

# STATS 305C: Practice Exam

**Name:**

**Problem 1: Gaussian models**

Consider the following model,

$$\begin{aligned} x_{n,d} &\stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_d^2) && \text{for } n = 1, \dots, N; d = 1, \dots, D \\ \sigma_d^2 &= \prod_{k=1}^d \lambda_k^{-1} && \text{for } d = 1, \dots, D \\ \lambda_d &\stackrel{\text{iid}}{\sim} \text{Ga}(\alpha, 1) && \text{for } d = 1, \dots, D \end{aligned}$$

- (a) Suppose  $\alpha > 1$ . Describe how this *multiplicative inverse gamma* prior affects the distribution of the data,  $x_{n,d}$ . For example, how does the distribution of  $x_{n,1}$  generally compare to that of  $x_{n,D}$ ?
- (b) Let  $\boldsymbol{\lambda} = \{\lambda_k\}_{k=1}^K$  and  $\mathbf{X} = \{x_{n,d}\}_{n=1}^N, d=1}^D$ . Derive a Gibbs sampler for the posterior distribution  $p(\boldsymbol{\lambda} \mid \mathbf{X}; \alpha)$ .

**Solution:**

- (a)  $\mathbb{E}[\lambda_d] = \alpha$  so when  $\alpha > 1$ , we the precisions are generally greater than one as well. Thus, the precision  $(\sigma_d^2)^{-1} = \prod_{k=1}^d \lambda_k$  grows with  $d$ , and equivalently, the variance shrinks as  $d$  increases. Thus, we expect  $x_{n,d}$  to be more concentrated around the prior mean of 0 than  $x_{n,1}$ .

- (b) The joint probability is,

$$p(\mathbf{X}, \boldsymbol{\lambda}; \alpha) = \prod_{d=1}^D \text{Ga}(\lambda_d; \alpha, 1) \prod_{n=1}^N \prod_{d=1}^D \mathcal{N}\left(x_{n,d} \mid 0, \prod_{k=1}^d \lambda_k^{-1}\right)$$

Note that  $\lambda_k$  appears in the likelihood for all  $x_{n,d}$  where  $d \geq k$ . Thus, the conditional distribution of  $\lambda_k$  is proportional to,

$$\begin{aligned} p(\lambda_k \mid -) &\propto \text{Ga}(\lambda_k; \alpha, 1) \prod_{n=1}^N \prod_{d=k}^D \mathcal{N}\left(x_{n,d} \mid 0, \prod_{j=1}^d \lambda_j^{-1}\right) \\ &\propto \lambda_k^{\alpha-1} e^{-\lambda_k} \prod_{n=1}^N \prod_{d=k}^D \lambda_k^{\frac{1}{2}} \exp\left\{-\frac{\lambda_k \left(\prod_{i=1}^{k-1} \lambda_i\right) \left(\prod_{j=k+1}^d \lambda_j\right) x_{n,d}^2}{2}\right\} \\ &\propto \text{Ga}(\lambda_k; \alpha', \beta') \end{aligned}$$

where

$$\begin{aligned} \alpha' &= \alpha + \frac{N(D-k+1)}{2} \\ \beta' &= 1 + \sum_{n=1}^N \sum_{d=k}^D \frac{\left(\prod_{i=1}^{k-1} \lambda_i\right) \left(\prod_{j=k+1}^d \lambda_j\right) x_{n,d}^2}{2} \end{aligned}$$

To implement a Gibbs sampler, iteratively sample each  $\lambda_k$  from its conditional, holding the others fixed.

**Problem 2: Hierarchical models.**

Recall the probability density function of the gamma distribution,  $p(\lambda; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda}$ , where  $\Gamma(\cdot)$  is the gamma function. Now consider the following hierarchical model,

$$\begin{aligned} \beta &\sim \text{Ga}(\alpha_0, \beta_0) \\ \lambda_g &\sim \text{Ga}(\alpha, \beta) && \text{for } g = 1, \dots, G \\ x_{g,n} &\sim \text{Po}(\lambda_g) && \text{for } g = 1, \dots, G; n = 1, \dots, N. \end{aligned}$$

Using the Poisson probability mass function  $p(x | \lambda) = \frac{1}{x!} \lambda^x e^{-\lambda}$ , derive a Gibbs sampling algorithm for this hierarchical model. Specifically, derive the conditional distributions,

- $p(\lambda_g | \{x_{g,n}\}_{n=1}^N, \beta_g; \alpha)$ ,
- $p(\beta | \{\lambda_g\}_{g=1}^G; \alpha_0, \beta_0)$ .

