



第5章：文本聚类

课件制作：王新宇、吴昊一、屠可伟
讲解人：王新宇

文本聚类

• 应用

• 新闻聚合网站

探索 DISCOVERY ▾

- **追光 | 来热带雨林研学，每天至少2万步**
 - 今日好稿 | 快看，藏在博物馆里的“龙”
 - 蛇、巨蜥、梅花鹿、白牦牛.....住在动物园的它们身上藏..
 - 一定要试试！8个在家轻松完成的科学小实验
 - 游天府看展览闹民俗 金沙太阳节春节假期接待观众较..
 - 普及天文知识，在孩子们心田播下“科学”的种子
- **肉眼可见！15日晚有木星伴月天象现身夜空**
 - 欢欢喜喜过大年 | 龙年说龙：来瞧瞧自博馆里的“中国龙..
 - 河北昌黎：15只黑嘴鸥现身七里海潟湖湿地
 - 美！玉兰花迎春惊艳绽放
 - 金钱豹频繁现身黄柏塬国家级自然保护区
 - 联播观察 | 龙年寻龙 成都“巨龙”藏在这！

文本聚类

- 应用
 - 客户评价分类



文本聚类：定义

- 输入：
 - 一系列文档 $\{d_1, d_2, \dots, d_n\}$
- 输出：
 - 分配后的簇
 - $C_1 = \{d_1, d_3, \dots\}$
 - $C_2 = \{d_2, d_6, \dots\}$
 - $C_3 = \{d_4, \dots\}$
 -

常用方法

- 用特征向量表示文本（第4章）
- 应用任意聚类算法
 - k 均值
 - 层次凝聚聚类
 - 基于高斯混合的最大期望值法
 -

需要测量向量之间的距离，例如 L2



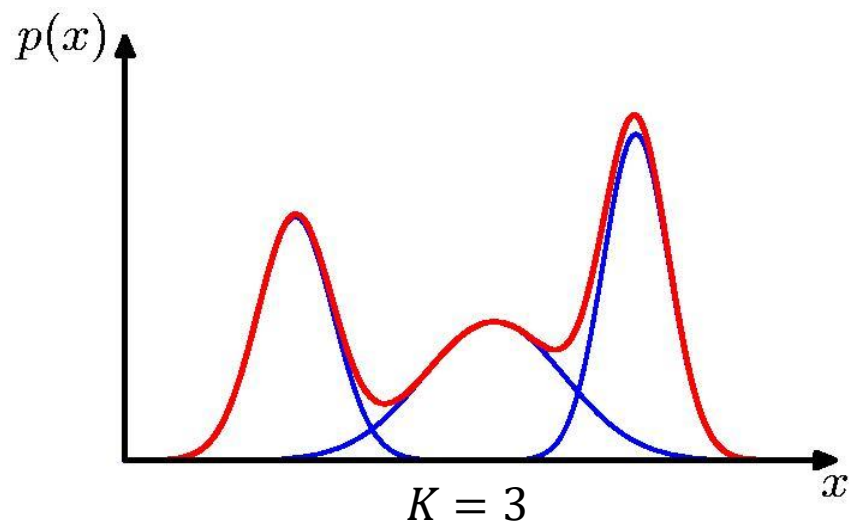
高斯混合 (MoG)

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

混合系数

高斯函数

$$\forall k : \pi_k \geq 0 \quad \sum_{k=1}^K \pi_k = 1$$

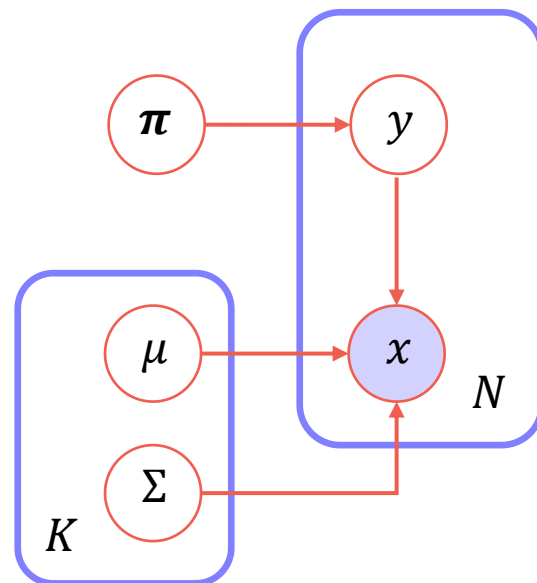


高斯混合 (MoG)

