

CATS



Crawl Data

Movie



Activity

Book



Sentiment

Province



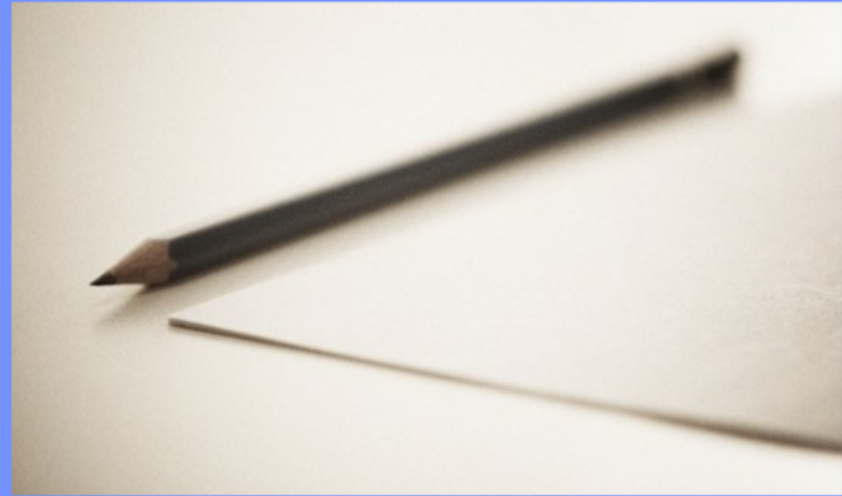
Result &
Appendix



DOGS

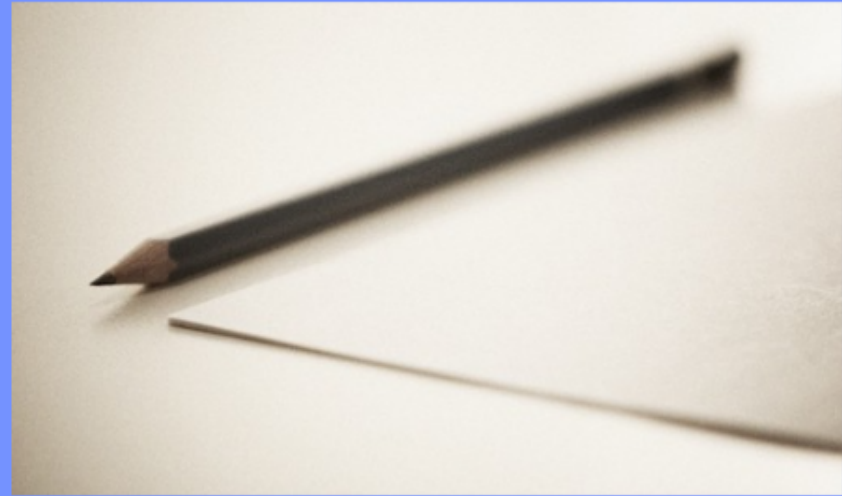
Data crawling

1. Two biggest groups in Douban -
cat:337,238 ; dog:156,360



Data crawling

1. Two biggest groups in Douban -
cat:337,238 ; dog:156,360
2. Location Data: 10k for each group



Data crawling

1. Two biggest groups in Douban -
cat:337,238 ; dog:156,360
2. Location Data: 10k for each group
3. Randomly pick 2k samples for
each group to analyse their book and
movie preference and activity



