



Latent-Variable Models & Its Gaussian Case

Qinliang Su (苏勤亮)

Sun Yat-sen University

suqliang@mail.sysu.edu.cn

Outline

- Latent-Variable Models
- Gaussian Latent-Variable Model
- Relation to the PCA

What are LVMs & Why They are Needed?

- In supervised learning, both regression and classification can be understood as learning *conditional probability distributions*

$$p(y|\mathbf{x}; \mathbf{w})$$

- In regression, the conditional pdf is assumed of the form

$$p(y|\mathbf{x}; \mathbf{w}) = \mathcal{N}(y; \mathbf{w}^T \mathbf{x}, \sigma^2)$$

- For classification, the conditional pdf is assumed of the form

$$p(y|\mathbf{x}) = (\sigma(\mathbf{xw}))^y \cdot (1 - \sigma(\mathbf{xw}))^{1-y}$$

$$p(\mathbf{y}|\mathbf{x}) = \prod_{k=1}^K [\text{softmax}_k(\mathbf{W}\mathbf{x})]^{y_k}$$

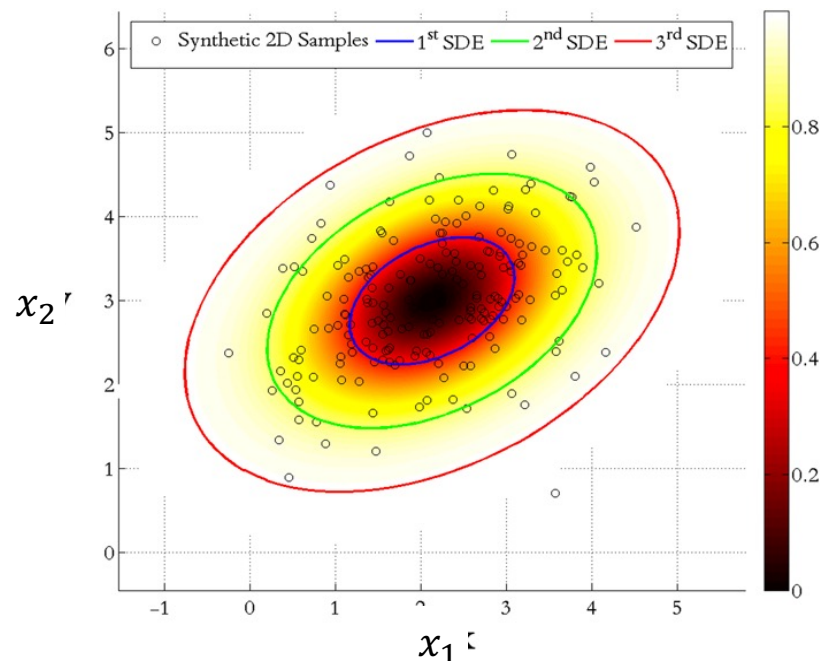
- Analogously, **unsupervised learning** can also be understood as learning a probability distribution, but it only concerns *input data x*

$$p(\mathbf{x}; \mathbf{w})$$

- Modeling x is much difficult than modeling the label y . The simplest way is to restrict $p(\mathbf{x}; \mathbf{w})$ to the Gaussian form

