

Identifying Influential Users' Professions via the Microblogs They Forward

Yuan Wang, Hangyu Mao, Zhen Xiao

Peking University, China

SocInf2017, 19 Aug 2017

Outline

- Background
- Data
- Our Method
- Evaluation
- Conclusion

Introduction

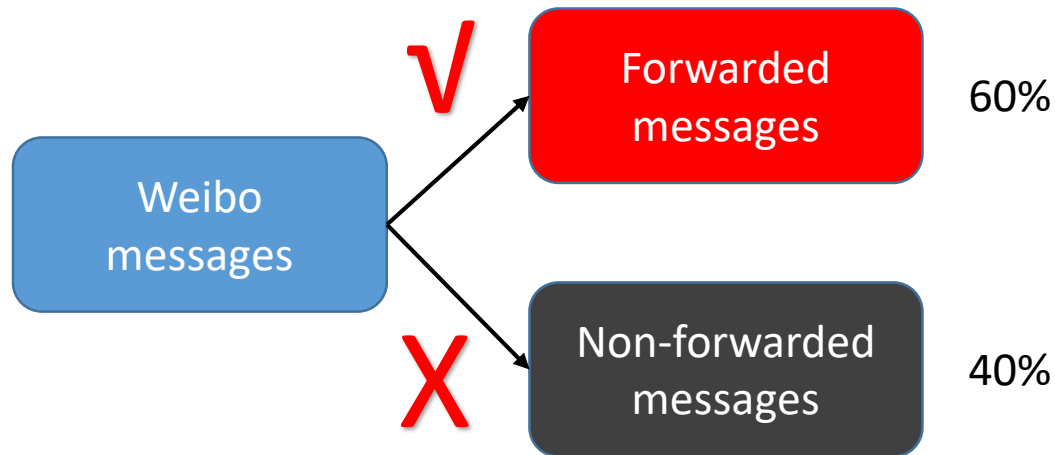


Profession



users are mainly organized by their professions

Message forwarding



马超 Terminal ★

3月7日 14:19 来自 iPhone 6

第二个工作我们也在做。

@王威廉 ✓

今天陪了谷歌Jeff Dean一天，总结几个有趣的事情：1) 谷歌邮箱自动回复功能早在2009年还是愚人节玩笑，但2015年后就被实现应用了。2) 谷歌正在用强化学习研究如何给GPU和CPU分别分配计算任务。3) 谷歌TPU论文刚被ISCA 2017接受了，论文马上放出来。🤔

