

Week 2: Linear Models

Matthew Caldwell

COMP0088 Introduction to Machine Learning • UCL Computer Science • Autumn 2025

Admin

- Labs
- Triage

Week 2 Recap

Drawing the Line

$$\theta^* = \operatorname*{argmin}_{\theta} L(f, \theta, \{X[, Y]\})$$

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} L(\mathbf{f}, \mathbf{w}, \{\mathbf{X}, \mathbf{Y}\})$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} L(\mathbf{f}, \mathbf{w}, \{\mathbf{X}, \mathbf{Y}\})$$

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} L(\mathbf{f}, \mathbf{w}, \mathbf{X}, \mathbf{y})$$

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} L(\mathbf{w}, \mathbf{X}, \mathbf{y})$$

$$X = \begin{bmatrix} x_1^\top \\ x_2^\top \\ \vdots \\ \vdots \\ x_n^\top \end{bmatrix}$$

$$X = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,d} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,d} \end{bmatrix}$$

↓ Feature dimensions ↓

$$X = \begin{bmatrix} x_1^\top \\ x_2^\top \\ \vdots \\ x_n^\top \end{bmatrix} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,d} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,d} \end{bmatrix}$$

↑ Samples ↑

“Design Matrix”

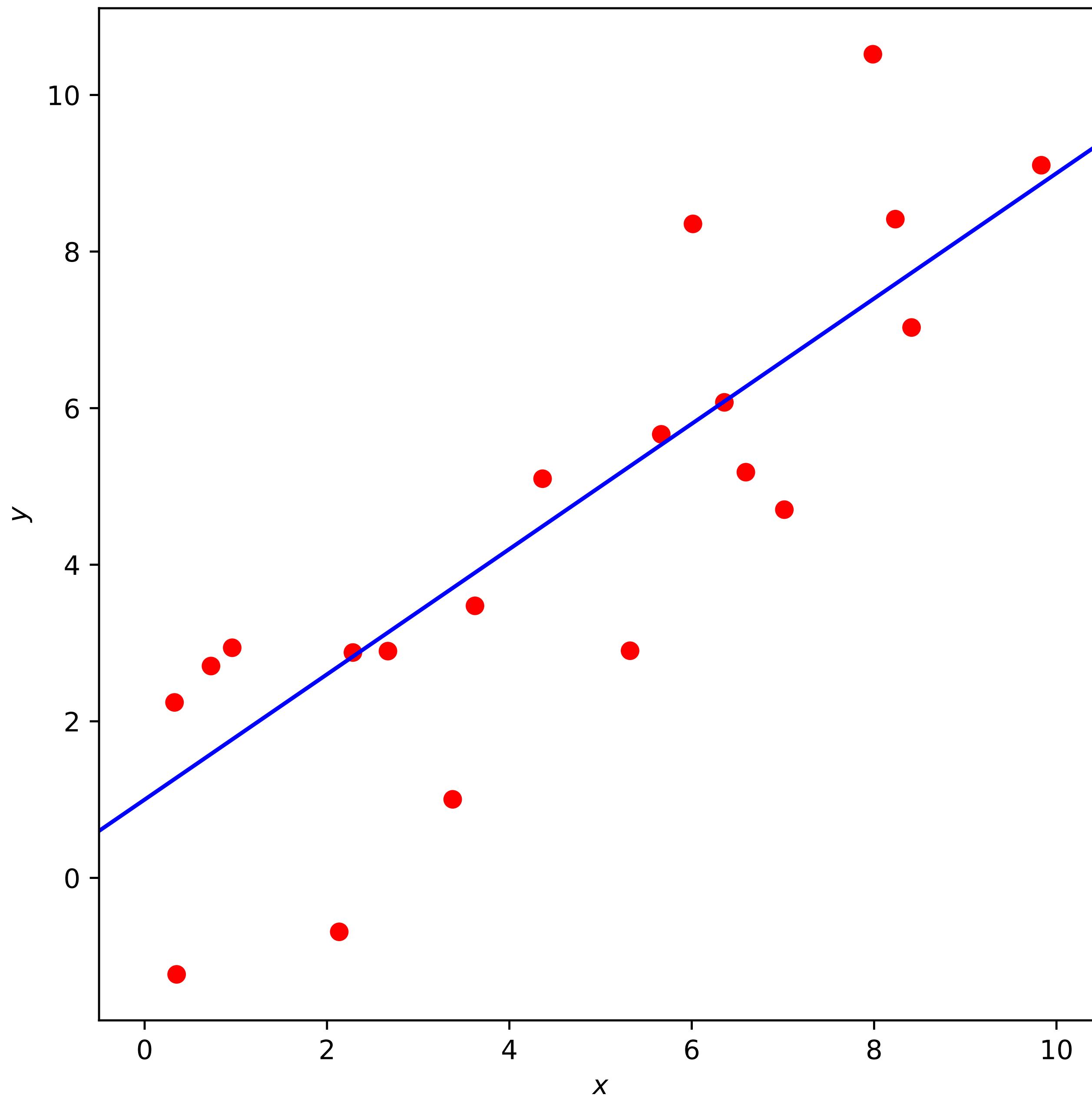


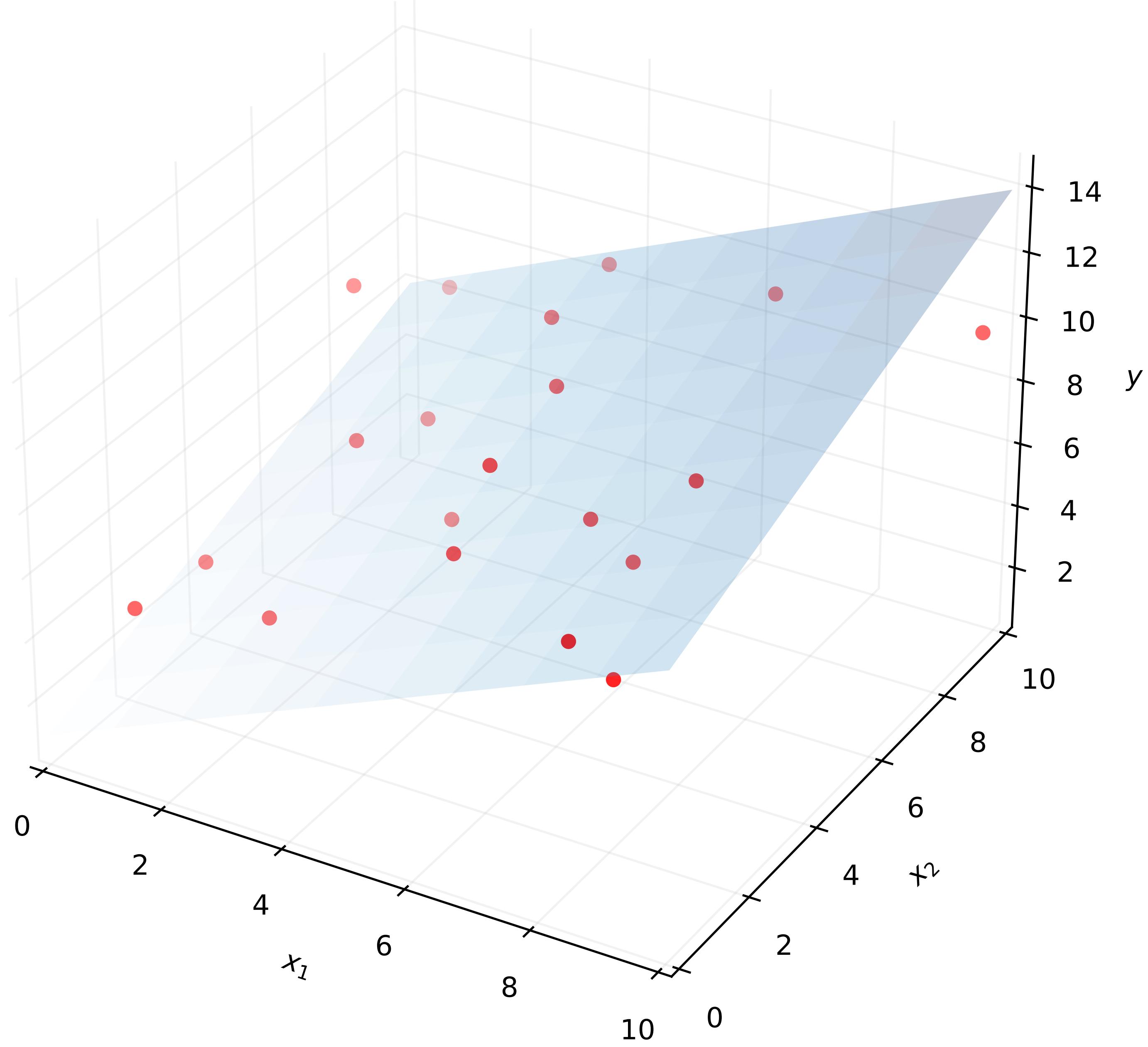
Meerkats

#	Age	Weight	Sex	Captive	...
1	11.8	815	1	0	
2	9.7	672	1	0	
3	8.9	446	1	0	
4	10.8	761	0	0	
5	8.3	1035	0	1	
6	11.7	930	1	1	
7	8.5	1027	0	1	
8	7.6	1234	0	1	
9	15.5	1461	0	1	
10	8.8	720	1	0	
11	12.8	1223	0	1	
12	7.7	711	1	0	
13	7.9	586			

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

\hat{y} = xw





Linear models are the best models!

\hat{y} = xw

**Comprehensible
Interpretable
Tractable
Versatile**

w =

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

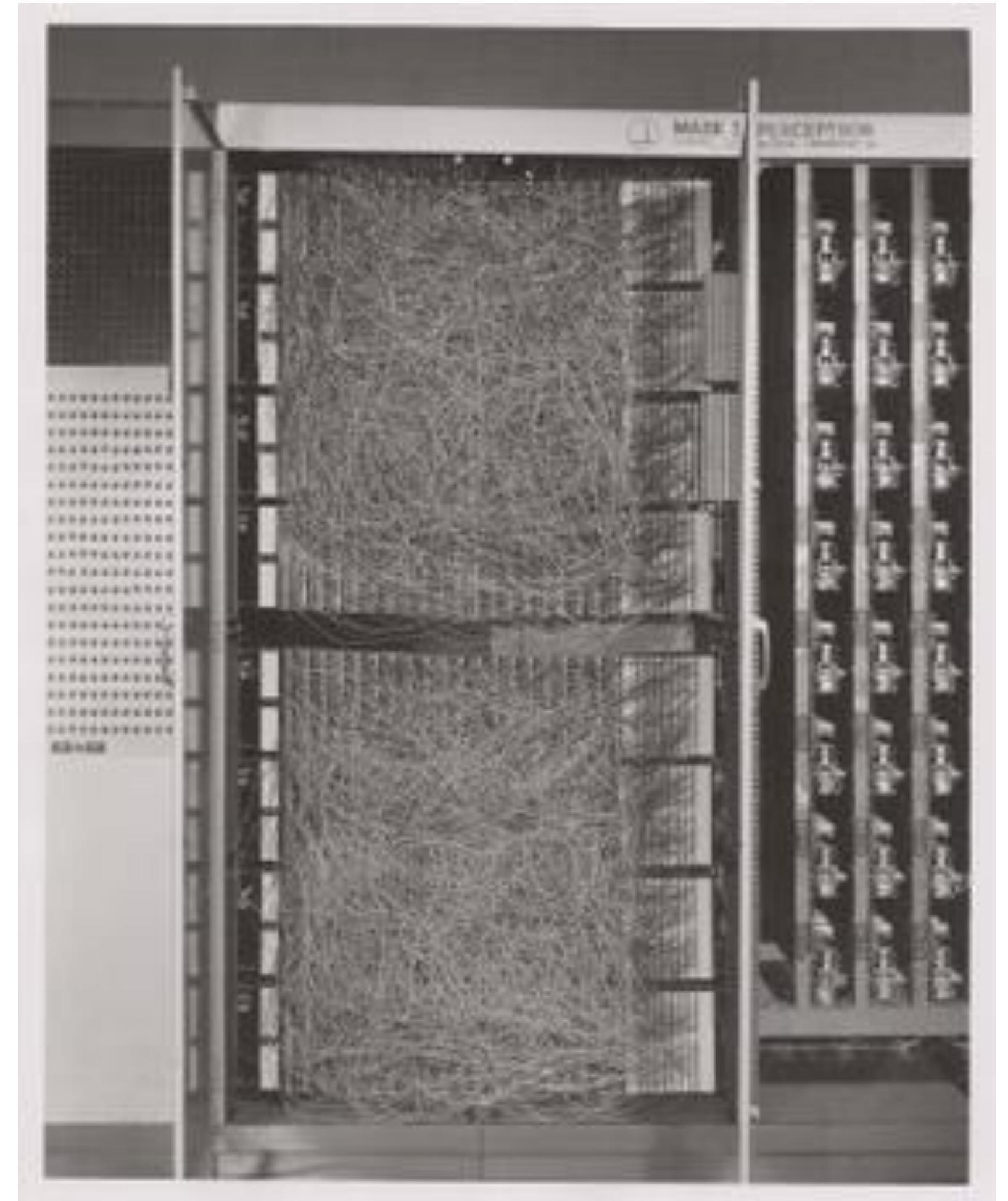
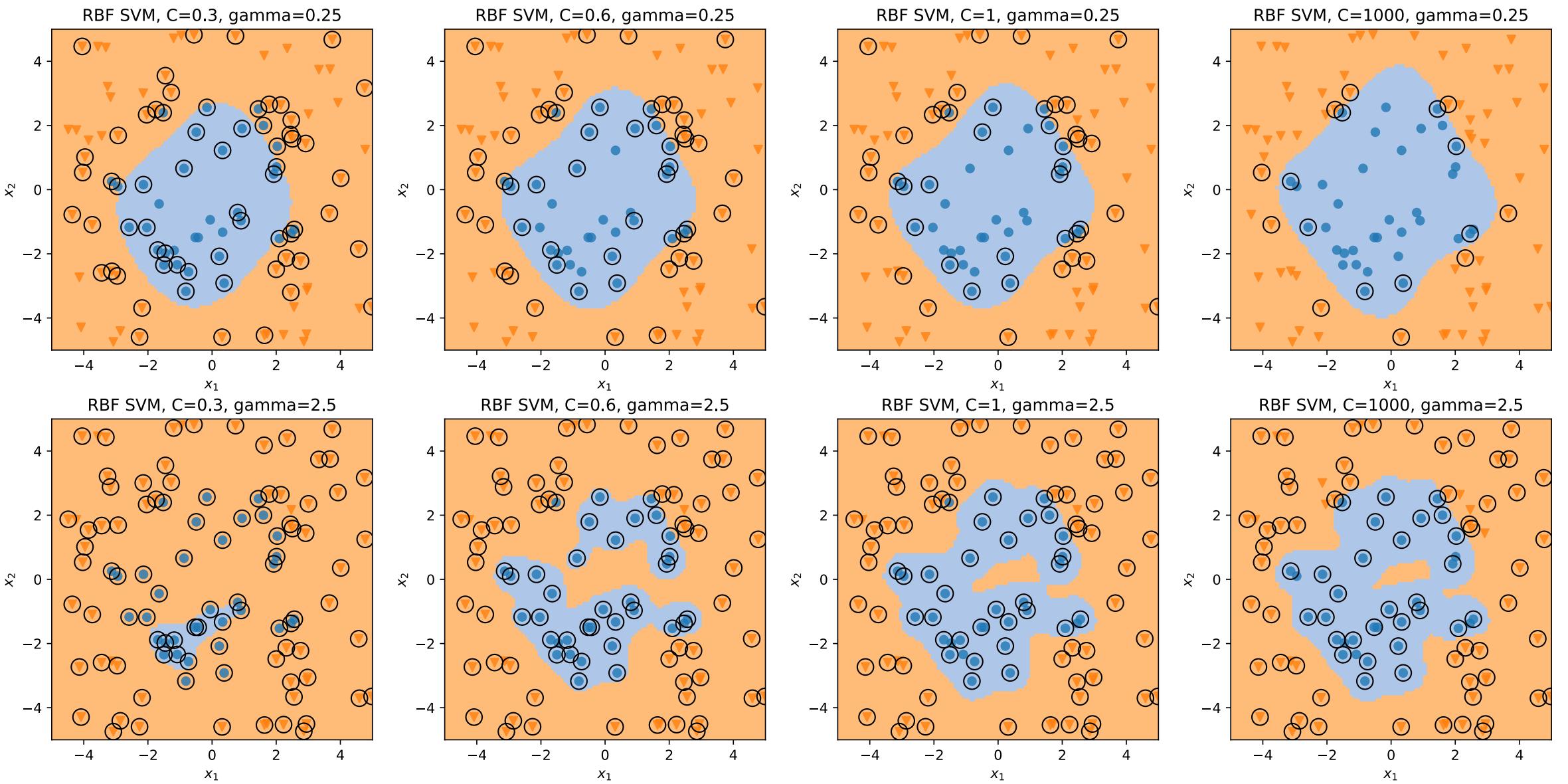
$$\mathbf{w} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{w} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

