# Introducción a la Inteligencia Artificial
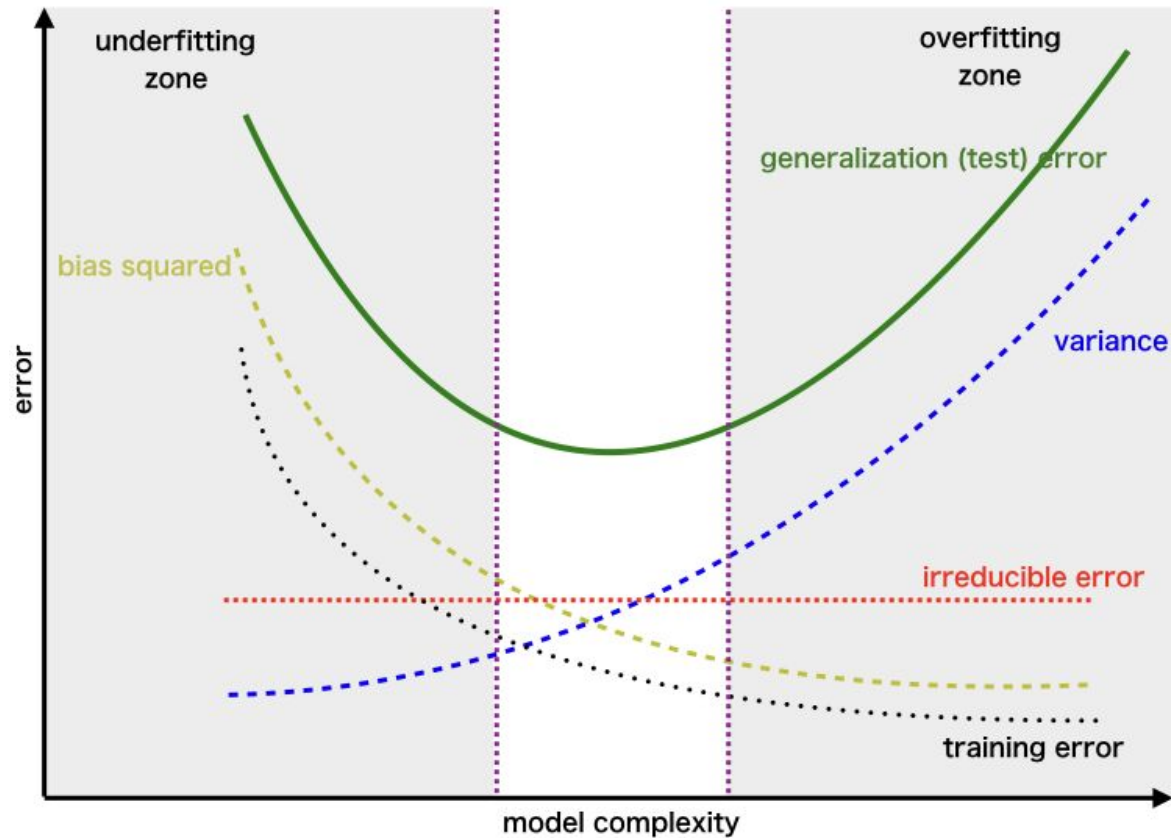## Clase 4
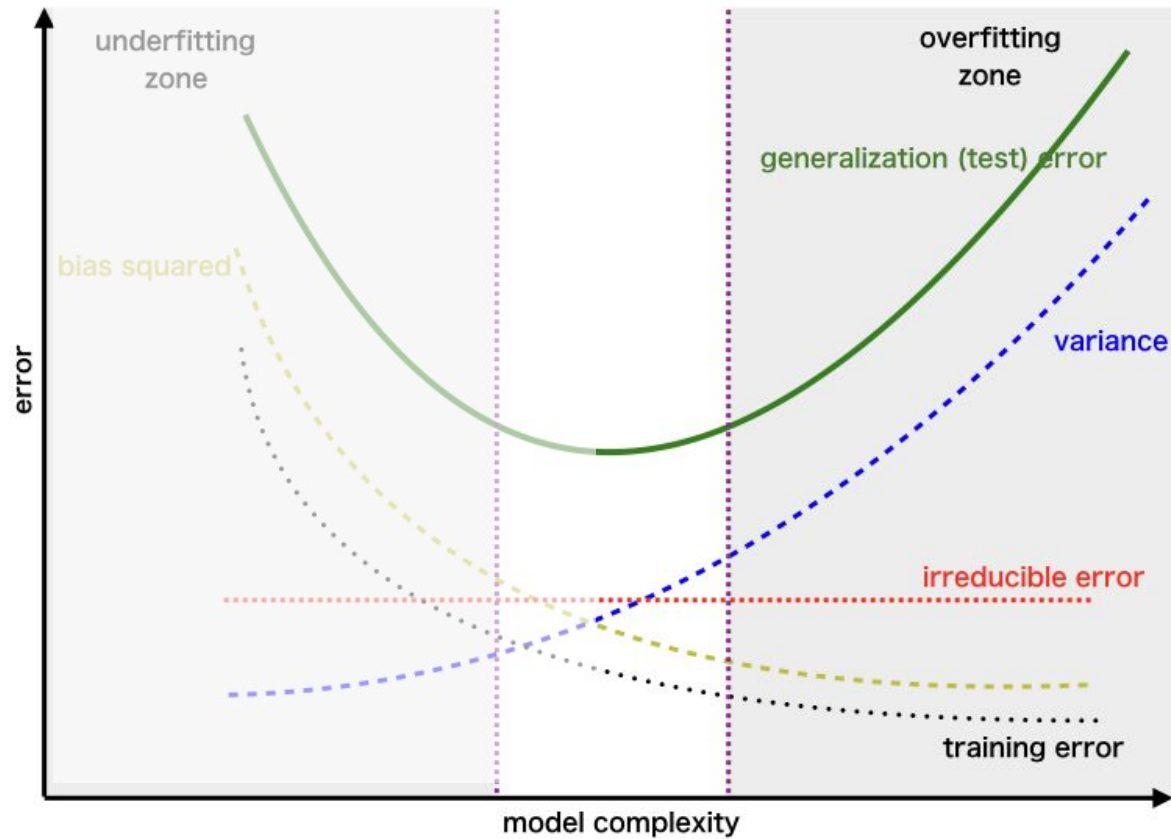
**Clase ~~5~~ 4**

0. Riesgo Empírico

.**UBA**fiuba
FACULTAD DE INGENIERÍA

## Regularización - Motivación



OLS
$\beta_1 = 1.022$, $\beta_0 = -0.17$

OLS with outlier
$\beta_1 = 5.402$, $\beta_0 = -13.15$

L2 Ridge with outlier
$\beta_1 = 1.369$, $\beta_0 = 7.02$

# Riesgo Empírico

$$R(f) = \mathbb{E}\left( \underbrace{L(f, \hat{f})}_{\text{función de pérdida (loss function)}} \right) \rightsquigarrow \quad L(f, \hat{f}) = (f - \hat{f})^2 \qquad \text{fn. de pérdida}$$

$$L2$$

función de
pérdida
(loss function)

outliers

regresión
esperada

⊗ — valor anómalo

¿Cómo podemos mejorar mi
reg. lineal?

