

HW5B

1) basic sampling (unif, rejection, inverse transform ...)

2) Gibbs sampling.

3) Dirichlet-sampling (Gibbs)

4) baby LDA inference, given data \Rightarrow find param

• 20NG $V = \text{vocab-size}$ $K = 53,000$

• by dataset 1000 docs, 500 unique words (debug)
 $K = 200$

LDA as generative process

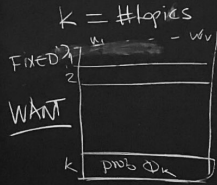
Given params \Rightarrow generate data (fake)

topic $\phi^k = \phi_k = \text{distribution (words)} = (\phi_1^k, \phi_2^k, \dots, \phi_v^k)$
 categories (words) probs over words

NOT FIXED
 topics = unknown?

$\phi^k = \text{R.V. topic}$

given prior Dirichlet (β)
 distributed $\beta = \text{prior params over } v \text{ word}$
 $(\beta^1, \dots, \beta^v)$



$d \in D$ documents $\theta_d = \text{distribution over topics} = (\theta_d^1, \theta_d^2, \dots, \theta_d^k)$
 prob over topics $\sum_{j=1}^k \theta_d^j = 1$

unknown $\theta_d = \text{R.V. doc}$
 given prior Dirichlet (α)
 distributed $\alpha = \text{prior params over } k \text{ topics}$
 $(\alpha^1, \alpha^2, \dots, \alpha^k)$



Generate data

- Sample k topics $\Phi^k \sim \text{Dirichlet}(\beta)$

For each doc $d \in D$

- Sample $\theta_d \sim \text{Dirichlet}(\alpha)$ $\theta_d = \text{dist}(\text{topics})$

• $L_d =$ ^(given) known length of doc d .

• For $i = 1: L_d$ $i =$ index of word i th word w_i

- Sample topic for word w_i $z_i \sim \text{Multinomial}(\theta_d)$
categorical?

- Sample word w from topic z_i : $w_i \sim \text{Multinomial}(\Phi_{z_i})$

caty dist
(discrete)

prob: (p_1, p_2, \dots, p_V)

$\sum p_i = 1$

prob($X=K$) = p_K

$V=2$ BINOMIAL (win)

N trials for V categories
output is $n_1, n_2, \dots, n_V = N$

$P(n_1, n_2, \dots, n_V | (p_1, p_2, \dots, p_V))$

$= \frac{N!}{n_1! n_2! \dots n_V!} \cdot p_1^{n_1} \cdot p_2^{n_2} \dots p_V^{n_V}$

costly (H03)
N=20 trials
 $P(\frac{5H}{10T}) = \frac{20!}{5^4 15^6} (0.5)^4 (0.5)^6$
MULTINOM

INFERENCE

start with initial $\bar{\alpha} = (\alpha_1, \dots, \alpha_K)$ $\bar{\beta} = (\beta_1, \dots, \beta_V)$

generative process $\rightarrow d_j = \text{bag-of-words-does}$ $d_j = (d_j^1, d_j^2, \dots, d_j^V)$ $\sum_j d_j^v = 1$

implicit $d_j = \Theta d = \text{dists over topics}$

likelihood $P(D = \text{docs} | \alpha, \beta, k, \text{Lang}) = \text{OB}$

params \rightarrow fixed

WANT: $\arg \max_{(\alpha, \beta)} P(D | \alpha, \dots)$

dist $\alpha = \text{prop}$ with $(\frac{1}{T}, \frac{1}{T}, \frac{2}{T}, \frac{2}{T}) \propto (1, 1, 2, 2)$

$$V=2 \text{ COIN} = R.V. (p, 1-p) \quad p \sim ?$$

of probab
 prior (initial guess)
 (belief about)

distns over possible (p, 1-p)

combin	50-50	60-40	70-30	90-10
prob	60%	10%	10%	20%

$$\text{Beta} \left(\begin{matrix} \text{bias} \\ p \\ 1-p \end{matrix} \middle| \begin{matrix} \text{params} \\ a, b \end{matrix} \right) = \frac{p^{a-1} \cdot (1-p)^{b-1} \cdot \Gamma(a+b)}{\Gamma(a) \Gamma(b)}$$

$$\Gamma(a) = (a-1)! \quad \Gamma(x+1) = (x+1)! = x! \cdot (x+1)$$

$$\Gamma(b) = (b-1)! \quad \Gamma(x) = (x-1)! = \frac{\Gamma(x+1)}{x}$$

$$\Gamma(a+b) = (a+b-1)!$$

a=2, b=1	n ₁ =6H	n ₂ =4T	a'=5, b'=5
a=2, b=1	n ₁ =900	n ₂ =900	a'=502, b'=501
a=200, b=100	500H	500T	a'=2000, b'=1500

$$P(N=n_1+n_2 \mid \text{prior beta } a, b) \stackrel{\text{Multi-Dirichlet}}{\propto} \text{Beta}(a+n_1, b+n_2)$$