

Chapter 13

Entangled Models

supplementary slides to
Machine Learning Fundamentals
© **Hui Jiang 2020**
published by Cambridge University Press

August 2020



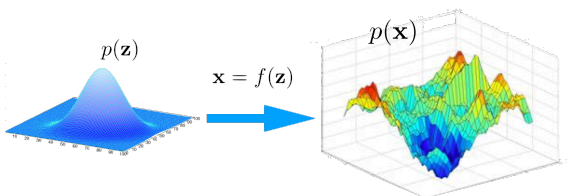
Outline

- 1 Formulation of Entangled Models
- 2 Linear Gaussian Models
- 3 Non-Gaussian Models
- 4 Deep Generative Models

Entangled Models (I)

Borel isomorphism theorem

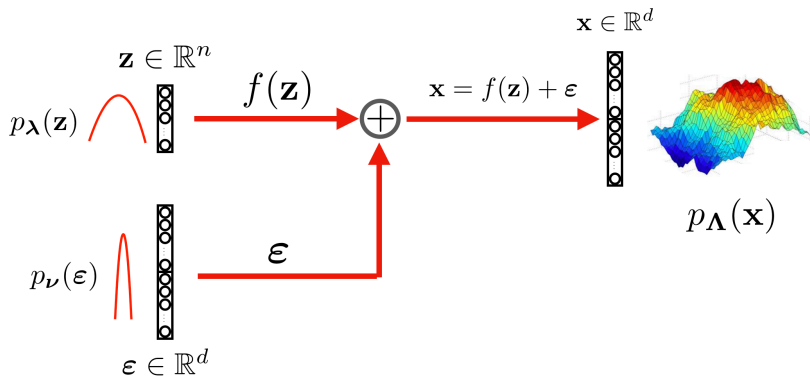
Given a normally distributed random variable $z \sim \mathcal{N}(0, 1)$, for any smooth probability distribution $p(\mathbf{x})$ ($\mathbf{x} \in \mathbb{R}^d$), there exist L^p functions: $f_1(z), \dots, f_d(z)$ to convert z into a vector $f(z) = [f_1(z) \ \dots \ f_d(z)]^\top$ so that $f(z)$ follows this distribution, i.e. $f(z) \sim p(\mathbf{x})$.



Entangled Models (II)

- **entangled models**: combine simple models with a complex transformation to derive complex models
- linear vs. non-linear transformations
- disentangled representation learning
 - entangling assumption: independent features \mathbf{z} , called *factors*, are entangled by a mixing function $\mathbf{x} = f(\mathbf{z})$ to observations \mathbf{x}
 - add residual noises to compute likelihood
 - disentangling: $\mathbf{x} \mapsto \mathbf{z}$
- three subgroups of entangled models:
 - linear Gaussian models: Gaussian + linear transformation
 - non-Gaussian models: non-Gaussian + linear transformation
 - deep generative models: Gaussian + neural networks

Entangled Models (III)



Entangled Models (IV)

| entangled models | factor $\mathbf{z} \sim p(\mathbf{z})$ | residual $\boldsymbol{\varepsilon} \sim p(\boldsymbol{\varepsilon})$ | mixing $f(\mathbf{z})$ |
|-------------------|--|---|---|
| probabilistic PCA | $\mathcal{N}(\mathbf{z} \mathbf{0}, \mathbf{I})$ | $\mathcal{N}(\boldsymbol{\varepsilon} \mathbf{0}, \sigma^2\mathbf{I})$ | $\mathbf{W}\mathbf{z}$ linear |
| factor analysis | $\mathcal{N}(\mathbf{z} \mathbf{0}, \mathbf{I})$ | $\mathcal{N}(\boldsymbol{\varepsilon} \mathbf{0}, \mathbf{D})$ \mathbf{D} : diagonal | $\mathbf{W}\mathbf{z}$ linear |
| ICA | $\prod_i p_i(z_i)$ non-Gaussian | — | $\mathbf{W}\mathbf{z}$ linear |
| IFA | $\prod_i p_i(z_i)$ factorial GMM | $\mathcal{N}(\boldsymbol{\varepsilon} \mathbf{0}, \boldsymbol{\Lambda})$ | $\mathbf{W}\mathbf{z}$ linear |
| HOPE | mixture model (movMF/GMM) | $\mathcal{N}(\boldsymbol{\varepsilon} \mathbf{0}, \sigma^2\mathbf{I})$ | $\mathbf{W}\mathbf{z}$ \mathbf{W} : orthogonal |
| VAE | $\mathcal{N}(\mathbf{z} \mathbf{0}, \mathbf{I})$ | $\mathcal{N}(\boldsymbol{\varepsilon} \mathbf{0}, \sigma^2\mathbf{I})$ | $f(\cdot) \in L^p$ neural nets |
| GAN | $\mathcal{N}(\mathbf{z} \mathbf{0}, \mathbf{I})$ | — | $f(\cdot) \in L^p$ neural nets |

