

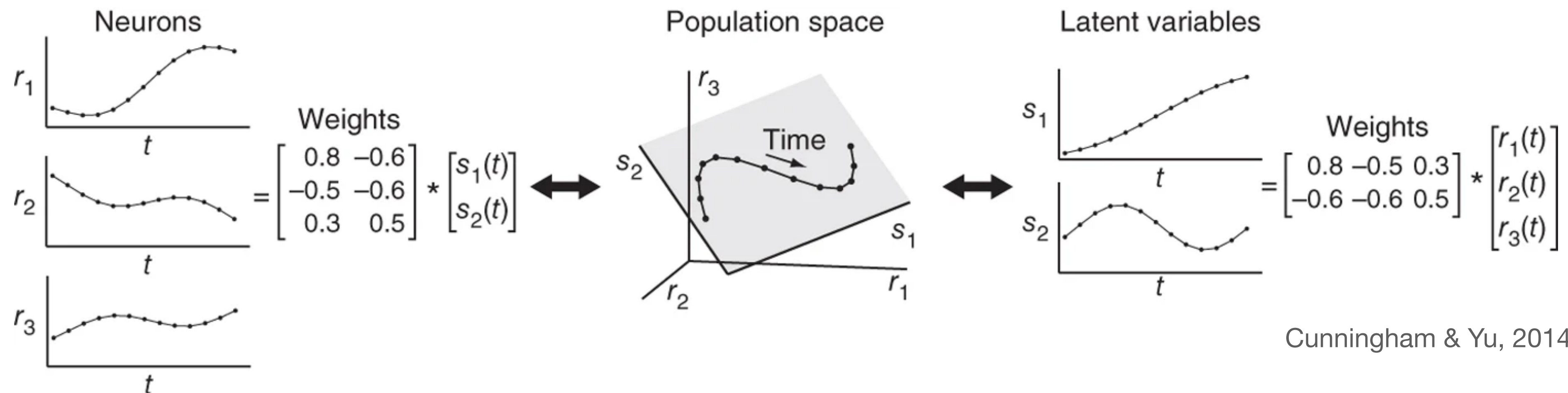
# Latent Variable Models

# Overview of Latent Variable Models

- **Motivation:** Find “latent” (unobserved) structure in neural population activity
- Latents can be discrete or continuous

# Overview of Latent Variable Models

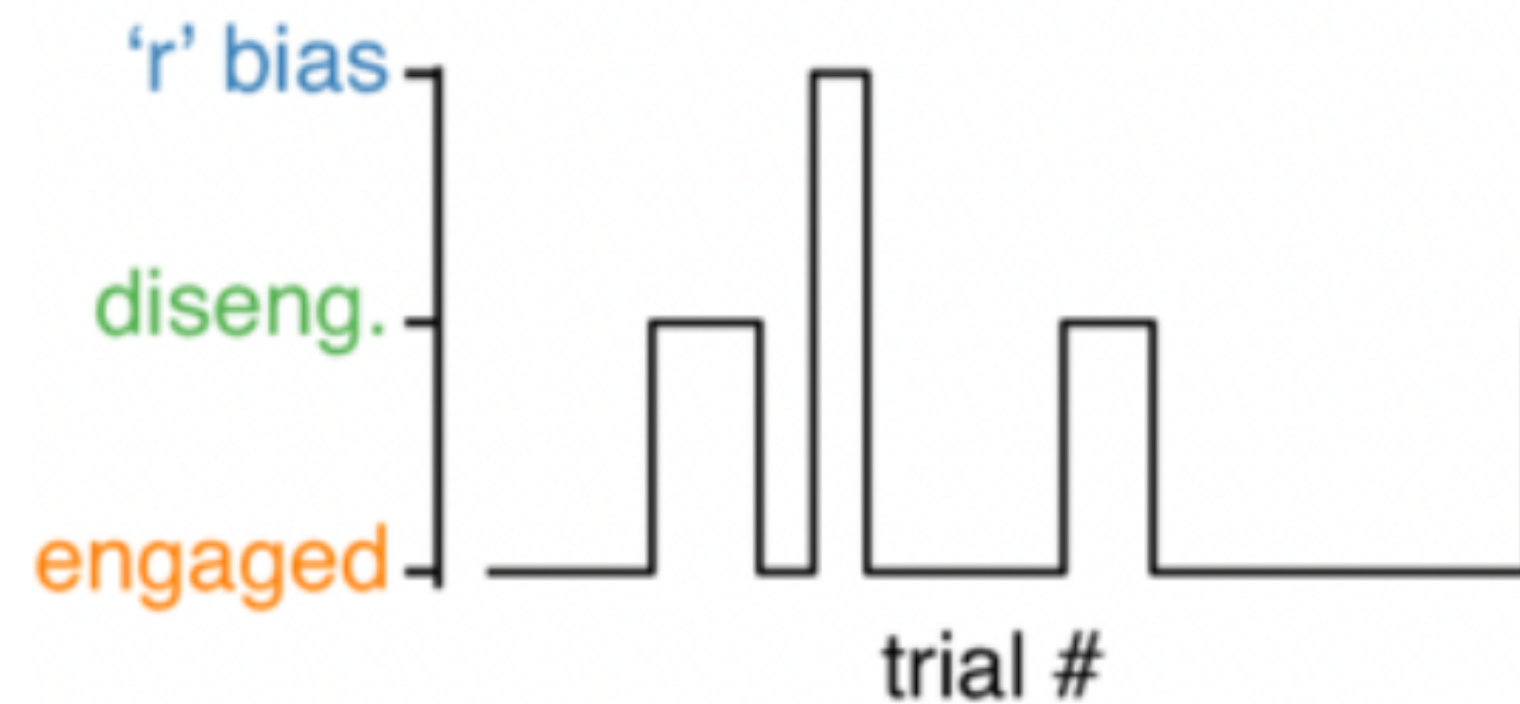
- **Motivation:** Find “latent” (unobserved) structure in neural population activity
- Latents can be discrete or **continuous**



Cunningham & Yu, 2014

# Overview of Latent Variable Models

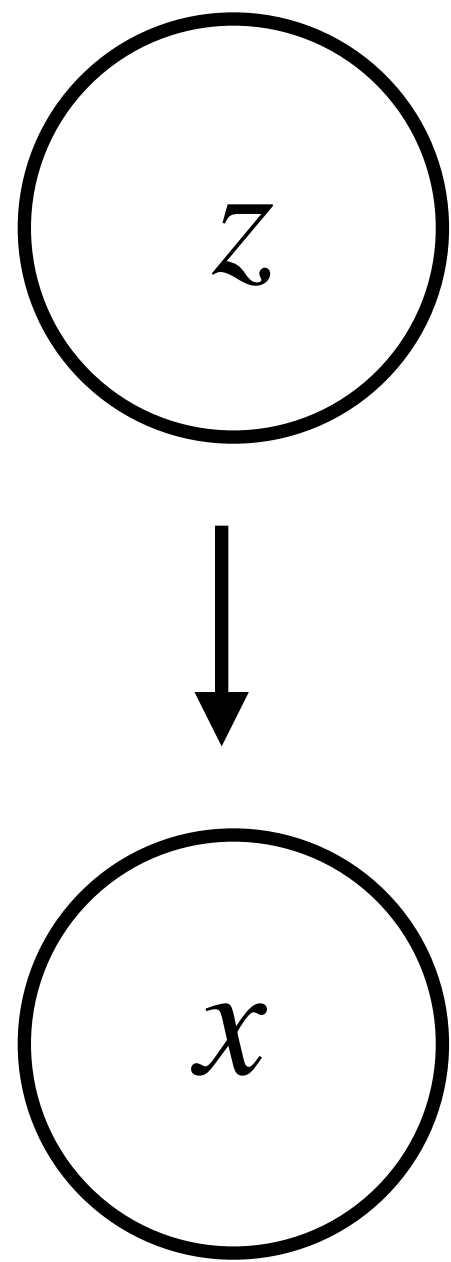
- **Motivation:** Find “latent” (unobserved) structure in neural population activity
- Latents can be **discrete** or continuous



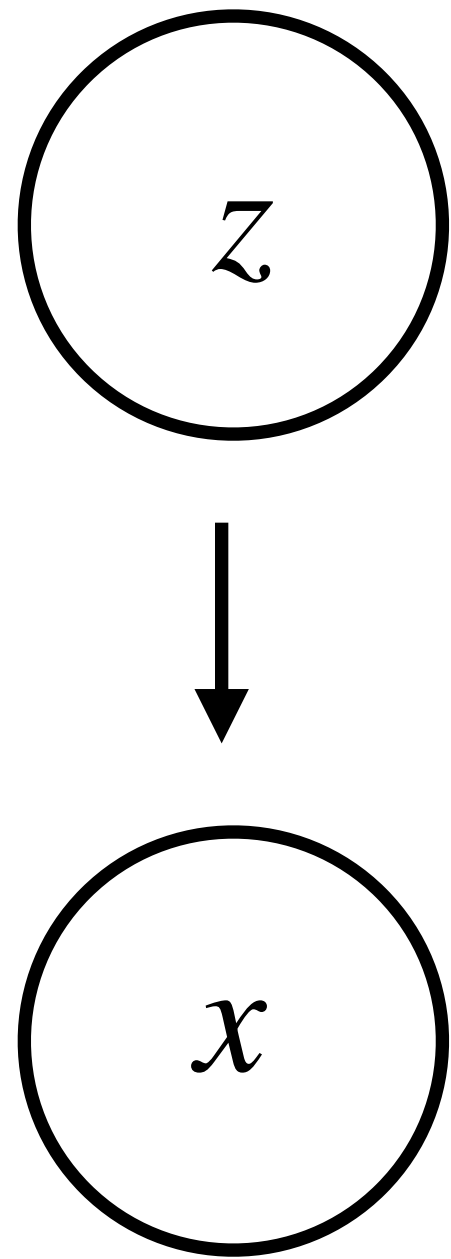
# Overview of Latent Variable Models

- Class 1: Intro to LVMs, Factor Analysis (FA), Gaussian Processes (GPs)
- Class 2: Hidden Markov Models (HMMs), Linear Dynamical Systems (LDS)
- Class 3: JC on GLM-HMMs and GPFA
- Class 4: Statistical inference: Expectation-Maximization, Markov Chain Monte Carlo (MCMC), Variational Inference
- Class 5: Switching Dynamical Systems, Poisson Linear Dynamical Systems, Variational Autoencoders (VAEs)
- Class 6: JC on LFADS (sequential VAEs), RSLDSs in neural data
- Anything else you'd like covered?

# Intro to Latent Variable Models

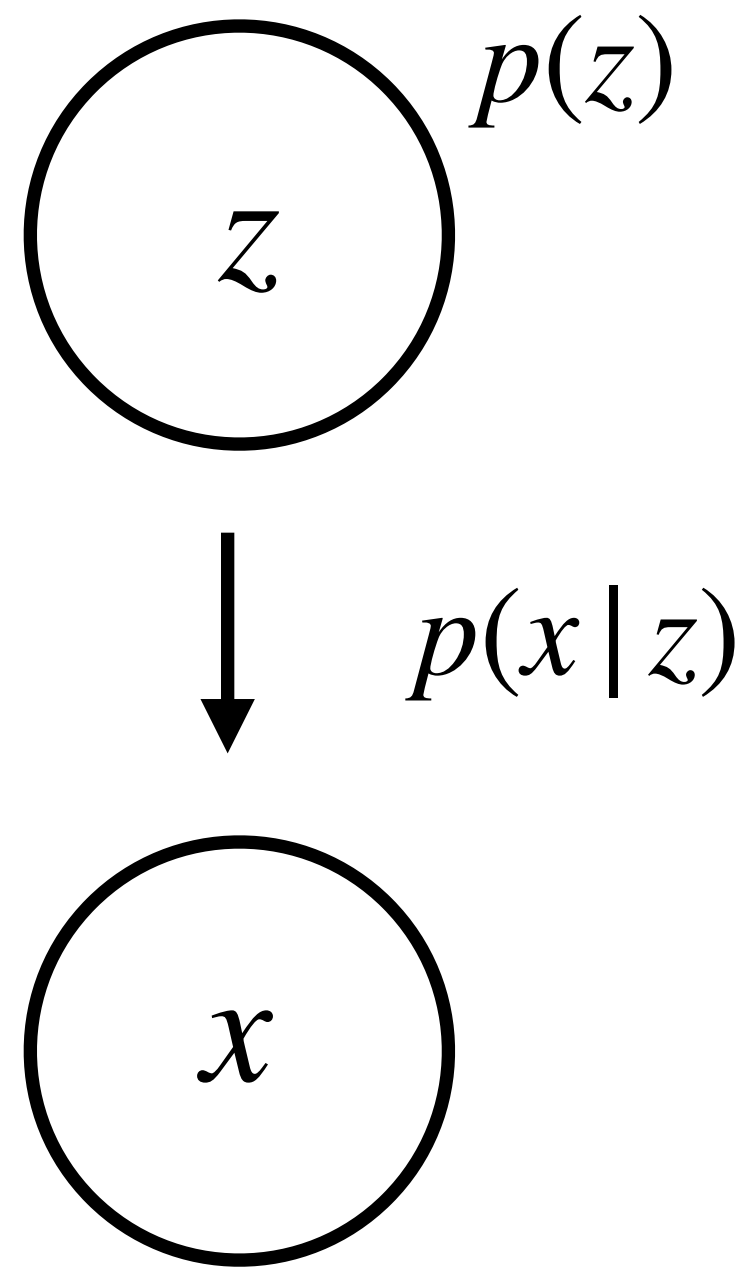


# Intro to Latent Variable Models



- Two parts of a LVM
  - Prior:  $z \sim p(z)$
- Conditional probability of observed data:  $x|z \sim p(x|z)$

