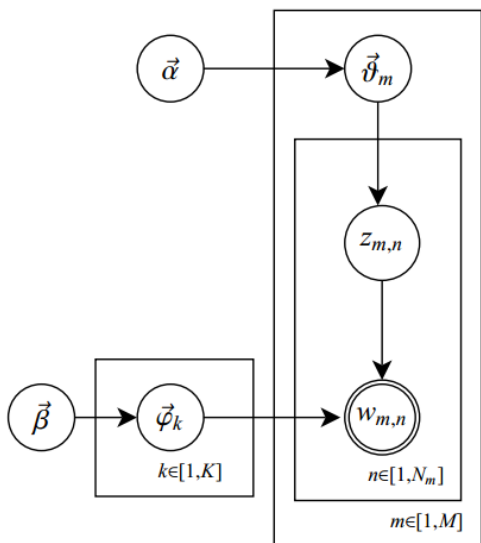


主题模型



22.4.petrol 22.4.petrol

每个主题的词分布:

主题#0: 溶脱 地面 下 发现 溜 吹 漏洞

概率: [0.03567895 0.0305515 0.02995125 0.02905559 0.02847266 0.02672661 0.02639478]

主题#1: 卫生 清扫 新增 本月 缺陷 无 空调

概率: [0.09127588 0.04710893 0.03864091 0.0385492 0.03507678 0.03243568 0.0211739]

主题#2: 过滤器 设备 高 差压 少 入口 17

概率: [0.07382152 0.04735673 0.03770307 0.03342033 0.03160451 0.02609018 0.02332794]

主题#3: 号牌 位 脱落 铁皮 北侧 塔 规范

概率: [0.06217475 0.06146815 0.04165677 0.03554779 0.03128229 0.02976478 0.02486769]

主题#4: 松 建议 液位 损坏 皮带 区域 断裂

概率: [0.03845036 0.03260667 0.03243315 0.031336 0.03074858 0.03062344 0.02835042]

主题#5: 盘根 漏 阀 出口 采样 引出 内

概率: [0.09527604 0.08630304 0.07457675 0.0508139 0.04410229 0.03406847 0.03309818]

主题#6: 日 月 8 红线 压力表 技术员 23

概率: [0.04320296 0.04268257 0.04081915 0.03381051 0.02013684 0.01981086 0.01615044]

主题#7: 地沟 错误 单向阀 螺丝 装车 旁 年

概率: [0.03169998 0.02547323 0.02350422 0.02196816 0.02146054 0.02122409 0.01970546]

主题#8: 皮带 断 不准 松动 一次 盖 被

概率: [0.15231247 0.10833394 0.05339595 0.04884857 0.04585281 0.03503216 0.03475644]

主题#9: 电机 杂音 接头 大盖 冲洗 有 活

概率: [0.07620952 0.0712979 0.06082735 0.03423004 0.03310958 0.02874415 0.02728715]

主题#10: 蒸汽 砂眼 管线 法兰 漏 前 有

概率: [0.04160303 0.03960238 0.03123862 0.0293527 0.02685108 0.02634925 0.0248503]

主题#11: 泄漏 伴热 量 牌 入口 底 法兰

概率: [0.07545736 0.05467109 0.05165556 0.03689506 0.03672955 0.03457064 0.03128606]

主题#12: 密封 平台 保温 脱开 缺失 堵头 汽提

概率: [0.04631792 0.04157294 0.0396999 0.03857663 0.03751675 0.03526295 0.03418964]

主题#13: 后端 长明灯 无法 号 灭火器 器 油站

概率: [0.04065264 0.03001566 0.02572 0.02520191 0.02510101 0.02370398 0.02267814]

主题#14: 坏 炉 压力表 润滑油 泵 安全阀 清理

概率: [0.04789169 0.03913219 0.03672254 0.03227133 0.03054908 0.02618981 0.02523409]

主题#15: 缺陷 新增 无 本月 交接班 胶带机 日志

概率: [0.07364203 0.07360244 0.06751279 0.06417833 0.03489421 0.02254108 0.02149432]

主要内容

□ LDA

- 隐Dirichlet分布
- Latent Dirichlet Allocation

□ 先验分布 – 共轭分布

□ Beta分布 – Dirichlet分布

□ 三层贝叶斯网络模型LDA

□ Gibbs采样和更新规则

LDA的应用方向

☐ 信息提取和搜索

☒ 语义分析

☐ 文档分类/聚类、文章摘要、社区挖掘

☐ 基于内容的图像聚类、目标识别

☒ 以及其他计算机视觉应用

☐ 生物信息数据的应用

朴素贝叶斯的分析

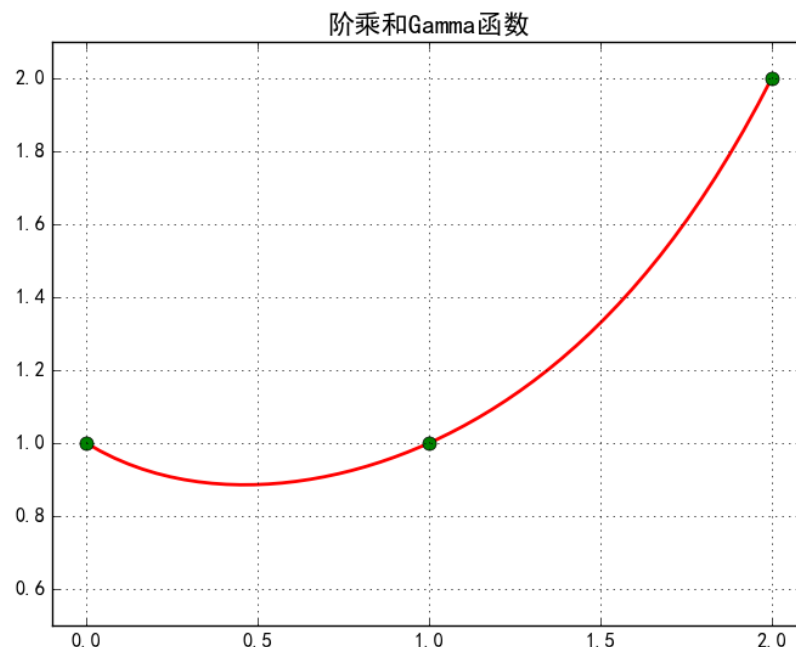
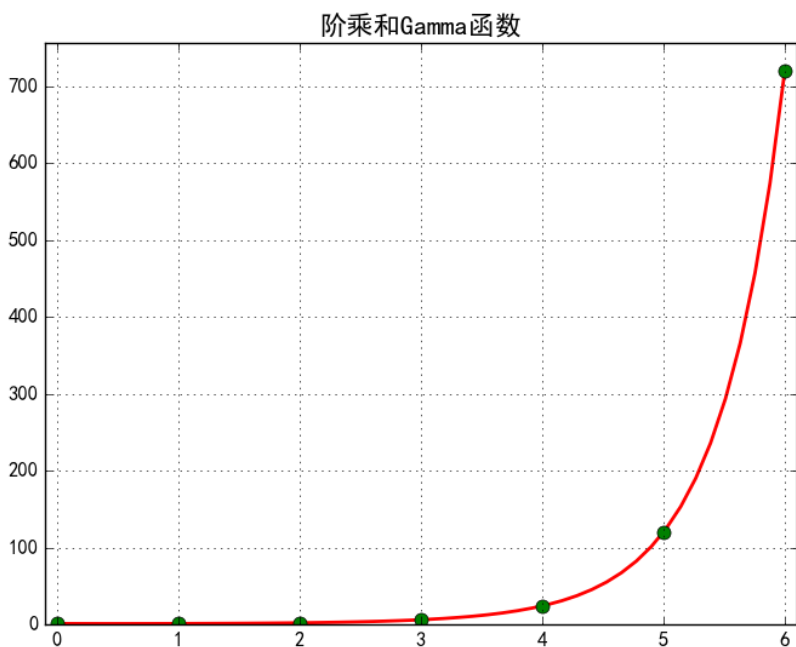
- 可以胜任许多文本分类问题。
- 无法解决语料中**一词多义**和**多词一义**的问题——它更像是词法分析，而非语义分析。
- 如果使用词向量作为文档的特征，**一词多义**和**多词一义**会造成计算文档间相似度的不准确性。
- 可以通过增加“主题”的方式，一定程度的解决上述问题：
 - 一个词可能被映射到多个主题中
 - ——**一词多义**
 - 多个词可能被映射到某个主题的概率很高
 - ——**多词一义**

引：Γ函数

$$\Gamma(x) = (x-1) \cdot \Gamma(x-1) \Rightarrow \frac{\Gamma(x)}{\Gamma(x-1)} = x-1$$

□ Γ函数是阶乘在实数上的推广

$$\Gamma(x) = \int_0^{+\infty} t^{x-1} e^{-t} dt = (x-1)!$$



Beta分布

□ Beta分布的概率密度：
$$f(x) = \begin{cases} \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, & x \in [0,1] \\ 0, & \text{其他} \end{cases}$$

□ 其中系数B为：

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$$

□ Gamma函数看成阶乘的实数域推广：

$$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt$$

$$\Rightarrow \Gamma(n) = (n-1)! \Rightarrow B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$$

Beta分布的期望

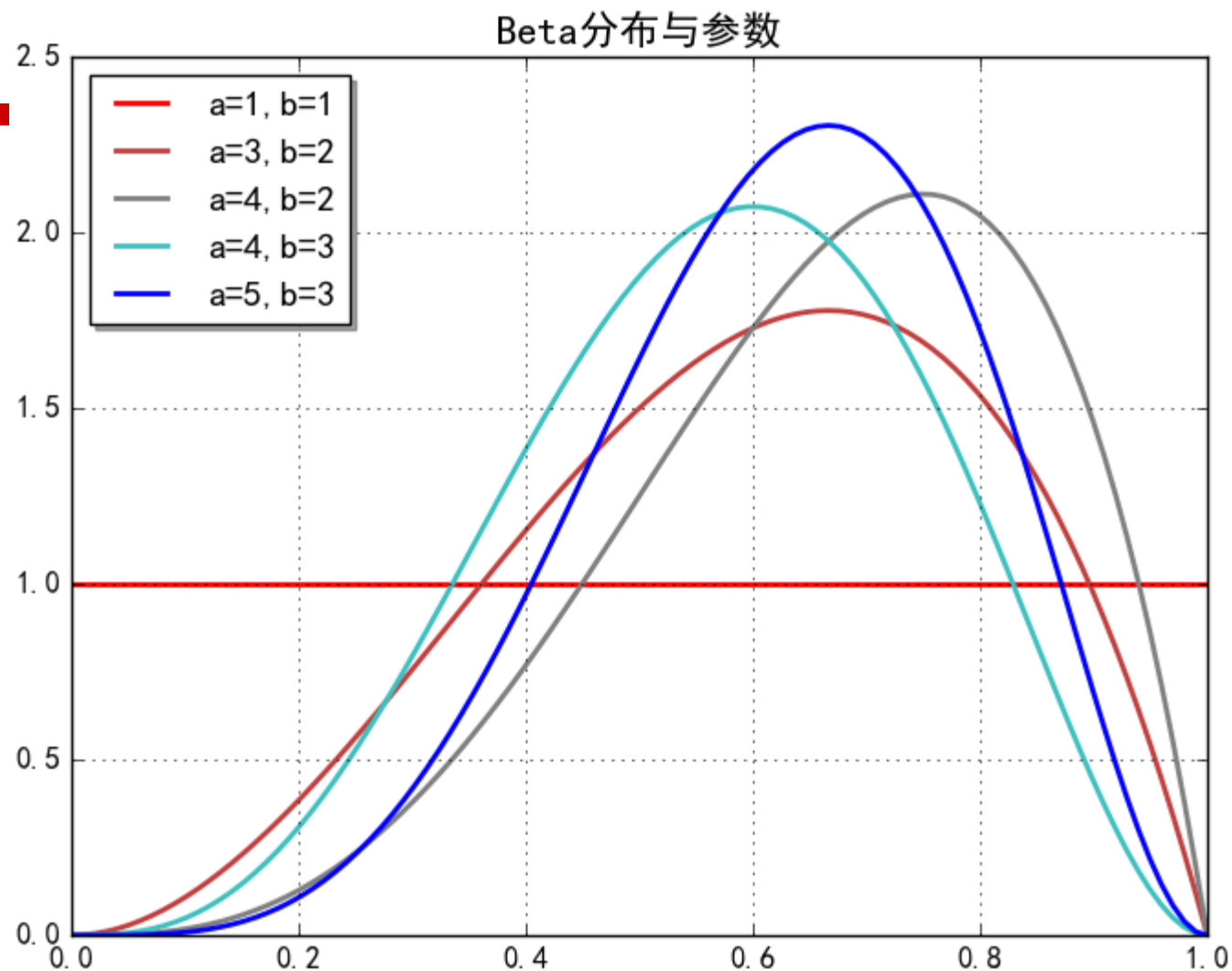
$$f(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, x \in [0,1]$$

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$$

□ 根据定义：

$$\begin{aligned} E(X) &= \int_0^1 x \cdot \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} dx \\ &= \frac{1}{B(\alpha, \beta)} \int_0^1 x^{(\alpha+1)-1} (1-x)^{\beta-1} dx \\ &= \frac{B(\alpha+1, \beta)}{B(\alpha, \beta)} = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \bigg/ \frac{\Gamma(\alpha+\beta+1)}{\Gamma(\alpha+1)\Gamma(\beta)} \\ &= \frac{\alpha}{\alpha+\beta} \end{aligned}$$

Beta分布

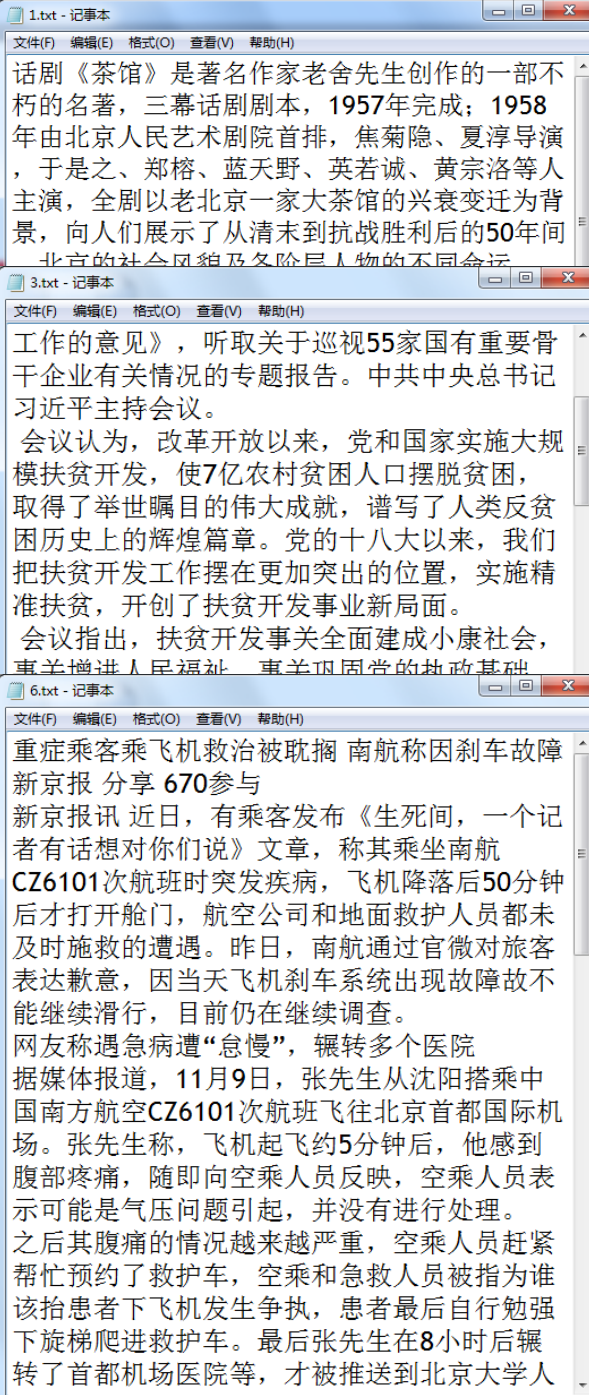


文档和主题

文档 1 : 茶馆 (0.0163591635916) 社会 (0.00528905289053) 王利发 (0.00528905289053)
文档 2 : 决议 (0.0138983050847) 打击 (0.00824858757062) 安理会 (0.00824858757062)
文档 3 : 会议 (0.0124491456469) 脱贫 (0.0124491456469) 党校 (0.0108218063466)
文档 4 : 美团 (0.0306066176471) 阿里 (0.0103860294118) 业务 (0.0103860294118)
文档 5 : 户口 (0.0221347331584) 登记 (0.0195100612423) 人口 (0.0142607174103)
文档 6 : 人员 (0.0111471861472) 飞机 (0.00898268398268) 称 (0.00681818181818)
文档 7 : 号线 (0.0328544061303) 站 (0.0194444444444) 14 (0.0184865900383)
文档 8 : 支付 (0.0198394495413) 腾讯 (0.0072247706422) 支付宝 (0.0072247706422)
文档 9 : 决议 (0.0138983050847) 打击 (0.00824858757062) 安理会 (0.00824858757062)
文档 10 : 足协 (0.0186473429952) 足球 (0.0138164251208) 佩兰 (0.011884057971)

=====

主题 1 : 美团 (0.0306066176471) 阿里 (0.0103860294118) 业务 (0.0103860294118)
主题 2 : 会议 (0.0124491456469) 脱贫 (0.0124491456469) 党校 (0.0108218063466)
主题 3 : 号线 (0.0328544061303) 站 (0.0194444444444) 14 (0.0184865900383)
主题 4 : 人物 (0.00214876033058) 民族 (0.00214876033058) 资本家 (0.00214876033058)
主题 5 : 足协 (0.0186473429952) 足球 (0.0138164251208) 佩兰 (0.011884057971)
主题 6 : 户口 (0.0221347331584) 登记 (0.0195100612423) 人口 (0.0142607174103)
主题 7 : 决议 (0.0138983050847) 打击 (0.00824858757062) 安理会 (0.00824858757062)
主题 8 : 人员 (0.0111471861472) 飞机 (0.00898268398268) 称 (0.00681818181818)
主题 9 : 茶馆 (0.0163591635916) 社会 (0.00528905289053) 王利发 (0.00528905289053)
主题 10 : 支付 (0.0198394495413) 腾讯 (0.0072247706422) 支付宝 (0.0072247706422)



LDA涉及的主要问题

- 共轭先验分布
- Dirichlet分布
- LDA模型
 - Gibbs采样算法学习参数

共轭先验分布

- 由于 x 为给定样本， $P(x)$ 有时被称为“证据”，仅仅是归一化因子，如果不关心 $P(\theta|x)$ 的具体值，只考察 θ 取何值时后验概率 $P(\theta|x)$ 最大，则可将分母省去。

$$P(\theta | x) = \frac{P(x | \theta)P(\theta)}{P(x)} \propto P(x | \theta)P(\theta)$$

- 在贝叶斯概率理论中，如果后验概率 $P(\theta|x)$ 和先验概率 $p(\theta)$ 满足同样的分布律，那么，先验分布和后验分布被叫做共轭分布，同时，先验分布叫做似然函数的共轭先验分布。
- In Bayesian probability theory, if the posterior distributions $p(\theta|x)$ are in the same family as the prior probability distribution $p(\theta)$, the prior and posterior are then called conjugate distributions, and the prior is called a conjugate prior for the likelihood function.

