概率主题模型及其应用

易磊,李晓东

大连理工大学软件学院

May 24, 2012

- Introduction to topic model
 - What are topic model?
 - LDA
- 2 Applications
 - Topic model in text mining
 - Topic model in image processing
- Model extensions
- 4 Conclusion
- 6 Reference
- 6 The end

Information overload



 As more information becomes available, it becomes more difficult to find and discover what we need.

Information overload



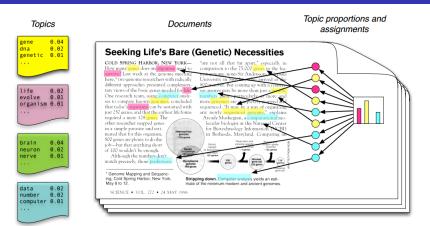
- As more information becomes available, it becomes more difficult to find and discover what we need.
- We need new tools to help us organize, search, and understand these vast amounts of information.

Information overload



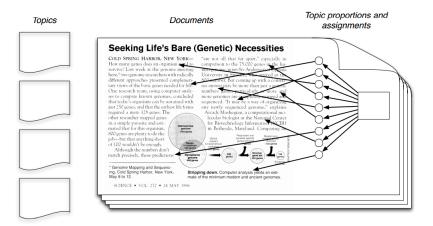
- As more information becomes available, it becomes more difficult to find and discover what we need.
- We need new tools to help us organize, search, and understand these vast amounts of information.
- Topic modeling provides methods for automatically organizing, understanding, searching, and summarizing large electronic archives.

Latent Dirichilet allocation



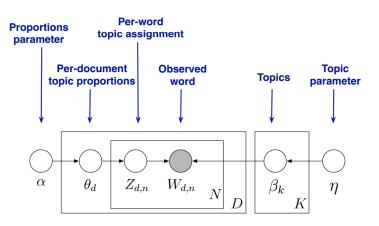
- Each topic is a distribution over words.
- Each document is a mixture of topics.
- Each word is drawn from one of those topics.

The posterior distribution



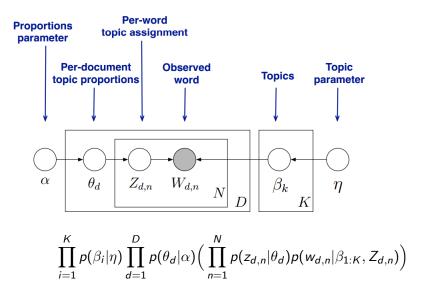
- In reality, we only observe the documents.
- The other structure are hidden variables
- Our goal is to infer the hidden variables.

LDA as a graphicial model

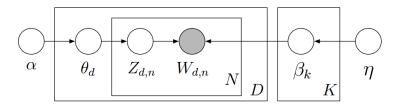


- Nodes are random variables; edges indicate dependence.
- Shaded nodes are observed.
- Plates indicate replicated variables.

LDA as a graphicial model



LDA



This joint defines a posterior approximate posterior inference algorithms:

- Gibbs sampling (Griffiths and Stryvers, 2002)
- variational inference (Teh et at., 2006)

- Introduction to topic model
 - What are topic model?
 - LDA
- 2 Applications
 - Topic model in text mining
 - Topic model in image processing
- Model extensions
- 4 Conclusion
- 6 Reference
- 6 The end

分析新浪微博知名微博主的微博主题(just a toy)

- 实验数据
 - 李开复(创新工场董事长兼首席执行官) 4000条最新微博
 - 杨幂(中国大陆知名女演员) 2000条最新微博
 - 36氪(关注互联网创业的科技博客) 7700条最新微博
- 预处理
 - 分词
 - 去除停用词
- LDA建模,使用LDA的开源实现gensim
- 实验环境及参数
 - 机器配置: Intel(R) Core(TM)2 Duo CPU E7400 @ 2.80GHz, 内存2G, 操作系统Ubuntu 11.04
 - 主题个数设为10

实验结果

• 李开复的微博主题

```
2012-05-22 21:42:11,659 : INFO : topic #8: 0.004*公司 + 0.004**恒国 + 0.003**创业 + 0.003**创业 + 0.003*英国 + 0.003*英国 + 0.003*工场 + 0.003*世 + 0.003*开 + 0.003*世 + 0.003*世 + 0.003*世 + 0.003*世 + 0.003*世 + 0.003*世 + 0.003*开 + 0.003*开 + 0.003*开 + 0.003*世 + 0.003*
```