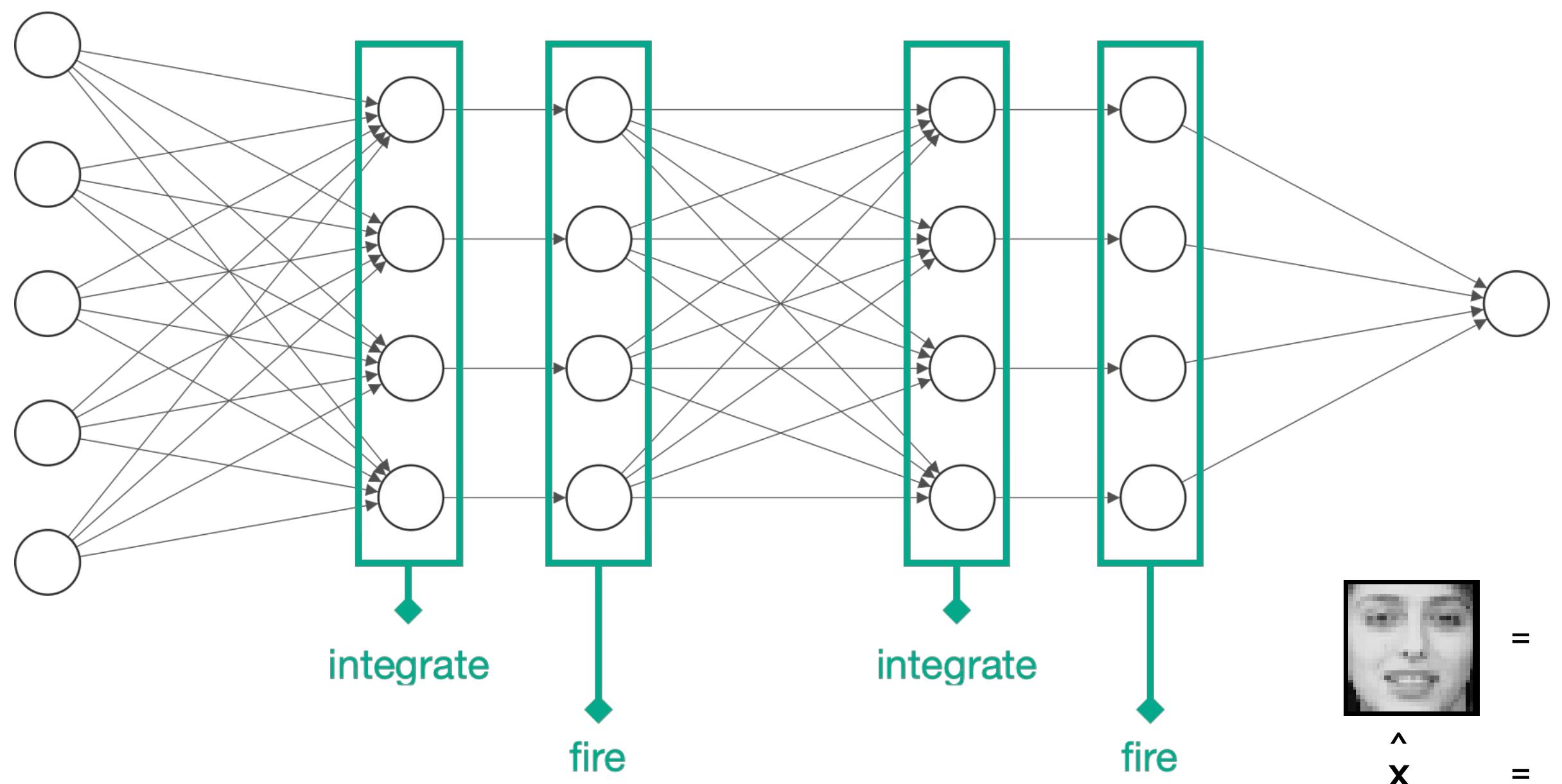
$$\mathbf{w} = \begin{bmatrix} 0.25 \\ 0.25 \\ 0.25 \\ 0.25 \end{bmatrix}$$

w =

[?
? .
?
?]

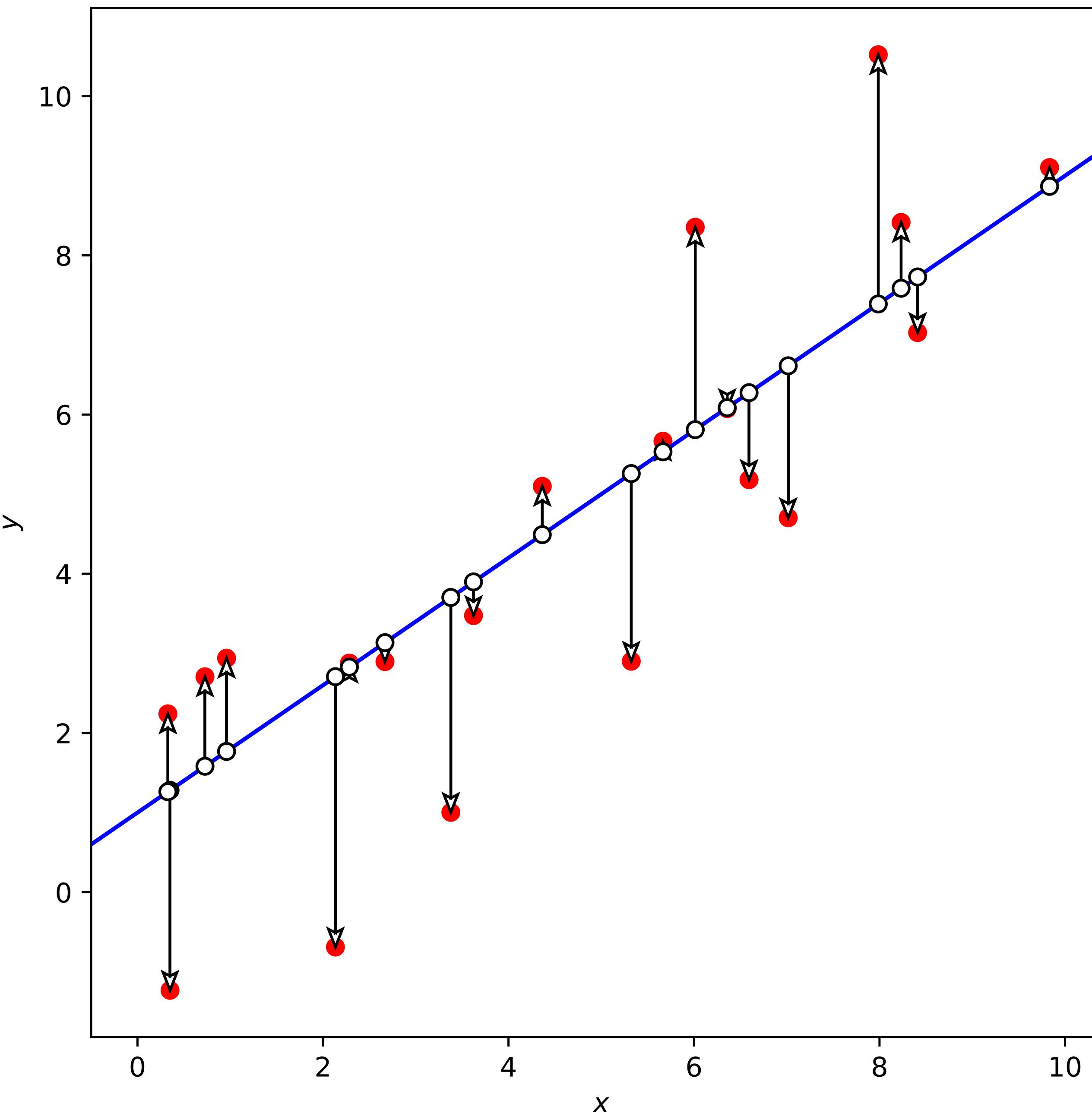


$$\hat{x} = \mu + w_1 u_1 + w_2 u_2 + w_3 u_3 + w_4 u_4 + \dots$$

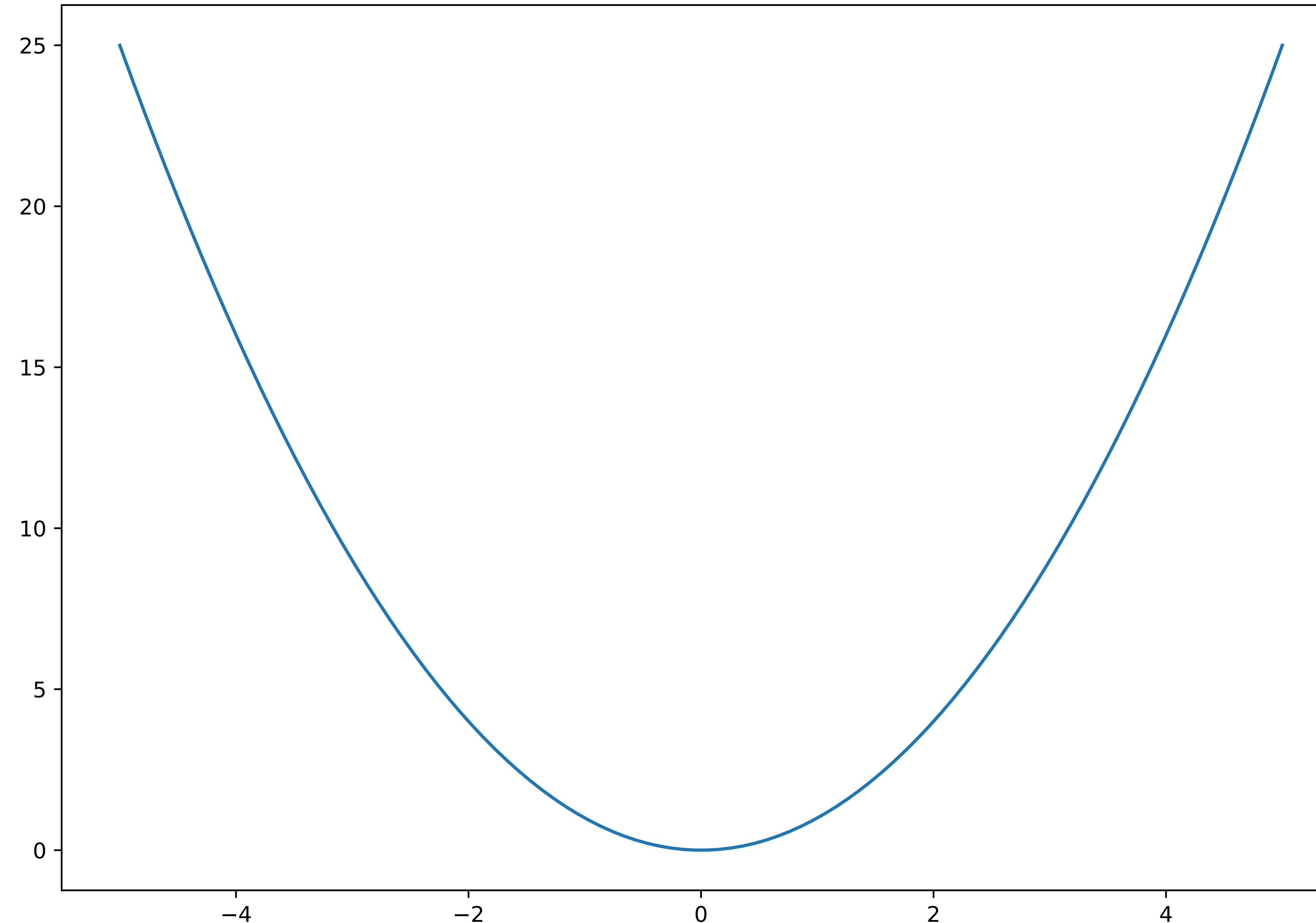


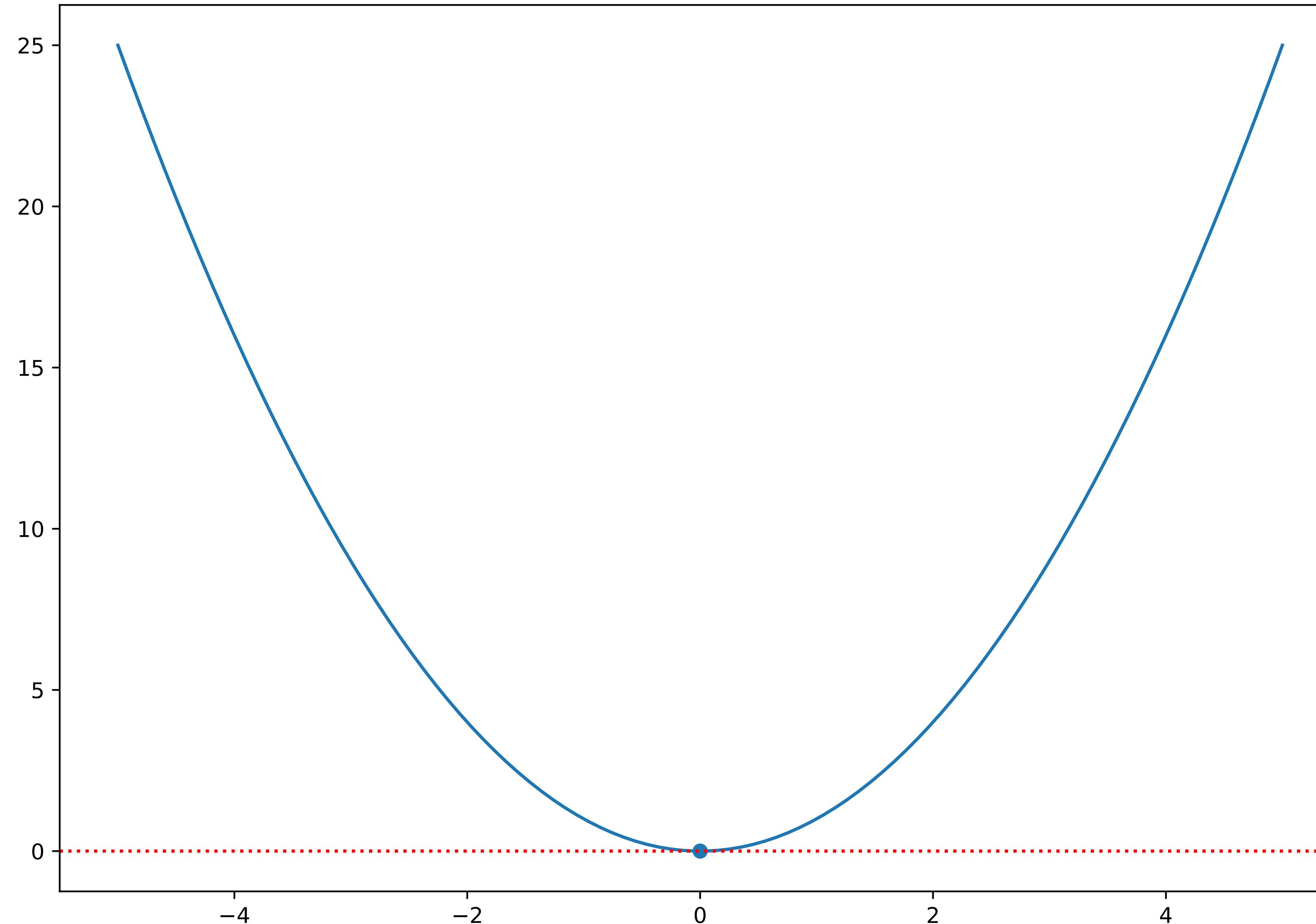
$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} L(\mathbf{w}, \mathbf{X}, \mathbf{y})$$

