

吉布斯采样和概率图模型介绍

钱烽

qf6101 at gmail.com

References

- Resnik, P., Resnik, P., Hardisty, E., & Hardisty, E. (2009). Gibbs Sampling for the Uninitiated. *Umiacs.Umd.Edu*, (June), 1–23.
- Heinrich, G. (2008). Parameter Estimation for Text Analysis.
- Advanced MCMC Methods (<http://mlg.eng.cam.ac.uk/zoubin/tutorials06.html>)
- Wikipedia

Part I

MCMC and Gibbs Sampling

Why use MCMC or Gibbs sampling?

- To approximate the value of an integral



Why integral?

- 离散场景并不需要 (点估计)

- Maximum likelihood estimate (MLE)

$$\tilde{\pi}_{MLE} = \operatorname{argmax}_{\pi} P(\mathcal{X}|\pi)$$

$$P(y|\mathcal{X}) \approx P(y|\tilde{\pi}_{MLE})$$

- Maximum a posteriori (MAP)

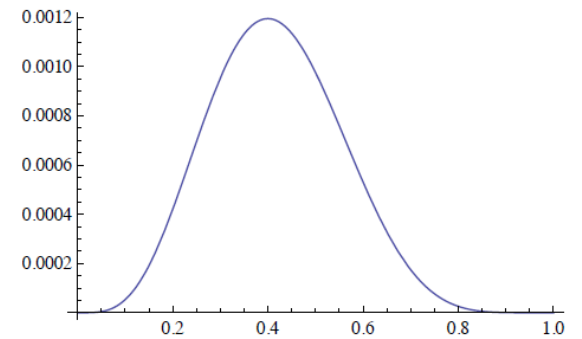
$$\tilde{\pi}_{MAP} = \operatorname{argmax}_{\pi} P(\pi|\mathcal{X})$$

$$= \operatorname{argmax}_{\pi} \frac{P(\mathcal{X}|\pi)P(\pi)}{P(\mathcal{X})}$$

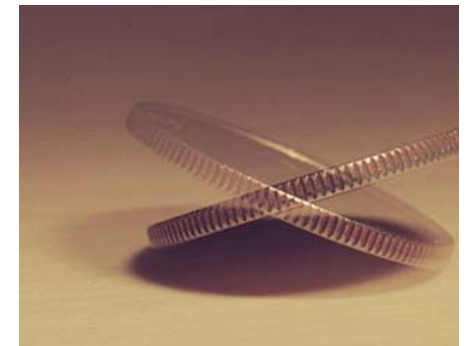
$$= \operatorname{argmax}_{\pi} P(\mathcal{X}|\pi)P(\pi)$$

$$P(y|\mathcal{X}) \approx P(y|\tilde{\pi}_{MAP})$$

$(\pi, P(X|\pi))$



HHHHTTTTTT



Why integral?

- 🤔: 真的不需要吗 (考虑整体分布, 例如求期望)?

$$E[f(z)] = \int f(z)p(z) dz.$$

$$P(y|\mathcal{X}) = \int P(y|\pi)P(\pi|\mathcal{X}) d\pi$$

$$P(\pi|\mathcal{X}) = \frac{P(\mathcal{X}|\pi)P(\pi)}{P(\mathcal{X})} = \frac{P(\mathcal{X}|\pi)P(\pi)}{\int_{\pi} P(\mathcal{X}|\pi)P(\pi) d\pi}.$$

$$E_{\pi \sim P(\pi|\mathcal{X})}[P(y|\pi)]$$



Why sampling?

- $f(z)$ 不可积。。

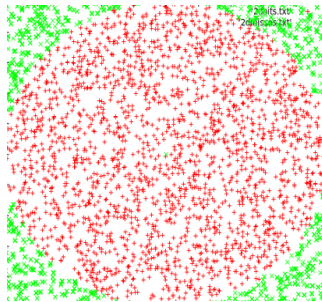


Monte Carlo Simulation

- Approximate π (Probabilistic Choice)

$$\frac{C}{S} \approx \frac{\pi(\frac{d}{2})^2}{d^2}$$

$$\pi \approx \frac{4C}{S}$$



撒点即采样

- Properties of Monte Carlo

- Estimator is unbiased
- Variance shrinks $\propto 1/N$
(N 是撒点数)



"Monte Carlo is an extremely bad method; it should be used only when all alternative methods are worse."

— Alan Sokal, 1996

Markov Chain Monte Carlo (MCMC)

- 回顾目标

$$E[f(z)] = \int f(z)p(z) dz.$$

- 基于 $p(z)$ 采样 $z \rightarrow$ 代入 $f(z) \rightarrow$ 求均值

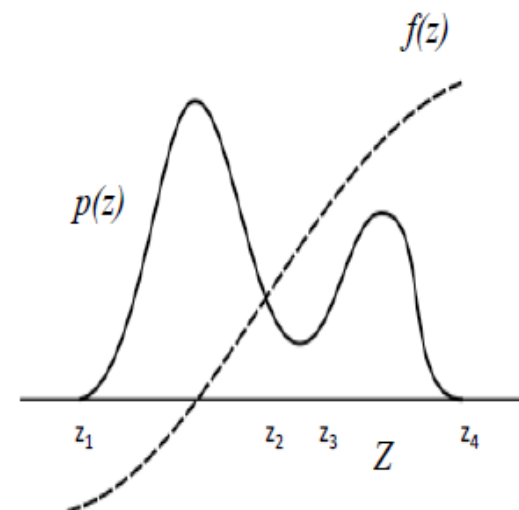
sample N points $z^{(0)}, z^{(1)}, z^{(2)}, \dots, z^{(N)}$ at random from the probability density $p(z)$.

$$E_{p(z)}[f(z)] = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=1}^N f(z^{(t)})$$

- Walking the right walk (找一个合适的 $g(z)$)

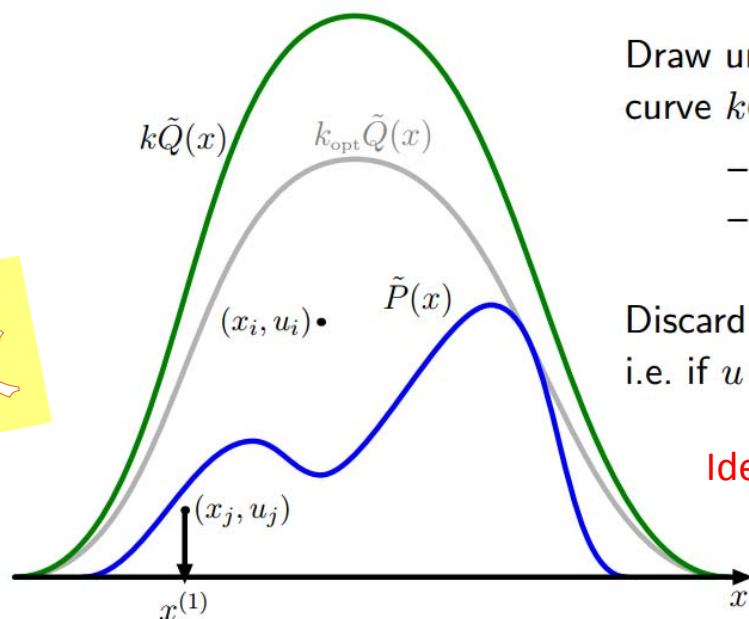
- 1: $z^{(0)} :=$ a random initial point
- 2: for $t = 1$ to T do
- 3: $z^{(t+1)} := g(z^{(t)})$
- 4: end for

$$P_{\text{trans}}(z^{(t+1)} | z^{(0)}, z^{(1)}, \dots, z^{(t)}) = P_{\text{trans}}(z^{(t+1)} | z^{(t)})$$