**Solution:** The first conditional is,

$$\begin{aligned} p(\lambda_g | \{x_{g,n}\}_{n=1}^N, \beta_g; \alpha) &\propto p(\lambda_g | \beta_g; \alpha) \prod_{n=1}^N p(x_{g,n} | \lambda_g) \\ &\propto \text{Ga}(\lambda_g; \alpha, \beta_g) \prod_{n=1}^N \text{Po}(x_{g,n} | \lambda_g) \\ &\propto \lambda_g^{\alpha-1} e^{-\beta_g \lambda_g} \prod_{n=1}^N \lambda_g^{x_{g,n}} e^{-\lambda_g} \\ &\propto \text{Ga}(\lambda_g; \alpha', \beta') \end{aligned}$$

where

$$\begin{aligned} \alpha' &= \alpha + \sum_{n=1}^N x_{g,n} \\ \beta' &= \beta_g + N \end{aligned}$$

The second is,

$$\begin{aligned} p(\beta | \{\lambda_g\}_{g=1}^G; \alpha_0, \beta_0, \alpha) &\propto \text{Ga}(\beta; \alpha_0, \beta_0) \prod_{g=1}^G \text{Ga}(\lambda_g; \alpha, \beta) \\ &\propto \beta^{\alpha_0-1} e^{-\beta_0 \beta} \prod_{g=1}^G \beta^\alpha e^{-\beta \lambda_g} \\ &\propto \text{Ga}(\beta; \alpha', \beta') \end{aligned}$$

where

$$\begin{aligned} \alpha' &= \alpha_0 + \alpha G \\ \beta' &= \beta_0 + \sum_{g=1}^G \lambda_g \end{aligned}$$

**Problem 3:** Graphical models.

(a) Draw the graphical model corresponding to this joint probability distribution,

$$p(\{x_n, y_n\}_{n=1}^N; \alpha, \beta, \gamma) = p(x_1 | \alpha) \left[ \prod_{n=2}^N p(x_n | x_{n-1}; \beta) \right] \left[ \prod_{n=1}^N p(y_n | x_n; \gamma) \right].$$

**Solution:**

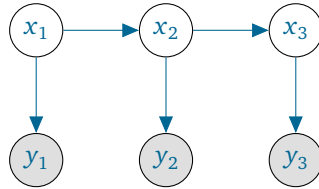


Figure 1: Note: This omits the hyperparameters, which could be drawn as black dots.

(b) Write the joint distribution corresponding to the graphical model in Figure 2.

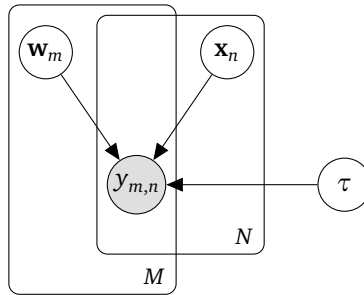


Figure 2

**Solution:**

$$p(\tau) \prod_{m=1}^M p(w_m) \prod_{n=1}^N p(x_n) \prod_{m=1}^M \prod_{n=1}^N p(y_{m,n} | w_m, x_n)$$

**Problem 4:** *Continuous latent variable models*

Canonical correlation analysis is a technique for paired datasets  $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$  where  $\mathbf{x}_n \in \mathbb{R}^{D_x}$  and  $\mathbf{y}_n \in \mathbb{R}^{D_y}$ . Like PCA, it can be viewed as a limiting case of a linear Gaussian model latent variable model,

$$\begin{aligned}\mathbf{z}_n &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ \mathbf{x}_n &\sim \mathcal{N}(\mathbf{W}_x \mathbf{z}_n + \mathbf{b}_x, \Sigma_x) \\ \mathbf{y}_n &\sim \mathcal{N}(\mathbf{W}_y \mathbf{z}_n + \mathbf{b}_y, \Sigma_y).\end{aligned}$$

Derive the conditional distribution  $p(\mathbf{y}_n | \mathbf{x}_n; \boldsymbol{\theta})$  where  $\boldsymbol{\theta} = (\mathbf{W}_x, \mathbf{W}_y, \mathbf{b}_x, \mathbf{b}_y, \Sigma_x, \Sigma_y)$ .

**Solution:** Note that

$$p(\mathbf{y}_n | \mathbf{x}_n; \boldsymbol{\theta}) = \int p(\mathbf{y}_n | \mathbf{z}_n; \boldsymbol{\theta}) p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta}) d\mathbf{z}_n$$

since  $\mathbf{y}_n \perp\!\!\!\perp \mathbf{x}_n | \mathbf{z}_n$ . The conditional in the integrand is,

$$\begin{aligned}p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta}) &\propto p(\mathbf{z}_n) p(\mathbf{x}_n | \mathbf{z}_n; \boldsymbol{\theta}) \\ &= \mathcal{N}(\mathbf{z}_n; \mathbf{0}, \mathbf{I}) \mathcal{N}(\mathbf{x}_n | \mathbf{W}_x \mathbf{z}_n + \mathbf{b}_x; \Sigma_x) \\ &\propto \exp\left\{-\frac{1}{2} \mathbf{z}_n^\top \mathbf{J}_z \mathbf{z}_n + \mathbf{h}_z^\top \mathbf{z}_n\right\}\end{aligned}$$

where

$$\begin{aligned}\mathbf{J}_z &= \mathbf{I} + \mathbf{W}_x^\top \Sigma_x^{-1} \mathbf{W}_x \\ \mathbf{h}_z &= \mathbf{W}_x^\top \Sigma_x^{-1} (\mathbf{x}_n - \mathbf{b}_x).\end{aligned}$$