$$p(\mathbf{x}; \mathbf{w}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

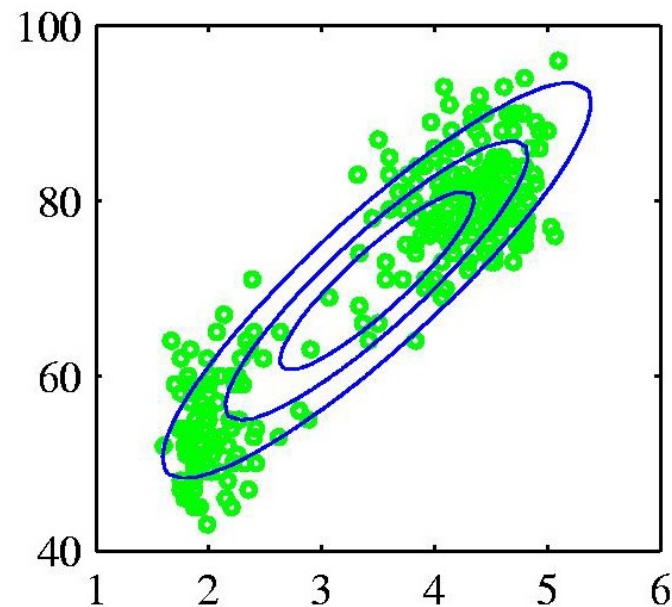
- $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are optimized to describe the data points $\{\mathbf{x}^{(n)}\}_{n=1}^N$ best



If x_1 is observed to be 3, what value x_2 would most likely be?

- The distribution of data \mathbf{x} is much more complex than what a simple Gaussian distribution can represent

For example, data points below cannot be fit well by a Gaussian distribution



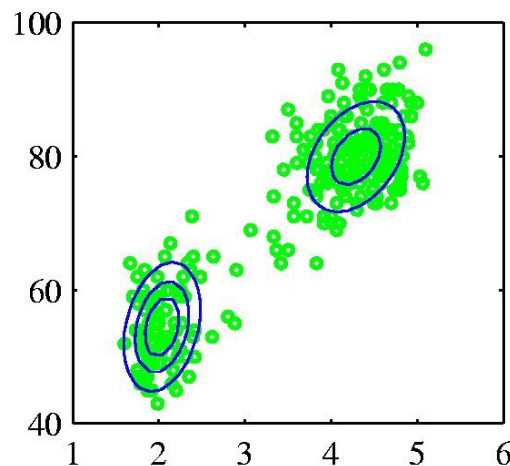
Obviously, to model the data well, the distribution $p(\mathbf{x})$ is required to be sufficiently expressive

- But how to obtain $p(\mathbf{x})$ that is able to *model complex data*?
- A possible way is to *composite 'many simple distributions'* to constitute *a expressive one*, that is,

$$p(\mathbf{x}) = \pi_1 p_1(\mathbf{x}) + \cdots + \pi_K p_K(\mathbf{x})$$

where $p_k(\mathbf{x})$ are simple distributions; to ensure $p(\mathbf{x})$ is valid, $\sum_{k=1}^K \pi_k = 1$ and $\pi_k \geq 0$

For example, the previous data can be modelled well *by two Gaussian distributions*



- The distribution $p(\mathbf{x}) = \pi_1 p_1(\mathbf{x}) + \cdots + \pi_K p_K(\mathbf{x})$ can be **derived from the following joint distribution**

$$p(\mathbf{x}, z) = p(\mathbf{x}|z)p(z)$$

where $p(\mathbf{x}|z = k) = p_k(\mathbf{x})$ and $p(z = k) = \pi_k$

- Using the marginal-joint relation $p(\mathbf{x}) = \sum_k p(\mathbf{x}, k)$, we know that

$$\begin{aligned} p(\mathbf{x}) &= \sum_{k=1}^K p_k(\mathbf{x}) \pi_k \\ &= \pi_1 p_1(\mathbf{x}) + \cdots + \pi_K p_K(\mathbf{x}) \end{aligned}$$

Thus, by defining **a simple joint distribution $p(\mathbf{x}, z)$** , it is possible to obtain a **flexible marginal distribution $p(\mathbf{x})$**

LVMs in General Form

- **LVMs**: a probabilistic model with latent variables

$$p(\mathbf{x}, \mathbf{z})$$

- \mathbf{x} is the **random variable of interest**
 - \mathbf{z} is the **latent variable** (nuisance variable)
- There sometimes exist multiple latent variables, *i.e.* $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K$