复习：二项分布的最大似然估计

- 投硬币试验中，进行N次独立试验，n次朝上，N-n次朝下。
- 假定朝上的概率为p，使用对数似然函数作为目标函数：

$$f(n | p) = \log(p^n (1-p)^{N-n}) \xrightarrow{\Delta} h(p)$$

$$\frac{\partial h(p)}{\partial p} = \frac{n}{p} - \frac{N-n}{1-p} \xrightarrow{\Delta} 0 \Rightarrow p = \frac{n}{N}$$

二项分布与先验举例

□ 在校门口统计一定时间段内出入的男女生数目分别为 N_B 和 N_G ，估算该校男女生比例。

$$\begin{cases} P_B = \frac{N_B}{N_B + N_G} \\ P_G = \frac{N_G}{N_B + N_G} \end{cases}$$

□ 若观察到4个女生和1个男生，可以得出该校女生比例是80%吗？

□ 修正公式：

$$\begin{cases} P_B = \frac{N_B + 5}{N_B + N_G + 10} \\ P_G = \frac{N_G + 5}{N_B + N_G + 10} \end{cases} \Rightarrow \begin{cases} P_B = \frac{1 + 5}{1 + 4 + 10} = 40\% \\ P_G = \frac{4 + 5}{1 + 4 + 10} = 60\% \end{cases}$$

上述过程的理论解释

- 投掷一个非均匀硬币，可以使用参数为 θ 的伯努利模型， θ 为硬币为正面的概率，那么结果 x 的分布形式为： $P(x|\theta) = C_n^k \cdot \theta^k \cdot (1-\theta)^{n-k}$
- 两点分布/二项分布的共轭先验是Beta分布，它具有两个参数 α 和 β ，Beta分布形式为

$$P(\theta | \alpha, \beta) = \begin{cases} \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}, & \theta \in [0,1] \\ 0, & \text{其他} \end{cases}$$

先验概率和后验概率的关系

□ 根据似然和先验：

$$P(x|\theta) = C_n^k \cdot \theta^k \cdot (1-\theta)^{n-k}$$

$$P(\theta|\alpha, \beta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

计算后验概率：

$$P(\theta|x) = \frac{P(x|\theta) \cdot P(\theta)}{P(x)} \propto P(x|\theta) \cdot P(\theta)$$

$$= \left(C_n^k \theta^k (1-\theta)^{n-k} \right) \cdot \left(\frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \right)$$

$$= \frac{C_n^k}{B(\alpha, \beta)} \theta^{(k+\alpha)-1} (1-\theta)^{(n-k+\beta)-1}$$

$$\propto \frac{1}{B(k+\alpha, n-k+\beta)} \theta^{(k+\alpha)-1} (1-\theta)^{(n-k+\beta)-1}$$

□ 后验概率是参数为 $(k+\alpha, n-k+\beta)$ 的Beta分布，即：伯努利分布/二项分布的共轭先验是Beta分布。

$$P(\theta) = \frac{1}{B(\alpha, \beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

伪计数

$$P(\theta|x) = \frac{1}{B(k+\alpha, n-k+\beta)} \theta^{(k+\alpha)-1} (1-\theta)^{(n-k+\beta)-1}$$

- 参数 α 、 β 是决定参数 θ 的参数，即超参数。
- 在后验概率的最终表达式中，参数 α 、 β 和 k 、 $n-k$ 一起作为参数 θ 的指数——后验概率的参数为 $(k+\alpha, n-k+\beta)$ 。
- 根据这个指数的实践意义：投币过程中，正面朝上的次数， α 和 β 先验性的给出了在没有任何实验的前提下，硬币朝上的概率分配；因此， α 和 β 可被称作“伪计数”。

共轭先验的直接推广

□ 从2到K:

- 二项分布 \rightarrow 多项分布
- Beta分布 \rightarrow Dirichlet分布

Dirichlet分布

□ Beta 分布: $f(x) = \begin{cases} \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, & x \in [0,1] \\ 0, & \text{其他} \end{cases}$

■ 其中: $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$

□ Dirichlet 分布: $f(\vec{p} | \vec{\alpha}) = \begin{cases} \frac{1}{\Delta(\vec{\alpha})} \prod_{k=1}^K p_k^{\alpha_k-1}, & p_k \in [0,1] \\ 0, & \text{其他} \end{cases}$

■ 简记: $Dir(\vec{p} | \vec{\alpha}) = \frac{1}{\Delta(\vec{\alpha})} \prod_{k=1}^K p_k^{\alpha_k-1}$ 其中: $\Delta(\vec{\alpha}) = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}$

Dirichlet分布的期望

□ 根据Beta分布的期望公式：

$$f(x) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}, x \in [0,1] \Rightarrow E(X) = \frac{\alpha}{\alpha + \beta}$$

□ 推广得到：

$$f(\vec{p} | \vec{\alpha}) = \frac{1}{\Delta(\vec{\alpha})} \prod_{k=1}^K p_k^{\alpha_k-1}, p \in [0,1] \Rightarrow E(p_i) = \frac{\alpha_i}{\sum_{k=1}^K \alpha_k}$$

Dirichlet分布分析

$$Dir(\vec{p} | \vec{\alpha}) = \frac{1}{\Delta(\vec{\alpha})} \prod_{k=1}^K p_k^{\alpha_k - 1}$$

□ α 是参数向量，共 K 个

□ 定义在 x_1, x_2, \dots, x_{K-1} 维上

■ $x_1 + x_2 + \dots + x_{K-1} + x_K = 1$

■ $x_1, x_2, \dots, x_{K-1} > 0$

■ 定义在 $(K-1)$ 维的 **单纯形** 上，其他区域的概率密度为 0

□ α 的取值对 $Dir(p|\alpha)$ 有什么影响？

Symmetric Dirichlet distribution

- A very common special case is the **symmetric Dirichlet distribution**, where all of the elements making up the parameter **vector** have the same value. Symmetric Dirichlet distributions are often used when a Dirichlet **prior** is called for, since there typically is no prior knowledge favoring one component over another. Since all elements of the parameter vector have the same value, the distribution alternatively can be parametrized by a single **scalar value** α , called the **concentration parameter**(**聚集参数**).

对称Dirichlet分布

□ Dirichlet 分布： $Dir(\vec{p} | \vec{\alpha}) = \frac{1}{\Delta(\vec{\alpha})} \prod_{k=1}^K p_k^{\alpha_k - 1}$

■ 其中： $\Delta(\vec{\alpha}) = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}$

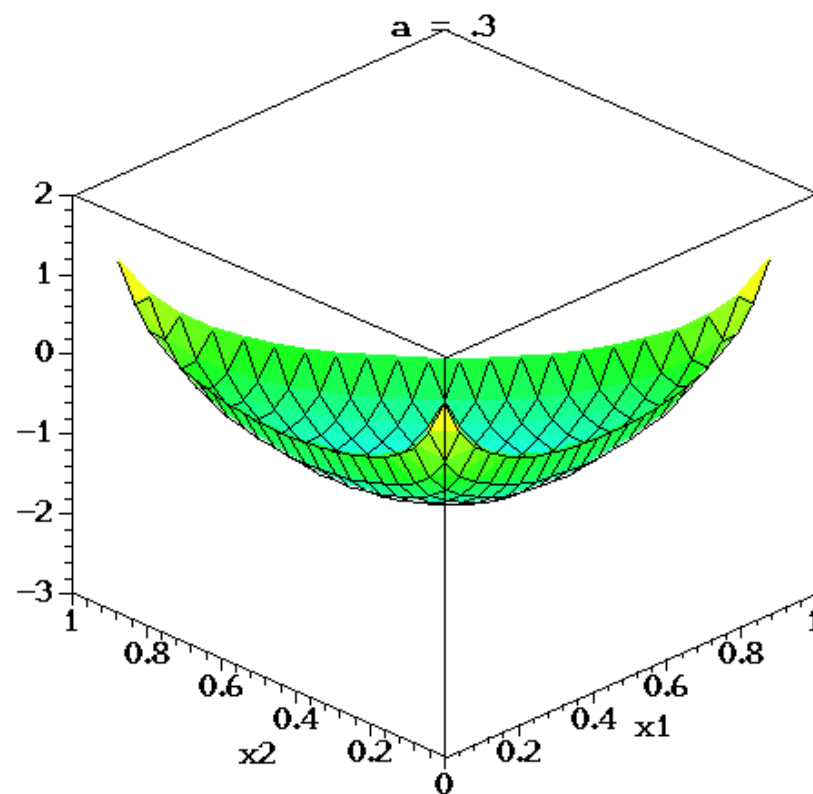
□ 对称Dirichlet分布： $Dir(\vec{p} | \alpha, K) = \frac{1}{\Delta_K(\alpha)} \prod_{k=1}^K p_k^{\alpha - 1}$

■ 其中： $\Delta_K(\alpha) = \frac{\Gamma^K(\alpha)}{\Gamma(K \cdot \alpha)}$

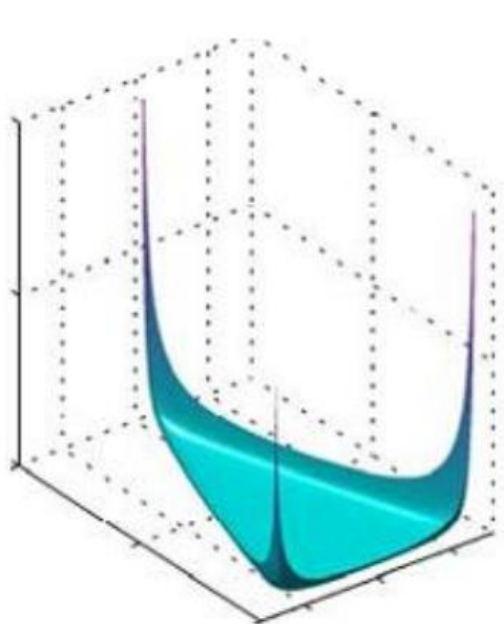
对称Dirichlet分布的参数分析

- $\alpha=1$ 时
 - 退化为均匀分布
- 当 $\alpha>1$ 时
 - $p_1=p_2=\dots=p_k$ 的概率增大
- 当 $\alpha<1$ 时
 - $p_i=1$, $p_{\text{非}i}=0$ 的概率增大

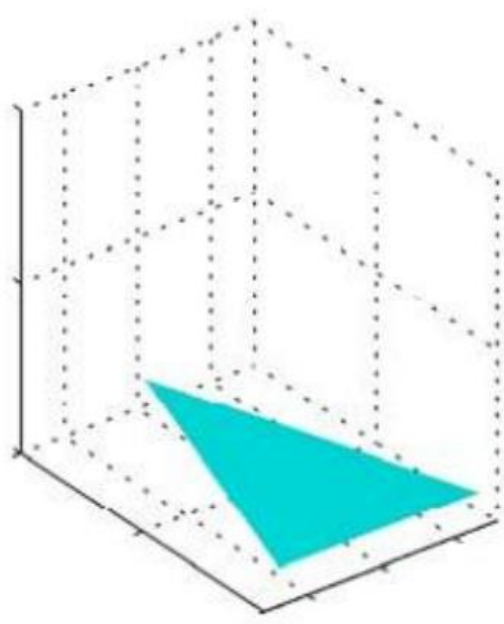
图像说明：将Dirichlet分布的概率密度函数取对数,绘制对称Dirichlet分布的图像,取 $K=3$,也就是有两个独立参数 x_1, x_2 , 分别对应图中的两个坐标轴, 第三个参数始终满足 $x_3=1-x_1-x_2$ 且 $\alpha_1=\alpha_2=\alpha_3=\alpha$, 图中反映的是 α 从0.3变化到2.0的概率对数值的变化情况。



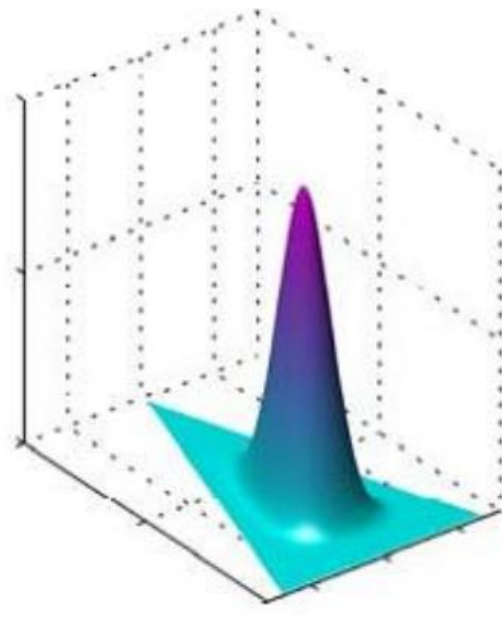
参数 α 对Dirichlet分布的影响



$$\{\alpha_k\} = 0.1$$



$$\{\alpha_k\} = 1$$



$$\{\alpha_k\} = 10$$

参数选择对对称Dirichlet分布的影响

- When $\alpha=1$, the symmetric Dirichlet distribution is equivalent to a uniform distribution over the open standard $(K-1)$ -simplex, i.e. it is uniform over all points in its support. Values of the concentration parameter above 1 prefer variants that are dense, evenly distributed distributions, i.e. all the values within a single sample are similar to each other. Values of the concentration parameter below 1 prefer sparse distributions, i.e. most of the values within a single sample will be close to 0, and the vast majority of the mass will be concentrated in a few of the values.

