

Factor Analysis, Probabilistic PCA, and Variational Inference

Statistical Machine Learning (ENGG*6600*02)

School of Engineering,
University of Guelph, ON, Canada

Course Instructor: Benjamin Ghojogh
Summer 2023

Variational Inference

Variational Inference

- Consider a dataset $\{\mathbf{x}_i\}_{i=1}^n$. Assume that every data point $\mathbf{x}_i \in \mathbb{R}^d$ is generated from a latent variable $\mathbf{z}_i \in \mathbb{R}^p$. This latent variable has a prior distribution $\mathbb{P}(\mathbf{z}_i)$. According to Bayes' rule, we have:

$$\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i) = \frac{\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i) \mathbb{P}(\mathbf{z}_i)}{\mathbb{P}(\mathbf{x}_i)}. \quad (1)$$

- Let $\mathbb{P}(\mathbf{z}_i)$ be an arbitrary distribution denoted by $q(\mathbf{z}_i)$. Suppose the parameter of conditional distribution of \mathbf{z}_i on \mathbf{x}_i is denoted by $\boldsymbol{\theta}$; hence, $\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i) = \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})$. Therefore, we can say:

$$\mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}) = \frac{\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) \mathbb{P}(\mathbf{z}_i | \boldsymbol{\theta})}{\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})}. \quad (2)$$

Variational Inference

- Consider the Kullback-Leibler (KL) divergence [1] between the prior probability of the latent variable and the posterior of the latent variable:

$$\begin{aligned}\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta})) &\stackrel{(a)}{=} \int q(\mathbf{z}_i) \log\left(\frac{q(\mathbf{z}_i)}{\mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta})}\right) d\mathbf{z}_i \\&= \int q(\mathbf{z}_i) (\log(q(\mathbf{z}_i)) - \log(\mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta}))) d\mathbf{z}_i \\&\stackrel{(2)}{=} \int q(\mathbf{z}_i) (\log(q(\mathbf{z}_i)) - \log(\mathbb{P}(\mathbf{x}_i \mid \mathbf{z}_i, \boldsymbol{\theta})) - \log(\mathbb{P}(\mathbf{z}_i \mid \boldsymbol{\theta})) + \log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta}))) d\mathbf{z}_i \\&\stackrel{(b)}{=} \log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) + \int q(\mathbf{z}_i) (\log(q(\mathbf{z}_i)) - \log(\mathbb{P}(\mathbf{x}_i \mid \mathbf{z}_i, \boldsymbol{\theta})) - \log(\mathbb{P}(\mathbf{z}_i \mid \boldsymbol{\theta}))) d\mathbf{z}_i \\&= \log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) + \int q(\mathbf{z}_i) \log\left(\frac{q(\mathbf{z}_i)}{\mathbb{P}(\mathbf{x}_i \mid \mathbf{z}_i, \boldsymbol{\theta})\mathbb{P}(\mathbf{z}_i \mid \boldsymbol{\theta})}\right) d\mathbf{z}_i \\&= \log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) + \int q(\mathbf{z}_i) \log\left(\frac{q(\mathbf{z}_i)}{\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \boldsymbol{\theta})}\right) d\mathbf{z}_i \\&= \log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) + \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \boldsymbol{\theta})),\end{aligned}$$

where (a) is for definition of KL divergence and (b) is because $\log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta}))$ is independent of \mathbf{z}_i and comes out of integral and $\int d\mathbf{z}_i = 1$.

- Hence:

$$\log(\mathbb{P}(\mathbf{x}_i \mid \boldsymbol{\theta})) = \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i \mid \mathbf{x}_i, \boldsymbol{\theta})) - \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i \mid \boldsymbol{\theta})). \quad (3)$$

Variational Inference

- We found:

$$\log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})) = \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})) - \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})).$$

- We define the **Evidence Lower Bound (ELBO)** as:

$$\mathcal{L}(q, \boldsymbol{\theta}) := -\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \boldsymbol{\theta})). \quad (4)$$

So:

$$\log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})) = \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})) + \mathcal{L}(q, \boldsymbol{\theta}).$$

- Therefore:

$$\mathcal{L}(q, \boldsymbol{\theta}) = \log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})) - \underbrace{\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}))}_{\geq 0}. \quad (5)$$

- As the second term is negative with its minus, the ELBO is a lower bound on the log likelihood of data:

$$\mathcal{L}(q, \boldsymbol{\theta}) \leq \log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})). \quad (6)$$

The likelihood $\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})$ is also referred to as the **evidence**.