CATS



Crawl Data

Movie



Activity

Book



Sentiment

Province



Result &
Appendix



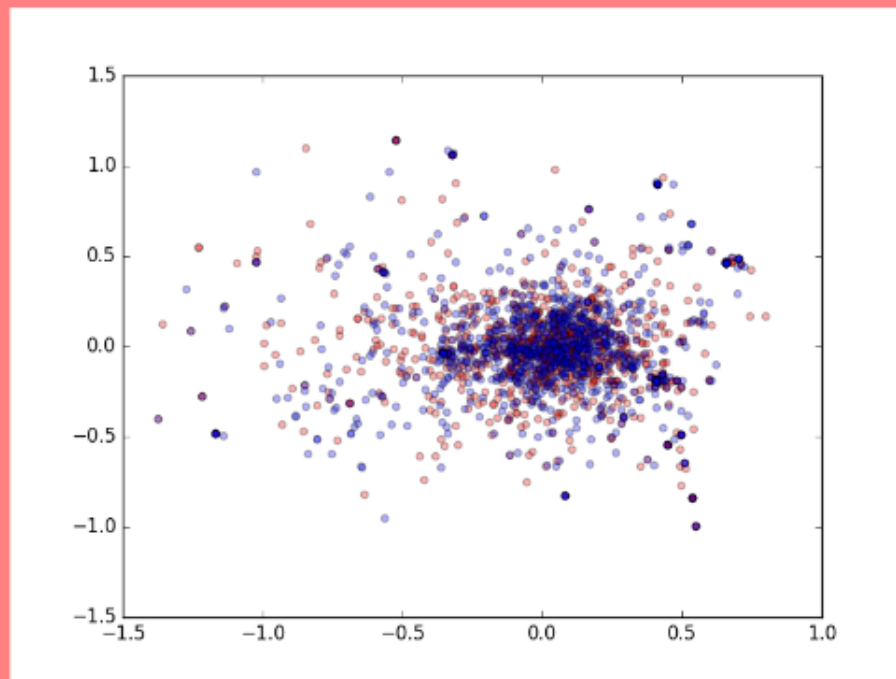
DOGS

Movie

- 37 types in total
- 2k people each group - build corresponding feature vector for each people
- dimension reduction - SVD - $\dim(23)$ - 90%
- LR ; decision tree ; SVM ; naïve bayes ; KNN ; boosting



Movie



A high overlap between them

Movie

Are they different?

Movie

Are they different? NOT

Movie

Are they different? NOT

Result:

logistic regression score: 0.505

lddecision tree score: 0.541

SVM score: 0.533

Naive bayes score: 0.530

k nearest neighbour score: 0.523

boosting score: 0.513

CATS



Crawl Data

Movie



Activity

Book



Sentiment

Province



Result &
Appendix



DOGS

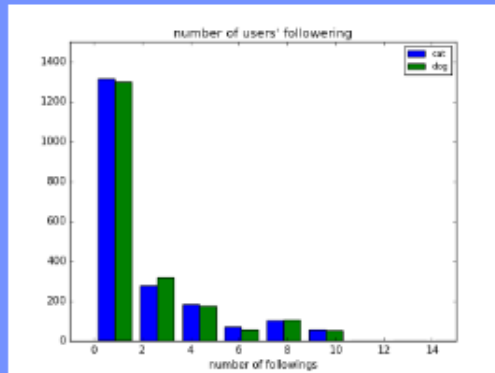
Activity

Following

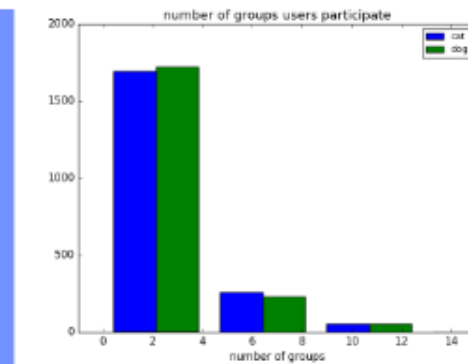
Follower

Group

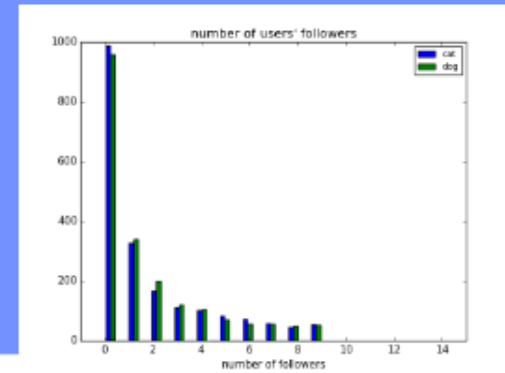
Activity



Following



Group



Follower

Activity

Are they different?

Activity

Are they different? STILL NOT

Activity

Are they different? STILL NOT

Result:

logistic regression score: 0.549

decision tree score: 0.503

SVM score: 0.528

Naive bayes score: 0.497

k nearest neighbour score: 0.537

boosting score: 0.555

CATS



Crawl Data

Movie



Activity

Book



Sentiment

Province



Result &
Appendix



DOGS

Book

There is no special reading interests of people who like cats and dogs.

