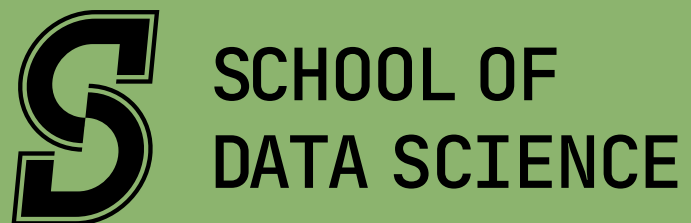


# Statistical Learning STA4042



INTERPRETABILITY  
VS.  
FLEXIBILITY



## A general regression example

Consider a general regression problem:

- $p$  random variables  $X^{(1)}, \dots, X^{(p)}$
- a target random variable  $Y = f(X^{(1)}, \dots, X^{(p)})$ .

# Statistical Learning STA4042



INTERPRETABILITY  
VS.  
FLEXIBILITY



SCHOOL OF  
DATA SCIENCE

## A general regression example

Consider a general regression problem:

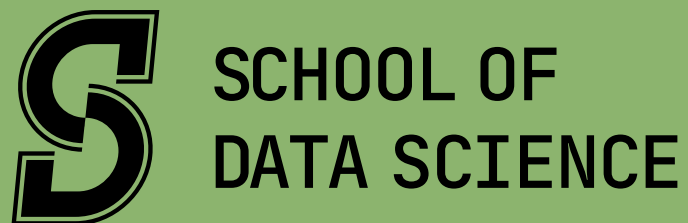
- $p$  random variables  $X^{(1)}, \dots, X^{(p)}$
- a target random variable  $Y = f(X^{(1)}, \dots, X^{(p)})$ .

*ex: independent drawing = individual,  
 $X^{(1)}, \dots, X^{(p)}$  are some characteristics of each individual (years  
of education, seniority...)  $Y$  is the level of income*

# Statistical Learning STA4042



INTERPRETABILITY  
VS.  
FLEXIBILITY



## A general regression example

Consider a general regression problem:

- $p$  random variables  $X^{(1)}, \dots, X^{(p)}$
- a target random variable  $Y = f(X^{(1)}, \dots, X^{(p)})$ .

*ex: independent drawing = individual,  
 $X^{(1)}, \dots, X^{(p)}$  are some characteristics of each individual (years of education, seniority...)  $Y$  is the level of income*

**Goal:** Given a training data set ( $= n$  independent drawings of  $Y$  and  $X = (X^{(1)}, \dots, X^{(p)})$ ):

$$D \equiv \{(x_1, y_1), \dots, (x_n, y_n)\}$$

→ Approximate  $f$  with a mapping  $\hat{f}_D$  s.t.  $Y \approx \hat{f}_D(X)$ .

**Simple solution:**

$$\text{Look for } \hat{f} : (x^{(1)}, \dots, x^{(p)}) \mapsto \beta_1 x^{(1)} + \dots + \beta_p x^{(p)}$$

# Statistical Learning STA4042



INTERPRETABILITY  
VS.  
FLEXIBILITY



SCHOOL OF  
DATA SCIENCE

## A general regression example

Consider a general regression problem:

- $p$  random variables  $X^{(1)}, \dots, X^{(p)}$
- a target random variable  $Y = f(X^{(1)}, \dots, X^{(p)})$ .

*ex: independent drawing = individual,  
 $X^{(1)}, \dots, X^{(p)}$  are some characteristics of each individual (years of education, seniority...)  $Y$  is the level of income*

**Goal:** Given a training data set ( $= n$  independent drawings of  $Y$  and  $X = (X^{(1)}, \dots, X^{(p)})$ ):

$$D \equiv \{(x_1, y_1), \dots, (x_n, y_n)\}$$

→ Approximate  $f$  with a mapping  $\hat{f}_D$  s.t.  $Y \approx \hat{f}_D(X)$ .

**Simple solution:**

$$\text{Look for } \hat{f} : (x^{(1)}, \dots, x^{(p)}) \mapsto \beta_1 x^{(1)} + \dots + \beta_p x^{(p)}$$

Parameters = Interpretation

(No parameter = no interpretation)

# Statistical Learning STA4042

## INTERPRETABILITY VS. FLEXIBILITY



# A general regression example

Consider a general regression problem:

- $p$  random variables  $X^{(1)}, \dots, X^{(p)}$
- a target random variable  $Y = f(X^{(1)}, \dots, X^{(p)})$ .

*ex: independent drawing = individual,  
 $X^{(1)}, \dots, X^{(p)}$  are some characteristics of each individual (years of education, seniority...)  $Y$  is the level of income*

**Goal:** Given a training data set ( $= n$  independent drawings of  $Y$  and  $X = (X^{(1)}, \dots, X^{(p)})$ ):

$$D \equiv \{(x_1, y_1), \dots, (x_n, y_n)\}$$

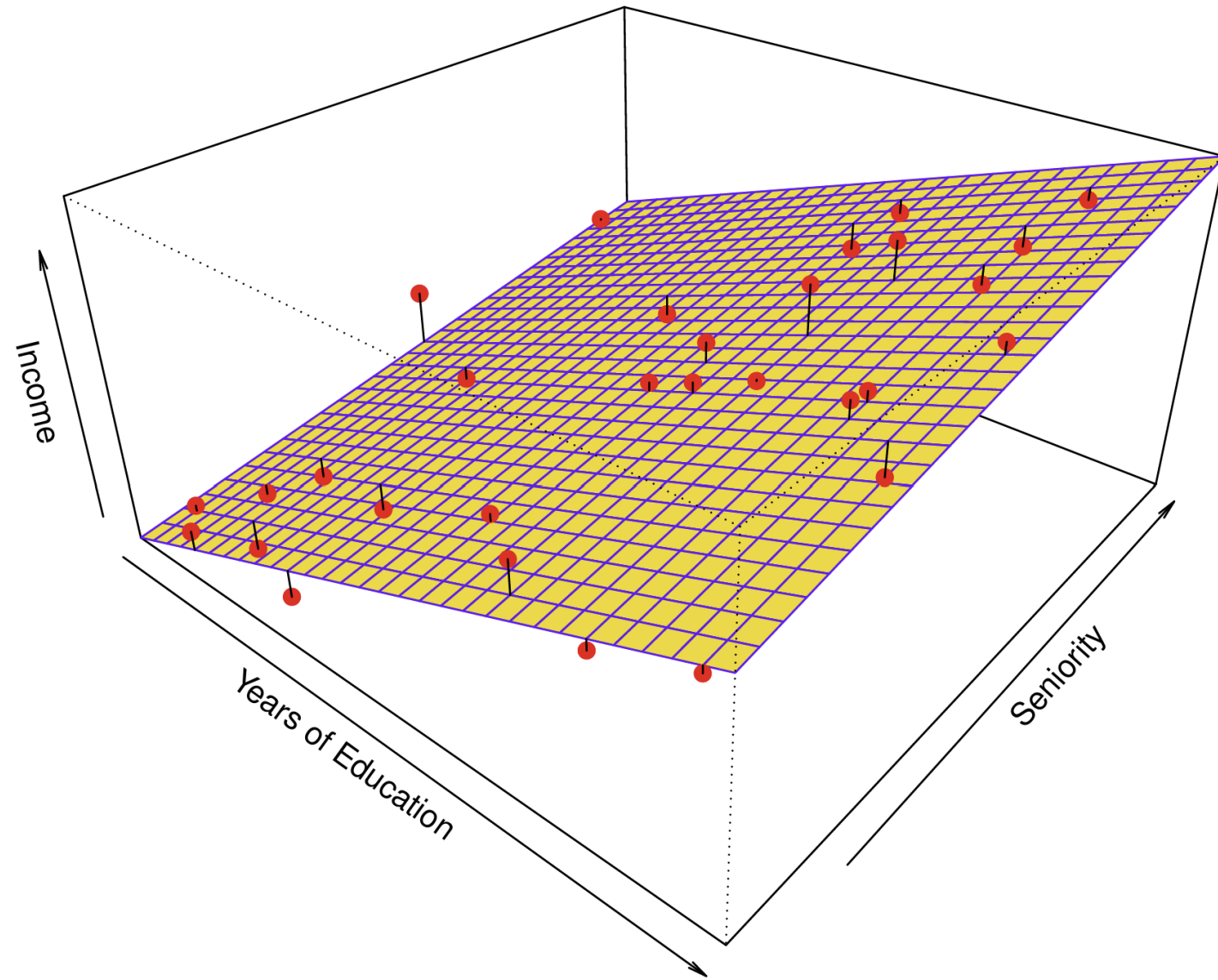
→ Approximate  $f$  with a mapping  $\hat{f}_D$  s.t.  $Y \approx \hat{f}_D(X)$ .

**Simple solution:**

$$\text{Look for } \hat{f} : (x^{(1)}, \dots, x^{(p)}) \mapsto \beta_1 x^{(1)} + \dots + \beta_p x^{(p)}$$

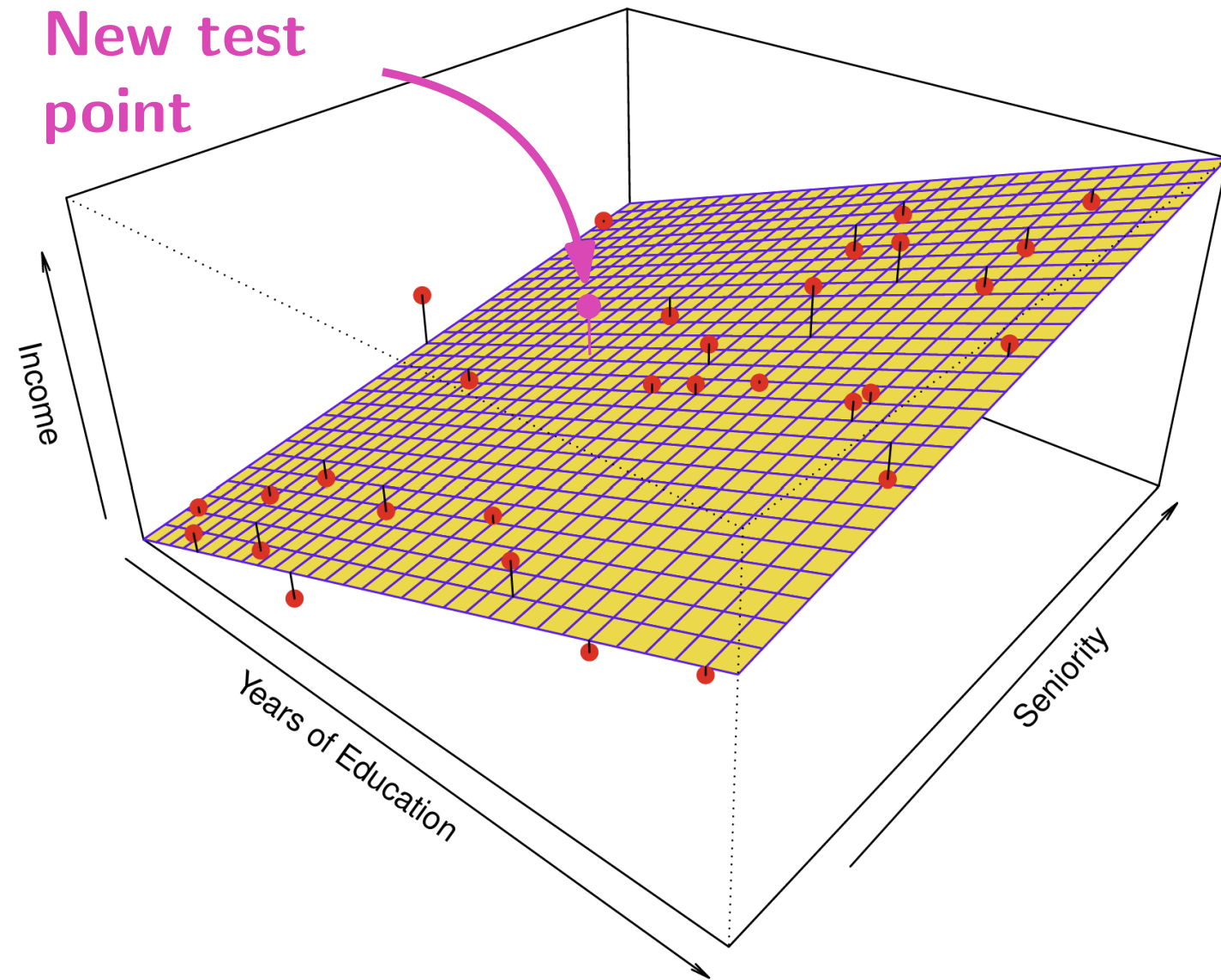
Parameters = Interpretation       $\beta_1, \dots, \beta_p$  provide dependence on  
(No parameter = no interpretation)      each predictor  $x^{(1)}, \dots, x^{(p)}$

# A general regression example



“Parametric model”: look for  $\beta_0, \beta_1, \beta_2$  s.t.:  
“Income” =  $\beta_0 + \beta_1$  “y. o. education” +  $\beta_2$  “seniority”