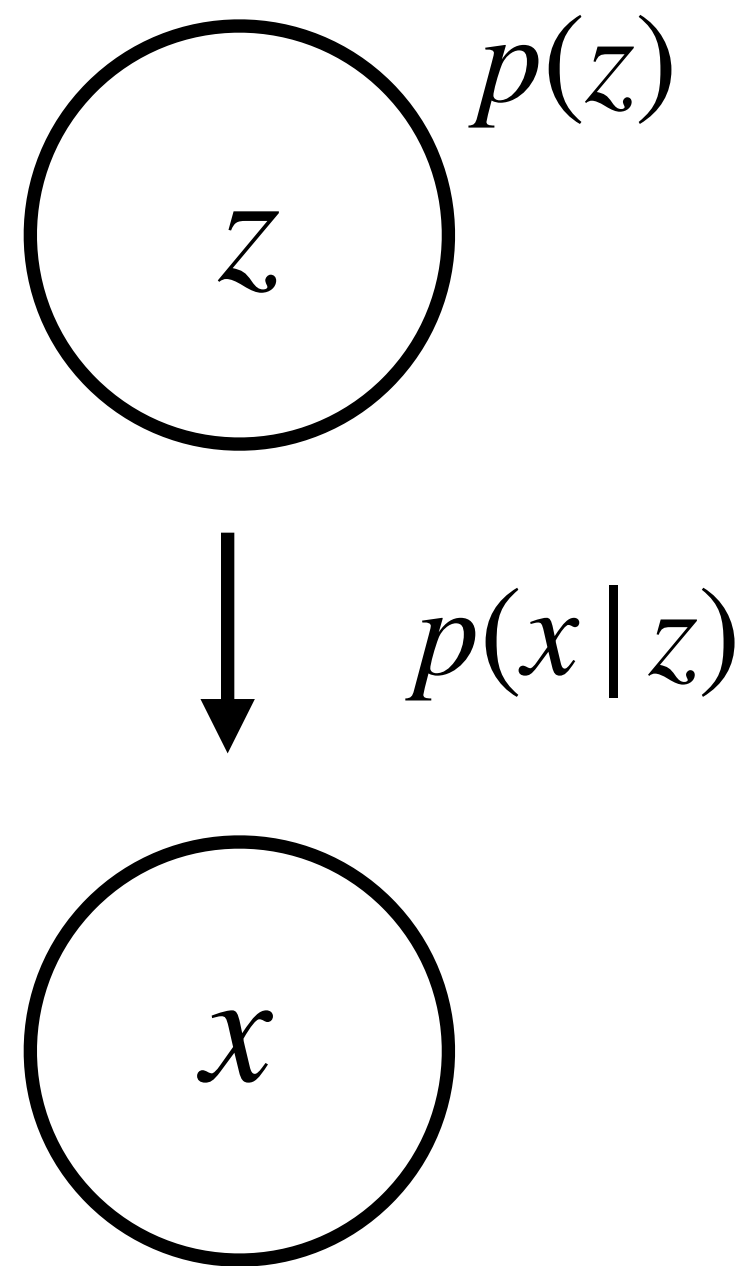
# Intro to Latent Variable Models



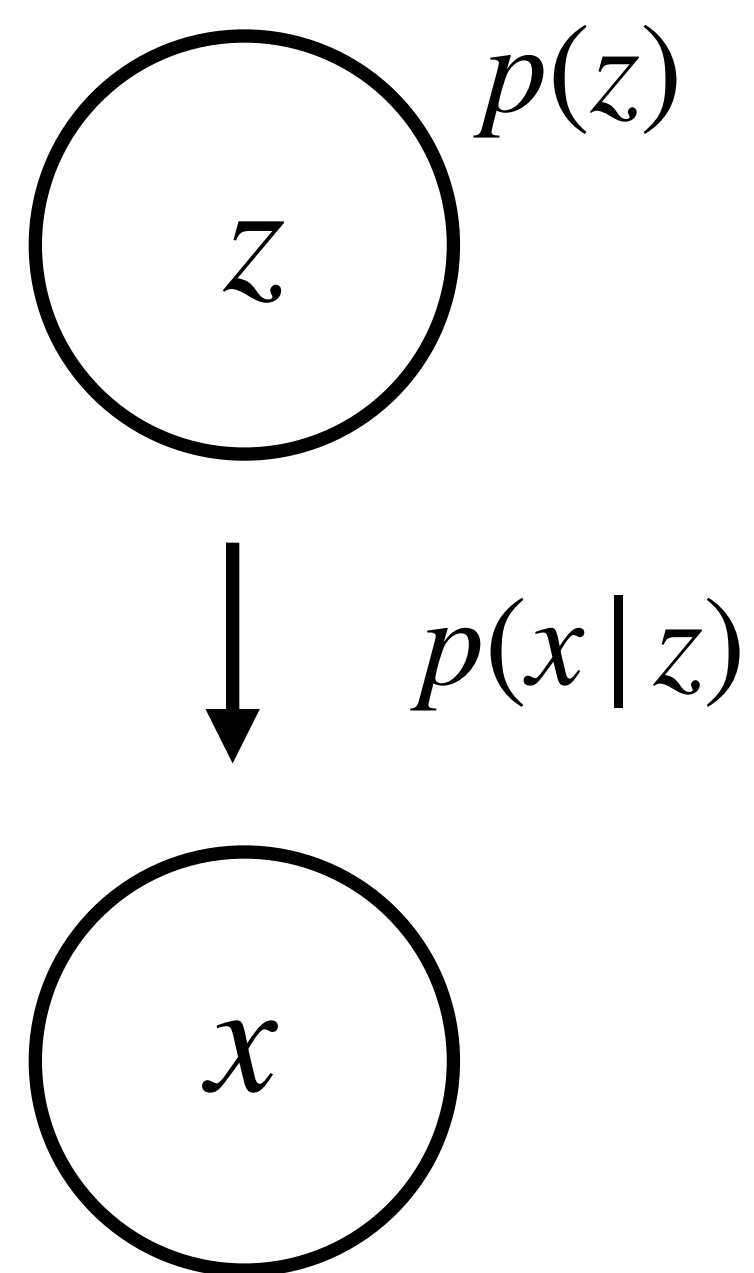


# Intro to Latent Variable Models

- Probability of observed data,  $p(x)$  , is:



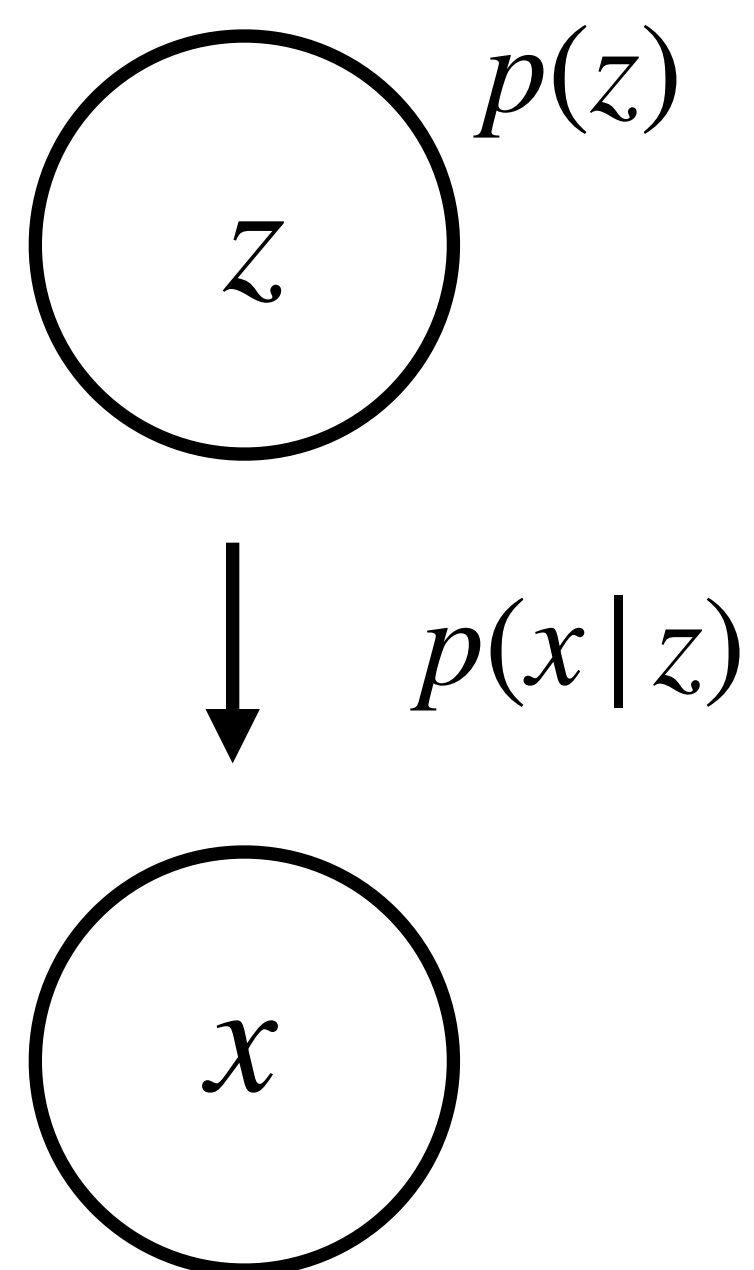
# Intro to Latent Variable Models



- Probability of observed data,  $p(x)$  , is:
- For discrete latents:

$$p(x) = \sum_{i=1}^m p(x|z = z_i)p(z = z_i)$$

# Intro to Latent Variable Models



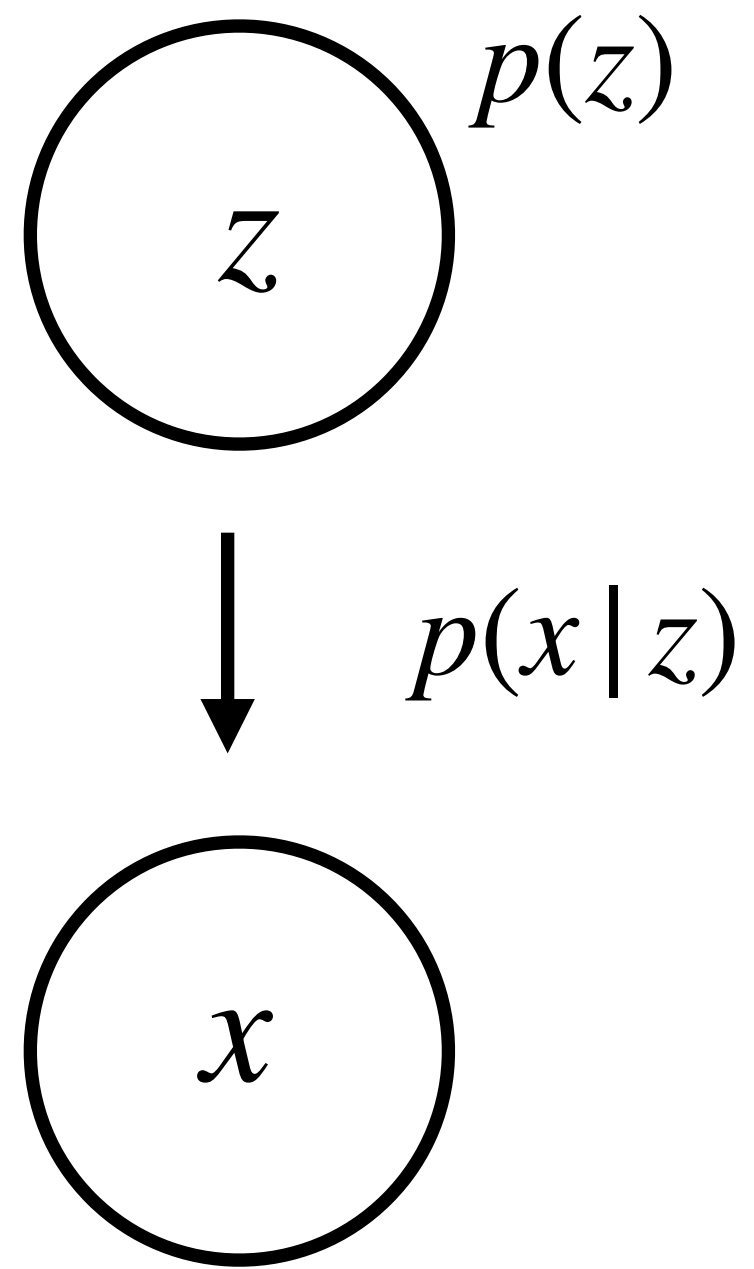
- Probability of observed data,  $p(x)$  , is:
  - For discrete latents:

$$p(x) = \sum_{i=1}^m p(x|z = z_i)p(z = z_i)$$

- For continuous latents:

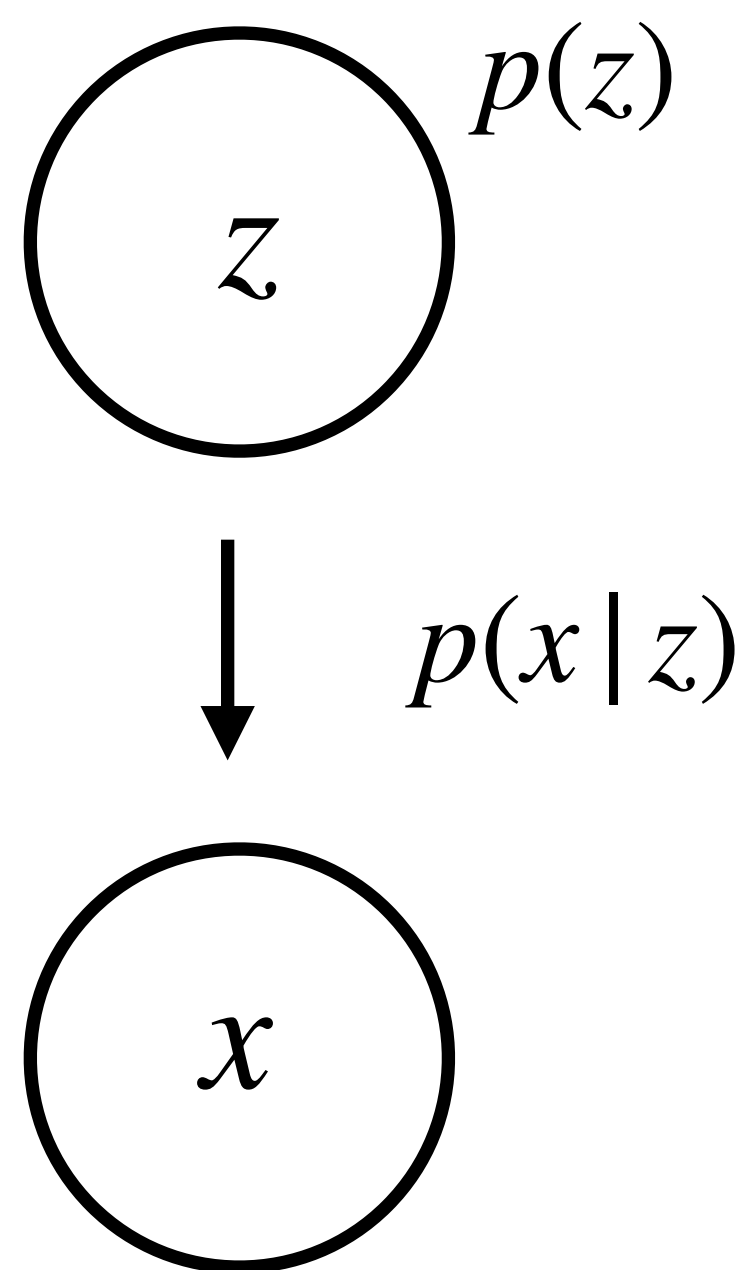
$$p(x) = \int p(x|z)p(z)dz$$

# Intro to Latent Variable Models: Goals



# Intro to Latent Variable Models: Goals

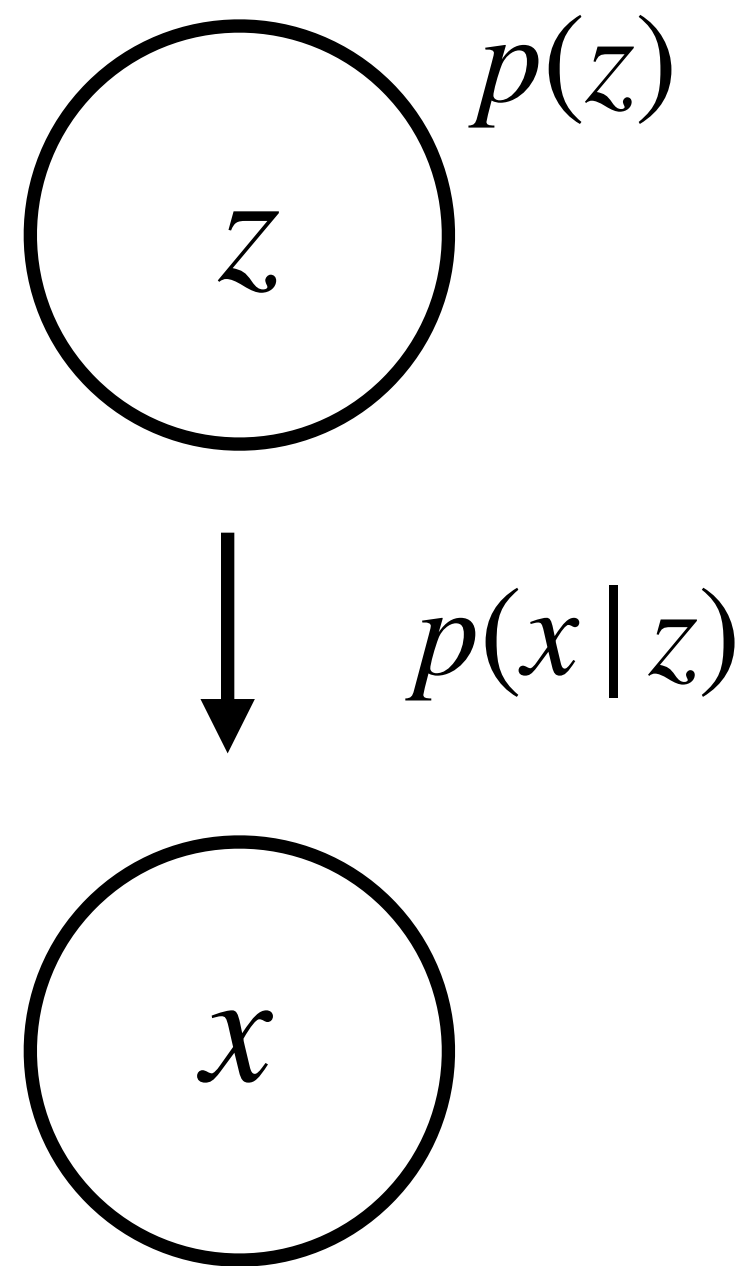
- Recognition/Inference



$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

# Intro to Latent Variable Models: Goals

- **Recognition/Inference**



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

- **Model Fitting**

- Model including parameters is actually:

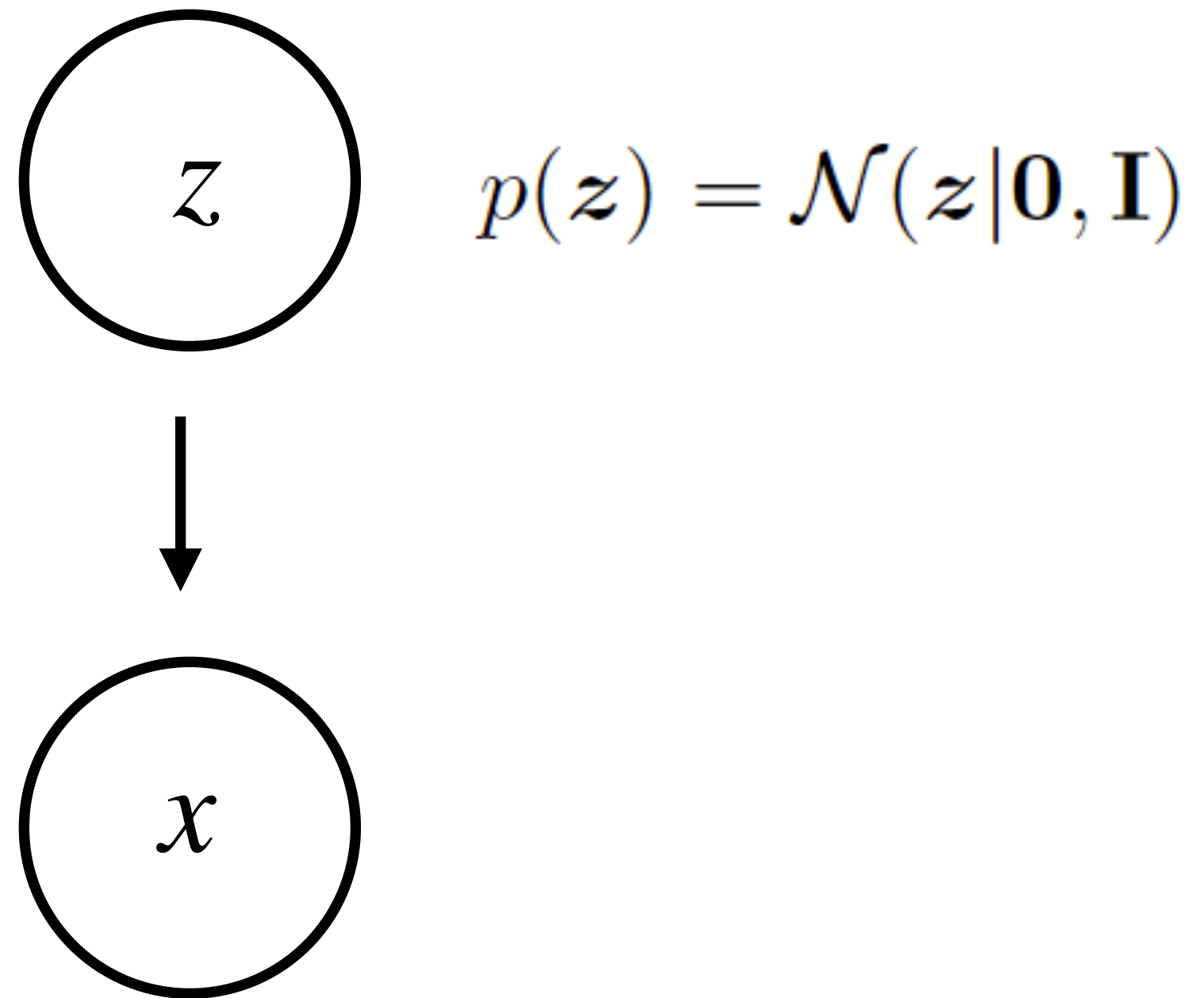
$$p(x, z|\theta) = p(x|z, \theta)p(z|\theta)$$

- Learning parameters by maximum likelihood:

$$\hat{\theta} = \arg \max_{\theta} p(x|\theta) = \arg \max_{\theta} \int p(x, z|\theta) dz.$$

# Factor Analysis

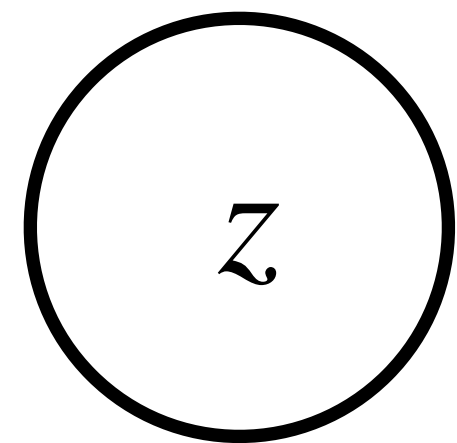
# Factor Analysis: Generative Model



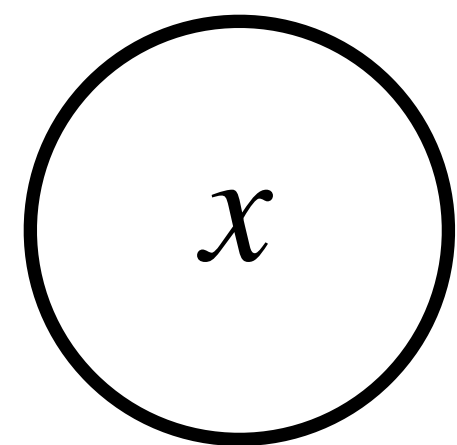
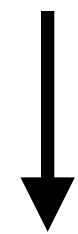
- Can be any Gaussian (see Murphy, book 1, section 20.2)