1. Cambiar $L$

— $L_1 = |f - \hat{f}|$

— $L_2 = (f - \hat{f})^2$

— $L_{\mathbb{I}} = \mathbb{I}_{|f - \hat{f}| \leq C}$

— $L_{huber}$  ⎫  Estimadores
         ⎬  robustos de
— $L_{Tukey}$  ⎭  la pérdida
              (multiparámetricos)

$$L' = \begin{cases} (f - \hat{f})^2 & \text{si } |f - \hat{f}| \leq C \\ a & \text{o.w.} \end{cases}$$
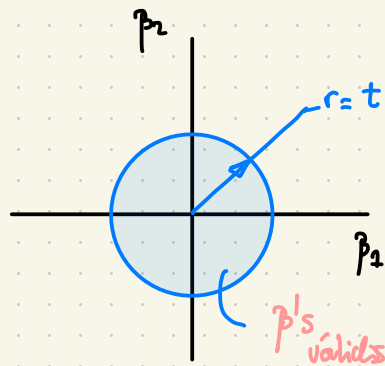
nosotros siempre buscamos minimizar R de la forma más sencillo, a veces cambiar L no es barato o peor no es útil $\Rightarrow$ probamos ==regularizar==

## 2. Regularización.

Con reg. buscamos min R al mismo tiempo que restringimos (limitamos) el comp. de los parámetros (parameter shrinking).

partimos de

$$\begin{cases} \hat{y} = \beta_0 + \sum_i \beta_i x_i \\[2mm] \hat{\beta} = \underset{\beta}{\arg\min} \left( \sum_i y_i - \sum_j \beta_j x_{ij} \right)^2 \\[2mm] \frac{1}{2} \sum_j \beta_j^{\,q} \leq t \qquad ( \|\beta\|^q \leq t ) \end{cases}$$



$\beta_2$

$r = t$

$\beta's$ válidos

Espacio de parámetros

$q = 2 \ \longrightarrow \ \hat{\beta} = \underset{\beta}{\arg\min} \left( \sum_i y_i - \sum_j \beta_j x_{ij} \right)^2 - \lambda \sum_j \beta_j^{\,2}$

parámetro de complejidad

$\big\}$ Termino de regularización (Weigth decay)

# Regularización

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{t_n - \mathbf{w}^{\mathrm{T}} \phi(\mathbf{x}_n)\}^2$$
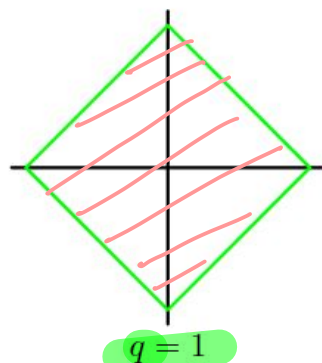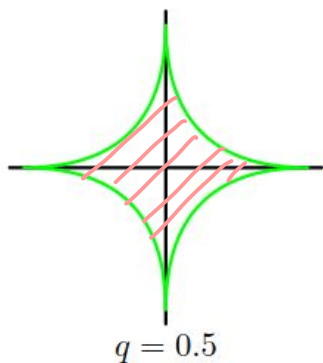
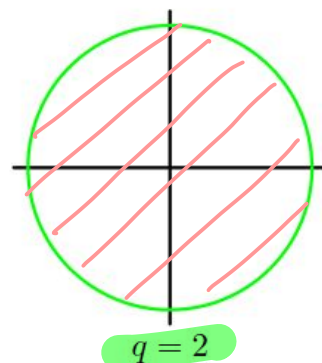Observado **-** Predicción

w está "libre"

$$\frac{1}{2}\sum_{n=1}^{N}\{t_n - \mathbf{w}^{\mathrm{T}}\phi(\mathbf{x}_n)\}^2 + \boxed{\frac{\lambda}{2}\sum_{j=1}^{M}|w_j|^q}$$

Término de regularización "weight decay" $\longrightarrow$ w afecta la pérdida
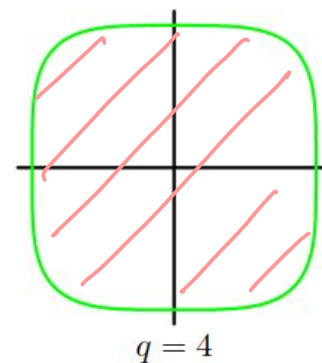


$q = 0.5$      $q = 1$      $q = 2$      $q = 4$

Lasso      Ridge

Valores validos de $\|w\|_2$

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

.UBAfiuba
FACULTAD DE INGENIERÍA

## Maximum A Posteriori como regularización

Distribución "a priori" de los parámetros $\longrightarrow$ Observar data $\longrightarrow$ Actualizar distribución (Posterior)

$$p(w) \sim D(\theta)$$

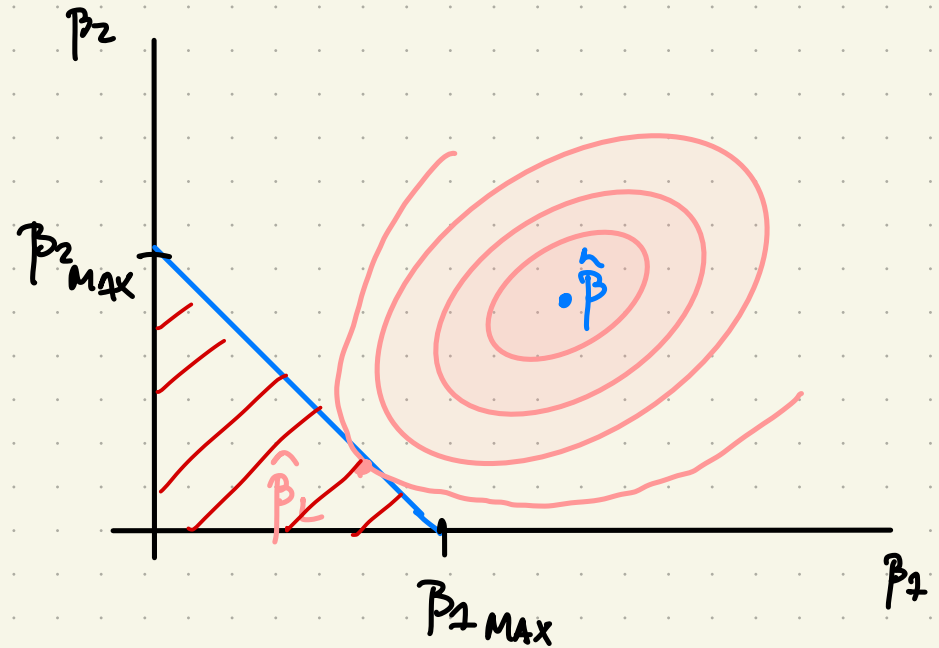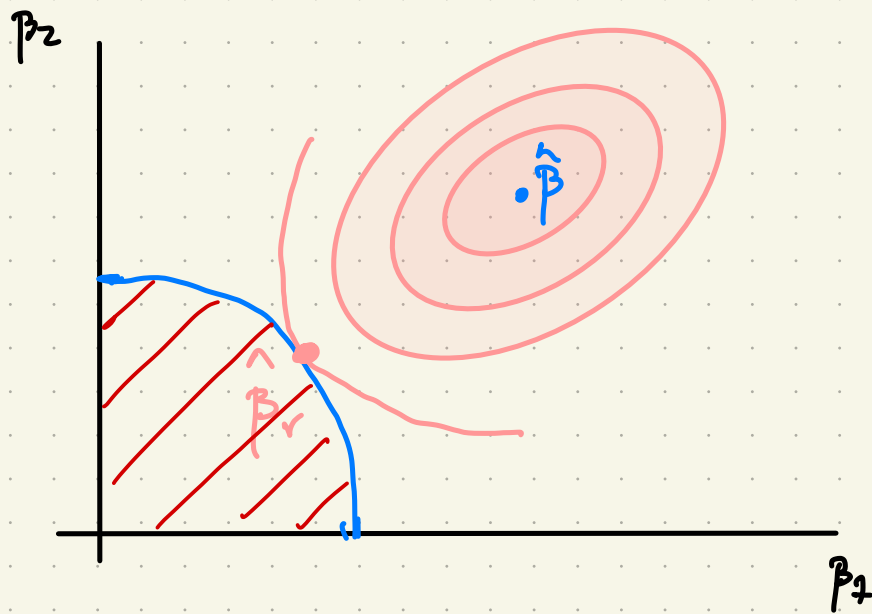$$(\mathcal{X}, \mathcal{Y})$$

