

4/22/20

Linear two-class Classification

P Input pairs: $\left\{ (\bar{x}_p, y_p) \right\}_{p=1}^P$

$\bar{x}_p \in \mathbb{R}^N$; $y_p \in \{0, +1\}$, OR $y_p \in \{-1, 1\}$

Two perspectives on classification

→ regression perspective (extended view)
→ Perceptron ↼ (view from "above")

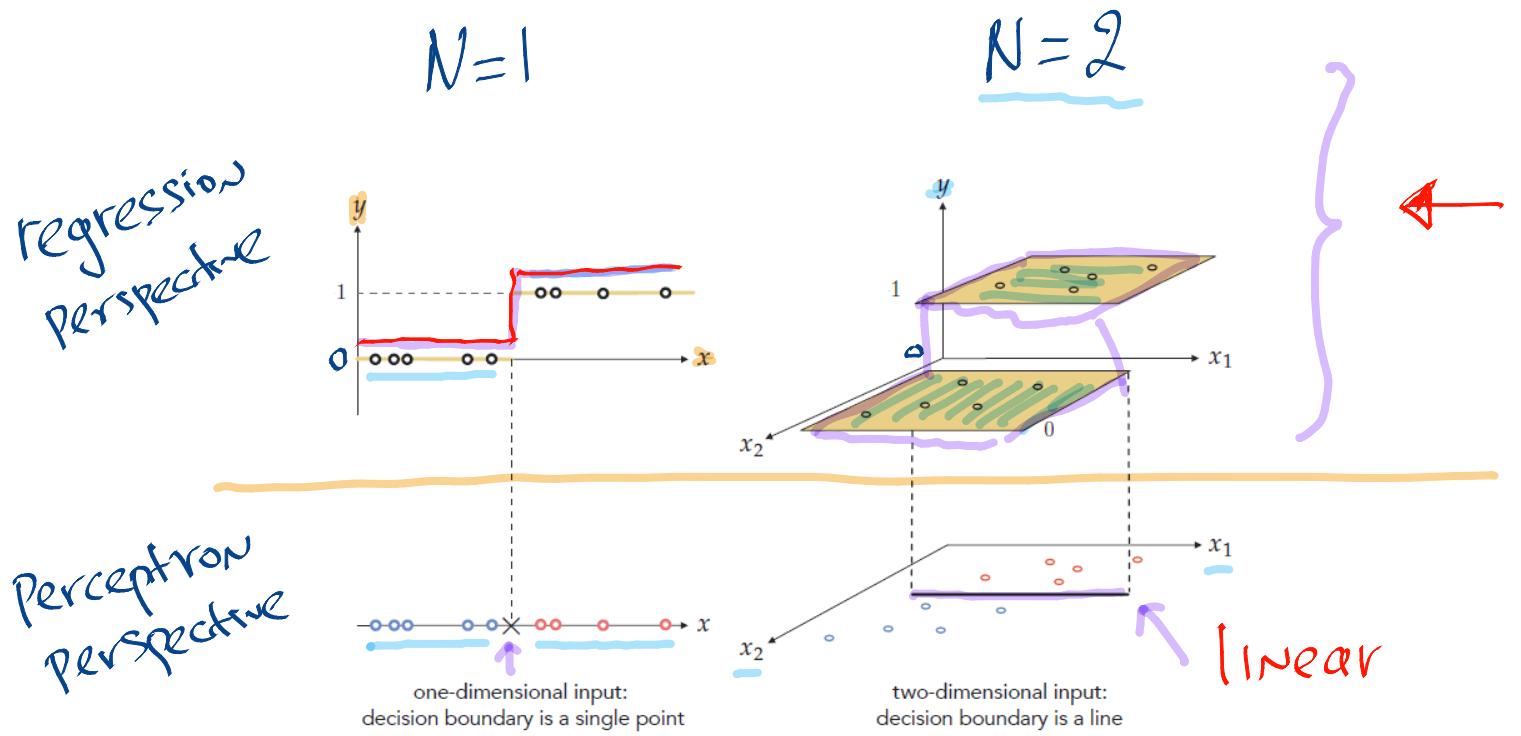
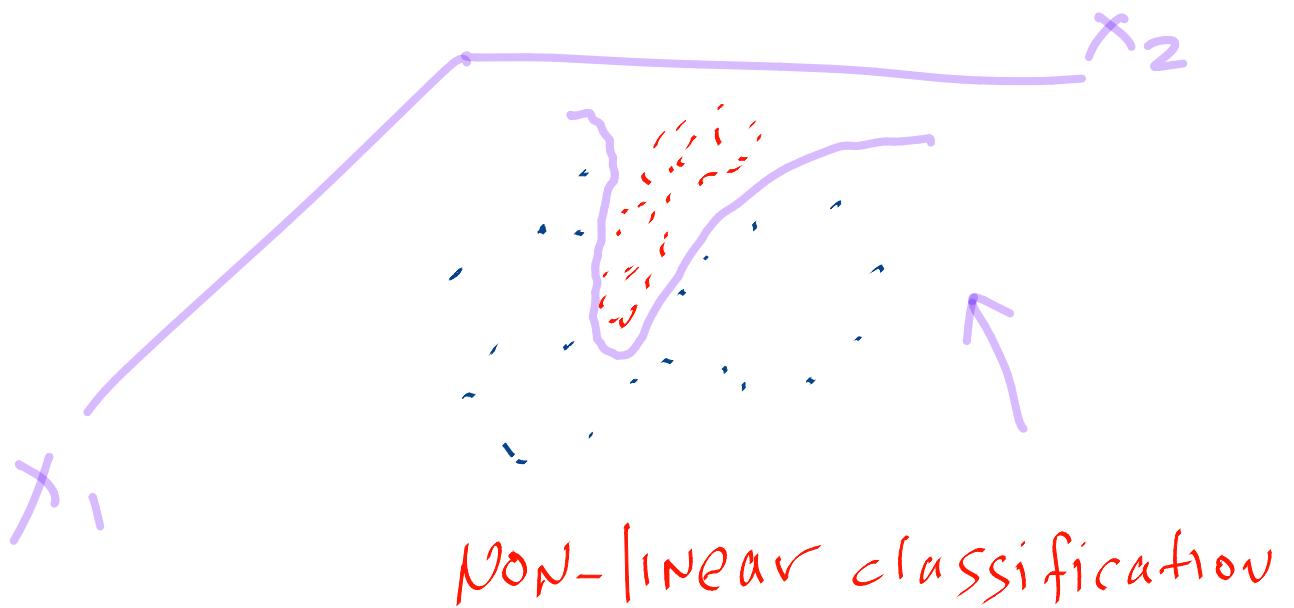


