

Sina Weibo as a Corpus for Studying Public Opinions

Christine Kuang, Siqu Wu, and Angie Zhu

Department of Statistics, UC Berkeley

May 3, 2012

Outline

1 Introduction

2 Processing

3 EDA

4 Analysis

- LASSO
- l_1 -Norm Support Vector Machine

5 Further Work

Introduction

- Opinions on microblogging and social networking websites
- Sina Weibo 新浪微博 is the largest microblogging website:
accounted for 65% of China's microblog market as of December 2011
- Study public opinions using Sina Weibo as a corpus for a given topic

Topic

- Internet censorship in China
- Time sensitive
- Processing is topic-dependent
- Hot topic is preferred
- Chosen topic: Han Han 韩寒

Background

- HAN Han 韩寒 (born 23 September 1982) is a Chinese best-selling author, professional rally driver, and wildly popular blogger
- Published his first novel *Triple Gate* 三重门 at age of 17
- High school dropout



Photograph by Tony Law / Redux. Source:

[http://www.time.com/time/magazine/article/](http://www.time.com/time/magazine/article/0,9171,1931619,00.html)

[0,9171,1931619,00.html](http://www.time.com/time/magazine/article/0,9171,1931619,00.html)

Background

- Ghostwriting allegation against Han from January 2012
- FANG Zhouzi 方舟子, a scientific author and anti-fraud crusader, created widespread debate on the internet
- 光明与磊落
- Han received a death threat on April 15, 2012

Data Collection

- Topic searching via API:
only the latest results are returned
up to 30 each time
- Collected on April 16 and 17, 2012

Characteristics of Chinese Language

- No explicit delimiter
- Ambiguities in phrases
 - Context ambiguiton: e.g., 他好吃
 - Word definition ambiguiton: e.g., 打
- Out-of-vocabulary words
- No 1-to-1 correspondence between traditional and simplified Chinese

Characteristics of Sina Weibo Posts



Pre-tagging Processing



Tagging

- Process: tagged 3000 total posts with four categories

- Examples:

Positive 支持韩寒! Support Han Han!

Negative 看到韩寒就恶心。 Feel nauseous when I see Han Han.

- Limitations:

- Subjective responses:

e.g., “that wasn’t too bad”

- Uncertain tags

- Quotes

- Posts without subjects

- Posts that just mention opposing author

Pre-segmentation Processing



Segmentation



Conjunction Rules



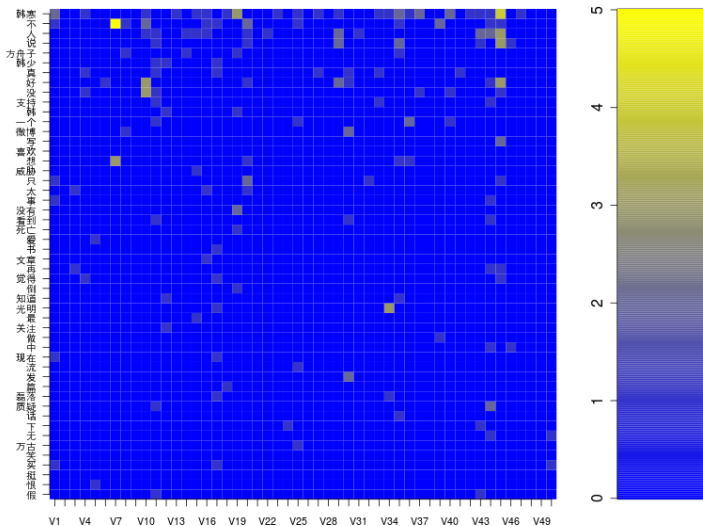
Stop Words and Punctuation Elimination

■ .