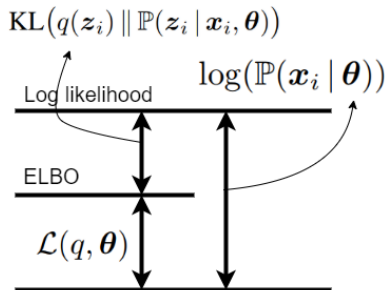
- Note that this lower bound gets tight when:

$$\begin{aligned} \mathcal{L}(q, \boldsymbol{\theta}) \approx \log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})) &\implies 0 \leq \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})) \stackrel{\text{set}}{=} 0 \\ &\implies q(\mathbf{z}_i) = \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta}). \end{aligned} \quad (7)$$

Variational Inference

- We found:

$$\log(\mathbb{P}(\mathbf{x}_i | \boldsymbol{\theta})) = \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \boldsymbol{\theta})) + \mathcal{L}(q, \boldsymbol{\theta}).$$



Expectation Maximization in Variational Inference

- According to MLE, we want to maximize the log-likelihood of data. According to Eq. (6):

$$\mathcal{L}(q, \theta) \leq \log(\mathbb{P}(\mathbf{x}_i | \theta)),$$

maximizing the ELBO will also maximize the log-likelihood.

- The Eq. (6) holds for any prior distribution q . We want to find the best distribution to maximize the lower bound.
- Hence, EM for variational inference is performed iteratively as:

$$\text{E-step: } q^{(t)} := \arg \max_q \mathcal{L}(q, \theta^{(t-1)}), \quad (8)$$

$$\text{M-step: } \theta^{(t)} := \arg \max_{\theta} \mathcal{L}(q^{(t)}, \theta), \quad (9)$$

where t denotes the iteration index.

Expectation Maximization in Variational Inference

- **E-step in EM for Variational Inference:** The E-step is:

$$\begin{aligned}\max_q \mathcal{L}(q, \theta^{(t-1)}) &\stackrel{(5)}{=} \max_q \log(\mathbb{P}(\mathbf{x}_i | \theta^{(t-1)})) + \max_q (-\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta^{(t-1)}))) \\ &= \max_q \log(\mathbb{P}(\mathbf{x}_i | \theta^{(t-1)})) + \min_q \text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta^{(t-1)})).\end{aligned}$$

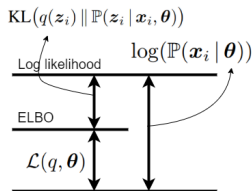
- The second term is always non-negative; hence, its minimum is zero:

$$\text{KL}(q(\mathbf{z}_i) \parallel \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta^{(t-1)})) \stackrel{\text{set}}{=} 0 \implies q(\mathbf{z}_i) = \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta^{(t-1)}),$$

which was already found in Eq. (7). Thus, the E-step assigns:

$$q^{(t)}(\mathbf{z}_i) \leftarrow \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \theta^{(t-1)}). \quad (10)$$

- In other words, in the figure, it pushes the middle line toward the above line by maximizing the ELBO.



Expectation Maximization in Variational Inference

- **M-step in EM for Variational Inference:** The M-step is:

$$\begin{aligned}\max_{\theta} \mathcal{L}(q^{(t)}, \theta) &\stackrel{(4)}{=} \max_{\theta} (-\text{KL}(q^{(t)}(z_i) \parallel \mathbb{P}(x_i, z_i \mid \theta))) \\ &\stackrel{(a)}{=} \max_{\theta} \left[- \int q^{(t)}(z_i) \log\left(\frac{q^{(t)}(z_i)}{\mathbb{P}(x_i, z_i \mid \theta)}\right) dz_i \right] \\ &= \max_{\theta} \int q^{(t)}(z_i) \log(\mathbb{P}(x_i, z_i \mid \theta)) dz_i - \max_{\theta} \int q^{(t)}(z_i) \log(q^{(t)}(z_i)) dz_i,\end{aligned}$$

where (a) is for definition of KL divergence.