Rejection Sampling and Importance Sampling

初级



Draw underneath a simple curve $k\tilde{Q}(x) \geq \tilde{P}(x)$:

- Draw $x \sim Q(x)$
- height $u \sim \text{Uniform}[0, k\tilde{Q}(x)]$

Discard the point if above \tilde{P} , i.e. if $u > \tilde{P}(x)$

Idea: 找一个容易采样的分布 Q 把 P 包起来

$$\int f(x)P(x) dx = \int f(x) \frac{P(x)}{Q(x)} Q(x) dx, \quad (Q(x) > 0 \text{ if } P(x) > 0)$$

$$\approx \frac{1}{S} \sum_{s=1}^S f(x^{(s)}) \frac{P(x^{(s)})}{Q(x^{(s)})}, \quad x^{(s)} \sim Q(x)$$

Monte Carlo

Metropolis–Hastings Sampling

- 算法过程

- Propose a move from the current state $Q(x'; x)$, e.g. $\mathcal{N}(x, \sigma^2)$
- Accept with probability $\min\left(1, \frac{P(x')Q(x; x')}{P(x)Q(x'; x)}\right)$
- Otherwise next state in chain is a copy of current state

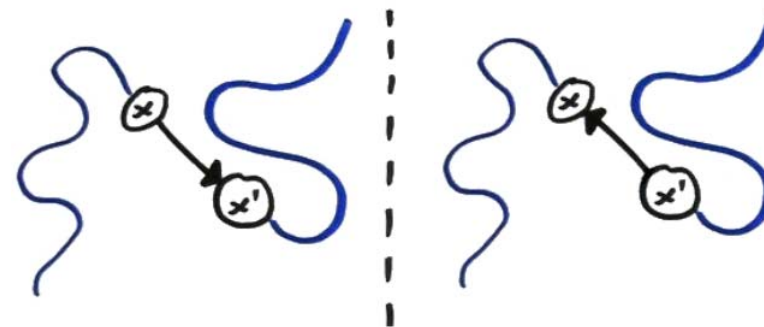
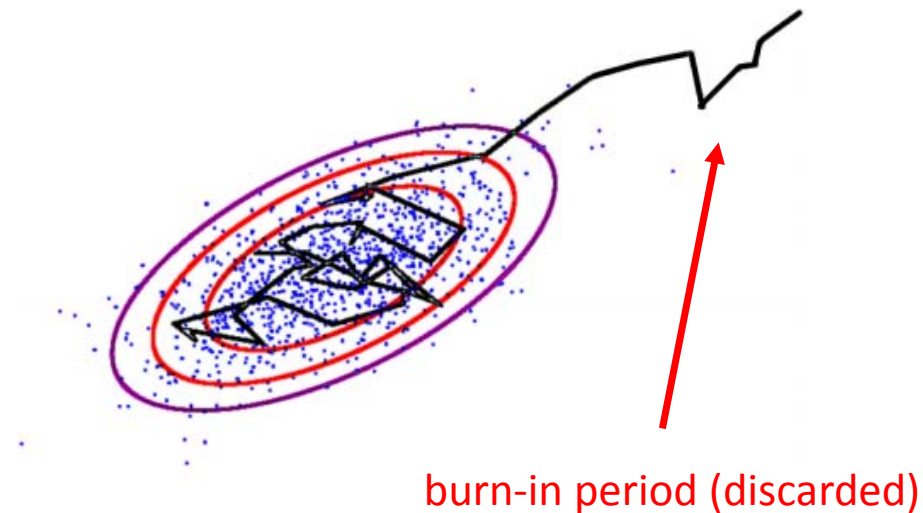
- 满足Detailed balance condition

$$\begin{aligned} P(x) \cdot T(x' \leftarrow x) &= P(x) \cdot Q(x'; x) \min\left(1, \frac{P(x')Q(x; x')}{P(x)Q(x'; x)}\right) = \min\left(P(x)Q(x'; x), P(x')Q(x; x')\right) \\ &= P(x') \cdot Q(x; x') \min\left(1, \frac{P(x)Q(x'; x)}{P(x')Q(x; x')}\right) = P(x') \cdot T(x \leftarrow x') \end{aligned}$$

中级

Detailed Balance意味着什么？

- Implies the invariant condition
 - 跳转 N 步之后，采样收敛到 $P(x)$



$$T(x' \leftarrow x)P^*(x) = T(x \leftarrow x')P^*(x')$$

还有些问题。。

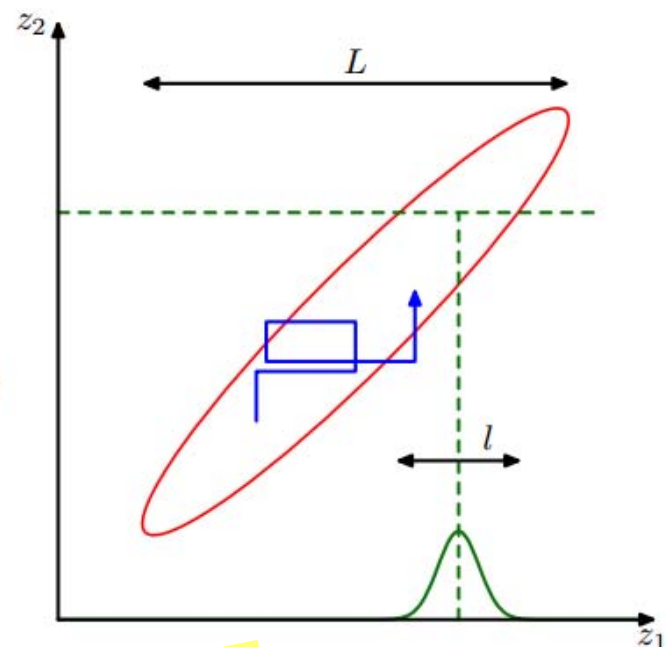
- 选择什么样的 Q 函数
- 采样高维数据很容易被拒绝



Gibbs Sampling

A method with no rejections:

- Initialize \mathbf{x} to some value
- Pick each variable in turn or randomly and resample $P(x_i | \mathbf{x}_{j \neq i})$



Proof of validity: a) check detailed balance for component update.
b) Metropolis–Hastings ‘proposals’ $P(x_i | \mathbf{x}_{j \neq i}) \Rightarrow$ accept with prob. 1
Apply a series of these operators. Don't need to check acceptance.

高级：从不拒绝

Gibbs Sampling

- 算法过程

```
1:  $z^{(0)} := \langle z_1^{(0)}, \dots, z_k^{(0)} \rangle$ 
2: for  $t = 1$  to  $T$  do
3:   for  $i = 1$  to  $k$  do
4:      $z_i^{(t+1)} \sim P(Z_i | z_1^{(t+1)}, \dots, z_{i-1}^{(t+1)}, z_{i+1}^{(t)}, \dots, z_k^{(t)})$ 
5:   end for
6: end for
```

$$P(Z_i | z_1^{(t+1)}, \dots, z_{i-1}^{(t+1)}, z_{i+1}^{(t)}, \dots, z_k^{(t)}) = \frac{P(z_1^{(t+1)}, \dots, z_{i-1}^{(t+1)}, z_i^{(t)}, z_{i+1}^{(t)}, \dots, z_k^{(t)})}{P(z_1^{(t+1)}, \dots, z_{i-1}^{(t+1)}, z_{i+1}^{(t)}, \dots, z_k^{(t)})}$$

