

AI 501 Machine Learning for AI, Spring 2021

Final Exam

1. (60 points) Consider the Probabilistic Principal Component Analysis (PPCA) model with generative process defined as

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}_h, \mathbf{I}_h), \quad \mathbf{x}|\mathbf{z} \sim \mathcal{N}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}_d), \quad (1)$$

where $\mathbf{z} \in \mathbb{R}^h$, $\mathbf{x} \in \mathbb{R}^d$, $\mathbf{W} \in \mathbb{R}^{d \times h}$, and $\boldsymbol{\mu} \in \mathbb{R}^d$. Let $\boldsymbol{\theta} := \{\mathbf{W}, \boldsymbol{\mu}, \sigma^2\}$ be the collection of parameters.

- (a) (20 points) Show that the conditional distribution $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})$ and the marginal distribution $p(\mathbf{x}; \boldsymbol{\theta})$ are given as follows:

$$p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}) = \mathcal{N}(\mathbf{z}|\mathbf{V}^{-1}\mathbf{W}^\top(\mathbf{x} - \boldsymbol{\mu}), \sigma^2\mathbf{V}^{-1}), \quad p(\mathbf{x}; \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_d), \quad (2)$$

where

$$\mathbf{V} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_h. \quad (3)$$

Hint: you may use the following matrix inverse identity:

$$(\mathbf{I} + \mathbf{A}\mathbf{B})^{-1} = \mathbf{I} - \mathbf{A}(\mathbf{I} + \mathbf{B}\mathbf{A})^{-1}\mathbf{B}. \quad (4)$$

- (b) (10 points) Let $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ be a training set. Using the expression of $p(\mathbf{x}; \boldsymbol{\theta})$ obtained above, compute the maximum-likelihood estimator of the parameter $\boldsymbol{\mu}$.
- (c) (30 points) Now consider the Expectation-Maximization (EM) procedure to learn parameters (\mathbf{W}, σ^2) via maximum likelihood while holding $\boldsymbol{\mu}$ fixed. Let $\boldsymbol{\theta}_t = (\mathbf{W}_t, \boldsymbol{\mu}, \sigma_t^2)$ be the parameter estimate at step t . Show that the M-step update equations are given as follows:

$$\mathbf{W}_{t+1} = \mathbf{S}\mathbf{W}_t(\sigma_t^2 \mathbf{I}_h + \mathbf{V}_t^{-1}\mathbf{W}_t^\top \mathbf{S}\mathbf{W}_t)^{-1}, \quad (5)$$

$$\sigma_{t+1}^2 = \frac{1}{d} \text{Tr}(\mathbf{S} - \mathbf{S}\mathbf{W}_t \mathbf{V}_t^{-1} \mathbf{W}_t^\top), \quad (6)$$

where

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top, \quad \mathbf{V}_t = \mathbf{W}_t^\top \mathbf{W}_t + \sigma_t^2 \mathbf{I}_h. \quad (7)$$

Hint: you may use the following identity,

$$\mathbf{x}^\top \mathbf{A} \mathbf{y} = \text{Tr}(\mathbf{A} \mathbf{y} \mathbf{x}^\top). \quad (8)$$

You can also refer to the [matrix cookbook](#) for the computation of derivatives.

2. (40 points) Consider the Mixture of Probabilistic Principal Component Analysers (MPPCA) defined as follows:

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}_h, \mathbf{I}_h), \quad c \sim \text{Cat}(\boldsymbol{\pi}), \quad \mathbf{x}|\mathbf{z}, c \sim \mathcal{N}(\mathbf{W}_c \mathbf{z} + \boldsymbol{\mu}_c, \sigma_c^2 \mathbf{I}_d), \quad (9)$$

where $\mathbf{z} \in \mathbb{R}^h$, $\mathbf{x} \in \mathbb{R}^d$, $c \in \{1, \dots, k\}$. Let $\boldsymbol{\theta} = \{\boldsymbol{\pi}, \{\mathbf{W}_j, \boldsymbol{\mu}_j, \sigma_j^2\}_{j=1}^k\}$ be the parameters with $\mathbf{W}_j \in \mathbb{R}^{d \times h}$, $\boldsymbol{\mu}_j \in \mathbb{R}^d$, and $\boldsymbol{\pi} \in [0, 1]^k$.

