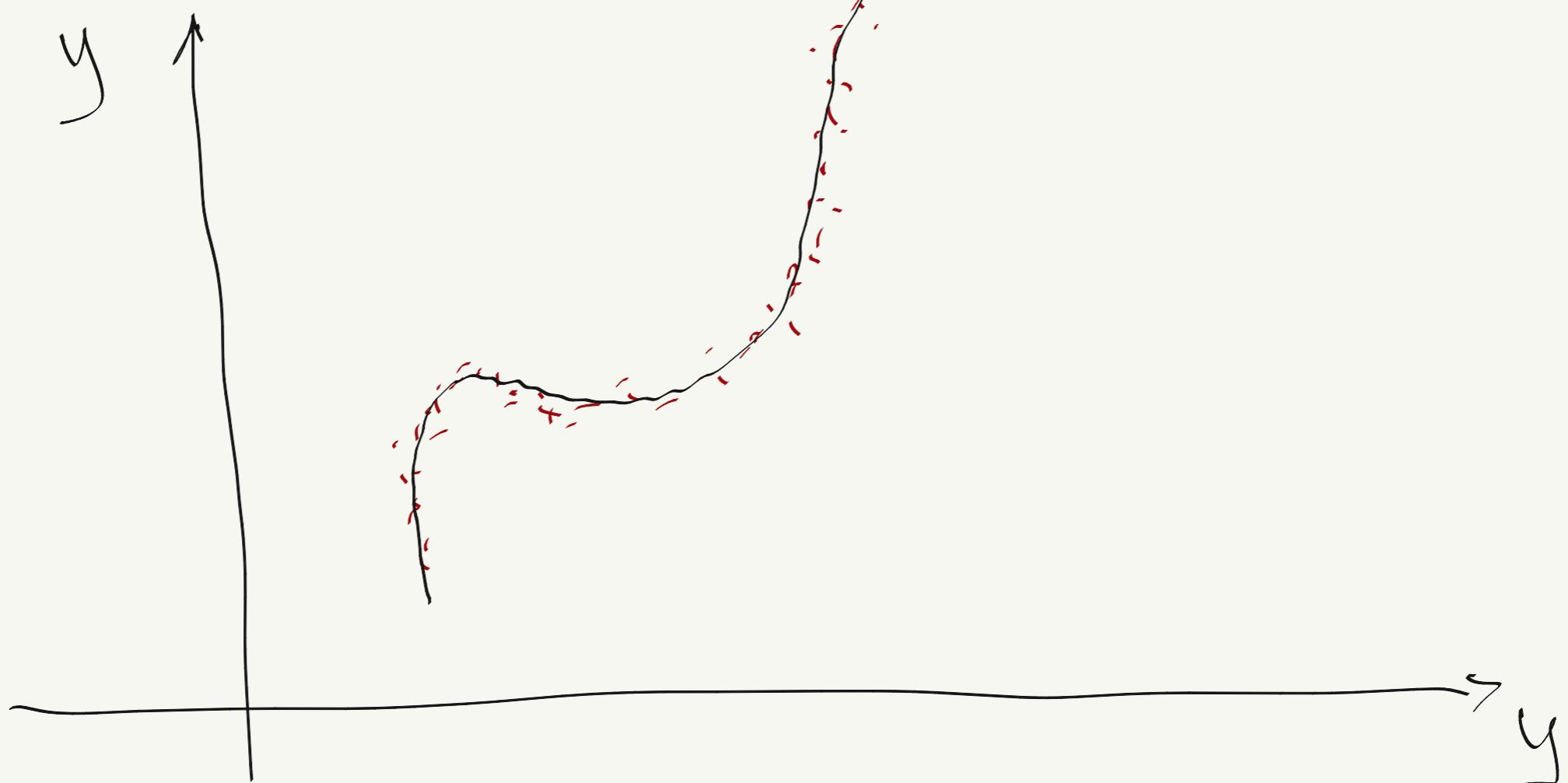
