

LEAST SQUARES LINEAR REGRESSION

In regression problems, we are given training data

$$(x_1, y_1), \dots, (x_n, y_n)$$

where

$$x_i \in \mathbb{R}^d$$

$$y_i \in \mathbb{R}.$$

We assume the (x_i, y_i) are realizations of a random pair (X, Y) . The goal of regression is to predict the response y associated to a new input x .

A regression model posits

$$Y = f(X) + \epsilon$$

where

$$\epsilon = \text{noise}$$

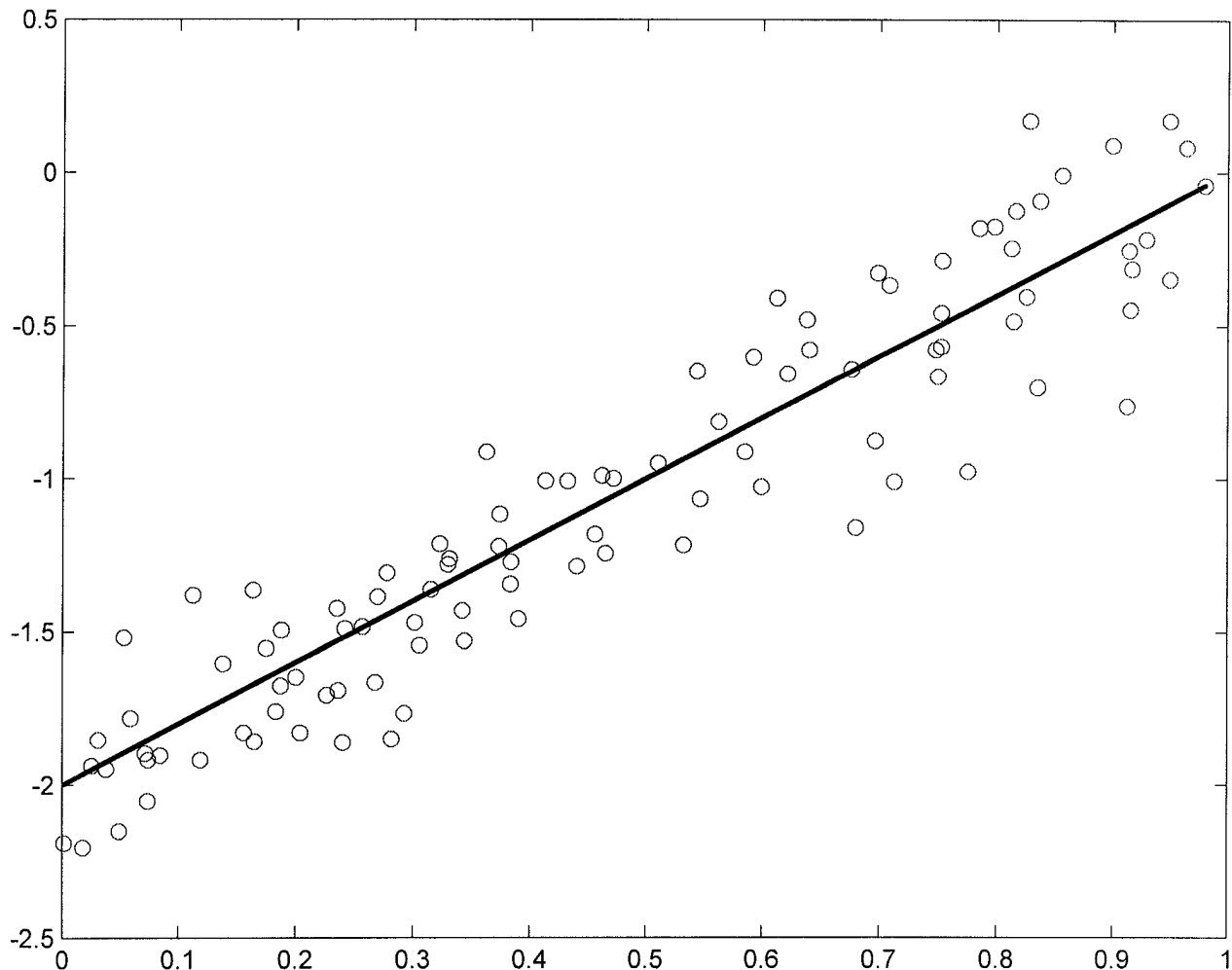
$f \in$ some class of functions.

In linear regression, we assume a linear model

$$f(x) = \beta^T x + \beta_0$$

where

$$\beta \in \mathbb{R}^d, \quad \beta_0 \in \mathbb{R}$$



The challenge in linear regression is to estimate the parameters β, β_0 , from training data.