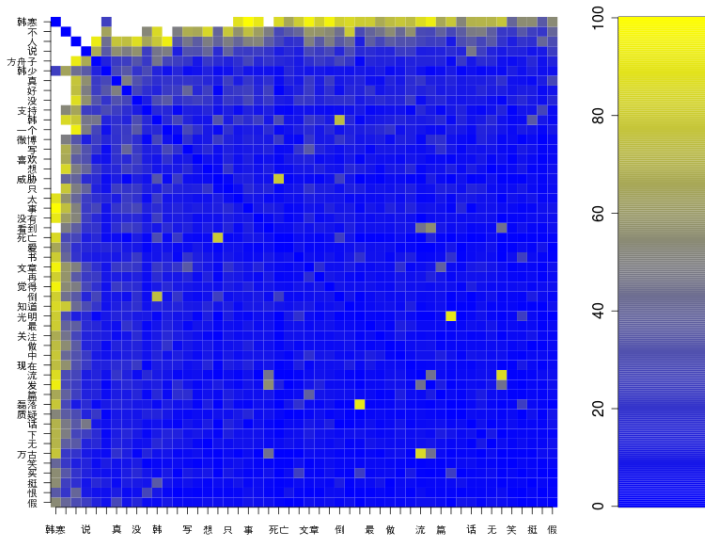
EDA

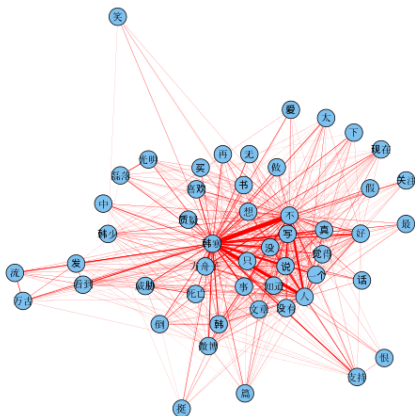
$$\blacksquare \geq 10$$

Frequency



Co-occurrence





Graphical LASSO

- Suppose the random vector $x \in R^p$ has a multivariate normal distribution with mean μ and covariance matrix Σ . The density function of x is

$$f(x|\mu, \Sigma) = \frac{1}{\sqrt{2\pi \det(\Sigma)}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}.$$

- Suppose we have IID $N(\nu, \Sigma)$ data x_1, x_2, \dots, x_n and we want to estimate the inverse covariance matrix Σ^{-1} . The joint likelihood of the data is

$$\begin{aligned} f(x_1, \dots, x_n | \mu, \Sigma) \\ = \frac{1}{(2\pi \det(\Sigma))^{n/2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu) \right\}. \end{aligned}$$

Graphical LASSO

- Taking logarithm to get the log-likelihood (and ignore constant terms):

$$l(\mu, \Sigma^{-1}) = \frac{1}{2 \log \det(\Sigma)} - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

- Do a maximum likelihood estimation (optimize over μ and $S = \Sigma^{-1}$; easy to see that the MLE for μ is \bar{X}):

$$\max_S \left\{ \frac{1}{2} \log \det(S) - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T S (x_i - \mu) \right\}$$

Graphical LASSO

- Here comes the trace trick $\sum_{i=1}^n (x_i - \mu)^T S (x_i - \mu) = \text{Tr}(\sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T S) = n \text{Tr}(\hat{\Sigma} S)$. We end up with the following optimization problem

$$\max_S \left\{ \log \det(S) - \text{Tr}(\hat{\Sigma} S) \right\}$$

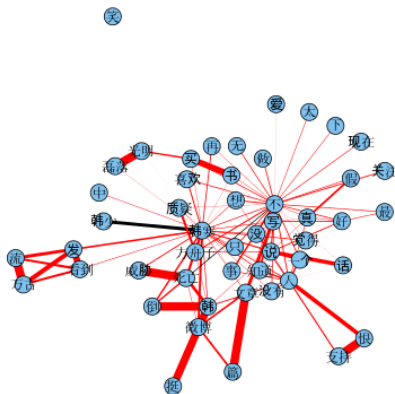
- Banerjee et al. (2007) propose the following optimization problem to recover the sparse structure in a gaussian graphical model

$$\max_S \left\{ \log \det S - \text{Tr}(\hat{\Sigma} S) - \lambda \|S\|_1 \right\}$$

where $\|S\|_1 = \sum_{i=1} \sum_{j=1} |s_{ij}|$.