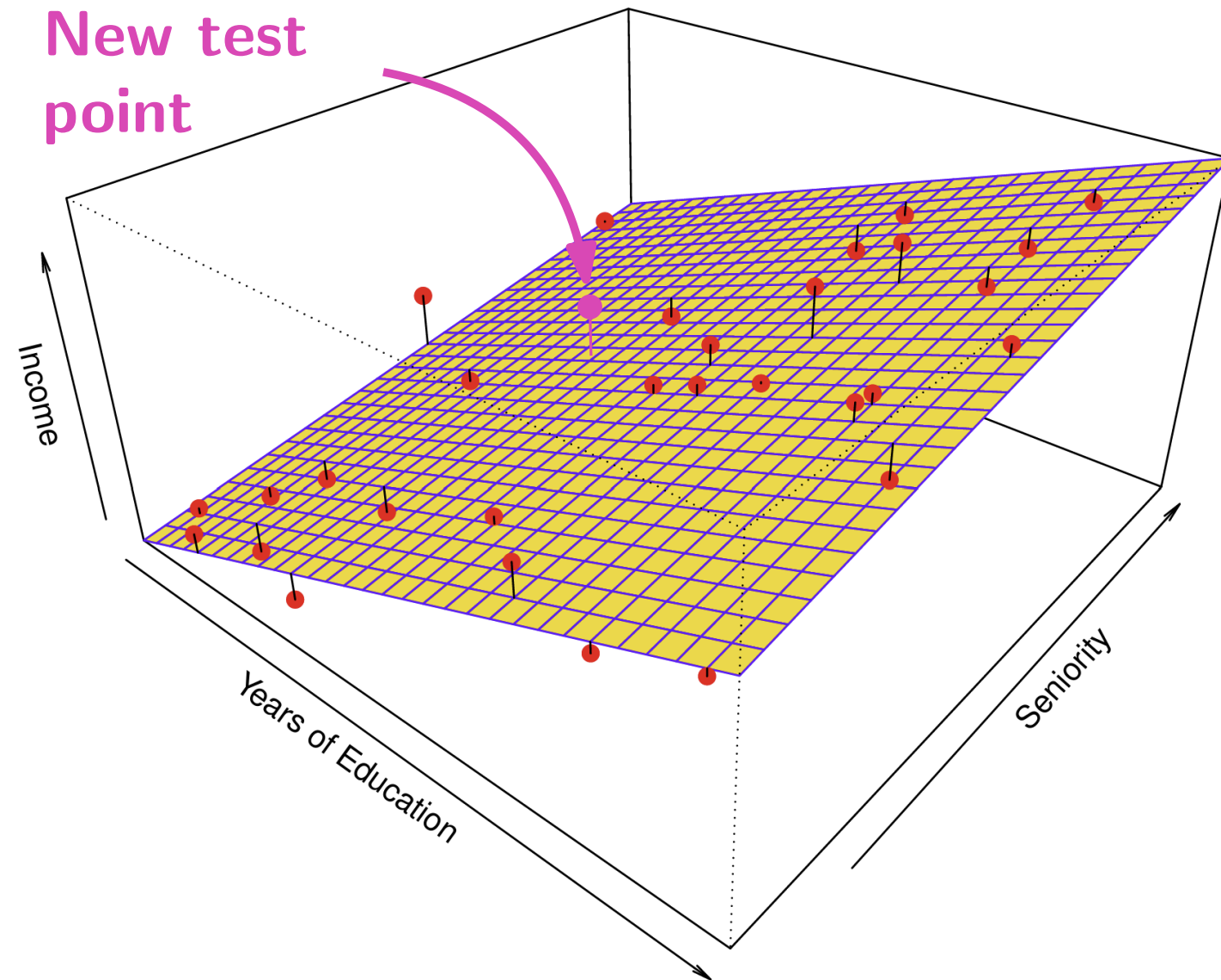
# A general regression example



“Parametric model”: look for  $\beta_0, \beta_1, \beta_2$  s.t.:  
“Income” =  $\beta_0 + \beta_1$  “y. o. education” +  $\beta_2$  “seniority”



# A general regression example

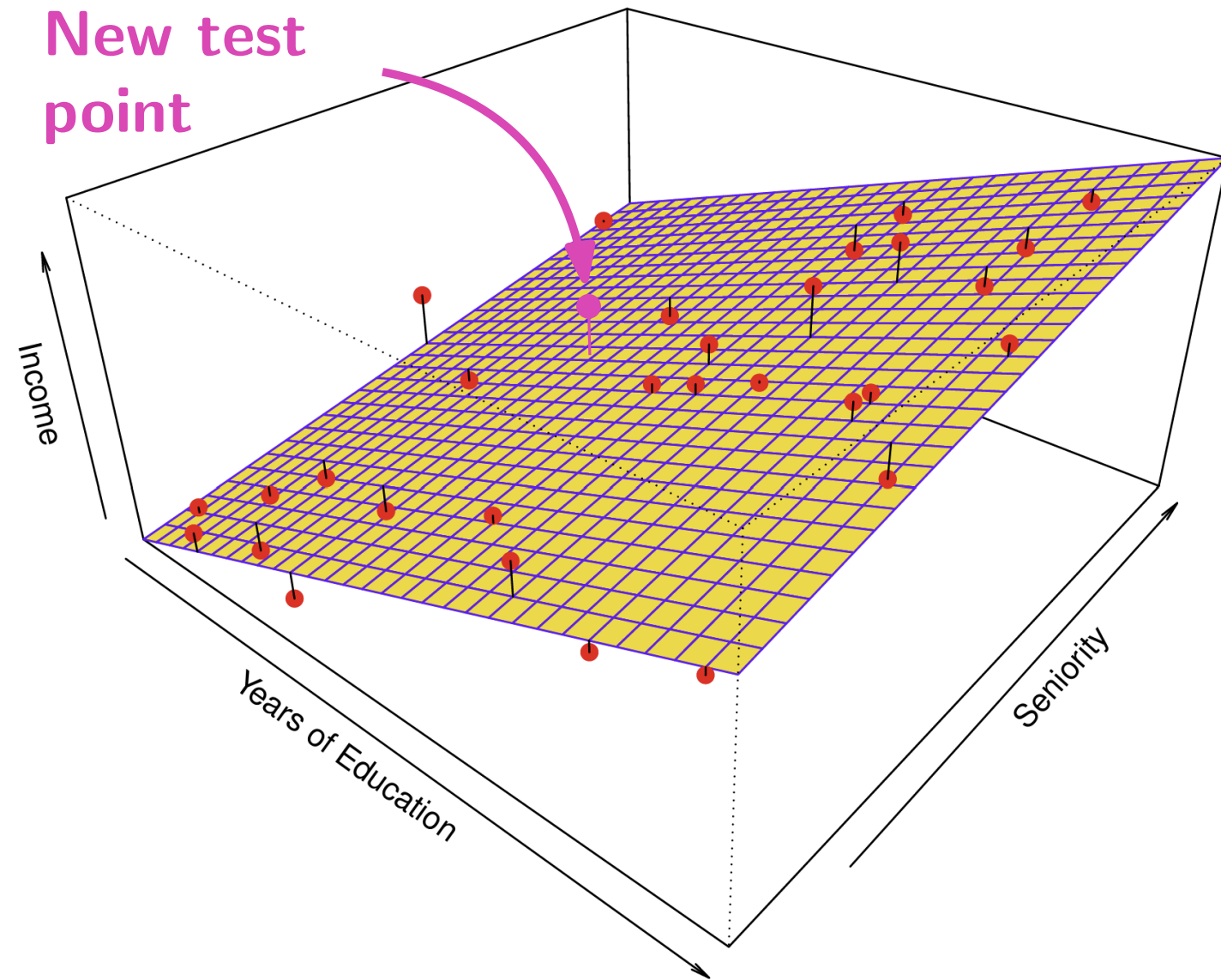


$$X = (x_1, \dots, x_n) \in \mathbb{R}^{p \times n}: \text{Input}$$

“Parametric model”: look for  $\beta_0, \beta_1, \beta_2$  s.t.:  
“Income” =  $\beta_0 + \beta_1$  “y. o. education” +  $\beta_2$  “seniority”



# A general regression example



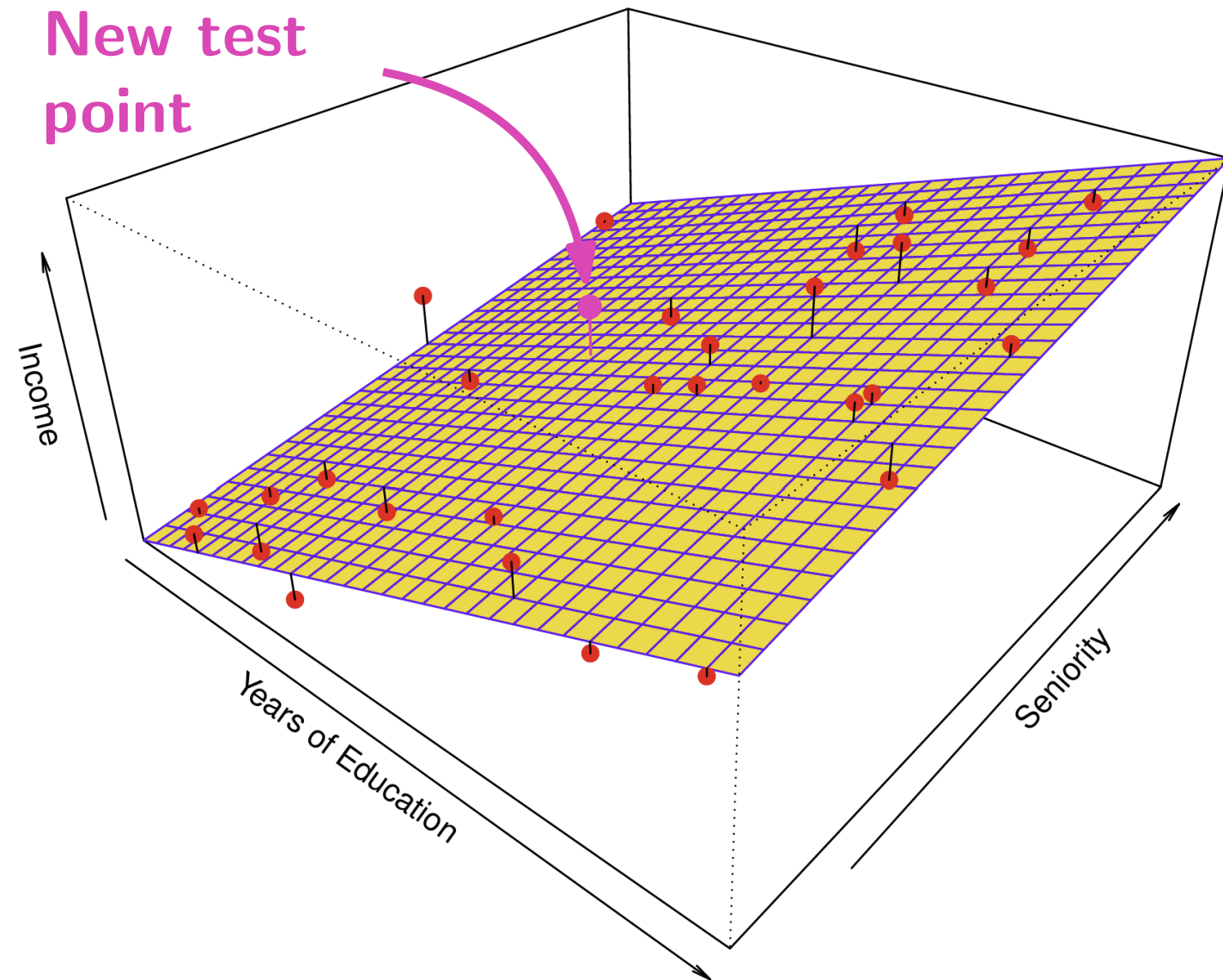
$X = (x_1, \dots, x_n) \in \mathbb{R}^{p \times n}$ : Input

$Y = (y_1, \dots, y_n) \in \mathbb{R}^n$ : Output

“Parametric model”: look for  $\beta_0, \beta_1, \beta_2$  s.t.:

“Income” =  $\beta_0 + \beta_1$  “y. o. education” +  $\beta_2$  “seniority”

# A general regression example



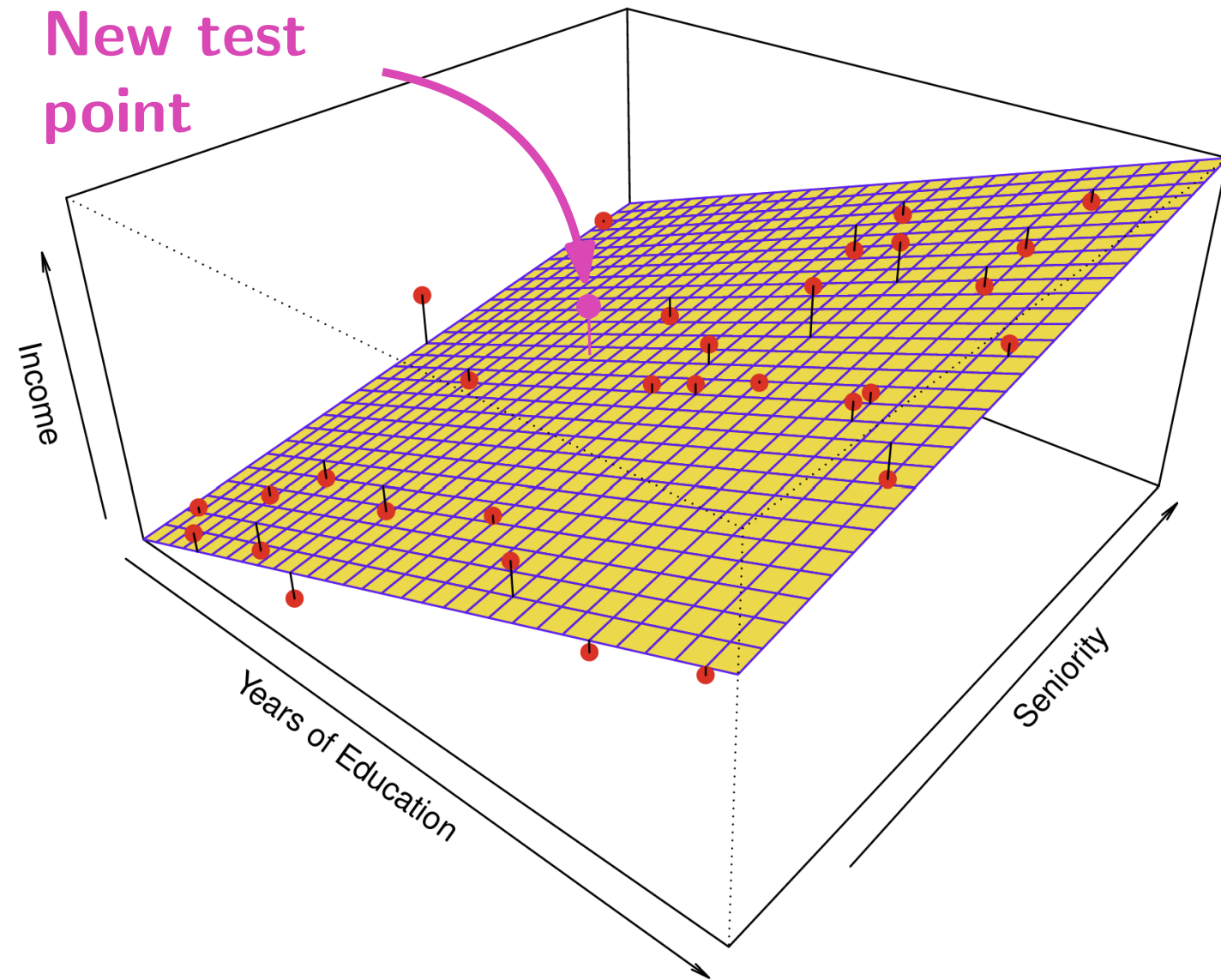
$X = (x_1, \dots, x_n) \in \mathbb{R}^{p \times n}$ : Input