3月7日 13:54 来自 微博 weibo.com

🔗 52 | 💬 9 | 👍 83

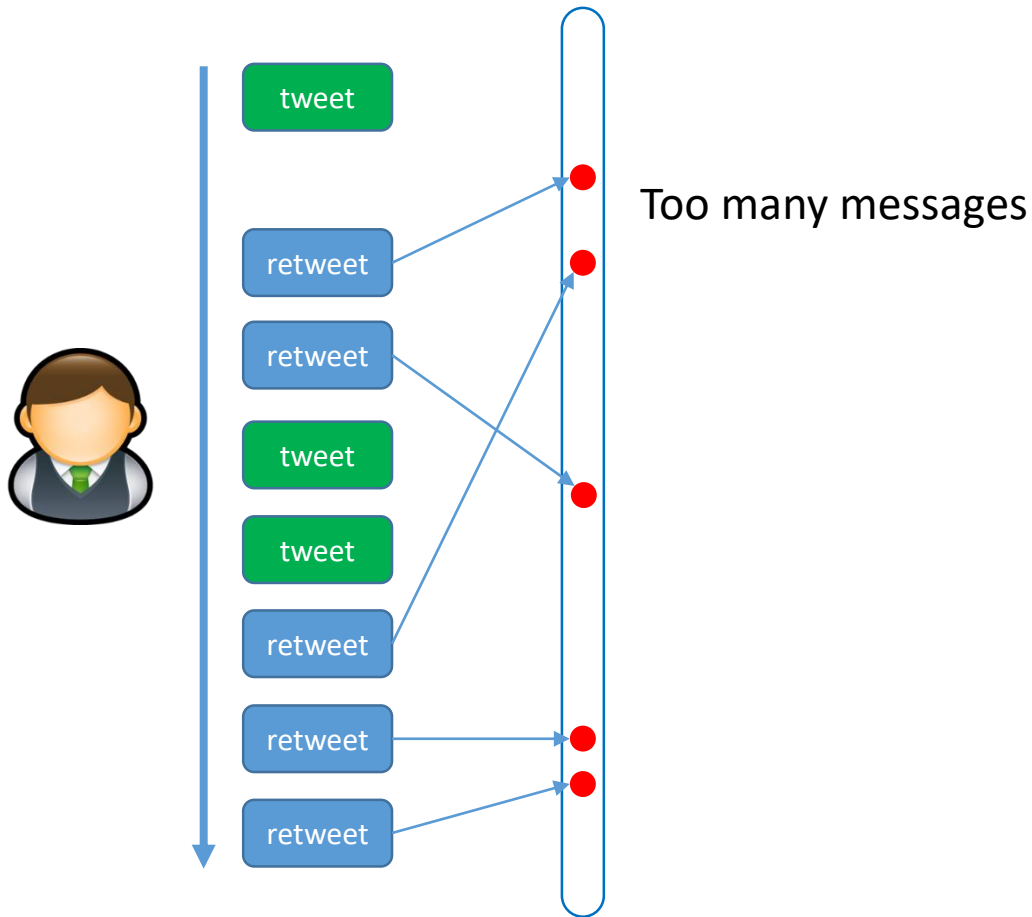
☆ 收藏

🔗 3

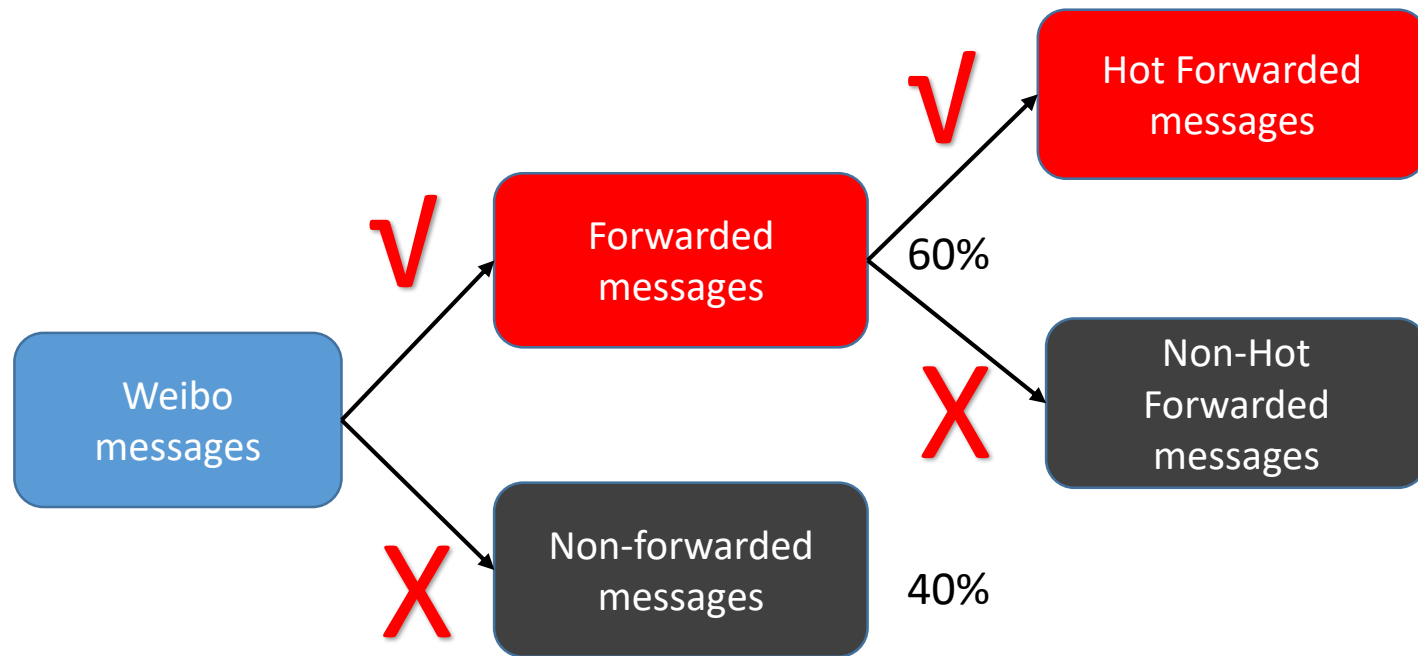
💬 评论

👍 3

Challenge



Challenge



Our Framework