Figure 6.1 Two perspectives on classification illustrated using single-input (left column) and two-input (right column) toy datasets. The regression perspective shown in the top panels is equivalent to the perceptron perspective shown in the bottom panels, where we look at each respective dataset from "above." In the Perceptron perspective we also mark the decision boundary. This is where the step function (colored in yellow in the top panels) transitions from its bottom to top step. See text for further details.



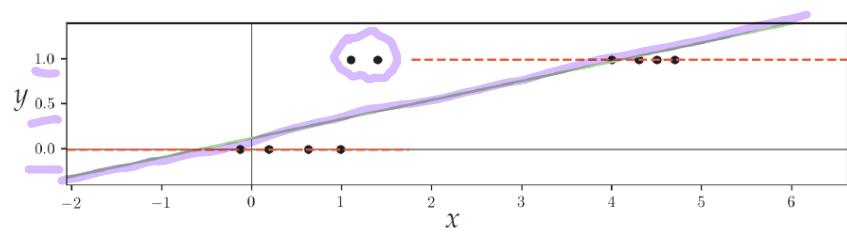
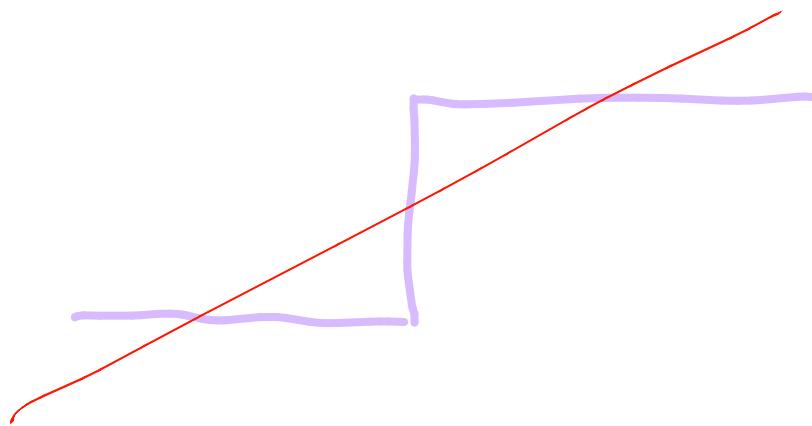
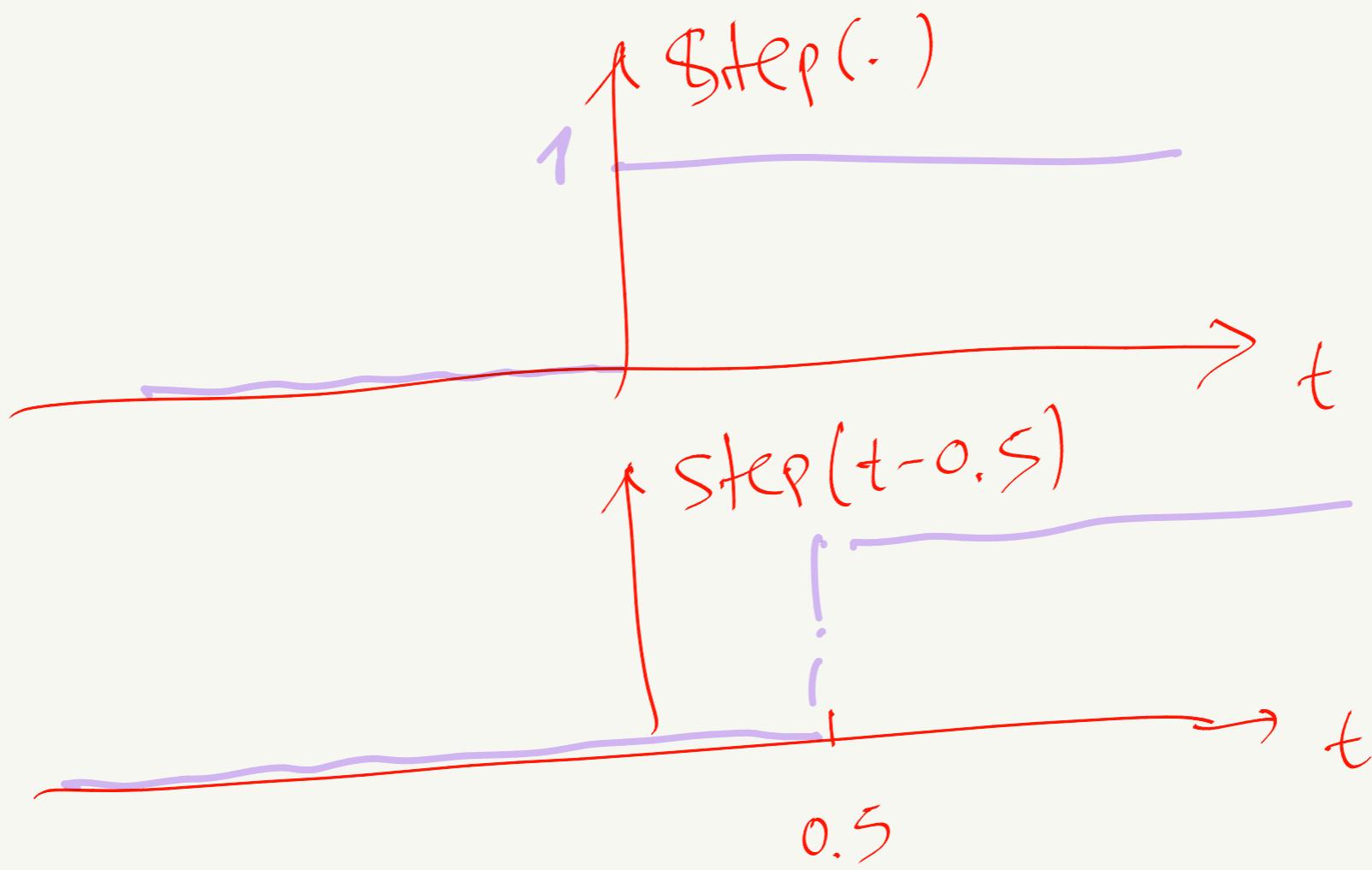


Figure 6.2 Figure associated with Example 6.1. See text for details.