Completing the square,  $p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}_z, \Sigma_z)$  where

$$\boldsymbol{\mu}_z = \mathbf{J}_z^{-1} \mathbf{h}_z = (\mathbf{I} + \mathbf{W}_x^\top \Sigma_x^{-1} \mathbf{W}_x)^{-1} \mathbf{W}_x^\top \Sigma_x^{-1} (\mathbf{x}_n - \mathbf{b}_x) \quad (1)$$

$$\Sigma_z = \mathbf{J}_z^{-1} = (\mathbf{I} + \mathbf{W}_x^\top \Sigma_x^{-1} \mathbf{W}_x)^{-1} \quad (2)$$

Finally, we obtain the predictive distribution by marginalizing over  $\mathbf{z}_n$  under this Gaussian conditional distribution,

$$p(\mathbf{y}_n | \mathbf{x}_n; \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}_n | \mathbf{W}_y \boldsymbol{\mu}_z + \mathbf{b}_y, \Sigma_y + \mathbf{W}_y \Sigma_z \mathbf{W}_y^\top).$$

An alternative way to obtain the same answer is to construct the joint multivariate normal distribution over  $(\mathbf{z}_n, \mathbf{x}_n, \mathbf{y}_n)$ , then use the rules of Gaussian marginalization and conditioning.

**Problem 5: The Bayesian Lasso**

The Lasso problem is an  $L_1$  penalized least squares problem,

$$\mathcal{L}(\mathbf{w}) = \sum_{n=1}^N \|y_n - \mathbf{x}_n^\top \mathbf{w}\|_2^2 + \lambda_0 \sum_{d=1}^D |w_d|. \quad (3)$$

From a Bayesian perspective, minimizing  $\mathcal{L}(\mathbf{w})$  is equivalent to *maximum a posteriori* (MAP) estimation in the following Bayesian model,

$$\begin{aligned} w_d &\stackrel{\text{iid}}{\sim} \text{Lap}(\lambda) \\ y_n &\stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{x}_n^\top \mathbf{w}, \sigma^2), \end{aligned} \quad (4)$$

where  $\text{Lap}(\lambda)$  denotes a Laplace distribution with density  $\text{Lap}(w; \lambda) = \frac{\lambda}{2} e^{-\lambda|w|}$ .

- (a) Find a setting of  $\lambda$  such that the MAP estimate of model (4) is the same as the minimizer of eq. 3. Your solution should be in terms of  $\lambda_0$  and  $\sigma^2$ .
- (b) The Laplace density can also be written as a *scale mixtures of Gaussians*,

$$\text{Lap}(w; \lambda) = \frac{\lambda}{2} e^{-\lambda|w|} = \int_0^\infty \mathcal{N}(w; 0, v) \cdot \text{Exp}(v; \frac{\lambda^2}{2}) dv = \int_0^\infty \frac{1}{\sqrt{2\pi v}} e^{-\frac{w^2}{2v}} \cdot \frac{\lambda^2}{2} e^{-\frac{\lambda^2 v}{2}} dv$$

Let  $\mathbf{y} = \{y_n\}_{n=1}^N$  and  $\mathbf{X} = \{\mathbf{x}_n\}_{n=1}^N$ . Use the integral representation above to write a joint distribution,

$$p(\mathbf{w}, \mathbf{v}, \mathbf{y} \mid \mathbf{X}; \lambda, \sigma^2)$$

on an extended space that includes the *augmentation variables*  $\mathbf{v} = (v_1, \dots, v_D)$ , such that the marginal distribution  $p(\mathbf{w}, \mathbf{y} \mid \mathbf{X}; \lambda, \sigma^2)$  matches that of the generative model described in eq. (4).

- (c) What algorithm would you use to perform Bayesian inference to approximate the posterior distribution  $p(\mathbf{w}, \mathbf{v} \mid \mathbf{X}, \mathbf{y}; \lambda, \sigma^2)$ ? Sketch out the steps involved.

**Solution:**

- (a) The log joint probability of the Bayesian model as a function of  $\mathbf{w}$  is,

$$\begin{aligned} \mathcal{J}(\mathbf{w}) &= \log p(\mathbf{w}, \mathbf{y} \mid \mathbf{X}) \\ &= \sum_{d=1}^D \log \text{Lap}(w_d; \lambda) + \sum_{n=1}^N \log \mathcal{N}(y_n; \mathbf{w}^\top \mathbf{x}_n, \sigma^2) + c \\ &= -\lambda \sum_{d=1}^D |w_d| - \sum_{n=1}^N \frac{1}{2\sigma^2} \|y_n - \mathbf{w}^\top \mathbf{x}_n\|_2^2 + c. \end{aligned}$$

Multiplying by  $2\sigma^2$  does not change the arg max of this objective (i.e. MAP estimate), so define,

$$\mathcal{J}'(\mathbf{w}) = -2\lambda\sigma^2 \sum_{d=1}^D |w_d| - \sum_{n=1}^N \|y_n - \mathbf{w}^\top \mathbf{x}_n\|_2^2 + c.$$

Note, if  $\lambda_0 = 2\lambda\sigma^2$  then  $\mathcal{J}'(\mathbf{w}) = -\mathcal{L}(\mathbf{w})$ . With this setting of  $\lambda_0$ , the Lasso estimate that minimizes  $\mathcal{L}(\mathbf{w})$  coincides the MAP estimate that maximizes  $\mathcal{J}'(\mathbf{w})$ .