• 杨幂的微博主题

2012-05-22 21:54:07.285: INFO: topic #1: 0.010+心 + 0.009+酸 + 0.009+物 + 0.005+節 + 0.005+野 + 0.005+野 前 + 0.005+野 前 + 0.004+宰福 + 0.004+常祖 + 0.005+教育 + 0.003+教育 - 0.003+物子 - 0.00	0.003*幸福	1111 0 .		0.011 az		. 0.000 /4			. 0.004 (6) (6)		
2012-05-22 21:54:07,294: INFO: topic #2: 0.008*数 + 0.007*心 + 0.007*心 + 0.007*前 + 0.005*常補 + 0.004*財的 + 0.004*財的 + 0.004*財制 + 0.004*財制 + 0.004*財		INFO:		0. 010*心	+ 0.009*爱	+ 0.006*@	+ 0.005*泪 +	0.005*好 +	0.005*谢谢 + (0.004*幸福	+ 0.004*" + 0.00
0.003*検系 2012-05-22 21:54:07,296: INFO: topic #3: 0.009*酸 + 0.007*心 + 0.006*® + 0.005*前 + 0.005*前 + 0.004*戦的 + 0.004*耶 + 0.004*谢谢 + 0.004*谢谢 + 0.003*今天 2012-05-22 21:54:07,298: INFO: topic #4: 0.009*酸 + 0.006*® + 0.005*心 + 0.005*前 + 0.005*技的 + 0.005*竹 + 0.004*谢谢 + 0.004*谢谢 + 0.003*你的 2012-05-22 21:54:07,398: INFO: topic #5: 0.009*酸 + 0.007*心 + 0.005*通 + 0.005*技的 + 0.004*谢谢 + 0.004*谢谢 + 0.004*附 + 0.004*好 + 0.004*好 + 0.004*报 + 0.005*技的 + 0.004*谢谢 + 0.004*附 + 0.004*好 + 0.004*	0.004*耶										
2012-05-22 21:54:07,296: INFO: topic #3: 0.009*量 + 0.007*心 + 0.006*争 + 0.005*报 + 0.005*报 + 0.004*银的 + 0.004*取 + 0.004*谢谢 + 0.00 .003*今天 2012-05-22 21:54:07,298: INFO: topic #4: 0.009*量 + 0.006*争 + 0.005*心 + 0.005*报 + 0.005*报的 + 0.005*籽	2012-05-22 21:54:07, 294 :	INFO:		0.008*爱	+ 0.007*@	+ 0.007*心	+ 0.007*泪 +	0.005*幸福	+ 0.005*我的 -	- 0.004*谢	谢 + 0.004*好 + 0.
0.003*今天 2012-05-22 21:54:07,298: INFO: topic #4: 0.005*酸 + 0.005*砂 + 0.005*池 + 0.005*預 + 0.005*預 + 0.005*野 + 0.005*学 + 0.004*今天 + 0.004*谢谢 + 0.003*你的 - 0.003*你的 - 0.003*你的 - 0.003*快乐 - 0.003*快乐 - 0.003*快乐 - 0.003*快乐 - 0.003*快乐 - 0.003*快乐 - 0.003*セス・07*心 + 0.006*製 + 0.005*園 + 0.005*園 + 0.005*園 + 0.004*場 + 0.004*場 + 0.004*財	0. 003*快乐										
2012-05-22 21:54:07,298: INFO: topic #4: 0.009*数 + 0.006*% + 0.005*心 + 0.005*我 + 0.005*我的 + 0.005*好 + 0.005*好 + 0.004*今天 + 0.004*谢谢 + 0.00.005*投 + 0.005*我的 + 0.005*我的 + 0.004*谢谢 + 0.004*今天 + 0.004*今天 + 0.005*投 + 0.005*投 + 0.005*投 + 0.005*设 + 0.005*投 + 0.005*投 + 0.005*投 + 0.005*设 + 0.0	2012-05-22 21:54:07, 296 :	INFO:	topic #3:	0.009*爱		+ 0.006*®	+ 0.005*好 +	0.005*泪 +	0.004*我的 + (0.004*耶 +	0.004*谢谢 + 0.00
+ 0.003*他的 2012-05-22 21:54:07,308: INFO: topic #5: 0,009*酸 + 0.007*心 + 0.005*油 + 0.005*独的 + 0.004*® + 0.004*園 + 0.004*財 + 0.004*労 + 0.004*今天 + 0. + 0.003*快乐 2012-05-22 21:54:07,309: INFO: topic #6: 0.007*心 + 0.006*製 + 0.005*像 + 0.005*園 + 0.005*港 + 0.005*港 + 0.004*財 + 0.004	0.003*今天										