Book - The process

Data source:

very sparse and no data for a lot of people

Data volume:

about 2k people for each group.

Implementation language:

Java



CATS



Crawl Data

Movie



Activity

Book



Sentiment

Province



Result &
Appendix

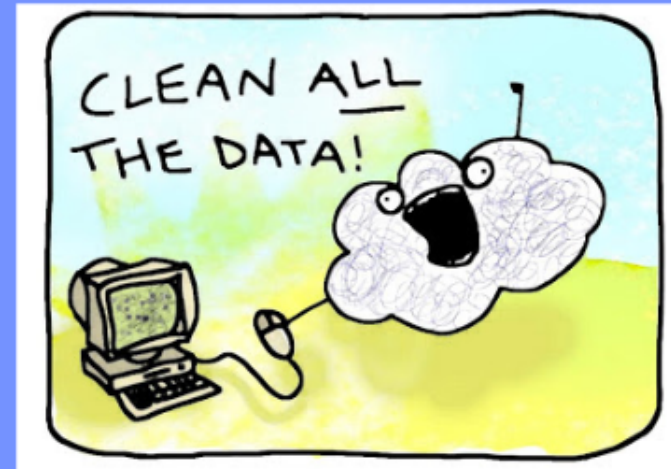


DOGS

Sentiment

1. Pre-processing Data

- 1.1 Get data with cat and dog tags by Twitter API
- 1.2 Extract words from data and remove stop words
- 1.3 Grouped tweets together with same tags
- 1.4 Count TF-IDF



Sentiment

1. Pre-processing Data

2. Count the sentiment of different groups

- 2.1 Use dictionary to divided words into three categories: positive, negative and neutral
- 2.2 Count the weight of positive and negative terms
- 2.3 Normalize the result

	Cat	Dog
Positive	0.60	0.61
Negative	0.40	0.39

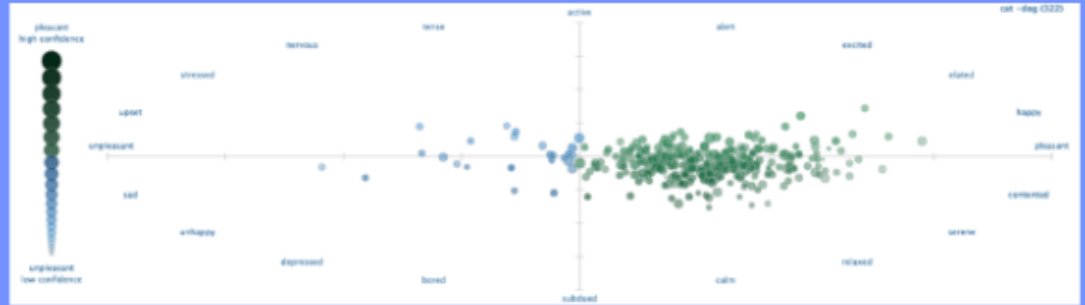
Sentiment

1. Pre-processing Data
2. Count the sentiment of different groups
3. Visualize the result



Sentiment

1. Pre-processing Data
2. Count the sentiment of different groups
3. Visualize the result
4. Another visualization result



