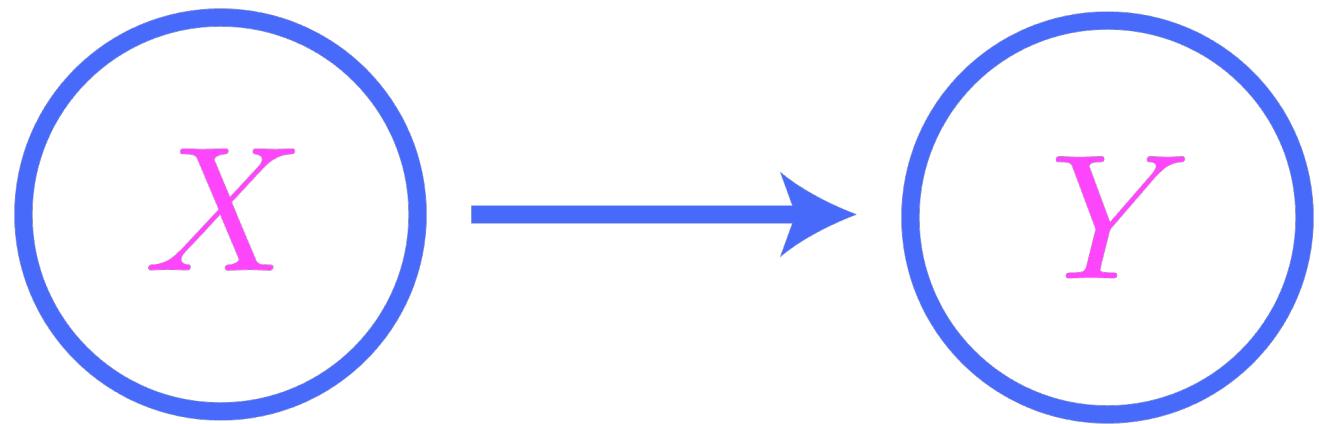
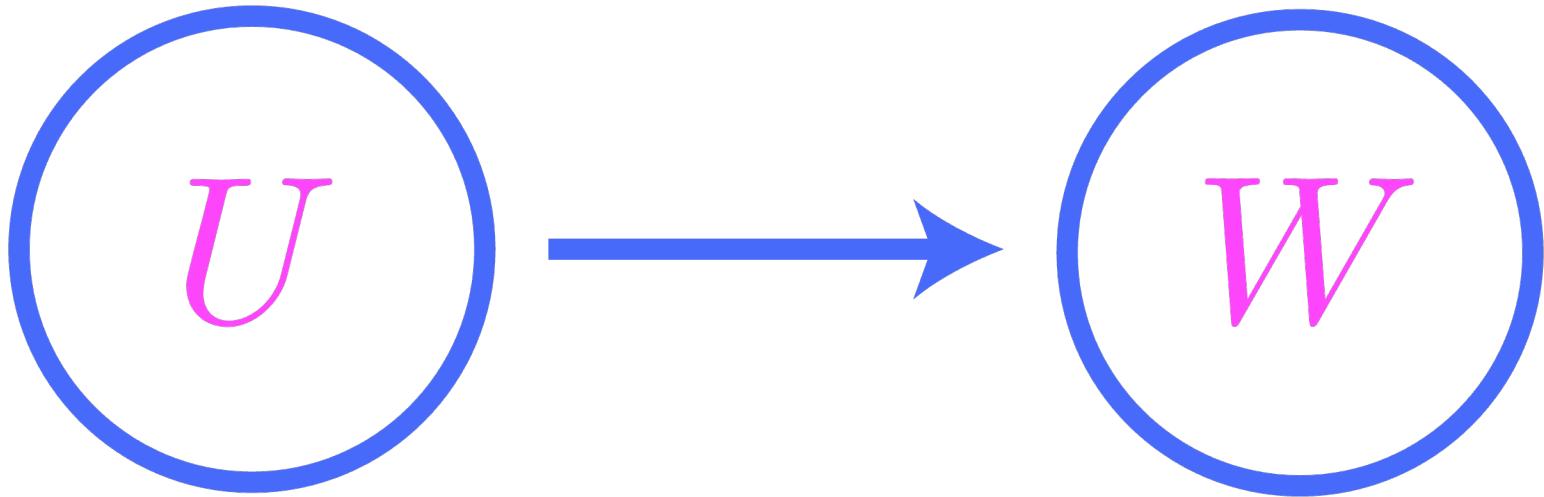
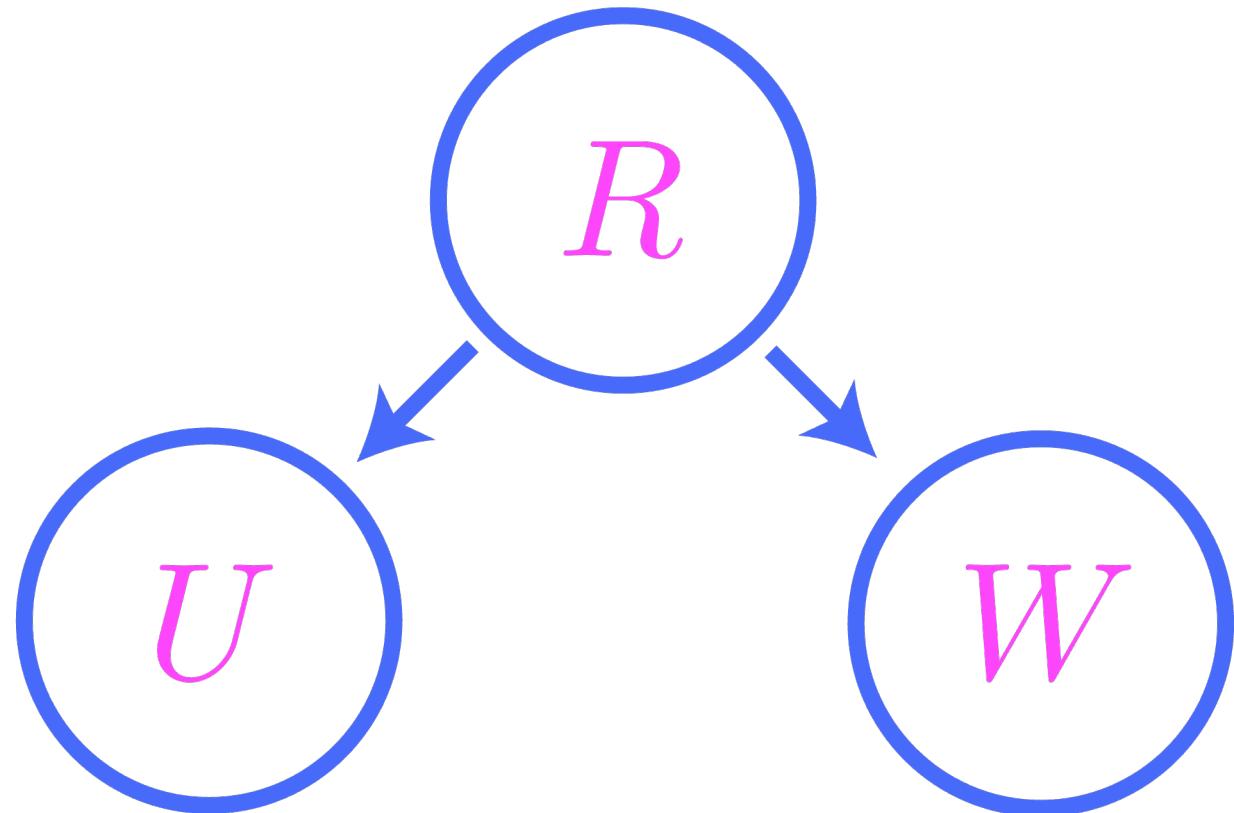


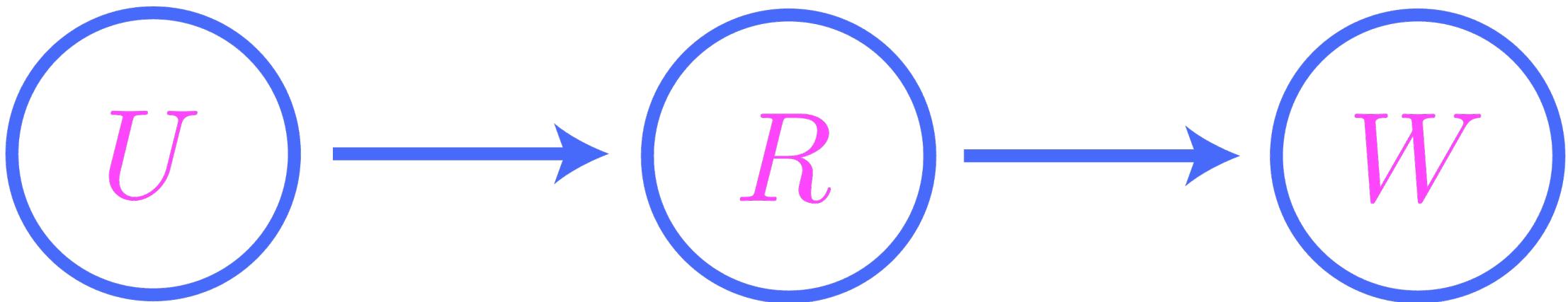
# 12. Probabilistic graphical models

## 12.1. A brief look at causal inference









## **Causal structures and probability.**

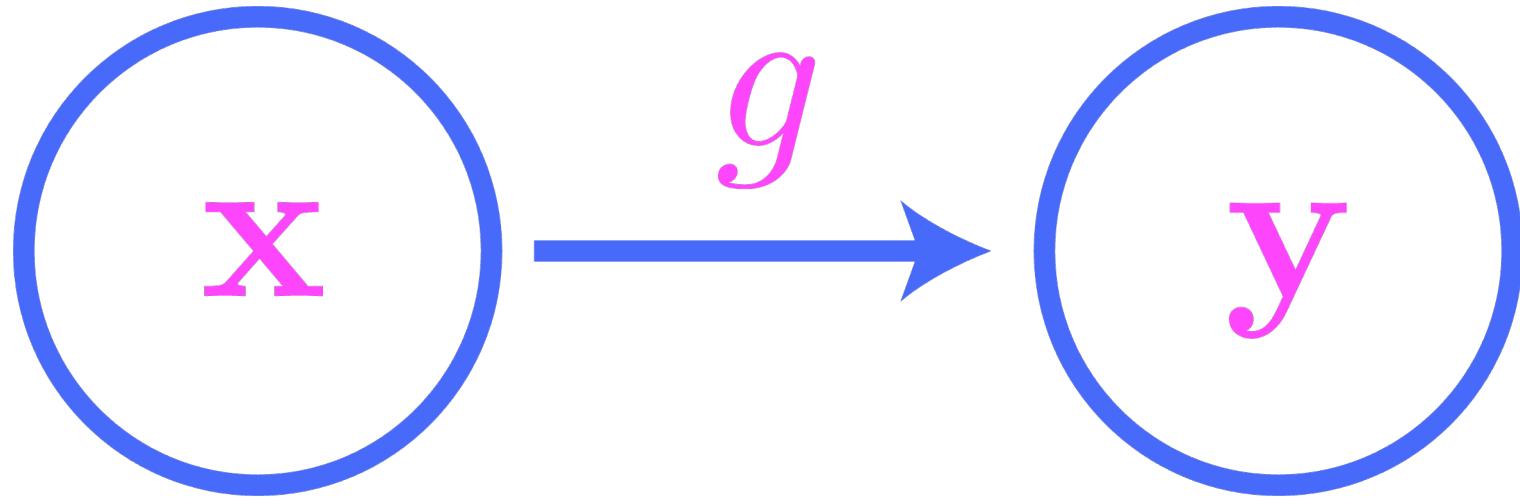
- Relationships of cause and effect represent strictly more structure than a joint probability distribution.
- A causal structure *refines* a joint probability distribution; it encodes *more* knowledge.
- The mapping from causal structures to joint probability distributions is many-to-one.