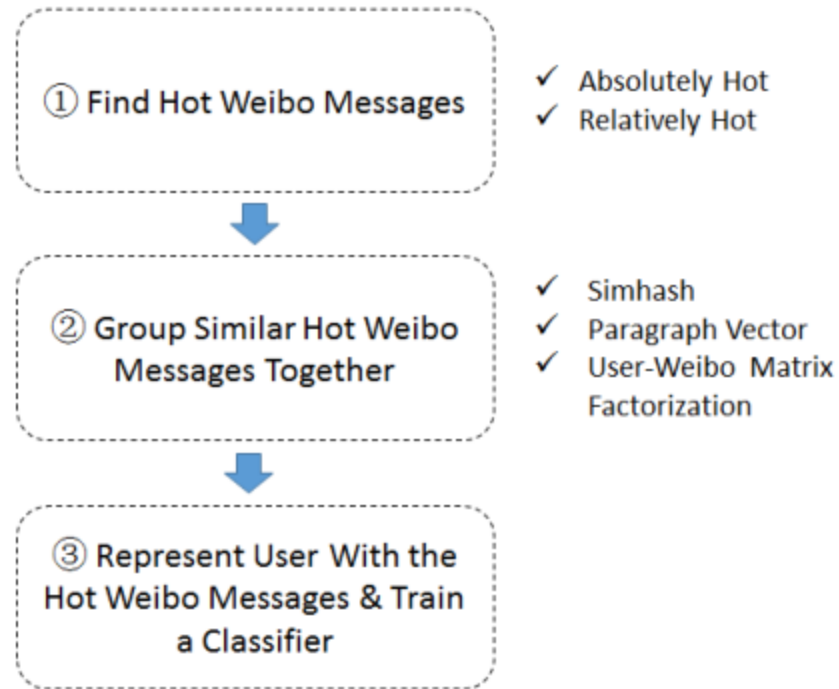


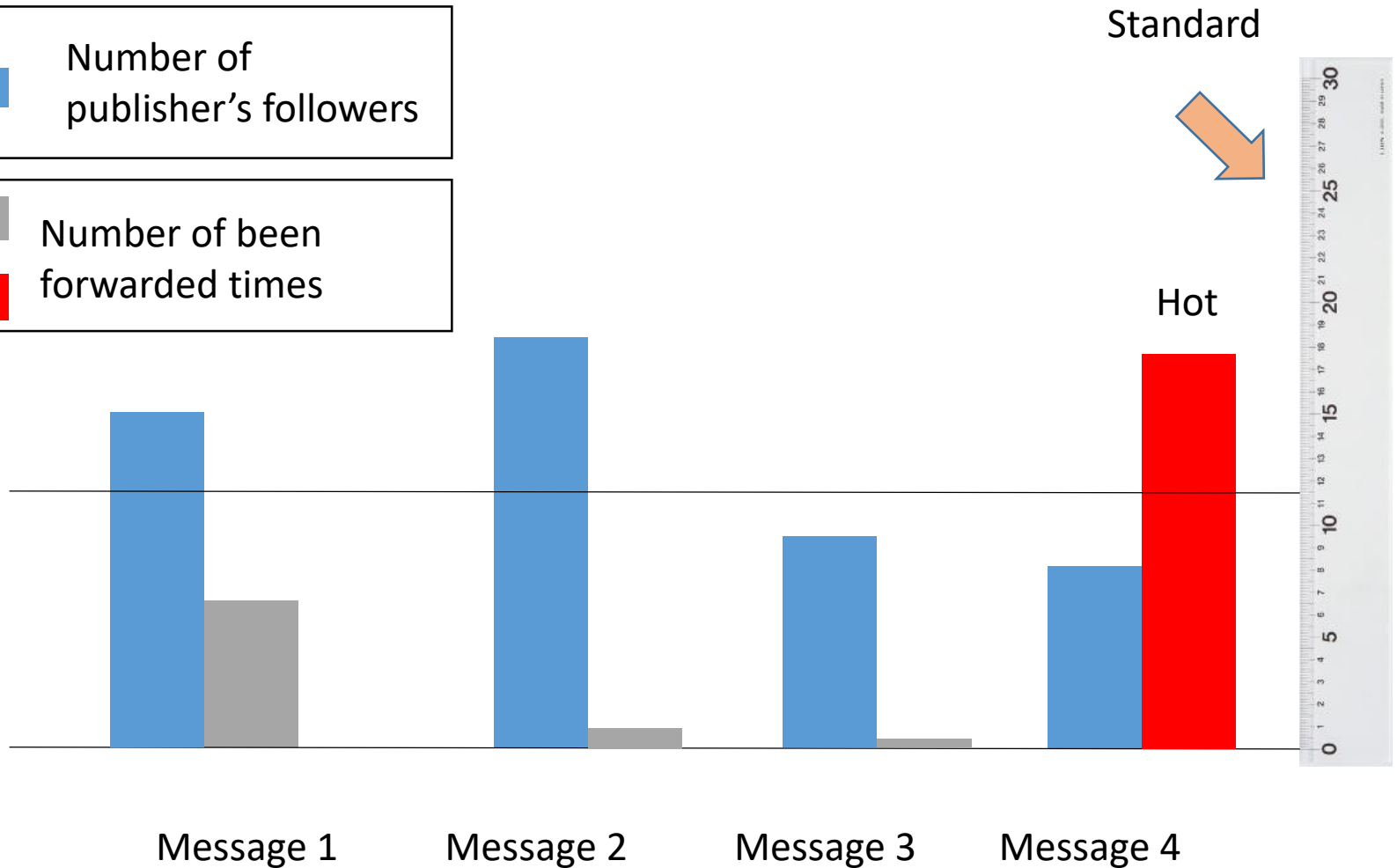
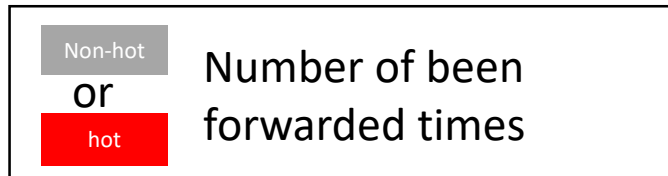
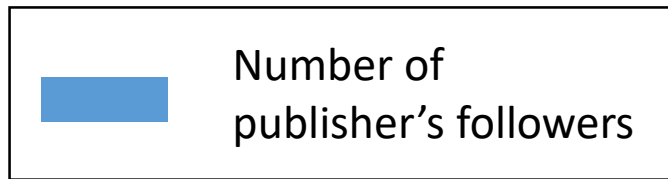
Fig. 1. The framework of PIFB