# Factor Analysis: Generative Model



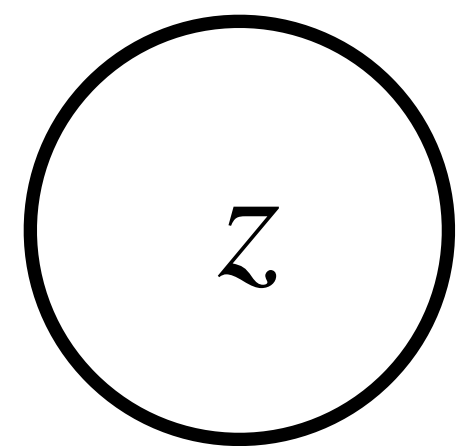
$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$$



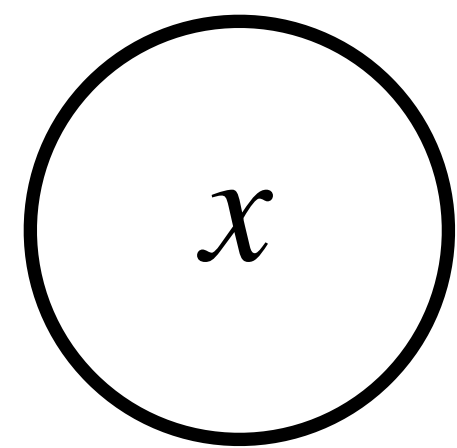
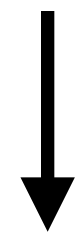
$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \boldsymbol{\Psi})$$

- Can be any Gaussian (see Murphy, book 1, section 20.2)
- Linear Gaussian model
- $\mathbf{z} : D \times T$  latent. dim x samples (timepoints)
- $\mathbf{x} : N \times T$  obs. dim (neurons) x samples (timepoints)
- $\mathbf{W} : N \times D$  obs. dim. (neurons) x latent dim.
- $\boldsymbol{\Psi} : D \times D$  diagonal covariance matrix

# Factor Analysis: Generative Model



$$p(z) = \mathcal{N}(z|\mathbf{0}, \mathbf{I})$$



$$p(x|z) = \mathcal{N}(x|\mathbf{W}z + \mu, \Psi)$$

$$p(x) = \int p(x|z)p(z)dz$$

$$p(x) = \mathcal{N}(x|\mu, \mathbf{W}\mathbf{W}^T + \Psi)$$

- Can be any Gaussian (see Murphy, book 1, section 20.2)
- Linear Gaussian model
- $z$  :  $D \times T$  latent. dim x samples (timepoints)
- $x$  :  $N \times T$  obs. dim (neurons) x samples (timepoints)
- $\mathbf{W}$  :  $N \times D$  obs. dim. (neurons) x latent dim.
- $\Psi$  :  $D \times D$  diagonal covariance matrix

# Factor Analysis: Generative Model

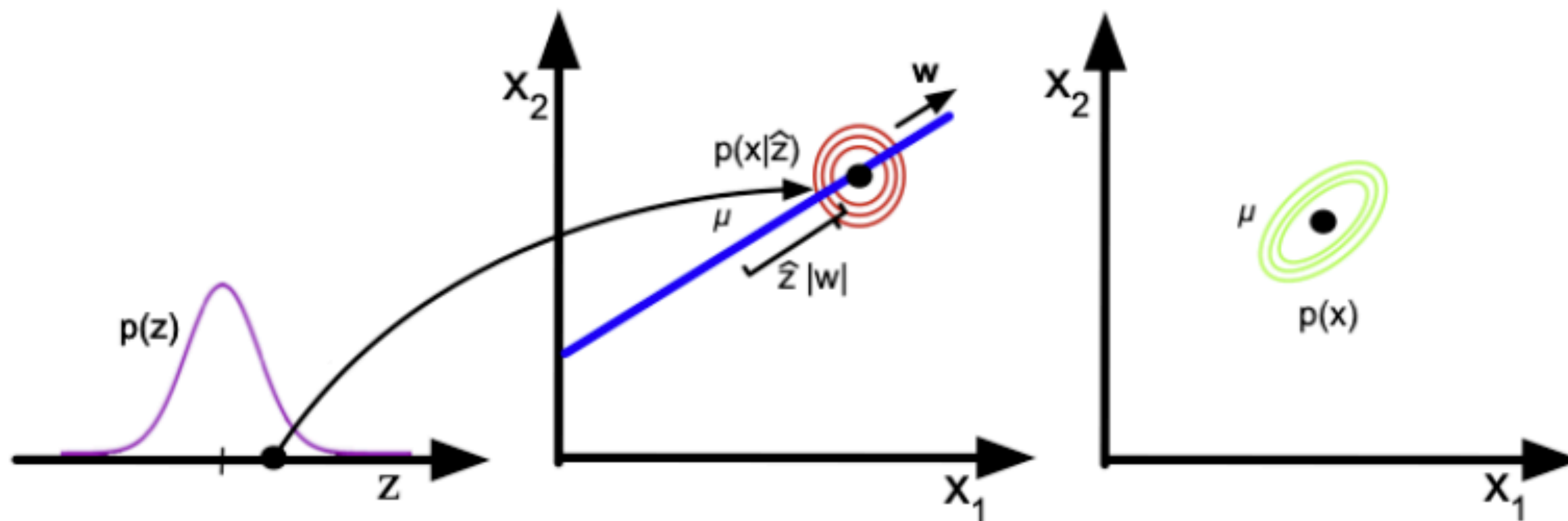
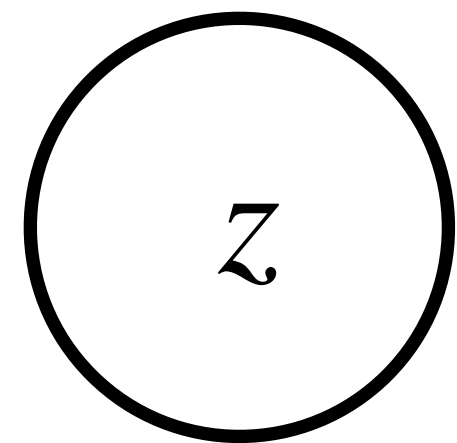
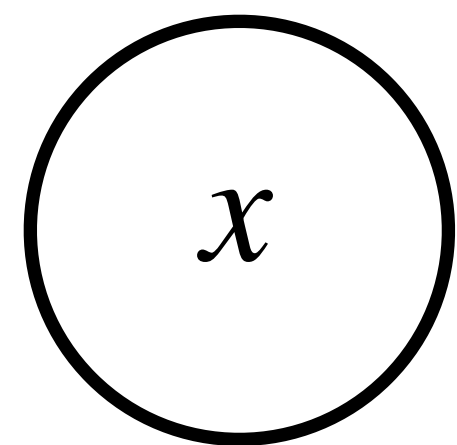
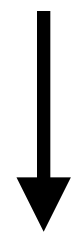


Figure 20.9: Illustration of the FA generative process, where we have  $L = 1$  latent dimension generating  $D = 2$  observed dimensions; we assume  $\Psi = \sigma^2 \mathbf{I}$ . The latent factor has value  $z \in \mathbb{R}$ , sampled from  $p(z)$ ; this gets mapped to a 2d offset  $\delta = zw$ , where  $w \in \mathbb{R}^2$ , which gets added to  $\mu$  to define a Gaussian  $p(x|z) = \mathcal{N}(x|\mu + \delta, \sigma^2 \mathbf{I})$ . By integrating over  $z$ , we “slide” this circular Gaussian “spray can” along the principal component axis  $w$ , which induces elliptical Gaussian contours in  $x$  space centered on  $\mu$ . Adapted from Figure 12.9 of [Bis06].

# Factor Analysis: Generative Model



$$p(z) = \mathcal{N}(z|\mathbf{0}, \mathbf{I})$$



$$p(x|z) = \mathcal{N}(x|\mathbf{W}z + \mu, \Psi)$$

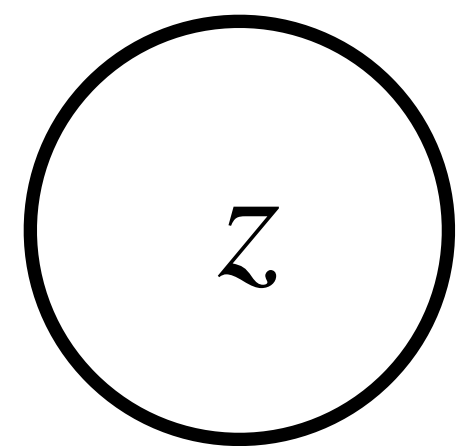
$$p(x) = \int p(x|z)p(z)dz$$

$$p(x) = \mathcal{N}(x|\mu, \mathbf{W}\mathbf{W}^T + \Psi)$$

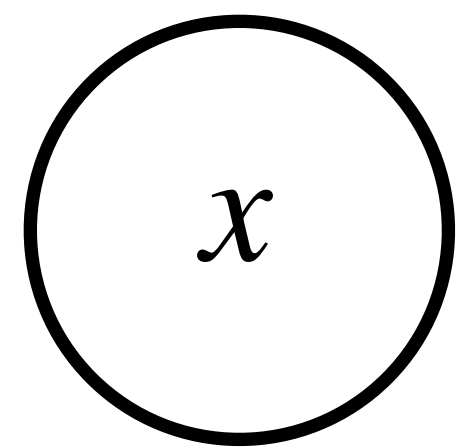
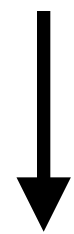
- Can be any Gaussian (see Murphy, book 1, section 20.2)
- Linear Gaussian model
- $z : D \times T$  latent. dim x samples (timepoints)
- $x : N \times T$  obs. dim (neurons) x samples (timepoints)
- $\mathbf{W} : N \times D$  obs. dim. (neurons) x latent dim.
- $\Psi : D \times D$  diagonal covariance matrix



# Factor Analysis: Generative Model



$$p(z) = \mathcal{N}(z|\mathbf{0}, \mathbf{I})$$



$$p(x|z) = \mathcal{N}(x|\mathbf{W}z + \mu, \Psi)$$

$$p(x) = \int p(x|z)p(z)dz$$

$$p(x) = \mathcal{N}(x|\mu, \mathbf{W}\mathbf{W}^T + \Psi)$$



Low Rank + Noise

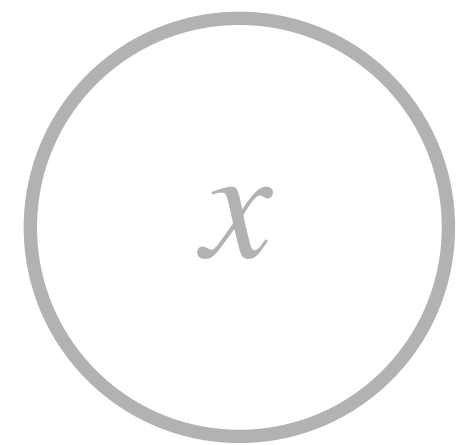
Shared and Unique

- Can be any Gaussian (see Murphy, book 1, section 20.2)
- Linear Gaussian model
- $z : D \times T$  latent. dim x samples (timepoints)
- $x : N \times T$  obs. dim (neurons) x samples (timepoints)
- $\mathbf{W} : N \times D$  obs. dim. (neurons) x latent dim.
- $\Psi : D \times D$  diagonal covariance matrix

