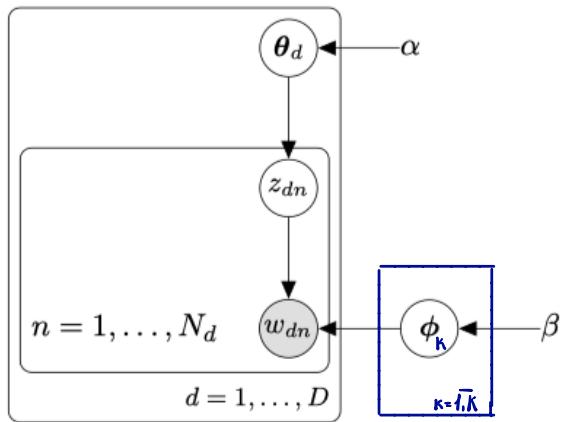


ML2 HOME ASSIGNMENT 6

id: 12179078

VOLODYMYR MEDENTSIV

PROBLEM 3



Topics $\varphi_\kappa \in \mathbb{R}^{|\mathcal{V}|} \sim \text{Dir}(\beta, \dots, \beta)$ if $\kappa = 1, K$

Documents $w_d \in \mathcal{D}$:

a) Topic distribution $\theta_d \sim \text{Dir}(\alpha, \dots, \alpha)$

b) Words w_{dn} in document d :

i) $z_{dn} \sim \text{Mult}(\theta_d)$

ii) $w_{dn} | z_{dn}, \varphi_{z_{dn}} \sim \text{Mult}(\varphi_{z_{dn}})$

$$A_{d\kappa} = \sum_{n=1}^{N_d} \delta(z_{dn} = \kappa)$$

$$B_{\kappa\omega} = \sum_{d=1}^D \sum_{i=n}^{N_d} \delta(w_{dn} = \omega) \delta(z_{dn} = \kappa)$$

$$M_\kappa = \sum_\omega B_{\kappa\omega}$$

1. The joint over observed data and latent:

$$\begin{aligned}
 p_{\text{joint}} &= p(\{\omega_{dn}\}_{n=1}^{N_d}, \zeta_{d=1}^D, \{\theta_d\}_{d=1}^D, \{\varphi_k\}_{k=1}^K, \{z_{dn}\}_{n=1}^{N_d}, \zeta_{d=1}^D) = \\
 &= \prod_{k=1}^K p(\varphi_k | \beta) \prod_{d=1}^D \prod_{n=1}^{N_d} p(\omega_{dn} | z_{dn}, \varphi_{z_{dn}}) \cdot p(z_{dn} | \theta_d) p(\theta_d | \alpha) = \\
 &= \prod_{k=1}^K p(\varphi_k | \beta) \cdot \prod_{d=1}^D p(\theta_d | \alpha) \prod_{n=1}^{N_d} p(\omega_{dn} | z_{dn}, \varphi_{z_{dn}}) \cdot p(z_{dn} | \theta_d) \quad \text{①}
 \end{aligned}$$

$$p(\varphi_k | \beta) = \frac{1}{B(\vec{\beta})} \cdot \prod_{\omega \in V} \varphi_{k\omega}^{\beta-1}$$

$$p(\theta_d | \alpha) = \frac{1}{B(\vec{\alpha})} \prod_{k=1}^K \theta_{dk}^{\alpha-1}$$

$$p(z_{dn} | \theta_d) = \prod_{k=1}^K \theta_{dk}^{[z_{dn}=k]}$$

$$p(\omega_{dn} | z_{dn}, \varphi_{z_{dn}}) = \prod_{\omega \in V} \varphi_{z_{dn}\omega}^{[\omega_{dn}=\omega]}$$