Our Data

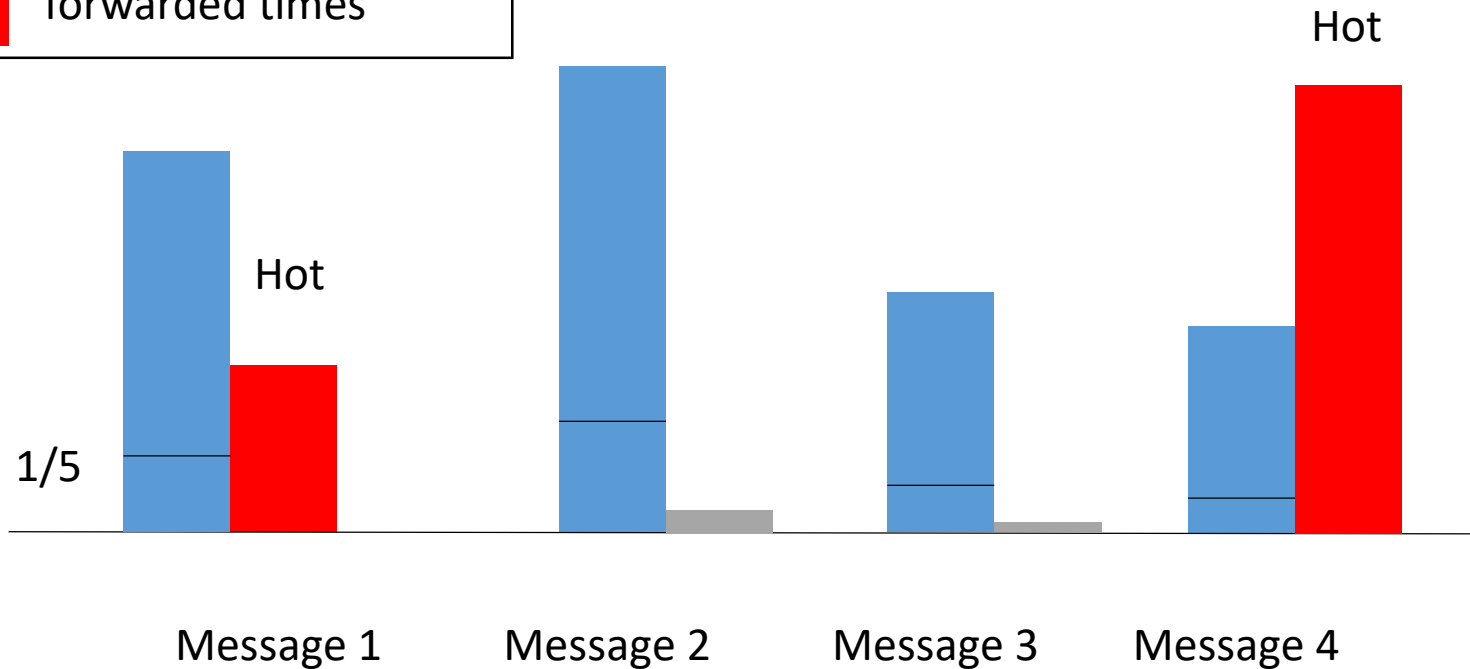
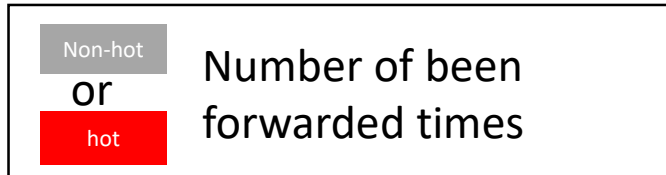
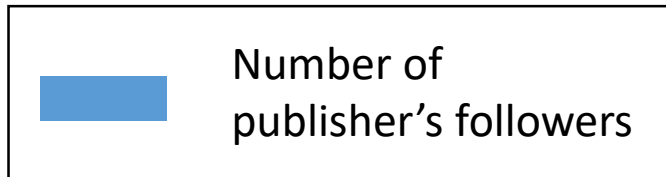
No.	Category	Ratio (%)
1	Media	26.3
2	Entertainment	10.1
3	Estate	9.1
4	Finance	8.6
5	Government	8.5
6	IT	8.4
7	Sports	6.4
8	Fashion	6.2
9	Education	5.9
10	Literature	5.4
11	Game	5.1