pass hyperplane thru a step function
 → perform optimization

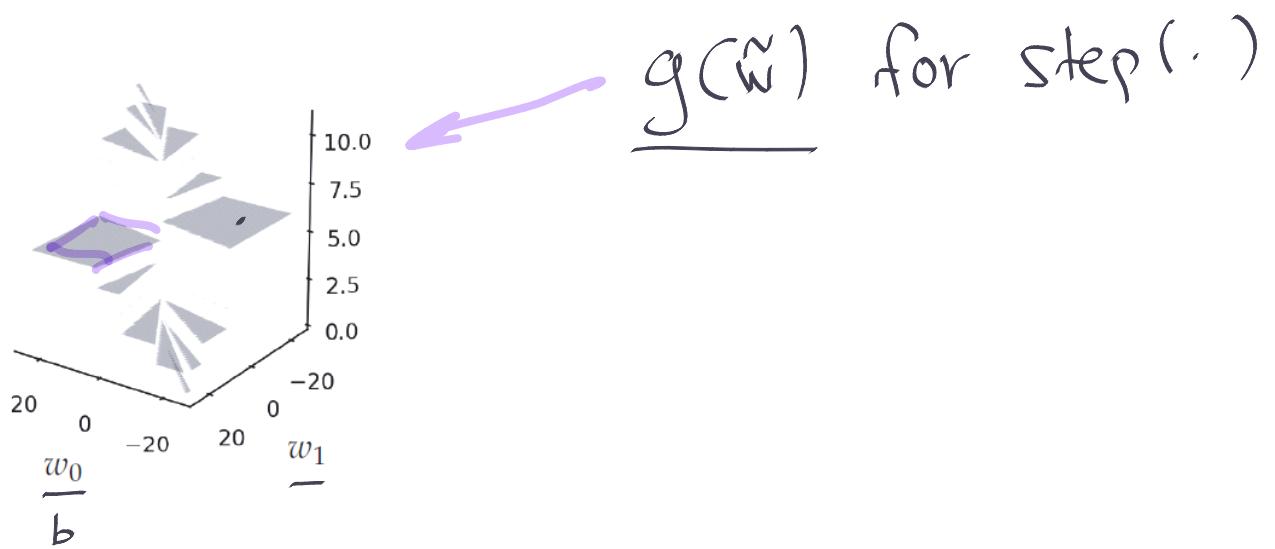
$$\text{Step}(t) = \begin{cases} 1, & t > 0 \\ 0, & t < 0 \end{cases} \quad (\text{undefined for } t=0)$$



$$\underbrace{N-d}_{\text{(text: } \tilde{x} \rightarrow \overset{\circ}{x})} \quad \tilde{X} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_N \end{bmatrix}_{(N+1) \times 1} \quad \tilde{w} = \begin{bmatrix} b(w_0) \\ w_1 \\ \vdots \\ w_N \end{bmatrix}$$

$$\begin{pmatrix} \tilde{x} \\ \tilde{w} \end{pmatrix} \rightarrow \overset{\circ}{x}$$

$$\text{hyperplane: } \tilde{X}^T \tilde{w}$$



$$\text{Step}(\tilde{x}_p^\top \tilde{w}) \approx y_p, \quad p=1, \dots, P$$

One approach

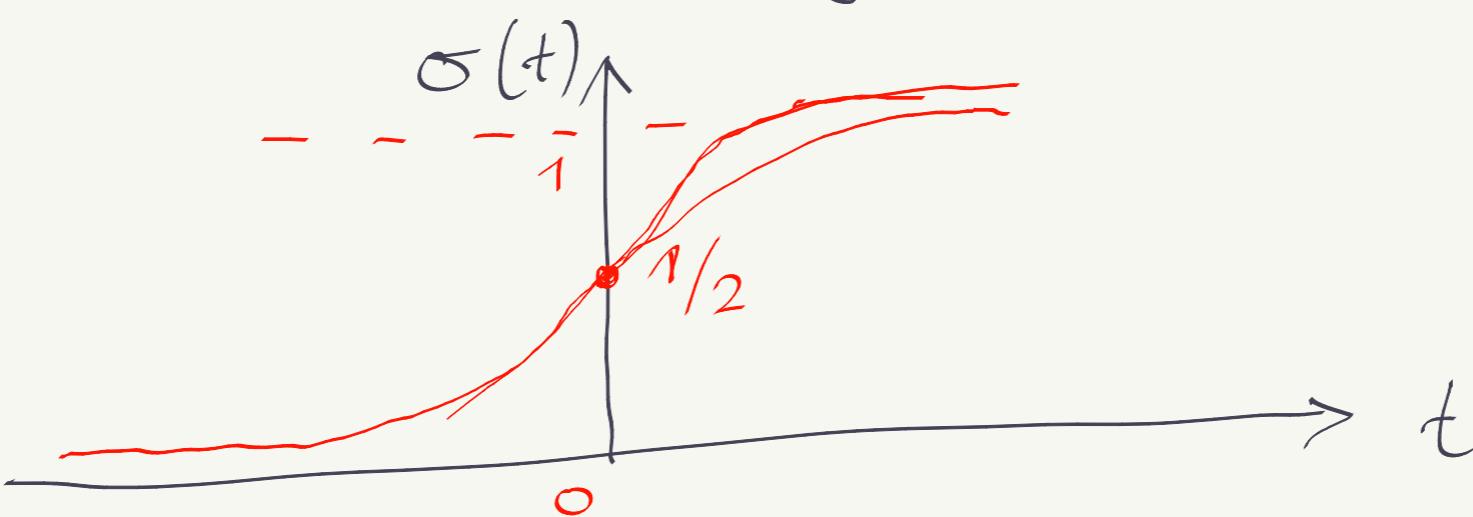
$$\Rightarrow \text{LS: } g(\tilde{w}) = \frac{1}{P} \sum_{p=1}^P (\text{Step}(\tilde{x}_p^\top \tilde{w}) - y_p)^2$$

Problematic

MAIN IDEA: approximate step function

Logistic Sigmoid function.

