

Problem 1

a) We shall assume that $cq(x) \geq \tilde{p}(x)$, i.e.

Until termination (e.g. exhausted comp. res.):

1) Draw $x_0 \sim q(x)$

2) Draw $v_0 \sim \text{Uni}(0, kq(x_0))$

3) If $v_0 > \tilde{p}(x_0)$ then reject the sample, accept otherwise with prob:

$$p(\text{accept}) = \frac{1}{k} \int \tilde{p}(x) dx$$

b) Yes, they are ^{indep} (Ian Murray, 2004). Steps are indep of prev. samples.

$$c) E[f] = \int f(z) p(z) dz = \int f(z) \frac{p(z)}{q(z)} q(z) dz \approx \frac{1}{L} \sum_{i=1}^L \frac{p(x_i)}{q(x_i)} f(x_i) =$$

$$= \frac{1}{L} \sum_{i=1}^L w_i f(x_i), \text{ where } w_i = \frac{p(x_i)}{q(x_i)}$$

$$d) d(x_{t+1}, x_t) = \min \left(1, \frac{\tilde{p}(x_{t+1})}{\tilde{p}(x_t)} \frac{q(x_t)}{q(x_{t+1})} \right) = \min \left(1, \frac{p(x_{t+1})}{p(x_t)} \frac{\tilde{p}^2 q(x_t)}{\tilde{p}^2 q(x_{t+1})} \right) = \\ = \min \left(1, \frac{p(x_{t+1}) q(x_t)}{p(x_t) q(x_{t+1})} \right)$$

e) By rejecting x_2 and x_5 we leave them as prev. samples for the next iteration, thus, the sequence will be: x_1, x_1, x_3, x_4, x_6

f) 1) Rejection sampler \rightarrow Does not scale well, because as dim. grows, we have a large space (volume) between \tilde{p} and q , thus many samples will be rejected.

2) Importance sampler \rightarrow Does not scale well because as # of dimensions increases it becomes hard to find $q(x)$ that fits to $p(x)$, and thus a poor choice will make it hard to generate samples such that they are both probable according to $p(x)$ and $q(x)$.

$$\text{as } w_i = \frac{p(x_i)}{q(x_i)} \quad E[f] \approx \frac{\sum_i w_i f(x_i)}{\sum_j w_j}$$

3) Independence sampler \rightarrow Scale much better because Markov perspective leads to efficient decompositions of high-dimensional problems in a sequence of smaller problems that are much easier to solve.

e) Even though we get samples $x \sim q(x)$ which do not depend on the prev. samples, our threshold d still takes into account prev. samples, and thus the decision becomes dependent on the previous draws. Therefore, generated samples are not independent.

Problem 5

a) $E[\mathbf{x}] = \begin{pmatrix} E[x_1] \\ \vdots \\ E[x_n] \end{pmatrix}$, since $x_i \sim \text{Bern}(\mu_i)$, $E[x_i] = \mu_i + 0(1-\mu_i) = \mu_i$; thus $E[\mathbf{x}] = \mathbf{\mu}$

b) Since Bernoulli trials are independent, we'll get a diagonal matrix with variances along diagonal.

So $\text{Var}[x_i] = E[x_i^2] - (E[x_i])^2 = \mu_i - \mu_i^2 = \mu_i(1-\mu_i)$

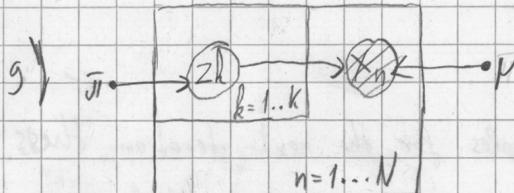
thus $\text{cov}[\mathbf{x}] = \text{diag}(\mu_i(1-\mu_i))$

c) $E[\mathbf{x}] = \sum_{k=1}^K \pi_k \mathbf{v}_k$; thus $E[\mathbf{x}] = \sum_{k=1}^K \pi_k \mathbf{v}_k$

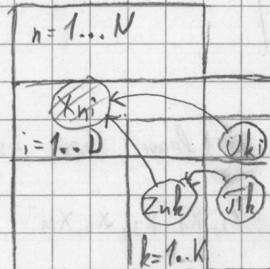
d) $L_I = \prod_x p(x|\nu, \pi)$, $\log L_I = \sum_x \left[\log \sum_{k=1}^K \left[\pi_k \prod_{i=1}^D \nu_{ki}^{x_i} (1-\nu_{ki})^{1-x_i} \right] \right]$

e) Because we are dealing with an incomplete log-likelihood, and we have a sum inside of the logarithm, there is no closed form solution and it is challenging to perform the optimization.