(b) With this integral representation of the Laplace distribution,

$$\begin{aligned}
 p(\mathbf{w}, \mathbf{y} | \mathbf{X}) &= \prod_{d=1}^D \left[ \int_0^\infty \mathcal{N}(w_d; 0, \nu_d) \text{Exp}(\nu_d; \frac{\lambda^2}{2}) d\nu_d \right] \prod_{n=1}^N \mathcal{N}(y_n; \mathbf{w}^\top \mathbf{x}_n, \sigma^2) \\
 &= \int_0^\infty \cdots \int_0^\infty \prod_{d=1}^D \mathcal{N}(w_d; 0, \nu_d) \text{Exp}(\nu_d; \frac{\lambda^2}{2}) \prod_{n=1}^N \mathcal{N}(y_n; \mathbf{w}^\top \mathbf{x}_n, \sigma^2) d\nu_1 \dots d\nu_D \\
 &= \int p(\mathbf{w}, \boldsymbol{\nu}, \mathbf{y} | \mathbf{X}) d\boldsymbol{\nu}
 \end{aligned}$$

where

$$p(\mathbf{w}, \boldsymbol{\nu}, \mathbf{y} | \mathbf{X}) = \prod_{d=1}^D \mathcal{N}(w_d; 0, \nu_d) \text{Exp}(\nu_d; \frac{\lambda^2}{2}) \prod_{n=1}^N \mathcal{N}(y_n; \mathbf{w}^\top \mathbf{x}_n, \sigma^2)$$

(c) The simple answer is you could do HMC, but you could have done that on the original model anyway. (There's a slight complication with the prior being non-differentiable at  $w_d = 0$ , but you could throw HMC at it anyway.) However, this augmentation scheme presents an opportunity to introduce a pretty cool trick. It turns out that in the augmented model, the conditional distributions of  $\mathbf{w}$  and  $\boldsymbol{\nu}$  are tractable, and we can use them to perform Gibbs sampling. We have,

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \boldsymbol{\nu}) = \mathcal{N}(\mathbf{w}; \mathbf{J}^{-1} \mathbf{h}, \mathbf{J}^{-1})$$

where

$$\begin{aligned}
 \mathbf{J} &= \text{diag}(\boldsymbol{\nu})^{-1} + \frac{1}{\sigma^2} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \\
 \mathbf{h} &= \frac{1}{\sigma^2} \sum_{n=1}^N \mathbf{x}_n y_n.
 \end{aligned}$$

The conditional distribution of  $\nu_d$  is,

$$\begin{aligned}
 p(\nu_d | w_d) &\propto \mathcal{N}(w_d; 0, \nu_d) \text{Exp}(\nu_d; \frac{\lambda^2}{2}) \\
 &\propto \frac{1}{\sqrt{\nu_d}} \exp \left\{ -\frac{w_d^2}{2\nu_d} - \frac{\lambda^2 \nu_d}{2} \right\}
 \end{aligned}$$

This is a strange looking conditional... in exponential family form, its sufficient statistics are  $\log \nu_d$ ,  $\nu_d^{-1}$ , and  $\nu_d$ .

You wouldn't be expected to know this for the actual final, but it turns out that this conditional distribution can be rewritten as a (somewhat) standard density. Instead of the variance  $\nu_d$ , let's do



a change of variables to consider the conditional distribution of the precision  $\tau_d = \nu_d^{-1}$ ,

$$\begin{aligned}
 p(\tau_d | w_d) &\propto \left| \frac{d\tau_d^{-1}}{d\tau_d} \right| \mathcal{N}(w_d; 0, \tau_d^{-1}) \text{Exp}(\tau_d^{-1}; \frac{\lambda^2}{2}) \\
 &\propto \frac{1}{\tau_d^2} \sqrt{\tau_d} \exp \left\{ -\frac{1}{2} \left( \tau_d w_d^2 + \frac{\lambda^2}{\tau_d} \right) \right\} \\
 &\propto \tau_d^{-\frac{3}{2}} \exp \left\{ -\frac{w_d^2}{2\tau_d} \left( \tau_d^2 + \frac{\lambda^2}{w_d^2} \right) \right\} \\
 &\propto \text{IG} \left( \tau_d; \frac{\lambda}{w_d}, \lambda^2 \right)
 \end{aligned}$$

where  $\text{IG}(x; \mu, \kappa)$  denotes the *inverse Gaussian* distribution with mean  $\mu$  and shape  $\kappa$ . This isn't a distribution that came up in class, but it appears in the context of stochastic processes. For example, the time a Brownian motion with positive drift takes to reach a fixed positive level is inverse-Gaussian distributed.