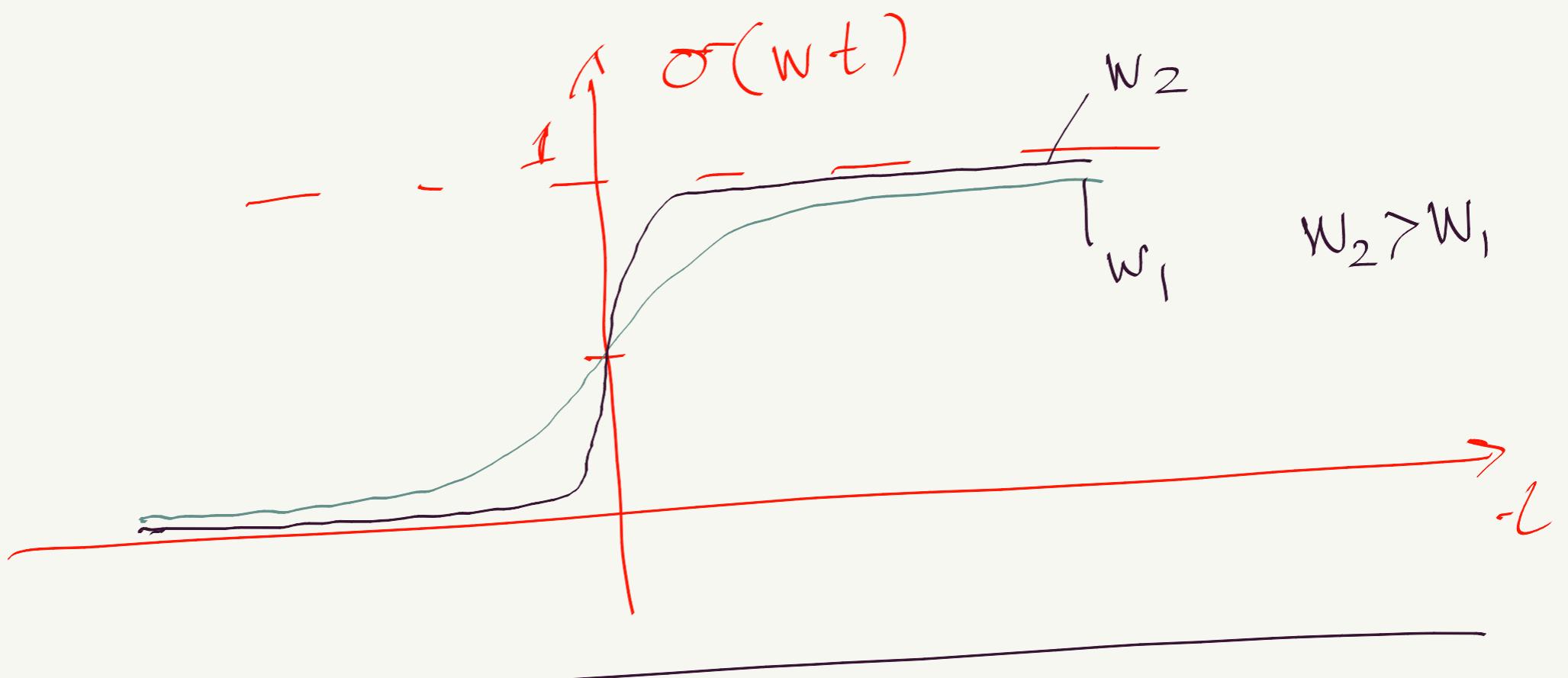
$$\sigma(t) = \frac{1}{1 + e^{-t}}$$



$$\lim_{t \rightarrow \infty} \sigma(t) = \frac{1}{1 + e^{-\infty}} = 1$$

$$\left. \begin{aligned} & t = 0 \\ & \sigma(0) = \frac{1}{1+1} = 1/2 \end{aligned} \right|$$