So $\prod_{d=1}^D \prod_{n=1}^{N_d} p(\omega_{dn} | z_{dn}, \varphi_{z_{dn}}) = \prod_{\omega \in V} \varphi_{z_{dn}\omega}^{B_{k\omega}}$, where $\varphi_{k\omega}$ is the ω -th component of the topic k

$$\prod_{d=1}^D \prod_{n=1}^{N_d} p(z_{dn} | \theta_d) = \prod_{d=1}^D \theta_{dk}^{A_{dk}}, \text{ where } \theta_{dk} \text{ is the } k\text{-th component in a distribution of topic } d.$$

$$\begin{aligned}
 &\textcircled{=} \prod_{k=1}^K \frac{1}{B(\vec{\beta})} \cdot \prod_{\omega \in V} \varphi_{k\omega}^{\beta-1} \cdot \prod_{d=1}^D \frac{1}{B(\vec{\alpha})} \cdot \prod_{k=1}^K \theta_{dk}^{\alpha-1} \cdot \prod_{\omega \in V} \varphi_{k\omega}^{B_{k\omega}} \theta_{dk}^{A_{dk}} = \\
 &= \frac{1}{B^K(\vec{\beta}) \cdot B^D(\vec{\alpha})} \cdot \prod_{k=1}^K \prod_{\omega \in V} \varphi_{k\omega}^{B_{k\omega} + \beta - 1} \cdot \prod_{d=1}^D \theta_{dk}^{A_{dk} + \alpha - 1}.
 \end{aligned}$$

$$2. p(\{w_{dn}\}_{n=1}^{Nd}, \{z_{dn}\}_{d=1}^D | \{z_{dn}\}_{n=1}^{Nd}, \{y_{d=1}^D\}) = \prod_{\omega} \prod_{\omega} \text{Joint d } \theta_{dk} d\varphi_{k\omega} =$$

$$\prod_{\omega} \frac{1}{B^K(\vec{\beta}) \cdot B^D(\vec{\alpha})} \cdot \prod_{k=1}^K \prod_{\omega \in U} \varphi_{k\omega}^{\beta_{k\omega} + \beta - 1} \cdot \prod_{d=1}^D \theta_{dk}^{\alpha_{dk} + \alpha - 1} \cdot d\theta_{dk} d\varphi_{k\omega}$$

$$= \frac{1}{B^K(\vec{\beta}) B^D(\vec{\alpha})} \cdot \prod_{k=1}^K \prod_{\omega \in U} \int_0^{\beta_{k\omega} + \beta - 1} \varphi_{k\omega} d\varphi_{k\omega} \prod_{d=1}^D \int_0^{\alpha_{dk} + \alpha - 1} \theta_{dk} d\theta_{dk} \quad \text{②}$$

$$\int_0^{\beta_{k\omega} + \beta - 1} \varphi_{k\omega} d\varphi_{k\omega} = \left[\frac{\varphi_{k\omega}}{\beta_{k\omega} + \beta} \right]_0^1 = \frac{1}{\beta_{k\omega} + \beta}$$

$$\int_0^1 \theta_{dk}^{\alpha_{dk} + \alpha - 1} d\theta_{dk} = \left[\frac{\theta_{dk}}{\alpha_{dk} + \alpha} \right]_0^1 = \frac{1}{\alpha_{dk} + \alpha}$$

$$\text{②} = \frac{1}{B^K(\vec{\beta}) B^D(\vec{\alpha}) \prod_{k=1}^K \prod_{d=1}^D (\alpha_{dk} + \alpha) \cdot \prod_{\omega \in U} (\beta_{k\omega} + \beta)}$$

$$3. p(z_{di} | \{w_{dn}\}_{n=1}^{Nd}, \{z_{dn}\}_{d=1}^D, \{y_{d=1}^D\}_{d=1}^D) =$$

