CEE 616: Probabilistic Machine Learning

# M5 Unsupervised Learning:
# 5B: Factor Analysis and Autoencoders

**Jimi Oke**

UMass Amherst

College of Engineering

Dec 4, 2025

# Outline

**①** Factor analysis

**②** FA Estimation

**③** Autoencoders

**④** AE variants

**⑤** Outlook

# Factor analysis model

Factor analysis model

- Basic idea: there are latent (hidden) **common factors $z$** underlying some multivariate observations $\boldsymbol{x}_n \in \mathbb{R}^D$

## Factor analysis model

- Basic idea: there are latent (hidden) **common factors $z$** underlying some multivariate observations $\boldsymbol{x}_n \in \mathbb{R}^D$

Factor analysis (FA) is a latent variable generative model specified as:

# Factor analysis model

- Basic idea: there are latent (hidden) **common factors $z$** underlying some multivariate observations $\boldsymbol{x}_n \in \mathbb{R}^D$

Factor analysis (FA) is a latent variable generative model specified as:

$$p(\boldsymbol{z}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \tag{1}$$

## Factor analysis model

- Basic idea: there are latent (hidden) **common factors $z$** underlying some multivariate observations $\boldsymbol{x}_n \in \mathbb{R}^D$

Factor analysis (FA) is a latent variable generative model specified as:

$$
\begin{aligned}
p(\boldsymbol{z}) &= \mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) && (1) \\
p(\boldsymbol{x}|\boldsymbol{z}, \boldsymbol{\theta}) &= \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) && (2)
\end{aligned}
$$

# Factor analysis model

- Basic idea: there are latent (hidden) **common factors $z$** underlying some multivariate observations $\boldsymbol{x}_n \in \mathbb{R}^D$

Factor analysis (FA) is a latent variable generative model specified as:

$$
\begin{align}
p(\boldsymbol{z}) &= \mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \tag{1} \\
p(\boldsymbol{x}|\boldsymbol{z}, \boldsymbol{\theta}) &= \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \tag{2}
\end{align}
$$

where:

- $\boldsymbol{z}$: latent vector

## Factor analysis model

- Basic idea: there are latent (hidden) **common factors $z$** underlying some multivariate observations $\boldsymbol{x}_n \in \mathbb{R}^D$

Factor analysis (FA) is a latent variable generative model specified as:

$$
\begin{aligned}
p(\boldsymbol{z}) &= \mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \qquad (1)\\
p(\boldsymbol{x}|\boldsymbol{z}, \boldsymbol{\theta}) &= \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \qquad (2)
\end{aligned}
$$

where:

- $\boldsymbol{z}$: latent vector
- $\boldsymbol{W}$: factor loading matrix, $D \times L$

# Factor analysis model

- Basic idea: there are latent (hidden) **common factors $z$** underlying some multivariate observations $\boldsymbol{x}_n \in \mathbb{R}^D$

Factor analysis (FA) is a latent variable generative model specified as:

$$
\begin{array}{rcl}
p(\boldsymbol{z}) &=& \mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad\quad (1) \\
p(\boldsymbol{x}|\boldsymbol{z}, \boldsymbol{\theta}) &=& \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \quad\quad (2)
\end{array}
$$

where:

- $\boldsymbol{z}$: latent vector
- $\boldsymbol{W}$: factor loading matrix, $D \times L$
- $\boldsymbol{\Psi}$: diagonal covariance matrix, $D \times D$

# Factor analysis model

- Basic idea: there are latent (hidden) **common factors $z$** underlying some multivariate observations $\boldsymbol{x}_n \in \mathbb{R}^D$

Factor analysis (FA) is a latent variable generative model specified as:

$$
\begin{align}
p(\boldsymbol{z}) &= \mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \tag{1} \\
p(\boldsymbol{x}|\boldsymbol{z}, \boldsymbol{\theta}) &= \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \tag{2}
\end{align}
$$

where:

- $\boldsymbol{z}$: latent vector
- $\boldsymbol{W}$: factor loading matrix, $D \times L$
- $\boldsymbol{\Psi}$: diagonal covariance matrix, $D \times D$

# Induced marginal distribution

# Induced marginal distribution

$$p(\boldsymbol{x}|\boldsymbol{\theta}) \quad =$$

# Induced marginal distribution

$$p(\boldsymbol{x}|\boldsymbol{\theta}) = \int \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z}+\boldsymbol{\mu}, \boldsymbol{\Psi})\mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)d\boldsymbol{z} \qquad (3)$$

# Induced marginal distribution

$$p(\boldsymbol{x}|\boldsymbol{\theta}) = \int \mathcal{N}(\boldsymbol{x}|\boldsymbol{Wz}+\boldsymbol{\mu},\boldsymbol{\Psi})\mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0,\boldsymbol{\Sigma}_0)d\boldsymbol{z} \tag{3}$$

$$= \mathcal{N}(\boldsymbol{x}|\underbrace{\boldsymbol{W}\boldsymbol{\mu}_0+\boldsymbol{\mu}}_{\text{mean}},\underbrace{\boldsymbol{\Psi}+\boldsymbol{W}\boldsymbol{\Sigma}_0\boldsymbol{W}^\top}_{\text{variance}}) \tag{4}$$