$$p(w|\mathcal{X}, \mathcal{Y}) = \frac{p(\mathcal{Y}|\mathcal{X}, w)p(w)}{p(\mathcal{Y}|\mathcal{X})}$$

$$w_{map} = (\Phi^T \Phi + \frac{\sigma^2}{b^2} I)^{-1} \Phi^T y$$

Gaussian prior con varianza b2

.UBAfiuba
FACULTAD DE INGENIERÍA

# Representación gráfica:

$\hat{\beta}$: Es el mejor $\hat{\beta}$ posible



☒ $\beta$'s válidos

## Ridge
↓
porque enrobustece mi sist.

$$-\lambda \sum_j |\beta_j)^2$$

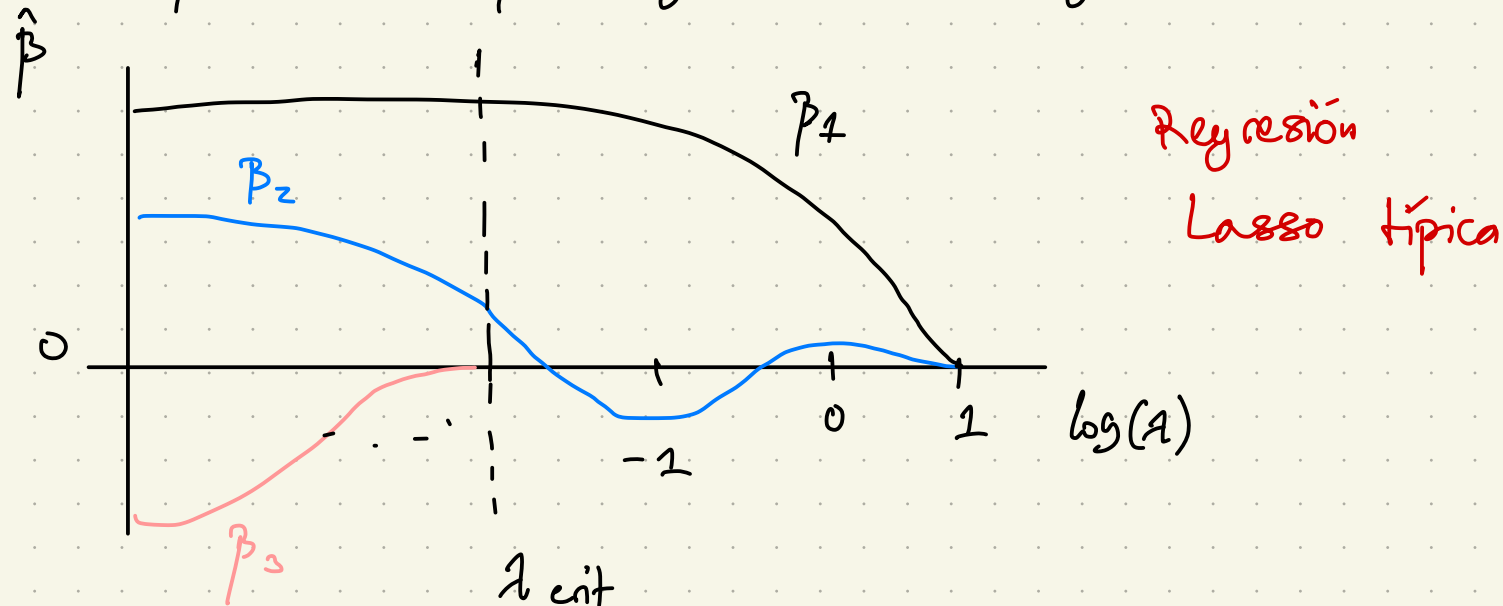## Lasso
↓
porque hace un proceso de selección de variables

Recordemos que $\hat{\beta}$ en gal tiene curvas de nivel elípticas (por construcción del estimador).

Pero, regularizar viene a costo de alejarnos del óptimo.

¿Cómo funciona?

1. Elegir $q$    (Lasso: 1, Ridge: 2)

2. Voy a elegir un vector de $\lambda$'s "apropiado"   ($\uparrow \lambda \rightsquigarrow \uparrow$ penalidad)

3. Optimizo para $\lambda_i \rightsquigarrow RSS(\lambda; q) = \| y - X\beta \|^q + \lambda \|\beta\|_q$

4. Calcular las metricas ( Error de representación, bondad, algún criterio externo).

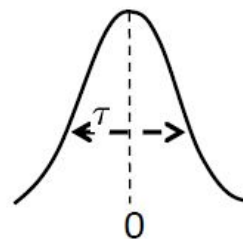5. Comparamos y elegimos el mejor:



Regresión
Lasso típica

$\hat{\beta}$

$\beta_1$

$\beta_2$

$\beta_3$

$0$   $1$   $\log(\lambda)$

$-1$

$\lambda$ crit

**Maximum A Posteriori como regularización - Ridge (L2)**

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\max_{\beta} \underbrace{\log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=}^n}_{\text{Conditional log likelihood}} + \underbrace{\log p(\beta)}_{\text{log prior}}$$

I) Gaussian Prior

$$\beta \sim \mathcal{N}(0, \tau^2 \mathbf{I}) \qquad p(\beta) \propto e^{-\beta^T \beta / 2\tau^2}$$



$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2$$

$\downarrow$

$\text{constant}(\sigma^2, \tau^2)$

Ridge Regression

$$\widehat{\beta}_{\mathsf{MAP}} = (\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^\top \mathbf{Y}$$