$P(z)$

$P(z)^{(-i)}$



所以。。只需要凑出联合概率形式就阔以啦

Part II

Probabilistic Graphical Model

Notations of Text Classification Task

priori params	V	number of words in the vocabulary.
	N	number of documents in the corpus.
	$\gamma_{\pi 1}, \gamma_{\pi 0}$	hyperparameters of the Beta distribution.
	γ_{θ}	hyperparameter vector for the multinomial prior.
set of docs	$\gamma_{\theta i}$	pseudocount for word i .
	\mathbb{C}_x	set of documents labeled x .
doc = set of words	\mathbb{C}	the set of all documents.
	C_0 (C_1)	number of documents labeled 0 (1).
	\mathbf{W}_j	document j 's frequency distribution.
doc label	W_{ji}	frequency of word i in document j .
	\mathbf{L}	vector of document labels.
	L_j	label for document j .
	R_j	number of words in document j .
word probs	θ_i	probability of word i .
	$\theta_{x,i}$	probability of word i from the distribution of class x .
without j^{th} doc	$\mathcal{N}_{\mathbb{C}_x}(i)$	number of times word i occurs in the set of all documents labeled x .
	$\mathbb{C}^{(-j)}$	set of all documents <i>except</i> \mathbf{W}_j
	$\mathbf{L}^{(-j)}$	vector of all document labels <i>except</i> L_j
	$C_0^{(-j)}$ ($C_1^{(-j)}$)	number of documents labeled 0 (1) <i>except</i> for \mathbf{W}_j
	μ	set of hyperparameters $\langle \gamma_{\pi 1}, \gamma_{\pi 0}, \gamma_{\theta} \rangle$



" j " indicates doc
and
" i " indicates word

Naïve Bayes Model

- 采用MAP求解(点估计)

$$\begin{aligned} L_j = \operatorname{argmax}_L P(L|W_j) &= \operatorname{argmax}_L \frac{P(W_j|L)P(L)}{P(W_j)} \\ &= \operatorname{argmax}_L P(W_j|L)P(L), \end{aligned}$$

- 属性之间条件独立

$$P(W_j | L) = \prod_{i=1}^n P(w_i | L)$$

$$W_j = \{w_1, w_2, \dots, w_n\}$$



基于PGM的Naïve Bayes Model

- 优点
 - 贝叶斯方法：引入不确定性，物理过程更平滑
- 理解为生成过程

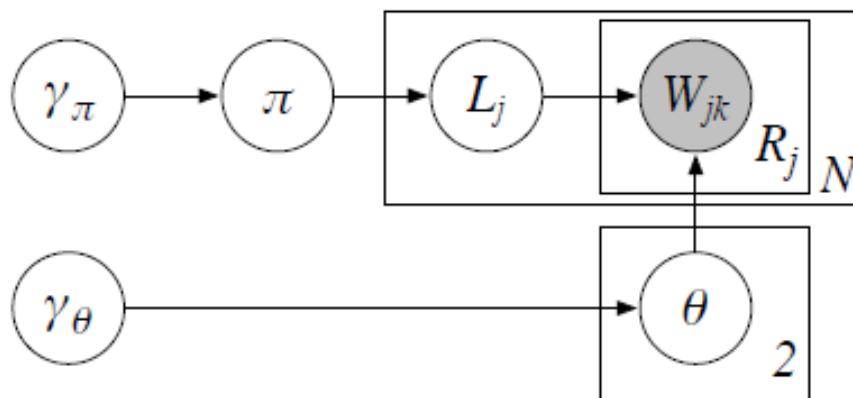
$$W_j \sim \text{Multinomial}(R_j, \theta_{L_j})$$

$$L_j \sim \text{Bernoulli}(\pi)$$

- 先验分布

$$\pi \sim \text{Beta}(\gamma_\pi)$$

$$\theta \sim \text{Dirichlet}(\gamma_\theta)$$



欢迎来到
二次元



等下。。Beta和Dirichlet是什么？



很久很久以前，在欧拉的时代。。

- 从Gamma函数说起

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

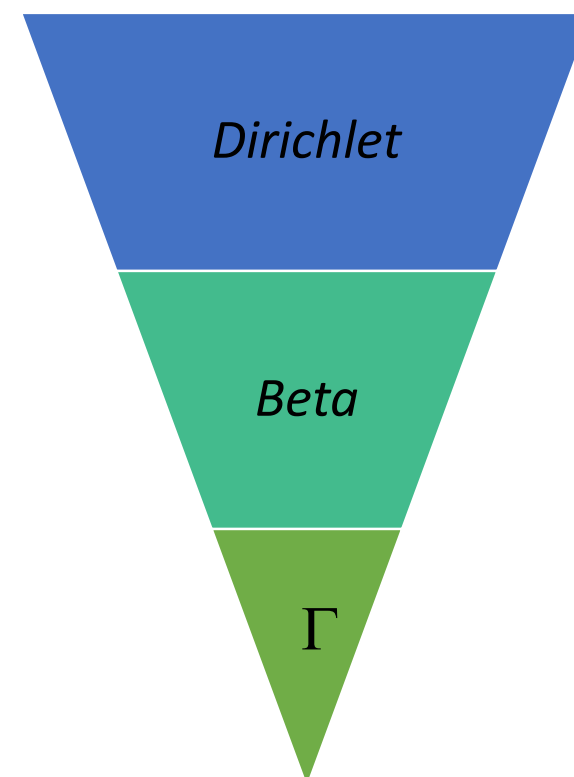
- 将阶乘扩展到实数范围

$$\Gamma(t+1) = t\Gamma(t)$$

$$\begin{aligned}\Gamma(-1) &= (-2)! & \Gamma(-\frac{3}{2}) &= \frac{4}{3}\sqrt{\pi} \\ \Gamma(0) &= (-1)! & \Gamma(-\frac{1}{2}) &= -2\sqrt{\pi} \\ \Gamma(1) &= 0! & \Gamma(\frac{1}{2}) &= \sqrt{\pi} \\ \Gamma(2) &= 1! & \Gamma(\frac{3}{2}) &= \frac{1}{2}\sqrt{\pi} \\ \Gamma(3) &= 2! & \Gamma(\frac{5}{2}) &= \frac{3}{4}\sqrt{\pi} \\ \Gamma(4) &= 3! & \Gamma(\frac{7}{2}) &= \frac{15}{8}\sqrt{\pi}\end{aligned}$$

Handwritten derivation of the Gamma function recurrence relation:

$$\begin{aligned}\Gamma(x) &= \int_0^{\infty} t^{x-1} e^{-t} dt \\ &= \int_0^{\infty} t^{x-1} \cdot (-e^{-t})' dt \\ &= -t^{x-1} e^{-t} \Big|_0^{\infty} + \int_0^{\infty} (x-1) t^{x-2} e^{-t} dt \\ &= -\frac{t^{x-1}}{e^t} \Big|_0^{\infty} + (x-1) \int_0^{\infty} t^{(x-1)-1} e^{-t} dt \\ &= (0-0) + (x-1) \Gamma(x-1) \\ \Gamma(x) &= (x-1) \Gamma(x-1) + \Gamma(1) \\ \Gamma(n) &= (n-1)! \Gamma(1)\end{aligned}$$



Gamma分布

- 被积项是密度函数

$$\int_0^{\infty} \frac{x^{\alpha-1} e^{-x}}{\Gamma(\alpha)} dx = 1$$

- Gamma分布

$$Gamma(x|\alpha) = \frac{x^{\alpha-1} e^{-x}}{\Gamma(\alpha)}$$

$\beta=1$ ↑

$$Gamma(t|\alpha, \beta) = \frac{\beta^{\alpha} t^{\alpha-1} e^{-\beta t}}{\Gamma(\alpha)}$$