cat sentiment analysis



dog sentiment analysis

CATS



Crawl Data

Movie



Activity

Book



Sentiment

Province



Result &
Appendix



DOGS

Province

1. Organize data

```
1 id, location
2 146543835, 上海
3 151632631, 东莞
4 65494746, 武汉
5 173427242, 合肥
6 77184885, 北京
7 176576581, 上海
8 142915595, 厦门
9 138708775, 上海
10 176575984, 沈阳
11 78234944, 重庆
12 175555322, 广州
13 158996233, 北京
14 175581878, 苏州
15 47877728, 杭州
16 176474997, 郑州
17 146394266, 北京
18 136709383, 北京
19 58342689, 杭州
20 148890271, 北京
21 176529476, 杭州
22 68913148, 克拉玛依
23 121183577, 北京
24 74216638, 上海
25 52358238, 北京
26 4721128, 成都
27 168769112, 上海
28 141582262, 上海
29 68288742, 广州
30 88129546, 北京
31 89722793, 杭州
32 44633627, 北京
33 131801253, 广州
34 141383965, 石家庄
35 138806062, 北京
```



```
1 ("北京", 13583)
2 ("上海", 18387)
3 ("广州", 5822)
4 ("深圳", 4358)
5 ("杭州", 3899)
6 ("成都", 3689)
7 ("南京", 2836)
8 ("武汉", 2633)
9 ("重庆", 2194)
10 ("西安", 2049)
11 ("苏州", 1625)
12 ("天津", 1622)
13 ("长沙", 1478)
14 ("郑州", 1883)
15 ("沈阳", 971)
16 ("青岛", 944)
17 ("厦门", 936)
18 ("济南", 889)
19 ("大连", 826)
20 ("合肥", 768)
21 ("宁波", 765)
22 ("无锡", 738)
23 ("哈尔滨", 712)
24 ("福州", 692)
25 ("昆明", 626)
26 ("长春", 553)
27 ("石家庄", 589)
28 ("漳州", 489)
29 ("南宁", 481)
30 ("佛山", 462)
```

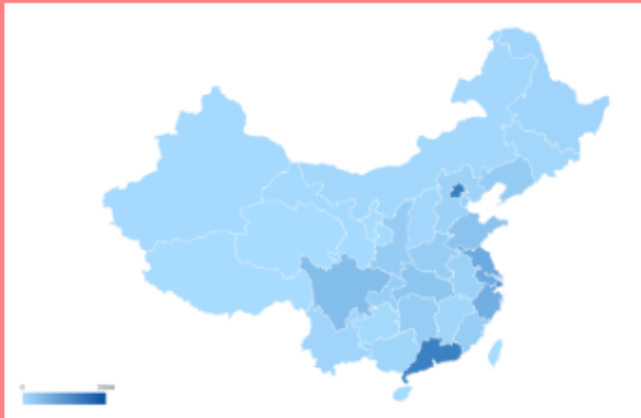


```
1 [{"name": "北京", "value": 13583},
2 [{"name": "上海", "value": 18387},
3 [{"name": "广东", "value": 13810},
4 [{"name": "浙江", "value": 6587},
5 [{"name": "四川", "value": 4581},
6 [{"name": "江苏", "value": 7896},
7 [{"name": "湖北", "value": 3329},
8 [{"name": "重庆", "value": 2194},
9 [{"name": "陕西", "value": 2489},
10 [{"name": "天津", "value": 1622},
11 [{"name": "湖南", "value": 2228},
12 [{"name": "河南", "value": 2218},
13 [{"name": "辽宁", "value": 2488},
14 [{"name": "山东", "value": 3687},
15 [{"name": "福建", "value": 2277},
16 [{"name": "安徽", "value": 1664},
17 [{"name": "黑龙江", "value": 1347},
18 [{"name": "云南", "value": 1134},
19 [{"name": "吉林", "value": 884},
20 [{"name": "河北", "value": 1528},
21 [{"name": "广西", "value": 1873},
22 [{"name": "江西", "value": 976},
23 [{"name": "山西", "value": 848},
24 [{"name": "香港", "value": 417},
25 [{"name": "贵州", "value": 558},
26 [{"name": "宁夏", "value": 429},
27 [{"name": "澳门", "value": 231},
28 [{"name": "新疆", "value": 468},
29 [{"name": "内蒙古", "value": 662},
30 [{"name": "海南", "value": 312},
31 [{"name": "宁夏", "value": 187},
32 [{"name": "台湾", "value": 195},
33 [{"name": "青海", "value": 144},
34 [{"name": "西藏", "value": 188}]
```

