

# A Semi-Supervised Bayesian Network Model for Microblog Topic Classification

Yan Chen<sup>1,2</sup>   Zhoujun Li<sup>1</sup>   Liqiang Nie<sup>2</sup>   Xia Hu<sup>3</sup>   Xiangyu  
Wang<sup>2</sup>   Tat-seng Chua<sup>2</sup>   Xiaoming Zhang<sup>1</sup>

<sup>1</sup>State Key Laboratory of Software Development Environment, Beihang University, China

<sup>2</sup>School of Computing, National University of Singapore, Singapore

<sup>3</sup>Arizona State University, United States

11-12-2012

# Outline

- 1 Background and Motivation
- 2 Related Work
- 3 Semi-Supervised Graphical Model
  - The General Framework
  - Probabilistic Graph Model Construction
  - Parameter Inference
- 4 Experiments
  - Experimental Settings
  - Analysis
  - Parameter Analysis
- 5 Conclusion and Future Work

# Outline

- 1 Background and Motivation
- 2 Related Work
- 3 Semi-Supervised Graphical Model
  - The General Framework
  - Probabilistic Graph Model Construction
  - Parameter Inference
- 4 Experiments
  - Experimental Settings
  - Analysis
  - Parameter Analysis
- 5 Conclusion and Future Work

# Background

- Microblogging services are becoming immensely popular in breaking-news disseminating, information sharing, and events participation.

# Background

- Microblogging services are becoming immensely popular in breaking-news disseminating, information sharing, and events participation.
- The most well known one is Twitter, which has more than 140 million active users with 1 billion Tweets every 3 days as of March 2012.

# Background

- Microblogging services are becoming immensely popular in breaking-news disseminating, information sharing, and events participation.
- The most well known one is Twitter, which has more than 140 million active users with 1 billion Tweets every 3 days as of March 2012.
- In China, Weibo ([www.weibo.com](http://www.weibo.com)) has accumulated more than 300 millions users in less than three years. Every second, more than 1000 Chinese tweets are posted in Weibo.

# Background

- Microblogging services are becoming immensely popular in breaking-news disseminating, information sharing, and events participation.
- The most well known one is Twitter, which has more than 140 million active users with 1 billion Tweets every 3 days as of March 2012.
- In China, Weibo ([www.weibo.com](http://www.weibo.com)) has accumulated more than 300 millions users in less than three years. Every second, more than 1000 Chinese tweets are posted in Weibo.

With the large volume and multi-aspect messages, how do users locate the specific messages that they are interested in?

# Motivation

## Example 1:





# Motivation

## Example 1:



DBS Bank

# Motivation

## Example 1:



DBS Bank

# Motivation

## Example 1:



# Motivation

## Example 2:

新浪微博 首页 广场 微吧 应用 游戏 搜索微博, 找人 Candace2

#航母style# 发布

### 热门话题

24小时 | 一周 | 一月

- 全国热点
- 区域热点
- 明星人物
- 兴趣
- 影视
- 情感
- 公益
- 行业
- 生活
- 族群
- 微博生态
- 大杂烩

**TOP 1**

**#请问玛雅人靠谱吗#** 519803

"玛雅体"未了, 玛雅人靠谱的话我就XXX

PK 你觉得玛雅人靠谱吗? (共有4412条观点PK)

红方: 靠谱啊, 所... (3004) VS 蓝方: 不靠谱, 老... (1408)

**TOP 2**

**#一九四二#** 265174

《1942》上映了, 来写一句话影评吧

微影评 (共有12条相关微博)

2012-11-11 13:25:16 **沈慕白** 发布了第一条事件相关微博

**TOP 3**

**#王的盛宴#** 98322

电影《王的盛宴》

微影评 (共有6条相关微博)