Machine Learning 10-315, Carnegie Mellon

.UBAfiuba
FACULTAD DE INGENIERÍA

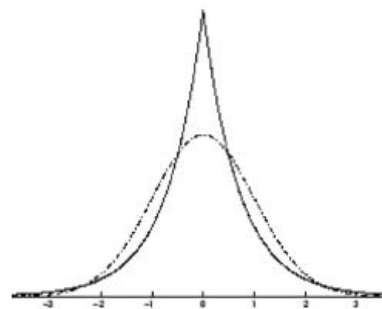**Maximum A Posteriori como regularización - LASSO (L1)**

$$\hat{\beta}_{\text{MAP}} = \arg\max_{\beta} \log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=}^n) + \log p(\beta)$$

$$\underbrace{\qquad\qquad}_{\text{Conditional log likelihood}} \quad \underbrace{\qquad}_{\text{log prior}}$$

II) Laplace Prior

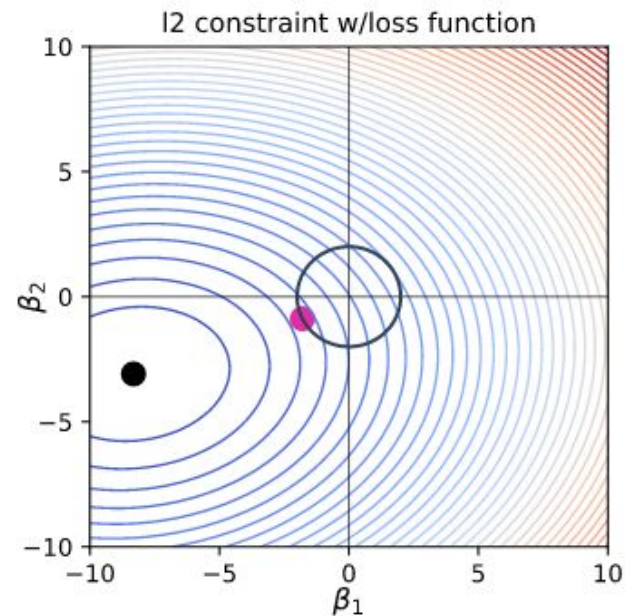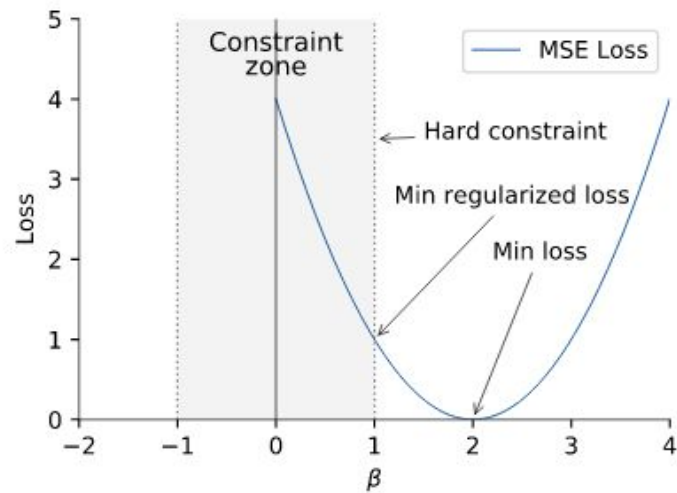$$\beta_i \overset{iid}{\sim} \text{Laplace}(0, t) \qquad p(\beta_i) \propto e^{-|\beta_i|/t}$$



$$\hat{\beta}_{\text{MAP}} = \arg\min_{\beta} \sum_{i=1}^n (Y_i - X_i\beta)^2 + \lambda\|\beta\|_1 \qquad \text{Lasso}$$