$$\lim_{t \rightarrow -\infty} \sigma(t) = \frac{1}{1 + e^{\infty}} = 0$$



19th century mathematician Verhulst
population growth in finite system

f : current population

1 : max capacity of system

$(1-f)$: remaining capacity

$$\rightarrow \frac{df}{dt} = f(1-f)$$

we can easily show that $\sigma(t)$ satisfies

$$\sigma(t) = \frac{1}{1+e^{-t}}$$

$$\sigma'(t) = \frac{-(-e^{-t})}{(1+e^{-t})^2} =$$

$$= \frac{1}{1+e^{-t}} \cdot \frac{e^{-t}}{1+e^{-t}}$$

$\underbrace{1+e^{-t}}_{\sigma(t)} \quad \underbrace{1+e^{-t}}_{1-\sigma(t)}$

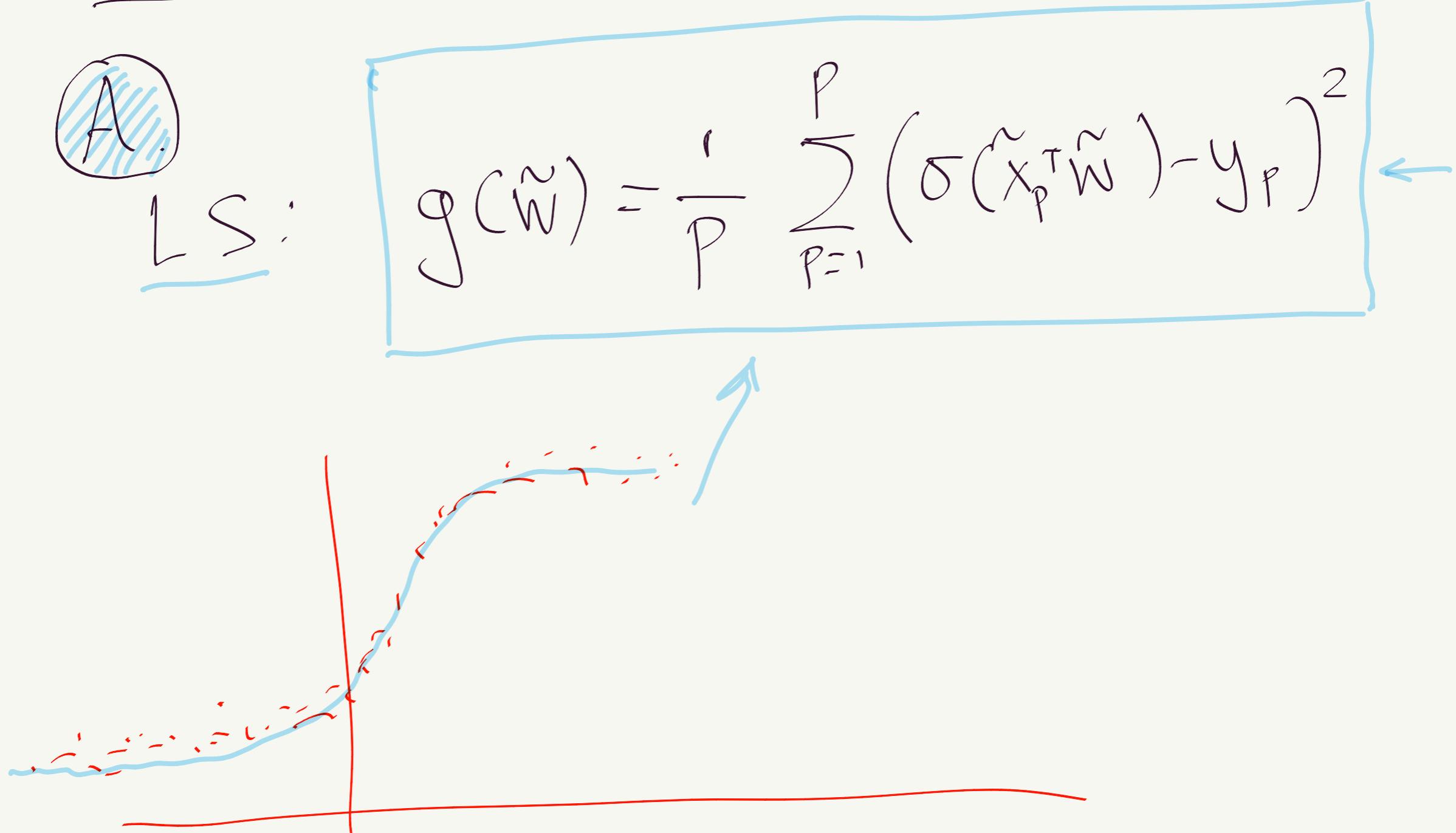
$$1-\sigma(t) = \frac{1+e^{-t}}{1+e^{-t}} - \frac{1}{1+e^{-t}} = \frac{e^{-t}}{1+e^{-t}} \quad \checkmark$$

