

## 1 Problem 1

For this problem, I will call the random variables  $A$ ,  $B$ ,  $C$  and  $D$  for simpler notation. The graph structure  $A - B - C - D - A$  implies that the distribution can be factorized as

$$p(a, b, c, d) \propto \phi_{AB}(a, b) \phi_{BC}(b, c) \phi_{CD}(c, d) \phi_{DA}(d, a) \quad . \quad (1)$$

Therefore, in order for  $p(0, 1, 1, 0) = 0$ , at least one of

$$\phi_{AB}(0, 1), \quad \phi_{BC}(1, 1), \quad \phi_{CD}(1, 0), \quad \text{or} \quad \phi_{DA}(0, 0) \quad (2)$$

must be equal to zero. This is clearly not true since

$$(i) p(0, 1, 1, 1) \neq 0, \quad (ii) p(1, 1, 1, 0) \neq 0 \quad \text{and} \quad (iii) p(0, 0, 0, 0) \neq 0 \quad . \quad (3)$$

Statement (i) implies that none of

$$\phi_{AB}(0, 1), \quad \phi_{BC}(1, 1), \quad \phi_{CD}(1, 1), \quad \text{or} \quad \phi_{DA}(1, 0) \quad (4)$$

can equal zero. Similarly, statement (ii) implies that none of

$$\phi_{AB}(1, 1), \quad \phi_{BC}(1, 1), \quad \phi_{CD}(1, 0), \quad \text{or} \quad \phi_{DA}(0, 1) \quad (5)$$

can equal zero. And finally, statement (iii) implies that none of

$$\phi_{AB}(0, 0), \quad \phi_{BC}(0, 0), \quad \phi_{CD}(0, 0), \quad \text{or} \quad \phi_{DA}(0, 0) \quad (6)$$

are zero. Therefore, by this contradicting example, the distribution cannot be factorized over the given graph structure.

## 2 Problem 2

(a) A probability distribution over three discrete random variables  $A$ ,  $B$  and  $C$  is parameterized as

$$p(a, b, c) \propto \exp(-\epsilon_1(a, b) - \epsilon_2(b, c)) \quad . \quad (7)$$

If we redefine

$$\epsilon'_1(a, B = b^i) \leftarrow \epsilon_1(a, B = b^i) + \lambda^i \quad \text{and} \quad \epsilon'_2(B = b^i, c) \leftarrow \epsilon_2(B = b^i, c) - \lambda^i \quad (8)$$

then the new distribution is

$$p(a, b^i, c) \propto \exp(-\epsilon'_1(a, b^i) - \epsilon'_2(b^i, c)) \quad (9)$$

$$= \exp(-\epsilon_1(a, b^i) - \lambda^i - \epsilon_2(b^i, c) + \lambda^i) \quad (10)$$

$$= \exp(-\epsilon_1(a, b^i) - \epsilon_2(b^i, c)) \quad (11)$$

which is equivalent to the original parameterization. Therefore, any symmetric reparameterization of the energy functions will leave the distribution unchanged.

(b) We would like to find  $w'_{ij}$  and  $u'_i$  such that

$$p_{\text{Ising}} \propto \exp \left( - \sum_{i < j \in E} w'_{ij} z_i z_j - \sum_i u'_i z_i \right) \quad (12)$$

for  $z_i = \pm 1$  is equivalent to

$$p_{\text{Boltzmann}} \propto \exp \left( - \sum_{i,j \in E} w_{ij} x_i x_j - \sum_i u_i x_i \right) \quad (13)$$

for  $x_i \in \{0, 1\}$ . This can be easily achieved with the mapping  $x_i \leftarrow (z_i + 1)/2$ . Substituting this into Equation (13), we find

$$p_{\text{Boltzmann}} \propto \exp \left( - \sum_{i,j \in E} \frac{w_{ij}}{4} (z_i z_j + z_i + z_j + 1) - \sum_i \frac{u_i}{2} (z_i + 1) \right) \quad (14)$$

$$= \exp \left( - \sum_{i,j \in E} \frac{w_{ij}}{4} z_i z_j - \sum_{i,j \in E} \frac{w_{ij}}{4} (z_i + z_j) - \sum_i \frac{u_i}{2} z_i - \sum_{i,j \in E} \frac{w_{ij}}{4} - \sum_i \frac{u_i}{2} \right) \quad (15)$$