2012-05-22 21:54:07,308: INFO: topic #5: 0.009*夏 + 0.007*心 + 0.005*周 + 0.005*預 + 0.004*⑨ + 0.004*⑨ + 0.004*谢谢 + 0.004*炒千 + 0.004*少天 + 0. + 0.003*版5. 2012-05-22 21:54:07,309: INFO: topic #6: 0.007*心 + 0.006*夏 + 0.005*⑩ + 0.005*洵 + 0.005*瑜福 + 0.004*琐的 + 0.004*琐的 + 0.004*琐的 + 0.004*琐的 + 0.004*琐的 + 0.003*谢谢 + 0. + 0.003*顺5-22 21:54:07,311: INFO: topic #7: 0.009*夏 + 0.008*心 + 0.005*⑩ + 0.005*洵 + 0.004*祥福 + 0.004*谢谢 + 0.004*琐谢 + 0.004*贵的 + 0.003*版5	2012-05-22 21:54:07, 298 :	INFO:		0.009*爱	+ 0.006*@		+ 0.005*泪 +	0.005*我的	+ 0.005*好 + 0	0.004*今天	+ 0.004*谢谢 + 0.
+ 0.003*検系 2012-05-22 1:54:07,309: INFO: topic #6: 0,007*心 + 0.006*載 + 0.005*® + 0.005*園 + 0.005*業福 + 0.004*教的 + 0.004*財 + 0.004*財 + 0.003*謝 樹 + 0. + 0.003*耶 2012-05-22 21:54:07,311: INFO: topic #7: 0.009*慶 + 0.008*心 + 0.005*® + 0.005*園 + 0.004*報福 + 0.004*財 樹 + 0.004*助 樹 + 0.004*助 か + 0.004*財	+ 0.003*你的										
2012-05-22 21:54:07,309: INFO: topic #6: 0.007+心 + 0.006+酸 + 0.005+® + 0.005+馏 + 0.005+增福 + 0.004+教的 + 0.004+教的 + 0.003+馏留 + 0.003+馏留 + 0.003+馏留 + 0.003+馏留 + 0.003+馏留 + 0.003+馏留 + 0.004+馏留 + 0.004+馏留 + 0.004+馏留 + 0.004+数的 + 0.003+收系 2012-05-22 21:54:07,311: INFO: topic #7: 0.003+腹 + 0.005+⑩ + 0.005+⑩ + 0.005+⑩ + 0.005+馏 + 0.005+૫ + 0.005+10.	2012-05-22 21:54:07,308 :	INFO:	topic #5:	0.009*爱		+ 0.005*泪	+ 0.005*我的	+ 0.004*@	+ 0.004*谢谢 -	- 0.004*好	+ 0.004*今天 + 0.
+ 0.003*取 2012-05-22 21:54:07,311 : INFO : topic #7: 0.009*数 + 0.008*心 + 0.005*参 + 0.005*语 + 0.004*堆福 + 0.004*坩樹 + 0.004*坩樹 + 0.004*サ + 0.004*サ 0 + 0.003*サ 5 - 0.003*サ 5 - 0.003*サ 5 - 0.003*サ 5 - 0.004*サ 7 - 0											
2012-05-22 21:54:07,311 : INFO: topic #7: 0.009*慶 + 0.008*心 + 0.005*® + 0.005*園 + 0.004*準福 + 0.004*謝 + 0.004*謝 + 0.004*散 + 0.003*快乐 2012-05-22 21:54:07,312 : INFO: topic #8: 0.008*夏 + 0.007*心 + 0.006*好 + 0.005*® + 0.005*1 + 0.005*谢 + 0.005*谢谢 + 0.004*快乐 + 0.004*我的 + 0.004**	2012-05-22 21:54:07,309 :	INFO:			+ 0.006*爱	+ 0.005*@	+ 0.005*泪 +	0.005*幸福	+ 0.004*我的 -	- 0.004*好	+ 0.003*谢谢 + 0.
0.003*快乐 2012-05-22 21:54:07,312: INFO: topic #8: 0.008*最 + 0.007*心 + 0.006*好 + 0.005*像 + 0.005*准 + 0.005*谢谢 + 0.004*快乐 + 0.004*投阶 + 0. 0.004**											
2012-05-22 21:54:07,312: INFO: topic #8: 0.008*夏 + 0.007*心 + 0.006*好 + 0.005*億 + 0.005*泪 + 0.005*谢谢 + 0.004*快乐 + 0.004*我的 + 0.004*	2012-05-22 21:54:07,311 :	INFO:	topic #7:	0.009*爱	+ 0.008*心	+ 0.005*@	+ 0.005*泪 +	0.004*幸福	+ 0.004*谢谢 -	- 0.004*耶	+ 0.004*我的 + 0.
0.004*"											
		INFO:		0.008*爱		+ 0.006*好	+ 0.005*@+	0.005*泪 +	0.005*谢谢 + (). 004*快乐	+ 0.004*我的 + 0.
2012-05-22 21:54:07,314: INFO: topic #9: 0.008*心 + 0.006*® + 0.005*泪 + 0.005*漫 + 0.005*晕禍 + 0.004*今天 + 0.004*耶 + 0.004*戼 + 0.00											
	2012-05-22 21:54:07,314 :	INFO:		0.008*心	+ 0.006*@	+ 0.005*泪	+ 0.005*愛 +	0.005*幸福	+ 0.004*今天	- 0.004*耶	+ 0.004*好 + 0.00