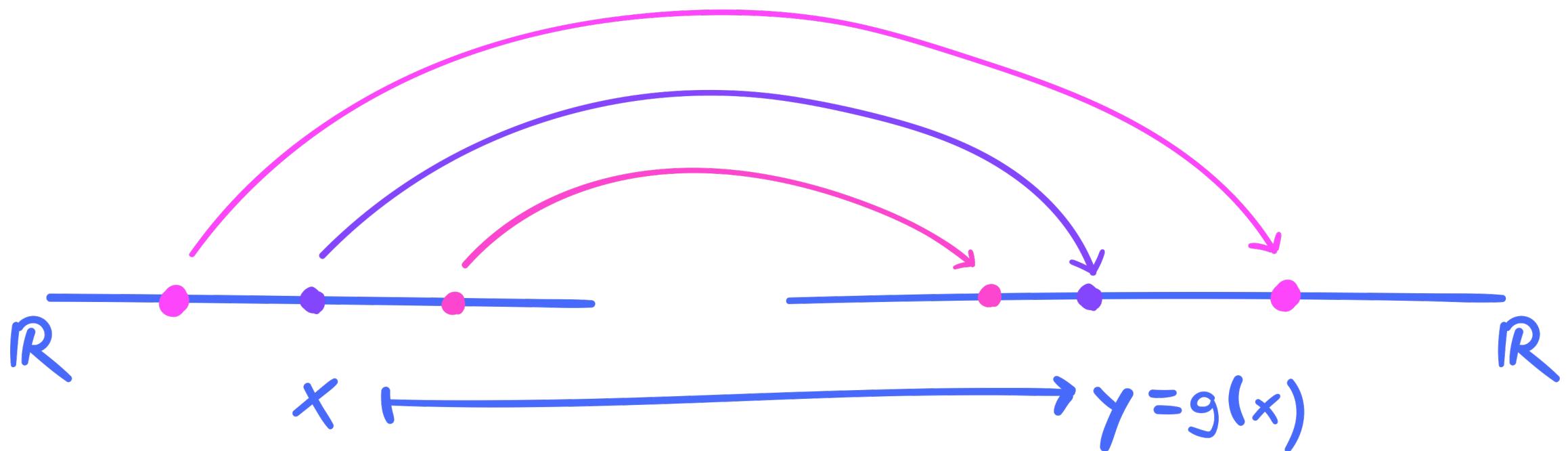
### Problem Prompt

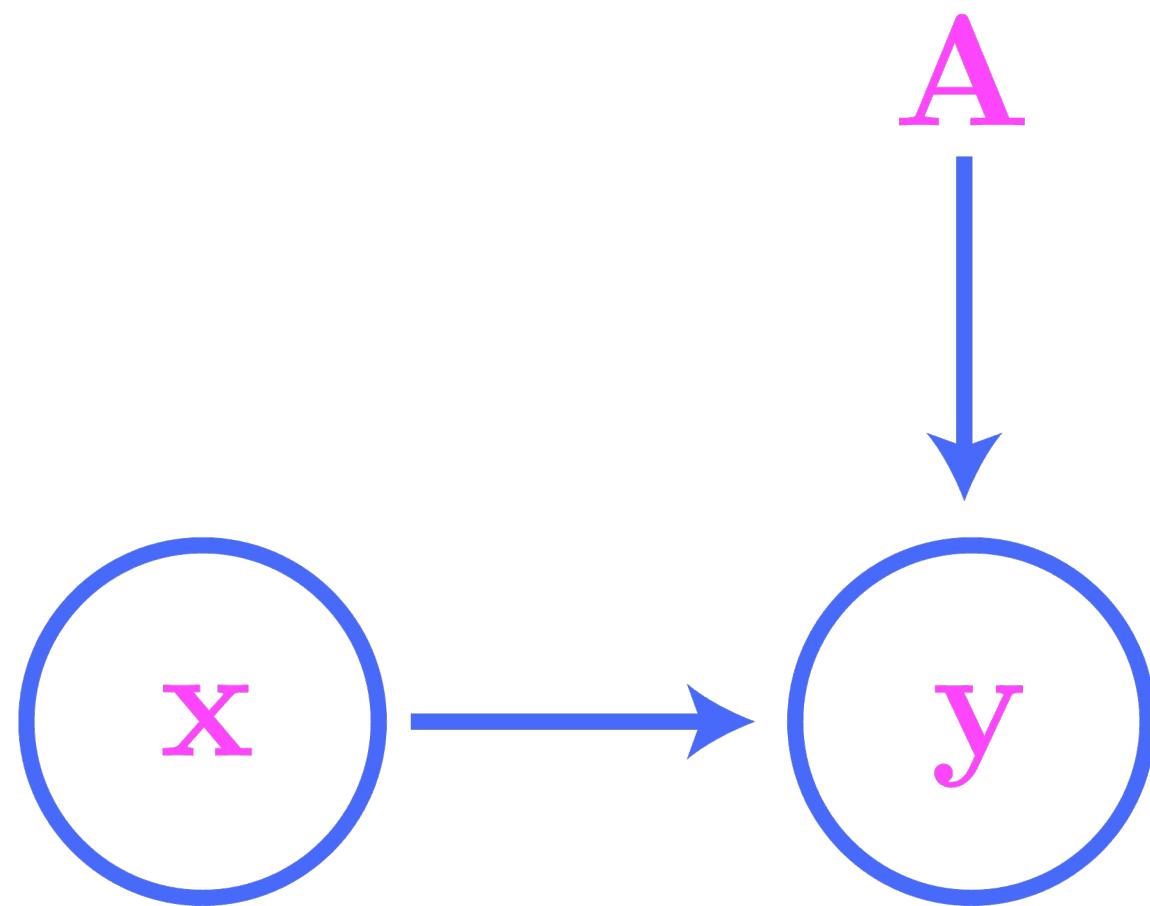
Do problem 1 on the worksheet.

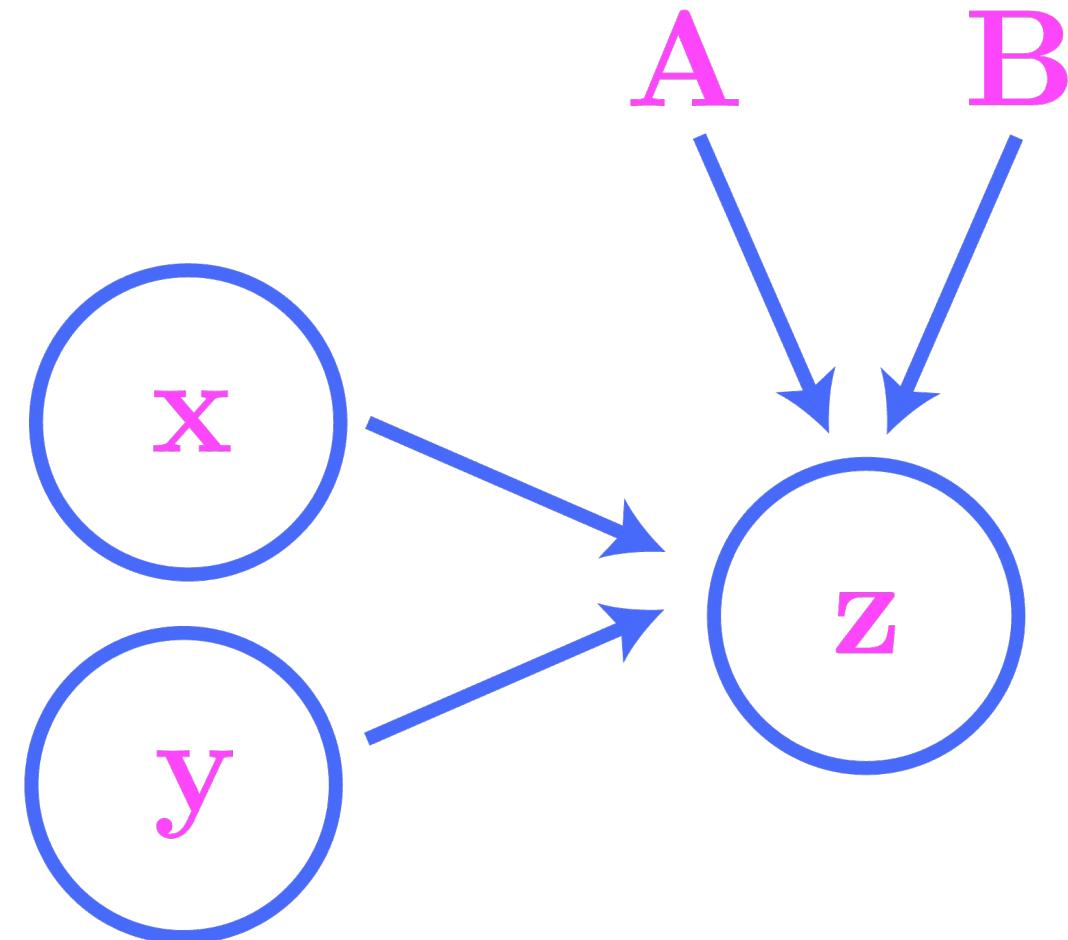
## 12.2. General probabilistic graphical models

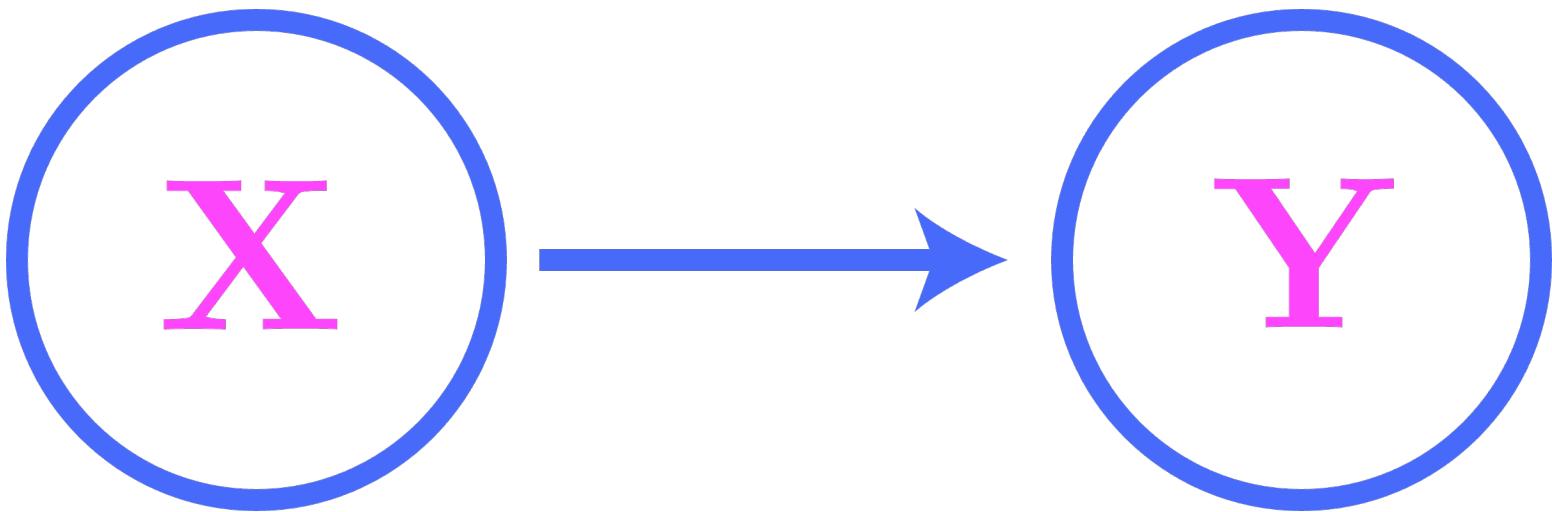


Deterministic link

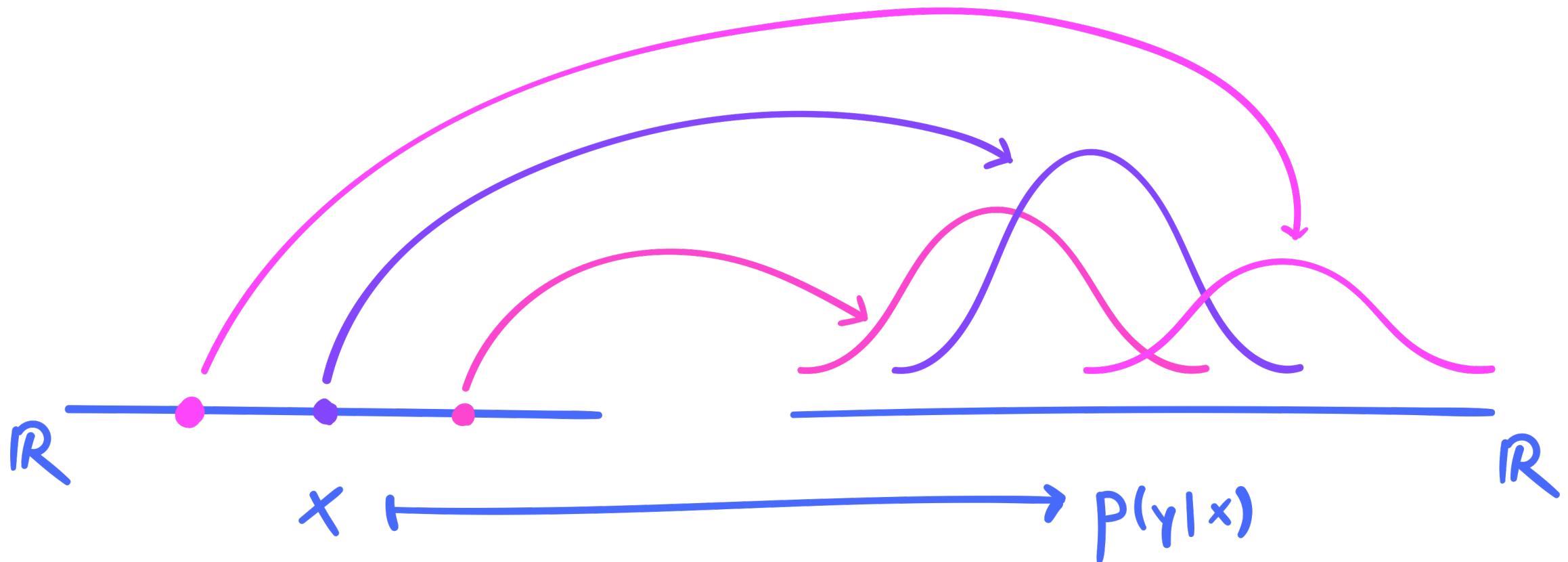








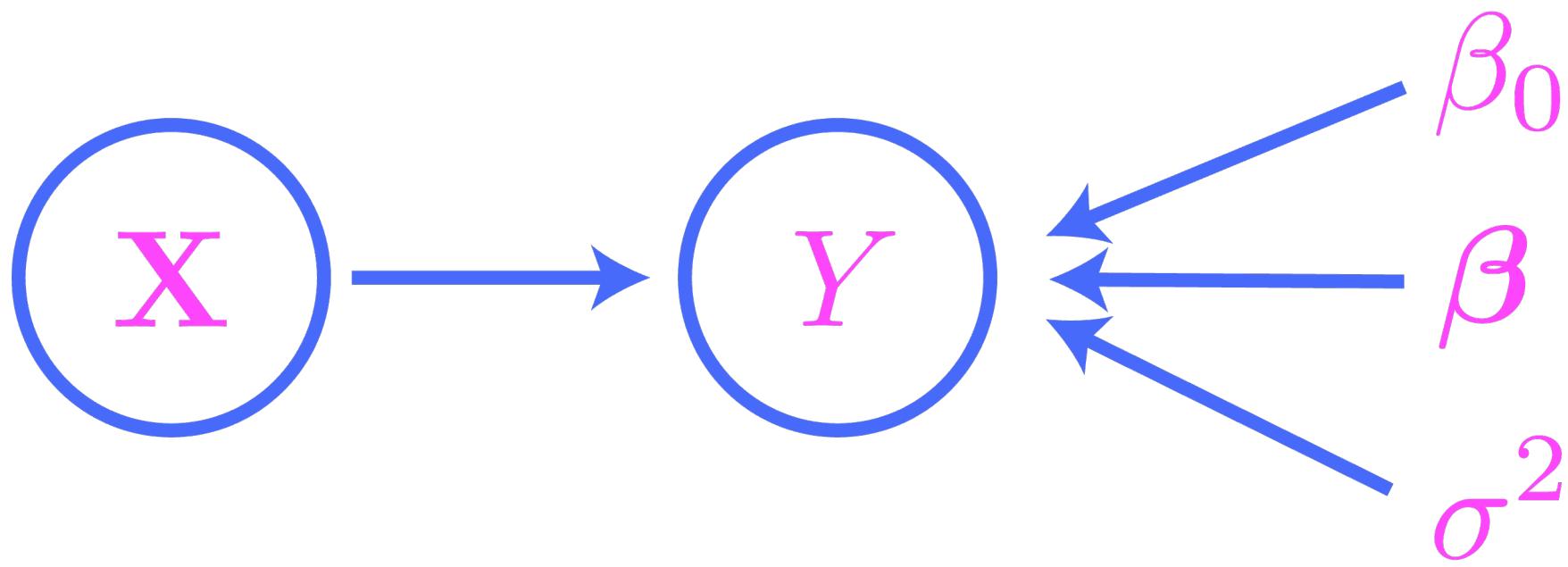
## Stochastic link (Markov Kernel)