嗯嗯。。分母就一阶乘

Beta分布

- 凑一个Beta分布

$$\begin{aligned}f(x; \alpha, \beta) &= \text{constant} \cdot x^{\alpha-1}(1-x)^{\beta-1} \\&= \frac{x^{\alpha-1}(1-x)^{\beta-1}}{\int_0^1 u^{\alpha-1}(1-u)^{\beta-1} du} \\&= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1} \\&= \frac{1}{B(\alpha, \beta)} x^{\alpha-1}(1-x)^{\beta-1}\end{aligned}$$

- Beta分布的期望

$$\begin{aligned}\mu = E[X] &= \int_0^1 x f(x; \alpha, \beta) dx \\&= \int_0^1 x \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)} dx \\&= \frac{\alpha}{\alpha + \beta} \\&= \frac{1}{1 + \frac{\beta}{\alpha}}\end{aligned}$$

数学题就一个字：“凑”



为什么用Beta做先验分布？

- 共轭先验(形式一样→容易凑)

- When the posterior probability distribution is of the same family as the prior probability distribution, it is said to be the conjugate prior of the posterior

$$p(\theta|x) = \frac{p(x|\theta) p(\theta)}{\int p(x|\theta') p(\theta') d\theta'}.$$

- Beta-Binomial共轭

Likelihood: Binomial $P(s, f|q = x) = \binom{s+f}{s} x^s (1-x)^f,$

Prior: Beta $P(x) = \frac{x^{\alpha-1} (1-x)^{\beta-1}}{B(\alpha, \beta)},$

Posterior: Beta $P(q = x|s, f) = \frac{P(s, f|x) P(x)}{\int P(s, f|x) P(x) dx}$

$$= \frac{\binom{s+f}{s} x^{s+\alpha-1} (1-x)^{f+\beta-1} / B(\alpha, \beta)}{\int_{y=0}^1 \left(\binom{s+f}{s} y^{s+\alpha-1} (1-y)^{f+\beta-1} / B(\alpha, \beta) \right) dy}$$

$$= \frac{x^{s+\alpha-1} (1-x)^{f+\beta-1}}{B(s+\alpha, f+\beta)},$$

By qfeng

Dirichlet的故事是一样的~

$$\text{Beta}(x|\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

$$\text{Dir}(\vec{p}|\vec{\alpha}) = \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K p_i^{\alpha_i-1} \quad \text{s.t. } \sum_{i=1}^K p_i = 1$$

$$\text{Binomial}(x|m_1, m_2) = \binom{m_1+m_2}{m_1} x^{m_1} (1-x)^{m_2}$$

$$\text{Multinomial}(\vec{p}|\vec{m}) = \binom{\sum_{i=1}^K m_i}{\vec{m}} \prod_{i=1}^K p_i^{m_i}$$

$$\text{Beta}(x|\alpha, \beta) + \text{Binomial}(x|m_1, m_2) \propto \text{Beta}(x|\alpha+m_1, \beta+m_2)$$

$$\text{Dir}(\vec{p}|\vec{\alpha}) + \text{Multinomial}(\vec{p}|\vec{m}) \propto \text{Dir}(\vec{p}|\vec{\alpha}+\vec{m})$$

在这里能找到更多的共轭先验
https://en.wikipedia.org/wiki/Conjugate_prior

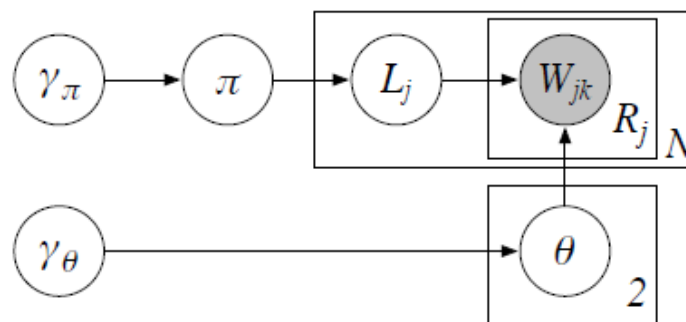
回到Naïve Bayes Model(基于PGM)

- 步骤一：凑出简单形式的联合概率

$$P(\mathbf{C}, \mathbf{L}, \pi, \theta_0, \theta_1; \gamma_{\pi 1}, \gamma_{\pi 0}, \gamma_{\theta}) = P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) P(\mathbf{L} | \pi) P(\theta_0 | \gamma_{\theta}) P(\theta_1 | \gamma_{\theta}) P(\mathbf{C}_0 | \theta_0, \mathbf{L}) P(\mathbf{C}_1 | \theta_1, \mathbf{L})$$

- 步骤二：用Gibbs sampling估计参数

$$\pi, \theta, L$$



步骤一：凑出简单形式的联合概率

- 分别凑每个因子

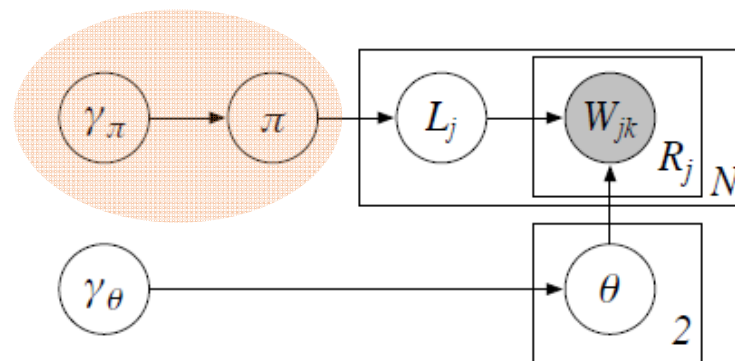
$$P(\mathbb{C}, \mathbf{L}, \pi, \theta_0, \theta_1; \gamma_{\pi 1}, \gamma_{\pi 0}, \gamma_{\theta}) = P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) P(\mathbf{L} | \pi) P(\theta_0 | \gamma_{\theta}) P(\theta_1 | \gamma_{\theta}) P(\mathbb{C}_0 | \theta_0, \mathbf{L}) P(\mathbb{C}_1 | \theta_1, \mathbf{L})$$

- 第一个因子(Beta)

$$P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) = \frac{\Gamma(\gamma_{\pi 1} + \gamma_{\pi 0})}{\Gamma(\gamma_{\pi 1}) \Gamma(\gamma_{\pi 0})} \pi^{\gamma_{\pi 1} - 1} (1 - \pi)^{\gamma_{\pi 0} - 1}$$



$$P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) \propto \pi^{\gamma_{\pi 1} - 1} (1 - \pi)^{\gamma_{\pi 0} - 1}$$



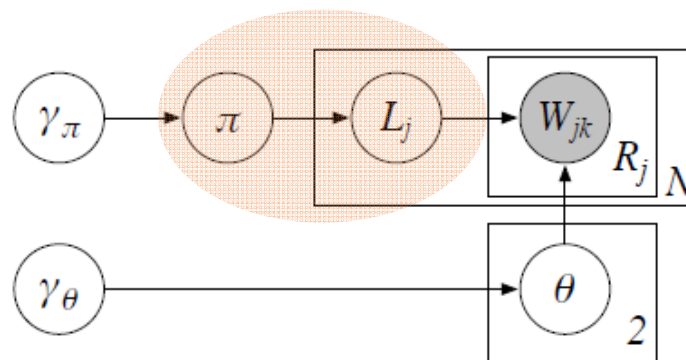
步骤一：凑出简单形式的联合概率

- 分别凑每个因子

$$P(\mathbb{C}, \mathbf{L}, \pi, \theta_0, \theta_1; \gamma_{\pi 1}, \gamma_{\pi 0}, \gamma_{\theta}) = P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) P(\mathbf{L} | \pi) P(\theta_0 | \gamma_{\theta}) P(\theta_1 | \gamma_{\theta}) P(\mathbb{C}_0 | \theta_0, \mathbf{L}) P(\mathbb{C}_1 | \theta_1, \mathbf{L})$$