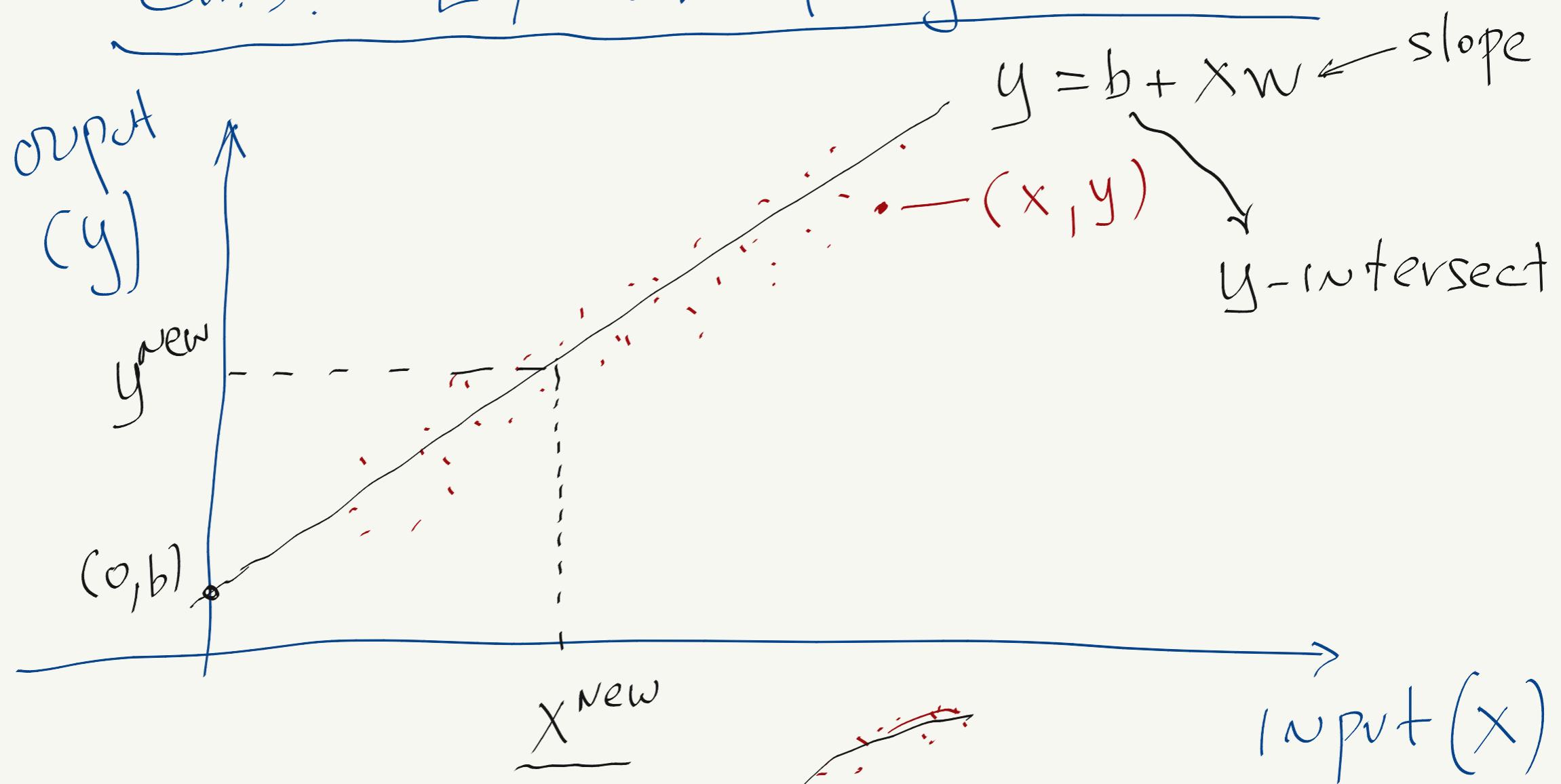