多项分布的共轭分布是Dirichlet分布

$\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$ = concentration hyperparameter

$\mathbf{p} \mid \boldsymbol{\alpha} = (p_1, \dots, p_K) \sim \text{Dir}(K, \boldsymbol{\alpha})$

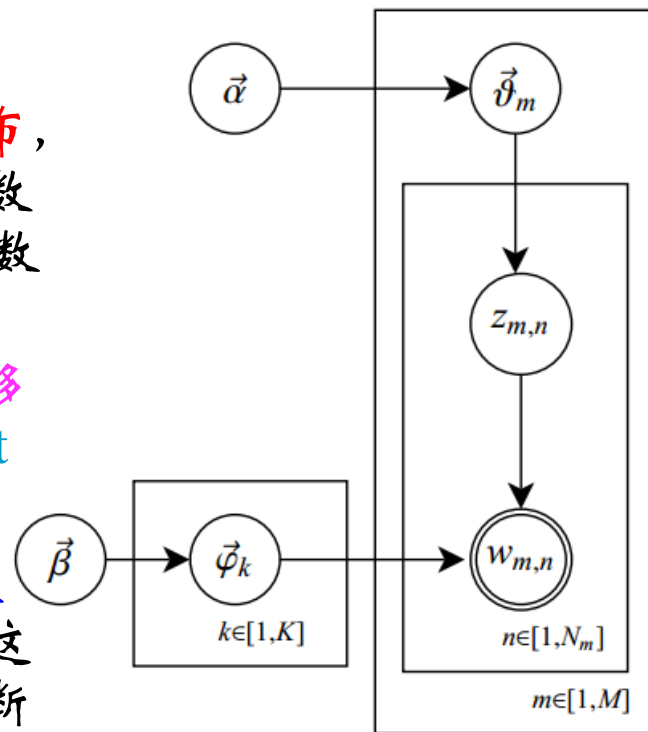
$\mathbb{X} \mid \mathbf{p} = (\mathbf{x}_1, \dots, \mathbf{x}_K) \sim \text{Cat}(K, \mathbf{p})$

$\mathbf{c} = (c_1, \dots, c_K)$ = number of occurrences of category i

$\mathbf{p} \mid \mathbb{X}, \boldsymbol{\alpha} \sim \text{Dir}(K, \mathbf{c} + \boldsymbol{\alpha}) = \text{Dir}(K, c_1 + \alpha_1, \dots, c_K + \alpha_K)$

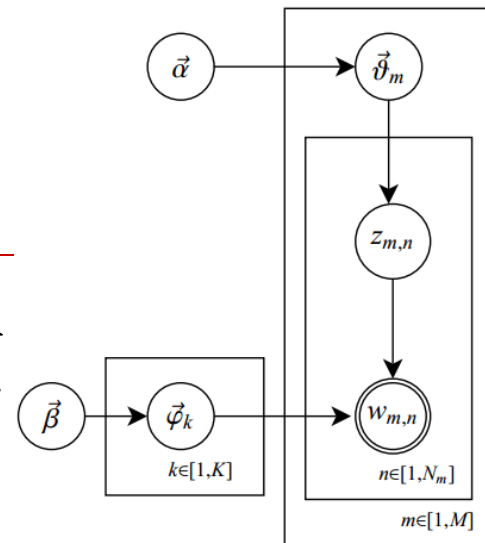
LDA的解释

- 共有 m 篇文章，一共涉及了 K 个主题；
- 每篇文章(长度为 N_m)都有各自的**主题分布**，**主题分布**是**多项分布**，该**多项分布**的参数服从**Dirichlet分布**，该**Dirichlet分布**的参数为 α ；
- 每个**主题**都有各自的**词分布**，**词分布**为**多项分布**，该**多项分布**的参数服从**Dirichlet分布**，该**Dirichlet分布**的参数为 β ；
- 对于某篇文章中的第 n 个**词**，首先从该文章的**主题分布**中采样一个**主题**，然后在这个**主题**对应的**词分布**中采样一个**词**。不断重复这个随机生成过程，直到 m 篇文章全部完成上述过程。



详细解释

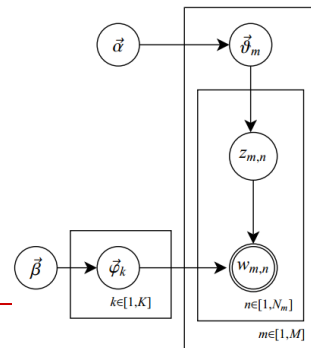
- 字典中共有 V 个term(不可重复), 这些term出现在具体的文章中, 就是word——在具体某文章中的word当然是有可能重复的。
- 语料库中共有 m 篇文档 $d_1, d_2 \dots d_m$;
- 对于文档 d_i , 由 N_i 个word组成, 可重复;
- 语料库中共有 K 个主题 $T_1, T_2 \dots T_k$;
- α 和 β 为先验分布的参数, 一般事先给定: 如取0.1的对称Dirichlet分布——表示在参数学习结束后, 期望每个文档的主题不会十分集中。
- θ 是每篇文档的**主题分布**
 - 对于第 i 篇文档 d_i 的主题分布是 $\theta_i = (\theta_{i1}, \theta_{i2}, \dots, \theta_{iK})$, 是长度为 K 的向量;
- 对于第 i 篇文档 d_i , 在主题分布 θ_i 下, 可以确定一个具体的主题 $z_{ij} = k$, $k \in [1, K]$,
- ϕ_k 表示第 k 个主题的**词分布**, $k \in [1, K]$
 - 对于第 k 个主题 T_k 的词分布 $\phi_k = (\phi_{k1}, \phi_{k2}, \dots, \phi_{kv})$, 是长度为 v 的向量
- 由 z_{ij} 选择 $\phi_{z_{ij}}$, 表示由词分布 $\phi_{z_{ij}}$ 确定term, 即得到观测值 w_{ij} 。



详细解释

- 图中 K 为主题个数， M 为文档总数， N_m 是第 m 个文档的单词总数。 β 是每个Topic下词的多项分布的Dirichlet先验参数， α 是每个文档下Topic的多项分布的Dirichlet先验参数。
 z_{mn} 是第 m 个文档中第 n 个词的主题， w_{mn} 是 m 个文档中的第 n 个词。两个隐含变量 θ 和 ϕ 分别表示第 m 个文档下的Topic分布和第 k 个Topic下词的分布，前者是 k 维(k 为Topic总数)向量，后者是 v 维向量(v 为词典中term总数)

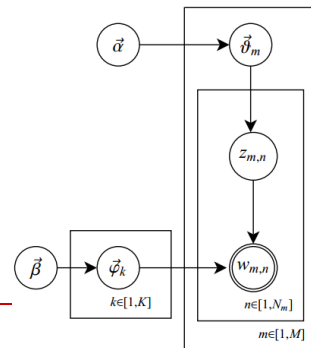
参数的学习



□ 给定一个文档集合， w_{mn} 是可以观察到的已知变量， α 和 β 是根据经验给定的先验参数，其他的变量 z_{mn} 、 θ 、 ϕ 都是未知的隐含变量，需要根据观察到的变量来学习估计的。根据LDA的图模型，可以写出所有变量的联合分布：

$$p(\vec{w}_m, \vec{z}_m, \vec{\theta}_m, \underline{\Phi} | \vec{\alpha}, \vec{\beta}) = \prod_{n=1}^{N_m} p(w_{m,n} | \vec{\phi}_{z_{m,n}}) p(z_{m,n} | \vec{\theta}_m) \cdot p(\vec{\theta}_m | \vec{\alpha}) \cdot p(\underline{\Phi} | \vec{\beta})$$

似然概率



□ 一个词 w_{mn} 初始化为一个词 t 的概率是

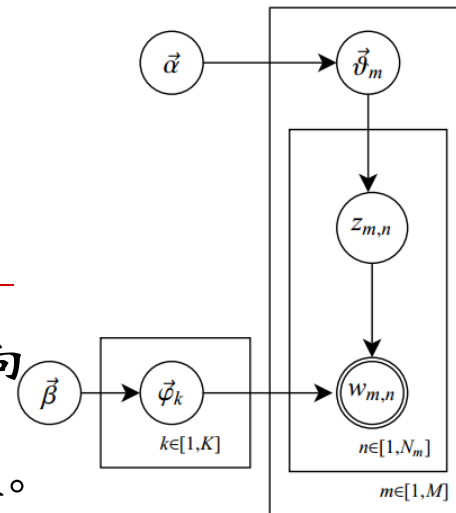
$$p(w_{m,n}=t|\vec{v}_m, \underline{\Phi}) = \sum_{k=1}^K p(w_{m,n}=t|\vec{\varphi}_k) p(z_{m,n}=k|\vec{v}_m)$$

□ 每个文档中出现主题 k 的概率乘以主题 k 下出现词 t 的概率，然后枚举所有主题求和得到。
整个文档集合的似然函数为：

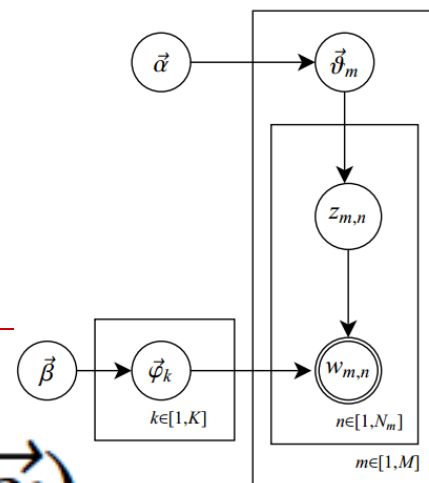
$$p(\mathcal{W}|\underline{\Theta}, \underline{\Phi}) = \prod_{m=1}^M p(\vec{w}_m|\vec{v}_m, \underline{\Phi}) = \prod_{m=1}^M \prod_{n=1}^{N_m} p(w_{m,n}|\vec{v}_m, \underline{\Phi})$$

Gibbs Sampling

- Gibbs Sampling算法的运行方式是每次选取概率向量的一个维度，给定其他维度的变量值采样当前维度的值。不断迭代直到收敛输出待估计的参数。
- 初始时随机给文本中的每个词分配主题 $z^{(0)}$ ，然后统计每个主题 z 下出现词 t 的数量以及每个文档 m 下出现主题 z 的数量，每一轮计算 $p(z_i|z_{-i}, \mathbf{d}, \mathbf{w})$ ，即排除当前词的主题分布：
 - 根据其他所有词的主题分布估计当前词分配各个主题的概率。
- 当得到当前词属于所有主题 z 的概率分布后，根据这个概率分布为该词采样一个新的主题。
- 用同样的方法更新下一个词的主题，直到发现每个文档的主题分布 θ_i 和每个主题的词分布 ϕ_j 收敛，算法停止，输出待估计的参数 θ 和 ϕ ，同时每个单词的主题 z_{mn} 也可同时得出。
- 实际应用中会设置最大迭代次数。每一次计算 $p(z_i|z_{-i}, \mathbf{d}, \mathbf{w})$ 的公式称为Gibbs updating rule。



联合分布

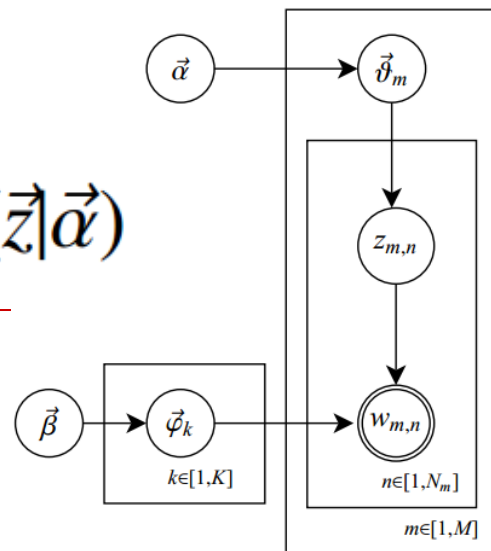


$$p(\vec{w}, \vec{z} | \vec{\alpha}, \vec{\beta}) = p(\vec{w} | \vec{z}, \vec{\beta}) p(\vec{z} | \vec{\alpha})$$

- 第一项因子是给定主题采样词的过程
- 后面的因子计算， $n_z^{(t)}$ 表示词 t 被观察到分配给主题 z 的次数， $n_m^{(k)}$ 表示主题 k 分配给文档 m 的次数。

计算因子 $p(\vec{w}, \vec{z} | \vec{\alpha}, \vec{\beta}) = p(\vec{w} | \vec{z}, \vec{\beta}) p(\vec{z} | \vec{\alpha})$

$$p(\vec{w} | \vec{z}, \vec{\beta}) = \int p(\vec{w} | \vec{z}, \underline{\Phi}) p(\underline{\Phi} | \vec{\beta}) d\underline{\Phi}$$



$$= \int \prod_{z=1}^K \frac{1}{\Delta(\vec{\beta})} \prod_{t=1}^V \varphi_{z,t}^{n_z^{(t)} + \beta_t - 1} d\vec{\varphi}_z$$

$$= \prod_{z=1}^K \frac{\Delta(\vec{n}_z + \vec{\beta})}{\Delta(\vec{\beta})}, \quad \vec{n}_z = \{n_z^{(t)}\}_{t=1}^V$$

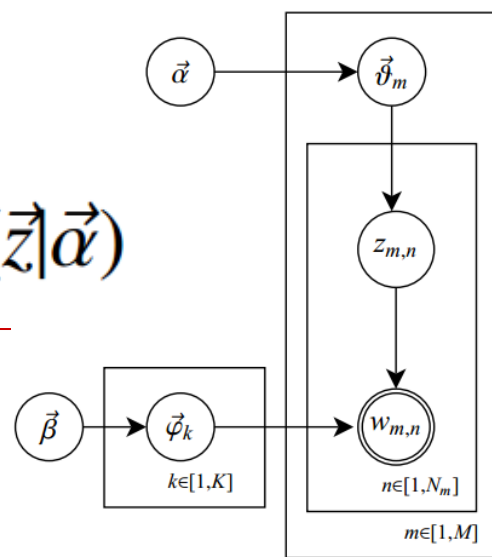
$$\int_{\vec{p}} \prod_{k=1}^K p_k^{\alpha_k - 1} d\vec{p} = \Delta(\vec{\alpha})$$

计算因子 $p(\vec{w}, \vec{z} | \vec{\alpha}, \vec{\beta}) = p(\vec{w} | \vec{z}, \vec{\beta}) p(\vec{z} | \vec{\alpha})$

$$p(\vec{z} | \vec{\alpha}) = \int p(\vec{z} | \underline{\Theta}) p(\underline{\Theta} | \vec{\alpha}) d\underline{\Theta}$$

$$= \int \prod_{m=1}^M \frac{1}{\Delta(\vec{\alpha})} \prod_{k=1}^K \vartheta_{m,k}^{n_m^{(k)} + \alpha_k - 1} d\vec{\vartheta}_m$$

$$= \prod_{m=1}^M \frac{\Delta(\vec{n}_m + \vec{\alpha})}{\Delta(\vec{\alpha})}, \quad \vec{n}_m = \{n_m^{(k)}\}_{k=1}^K$$



$$\int_{\vec{p}} \prod_{k=1}^K p_k^{\alpha_k - 1} d\vec{p} = \Delta(\vec{\alpha})$$

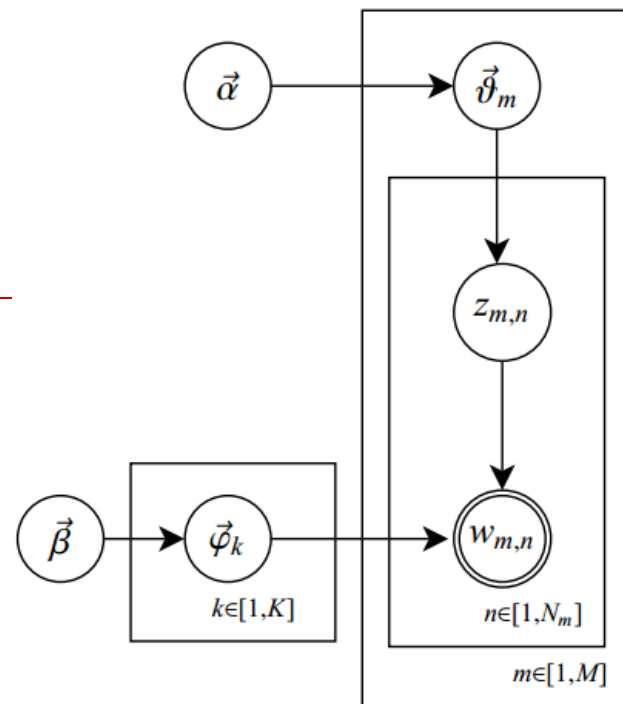
Gibbs updating rule

$$\begin{aligned} p(z_i=k|\vec{z}_{\neg i}, \vec{w}) &= \frac{p(\vec{w}, \vec{z})}{p(\vec{w}, \vec{z}_{\neg i})} = \frac{p(\vec{w}|\vec{z})}{p(\vec{w}_{\neg i}|\vec{z}_{\neg i})p(w_i)} \cdot \frac{p(\vec{z})}{p(\vec{z}_{\neg i})} \\ &\propto \frac{\Delta(\vec{n}_z + \vec{\beta})}{\Delta(\vec{n}_{z,\neg i} + \vec{\beta})} \cdot \frac{\Delta(\vec{n}_m + \vec{\alpha})}{\Delta(\vec{n}_{m,\neg i} + \vec{\alpha})} \\ &= \frac{\Gamma(n_k^{(t)} + \beta_t) \Gamma(\sum_{t=1}^V n_{k,\neg i}^{(t)} + \beta_t)}{\Gamma(n_{k,\neg i}^{(t)} + \beta_t) \Gamma(\sum_{t=1}^V n_k^{(t)} + \beta_t)} \cdot \frac{\Gamma(n_m^{(k)} + \alpha_k) \Gamma(\sum_{k=1}^K n_{m,\neg i}^{(k)} + \alpha_k)}{\Gamma(n_{m,\neg i}^{(k)} + \alpha_k) \Gamma(\sum_{k=1}^K n_m^{(k)} + \alpha_k)} \\ &= \frac{n_{k,\neg i}^{(t)} + \beta_t}{\sum_{t=1}^V n_{k,\neg i}^{(t)} + \beta_t} \cdot \frac{n_{m,\neg i}^{(k)} + \alpha_k}{[\sum_{k=1}^K n_m^{(k)} + \alpha_k] - 1} \\ &\propto \frac{n_{k,\neg i}^{(t)} + \beta_t}{\sum_{t=1}^V n_{k,\neg i}^{(t)} + \beta_t} (n_{m,\neg i}^{(k)} + \alpha_k) \end{aligned}$$

