

# Aprendizaje automatizado

## MÍNIMOS CUADRADOS

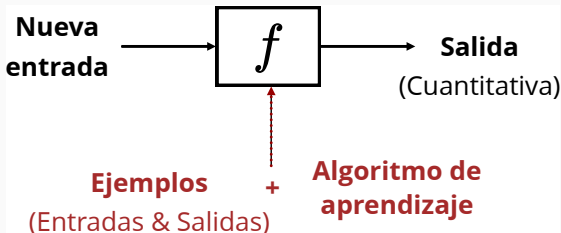
---

Gibran Fuentes Pineda

Febrero 2023

# Regresión

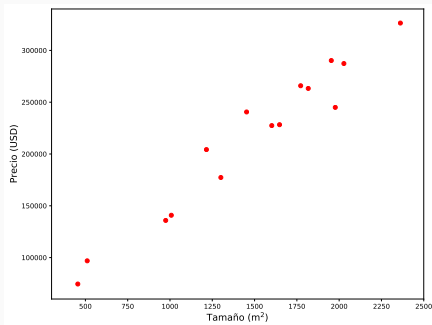
- Salida continua (cuantitativa)
- Ejemplos: predicción de temperatura de un cuarto, etc.



# Prediciendo el precio de casas

- ¿Cómo podemos ajustar nuestra función  $f$  para modelar la relación entre el tamaño y el precio de casas?

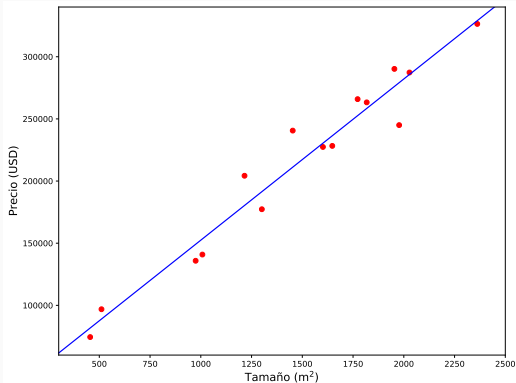
Tamaño (m <sup>2</sup> )	Precio (USD)
489.59	489.59
556.08	556.08
570.35	570.35
772.84	772.84
970.95	970.95
1162.00	1162.00
1263.10	1263.10
⋮	⋮



# Prediciendo el precio de casas

- Podemos hacer presuposiciones sobre  $f$ , por ejemplo que la relación es lineal:

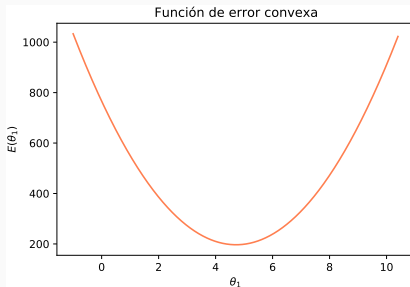
$$f_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$



## ¿Cómo medimos la calidad del ajuste?

- Definimos una función de error, por ejemplo la suma de errores cuadráticos:

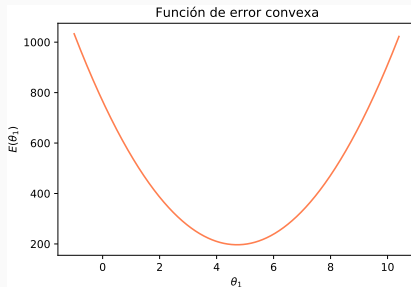
$$E(\boldsymbol{\theta}) = \frac{1}{2} \cdot \sum_{i=1}^n \left\{ f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)} \right\}^2$$



## ¿Cómo medimos la calidad del ajuste?

- Definimos una función de error, por ejemplo la suma de errores cuadráticos:

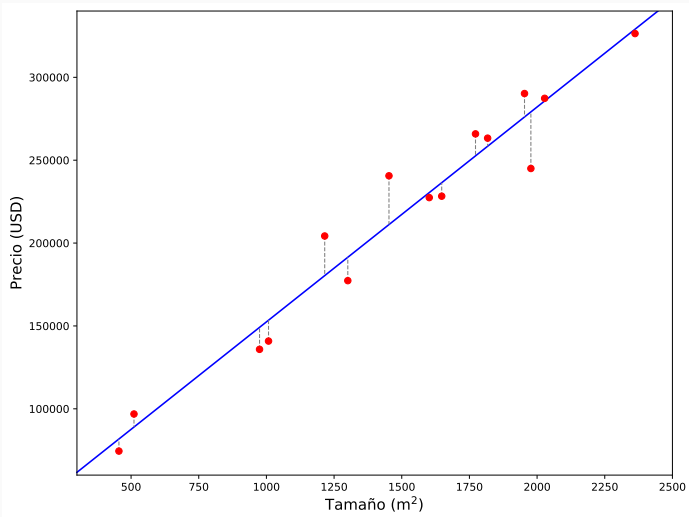
$$E(\boldsymbol{\theta}) = \frac{1}{2} \cdot \sum_{i=1}^n \left\{ f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)} \right\}^2$$