Least Squares

In least squares linear regression, we select β, β_0 to minimize the sum of squared errors,

$$SSE(\beta, \beta_0) := \sum_{i=1}^n (y_i - \beta^T x_i - \beta_0)^2$$

Example | Suppose $d=1$, so x_i, β are scalars.

(A)

$$\frac{\partial SSE}{\partial \beta_0} = \dots = 0$$

$$\Rightarrow \beta_0 =$$

$$\frac{\partial SSE}{\partial \beta} =$$

$$\Rightarrow \beta =$$

In matrix form,

$$\begin{bmatrix} \dots \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix} = \begin{bmatrix} \dots \end{bmatrix}$$

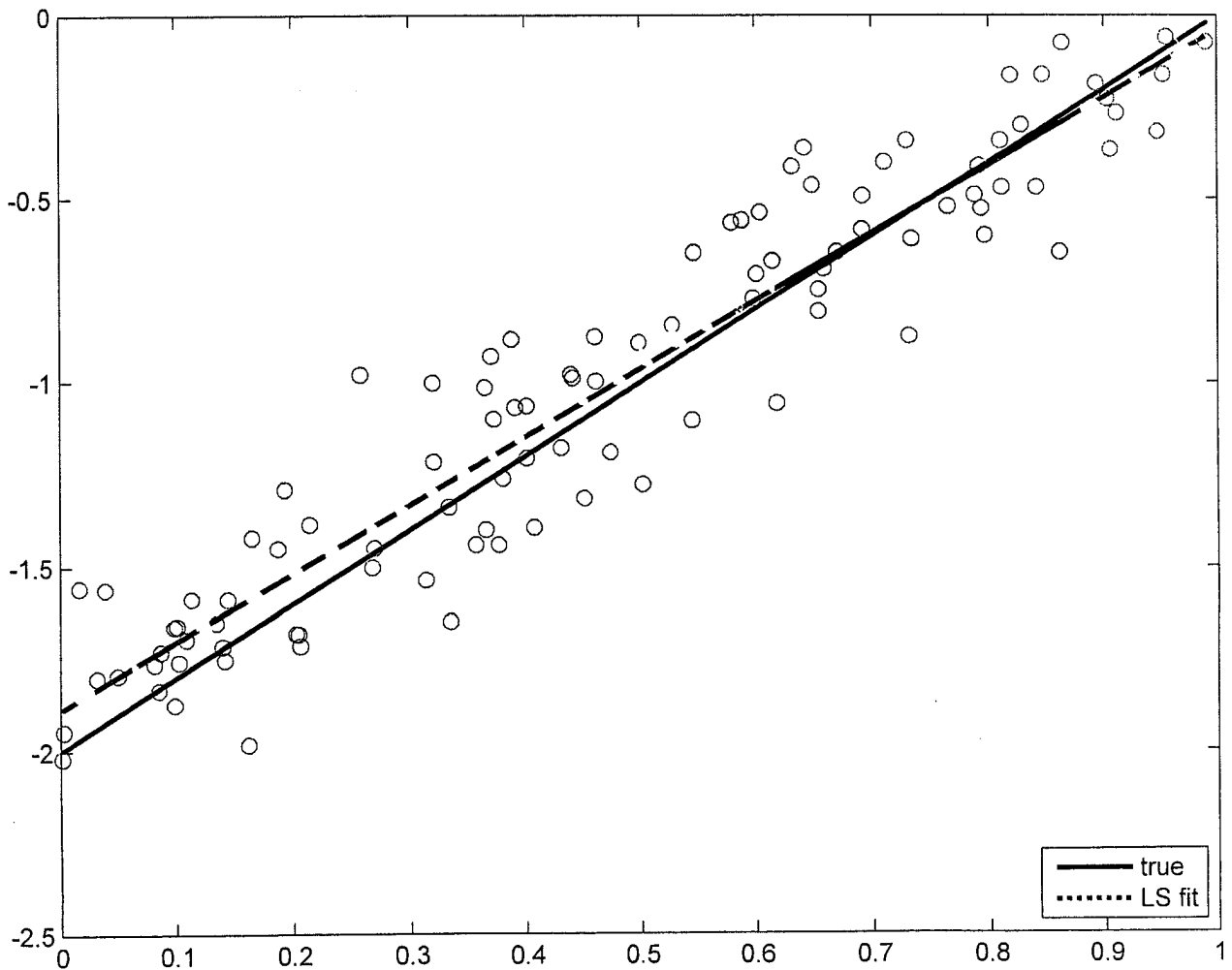
Inverting the matrix,

$$\begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta} \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\bar{y} (\sum x_i^2) - \bar{x} \sum x_i y_i}{\sum x_i^2 - n \bar{x}^2} \\ \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} \end{bmatrix}$$