词分布和主题分布

$$\varphi_{k,t} = \frac{n_k^{(t)} + \beta_t}{\sum_{t=1}^V n_k^{(t)} + \beta_t}$$

$$\vartheta_{m,k} = \frac{n_m^{(k)} + \alpha_k}{\sum_{k=1}^K n_m^{(k)} + \alpha_k}$$



$$p(\vec{\vartheta}_m | \vec{z}_m, \vec{\alpha}) = \frac{1}{Z_{\vartheta_m}} \prod_{n=1}^{N_m} p(z_{m,n} | \vec{\vartheta}_m) \cdot p(\vec{\vartheta}_m | \vec{\alpha}) = \text{Dir}(\vec{\vartheta}_m | \vec{n}_m + \vec{\alpha})$$

$$p(\vec{\varphi}_k | \vec{z}, \vec{w}, \vec{\beta}) = \frac{1}{Z_{\varphi_k}} \prod_{\{i: z_i=k\}} p(w_i | \vec{\varphi}_k) \cdot p(\vec{\varphi}_k | \vec{\beta}) = \text{Dir}(\vec{\varphi}_k | \vec{n}_k + \vec{\beta})$$

Gibbs采样算法

Algorithm LdaGibbs($\{\vec{w}\}, \alpha, \beta, K$)

Input: word vectors $\{\vec{w}\}$, hyperparameters α, β , topic number K

Global data: count statistics $\{n_m^{(k)}\}, \{n_k^{(i)}\}$ and their sums $\{n_m\}, \{n_k\}$, memory for full conditional array $p(z_i|\cdot)$

Output: topic associations $\{\vec{z}\}$, multinomial parameters $\underline{\Phi}$ and $\underline{\Theta}$, hyperparameter estimates α, β

// initialisation

zero all count variables, $n_m^{(k)}, n_m, n_k^{(i)}, n_k$

for all documents $m \in [1, M]$ **do**

for all words $n \in [1, N_m]$ in document m **do**

 sample topic index $z_{m,n}=k \sim \text{Mult}(1/K)$

 increment document–topic count: $n_m^{(k)} += 1$

 increment document–topic sum: $n_m += 1$

 increment topic–term count: $n_k^{(i)} += 1$

 increment topic–term sum: $n_k += 1$

// Gibbs sampling over burn-in period and sampling period

while not finished **do**

for all documents $m \in [1, M]$ **do**

for all words $n \in [1, N_m]$ in document m **do**

 // for the current assignment of k to a term t for word $w_{m,n}$:

 decrement counts and sums: $n_m^{(k)} -= 1; n_m -= 1; n_k^{(i)} -= 1; n_k -= 1$

 // multinomial sampling acc. to Eq. 78 (decrements from previous step):

 sample topic index $\tilde{k} \sim p(z_i|\vec{z}_{-i}, \vec{w})$

 // for the new assignment of $z_{m,n}$ to the term t for word $w_{m,n}$:

 increment counts and sums: $n_m^{(\tilde{k})} += 1; n_m += 1; n_k^{(i)} += 1; n_{\tilde{k}} += 1$

 // check convergence and read out parameters

if converged and L sampling iterations since last read out **then**

 // the different parameters read outs are averaged.

 read out parameter set $\underline{\Phi}$ according to Eq. 81

 read out parameter set $\underline{\Theta}$ according to Eq. 82

代码实现

□ 数目：

- 文档数目： M

- 词数目： V (非重复的, “term”)

- 主题数目： K

□ 记号：

- 用 d 表述第几个文档， k 表示主题， w 表示词汇(term)， n 表示词(word)

三个矩阵和三个向量

- $z[d][w]$: 第 d 篇文档的第 w 个词来自哪个主题。M行, X列, X为相应文档长度: 即词(可重复)的数目。
- $nw[w][t]$: 第 w 个词是第 t 个主题的次数。word-topic矩阵, 列向量 $nw[][t]$ 表示主题 t 的词频数分布; V行K列
- $nd[d][t]$: 第 d 篇文档中第 t 个主题出现的次数, doc-topic矩阵, 行向量 $nd[d]$ 表示文档 d 的主题频数分布。M行, K列。
- 辅助向量:
 - $ntSum[t]$: 第 t 个主题在所有语料出现的次数, K维
 - $ndSum[d]$: 第 d 篇文档中词的数目(可重复), M维;
 - $P[t]$: 对于当前计算的某词属于主题 t 的概率, K维。

Code

```
if __name__ == "__main__":
    doc_num = 10 # 文档数目
    # 载入停止词库
    stop_words = load_stopwords()
    dic = {}
    doc = read_document(doc_num, stop_words, dic)

    # LDA
    term_num = len(dic) # 词汇的数目
    # nt[w][t]: 第term个词属于第t个主题的次数
    nt = [[0 for t in range(topic_number)] for term in range(term_num)]
    # nd[d][t]: 第d个文档中出现第t个主题的次数
    nd = [[0 for t in range(topic_number)] for d in range(doc_num)]
    # nt_sum[t]: 第t个主题出现的次数(nt矩阵的第t列)
    nt_sum = [0 for t in range(topic_number)]
    # nd_sum[d]: 第d个文档的长度(nd矩阵的第d行)
    nd_sum = [0 for d in range(doc_num)]
    z = init_topic(doc, nt, nd, nt_sum, nd_sum, dic)
    theta, phi = lda(z, nt, nd, nt_sum, nd_sum, dic, doc)
    show_result(theta, phi, dic) # 输出每个文档的主题和每个主题的关键字
```

Code

```
def lda(z, nt, nd, nt_sum, nd_sum, dic, doc):
    doc_num = len(z)
    for time in range(50):
        for m in range(doc_num):
            doc_length = len(z[m])
            for i in range(doc_length):
                term = dic[doc[m][i]] # 词语 -> 词汇
                gibbs_sampling(z, m, i, nt, nd, nt_sum, nd_sum, term)
    theta = calc_theta(nd, nd_sum) # 计算每个文档的主题分布
    phi = calc_phi(nt, nt_sum) # 计算每个主题的词分布
    return theta, phi
```

Code

```
def calc_theta(nd, nd_sum):  # 每个文档的主题分布
    doc_num = len(nd)
    topic_alpha = topic_number * alpha
    theta = [[0 for t in range(topic_number)] for d in range(doc_num)]
    for m in range(doc_num):
        for k in range(topic_number):
            theta[m][k] = (nd[m][k] + alpha) / (nd_sum[m] + topic_alpha)
    return theta

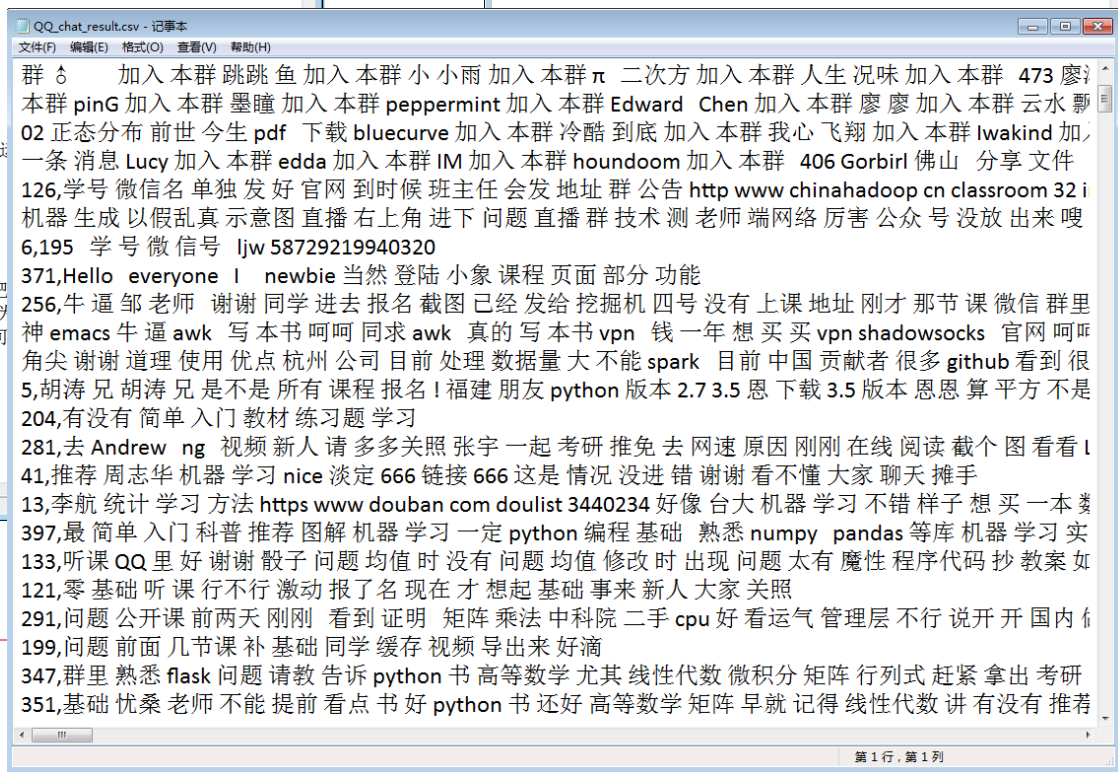
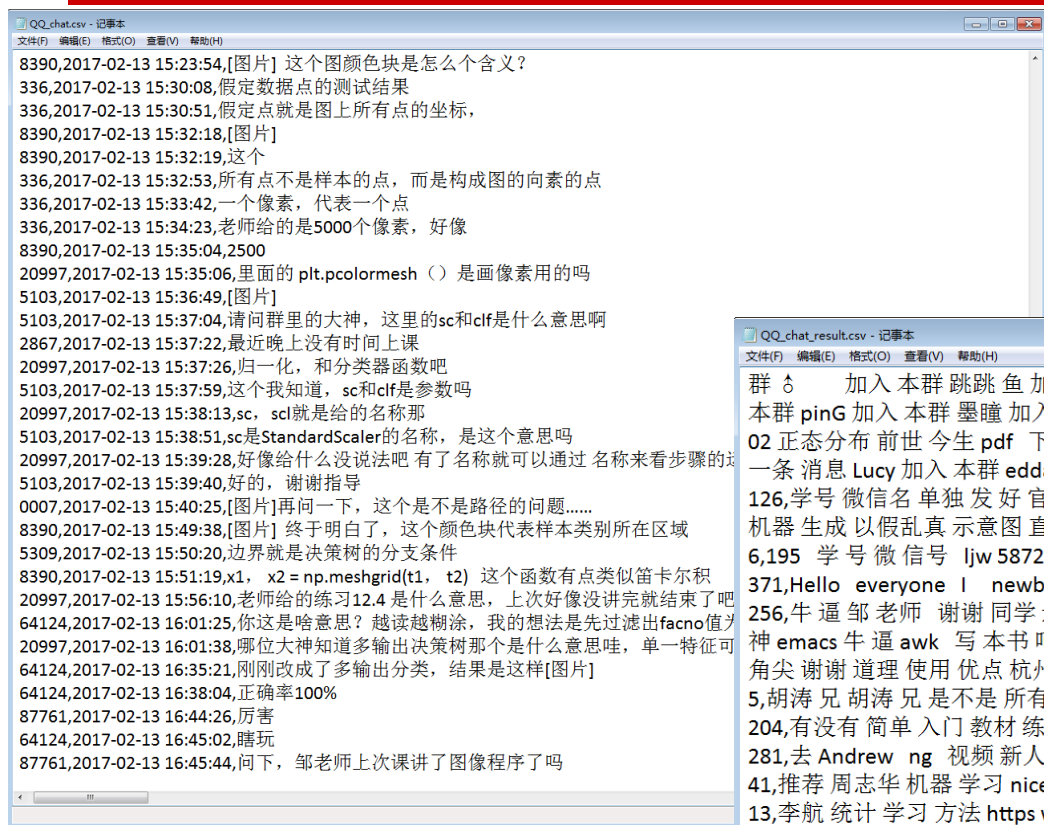
def calc_phi(nt, nt_sum):  # 每个主题的词分布
    term_num = len(nt)
    term_beta = term_num * beta
    phi = [[0 for w in range(term_num)] for t in range(topic_number)]
    for k in range(topic_number):
        for term in range(term_num):
            phi[k][term] = (nt[term][k] + beta) / (nt_sum[k] + term_beta)
    return phi
```


文档和主题

```
文档 1 : 茶馆 ( 0.0163591635916 ) 社会 ( 0.00528905289053 ) 王利发 ( 0.00528905289053 )
文档 2 : 决议 ( 0.0138983050847 ) 打击 ( 0.00824858757062 ) 安理会 ( 0.00824858757062 )
文档 3 : 会议 ( 0.0124491456469 ) 脱贫 ( 0.0124491456469 ) 党校 ( 0.0108218063466 )
文档 4 : 美团 ( 0.0306066176471 ) 阿里 ( 0.0103860294118 ) 业务 ( 0.0103860294118 )
文档 5 : 户口 ( 0.0221347331584 ) 登记 ( 0.0195100612423 ) 人口 ( 0.0142607174103 )
文档 6 : 人员 ( 0.0111471861472 ) 飞机 ( 0.00898268398268 ) 称 ( 0.00681818181818 )
文档 7 : 号线 ( 0.0328544061303 ) 站 ( 0.0194444444444 ) 14 ( 0.0184865900383 )
文档 8 : 支付 ( 0.0198394495413 ) 腾讯 ( 0.0072247706422 ) 支付宝 ( 0.0072247706422 )
文档 9 : 决议 ( 0.0138983050847 ) 打击 ( 0.00824858757062 ) 安理会 ( 0.00824858757062 )
文档 10 : 足协 ( 0.0186473429952 ) 足球 ( 0.0138164251208 ) 佩兰 ( 0.011884057971 )

=====
主题 1 : 美团 ( 0.0306066176471 ) 阿里 ( 0.0103860294118 ) 业务 ( 0.0103860294118 )
主题 2 : 会议 ( 0.0124491456469 ) 脱贫 ( 0.0124491456469 ) 党校 ( 0.0108218063466 )
主题 3 : 号线 ( 0.0328544061303 ) 站 ( 0.0194444444444 ) 14 ( 0.0184865900383 )
主题 4 : 人物 ( 0.00214876033058 ) 民族 ( 0.00214876033058 ) 资本家 ( 0.00214876033058 )
主题 5 : 足协 ( 0.0186473429952 ) 足球 ( 0.0138164251208 ) 佩兰 ( 0.011884057971 )
主题 6 : 户口 ( 0.0221347331584 ) 登记 ( 0.0195100612423 ) 人口 ( 0.0142607174103 )
主题 7 : 决议 ( 0.0138983050847 ) 打击 ( 0.00824858757062 ) 安理会 ( 0.00824858757062 )
主题 8 : 人员 ( 0.0111471861472 ) 飞机 ( 0.00898268398268 ) 称 ( 0.00681818181818 )
主题 9 : 茶馆 ( 0.0163591635916 ) 社会 ( 0.00528905289053 ) 王利发 ( 0.00528905289053 )
主题 10 : 支付 ( 0.0198394495413 ) 腾讯 ( 0.0072247706422 ) 支付宝 ( 0.0072247706422 )
```

聊天记录分析感兴趣话题



数据处理流程

□ 获取QQ聊天记录：txt文本格式(图1)

□ 整理成“QQ号/时间/留言”的规则形

■ 正则表达式

■ 清洗特定词：表情、@XX

■ 使用停止词库

■ 获得CSV表格数据(图2)

□ 合并相同QQ号的留言

■ 长文档利于计算每人感兴趣话题(图3)