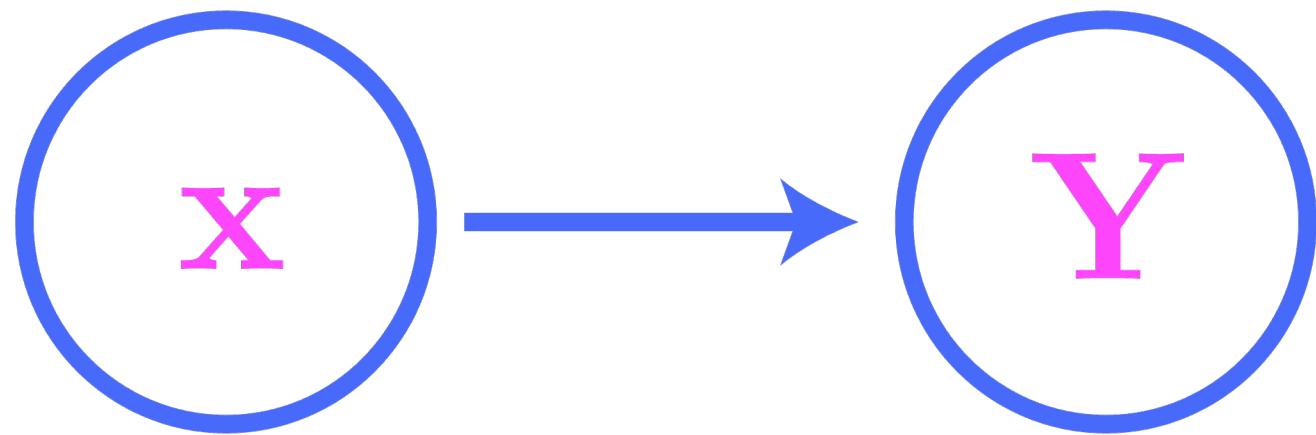


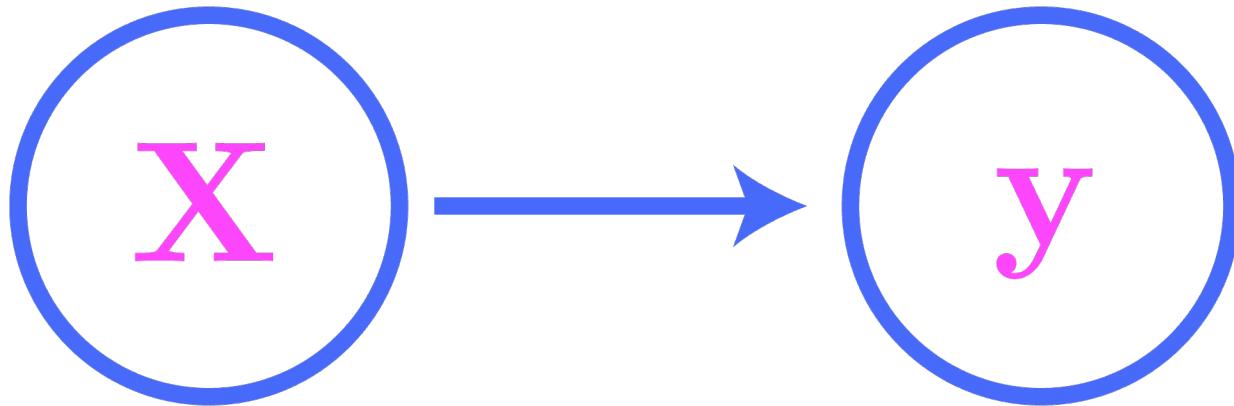


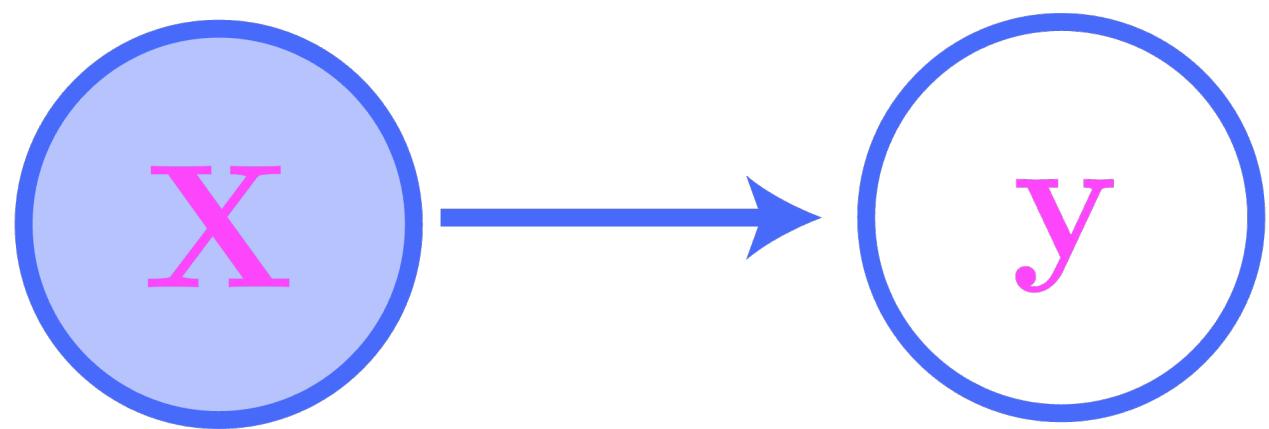
### Problem Prompt

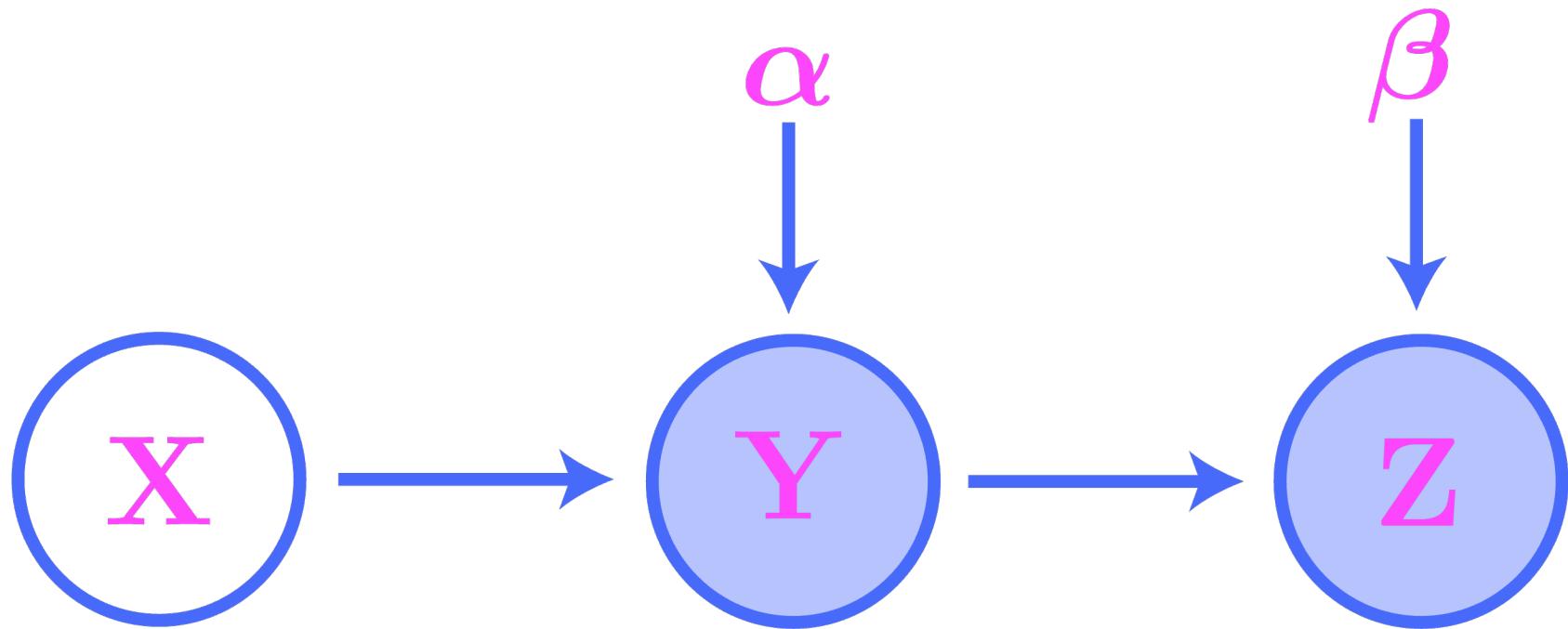
Do problem 2 on the worksheet.