$$\sigma'(t) = \sigma(t) [1 - \sigma(t)]$$

Logistic regression

$$\sigma(\tilde{x}_p^\top \tilde{w}) \approx y_p, \quad p=1, \dots, P$$

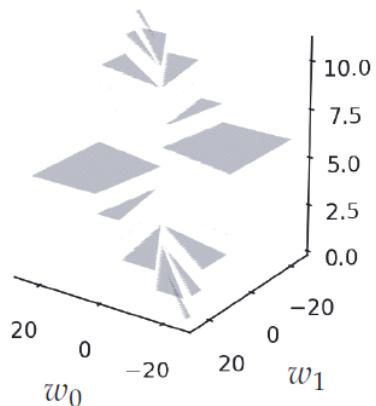
Cost function



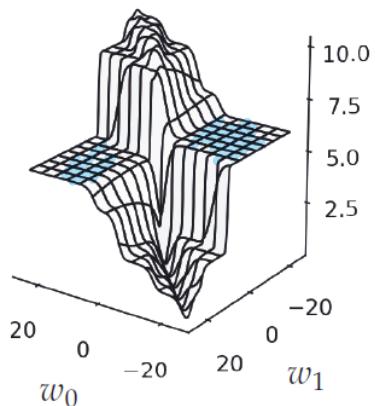
$g(\tilde{w})$:

- NON-CONVEX
- can be optimized using GD, Newton (w/ variations)

$g(\tilde{w})$



$g(\tilde{w})$



↑
step

↑
 $\sigma(t)$

Logistic regression

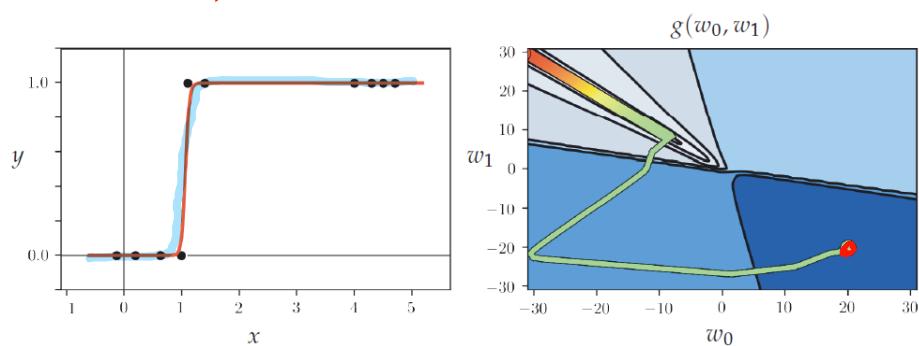


Figure 6.5 Figure associated with Example 6.3. See text for details.

GD

$$w_0 = -w_1 = 20$$

Normalized GD

Logistic Regression using Cross Entropy Cost

$$y_p \in \{0, 1\}$$

log error

$$g_p(\tilde{w}) = \begin{cases} -\log(\sigma(\tilde{x}_p^T \tilde{w})), & y_p = 1 \\ -\log(1 - \sigma(\tilde{x}_p^T \tilde{w})), & y_p = 0 \end{cases}$$

(NON-Negative)
 $\min = 0$

$$g_p(\tilde{w}) = -y_p \cdot \log \sigma(\tilde{x}_p^T \tilde{w}) - (1-y_p) \log(1 - \sigma(\tilde{x}_p^T \tilde{w}))$$