**Problem 6: Mixture Models**

Consider the following *zero-inflated Poisson regression* model where  $w, x_n \in \mathbb{R}_+$ ,  $y_n \in \mathbb{N}$ , and  $z_n \in \{0, 1\}$ ,

$$\begin{aligned} w &| \alpha, \beta \sim \text{Gamma}(\alpha, \beta) \\ z_n &| \gamma \stackrel{\text{iid}}{\sim} \text{Bern}(\gamma) \\ y_n &| x_n, z_n, w \stackrel{\text{iid}}{\sim} \text{Poisson}(wx_n z_n). \end{aligned}$$

- (a) Sketch the probability mass function of the marginal distribution  $p(y_n | x_n, w, \gamma)$  for  $\gamma \in \{0, 0.5, 1\}$ , assuming  $wx_n = 5$ . What is  $p(y_n = 0 | x_n, w, \gamma)$ ? (Note:  $0! = 1$  and  $0^0 = 1$ .)
- (b) Compute the conditional distribution  $p(z_n = 1 | y_n, x_n, w, \gamma)$ .
- (c) Compute the expected log probability,

$$\mathcal{L}(w) = \mathbb{E}_{p(z|y,x,w',\gamma)} \left[ \log p(\{y_n, x_n, z_n\}_{n=1}^N, w | \alpha, \beta, \gamma) \right],$$

where  $w'$  denotes a fixed weight. For notational simplicity, let  $q_n \triangleq p(z_n = 1 | y_n, x_n, w', \gamma)$  denote the solution to part (c), and drop terms in  $\mathcal{L}(w)$  that are constant with respect to  $w$ .

- (d) Assume  $\alpha > 1$ . Solve for  $w^* = \arg \max \mathcal{L}(w)$  using the fact that the mode of the  $\text{Gamma}(a, b)$  distribution is at  $(a - 1)/b$  when  $a > 1$ .

**Solution:**

- (a) The marginal distribution is,

$$\begin{aligned} p(y_n | wx_n = 5; \gamma) &= \sum_{z_n} \text{Po}(y_n; wx_n z_n) p(z_n; \gamma) \\ &= \gamma \text{Po}(y_n; 5) + (1 - \gamma) \text{Po}(y_n; 0) \\ &= \gamma \text{Po}(y_n; 5) + (1 - \gamma) \mathbb{I}[y_n = 0] \end{aligned}$$

It looks like a re-scaled Poisson pmf with an extra mass of  $(1 - \gamma)$  added to  $p(y_n = 0 | -)$ .

- (b) The conditional is,

$$\begin{aligned} p(z_n = 1 | y_n, x_n, w; \gamma) &\propto p(z_n = 1; \gamma) p(y_n | z_n = 1, x_n, w) \\ &\propto \gamma \text{Po}(y_n; wx_n) \end{aligned}$$

Likewise,

$$\begin{aligned} p(z_n = 0 | y_n, x_n, w; \gamma) &\propto p(z_n = 0; \gamma) p(y_n | z_n = 0, x_n, w) \\ &\propto (1 - \gamma) \text{Po}(y_n; 0) \\ &\propto (1 - \gamma) \mathbb{I}[y_n = 0]. \end{aligned}$$

Normalizing yields,

$$p(z_n = 1 | y_n, x_n, w; \gamma) = \frac{\gamma \text{Po}(y_n; wx_n)}{\gamma \text{Po}(y_n; wx_n) + (1 - \gamma) \mathbb{I}[y_n = 0]}$$

(c) The expected log probability is,

$$\begin{aligned}
\mathcal{L}(w) &= \mathbb{E}_{p(z|y,x,w',\gamma)} \left[ \log p(\{y_n, x_n, z_n\}_{n=1}^N, w \mid \alpha, \beta, \gamma) \right] \\
&= \mathbb{E}_{p(z|y,x,w',\gamma)} \left[ \log p(w; \alpha, \beta) + \sum_{n=1}^N \log p(y_n \mid w, x_n, z_n) \right] + c \\
&= \log \text{Ga}(w; \alpha, \beta) + \sum_{n=1}^N \mathbb{E}_{p(z_n|-)} [\log \text{Po}(y_n \mid wx_n z_n)] + c \\
&= (\alpha - 1) \log w - \beta w + \sum_{n=1}^N \mathbb{E}_{p(z_n|-)} [y_n \log(wx_n z_n) - wx_n z_n] + c \\
&= (\alpha - 1) \log w - \beta w + \sum_{n=1}^N y_n \log w - wx_n q_n + c \\
&= \log \text{Ga}(w; \alpha', \beta')
\end{aligned}$$

where

$$\begin{aligned}
\alpha' &= \alpha + \sum_{n=1}^N y_n \\
\beta' &= \beta + \sum_{n=1}^N x_n q_n.
\end{aligned}$$

(d) The mode of the gamma is at

$$w^* = \frac{\alpha' - 1}{\beta'} = \frac{\alpha + \sum_{n=1}^N y_n - 1}{\beta + \sum_{n=1}^N x_n q_n}$$

As a sanity check, consider the case where  $x_n = 1$  for all  $n$ . Then the maximum likelihood estimate of  $w$  under a standard Poisson model (equivalently, with  $z_n = 1$  for all  $n$ ), would be  $w_{\text{MLE}} = \frac{1}{N} \sum_{n=1}^N y_n$ . With the zero-inflated model, some of the observations where  $y_n = 0$  are attributed to  $z_n$  being zero. Thus, we should expect our estimate of  $w$  to be a bit higher. Indeed, under an uninformative prior ( $\alpha = 1, \beta = 0$ ), we would have  $w^* = (\sum_n y_n) / (\sum_n q_n) \geq w_{\text{MLE}}$ , since the denominator is at most  $N$ .