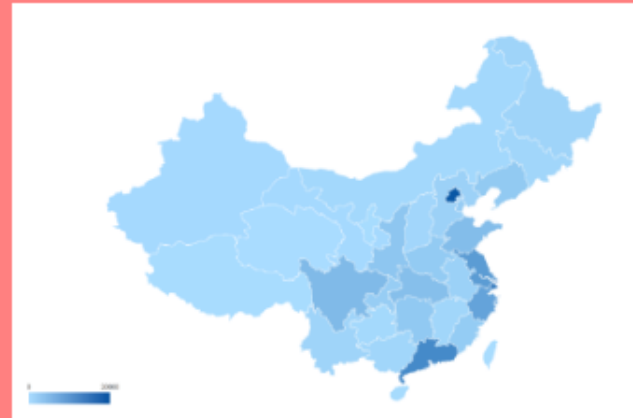
Province

1. Organize data

2. Visualize data Demo: morikka.me:8000/tools/cats.html or [dogs.html](http://morikka.me:8000/tools/dogs.html)



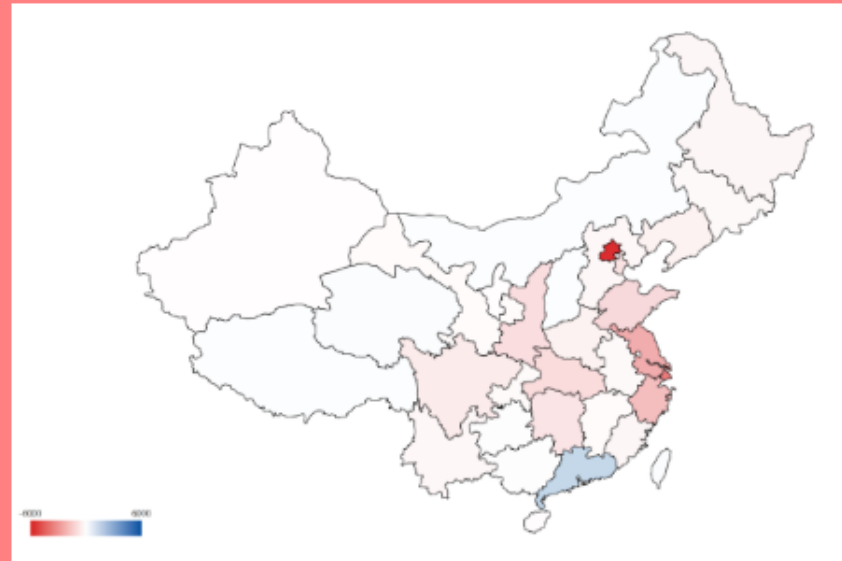
Cat



Dog

Province

1. Organize data
2. Visualize data
3. The difference



morikka.me:8000/tools/diff.html

CATS



Crawl Data

Movie



Activity

Book



Sentiment

Province



Result &
Appendix



DOGS

Result

Similar

Movie & Book choice

Positive

Friend activity

Difference

Cat is more popular online

Cat people are more negative

Cat people love Guangzhou
Dog people love Beijing

Appendix

Community independence

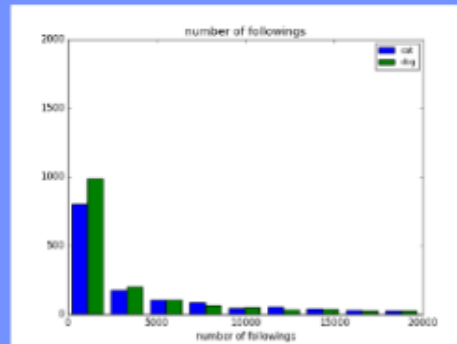
Activity

Following

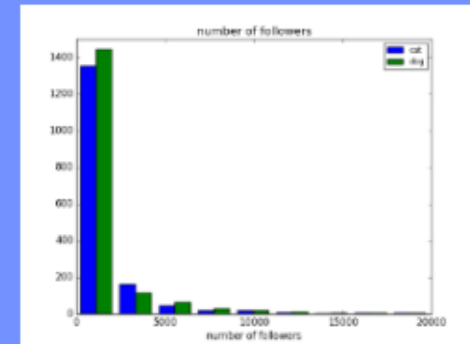
Follower

Group

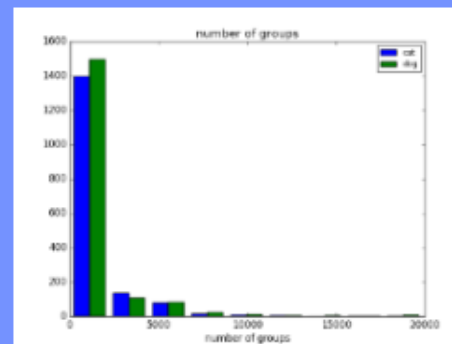
Activity



Following



Follower



Group

CATS



Crawl Data

Movie



Activity

Book



Sentiment

Province



Result &
Appendix



DOGS