- $P(Y)$: k 个簇上的分布
- $P(X|Y)$: 每个簇从具有均值 μ_i 和协方差矩阵的 Σ_i 的混合高斯中生成数据点的分布

每个数据点都是从**生成过程**中采样的:

1. 以概率 π_i 选择 $y = i$
2. 从高斯函数 $\mathcal{N}(\mathbf{x}|\mu_i, \Sigma_i)$ 生成数据点

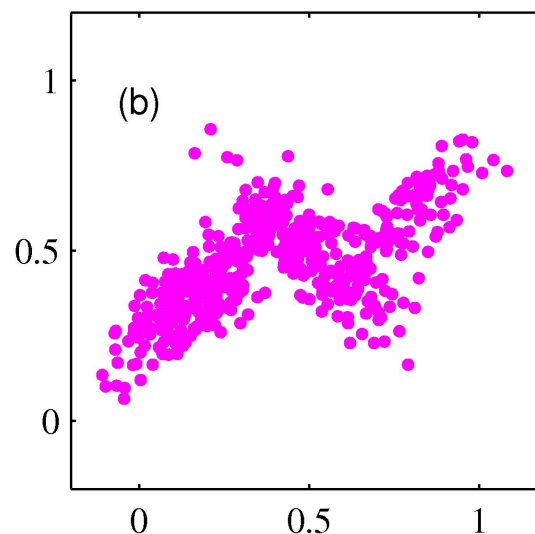


高斯混合 (MoG)

- 在聚类中，我们不知道标签 Y 的分布
- 最大化边际似然：

$$\prod_j P(\mathbf{x}_j) = \prod_j \sum_i P(y_j = i, \mathbf{x}_j) = \prod_j \sum_i \pi_i N(\mathbf{x}_j | \mu_i, \Sigma_i)$$

- 我们如何优化它？



期望最大值法 (EM)

- 选择 K 个随机聚类模型（高斯）
- 交替执行：
 - [E步骤] 将数据点按比例分配给不同的聚类模型
 - [M步骤] 根据（按比例）分配的点修改每个聚类模型
- 当边际概率没有显著变化时停止

- EM = Expectation-Maximization 通过坐标上升最大化边际似然

E步骤

- [E步骤] 将数据点按**比例分配**给不同的模型

- 计算每个数据点的标签分布

$$P(y_j = i | \mathbf{x}_j, \theta^{(t)}) \propto \pi_i^{(t)} N(\mathbf{x}_j | \mu_i^{(t)}, \Sigma_i^{(t)})$$

只需评价高斯分布

\mathbf{x}_j

M步骤

- [E步骤] 将数据实例按比例分配给不同的模型

- 计算每个数据点的标签分布

$$P(y_j = i | \mathbf{x}_j, \theta^{(t)}) \propto \pi_i^{(t)} N(\mathbf{x}_j | \mu_i^{(t)}, \Sigma_i^{(t)})$$

- [M 步骤] 根据（按比例）分配的点修改每个聚类模型

- 计算给定标签分布的参数的加权最大似然估计

$$\mu_i^{(t+1)} = \frac{\sum_j P(y_j = i | \mathbf{x}_j, \theta^{(t)}) \mathbf{x}_j}{\sum_{j'} P(y_{j'} = i | \mathbf{x}_{j'}, \theta^{(t)})}$$

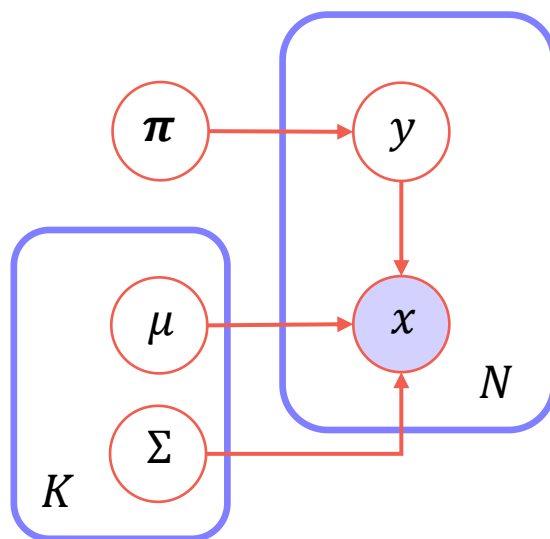
$$\pi_i^{(t+1)} = \frac{\sum_j P(y_j = i | \mathbf{x}_j, \theta^{(t)})}{m}$$

$$\Sigma_i^{(t+1)} = \frac{\sum_j P(y_j = i | \mathbf{x}_j, \theta^{(t)}) [\mathbf{x}_j - \mu_i^{(t+1)}][\mathbf{x}_j - \mu_i^{(t+1)}]^T}{\sum_{j'} P(y_{j'} = i | \mathbf{x}_{j'}, \theta^{(t)})}$$