$$\downarrow$$

$$\text{constant}(\sigma^2, t)$$

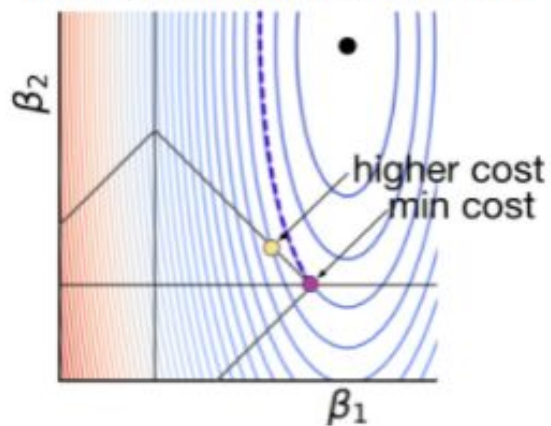Machine Learning 10-315, Carnegie Mellon

.UBAfiuba
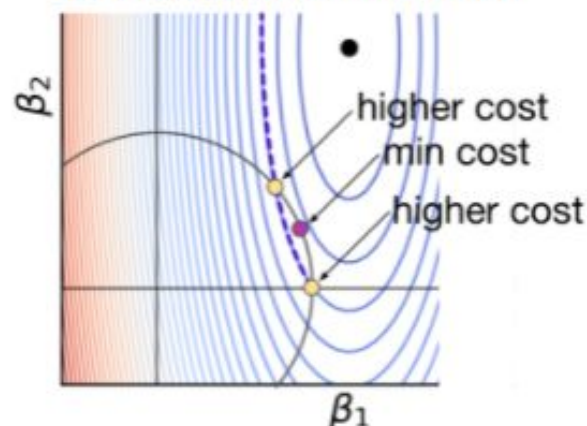FACULTAD DE INGENIERÍA

# Regularización

## Regularización
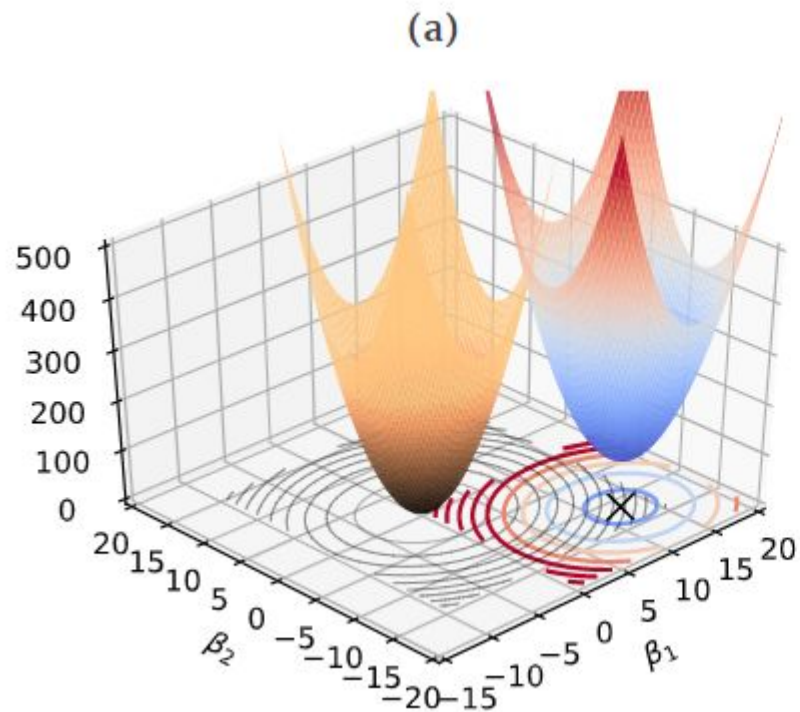


(a) L1 Constraint Diamond

(b) L2 Constraint Circle

## ElasticNet

$$(\alpha\lambda||\beta||_1 + \tfrac{1}{2}(1-\alpha)||\beta||_2^2)$$

¿Qué β se reduce más?
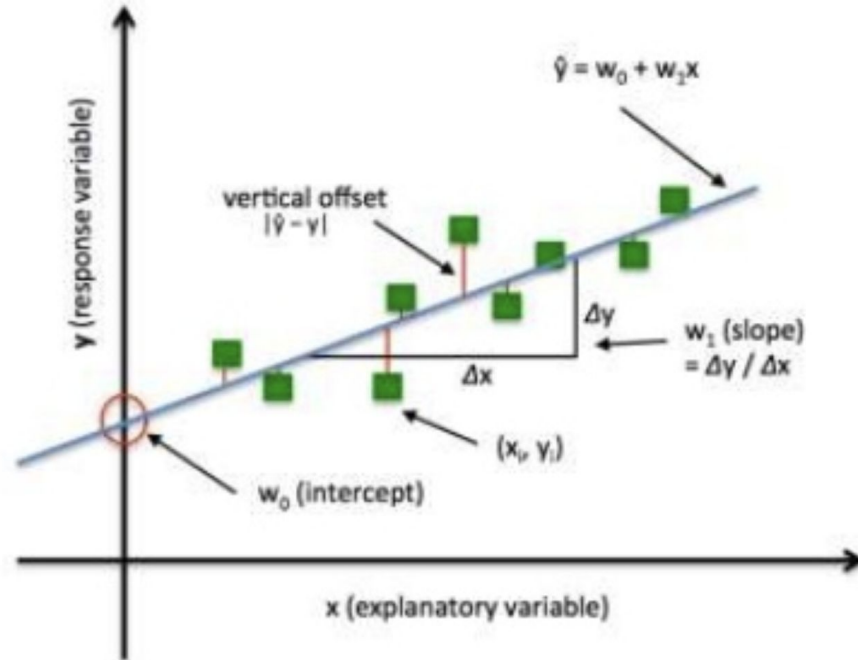
.UBAfiuba
FACULTAD DE INGENIERÍA

## Regularización

**Gradiente Descendente**

## Implementación de Gradiente Descendente

Solucion analitica

$$\min_{W} \|Y - XW\|_2^2$$

$$W = (X^T X)^{-1} X^T Y$$



$\hat{y} = w_0 + w_1 x$

y (response variable)

vertical offset
$|\hat{y} - y|$

$\Delta y$

$w_1$ (slope)
$= \Delta y / \Delta x$

$\Delta x$

$(x_i, y_i)$

$w_0$ (intercept)

x (explanatory variable)
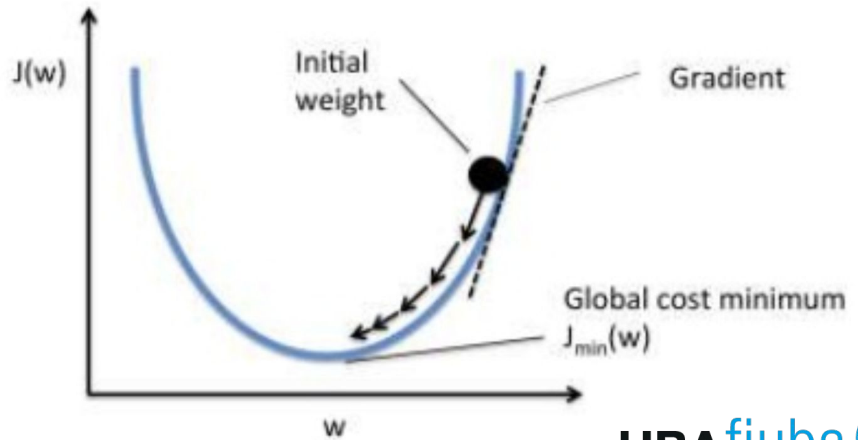
.UBAfiuba
FACULTAD DE INGENIERÍA
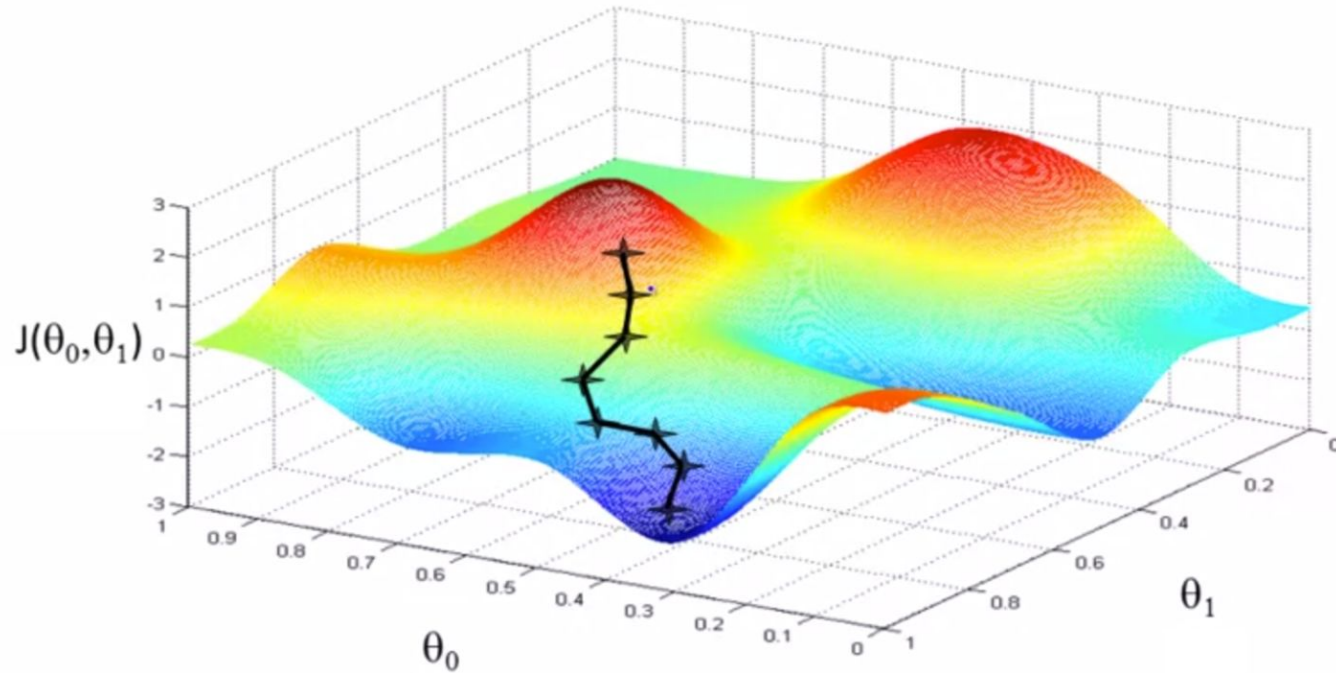
## Implementación de Gradiente Descendente