$Y = (y_1, \dots, y_n) \in \mathbb{R}^n$ : Output

$\beta$  minimizes  $\|\beta X - Y\|^2$

“Parametric model”: look for  $\beta_0, \beta_1, \beta_2$  s.t.:  
“Income” =  $\beta_0 + \beta_1$  “y. o. education” +  $\beta_2$  “seniority”

# A general regression example



$X = (x_1, \dots, x_n) \in \mathbb{R}^{p \times n}$ : Input

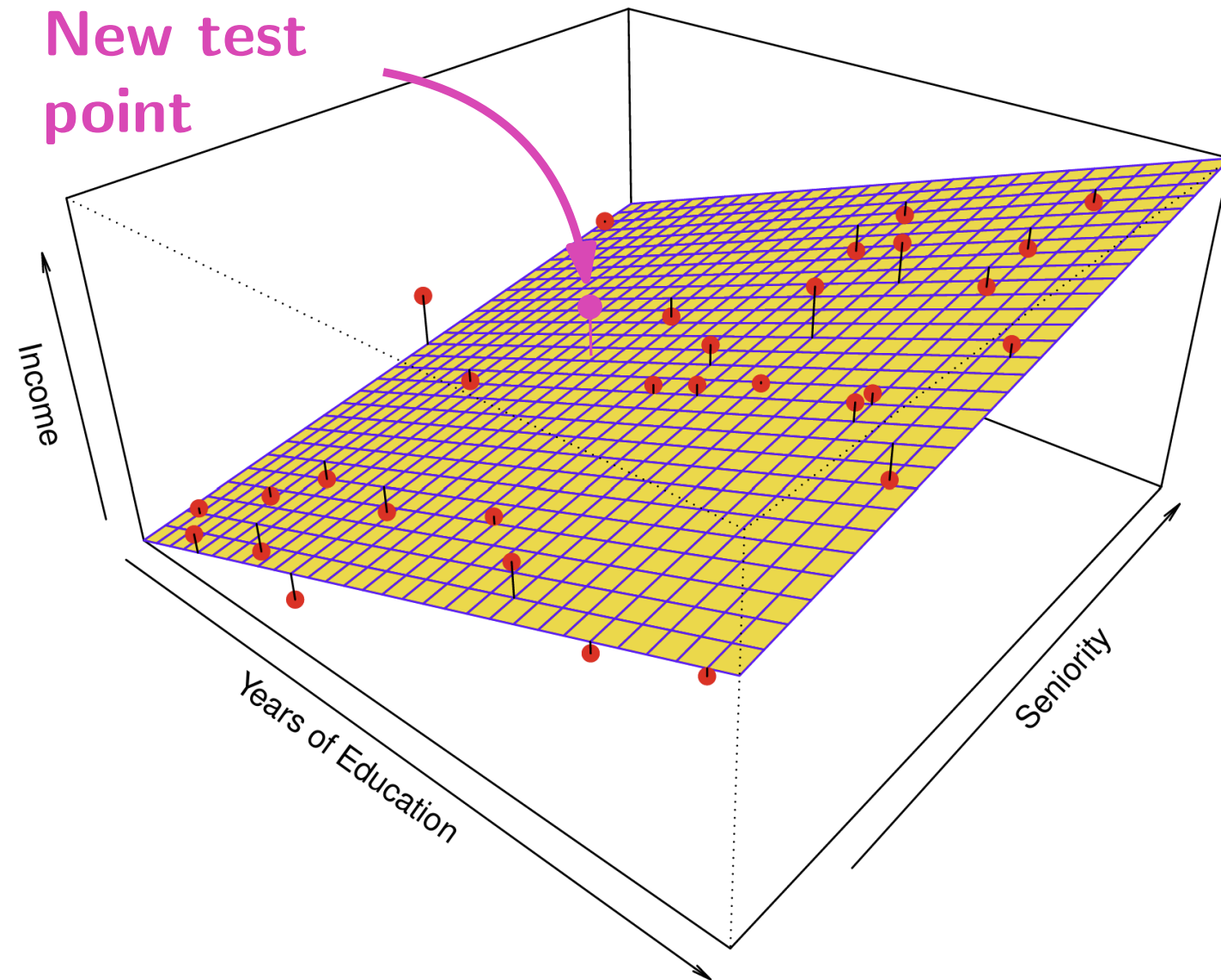
$Y = (y_1, \dots, y_n) \in \mathbb{R}^n$ : Output

$\beta$  minimizes  $\|\beta X - Y\|^2$

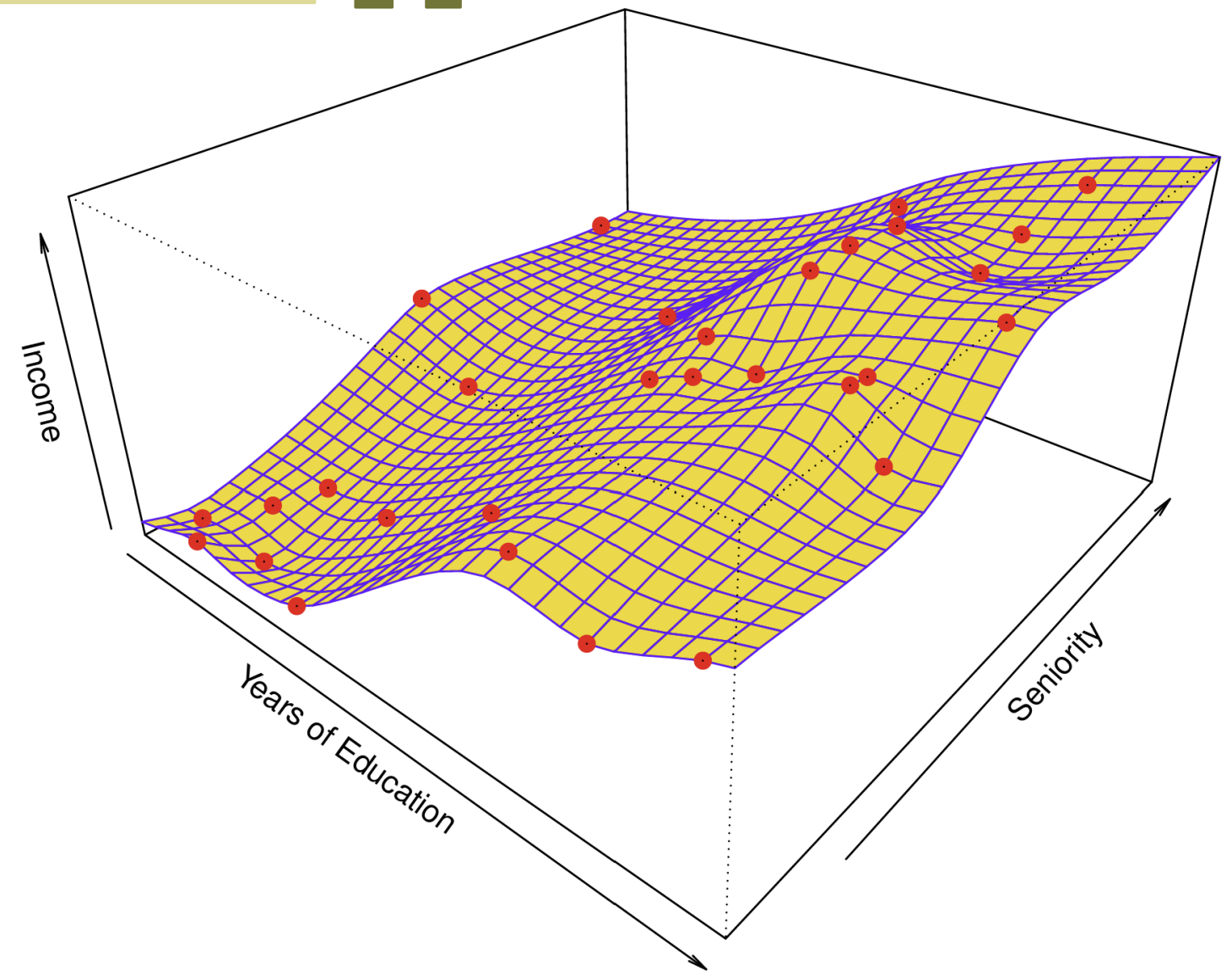
Solution:  $\beta = (XX^T)^{-1}XY \in \mathbb{R}^p$

“Parametric model”: look for  $\beta_0, \beta_1, \beta_2$  s.t.:  
“Income” =  $\beta_0 + \beta_1$  “y. o. education” +  $\beta_2$  “seniority”

# A general regression example



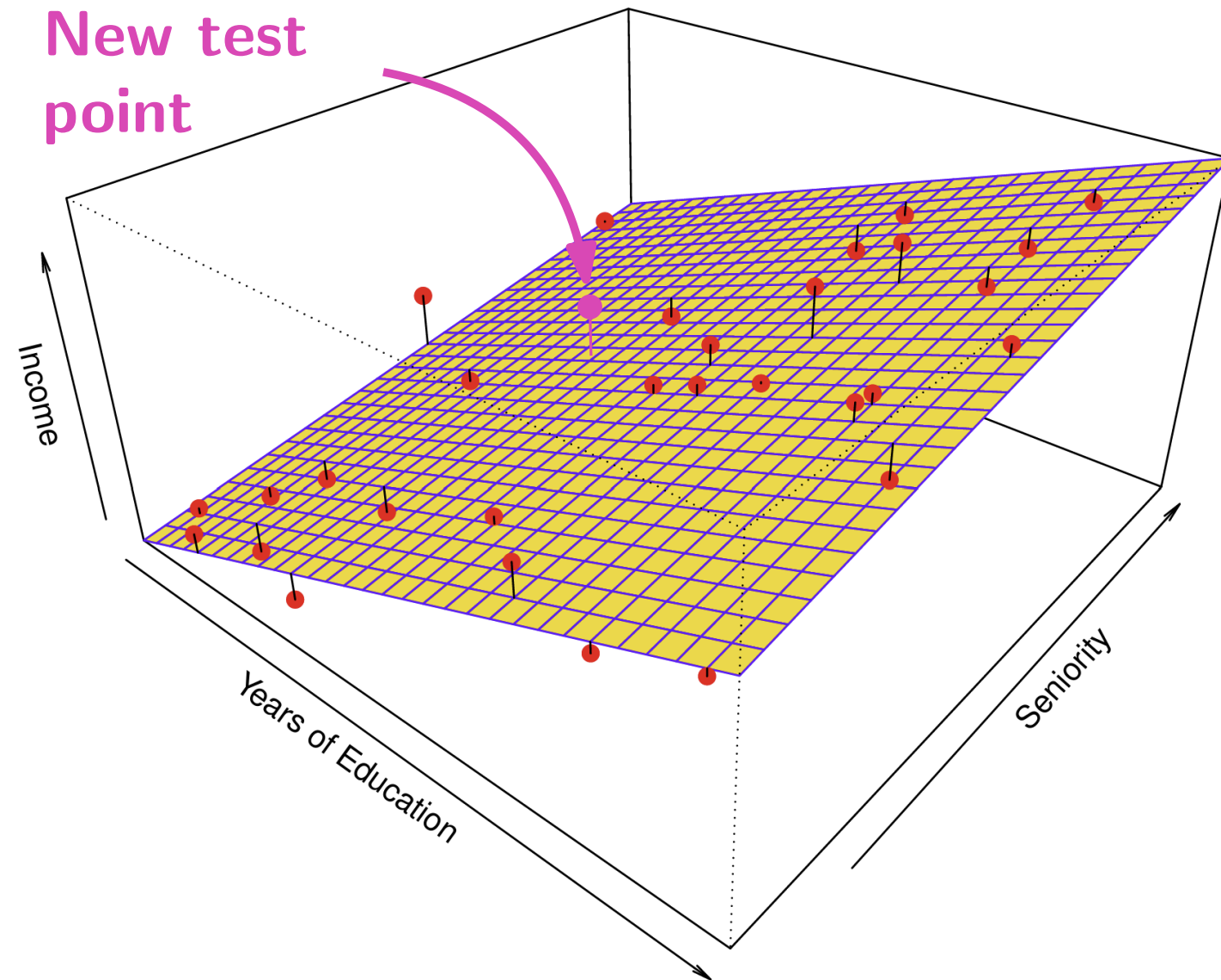
“Parametric model”: look for  $\beta_0, \beta_1, \beta_2$  s.t.:  
“Income” =  $\beta_0 + \beta_1$  “y. o. education” +  $\beta_2$  “seniority”