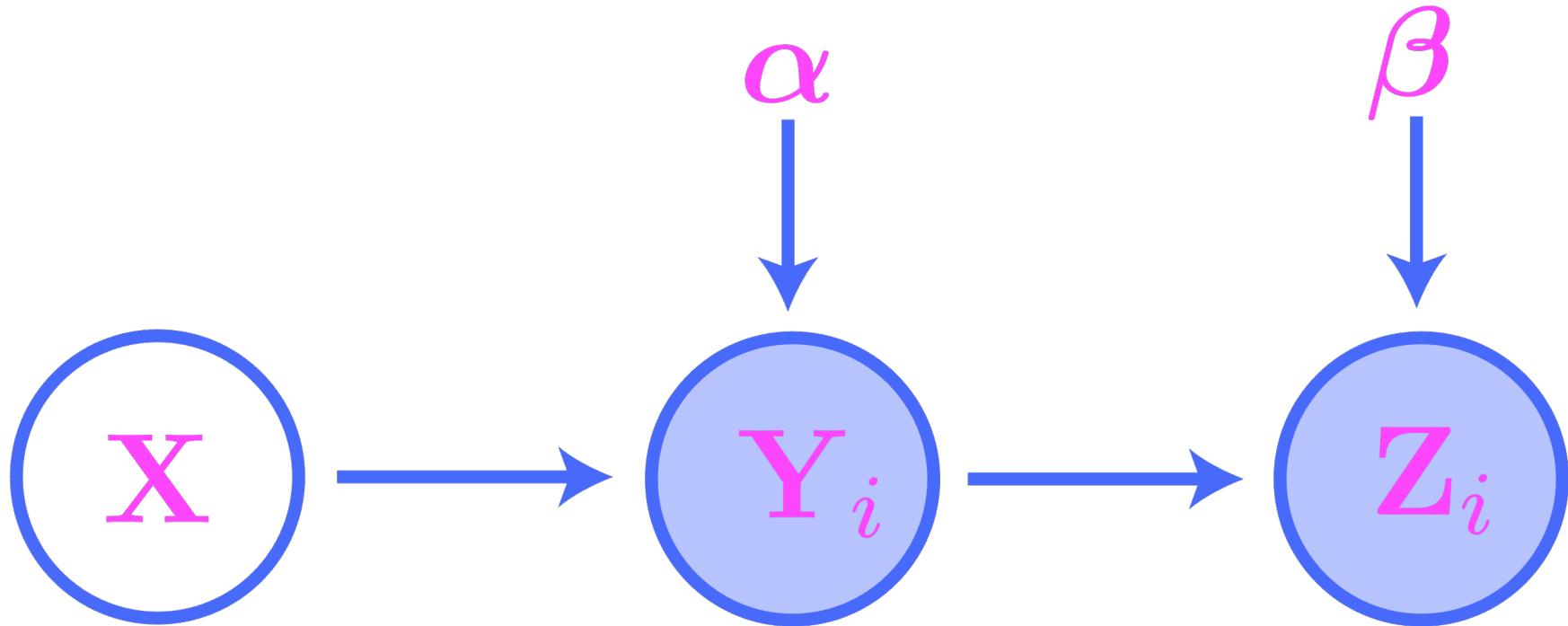


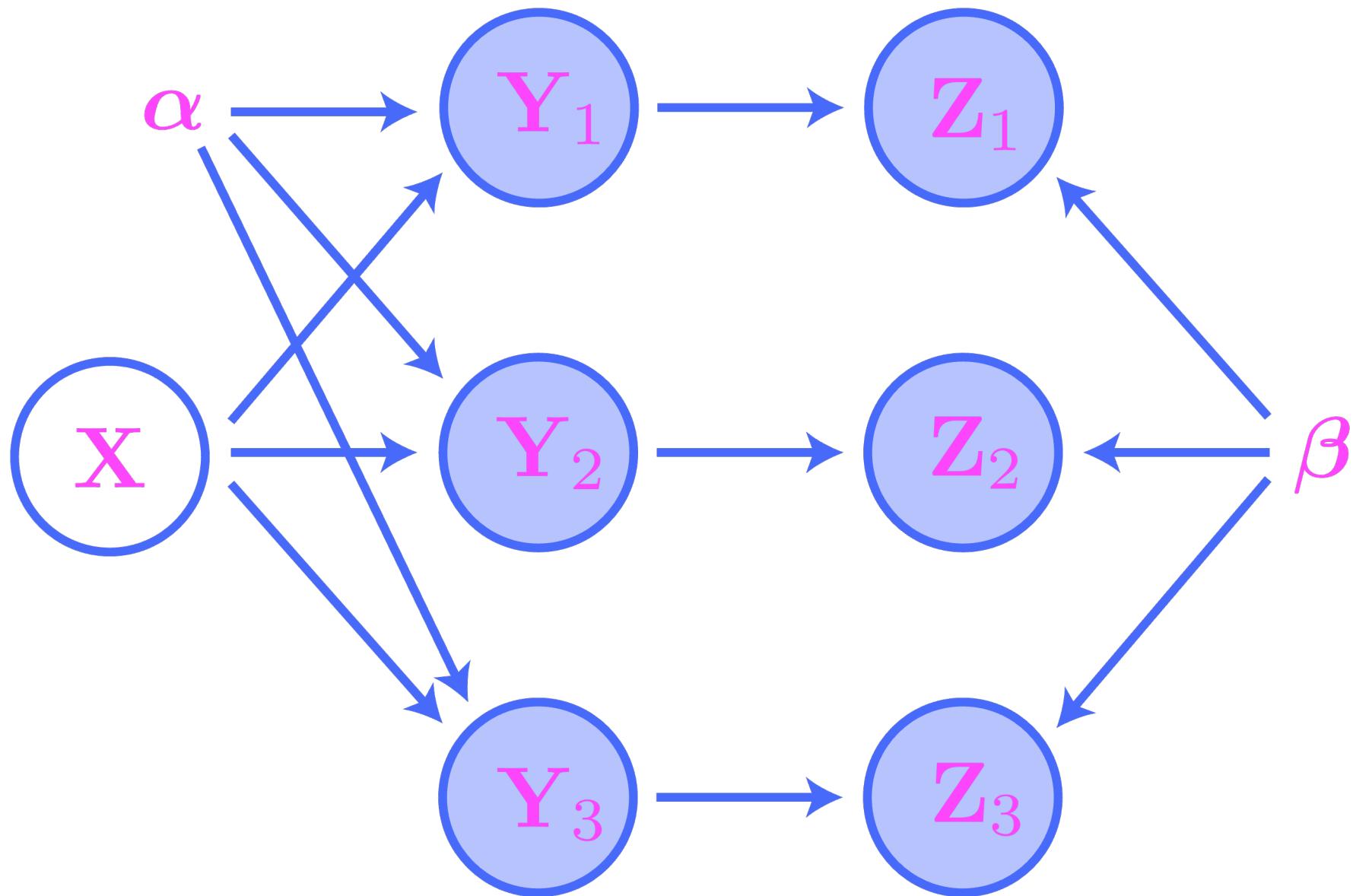


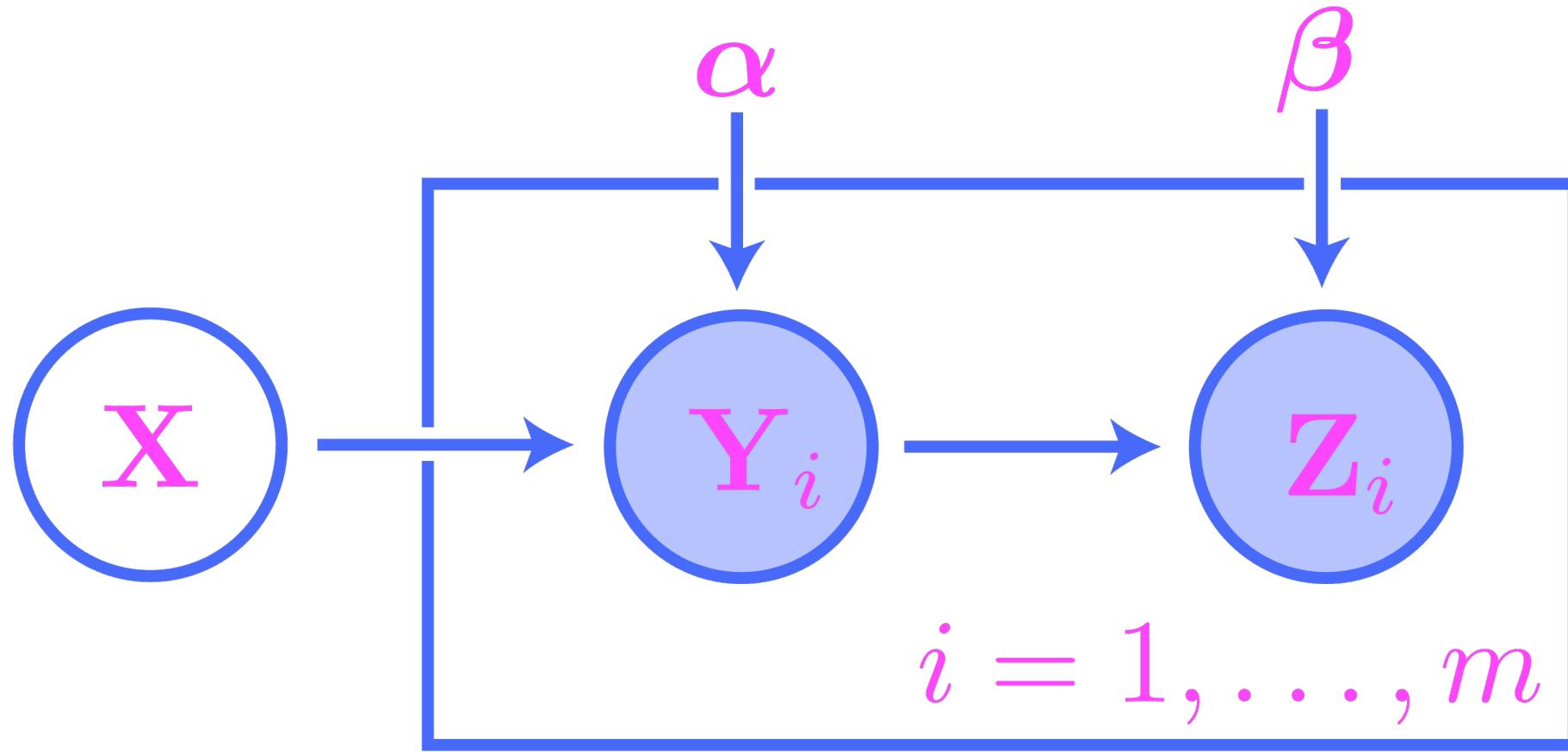














### Problem Prompt

Do problem 3 on the worksheet.

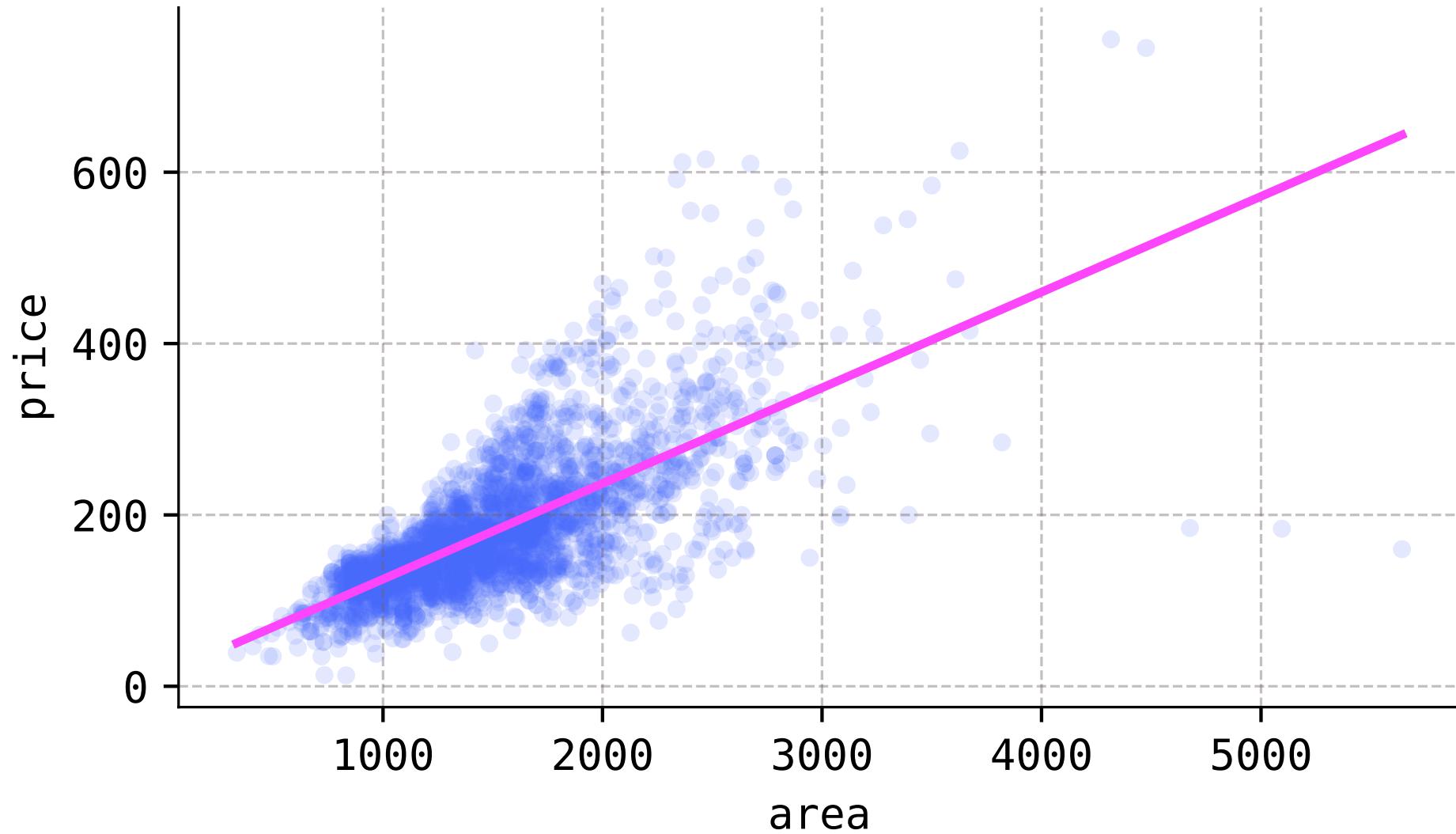


## Definition 12.1

A *probabilistic graphical model (PGM)* consists of the following:

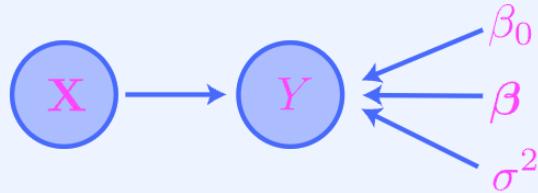
1. A set of vectors, some random and some deterministic, and some marked as observed and all others as hidden.
2. A graphical structure depicting the vectors as nodes and flows of influence (or information) as arrows between the nodes. If any of these flows are parametrized, then the graphical structure also has (un-circled) nodes for the parameters.
3. Mathematical descriptions of the flows as (possibly parametrized) link functions.

## 12.3. Linear regression models



### Definition 12.2

A *linear regression model* is a probabilistic graphical model whose underlying graph is of the form



where  $\mathbf{X} \in \mathbb{R}^n$ . The model has the following parameters:

- A real parameter  $\beta_0 \in \mathbb{R}$ .
- A parameter vector  $\boldsymbol{\beta} \in \mathbb{R}^n$ .
- A positive real parameter  $\sigma^2 > 0$ .

The link function at  $Y$  is given by

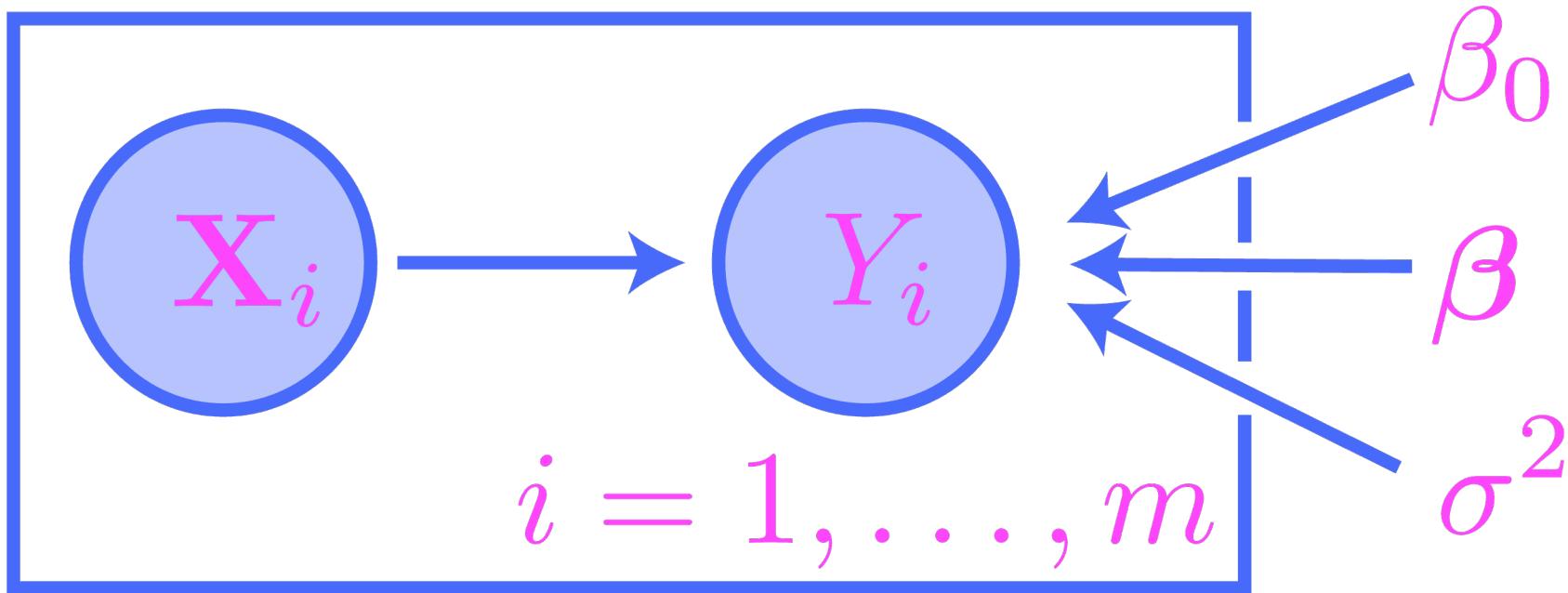
$$Y | \mathbf{X} \sim \mathcal{N}(\mu, \sigma^2), \quad \text{where } \mu = \beta_0 + \mathbf{x}^\top \boldsymbol{\beta}.$$

### Definition 12.3

For fixed  $\mathbf{x} \in \mathbb{R}^n$  and  $y \in \mathbb{R}$ , the *model likelihood function* for a linear regression model is the function

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}, \sigma^2; y | \mathbf{x}) \stackrel{\text{def}}{=} f(y | \mathbf{x}; \beta_0, \boldsymbol{\beta}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{1}{2\sigma^2} (y - \mu)^2 \right] \quad (12.2)$$

of the parameters  $\beta_0, \boldsymbol{\beta}, \sigma^2$ , where  $\mu = \beta_0 + \mathbf{x}^\top \boldsymbol{\beta}$ .



#### 🔔 Definition 12.4

Given an observed dataset

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^n \times \mathbb{R},$$

the *data likelihood function* for a linear regression model is the function

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}, \sigma^2; y_1, \dots, y_m \mid \mathbf{x}_1, \dots, \mathbf{x}_m) \stackrel{\text{def}}{=} f(y_1, \dots, y_m \mid \mathbf{x}_1, \dots, \mathbf{x}_m; \beta_0, \boldsymbol{\beta}, \sigma^2)$$

of the parameters  $\beta_0, \boldsymbol{\beta}, \sigma^2$ .

### 🔔 Theorem 12.1 (Data likelihood functions of linear regression models)

Given an observed dataset

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^n \times \mathbb{R},$$

the data likelihood function for a linear regression model is given by

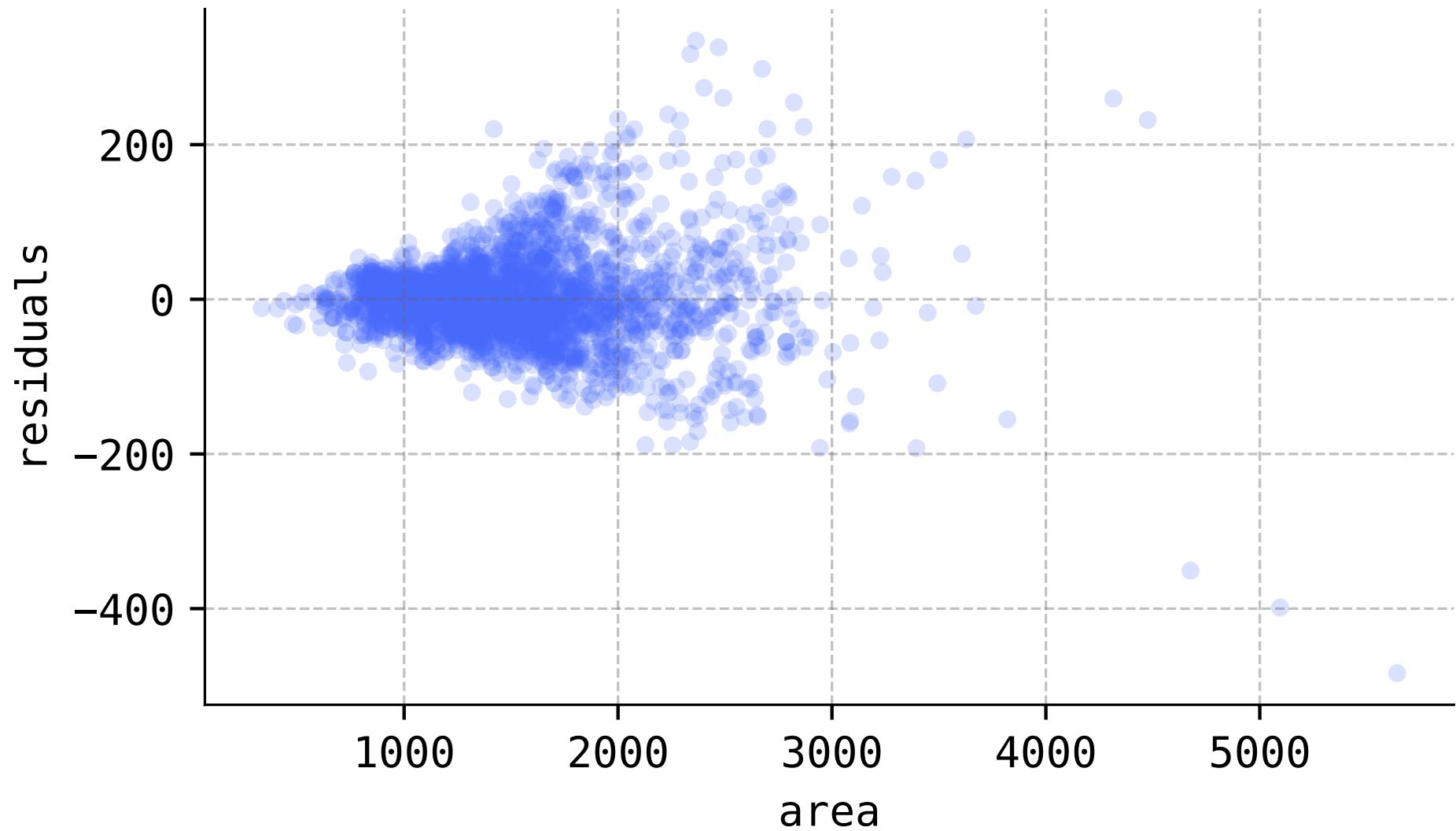
$$\begin{aligned}\mathcal{L}_{\text{data}}(\beta_0, \boldsymbol{\beta}, \sigma^2) &= \prod_{i=1}^m \mathcal{L}(\beta_0, \boldsymbol{\beta}, \sigma^2; y_i | \mathbf{x}_i) \\ &= \frac{1}{(2\pi\sigma^2)^{m/2}} \exp \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - \mu_i)^2 \right],\end{aligned}$$