41,531 manually annotated influential users

Absolutely Hot Message



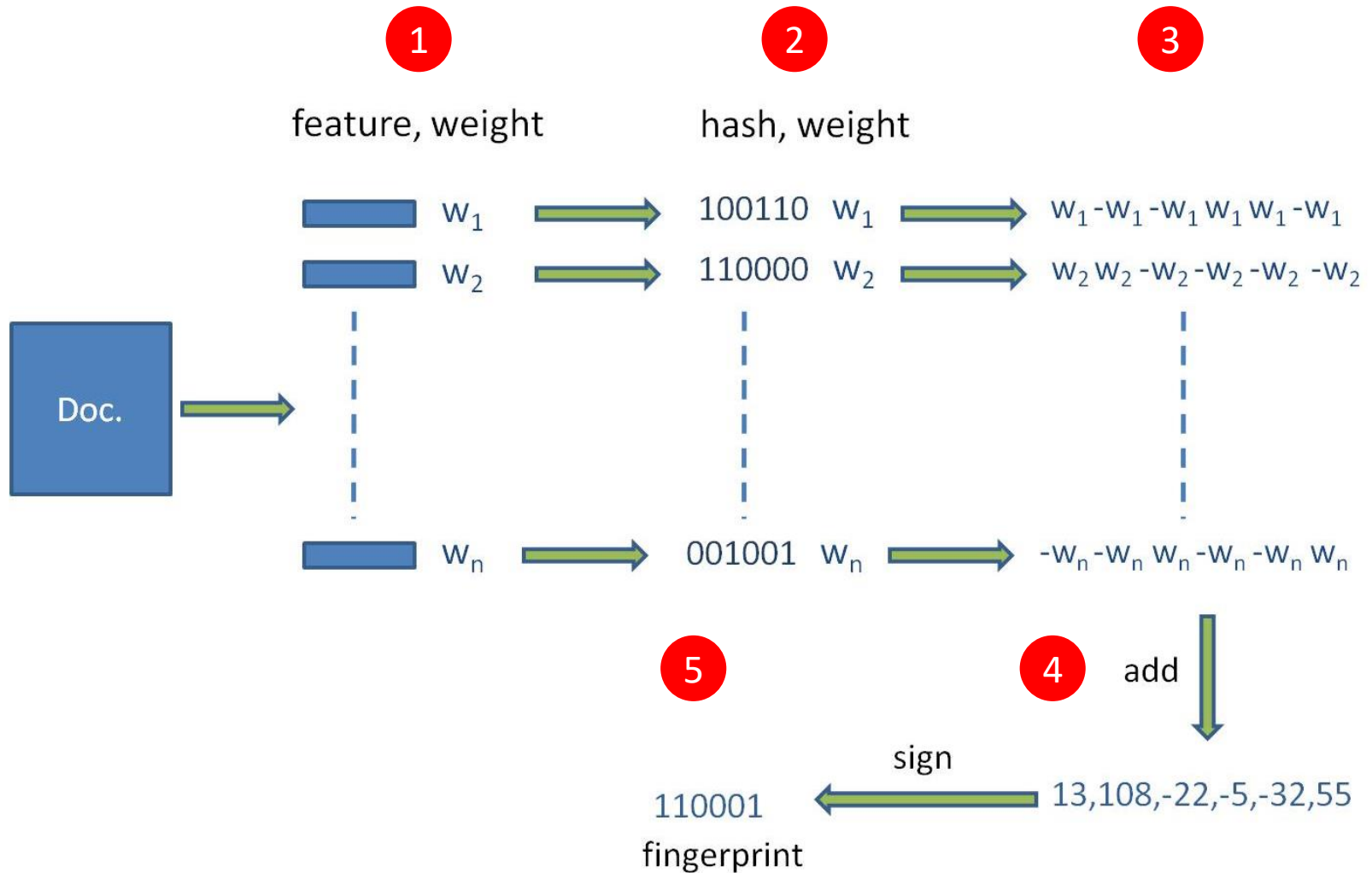
Relatively Hot Message



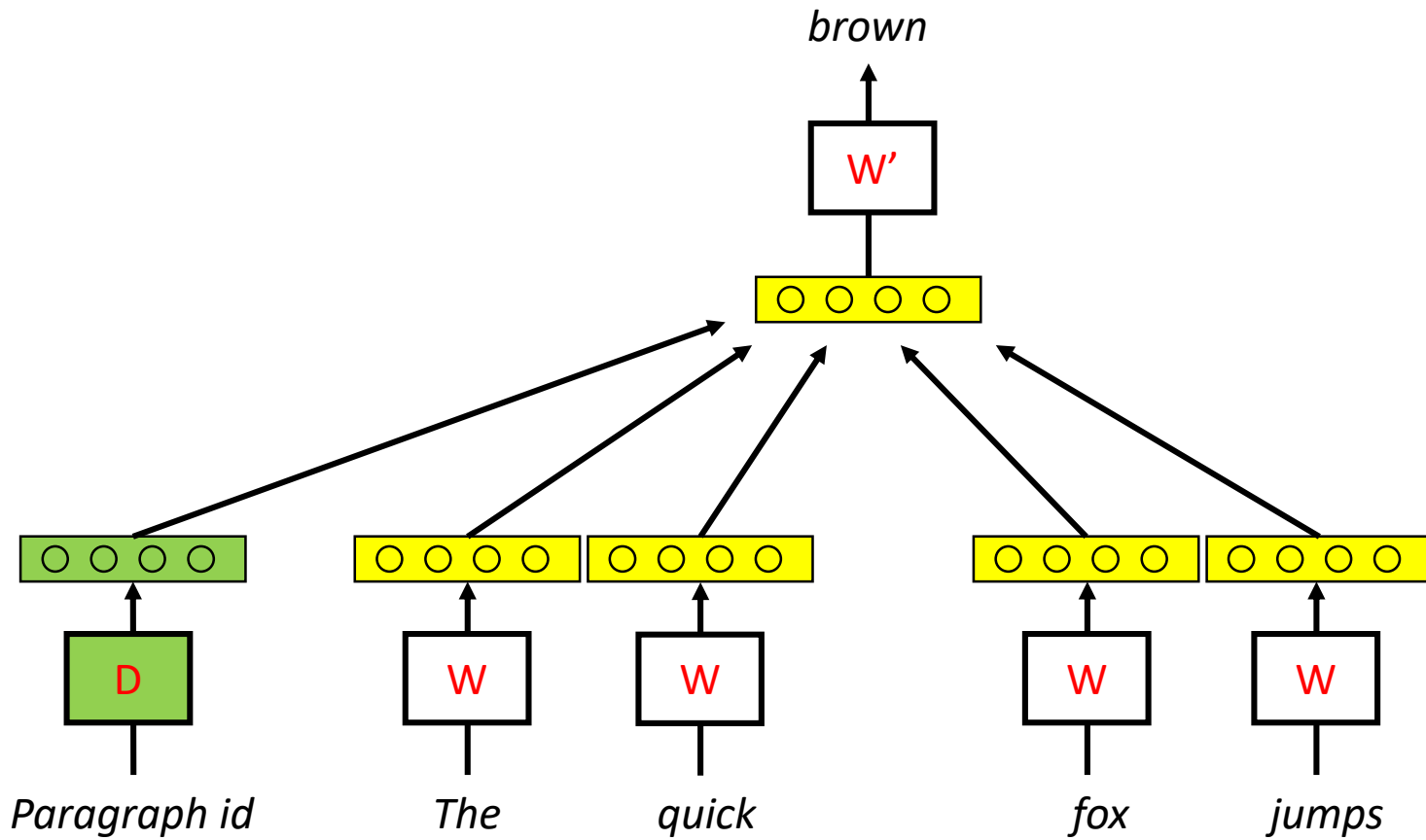
Group Similar Hot Messages Together