\hat{y} = xw



$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$





$$\frac{dy}{dx} = 0$$

$$\frac{dL}{d\mathbf{w}} = 0$$

$$\nabla_{\mathbf{w}} L = 0$$

$$\nabla_{\mathbf{w}} L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \\ \frac{\partial L}{\partial w_d} \end{bmatrix}$$

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

$$\nabla_w \left\| Xw - y \right\|^2 = 0$$

$$\begin{aligned}
\|\mathbf{Xw} - \mathbf{y}\|^2 &= (\mathbf{Xw} - \mathbf{y})^\top (\mathbf{Xw} - \mathbf{y}) \\
&= (\mathbf{Xw})^\top \mathbf{Xw} - (\mathbf{Xw})^\top \mathbf{y} - \mathbf{y}^\top (\mathbf{Xw}) + \mathbf{y}^\top \mathbf{y} \\
&= (\mathbf{Xw})^\top \mathbf{Xw} - 2\mathbf{y}^\top (\mathbf{Xw}) + \mathbf{y}^\top \mathbf{y} \\
&= (\mathbf{Xw})^\top \mathbf{Xw} - 2(\mathbf{y}^\top \mathbf{X})\mathbf{w} + \mathbf{y}^\top \mathbf{y} \\
&= \mathbf{w}^\top \mathbf{X}^\top \mathbf{Xw} - 2(\mathbf{X}^\top \mathbf{y})^\top \mathbf{w} + \mathbf{y}^\top \mathbf{y} \\
&= \mathbf{w}^\top (\mathbf{X}^\top \mathbf{X})\mathbf{w} - 2(\mathbf{X}^\top \mathbf{y})^\top \mathbf{w} + \mathbf{y}^\top \mathbf{y}
\end{aligned}$$

$$\begin{aligned}
\nabla_{\mathbf{w}} L &= \nabla_{\mathbf{w}} \left(\mathbf{w}^T (\mathbf{X}^T \mathbf{X}) \mathbf{w} - 2(\mathbf{X}^T \mathbf{y})^T \mathbf{w} + \mathbf{y}^T \mathbf{y} \right) \\
&= \nabla_{\mathbf{w}} \left(\mathbf{w}^T (\mathbf{X}^T \mathbf{X}) \mathbf{w} - 2(\mathbf{X}^T \mathbf{y})^T \mathbf{w} \right) \\
&= 2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y}
\end{aligned}$$

$$\mathbf{X}^\top \mathbf{X} \mathbf{w}^* = \mathbf{X}^\top \mathbf{y}$$

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

“Normal Equations”

$$(X^T X + \lambda I) w^* = X^T y$$

$$w^* = (X^T X + \lambda I)^{-1} X^T y$$

“Revised Normal Equations”

$$(X^T X + \lambda I) w^* = X^T y$$

$$w^* = (X^T X + \lambda I)^{-1} X^T y$$

“Revised Normal Equations”



$$y = f(x, \theta)$$

$$y \equiv f_{\theta}(x)$$

$$\theta^* = \operatorname*{argmin}_{\theta} L(\mathbf{f}, \theta, \{\mathbf{X}, \mathbf{Y}\})$$

$$\theta^* = \operatorname*{argmin}_{\theta} L_{f,X,Y}(\theta)$$

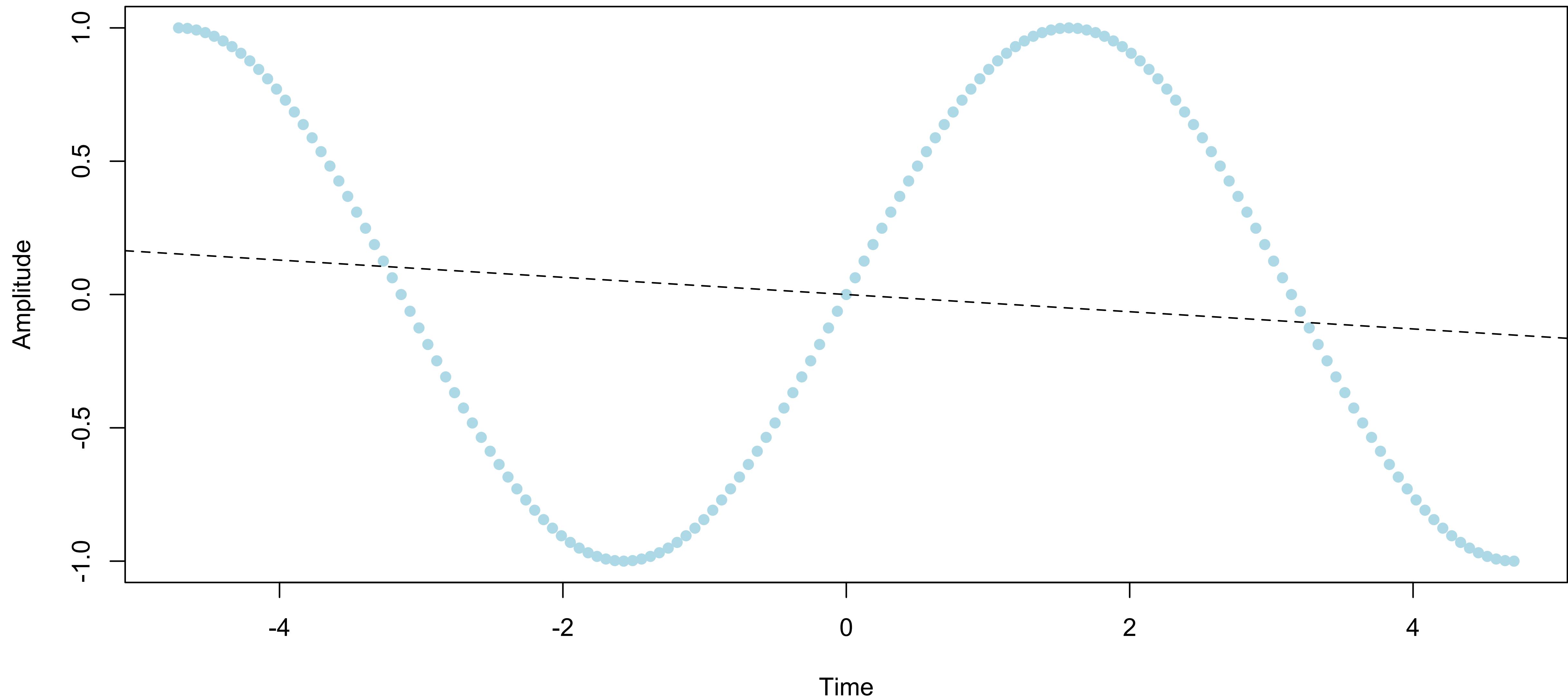
$$P(y|x; \theta) = f(x, y, \theta)$$

$$L(\theta) = f(x, y, \theta)$$

$$\theta^* = \operatorname*{argmax}_{\theta} L_{f,X,Y}(\theta)$$

training \neq test

**Given the training set, loss (or likelihood)
is solely a function of the parameters**



Don't fit a straight line to curvy data



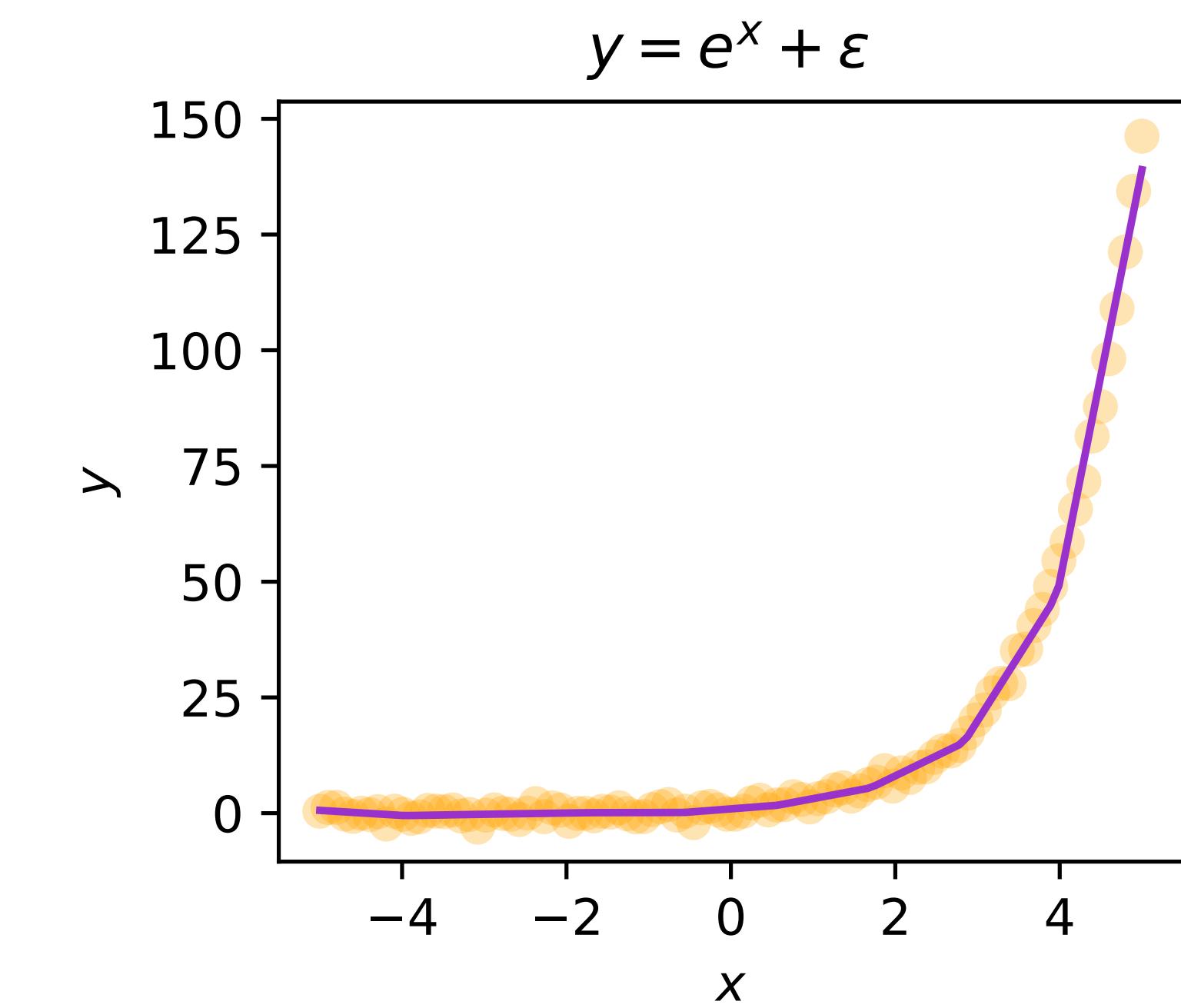
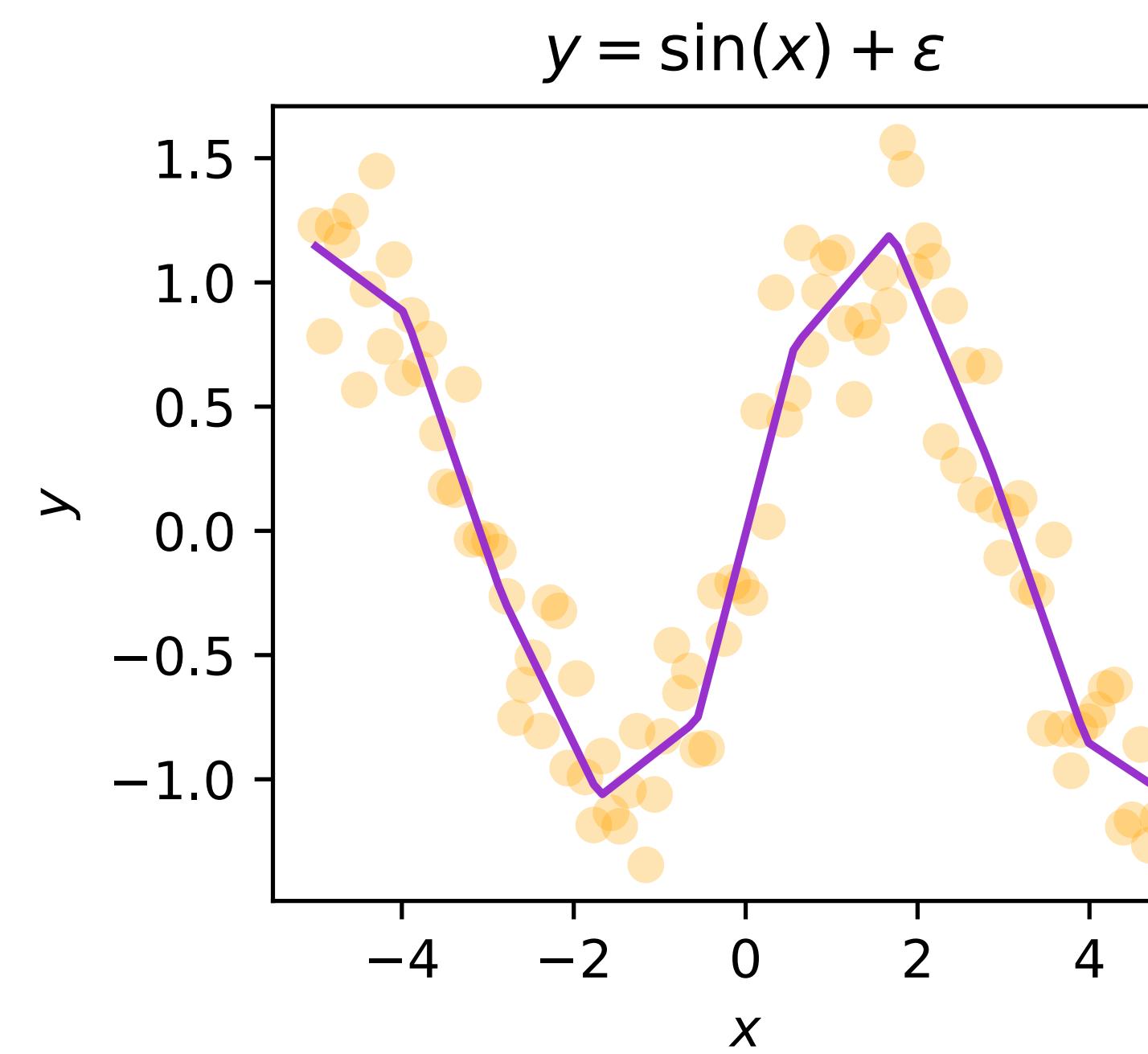
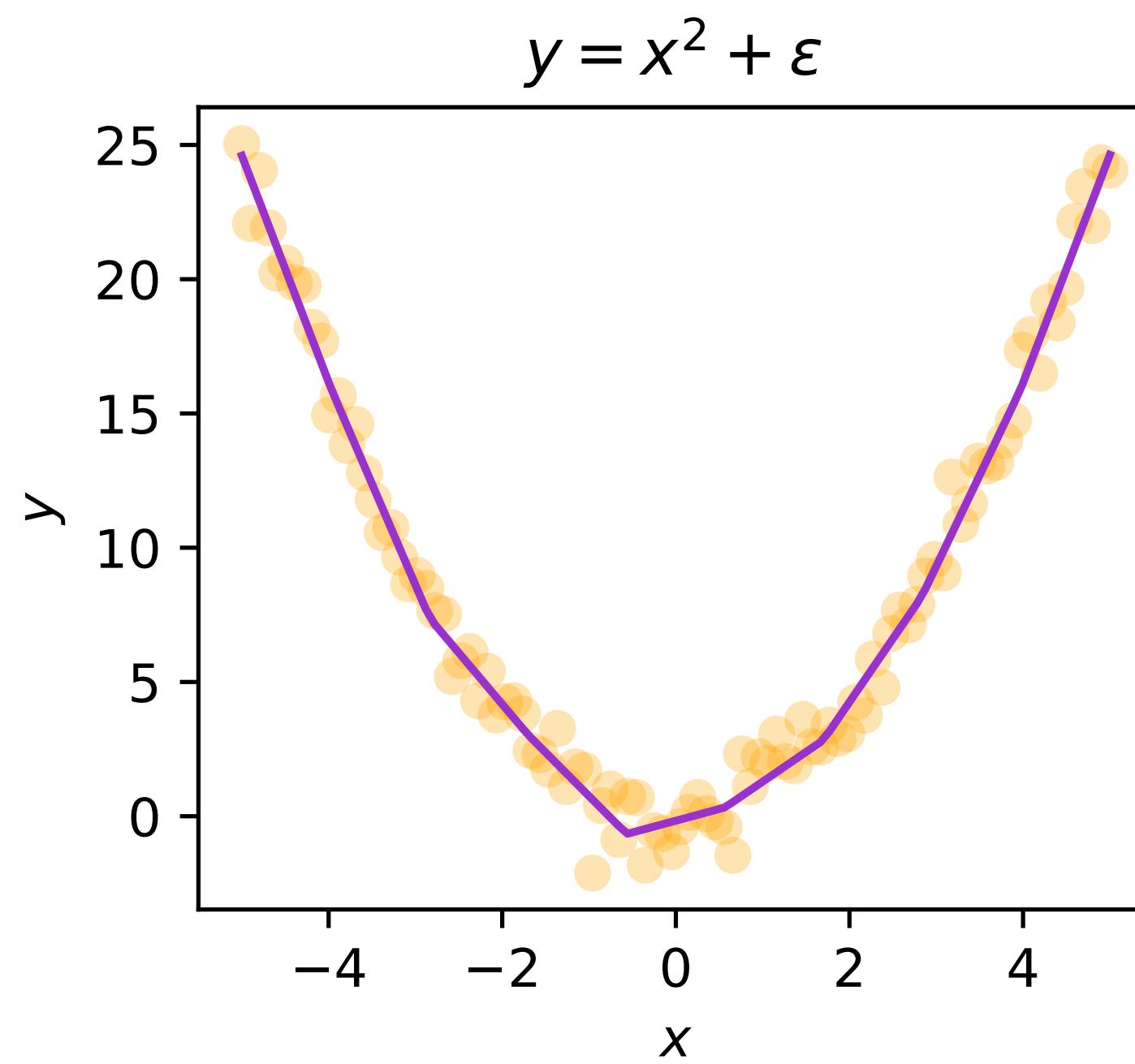
$$\hat{y} = \sum_i^k w_i h_i(\mathbf{x}) = \mathbf{h}(\mathbf{x}) \cdot \mathbf{w} = \mathbf{x}' \cdot \mathbf{w}$$