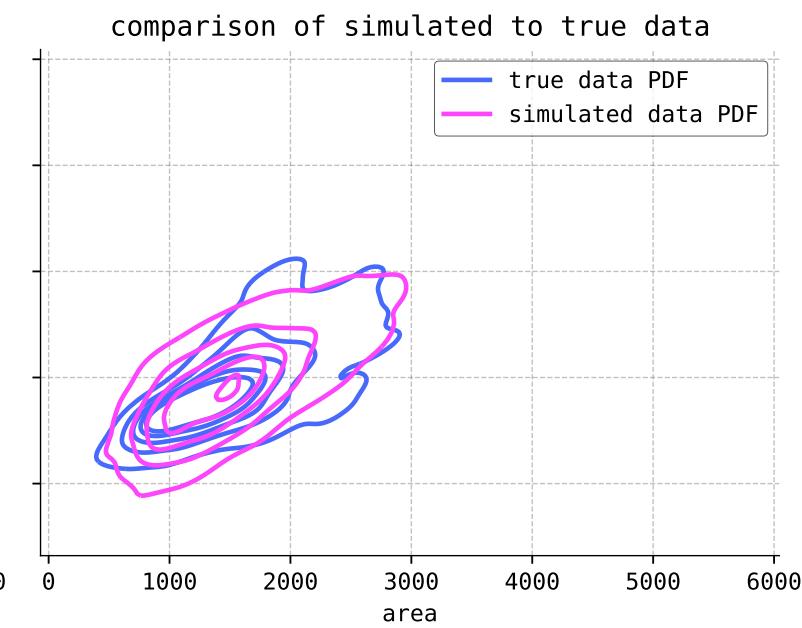
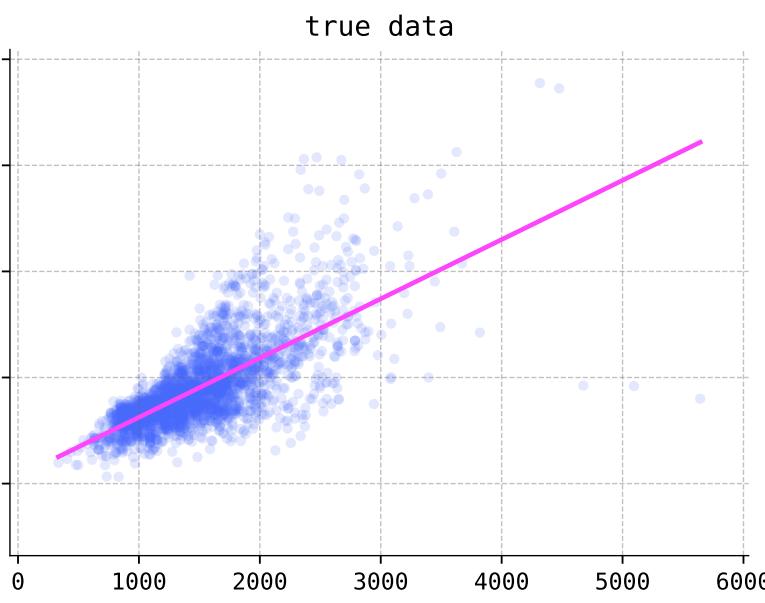
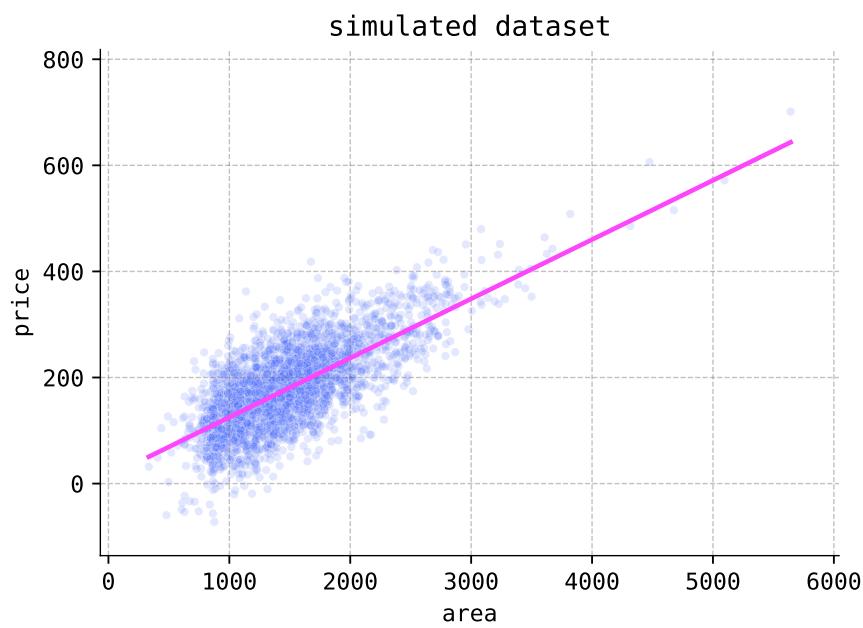
where  $\mu_i = \beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}$  for each  $i = 1, \dots, m$ .



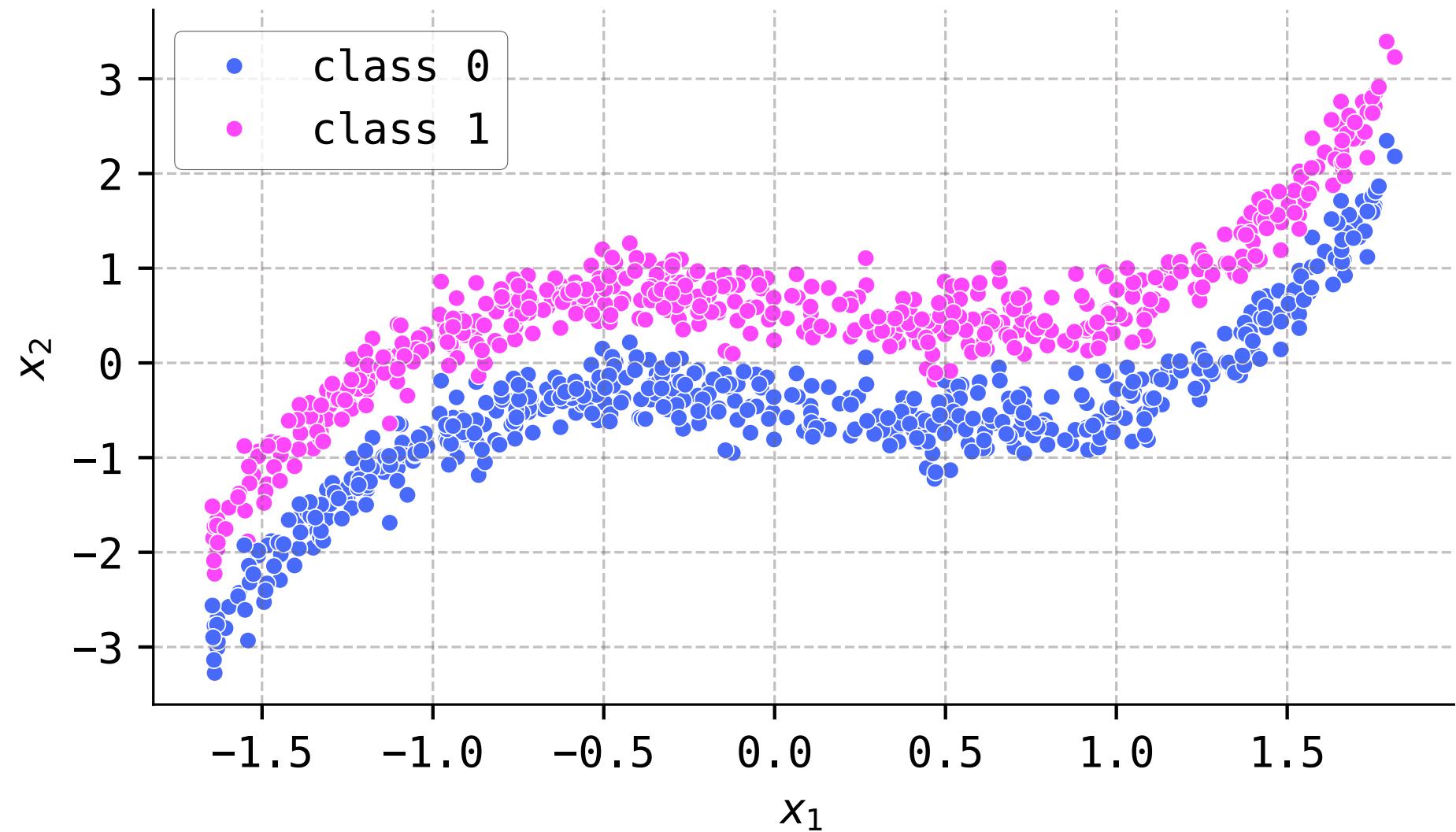
### Problem Prompt

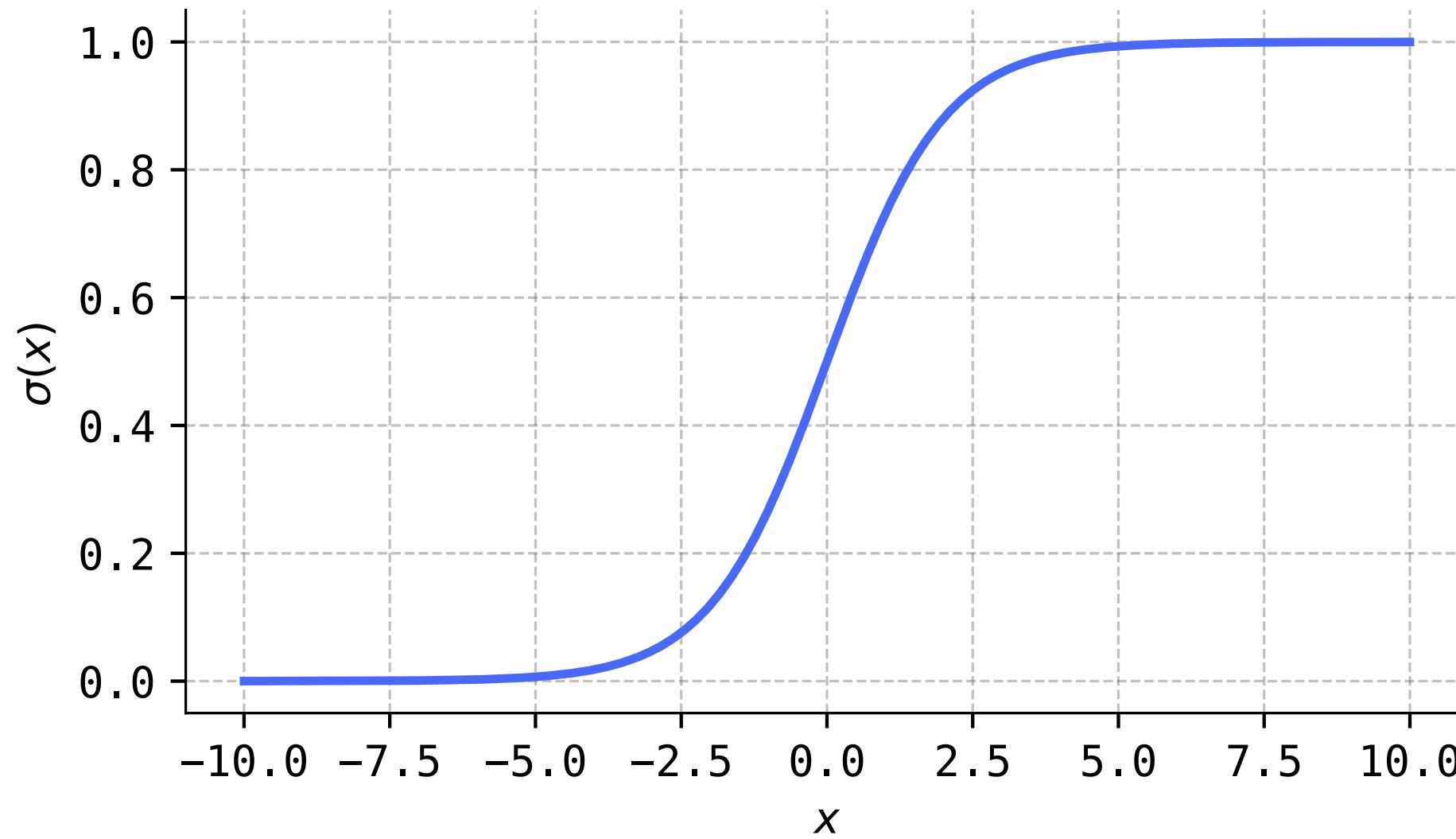
Do problems 3 and 4 on the worksheet.





## 12.4. Logistic regression models





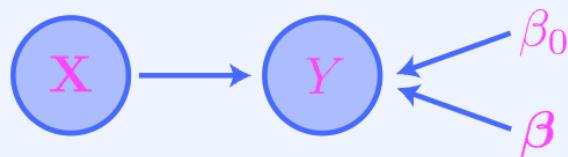


### Problem Prompt

Do problem 5 on the worksheet.

### 🔔 Definition 12.5

A *logistic regression model* is a probabilistic graphical model whose underlying graph is of the form



where  $\mathbf{X} \in \mathbb{R}^n$ . The model has the following parameters:

- A real parameter  $\beta_0 \in \mathbb{R}$ .
- A parameter vector  $\boldsymbol{\beta} \in \mathbb{R}^n$ .