# Induced marginal distribution

$$p(\boldsymbol{x}|\boldsymbol{\theta}) = \int \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) d\boldsymbol{z} \tag{3}$$

$$= \mathcal{N}(\boldsymbol{x}|\underbrace{\boldsymbol{W}\boldsymbol{\mu}_0 + \boldsymbol{\mu}}_{\text{mean}}, \underbrace{\boldsymbol{\Psi} + \boldsymbol{W}\boldsymbol{\Sigma}_0\boldsymbol{W}^\top}_{\text{variance}}) \tag{4}$$

Simplifications:

- $\boldsymbol{\mu}_0 \to \boldsymbol{0}$

# Induced marginal distribution

$$
\begin{aligned}
p(\boldsymbol{x}|\boldsymbol{\theta}) &= \int \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z}+\boldsymbol{\mu}, \boldsymbol{\Psi})\mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)d\boldsymbol{z} && (3) \\
&= \mathcal{N}(\boldsymbol{x}|\underbrace{\boldsymbol{W}\boldsymbol{\mu}_0+\boldsymbol{\mu}}_{\text{mean}}, \underbrace{\boldsymbol{\Psi} + \boldsymbol{W}\boldsymbol{\Sigma}_0\boldsymbol{W}^\top}_{\text{variance}}) && (4)
\end{aligned}
$$

Simplifications:

- $\boldsymbol{\mu}_0 \rightarrow \boldsymbol{0}$
- $\boldsymbol{\Sigma}_0 \rightarrow \boldsymbol{I}$

## Induced marginal distribution

$$p(\boldsymbol{x}|\boldsymbol{\theta}) = \int \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi})\mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)d\boldsymbol{z} \tag{3}$$

$$= \mathcal{N}(\boldsymbol{x}|\underbrace{\boldsymbol{W}\boldsymbol{\mu}_0 + \boldsymbol{\mu}}_{\text{mean}}, \underbrace{\boldsymbol{\Psi} + \boldsymbol{W}\boldsymbol{\Sigma}_0\boldsymbol{W}^{\top}}_{\text{variance}}) \tag{4}$$

Simplifications:

- $\boldsymbol{\mu}_0 \to \boldsymbol{0}$
- $\boldsymbol{\Sigma}_0 \to \boldsymbol{I}$

The simplified marginal distribution then becomes:

## Induced marginal distribution

$$
\begin{aligned}
p(\boldsymbol{x}|\boldsymbol{\theta}) &= \int \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z}+\boldsymbol{\mu}, \boldsymbol{\Psi})\mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)d\boldsymbol{z} \quad (3)\\
&= \mathcal{N}(\boldsymbol{x}|\underbrace{\boldsymbol{W}\boldsymbol{\mu}_0+\boldsymbol{\mu}}_{\text{mean}}, \underbrace{\boldsymbol{\Psi}+\boldsymbol{W}\boldsymbol{\Sigma}_0\boldsymbol{W}^\top}_{\text{variance}}) \quad (4)
\end{aligned}
$$

Simplifications:

- $\boldsymbol{\mu}_0 \rightarrow \boldsymbol{0}$
- $\boldsymbol{\Sigma}_0 \rightarrow \boldsymbol{I}$

The simplified marginal distribution then becomes:

$$
p(\boldsymbol{x}|\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{W}\boldsymbol{W}^\top + \boldsymbol{\Psi}) \quad (5)
$$

## Generative model simplified

$$[\text{Prior}] \quad p(z) \;\; = \;\; \mathcal{N}(z|\mathbf{0}, I) \tag{6}$$

## Generative model simplified

$$[\text{Prior}] \quad p(\boldsymbol{z}) \;=\; \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \tag{6}$$
$$[\text{Likelihood}] \quad p(\boldsymbol{x}|\boldsymbol{z})$$

## Generative model simplified

$$[\text{Prior}] \quad p(\boldsymbol{z}) \;=\; \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \tag{6}$$

$$[\text{Likelihood}] \quad p(\boldsymbol{x}|\boldsymbol{z}) \;=\; \mathcal{N}(\boldsymbol{x}|\boldsymbol{Wz} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \tag{7}$$

# Generative model simplified

$$
\begin{array}{rcll}
[\text{Prior}] & p(\mathbf{z}) & = & \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}) \hspace{2em} (6) \\
[\text{Likelihood}] & p(\mathbf{x}|\mathbf{z}) & = & \mathcal{N}(\mathbf{x}|\mathbf{Wz} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \hspace{2em} (7) \\
[\text{Evidence/marginal}] & p(\mathbf{x}) & &
\end{array}
$$

## Generative model simplified

$$
\begin{array}{rrcl}
[\text{Prior}] & p(\boldsymbol{z}) & = & \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \quad\quad (6) \\
[\text{Likelihood}] & p(\boldsymbol{x}|\boldsymbol{z}) & = & \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \quad\quad (7) \\
[\text{Evidence/marginal}] & p(\boldsymbol{x}) & = & \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{W}\boldsymbol{W}^{\top} + \boldsymbol{\Psi}) \quad\quad (8)
\end{array}
$$