**Problem 7: Mixed Membership Models**

Latent Dirichlet allocation (LDA) corresponds to the following generative model,

$$\begin{aligned}
 \boldsymbol{\eta}_k &\sim \text{Dir}(\boldsymbol{\phi}) && \text{for } k = 1, \dots, K \\
 \boldsymbol{\pi}_n &\sim \text{Dir}(\boldsymbol{\alpha}) && \text{for } n = 1, \dots, N \\
 z_{n,\ell} &\sim \text{Cat}(\boldsymbol{\pi}_n) && \text{for } n = 1, \dots, N; \ell = 1, \dots, L \\
 x_{n,\ell} &\sim \boldsymbol{\eta}_{z_{n,\ell}} && \text{for } n = 1, \dots, N; \ell = 1, \dots, L
 \end{aligned}$$

where  $\boldsymbol{\eta}_k \in \Delta_V$  are the *topics* (i.e. distributions over words) and  $\boldsymbol{\pi}_n \in \Delta_K$  are the *topic proportions* (i.e. distributions over topics).

However, this model fails to capture correlations in the topic proportions; for example, that a “finance” topic and a “government” topic may often co-occur in the same document. *Correlated topic models* address this limitation by replacing the Dirichlet prior on  $\boldsymbol{\pi}_n$  with a logistic normal prior,

$$\begin{aligned}
 \boldsymbol{\pi}_n &= \text{softmax}(\mathbf{u}_n) = \left[ \frac{e^{u_{n1}}}{1 + \sum_{k=1}^{K-1} e^{u_{nk}}}, \dots, \frac{e^{u_{n,K-1}}}{1 + \sum_{k=1}^{K-1} e^{u_{nk}}}, \frac{1}{1 + \sum_{k=1}^{K-1} e^{u_{nk}}} \right]^\top \\
 \mathbf{u}_n &\sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})
 \end{aligned}$$

where  $\mathbf{u}_n \in \mathbb{R}^{K-1}$ . The correlations in  $\mathbf{u}_n$  due to the multivariate normal prior induce correlations in  $\boldsymbol{\pi}_n$  as well.

- (a) Without doing any math, sketch the density of  $\boldsymbol{\pi}_n \in \Delta_3$  when  $\boldsymbol{\mu} = [0, 0]^\top$  and  $\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ . Do the same for  $\boldsymbol{\mu} = [0, 0]^\top$  and  $\boldsymbol{\Sigma} = \begin{bmatrix} 1 & \frac{1}{2} \\ \frac{1}{2} & 1 \end{bmatrix}$ . Explain your reasoning.
- (b) Try to derive CAVI updates for this model. Where do you run into trouble and why?

**Solution:**

- (a) See the following figure from from Blei, David and John Lafferty. "Correlated topic models." *Advances in Neural Information Processing Systems* 18 (2005). When the first two coordinates are

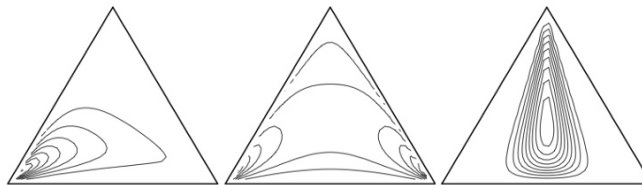


Figure 1: Top: Graphical model representation of the correlated topic model. The logistic normal distribution, used to model the latent topic proportions of a document, can represent correlations between topics that are impossible to capture using a single Dirichlet. Bottom: Example densities of the logistic normal on the 2-simplex. From left: diagonal covariance and nonzero-mean, negative correlation between components 1 and 2, positive correlation between components 1 and 2.

uncorrelated, the distribution is essentially symmetric. (The figure here shows what happens when the mean is nonzero — it's pulled toward one of the simplex vertices.) When the two coordinates of the Gaussian are anti-correlated, you end up with a multi-modal distribution on the simplex with modes toward one vertex or another. When the two coordinates are positively correlated, the mass is pulled toward the midline. That's because both Gaussian coordinates are large and positive, it maps to  $\pi_n = [0.5, 0.5, 0]$ ; when both coordinates are negative, it maps to  $\pi_n = [0, 0, 1]$ .

- (b) We can work with either  $\mathbf{u}_n$  or  $\pi_n$ . Since the prior is specified on  $\mathbf{u}_n$ , let's go with that. The problem is that we have a Gaussian prior and a categorical likelihood. The optimal CAVI updates are still,

$$\begin{aligned} q(\mathbf{u}_n) &\propto \exp \left\{ \mathbb{E}_{q(z_n)} [\log p(\mathbf{u}_n) + \log p(\mathbf{z}_n | \mathbf{u}_n)] \right\} \\ &\propto \mathcal{N}(\mathbf{u}_n; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \exp \left\{ \sum_{l=1}^L \mathbb{E}_{q(z_{n,l})} [\log \text{Cat}(z_{n,l}; \text{softmax}(\mathbf{u}_n))] \right\} \\ &\propto \mathcal{N}(\mathbf{u}_n; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \exp \left\{ \sum_{l=1}^L \sum_{v=1}^V q_{n,l,v} \left[ u_{n,v} - \log \left( 1 + \sum_{v'=1}^V e^{u_{n,v'}} \right) \right] \right\} \end{aligned}$$

where  $q_{n,l,v} = q(z_{n,l} = v)$ . The problem is that the log-sum-exp of  $\mathbf{u}_n$  in the exponent is not conjugate with the Gaussian prior.