The link function at  $Y$  is given by

$$Y | \mathbf{X} \sim \text{Ber}(\phi), \quad \text{where} \quad \phi = \sigma(\beta_0 + \mathbf{x}^\top \boldsymbol{\beta}),$$

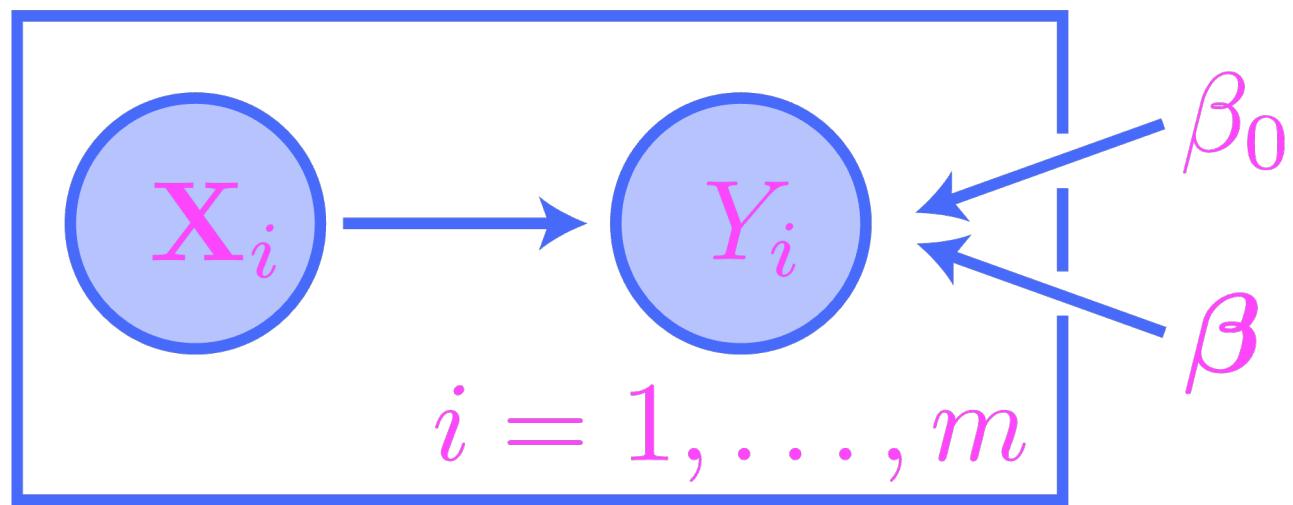
and where  $\sigma$  is the sigmoid function.

### Definition 12.6

For fixed  $\mathbf{x} \in \mathbb{R}^n$  and  $y \in \{0, 1\}$ , the *model likelihood function* for a logistic regression model is the function

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}; y | \mathbf{x}) \stackrel{\text{def}}{=} p(y | \mathbf{x}; \beta_0, \boldsymbol{\beta}) = \phi^y(1 - \phi)^{1-y}$$

of the parameters  $\beta_0, \boldsymbol{\beta}$ , where  $\phi = \sigma(\beta_0 + \mathbf{x}^\top \boldsymbol{\beta})$ .



### Definition 12.7

Given an observed dataset

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^n \times \{0, 1\},$$

the *data likelihood function* for a logistic regression model is the function

$$\mathcal{L}(\beta_0, \boldsymbol{\beta}; y_1, \dots, y_m \mid \mathbf{x}_1, \dots, \mathbf{x}_m) \stackrel{\text{def}}{=} p(y_1, \dots, y_m \mid \mathbf{x}_1, \dots, \mathbf{x}_m; \beta_0, \boldsymbol{\beta})$$

of the parameters  $\beta_0, \boldsymbol{\beta}$ .

## 🔔 Theorem 12.2 (Data likelihood functions of logistic regression models)

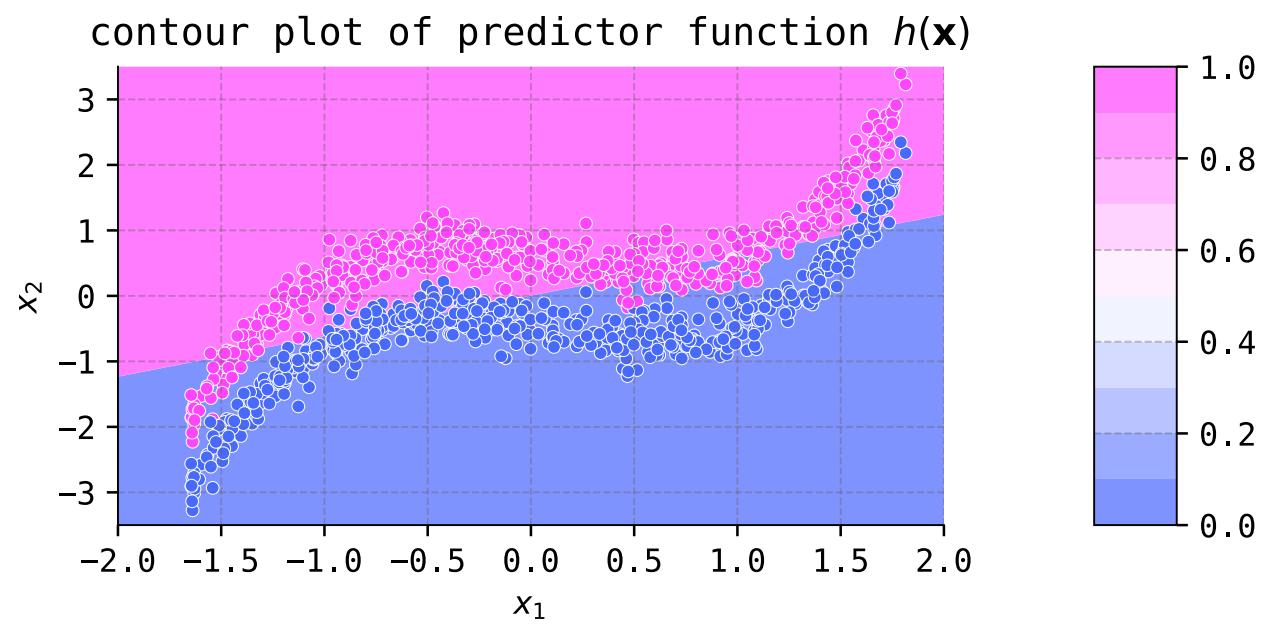
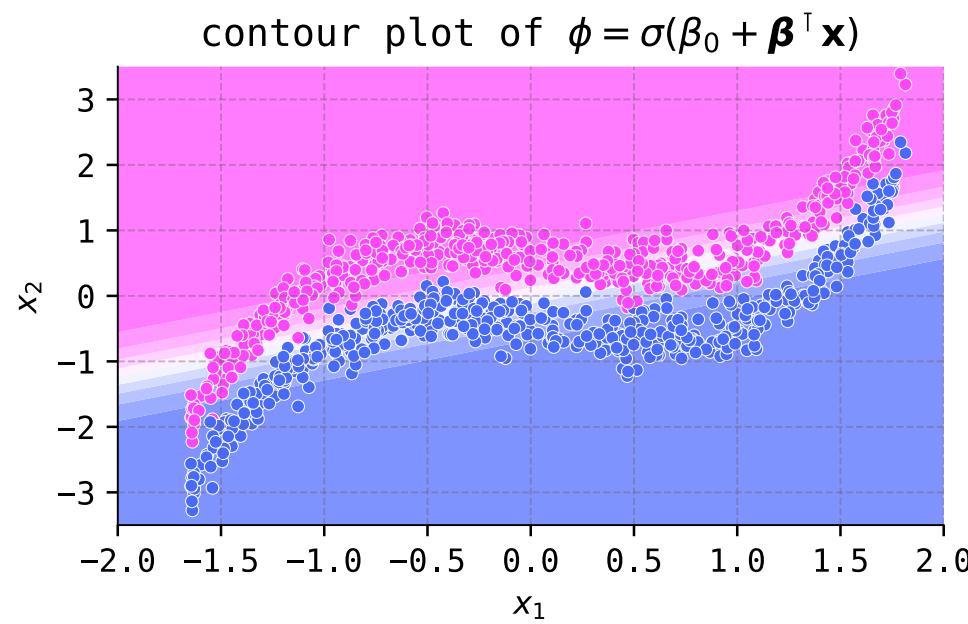
Given an observed dataset

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^n \times \{0, 1\},$$

the data likelihood function for a logistic regression model is given by

$$\mathcal{L}_{\text{data}}(\beta_0, \boldsymbol{\beta}) = \prod_{i=1}^m \mathcal{L}(\beta_0, \boldsymbol{\beta}; y_i | \mathbf{x}_i) = \prod_{i=1}^m \phi_i^{y_i} (1 - \phi_i)^{1-y_i},$$

where  $\phi_i = \sigma(\beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta})$  for each  $i = 1, \dots, m$ .

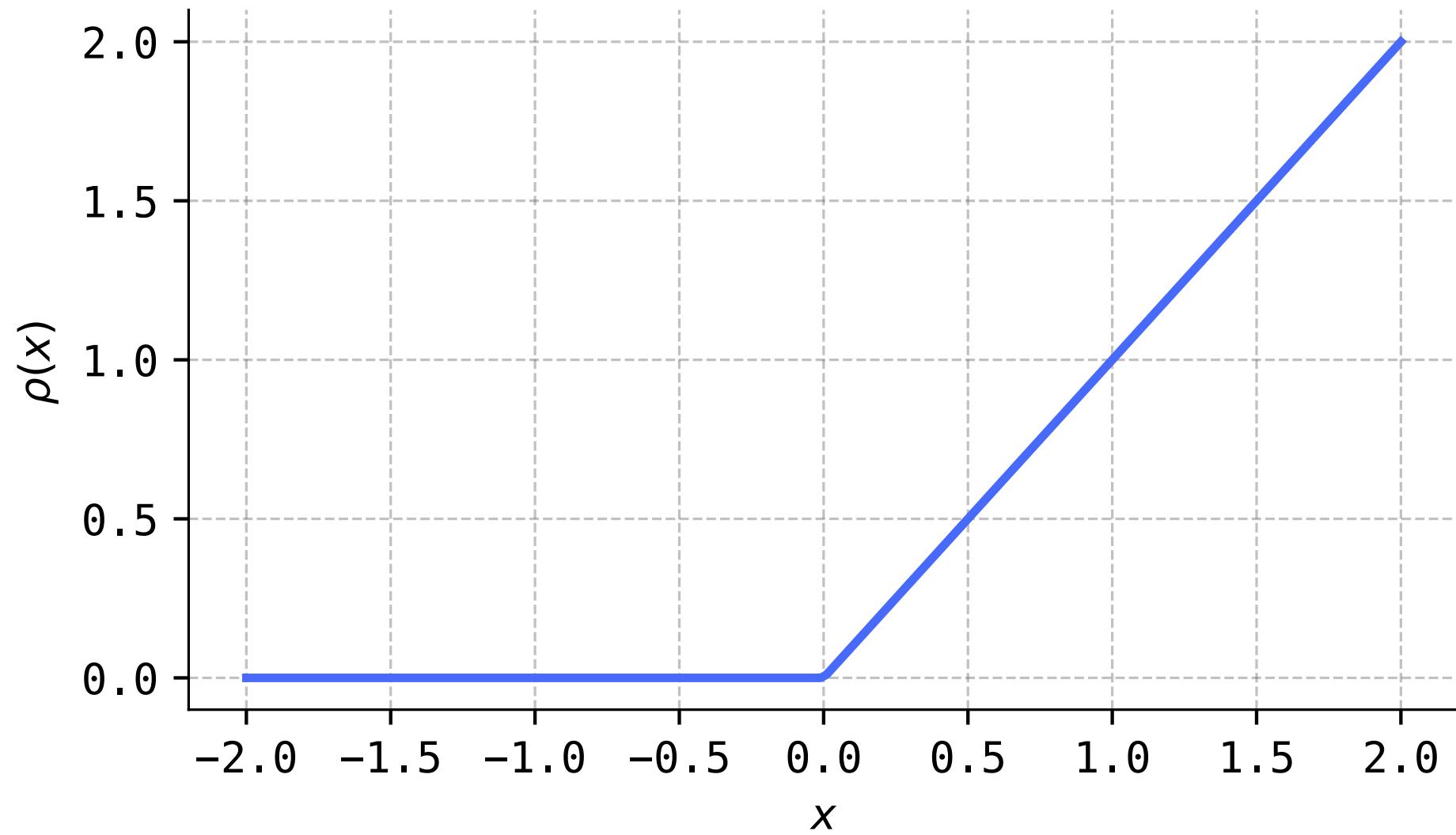




### Problem Prompt

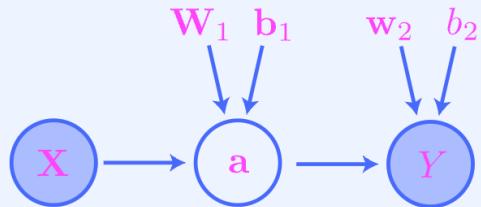
Do problem 6 on the worksheet.

## 12.5. Neural network models



### Definition 12.8

A (*fully-connected, feedforward*) neural network with one hidden layer is a probabilistic graphical model whose underlying graph is of the form



where  $\mathbf{X} \in \mathbb{R}^n$  and  $\mathbf{a} \in \mathbb{R}^p$ . The model has the following parameters:

- A parameter matrix  $\mathbf{W}_1 \in \mathbb{R}^{n \times p}$ .
- A parameter vector  $\mathbf{b}_1 \in \mathbb{R}^p$ .
- A parameter vector  $\mathbf{w}_2 \in \mathbb{R}^p$ .
- A real parameter  $b_2 \in \mathbb{R}$ .

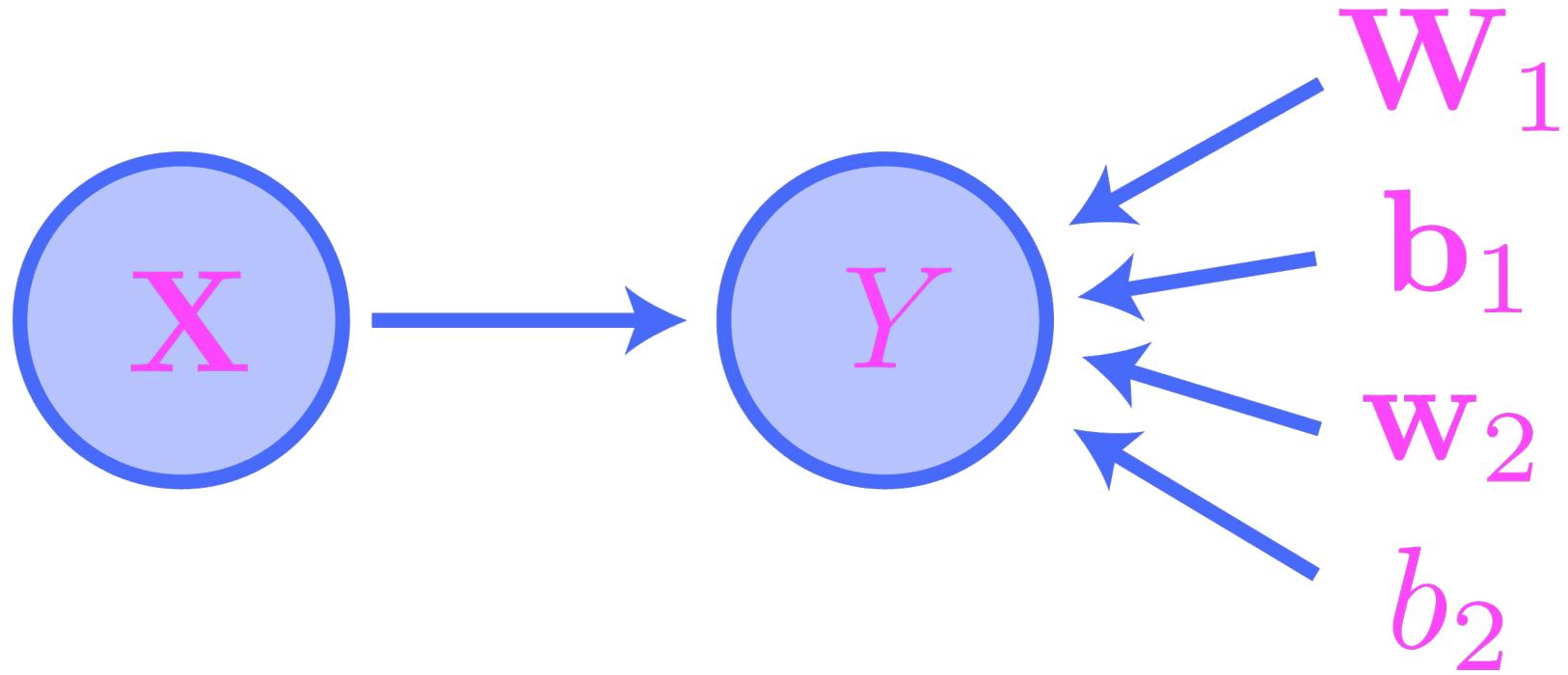
The link function at  $\mathbf{a}$  is given by

$$\mathbf{a}^\top = \rho(\mathbf{x}^\top \mathbf{W}_1 + \mathbf{b}_1^\top),$$

while the link function at  $Y$  is given by

$$Y | \mathbf{X} \sim \text{Ber}(\phi), \quad \text{where } \phi = \sigma(\mathbf{a}^\top \mathbf{w}_2 + b_2). \quad (12.5)$$

Here,  $\rho$  is the ReLU function and  $\sigma$  is the sigmoid function.