# Motivation

## Example 2:

The screenshot shows the Sina Weibo website interface. A red box highlights the left sidebar menu, which contains the following items: 热门话题 (Hot Topics), 全国热点 (National Hotspots), 区域热点 (Regional Hotspots), 明星人物 (Celebrity Figures), 兴趣 (Interests), 影视 (Movies/TV), 情感 (Relationships), 公益 (Public Welfare), 行业 (Industry), 生活 (Life), 族群 (Groups), 微博生态 (Weibo Ecosystem), and 大杂烩 (Miscellaneous). The main content area displays a list of trending topics:

- TOP 1 #请问玛雅人靠谱吗#** (519803)
  - “玛雅体”来了，玛雅人靠谱的话我就XXX
  - PK 你觉得玛雅人靠谱吗？（共有4412条观点PK）
  - 红方：靠谱啊，所... (3004) VS 蓝方：不靠谱，老... (1408)
- TOP 2 #一九四二#** (265174)
  - 《1942》上映了，来写一句话影评吧
  - 微影评（共有12条相关微博）
  - 2012-11-11 13:25:16 \_沈慕白 发布了第一条事件相关微博
- TOP 3 #王的盛宴#** (98322)
  - 电影《王的盛宴》
  - 微影评（共有6条相关微博）

# Motivation

## Example 2:

新浪微博 首页 广场 微吧 应用 游戏 搜索微博, 找人 Candace2

#航母style# 发布

热门话题

全国热点 24小时 | 一周 | 一月

区域热点

明星人物

兴趣

影视

情感

公益

行业

生活

族群

微博生态

大杂烩

TOP 1 #请问玛雅人靠谱吗# 519803

"玛雅体"未了, 玛雅人靠谱的话我就XXX

PK 你觉得玛雅人靠谱吗? (共有4412条观点PK)

TOP 2 #一九四二# 265174

《1942》上映了, 来写一句话影评吧

微影评 (共有12条相关微博)

2012-11-11 13:25:16 \_沈慕白 发布了第一条事件相关微博

TOP 3 #王的盛宴# 98322

电影《王的盛宴》

微影评 (共有6条相关微博)

# Motivation

## Example 2:

话题 > 热点

#一九四二# 分享 3767 全部讨论: 3790939

**《1942》上映了，来写一句话影评吧**

由冯小刚执导，根据刘震云小说改编的电影《一九四二》将于11月29日零点首映。本片讲述了1942年发生在河南的大饥荒。陈道明、张国立、李雪健、张涵予、阿德里安·布罗迪、希姆·罗宾斯等参演。已看过试映场的@张泉灵 @任志强 @马未都 等已经在微博上发表了对此片的感想。你会去看这部片吗？快来写下你的微影评吧！[\[专题\]电影吧热议中](#)

你看《1942》了吗？来写一句话影评吧` 发布

微影评

 **张泉灵**：刚在上海采访冯小刚，他讲述了他和「1942」18年的缘分，讲他当年和刘震云沿着灾民逃难的路径走的那几个月，讲这回演员们抛开自己的演技，生生冻着饿着自己。冯小刚、刘震云都是少见的聪明人，带着善的底子、较真的样子。可他们用最笨的办法拼出了一部戏。这世界的好东西常常是聪明人用笨办法做出来的。

11月23日 17:59 来自iPhone客户端 转发(5453) | 评论(1800)

网友投票

你给电影《一九四二》打几星？

单选 查看详情

- ☐ ★★★★★ (猛赞)
- ☐ ★★★★ (顶一下)
- ☐ ★★★ (马马虎虎)
- ☐ ★★ (就是个烂片)
- ☐ ★ (烂片中的奇葩)
- ☐ ☆ (我根本不想看)

投票 ☐ 匿名 ☒ 转发到微博

# Motivation

## Example 2:



How do we provide users an overviews of search results based on meaningful and structural categories.



# Motivation

## Example 2:

The screenshot shows a Weibo post from user 'huaili' about the movie '1942'. A red box highlights the main text of the post. The post includes a movie poster, a title, and a detailed description. To the right of the post is a '网友投票' (User Poll) section with five options for rating the movie.