# FA vs. Probabilistic PCA vs. PCA

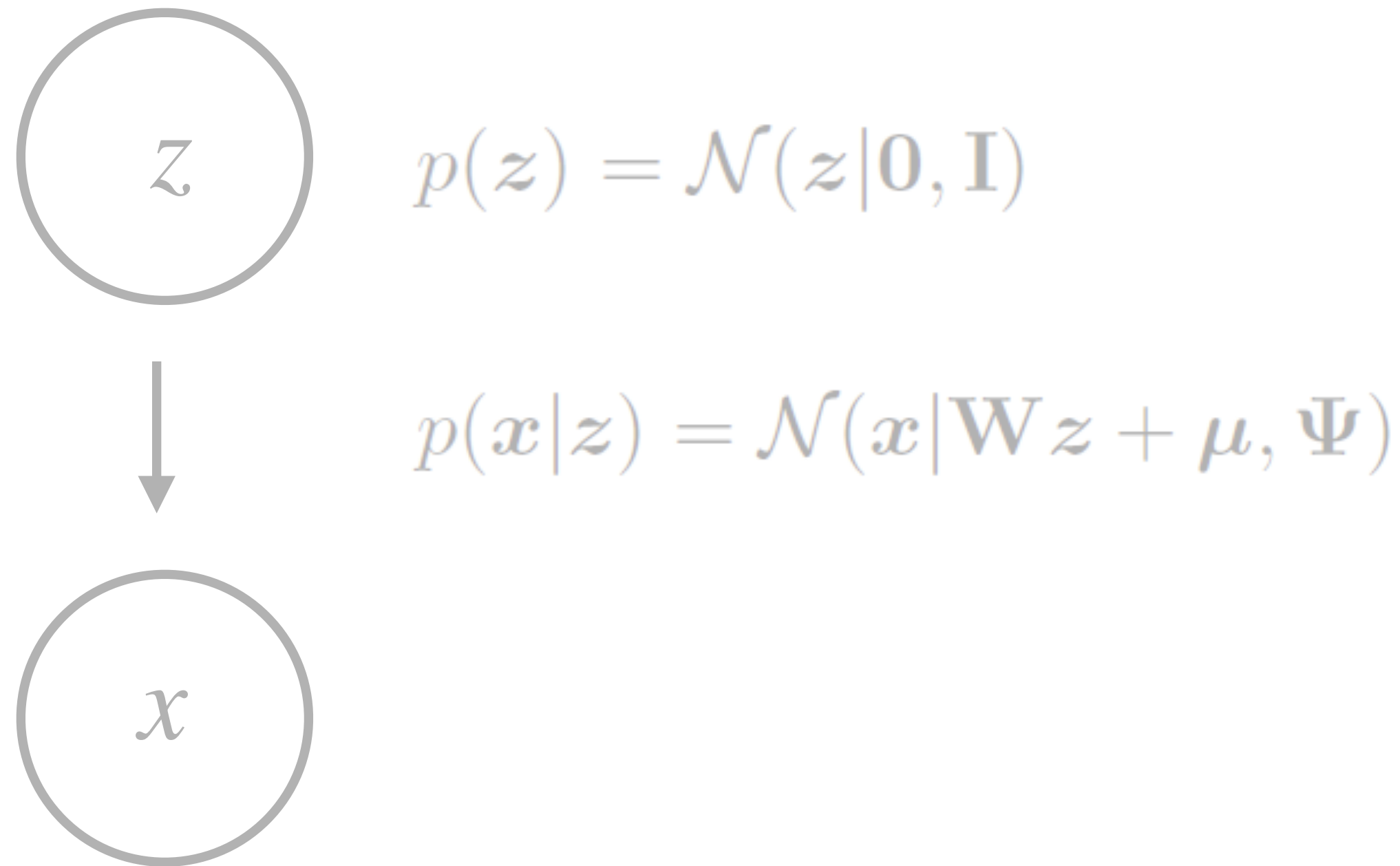


$$p(z) = \mathcal{N}(z|\mathbf{0}, \mathbf{I})$$



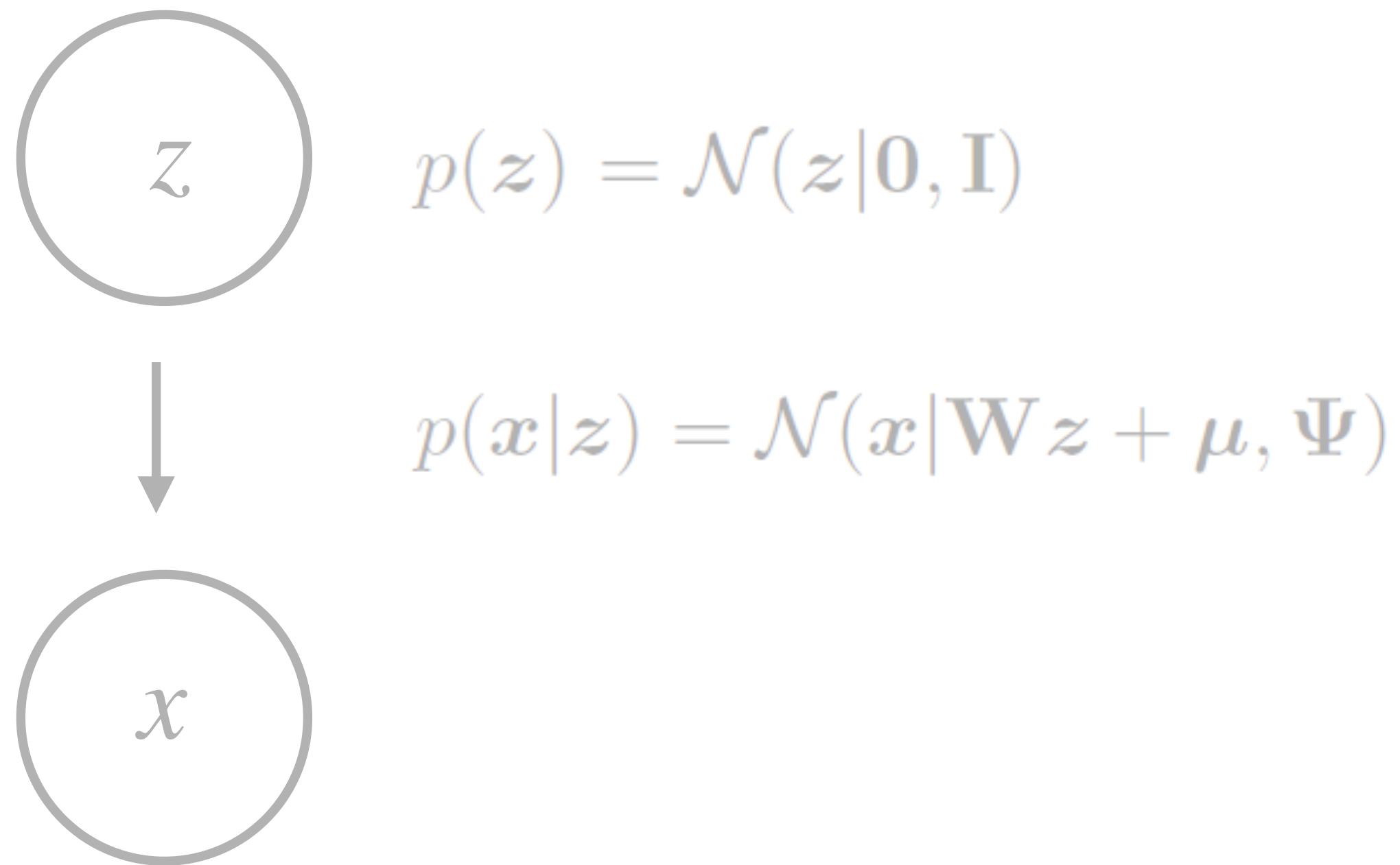
$$p(x|z) = \mathcal{N}(x|\mathbf{W}z + \mu, \Psi)$$

# FA vs. Probabilistic PCA vs. PCA



- Probabilistic PCA is Factor Analysis where  $\Psi$  is the identity matrix (all observations have the same independent noise)

# FA vs. Probabilistic PCA vs. PCA



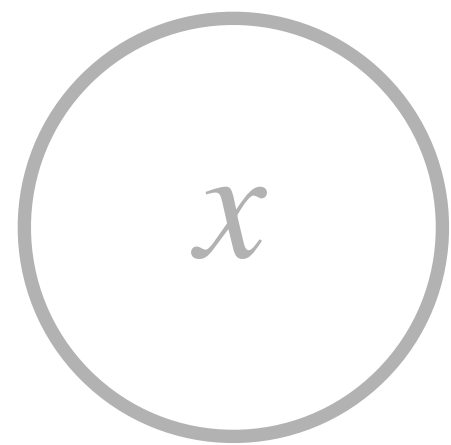
- Probabilistic PCA is Factor Analysis where  $\Psi$  is the identity matrix (all observations have the same independent noise)
- PPCA when  $\Psi \rightarrow 0$  becomes PCA



# FA vs. PCA: An Example



$$p(z) = \mathcal{N}(z|\mathbf{0}, \mathbf{I})$$

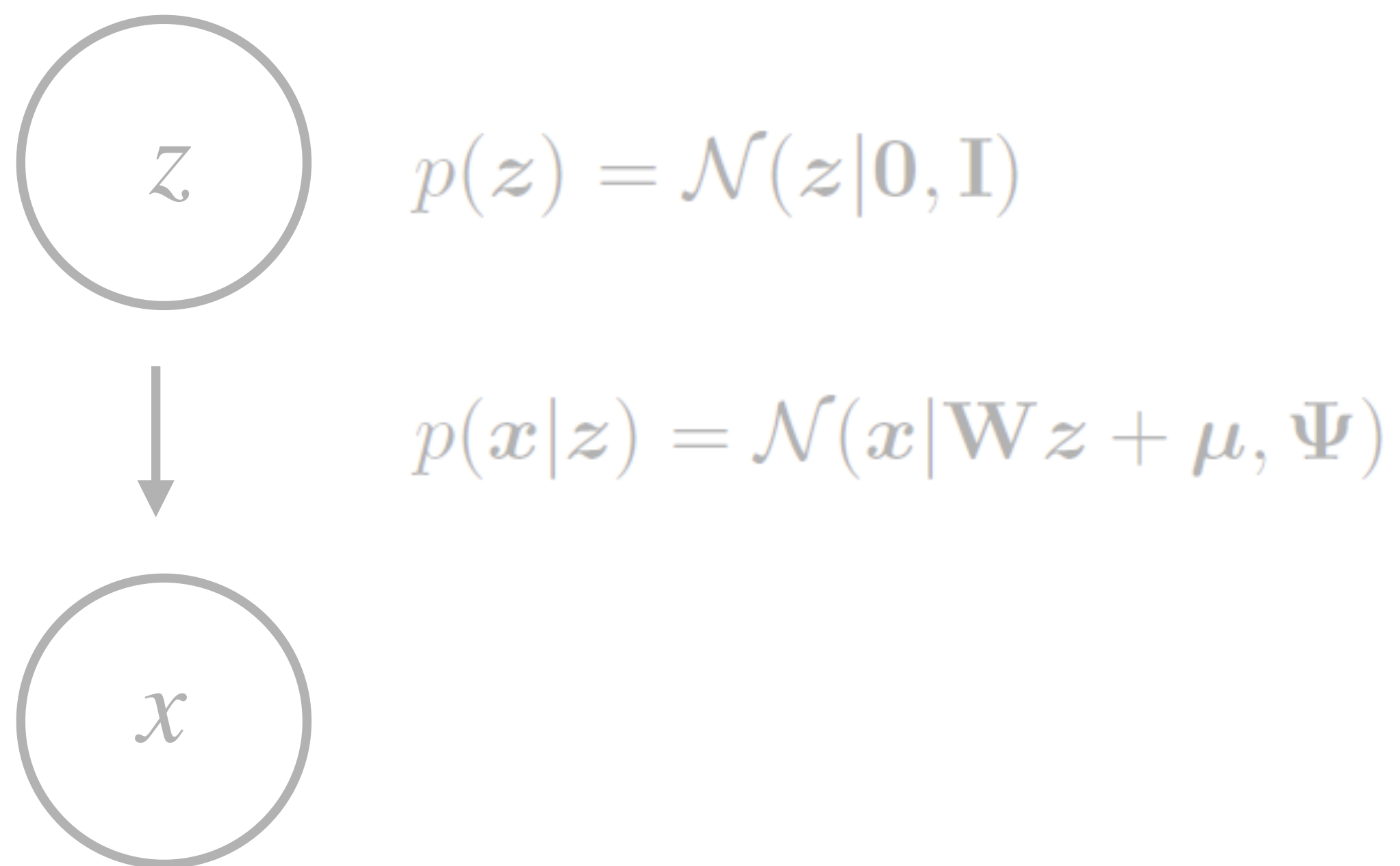


$$p(x|z) = \mathcal{N}(x|\mathbf{W}z + \mu, \Psi)$$

$$\mathbf{W} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\Psi = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