$m = \#$ 训练样本

生成模型

- MoG生成具有高斯分布
- 我们可以直接生成文档（单词序列）吗？
 - 可以使用离散分布



生成模型

- MoG是用K个高斯分布中的一个来生成一个文档的特征向量
- 我们可以直接生成文档（单词序列）吗？
 - 可以使用离散分布

每个主题的单词分布

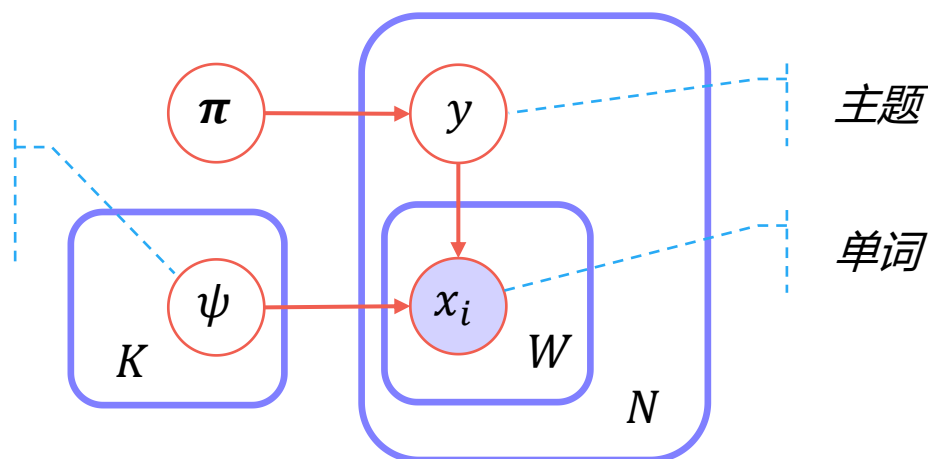
例如, 单词- [如何、怎么、得到.....]

体育: [0.1, 0.3, 0.2, ...]

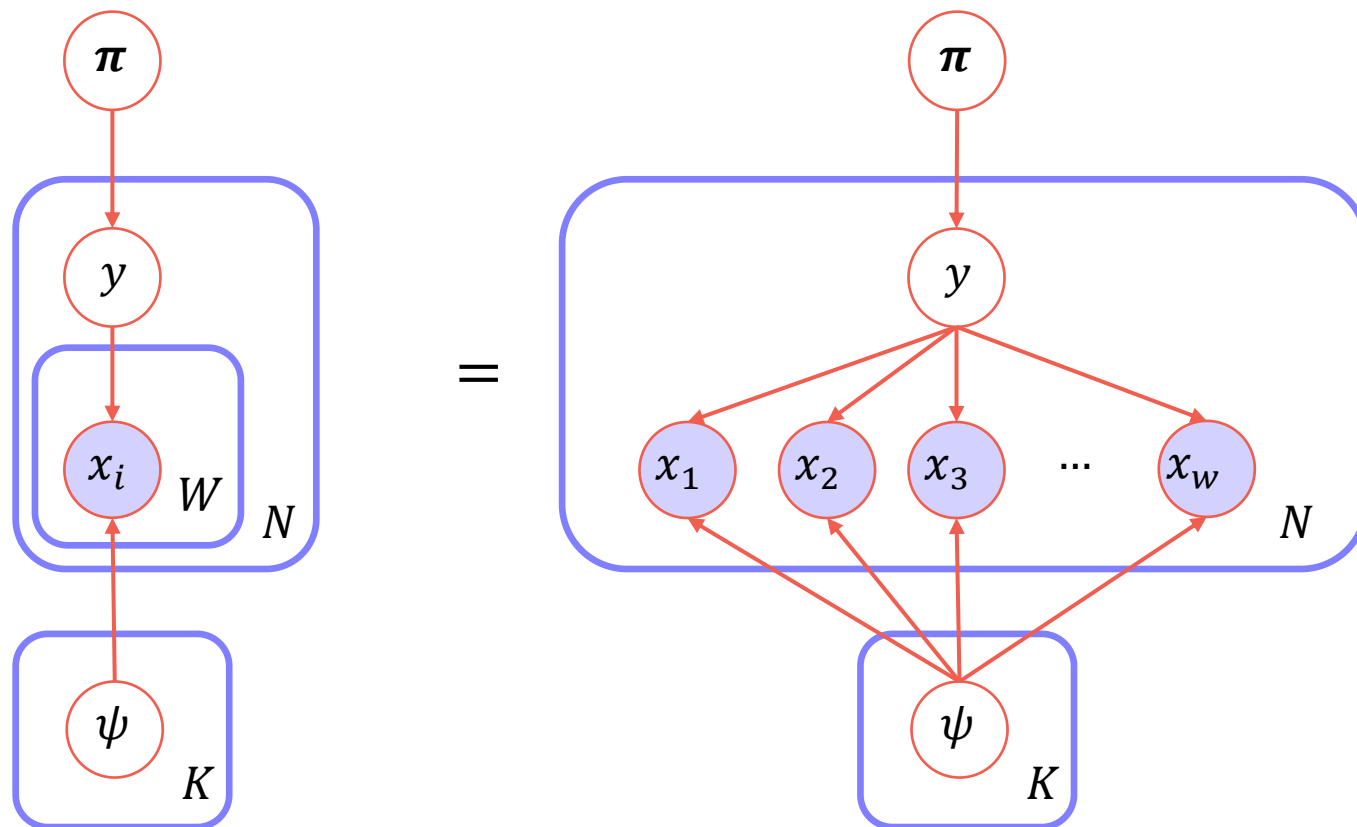
艺术: [0.2, 0.3, 0.1, ...]