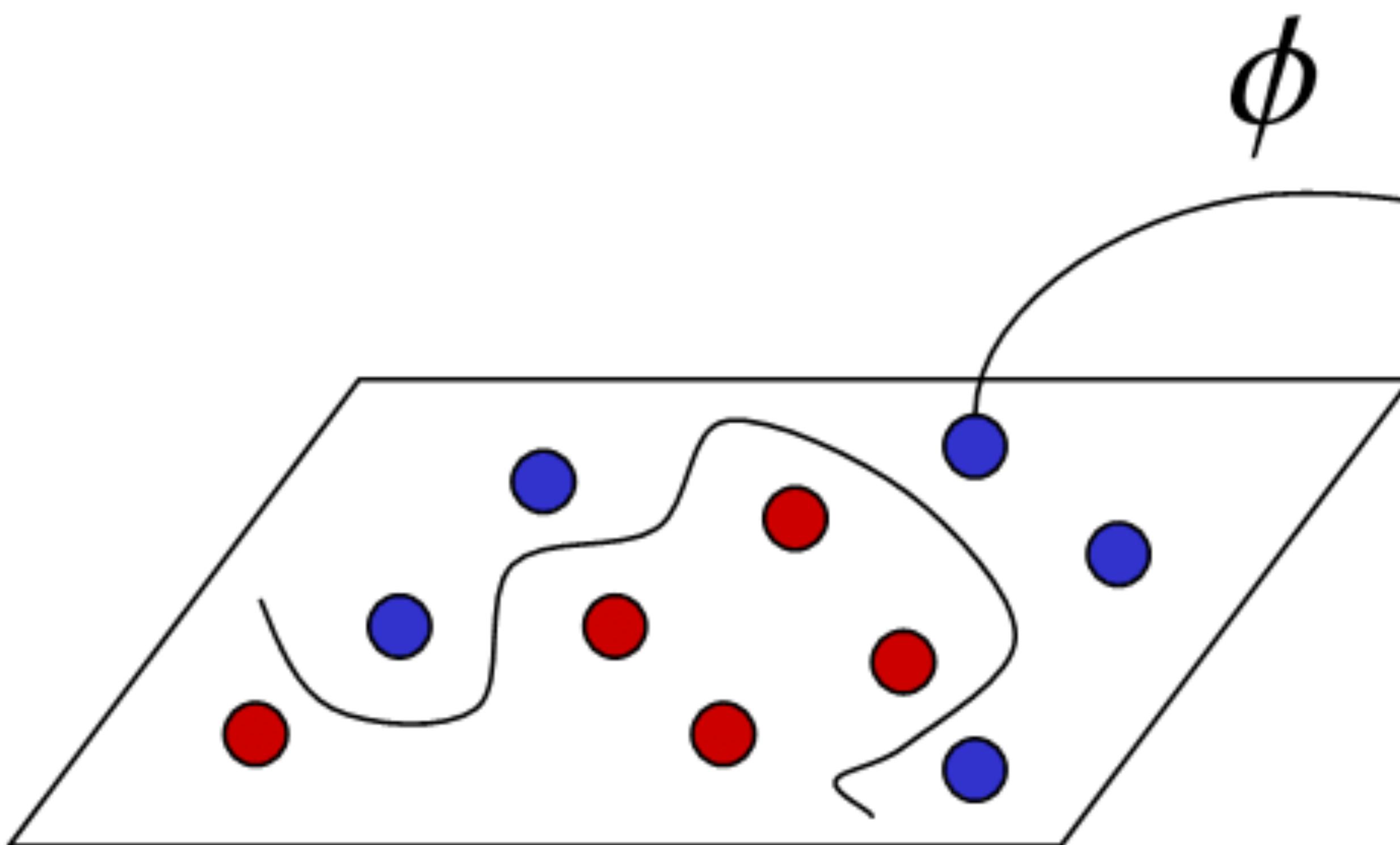
$$\mathbf{x}' = \mathbf{h}(\mathbf{x}) = [h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_k(\mathbf{x})]$$

$$h_1(x) = 1$$

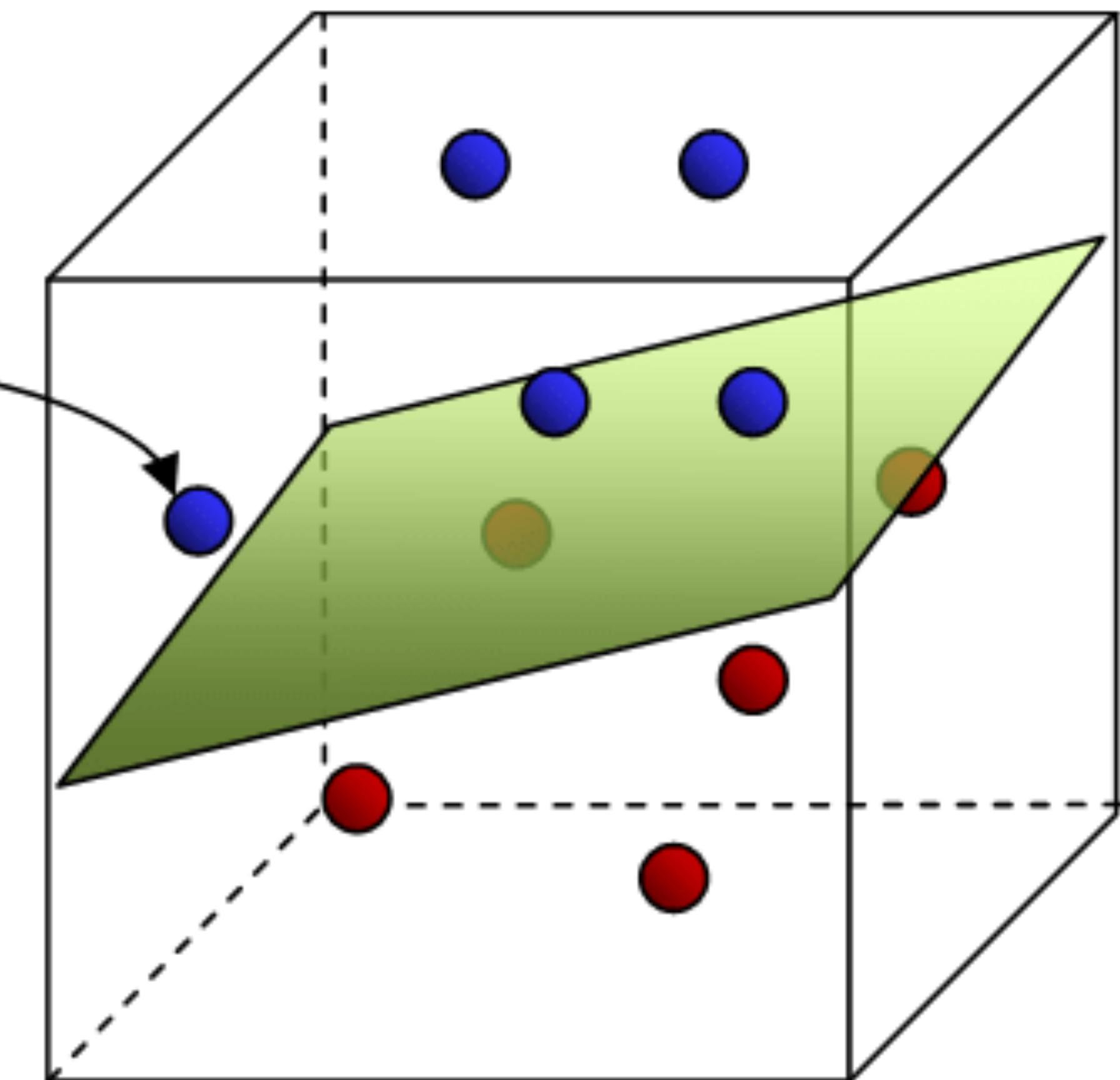
$$h_2(x) = x$$

$$h_{i+2}(x) = \max(0, x - t_i), \quad i \in \{1, 2, \dots, m-1\}$$





a) Input Space



b) Feature Space



DEFEND · THE · CHILDREN · OF · THE ·
POOR · & · PUNISH · THE · WRONGDOER ·

$$L(\mathbf{X}, \mathbf{y}; \theta) = D(f_{\theta}(\mathbf{X}, \mathbf{y}), \mathbf{y}) + \lambda P(\theta)$$

$$L(\mathbf{X}, \mathbf{y}; \mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda P(\mathbf{w})$$