- Simhash
- Paragraph2vec
- Matrix factorization

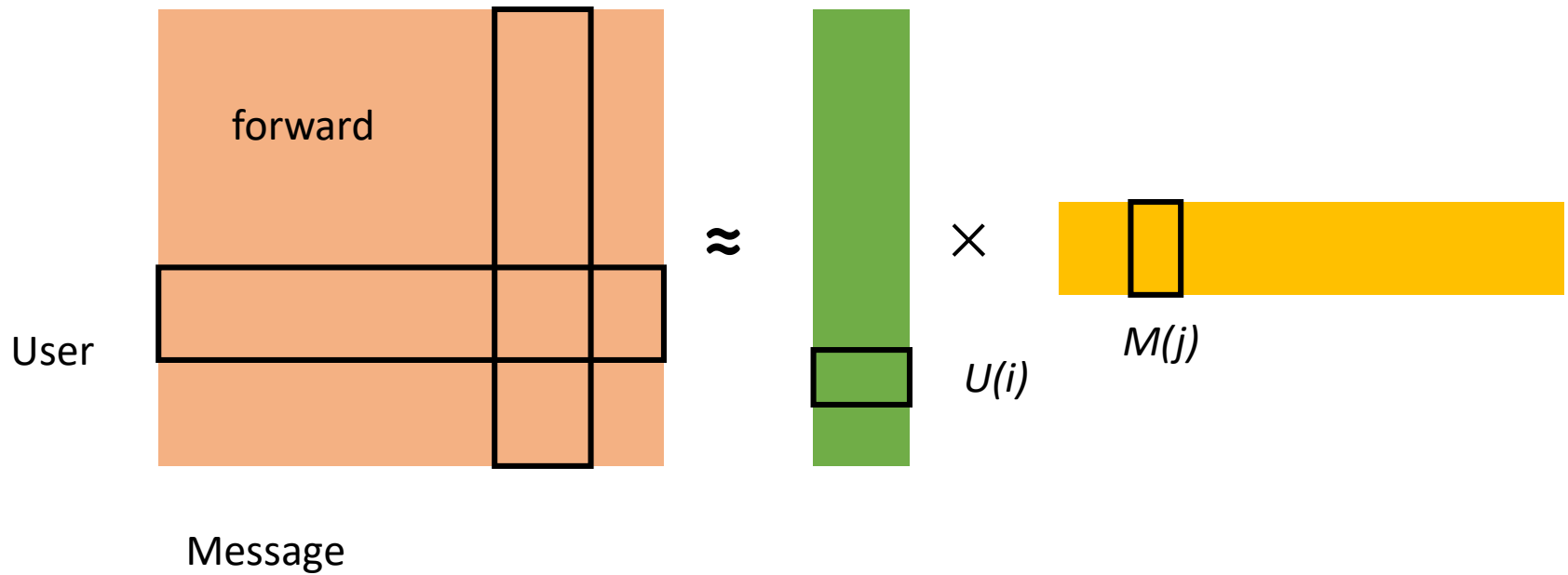
Simhash



Paragraph2Vec

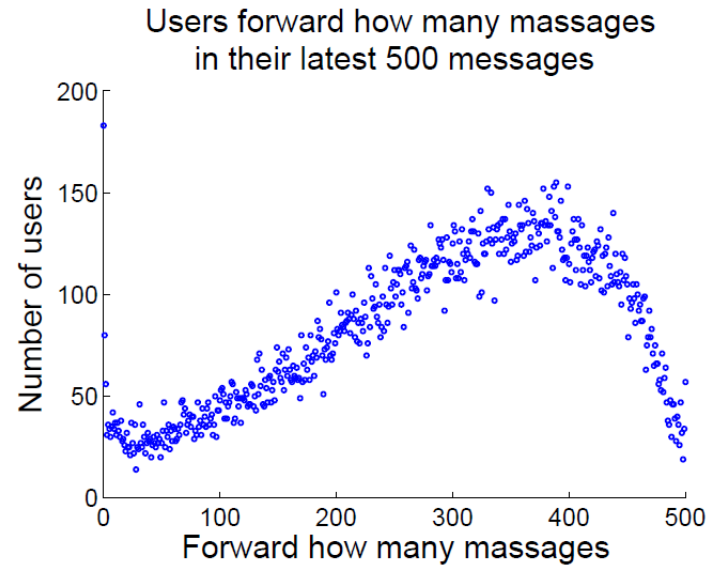
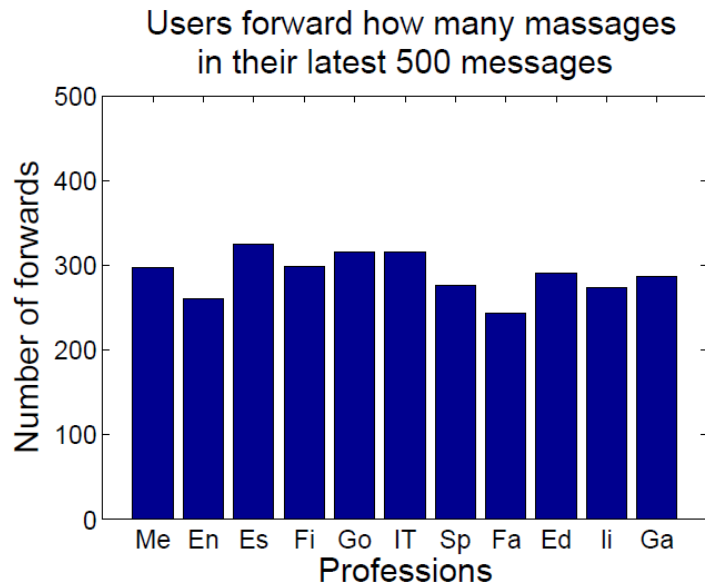