f) $\log L_c = \sum_{n=1}^N \sum_{k=1}^K \log p(x_n, z_{nk} | \nu, \pi) = \sum_{n=1}^N \sum_{k=1}^K \left[z_{nk} (\log p(z_{nk} | \pi) + \log p(x_n | z_{nk}, \nu)) \right] =$
 $= \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left(\log \pi_k + \sum_{i=1}^D [x_{ni} \log \nu_{ki} + (1-x_{ni}) \log (1-\nu_{ki})] \right)$



or more detailed



h) To find $B(q, \theta)$ (I shall use $L(q, \theta)$ instead as Bishop does) we take expectation of the complete log-likelihood, which is going to be our lower bound. $\theta = \{\nu, \pi\}$

$$L(q, \theta) = E_q [\log L_c] = \sum_{n=1}^N \sum_{k=1}^K q_n(k) \log p(x_n, z_{nk} | \nu, \pi) =$$

$$= \sum_{n=1}^N \sum_{k=1}^K q_n(k) \left[\log \pi_k + \sum_{i=1}^D x_{ni} \log \nu_{ki} + (1-x_{ni}) \log (1-\nu_{ki}) - \log q_n(k) \right] =$$

$$= \sum_{n=1}^N \sum_{k=1}^K q_n(k) \left(\log \pi_k + \sum_{i=1}^D (x_{ni} \log \nu_{ki} + (1-x_{ni}) \log (1-\nu_{ki})) - \log q_n(k) \right)$$

Problem 5 (continued)

$$i) \tilde{\mathcal{L}}(q, \tilde{\theta}) = \mathcal{L}(q, \theta) + \lambda \left(\sum_{k=1}^K \bar{\pi}_k - 1 \right) + \sum_{n=1}^N \beta_n \left(\sum_{k=1}^K q_n(k) - 1 \right)$$

$$\tilde{\theta} = \{\mu, \bar{\pi}, \lambda, \beta\}$$

$$j) \frac{d \tilde{\mathcal{L}}(q, \tilde{\theta})}{d q(z_{nk})} = \log \bar{\pi}_k + \left[\sum_{i=1}^D x_{ni} \log \nu_{ki} + (1-x_{ni}) \log (1-\nu_{ki}) \right] - \log q_n(k) - 1 + \beta_n = 0$$

$$q_n(k) = \exp(\beta_n - 1) \cdot \bar{\pi}_k \cdot \prod_{i=1}^D \nu_{ki}^{x_{ni}} (1-\nu_{ki})^{1-x_{ni}}$$

$q_n(k)$ is a normalized posterior probability of assigning datapoint x_n to class k .
 which by investigation is proportional to: $p(z_{nk} | \bar{\pi}) p(x_n | \nu_k)$

$$k) \frac{d \tilde{\mathcal{L}}(q, \theta)}{d \bar{\pi}_k} = \sum_{n=1}^N \frac{q_n(k)}{\bar{\pi}_k} + \lambda = 0 \Rightarrow \bar{\pi}_k = -\frac{1}{\lambda} \sum_{n=1}^N q_n(k)$$

$$-\bar{\pi}_k \lambda = \sum_{n=1}^N q_n(k) \Rightarrow -\lambda = \sum_{n=1}^N \sum_{k=1}^K q_n(k) = -N, \text{ thus we finally obtain:}$$

$$\bar{\pi}_k = \frac{N_k}{N}, \text{ where } N_k = \sum_{n=1}^N q_n(k)$$

Problem 3

Because we are sampling from $p(\nu, \tau | x)$, we need $p(\nu | x, \tau)$ and $p(\tau | x, \nu)$

$p(\nu | x, \tau) \propto p(x | \nu, \tau) p(\nu)$ since $x \sim \text{Gauss}$ and $\nu \sim \text{Gauss}$, the product will also be gaussian.

$$\begin{aligned} p(\nu | x, \tau) &\propto C_1 C_2 \exp\left(-\frac{1}{2} \left(\frac{(x-\nu)^2}{\tau} + (\nu - \mu_0)^2 s_0^{-2} \right)\right) = \\ &= C_1 C_2 \exp\left(-\frac{1}{2} \left(\frac{x^2}{\tau} - 2x\nu\tau + \frac{\nu^2}{\tau} + \nu^2 s_0^{-2} - 2\nu\nu_0 s_0^{-2} + \mu_0^2 s_0^{-2} \right)\right) = \\ &= C_1 C_2 \exp\left(-\frac{1}{2} \left(\underbrace{\nu^2(\tau + s_0^{-2})}_{S_N^{-1}} - 2\nu \underbrace{(x\tau + \nu_0 s_0^{-2})}_{S_N^{-1} \bar{\nu}_N} + \text{const} \right)\right) = \end{aligned}$$

And we can instantly see that terms have quadratic structure and thus,

$$p(\nu | x, \tau) = N\left(\nu | \frac{x\tau + \nu_0 s_0^{-2}}{\tau + s_0^{-2}}, \frac{1}{\tau + s_0^{-2}}\right) \quad \text{as } \nu_N = (\tau + s_0^{-2})^{-1} (x\tau + \nu_0 s_0^{-2})$$

$$\begin{aligned} p(\tau | x, \nu) &= \frac{p(x | \nu, \tau) p(\tau)}{p(x)} \propto p(x | \nu, \tau) p(\tau) = \\ &= \frac{1}{\Gamma(a)} b^a \tau^{a-1} \exp(-b\tau) \tau^{\frac{1}{2}a-1} \exp\left(-\frac{1}{2}(x-\nu)^2 \tau\right) \propto \tau^{a+\frac{1}{2}-1} \exp\left(-\left(b + \frac{(x-\nu)^2}{2}\right)\tau\right) \end{aligned}$$

Which by inspection we can see is Gamma distribution.

$$\text{Gamma}\left(\tau | a + \frac{1}{2}, b + \frac{(x-\nu)^2}{2}\right)$$