科学: [0.3, 0.4, 0.1, ...]

.....



生成模型



这正是朴素贝叶斯模型!

无监督朴素贝叶斯

- 我们可以运行 EM 来进行朴素贝叶斯的无监督学习
 - 即基于单词而不是特征的文本聚类
- [E 步骤] 根据不同的主题按比例分配文档
 - 计算每个文档的主题分布

$$P(y_j = i | x_{j,1:w}, \theta^{(t)}) \propto \pi_i^{(t)} \prod_{k=1}^w P(x_{j,k} | \psi_i^{(t)})$$

无监督朴素贝叶斯

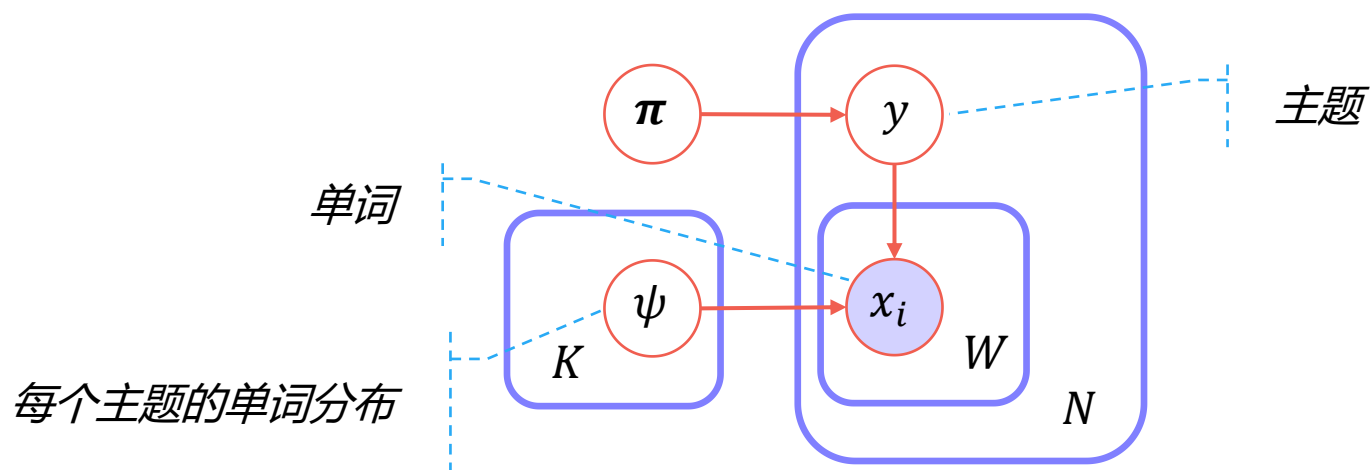
- 我们可以运行 EM 来进行朴素贝叶斯的无监督学习
 - 即基于单词而不是特征的文本聚类
- [M 步骤] 根据（按比例）指定的单词修改每个主题
 - 计算给定主题分布的参数的加权最大似然估计
 - 表示 $\psi_i = \{p_{i,1}, p_{i,2}, \dots, p_{i,v}\}$

