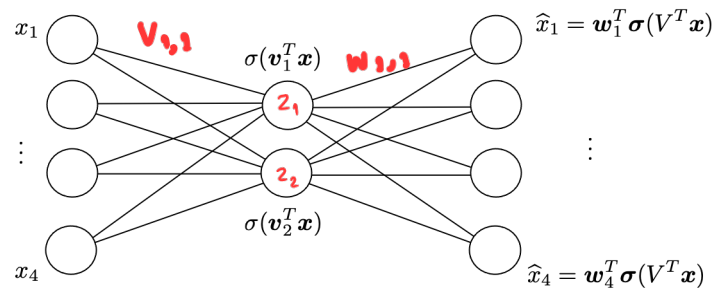


Autoencoders, Multivariate Normal KL

Submit a PDF of your answers to Canvas

1. Consider the undercomplete autoencoder shown in the figure below. In this problem, there is no activation function on the output layer (the first output is given as $\hat{x}_1 = \mathbf{w}_1^T \boldsymbol{\sigma}(V^T \mathbf{x})$).



$\Rightarrow V$ and W $V \rightarrow 2 \times 4, W \rightarrow 4 \times 2 \rightarrow 16$ parameters in total

- a) How many parameters do you need to learn to train your autoencoder?
 - b) Let $w_{1,1}$ denote the weight associated with the connection between the first (top) hidden node and the first output node, and $w_{2,1}$ denote the weight associated with the bottom hidden node and the first output. Compute the partial derivative of the output \hat{x}_1 with respect to weight $w_{1,1}$ when $\sigma(\cdot)$ is the logistic activation function. $\sigma(v_1^T x)$
 - c) Let $v_{1,1}$ denote the weight associated with the connection between the first (top) input node and the first hidden node. Compute the partial derivative of the output \hat{x}_1 with respect to weight $v_{1,1}$ when $\sigma(\cdot)$ is the logistic activation function.
 $\frac{d}{dv_{1,1}} \hat{x}_1 = \frac{d}{dz_1} \hat{x}_1 \frac{dz_1}{dv_{1,1}} = w_{1,1} (\sigma(z_1) (1 - \sigma(z_1))) x_1$
 - d) You decide to use the squared error loss function, i.e. $\ell(\hat{\mathbf{x}}, \mathbf{x}) = \|\hat{\mathbf{x}} - \mathbf{x}\|^2$. Compute the partial derivative of $\ell(\hat{\mathbf{x}}, \mathbf{x})$ first with respect to $w_{1,1}$, and then with respect to $v_{1,1}$.
 $\frac{\partial}{\partial v_{1,1}} \ell(\hat{\mathbf{x}}, \mathbf{x}) = \sum_{i=1}^4 \left(\frac{\partial}{\partial v_{1,1}} \hat{x}_i \right) 2(x_i - \hat{x}_i) = 2 \sum_{i=1}^4 (x_i - \hat{x}_i) w_{1,i} (\sigma(z_1) (1 - \sigma(z_1))) x_1$
2. In this problem we will compute the Kullback–Leibler divergence between two multivariate normal distributions (a quantity we will use in a later lecture). For reference, recall the general form of a multivariate normal distribution:

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{n}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

and the definition of KL divergence:

$$D(p(\mathbf{x}) || q(\mathbf{x})) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}.$$

Let $p(\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $q(\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$. Describe and justify each step in the derivation of the KL divergence between the two multivariate normal distributions:

$$\begin{aligned}
 & D(p(\mathbf{x})||q(\mathbf{x})) \\
 & \stackrel{(1)}{=} \int p(\mathbf{x}) \left(\frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) \right) d\mathbf{x} \\
 & \stackrel{(2)}{=} \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} + \frac{1}{2} E_p \left[(\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \right] \quad \text{defn. of expectation} \\
 & \quad - \frac{1}{2} E_p \left[(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) \right] \quad E(g(\mathbf{x})) = \int g(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\
 & \stackrel{(3)}{=} \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} + \frac{1}{2} E_p \left[\text{tr} \left((\mathbf{x} - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \right) \right] \quad \text{is a scalar,} \\
 & \quad - \frac{1}{2} E_p \left[\text{tr} \left((\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) \right) \right] \quad \text{trace is the identity transformation!} \\
 & \stackrel{(4)}{=} \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} + \frac{1}{2} \text{tr} \left(\boldsymbol{\Sigma}_2^{-1} E_p \left[(\mathbf{x} - \boldsymbol{\mu}_2)(\mathbf{x} - \boldsymbol{\mu}_2)^T \right] \right) \quad \text{tr}(AB) = \text{tr}(BA) \\
 & \quad - \frac{1}{2} \text{tr} \left(\boldsymbol{\Sigma}_1^{-1} E_p \left[(\mathbf{x} - \boldsymbol{\mu}_1)(\mathbf{x} - \boldsymbol{\mu}_1)^T \right] \right) \quad \text{we may remove } \boldsymbol{\Sigma}_i^{-1} \text{ out of expectation} \\
 & \quad \quad \text{bc. } A = (\mathbf{x} - \boldsymbol{\mu}_1)^T, B = \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) \\
 & \stackrel{(5)}{=} \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} + \frac{1}{2} \text{tr} \left(\boldsymbol{\Sigma}_2^{-1} E_p \left[(\mathbf{x} - \boldsymbol{\mu}_2)(\mathbf{x} - \boldsymbol{\mu}_2)^T \right] \right) - \frac{1}{2} \text{tr} \left(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_1 \right) \quad \text{definition of covariance} \\
 & \quad \quad \boldsymbol{\Sigma}_i = E_p \left[(\mathbf{x}_i - \boldsymbol{\mu}_i)^T (\mathbf{x}_i - \boldsymbol{\mu}_i) \right] \\
 & \stackrel{(6)}{=} \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} + \frac{1}{2} \text{tr} \left(\boldsymbol{\Sigma}_2^{-1} E_p \left[\mathbf{x}\mathbf{x}^T - 2\mathbf{x}\boldsymbol{\mu}_2^T + \boldsymbol{\mu}_2\boldsymbol{\mu}_2^T \right] \right) - \frac{1}{2} \text{tr} \left(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_1 \right) \quad \text{expansion} \\
 & \stackrel{(7)}{=} \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} + \frac{1}{2} \text{tr} \left(\boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\Sigma}_1 + \boldsymbol{\mu}_1\boldsymbol{\mu}_1^T - 2\boldsymbol{\mu}_1\boldsymbol{\mu}_2^T + \boldsymbol{\mu}_2\boldsymbol{\mu}_2^T) \right) - \frac{1}{2} n \quad E_p[\mathbf{x}] = \boldsymbol{\mu}_1, E_p[\mathbf{x}^T \mathbf{x}] = \boldsymbol{\Sigma}_1 \\
 & \stackrel{(8)}{=} \frac{1}{2} \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} + \frac{1}{2} \text{tr} \left(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1 \right) + \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - \frac{1}{2} n \quad \text{arrangement of terms}
 \end{aligned}$$