**Highlighted Text:**

**#一九四二#** 分享 **3767** 全部讨论: 3790939

**《1942》上映了，来写一句话影评吧**

由冯小刚执导，根据刘震云小说改编的电影《一九四二》将于11月29日零点首映。本片讲述了1942年发生在河南的大饥荒。陈道明、张国立、李雪健、张涵予、阿德里安·布罗迪、希姆·罗宾斯等参演。已看过试映场的@张泉灵 @任志强 @马未都 等已经在微博上发表了对此片的感想。你会去看这部片吗？快来写下你的微影评吧！[专题]电影吧热议中]

你看《1942》了吗？来写一句话影评吧~ 发布

**微影评**

**张泉灵**：刚在上海采访冯小刚，他讲述了他和「1942」18年的缘分，讲他当年和刘震云沿着灾民逃难的路径走的那几个月，讲这演员们抛开自己的演技，生生冻着饿着自己。冯小刚、刘震云都是少见的聪明人，带着善的底子、较真的样子。可他们用最笨的办法拍出了一部戏。这世界的好东西常常是聪明人用笨办法做出来的。

11月23日 17:59 来自iPhone客户端 转发(5453) | 评论(1800)

**网友投票**

你给电影《一九四二》打几星？

单选 查看详情

- ☐ ★★★★★ (猛赞)
- ☐ ★★★★ (顶一下)
- ☐ ★★★ (马马虎虎)
- ☐ ★★ (就是个烂片)
- ☐ ★ (烂片中的奇葩)
- ☐ ☆ (我根本不想看)

投票 ☐ 匿名 ☒ 转发到微博

## Topic Classification!

# Outline

- 1 Background and Motivation
- 2 Related Work
- 3 Semi-Supervised Graphical Model
  - The General Framework
  - Probabilistic Graph Model Construction
  - Parameter Inference
- 4 Experiments
  - Experimental Settings
  - Analysis
  - Parameter Analysis
- 5 Conclusion and Future Work

# Related Work

## 1 Topic Model based Methods

- [Hong and Davison, 2010] employs latent dirichlet allocation (LDA) [Blei et al., 2003] and author-topic model [Rosen-Zvi et al., 2010] to deeply investigate to automatically find hidden topic structures on Twitter.
- Several variants of LDA to incorporate supervision have been proposed by [Ramage et al., 2009, Ramage et al., 2010], and have been shown to be competitive with strong baselines in the microblogging environment.

## 2 Traditional Classification Methods

- [Lee et al., 2011] classified tweets into pre-defined categories such as sports, technology, politics, etc. They constructed word vectors with tf-idf weights and utilized a Naive Bayesian Multinomial classifier to classify tweets.
- [Sriram et al., 2010] proposed to use a small set of domain-specific features extracted from the author's profile and text to represent short messages. Their method requires extensive pre-processing to conduct effectively feature analysis.

# Challenges and Contribution

## 1 Challenges

- Sparseness: lack sufficient word co-occurrence or shared contexts for effective similarity measure-[Hu et al., 2009].
- Informal: not well conformed as standard structures of documents.
- Lack of label information. It is time and labor consuming to label the huge amount of messages.

# Challenges and Contribution

## 1 Challenges

- Sparseness: lack sufficient word co-occurrence or shared contexts for effective similarity measure-[Hu et al., 2009].
- Informal: not well conformed as standard structures of documents.
- Lack of label information. It is time and labor consuming to label the huge amount of messages.

## 2 Contribution

- to handle data sparseness problem, we employ query related external resources from Google Search Engine to enrich the short messages.
- to alleviate negative effect brought by informal words, we utilize linguistic corpus to detect informal words and correct them.
- to require less labelled data, we attempt to use a semi-supervised learning approach for microblog categorization task.