$$= \frac{p(\{w_{dn}\}_{n=1}^{Nd}, \{z_{dn}\}_{d=1}^D, \{y_{d=1}^D\}_{d=1}^D)}{p(\{w_{dn}\}_{n=1}^{Nd}, \{z_{dn}\}_{d=1}^D, \{y_{d=1}^D\}_{d=1}^D)} = \begin{cases} p(\{w_{dn}\}_{n=1}^{Nd}, \{z_{dn}\}_{d=1}^D, \{y_{d=1}^D\}_{d=1}^D) \\ = p(w_{di}) \cdot p(\{w_{dn}\}_{n=1}^{Nd}, \{z_{dn}\}_{d=1}^D, \\ \{y_{d=1}^D\}_{d=1}^D | w_{di}) \end{cases}$$

$$= \frac{B^K(\vec{\beta}) B^D(\vec{\alpha}) \prod_{k=1}^K \prod_{d=1}^D (\alpha_{dk} + \alpha) \prod_{\omega \in U \setminus \{w_i\}} (\beta_{k\omega} + \beta)}{B^K(\vec{\beta}) B^D(\vec{\alpha}) \prod_{k=1}^K \prod_{d=1}^D (\alpha_{dk} + \alpha) \prod_{\omega \in U} (\beta_{k\omega} + \beta) \cdot p(w_{di})}$$

$$p(w_{di}) = \sum_{z_{di}} p(w_{di} | z_{di}, q_{di}) \cdot p(z_{di}) = \sum_{z_i} \text{Multi}(w_{di} | q_{di}) \cdot \text{Multi}(z_{di} | \theta_d)$$

PROBLEM 1

a)

$$\tilde{q}(x) = c q(x)$$

$$\tilde{q}(x) \geq \tilde{p}(x) \quad \forall x$$

Pseudocode of Rejection Sampling:

1. Set number of samples N

2. $i = 0$

3. do while $i < N$

3.1. $z \sim q(x)$

3.2. $u \sim U[0; \tilde{q}(z)]$

3.3. if $u > \tilde{p}(z)$:

reject z

else:

accept z as a sample from $p(x)$

$i += 1$

b) Yes, samples are independent.

$$p(x) = c_p \tilde{p}(x)$$

$$q(x) = c_q \tilde{q}(x)$$

Using Bishop's 11.23:

$$w_e = \frac{\tilde{z}_e}{\sum_m \tilde{z}_m} = \frac{\tilde{p}(z^{(e)}) / q(z^{(e)})}{\sum_m \tilde{p}(z^{(m)}) / q(z^{(m)})} = \frac{p(z^{(e)})}{c_p \cdot q(z^{(e)}) \cdot \sum_m \frac{p(z^{(m)})}{c_p \cdot q(z^{(m)})}} =$$

$$= \frac{p(z^{(l)}) / g(z^{(l)})}{\sum_m p(z^{(m)}) / g(z^{(m)})}.$$

d) Independence Sampler with $g(x_{t+1}|x_t) = g(x_t)$

$$\alpha = \min \left(1, \frac{p(x_{t+1})}{p(x_t)} \frac{g(x_t|x_{t+1})}{g(x_{t+1}|x_t)} \right) =$$

$$= \left[g(x_t|x_{t+1}) = \frac{g(x_{t+1}|x_t) g(x_t)}{g(x_{t+1})} = \frac{g(x_{t+1}) g(x_t)}{g(x_{t+1})} = g(x_t) \right] =$$

$$= \min \left(1, \frac{p(x_{t+1}) g(x_t)}{p(x_t) g(x_{t+1})} \right).$$

e) Samples are drawn from a probability distribution, which is independent from the previous sample ($x_{t+1} \sim g(x_{t+1})$) but the acceptance probability depends on x_t , thus samples are not independent.

f) We will obtain such a sequence:

$$x_1, x_1, x_3, x_4, x_4.$$

In practical applications we will most probably retain only:
 x_1, x_3, x_4 , so no duplicates.

g) Rejection and Importance sampler work very bad in high dimensional spaces. I think there are two reasons:

