

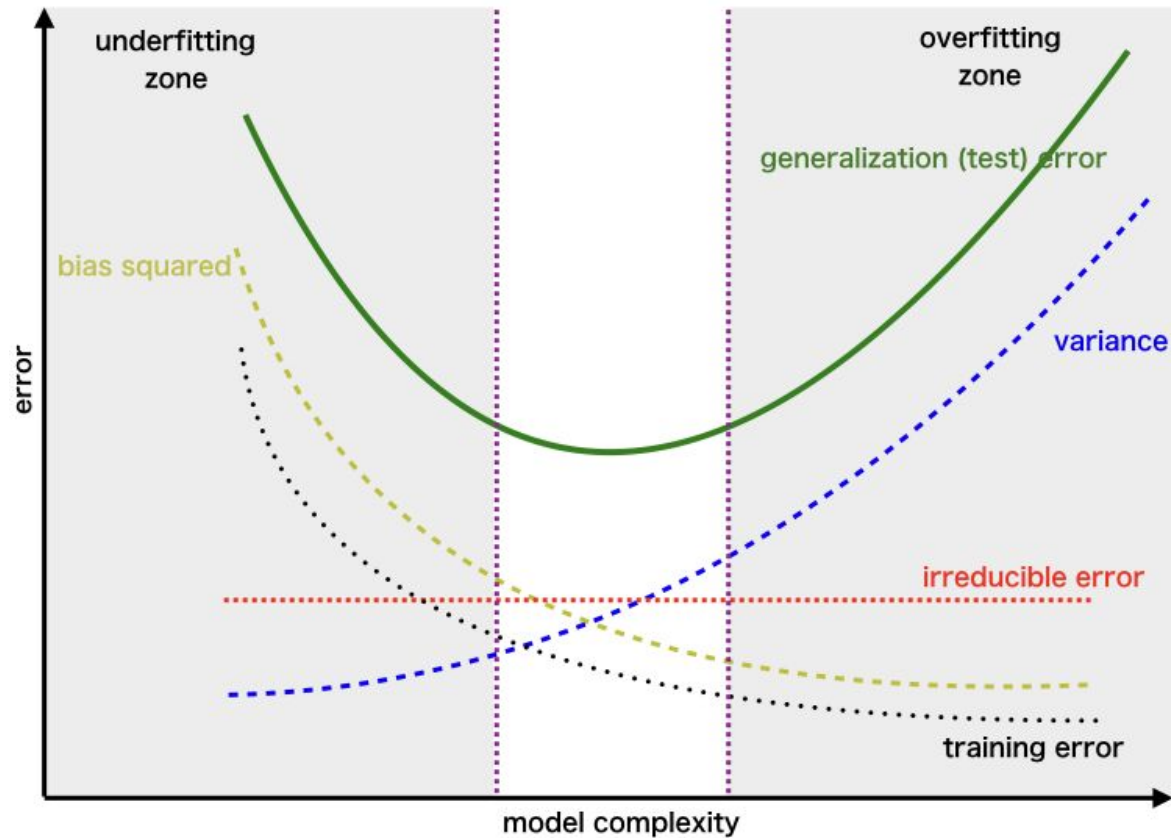
Introducción a la Inteligencia Artificial