$$P(\mathbf{w}) = \|\mathbf{w}\|^2$$

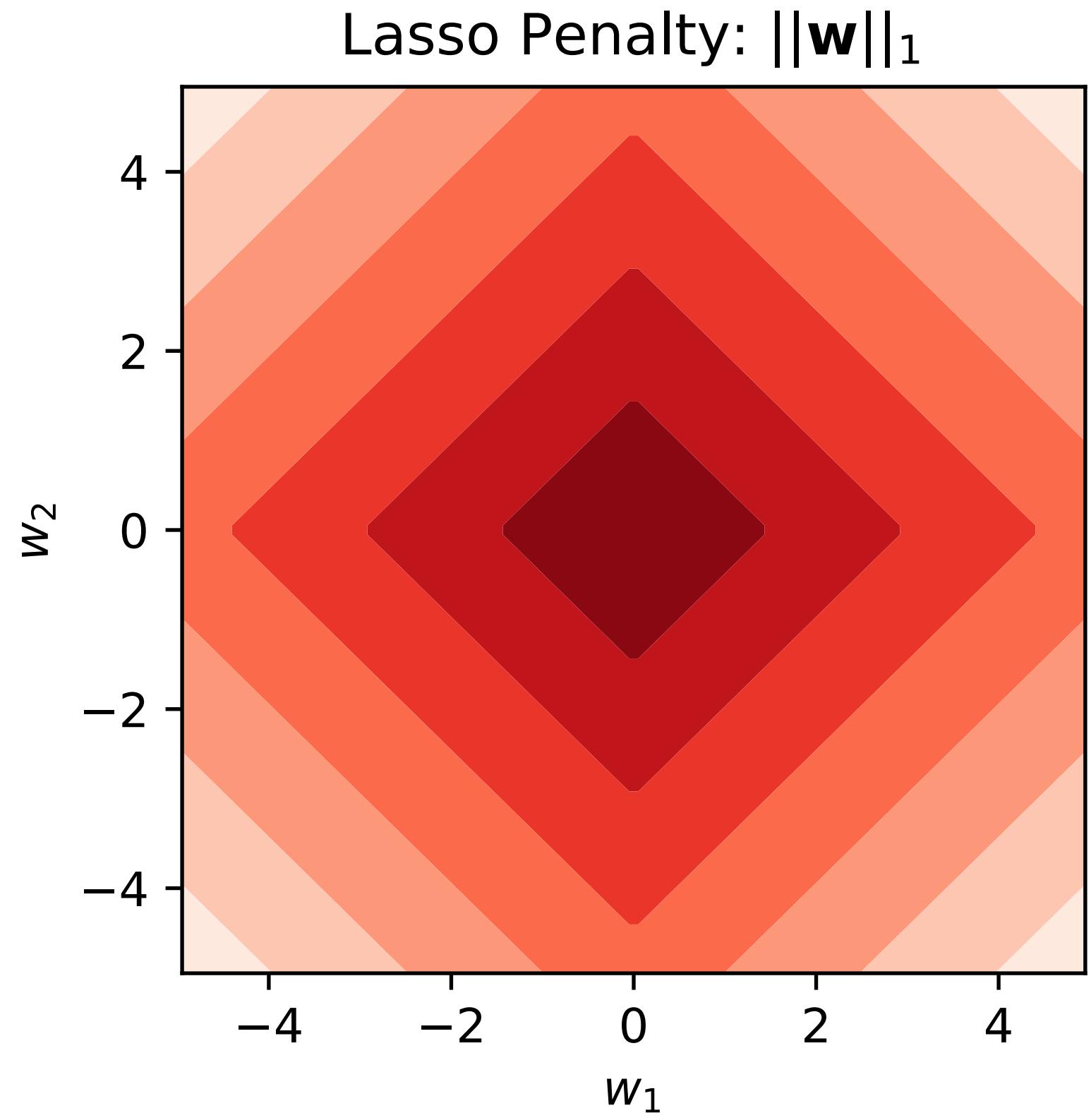
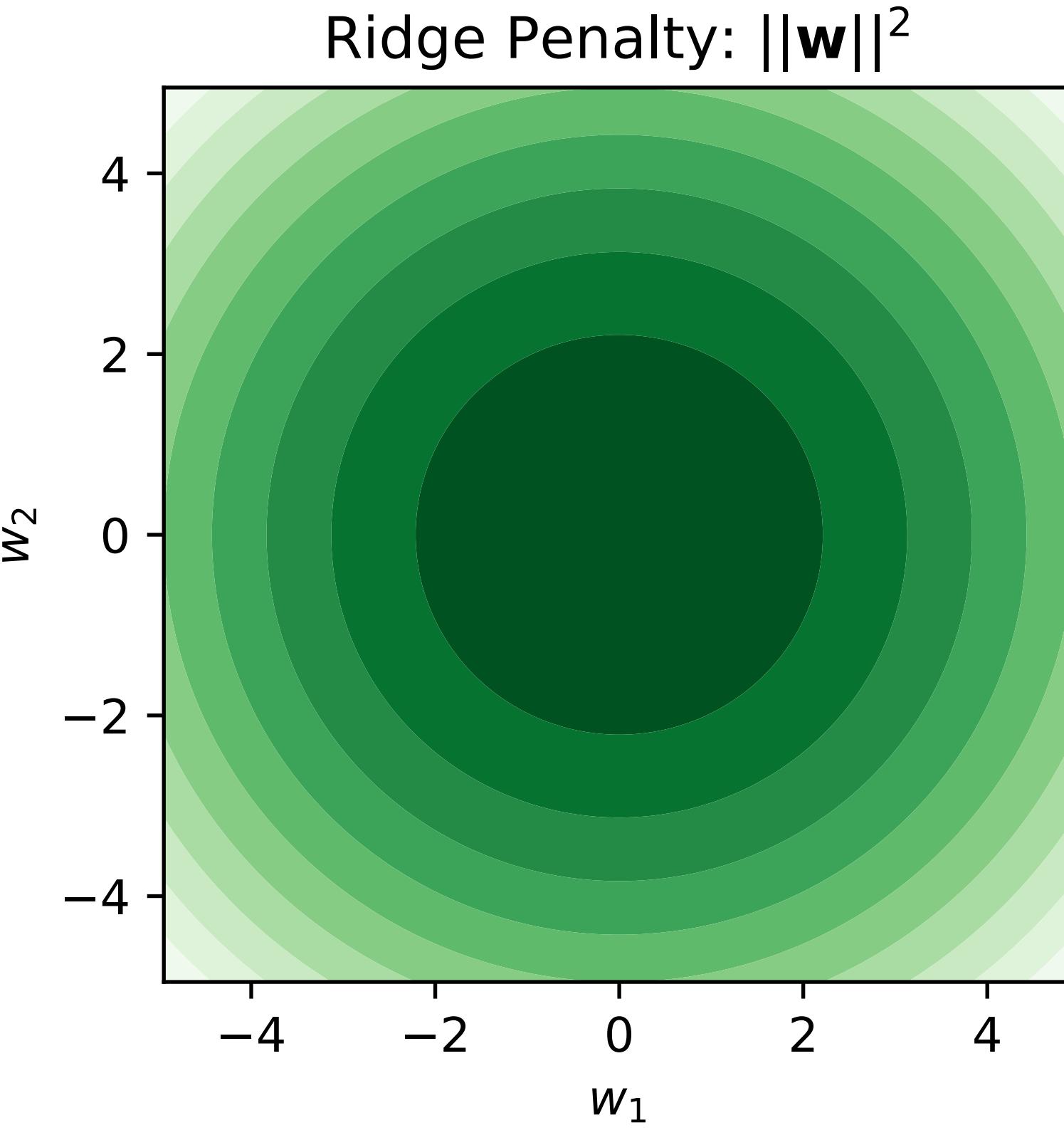
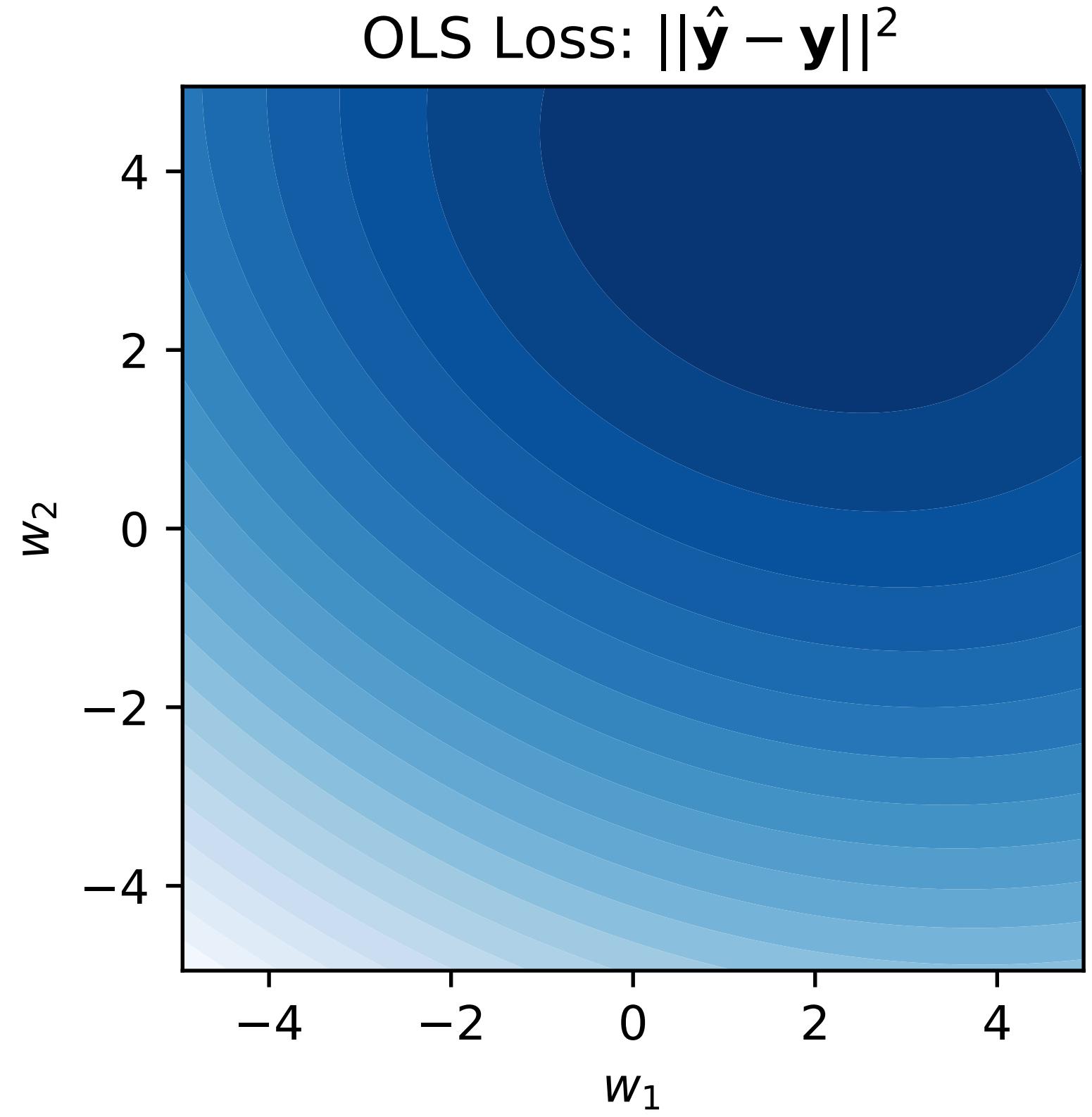
$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \|\mathbf{Xw} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2$$

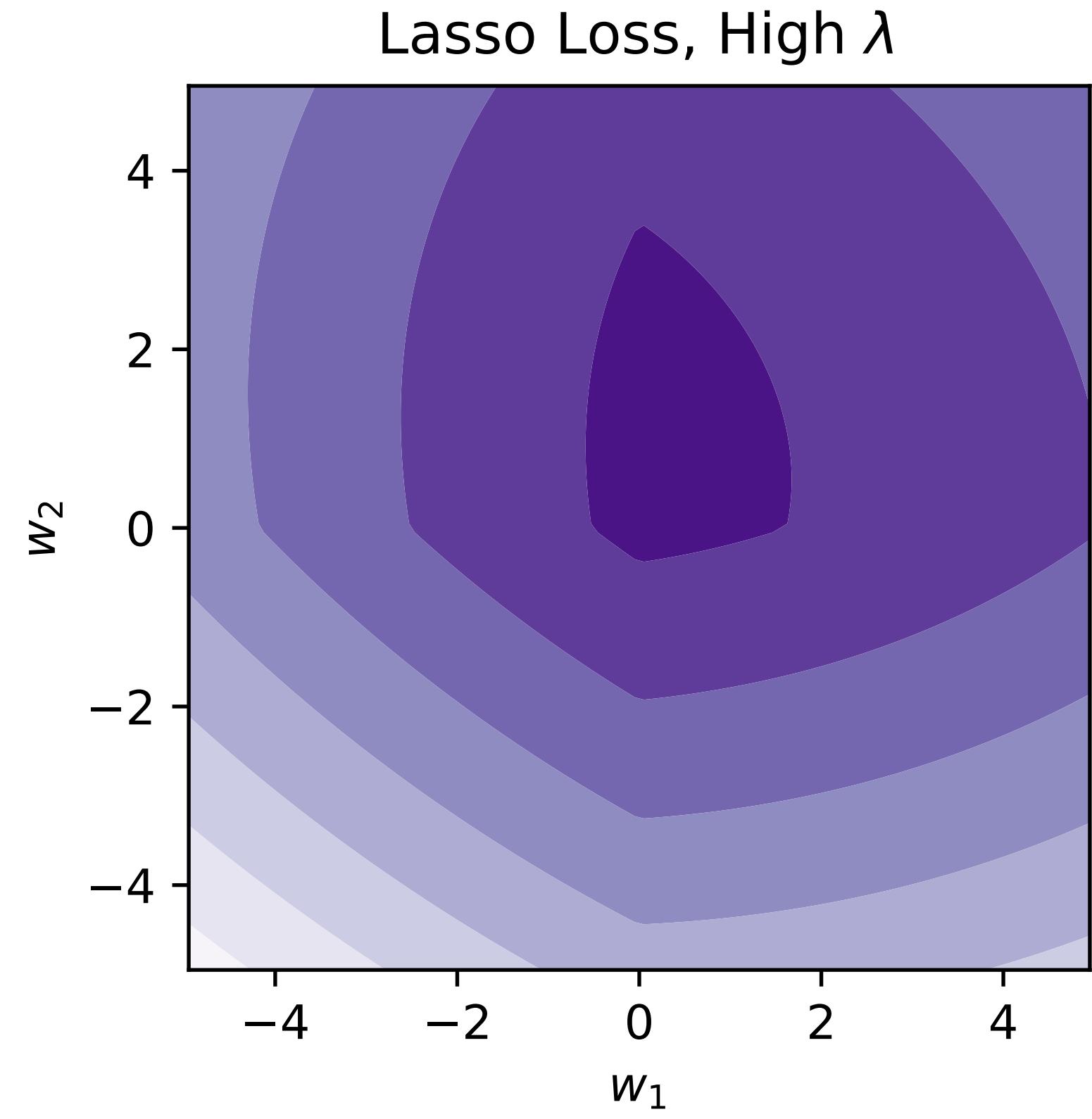
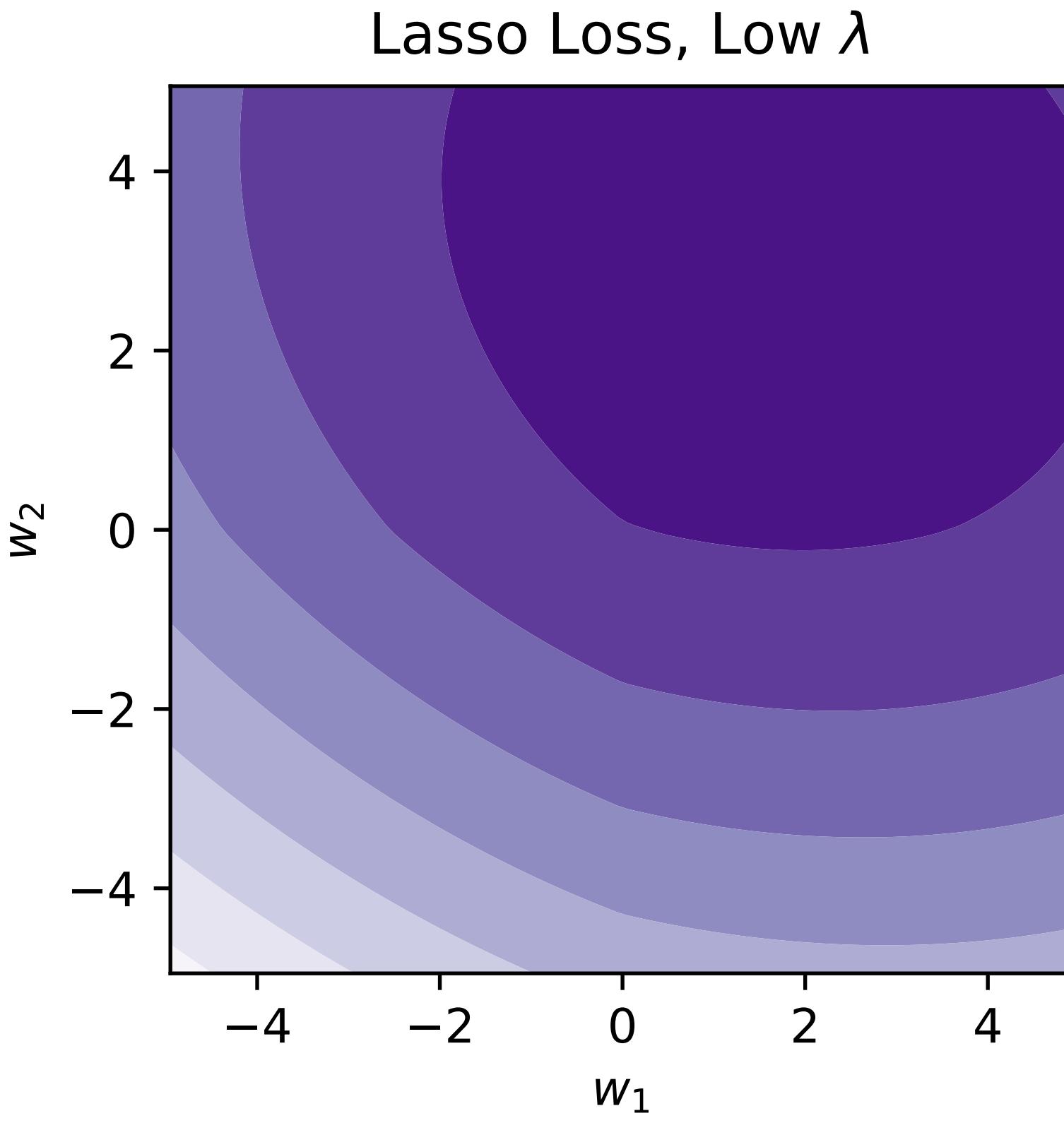
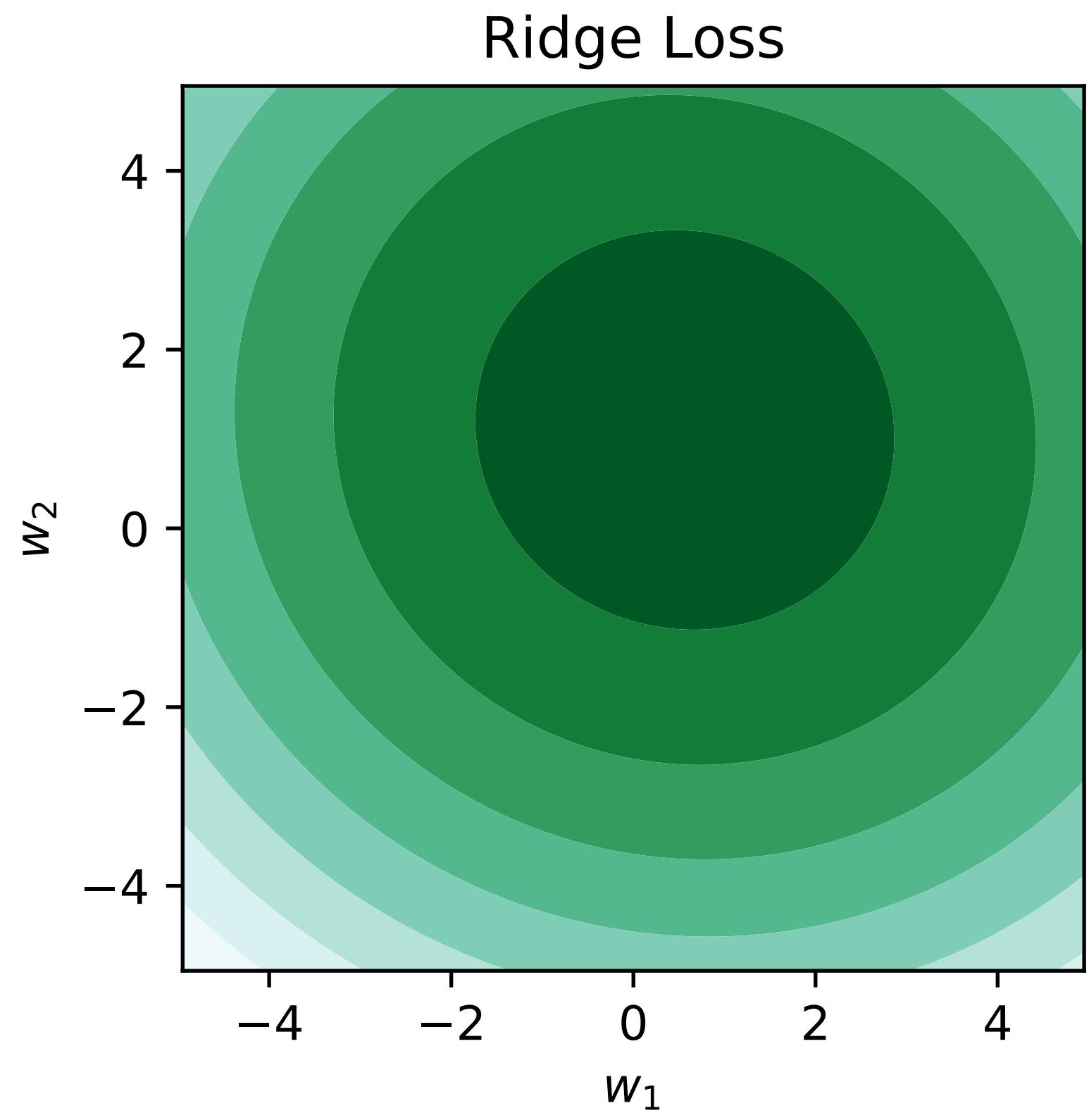
“Ridge Regression”