The last two terms in this equation are constant so they can be absorbed into the partition function, leaving

$$p_{\text{Boltzmann}} \propto \exp \left( - \sum_{i,j \in E} \frac{w_{ij}}{4} z_i z_j - \sum_i \left[ \frac{u_i}{2} + \sum_{j \in E_i} \frac{w_{ij}}{4} \right] z_i \right) \quad (16)$$

where  $E_i$  is the set of edges containing node  $i$ . Therefore, the correct mapping between the two distributions is

$$w'_{ij} = \frac{w_{ij}}{4} \quad \text{and} \quad u'_i = \frac{u_i}{2} + \sum_{j \in E_i} \frac{w_{ij}}{4} \quad . \quad (17)$$

### 3 Problem 3

For a simple (non-pairwise) distribution on 3 random variables  $A$ ,  $B$  and  $C$ , factored according to

$$p(a, b, c) \propto \phi_{ABC}(a, b, c) \quad , \quad (18)$$

we can introduce a new variable  $D$  to convert it to the pairwise

$$p(a, b, c, d) \propto \phi_{AD}(a, d) \phi_{BD}(b, d) \phi_{CD}(c, d) \quad (19)$$

where

$$p(a, b, c) = \sum_i p(a, b, c, D = d^i) \quad (20)$$

$$\phi_{ABC}(a, b, c) \propto \sum_i \phi_{AD}(a, D = d^i) \phi_{BD}(b, D = d^i) \phi_{CD}(c, D = d^i) \quad . \quad (21)$$

To determine the forms of the pairwise potentials, first, we can assert that  $D$  assumes a tuple value  $(d_1, d_2, \dots)$  with one entry for each connected node (i.e. three in this example for  $A$ ,  $B$ , and  $C$ ). Then, the potentials will be

$$\phi_{AD}(a, d) = \phi_{ABC}(a, d_2, d_3) \quad (22)$$

with similar forms for the other edges.

Following this example, the general procedure for *one* particular non-pairwise potential will be

$$\phi_X(\mathbf{x}) \rightarrow \prod_{i=1}^N \phi_{X_i, Y}(x_i, Y) \quad (23)$$

where  $Y$  is an  $N$ -tuple and

$$\phi_{X_i, Y}(x_i, Y) = \phi_X(y_1, \dots, y_{i-1}, x_i, y_{i+1}, \dots, y_N) \quad . \quad (24)$$

## 4 Problem 4: Exponential Families

(a)

1. A multivariate Gaussian with identity covariance in  $K$  dimensions is part of the exponential family:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, I) = \frac{1}{(2\pi)^{K/2}} \exp \left( -\frac{1}{2} \sum_{i=1}^K [x_i - \mu_i]^2 \right) \quad (25)$$

$$= \frac{1}{(2\pi)^{K/2}} \exp \left( \sum_{i=1}^K \mu_i x_i - \frac{1}{2} \sum_{i=1}^K x_i^2 - \frac{1}{2} \sum_{i=1}^K \mu_i^2 \right) \quad . \quad (26)$$

Therefore, setting  $\mathbf{f}(\mathbf{x}) = (x_1, x_1^2, \dots, x_K, x_K^2)^T$ ,  $\boldsymbol{\eta} = (\mu_1, -1/2, \dots, \mu_K, -1/2)^T$ ,  $\ln Z = \sum \mu_i^2/2 + K \ln(2\pi)/2$  and  $h(\mathbf{x}) = 1$  puts this distribution in the correct form.