- (a) (10 points) Compute the distributions $p(c|\mathbf{x}; \boldsymbol{\theta})$, $p(\mathbf{z}|c, \mathbf{x}; \boldsymbol{\theta})$ and $p(\mathbf{x}; \boldsymbol{\theta})$.
- (b) (15 points) Let $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n$ be a training set to fit MPPCA. Let us first consider the parameters $\boldsymbol{\pi}$ and $\{\boldsymbol{\mu}_j\}_{j=1}^k$. Given a current parameter estimate $\boldsymbol{\theta}_t = \{\{\mathbf{W}_{t,j}, \boldsymbol{\mu}_{t,j}, \sigma_{t,j}^2\}_{j=1}^k, \boldsymbol{\pi}_j\}$, we first compute the expected complete-data log-likelihood *only* for c as follows:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}) &= \sum_{i=1}^n \mathbb{E}_q[\log p(\mathbf{x}_i, c_i; \boldsymbol{\theta})] \\ &= \sum_{i=1}^n \mathbb{E}_q \left[\log \prod_{j=1}^k \left(\pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \mathbf{W}_j \mathbf{W}_j^\top + \sigma_j^2 \mathbf{I}_d) \right)^{\mathbb{1}_{\{c_i=j\}}} \right] \\ &= \sum_{i=1}^n \sum_{j=1}^k r_{t,i,j} \left(\log \pi_j + \log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \mathbf{W}_j \mathbf{W}_j^\top + \sigma_j^2 \mathbf{I}_d) \right),\end{aligned}\tag{10}$$

where

$$r_{t,i,j} := \frac{\pi_{t,j} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{t,j}, \mathbf{W}_{t,j} \mathbf{W}_{t,j}^\top + \sigma_{t,j}^2 \mathbf{I}_d)}{\sum_{\ell=1}^k \pi_{t,\ell} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_{t,\ell}, \mathbf{W}_{t,\ell} \mathbf{W}_{t,\ell}^\top + \sigma_{t,\ell}^2 \mathbf{I}_d)}.\tag{11}$$

Show that the M-step update for $\boldsymbol{\pi}$ and $\{\boldsymbol{\mu}_j\}_{j=1}^k$ are given as

$$\pi_{t+1,j} = \frac{\sum_{i=1}^n r_{t,i,j}}{n}, \quad \boldsymbol{\mu}_{t+1,j} = \frac{\sum_{i=1}^n r_{t,i,j} \mathbf{x}_i}{\sum_{i=1}^n r_{t,i,j}}.\tag{12}$$

- (c) (15 points) Consider the parameters $\{\mathbf{W}_j, \sigma_j\}_{j=1}^k$. While holding $\boldsymbol{\pi}, \{\boldsymbol{\mu}_j\}_{j=1}^k$ fixed, we would like to apply M-step with the expected complete-data log-likelihood with both \mathbf{z}_i and c_i ,

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}) = \sum_{i=1}^n \mathbb{E}_q[\log p(\mathbf{x}_i, \mathbf{z}_i, c_i; \boldsymbol{\theta})],\tag{15}$$

where

$$q(\mathbf{z}_i, c_i; \boldsymbol{\theta}_t) = p(\mathbf{z}_i | c_i, \mathbf{x}_i; \boldsymbol{\theta}_t) p(c_i | \mathbf{x}_i; \boldsymbol{\theta}_t).\tag{16}$$

Derive the M-step update equations for the parameters $\{\mathbf{W}_j, \sigma_j\}_{j=1}^k$.