$$p_{i,l}^{(t+1)} = \frac{\sum_j P(y_j = i | x_{j,1:w}, \theta^{(t)}) \sum_k \mathbf{1}(x_{j,k} = l)}{\sum_j P(y_j = i | x_{j,1:w}, \theta^{(t)}) \cdot w_j} \quad \pi_i^{(t+1)} = \frac{\sum_j P(y_j = i | x_{j,1:w}, \theta^{(t)})}{m}$$

其中 v 是词汇量， m 是训练文档的数量。

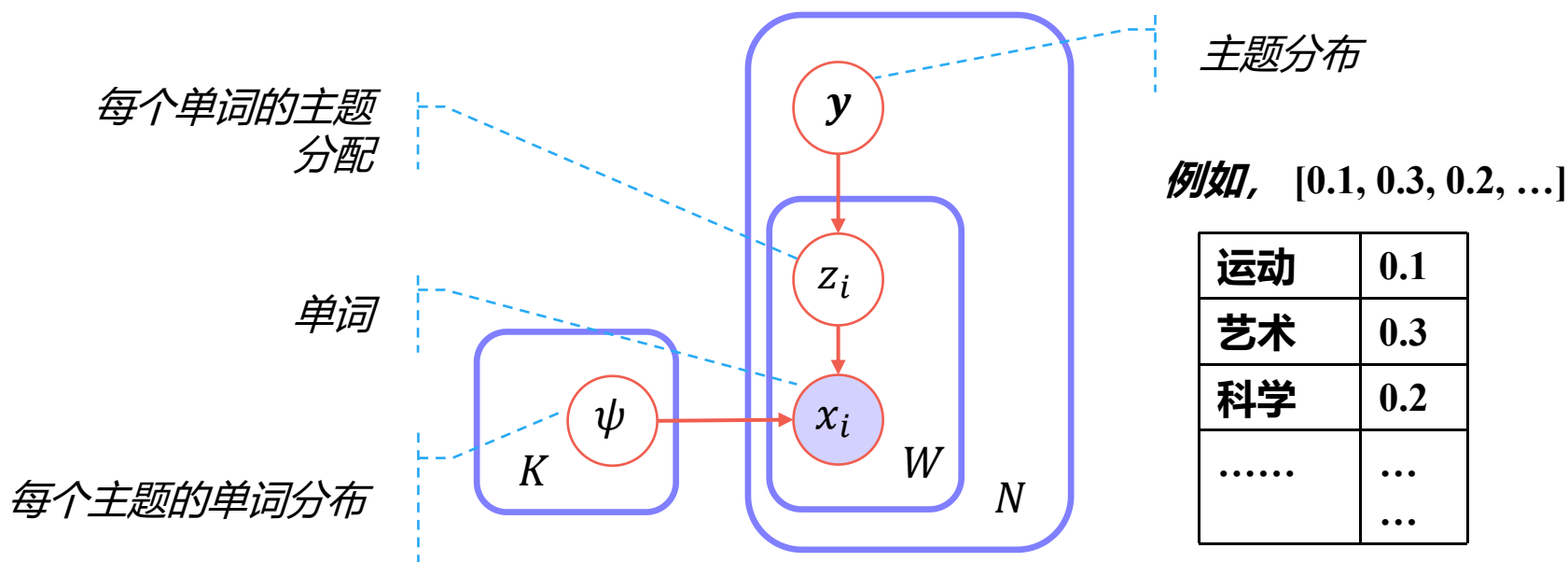