1. it's hard to find approximating distribution $g(x)$ such that its high-density regions will coincide with such regions in $p(x)$.
2. the approximation error, volume between $g(x)$ and $p(x)$, increases as the dimensionality of x increases (curse of dimensionality)

MCMC works better in high-dimensional spaces than rejection and importance sampling, but still its performance becomes worse with the increase in dimensionality. For example, let's consider

Metropolis-Hastings algorithm, with symmetric proposal and assume that $p(x) = p(x^0) \dots p(x^{(D)})$, $x \in \mathbb{R}^D$ then we can show how rapidly decreases acceptance probability with increase of D :

if $\frac{p(x_{t+1}^{(i)})}{p(x_t^{(i)})} = 0,9$ \Rightarrow acceptance probability for each component

of x is 0,9 but acceptance probability for the whole x is $0,9^D$. which could be < 1 as D increases. This example also illustrates curse of dimensionality.

PROBLEM 2

$$x \sim N(\mu | \mu, \tau^{-1})$$

$$\mu \sim N(\mu | \mu_0, s_0)$$

$$\tau \sim \text{Gamma}(\gamma | \alpha, \beta)$$

Derive Gibbs sampling distribution for the posterior $p(\mu, \tau | x)$.

1. The sampling procedure:

for $t = 1, T$

$$\mu^{(t)} \sim p(\mu | x, \tau^{(t-1)})$$

$$\tau^{(t)} \sim p(\tau | x, \mu^{(t)})$$

Let's find $p(\mu | x, \tau^{(t-1)})$ and $p(\tau | x, \mu^{(t)})$ using 2.139-2.143 and 2.149-2.151 in Bishop.

$$p(\mu | x, \tau^{(t-1)}) = N(\mu | \tilde{\mu}, \tilde{s}^2), \text{ assume that we observe all } x \text{ and } \tau^{(t-1)} \text{ is known}$$

$$\text{with } \tilde{\mu} = \frac{(\tau^{(t-1)})^{-1}}{s_0 + (\tau^{(t-1)})^{-1}} \mu_0 + \frac{s_0}{s_0 + (\tau^{(t-1)})^{-1}} \cdot x$$

$$\tilde{s}^2 = \left(\frac{1}{s_0} + \tau^{(t-1)} \right)^{-1} = \frac{s_0}{1 + s_0 \cdot \tau^{(t-1)}}$$

$$p(\tau | x, \mu^{(t)}) = \text{Gam}(\tau | \tilde{\alpha}, \tilde{\beta}), \text{ assume that we observe all } x \text{ and } \mu^{(t)} \text{ is known.}$$

$$\text{with } \tilde{\alpha} = \alpha_0 + 1/2$$

$$\tilde{\beta} = \beta_0 + 1/2 \cdot (x - \mu^{(t)})$$

PROBLEM 4

Given: $p(x|\mu) = \prod_{i=1}^D \mu_i^{x_i} (1-\mu_i)^{1-x_i}$
 $x \in \{0,1\}^D, \mu \in [0,1]^D$

- a) $E p(x|\mu) = \mu$.
- b) $\text{Cov } p(x|\mu) = E(x - \mu)(x - \mu)^T :$

$$E(x_i - \mu_i)(x_j - \mu_j) = E(x_i x_j - \mu_i x_j - x_i \mu_j + \mu_i \mu_j) \quad \text{①}$$

if $i \neq j$ $\stackrel{1}{=} E x_i E x_j - \mu_i E x_j - E x_i \cdot \mu_j + \mu_i \mu_j = 0$.

if $i = j$ $\stackrel{2}{=} \mu_i (1 - \mu_i)$

So $\text{Cov } x = \text{diag}(\mu(1-\mu))$

Given: $p(x|\mu, \pi) = \sum_{n=1}^K \pi_n p(x|\mu_n)$ with $\pi = (\pi_1, \dots, \pi_K)$, $\mu = (\mu_1, \dots, \mu_K)$
 and $\forall i: \mu_i \in [0,1]^D$. And $p(x|\mu_n) = \prod_{i=1}^D \mu_{ni}^{x_i} (1-\mu_{ni})^{1-x_i}$

c) $E x = E \sum_{n=1}^K \pi_n p(x|\mu_n) = \sum_{n=1}^K \pi_n E p(x|\mu_n) = \sum_{n=1}^K \pi_n \cdot \mu_n$.