$$P(\mathbf{w}) = \|\mathbf{w}\|_1$$

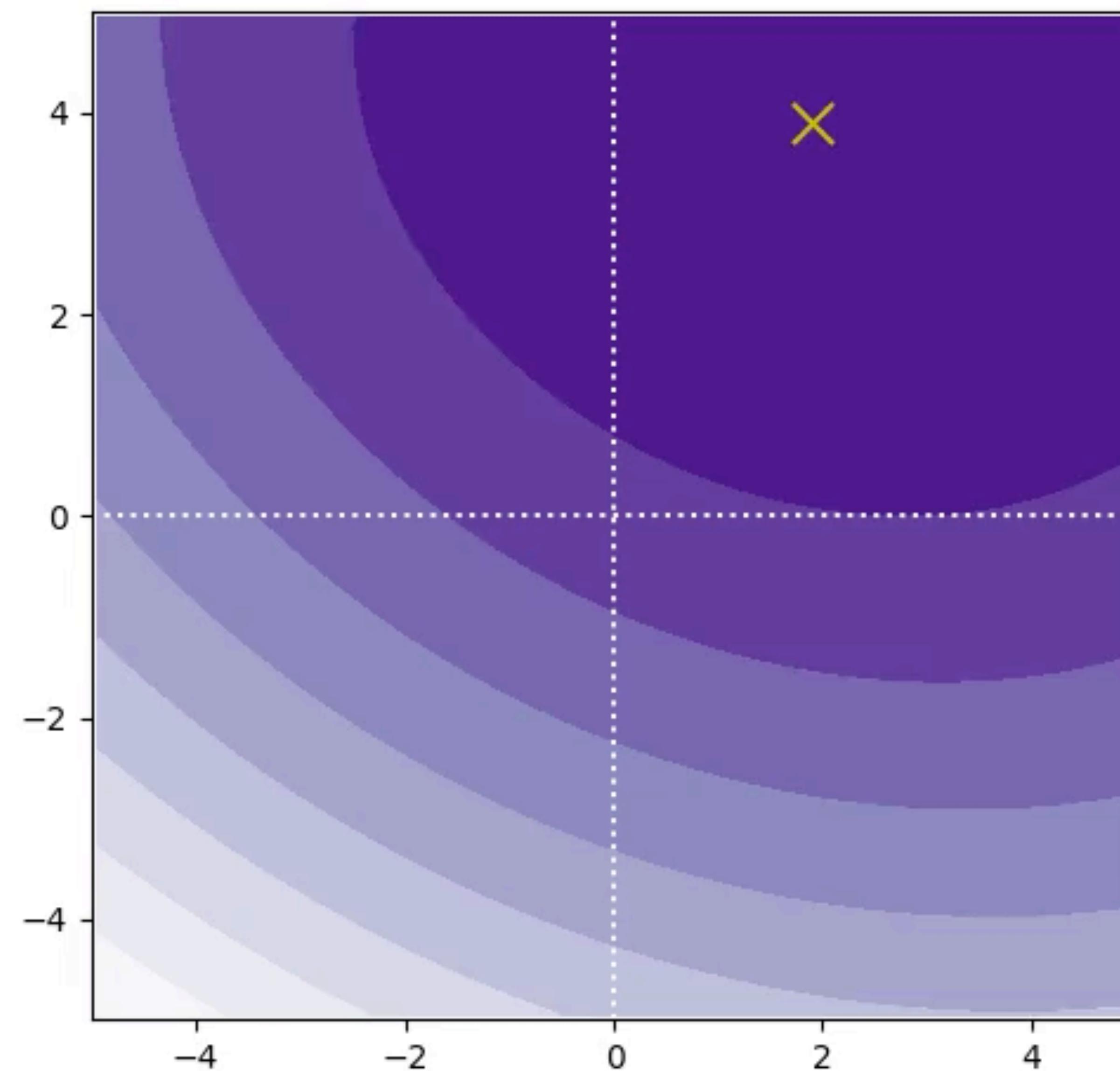
$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \|\mathbf{Xw} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1$$

“Lasso”





Lasso Loss



$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|^2$$

$$\mathbf{w}^* = \operatorname*{argmin}_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1$$





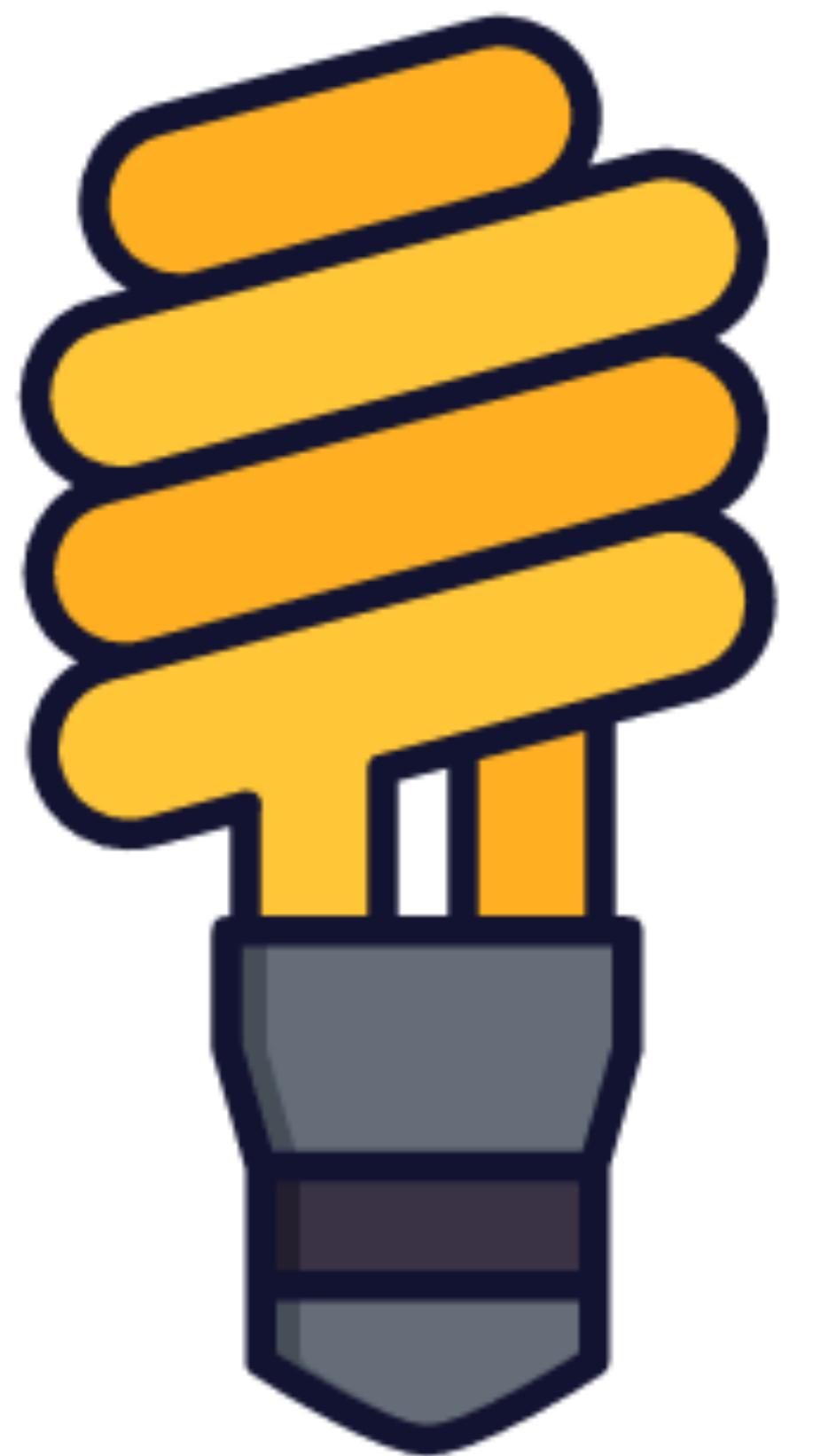
$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} L(\mathbf{f}, \mathbf{X}, \mathbf{y}, \mathbf{w})$$

“Gradient Descent”

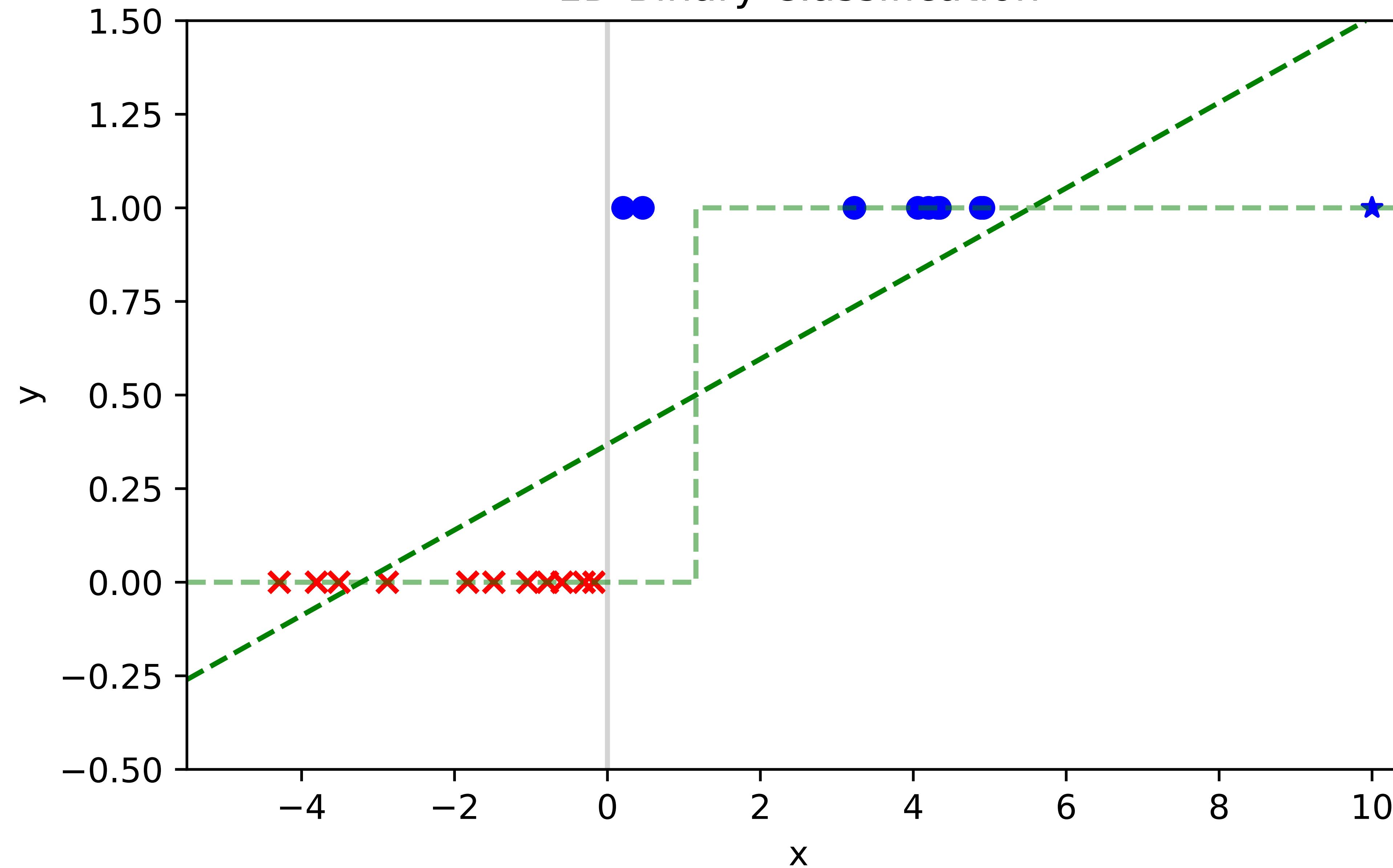
$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} L(\mathbf{f}, \mathbf{X}, \mathbf{y}, \mathbf{w})$$