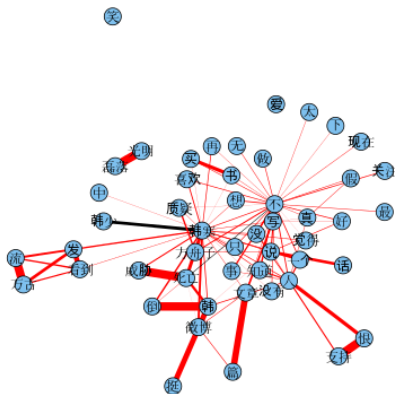


Graphical LASSO



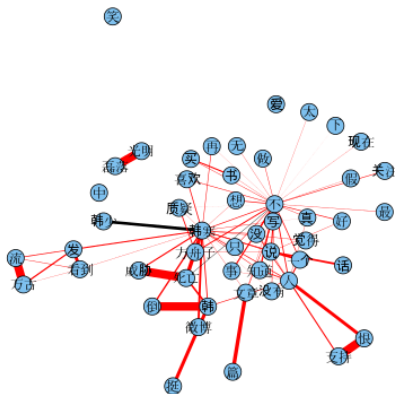
$\lambda = 0.01$

Graphical LASSO



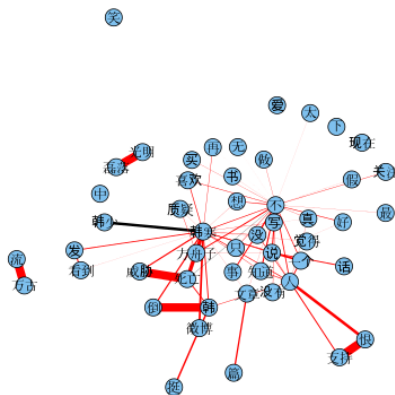
$\lambda = 0.01222$

Graphical LASSO



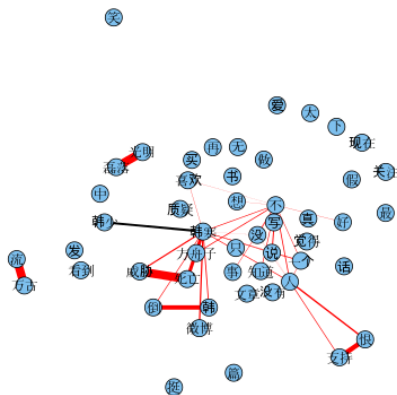
$\lambda = 0.01444$

Graphical LASSO



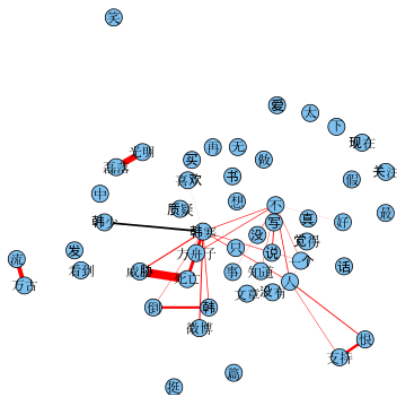
$\lambda = 0.01667$

Graphical LASSO



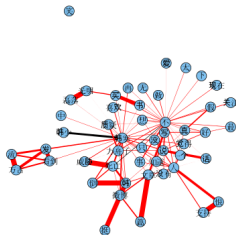
$\lambda = 0.02556$

Graphical LASSO

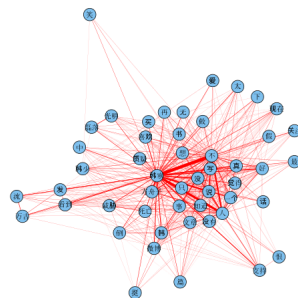


$\lambda = 0.03$

Graphical LASSO



lambda = 0.01



LASSO

$$\hat{\beta}(\lambda) = \arg \min_{\beta} \frac{1}{2} \|y - (\beta_0 + X\beta)\|_2^2 + \lambda \|\beta\|_1$$

- Four models for each category for classification
- General overview of method
- General overview of application to data
 - for 4 categories
 - 10 fold CV
 - Frequency matrix is 3000×795
 - classification error