Best polynomial fit to the data  $Y = P(X)$

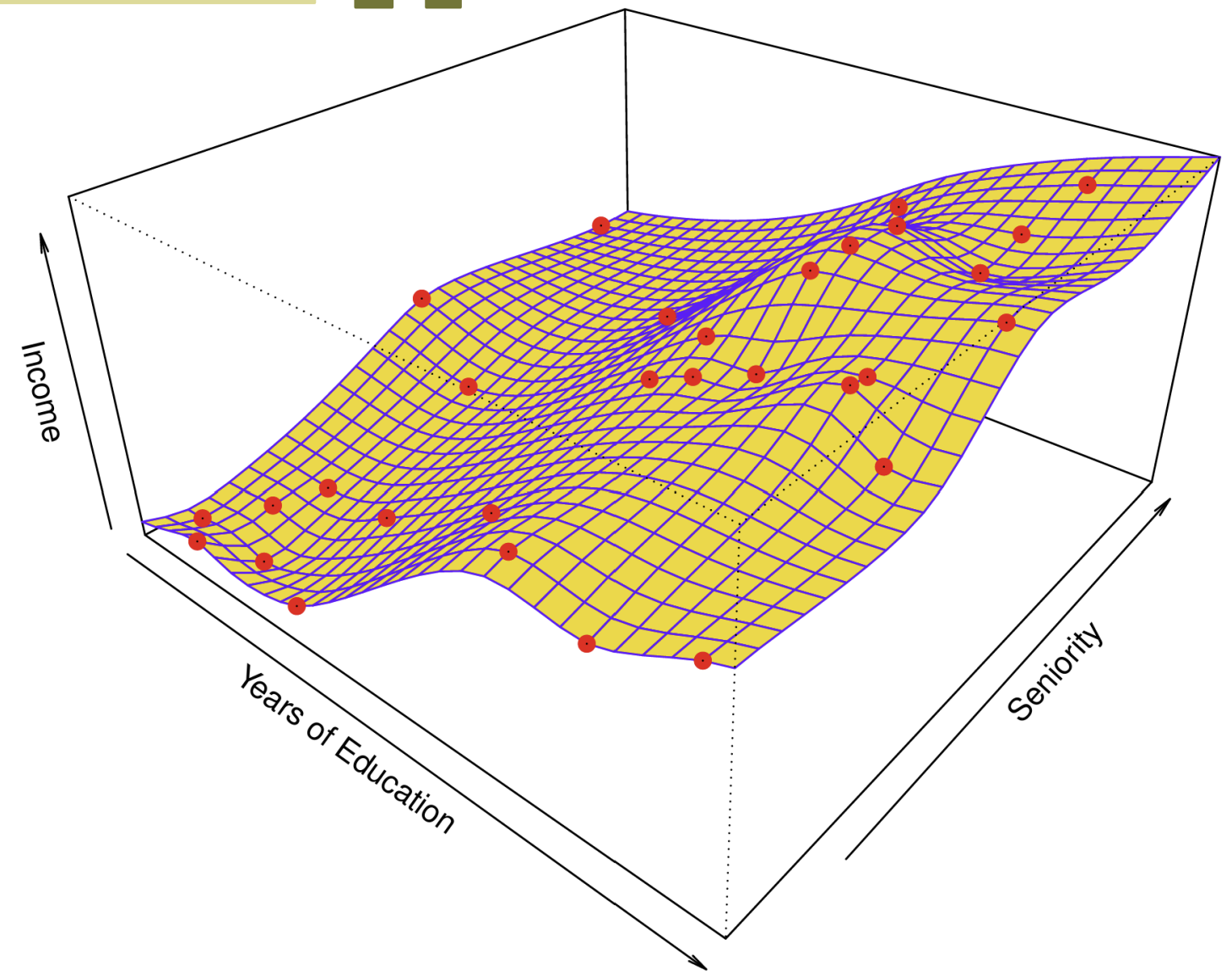


# A general regression example



“Parametric model”: look for  $\beta_0, \beta_1, \beta_2$  s.t.:  
“Income” =  $\beta_0 + \beta_1$  “y. o. education” +  $\beta_2$  “seniority”

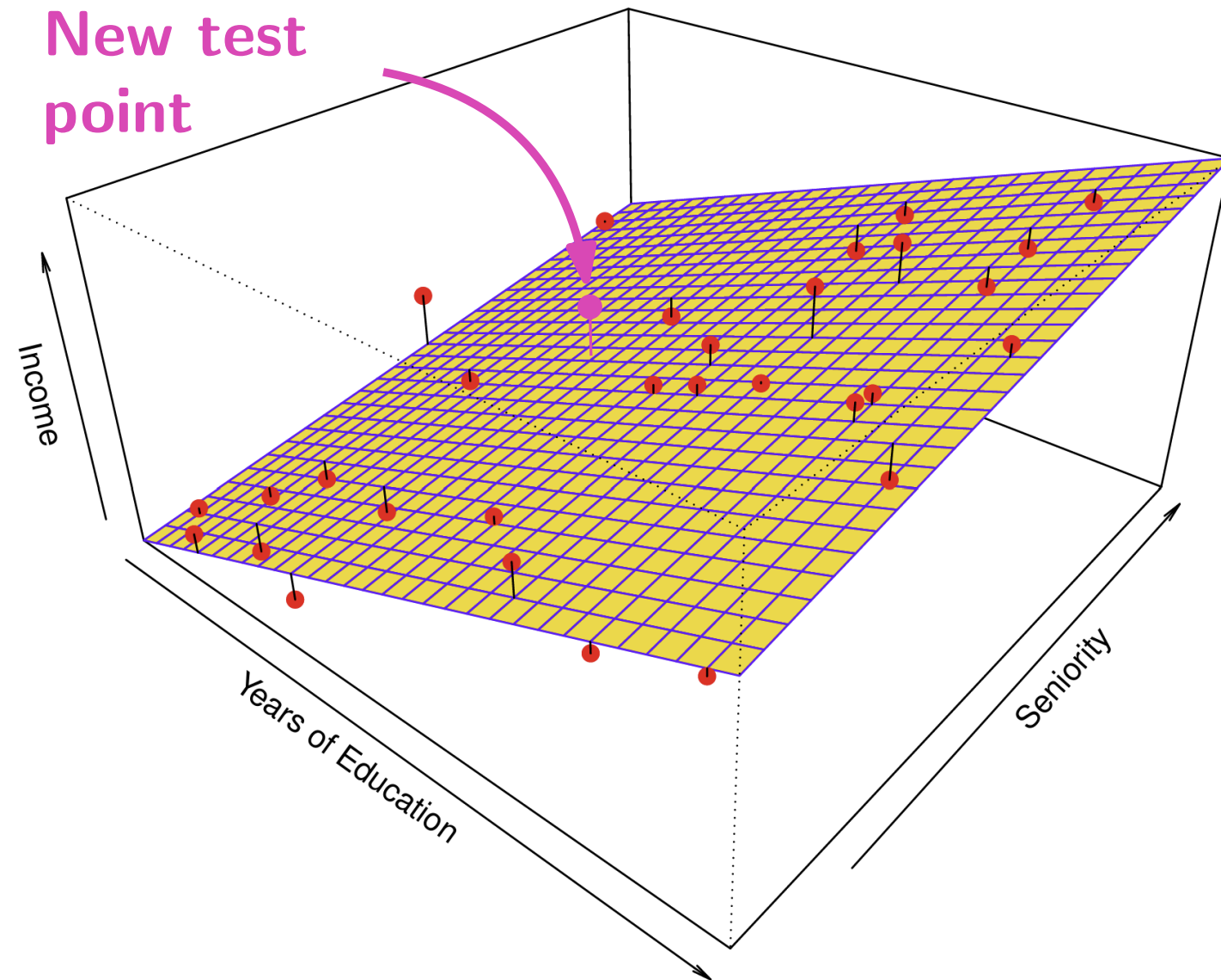
INTERPRETABLE



Best polynomial fit to the data  $Y = P(X)$

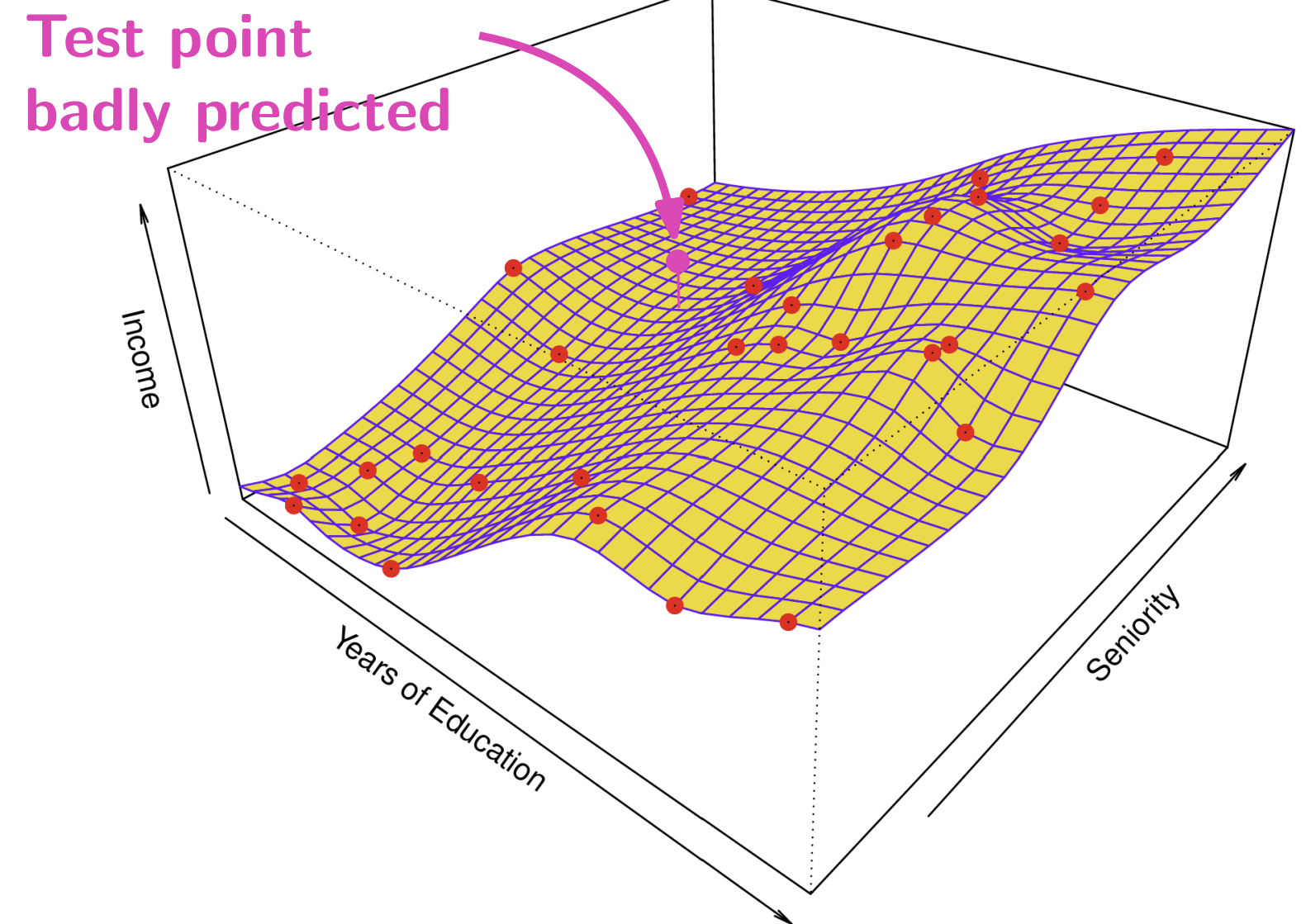
FLEXIBLE

# A general regression example



“Parametric model”: look for  $\beta_0, \beta_1, \beta_2$  s.t.:  
“Income” =  $\beta_0 + \beta_1$  “y. o. education” +  $\beta_2$  “seniority”

INTERPRETABLE



Best polynomial fit to the data  $Y = P(X)$

FLEXIBLE

# A non parametric method instead of regression: KNN

**$k$ -nearest neighbors** method:

Supervised method:

$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$  given



# A non parametric method instead of regression: KNN

**$k$ -nearest neighbors** method:

Supervised method:

$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$  given

New data test  $x \in \mathbb{R}^p$ :

- Look for  $i_1, \dots, i_k$  s.t.:  $x_{i_1}, \dots, x_{i_k}$  are the  $k$  closest neighbors of  $x$ :

$$\forall i \in [n] : \|x - x_i\| \geq \sup_{j \in [k]} \|x - x_{i_j}\|.$$

- $\hat{f}_D(x) \equiv \frac{1}{k} \sum_{j=1}^k y_{i_j}.$

# A non parametric method instead of regression: KNN

**$k$ -nearest neighbors method:**

Supervised method:

$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$  given

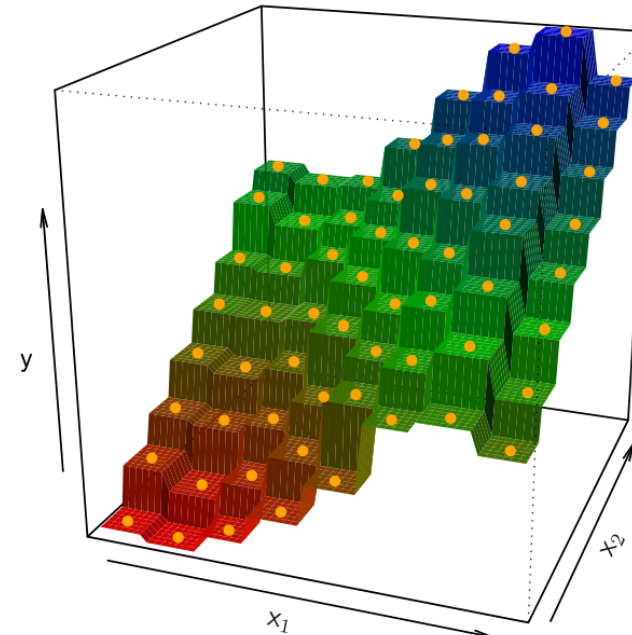
New data test  $x \in \mathbb{R}^p$ :