Solución numérica

$$\min_W \|Y - XW\|_2^2 \implies \min_W \sum_i (y_i - X_i \cdot W)^2$$

$$W \;\leftarrow\; W - \alpha \nabla \left( \sum_i (y_i - X_i \cdot W)^2 \right)$$
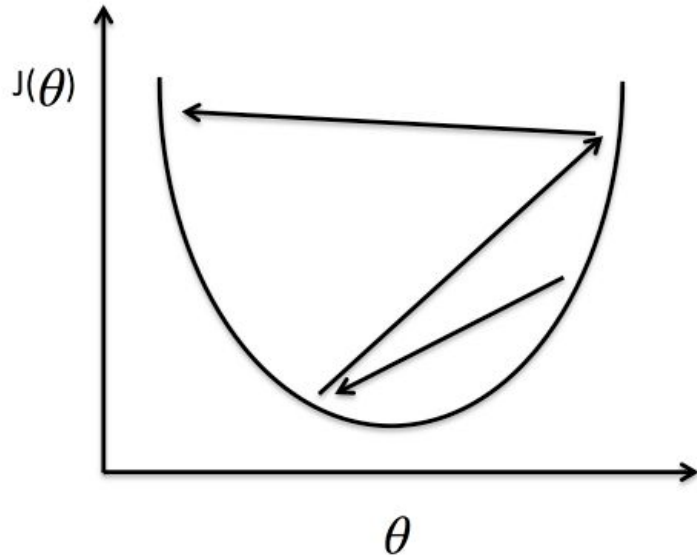
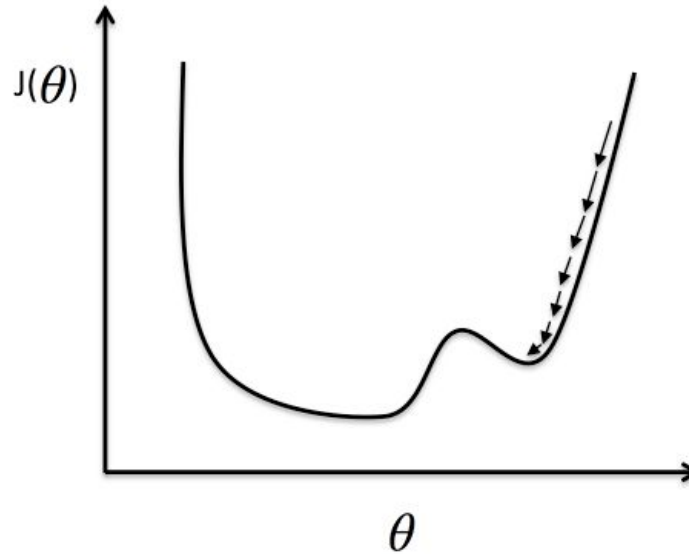

.UBAfiuba
FACULTAD DE INGENIERÍA

## Gradiente Descendente

## Gradiente Descendente



Large learning rate: Overshooting.

Small learning rate: Many iterations until convergence and trapping in local minima.

## Implementación de Gradiente Descendente

Solución numérica

$$\nabla_w J(w) = \nabla_w \Big( \sum_i (y_i - X_i W)^2 \Big)$$

$$= \sum_i \Big( \nabla_w (y_i - X_i W)^2 \Big)$$

$$= \sum_i \Big( \nabla_w (y_i - (x_{i1}w_1 + x_{i2}w_2 + \cdots + x_{im}w_m))^2 \Big)$$

$$= \sum_i \Big( -2(y_i - \hat{y}_i)x_{ij} \Big) \qquad \forall j \in (1 \cdots m)$$

**Implementación de Gradiente Descendente**

Solución numérica

$$\nabla \left( \sum_{\text{all samples}} (y_i - f_W(X_i))^2 \right)$$

## Gradient Descent algorithm

for epoch in n_epochs:
- compute the predictions for **all the samples**
- compute the error between truth and predictions
- compute the gradient using **all the samples**
- update the parameters of the model

.UBAfiuba
FACULTAD DE INGENIERÍA

**Implementación de Gradiente Descendente Estocástico**

Solución numérica
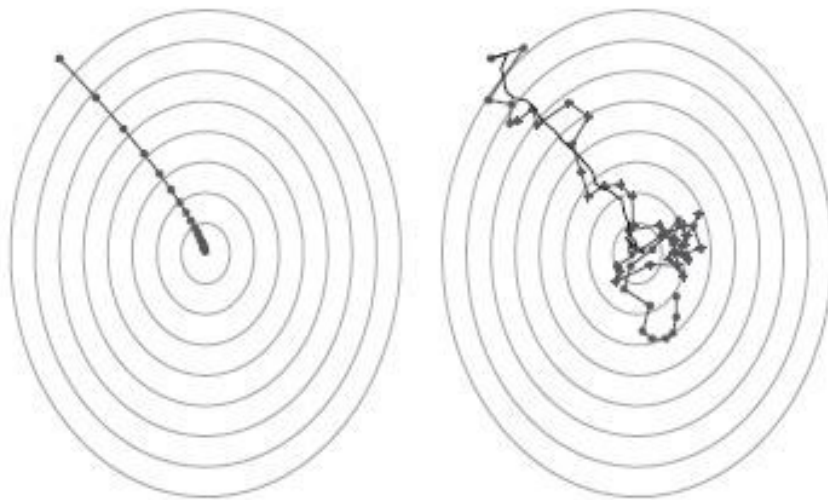
$$\nabla\left(\left(y_i - f_W(X_i)\right)^2\right)$$

## Stochastic Gradient Descent algorithm

for epoch in n_epochs:
- shuffle the samples
- **for sample in n_samples:**
  - compute the predictions for the sample
  - compute the error between truth and predictions
  - compute the gradient using the sample
  - update the parameters of the model