□ LDA模型计算主题

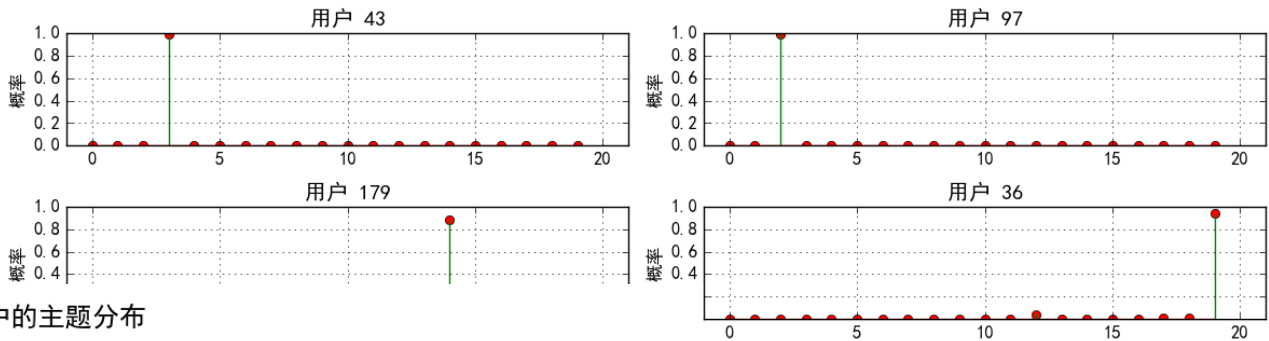
■ 调参与可视化

□ 计算每个QQ号及众人感兴趣话题

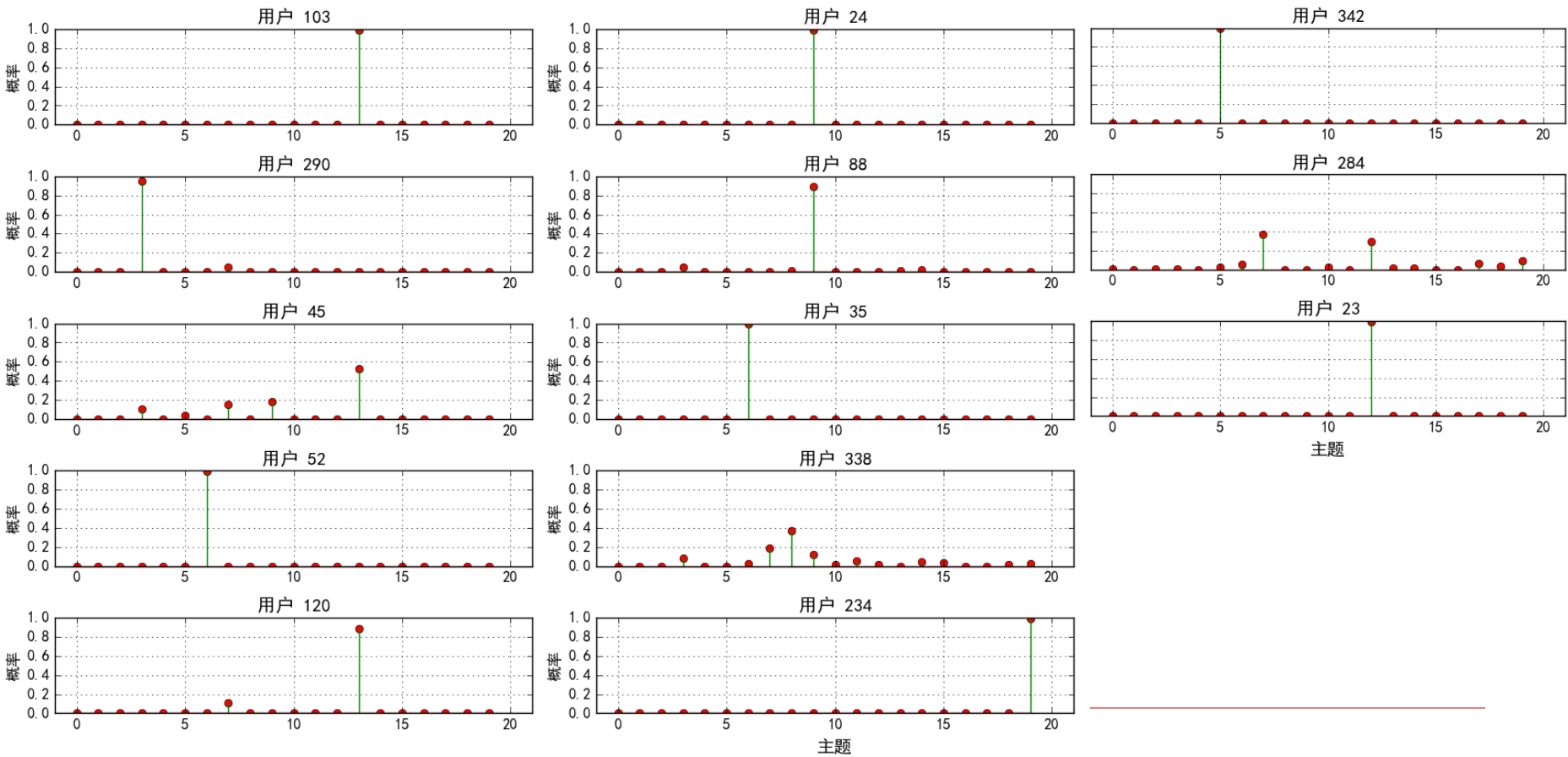


主题分布

用户的主题分布

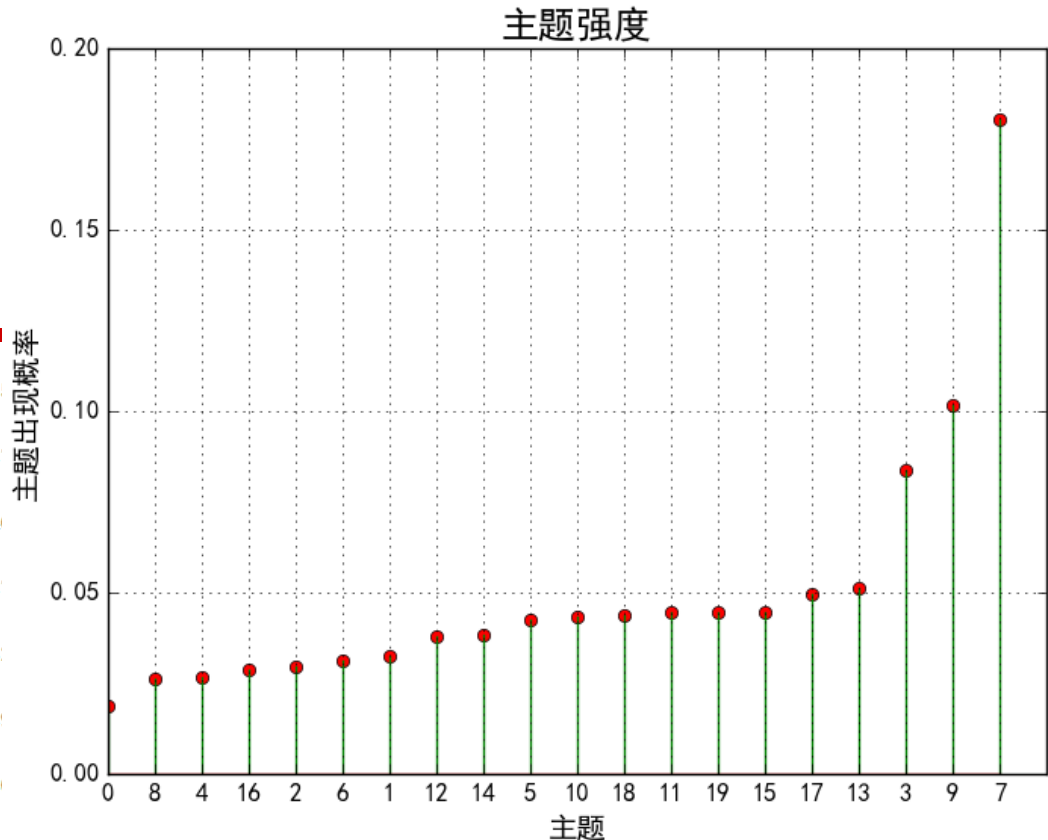


用户的主题分布



感兴趣话题

主题#1:	下周 用到 绑定 matlab 手册 详细描述 666
概率:	[0.00386318 0.00379716 0.00275161 0.00267257 0.0025779 0.002
主题#2:	决策树 下载 代码 新入 记得 无法 if
概率:	[0.00459593 0.00368178 0.00332368 0.00322059 0.00282207 0.002
主题#3:	下载 视频 app 中 cn 电脑
概率:	[0.01650318 0.00643555 0.00534376 0.00493769 0.00477131 0.004
主题#4:	互扫 上装 停留 一份 层次 缺 找点
概率:	[0.00457096 0.00322019 0.00314435 0.00283468 0.0027911 0.002
主题#5:	同问 变量 极限 训练 毕竟 问
概率:	[0.00799138 0.00796053 0.00353821 0.00350646 0.00323556 0.003
主题#6:	上网 元素 等待 白屏 中国 在线看
概率:	[0.00678803 0.00387909 0.00340132 0.00335809 0.00303858 0.002
主题#7:	好 大家 老师 学习 没有 谢谢
概率:	[0.02435813 0.01137977 0.01116381 0.01048067 0.01000842 0.008
主题#8:	参考书 需不需要 离散 型 基础知识 组成
概率:	[0.00909812 0.00337894 0.00242403 0.00229328 0.00220396 0.00219262 0.00209999]
主题#9:	com 9lpintuan wxe9c061d2aed9ae09 group https 请
概率:	[0.01010167 0.00642439 0.00517232 0.00517232 0.00495921 0.00485321 0.00467808]
主题#10:	画面 端 找回 QQ 群 账号
概率:	[0.00641685 0.00384681 0.00359311 0.00355134 0.00343251 0.00288641 0.00272948]
主题#11:	直播 今天 一片空白 没有 进不去 问题
概率:	[0.01035459 0.00352562 0.00291557 0.00276692 0.00275323 0.00259631 0.00257422]
主题#12:	太过分 编程 讲完 闪照 配置 这部分
概率:	[0.00510845 0.00290508 0.00289446 0.00273205 0.00265899 0.00251825 0.00249154]
主题#13:	pydotplus 学号 积分 浏览器 弄 请问
概率:	[0.01348479 0.00478899 0.00342937 0.00260794 0.00236339 0.00226667 0.00213521]
主题#14:	分享 两个 满足 话 京东 需求
概率:	[0.0143761 0.00353596 0.00270139 0.00264514 0.00255197 0.00249422 0.00244258]
主题#15:	相等 听不见 openCV 版 数据库 设置
概率:	[0.01212347 0.00449718 0.0030489 0.0030018 0.00285451 0.00281434 0.00263035]
主题#16:	回放 课程 OpenStack 牛云 投放 点哇
概率:	[0.01010629 0.00425028 0.0026812 0.00238546 0.00221883 0.00221883 0.00221883]
主题#17:	福 敬业 早 断 厉害 歌



附：正则表达式

语法	说明	表达式实例	完整匹配的字符串
字符			
一般字符	匹配自身	abc	abc
.	匹配任意除换行符"\n"外的字符。 在DOTALL模式中也能匹配换行符。	a.c	abc
\	转义字符，使后一个字符改变原来的意思。 如果字符串中有字符*需要匹配，可以使用\"或者字符集[*]。	a\\.c a\\c	a.c a\\c
[...]	字符集（字符类）。对应的位置可以是字符集中任意字符。 字符集中的字符可以逐个列出，也可以给出范围，如[abc]或[a-c]。第一个字符如果是^则表示取反，如[^abc]表示不是abc的其他字符。 所有的特殊字符在字符集中都失去其原有的特殊含义。在字符集中如果要使用]、-或^，可以在前面加上反斜杠，或把]、-放在第一个字符，把^放在非第一个字符。	a[bcd]e	abe ace ade
预定义字符集（可以写在字符集[...]中）			
\\d	数字：[0-9]	a\\dc	a1c
\\D	非数字：[^\\d]	a\\Dc	abc
\\s	空白字符：[<空格>\\t\\r\\n\\f\\v]	a\\sc	a c
\\S	非空白字符：[^\\s]	a\\Sc	abc
\\w	单词字符：[A-Za-z0-9_]	a\\wc	abc
\\W	非单词字符：[^\\w]	a\\Wc	a c
数量词（用在字符或(...)之后）			
*	匹配前一个字符0或无限次。	abc*	ab abccc
+	匹配前一个字符1次或无限次。	abc+	abc abccc
?	匹配前一个字符0次或1次。	abc?	ab abc
{m}	匹配前一个字符m次。	ab{2}c	abbc
{m,n}	匹配前一个字符m至n次。 m和n可以省略：若省略m，则匹配0至n次；若省略n，则匹配m至无限次。	ab{1,2}c	abc abbc
*? +? ?? {m,n}?	使 * + ? {m,n} 变成非贪婪模式。	示例将在下文中介。	

边界匹配（不消耗待匹配字符串中的字符）			
^	匹配字符串开头。 在多行模式中匹配每一行的开头。	^abc	abc
\$	匹配字符串末尾。 在多行模式中匹配每一行的末尾。	abc\$	abc
\\A	仅匹配字符串开头。	\\Aabc	abc
\\Z	仅匹配字符串末尾。	abc\\Z	abc
\\b	匹配\\w和\\W之间。	a\\b!bc	a!bc
\\B	[^\\b]	a\\Bbc	abc
逻辑、分组			
	代表左右表达式任意匹配一个。 它总是先尝试匹配左边的表达式，一旦成功匹配则跳过匹配右边的表达式。 如果 没有被包括在()中，则它的范围是整个正则表达式。	abc def	abc def
(...)	被括起来的表达式将作为分组，从表达式左边开始每遇到一个分组的左括号'('，编号+1。 另外，分组表达式作为一个整体，可以后接数量词。表达式中的 仅在该组中有效。	(abc){2} a(123 456)c	abcaabc a456c
(?P<name>...)	分组，除了原有的编号外再指定一个额外的别名。	(?P<id>abc){2}	abcaabc
\\<number>	引用编号为<number>的分组匹配到的字符串。	(\\d)abc\\1	1abc1 5abc5
(?P=name)	引用别名为<name>的分组匹配到的字符串。	(?P<id>\\d)abc(?P=id)	1abc1 5abc5
特殊构造（不作为分组）			
(?...)	(...)的不分组版本，用于使用' '或后接数量词。	(?:abc){2}	abcaabc
(?i ms ux)	i ms ux的每个字符代表一个匹配模式，只能用在正则表达式的开头，可选多个。匹配模式将在下文中介。	(?i)abc	AbC
(?#...)	#后的内容将作为注释被忽略。	abc(?#comment)123	abc123
(?=...)	之后的字符串内容需要匹配表达式才能成功匹配。 不消耗字符串内容。	a(?=\\d)	后面是数字的a
(?!...)	之后的字符串内容需要不匹配表达式才能成功匹配。 不消耗字符串内容。	a(?!\\d)	后面不是数字的a
(?<=...)	之前的字符串内容需要匹配表达式才能成功匹配。 不消耗字符串内容。	(?<=\\d)a	前面是数字的a
(?<!=...)	之前的字符串内容需要不匹配表达式才能成功匹配。 不消耗字符串内容。	(?<!=\\d)a	前面不是数字的a
(?(id/name) yes-pattern no-pattern)	如果编号为id/别名为name的组匹配到字符，则需要匹配yes-pattern，否则需要匹配no-pattern。 no-pattern可以省略。	(\\d)abc(?:1 \\d abc)	1abc2 abcabc

常用正则表达式

- ❑ 匹配中文字符: `[\u4e00-\u9fa5]`
- ❑ 匹配双字节字符(包括汉字在内): `[^\x00-\xff]`
- ❑ 匹配空白行: `\n\s*\r`
- ❑ 匹配HTML标记: `<(\S*?)[^>]*>.*?</\1>|<.*? />`
- ❑ 匹配首尾空白字符: `^\s*|\s*$`
- ❑ 匹配Email地址: `\w+([-+.] \w+)*@\w+([-.] \w+)*\.\w+([-.] \w+)*`
- ❑ 匹配网址URL: `[a-zA-z]+://[^\s]*`
- ❑ 匹配帐号合法(5-16位, 字母开头, 允许字母数字下划线): `^[a-zA-Z][a-zA-Z0-9_]{4,15}$`
- ❑ 匹配国内电话号码: `\d{3}-\d{8}|\d{4}-\d{7}`
- ❑ 匹配腾讯QQ号: `[1-9][0-9]{4,}`
- ❑ 匹配中国邮政编码: `[1-9]\d{5}(?! \d)`
- ❑ 匹配身份证: `\d{15}|\d{18}|\d{17}[xX]`
- ❑ 匹配ip地址: `\d{1,3}\.\d{1,3}\.\d{1,3}\.\d{1,3}`

常用正则表达式

□ 匹配特定数字：

- 匹配正整数：`^[1-9]\d*$`
- 匹配负整数：`^-[1-9]\d*$`
- 匹配整数：`^-?[1-9]\d*$`
- 匹配非负整数(正整数 + 0)：`^[1-9]\d*|0$`
- 匹配非正整数(负整数 + 0)：`^-[1-9]\d*|0$`
- 匹配正浮点数：`^[1-9]\d*\.\d*|0\.\d*[1-9]\d*$`
- 匹配负浮点数：`^-([1-9]\d*\.\d*|0\.\d*[1-9]\d*)$`
- 匹配浮点数：`^-?([1-9]\d*\.\d*|0\.\d*[1-9]\d*|0?\.\d+|0)$`
- 匹配非负浮点数(正浮点数 + 0)：`^[1-9]\d*\.\d*|0\.\d*[1-9]\d*|0?\.\d+|0$`
- 匹配非正浮点数(负浮点数 + 0)：`^-([1-9]\d*\.\d*|0\.\d*[1-9]\d*)|0?\.\d+|0$`

□ 匹配特定字符串：

- 匹配由26个英文字母组成的字符串：`^[A-Za-z]+$`
- 匹配由26个英文字母的大写组成的字符串：`^[A-Z]+$`
- 匹配由26个英文字母的小写组成的字符串：`^[a-z]+$`
- 匹配由数字和26个英文字母组成的字符串：`^[A-Za-z0-9]+$`
- 匹配由数字26个英文字母或下划线组成的字符串：`^\w+$`

超参数的确定

- 交叉验证
- α 表达了不同文档间主题是否鲜明， β 度量了有多少近义词能够属于同一个类别。
- 主题数目 K ，词项数目为 W ，可以使用：
 - $\alpha=50/K$
 - $\beta=200/W$
 - 注：不一定普遍适用

一种迭代求超参数的方法

□ Digamma函数: $\Psi(x) = \frac{d \ln \Gamma(x)}{dx} = \frac{\Gamma'(x)}{\Gamma(x)}$