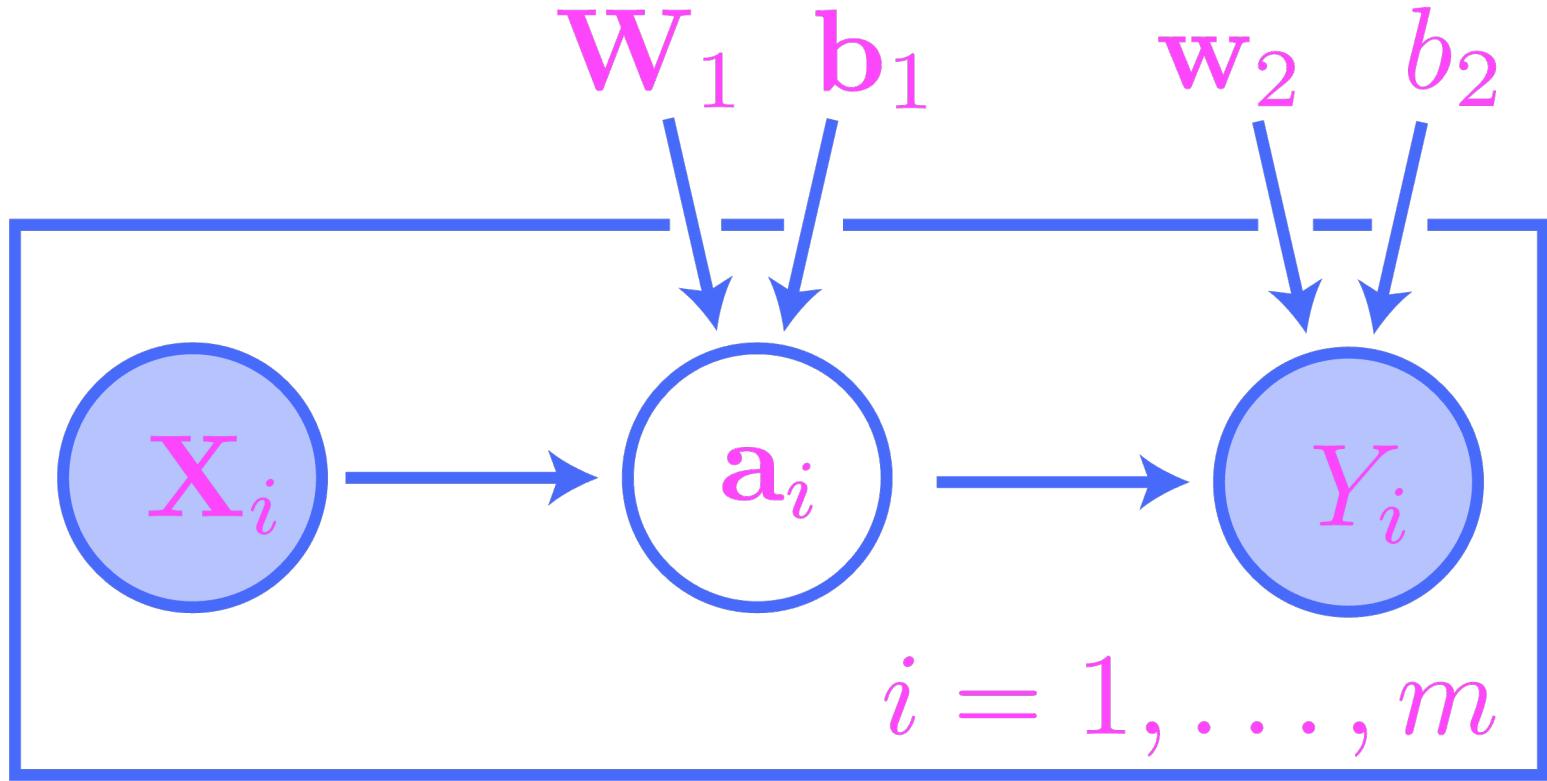
### Definition 12.9

For fixed  $\mathbf{x} \in \mathbb{R}^n$  and  $y \in \{0, 1\}$ , the *model likelihood function* for a neural network model is the function

$$\mathcal{L}(\mathbf{W}_1, \mathbf{b}_1, \mathbf{w}_2, b_2; y | \mathbf{x}) \stackrel{\text{def}}{=} p(y | \mathbf{x}; \mathbf{W}_1, \mathbf{b}_1, \mathbf{w}_2, b_2) = \phi^y(1 - \phi)^{1-y}$$

of the parameters  $\mathbf{W}_1, \mathbf{b}_1, \mathbf{w}_2, b_2$ , where

$$\begin{aligned}\mathbf{a}^\top &= \rho(\mathbf{x}^\top \mathbf{W}_1 + \mathbf{b}_1^\top), \\ \phi &= \sigma(\mathbf{a}^\top \mathbf{w}_2 + b_2).\end{aligned}$$



### Definition 12.10

Given an observed dataset

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^n \times \{0, 1\},$$

the *data likelihood function* for a neural network model is the function

$$\mathcal{L}(\mathbf{W}_1, \mathbf{b}_1, \mathbf{w}_2, b_2; y_1, \dots, y_m \mid \mathbf{x}_1, \dots, \mathbf{x}_m) \stackrel{\text{def}}{=} p(y_1, \dots, y_m \mid \mathbf{x}_1, \dots, \mathbf{x}_m; \mathbf{W}_1, \mathbf{b}_1, \mathbf{w}_2, b_2)$$

of the parameters  $\mathbf{W}_1, \mathbf{b}_1, \mathbf{w}_2, b_2$ .

🔔 **Theorem 12.3 (Data likelihood functions of neural network models)**

Given an observed dataset

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m) \in \mathbb{R}^n \times \{0, 1\},$$

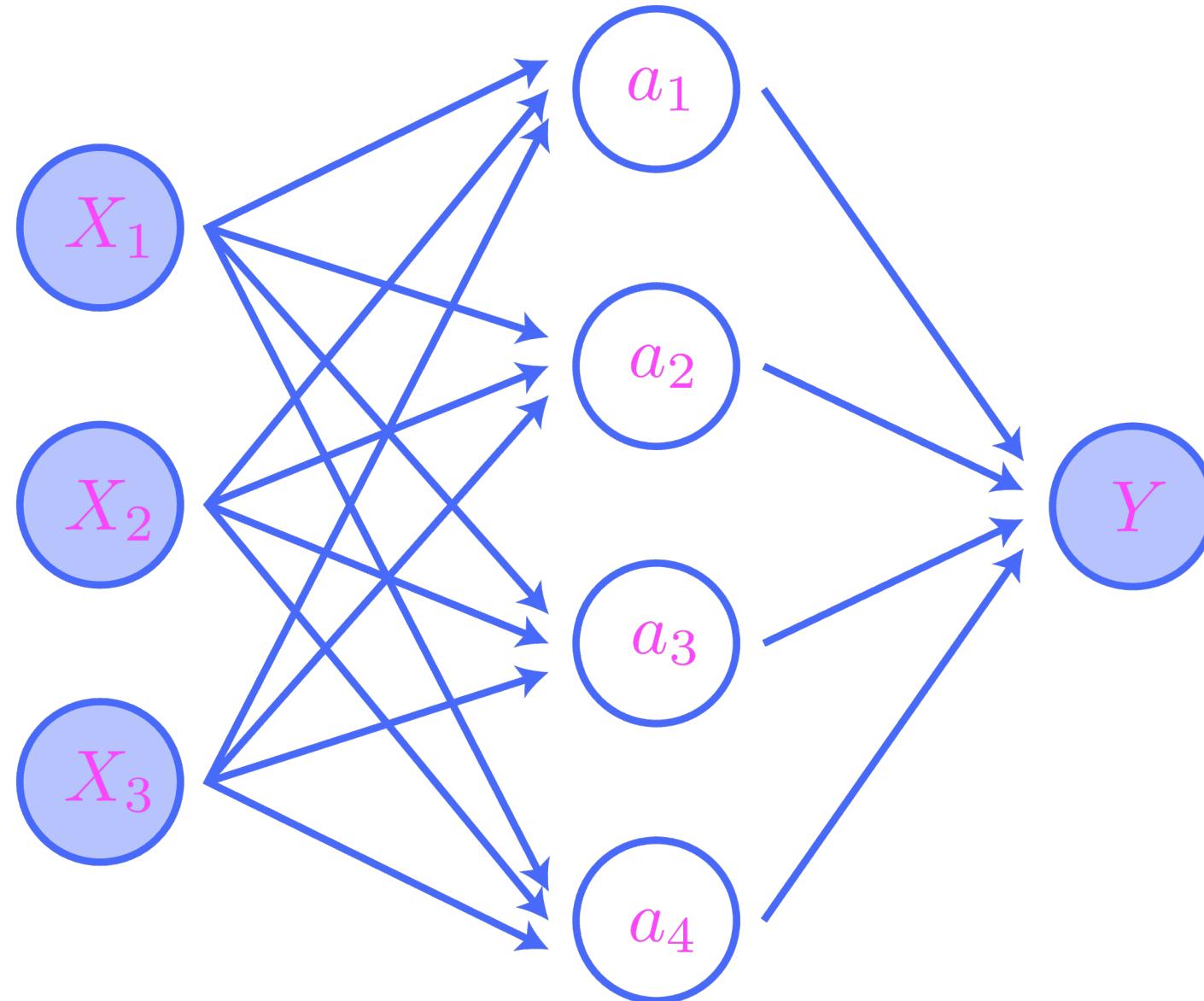
the data likelihood function for a neural network model is given by

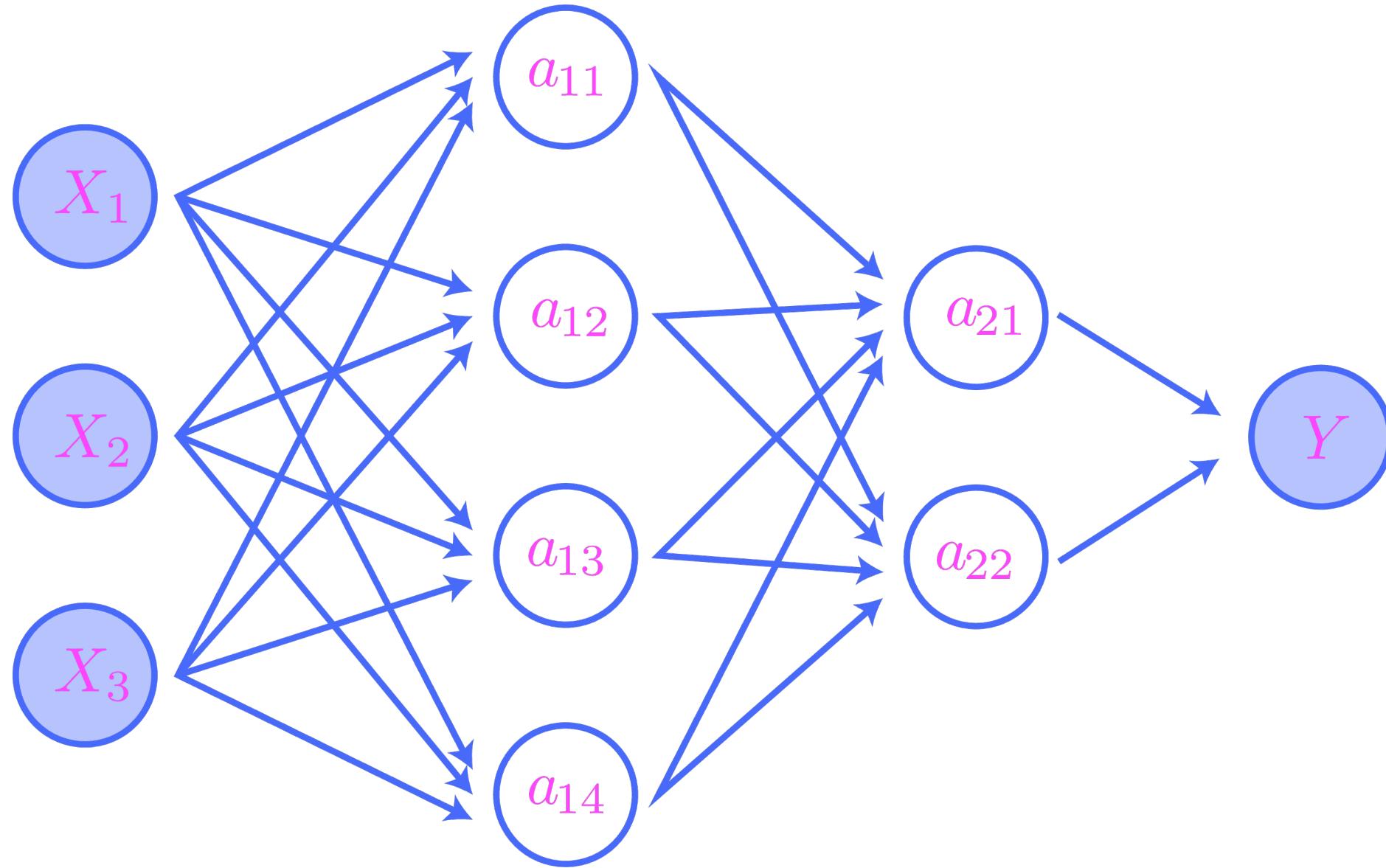
$$\mathcal{L}_{\text{data}}(\mathbf{W}_1, \mathbf{b}_1, \mathbf{w}_2, b_2) = \prod_{i=1}^m \mathcal{L}(\mathbf{W}_1, \mathbf{b}_1, \mathbf{w}_2, b_2; y_i \mid \mathbf{x}_i) = \prod_{i=1}^m \phi_i^{y_i} (1 - \phi_i)^{1-y_i}$$