- Look for  $i_1, \dots, i_k$  s.t.:  $x_{i_1}, \dots, x_{i_k}$  are the  $k$  closest neighbors of  $x$ :

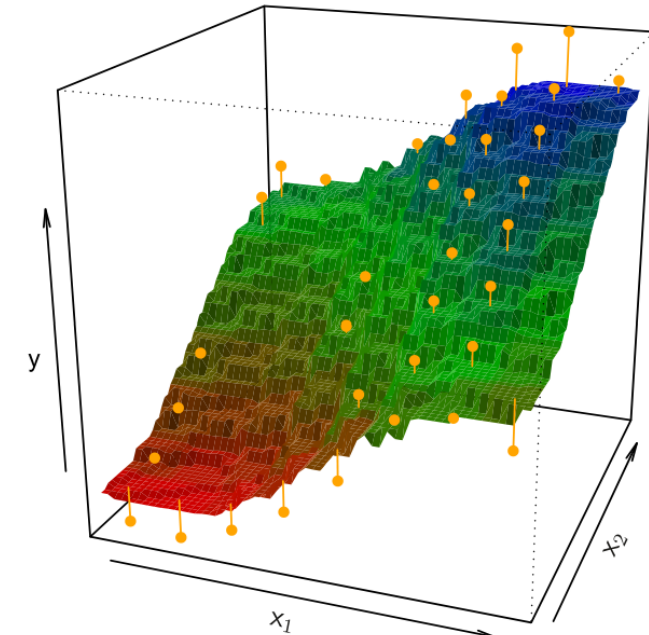
$$\forall i \in [n] : \|x - x_i\| \geq \sup_{j \in [k]} \|x - x_{i_j}\|.$$

$$\bullet \hat{f}_D(x) \equiv \frac{1}{k} \sum_{j=1}^k y_{i_j}.$$

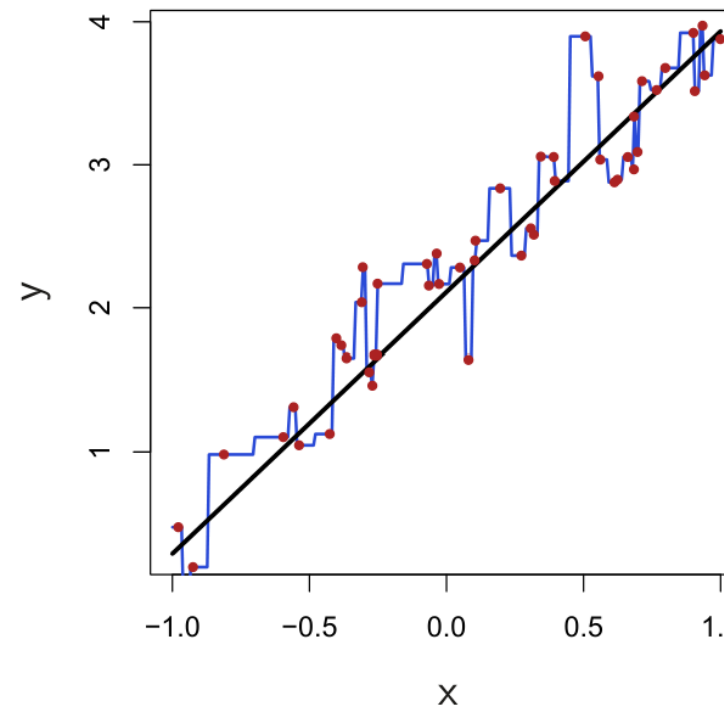
Method highly relies on the number of neighbors chosen



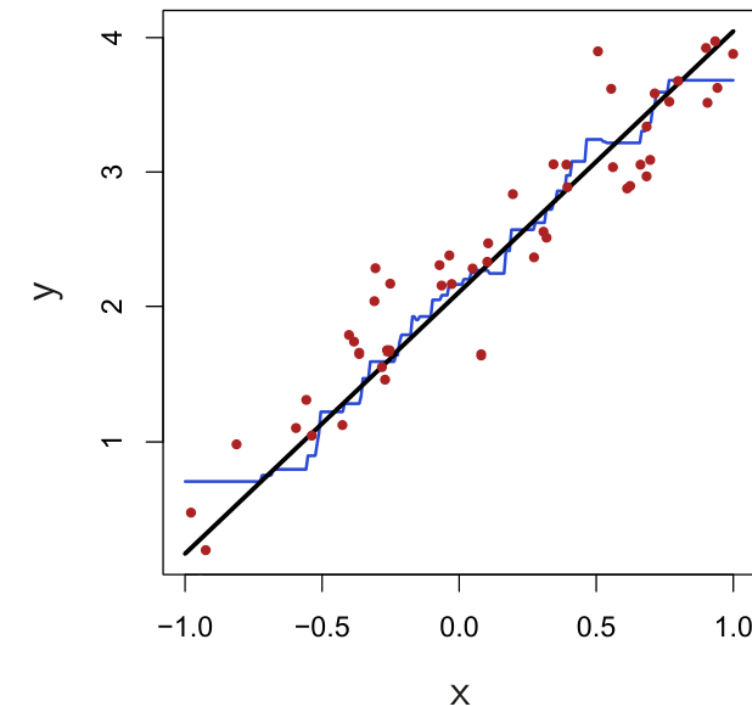
$k = 1$



$k = 9$



$k = 1$



$k = 9$

# A non parametric method instead of regression: KNN

Training dataset:  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$  given

Least squared regresion (**LSR**):

$$\hat{f}_D(x) \equiv \hat{\beta}_0 x^{(1)} + \dots + \hat{\beta}_p x^{(p)}$$

where  $\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n \|\beta^T x - y\|^2.$

# A non parametric method instead of regression: KNN

Training dataset:  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$  given

Least squared regresion (**LSR**):

$$\hat{f}_D(x) \equiv \hat{\beta}_0 x^{(1)} + \dots + \hat{\beta}_p x^{(p)}$$

where  $\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n \|\beta^T x - y\|^2.$

k-nearest neighbors method (**KNN**):

$x_{i_1}, \dots, x_{i_k}$  k nearest neighbors of  $x$  and:

$$\hat{f}_D(x) \equiv \frac{1}{k} \sum_{j=1}^k y_{i_j}.$$

# A non parametric method instead of regression: KNN

Training dataset:  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$  given

Least squared regresion (**LSR**):

$$\hat{f}_D(x) \equiv \hat{\beta}_0 x^{(1)} + \dots + \hat{\beta}_p x^{(p)}$$

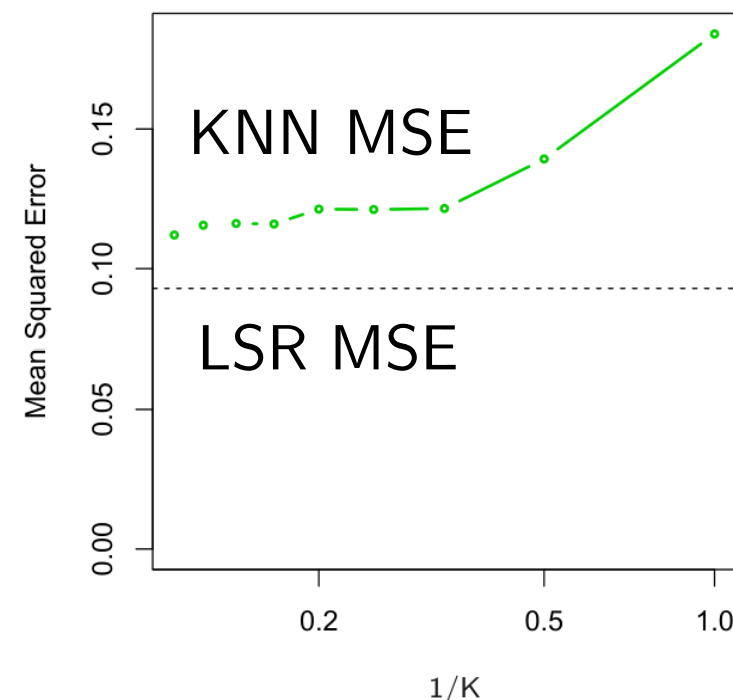
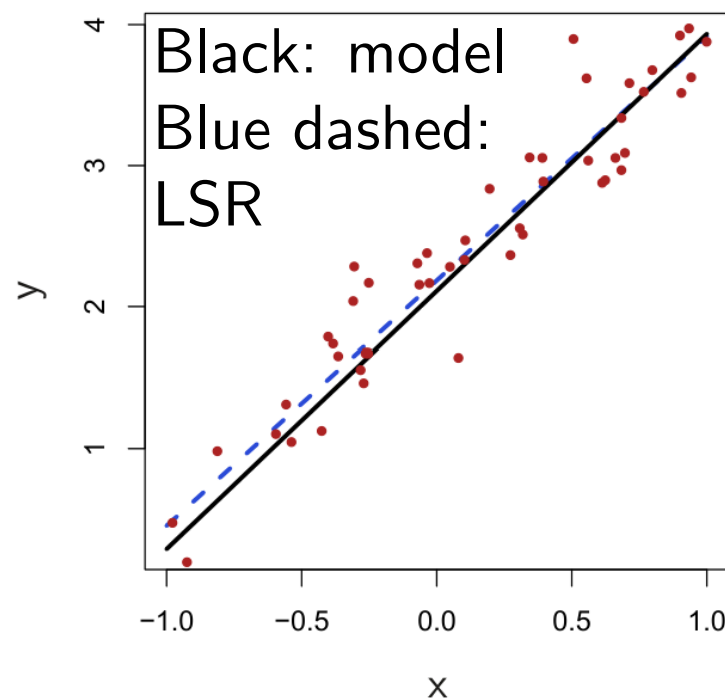
where  $\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n \|\beta^T x - y\|^2$ .

k-nearest neighbors method (**KNN**):

$x_{i_1}, \dots, x_{i_k}$  k nearest neighbors of  $x$  and:

$$\hat{f}_D(x) \equiv \frac{1}{k} \sum_{j=1}^k y_{i_j}.$$

- With linear model, **LSR** better



# A non parametric method instead of regression: KNN

Training dataset:  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$  given

Least squared regresion (**LSR**):

$$\hat{f}_D(x) \equiv \hat{\beta}_0 x^{(1)} + \dots + \hat{\beta}_p x^{(p)}$$

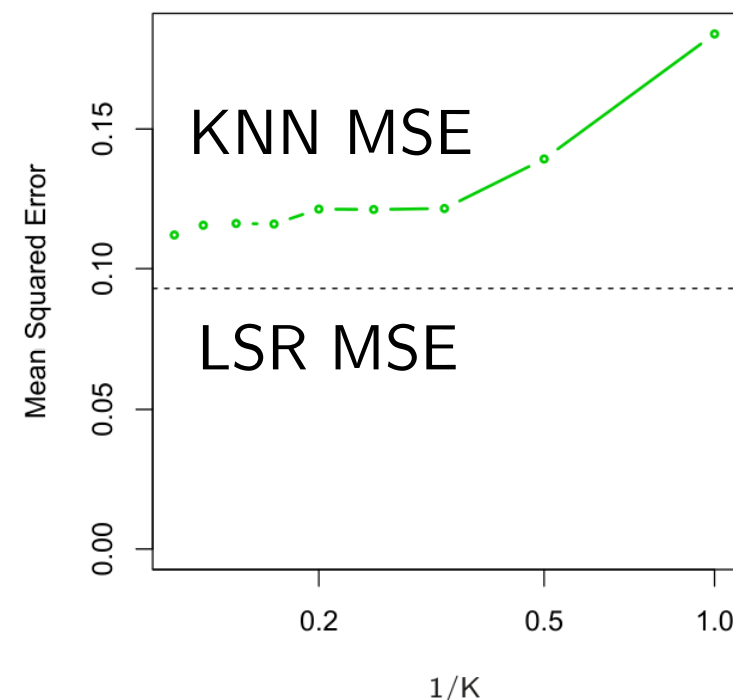
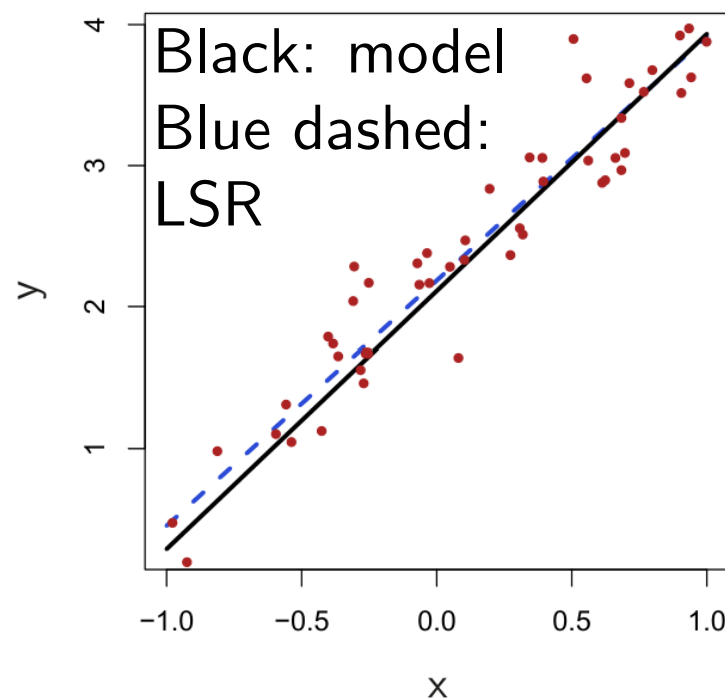
where  $\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n \|\beta^T x - y\|^2$ .

k-nearest neighbors method (**KNN**):