# What FA does

## What FA does

It approximates the covariance matrix of the visible/observed vector $\boldsymbol{x}$ using a low-rank decomposition:

## What FA does

It approximates the covariance matrix of the visible/observed vector $\boldsymbol{x}$ using a low-rank decomposition:

$$\boldsymbol{C} = \mathrm{Cov}[\boldsymbol{x}] = \underbrace{\boldsymbol{W}\boldsymbol{W}^\top + \boldsymbol{\Psi}}_{\text{low-rank decomp}} \tag{9}$$

## What FA does

It approximates the covariance matrix of the visible/observed vector $\boldsymbol{x}$ using a low-rank decomposition:

$$\boldsymbol{C} = \text{Cov}[\boldsymbol{x}] = \underbrace{\boldsymbol{W}\boldsymbol{W}^\top + \boldsymbol{\Psi}}_{\text{low-rank decomp}} \tag{9}$$

- $\boldsymbol{W}\boldsymbol{W}^\top$ is $D \times D$ (recall: $\boldsymbol{W} \in \mathbb{R}^{D \times L}$)

## What FA does

It approximates the covariance matrix of the visible/observed vector $\boldsymbol{x}$ using a low-rank decomposition:

$$\boldsymbol{C} = \text{Cov}[\boldsymbol{x}] = \underbrace{\boldsymbol{W}\boldsymbol{W}^\top + \boldsymbol{\Psi}}_{\text{low-rank decomp}} \tag{9}$$

- $\boldsymbol{W}\boldsymbol{W}^\top$ is $D \times D$ (recall: $\boldsymbol{W} \in \mathbb{R}^{D \times L}$)
- $\boldsymbol{\Psi}$ is $D \times D$ (restricted to be diagonal)

## What FA does

It approximates the covariance matrix of the visible/observed vector $\boldsymbol{x}$ using a low-rank decomposition:

$$\boldsymbol{C} = \text{Cov}[\boldsymbol{x}] = \underbrace{\boldsymbol{W}\boldsymbol{W}^{\top} + \boldsymbol{\Psi}}_{\text{low-rank decomp}} \tag{9}$$

- $\boldsymbol{W}\boldsymbol{W}^{\top}$ is $D \times D$ (recall: $\boldsymbol{W} \in \mathbb{R}^{D \times L}$)
- $\boldsymbol{\Psi}$ is $D \times D$ (restricted to be diagonal)
- For each variable $x_d$, the **marginal variance** (each diagonal term in $\boldsymbol{C}$) is given by:

## What FA does

It approximates the covariance matrix of the visible/observed vector $\boldsymbol{x}$ using a low-rank decomposition:

$$\boldsymbol{C} = \text{Cov}[\boldsymbol{x}] = \underbrace{\boldsymbol{W}\boldsymbol{W}^\top + \boldsymbol{\Psi}}_{\text{low-rank decomp}} \tag{9}$$

- $\boldsymbol{W}\boldsymbol{W}^\top$ is $D \times D$ (recall: $\boldsymbol{W} \in \mathbb{R}^{D \times L}$)
- $\boldsymbol{\Psi}$ is $D \times D$ (restricted to be diagonal)
- For each variable $x_d$, the **marginal variance** (each diagonal term in $\boldsymbol{C}$) is given by:

$$\mathbb{V}[x_d] = \sum_{k=1}^{L} \underbrace{w_{dk}^2}_{\text{common}} + \underbrace{\psi_d}_{\text{unique}} \tag{10}$$

# Parameters to be estimated

## Parameters to be estimated

The unknown parameters in FA are:

## Parameters to be estimated

The unknown parameters in FA are:

- $W$: factor loading matrix

## Parameters to be estimated

The unknown parameters in FA are:

- $W$: factor loading matrix
- $\Psi$: covariance matrix

# Parameters to be estimated

The unknown parameters in FA are:

- $W$: factor loading matrix
- $\Psi$: covariance matrix

These can be estimated via:

- Maximum likelihood estimation (MLE)

# Parameters to be estimated

The unknown parameters in FA are:

- $W$: factor loading matrix
- $\Psi$: covariance matrix

These can be estimated via:

- Maximum likelihood estimation (MLE)
- Expectation-maximization (EM) algorithm

# Parameters to be estimated

The unknown parameters in FA are:

- $W$: factor loading matrix
- $\Psi$: covariance matrix

These can be estimated via:

- Maximum likelihood estimation (MLE)
- Expectation-maximization (EM) algorithm
- Bayesian methods (e.g., variational inference, MCMC)

# Parameters to be estimated

The unknown parameters in FA are:

- $W$: factor loading matrix
- $\Psi$: covariance matrix

These can be estimated via:

- Maximum likelihood estimation (MLE)
- Expectation-maximization (EM) algorithm
- Bayesian methods (e.g., variational inference, MCMC)

Once estimated, the **posterior** of latent embeddings is given by:

## Parameters to be estimated

The unknown parameters in FA are:

- $W$: factor loading matrix
- $\Psi$: covariance matrix

These can be estimated via:

- Maximum likelihood estimation (MLE)
- Expectation-maximization (EM) algorithm
- Bayesian methods (e.g., variational inference, MCMC)