- The second term is constant w.r.t. θ . Hence:

$$\begin{aligned}\max_{\theta} \mathcal{L}(q^{(t)}, \theta) &= \max_{\theta} \int q^{(t)}(z_i) \log(\mathbb{P}(x_i, z_i \mid \theta)) dz_i \\ &\stackrel{(a)}{=} \max_{\theta} \mathbb{E}_{\sim q^{(t)}(z_i)} [\log \mathbb{P}(x_i, z_i \mid \theta)],\end{aligned}$$

where (a) is because of definition of expectation. Thus, the M-step assigns:

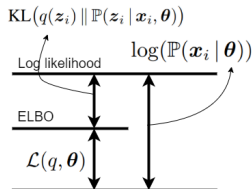
$$\theta^{(t)} \leftarrow \arg \max_{\theta} \mathbb{E}_{\sim q^{(t)}(z_i)} [\log \mathbb{P}(x_i, z_i \mid \theta)]. \quad (11)$$

Expectation Maximization in Variational Inference

- We found:

$$\theta^{(t)} \leftarrow \arg \max_{\theta} \mathbb{E}_{\sim q^{(t)}(z_i)} [\log \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta)].$$

- In other words, in the figure, it pushes the above line higher.



- The E-step and M-step together somehow play a **game** where the E-step tries to reach the middle line (or the ELBO) to the log-likelihood and the M-step tries to increase the above line (or the log-likelihood). This procedure is done repeatedly so the two steps help each other improve to higher values.
- To summarize, the EM in variational inference is:

$$q^{(t)}(z_i) \leftarrow \mathbb{P}(z_i | \mathbf{x}_i, \theta^{(t-1)}), \quad (12)$$

$$\theta^{(t)} \leftarrow \arg \max_{\theta} \mathbb{E}_{\sim q^{(t)}(z_i)} [\log \mathbb{P}(\mathbf{x}_i, \mathbf{z}_i | \theta)]. \quad (13)$$

Expectation Maximization in Variational Inference

- It is noteworthy that, in variational inference, sometimes, the parameter θ is absorbed into the latent variable \mathbf{z}_i .
- According to the chain rule, we have:

$$\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i, \theta) = \mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \theta) \mathbb{P}(\mathbf{z}_i | \theta) \mathbb{P}(\theta).$$

- Considering the term $\mathbb{P}(\mathbf{z}_i | \theta) \mathbb{P}(\theta)$ as one probability term, we have:

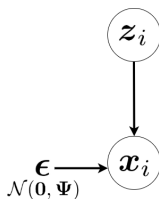
$$\mathbb{P}(\mathbf{x}_i, \mathbf{z}_i) = \mathbb{P}(\mathbf{x}_i | \mathbf{z}_i) \mathbb{P}(\mathbf{z}_i),$$

where the parameter θ disappears because of absorption.

Factor Analysis

Factor Analysis

- Factor analysis [2, 3, 4, 5] is one of the simplest and most fundamental **generative** models.
- Factor analysis assumes that every data point $\mathbf{x}_i \in \mathbb{R}^d$ is generated from a latent variable $\mathbf{z}_i \in \mathbb{R}^p$. The **latent variable** is also referred to as the **latent factor**; hence, the name of factor analysis comes from the fact that it analyzes the latent factors.
- In factor analysis, we assume that the data point \mathbf{x}_i is obtained through the following steps: (1) by **linear projection** of the p -dimensional \mathbf{z}_i onto a d -dimensional space by projection matrix $\mathbf{A} \in \mathbb{R}^{d \times p}$, then (2) applying some **linear translation**, and finally (3) **adding a Gaussian noise** $\epsilon \in \mathbb{R}^d$ with covariance matrix $\Psi \in \mathbb{R}^{d \times d}$.
- Note that as the noises in different dimensions are **independent**, the covariance matrix Ψ is diagonal.
- Factor analysis can be illustrated as a graphical model [6] where the visible data variable is conditioned on the latent variable and the noise random variable.



Factor Analysis

- For simplicity, the prior distribution of the latent variable can be assumed to be a multivariate Gaussian distribution:

$$\mathbb{P}(\mathbf{z}_i) = \mathcal{N}(\mathbf{z}_i | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) = \frac{1}{\sqrt{(2\pi)^p |\boldsymbol{\Sigma}_0|}} \exp\left(-\frac{(\mathbf{z}_i - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_0^{-1} (\mathbf{z}_i - \boldsymbol{\mu}_0)}{2}\right), \quad (14)$$

where $\boldsymbol{\mu}_0 \in \mathbb{R}^p$ and $\boldsymbol{\Sigma}_0 \in \mathbb{R}^{p \times p}$ are the mean and the covariance matrix of \mathbf{z}_i and $|\cdot|$ is the determinant of matrix.