- 第二个因子(Binomial)

$$\begin{aligned} P(\mathbf{L} | \pi) &= \prod_{n=1}^N \pi^{L_n} (1 - \pi)^{(1 - L_n)} \\ &= \pi^{C_1} (1 - \pi)^{C_0} \end{aligned}$$



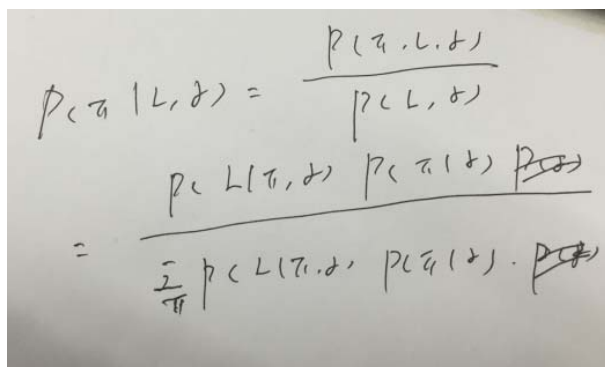
步骤一：凑出简单形式的联合概率

- 分别凑每个因子

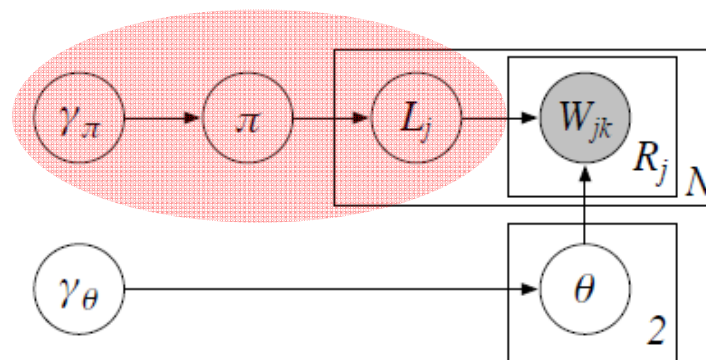
$$P(\mathbf{C}, \mathbf{L}, \pi, \theta_0, \theta_1; \gamma_{\pi 1}, \gamma_{\pi 0}, \gamma_{\theta}) = P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) P(\mathbf{L} | \pi) P(\theta_0 | \gamma_{\theta}) P(\theta_1 | \gamma_{\theta}) P(\mathbf{C}_0 | \theta_0, \mathbf{L}) P(\mathbf{C}_1 | \theta_1, \mathbf{L})$$

- 第一个因子(Beta)*第二个因子(Binomial)

$$\begin{aligned} P(\pi | \mathbf{L}; \gamma_{\pi 1}, \gamma_{\pi 0}) &\propto P(\mathbf{L} | \pi) P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) \\ &\propto [\pi^{C_1} (1 - \pi)^{C_0}] [\pi^{\gamma_{\pi 1} - 1} (1 - \pi)^{\gamma_{\pi 0} - 1}] \\ &\propto \pi^{C_1 + \gamma_{\pi 1} - 1} (1 - \pi)^{C_0 + \gamma_{\pi 0} - 1} \end{aligned}$$



$$\begin{aligned} p(\pi | \mathbf{L}, \gamma) &= \frac{p(\pi, \mathbf{L}, \gamma)}{p(\mathbf{L}, \gamma)} \\ &= \frac{p(\mathbf{L} | \pi, \gamma) p(\pi | \gamma)}{\sum_{\pi} p(\mathbf{L} | \pi, \gamma) p(\pi | \gamma)} \end{aligned}$$



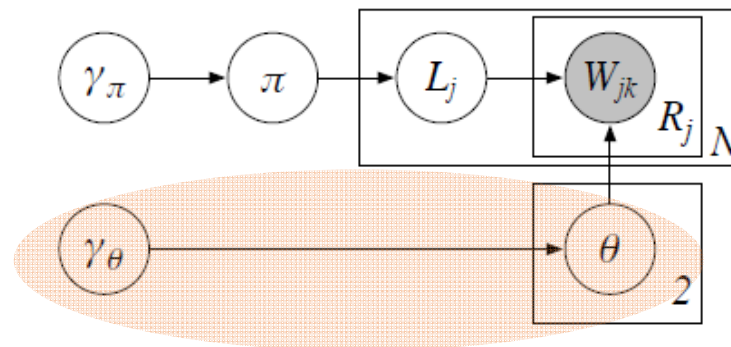
步骤一：凑出简单形式的联合概率

- 分别凑每个因子

$$P(\mathbb{C}, \mathbf{L}, \pi, \theta_0, \theta_1; \gamma_{\pi 1}, \gamma_{\pi 0}, \gamma_{\theta}) = P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) P(\mathbf{L} | \pi) P(\theta_0 | \gamma_{\theta}) P(\theta_1 | \gamma_{\theta}) P(\mathbb{C}_0 | \theta_0, \mathbf{L}) P(\mathbb{C}_1 | \theta_1, \mathbf{L})$$

- 第三、四个因子(Dirichlet)

$$\begin{aligned} P(\theta | \gamma_{\theta}) &= \frac{\Gamma(\sum_{i=1}^V \gamma_{\theta i})}{\prod_{i=1}^V \Gamma(\gamma_{\theta i})} \prod_{i=1}^V \theta_i^{\gamma_{\theta i} - 1} \\ &= c' \prod_{i=1}^V \theta_i^{\gamma_{\theta i} - 1} \\ &\propto \prod_{i=1}^V \theta_i^{\gamma_{\theta i} - 1} \end{aligned}$$



步骤一：凑出简单形式的联合概率

- 分别凑每个因子

$$P(\mathbb{C}, \mathbf{L}, \pi, \theta_0, \theta_1; \gamma_{\pi 1}, \gamma_{\pi 0}, \gamma_{\theta}) = P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) P(\mathbf{L} | \pi) P(\theta_0 | \gamma_{\theta}) P(\theta_1 | \gamma_{\theta}) P(\mathbb{C}_0 | \theta_0, \mathbf{L}) P(\mathbb{C}_1 | \theta_1, \mathbf{L})$$

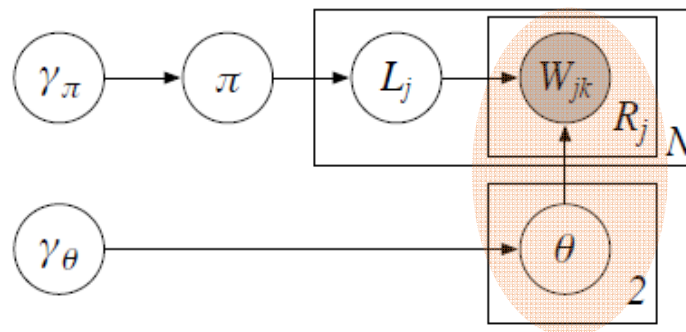
- 第五、六个因子(Multinomial)

$$P(\mathbf{W}_n | \mathbf{L}, \theta_{\mathbf{L}_n}) = \prod_{i=1}^V \theta_i^{W_{ni}}$$