Learning of Entangled Models

1 Jacobian method: if the mixing function is invertible and differentiable

- entangled models \implies maximum likelihood estimation

$$p_{\Lambda}(\mathbf{x}) = |\mathbf{J}| p_{\lambda}(f_1^{-1}(\mathbf{x})) p_{\nu}(f_2^{-1}(\mathbf{x}))$$

- disentangling: $\mathbf{z} = f_1^{-1}(\mathbf{x})$

2 marginalization method

- entangled models \implies maximum likelihood estimation

$$p_{\Lambda}(\mathbf{x}) = \int_{\mathbf{z}} p_{\lambda}(\mathbf{z}) p_{\nu}(\mathbf{x} - f(\mathbf{z}; \mathbf{W})) d\mathbf{z}$$

- disentangling:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{z}, \mathbf{x})}{p(\mathbf{x})} = \frac{p_{\lambda}(\mathbf{z}) p_{\nu}(\mathbf{x} - f(\mathbf{z}; \mathbf{W}))}{\int_{\mathbf{z}} p_{\lambda}(\mathbf{z}) p_{\nu}(\mathbf{x} - f(\mathbf{z}; \mathbf{W})) d\mathbf{z}}$$

Linear Gaussian Models

- linear Gaussian models are a group of simple entangled models
- factors \mathbf{z} follow a zero-mean multivariate Gaussian:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \Sigma_1)$$

- residual ε follows another multivariate Gaussian

$$p(\varepsilon) = \mathcal{N}(\varepsilon \mid \boldsymbol{\mu}, \Sigma_2)$$

- the mixing function is linear:

$$f(\mathbf{z}; \mathbf{W}) = \mathbf{W}\mathbf{z}$$

- linear Gaussian models:

$$p_{\Lambda}(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{W}\Sigma_1\mathbf{W}^{\top} + \Sigma_2)$$

Probabilistic PCA

- probabilistic PCA: a special case of linear Gaussian models
- factor distribution: $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I})$
- residual distribution: an isotropic covariance Gaussian
 $p_\sigma(\boldsymbol{\varepsilon}) = \mathcal{N}(\boldsymbol{\varepsilon} \mid \boldsymbol{\mu}, \sigma^2 \mathbf{I})$
- probabilistic PCA models: $p_\Lambda(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I})$
- the log-likelihood function

$$l(\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = C - \frac{N}{2} \ln |\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}| - \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I})^{-1} (\mathbf{x}_i - \boldsymbol{\mu})$$

- maximum likelihood estimation $\mathbf{W}_{\text{MLE}} \implies$ a rotation of principle components in regular PCA
- disentangling: $p(\mathbf{z} \mid \mathbf{x}) = \mathcal{N}(\mathbf{z} \mid \mathbf{M}^{-1} \mathbf{W}^\top (\mathbf{x} - \boldsymbol{\mu}), \sigma^{-2} \mathbf{M})$
where $\mathbf{M} = \mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}$

Factor Analysis (I)

- probabilistic PCA: a special case of linear Gaussian models
- factor distribution: $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \mathbf{0}, \mathbf{I})$
- residual distribution: a diagonal covariance Gaussian

$$p(\boldsymbol{\varepsilon}) = \mathcal{N}(\boldsymbol{\varepsilon} \mid \boldsymbol{\mu}, \mathbf{D})$$

where $\mathbf{D} \in \mathbb{R}^{d \times d}$ is a diagonal covariance matrix

- factor analysis models:

$$p_{\Lambda}(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{W}\mathbf{W}^{\top} + \mathbf{D})$$