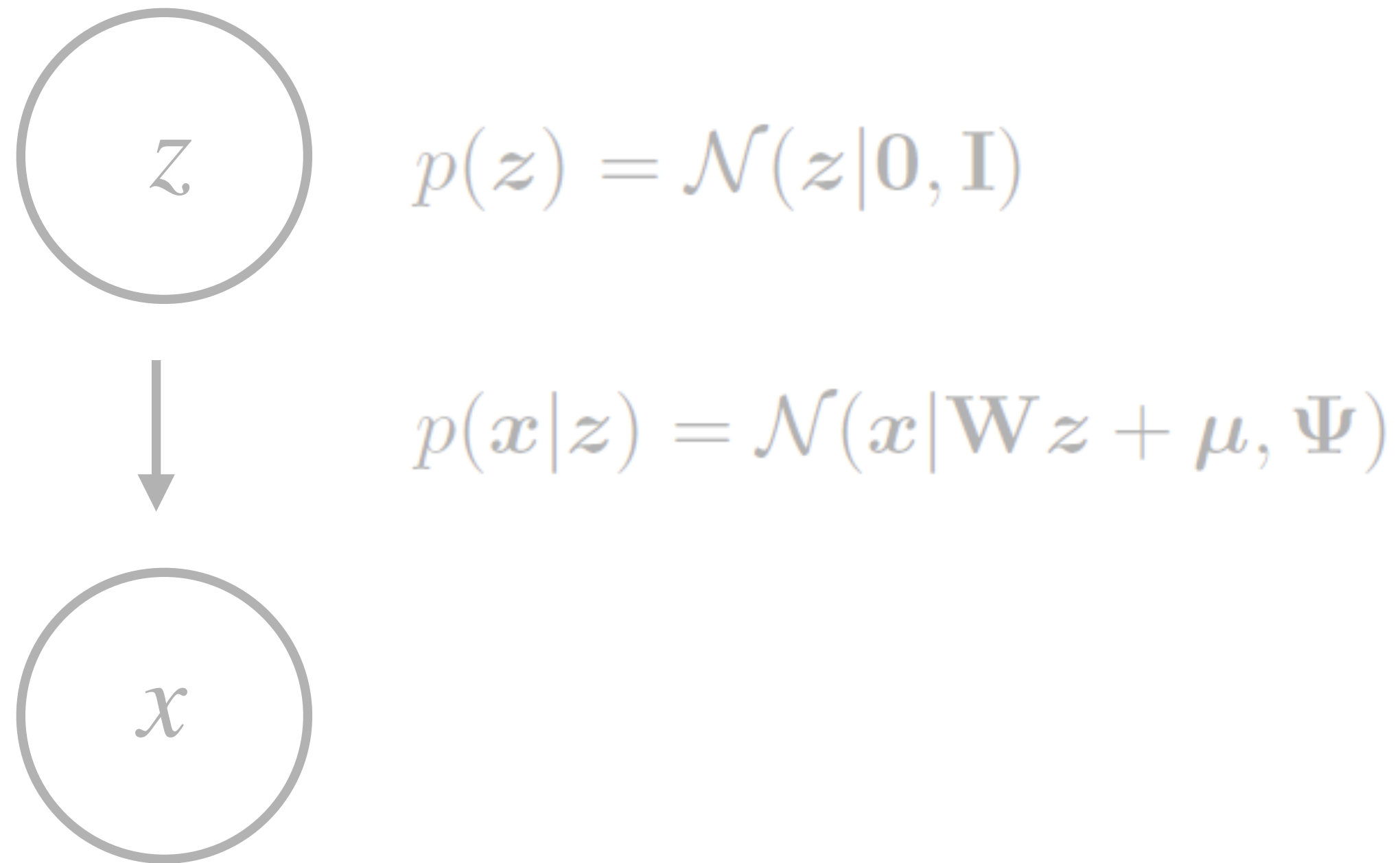
# FA vs. PCA: An Example



$$W = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \Psi = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\text{cov}(X) = WW^T + \Psi = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$$

# FA vs. PCA: An Example



$$W = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \Psi = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\text{cov}(X) = WW^T + \Psi = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$$

- PCA would give a top eigenvector primarily lying along the first dimension

# Factor Analysis: Inferring the latents

$$\begin{aligned} p(z | x) &\propto p(x | z)p(z) \\ &= \mathcal{N}(x | Wz, \Psi) \cdot \mathcal{N}(z | 0, I) \end{aligned}$$

# Factor Analysis: Inferring the latents

$$p(z|x) \propto p(x|z)p(z)$$

$$= \mathcal{N}(x|Wz, \Psi) \cdot \mathcal{N}(z|0, I)$$

$$\vdots$$

$$= \mathcal{N}(\Lambda W^T \Psi^{-1} x, \Lambda) \quad \text{where} \quad \Lambda = (W^T \Psi^{-1} W + I)^{-1}$$

# Factor Analysis: Inferring the latents

$$p(z | x) \propto p(x | z)p(z)$$

$$= \mathcal{N}(x | Wz, \Psi) \cdot \mathcal{N}(z | 0, I)$$

$$\vdots$$

$$= \mathcal{N}(\Lambda W^T \Psi^{-1} x, \Lambda) \quad \text{where} \quad \Lambda = (W^T \Psi^{-1} W + I)^{-1}$$

- When inferring the latent, the components of  $x$  are downweighted in proportion to their amount of independent noise (value in  $\Psi$ ).

# EM for Factor Analysis

- E step: Estimate the posterior,  $p(z|x)$  , given set parameters
- M step: Estimate the parameters,  $[W, \Psi]$ , given the expectations of the latents



# EM for ~~Factor Analysis~~ PPCA

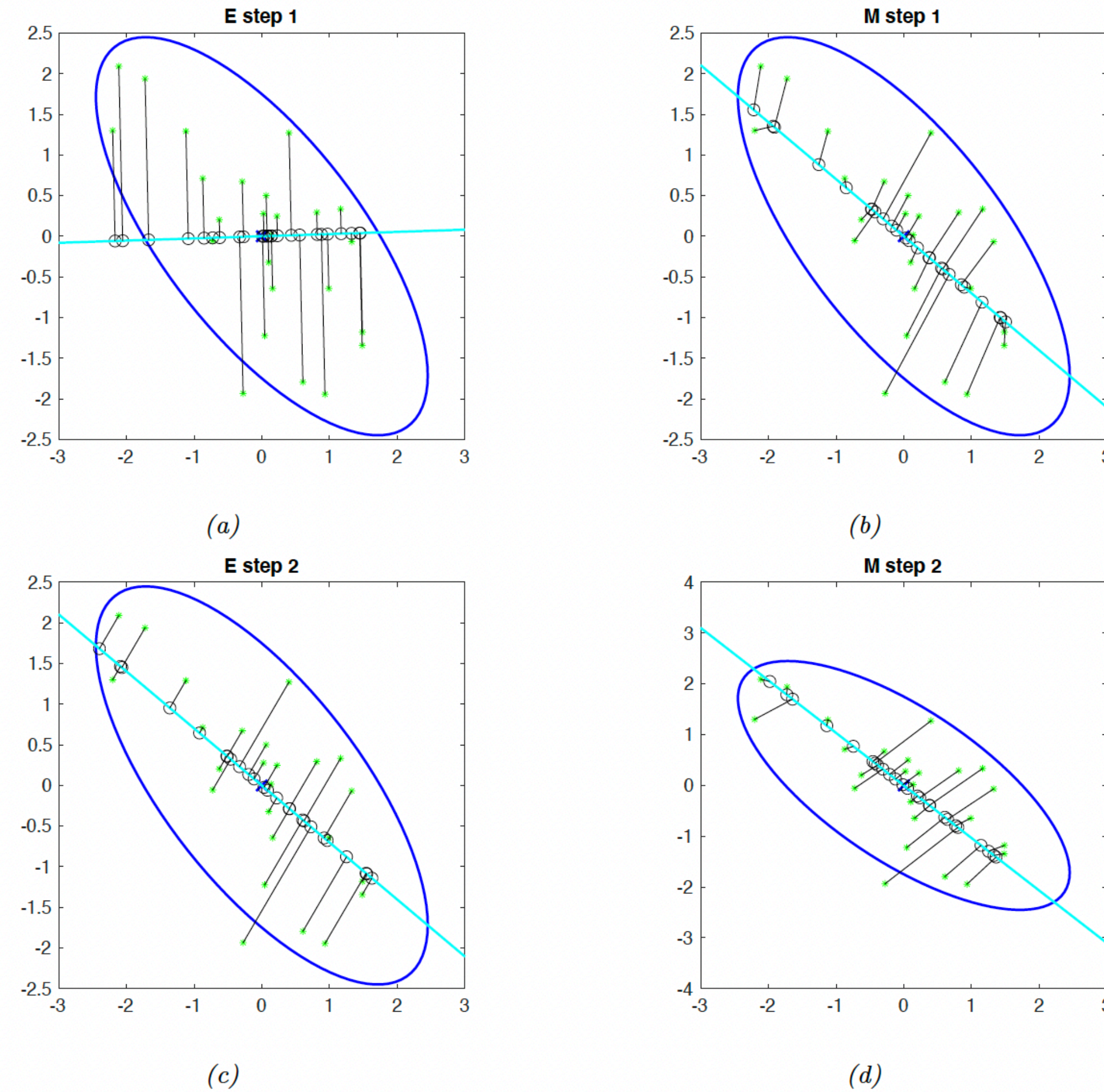


Figure 20.10: Illustration of EM for PCA when  $D = 2$  and  $L = 1$ . Green stars are the original data points, black circles are their reconstructions. The weight vector  $w$  is represented by blue line. (a) We start with a random initial guess of  $w$ . The E step is represented by the orthogonal projections. (b) We update the rod  $w$  in the M step, keeping the projections onto the rod (black circles) fixed. (c) Another E step. The black circles can 'slide' along the rod, but the rod stays fixed. (d) Another M step. Adapted from Figure 12.12 of [Bis06].



# Why probabilistic models, versus PCA?

- Allows having more sophisticated, and more accurate models
  - Different noise models (FA vs PPCA), mixture of factor analyzers, etc...
- Principled
- Better for missing data, or streaming data

# Gaussian Processes

# Gaussian Processes

Now consider a function  $f : \mathcal{X} \rightarrow \mathbb{R}$  evaluated at a set of inputs,  $\mathbf{X} = \{x_n \in \mathcal{X}\}_{n=1}^N$ . Let  $\mathbf{f}_X = [f(x_1), \dots, f(x_N)]$  be the set of unknown function values at these points.

# Gaussian Processes

Now consider a function  $f : \mathcal{X} \rightarrow \mathbb{R}$

evaluated at a set of inputs,  $\mathbf{X} = \{x_n \in \mathcal{X}\}_{n=1}^N$ . Let  $\mathbf{f}_X = [f(x_1), \dots, f(x_N)]$  be the set of unknown function values at these points. If  $\mathbf{f}_X$  is jointly Gaussian for any set of  $N \geq 1$  points, then we say that  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a **Gaussian process**.

# Gaussian Processes

Now consider a function  $f : \mathcal{X} \rightarrow \mathbb{R}$  evaluated at a set of inputs,  $\mathbf{X} = \{x_n \in \mathcal{X}\}_{n=1}^N$ . Let  $\mathbf{f}_X = [f(x_1), \dots, f(x_N)]$  be the set of unknown function values at these points. If  $\mathbf{f}_X$  is jointly Gaussian for any set of  $N \geq 1$  points, then we say that  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a **Gaussian process**. Such a process is defined by its **mean function**  $m(x) \in \mathbb{R}$  and a **covariance function**,  $\mathcal{K}(x, x') \geq 0$ , which is any positive definite **Mercer kernel**

# Gaussian Processes

Now consider a function  $f : \mathcal{X} \rightarrow \mathbb{R}$

evaluated at a set of inputs,  $\mathbf{X} = \{\mathbf{x}_n \in \mathcal{X}\}_{n=1}^N$ . Let  $\mathbf{f}_X = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]$  be the set of unknown function values at these points. If  $\mathbf{f}_X$  is jointly Gaussian for any set of  $N \geq 1$  points, then we say that  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a **Gaussian process**. Such a process is defined by its **mean function**  $m(\mathbf{x}) \in \mathbb{R}$  and a **covariance function**,  $\mathcal{K}(\mathbf{x}, \mathbf{x}') \geq 0$ , which is any positive definite **Mercer kernel**

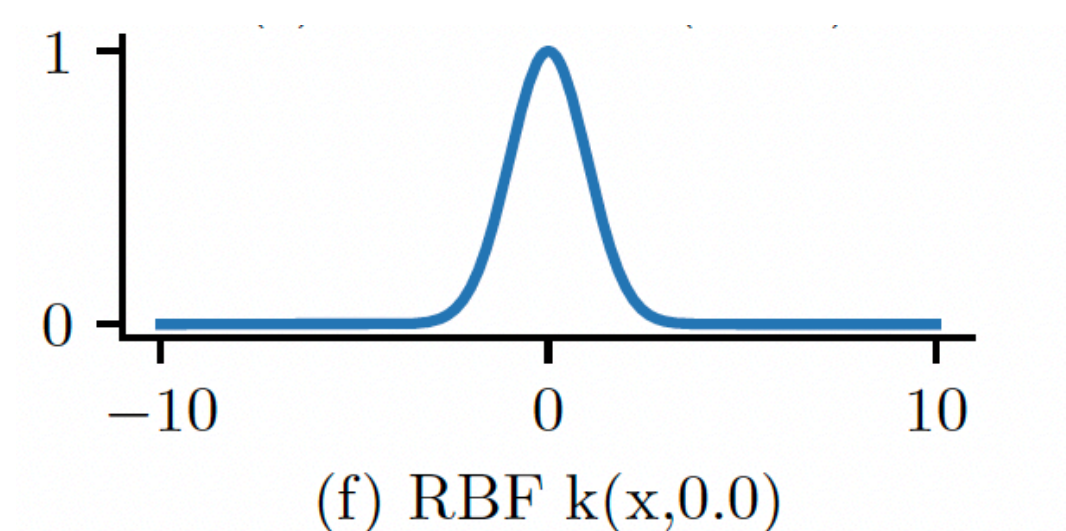
- Example Kernel (“Radial Basis Function”):  $\mathcal{K}(\mathbf{x}, \mathbf{x}'; \ell) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right)$