↓

$$P(\mathbb{C}_x | \mathbf{L}, \theta_x) = \prod_{n \in \mathbb{C}_x} \prod_{i=1}^V \theta_{x,i}^{W_{ni}}$$

$$= \prod_{i=1}^V \theta_{x,i}^{N_{\mathbb{C}_x}(i)}$$



步骤一：凑出简单形式的联合概率

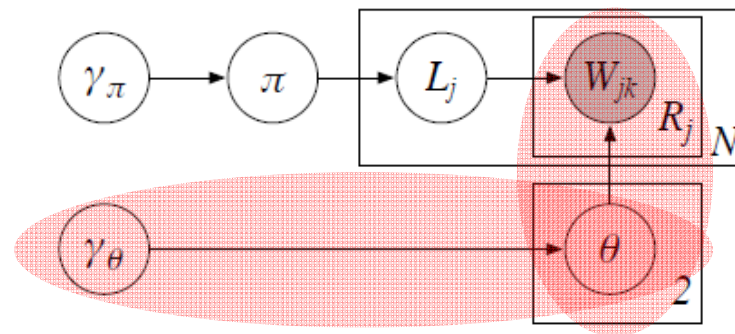
- 分别凑每个因子

$$P(\mathbb{C}, \mathbf{L}, \pi, \theta_0, \theta_1; \gamma_{\pi 1}, \gamma_{\pi 0}, \gamma_{\theta}) = P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) P(\mathbf{L} | \pi) P(\theta_0 | \gamma_{\theta}) P(\theta_1 | \gamma_{\theta}) P(\mathbb{C}_0 | \theta_0, \mathbf{L}) P(\mathbb{C}_1 | \theta_1, \mathbf{L})$$

- 第三、四个因子(Dirichlet)*第五、六个因子(Multinomial)

$$\begin{aligned} P(\theta | \mathbf{W}_n; \gamma_{\theta}) &= P(\mathbf{W}_n | \theta) P(\theta | \gamma_{\theta}) \\ &\propto \prod_{i=1}^V \theta_i^{W_{ni}} \prod_{i=1}^V \theta_i^{\gamma_{\theta i} - 1} \\ &\propto \prod_{i=1}^V \theta_i^{W_{ni} + \gamma_{\theta i} - 1} \end{aligned}$$

$$P(\theta_x | \mathbb{C}_x; \gamma_{\theta}) \propto \prod_{i=1}^V \theta_{x,i}^{N_{\mathbb{C}_x}(i) + \gamma_{\theta i} - 1}$$



步骤一：凑出简单形式的联合概率

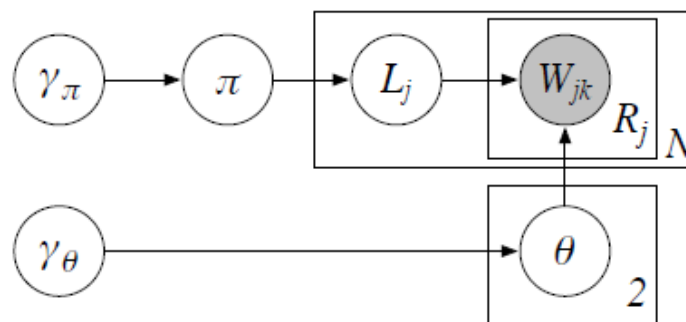
• 凑出的结果

$$P(\mathbb{C}, \mathbf{L}, \pi, \theta_0, \theta_1; \gamma_{\pi 1}, \gamma_{\pi 0}, \gamma_{\theta}) = P(\pi | \gamma_{\pi 1}, \gamma_{\pi 0}) P(\mathbf{L} | \pi) P(\theta_0 | \gamma_{\theta}) P(\theta_1 | \gamma_{\theta}) P(\mathbb{C}_0 | \theta_0, \mathbf{L}) P(\mathbb{C}_1 | \theta_1, \mathbf{L})$$

$$\propto \pi^{C_1 + \gamma_{\pi 1} - 1} (1 - \pi)^{C_0 + \gamma_{\pi 0} - 1} \prod_{i=1}^V \theta_{0,i}^{\mathcal{N}_{\mathbb{C}_0}(i) + \gamma_{\theta 0} - 1} \theta_{1,i}^{\mathcal{N}_{\mathbb{C}_1}(i) + \gamma_{\theta 1} - 1}$$



搞定



步骤一：凑出简单形式的联合概率

- 还有更简单的形式吗(**COLLAPSED Gibbs sampling: 积掉 π**)

$$\begin{aligned}
 P(\mathbf{L}, \mathbb{C}, \boldsymbol{\theta}_0, \boldsymbol{\theta}_1; \boldsymbol{\mu}) &= \int_{\pi} P(\mathbf{L}, \mathbb{C}, \boldsymbol{\theta}_0, \boldsymbol{\theta}_1, \pi; \boldsymbol{\mu}) \, d\pi \\
 &= P(\boldsymbol{\theta}_0 | \boldsymbol{\gamma}_{\boldsymbol{\theta}}) P(\boldsymbol{\theta}_1 | \boldsymbol{\gamma}_{\boldsymbol{\theta}}) P(\mathbb{C}_0 | \boldsymbol{\theta}_0, \mathbf{L}) P(\mathbb{C}_1 | \boldsymbol{\theta}_1, \mathbf{L}) \int_{\pi} P(\pi | \boldsymbol{\gamma}_{\pi 1}, \boldsymbol{\gamma}_{\pi 0}) P(\mathbf{L} | \pi) \, d\pi \\
 P(\mathbf{L}, \mathbb{C}, \boldsymbol{\theta}_0, \boldsymbol{\theta}_1; \boldsymbol{\mu}) &\propto \frac{\Gamma(\boldsymbol{\gamma}_{\pi 1} + \boldsymbol{\gamma}_{\pi 0})}{\Gamma(\boldsymbol{\gamma}_{\pi 1}) \Gamma(\boldsymbol{\gamma}_{\pi 0})} \frac{\Gamma(C_1 + \boldsymbol{\gamma}_{\pi 1}) \Gamma(C_0 + \boldsymbol{\gamma}_{\pi 0})}{\Gamma(N + \boldsymbol{\gamma}_{\pi 1} + \boldsymbol{\gamma}_{\pi 0})} \prod_{i=1}^V \theta_{0,i}^{\mathcal{N}_{\mathbb{C}_0}(i) + \boldsymbol{\gamma}_{\boldsymbol{\theta}_0} - 1} \theta_{1,i}^{\mathcal{N}_{\mathbb{C}_1}(i) + \boldsymbol{\gamma}_{\boldsymbol{\theta}_1} - 1}
 \end{aligned}$$



$$\begin{aligned}
 \int_{\pi} P(\pi | \boldsymbol{\gamma}_{\pi 1}, \boldsymbol{\gamma}_{\pi 0}) P(\mathbf{L} | \pi) \, d\pi &= \int_{\pi} \frac{\Gamma(\boldsymbol{\gamma}_{\pi 1} + \boldsymbol{\gamma}_{\pi 0})}{\Gamma(\boldsymbol{\gamma}_{\pi 1}) \Gamma(\boldsymbol{\gamma}_{\pi 0})} \pi^{\boldsymbol{\gamma}_{\pi 1} - 1} (1 - \pi)^{\boldsymbol{\gamma}_{\pi 0} - 1} \pi^{C_1} (1 - \pi)^{C_0} \, d\pi \\
 &= \frac{\Gamma(\boldsymbol{\gamma}_{\pi 1} + \boldsymbol{\gamma}_{\pi 0})}{\Gamma(\boldsymbol{\gamma}_{\pi 1}) \Gamma(\boldsymbol{\gamma}_{\pi 0})} \int_{\pi} \pi^{C_1 + \boldsymbol{\gamma}_{\pi 1} - 1} (1 - \pi)^{C_0 + \boldsymbol{\gamma}_{\pi 0} - 1} \, d\pi
 \end{aligned}$$