Given: data set $X = (x_1, \dots, x_N)$

d) $\text{Loglikelihood} = \log p(X|\mu, \pi) = \log \prod_{n=1}^N p(x_n|\mu, \pi) = \sum_{i=1}^N \log p(x_i|\mu, \pi)$

$$= \sum_{i=1}^N \log \sum_{n=1}^K \pi_n p(x_i|\mu_n) = \sum_{i=1}^N \log \sum_{n=1}^K \left(\pi_n \cdot \prod_{d=1}^D \mu_{nd}^{x_i^d} (1-\mu_{nd})^{1-x_i^d} \right),$$

with x_i^d - d-th component of observed data point x_i .

e) We use loglikelihood because it decomposes $\log \prod$ into $\sum \log$ and

there we can then easily derive closed-form max-likelihood solutions.
 However in this mixture model we have $\sum \log \sum$ so we can't further decompose $\log \sum$ to easily derive analytical solution and will need numerical methods to optimize the log-likelihood.

Given: we introduce latent variables z_n - all-hot encoded vector

$$p(x_n, z_n | \mu, \pi) = p(z_n | \pi) p(x_n | z_n, \mu) = \prod_{k=1}^K \pi_k^{z_{nk}} p(x_n | \mu_k)^{z_{nk}}$$

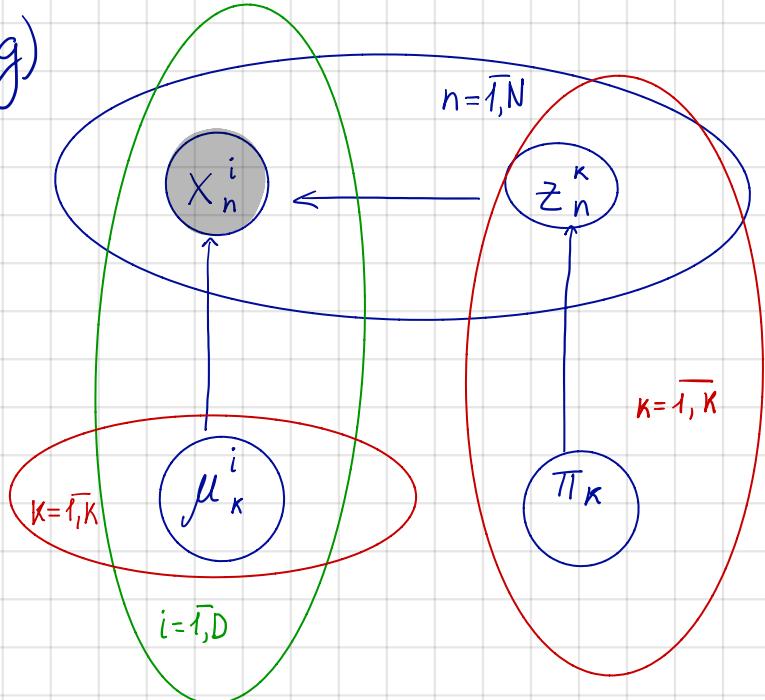
f) $\log p(X, Z | \mu, \pi) = \log \prod_{n=1}^N p(x_n, z_n | \mu, \pi) =$

$$= \sum_{n=1}^N \log p(x_n, z_n | \mu, \pi) = \sum_{n=1}^N \log p(z_n | \pi) p(x_n | z_n, \mu) =$$

$$= \sum_{n=1}^N \log \prod_{k=1}^K \pi_k^{z_{nk}} \left(\prod_{i=1}^D \mu_{xi}^{x_n^i} (1-\mu_{xi})^{1-x_n^i} \right)^{z_{nk}} =$$