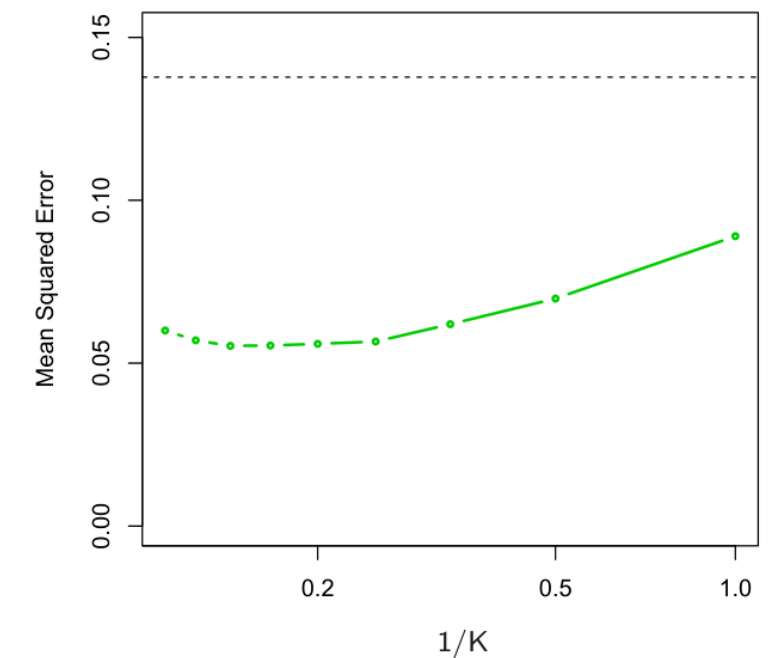
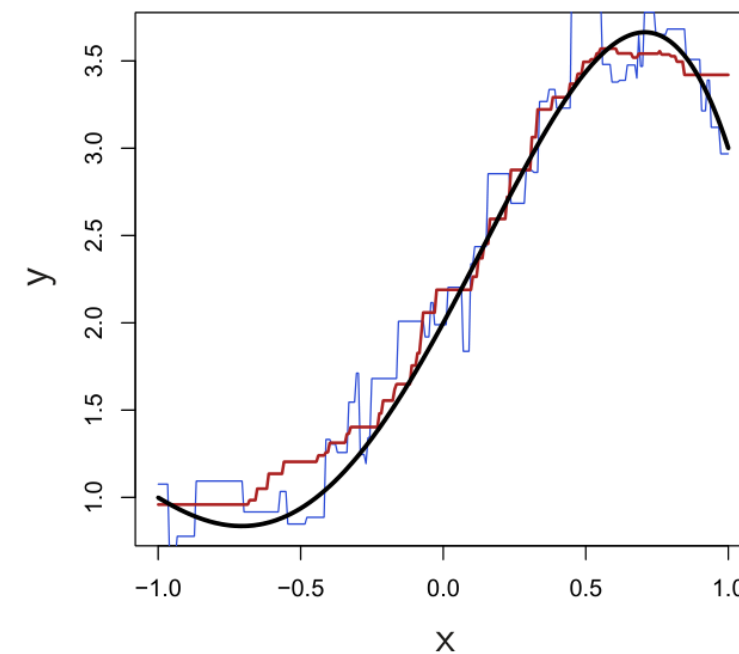
$x_{i_1}, \dots, x_{i_k}$  k nearest neighbors of  $x$  and:

$$\hat{f}_D(x) \equiv \frac{1}{k} \sum_{j=1}^k y_{i_j}.$$

- With linear model, **LSR** better



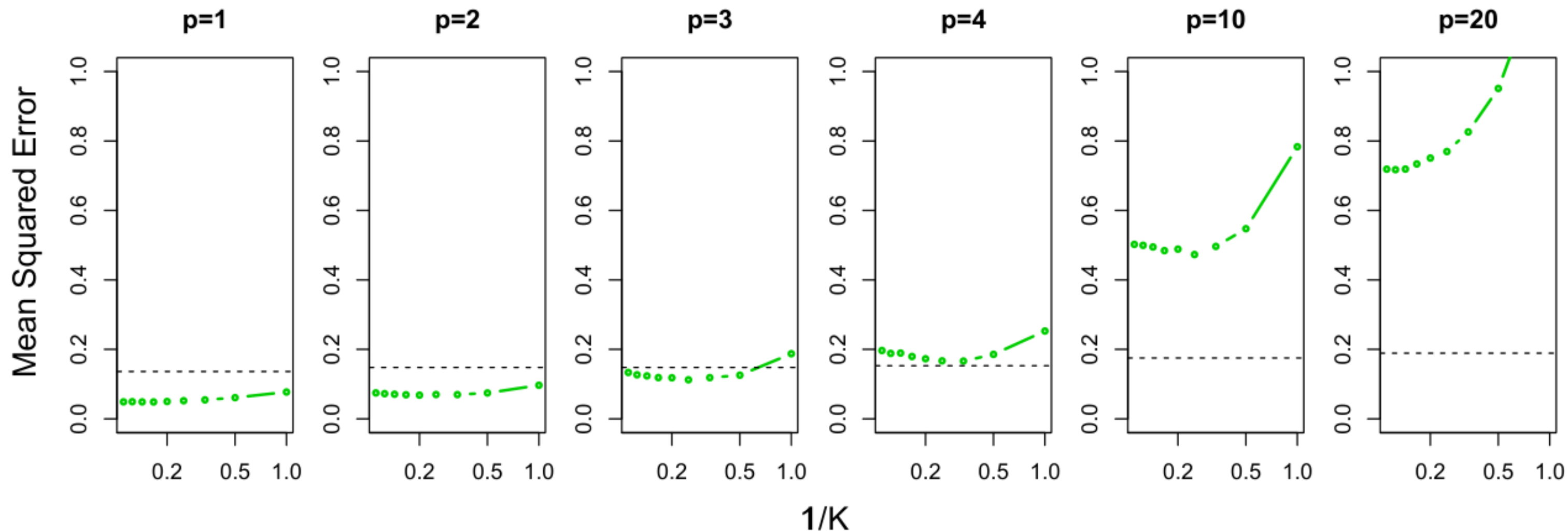
- More complex model of data: **KNN** better



# A non parametric method instead of regression: KNN

When dimension grows, with fixed amount of data in training data set ( $n$  constant)  
KNN highly dependent on the variance of the data: lack of interpretability.

(Here model of data non-linear because otherwise least squares regression would be always performing better)





# Maximum likelihood (interpretable parametric method)

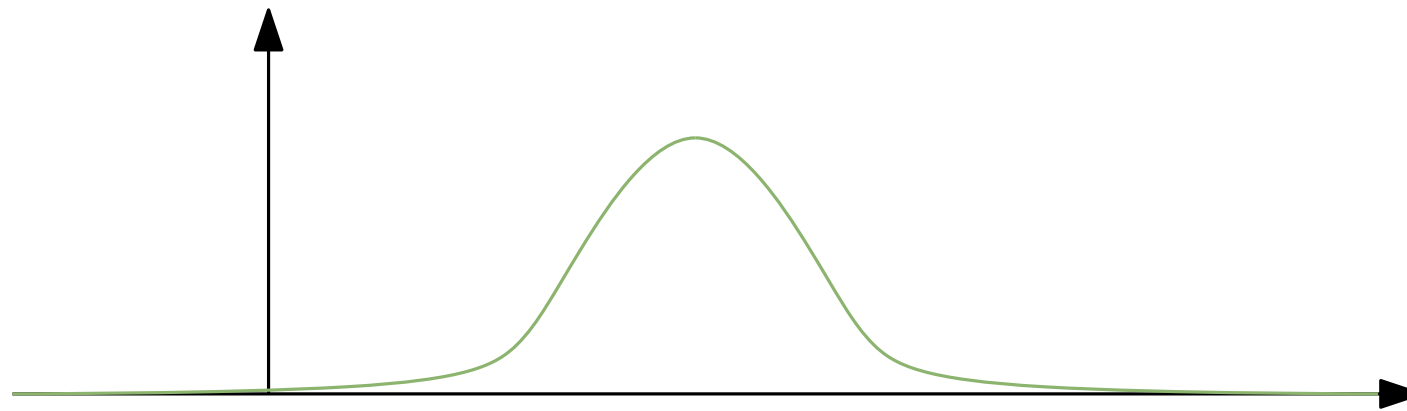
We know the model of  $X$ , we want to “fit the model”

# Maximum likelihood (interpretable parametric method)

We know the model of  $X$ , we want to “fit the model”

Find the parameters defining the model

Ex:  $X$  follows a Gaussian distribution

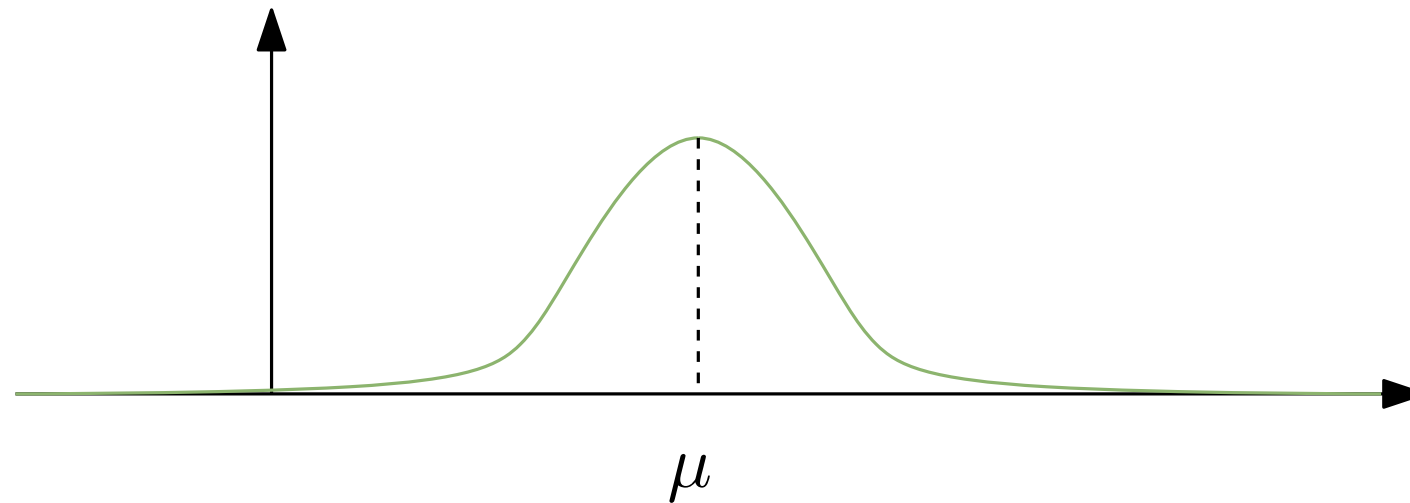


# Maximum likelihood (interpretable parametric method)

We know the model of  $X$ , we want to “fit the model”

Find the parameters defining the model

Ex:  $X$  follows a Gaussian distribution

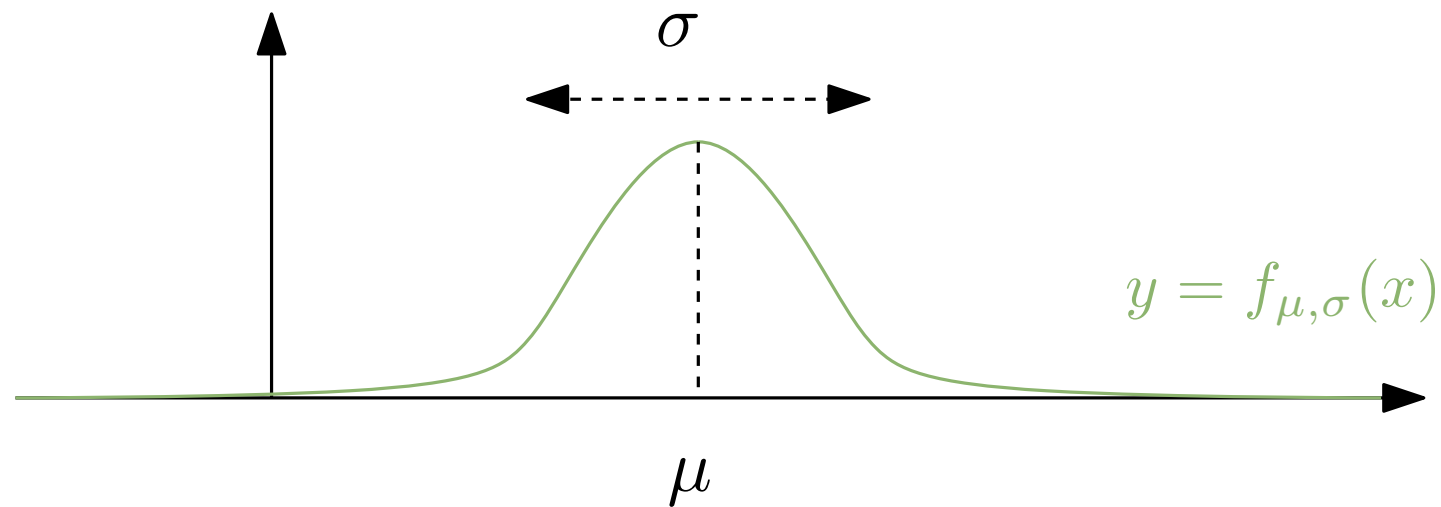


# Maximum likelihood (interpretable parametric method)

We know the model of  $X$ , we want to “fit the model”