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla_{\mathbf{w}} L(\mathbf{f}, \mathbf{x}, y, \mathbf{w})$$

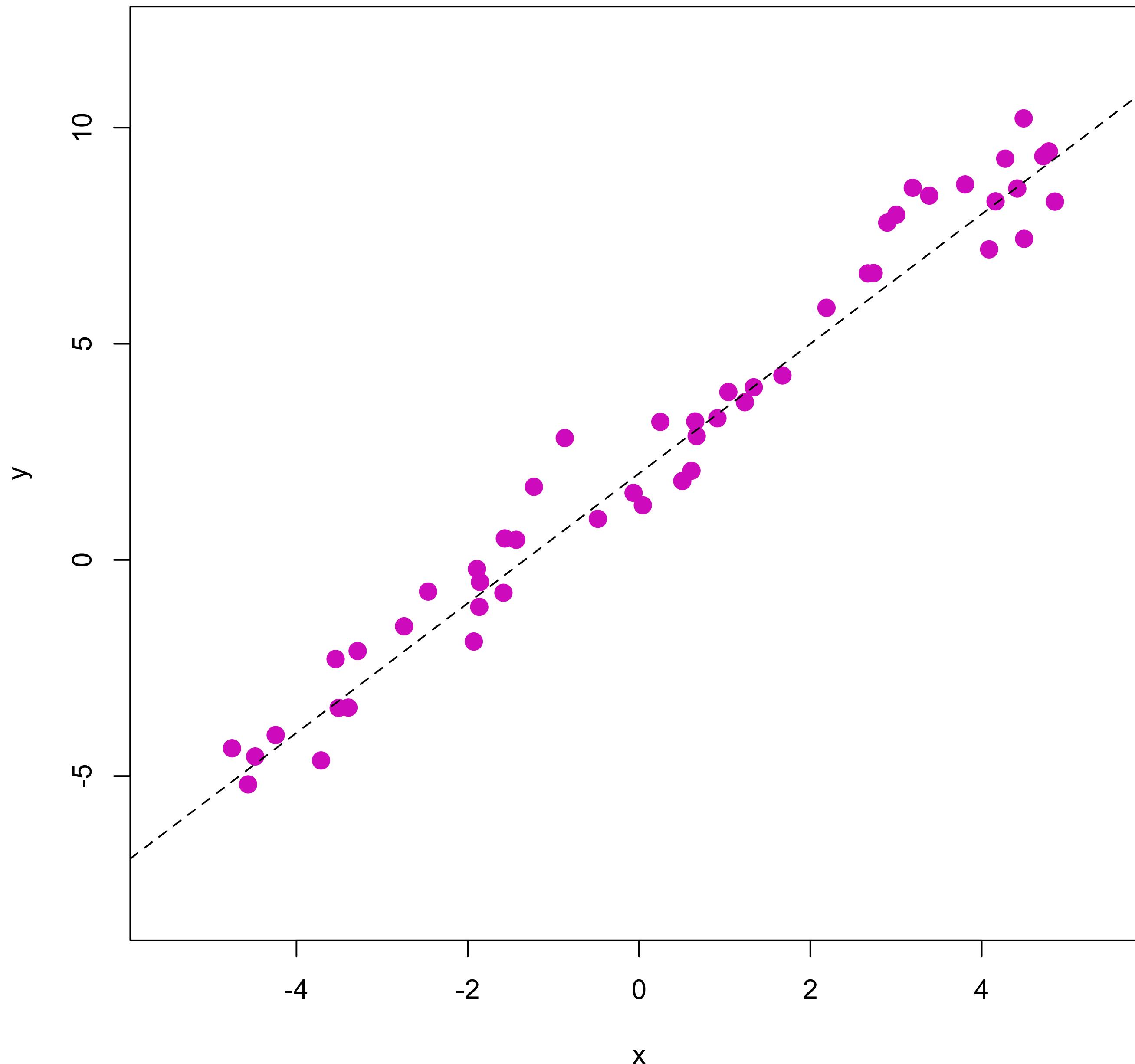
“Stochastic Gradient Descent”



