



Latent-Variable Models

Qinliang Su (苏勤亮)

Sun Yat-sen University

suqliang@mail.sysu.edu.cn

Outline

- Introduction of Latent-Variable Models
- Gaussian Latent-Variable Model
- Gaussian Mixture Model
- Examples of other LVMs

Unsupervised Probabilistic Modeling

- 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

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

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

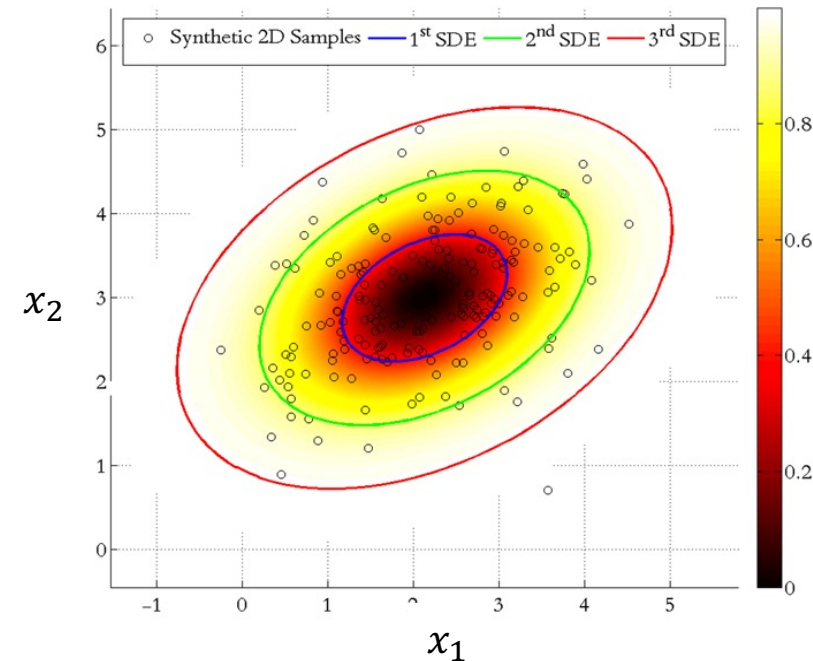
- **Unsupervised learning** can also be understood from the perspective of learning probability distributions. But it only concerns the distribution of *input data x*

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

- Modeling x is much difficult than modeling the label y . A naïve 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

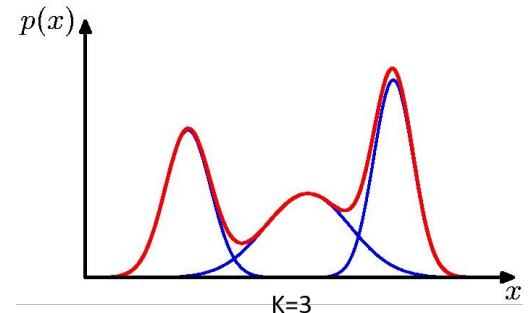


Obviously, the representational ability of the model is very limited

How does Latent Variables Arise?

- **Reason 1:** Building expressive models using the composition of simple models
 - Suppose there exists a simple categorical distribution $p(z) = \text{Cat}(K, \pi)$ and a Gaussian distribution $p(x) = \mathcal{N}(x|\mu, \sigma^2)$
 - By using them separately, only simple statistical relations can be modelled
 - But if we composite them as $p(x, z) = p(x|z)p(z)$, the induced marginal distribution $p(x)$ could be much more expressive

$$p(x) = \sum_z p(x|z)p(z) = \sum_{k=1}^K \pi_z \mathcal{N}(x|\mu_z, \sigma_z^2)$$



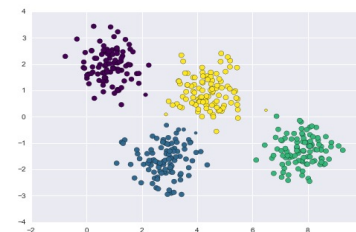
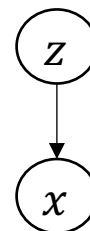
- Theoretically, it is able to represent any complex distribution