4/17/20

Ch. 5. Linear Regression.



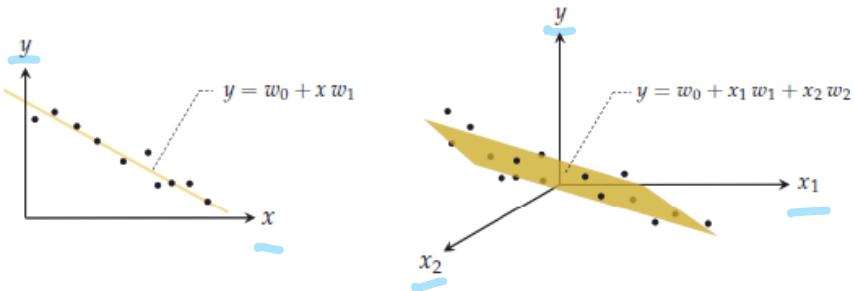


Figure 5.1 (left panel) A simulated dataset in two dimensions along with a well-fitting line. A line in two dimensions is defined as $w_0 + xw_1 = y$, where w_0 is referred to as the bias and w_1 the slope, and a point (x_p, y_p) lies close to it if $w_0 + x_p w_1 \approx y_p$. (right panel) A simulated three-dimensional dataset along with a well-fitting hyperplane. A hyperplane in general is defined as $w_0 + x_1w_1 + x_2w_2 + \cdots + x_Nw_N = y$, where again w_0 is the bias and w_1, w_2, \dots, w_N the hyperplane's coordinate-wise slopes, and a point (\mathbf{x}_p, y_p) lies close to it if $w_0 + x_{1,p}w_1 + x_{2,p}w_2 + \cdots + x_{N,p}w_N \approx y_p$. Here $N = 2$.