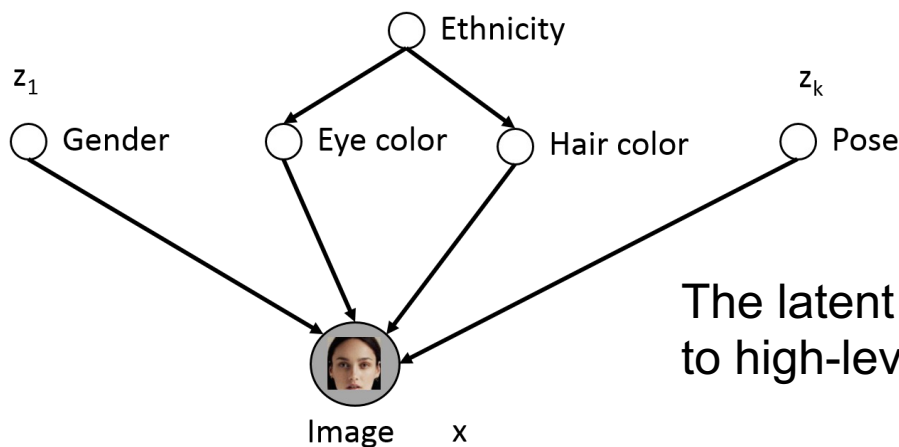
$$p(\mathbf{x}, \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K)$$

- The probabilistic model *w.r.t.* the interested variable \mathbf{x}

$$p(\mathbf{x}) = \int_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad \text{or} \quad p(\mathbf{x}) = \int_{\mathbf{z}_1 \cdots \mathbf{z}_K} p(\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_K) d\mathbf{z}_1 \cdots d\mathbf{z}_K$$

Advantages of LVMs

- **Representation ability:** yielding expressive distributions by compositing simple ones
- **Interpretability:** latent variables can sometimes be associated with some physical meanings



The latent variables z often correspond to high-level features

- **Integrating prior knowledges:** injecting our prior knowledges about a task or data into the model
- **Low-dimensional feature:** the latent vector can often be used as the feature of a data instance

Outline

- Latent-Variable Models
- Gaussian Latent-Variable Model
- Relation to the PCA

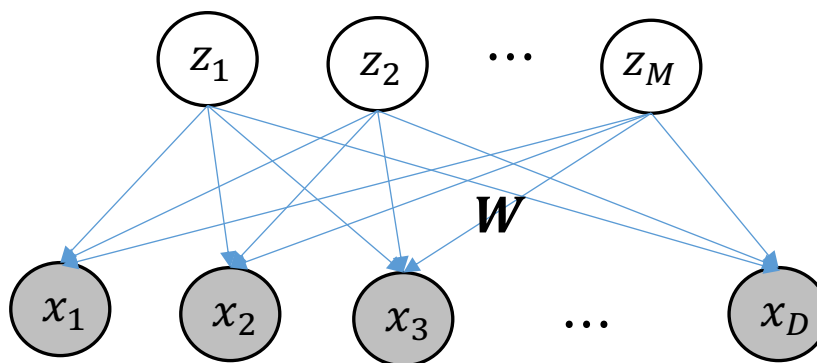
Modelling

- Assuming both of the prior and conditional pdfs are **independent Gaussian**

Prior distribution: $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, I)$

Likelihood function: $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 I)$

- Actually, the model describes how data samples \mathbf{x} are generated



$$\mathbf{z} = [z_1, \dots, z_M] \text{ \& } \mathbf{x} = [x_1, \dots, x_D]$$

Training Objective

- Given a set of samples $\{\mathbf{x}_n\}_{n=1}^N$, the question now becomes how to train the model $p(\mathbf{x}, \mathbf{z})$ so that it can describe the data best
- The model parameter \mathbf{W} can be learned by maximizing the log-likelihood

$$\max_{\mathbf{W}} \sum_{n=1}^N \log p(\mathbf{x}_n)$$

- In LVMs, what we have is the joint pdf

$$\begin{aligned} p(\mathbf{x}_n, \mathbf{z}_n) &= p(\mathbf{x}_n | \mathbf{z}_n) p(\mathbf{z}_n) \\ &= \mathcal{N}(\mathbf{x}_n; \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{z}_n; \mathbf{0}, \mathbf{I}), \end{aligned}$$