Once estimated, the **posterior** of latent embeddings is given by:

$$p(z|x) = \mathcal{N}(z|W^\top C^{-1}(x - \mu), I - W^\top C^{-1} W) \tag{11}$$

## Parameters to be estimated

The unknown parameters in FA are:

- $W$: factor loading matrix
- $\Psi$: covariance matrix

These can be estimated via:

- Maximum likelihood estimation (MLE)
- Expectation-maximization (EM) algorithm
- Bayesian methods (e.g., variational inference, MCMC)

Once estimated, the **posterior** of latent embeddings is given by:

$$p(z|x) = \mathcal{N}(z|W^\top C^{-1}(x - \mu), I - W^\top C^{-1}W) \tag{11}$$

- Posterior has closed-form solution under Gaussian distribution

Factor analysis
000000●
FA Estimation
00
Autoencoders
000
AE variants
0000
Outlook
0

# FA model summary

## FA model summary

$$[\text{Prior}] \quad p(\boldsymbol{z}) \;=\; \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \tag{12}$$

## FA model summary

$$[\text{Prior}] \quad p(\boldsymbol{z}) \;=\; \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \tag{12}$$

$$[\text{Likelihood}] \quad p(\boldsymbol{x}|\boldsymbol{z})$$

# FA model summary

$$[\text{Prior}] \quad p(\boldsymbol{z}) \quad = \quad \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \tag{12}$$

$$[\text{Likelihood}] \quad p(\boldsymbol{x}|\boldsymbol{z}) \quad = \quad \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \tag{13}$$

## FA model summary

$$
\begin{align}
[\text{Prior}] \quad p(\boldsymbol{z}) &= \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \tag{12} \\
[\text{Likelihood}] \quad p(\boldsymbol{x}|\boldsymbol{z}) &= \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \tag{13} \\
[\text{Evidence/marginal}] \quad p(\boldsymbol{x}) & 
\end{align}
$$

# FA model summary

$$[\text{Prior}] \quad p(\boldsymbol{z}) \;=\; \mathcal{N}(\boldsymbol{z}|\boldsymbol{0},\boldsymbol{I}) \tag{12}$$

$$[\text{Likelihood}] \quad p(\boldsymbol{x}|\boldsymbol{z}) \;=\; \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z}+\boldsymbol{\mu},\boldsymbol{\Psi}) \tag{13}$$

$$[\text{Evidence/marginal}] \quad p(\boldsymbol{x}) \;=\; \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{C}) \tag{14}$$

## FA model summary

$$\begin{align}
[\text{Prior}] \quad p(\boldsymbol{z}) &= \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \tag{12} \\
[\text{Likelihood}] \quad p(\boldsymbol{x}|\boldsymbol{z}) &= \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \tag{13} \\
[\text{Evidence/marginal}] \quad p(\boldsymbol{x}) &= \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{C}) \tag{14} \\
[\text{Posterior}] \quad p(\boldsymbol{z}|\boldsymbol{x}) &= \mathcal{N}(\boldsymbol{z}|\boldsymbol{W}^\top \boldsymbol{C}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}), \boldsymbol{I} - \boldsymbol{W}^\top \boldsymbol{C}^{-1} \boldsymbol{W}) \tag{15}
\end{align}$$

## FA model summary

$$
\begin{array}{rrcl}
\text{[Prior]} & p(\boldsymbol{z}) & = & \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \qquad (12) \\
\text{[Likelihood]} & p(\boldsymbol{x}|\boldsymbol{z}) & = & \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \qquad (13) \\
\text{[Evidence/marginal]} & p(\boldsymbol{x}) & = & \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{C}) \qquad (14) \\
\text{[Posterior]} & p(\boldsymbol{z}|\boldsymbol{x}) & = & \mathcal{N}(\boldsymbol{z}|\boldsymbol{W}^\top \boldsymbol{C}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}), \boldsymbol{I} - \boldsymbol{W}^\top \boldsymbol{C}^{-1}\boldsymbol{W}) \quad (15)
\end{array}
$$

- $\boldsymbol{z}$: latent vector, length $L$ (assumed to be zero-mean, unit variance)

## FA model summary

$$[\text{Prior}] \quad p(\boldsymbol{z}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \tag{12}$$

$$[\text{Likelihood}] \quad p(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \tag{13}$$

$$[\text{Evidence/marginal}] \quad p(\boldsymbol{x}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{C}) \tag{14}$$

$$[\text{Posterior}] \quad p(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{W}^{\top}\boldsymbol{C}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}), \boldsymbol{I} - \boldsymbol{W}^{\top}\boldsymbol{C}^{-1}\boldsymbol{W}) \tag{15}$$

- $\boldsymbol{z}$: latent vector, length $L$ (assumed to be zero-mean, unit variance)
- $\boldsymbol{x}$: observed vector

## FA model summary

$$
\begin{align}
[\text{Prior}] \quad p(z) &= \mathcal{N}(z|\mathbf{0}, I) \tag{12} \\
[\text{Likelihood}] \quad p(x|z) &= \mathcal{N}(x|Wz + \mu, \Psi) \tag{13} \\
[\text{Evidence/marginal}] \quad p(x) &= \mathcal{N}(x|\mu, C) \tag{14} \\
[\text{Posterior}] \quad p(z|x) &= \mathcal{N}(z|W^\top C^{-1}(x - \mu), I - W^\top C^{-1}W) \tag{15}
\end{align}
$$