- Objetivo: encontrar el valor de  $\boldsymbol{\theta}$  que minimice  $E(\boldsymbol{\theta})$

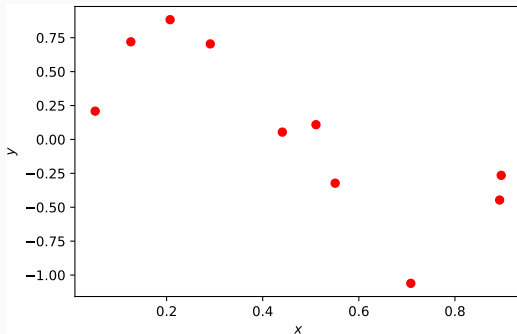
$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} E(\boldsymbol{\theta})$$

## ¿Cómo medimos la calidad del ajuste?



# Modelando relaciones no lineales

- ¿Qué función se ajusta a estos datos?





- Podemos ajustar un polinomio de la siguiente forma<sup>1</sup>

$$f_{\theta}(x) = \theta_0 + \theta_1 \cdot x + \theta_1 \cdot x^2 + \dots + \theta_d \cdot x^d$$

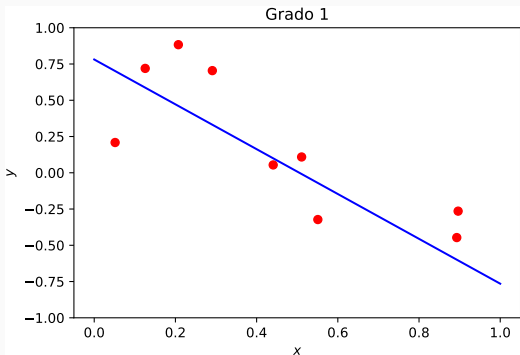
---

<sup>1</sup>Nota que esta forma no está considerando interacciones

## ¿Qué grado del polinomio es adecuado?

- Podemos usar uno lineal nuevamente

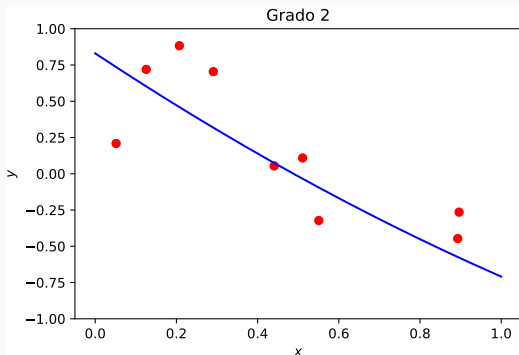
$$f_{\theta}(x) = \theta_0 + \theta_1 \cdot x$$



## ¿Qué grado del polinomio es adecuado?

- O uno cuadrático

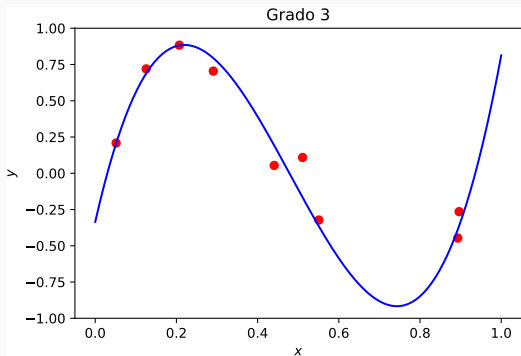
$$f_{\theta}(x) = \theta_0 + \theta_1 \cdot x + \theta_2 \cdot x^2$$



# ¿Qué grado del polinomio es adecuado?

- Grado 3

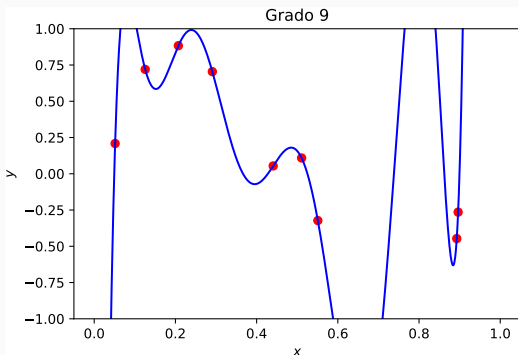
$$f_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 \cdot x^2 + \theta_3 \cdot x^3$$