- Reason 2: hidden structures in the data

- 1) Data with hidden cluster structure

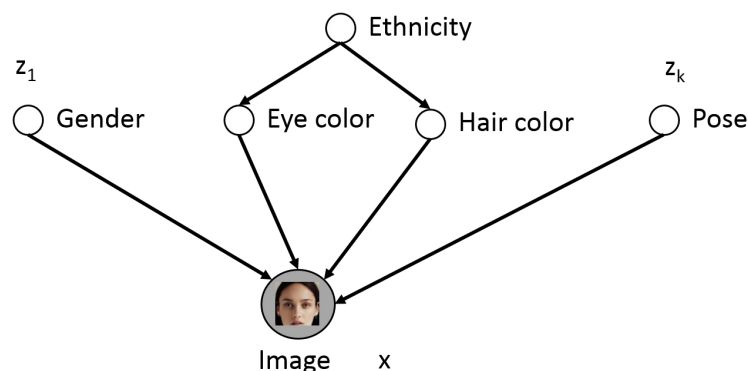
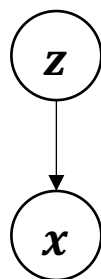
$$z_n \sim \text{Cat}(K, \boldsymbol{\pi})$$

$$x_n \sim \mathcal{N}(x | \boldsymbol{\mu}_{z_n}, \boldsymbol{\Sigma}_{z_n})$$



- 2) Topic model for documents

- 3) Image Modeling



- In the examples above, the latent variables z often correspond to high-level features
- *If the latent structure is respected, more interpretable models could be obtained*

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)
- Sometimes, there may exist multiple latent variables $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_K$

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

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

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

Outline

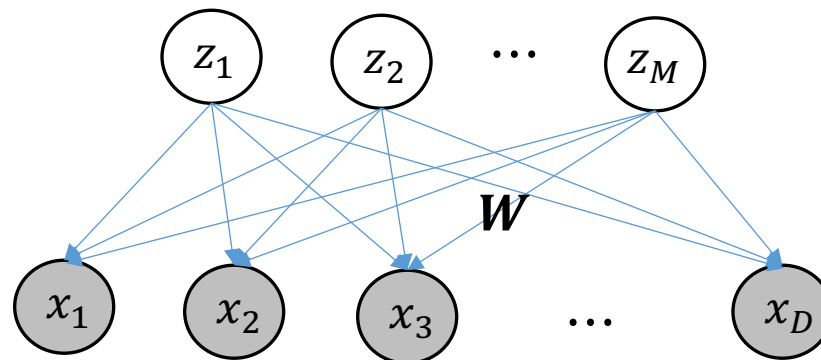
- Introduction of Latent-Variable Models
- **Gaussian Latent-Variable Model**
- Gaussian Mixture Model
- Examples of other LVMs

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

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

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

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



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

Training Objective

- Given the samples $\{\mathbf{x}_n\}_{n=1}^N$, the question becomes how to train the model $p(\mathbf{x}, \mathbf{z})$ to make it able to 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 from

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

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

Theorem: A linear combination of Gaussian random variables also follows a Gaussian distribution

- 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 Maximizing $\log p(\mathbf{x})$

- 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}

Relation to PCA

- Comparing the expression

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

to the principle components of PCA, which are the matrix U , we can see that

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

Gaussian latent-variable models are called *probabilistic PCA*

Outline

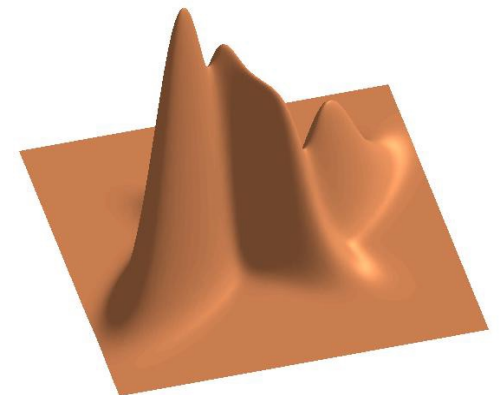
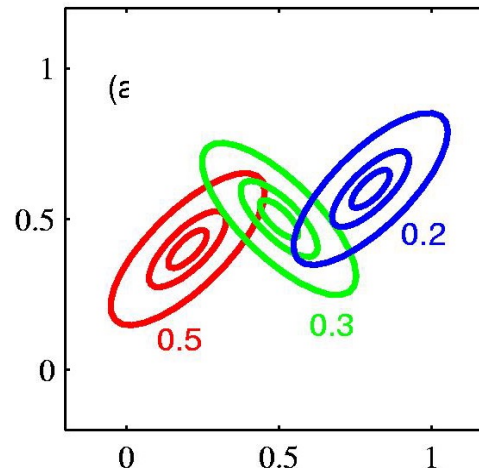
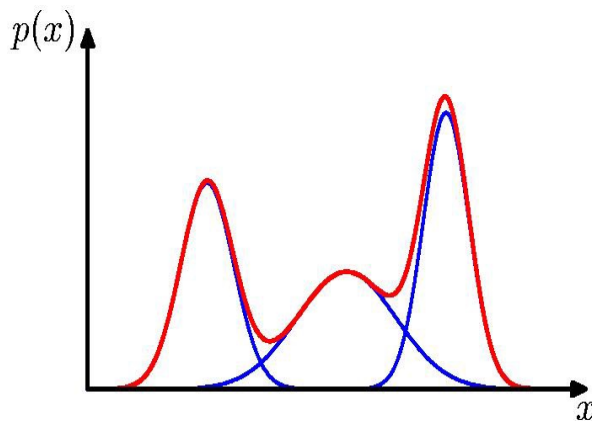
- Introduction of Latent-Variable Models
- Gaussian Latent-Variable Model
- **Gaussian Mixture Model**
- Examples of other LVMs