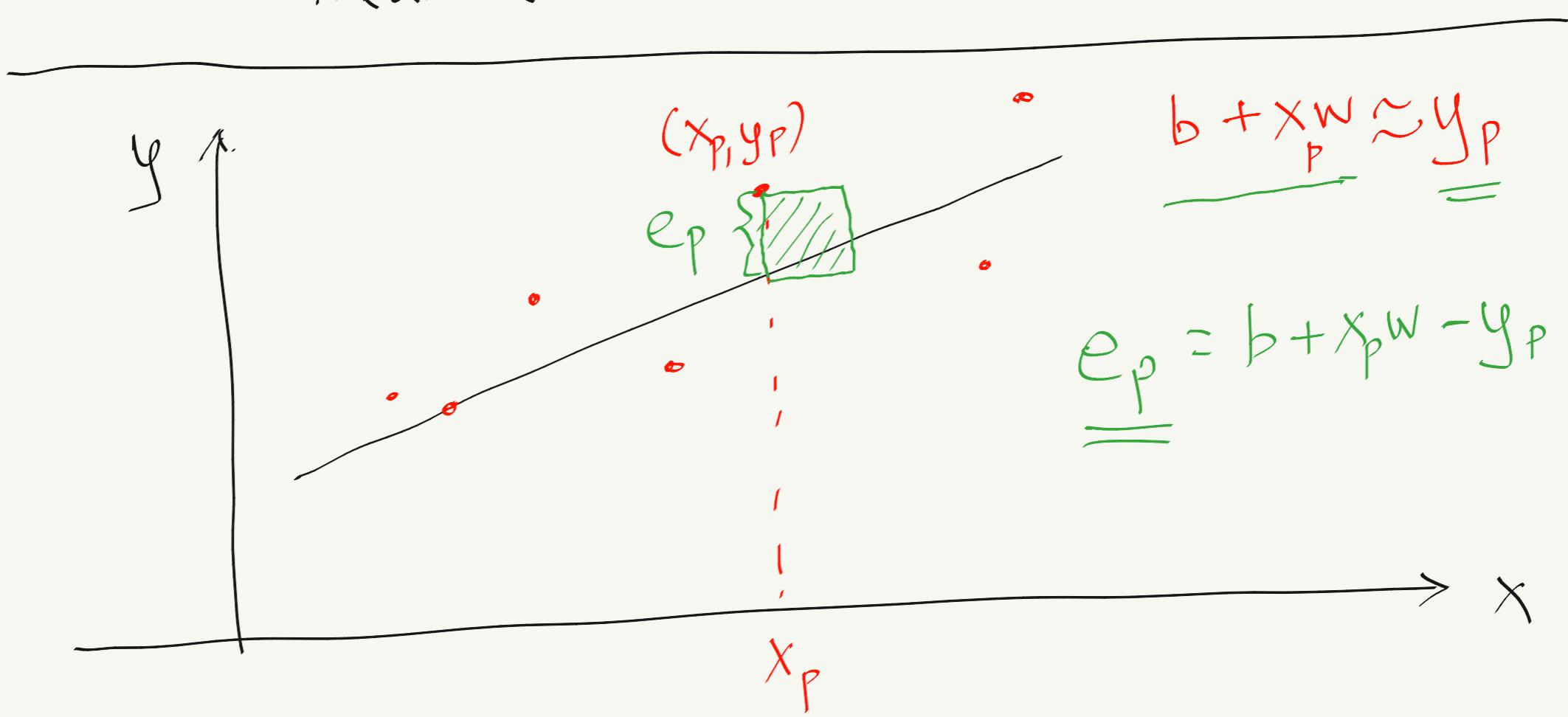
Regression data: P input/output pairs

$$\left\{ (\bar{x}_1, y_1), (\bar{x}_2, y_2), \dots, (\bar{x}_P, y_P) \right\}$$

\downarrow \downarrow
 $\in \mathbb{R}^N$ $\in \mathbb{R}^1$

$$\bar{x}_P = \begin{bmatrix} x_{1P} & \dots & x_{NP} \end{bmatrix}^T$$

\uparrow \uparrow
1st input Nth input feature
feature



error for p -th data point

$$e_p = (b + x_p w - y_p)$$

Least-Squares Regression

$$g(b, w) = \frac{1}{P} \sum_{p=1}^P e_p^2$$

$$= \frac{1}{P} \sum_{p=1}^P (b +$$

$$g(b, \bar{w}) = \left(\frac{1}{P} \sum_{p=1}^P \left(b + \bar{x}_p^T \bar{w} - y_p \right)^2 \right)$$

↓
 feature
 forcing weights Unknown $|XN$ $N \times 1$
 forcing weights Unknowns

$$N\text{-dim} : b + \bar{x}^T w = y \quad \text{equation of hyperplane}$$

Optimization Problem

Find b, \bar{w} so that $g(b, \bar{w})$ is as small as possible

Modify Notation

$$\tilde{x}_p = \begin{bmatrix} 1 \\ \bar{x}_p \end{bmatrix} \quad \tilde{w} = \begin{bmatrix} b \\ \bar{w} \end{bmatrix}$$

$(N+1) \times 1 \qquad \qquad \qquad (N+1) \times 1$

$$\bar{x}_p^T \cdot \tilde{w} = [1 \ \bar{x}_p^T] \begin{bmatrix} b \\ \bar{w} \end{bmatrix} = b + \bar{x}_p^T \bar{w}$$

Rewrite cost function

$$g(\tilde{w}) = \frac{1}{P} \sum_{p=1}^P \left(\underbrace{\tilde{x}_p^T \tilde{w} - y_p}_{\text{error}} \right)^2$$

$\left(\bar{x}^T \bar{B} \bar{x} + \bar{x}^T b \right)$

{ Optimization problem

$$\min_{\tilde{w}} g(\tilde{w}) \Rightarrow \tilde{w}^* = \arg \min_{\tilde{w}} g(\tilde{w})$$

1st order optimality condition

stationary points

$$\nabla_{\tilde{w}} g(\tilde{w}) = 0$$

$$\therefore \nabla \frac{1}{P} \sum_{p=1}^P (\tilde{x}_p^\top \tilde{w} - y_p)^2 = 0$$

$$\Rightarrow \frac{1}{P} \sum_{p=1}^P \nabla (\tilde{x}_p^\top \tilde{w} - y_p)^2 = 0$$

$$\Rightarrow \frac{1}{P} \sum_{p=1}^P 2 (\tilde{x}_p^\top \tilde{w} - y_p) \cdot \underbrace{\nabla_{\tilde{w}} (\tilde{x}_p^\top \tilde{w} - y_p)}_0 = 0$$