- $z$: latent vector, length $L$ (assumed to be zero-mean, unit variance)
- $x$: observed vector
- $W$: $D \times L$ factor loading matrix

## FA model summary

$$
\begin{array}{rcll}
[\text{Prior}] & p(\boldsymbol{z}) & = & \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) & (12) \\
[\text{Likelihood}] & p(\boldsymbol{x}|\boldsymbol{z}) & = & \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) & (13) \\
[\text{Evidence/marginal}] & p(\boldsymbol{x}) & = & \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{C}) & (14) \\
[\text{Posterior}] & p(\boldsymbol{z}|\boldsymbol{x}) & = & \mathcal{N}(\boldsymbol{z}|\boldsymbol{W}^\top \boldsymbol{C}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}), \boldsymbol{I} - \boldsymbol{W}^\top \boldsymbol{C}^{-1} \boldsymbol{W}) & (15)
\end{array}
$$

- $\boldsymbol{z}$: latent vector, length $L$ (assumed to be zero-mean, unit variance)
- $\boldsymbol{x}$: observed vector
- $\boldsymbol{W}$: $D \times L$ factor loading matrix
- $\boldsymbol{\Psi}$: $D \times D$ diagonal covariance matrix or **matrix of unique variances**

# FA model summary

$$[\text{Prior}] \quad p(\boldsymbol{z}) \ = \ \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \tag{12}$$

$$[\text{Likelihood}] \quad p(\boldsymbol{x}|\boldsymbol{z}) \ = \ \mathcal{N}(\boldsymbol{x}|\boldsymbol{W}\boldsymbol{z} + \boldsymbol{\mu}, \boldsymbol{\Psi}) \tag{13}$$

$$[\text{Evidence/marginal}] \quad p(\boldsymbol{x}) \ = \ \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}, \boldsymbol{C}) \tag{14}$$

$$[\text{Posterior}] \quad p(\boldsymbol{z}|\boldsymbol{x}) \ = \ \mathcal{N}(\boldsymbol{z}|\boldsymbol{W}^\top \boldsymbol{C}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}), \boldsymbol{I} - \boldsymbol{W}^\top \boldsymbol{C}^{-1}\boldsymbol{W}) \tag{15}$$

- $\boldsymbol{z}$: latent vector, length $L$ (assumed to be zero-mean, unit variance)
- $\boldsymbol{x}$: observed vector
- $\boldsymbol{W}$: $D \times L$ factor loading matrix
- $\boldsymbol{\Psi}$: $D \times D$ diagonal covariance matrix or **matrix of unique variances**
- $\boldsymbol{C} = \boldsymbol{W}\boldsymbol{W}^\top + \boldsymbol{\Psi}$

# Unidentifiability of FA parameters

Factor analysis
000000

FA Estimation
●○

Autoencoders
000

AE variants
0000

Outlook
○

# Unidentifiability of FA parameters

The parameters $W$ and $\Psi$ are unidentifiable.

Factor analysis
000000

FA Estimation
●○

Autoencoders
000

AE variants
0000

Outlook
○

Unidentifiability of FA parameters

The parameters $W$ and $\Psi$ are unidentifiable. This can be addressed via:

Factor analysis
000000
FA Estimation
●0
Autoencoders
000
AE variants
0000
Outlook
0

# Unidentifiability of FA parameters

The parameters $W$ and $\Psi$ are unidentifiable. This can be addressed via:

- Constraining $W$ to have orthonormal columns [PCA]

Factor analysis
000000

FA Estimation
●○

Autoencoders
000

AE variants
0000

Outlook
○

# Unidentifiability of FA parameters

The parameters $W$ and $\Psi$ are unidentifiable. This can be addressed via:

- Constraining $W$ to have orthonormal columns [PCA]
- Constraining $W$ to be lower triangular

Factor analysis
000000

FA Estimation
●○

Autoencoders
000

AE variants
0000

Outlook
○

## Unidentifiability of FA parameters

The parameters $\boldsymbol{W}$ and $\boldsymbol{\Psi}$ are unidentifiable. This can be addressed via:

- Constraining $\boldsymbol{W}$ to have orthonormal columns [PCA]
- Constraining $\boldsymbol{W}$ to be lower triangular
- Informative rotation: $\tilde{\boldsymbol{W}} = \boldsymbol{W}\boldsymbol{R}$, where $\boldsymbol{R}$ is the rotation matrix

# Unidentifiability of FA parameters

The parameters $\boldsymbol{W}$ and $\boldsymbol{\Psi}$ are unidentifiable. This can be addressed via:

- Constraining $\boldsymbol{W}$ to have orthonormal columns [PCA]
- Constraining $\boldsymbol{W}$ to be lower triangular
- Informative rotation: $\tilde{\boldsymbol{W}} = \boldsymbol{WR}$, where $\boldsymbol{R}$ is the rotation matrix
    - Commonly used rotations: Varimax, Promax, Oblimin, Geomin, Thurston, Equamax