# ¿Qué grado del polinomio es adecuado?

- 0 grado 9

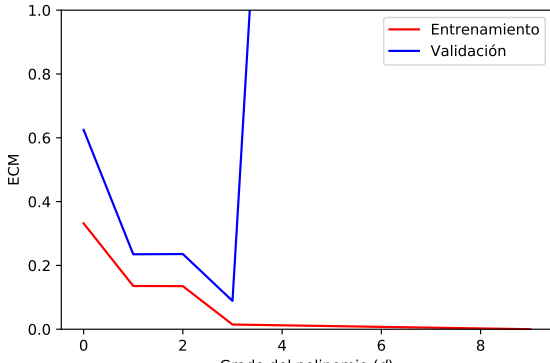
$$f_{\theta}(x) = \theta_0 + \theta_1 + \theta_2 \cdot x^2 + \cdot x + \cdots + \theta_9 \cdot x^9$$



# El problema de la generalización

- Comparamos los desempeños con distintos grados de polinomio usando el error cuadrático medio (ECM)

$$E(\boldsymbol{\theta}) = \frac{1}{n} \cdot \sum_{i=1}^n \left\{ f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)} \right\}^2$$



## ¿Por qué está sobreajustando?

	$d = 0$	$d = 1$	$d = 3$	$d = 9$
$\theta_0$	0.05	0.78	-0.33	-17.62
$\theta_1$		-1.54	12.32	762.18
$\theta_2$			-36.32	12071.82
$\theta_3$			25.14	98135.73
$\theta_4$				-459092.41
$\theta_5$				1301097.36
$\theta_6$				-2263938.71
$\theta_7$				2358449.27
$\theta_8$				-1347197.15
$\theta_9$				324015.43

## ¿Cómo evito el sobreajuste?

- Penalizando parámetros con valores grandes

$$\tilde{E}(\boldsymbol{\theta}) = \frac{1}{2} \cdot \sum_{i=1}^n \left\{ f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)} \right\}^2 + \frac{\lambda}{2} \cdot \|\boldsymbol{\theta}\|_2^2$$

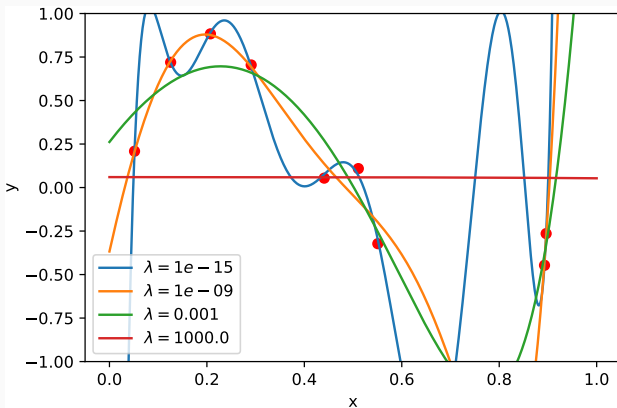
- $\lambda$  determina la ponderación que se le da al término de penalización



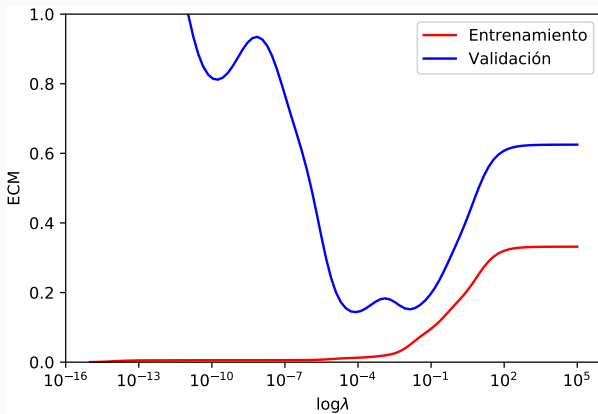
## ¿Cómo evito el sobreajuste?