Find the parameters defining the model

Ex:  $X$  follows a Gaussian distribution



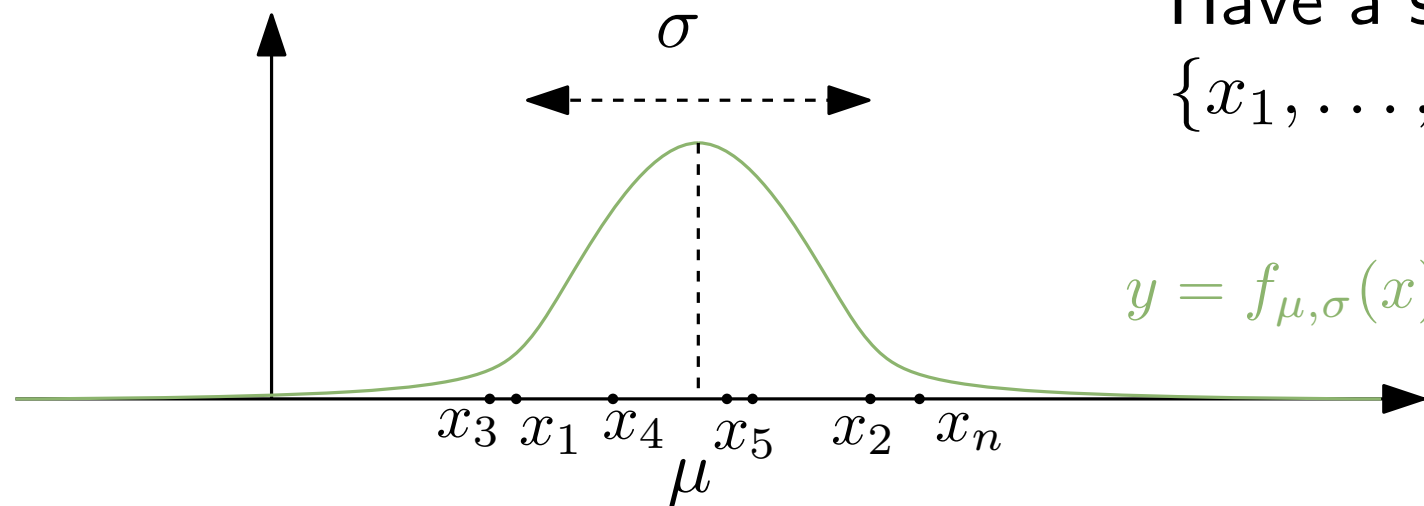
# Maximum likelihood (interpretable parametric method)

We know the model of  $X$ , we want to “fit the model”

Find the parameters defining the model

Ex:  $X$  follows a Gaussian distribution

Have a set of drawings  
 $\{x_1, \dots, x_n\}$



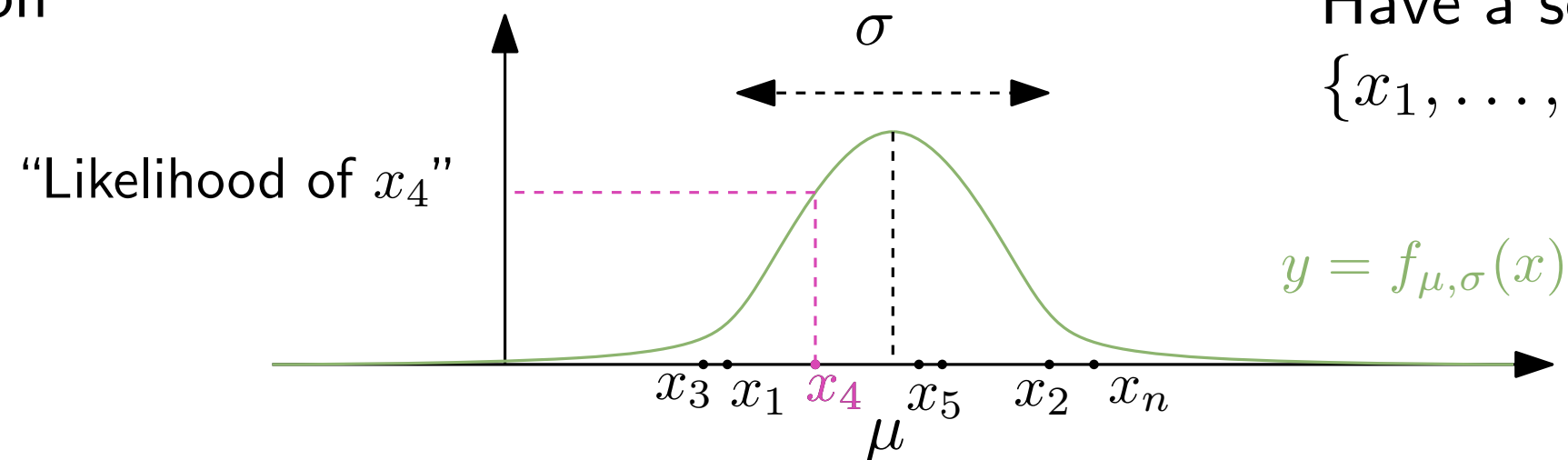
# Maximum likelihood (interpretable parametric method)

We know the model of  $X$ , we want to “fit the model”

Find the parameters defining the model

Ex:  $X$  follows a Gaussian distribution

Have a set of drawings  
 $\{x_1, \dots, x_n\}$



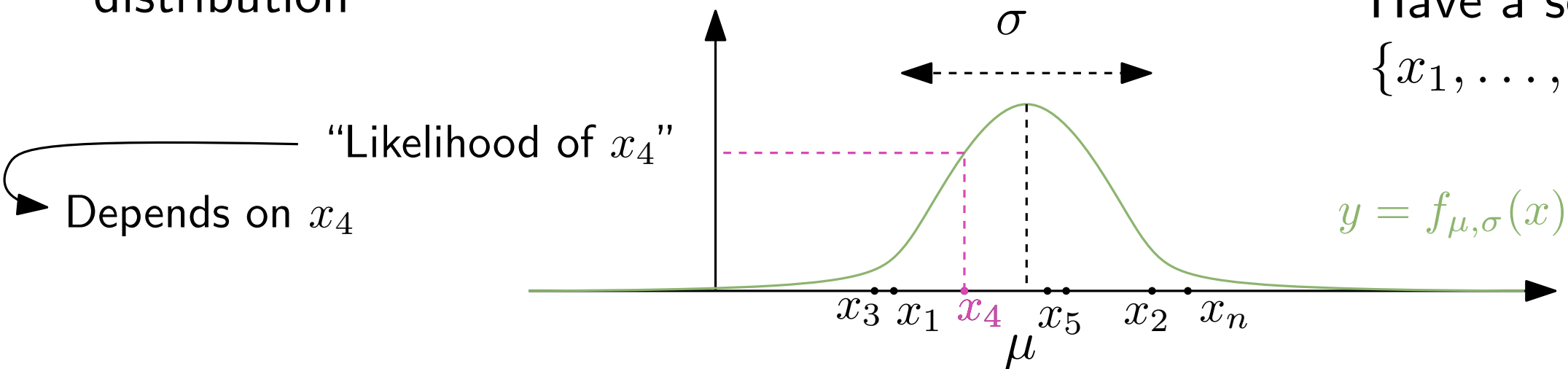
# Maximum likelihood (interpretable parametric method)

We know the model of  $X$ , we want to “fit the model”

Find the parameters defining the model

Ex:  $X$  follows a Gaussian distribution

Have a set of drawings  
 $\{x_1, \dots, x_n\}$





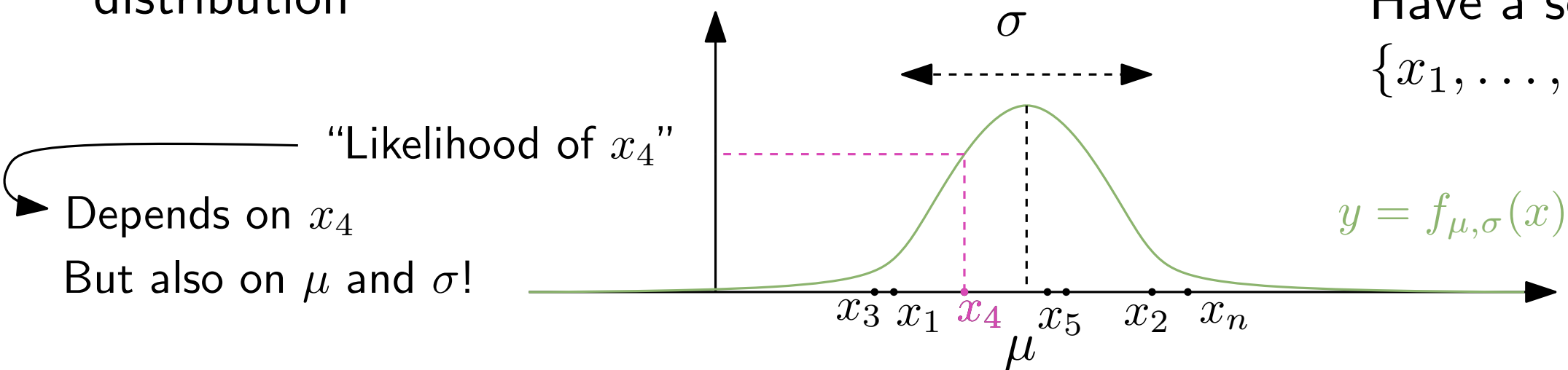
# Maximum likelihood (interpretable parametric method)

We know the model of  $X$ , we want to “fit the model”

Find the parameters defining the model

Ex:  $X$  follows a Gaussian distribution

Have a set of drawings  
 $\{x_1, \dots, x_n\}$



Maximize:  $\prod_{i=1}^n f_{\mu, \sigma}(x_i), \quad \mu, \sigma > 0$

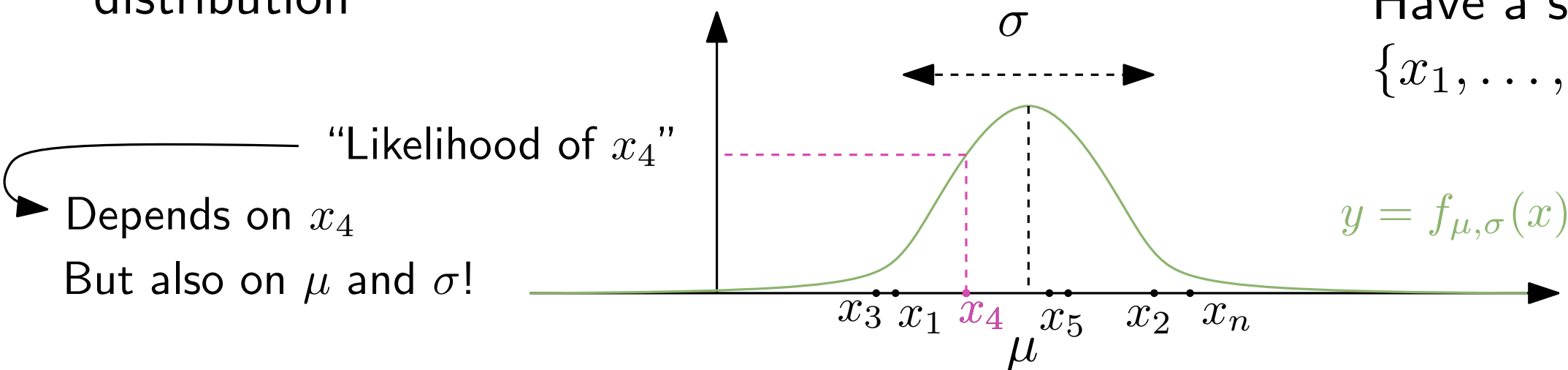
# Maximum likelihood (interpretable parametric method)

We know the model of  $X$ , we want to “fit the model”

Find the parameters defining the model

Ex:  $X$  follows a Gaussian distribution

Have a set of drawings  
 $\{x_1, \dots, x_n\}$



$$\text{Maximize: } \prod_{i=1}^n f_{\mu, \sigma}(x_i), \quad \mu, \sigma > 0 \quad \Longleftrightarrow \quad \text{Maximize: } \sum_{i=1}^n \log(f_{\mu, \sigma}(x_i)) \quad \mu, \sigma > 0$$

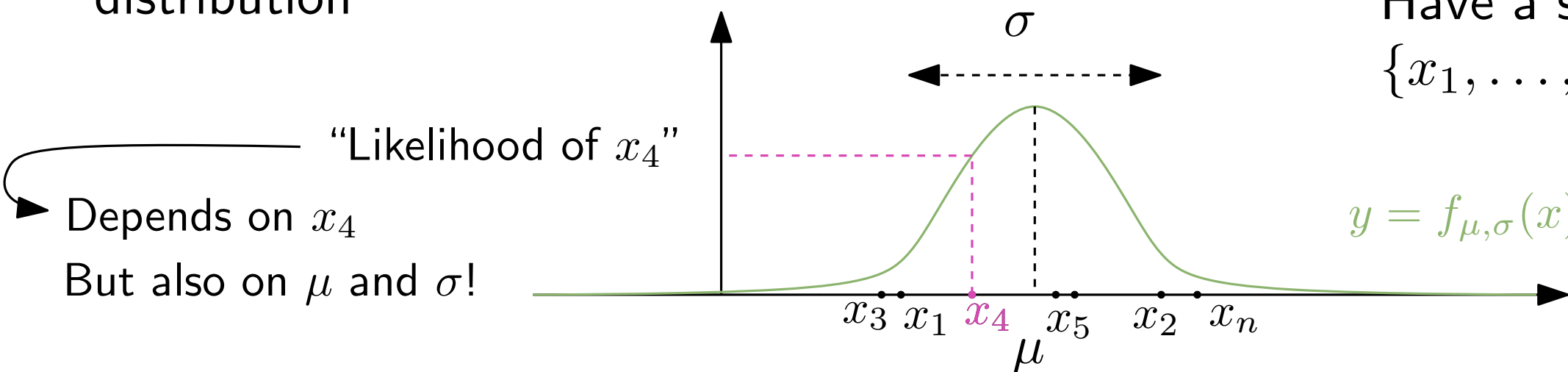
# Maximum likelihood (interpretable parametric method)

We know the model of  $X$ , we want to “fit the model”

Find the parameters defining the model

Ex:  $X$  follows a Gaussian distribution

Have a set of drawings  
 $\{x_1, \dots, x_n\}$



“loglikelihood minimization”

$$\text{Maximize: } \prod_{i=1}^n f_{\mu, \sigma}(x_i), \quad \mu, \sigma > 0 \quad \Longleftrightarrow \quad \text{Maximize: } \sum_{i=1}^n \log(f_{\mu, \sigma}(x_i)) \quad \mu, \sigma > 0$$