□ 迭代公式: (T. Minka)

$$\alpha_k = \frac{\left(\left(\sum_{m=1}^M \Psi(n_m^{(k)} + \alpha_k) \right) - M \cdot \Psi(\alpha_k) \right)}{\left(\sum_{m=1}^M \Psi\left(n_m + \sum_{j=1}^K \alpha_j\right) \right) - M \cdot \Psi\left(\sum_{j=1}^K \alpha_j\right)} \cdot \alpha_k$$

主题个数的确定

- 相似度最小
- 选取初始的主题个数 K ，训练LDA模型，计算各主题之间的相似度
- 增加或减少 K 的值，重新训练LDA模型，再次计算topic之间的相似度
- 选择相似度最小的模型所对应的 K 作为主题个数。

概率分布的困惑度/复杂度Perplexity

□ 某离散概率分布 p 的困惑度为：

$$Perplexity = 2^{H(p)} = 2^{-\sum_x p(x) \log_2 p(x)}$$

□ 样本集 $x_1, x_2 \dots x_n$ 的估计模型 q 的困惑度为：

$$Perplexity = a^{-\frac{1}{N} \sum_{i=1}^n \log_a q(x_i)}, \quad a \text{ 为任意整数}$$

■ 交叉熵为： $H(p, q) = -\sum_x p(x) \log_2 q(x)$

困惑度Perplexity与主题模型

□ 使用训练数据得到无监督模型，在测试数据集中计算所有token似然值几何平均数的倒数

■ 测试数据集中词典大小的期望

$$P(W | Model) = \prod_{i=1}^V p(w_i | Model)^{-\frac{1}{V}} = \exp\left(-\frac{1}{V} \cdot \sum_{i=1}^V \log p(w_i | Model)\right)$$

□ 其中，LDA中词的似然概率为：

$$P(\vec{w}_m | Model) = \prod_{n=1}^{N_m} \sum_{k=1}^K p(z = k | d = m) \cdot p(w = t | z = k) = \prod_{t=1}^V \left(\sum_{k=1}^K \vartheta_{m,k} \cdot \varphi_{k,t} \right)^{n_m^{(t)}}$$

附：PageRank

□ 一个网页*i*的重要度可以使用指向网页*i*的其他网页*j*的重要度加权得到。

■ 权值不妨取网页*j*包含的链接数目。

$$D(P_i) = (1-d) + d \cdot \sum_{j \in In(P_i)} \frac{1}{|Out(P_j)|} \cdot D(P_j)$$

□ 参数的意义为：

■ 网页*i*的中重要性 $D(P_i)$

■ 阻尼系数*d*，如设置为常系数0.85

■ 指向网页*i*的网页集合 $In(P_i)$

■ 网页*j*指向的网页集合 $Out(P_j)$

TextRank

□ 将PageRank中的“网页”换成“词”，结论仍成立。

■ 选择合适的窗口大小，认为窗口内的词相互指向。

$$D(w_i) = (1-d) + d \cdot \sum_{j \in \text{In}(w_i)} \frac{1}{|\text{Out}(w_j)|} \cdot D(w_j)$$

□ 句子 S_i 和 S_j 的相似度：

$$\text{Similar}(S_i, S_j) = \frac{|\{w_k \mid w_k \in S_i \ \& \ w_k \in S_j\}|}{\log|S_i| + \log|S_j|}$$

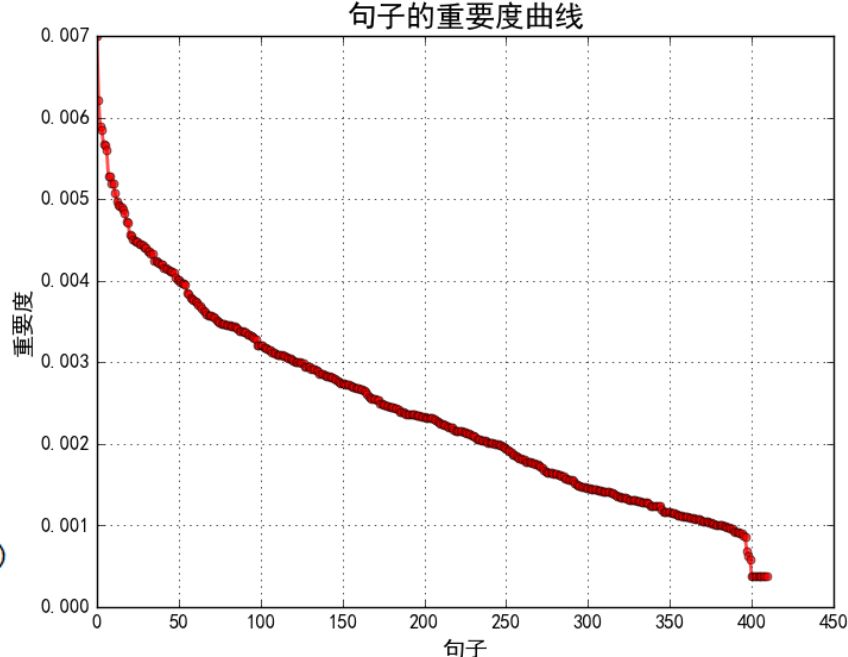
□ 将PageRank中“网页”换成“句子”，结论仍然基本成立，只需考虑将“链接”加权：

$$D(S_i) = (1-d) + d \cdot \sum_{j \in \text{In}(S_i)} \frac{\text{similar}(S_j, S_i)}{\sum_{k \in \text{Out}(S_j)} \text{similar}(S_j, S_k)} \cdot D(S_j)$$

Text Rank

```
tr4s = TextRank4Sentence()
tr4s.analyze(text=text, lower=True, source = 'no_stop_words')
data = pd.DataFrame(data=tr4s.key_sentences)
mpl.rcParams['font.sans-serif'] = [u'SimHei']
mpl.rcParams['axes.unicode_minus'] = False
plt.figure(facecolor='w')
plt.plot(data['weight'], 'ro-', lw=2, ms=5, alpha=0.7)
plt.grid(b=True)
plt.xlabel(u'句子', fontsize=14)
plt.ylabel(u'重要度', fontsize=14)
plt.title(u'句子的重要度曲线', fontsize=18)
plt.show()

key_sentences = tr4s.get_key_sentences(num=20, sentence_min_len=4)
for sentence in key_sentences:
    print sentence['weight'], sentence['sentence']
```



0.00699560759634 她知道我心里的苦闷，知道不该阻止我出去走走，知道我要是老呆在家里结果会更糟，但她又担心我一个人在那荒僻的园子里整天都想些什么

0.00621160375013 这一来你中了魔了，整天都在想哪一件事可以写，哪一个人可以让你写成小说

0.00588860912528 那时她的儿子，还太年轻，还来不及为母亲想，他被命运击昏了头，一心以为自己是世上最不幸的一个，不知道儿子的不幸在母亲那儿总是要加倍的

0.00584459738866 她想，只要儿子能活下去哪怕自己去死呢也行，可她又确信一个人不能仅仅是活着，儿子得有一条路走向自己的幸福

0.00567083997126 我奇怪这么小的孩子怎么一个人跑来这园子里

0.00565208315006 我那时脾气坏到极点，经常是发了疯一样地离开家，从那园子里回来又中了魔似的什么话都不说

0.00559372837107 如今我摇着车在这园子里慢慢走，常常有一种感觉，觉得我一个人跑出来已经玩得太久了

0.00527989619912 而且我想，他的母亲也比我的母亲运气好，他的母亲没有一个双腿残废的儿子，否则事情就不这么简单

0.00527906358787 年年月月我都到这园子里来，年年月月我都要想，母亲盼望我找到的那条路到底是什么

0.00519622569726 那天你又说你不如死了好，你的一个朋友劝你：你不能死，你还得写呢，还有好多好作品等着你去写呢

0.00519145626625 他的衣着过分随便，走路的姿态也不慎重，走上五六十米路便选定一处地方，一只脚踏在石凳上或土埂上或树墩上，解下腰间的酒瓶，解酒瓶的当儿

0.00507970004724 她一个人在园子里走，走过我的身旁，走过我经常呆的一些地方，步履茫然又急迫

0.00497335014554 “在那段日子里——那是好几年长的一段日子，我想我一定使母亲作过了最坏的准备了，但她从来没有对我说过：“你为我想想”

0.0049360646412 我便又不能在家里呆了，又整天整天独自跑到地坛去，心里是没头没尾的沉郁和哀怨，走遍整个园子却怎么也想不通：母亲为什么就不能再多活两年

0.00491815078362 是中魔了，我走到哪儿想到哪儿，在人山人海只寻找小说，要是有一种小说试剂就好了，见人就滴两滴看他是不是一篇小说，要是有一种小说显

0.00490464531034 我在这园子里坐着，我听见园神告诉我，每一个有激情的演员都难免是一个人质

0.00486932833768 十五年前的一个下午，我摇着轮椅进入园中，它为一个失魂落魄的人把一切都准备好了

0.00483241065578 有一天我在这园子碰见一个老太太，她说：“哟，你还在这儿哪

0.0047216969869 我才想到，当年我总是独自跑到地坛去，曾经给母亲出了一个怎样的难题

0.00470900507196 我带着本子和笔，到园中找一个最不为人打扰的角落，偷偷地写

LDA的实现

- LDA-C: David Blei, C实现, VBEM参数估计
 - <http://www.cs.princeton.edu/~blei/lda-c/index.html>
- GibbsLDA++/JGibbLDA : C/C++实现/Java实现
 - <http://gibbslda.sourceforge.net/> <http://jgibblda.sourceforge.net/>
 - Xuan-Hieu Phan/Cam-Tu Nguyen, 输入输出一致
- Matlab Topic Modeling Toolbox 1.4, Mark Steyvers, Gibbs采样
 - http://psiexp.ss.uci.edu/research/programs_data/toolbox.htm
- Gensim: Online VB
 - 官网: <http://radimrehurek.com/gensim/index.html>
 - github: http://www.cs.columbia.edu/~blei/topicmodeling_software.html
- Scikit-learn: sklearn.decomposition.LatentDirichletAllocation/Online VB
- LDA/Online VB: <https://pypi.python.org/pypi/lda>
- LDA不完全列表:
 - http://www.cs.columbia.edu/~blei/topicmodeling_software.html

例：Gensim的安装

```
C:\Users\zou>pip install gensim
```

Downloading gen

100%  4.2MB 141kB/s

Requirement already satisfied (use --upgrade to upgrade): scipy>=0.7.0 in c:\python27\lib\site-packages (from gensim)

Requirement already satisfied (use --upgrade to upgrade): six>=1.5.0 in c:\python27\lib\site-packages (from gensim)

Requirement already satisfied (use --upgrade to upgrade): numpy>=1.3 in c:\python27\lib\site-packages (from gensim)

```
Downloading smart open-1.3.4.tar.gz
```

Collecting boto>=2.3.2 (from smart-open>=1.2.1->gensim)

100% ██████████ 1.4MB 249kB/s

```
Collecting bz2file (from smart-open>=1.2.1->gensim)
```

Downloading bz2file-0.98.tar.gz

```
downloading requests-2.11.1-py2.py3-none-any.whl (514kB)
```

100% 

Running setup.py install for smart-open ... done

TF-IDF

```
Text =  
[['human', 'machine', 'interface', 'lab', 'abc', 'computer', 'applications'],  
 ['survey', 'user', 'opinion', 'computer', 'system', 'response', 'time'],  
 ['eps', 'user', 'interface', 'management', 'system'],  
 ['system', 'human', 'system', 'engineering', 'testing', 'eps'],  
 ['relation', 'user', 'perceived', 'response', 'time', 'error', 'measurement'],  
 ['generation', 'random', 'binary', 'unordered', 'trees'],  
 ['intersection', 'graph', 'paths', 'trees'],  
 ['graph', 'minors', 'iv', 'widths', 'trees', 'well', 'quasi', 'ordering'],  
 ['graph', 'minors', 'survey']]
```

TF-IDF:

```
[(0, 0.4301019571350565), (1, 0.4301019571350565), (2, 0.4301019571350565), (3, 0.4301019571350565), (4, 0.2944198962221451), (5,  
 [4, 0.3726494271826947), (7, 0.27219160459794917), (8, 0.3726494271826947), (9, 0.27219160459794917), (10, 0.3726494271826947),  
 [6, 0.438482464916089), (7, 0.32027755044706185), (9, 0.32027755044706185), (13, 0.6405551008941237), (14, 0.438482464916089)]  
 [(5, 0.3449874408519962), (7, 0.5039733231394895), (14, 0.3449874408519962), (15, 0.5039733231394895), (16, 0.5039733231394895)]  
 [(9, 0.21953536176370683), (10, 0.30055933182961736), (12, 0.30055933182961736), (17, 0.43907072352741366), (18, 0.4390707235274  
 [(21, 0.48507125007266594), (22, 0.48507125007266594), (23, 0.48507125007266594), (24, 0.48507125007266594), (25, 0.242535625036  
 [(25, 0.31622776601683794), (26, 0.31622776601683794), (27, 0.6324555320336759), (28, 0.6324555320336759)]  
 [(25, 0.20466057569885868), (26, 0.20466057569885868), (29, 0.2801947048062438), (30, 0.40932115139771735), (31, 0.4093211513977  
 [(8, 0.6282580468670046), (26, 0.45889394536615247), (29, 0.6282580468670046)]
```

LSI

LSI Model:

```
[(0, 0.34057117986841989), (1, -0.20602251622679696)],  
[(0, 0.69330400021715577), (1, 0.0072327583903918488)],  
[(0, 0.59026076703897357), (1, -0.35260469490855789)],  
[(0, 0.52149018218251453), (1, -0.33887976154055377)],  
[(0, 0.39533193176354431), (1, -0.059192853366596486)],  
[(0, 0.036353173528493307), (1, 0.18146550208818862)],  
[(0, 0.14709012328778862), (1, 0.49432948127822229)],  
[(0, 0.21407117317565286), (1, 0.640645666445394)],  
[(0, 0.40066568318170664), (1, 0.64131082990940158)]]
```

LSI Topics:

```
[(0,  
  u' 0.400*"system" + 0.318*"survey" + 0.290*"user" + 0.274*"eps" + 0.236*"management"',  
  (1,  
    u' 0.421*"minors" + 0.420*"graph" + 0.293*"survey" + 0.239*"trees" + 0.226*"intersection"')]
```

思考

- LSI/LFM/ICA 的关系
- LSI和pLSA的关系

相似度

Similarity:

```
[array([ 1.          ,  0.85017949,  0.99998462,  0.99948108,  0.92283762,
        -0.33944285, -0.2520774 , -0.21974573,  0.01438823], dtype=float32),
 array([ 0.85017949,  1.          ,  0.85309052,  0.83277911,  0.98737705,
        0.20664607,  0.29518002,  0.32680073,  0.53867108], dtype=float32),
 array([ 0.99998462,  0.85309052,  1.          ,  0.99928677,  0.92496276,
        -0.33421332, -0.24669874, -0.214324  ,  0.01994151], dtype=float32),
 array([ 0.99948108,  0.83277911,  0.99928677,  1.          ,  0.90995121,
        -0.36956567, -0.28311783, -0.25105584, -0.01782739], dtype=float32),
 array([ 0.92283762,  0.98737705,  0.92496276,  0.90995121,  1.          ,
        0.04906873,  0.14012395,  0.1729846 ,  0.39842743], dtype=float32),
 array([-0.33944285,  0.20664607, -0.33421332, -0.36956567,  0.04906873,
         1.          ,  0.99581695,  0.99222624,  0.93564534], dtype=float32),
 array([-0.2520774 ,  0.29518002, -0.24669874, -0.28311783,  0.14012395,
         0.99581695,  1.          ,  0.99944651,  0.96397996], dtype=float32),
 array([-0.21974573,  0.32680073, -0.214324  , -0.25105584,  0.1729846 ,
         0.99222624,  0.99944651,  0.99999994,  0.97229445], dtype=float32),
 array([ 0.01438823,  0.53867108,  0.01994151, -0.01782739,  0.39842743,
         0.93564534,  0.96397996,  0.97229445,  1.          ], dtype=float32)]
```

主题和主题分布

LDA Model:

Document-Topic:

```
[[ (0, 0.68548441915170544), (1, 0.31451558084829462)],  
 [ (0, 0.65732202058761513), (1, 0.34267797941238493)],  
 [ (0, 0.67101883898793013), (1, 0.32898116101206987)],  
 [ (0, 0.29774557750241137), (1, 0.70225442249758874)],  
 [ (0, 0.55150516193766697), (1, 0.44849483806233303)],  
 [ (0, 0.25456933670287446), (1, 0.7454306632971256)],  
 [ (0, 0.67476418767307922), (1, 0.32523581232692073)],  
 [ (0, 0.29509659300584296), (1, 0.7049034069941571)],  
 [ (0, 0.69445879658152987), (1, 0.30554120341847024)]]
```

Topic 0

```
[(u' survey', 0.042573497130974247),  
 (u' minors', 0.03943557671036535),  
 (u' graph', 0.038776707760135178),  
 (u' system', 0.034575198665359616),  
 (u' trees', 0.032742027152788719),  
 (u' opinion', 0.031680224783503845),  
 (u' generation', 0.031141365123546434),  
 (u' unordered', 0.030981049002428096),  
 (u' time', 0.030911535753312992),  
 (u' random', 0.03090147631201922)]
```