Gaussian Mixture Distributions

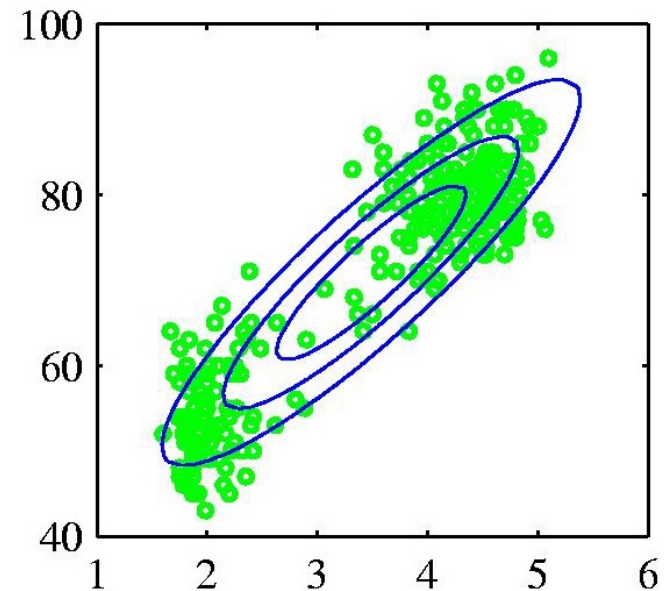
- The distribution expression

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

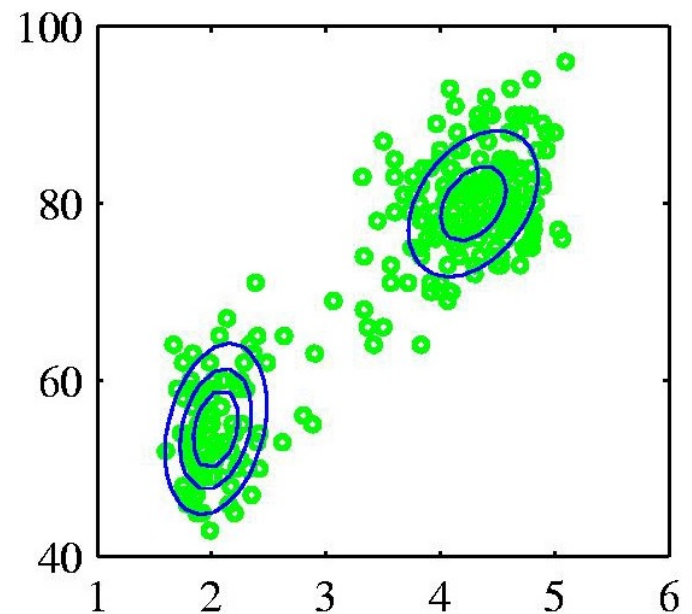
- K is the number of Gaussian distributions
- π_k is the weight of the k -th distribution with $\sum_{k=1}^K \pi_k = 1$
- $\boldsymbol{\mu}_k$ and $\boldsymbol{\Sigma}_k$ are the mean vector and covariance matrix of the k -th Gaussian distribution



- It is very difficult to model the green points by a Gaussian distribution



- But if we model it with the mixture of two Gaussian distributions, it looks much better



Representing Gaussian Mixture Distribution as LVM

- For a latent-variable model $p(\mathbf{x}, \mathbf{z})$, if we set its conditional distribution $p(\mathbf{x}|\mathbf{z})$ and prior distribution $p(\mathbf{z})$ as