\equiv

$$y_p = 1 \rightarrow 1$$

$$y_p = 0 \rightarrow 0$$

0

1

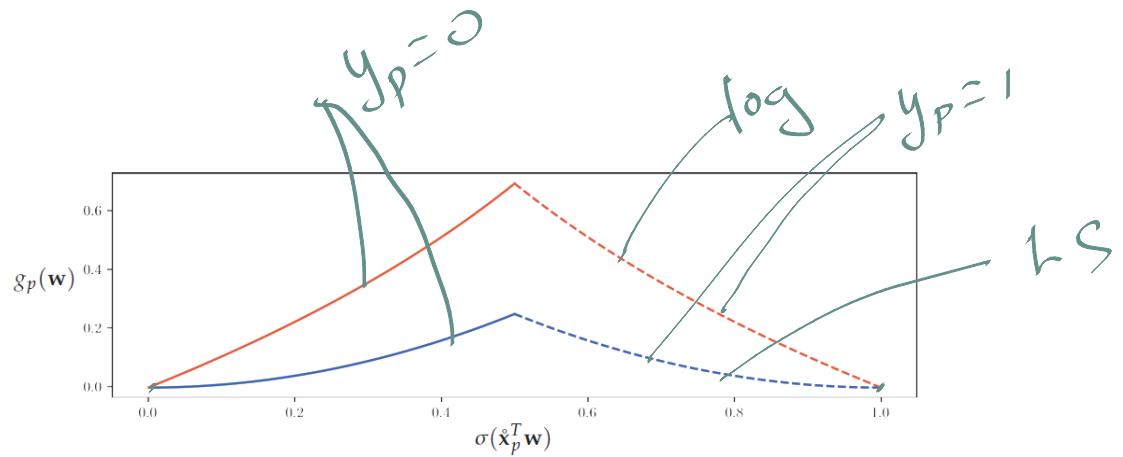


Figure 6.6 Visual comparison of the squared error (in blue) and the log error (in red) for two cases: $y_p = 0$ (solid curves) and $y_p = 1$ (dashed curves). In both cases the log error penalizes deviation from the true label value to a greater extent than the squared error.

Cross-entropy Cost

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

logistic regression classifier
using cross-entropy

(CONVEX)

$$\begin{aligned} & \nabla_{\tilde{w}} \left(y_p \log(\sigma(\tilde{x}_p^T \tilde{w})) \right) \\ &= y_p \underbrace{\frac{1}{\sigma(\tilde{x}_p^T \tilde{w})}}_{\text{Scalar}} \cdot \cancel{\sigma(\cdot)} [1-\sigma(\cdot)] \cdot \underbrace{\nabla_{\tilde{w}} (\tilde{x}_p^T \tilde{w})}_{\tilde{x}_p} \end{aligned}$$

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

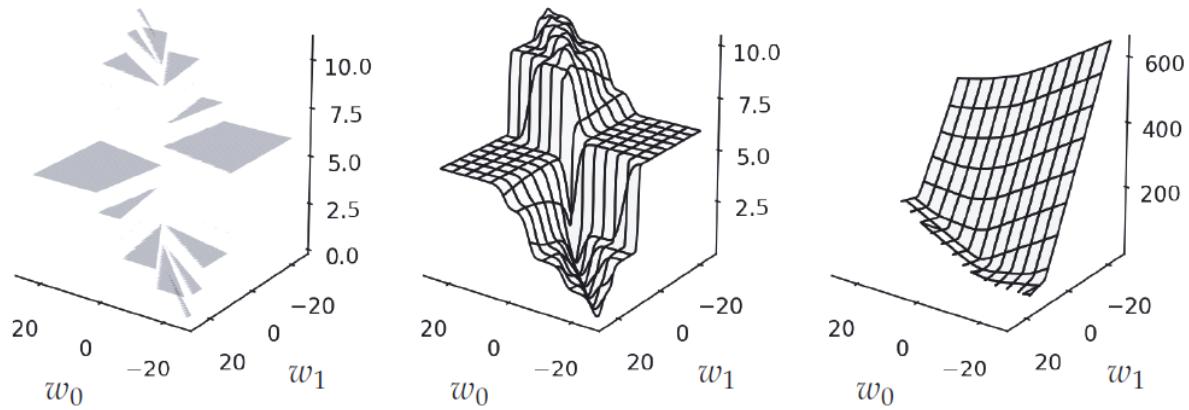


Figure 6.3 Figure associated with Example 6.2. See text for details.

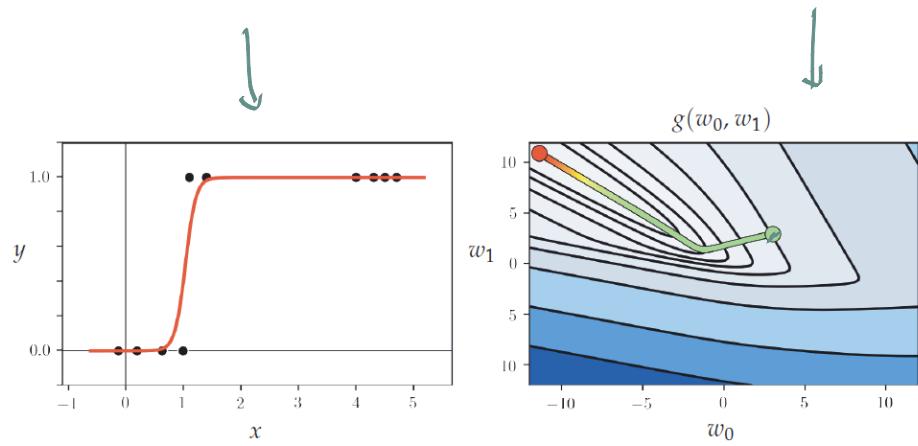


Figure 6.7 Figure associated with Example 6.4. See text for details.