service', 'quality', 'staff', 'waiter', 'waitress', \
                            'prepare', 'price']
service_good_notext_list = ['24hour', 'welcoming', 'fantastic', 'nice', \
                            'communicative', 'helpful', 'quick', 'fast', 'super']
service_good_notext_list.extend(uni_good_notext_list)
service_bad_notext_list = ['bad', 'worse', 'worst', 'confusing', 'improper', 'late', \
                            'disrespect', 'tragic']

taste_noun_notext_list = ['taste', 'food', 'drink', 'appetizer', 'dessert']
taste_good_notext_list = ['great', 'good', 'fantastic', 'pleasant', 'pleased', 'quick', \
                            'fast', 'decent', 'organic', 'inexpensive', 'cheap', 'fresh']
taste_bad_notext_list = uni_bad_notext_list
```

# Error Analysis

- City economic status feature doesn't work
- Not all city data are released