- \mathbf{x}_i is obtained through (1) the linear projection of \mathbf{z}_i by $\boldsymbol{\Lambda} \in \mathbb{R}^{d \times p}$, (2) applying some linear translation, and (3) adding a Gaussian noise $\boldsymbol{\epsilon} \in \mathbb{R}^d$ with covariance $\boldsymbol{\Psi} \in \mathbb{R}^{d \times d}$.
- Hence, the data point \mathbf{x}_i has a **conditional multivariate Gaussian distribution given the latent variable**; its conditional likelihood is:

$$\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i) = \mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\Lambda}, \boldsymbol{\mu}, \boldsymbol{\Psi}) = \mathcal{N}(\boldsymbol{\Lambda} \mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}), \quad (15)$$

where $\boldsymbol{\mu}$, which is the translation vector, is the mean of data $\{\mathbf{x}_i\}_{i=1}^n$:

$$\mathbb{R}^d \ni \boldsymbol{\mu} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i. \quad (16)$$

Factor Analysis

- The marginal distribution of \mathbf{x}_i is:

$$\mathbb{P}(\mathbf{x}_i) = \int \mathbb{P}(\mathbf{x}_i | \mathbf{z}_i) \mathbb{P}(\mathbf{z}_i) d\mathbf{z}_i \implies$$

$$\mathbb{P}(\mathbf{x}_i | \mathbf{\Lambda}, \boldsymbol{\mu}, \boldsymbol{\Psi}) = \int \mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \mathbf{\Lambda}, \boldsymbol{\mu}, \boldsymbol{\Psi}) \mathbb{P}(\mathbf{z}_i | \boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) d\mathbf{z}_i$$

$$\stackrel{(a)}{=} \mathcal{N}(\mathbf{\Lambda}\boldsymbol{\mu}_0 + \boldsymbol{\mu}, \boldsymbol{\Psi} + \mathbf{\Lambda}\boldsymbol{\Sigma}_0\mathbf{\Lambda}^\top) \quad (17)$$

$$= \mathcal{N}(\hat{\boldsymbol{\mu}}, \boldsymbol{\Psi} + \hat{\boldsymbol{\Lambda}}\hat{\boldsymbol{\Lambda}}^\top), \quad (18)$$

where $\mathbb{R}^d \ni \hat{\boldsymbol{\mu}} := \mathbf{\Lambda}\boldsymbol{\mu}_0 + \boldsymbol{\mu}$, $\mathbb{R}^{d \times d} \ni \hat{\boldsymbol{\Lambda}} := \mathbf{\Lambda}\boldsymbol{\Sigma}_0^{(1/2)}$, and (a) is because mean is linear and variance is quadratic so the mean and variance of projection are applied linearly and quadratically, respectively.

- As the mean $\hat{\boldsymbol{\mu}}$ and covariance $\hat{\boldsymbol{\Lambda}}$ are needed to be learned, we can absorb $\boldsymbol{\mu}_0$ and $\boldsymbol{\Sigma}_0$ into $\boldsymbol{\mu}$ and $\mathbf{\Lambda}$ and assume that $\boldsymbol{\mu}_0 = \mathbf{0}$ and $\boldsymbol{\Sigma}_0 = \mathbf{I}$.
- In summary, factor analysis assumes every data point $\mathbf{x}_i \in \mathbb{R}^d$ is obtained by projecting a latent variable $\mathbf{z}_i \in \mathbb{R}^p$ onto a d -dimensional space by projection matrix $\mathbf{\Lambda} \in \mathbb{R}^{d \times p}$ and translating it by $\boldsymbol{\mu} \in \mathbb{R}^d$ and finally adding some Gaussian noise $\boldsymbol{\epsilon} \in \mathbb{R}^d$ (whose dimensions are independent) as:

$$\mathbf{x}_i := \mathbf{\Lambda}\mathbf{z}_i + \boldsymbol{\mu} + \boldsymbol{\epsilon}, \quad (19)$$

$$\mathbb{P}(\mathbf{z}_i) = \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (20)$$

$$\mathbb{P}(\boldsymbol{\epsilon}) = \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi}). \quad (21)$$

Factor Analysis

- The joint distribution of \mathbf{x}_i and \mathbf{z}_i is:

$$\mathbf{y}_i := \begin{bmatrix} \mathbf{x}_i \\ \mathbf{z}_i \end{bmatrix} \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y). \quad (22)$$