Topic 1

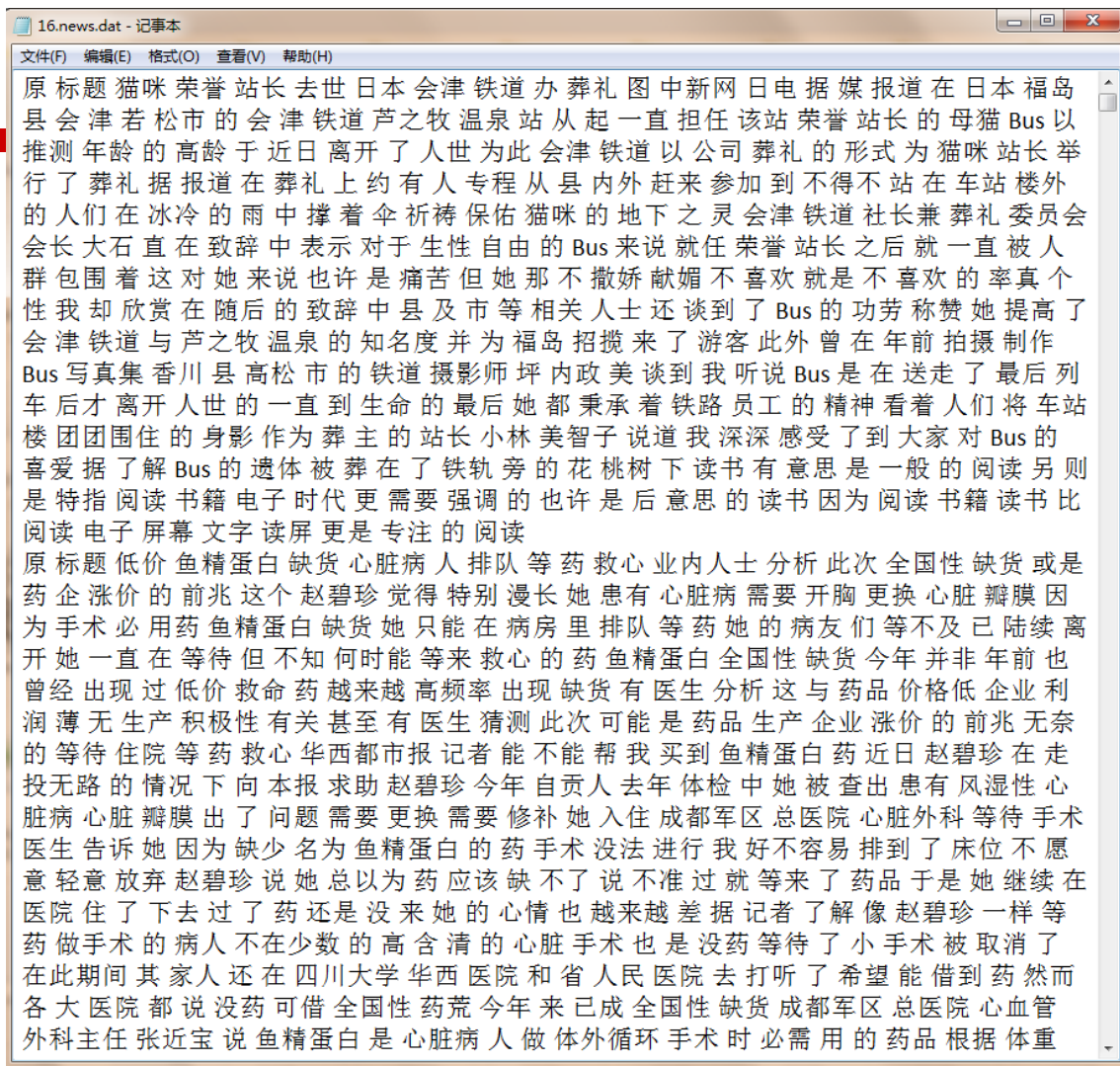
```
[(u' system', 0.037724259436260198),  
 (u' eps', 0.03524885080697393),  
 (u' interface', 0.034303635122775261),  
 (u' intersection', 0.03398428810730824),  
 (u' user', 0.033982740385072041),  
 (u' management', 0.033477115230294417),  
 (u' human', 0.032957111835837112),  
 (u' paths', 0.032333361709319365),  
 (u' engineering', 0.030715385582341159),  
 (u' computer', 0.030706245324429286)]
```

LDA计算的相似度

Similarity:

```
[array([ 0.99999994,  0.79683411,  0.99871153,  0.9988395 ,  0.99509394,
         0.68600154,  0.98457313,  0.66293609,  0.71091771], dtype=float32),
 array([ 0.79683411,  1.          ,  0.82646847,  0.82500935,  0.85270071,
         0.98624408,  0.89025998,  0.98059875,  0.99140108], dtype=float32),
 array([ 0.99871153,  0.82646847,  1.          ,  0.99999666,  0.9988324 ,
         0.72204095,  0.99218392,  0.70007479,  0.74569058], dtype=float32),
 array([ 0.9988395 ,  0.82500935,  0.99999666,  1.          ,  0.99870408,
         0.72024882,  0.99185783,  0.69822526,  0.74396449], dtype=float32),
 array([ 0.99509394,  0.85270071,  0.9988324 ,  0.99870408,  0.99999994,
         0.75462055,  0.99705368,  0.7337535 ,  0.77700794], dtype=float32),
 array([ 0.68600154,  0.98624408,  0.72204095,  0.72024882,  0.75462055,
         1.          ,  0.80272919,  0.9995119 ,  0.9993937 ], dtype=float32),
 array([ 0.98457313,  0.89025998,  0.99218392,  0.99185783,  0.99705368,
         0.80272919,  1.          ,  0.78370738,  0.82300484], dtype=float32),
 array([ 0.66293609,  0.98059875,  0.70007479,  0.69822526,  0.7337535 ,
         0.9995119 ,  0.78370738,  1.          ,  0.99781829], dtype=float32),
 array([ 0.71091771,  0.99140108,  0.74569058,  0.74396449,  0.77700794,
         0.9993937 ,  0.82300484,  0.99781829,  1.          ], dtype=float32)]
```

网易新闻语料



LDA

```
初始化停止词列表 --
开始读入语料数据 --
读入语料数据完成，用时9.256秒
文本数目：2043个
正在建立词典 --
正在计算文本向量 --
正在计算文档TF-IDF --
建立文档TF-IDF完成，用时0.185秒
LDA模型拟合推断 --
LDA模型完成，训练时间为 37.687秒
10个文档的主题分布：
第532个文档的前10个主题： [20 18 25  4 13  6 19 28 22 24]
[ 0.50757285  0.10239849  0.07296044  0.05082813  0.02763301  0.02729199
  0.02345142  0.02105525  0.01749429  0.01665937]
第1043个文档的前10个主题： [23 14  4 18 10 24  6 15  0 20]
[ 0.4981378  0.06441008  0.05225744  0.05120348  0.0392068  0.03119684
  0.02928527  0.02618645  0.02403664  0.02328459]
第1035个文档的前10个主题： [19 25  4 11 15  0 16 28  6  7]
[ 0.26742334  0.16533452  0.08484096  0.07141483  0.0688248  0.05389866
  0.05103031  0.03916642  0.03554837  0.027476  ]
第588个文档的前10个主题： [ 7 20  5 12 19 21 15 17 23 14]
[ 0.26408634  0.20762942  0.1160332  0.10415797  0.06068137  0.05660975
  0.02997539  0.01992539  0.01928632  0.01816658]
第1412个文档的前10个主题： [ 6 25  3 22 26 16 19  4 18  7]
[ 0.16465983  0.15589012  0.15210117  0.1234063  0.08512253  0.0831406
  0.04052934  0.03234385  0.0246687  0.02238315]
第805个文档的前10个主题： [ 1 25 19  4 10 15 23 28 26 18]
[ 0.33525038  0.23190863  0.09825045  0.06684136  0.05441141  0.0435245
  0.03313123  0.01985919  0.01859455  0.01445226]
```


主题

每个主题的词分布:

主题#0:

词: 村民 乘客 云南 旅客 地上 裤子 妈妈

概率: [0.00682393 0.0042878 0.00323379 0.00318589 0.00316816 0.00306
0.0029545]

主题#1:

词: 广东省 王某 刘某 皋丸 参议院 榆阳区 陈满

概率: [0.00751067 0.00640128 0.00602311 0.00560491 0.00496156 0.00429384
0.00428849]

主题#2:

词: 工匠 台当局 失误 退役 假如 暴力事件 其一

概率: [0.00298584 0.00257124 0.00240675 0.002152 0.0019788 0.00188708
0.00166307]

主题#3:

词: 李某 充值 工资 毫米 小杰 平均工资 徐某

概率: [0.01106934 0.00335774 0.00318256 0.00301278 0.00299619 0.00285581
0.00280268]

主题#4:

词: 阅读 读书 李 女子 视频 书籍 电子

概率: [0.00996042 0.00583026 0.00567708 0.00562837 0.00477045 0.00399124
0.0039597]

主题#5:

词: 普京 伦敦 俄 会谈 安倍 身份证 被捕

概率: [0.00617236 0.00446138 0.00441558 0.0041976 0.00326962 0.00310108
0.00297686]

主题#6:

词: 企业 政府 患者 公司 医院 建设 医疗

概率: [0.00433506 0.00424583 0.0039865 0.00391137 0.00326307 0.0031044
0.00304006]

路透社数据

159 0:1 2:1 6:1 9:1 12:5 13:2 20:1 21:4 24:2 29:1 35:1 38:2 39:7 48:1 49:1 54:1 59:2 60:1 61:7 66:1
107 0:7 2:2 7:1 16:1 17:1 20:1 24:1 38:3 42:1 59:1 62:1 65:2 70:1 76:2 84:1 87:1 90:2 101:1 107:1 1
153 3:1 4:10 6:4 7:1 8:1 11:9 13:1 20:1 31:3 32:1 33:1 35:2 44:5 45:3 48:5 49:1 62:1 64:1 68:1 71:2
156 0:6 2:1 6:1 7:1 8:1 12:7 18:3 19:1 21:3 22:1 24:3 26:3 27:1 37:1 39:2 40:1 45:1 57:2 60:2 61:3
192 3:2 4:14 5:1 6:1 8:2 9:1 11:11 13:2 14:1 15:3 20:1 26:1 30:1 31:5 33:1 34:1 35:2 37:1 41:1 43:1
180 2:2 3:2 4:24 6:2 8:2 9:1 11:16 13:2 15:2 26:1 31:3 33:3 34:1 35:2 37:3 44:1 48:4 49:1 57:3 64:1
147 3:2 4:7 5:1 6:1 8:1 11:5 13:1 14:1 15:1 31:1 32:1 33:2 34:1 35:2 37:1 41:1 44:4 45:1 48:2 49:2
0 UK: Prince Charles spearheads British royal revolution. LONDON 1996-08-20
184 2:2 3:2 4:20 6:2 8:3 9:1 11:15 13:1 15:1 21
1 GERMANY: Historic Dresden church rising from WW2 ashes. DRESDEN, Germany 1996-08-21
163 1:1 3:2 4:17 5:2 6:2 11:14 13:2 14:2 26:1 2
2 INDIA: Mother Teresa's condition said still unstable. CALCUTTA 1996-08-23
187 0:2 2:2 5:2 7:1 9:3 12:11 14:1 16:1 18:1 13
3 UK: Palace warns British weekly over Charles pictures. LONDON 1996-08-25
170 0:2 3:1 4:1 7:1 12:15 15:2 18:3 19:1 20:1 4
4 INDIA: Mother Teresa, slightly stronger, blesses nuns. CALCUTTA 1996-08-25
224 0:1 2:4 4:3 5:1 6:2 7:2 8:2 9:1 10:3 13:1 5
5 INDIA: Mother Teresa's condition unchanged, thousands pray. CALCUTTA 1996-08-25
193 0:1 1:1 2:1 3:1 4:10 6:2 7:3 8:2 11:10 13:7
6 INDIA: Mother Teresa shows signs of strength, blesses nuns. CALCUTTA 1996-08-26
180 0:1 1:1 2:1 3:1 4:12 6:1 7:2 8:2 11:12 13:8
7 INDIA: Mother Teresa's condition improves, many pray. CALCUTTA, India 1996-08-25
237 0:1 2:4 6:2 7:1 8:1 9:2 10:1 12:11 15:1 18
9 UK: Charles under fire over prospect of Queen Camilla. LONDON 1996-08-26
195 0:1 2:1 4:1 12:5 15:1 18:5 19:1 21:3 22:2 10
10 UK: Britain tells Charles to forget Camilla. LONDON 1996-08-27
194 0:2 2:3 3:1 5:1 6:1 7:3 9:1 12:17 15:4 18:11
11 COTE D'IVOIRE: FEATURE - Quiet homecoming for reprieved Ivory Coast maid. ABIDJAN 1996-08-28
165 0:1 3:1 5:1 7:3 12:5 15:3 19:1 20:1 21:2 213
12 INDIA: Mother Teresa (I want to go home) sits and prays. CALCUTTA 1996-08-28
134 0:4 5:1 6:2 9:1 15:1 18:1 19:1 23:3 26:1 314
13 INDIA: Mother Teresa nears end of crisis, nuns rejoice. CALCUTTA 1996-08-28
193 0:3 1:1 2:1 3:1 6:3 8:1 9:2 10:1 13:2 14:215
14 UK: Prosaic end for marriage of Charles and Diana. LONDON 1996-08-28
177 0:4 2:2 5:1 6:3 8:1 9:3 13:1 14:1 15:2 16:16
15 UK: No respite for British royals despite divorce. LONDON 1996-08-28
180 2:2 3:1 8:1 14:1 17:6 19:1 34:1 36:3 41:117
16 UK: Camilla, love of Charles' life, an unlikely queen. LONDON 1996-08-28
113 0:1 3:1 5:1 6:1 9:1 15:1 30:1 36:1 37:2 4219
17 UK: Diana sets out on new life as single woman. LONDON 1996-08-28
93 0:3 4:1 5:1 7:4 9:1 14:1 15:1 19:1 20:1 24:20
18 USA: O.J. Simpson attacks media, hints at lawsuits. WASHINGTON 1996-08-28
166 0:2 1:11 3:2 5:2 6:2 7:2 10:1 14:5 15:1 1821
19 USA: U.S. Cardinal Bernardin has one year or less to live. CHICAGO 1996-08-30
20 USA: U.S. Cardinal Bernardin says has terminal cancer. CHICAGO 1996-08-30
21 ROMANIA: German architect wins Bucharest rebuilding prize. BUCHAREST 1996-09-02
22 ARGENTINA: Argentina's "Blond Angel" finally quits Navy. BUENOS AIRES, Argentina 1996-09-02
23 UK: Disney lights up Pocahontas resting place. GRAVESEND, England 1996-09-06
24 HUNGARY: POPE LEAVES HUNGARY AFTER DEMANDING TWO-DAY VISIT. BUDAPEST 1996-09-07
25 HUNGARY: Pope says mass in Hungary, health in spotlight. GYOR, Hungary 1996-09-07
26 UK: Prince Charles' love will not wed him, paper says. LONDON 1996-09-09

church
pope
years
people
mother
last
told
first
world
year
president
teresa
charles
catholic
during
life
u. s
city
public
time
since
family
king
former
british
harriman
against
country
vatican
made
three
hospital

LDA

```
C:\Python27\python.exe D:/Python/16.3.reuters.py
```

```
type(X): <type 'numpy.ndarray'>
```

```
shape: (395, 4258)
```

```
[[ 1  0  1  0  0  0  1  0  0  1]
 [ 7  0  2  0  0  0  0  1  0  0]
 [ 0  0  0  1 10  0  4  1  1  0]
 [ 6  0  1  0  0  0  1  1  1  0]
 [ 0  0  0  2 14  1  1  0  2  1]
 [ 0  0  2  2 24  0  2  0  2  1]
 [ 0  0  0  2  7  1  1  0  1  0]
 [ 0  0  2  2 20  0  2  0  3  1]
 [ 0  1  0  2 17  2  2  0  0  0]
 [ 2  0  2  0  0  2  0  1  0  3]]
```

```
type(vocab): <type 'tuple'>
```

```
len(vocab): 4258
```

```
('church', 'pope', 'years', 'people', 'mother', 'last', 'told', 'first', 'world', 'year')
```

```
type(titles): <type 'tuple'>
```

```
len(titles): 395
```

```
('0 UK: Prince Charles spearheads British royal revolution. LONDON 1996-08-20', '1 GERMANY:
```

```
LDA start ----
```

```
INFO:lda:n_documents: 395
```

```
INFO:lda:vocab_size: 4258
```

```
INFO:lda:n_words: 84010
```

```
INFO:lda:n_topics: 20
```

```
INFO:lda:n_iter: 500
```