2. The Dirichlet distribution in  $K$  dimensions is

$$D(\boldsymbol{\theta}; \boldsymbol{\alpha}) = \frac{1}{Z(\boldsymbol{\alpha})} \prod_{i=1}^K \theta_i^{\alpha_i - 1} \quad (27)$$

where

$$Z(\boldsymbol{\alpha}) = \frac{\prod_{i=1}^K \Gamma(\alpha_i)}{\Gamma \left( \sum_{i=1}^K \alpha_i \right)} \quad . \quad (28)$$

This can be rewritten as

$$D(\boldsymbol{\theta}; \boldsymbol{\alpha}) = \exp \left( \sum_{i=1}^K [1 - \alpha_i] \ln \theta_i - \ln Z(\boldsymbol{\alpha}) \right) \quad . \quad (29)$$

This is clearly a member of the exponential family with  $\mathbf{f}(\boldsymbol{\theta}) = (\ln \theta_1, \dots, \ln \theta_K)^T$ ,  $\boldsymbol{\eta} = (1 - \alpha_1, \dots, 1 - \alpha_K)^T$  and  $h(\boldsymbol{\theta}) = 1$ .

3. The log-normal distribution is parameterized as

$$\mathcal{L}(y; 0, \sigma^2) = \left| \frac{d \ln y}{dy} \right| \mathcal{N}(\ln y; 0, \sigma^2) \quad (30)$$

$$= \frac{1}{y\sqrt{2\pi\sigma^2}} \exp\left(-\frac{[\ln y]^2}{2\sigma^2}\right) . \quad (31)$$

Setting  $f(y) = (\ln y)^2$ ,  $\eta = -1/2\sigma^2$ ,  $h(y) = y^{-1}$  and  $Z = 1/\sqrt{2\pi\sigma^2}$  shows that this is also a member of the exponential family.

4. The Boltzmann distribution can be written as

$$p(\mathbf{x}) = \frac{1}{Z} \exp\left(\sum_i u_i x_i + \sum_{i,j \in E} w_{ij} x_i x_j\right) \quad (32)$$

$$= \frac{1}{Z} \exp\left(\sum_i u_i x_i + \sum_{i,j \in E} w_{ij} x_i x_j - \ln Z\right) . \quad (33)$$

Therefore, we can set  $\boldsymbol{\eta} = \{\mathbf{u}, \mathbf{w}\}$  where  $\mathbf{w} = \{w_{ij}, \forall (i, j) \in E\}$  and  $\mathbf{f}(\mathbf{x}) = \{\mathbf{x}, \boldsymbol{\xi}\}$  where  $\boldsymbol{\xi} = \{x_i x_j, \forall (i, j) \in E\}$  to show that this distribution can also be written in the form of a member of the exponential family.

(b) For a continuous distribution, the partition function is given by

$$Z(\boldsymbol{\eta}) = \int h(\mathbf{x}) \exp(\boldsymbol{\eta} \cdot \mathbf{f}(\mathbf{x})) d\mathbf{x} \quad (34)$$

and the gradient with respect to  $\boldsymbol{\eta}$  is

$$\nabla_{\boldsymbol{\eta}} Z(\boldsymbol{\eta}) = \int \mathbf{f}(\mathbf{x}) h(\mathbf{x}) \exp(\boldsymbol{\eta} \cdot \mathbf{f}(\mathbf{x})) d\mathbf{x} . \quad (35)$$

Therefore, since

$$\nabla_{\boldsymbol{\eta}} \ln Z(\boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} \nabla_{\boldsymbol{\eta}} Z(\boldsymbol{\eta}) \quad (36)$$

$$= \int \mathbf{f}(\mathbf{x}) p(\mathbf{x}; \boldsymbol{\eta}) \quad (37)$$

where

$$p(\mathbf{x}; \boldsymbol{\eta}) = h(\mathbf{x}) \exp(\boldsymbol{\eta} \cdot \mathbf{f}(\mathbf{x}) - \ln Z(\boldsymbol{\eta})) , \quad (38)$$

the result is immediately clear

$$\nabla_{\boldsymbol{\eta}} \ln Z(\boldsymbol{\eta}) = \int \mathbf{f}(\mathbf{x}) p(\mathbf{x}; \boldsymbol{\eta}) = E_{p(\mathbf{x}; \boldsymbol{\eta})} [\mathbf{f}(\mathbf{x})] . \quad (39)$$