- the log-likelihood function:

$$l(\mathbf{W}, \mathbf{D}) = C - \frac{N}{2} \left[\ln |\mathbf{W}\mathbf{W}^{\top} + \mathbf{D}| + \text{tr} \left((\mathbf{W}\mathbf{W}^{\top} + \mathbf{D})^{-1} \mathbf{S} \right) \right]$$

- no closed-form solution for maximum likelihood estimation

Factor Analysis (II)

alternating MLE method

Input: the sample covariance matrix \mathbf{S}

Output: \mathbf{W} and \mathbf{D}

randomly initialize \mathbf{D}_0 ; set $t = 1$

while not converged **do**

1. construct \mathbf{P}_t using the n leading eigenvectors of $\mathbf{D}_{t-1}^{-\frac{1}{2}} \mathbf{S} \mathbf{D}_{t-1}^{-\frac{1}{2}}$
2. $\mathbf{W}_t = \mathbf{D}_t^{\frac{1}{2}} \mathbf{P}_t$
3. $\mathbf{D}_t = \text{diag}(\mathbf{S} - \mathbf{W}_t \mathbf{W}_t^\top)$
4. $t = t + 1$

end while

Non-Gaussian Entangled Models

use a non-Gaussian factor distribution and a linear mixing function

■ Independent Component Analysis (ICA)

$$p(\mathbf{z}) = \prod_{j=1}^n p(z_j) = \prod_{j=1}^n \frac{4}{\pi(e^{z_j} + e^{-z_j})}$$

■ Independent Factor Analysis (IFA)

$$p(z_j) = \sum_{m=1}^M w_{jm} \mathcal{N}(z_j | \mu_{jm}, \sigma_{jm}^2)$$

■ Hybrid Orthogonal Projection and Estimation (HOPE)

- $p(\mathbf{z})$ is a mixture model in \mathbb{R}^n
- $p(\boldsymbol{\varepsilon})$ is a zero-mean isotropic covariance Gaussian in \mathbb{R}^{d-n} : $p(\boldsymbol{\varepsilon}) = \mathcal{N}(\boldsymbol{\varepsilon} | \mathbf{0}, \sigma^2 \mathbf{I})$
- \mathbf{W} is an orthogonal matrix

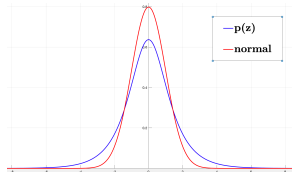


Figure: the normal distribution vs. the heavy-tail distribution $p(z)$ in ICA

Independent Component Analysis (ICA)

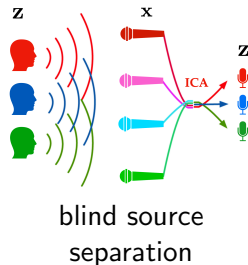
- ICA is used for blind source separation
- use a linear transformation to recover \mathbf{z} from \mathbf{x}

$$\mathbf{z} = \mathbf{W}^{-1}\mathbf{x}$$

- the log-likelihood function:

$$l(\mathbf{W}^{-1}) = \sum_{i=1}^N \sum_{j=1}^n \ln p(\mathbf{w}_j^T \mathbf{x}_i) + N \ln |\mathbf{W}^{-1}|$$

- maximum likelihood estimation: use any gradient descent methods to maximize $l(\mathbf{W}^{-1})$ with respect to \mathbf{W}^{-1}



Deep Generative Models

- combine a simple Gaussian with a complex non-linear mixing function
- use a deep neural network \mathbb{W} for the mixing function:

$$\mathbf{x} = f(\mathbf{z}; \mathbb{W})$$

- cannot explicitly evaluate the likelihood function
- no direct way to disentangle: $p(\mathbf{z}|\mathbf{x})$
- use some tricks to bypass these difficulties:
 - 1 variational autoencoders (VAE)
 - 2 generative adversarial nets (GAN)
 - 3 normalizing flows