# Outline

- 1 Background and Motivation
- 2 Related Work
- 3 Semi-Supervised Graphical Model
  - The General Framework
  - Probabilistic Graph Model Construction
  - Parameter Inference
- 4 Experiments
  - Experimental Settings
  - Analysis
  - Parameter Analysis
- 5 Conclusion and Future Work

# the General Framework

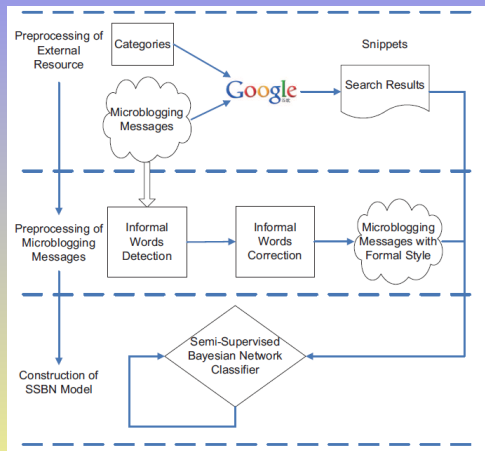
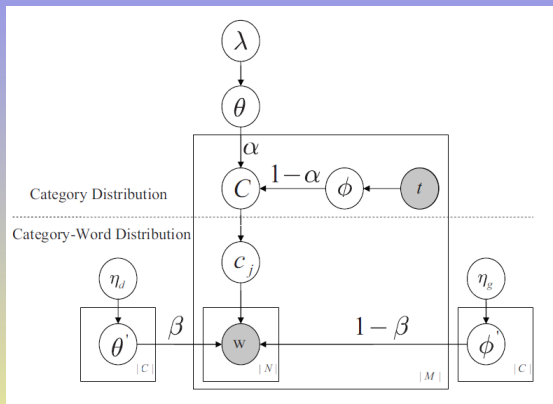


Figure: The General Framework.

# Semi-Supervised Bayesian Network Graph Model



**Figure:** Probabilistic graphical representation of semi-supervised Bayesian network model.



# Parameter Inference

The maximum likelihood category label for a given message  $m_i$  is,

$$y_i = \arg \max_{c_j} P(c_j | m_i, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') = \frac{P(c_j | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') P(m_i | c_j, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}')}{P(m_i | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}')}$$

# Parameter Inference

The maximum likelihood category label for a given message  $m_i$  is,

$$y_i = \arg \max_{c_j} P(c_j | m_i, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') = \frac{P(c_j | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') P(m_i | c_j, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}')}{P(m_i | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}')}$$

$$P(c_j | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') = P(c_j | \hat{\theta}, \hat{\phi}) = \hat{\alpha} P(c_j | \hat{\theta}) + (1 - \hat{\alpha}) P(c_j | \hat{\phi})$$

# Parameter Inference

The maximum likelihood category label for a given message  $m_i$  is,

$$y_i = \arg \max_{c_j} P(c_j | m_i, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') = \frac{P(c_j | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') P(m_i | c_j, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}')}{P(m_i | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}')}$$

$$P(c_j | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') = P(c_j | \hat{\theta}, \hat{\phi}) = \hat{\alpha} P(c_j | \hat{\theta}) + (1 - \hat{\alpha}) P(c_j | \hat{\phi})$$

$$P(m_i | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') = \sum_{c_j} P(c_j | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') P(m_i | c_j, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}')$$

# Parameter Inference

The maximum likelihood category label for a given message  $m_i$  is,

$$y_i = \arg \max_{c_j} P(c_j | m_i, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') = \frac{P(c_j | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') P(m_i | c_j, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}')}{P(m_i | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}')}$$

$$P(c_j | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') = P(c_j | \hat{\theta}, \hat{\phi}) = \hat{\alpha} P(c_j | \hat{\theta}) + (1 - \hat{\alpha}) P(c_j | \hat{\phi})$$

$$P(m_i | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') = \sum_{c_j} P(c_j | \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') P(m_i | c_j, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}')$$

$$\begin{aligned}
 P(m_i | c_j, \hat{\theta}, \hat{\phi}, \hat{\theta}', \hat{\phi}') &= P(m_i | c_j, \hat{\theta}', \hat{\phi}') = \prod_{k=1}^{|m_i|} P(w_k | c_j, \hat{\theta}', \hat{\phi}') \\
 &= \prod_{k=1}^{|m_i|} \{ \beta P(w_k | c_j, \hat{\theta}') + (1 - \beta) P(w_k | c_j, \hat{\phi}') \}
 \end{aligned}$$