(c) For the multivariate Gaussian in example 1, the log-partition function becomes

$$\ln Z = \frac{1}{2} \sum_{i=1}^K \left[ \frac{\mu_i^2}{\sigma_i^2} + \ln \sigma_i^2 + \ln 2\pi \right] \quad (40)$$

when we introduce a diagonal covariance tensor  $\Sigma = \text{Diag}(\sigma_1^2, \sigma_2^2, \dots)$ . Also,  $\boldsymbol{\eta} \leftarrow \{\mu_i/\sigma_i^2, -1/2\sigma_i^2\} = \{\alpha_i, \beta_i\}$  so the log partition function can be re-written

$$\ln Z = \sum_i \left[ -\frac{\alpha_i^2}{4\beta_i} + \ln \left( -\frac{1}{2\beta_i} \right) + \ln 2\pi \right] \quad (41)$$

Therefore, the derivative of  $\ln Z$  with respect to  $\alpha_i = \mu_i/\sigma_i^2$  is

$$\frac{d \ln Z}{d \alpha_i} = -\frac{\alpha_i}{2\beta_i} = \mu_i = E[x_i] \quad (42)$$

as expected and

$$\frac{d \ln Z}{d \beta_i} = \frac{\alpha_i^2}{4\beta_i^2} - \frac{1}{2\beta_i} = \mu_i^2 + \sigma_i^2 = E[x_i^2] \quad (43)$$

This verifies the claim that the gradient of the log partition function gives the expectation values of  $\mathbf{f}(\mathbf{x})$ .

(d) Since

$$p(Y = 1|\mathbf{x}; \boldsymbol{\alpha}) = \frac{1}{1 + e^{-\boldsymbol{\alpha} \cdot \mathbf{x}}} \quad (44)$$

for  $\mathbf{x} = (1, x_1, x_2, \dots, x_n)$ ,  $\boldsymbol{\alpha} = (\alpha_0, \dots, \alpha_n)$  and binary  $Y$ ,

$$p(Y = 1|\mathbf{x}; \boldsymbol{\alpha}) = 1 - \frac{1}{1 + e^{-\boldsymbol{\alpha} \cdot \mathbf{x}}} = \frac{e^{-\boldsymbol{\alpha} \cdot \mathbf{x}}}{1 + e^{-\boldsymbol{\alpha} \cdot \mathbf{x}}} \quad (45)$$

Therefore,

$$p(Y = y|\mathbf{x}; \boldsymbol{\alpha}) = \frac{e^{(1-y)\boldsymbol{\alpha} \cdot \mathbf{x}}}{Z(\boldsymbol{\alpha}, \mathbf{x})} \quad (46)$$

where

$$Z(\boldsymbol{\alpha}, \mathbf{x}) = 1 + e^{-\boldsymbol{\alpha} \cdot \mathbf{x}} \quad (47)$$

Therefore, setting  $\mathbf{f}(y, \mathbf{x}) = (1 - y)\mathbf{x}$  and  $h(\mathbf{x}, \mathbf{y}) = 1$ , we see that this conditional distribution is part of the exponential family.

## 5 Problem 5: Conjugacy and Prediction

(a) The Dirichlet distribution is

$$\text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha}) \propto \prod_k \theta_k^{\alpha_k - 1} \quad (48)$$

and the Multinomial distribution is

$$\text{Mult}(\mathbf{x}|\boldsymbol{\theta}) \propto \prod_k \theta_k^{x_k} \quad (49)$$

Therefore, the posterior on  $\boldsymbol{\theta}$  is

$$p(\boldsymbol{\theta}|\mathbf{x}; \boldsymbol{\alpha}) \propto \left( \prod_k \theta_k^{x_k} \right) \left( \prod_k \theta_k^{\alpha_k - 1} \right) = \prod_k \theta_k^{\alpha_k - 1 + x_k} \propto \text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha} + \mathbf{x}) \quad . \quad (50)$$