$$= \sum_{n=1}^N \sum_{k=1}^K \left(z_{nk} \log \pi_k + z_{nk} \sum_{i=1}^D (x_n^i \log \mu_{xi} + (1-x_n^i) \log (1-\mu_{xi})) \right)$$

g)



h) VEM objective function $B(\{q_n(z_n)\}, \mu, \pi)$

$$B(\{q_n(z_n)\}, \mu, \pi) = \sum_n \sum_{z_n} q_n(z_n) \log p(x_n, z_n | \mu, \pi) - \sum_n \sum_{z_n} q_n(z_n) \log q_n(z_n) =$$

$$= \sum_n \sum_{z_n} q_n(z_n) \left[\sum_{k=1}^K \left(z_{nk} \log \pi_k + z_{nn} \sum_{i=1}^D (x_n^i \log \mu_{xi} + (1-x_n^i) \log (1-\mu_{xi})) \right) - \log q_n(z_n) \right]$$

i) Constraints: $\sum_k \pi_k = 1$ and $H_n \sum_{z_n} q(z_n) = 1$.

Then

$$\tilde{B}(\{q_n(z_n)\}, \mu, \pi, \tilde{\lambda}, \{\lambda_n\}) = B(\{q_n(z_n)\}, \mu, \pi) + \tilde{\lambda} \left(\sum_k \pi_k - 1 \right) + \sum_n \lambda_n (\sum_{z_n} q_n(z_n) - 1)$$

j) $q_n(z_n)$ approximates the posterior $p(z_n | x_n, \mu, \pi, \{\lambda_n\})$ and can be viewed as the responsibility which z_n takes in explaining x_n .

$$\frac{\partial \tilde{B}(\{q_n(z_n)\}, \mu, \pi, \tilde{\lambda}, \{\lambda_n\})}{\partial q_n(z_n)} =$$

$$= \sum_{k=1}^K \left(z_{nk} \log \pi_k + z_{nn} \sum_{i=1}^D (x_n^i \log \mu_{xi} + (1-x_n^i) \log (1-\mu_{xi})) \right) -$$

$$- \log q_n(z_n) - 1 + \lambda_n = 0 \Rightarrow$$

$$\Rightarrow \log q_n(z_n) = \log \prod_{k=1}^K \pi_k^{z_{nk}} \circ \left(\prod_{i=1}^D \mu_{xi}^{x_n^i} \cdot (1-\mu_{xi})^{1-x_n^i} \right)^{z_{nn}} - 1 + \lambda_n \Rightarrow$$

$$\Rightarrow q_n(z_n) = \exp(\lambda_n - 1) \circ \prod_{k=1}^K \pi_k^{z_{nk}} \circ \left(\prod_{i=1}^D \mu_{xi}^{x_n^i} \cdot (1-\mu_{xi})^{1-x_n^i} \right)^{z_{nn}}$$

k) M-step:

$$\frac{\partial \tilde{B}(\{f_n(z_n)\}, \mu, \pi, \tilde{\lambda}, \{\lambda_n\})}{\partial \pi_k} = \sum_n \sum_{z_n} f_n(z_n) z_{nk} \cdot \frac{1}{\pi_k} + \tilde{\lambda} = 0$$

$$\Rightarrow \pi_k = - \frac{\sum_n \sum_{z_n} f_n(z_n) z_{nk}}{\tilde{\lambda}}$$

Let's find $\tilde{\lambda}$:

$$\tilde{\lambda} \underbrace{\sum_k \pi_k}_{1} = - \sum_k \sum_n \sum_{z_n} f_n(z_n) z_{nk} \Rightarrow \tilde{\lambda} = -N$$

$$\Rightarrow \pi_k = \frac{\sum_n \sum_{z_n} f_n(z_n) z_{nk}}{N}$$