# Estimating

## 1 Estimating $\theta$ :

$$\hat{\theta}_{c_j} \equiv P(c_j | \hat{\theta}) = \frac{1 + \sum_{i=1}^{|M|} \Lambda(i) P(y_i = c_j | m_i)}{|C| + |M^l| + \lambda |M^u|} \quad (1)$$

# Estimating

① Estimating  $\theta$ :

$$\hat{\theta}_{c_j} \equiv P(c_j | \hat{\theta}) = \frac{1 + \sum_{i=1}^{|M|} \Lambda(i) P(y_i = c_j | m_i)}{|C| + |M^l| + \lambda |M^u|} \quad (1)$$

② Estimating  $\phi$ :

$$\hat{\phi}_{c_j} \equiv P(c_j | \hat{\phi}) = \frac{\frac{1}{NGD(t, c_j)} + \mu}{\sum_{j=1}^{|C|} \frac{1}{NGD(t, c_j)} + |C| \mu} \quad (2)$$

# Estimating

① Estimating  $\theta$ :

$$\hat{\theta}_{c_j} \equiv P(c_j | \hat{\theta}) = \frac{1 + \sum_{i=1}^{|M|} \Lambda(i) P(y_i = c_j | m_i)}{|C| + |M^l| + \lambda |M^u|} \quad (1)$$

② Estimating  $\phi$ :

$$\hat{\phi}_{c_j} \equiv P(c_j | \hat{\phi}) = \frac{\frac{1}{NGD(t, c_j)} + \mu}{\sum_{j=1}^{|C|} \frac{1}{NGD(t, c_j)} + |C| \mu} \quad (2)$$

③ Estimating  $\theta'$  and  $\phi'$ :

$$\hat{\theta}'_{c_j}{}^{w_k} \equiv P(w_k | c_j, \hat{\theta}') = \frac{n_{dc_j}^{w_k} + \eta_d}{\sum_{p'=1}^{|N|} n_{dc_j}^{w_{p'}} + |N| \eta_d} \quad (3)$$

$$\hat{\phi}'_{c_j}{}^{w_k} \equiv P(w_k | c_j, \hat{\phi}') = \frac{n_{gc_j}^{w_k} + \eta_g}{\sum_{q'=1}^{|N|} n_{gc_j}^{w_{q'}} + |N| \eta_g} \quad (4)$$

# Outline

- 1 Background and Motivation
- 2 Related Work
- 3 Semi-Supervised Graphical Model
  - The General Framework
  - Probabilistic Graph Model Construction
  - Parameter Inference
- 4 Experiments**
  - Experimental Settings
  - Analysis
  - Parameter Analysis
- 5 Conclusion and Future Work



# Datasets and Evaluation Metrics

Twitter		Sina Weibo	
Total	16935	Total	15811
Sports	2720	Sports	2602
Entertainment	2816	Movies	2694
Business	2912	Games	2605
Science&Tech	2827	Science&Tech	2647
Politics	2937	Politics	2654
Education	2723	Music	2609

Table: The distribution of different categories over two datasets.

# Datasets and Evaluation Metrics

Twitter		Sina Weibo	
Total	16935	Total	15811
Sports	2720	Sports	2602
Entertainment	2816	Movies	2694
Business	2912	Games	2605
Science&Tech	2827	Science&Tech	2647
Politics	2937	Politics	2654
Education	2723	Music	2609

**Table:** The distribution of different categories over two datasets.

## ① Apple, stock business

- iBenApple Mon Jan 24 13:50:42 +0000 2011 #IHateItWhen Apple's stock continue to fall!