主题建模

- 文本簇可能对应不同的主题
- 到目前为止，我们假设每个文档都有一个簇标签
- 但是，一篇文章可能涵盖多个主题
 - 如何进行学习？



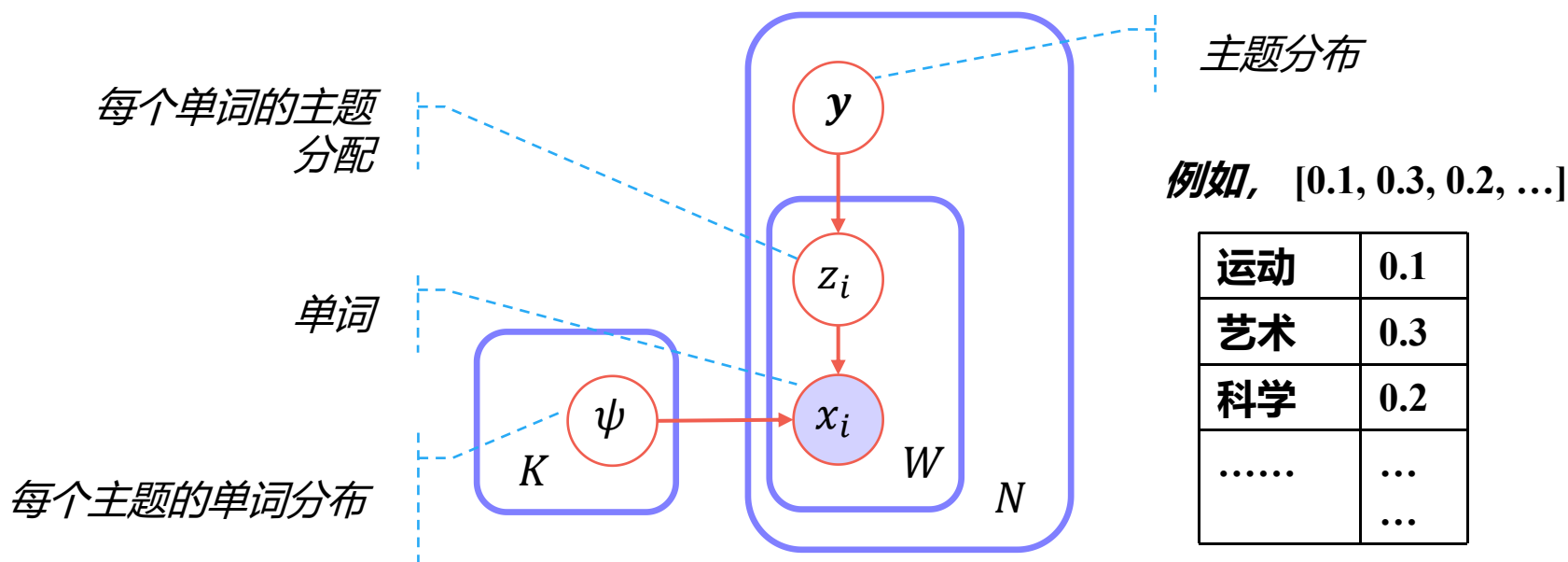
主题建模

- 文本簇可能对应不同的主题
- 到目前为止，我们假设每个文档都有一个簇标签
- 但是，一篇文章可能涵盖多个主题
 - 如何进行学习？



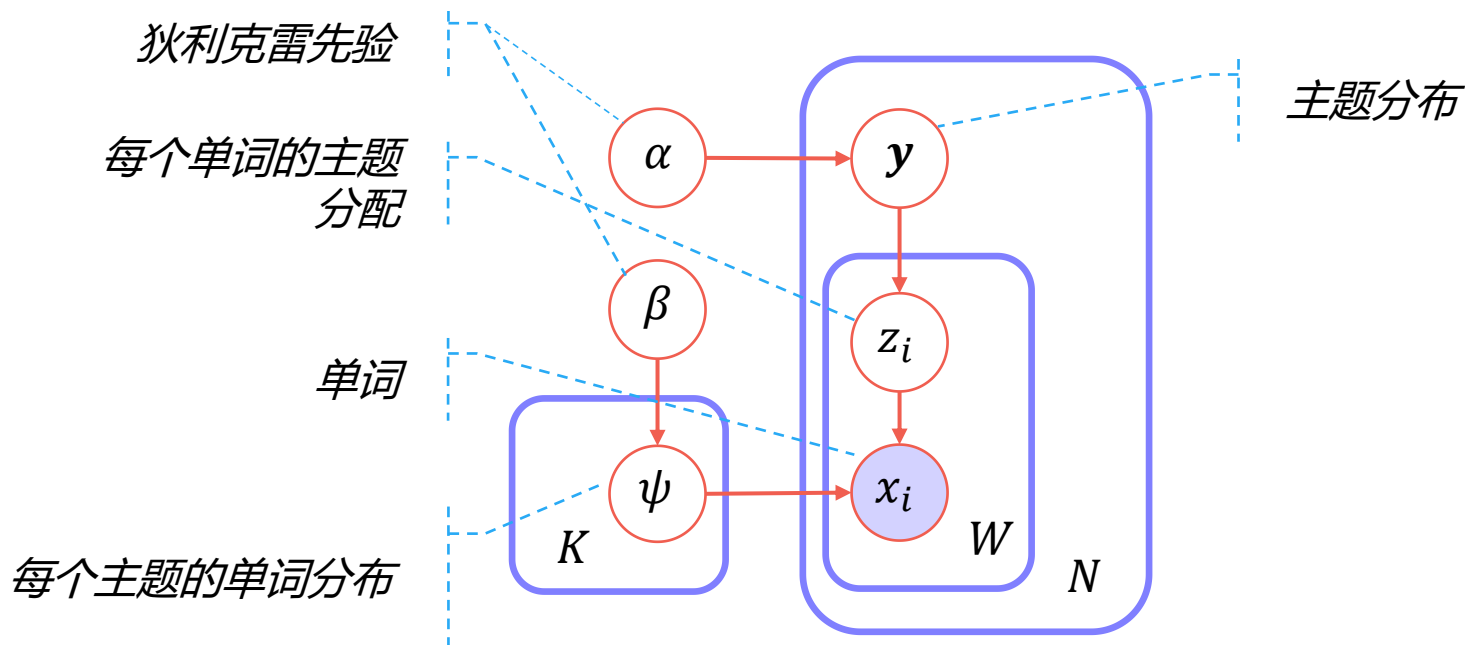
主题建模

- 这称为概率潜在语义分析 (pLSA)