User-Weibo Matrix Factorization



User-weibo relationship matrix M

Evaluation



Evaluation

- Absolute hot
- Relative hot

Table 2. Number of absolutely hot messages

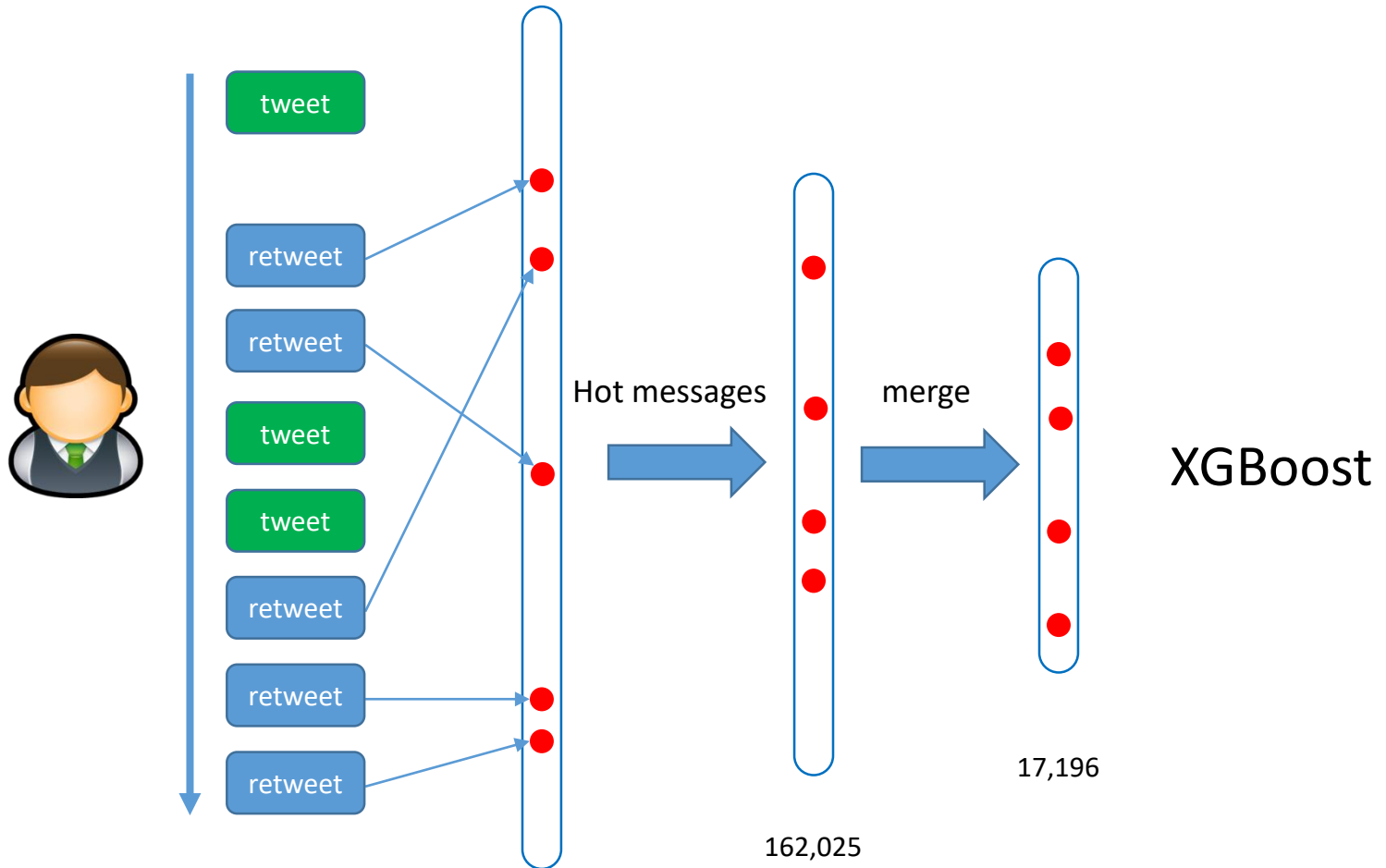
No.	Threshold	# before filter	# after filter
1	500	731,150	100,219
2	2000	426,019	82,339
3	10000	74,308	32,955