步骤二：用Gibbs sampling估计参数

• Sampling for Document Labels

$$\textcircled{1} \quad P(L_j | L^{(-j)}, C^{(-j)}, \theta_0, \theta_1; \mu) = \frac{P(L_j, W_j, L^{(-j)}, C^{(-j)}, \theta_0, \theta_1; \mu)}{P(L^{(-j)}, C^{(-j)}, \theta_0, \theta_1; \mu)}$$

$$= \frac{P(L, C, \theta_0, \theta_1; \mu)}{P(L^{(-j)}, C^{(-j)}, \theta_0, \theta_1; \mu)}$$

$$\textcircled{2} \quad P(L, C, \theta_0, \theta_1; \mu) \propto \frac{\Gamma(\gamma_{\pi 1} + \gamma_{\pi 0})}{\Gamma(\gamma_{\pi 1})\Gamma(\gamma_{\pi 0})} \frac{\Gamma(C_1 + \gamma_{\pi 1})\Gamma(C_0 + \gamma_{\pi 0})}{\Gamma(N + \gamma_{\pi 1} + \gamma_{\pi 0})} \prod_{i=1}^V \theta_{0,i}^{N_{C_0}(i) + \gamma_{\theta 0} - 1} \theta_{1,i}^{N_{C_1}(i) + \gamma_{\theta 1} - 1}$$

与 L 无关 \rightarrow cancel out

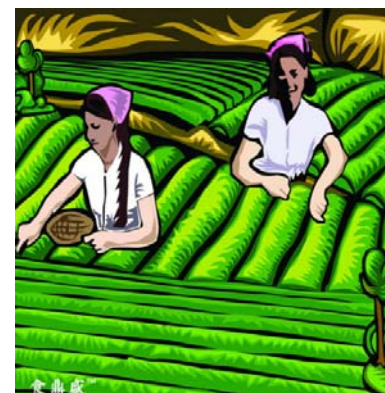
$$\frac{\Gamma(C_1 + \gamma_{\pi 1})\Gamma(C_0 + \gamma_{\pi 0})}{\Gamma(N + \gamma_{\pi 1} + \gamma_{\pi 0})}$$

$$\frac{\Gamma(C_0^{(-j)} + \gamma_{\pi 0})\Gamma(C_1^{(-j)} + \gamma_{\pi 1})}{\Gamma(N + \gamma_{\pi 1} + \gamma_{\pi 0} - 1)}$$

$$\frac{C_x + \gamma_{\pi x} - 1}{N + \gamma_{\pi 1} + \gamma_{\pi 0} - 1}$$

$$\prod_{i=1}^V \frac{\theta_{x,i}^{N_{C_x}(i) + \gamma_{\theta x} - 1}}{\theta_{x,i}^{N_{C_x^{(-j)}}(i) + \gamma_{\theta x} - 1}} = \prod_{i=1}^V \theta_{x,i}^{W_{jx}}$$

$$\textcircled{3} \quad P(L_j = x | L^{(-j)}, C^{(-j)}, \theta_0, \theta_1; \mu) = \frac{C_x + \gamma_{\pi x} - 1}{N + \gamma_{\pi 1} + \gamma_{\pi 0} - 1} \prod_{i=1}^V \theta_{x,i}^{W_{jx}}$$



采呀采呀快采样。
采样采得心花开。

步骤二：用Gibbs sampling估计参数

- Sampling for Document Labels

$$\Pr(L_j = x | \mathbf{L}^{(-j)}, \mathbf{C}^{(-j)}, \theta_0, \theta_1; \mu) = \frac{C_x + \gamma_{\pi x} - 1}{N + \gamma_{\pi 1} + \gamma_{\pi 0} - 1} \prod_{i=1}^V \theta_{x,i}^{W_{ji}}$$

已知Label的Docs，
按真实Label计数。



1. Let value0 = expression (49) with $x = 0$
2. Let value1 = expression (49) with $x = 1$
3. Let the distribution be $\langle \frac{\text{value0}}{\text{value0} + \text{value1}}, \frac{\text{value1}}{\text{value0} + \text{value1}} \rangle$
4. Select the value of $L_j^{(t+1)}$ as the result of a Bernoulli trial (weighted coin flip) according to this distribution.

步骤二：用Gibbs sampling估计参数

- Sampling for θ

- C, L 与 μ 独立

$$P(\theta_x | \mathbb{C}_x; \gamma_\theta) \propto \prod_{i=1}^V \theta_{x,i}^{N_{\mathbb{C}_x}(i) + \gamma_{\theta i} - 1}$$

$$P(C, L | \theta; \mu) = P(C, L | \theta)$$



$$P(\theta | \mathbb{C}, \mathbb{L}; \mu) \propto P(\mathbb{C}, \mathbb{L} | \theta) P(\theta | \mu)$$



Dirichlet

Multinomial

Dirichlet

- So, 直接按狄利克雷分布采样

$$\theta \sim \text{Dir}(N_{\mathbb{C}_x}(i) + \gamma_{\theta i})$$

问：怎么采样Dirichlet分布？

with $\mathbf{a} = \langle a_1, \dots, a_V \rangle$

答：分别采样每个因子

$$\text{Gamma}(\alpha_i, 1) = \frac{y_i^{\alpha_i - 1} e^{-y_i}}{\Gamma(\alpha_i)}$$

$$a_i = y_i / \sum_{j=1}^V y_j$$

PGM based Naïve Bayes Model

- Gibbs sampling algorithm

```
1: for  $t := 1$  to  $T$  do
2:   for  $j := 1$  to  $N$  do
3:     if  $j$  is not a training document then
4:       Subtract  $j$ 's word counts from the total word counts of whatever class it's currently a member of
5:       Subtract 1 from the count of documents with label  $L_j$ 
6:       Assign a new label  $L_j^{(t+1)}$  to document  $j$  as described at the end of Section 2.5.1
7:       Add 1 to the count of documents with label  $L_j^{(t+1)}$ 
8:       Add  $j$ 's word counts to the total word counts for class  $L_j^{(t+1)}$ 
9:     end if
10:  end for
11:   $\mathbf{t}_0 :=$  vector of total word counts from class 0, including pseudocounts
12:   $\theta_0 \sim \text{Dirichlet}(\mathbf{t}_0)$ , as described in Section 2.5.2
13:   $\mathbf{t}_1 :=$  vector of total word counts from class 1, including pseudocounts
14:   $\theta_1 \sim \text{Dirichlet}(\mathbf{t}_1)$ , as described in Section 2.5.2
15: end for
```



Part III

PLSA and LDA

Probabilistic latent semantic analysis (PLSA)

PLSA: $\alpha \rightarrow z \rightarrow w$

z : 隐含变量
 x : 文档特征
 z : 隐含变量
 w : 词特征
 z : $q(z) = p(z|x, \theta)$
 w : $J(\theta) = \sum_z q(z) \ln p(x, z|\theta)$
 w : $\theta^{new} = \arg \max_{\theta} J(\theta)$

公式太多，我手写了奥~



PLSA的EM求解过程 (E-Step)

$$\begin{aligned} \bar{c} &= p(z_k | w_j, d_i) = \frac{p(d_i, w_j, z_k)}{p(w_j, d_i)} \\ &= \frac{p(w_j | d_i, z_k) \cdot p(z_k | d_i) \cdot p(d_i)}{\sum_k p(w_j | d_i, z_k) \cdot p(z_k | d_i) \cdot p(d_i)} \\ &= \frac{p(w_j | z_k) p(z_k | d_i)}{\sum_k p(w_j | z_k) p(z_k | d_i)} \end{aligned}$$

令 $\phi_{kj} = p(w_j | z_k)$ 则参数 $\phi = \begin{pmatrix} \phi_{kj} \\ \theta_{ik} \end{pmatrix}$

$\theta_{ik} = p(z_k | d_i)$

参数化 令 $\phi_{kj} = p(w_j | z_k)$ 则参数 $\phi = \begin{pmatrix} \phi_{kj} \\ \theta_{ik} \end{pmatrix}$