.UBA*fiuba*
FACULTAD DE INGENIERÍA

## Implementación de Gradiente Descendente Estocástico

Solución numérica

**Implementación de Gradiente Descendente Mini-Batch**

Solución numérica

$$\nabla \left( \sum_{\text{batch samples}} (y_i - f_W(X_i))^2 \right)$$

Mini-Batch Gradient Descent algorithm
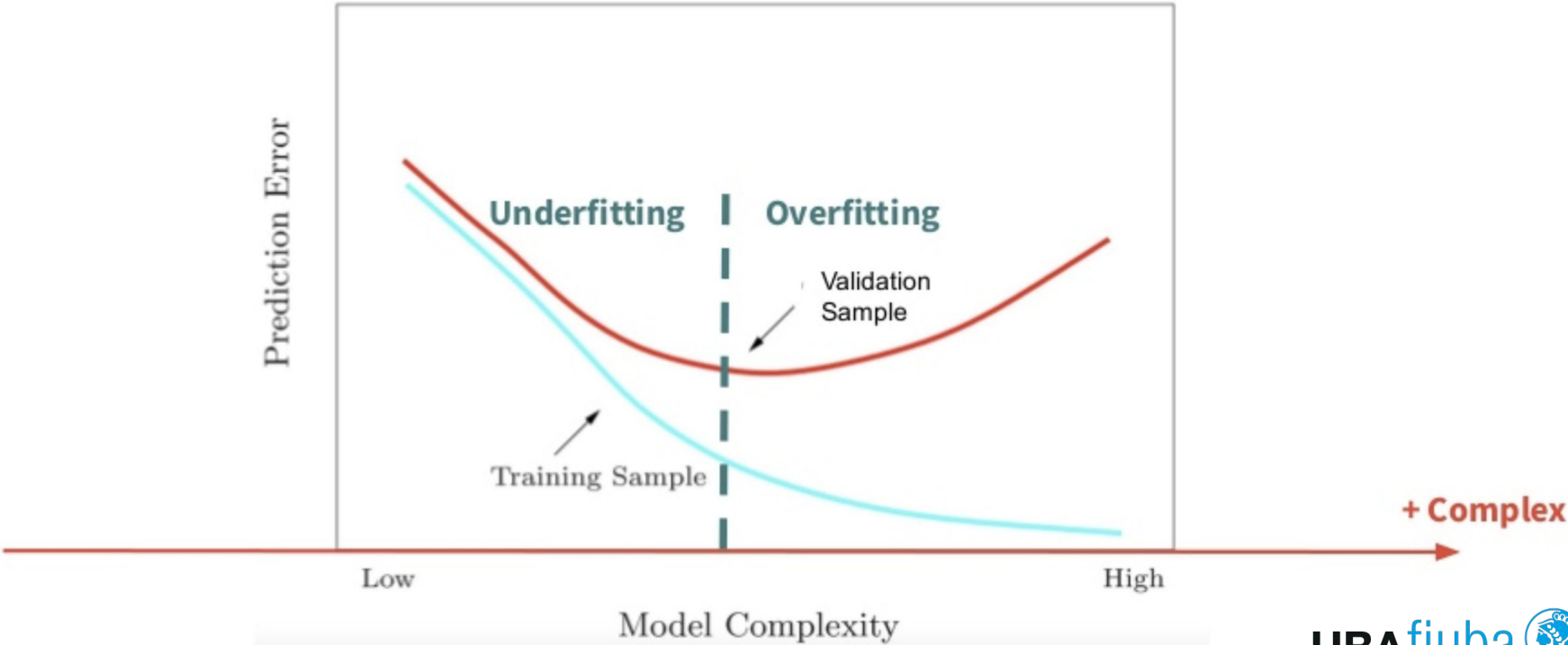for epoch in n_epochs:
- shuffle the batches
- for batch in n_batches:
  - compute the predictions for the batch
  - compute the error for the batch
  - compute the gradient for the batch
  - update the parameters of the model

.**UBA**fiuba
FACULTAD DE INGENIERÍA
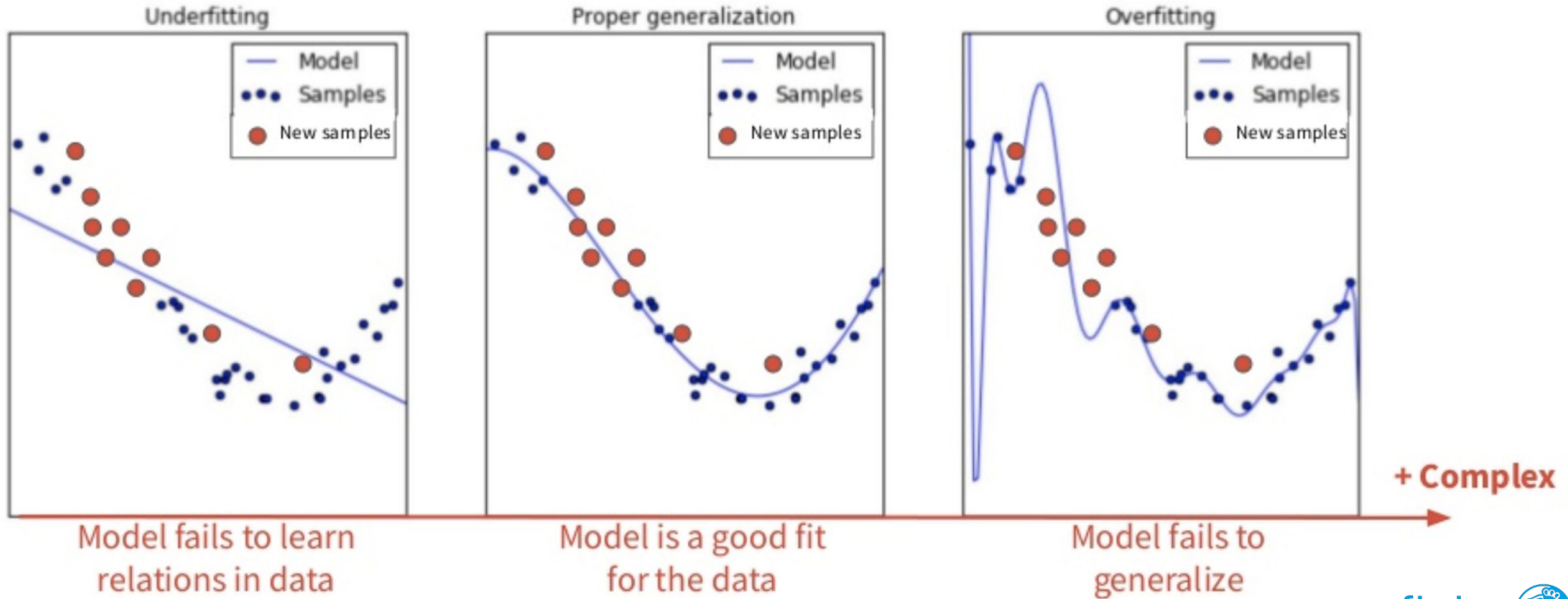
## Comparativa de gradientes

| | Gradient Descent | Stochastic Gradient Descent | Mini-Batch Gradient Descent |
|---|---|---|---|
| **Gradient** | $\nabla \left( \sum_{\text{all samples}} (y_i - f_W(X_i))^2 \right)$ | $\nabla \left( (y_i - f_W(X_i))^2 \right)$ | $\nabla \left( \sum_{\text{batch samples}} (y_i - f_W(X_i))^2 \right)$ |
| **Speed** | Very Fast (vectorized) | Slow (compute sample by sample) | Fast (vectorized) |
| **Memory** | O(dataset) | O(1) | O(batch) |
| **Convergence** | Needs more epochs | Needs less epochs | Middle point between GD and SGD |
| **Gradient Stability** | Smooth updates in params | Noisy updates in params | Middle point between GD and SGD |