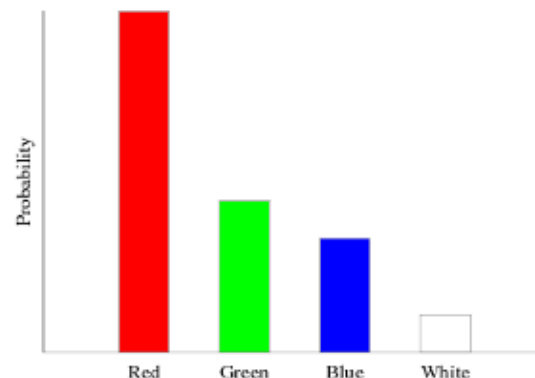
$$p(\mathbf{x}|\mathbf{z} = \mathbf{1}_k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$p(\mathbf{z} = \mathbf{1}_k) = \pi_k$$

- \mathbf{z} can only be a **one-hot vector**, with $\mathbf{1}_k$ denoting the k -th element to be 1
- $p(\mathbf{z} = \mathbf{1}_k) = \pi_k$ actually denotes a categorical distribution, that is,

$$p(\mathbf{z}) = \text{Cat}(\mathbf{z}; \boldsymbol{\pi})$$

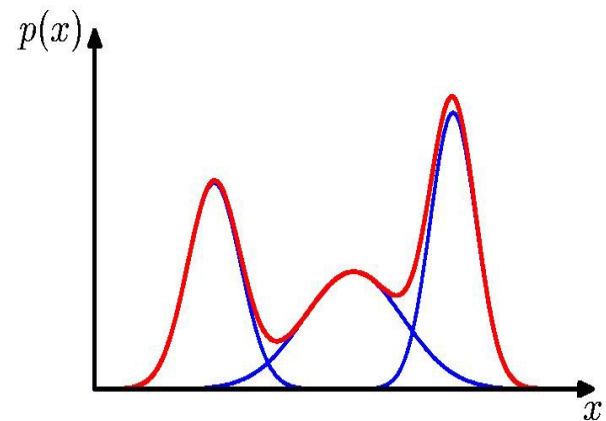
with $\text{Cat}(\mathbf{z} = \mathbf{1}_k; \boldsymbol{\pi}) = \pi_k$ and
 $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_K]$



- Due to $p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$, we can easily see that

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k),$$

which is exactly the Gaussian mixture distribution



Gaussian mixture distributions can be equivalently represented by the latent-variable model

$$p(\mathbf{x}, \mathbf{z}) = \prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{z_k}$$

Training by Maximizing the Marginal

- Given a set of training data $\{\mathbf{x}^{(n)}\}_{n=1}^N$, the goal is to learn the distribution parameters

$$\{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K \triangleq \boldsymbol{\theta}$$

- The data points $\mathbf{x}^{(n)}$ are assumed *i.i.d*, thus we can write the joint distribution as

$$p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \prod_{n=1}^N \underbrace{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}^{(n)}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{p(\mathbf{x}^n)}$$

The latent-variable form is not used

- For probabilistic models, the training objective is *to maximize the log-likelihood function*, that is,

$$\log p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}^{(n)}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Maximizing $\log p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)})$

- Substituting the expression of $\mathcal{N}(\mathbf{x}^{(n)}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ into it gives

$$\begin{aligned} \log p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) \\ = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_k) \right\} \right) \end{aligned}$$

- To optimize it, we require the *derivatives* of $\log p(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)})$ w.r.t. the model parameters $\{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$

How to Use the Learned Model?

- After learning the parameters $\{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$, that is, the distribution

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

is known, we can *use it to complete a lot of tasks*

- **Example:** Given a testing data point \mathbf{x} , can we use it to determine the probability that an \mathbf{x} belongs to the k -th cluster?

$$p(\mathbf{x} \in k\text{-th cluster}) = \frac{\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{i=1}^K \pi_i \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}$$

- Can we explain the probability in a more principled way?

$$p(\mathbf{z} = \mathbf{1}_k | \mathbf{x}) = ?$$

$$\begin{aligned} p(\mathbf{z} = \mathbf{1}_k | \mathbf{x}) &= \frac{p(\mathbf{x}, \mathbf{z} = \mathbf{1}_k)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}, \mathbf{z} = \mathbf{1}_k)}{\sum_{i=1}^K p(\mathbf{x}, \mathbf{1}_i)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{i=1}^K \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)} \end{aligned}$$