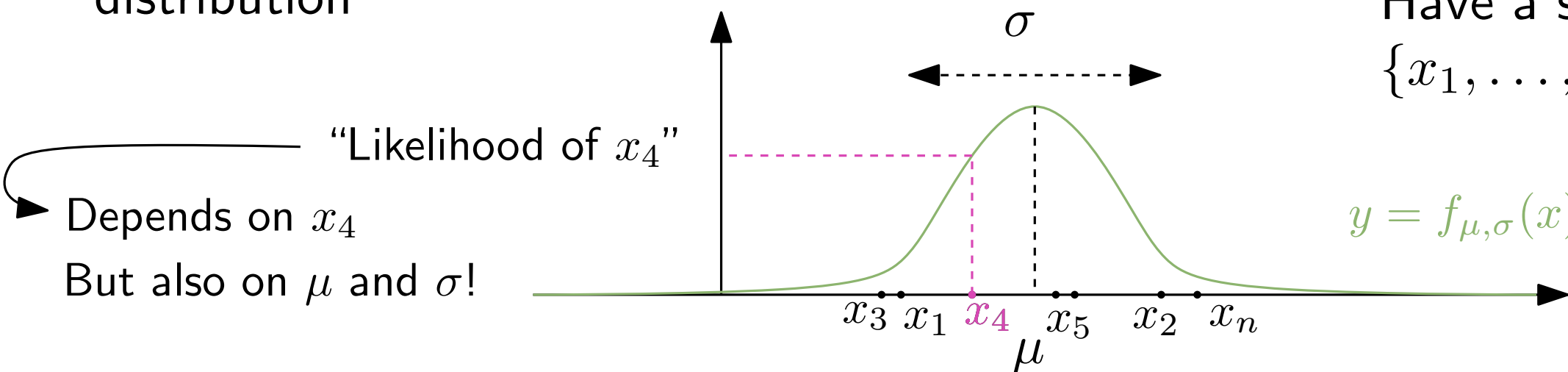
# Maximum likelihood (interpretable parametric method)

We know the model of  $X$ , we want to “fit the model”

Find the parameters defining the model

Ex:  $X$  follows a Gaussian distribution

Have a set of drawings  
 $\{x_1, \dots, x_n\}$



“loglikelihood minimization”

$$\text{Maximize: } \prod_{i=1}^n f_{\mu, \sigma}(x_i), \quad \mu, \sigma > 0 \quad \Longleftrightarrow \quad \text{Maximize: } \sum_{i=1}^n \log(f_{\mu, \sigma}(x_i)) \quad \mu, \sigma > 0$$

# Gaussian likelihood - Ordinary least square regression (LSR).

---

# Gaussian likelihood - Ordinary least square regression (LSR).

Model:  $Y = \beta^T X + \varepsilon$ , with:

$Y \in \mathbb{R}$ : Target

$X \in \mathbb{R}^p$ : data

$\varepsilon \in \mathbb{R}$ ,  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ : noise

$\beta \in \mathbb{R}^p$ : regression vector

# Gaussian likelihood - Ordinary least square regression (LSR).

Model:  $Y = \beta^T X + \varepsilon$ , with:

$Y \in \mathbb{R}$ : Target

$X \in \mathbb{R}^p$ : data

$\varepsilon \in \mathbb{R}$ ,  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ : noise

$\beta \in \mathbb{R}^p$ : regression vector

Given data set  $(x_1, y_1), \dots, (x_n, y_n)$ ,  
→ Estimate  $\beta$





# Gaussian likelihood - Ordinary least square regression (LSR).

Model:  $Y = \beta^T X + \varepsilon$ , with:

$Y \in \mathbb{R}$ : Target

$X \in \mathbb{R}^p$ : data

$\varepsilon \in \mathbb{R}, \varepsilon \sim \mathcal{N}(0, \sigma^2)$ : noise

$\beta \in \mathbb{R}^p$ : regression vector

Given data set  $(x_1, y_1), \dots, (x_n, y_n)$ ,

→ Estimate  $\beta$

$$y_i \sim \mathcal{N}(\beta^T x_i, \sigma^2) \text{ thus } f_\beta(y_i) = \frac{\exp(-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2})}{\sqrt{2\pi}\sigma}.$$

# Gaussian likelihood - Ordinary least square regression (LSR).

Model:  $Y = \beta^T X + \varepsilon$ , with:

$Y \in \mathbb{R}$ : Target

$X \in \mathbb{R}^p$ : data

$\varepsilon \in \mathbb{R}, \varepsilon \sim \mathcal{N}(0, \sigma^2)$ : noise

$\beta \in \mathbb{R}^p$ : regression vector

Given data set  $(x_1, y_1), \dots, (x_n, y_n)$ ,

→ Estimate  $\beta$

$y_i \sim \mathcal{N}(\beta^T x_i, \sigma^2)$  thus  $f_\beta(y_i) = \frac{\exp(-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2})}{\sqrt{2\pi}\sigma}$ .

Negative Log-likelihood:

$$L(\beta) = \sum_{i=1}^n \frac{(y_i - \beta^T x_i)^2}{2\sigma^2} + \frac{\log(2\pi)}{2} + \log(\sigma)$$

# Gaussian likelihood - Ordinary least square regression (LSR).

Model:  $Y = \beta^T X + \varepsilon$ , with:

$Y \in \mathbb{R}$ : Target

$X \in \mathbb{R}^p$ : data

$\varepsilon \in \mathbb{R}$ ,  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ : noise

$\beta \in \mathbb{R}^p$ : regression vector

Given data set  $(x_1, y_1), \dots, (x_n, y_n)$ ,

→ Estimate  $\beta$

Negative Log-likelihood:

$$y_i \sim \mathcal{N}(\beta^T x_i, \sigma^2) \text{ thus } f_\beta(y_i) = \frac{\exp(-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2})}{\sqrt{2\pi}\sigma}. \quad L(\beta) = \sum_{i=1}^n \frac{(y_i - \beta^T x_i)^2}{2\sigma^2} + \frac{\log(2\pi)}{2} + \log(\sigma)$$

$$\text{Minimize } L(\beta) \quad \Longleftrightarrow \quad \text{Minimize } \|Y - \beta^T X\|^2$$

With

$$Y = (y_1, \dots, y_n) \in \mathbb{R}^n,$$

$$X = (x_1, \dots, x_n) \in \mathbb{R}^{p \times n}$$

# Gaussian likelihood - Ordinary least square regression (LSR).

Model:  $Y = \beta^T X + \varepsilon$ , with:

$Y \in \mathbb{R}$ : Target

$X \in \mathbb{R}^p$ : data

$\varepsilon \in \mathbb{R}, \varepsilon \sim \mathcal{N}(0, \sigma^2)$ : noise

$\beta \in \mathbb{R}^p$ : regression vector

Given data set  $(x_1, y_1), \dots, (x_n, y_n)$ ,

→ Estimate  $\beta$

Negative Log-likelihood:

$$y_i \sim \mathcal{N}(\beta^T x_i, \sigma^2) \text{ thus } f_\beta(y_i) = \frac{\exp(-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2})}{\sqrt{2\pi}\sigma}. \quad L(\beta) = \sum_{i=1}^n \frac{(y_i - \beta^T x_i)^2}{2\sigma^2} + \frac{\log(2\pi)}{2} + \log(\sigma)$$

$$\text{Minimize } L(\beta) \iff \text{Minimize } \|Y - \beta^T X\|^2$$

With

$$Y = (y_1, \dots, y_n) \in \mathbb{R}^n,$$

$$X = (x_1, \dots, x_n) \in \mathbb{R}^{p \times n}$$

“Ordinary least square regression”

# Poisson likelihood - Poisson Regression.

---

Poisson distribution is used to model number of apperence of a random event in a given interval.

# Poisson likelihood - Poisson Regression.

---

Poisson distribution is used to model number of apperence of a random event in a given interval.

Ex: number of accident in one year, number of travels in one year, number of calls in one hour, number of day cases during an epidemic...

# Poisson likelihood - Poisson Regression.

Poisson distribution is used to model number of apperence of a random event in a given interval.

Ex: number of accident in one year, number of travels in one year, number of calls in one hour, number of day cases during an epidemic...

Law:  $f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$

# Poisson likelihood - Poisson Regression.

Poisson distribution is used to model number of apperence of a random event in a given interval.

Ex: number of accident in one year, number of travels in one year, number of calls in one hour, number of day cases during an epidemic...

Law:  $f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$



# Poisson likelihood - Poisson Regression.

Poisson distribution is used to model number of appearance of a random event in a given interval.

Ex: number of accident in one year, number of travels in one year, number of calls in one hour, number of day cases during an epidemic...

Law:  $f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$        $\lambda$ : expected number of appearance.

# Poisson likelihood - Poisson Regression.

Poisson distribution is used to model number of appearance of a random event in a given interval.

Ex: number of accident in one year, number of travels in one year, number of calls in one hour, number of day cases during an epidemic...

Law:  $f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$        $\lambda$ : expected number of appearance.

For the Poisson regression the mean is not  $\lambda = \beta^T x$  but  $\lambda = e^{\beta^T x}$

# Poisson likelihood - Poisson Regression.

Poisson distribution is used to model number of appearance of a random event in a given interval.

Ex: number of accident in one year, number of travels in one year, number of calls in one hour, number of day cases during an epidemic...

Law:  $f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$        $\lambda$ : expected number of appearance.

For the Poisson regression the mean is not  $\lambda = \beta^T x$  but  $\lambda = e^{\beta^T x}$

Given data set  $(x_1, y_1), \dots, (x_n, y_n)$ ,  
→ Estimate  $\beta$

# Poisson likelihood - Poisson Regression.

Poisson distribution is used to model number of appearance of a random event in a given interval.

Ex: number of accident in one year, number of travels in one year, number of calls in one hour, number of day cases during an epidemic...

Law:  $f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$        $\lambda$ : expected number of appearance.

For the Poisson regression the mean is not  $\lambda = \beta^T x$  but  $\lambda = e^{\beta^T x}$

Given data set  $(x_1, y_1), \dots, (x_n, y_n)$ ,  
→ Estimate  $\beta$

Negative log likelihood:  $L(\beta) = \sum_{i=1}^n e^{\beta^T x_i} - y_i \beta^T x_i + \log(y_i!)$

# Poisson likelihood - Poisson Regression.

Poisson distribution is used to model number of appearance of a random event in a given interval.

Ex: number of accident in one year, number of travels in one year, number of calls in one hour, number of day cases during an epidemic...

Law:  $f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$        $\lambda$ : expected number of appearance.

For the Poisson regression the mean is not  $\lambda = \beta^T x$  but  $\lambda = e^{\beta^T x}$

Given data set  $(x_1, y_1), \dots, (x_n, y_n)$ ,  
→ Estimate  $\beta$

Negative log likelihood:  $L(\beta) = \sum_{i=1}^n e^{\beta^T x_i} - y_i \beta^T x_i + \log(y_i!)$

No closed-form solution for  $\beta = \text{Argmin}_{\beta \in \mathbb{R}^p} L(\beta)$ .

# Poisson likelihood - Poisson Regression.

Poisson distribution is used to model number of appearance of a random event in a given interval.

Ex: number of accident in one year, number of travels in one year, number of calls in one hour, number of day cases during an epidemic...

Law:  $f(k) = \frac{\lambda^k}{k!} e^{-\lambda}$        $\lambda$ : expected number of appearance.

For the Poisson regression the mean is not  $\lambda = \beta^T x$  but  $\lambda = e^{\beta^T x}$

Given data set  $(x_1, y_1), \dots, (x_n, y_n)$ ,

→ Estimate  $\beta$

Negative log likelihood:  $L(\beta) = \sum_{i=1}^n e^{\beta^T x_i} - y_i \beta^T x_i + \log(y_i!)$

No closed-form solution for  $\beta = \text{Argmin}_{\beta \in \mathbb{R}^p} L(\beta)$ .

→ Use an optimizer to compute  $\beta$ .

# Interpretability or flexibility?

---

No good answer!

Most of the techniques using loglikelihood can be classified under the larger class of “**Empirical risk minimization**” method:



# Interpretability or flexibility?

No good answer!

Most of the techniques using loglikelihood can be classified under the larger class of “**Empirical risk minimization**” method:

$$\text{Minimize: } R(h) = \sum_{i=1}^n L(h(x_i), y_i), \quad h : \mathbb{R}^p \rightarrow \mathbb{R}^q$$



# Interpretability or flexibility?

No good answer!

Most of the techniques using loglikelihood can be classified under the larger class of “**Empirical risk minimization**” method:

$$\text{Minimize: } R(h) = \sum_{i=1}^n L(h(x_i), y_i), \quad h : \mathbb{R}^p \rightarrow \mathbb{R}^q$$

$h$ : “**hypothesis**”, in regression tasks,  $h_\beta : x \mapsto \beta^T x$

# Interpretability or flexibility?

No good answer!

Most of the techniques using loglikelihood can be classified under the larger class of “**Empirical risk minimization**” method:

$$\text{Minimize: } R(h) = \sum_{i=1}^n L(h(x_i), y_i), \quad h : \mathbb{R}^p \rightarrow \mathbb{R}^q$$

$h$ : “**hypothesis**”, in regression tasks,  $h_\beta : x \mapsto \beta^T x$

Multiple choices of  $L$ ,  $h \rightarrow$  more flexible methods.

# Interpretability or flexibility?

No good answer!

Most of the techniques using loglikelihood can be classified under the larger class of “**Empirical risk minimization**” method:

$$\text{Minimize: } R(h) = \sum_{i=1}^n L(h(x_i), y_i), \quad h : \mathbb{R}^p \rightarrow \mathbb{R}^q$$

$h$ : “**hypothesis**”, in regression tasks,  $h_\beta : x \mapsto \beta^T x$

Multiple choices of  $L$ ,  $h \rightarrow$  more flexible methods.

When interpretation is available, always prefer interpretable methods.