## ② Apple, ipad science

- Kericox3 Tue Feb 01 12:34:55 +0000 2011 Apple iphone 4g 32gb and blackberry bold 9700 Unlocked. - Anything ...: Apple Tablet iPad 64GB (Wi-Fi + 3G) .....  
<http://bit.ly/gbbW1J>

# Datasets and Evaluation Metrics

Twitter		Sina Weibo	
Total	16935	Total	15811
Sports	2720	Sports	2602
Entertainment	2816	Movies	2694
Business	2912	Games	2605
Science&Tech	2827	Science&Tech	2647
Politics	2937	Politics	2654
Education	2723	Music	2609

**Table:** The distribution of different categories over two datasets.

- accuracy
- precision
- recall
- $F_1$

# SSBN Model Performance

Twitter				Sina Weibo			
<i>Category</i>	<i>Precision</i>	<i>Recall</i>	<i>F1</i>	<i>Category</i>	<i>Precision</i>	<i>Recall</i>	<i>F1</i>
Sports	<b>0.9322</b>	0.9483	<b>0.9402</b>	Sports	<b>0.9318</b>	0.8747	<b>0.9023</b>
Entertainment	0.9000	0.5625	0.6923	Movies	0.8848	0.8207	0.8515
Business	0.8043	0.5323	0.6382	Games	0.8090	0.9283	0.8646
Science&Tech	0.6937	<b>0.9801</b>	0.8124	Science&Tech	0.8688	0.8323	0.8502
Politics	0.9096	0.9640	0.9360	Politics	0.8661	<b>0.9324</b>	0.8980
Education	0.5000	0.5519	0.5165	Music	0.8819	0.8699	0.8759
Micro-average	0.7979	0.7979	0.7979	Micro-average	0.8798	0.8798	0.8798
Macro-average	0.7934	0.6043	0.6128	Macro-average	0.8737	0.8764	0.8738

**Table:** Performance of SSBN model on two datasets with 5% training data and 95% testing data, respectively.

# Baselines

- SVM
- Naive Bayesian
- K Nearest Neighbors
- Rocchio
- Labeled LDA
- Transductive SVM
- Semi-Naive Bayesian classifier

# Comparison Performance

<i>Classifier</i>	<i>Accuracy</i>	<i>MicroP</i>	<i>MicroR</i>	<i>MicroF1</i>	<i>MacroP</i>	<i>MacroR</i>	<i>MacroF1</i>
SSBN	<b>0.8875</b>	<b>0.8875</b>	<b>0.8875</b>	<b>0.8875</b>	0.8282	0.7627	0.7845
SVM	0.8670	0.8670	0.8670	0.8670	0.8768	0.7611	<b>0.7860</b>
NB	0.8722	0.8696	0.8722	0.8722	<b>0.8879</b>	0.7329	0.7587
KNN	0.7268	0.7268	0.7268	0.7268	0.6721	0.6471	0.6516
Rocchio	0.8180	0.8204	0.8180	0.8192	0.7361	<b>0.8384</b>	0.7605
L-LDA	0.8605	0.8605	0.8605	0.8605	0.8467	0.7223	0.7532

**Table:** Performance comparison among SSBN and other supervised baseline methods on twitter with 90% training data.

<i>Classifier</i>	<i>Accuracy</i>	<i>MicroP</i>	<i>MicroR</i>	<i>MicroF1</i>	<i>MacroP</i>	<i>MacroR</i>	<i>MacroF1</i>
SSBN	<b>0.7979</b>	<b>0.7979</b>	<b>0.7979</b>	<b>0.7979</b>	<b>0.7934</b>	<b>0.6043</b>	<b>0.6128</b>
Trans-SVM	0.6707	0.6707	0.6707	0.6707	0.6602	0.5108	0.4491
Semi-NB	0.7156	0.7156	0.7156	0.7156	0.7308	0.5653	0.549

**Table:** Performance comparison among SSBN and other semi-supervised baseline methods on Twitter with 5% training data.