- The expectation of \mathbf{x}_i is:

$$\mathbb{E}[\mathbf{x}_i] \stackrel{(19)}{=} \mathbb{E}[\boldsymbol{\Lambda}\mathbf{z}_i + \boldsymbol{\mu} + \boldsymbol{\epsilon}] = \boldsymbol{\Lambda}\mathbb{E}[\mathbf{z}_i] + \boldsymbol{\mu} + \mathbb{E}[\boldsymbol{\epsilon}] \stackrel{(a)}{=} \boldsymbol{\mu}, \quad (23)$$

where (a) is because of Eqs. (20) and (21).

- Hence:

$$\boldsymbol{\mu}_y := \begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_z \end{bmatrix} \stackrel{(a)}{=} \begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{0} \end{bmatrix}, \quad (24)$$

where (a) is because of Eqs. (20) and (23).

Factor Analysis

- Lemma:

Lemma

Consider two random variables $\mathbf{x}_i \in \mathbb{R}^d$ and $\mathbf{z}_i \in \mathbb{R}^p$ and let $\mathbf{y}_i := [\mathbf{x}_i^\top, \mathbf{z}_i^\top]^\top \in \mathbb{R}^{d+p}$. Assume that \mathbf{x}_i and \mathbf{z}_i are jointly multivariate Gaussian; hence, the variable \mathbf{y}_i has a multivariate Gaussian distribution, i.e., $\mathbf{y}_i \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$. The mean and covariance can be decomposed as:

$$\boldsymbol{\mu}_y = [\boldsymbol{\mu}^\top, \boldsymbol{\mu}_0^\top]^\top \in \mathbb{R}^{d+p}, \quad (25)$$

$$\boldsymbol{\Sigma}_y = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \in \mathbb{R}^{(d+p) \times (d+p)}, \quad (26)$$

where $\boldsymbol{\mu} \in \mathbb{R}^d$, $\boldsymbol{\mu}_0 \in \mathbb{R}^p$, $\boldsymbol{\Sigma}_{11} \in \mathbb{R}^{d \times d}$, $\boldsymbol{\Sigma}_{22} \in \mathbb{R}^{p \times p}$, $\boldsymbol{\Sigma}_{12} \in \mathbb{R}^{d \times p}$, and $\boldsymbol{\Sigma}_{21} = \boldsymbol{\Sigma}_{12}^\top \in \mathbb{R}^{p \times d}$.

Factor Analysis

- Lemma [7]:

Lemma

$$\mathbb{R}^d \ni \boldsymbol{\mu}_{x|z} := \boldsymbol{\mu} + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{z}_i - \boldsymbol{\mu}_0), \quad (27)$$

$$\mathbb{R}^{d \times d} \ni \boldsymbol{\Sigma}_{x|z} := \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}, \quad (28)$$

and likewise for $\mathbf{z}_i|\mathbf{x}_i \sim \mathcal{N}(\boldsymbol{\mu}_{z|x}, \boldsymbol{\Sigma}_{z|x})$:

$$\mathbb{R}^p \ni \boldsymbol{\mu}_{z|x} := \boldsymbol{\mu}_0 + \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}(\mathbf{x}_i - \boldsymbol{\mu}), \quad (29)$$

$$\mathbb{R}^{p \times p} \ni \boldsymbol{\Sigma}_{z|x} := \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21}\boldsymbol{\Sigma}_{11}^{-1}\boldsymbol{\Sigma}_{12}. \quad (30)$$

Factor Analysis

- According to Eq. (20), we have $\Sigma_{22} = \Sigma_z = I$. According to Eq. (19), we have:

$$\begin{aligned}\Sigma_{11} &= \Sigma_x = \mathbb{E}[(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top] \\ &= \mathbb{E}[(\boldsymbol{\Lambda}\mathbf{z}_i + \boldsymbol{\mu} + \boldsymbol{\epsilon} - \boldsymbol{\mu})(\boldsymbol{\Lambda}\mathbf{z}_i + \boldsymbol{\mu} + \boldsymbol{\epsilon} - \boldsymbol{\mu})^\top] \\ &= \mathbb{E}[\boldsymbol{\Lambda}\mathbf{z}_i\mathbf{z}_i^\top\boldsymbol{\Lambda}^\top + \boldsymbol{\epsilon}\mathbf{z}_i^\top\boldsymbol{\Lambda}^\top + \boldsymbol{\Lambda}\mathbf{z}_i\boldsymbol{\epsilon}^\top + \boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top] \\ &= \boldsymbol{\Lambda}\mathbb{E}[\mathbf{z}_i\mathbf{z}_i^\top]\boldsymbol{\Lambda}^\top + \mathbb{E}[\boldsymbol{\epsilon}]\mathbb{E}[\mathbf{z}_i]^\top\boldsymbol{\Lambda}^\top + \boldsymbol{\Lambda}\mathbb{E}[\mathbf{z}_i]\mathbb{E}[\boldsymbol{\epsilon}]^\top + \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top] \\ &\stackrel{(a)}{=} \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \mathbf{0} + \mathbf{0} + \boldsymbol{\Psi} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi},\end{aligned}\tag{31}$$

where (a) is because of Eqs. (20) and (21).

- Moreover, we have:

$$\begin{aligned}\Sigma_{12} &= \Sigma_{xz} = \mathbb{E}[(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{z}_i - \boldsymbol{\mu}_0)^\top] \\ &\stackrel{(a)}{=} \mathbb{E}[(\boldsymbol{\Lambda}\mathbf{z}_i + \boldsymbol{\mu} + \boldsymbol{\epsilon} - \boldsymbol{\mu})(\mathbf{z}_i - \mathbf{0})^\top] \\ &\stackrel{(b)}{=} \boldsymbol{\Lambda}\mathbb{E}[\mathbf{z}_i\mathbf{z}_i^\top] + \mathbb{E}[\boldsymbol{\epsilon}]\mathbb{E}[\mathbf{z}_i]^\top = \boldsymbol{\Lambda}I + (\mathbf{0}\mathbf{0}^\top) = \boldsymbol{\Lambda},\end{aligned}\tag{32}$$

where (a) is because of Eqs. (19) and (20) and (b) is because \mathbf{z}_i and $\boldsymbol{\epsilon}$ are independent.

- We also have $\Sigma_{21} = \Sigma_{12}^\top = \boldsymbol{\Lambda}^\top$. Therefore:

$$\begin{bmatrix} \mathbf{x}_i \\ \mathbf{z}_i \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi} & \boldsymbol{\Lambda} \\ \boldsymbol{\Lambda}^\top & I \end{bmatrix}\right).\tag{33}$$

Factor Analysis

- Hence, the marginal distribution of data point \mathbf{x}_i is:

$$\mathbb{P}(\mathbf{x}_i) = \mathbb{P}(\mathbf{x}_i | \mathbf{\Lambda}, \boldsymbol{\mu}, \boldsymbol{\Psi}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{\Lambda}\mathbf{\Lambda}^\top + \boldsymbol{\Psi}). \quad (34)$$

According to Eqs. (29) and (30) [Lemma], the posterior or the conditional distribution of latent variable given data is:

$$\begin{aligned} q(\mathbf{z}_i) &\stackrel{(12)}{=} \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i) = \mathbb{P}(\mathbf{z}_i | \mathbf{x}_i, \mathbf{\Lambda}, \boldsymbol{\mu}, \boldsymbol{\Psi}) \\ &= \mathcal{N}(\boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}}), \end{aligned} \quad (35)$$

where:

$$\mathbb{R}^p \ni \boldsymbol{\mu}_{\mathbf{z}|\mathbf{x}} := \mathbf{\Lambda}^\top (\mathbf{\Lambda}\mathbf{\Lambda}^\top + \boldsymbol{\Psi})^{-1} (\mathbf{x}_i - \boldsymbol{\mu}), \quad (36)$$

$$\mathbb{R}^{p \times p} \ni \boldsymbol{\Sigma}_{\mathbf{z}|\mathbf{x}} := \mathbf{I} - \mathbf{\Lambda}^\top (\mathbf{\Lambda}\mathbf{\Lambda}^\top + \boldsymbol{\Psi})^{-1} \mathbf{\Lambda}. \quad (37)$$

- Recall that the conditional distribution of data given the latent variable, i.e. $\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i)$, was introduced in Eq. (15):

$$\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i) = \mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \mathbf{\Lambda}, \boldsymbol{\mu}, \boldsymbol{\Psi}) = \mathcal{N}(\mathbf{\Lambda}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}).$$