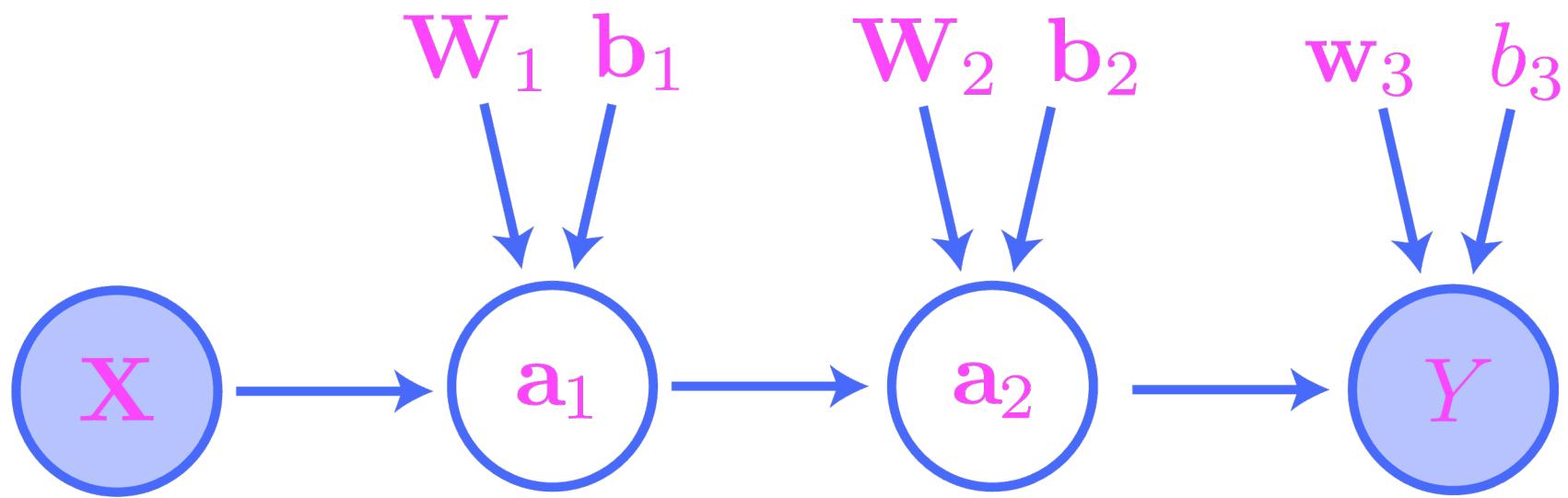
where

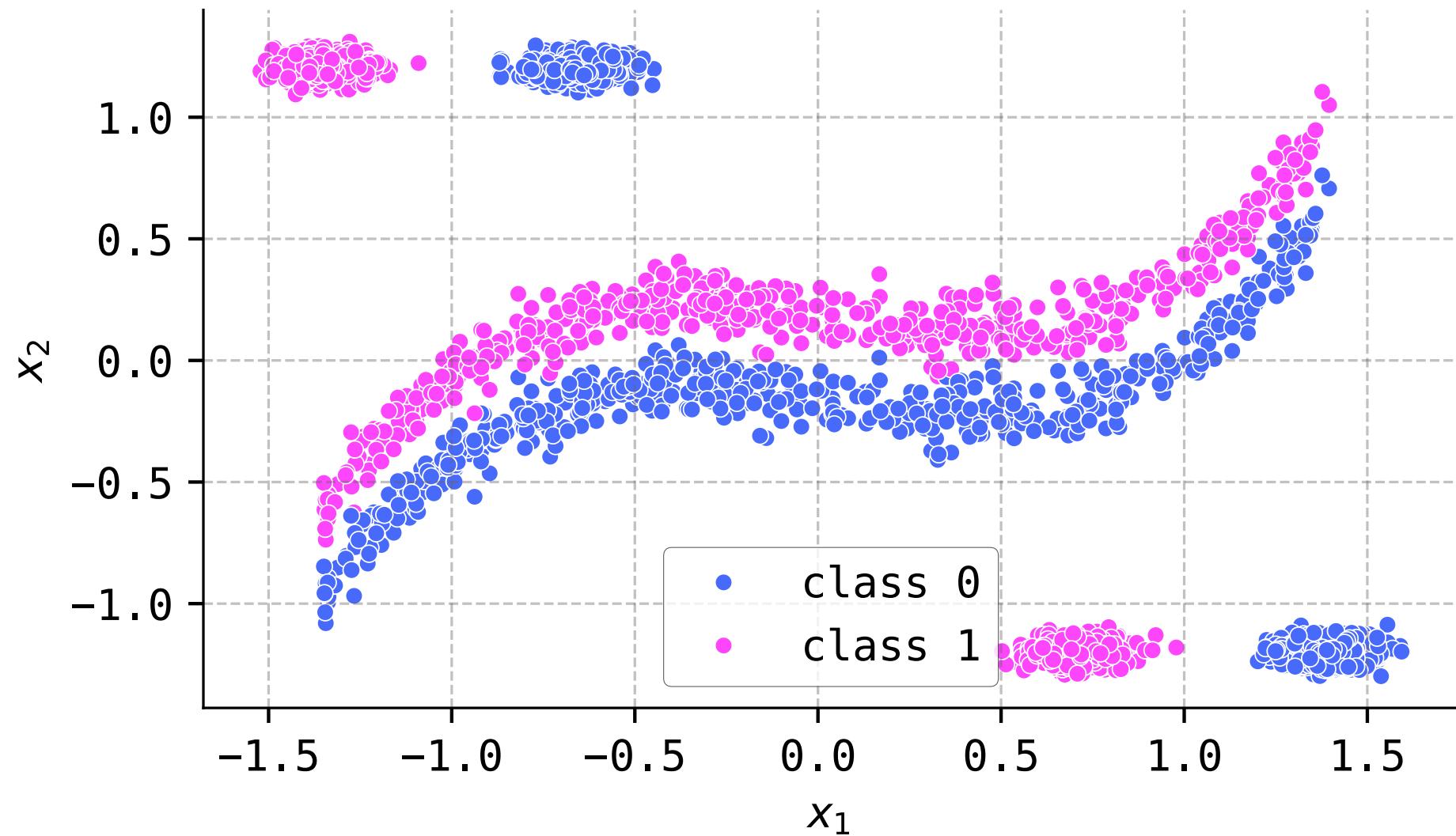
$$\begin{aligned}\mathbf{a}_i^\top &= \rho(\mathbf{x}_i^\top \mathbf{W}_1 + \mathbf{b}_1^\top), \\ \phi_i &= \sigma(\mathbf{a}_i^\top \mathbf{w}_2 + b_2),\end{aligned}$$

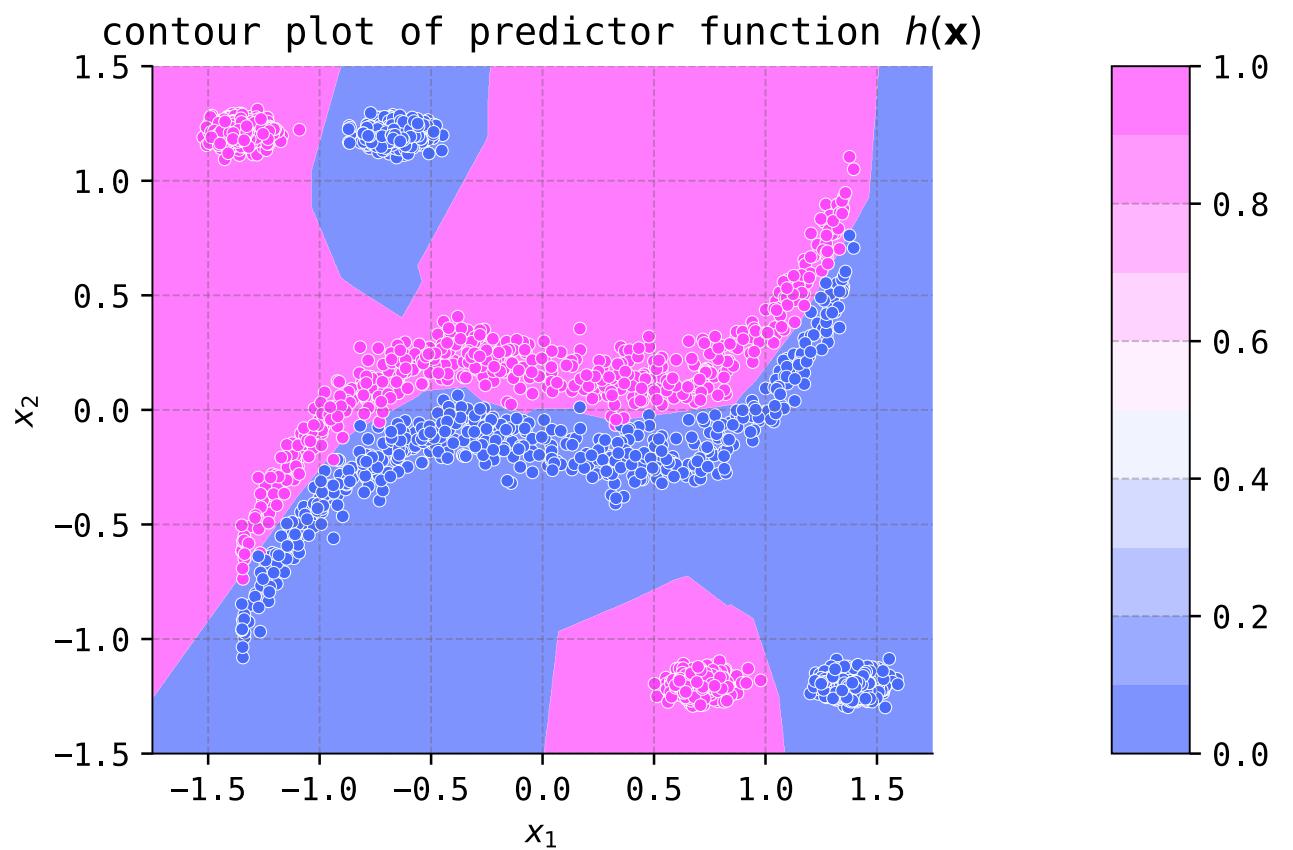
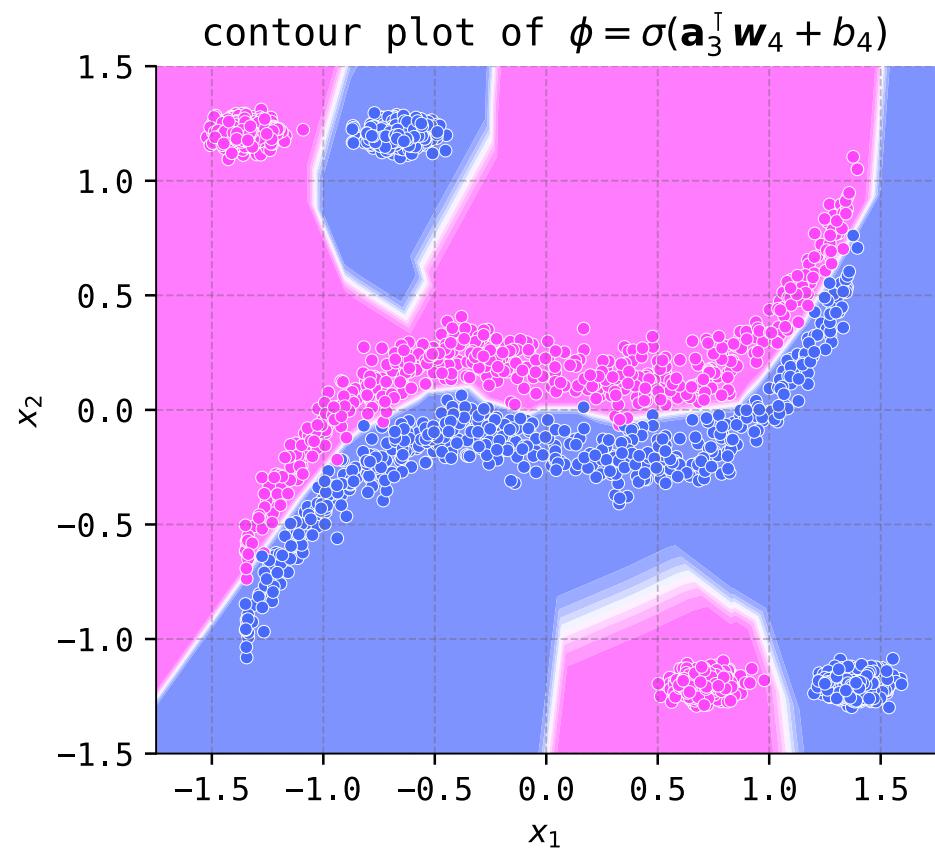
for each  $i = 1, \dots, m$ .



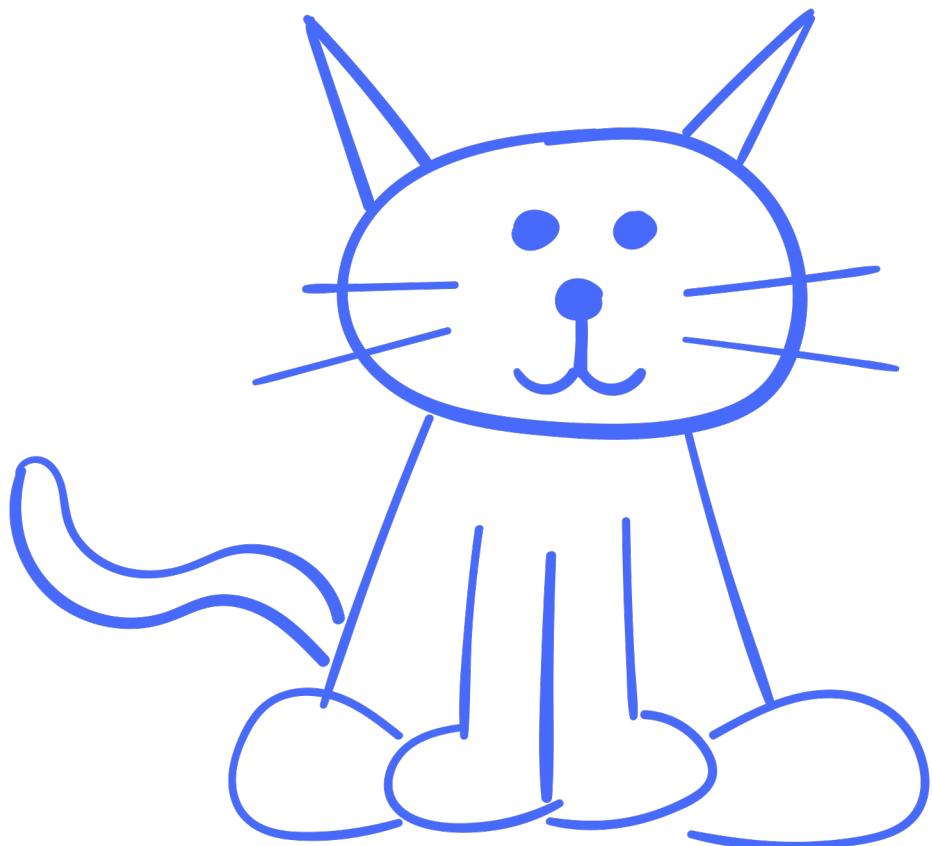








real picture of  
the author's real cat



representation

$$\begin{bmatrix} -1 \\ 2 \\ -3 \\ 4 \end{bmatrix}$$