实验结果(续)

• 36氪的微博主题

```
2012-05-22 22:28:55,374: 1NFO 1 topic = 80: 0.010*公司 + 0.009*周用 + 0.009*周元 + 0.009*周户 + 0.006*創立 + 0.005*報本 + 0.004*年東 + 0.005*年東 + 0.004*年東 + 0.005*年東 + 0.005*日 + 0.005*
```

主题归纳



■ 創新,中国,创业,投资,用户,硅谷,苹果



爱,心,幸福,谢谢,快乐,耶



公司,应用,美元,用户,创业,发布,移动,苹果,手机

Topic model in image processing

- 实验数据: youtube22数据集,包含22个概念的数据,实验只使用其中的一个概念,篮球。一共使用了50个与篮球相关的视频,提取的关键帧包括2632张图片。
- 特征提取:我了提取SIFT特征,SIFT算法是一种提取局部特征的算法,在尺度空间寻找极值点,提取位置,尺度,旋转不变量。
- 关键帧的向量表示.使用BOF模型。
- 实验参数及环境设置:SIFT的采样数目设为1万个,聚类为500类, 每张图片用一个500维的特征向量表示。然后通过LDA进行主题分析,主题个数我们设置为10,

实验结果



实验结果(续)



• Topic 4:

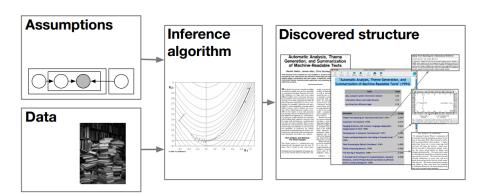
- Introduction to topic model
 - What are topic model?
 - LDA
- 2 Applications
 - Topic model in text mining
 - Topic model in image processing
- Model extensions
- 4 Conclusion
- 6 Reference
- 6 The end

Beyond latent Dirichlet allocation

- Modeling richer assumptions
 - Dynamic topic models
 - Correlated topic models
- Supervised topic models
 - Supervised LDA
 - Relational topic models
- Bayesian nonparametric topic models

- Introduction to topic model
 - What are topic model?
 - LDA
- 2 Applications
 - Topic model in text mining
 - Topic model in image processing
- Model extensions
- 4 Conclusion
- 6 Reference
- 6 The end

Summary



- Introduction to topic model
 - What are topic model?
 - LDA
- 2 Applications
 - Topic model in text mining
 - Topic model in image processing
- Model extensions
- 4 Conclusion
- Seference
- 6 The end

Reference

- David M. Blei, Andrew Ng, Michael Jordan. Latent Dirichlet allocation. JMLR (3) 2003 pp. 993-1022.
- David M. Blei. Introduction to Probabilistic Topic Models.
 Communications of the ACM.2011
- David M. Blei, Jon D. McAuliffe. Supervised Topic Models. NIPS (2007)
- David M. Blei, John D. Lafferty. Dynamic Topic Models. ICML (2006)
- Jonathan Chang, David Blei. Relational Topic Models for Document Networks. AlStats (2009)

- Introduction to topic model
 - What are topic model?
 - LDA
- 2 Applications
 - Topic model in text mining
 - Topic model in image processing
- Model extensions
- 4 Conclusion
- Reference
- 6 The end

Thanks!