Therefore, the posterior given a single “observation”  $\mathbf{x}$  is given by a Dirichlet with new hyperparameters. Given multiple, independent observations, this becomes

$$p(\boldsymbol{\theta}|\{\mathbf{x}\}; \boldsymbol{\alpha}) \propto \prod_{i=1}^N p(\boldsymbol{\theta}|\mathbf{x}^{(i)}; \boldsymbol{\alpha}) = \prod_{i=1}^N \text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha} + \mathbf{x}^{(i)}) \quad . \quad (51)$$

Now, to give the result in the notation of the problem, since  $\mathbf{x}^{(i)}$  is zero everywhere except in one component where it equals one, the posterior can be written as a Dirichlet distribution with hyperparameters  $\boldsymbol{\alpha}'$  given by

$$\alpha'_k = \alpha_k + \sum_{i=1}^N \begin{cases} 1, & \text{if } \mathbf{x}_k^{(i)} = 1 \\ 0, & \text{otherwise} \end{cases} \quad (52)$$

(b) The joint posterior on  $\mathbf{x}_{\text{new}}$  and  $\boldsymbol{\theta}$  is

$$p(\mathbf{x}_{\text{new}}, \boldsymbol{\theta}|\{\mathbf{x}^{(i)}\}; \boldsymbol{\alpha}) = p(\mathbf{x}_{\text{new}}|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\{\mathbf{x}^{(i)}\}; \boldsymbol{\alpha}) \quad (53)$$

and this is given by

$$\text{Mult}(\mathbf{x}_{\text{new}}|\boldsymbol{\theta}) \text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha}') \sim \text{Dir}(\boldsymbol{\theta}|\boldsymbol{\alpha}'') \quad . \quad (54)$$

Integrating over  $\boldsymbol{\theta}$ , we get

$$p(\mathbf{x}_{\text{new}}|\{\mathbf{x}^{(i)}\}; \boldsymbol{\alpha}) = \int d\boldsymbol{\theta} p(\mathbf{x}_{\text{new}}, \boldsymbol{\theta}|\{\mathbf{x}^{(i)}\}; \boldsymbol{\alpha}) = \frac{\prod_k \Gamma(\alpha'_k + x_{\text{new},k})}{\Gamma(\sum_k [\alpha'_k + x_{\text{new},k}])} \quad . \quad (55)$$

## 6 Problem 6: Kullback-Leibler divergence

(a) For a convex function  $f(x)$ , Jensen’s inequality is

$$E[f(x)] \geq f(E[x]) \quad . \quad (56)$$

For our problem, we can define  $y = q/p$  and  $f(x) = -\log x$ . Therefore,  $E_p[y] = \sum p \frac{q}{p} = \sum q = 1$  and

$$E_p[f(y)] = -\sum p \log \frac{q}{p} = \sum p \log \frac{p}{q} \geq -\log(E_p[y]) = -\log(1) = 0 \quad . \quad (57)$$

This proves that  $D(p||q) \geq 0$ . If  $p = q$  then

$$D(p||q) = \sum p \log \frac{p}{q} = \sum p \log 1 = \sum 0 = 0 \quad . \quad (58)$$

If  $p \neq q$  then  $D(p||q) > 0$  since equality in Jensen’s inequality holds only when  $f(x)$  is not strictly convex (i.e. only when  $p = q$ ). Therefore,  $D(p||q) = 0$  if and only if  $p = q$ .

(b) The K-L divergence can be rewritten as

$$D(p||q) = \sum p \log \frac{p}{q} = \sum (p \log p - p \log q) = -H(p) - \sum p \log q \quad . \quad (59)$$

Then, we can choose the uniform distribution  $q = 1/k$  and find

$$D(p||q) = -H(p) + \log k \sum p = -H(p) + \log k \geq 0 \rightarrow \log k \geq H(p) \quad (60)$$

with equality when  $p = q = 1/k$ .