# Comparison Performance

<i>Classifier</i>	<i>Accuracy</i>	<i>MicroP</i>	<i>MicroR</i>	<i>MicroF1</i>	<i>MacroP</i>	<i>MacroR</i>	<i>MacroF1</i>
SSBN	<b>0.9020</b>	<b>0.9020</b>	<b>0.9020</b>	<b>0.9020</b>	0.8976	<b>0.9045</b>	<b>0.9004</b>
SVM	0.8991	0.8991	0.8991	0.8991	<b>0.9017</b>	0.8971	0.8991
NB	0.9015	0.9015	0.9015	0.9015	0.8990	0.9024	0.9003
KNN	0.8565	0.8565	0.8565	0.8565	0.8589	0.8486	0.8526
Rocchio	0.8802	0.8803	0.8802	0.8802	0.8769	0.8832	0.8781
L-LDA	0.8905	0.8905	0.8905	0.8905	0.8876	0.8989	0.8932

**Table:** Performance comparison among SSBN and other supervised baseline methods on Sina Weibo with 90% training data.

<i>Classifier</i>	<i>Accuracy</i>	<i>MicroP</i>	<i>MicroR</i>	<i>MicroF1</i>	<i>MacroP</i>	<i>MacroR</i>	<i>MacroF1</i>
SSBN	<b>0.8798</b>	<b>0.8798</b>	<b>0.8798</b>	<b>0.8798</b>	<b>0.8737</b>	<b>0.8764</b>	<b>0.8738</b>
Trans-SVM	0.8084	0.8084	0.8084	0.8084	0.8049	0.8085	0.8052
Semi-NB	0.8198	0.8198	0.8198	0.8198	0.8225	0.8217	0.8204

**Table:** Performance comparison among SSBN and other semi-supervised baseline methods on Sina Weibo with 5% training data.

# On the Sensitivity of Training Data Size

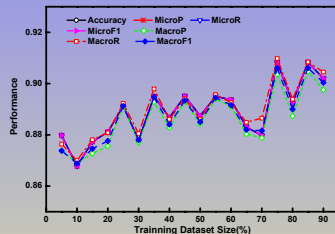
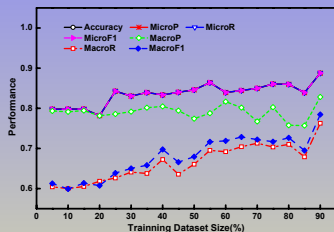
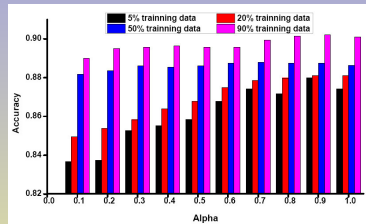
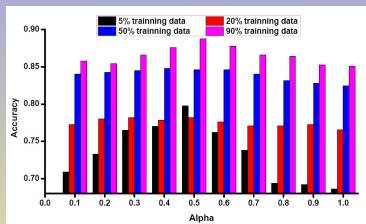


Figure: Performance sensitivity of training set size on Twitter and Sina Weibo



# Effect of $\alpha$

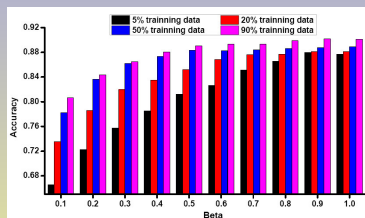
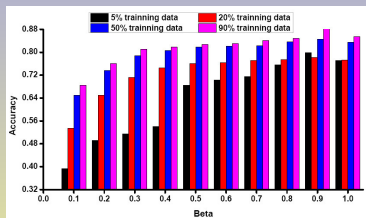
The trade-off parameter  $\alpha$  is used to balance the effects of two kinds of prior knowledge at category level: microblogging data collection and external resources.