where

$$\bar{x} = \frac{1}{n} \sum x_i, \quad \bar{y} = \frac{1}{n} \sum y_i$$



More generally, suppose d is arbitrary. Set

$$\theta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}$$

Then

$$\begin{aligned} \text{SSE}(\theta) &= \sum_{i=1}^n (y_i - \beta^T x_i - \beta_0)^2 \\ &= \| \underline{y} - A\theta \|^2 \end{aligned}$$

where

(B)

$$\underline{y} = \begin{bmatrix} \\ \\ \\ \\ \end{bmatrix}, \quad A = \begin{bmatrix} \phantom{x_{11}} & \phantom{x_{12}} & \phantom{x_{13}} & \phantom{x_{14}} & \phantom{x_{15}} \\ \phantom{x_{21}} & \phantom{x_{22}} & \phantom{x_{23}} & \phantom{x_{24}} & \phantom{x_{25}} \\ \phantom{x_{31}} & \phantom{x_{32}} & \phantom{x_{33}} & \phantom{x_{34}} & \phantom{x_{35}} \\ \phantom{x_{41}} & \phantom{x_{42}} & \phantom{x_{43}} & \phantom{x_{44}} & \phantom{x_{45}} \\ \phantom{x_{51}} & \phantom{x_{52}} & \phantom{x_{53}} & \phantom{x_{54}} & \phantom{x_{55}} \end{bmatrix}$$

$n \times (d+1)$

The minimizer $\hat{\theta}$ of this quadratic objective function is

$$\hat{\theta} = (A^T A)^{-1} A^T \underline{y}$$

provided

To see this, write

$$\textcircled{C} \quad \|\underline{y} - A\theta\|^2 =$$
$$=$$

Linear regression is easy and works well when the true f is indeed linear. But often it is not. What can we do?

Sometimes f is linear in variables $\phi_1(x), \dots, \phi_k(x)$, where ϕ_k is possibly nonlinear. In such a case we can model

$$f(x) = \sum_{k=1}^K \beta_k \phi_k(x) + \beta_0$$

and estimate $\beta_0, \beta_1, \dots, \beta_K$ using least squares applied to the "data"

$$(\phi(x_1), y_1), \dots, (\phi(x_n), y_n).$$

Exercise | Suppose $f(x)$ is a cubic polynomial.
Determine the least-squares estimate of f
given $(x_1, y_1), \dots, (x_n, y_n)$.