 - 同样，我们可以运行 EM 来学习它
- 我们可以进一步在主题和单词分布上添加狄利克雷先验

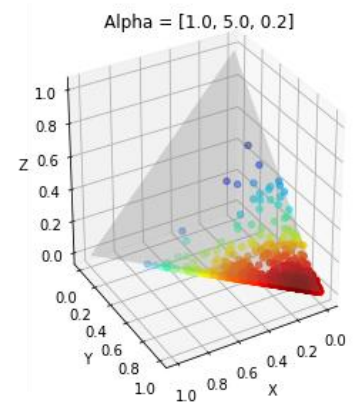
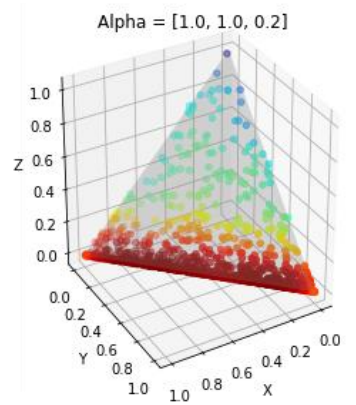
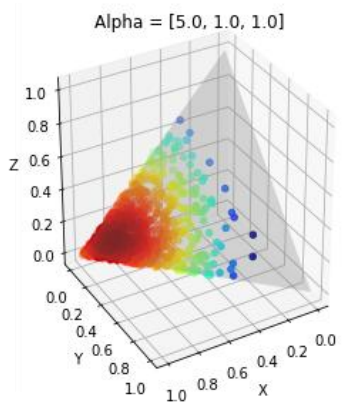
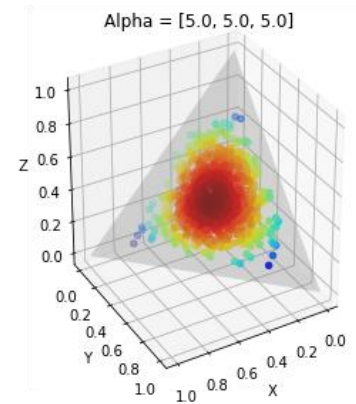
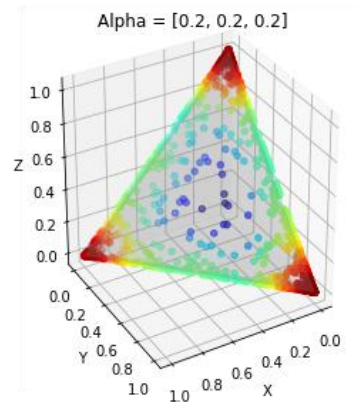
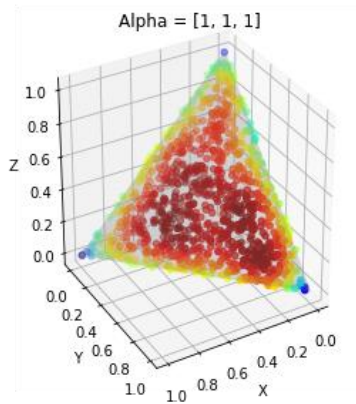