Variational Autoencoders (VAE): Formulation

- use a Gaussian $q(\mathbf{z}|\mathbf{x})$ to approximate the intractable distribution $p(\mathbf{z}|\mathbf{x})$ as:

$$p(\mathbf{z}|\mathbf{x}) \approx q(\mathbf{z}|\mathbf{x})$$

- $q(\mathbf{z}|\mathbf{x}) \triangleq \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{x}})$ has \mathbf{x} -dependent mean and covariance
- introduce another deep neural network \mathbb{V} , called *encoder*:

$$[\boldsymbol{\mu}_{\mathbf{x}} \ \boldsymbol{\Sigma}_{\mathbf{x}}] = h(\mathbf{x}; \mathbb{V})$$

- derive a lower bound for the intractable log-likelihood:

$$\underbrace{\ln p(\mathbf{x})}_{l(\mathbb{W}, \sigma|\mathbf{x})} = \underbrace{\text{KL}(q(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}|\mathbf{x}))}_{\geq 0} + \underbrace{\left\{ \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\ln p(\mathbf{x}|\mathbf{z})] - \text{KL}(q(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z})) \right\}}_{L(\mathbb{W}, \mathbb{V}, \sigma|\mathbf{x})}$$

Variational Autoencoders (VAE): Optimization

- VAE aims to learn both \mathbb{W} and \mathbb{V} :

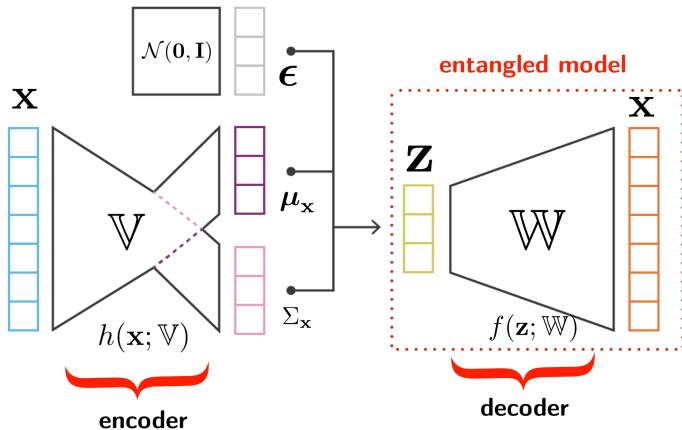
$$\arg \max_{\mathbb{W}, \mathbb{V}, \sigma} \sum_{i=1}^N L(\mathbb{W}, \mathbb{V}, \sigma | \mathbf{x}_i)$$

$$\implies \arg \max_{\mathbb{W}, \mathbb{V}, \sigma} \sum_{i=1}^N \left\{ \mathbb{E}_{q(\mathbf{z} | \mathbf{x}_i)} [\ln p(\mathbf{x}_i | \mathbf{z})] - \text{KL} \left(q(\mathbf{z} | \mathbf{x}_i) \parallel p(\mathbf{z}) \right) \right\}$$

- use sampling for expectation: $\mathbf{z}_j \sim q(\mathbf{z} | \mathbf{x}_i) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}_{\mathbf{x}_i}, \boldsymbol{\Sigma}_{\mathbf{x}_i})$, then $\mathbb{E}_{q(\mathbf{z} | \mathbf{x}_i)} [\ln p(\mathbf{x}_i | \mathbf{z})] \approx \frac{1}{G} \sum_{j=1}^G \ln p(\mathbf{x}_i | \mathbf{z}_j)$
- a reparameterization trick: $\boldsymbol{\epsilon}_j \sim \mathcal{N}(\boldsymbol{\epsilon} | \mathbf{0}, \mathbf{I})$, then

$$\mathbb{E}_{q(\mathbf{z} | \mathbf{x}_i)} [\ln p(\mathbf{x}_i | \mathbf{z})] \approx \frac{1}{G} \sum_{j=1}^G \ln p(\mathbf{x}_i | \boldsymbol{\Sigma}_{\mathbf{x}_i}^{\frac{1}{2}} \boldsymbol{\epsilon}_j + \boldsymbol{\mu}_{\mathbf{x}_i})$$

Variational Autoencoders (VAE)



Generative Adversarial Nets (GAN)

- introduce a discriminator \mathbb{V} to replace the intractable likelihood function
- use a pure sampling-based training procedure
- *equilibrium*: \mathbb{V} cannot distinguish fake samples from the true samples
- equilibrium \implies \mathbb{W} is a good entangled model for the data distribution