Solution

$$f(x) = \beta_3 x^3 + \beta_2 x^2 + \beta_1 x + \beta_0$$

$$\Rightarrow \phi_k(x) = x^k$$

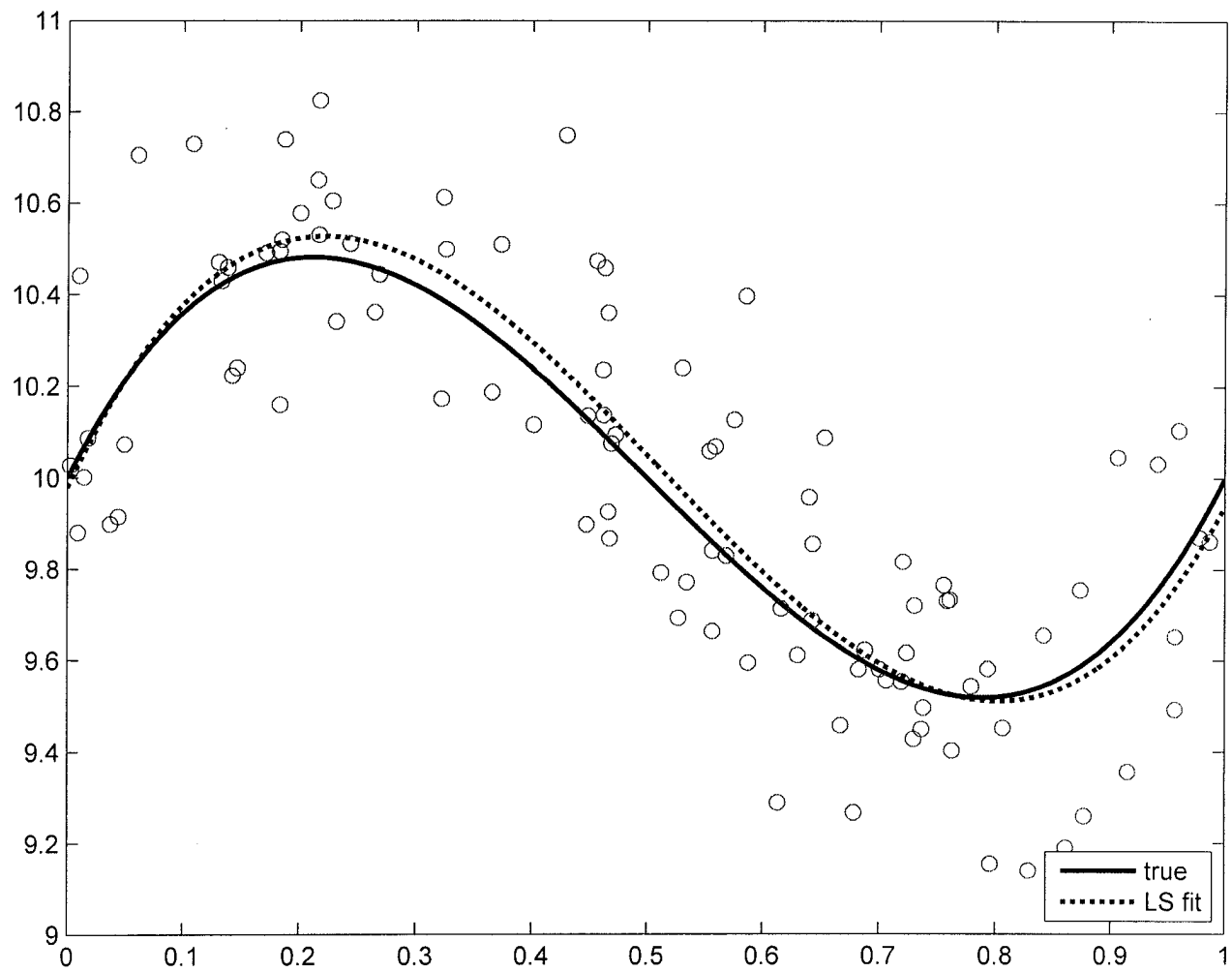
(D) \Rightarrow

$$A = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}$$

$$\Rightarrow \hat{\theta} = (A^T A)^{-1} A^T \underline{y}$$

gives the LS cubic polynomial fit.

What if a polynomial model is also not appropriate, or the degree is unknown? We'll address these and other issues later in the course.



Key

$$A. \quad \frac{\partial SSE}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta x_i - \beta_0) = 0$$

$$\Rightarrow \beta_0 = \frac{1}{n} \sum_i (y_i - \beta x_i)$$

$$\frac{\partial SSE}{\partial \beta} = -2 \sum_i x_i (y_i - \beta x_i - \beta_0) = 0$$

$$\Rightarrow \beta = \frac{\sum x_i (y_i - \beta_0)}{\sum x_i^2}$$

$$\begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

$$B. \quad \underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad A = \begin{bmatrix} 1 & x_{11} & \dots & x_{1d} \\ 1 & x_{21} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{nd} \end{bmatrix}$$

$$\hat{\theta} = (A^T A)^{-1} A^T \underline{y}$$

provided $A^T A$ is nonsingular $\Leftrightarrow \text{rank}(A) = d+1$

$$\begin{aligned}
 C. \quad \|\underline{y} - A\theta\|^2 &= (\underline{y} - A\theta)^T (\underline{y} - A\theta) \\
 &= \theta^T A^T A \theta - 2 \underline{y}^T A \theta + \underline{y}^T \underline{y}
 \end{aligned}$$

$$\text{proof 1: } \frac{\partial}{\partial \theta} (\checkmark) = 2 A^T A \theta - 2 A^T \underline{y} = 0$$

$$\Rightarrow \hat{\theta} = (A^T A)^{-1} A^T \underline{y}$$

$$\text{proof 2: } \theta^T B \theta + \underline{c}^T \theta + d$$

$$= \left(\theta + \frac{1}{2} B^{-1} \underline{c} \right)^T B \left(\theta + \frac{1}{2} B^{-1} \underline{c} \right)$$

$$+ \left(d - \frac{1}{4} \underline{c}^T B^{-1} \underline{c} \right)$$

If B is positive definite, then the unique minimizer is

$$\theta = -\frac{1}{2} B^{-1} \underline{c}.$$

Apply this with $B = A^T A$, $\underline{c} = -2 A^T \underline{y}$, $d = \underline{y}^T \underline{y}$

$$D. \quad A = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 \end{bmatrix}$$