But what we need is to optimize $p(\mathbf{x}_n)$

Marginal Distribution $p(\mathbf{x})$

- The most direct method is to compute the marginal pdf first

$$p(\mathbf{x}_n) = \int_{\mathbf{z}_n} p(\mathbf{x}_n, \mathbf{z}_n) d\mathbf{z}_n$$

- Deriving the analytical expression for $p(\mathbf{x}_n)$ is impossible in most scenarios due to existence of the integration
- But for the **Gaussian case**, we can easily obtain it as

$$p(\mathbf{x}_n) = \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2 \mathbf{I})$$

A simple method to derive the marginal distribution

- From the model

$$\mathcal{N}(\mathbf{x}_n; \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{z}_n; \mathbf{0}, \mathbf{I}),$$

the data point \mathbf{x}_n can be understood as generated as

$$\mathbf{x}_n = \boldsymbol{\mu} + \mathbf{W}\mathbf{z}_n + \boldsymbol{\epsilon}_n$$

where $\mathbf{z}_n \sim \mathcal{N}(\mathbf{z}_n; \mathbf{0}, \mathbf{I})$ and $\boldsymbol{\epsilon}_n \sim \mathcal{N}(\mathbf{z}_n; \mathbf{0}, \sigma^2 \mathbf{I})$

- That is, data \mathbf{x}_n can be understood as generated from \mathbf{z}_n and $\boldsymbol{\epsilon}_n$ as $\mathbf{x}_n = \boldsymbol{\mu} + \mathbf{W}\mathbf{z}_n + \boldsymbol{\epsilon}_n$

Lemma: Any linear combination of Gaussian random variables is also Gaussian

- Therefore, \mathbf{x}_n also follows a Gaussian distribution

How can a Gaussian distribution be determined?

⇒ Mean & Covariance

- Mean & Covariance

Mean: $\mathbb{E}[\mathbf{x}_n] = \boldsymbol{\mu} + \mathbf{W}\mathbb{E}[\mathbf{z}_n] + \mathbb{E}[\boldsymbol{\epsilon}_n] = \boldsymbol{\mu}$

Covariance: $\mathbb{E}[(\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T] = \mathbf{W}\mathbb{E}[\mathbf{z}_n\mathbf{z}_n^T]\mathbf{W}^T + \mathbb{E}[\boldsymbol{\epsilon}_n\boldsymbol{\epsilon}_n^T]$
 $= \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}$

- Thus, the marginal distribution of \mathbf{x}_n is

$$p(\mathbf{x}_n) = \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})$$

Training by Optimizing $\log p(\mathbf{x})$

- Given the training dataset $\{\mathbf{x}_n\}_{n=1}^N$, to learn \mathbf{W} , $\boldsymbol{\mu}$ and σ^2 , what we need to do is to optimize the log-probability

$$\log p(\mathbf{x}_1, \dots, \mathbf{x}_N)$$

- Due to $p(\mathbf{x}_n) = \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})$, we have

$$\log p(\mathbf{x}_1, \dots, \mathbf{x}_N) = \sum_{n=1}^N \log \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})$$

- It can be further written as

$$\begin{aligned} \log p(\mathbf{x}_1, \dots, \mathbf{x}_N) = & -\frac{ND}{2} \log 2\pi - \frac{N}{2} \log \det(\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}) \\ & - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T (\mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I})^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned}$$