$$\nabla_{\tilde{w}} (\tilde{x}_p^\top \tilde{w}) - \nabla_{\tilde{w}} y_p = 0$$

$$\nabla_{\tilde{w}} (\tilde{w}^\top \tilde{x}_p) = \tilde{x}_p$$

$$\Rightarrow \left(\frac{2}{P} \right) \sum_{p=1}^P \underbrace{(\tilde{x}_p^\top \tilde{w} - y_p)}_{\text{scalar}} \tilde{x}_p = 0$$

$$\Rightarrow \sum_{p=1}^P \tilde{x}_p \cdot (\tilde{x}_p^\top \tilde{w} - y_p) = 0$$

$$\Rightarrow \boxed{\sum_{p=1}^P \tilde{x}_p \cdot \tilde{x}_p^\top \tilde{w} = \sum_{p=1}^P \tilde{x}_p \cdot y_p}$$

Normal Equations

$$\sum_{P=1}^P \underbrace{\begin{pmatrix} \tilde{x}_P & \tilde{x}_P^T \end{pmatrix}}_A \tilde{w} = \sum_{P=1}^P \underbrace{\tilde{x}_P \cdot y_P}_{\bar{b}}$$

outer product $\left[(N+1) \times 1 \right] \left[1 \times (N+1) \right]$
 $\underbrace{(N+1) \times (N+1)}$
 $(N+1) \times 1$

$A \tilde{w} = \bar{b}$

If A is invertible:

$$\tilde{w}^* = A^{-1} \bar{b}$$

if A is non-invertible:

$$\tilde{w}^* = A^+ \bar{b}$$

generalized inverse

Iterative Optimization

- Newton's method : In one step we obtain the solution, since Hessian = A
 \Rightarrow solve normal equations in one-step

• GD

$$\tilde{w}^{k+1} = \tilde{w}^k - \alpha \nabla g(\tilde{w}^k)$$

$$= \tilde{w}^k - \alpha (A\tilde{w}^k - \bar{b})$$

$$= (\mathbb{I} - \alpha A) \tilde{w}^k + \alpha \bar{b}$$

need to be computed once

\Downarrow

$$\boxed{\tilde{w}^{k+1} = B \tilde{w}^k + \bar{c}}$$

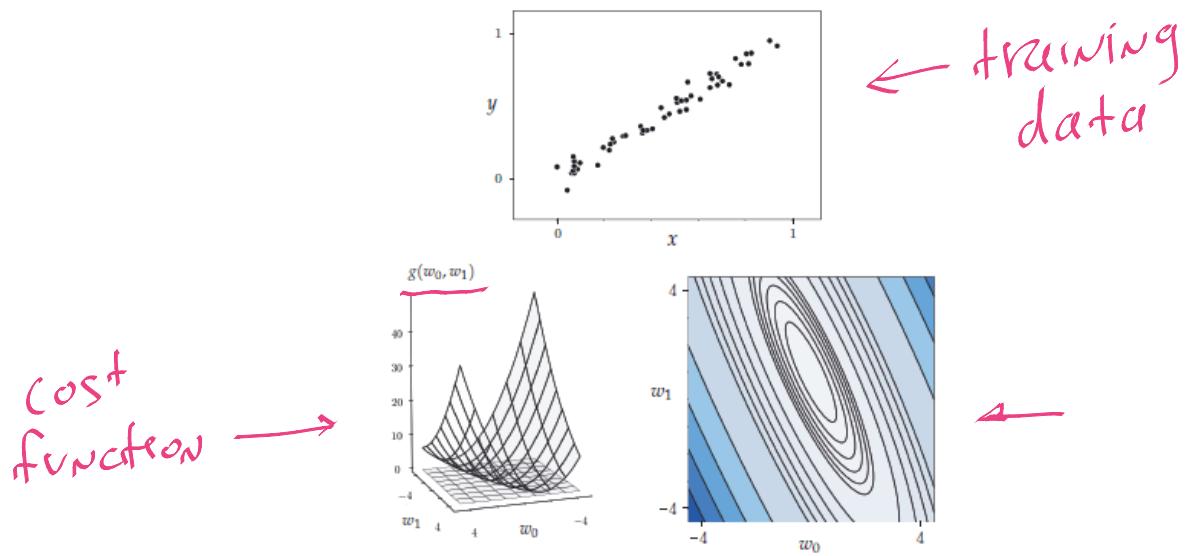


Figure 5.3 Figure associated with Example 5.1. See text for details.