# Gaussian Processes

Now consider a function  $f : \mathcal{X} \rightarrow \mathbb{R}$  evaluated at a set of inputs,  $\mathbf{X} = \{\mathbf{x}_n \in \mathcal{X}\}_{n=1}^N$ . Let  $\mathbf{f}_X = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]$  be the set of unknown function values at these points. If  $\mathbf{f}_X$  is jointly Gaussian for any set of  $N \geq 1$  points, then we say that  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a **Gaussian process**. Such a process is defined by its **mean function**  $m(\mathbf{x}) \in \mathbb{R}$  and a **covariance function**,  $\mathcal{K}(\mathbf{x}, \mathbf{x}') \geq 0$ , which is any positive definite **Mercer kernel**

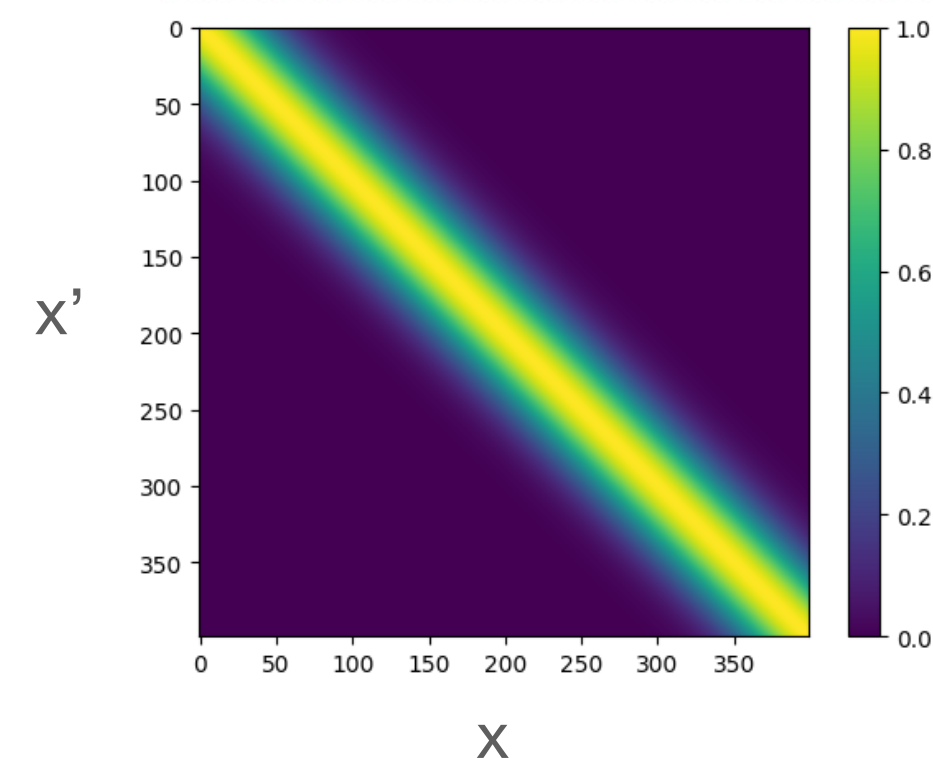
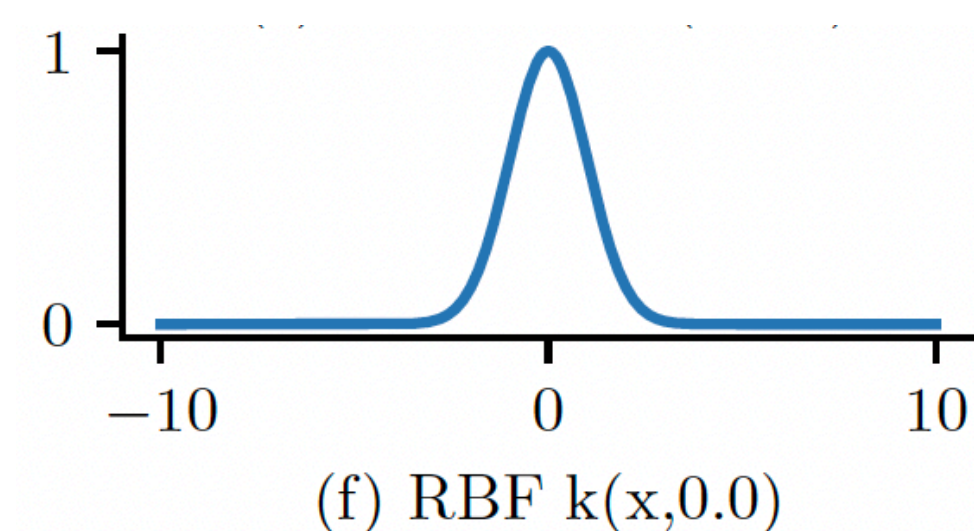
- Example Kernel (“Radial Basis Function”):  $\mathcal{K}(\mathbf{x}, \mathbf{x}'; \ell) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\ell^2}\right)$



# Gaussian Processes

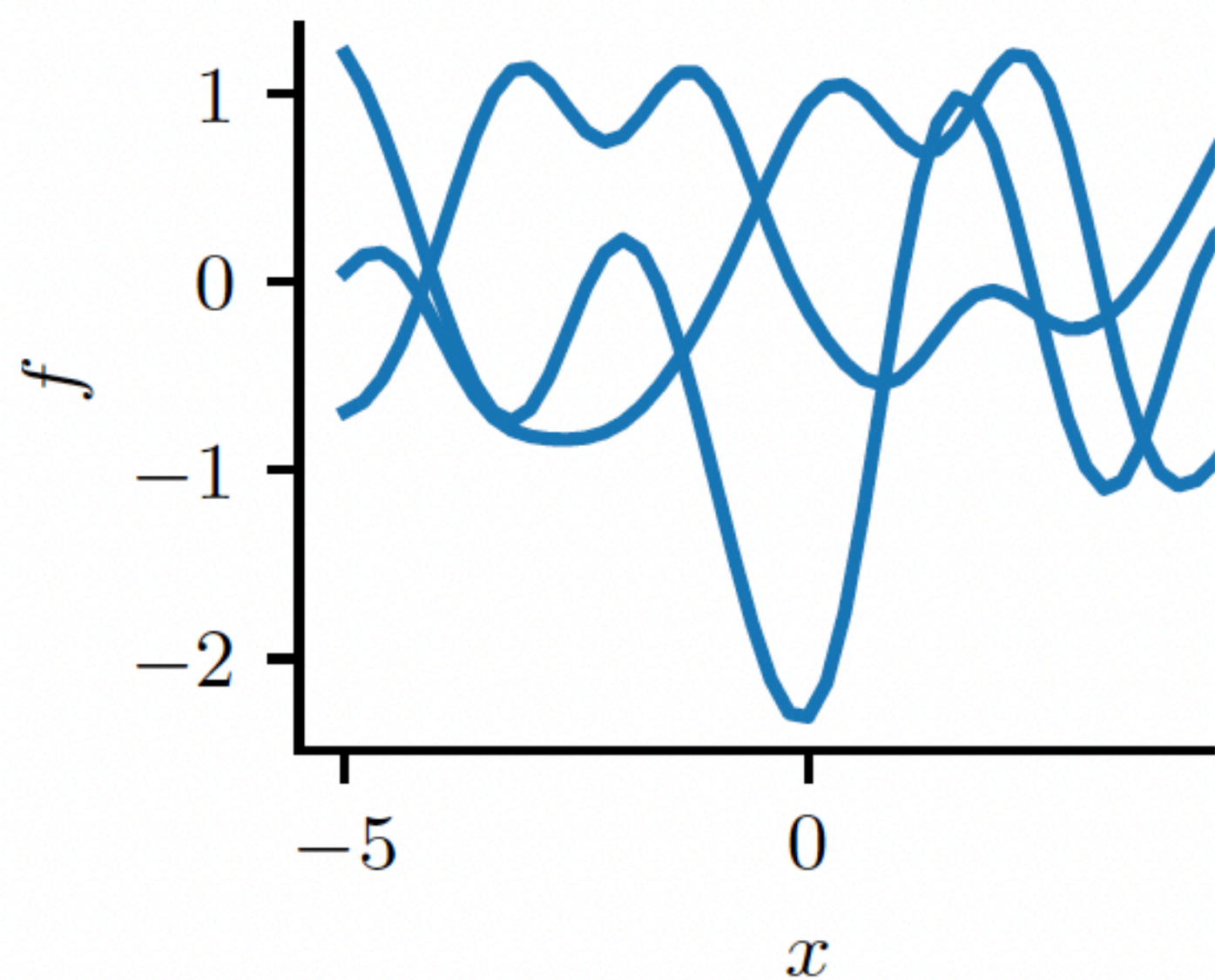
Now consider a function  $f : \mathcal{X} \rightarrow \mathbb{R}$  evaluated at a set of inputs,  $\mathbf{X} = \{x_n \in \mathcal{X}\}_{n=1}^N$ . Let  $\mathbf{f}_X = [f(x_1), \dots, f(x_N)]$  be the set of unknown function values at these points. If  $\mathbf{f}_X$  is jointly Gaussian for any set of  $N \geq 1$  points, then we say that  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a **Gaussian process**. Such a process is defined by its **mean function**  $m(x) \in \mathbb{R}$  and a **covariance function**,  $\mathcal{K}(x, x') \geq 0$ , which is any positive definite **Mercer kernel**

- Example Kernel (“Radial Basis Function”):  $\mathcal{K}(x, x'; \ell) = \exp\left(-\frac{\|x - x'\|^2}{2\ell^2}\right)$





# Gaussian Processes - sampling from the prior



(a)

Figure 18.7: Left: some functions sampled from a GP prior with RBF kernel. Middle: some samples from a GP posterior, after conditioning on 5 noise-free observations. Right: some samples from a GP posterior, after conditioning on 5 noisy observations. The shaded area represents  $\mathbb{E}[f(\mathbf{x})] \pm 2\sqrt{\mathbb{V}[f(\mathbf{x})]}$ . Adapted from Figure 2.2 of [RW06]. Generated by [gpr\\_demo\\_noise\\_free.ipynb](#).