Factor analysis
000000

FA Estimation
●○

Autoencoders
000

AE variants
0000

Outlook
○

## Unidentifiability of FA parameters

The parameters $W$ and $\Psi$ are unidentifiable. This can be addressed via:

- Constraining $W$ to have orthonormal columns [PCA]
- Constraining $W$ to be lower triangular
- Informative rotation: $\tilde{W} = WR$, where $R$ is the rotation matrix
  - Commonly used rotations: Varimax, Promax, Oblimin, Geomin, Thurstone, Equamax
- Sparsity-promoting priors on $W$

Factor analysis
000000

FA Estimation
●○

Autoencoders
000

AE variants
0000

Outlook
○

# Unidentifiability of FA parameters

The parameters $W$ and $\Psi$ are unidentifiable. This can be addressed via:

- Constraining $W$ to have orthonormal columns [PCA]
- Constraining $W$ to be lower triangular
- Informative rotation: $\tilde{W} = WR$, where $R$ is the rotation matrix
  - Commonly used rotations: Varimax, Promax, Oblimin, Geomin, Thurstone, Equamax
- Sparsity-promoting priors on $W$
- Non-Gaussian priors for latent factors

Factor analysis
000000

FA Estimation
○●

Autoencoders
000

AE variants
0000

Outlook
○

PCA as a special case of FA

# PCA as a special case of FA

Principal components analysis (PCA) is a special case of FA with:

# PCA as a special case of FA

Principal components analysis (PCA) is a special case of FA with:

$$\mathbf{\Psi} = \sigma^2 \mathbf{I} \tag{16}$$

# PCA as a special case of FA

Principal components analysis (PCA) is a special case of FA with:

$$\mathbf{\Psi} = \sigma^2 \mathbf{I} \tag{16}$$

where $\sigma^2$ is the isotropic noise variance.

# PCA as a special case of FA

Principal components analysis (PCA) is a special case of FA with:

$$\mathbf{\Psi} = \sigma^2 \mathbf{I} \tag{16}$$

where $\sigma^2$ is the isotropic noise variance. In PCA, the covariance of the observed vector is given by:

Factor analysis
000000

FA Estimation
00●

Autoencoders
000

AE variants
0000

Outlook
0

# PCA as a special case of FA

Principal components analysis (PCA) is a special case of FA with:

$$\mathbf{\Psi} = \sigma^2 \mathbf{I} \tag{16}$$

where $\sigma^2$ is the isotropic noise variance. In PCA, the covariance of the observed vector is given by:

$$\mathbf{C} = \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I} \tag{17}$$

# Autoencoders as nonlinear PCA/FA

In PCA/FA, we learn a **linear mapping** from a high-dimensional observed space $\boldsymbol{x} \in \mathbb{R}^D$ to a low-dimensional latent space $\boldsymbol{z} \in \mathbb{R}^L$ and vice-versa.

## Autoencoders as nonlinear PCA/FA

In PCA/FA, we learn a **linear mapping** from a high-dimensional observed space $\mathbf{x} \in \mathbb{R}^D$ to a low-dimensional latent space $\mathbf{z} \in \mathbb{R}^L$ and vice-versa.

- **Encoder** $f_e$: mapping from $\mathbf{x} \rightarrow \mathbf{z}$

Factor analysis
000000

FA Estimation
00

Autoencoders
●00

AE variants
0000

Outlook
0

## Autoencoders as nonlinear PCA/FA

In PCA/FA, we learn a **linear mapping** from a high-dimensional observed space $\mathbf{x} \in \mathbb{R}^D$ to a low-dimensional latent space $\mathbf{z} \in \mathbb{R}^L$ and vice-versa.

- **Encoder** $f_e$: mapping from $\mathbf{x} \to \mathbf{z}$
- **Decoder** $f_d$: mapping from $\mathbf{z} \to \mathbf{x}$

## Autoencoders as nonlinear PCA/FA

In PCA/FA, we learn a **linear mapping** from a high-dimensional observed space $\mathbf{x} \in \mathbb{R}^D$ to a low-dimensional latent space $\mathbf{z} \in \mathbb{R}^L$ and vice-versa.

- **Encoder** $f_e$: mapping from $\mathbf{x} \to \mathbf{z}$
- **Decoder** $f_d$: mapping from $\mathbf{z} \to \mathbf{x}$

In PCA, for example, $f_e$ is given by:

Factor analysis
000000

FA Estimation
00

Autoencoders
●○○

AE variants
0000

Outlook
○

## Autoencoders as nonlinear PCA/FA

In PCA/FA, we learn a **linear mapping** from a high-dimensional observed space $\boldsymbol{x} \in \mathbb{R}^D$ to a low-dimensional latent space $\boldsymbol{z} \in \mathbb{R}^L$ and vice-versa.

- **Encoder** $f_e$: mapping from $\boldsymbol{x} \to \boldsymbol{z}$
- **Decoder** $f_d$: mapping from $\boldsymbol{z} \to \boldsymbol{x}$

In PCA, for example, $f_e$ is given by:

$$\boldsymbol{z} = \boldsymbol{W}^\top \boldsymbol{x} \equiv f_e(\boldsymbol{x}) \tag{18}$$

## Autoencoders as nonlinear PCA/FA

In PCA/FA, we learn a **linear mapping** from a high-dimensional observed space $\boldsymbol{x} \in \mathbb{R}^D$ to a low-dimensional latent space $\boldsymbol{z} \in \mathbb{R}^L$ and vice-versa.