**Problem 8: Variational autoencoders**

Consider the following *deep mixture model*,

$$\begin{aligned}z_n &\sim \pi \\ \mathbf{x}_n &\sim \mathcal{N}(\boldsymbol{\mu}_{z_n}, \boldsymbol{\Sigma}_{z_n}) \\ \mathbf{y}_n &\sim \mathcal{N}(f(\mathbf{x}_n; \mathbf{w}), \sigma^2 \mathbf{I})\end{aligned}$$

where  $z_n \in \{1, \dots, K\}$  is a discrete latent variable,  $\mathbf{x}_n \in \mathbb{R}^M$  is a continuous latent variable,  $\mathbf{y}_n \in \mathbb{R}^D$  is an observed data point, and  $f : \mathbb{R}^M \mapsto \mathbb{R}^D$  is a neural network with weights  $\mathbf{w}$ . The generative model parameters are  $\boldsymbol{\theta} = (\pi, \{\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K, \mathbf{w})$ .

- (a) Suppose you wanted to perform fixed form variational inference to approximate the posterior,  $p(z_n, \mathbf{x}_n \mid \mathbf{y}_n; \boldsymbol{\theta}) \approx q(z_n, \mathbf{x}_n; \boldsymbol{\phi})$ , with variational parameters  $\boldsymbol{\phi}$ . What challenges might you encounter when trying to maximize the local ELBO,  $\mathcal{L}_n(\boldsymbol{\theta}, \boldsymbol{\phi})$ , using stochastic gradient ascent and the reparameterization trick (i.e. the pathwise gradient estimator)?
- (b) Suggest an alternative to the reparameterization trick that could allow you to fit  $\boldsymbol{\theta}$  and  $\boldsymbol{\phi}$ . What challenges might this alternative present?
- (c) Rewrite the generative model by marginalizing over  $z_n$  to obtain a collapsed model  $p(\mathbf{x}_n, \mathbf{y}_n; \boldsymbol{\theta})$ , and assume a variational posterior  $q(\mathbf{x}_n; \boldsymbol{\phi})$ . Can you use the reparameterization trick now?

**Solution:**

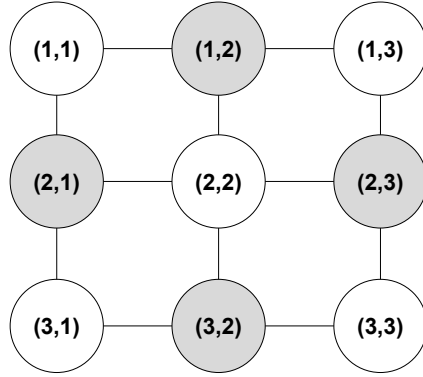
- (a) Since  $z_n$  is a discrete variable, we cannot reparameterize it as a differentiable transformation of “noise.” That prevents us from using the standard pathwise gradient estimator/reparameterization tricks.
- (b) We could use the score-function gradient estimator, since that works for discrete latent variables. However, without good control variates it can lead to higher variance estimates. Alternatively, we could use a continuous relaxation of the discrete variables, as in the Concrete or Gumbel-softmax approximations for discrete VAEs.
- (c) Marginalizing out  $z_n$  yields,

$$p(\mathbf{x}_n, \mathbf{y}_n; \boldsymbol{\theta}) = \left( \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \mathcal{N}(\mathbf{y}_n \mid f(\mathbf{x}_n; \mathbf{w}), \sigma^2 \mathbf{I}).$$

If we then choose a Gaussian variational posterior for  $q(\mathbf{x}_n; \boldsymbol{\phi})$ , then we **can** use the reparameterization trick since there are no longer any discrete variables to worry about.

**Problem 9: State space models**

In class we studied state space models for sequential data, like hidden Markov models and linear dynamical systems. Here we will consider similar models for 2-dimensional data. Suppose we observe an image  $\mathbf{y} \in \mathbb{R}^{H \times W}$  which we believe to be a noisy version of an underlying binary image  $\mathbf{x} \in \{0, 1\}^{H \times W}$ . Given  $\mathbf{y}$ , we wish to recover the true image  $\mathbf{x}$  which it was derived from. We formulate this as a probabilistic inference problem. We will assume the image is square and start by constructing a graph which connects neighboring pixels. The graph for  $H = W = 3$  is shown below, with the node labels corresponding to the indices in the vectors  $\mathbf{y}$  and  $\mathbf{x}$ .