	$\log \lambda = -\infty$	$\log \lambda = -18$	$\log \lambda = 0$
$\theta_0$	0.35	0.35	-17.62
$\theta_1$	232.37	4.74	-0.05
$\theta_2$	-5321.83	-0.77	-0.06
$\theta_3$	48568	-31.97	-0.05
$\theta_4$	-231639.30	-3.89	-0.03
$\theta_5$	640042.26	55.28	-0.02
$\theta_6$	-1061800.52	41.32	-0.01
$\theta_7$	1042400.18	-45.95	-0.00
$\theta_8$	-557682.99	-91.53	0.00
$\theta_9$	125201.43	72.68	0.01

# Mínimos cuadrados penalizados



# Mínimos cuadrados penalizados



- Modelo lineal

$$f_{\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{x} = \sum_{i=1}^d \theta_i \cdot x_i$$

- Modelo lineal

$$f_{\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{x} = \sum_{i=1}^d \theta_i \cdot x_i$$

- Con expansión de funciones base  $\phi$

$$f_{\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \boldsymbol{\phi}(\mathbf{x}) = \sum_{i=1}^d \theta_i \cdot \phi(\mathbf{x})_i$$

- Modelo lineal

$$f_{\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \mathbf{x} = \sum_{i=1}^d \theta_i \cdot x_i$$

- Con expansión de funciones base  $\phi$

$$f_{\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}^{\top} \boldsymbol{\phi}(\mathbf{x}) = \sum_{i=1}^d \theta_i \cdot \phi(\mathbf{x})_i$$

- Lineal en los parámetros  $\boldsymbol{\theta}$

- Asumiendo ruido  $\epsilon$  con distribución normal en el modelo

$$y = f_{\theta}(\mathbf{x}, \theta) + \epsilon$$

- Asumiendo ruido  $\epsilon$  con distribución normal en el modelo

$$y = f_{\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta}) + \epsilon$$

- Tratamos de modelar la probabilidad condicional de la salida dados los datos y parámetros

$$P(y|\mathbf{x}, \boldsymbol{\theta}, \sigma^2) = \mathcal{N}(y|f_{\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta}), \sigma^2)$$



# Obteniendo el estimador de máxima verosimilitud

- Se busca minimizar el negativo de la verosimilitud logarítmica

$$\begin{aligned} NVL(\boldsymbol{\theta}) &= - \sum_{i=1}^n \log P(y^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}) \\ &= - \sum_{i=1}^n \log \mathcal{N}(y^{(i)} | f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}), \sigma^2) \\ &= - \frac{1}{2\sigma^2} \sum_{i=1}^n \left\{ f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)} \right\}^2 - \frac{n}{2} \log 2\pi\sigma^2 \end{aligned}$$

# Obteniendo el estimador de máxima verosimilitud

- Se busca minimizar el negativo de la verosimilitud logarítmica

$$\begin{aligned} NVL(\boldsymbol{\theta}) &= - \sum_{i=1}^n \log P(y^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}) \\ &= - \sum_{i=1}^n \log \mathcal{N}(y^{(i)} | f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}), \sigma^2) \\ &= - \frac{1}{2\sigma^2} \sum_{i=1}^n \left\{ f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)} \right\}^2 - \frac{n}{2} \log 2\pi\sigma^2 \end{aligned}$$

- Equivalente a minimizar suma de errores cuadráticos

$$E(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n \left\{ f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)} \right\}^2$$

## Obteniendo el estimador de máxima verosimilitud

- Reformulando NVL

$$\begin{aligned} NVL(\boldsymbol{\theta}) &= \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= \frac{1}{2}\mathbf{y}^\top \mathbf{y} - \frac{1}{2}\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} - \frac{1}{2}\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} + \frac{1}{2}\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} \\ &= \frac{1}{2}\mathbf{y}^\top \mathbf{y} - \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} + \frac{1}{2}\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} \end{aligned}$$

## Obteniendo el estimador de máxima verosimilitud

- Reformulando NVL

$$\begin{aligned} NVL(\boldsymbol{\theta}) &= \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= \frac{1}{2}\mathbf{y}^\top \mathbf{y} - \frac{1}{2}\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} - \frac{1}{2}\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} + \frac{1}{2}\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} \\ &= \frac{1}{2}\mathbf{y}^\top \mathbf{y} - \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} + \frac{1}{2}\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} \end{aligned}$$

- Derivando con respecto a  $\boldsymbol{\theta}$  e igualando a cero

$$\mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} = \mathbf{X}^\top \mathbf{y}$$

# Obteniendo el estimador de máxima verosimilitud

- Reformulando NVL

$$\begin{aligned} NVL(\boldsymbol{\theta}) &= \frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= \frac{1}{2}\mathbf{y}^\top \mathbf{y} - \frac{1}{2}\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} - \frac{1}{2}\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} + \frac{1}{2}\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} \\ &= \frac{1}{2}\mathbf{y}^\top \mathbf{y} - \boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{y} + \frac{1}{2}\boldsymbol{\theta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} \end{aligned}$$