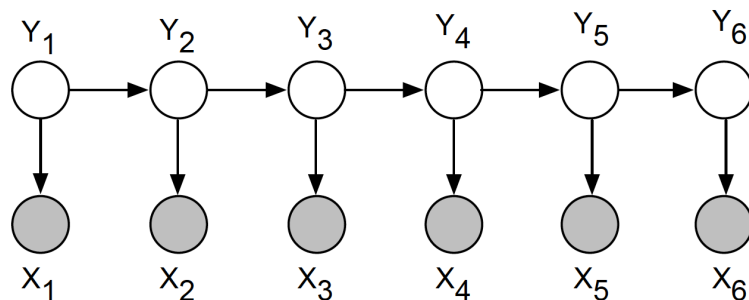
Thus, in the latent-variable model, the **posteriori** $p(\mathbf{z} | \mathbf{x})$ indicates the probability that a data instance belongs to different clusters

Outline

- Introduction of Latent-Variable Models
- Gaussian Latent-Variable Model
- Gaussian Mixture Model
- Examples of other LVMs

Application: Hidden Markov Model

- Hidden Markov Model (HMM)



- It is widely used in speech recognition, part-of-speech tagging, localization *etc.*
- Joint distribution

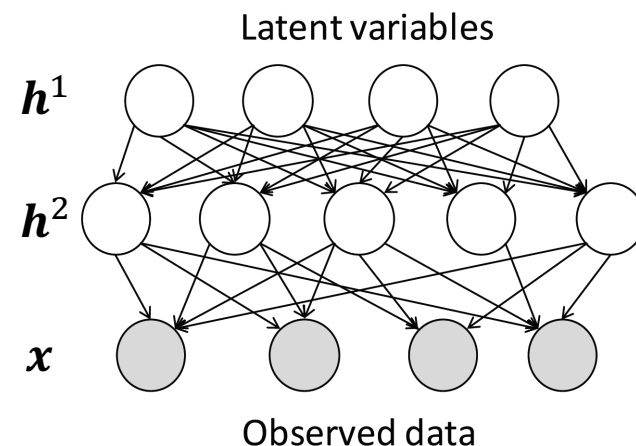
$$p(\mathbf{y}, \mathbf{x}) = p(y_1)p(x_1|y_1) \prod_{t=2}^T p(y_t|y_{t-1})p(x_t|y_t)$$

where $p(y_t|y_{t-1})$ is the transition probability; $p(x_t|y_t)$ is the emission probability

Application: Image Modeling

- Sigmoid belief networks (SBN)

- $h_i^1 \sim \text{Bernoulli}(0.5)$
- $h_j^2 \sim \text{Bernoulli}(\sigma([W_1 h^1 + b_1]_j))$
- $x_k \sim \text{Bernoulli}(\sigma([W_2 h^2 + b_2]_k))$



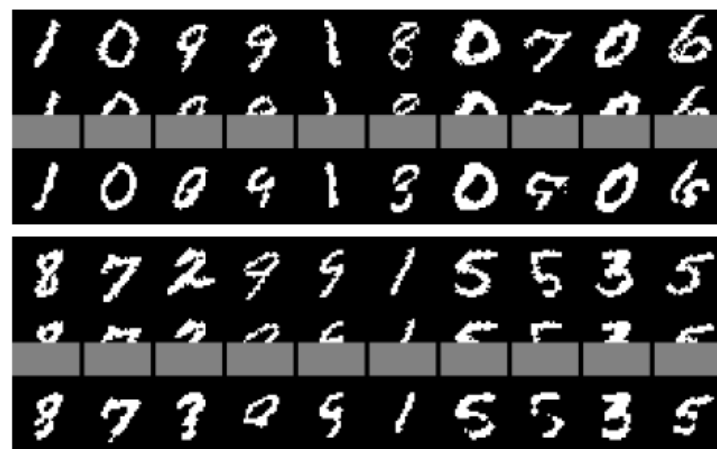
Joint pdf: $p(x, h^2, h^1) = p(x|h^2)p(h^2|h^1)p(h^1)$



Original



Generating



In-painting

Application: Text Modeling

- Topic Model: Latent Dirichlet Allocation (LDA)

- $\theta \sim Dir(\alpha)$: the distribution of different topics

- $\varphi_k \sim Dir(\beta)$: the distribution of words for topic

- $z_n \sim Multinomial(\theta)$: the topic of n -th word

- $w_n \sim Multinomial(\varphi_{z_n})$: the n -th word