Our prior on  $\mathbf{x}$  will be given as an *Ising model*, which encodes our belief that nearby pixels are likely to be similar:

$$p(\mathbf{x}) = \frac{1}{Z(\theta)} \prod_{(ij,kl) \in \mathcal{E}} \psi_{\theta}(x_{ij}, x_{kl})$$

Here,  $\mathcal{E}$  is the edge set of the pixel graph,  $Z(\theta)$  is a normalizing constant, and  $\psi_{\theta} : \{0, 1\} \times \{0, 1\} \rightarrow \mathbb{R}_{++}$  is defined by:

$$\psi_{\theta}(x_{ij}, x_{kl}) = \begin{cases} e^{\theta} & , x_{ij} = x_{kl} \\ 1 & , x_{ij} \neq x_{kl} \end{cases}$$

where  $\theta > 0$  is a hyperparameter. We assume a Gaussian noise model, which gives us a likelihood over  $\mathbf{y}$  given  $\mathbf{x}$  as:

$$p(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^H \prod_{j=1}^W \mathcal{N}(y_{ij} | x_{ij}, \sigma^2)$$

where  $\sigma^2$  is a hyperparameter. Given  $\mathbf{y}$ , we will obtain our denoised image by sampling from the posterior  $p(\mathbf{x} | \mathbf{y})$  using Gibbs sampling

- (a) Given a pixel  $(i, j)$ , let  $\mathcal{E}(i, j)$  denote its neighbors in the pixel graph. Similarly, given  $\mathbf{x} \in \{0, 1\}^{H \times W}$ , let  $N_1(i, j, \mathbf{x}_{-ij}) = \sum_{(k,l) \in \mathcal{E}(i,j)} x_{k,l}$  denote the number of neighbors of pixel  $(i, j)$  set to 1 and let  $N_0(i, j, \mathbf{x}_{-ij}) = \sum_{(k,l) \in \mathcal{E}(i,j)} 1 - x_{kl}$  denote the number of neighbors of pixel  $(i, j)$  set to 0.

Show that the complete conditional of  $x_{ij}$  is given by:

$$p(x_{ij} = 1 \mid \mathbf{x}_{-ij}, \mathbf{y}) = \frac{e^{\phi_1}}{e^{\phi_0} + e^{\phi_1}}$$

where

$$\begin{aligned}\phi_1 &= \theta N_1(i, j, \mathbf{x}_{-ij}) + \log \mathcal{N}(y_{ij} \mid x_{ij} = 1, \sigma^2) \\ \phi_0 &= \theta N_0(i, j, \mathbf{x}_{-ij}) + \log \mathcal{N}(y_{ij} \mid x_{ij} = 0, \sigma^2)\end{aligned}$$

- (b) Suppose we also incorporate a prior  $p(\theta)$  on  $\theta$ , e.g.  $p(\theta) = \text{Gamma}(\theta; \alpha, \beta)$ . It is not possible to derive a closed form for  $\theta$ 's complete conditional  $p(\theta \mid \mathbf{x}, \mathbf{y})$ . Explain what we may do instead to approximately sample from this conditional. Why might this be computationally challenging for large images (i.e. when  $D$  is large)?
- (c) [Bonus] Consider the pixel graph, and let  $\mathcal{S}$  be a maximal set of nodes such that  $\mathcal{E}(i, j) \cap \mathcal{S} = \emptyset$  for all  $(i, j) \in \mathcal{S}$ . For the example graph, we could use  $\mathcal{S}$  as the shaded set of nodes, so  $\mathcal{S} = \{(1, 2), (2, 1), (2, 3), (3, 2)\}$ . Explain why the random variables  $\{x_{ij} : (i, j) \in \mathcal{S}\}$  are independent given  $\mathbf{x}_{-\mathcal{S}}, \mathbf{y}, \theta$  and how we can exploit this for an efficient parallel block Gibbs update.

**Solution:**

(a) We have,

$$\begin{aligned}p(x_{ij} \mid \mathbf{x}_{-ij}, \mathbf{y}) &\propto \left( \prod_{k,l \in \mathcal{E}(i,j)} \psi_{\theta}(x_{ij}, x_{kl}) \right) \mathcal{N}(y_{ij} \mid x_{ij}, \sigma^2) \\ &\propto \left( \prod_{k,l \in \mathcal{E}(i,j)} e^{\theta \mathbb{I}[x_{ij}=x_{kl}]} \right) \mathcal{N}(y_{ij} \mid x_{ij}, \sigma^2) \\ &\propto e^{\theta \sum_{k,l \in \mathcal{E}(i,j)} \mathbb{I}[x_{ij}=x_{kl}]} \mathcal{N}(y_{ij} \mid x_{ij}, \sigma^2)\end{aligned}$$

For the two cases we have,

$$\begin{aligned}p(x_{ij} = 1 \mid \mathbf{x}_{-ij}, \mathbf{y}) &\propto e^{\theta N_1(i,j,\mathbf{x}_{-ij})} \mathcal{N}(y_{ij} \mid 1, \sigma^2) \\ p(x_{ij} = 0 \mid \mathbf{x}_{-ij}, \mathbf{y}) &\propto e^{\theta N_0(i,j,\mathbf{x}_{-ij})} \mathcal{N}(y_{ij} \mid 0, \sigma^2)\end{aligned}$$

Normalizing and writing in terms of the logs yields the desired answer.

(b) The hard part about inferring  $\theta$  is that it appears in the log normalizer,

$$Z(\theta) = \sum_{\mathbf{x} \in \{0,1\}^{H \times W}} \prod_{(ij,kl) \in \mathcal{E}} \psi_{\theta}(x_{ij}, x_{kl}),$$

and the log normalizer consists of a sum with  $2^{HW}$  terms. For small lattices we can enumerate all these terms, evaluate  $Z(\theta)$ , and do inference with algorithms like HMC. However, this won't scale to even moderately sized grids.