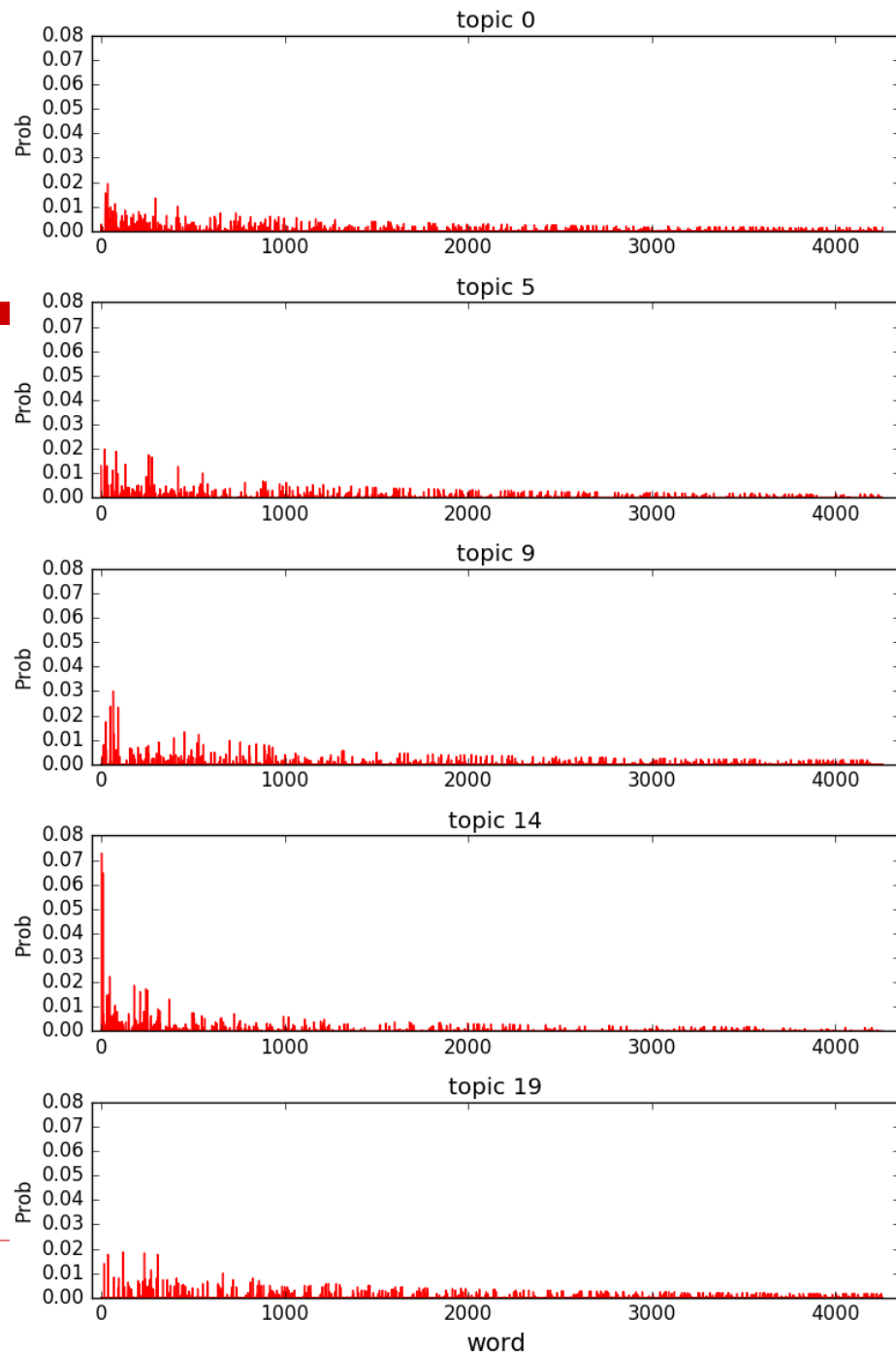
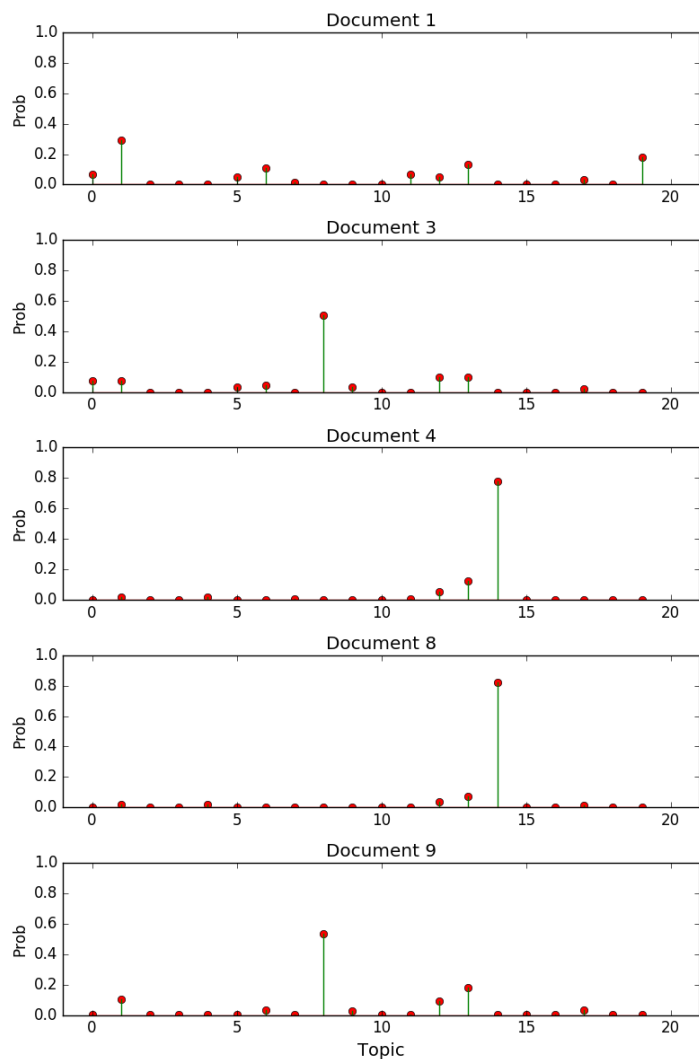
```
INFO:lda:<0> log likelihood: -1051748
```

```
INFO:lda:<10> log likelihood: -719800
```

```
INFO:lda:<20> log likelihood: -699115
```

```
INFO:lda:<30> log likelihood: -689370
```


主题和主题分布

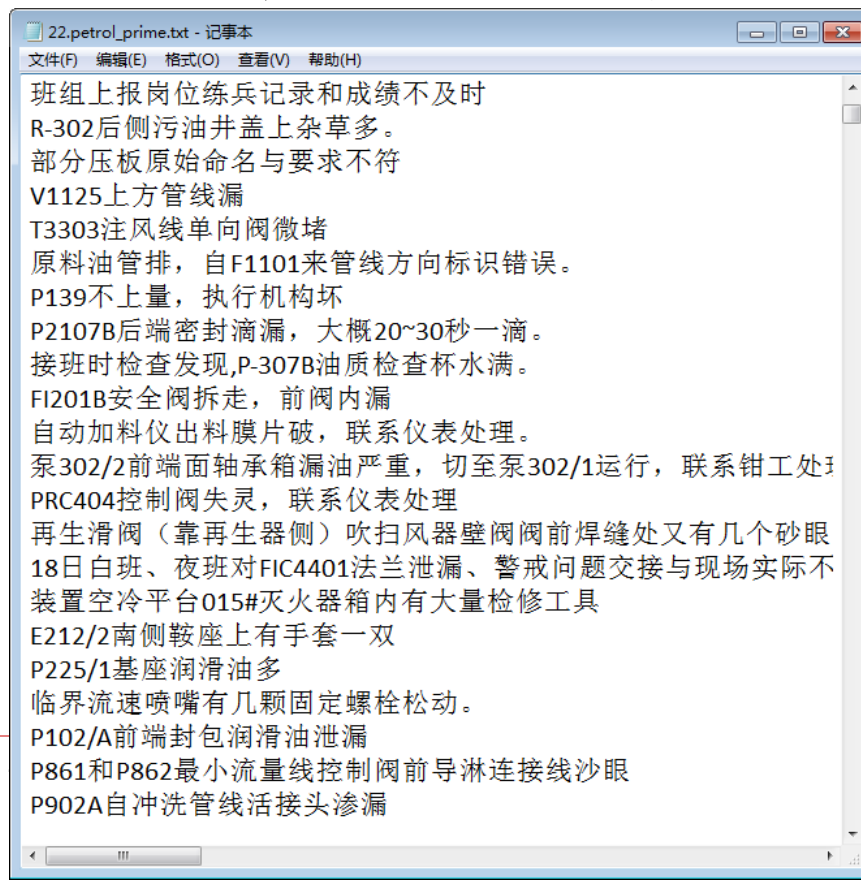


石油例检结果处理

□ 针对国内某石油企业的例行检查处理结果，
试通过主题模型方案，分析例检结果中最突出的问题是什么？

■ 文本共4700个，

■ 单个文档十数字



聚类 “主主题”

22.4.petrol 22.4.petrol

每个主题的词分布:

主题#0: 溶脱 地面 下 发现 溜 吹 漏洞

概率: [0.03567895 0.0305515 0.02995125 0.02905559 0.02847266 0.02672661 0.026

主题#1: 卫生 清扫 新增 本月 缺陷 无 空调

概率: [0.09127588 0.04710893 0.03864091 0.0385492 0.03507678 0.03243568 0.021

主题#2: 过滤器 设备 高 差压 少 入口 17

概率: [0.07382152 0.04735673 0.03770307 0.03342033 0.03160451 0.02609018 0.023

主题#3: 号牌 位 脱落 铁皮 北侧 塔 规范

概率: [0.06217475 0.06146815 0.04165677 0.03554779 0.03128229 0.02976478 0.024

主题#4: 松 建议 液位 损坏 皮带 区域 断裂

概率: [0.03845036 0.03260667 0.03243315 0.031336 0.03074858 0.03062344 0.028

主题#5: 盘根 漏 阀 出口 采样 引出 内

概率: [0.09527604 0.08630304 0.07457675 0.0508139 0.04410229 0.03406847 0.033

主题#6: 日 月 8 红线 压力表 技术员 23

概率: [0.04320296 0.04268257 0.04081915 0.03381051 0.02013684 0.01981086 0.016

主题#7: 地沟 错误 单向阀 螺丝 装车 旁 年

概率: [0.03169998 0.02547323 0.02350422 0.02196816 0.02146054 0.02122409 0.019

主题#8: 皮带 断 不准 松动 一次 盖 被

概率: [0.15231247 0.10833394 0.05339595 0.04884857 0.04585281 0.03503216 0.034

主题#9: 电机 杂音 接头 大盖 冲洗 有 活

概率: [0.07620952 0.0712979 0.06082735 0.03423004 0.03310958 0.02874415 0.027

主题#10: 蒸汽 砂眼 管线 法兰 漏 前有

概率: [0.04160303 0.03960238 0.03123862 0.0293527 0.02685108 0.02634925 0.024

主题#11: 泄漏 伴热 量 牌 入口 底 法兰

概率: [0.07545736 0.05467109 0.05165556 0.03689506 0.03672955 0.03457064 0.031

主题#12: 密封 平台 保温 脱开 缺失 堵头 汽提

概率: [0.04631792 0.04157294 0.0396999 0.03857663 0.03751675 0.03526295 0.034

主题#13: 后端 长明灯 无法 号 灭火器 器 油站

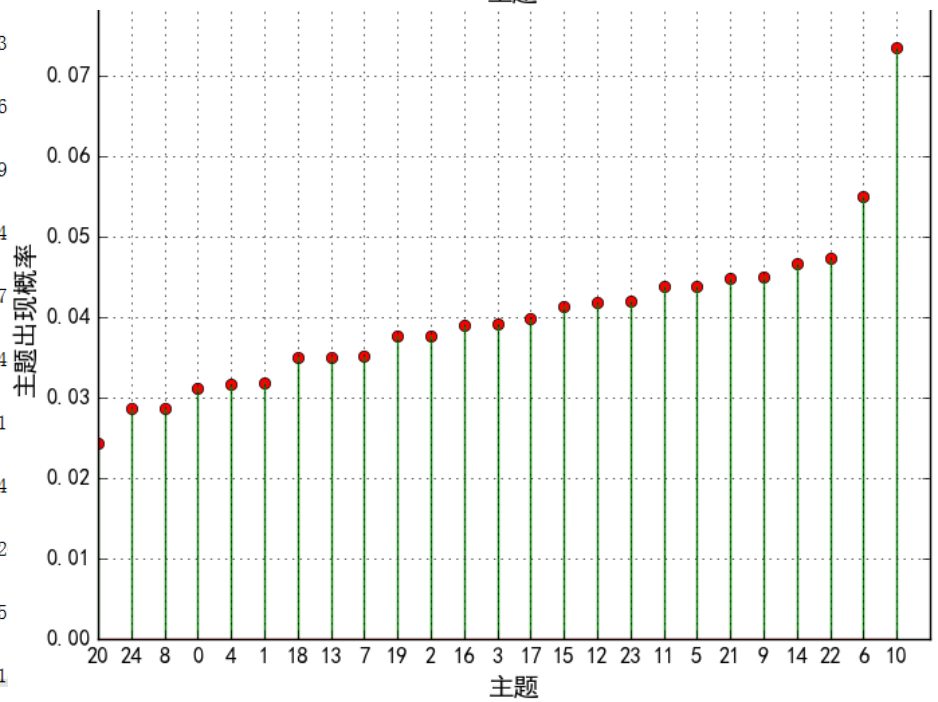
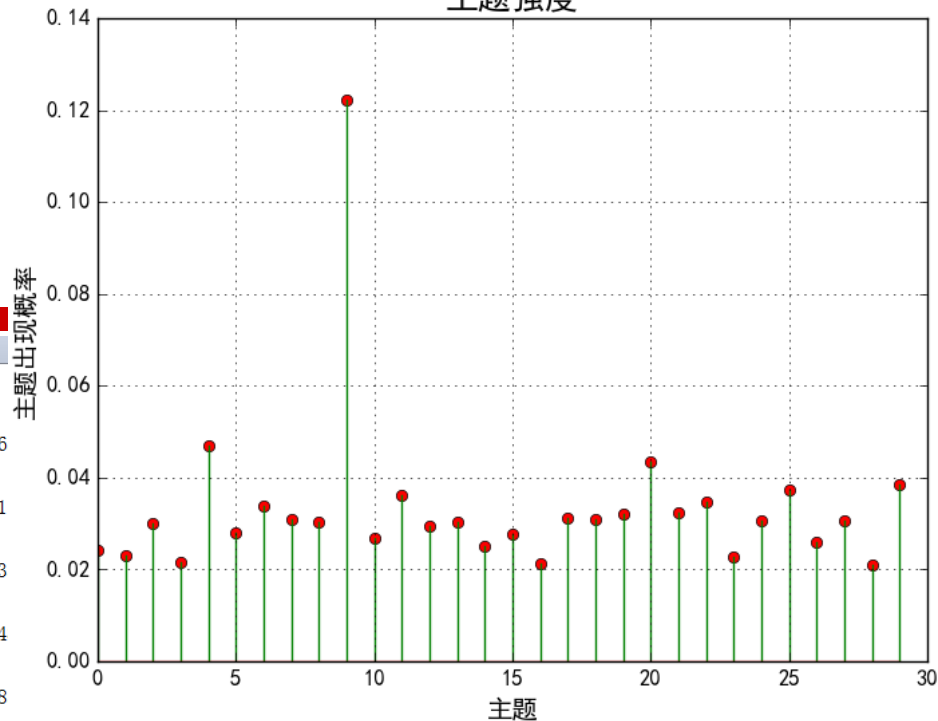
概率: [0.04065264 0.03001566 0.02572 0.02520191 0.02510101 0.02370398 0.022

主题#14: 坏 炉 压力表 润滑油 泵 安全阀 清理

概率: [0.04789169 0.03913219 0.03672254 0.03227133 0.03054908 0.02618981 0.025

主题#15: 缺陷 新增 无 本月 交接班 胶带机 日志

概率: [0.07364203 0.07360244 0.06751279 0.06417833 0.03489421 0.02254108 0.021



LDA总结

- 由于在词和文档之间加入的**主题**的概念，可以较好的解决**一词多义**和**多词一义**的问题。
- 在实践中发现，LDA用于**短文档**往往效果不明显——这是可以解释的：因为一个词被分配给某个主题的次数和一个主题包括的词数目尚未收敛。往往需要通过其他方案“**连接**”成长文档。
 - 用户评论/Twitter/微博
- LDA可以和其他算法**相结合**。首先使用LDA将长度为 N_i 的文档**降维**到 K 维(主题的数目)，同时给出每个主题的概率(主题分布)，从而可以使用if-idf继续分析或者直接作为文档的特征进入**聚类**或者**标签传播算法**——用于**社区发现**等问题。

参考文献

- ❑ David M. Blei, Andrew Y. Ng, Michael I. Jordan, *Latent Dirichlet Allocation*, 2003
- ❑ Gregor Heinrich, *Parameter estimation for text analysis*. 2008
- ❑ Matthew D. Hoffman, David M. Blei, Francis Bach. *Online learning for Latent Dirichlet Allocation*. 2010
- ❑ Mihalcea R, Tarau P. TextRank, *Bringing order into texts*. Association for Computational Linguistics, 2004.
- ❑ http://en.wikipedia.org/wiki/Dirichlet_distribution
- ❑ http://en.wikipedia.org/wiki/Conjugate_prior
- ❑ <https://en.wikipedia.org/wiki/Perplexity>
- ❑ <https://github.com/letiantian/TextRank4ZH>

感谢大家！

恳请大家批评指正！