- Derivando con respecto a  $\boldsymbol{\theta}$  e igualando a cero

$$\mathbf{X}^\top \mathbf{X}\boldsymbol{\theta} = \mathbf{X}^\top \mathbf{y}$$

- El estimador de máxima verosimilitud es

$$\hat{\boldsymbol{\theta}}_{EMV} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

## ¿Y si tenemos múltiples variables de salida?

- Solución de mínimos cuadrados

$$\hat{\Theta}_{EMV} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

- Equivalente a

$$\hat{\theta}_{kEMV} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}_k$$

## Obteniendo el estimador de máximo a posteriori

- Asumiendo distribución a priori normal sobre  $\boldsymbol{\theta}$

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{MAP} = \arg \max_{\boldsymbol{\theta}} & \sum_{i=1}^n \log \mathcal{N}(y^{(i)} | \theta_0 + \boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}^{(i)}), \sigma^2) \\ & + \sum_{j=0}^d \log \mathcal{N}(\theta_j | 0, \tau^2)\end{aligned}$$

## Obteniendo el estimador de máximo a posteriori

- Asumiendo distribución a priori normal sobre  $\boldsymbol{\theta}$

$$\hat{\boldsymbol{\theta}}_{MAP} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \log \mathcal{N}(y^{(i)} | \theta_0 + \boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}^{(i)}), \sigma^2) \\ + \sum_{j=0}^d \log \mathcal{N}(\theta_j | 0, \tau^2)$$

- Equivalente a minimizar suma de errores cuadráticos con los parámetros penalizados con la norma  $\ell_2$

$$\tilde{E}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n \{f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)}\}^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$



## Obteniendo el estimador de máximo a posteriori

- Asumiendo distribución a priori normal sobre  $\boldsymbol{\theta}$

$$\begin{aligned}\hat{\boldsymbol{\theta}}_{MAP} = \arg \max_{\boldsymbol{\theta}} & \sum_{i=1}^n \log \mathcal{N}(y^{(i)} | \theta_0 + \boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}^{(i)}), \sigma^2) \\ & + \sum_{j=0}^d \log \mathcal{N}(\theta_j | 0, \tau^2)\end{aligned}$$

- Equivalente a minimizar suma de errores cuadráticos con los parámetros penalizados con la norma  $\ell_2$

$$\tilde{E}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n \{f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)}\}^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

- Derivando  $\tilde{E}(\boldsymbol{\theta})$  con respecto a  $\boldsymbol{\theta}$  e igualando a cero

$$\hat{\boldsymbol{\theta}}_{ridge} = (\lambda \cdot \mathbf{I}_D + \mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

- Cuando la regularización es por norma  $\ell_1$  se conoce como LASSO

$$\hat{\boldsymbol{\theta}}_{\text{LASSO}} = \arg \min_{\boldsymbol{\theta}} \left[ \frac{1}{2} \cdot \sum_{i=1}^n \left\{ f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)} \right\}^2 + \frac{\lambda}{2} \cdot \|\boldsymbol{\theta}\|_1 \right]$$

- Cuando la regularización es por norma  $\ell_1$  se conoce como LASSO

$$\hat{\boldsymbol{\theta}}_{\text{LASSO}} = \arg \min_{\boldsymbol{\theta}} \left[ \frac{1}{2} \cdot \sum_{i=1}^n \left\{ f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)} \right\}^2 + \frac{\lambda}{2} \cdot \|\boldsymbol{\theta}\|_1 \right]$$

- Optimización cuadrática: no existe solución cerrada pero existen algoritmos eficientes

# Regularización con diferentes normas

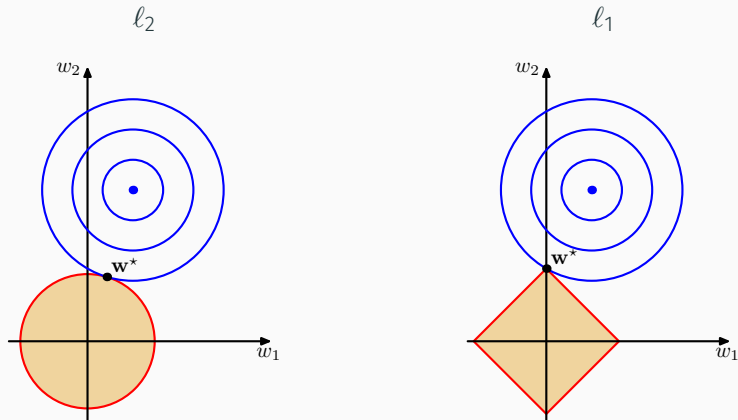


Imagen tomada de C. Bishop. PRML, 2009