- **Encoder** $f_e$: mapping from $\boldsymbol{x} \to \boldsymbol{z}$
- **Decoder** $f_d$: mapping from $\boldsymbol{z} \to \boldsymbol{x}$

In PCA, for example, $f_e$ is given by:

$$\boldsymbol{z} = \boldsymbol{W}^\top \boldsymbol{x} \equiv f_e(\boldsymbol{x}) \tag{18}$$

And $f_d$ is given by:

## Autoencoders as nonlinear PCA/FA

In PCA/FA, we learn a **linear mapping** from a high-dimensional observed space $\boldsymbol{x} \in \mathbb{R}^D$ to a low-dimensional latent space $\boldsymbol{z} \in \mathbb{R}^L$ and vice-versa.

- **Encoder** $f_e$: mapping from $\boldsymbol{x} \to \boldsymbol{z}$
- **Decoder** $f_d$: mapping from $\boldsymbol{z} \to \boldsymbol{x}$

In PCA, for example, $f_e$ is given by:

$$\boldsymbol{z} = \boldsymbol{W}^\top \boldsymbol{x} \equiv f_e(\boldsymbol{x}) \tag{18}$$

And $f_d$ is given by:

$$\hat{\boldsymbol{x}} = \boldsymbol{W}\boldsymbol{z} \equiv f_d(\boldsymbol{z}) \tag{19}$$

To introduce flexibility, we can specify $f_e$ and $f_e$ are nonlinear/more complex functions. This is best accomplished via neural network, resulting in an **autoencoder** (AE).

Reconstruction loss

The reconstruction function is the approximation of the observation from the decoder:

The reconstruction function is the approximation of the observation from the decoder:

$$\hat{\boldsymbol{x}} \equiv r(\boldsymbol{x}) = f_d(f_e(\boldsymbol{x})) \tag{20}$$

Reconstruction loss

The reconstruction function is the approximation of the observation from the decoder:

$$\hat{\boldsymbol{x}} \equiv r(\boldsymbol{x}) = f_d(f_e(\boldsymbol{x})) \tag{20}$$

An autoencoder is thus trained to minimize the reconstruction loss

Reconstruction loss

The reconstruction function is the approximation of the observation from the decoder:

$$\hat{\boldsymbol{x}} \equiv r(\boldsymbol{x}) = f_d(f_e(\boldsymbol{x})) \tag{20}$$

An autoencoder is thus trained to minimize the reconstruction loss

$$\mathcal{L}(\boldsymbol{\theta}) = ||r(\boldsymbol{x}) - \boldsymbol{x}||_2^2 \tag{21}$$

or equivalently, the negative log-likelihood:

# Reconstruction loss

The reconstruction function is the approximation of the observation from the decoder:

$$\hat{\boldsymbol{x}} \equiv r(\boldsymbol{x}) = f_d(f_e(\boldsymbol{x})) \tag{20}$$

An autoencoder is thus trained to minimize the reconstruction loss

$$\mathcal{L}(\boldsymbol{\theta}) = ||r(\boldsymbol{x}) - \boldsymbol{x}||_2^2 \tag{21}$$

or equivalently, the negative log-likelihood:

$$\mathcal{L}(\boldsymbol{\theta}) = -\log p(\boldsymbol{x}|r(\boldsymbol{x})) \tag{22}$$

# Basic autoencoder (AE) architecture

# Basic autoencoder (AE) architecture

Autoencoder with 2 single-layer MLPS: input layer, hidden layer (latent representation) and output layer (reconstruction)
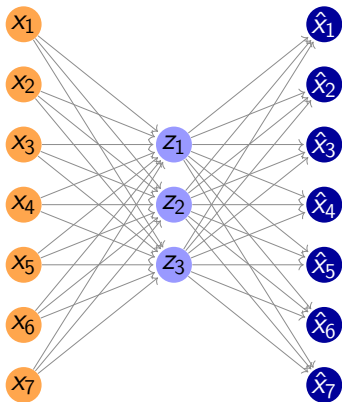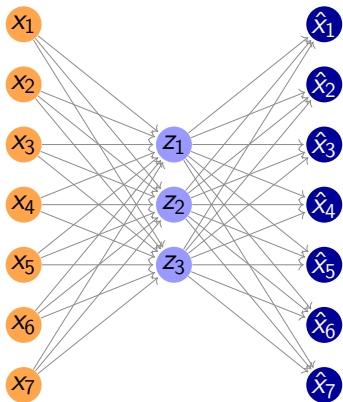
# Basic autoencoder (AE) architecture

Autoencoder with 2 single-layer MLPS: input layer, hidden layer (latent representation) and output layer (reconstruction)



- Hidden layer (size $L$) is a low-dimensional **bottleneck** between input and reconstruction

# Basic autoencoder (AE) architecture

Autoencoder with 2 single-layer MLPS: input layer, hidden layer (latent representation) and output layer (reconstruction)



- Hidden layer (size $L$) is a low-dimensional **bottleneck** between input and reconstruction
- $L \ll D$: undercomplete representation