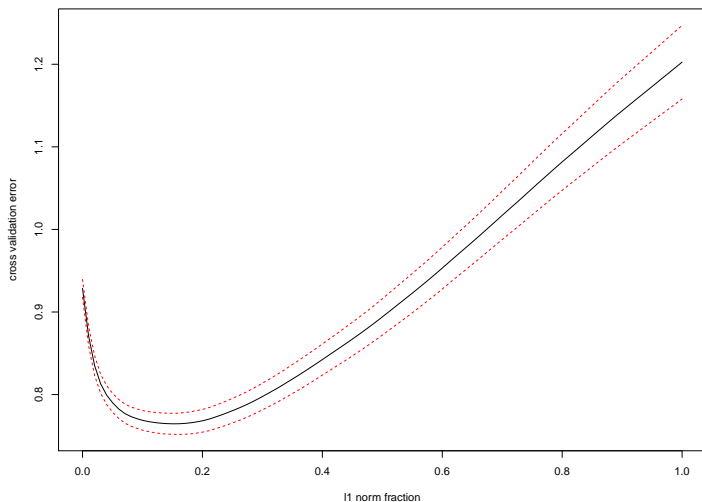
LASSO

LASSO Results

- Three different ways to look at coefficients
- Why:
 - Absolute value: most relevant
 - Positive: more likely to be in category
 - Negative: less likely to be in category

LASSO

LASSO Results: Positive Category



LASSO

LASSO Results: Positive Category

Word	Absolute Coef.	Word	Positive Coef.	Word	Negative Coef.
加油 (keep going)	0.820	加油 (keep going)	0.820	样子 (manner)	-0.396
韩少 (Master Han)	0.644	韩少 (Master Han)	0.644	恋 (love)	-0.344
成熟 (mature)	0.546	成熟 (mature)	0.546	发表 (announce)	-0.336
顶 (support)	0.533	顶 (support)	0.533	道理 (rational)	-0.336
宽容 (tolerant)	0.518	宽容 (tolerant)	0.518	利益 (benefit)	-0.335

LASSO

LASSO Results: Negative Category

Word	Absolute Coef.	Word	Positive Coef.	Word	Negative Coef.
讨厌 (hate)	0.481	讨厌 (hate)	0.481	支持 (support)	-0.008
无耻 (shameless)	0.412	无耻 (shameless)	0.412	不 (no)	0.000
恶心 (disgusting)	0.395	恶心 (disgusting)	0.395	人 (people/person)	0.000
骗子 (liar)	0.380	骗子 (liar)	0.380	说 (say)	0.000
扁 (beat up)	0.353	扁 (beat up)	0.353	方舟子 (FangZhouZi)	0.000

l_1 -Norm Support Vector Machine

- Support vector machine (Vapnik 1996) is another commonly-used machine learning method to classify data points into two categories. Consider again the linear decision function $f(x) = \beta_0 + \beta x$ and the classifier $Class(x) = \mathbf{sign}(f(x))$. To “learn” the parameters, we want the training misclassification rate to be small and the margin of the decision boundary (which can be shown to be $1/||\beta||_2$) to be wide. Hence we consider the following optimization problem (give a 2d illustration here...):

$$\min_{\beta_0, \beta} \sum_{i=1}^n (1 - y_i(\beta_0 + \beta^T x_i))_+ + \frac{\lambda}{2} ||\beta||_2,$$

where $z_+ = \max(0, z)$ (the function $h(z) = (1 - z)_+$ is also known as the hinge loss function)

l_1 -Norm Support Vector Machine

- Similarly, the sparse version of SVM simply replaces the l_2 -norm by l_1 -norm (which is just another measure of the wideness of the margin):

$$\min_{\beta_0, \beta} \left\{ \sum_{i=1}^n (1 - y_i(\beta_0 + \beta^T \mathbf{x}_i))_+ + \lambda \|\beta\|_1 \right\}.$$

We repeat the same data analysis as we did for the LASSO method. The results are summarized in table XX. For efficiently fitting the sparse svm, we use the matlab package by Fung and Mangasarian (2004). Again, 10-fold cross validations are performed in order to select the “best” λ .

Further Work