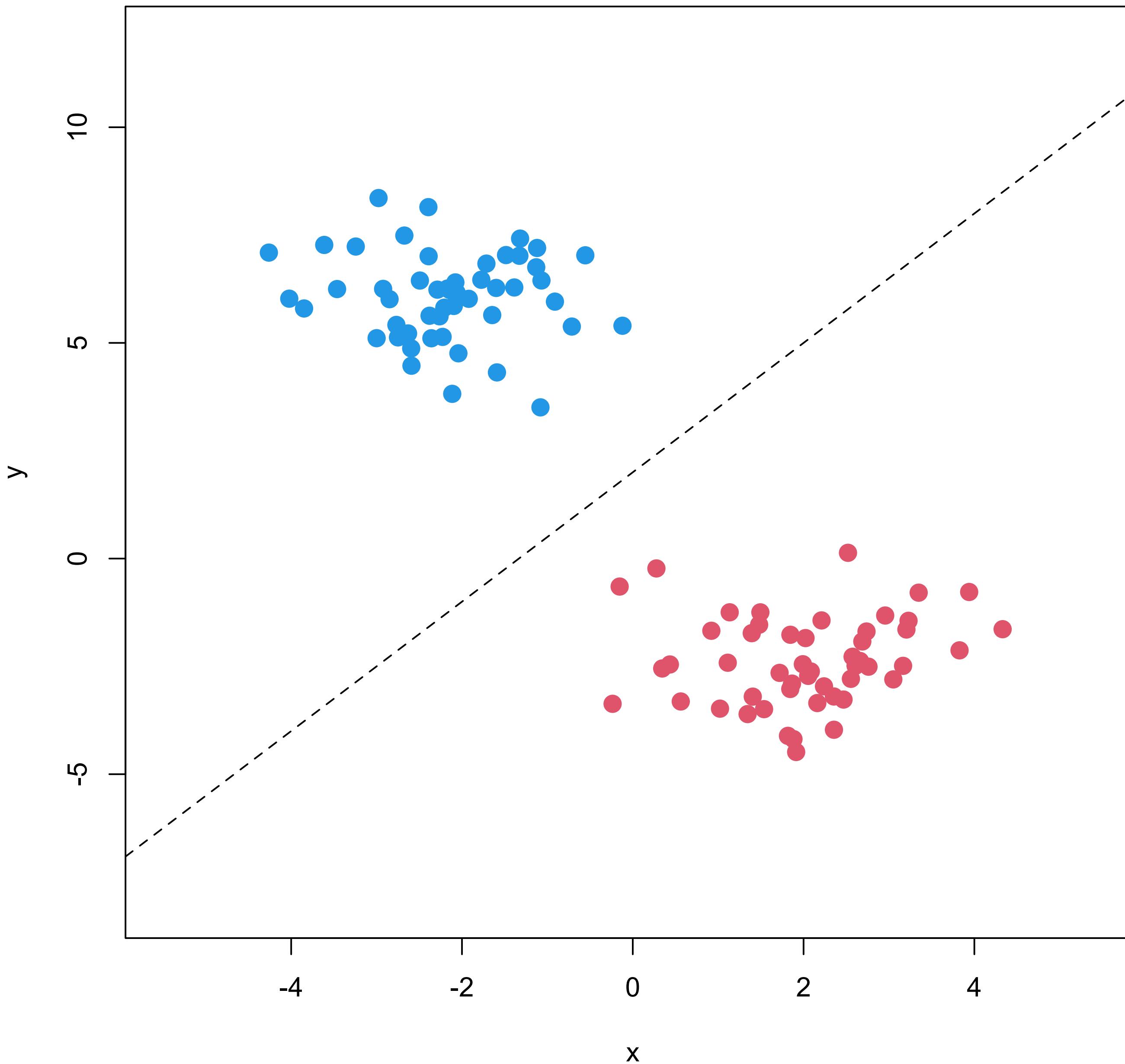
1D Binary Classification



Regression

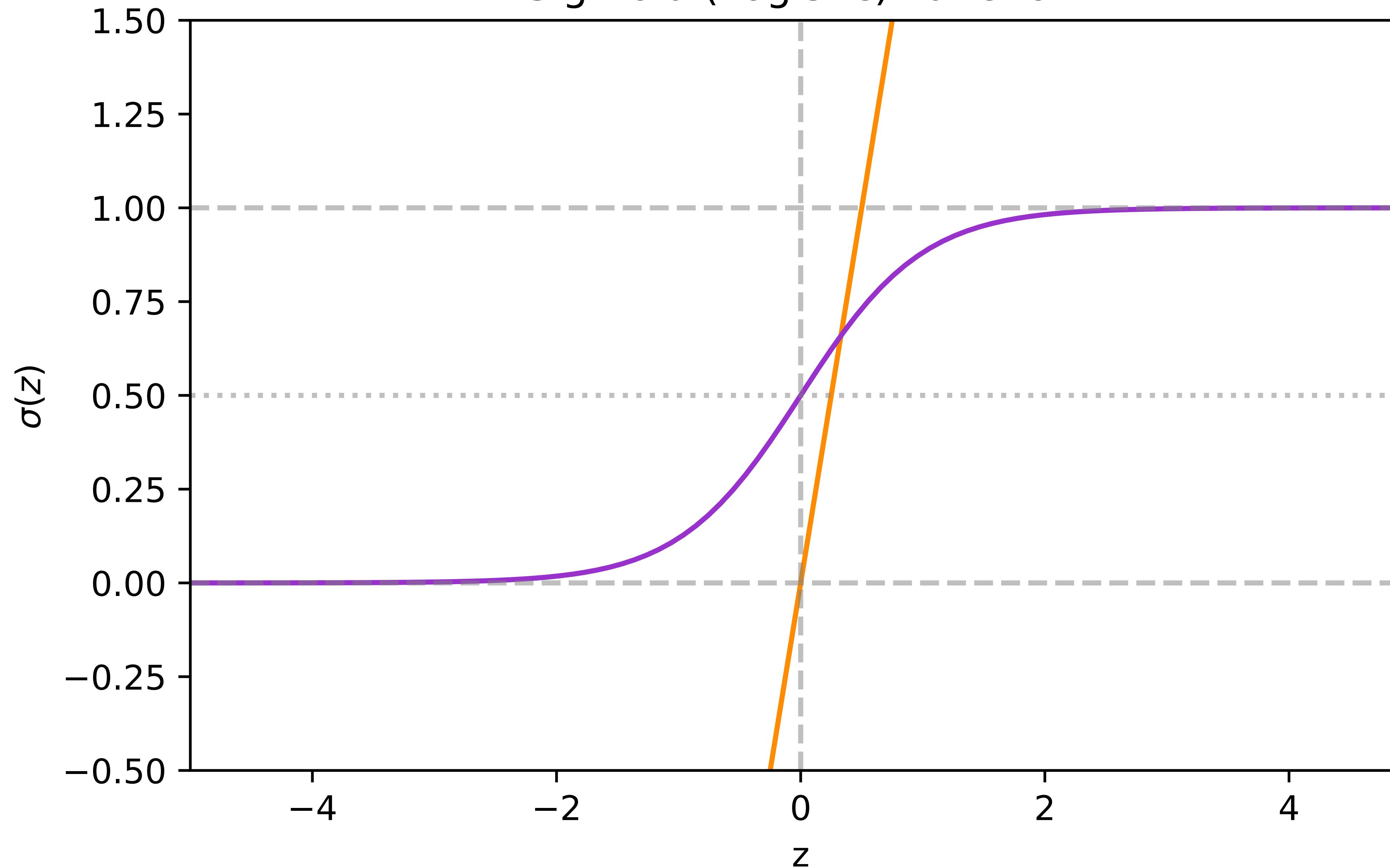


Classification



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

Sigmoid (Logistic) Function

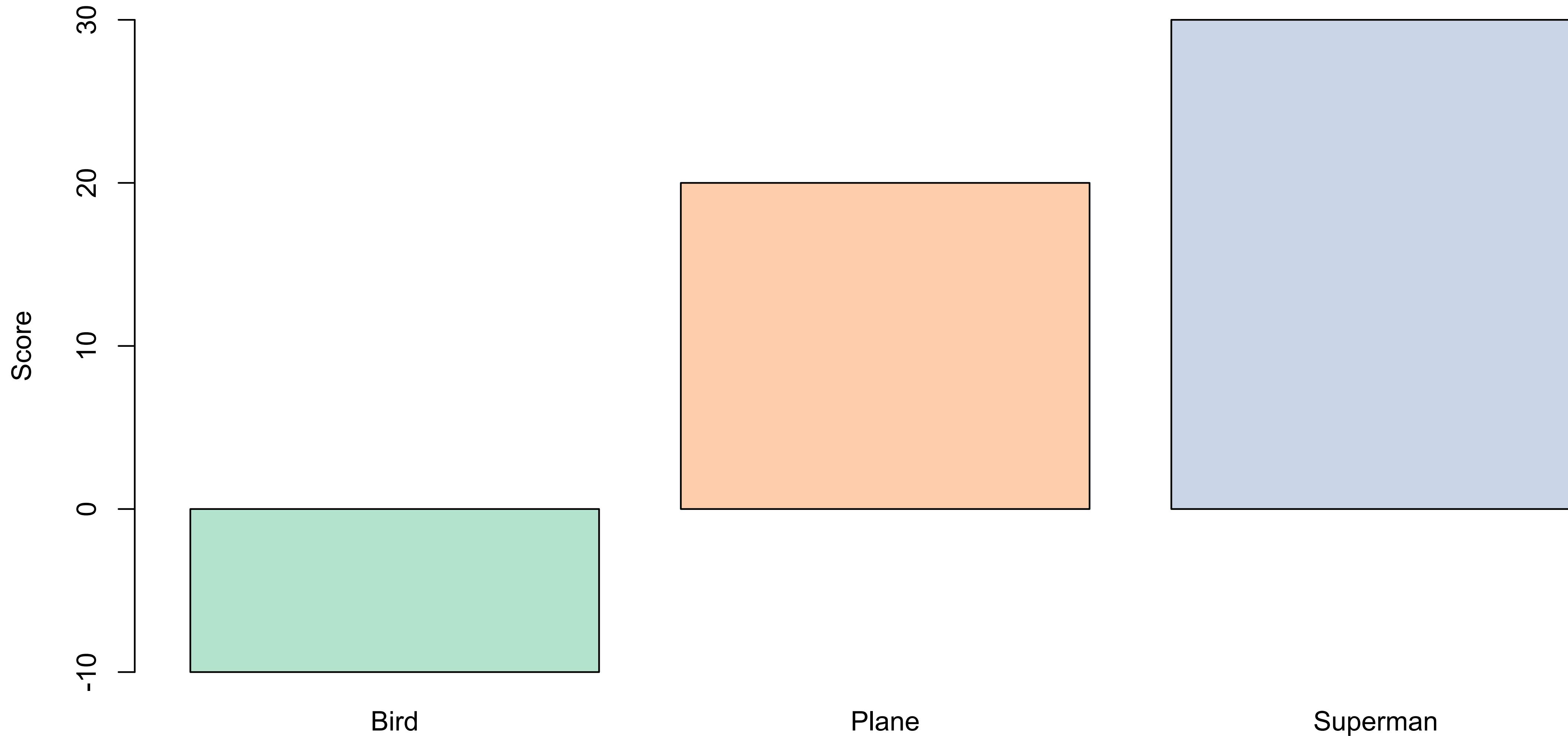


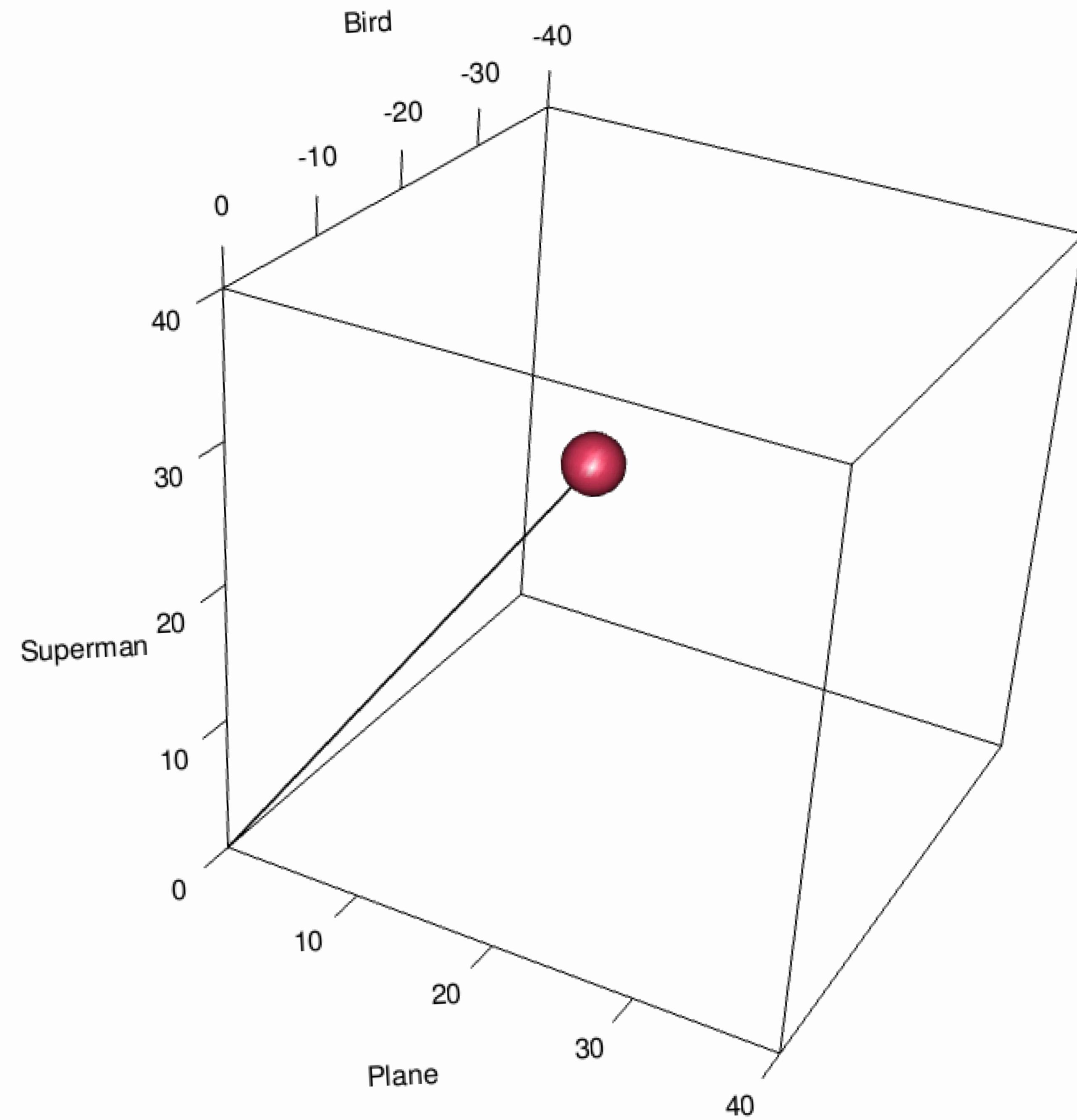
$$\nabla_{\mathbf{w}} l = \mathbf{x}^T (\mathbf{y} - \hat{\mathbf{y}})$$

$$= \mathbf{x}^T (\mathbf{y} - \sigma(\mathbf{x}_{\mathbf{w}}))$$

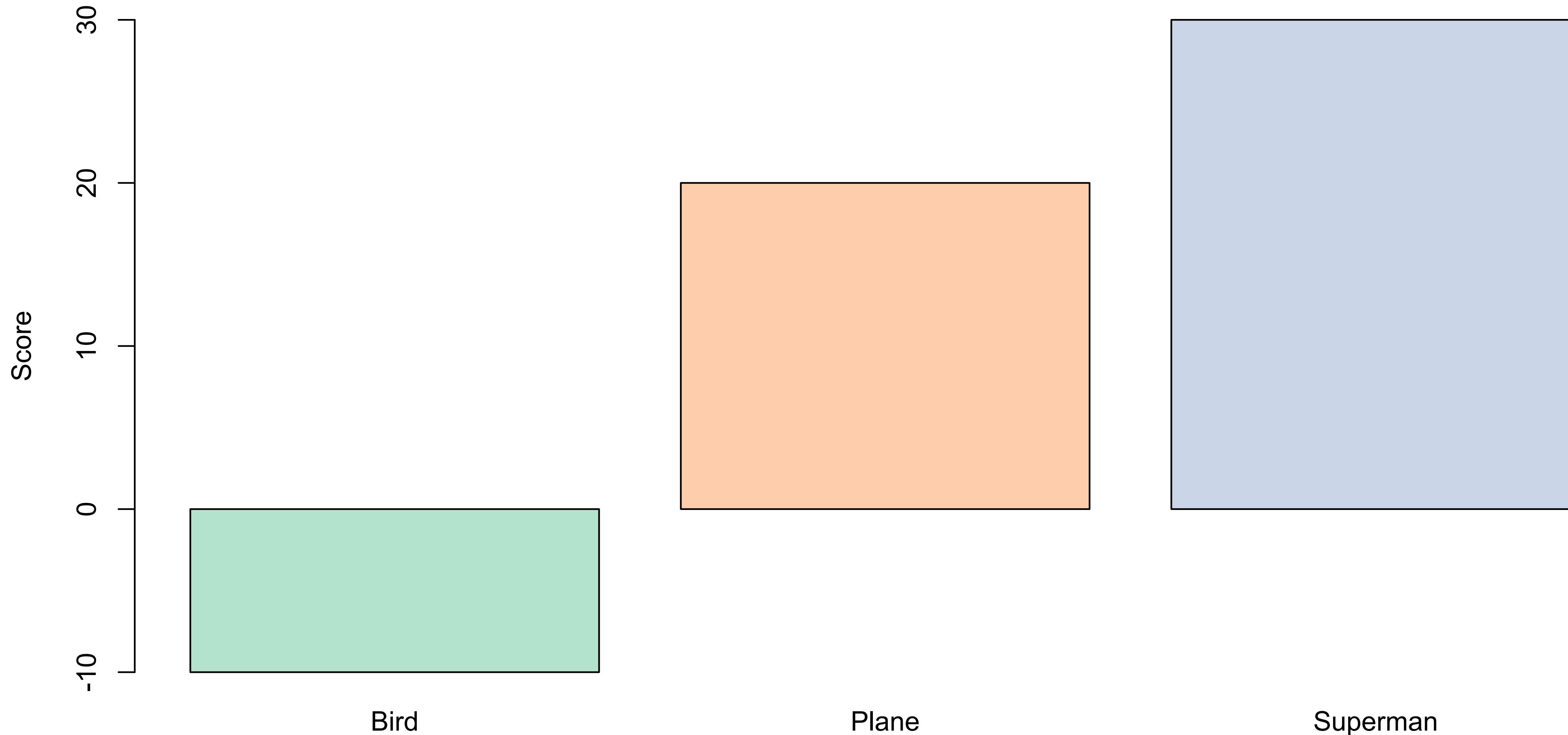
$$= \mathbf{x}^T \left(\mathbf{y} - \frac{1}{1 + e^{-\mathbf{x}_{\mathbf{w}}}} \right)$$

Class Scores (Logits)



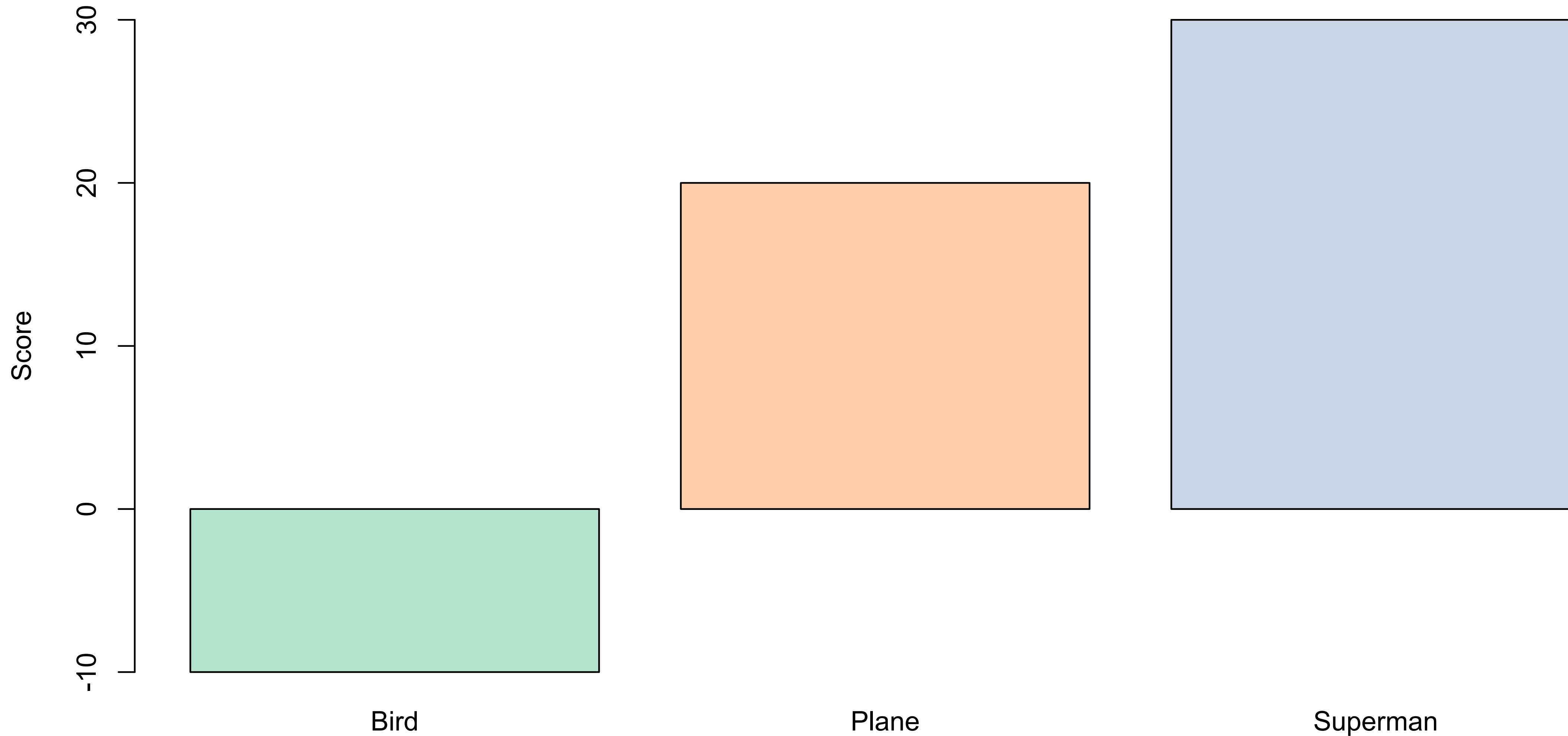


Class Scores (Logits)

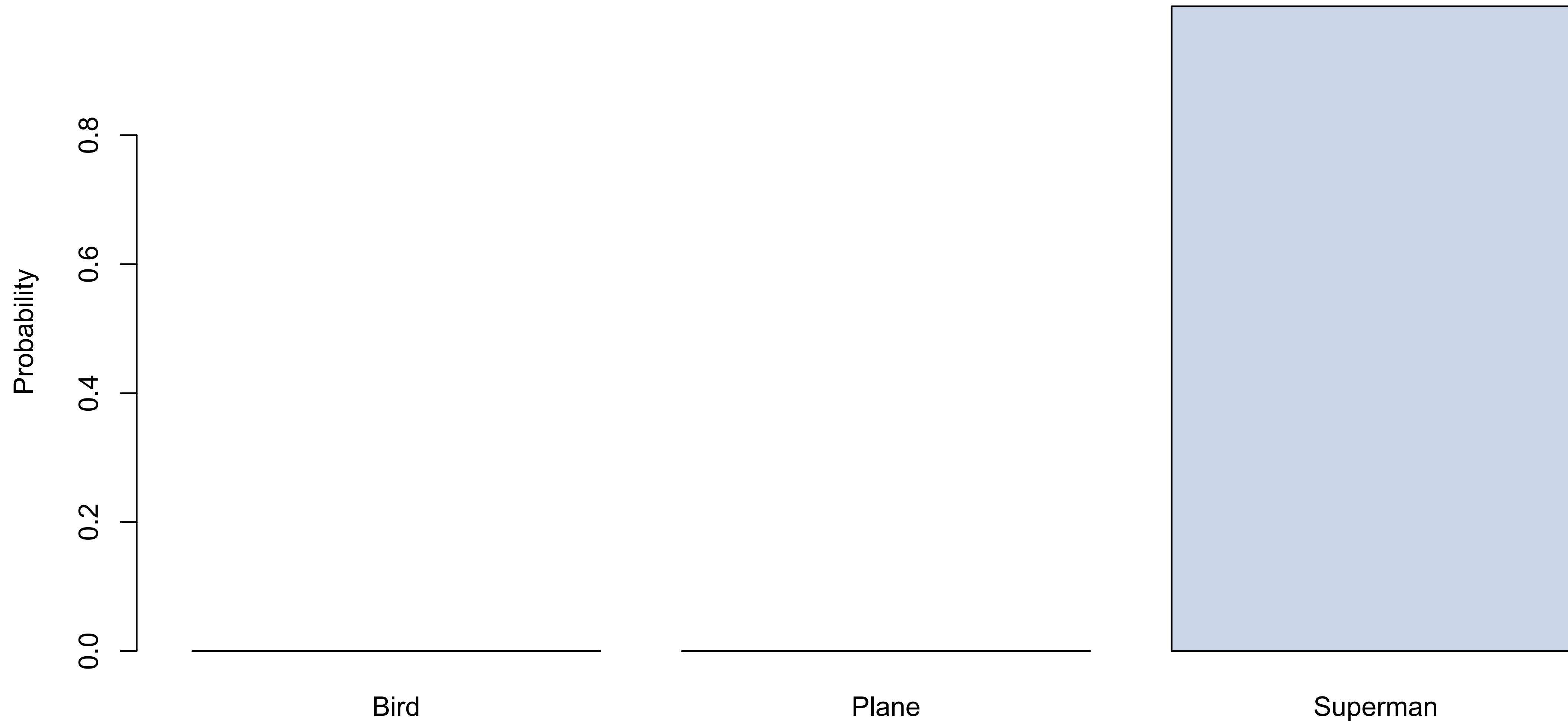


$$\varsigma(z) = \frac{e^z}{\sum e^z}$$

Class Scores (Logits)



Class Probabilities



$$\nabla_{\mathbf{W}} L = \mathbf{X}^T (\hat{\mathbf{Y}} - \mathbf{Y})$$

Questions?

Next: Non-Parametric Models