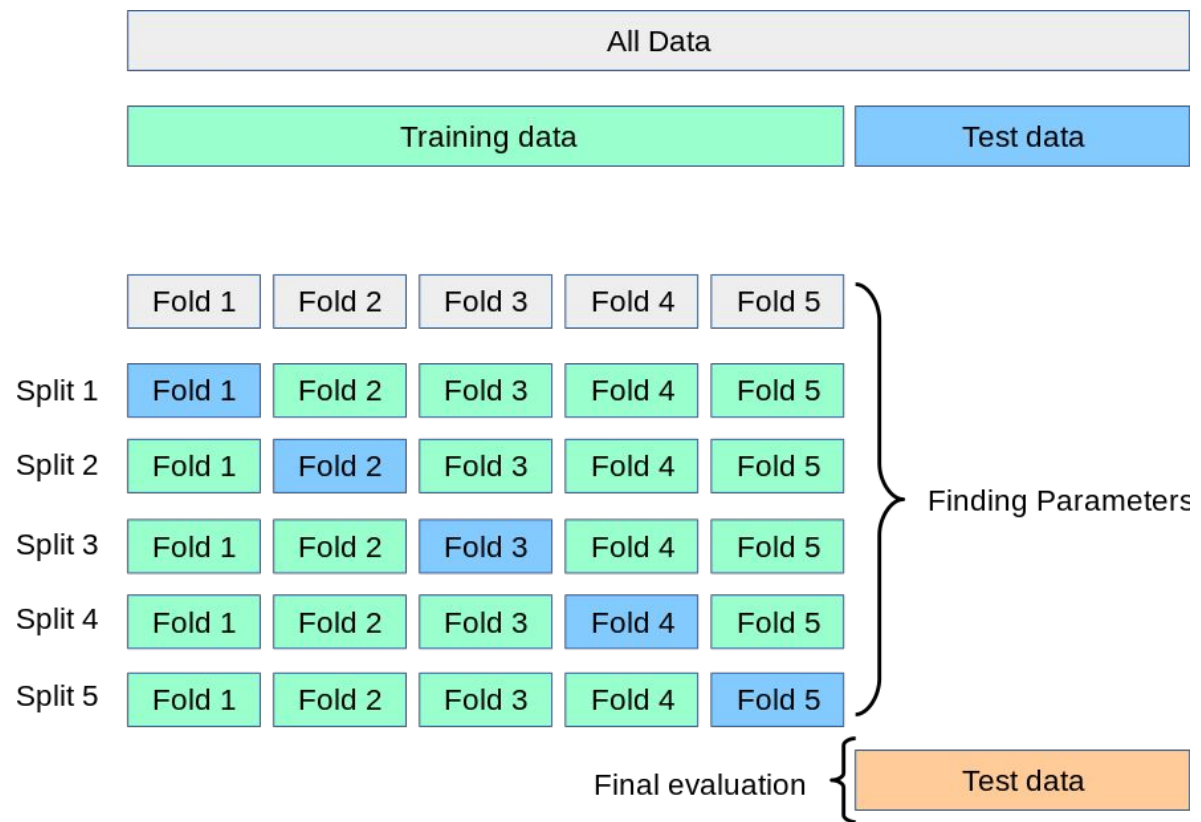.UBAfiuba
FACULTAD DE INGENIERÍA

**Selección de modelos**

## Selección de modelos

## Selección de modelos

## Cross-Validation

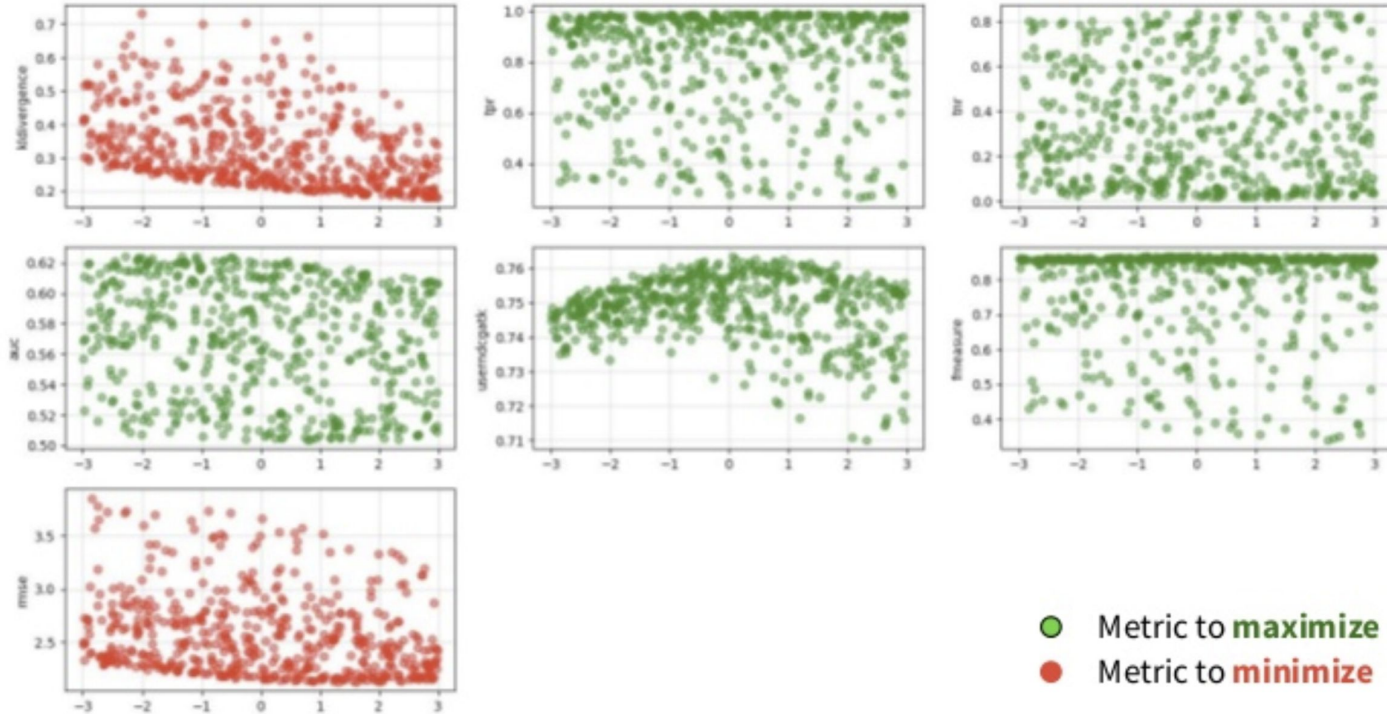## Entrenamiento numérico del modelo seleccionado - Obtención de parámetros

## Selección de los hiper parámetros



**Grid Search**

**Random Search**

Bibliografía

- The Elements of Statistical Learning | Trevor Hastie | Springer
- An Introduction to Statistical Learning | Gareth James | Springer
- Deep Learning | Ian Goodfellow | https://www.deeplearningbook.org/
- Mathematics for Machine Learning | Deisenroth, Faisal, Ong
- Artificial Intelligence, A Modern Approach | Stuart J. Russell, Peter Norvig
- A visual explanation for regularization of linear models - Terence Parr
- A Complete Tutorial on Ridge and Lasso Regression in Python - Aarshay Jain