If data $\{\mathbf{x}_i\}_{i=1}^n$ are centered, i.e. $\boldsymbol{\mu} = \mathbf{0}$, the marginal of data, Eq. (34), and the likelihood of data, Eq. (15), become:

$$\mathbb{P}(\mathbf{x}_i | \mathbf{\Lambda}, \boldsymbol{\Psi}) = \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi} + \mathbf{\Lambda}\mathbf{\Lambda}^\top), \quad (38)$$

$$\mathbb{P}(\mathbf{x}_i | \mathbf{z}_i, \mathbf{\Lambda}, \boldsymbol{\Psi}) = \mathcal{N}(\mathbf{\Lambda}\mathbf{z}_i, \boldsymbol{\Psi}), \quad (39)$$

respectively. In some works, people center the data as a pre-processing to factor analysis.

Factor Analysis

- We can find the parameters Λ and Ψ using Expectation Maximization.
- See our tutorial “Factor analysis, probabilistic principal component analysis, variational inference, and variational autoencoder: Tutorial and survey” [8] for the details of EM steps in factor analysis.

Probabilistic Principal Component Analysis

Probabilistic Principal Component Analysis

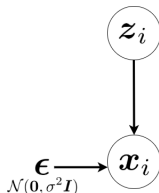
- **Probabilistic PCA (PPCA)** (1997-1999) [9, 10] is a **special case of factor analysis** where the variance of noise is equal in all dimensions of data space with covariance between dimensions, i.e.:

$$\Psi = \sigma^2 I. \quad (40)$$

- In other words, PPCA considers an **isotropic noise** in its formulation. Therefore, Eq. (21) is simplified to:

$$\mathbb{P}(\epsilon) = \mathcal{N}(\mathbf{0}, \sigma^2 I). \quad (41)$$

- Because of having zero covariance of noise between different dimensions, PPCA assumes that the data points are **independent** of each other given latent variables.
- PPCA can be illustrated as a graphical model, where the visible data variable is conditioned on the latent variable and the isotropic noise random variable.



Probabilistic Principal Component Analysis

- As PPCA is a special case of factor analysis, it also is solved using EM. Similar to factor analysis, it can be solved **iteratively using EM** [9].
- However, one can also find a **closed-form solution to its EM approach** [10]. Hence, by restricting the noise covariance to be isotropic, its **solution becomes simpler and closed-form**.
- We can find the parameters Λ and σ using Expectation Maximization.
- See our tutorial “Factor analysis, probabilistic principal component analysis, variational inference, and variational autoencoder: Tutorial and survey” [8] for the details of EM steps in PPCA.

Acknowledgment

- Some slides are based on our tutorial paper: “Factor analysis, probabilistic principal component analysis, variational inference, and variational autoencoder: Tutorial and survey” [8]
- Some slides of this slide deck are inspired by teachings of deep learning course at the Carnegie Mellon University (you can see their YouTube channel).
- Factor analysis in sklearn: <https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.FactorAnalysis.html>

References

- [1] S. Kullback and R. A. Leibler, "On information and sufficiency," *The annals of mathematical statistics*, vol. 22, no. 1, pp. 79–86, 1951.
- [2] B. Fruchter, *Introduction to factor analysis*.
Van Nostrand, 1954.
- [3] R. B. Cattell, "A biometrics invited paper. factor analysis: An introduction to essentials i. the purpose and underlying models," *Biometrics*, vol. 21, no. 1, pp. 190–215, 1965.
- [4] H. H. Harman, *Modern factor analysis*.
University of Chicago press, 1976.
- [5] D. Child, *The essentials of factor analysis*.
Cassell Educational, 1990.
- [6] Z. Ghahramani and G. E. Hinton, "The EM algorithm for mixtures of factor analyzers," tech. rep., Technical Report CRG-TR-96-1, University of Toronto, 1996.
- [7] A. Ng, "CS229 lecture notes for factor analysis," tech. rep., Stanford University, 2018.
- [8] B. Ghojogh, A. Ghodsi, F. Kararray, and M. Crowley, "Factor analysis, probabilistic principal component analysis, variational inference, and variational autoencoder: Tutorial and survey," *arXiv preprint arXiv:2101.00734*, 2021.
- [9] S. Roweis, "EM algorithms for PCA and SPCA," *Advances in neural information processing systems*, vol. 10, pp. 626–632, 1997.

References (cont.)

- [10] M. E. Tipping and C. M. Bishop, "Probabilistic principal component analysis," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 61, no. 3, pp. 611–622, 1999.