Evaluation

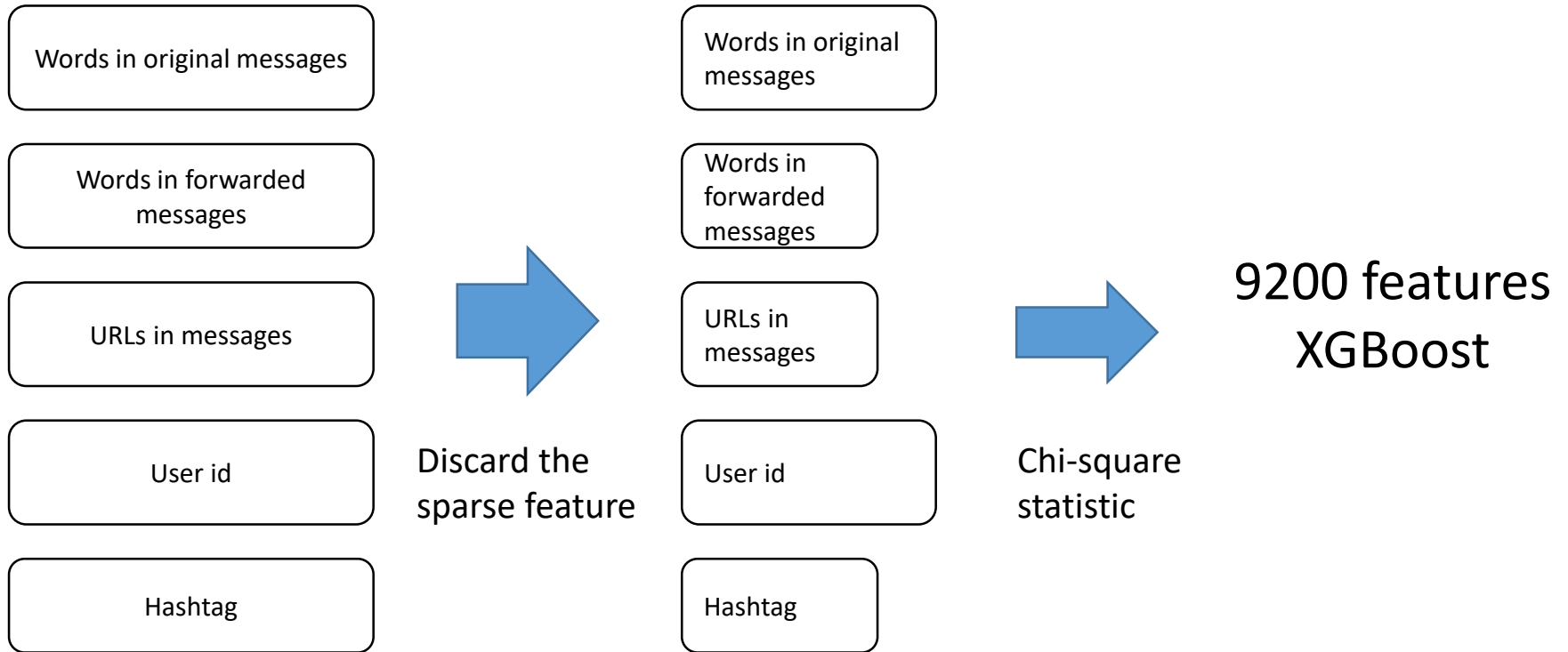
Table 3. Number of messages under different merging strategies

No.	Merging Strategy	# before	# after
1	Simhash	162,025	57,624
2	P2V	162,025	32,118
3	UWMF	162,025	27,129
4	Simhash+P2V+UWMF	162,025	17,196

Evaluation



Baseline



Evaluation

No.	Method	Accuracy	Precision	Recall	F1
1	Baseline	62.38	64.03	60.29	62.10
2	Simhash	69.24 \uparrow 6.86%	70.88	67.61	69.21 \uparrow 7.11%
3	Simhash+P2V	73.79 \uparrow 11.41%	73.90	71.28	72.57 \uparrow 10.47%
4	Simhash+P2V+UWMF	73.98 \uparrow 11.60%	74.81	72.95	73.87 \uparrow 11.77%

Table 4: Evaluation results for various features and combinations. (%)

Evaluation

	Me	En	Es	Fi	Go	IT	Sp	Fa	Ed	Li	Ga
Me	76.7	5.6	2.7	3.5	4.1	3.3	2.2	0.9	0.6	0.2	0.2
En	7.2	74.5	0.2	3.3	0.7	1.4	4.4	5.1	0.2	1.3	1.7
Es	7.4	2.0	72.9	8.5	5.3	2.2	0.4	0.9	0.1	0.0	0.3
Fi	8.4	0.1	6.4	70.2	5.3	6.2	0.2	1.3	1.7	0.1	0.1
Go	4.9	2.2	0.4	4.2	78.2	2.9	4.1	0.4	2.5	0.2	0.0
IT	6.1	0.7	3.9	4.3	1.3	76.3	0.2	0.1	2.6	0.7	3.8
Sp	5.1	2.9	0.0	0.3	0.3	1.0	86.2	2.2	0.7	0.0	1.3
Fa	9.7	14.9	1.0	6.2	0.2	0.0	3.3	61.5	0.9	1.2	1.1
Ed	5.2	3.9	3.3	4.6	2.0	3.2	0.7	1.8	68.4	4.2	2.7
Li	13.7	7.2	0.7	1.3	0.6	1.4	0.4	3.3	9.8	60.9	0.7
Ga	5.2	3.9	0.8	0.0	0.4	7.1	1.2	4.3	0.1	0.3	76.7

Table 5: Distribution of identified professions in each profession.

Conclusion

- find the hot weibo messages
- group similar hot messages together
- design a classifier to predict users' professions

Questions