# Gaussian Processes - Example kernels

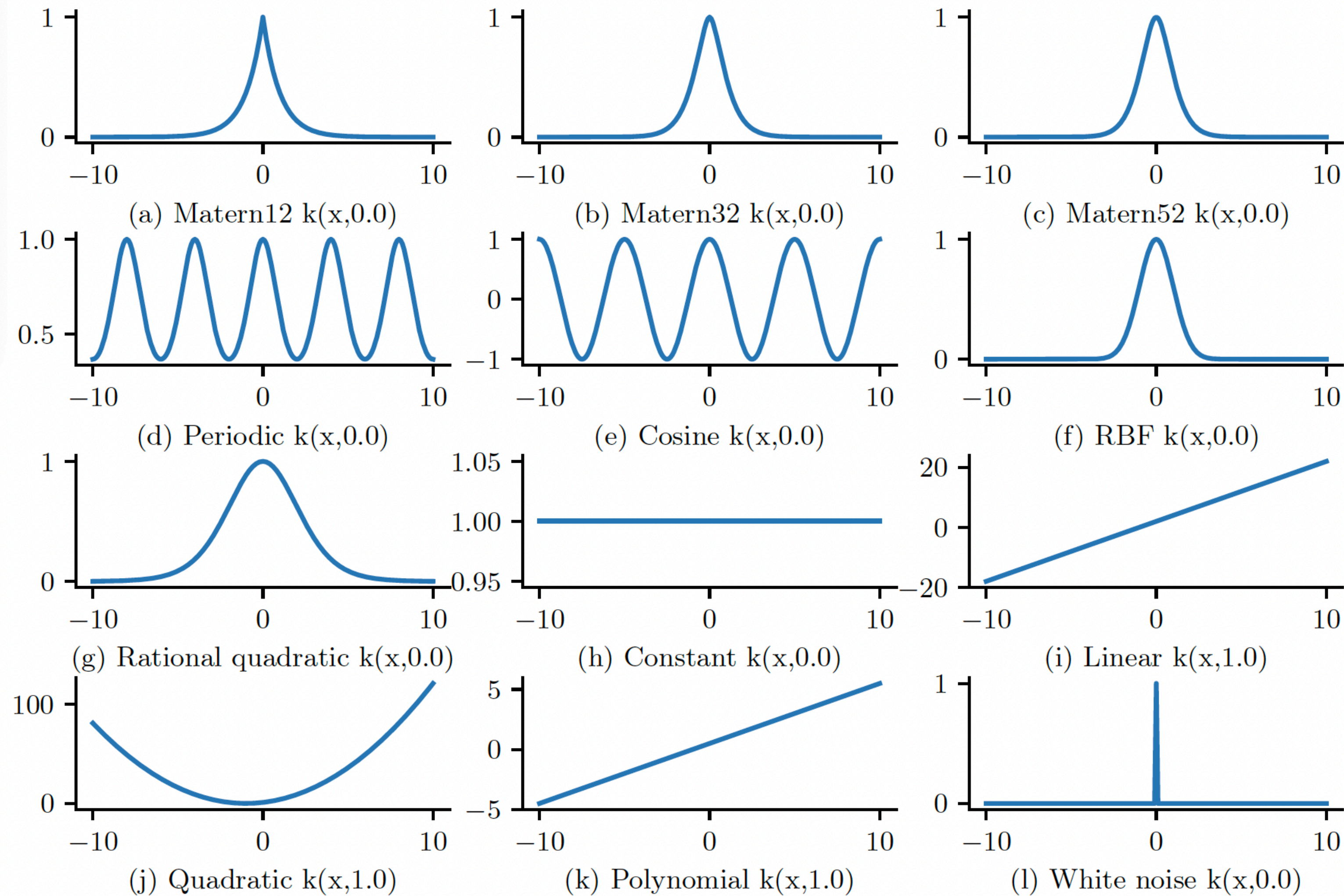
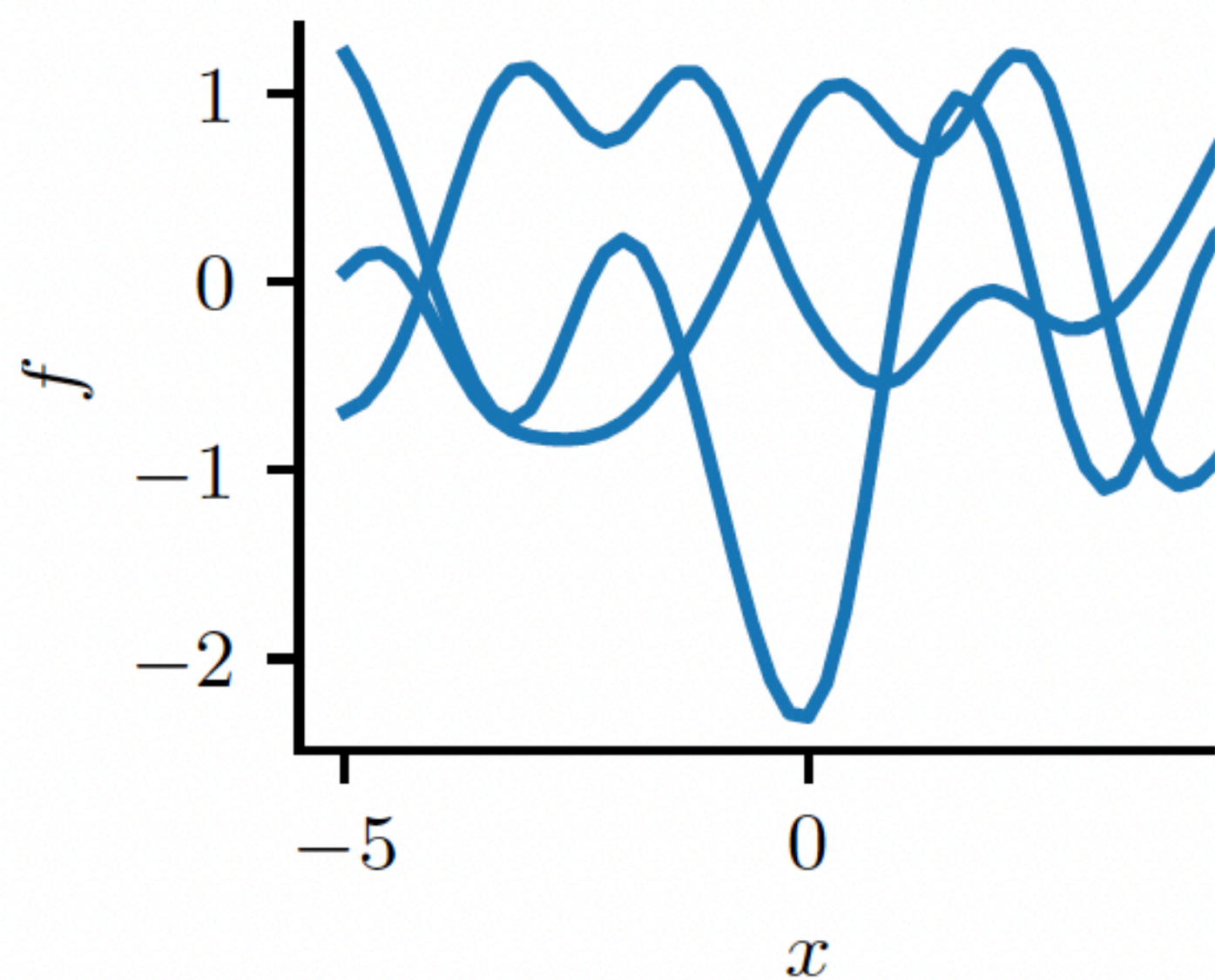


Figure 18.3: GP kernels evaluated at  $k(x, 0)$  as a function of  $x$ . Generated by [gpKernelPlot.ipynb](#).



# Gaussian Processes - estimating a posterior



(a)

Figure 18.7: Left: some functions sampled from a GP prior with RBF kernel. Middle: some samples from a GP posterior, after conditioning on 5 noise-free observations. Right: some samples from a GP posterior, after conditioning on 5 noisy observations. The shaded area represents  $\mathbb{E}[f(\mathbf{x})] \pm 2\sqrt{\mathbb{V}[f(\mathbf{x})]}$ . Adapted from Figure 2.2 of [RW06]. Generated by [gpr\\_demo\\_noise\\_free.ipynb](#).



# Gaussian Processes - estimating a posterior

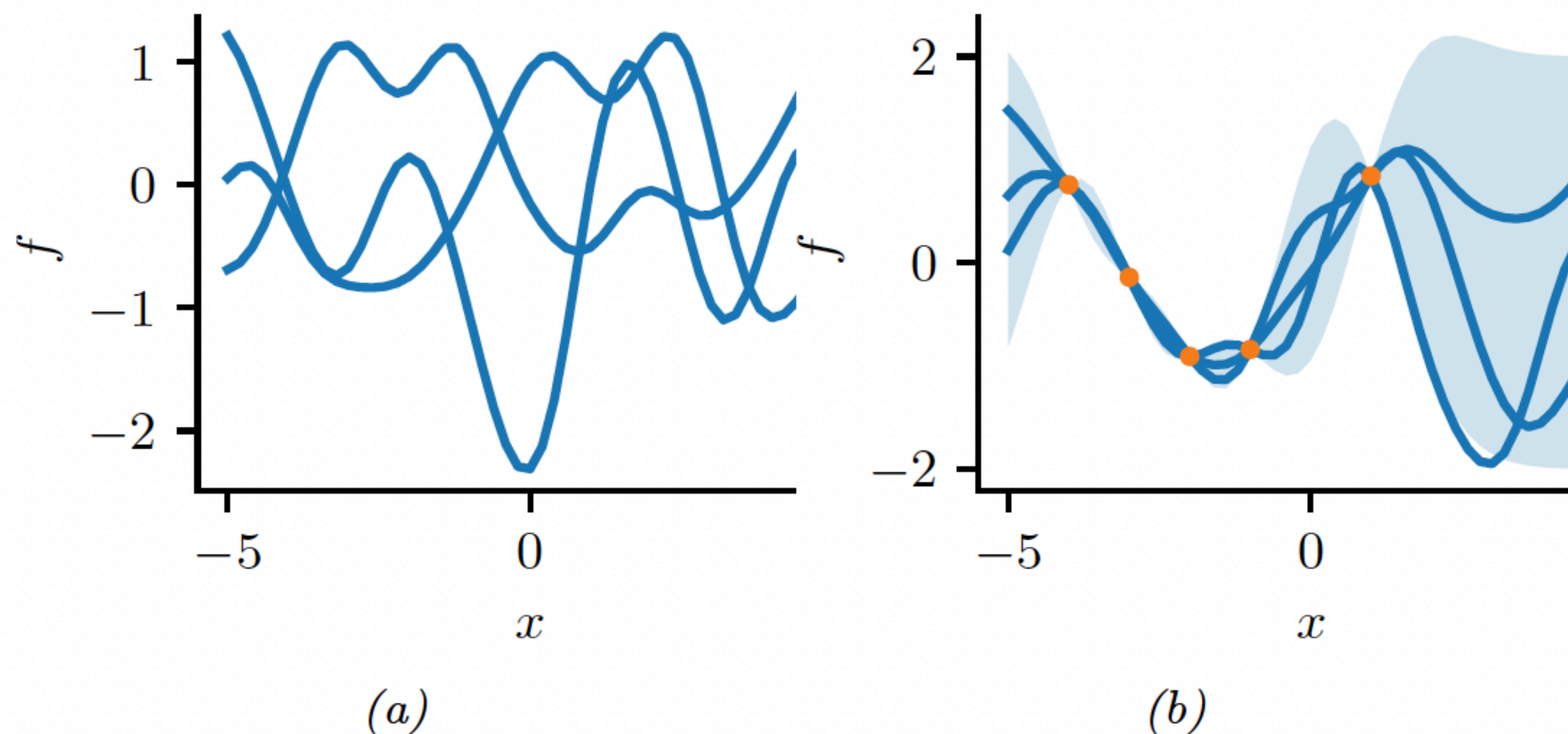


Figure 18.7: Left: some functions sampled from a GP prior with RBF kernel. Middle: some samples from a GP posterior, after conditioning on 5 noise-free observations. Right: some samples from a GP posterior, after conditioning on 5 noisy observations. The shaded area represents  $\mathbb{E}[f(\mathbf{x})] \pm 2\sqrt{\mathbb{V}[f(\mathbf{x})]}$ . Adapted from Figure 2.2 of [RW06]. Generated by [gpr\\_demo\\_noise\\_free.ipynb](#).



# Gaussian Processes - estimating a posterior

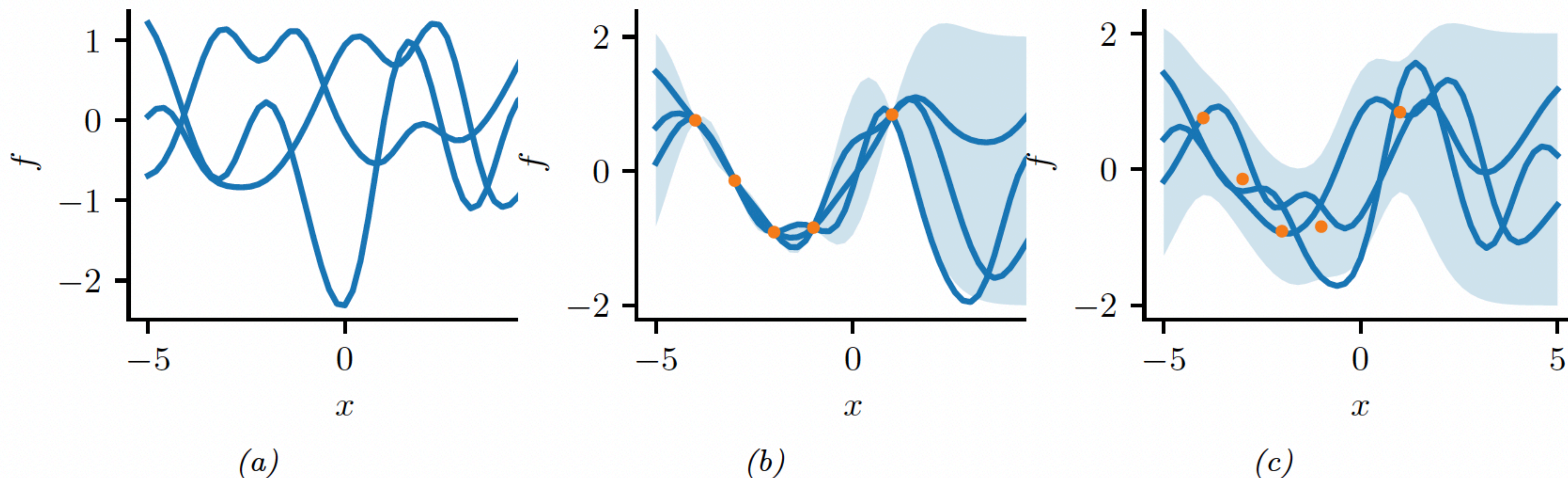


Figure 18.7: Left: some functions sampled from a GP prior with RBF kernel. Middle: some samples from a GP posterior, after conditioning on 5 noise-free observations. Right: some samples from a GP posterior, after conditioning on 5 noisy observations. The shaded area represents  $\mathbb{E}[f(\mathbf{x})] \pm 2\sqrt{\mathbb{V}[f(\mathbf{x})]}$ . Adapted from Figure 2.2 of [RW06]. Generated by [gpr\\_demo\\_noise\\_free.ipynb](#).