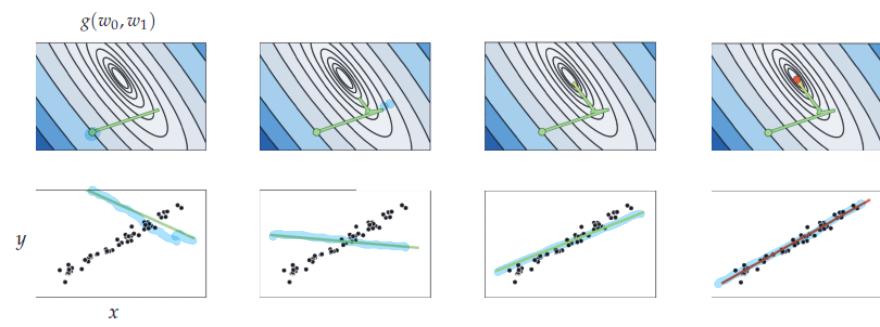
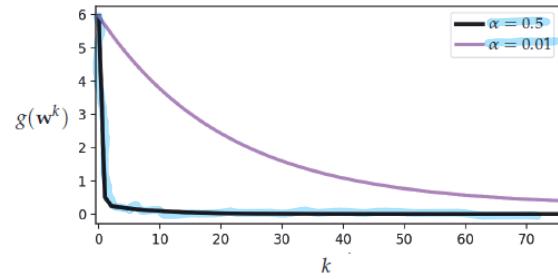


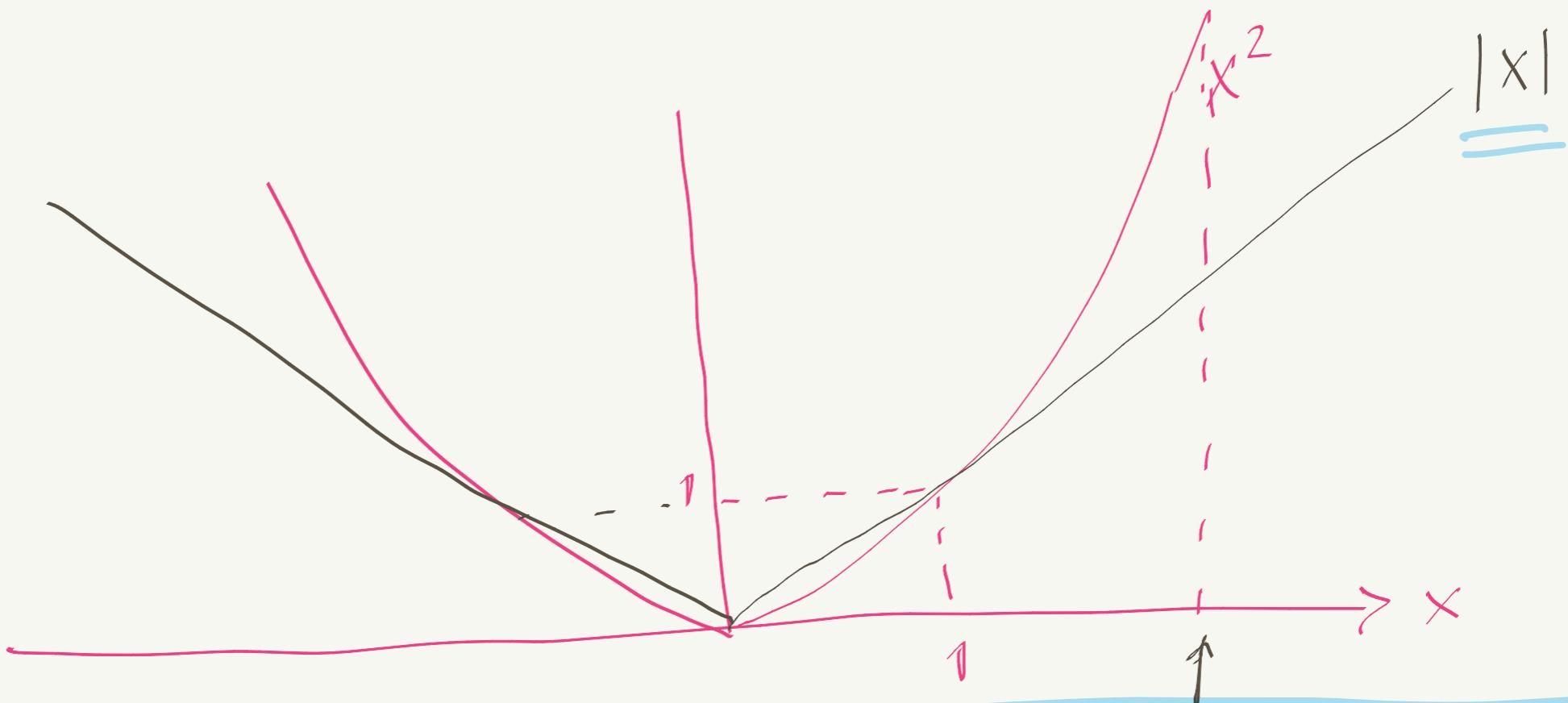
Figure 5.4 Figure associated with Example 5.2. See text for details.