## Basic autoencoder (AE) architecture

Autoencoder with 2 single-layer MLPS: input layer, hidden layer (latent representation) and output layer (reconstruction)



- Hidden layer (size $L$) is a low-dimensional **bottleneck** between input and reconstruction
- $L \ll D$: undercomplete representation
- $L \gg D$: overcomplete representation (regularize to prevent identity learning)

# Denoising autoencoders

Factor analysis
000000

FA Estimation
00

Autoencoders
000

AE variants
●000

Outlook
○

# Denoising autoencoders

In denoising autoencoders (DAEs), the input is corrupted ($\tilde{\boldsymbol{x}}$) by:

- Gaussian noise: $p_c(\tilde{\boldsymbol{x}}|\boldsymbol{x}) = \mathcal{N}(\tilde{\boldsymbol{x}}|\boldsymbol{x}, \sigma^2 \boldsymbol{I})$

# Denoising autoencoders

In denoising autoencoders (DAEs), the input is corrupted ($\tilde{\boldsymbol{x}}$) by:

- Gaussian noise: $p_c(\tilde{\boldsymbol{x}}|\boldsymbol{x}) = \mathcal{N}(\tilde{\boldsymbol{x}}|\boldsymbol{x}, \sigma^2\boldsymbol{I})$
- Bernoulli dropout: randomly setting a proportion of input nodes to zero



Schematic of a DAE.

Source: https://lilianweng.github.io/posts/2018-08-12-vae/

The model is then trained to minimize the loss between the reconstructed input $r(\tilde{\boldsymbol{x}})$ and its uncorrupted version $\boldsymbol{x}$
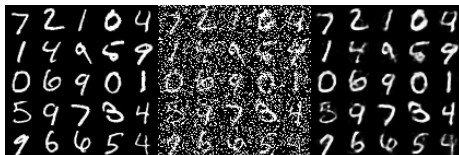
Factor analysis
000000

FA Estimation
00

Autoencoders
000

AE variants
0●00

Outlook
0

# Uses of DAE

Factor analysis
○○○○○○

FA Estimation
○○

Autoencoders
○○○

AE variants
○●○○

Outlook
○

Uses of DAE

- DAEs are used for denoising images

Factor analysis
000000

FA Estimation
00

Autoencoders
000

AE variants
0●00

Outlook
O

# Uses of DAE

- DAEs are used for denoising images



Original, corrupted and reconstructed images from MNIST dataset.

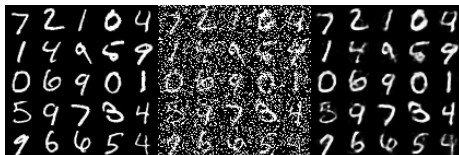Source: http://www.opendeep.org/v0.0.5/docs/tutorial-your-first-model

- They can also learn vector fields of input data

Factor analysis
000000

FA Estimation
00

Autoencoders
000

AE variants
0●00

Outlook
0

Uses of DAE

- DAEs are used for denoising images



Original, corrupted and reconstructed images from MNIST dataset.

Source: http://www.opendeep.org/v0.0.5/docs/tutorial-your-first-model

- They can also learn vector fields of input data
- Popular for their simplicity

# Sparse autoencoder (SAE)

## Sparse autoencoder (SAE)

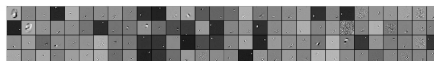**Sparse autoencoder** (SAE): sparsity penalty on latent activations
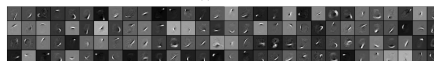
# Sparse autoencoder (SAE)

**Sparse autoencoder** (SAE): sparsity penalty on latent activations

$$\Omega(\boldsymbol{z}) = \lambda ||\boldsymbol{z}||_1 \qquad (23)$$

- *k*-Sparse autoencoder: use only *k* largest activations in training



(a) $k = 70$

(b) $k = 40$

(c) $k = 25$

(d) $k = 10$

Filters of the *k*-sparse autoencoder for different sparsity levels k, learnt from MNIST with 1000 hidden units.

Source: https://arxiv.org/pdf/1312.5663.pdf

Useful for interpretability

Factor analysis
000000
FA Estimation
00
Autoencoders
000
AE variants
000●
Outlook
0

Other AEs

- **Contractive autoencoder** (CAE): regularizes via penalty on reconstruction
  loss

Other AEs

- **Contractive autoencoder** (CAE): regularizes via penalty on reconstruction loss

$$\Omega(\boldsymbol{z}, \boldsymbol{x}) = \lambda \left\| \frac{\partial f_e(\boldsymbol{x})}{\partial \boldsymbol{x}} \right\|_F^2 \tag{24}$$

## Other AEs

- **Contractive autoencoder** (CAE): regularizes via penalty on reconstruction loss

$$\Omega(\boldsymbol{z}, \boldsymbol{x}) = \lambda \left\| \frac{\partial f_e(\boldsymbol{x})}{\partial \boldsymbol{x}} \right\|_F^2 \tag{24}$$

- Variational autoencoder (VAE): probablistic version of AE/generative model

# Reading

- **PMLI** 20.3
- **DL** 20