$\theta_{ik} = p(z_k | d_i)$

(注: 当文档 d_i 中, 词 w_j 出现 x_{ij} 次, 则 $\sum_k \phi_{kj} x_{ij} = 1$ 且 $\sum_k \theta_{ik} = 1$)

$\phi_{kj} \in K \times V$
 $\theta_{ik} \in M \times K$

根据EM原理, 有以下约束

$\sum_j \phi_{kj} = 1$
 $\sum_k \theta_{ik} = 1$

PLSA的EM求解过程 (M-Step)

M:

$$\begin{aligned} \frac{1}{2} q(z_k) &= p(z_k | w_j, d_i) \\ L(0) &= \sum_k q(z_k) \ln \prod_{i,j} p(d_i, w_j, z_k | 0) \\ &= \sum_k q(z_k) \ln \prod_{i,j} p(d_i, w_j, z_k | 0)^{n(d_i, w_j)} \\ &= \sum_k q(z_k) \sum_i \sum_j \ln p(d_i, w_j, z_k | 0) \\ &= \sum_k q(z_k) \sum_i \sum_j \ln p(d_i, w_j, z_k | 0) \\ &= \sum_i \sum_j n(d_i, w_j) \sum_k q(z_k) \ln p(d_i, w_j, z_k | 0) \\ &= \sum_i \sum_j n(d_i, w_j) \sum_k q(z_k) \ln \phi_{kj} \theta_{ik} \cdot p(d_i) \end{aligned}$$

求参数

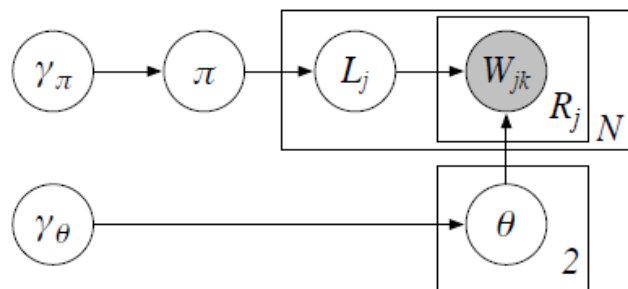
$$\begin{aligned} L(0) &= \sum_i \sum_j n(d_i, w_j) \sum_k q(z_k) \ln \phi_{kj} \theta_{ik} + \lambda_1 (1 - \sum_j \phi_{kj}) + \lambda_2 (1 - \sum_k \theta_{ik}) \\ \frac{\partial L(0)}{\partial \phi_{kj}} &= \sum_i n(d_i, w_j) q(z_k) \cdot \frac{1}{\phi_{kj}} - \lambda_1 = 0 \\ \frac{\partial L(0)}{\partial \theta_{ik}} &= \sum_j n(d_i, w_j) q(z_k) \cdot \frac{1}{\theta_{ik}} - \lambda_2 = 0 \end{aligned}$$

$$\begin{aligned} \sum_j \phi_{kj} &= \frac{1}{\lambda_1} \sum_i n(d_i, w_j) q(z_k) = 1 \\ \sum_k \theta_{ik} &= \frac{1}{\lambda_2} \sum_j n(d_i, w_j) q(z_k) = 1 \\ \text{求解: } \lambda_1 &= \sum_i n(d_i, w_j) q(z_k) \\ \lambda_2 &= \sum_j n(d_i, w_j) q(z_k) \\ &= \sum_j n(d_i, w_j) = n(d_i) \\ \text{求解参数: } \phi_{kj} &= \frac{\sum_i n(d_i, w_j) q(z_k)}{\sum_i n(d_i, w_j) q(z_k)} \\ \theta_{ik} &= \frac{\sum_j n(d_i, w_j) q(z_k)}{\sum_j n(d_i, w_j) q(z_k)} \end{aligned}$$

求解参数

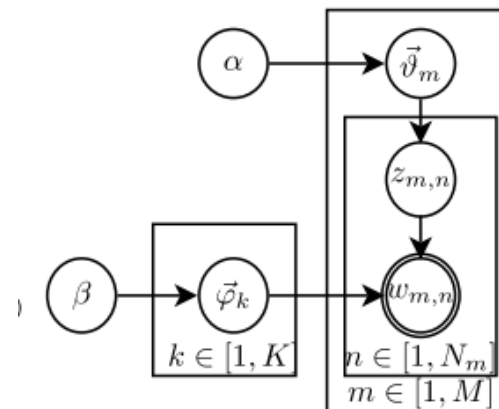
Latent Dirichlet allocation (LDA)

- 对比差异（Label的个数从2扩展到K）
 - π 从Beta分布扩展到Dirichlet分布
 - θ 被生成的次数从2扩展到K次



Naive Bayes Model

对比
↔



LDA Model

LDA Notations

M number of documents to generate (const scalar).

K number of topics / mixture components (const scalar).

V number of terms t in vocabulary (const scalar).

$\vec{\alpha}$ hyperparameter on the mixing proportions (K -vector or scalar if symmetric).

$\vec{\beta}$ hyperparameter on the mixture components (V -vector or scalar if symmetric).

$\vec{\vartheta}_m$ parameter notation for $p(z|d=m)$, the topic mixture proportion for document m . One proportion for each document, $\underline{\Theta} = \{\vec{\vartheta}_m\}_{m=1}^M$ ($M \times K$ matrix).

$\vec{\varphi}_k$ parameter notation for $p(t|z=k)$, the mixture component of topic k . One component for each topic, $\underline{\Phi} = \{\vec{\varphi}_k\}_{k=1}^K$ ($K \times V$ matrix).

N_m document length (document-specific), here modelled with a Poisson distribution [BNJ02] with constant parameter ξ .

$z_{m,n}$ mixture indicator that chooses the topic for the n th word in document m .

$w_{m,n}$ term indicator for the n th word in document m .

Solving LDA using Gibbs Sampling

- 凑出联合概率

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

$$p(\vec{z}, \vec{w} | \vec{\alpha}, \vec{\beta}) = \prod_{z=1}^K \frac{\Delta(\vec{n}_z + \vec{\beta})}{\Delta(\vec{\beta})} \cdot \prod_{m=1}^M \frac{\Delta(\vec{n}_m + \vec{\alpha})}{\Delta(\vec{\alpha})}.$$

- 凑出条件概率

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

topic part document part

数数而已



$$\begin{aligned} \varphi_{k,t} &= \frac{n_k^{(t)} + \beta_t}{\sum_{t=1}^V n_k^{(t)} + \beta_t}, \\ \vartheta_{m,k} &= \frac{n_m^{(k)} + \alpha_k}{\sum_{k=1}^K n_m^{(k)} + \alpha_k}. \end{aligned}$$

Remaining Problems in LDA

- hyper parameters
 - symmetric Dirichlet priors
 - $\alpha = 50/K$ and $\beta = 0.01$
- Querying (波浪表示查询文档)

$$p(\tilde{z}_i=k|\tilde{w}_i=t, \tilde{\mathbf{z}}_{-i}, \tilde{\mathbf{w}}_{-i}; \mathcal{M}) = \frac{n_k^{(t)} + \tilde{n}_{k,\neg i}^{(t)} + \beta_t}{\sum_{l=1}^V n_k^{(l)} + \tilde{n}_{k,\neg i}^{(l)} + \beta_t} \cdot \frac{n_{\tilde{m},\neg i}^{(k)} + \alpha_k}{[\sum_{k=1}^K n_{\tilde{m}}^{(k)} + \alpha_k] - 1}$$

大规模语料下数值很稳定

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

Further Reading

- Gershman, S. J., & Blei, D. M. (2012). A tutorial on Bayesian nonparametric models. *Journal of Mathematical Psychology*, 56(1), 1–12.
- Knight, K. (2009). Bayesian Inference with Tears. *A Tutorial Workbook for Natural Language Researchers*, 4(September), 388–395.