**Figure:** The Performance with varying  $\alpha$  and training data size when other parameters are fixed.

# Effect of $\beta$

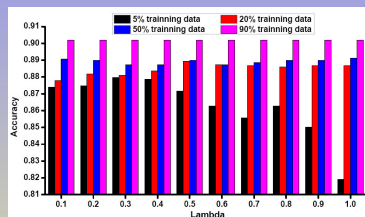
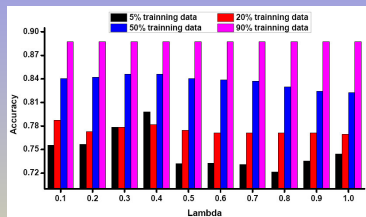
There are two category-word distributions,  $\theta'$  and  $\phi'$ , which are respectively generated from our data collection and google search results; and parameter  $\beta$  is utilized to adjust the contribution between these two different resources in category-word level.



**Figure:** The Performance with varying  $\beta$  and training data size when other parameters are fixed.

# Effect of $\lambda$

$\lambda$  indicates the contribution from unlabeled data points, between 0 and 1.



**Figure:** The Performance with varying  $\lambda$  and training data size when other parameters are fixed.

# Outline

- 1 Background and Motivation
- 2 Related Work
- 3 Semi-Supervised Graphical Model
  - The General Framework
  - Probabilistic Graph Model Construction
  - Parameter Inference
- 4 Experiments
  - Experimental Settings
  - Analysis
  - Parameter Analysis
- 5 Conclusion and Future Work

## ① Conclusion

- the incorporation of external resources to supplement the short microblogs well compensates the data sparseness issue;
- the semi-supervised classifier seamlessly fuse labeled data structure and external resources into the training process, which reduced the requirement for manually labeling to a certain degree;
- we model the category probability of a given message based on the category-word distribution, and this successfully avoided the difficulty brought about by the spelling errors that are common in microblogging messages.

## ① Conclusion

- the incorporation of external resources to supplement the short microblogs well compensates the data sparseness issue;
- the semi-supervised classifier seamlessly fuse labeled data structure and external resources into the training process, which reduced the requirement for manually labeling to a certain degree;
- we model the category probability of a given message based on the category-word distribution, and this successfully avoided the difficulty brought about by the spelling errors that are common in microblogging messages.

## ② Future Work

- the incorporation of social network structure can improve the performance of microblogging classification;
- the use of external resources such as Wikipedia and WordNet might be valuable for understanding microblogging messages;
- the provision of category summarization can help to organize microblogging messages.

# References I



Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003).  
Latent dirichlet allocation.  
*Journal of Machine Learning Research*, 3:993–1022.



Hong, L. and Davison, B. D. (2010).  
Empirical study of topic modeling in twitter.  
*In Proceedings of KDD Workshop on Social Media Analytics*.



Hu, X., Sun, N., Zhang, C., and Chua, T.-S. (2009).  
Exploiting internal and external semantics for the clustering of short texts using world knowledge.  
*In Proceedings of the ACM conference on Information and knowledge management*.



Lee, K., Palsetia, D., Narayanan, R., Patwary, M. M. A., Agrawal, A., and Choudhary, A. (2011).  
Twitter trending topic classification.  
*In Proceedings of ICDM Workshop on Optimization Based Methods for Emerging Data Mining Problems*.



Ramage, D., Dumais, S., and Liebling, D. (2010).  
Charaterizing microblog with topic models.  
*In Proceedings of International AAAI Conference on Weblogs and Social Media*.

# References II



Ramage, D., Hall, D., Nallapati, R., and Manning, C. D. (2009).  
Labeled LDA: a supervised topic model for credit attribution in multi-labeled corpora.  
In *Proceedings of International Conference on Empirical Methods in Natural Language Processing*.



Rosen-Zvi, M., Chemudugunta, C., Griffiths, T., Smyth, P., and Steyvers, M. (2010).  
Learning author-topic models from text corpora.  
*ACM Transactions on Information Systems*, 28:1–38.



Sriram, B., Fuhry, D., Demir, E., Ferhatosmanoglu, H., and Demirbas, M. (2010).  
Short text classification in twitter to improve information filtering.  
In *Proceedings of Annual ACM Conference on Research and Development in Information Retrieval*.



# Thank you!