Figure 5.5 Figure associated with Example 5.2. See text for details.



Squaring the error increases the importance of large errors

⇒ in trying to minimize them we overfit to outliers

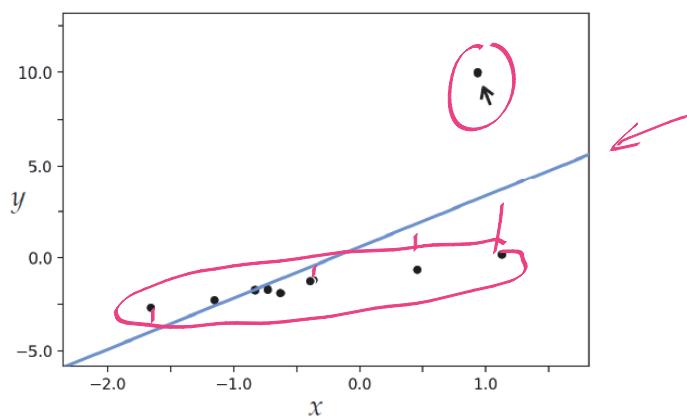


Least Absolute Deviations Regressor

$$g(\tilde{w}) = \frac{1}{P} \sum_{p=1}^P |\tilde{x}^T \tilde{w} - y_p|$$

error

Figure 5.6 Figure associated with Example 5.3. See text for details.



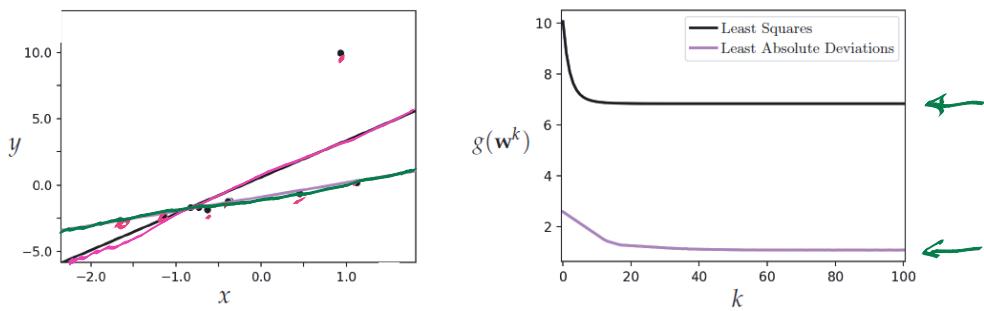


Figure 5.7 Figure associated with Example 5.4. See text for details.