- By setting $\frac{\partial \log p(\mathbf{x}_1, \dots, \mathbf{x}_N)}{\partial \boldsymbol{\mu}} = 0$, we obtain

$$\boldsymbol{\mu} = \frac{\sum_{n=1}^N \mathbf{x}_n}{N}$$

$$\begin{aligned} \frac{\partial \ln \det(\mathbf{X})}{\partial \mathbf{X}} &= (\mathbf{X}^{-1})^T \\ \frac{\partial \ln \text{trace}(\mathbf{X}^{-1} \mathbf{B})}{\partial \mathbf{X}} &= -(\mathbf{X}^{-1} \mathbf{B} \mathbf{X}^{-1})^T \end{aligned}$$

- By denoting $\boldsymbol{\Sigma} = \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}$, we have

$$\begin{aligned} \frac{\partial \log p(\mathbf{x}_1, \dots, \mathbf{x}_N)}{\partial \boldsymbol{\Sigma}} &= -\frac{N}{2} \boldsymbol{\Sigma}^{-1} + \frac{1}{2} \sum_{n=1}^N \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_n) (\mathbf{x}_n - \boldsymbol{\mu}_n)^T \boldsymbol{\Sigma}^{-1} \\ &= -\frac{N}{2} \boldsymbol{\Sigma}^{-1} + \frac{N}{2} \boldsymbol{\Sigma}^{-1} \mathbf{S} \boldsymbol{\Sigma}^{-1} \end{aligned}$$

\Rightarrow Thus, it can be derived that $\boldsymbol{\Sigma} = \mathbf{S}$

- When $\boldsymbol{\Sigma}$ is restricted to the form $\boldsymbol{\Sigma} = \mathbf{W} \mathbf{W}^T + \sigma^2 \mathbf{I}$, it can be derived that

$$\mathbf{W} = \mathbf{U}(\boldsymbol{\Lambda} - \sigma^2 \mathbf{I})^{\frac{1}{2}}$$

- \mathbf{U} consists of the top- M eigenvectors of \mathbf{S}
- $\boldsymbol{\Lambda}$ is a diagonal matrix with the top- M eigenvalues of \mathbf{S}

Outline

- Latent-Variable Models
- Gaussian Latent-Variable Model
- Relation to the PCA

- Comparing the expression

$$\mathbf{W} = \mathbf{U}(\mathbf{\Lambda} - \sigma^2 \mathbf{I})^{\frac{1}{2}}$$

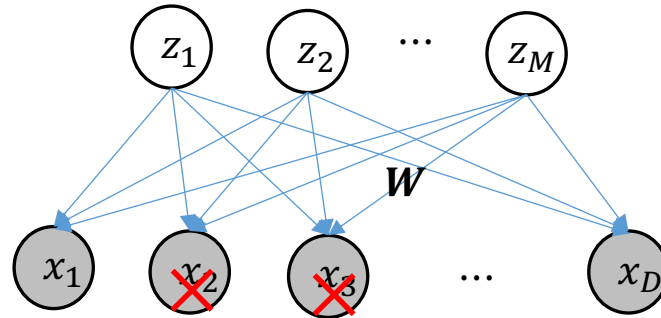
to the principle components of PCA, which are the matrix \mathbf{U} , we can see that

- \mathbf{W} can be viewed as un-normalized principle components of data \mathbf{x}_n , with the i -th component scaled by a coefficient $\sqrt{\lambda_i - \sigma^2}$

Gaussian latent-variable models are called *probabilistic PCA*

Advantages of Probabilistic PCA over PCA

- 1) Being able to deal with incomplete data observation



- 2) Easily extend to more complex models, *e.g.*, let $p(\mathbf{z})$ be a mixture distribution, or each x_i is assigned a distinctive σ_i^2

$$p(\mathbf{z}) = \sum_{k=1}^K \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$$

- 3) More computationally efficient: the optimal W can be learned by a mini-batch based SGD, rather than working on the whole dataset to compute the similarity matrix S