主题建模

- 狄利克雷先验：鼓励主题和单词分布稀疏
 - 文档应仅涵盖几个主题
 - 在一个主题中，只有一部分单词具有高频率

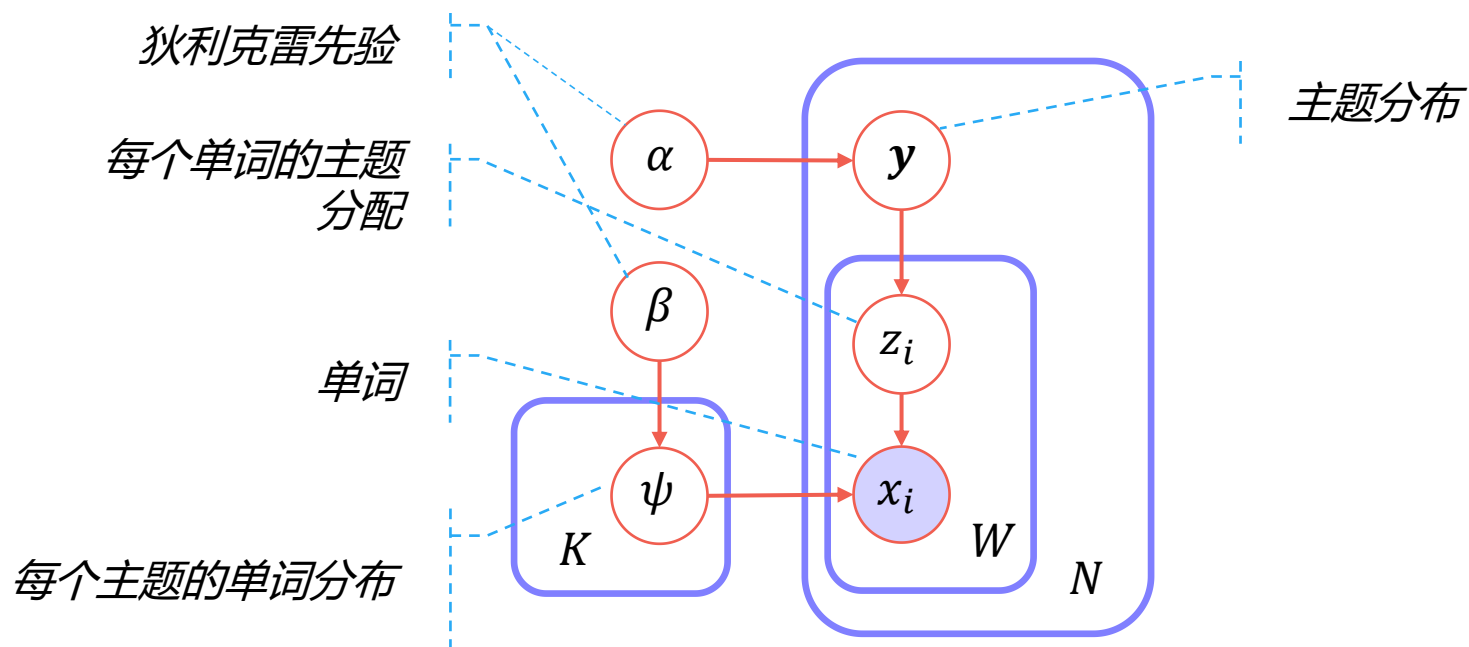


狄利克雷分布



主题建模

- 这称为潜在狄利克雷分配 (LDA)
 - 学习: 变分推理或MCMC



通用EM算法

- 可用于学习任何具有隐变量（缺失数据）的模型
- 交替计算：
 - 根据当前参数值计算隐变量的分布
 - 根据隐变量的分布计算新参数值以最大化预期对数似然
- 当没有变化时停止

- 可以达到局部最优，但不一定达到全局最优

文本聚类

- 高斯混合
- 无监督朴素贝叶斯
- 主题模型
- 学习
 - 期望最大值算法