Clase 5

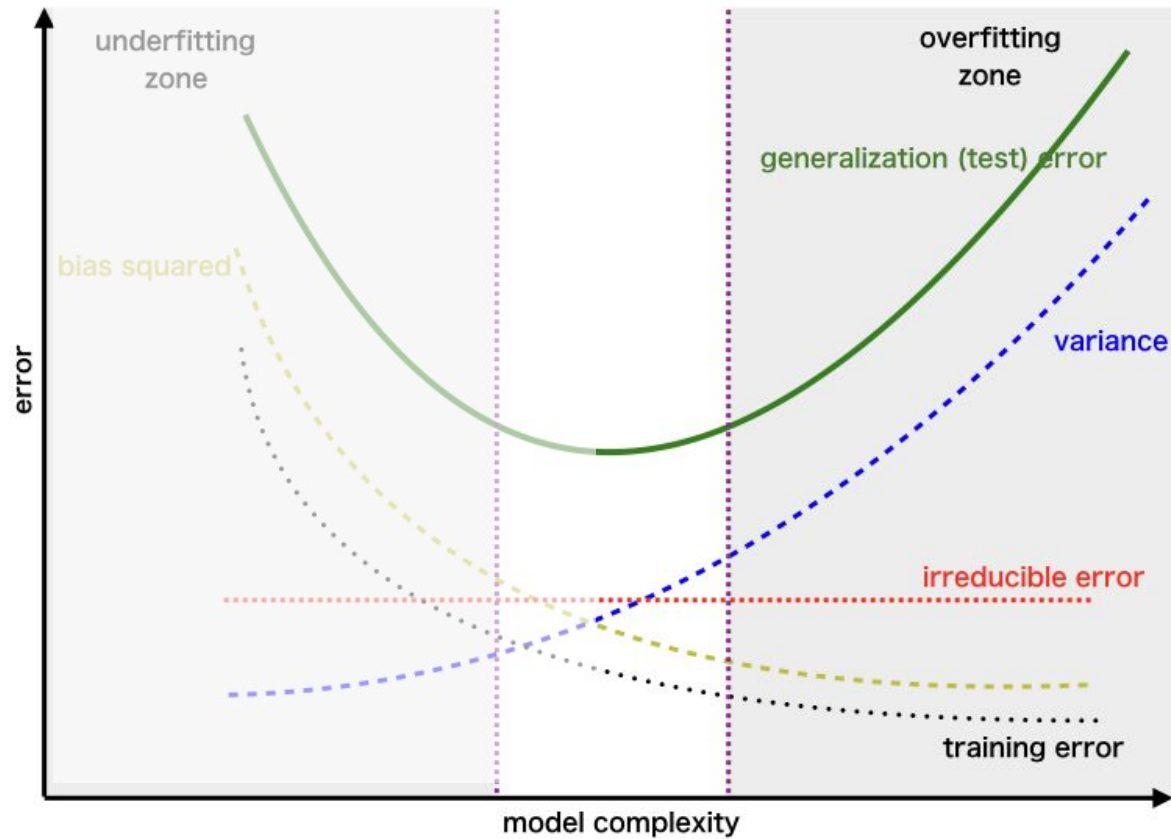


Clase 5

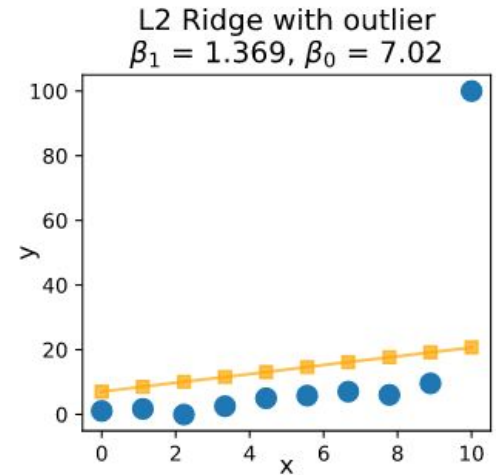
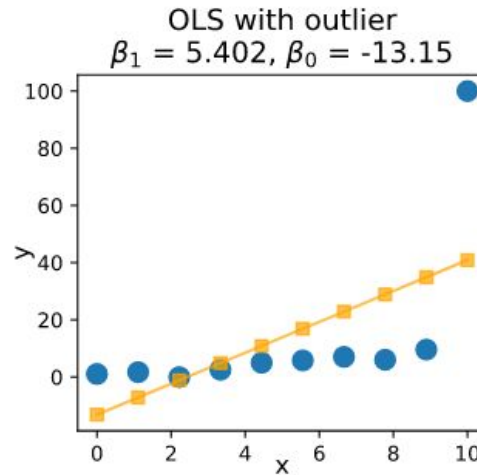
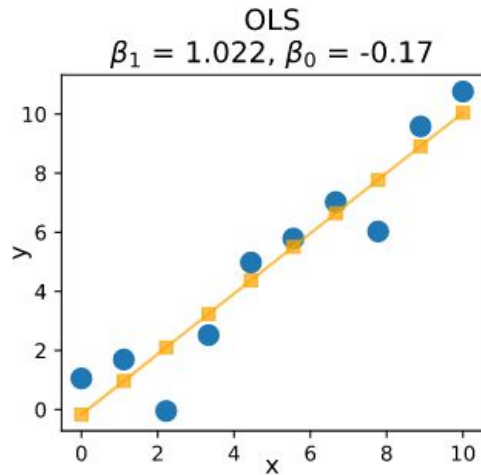
1. Regularización
 - a. Caso general
 - b. Ridge
 - c. Lasso
2. Gradient descent
 - a. GD
 - b. GD Estocástico
 - c. GD Mini-Batch
3. Entrenamiento de modelos
 - a. Selección de modelos
 - b. Cross-Validation



Regularización



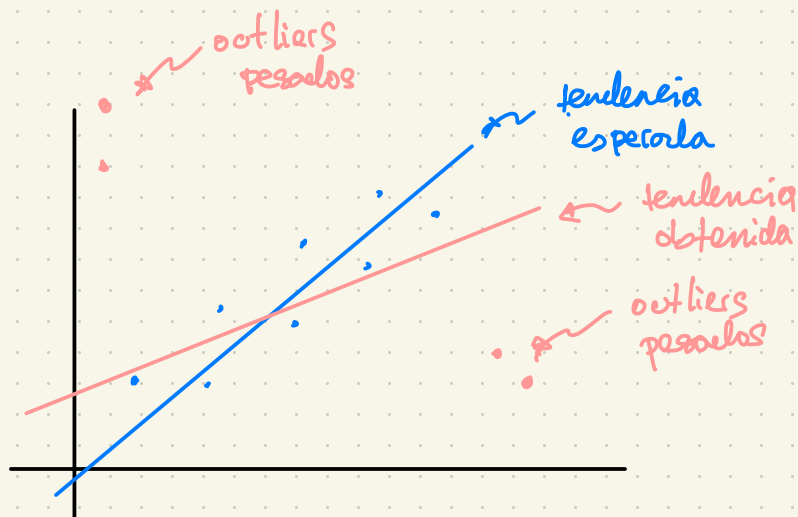
Regularización - Motivación



Riesgo Empírico

$$R(f) = \mathbb{E}(L(f, \hat{f})) \leadsto L(f, \hat{f}) = (f - \hat{f})^2 \rightarrow \text{pérdida cuadrática} \\ \text{penalidad } l_2$$

hasta
ahora
usamos
esto



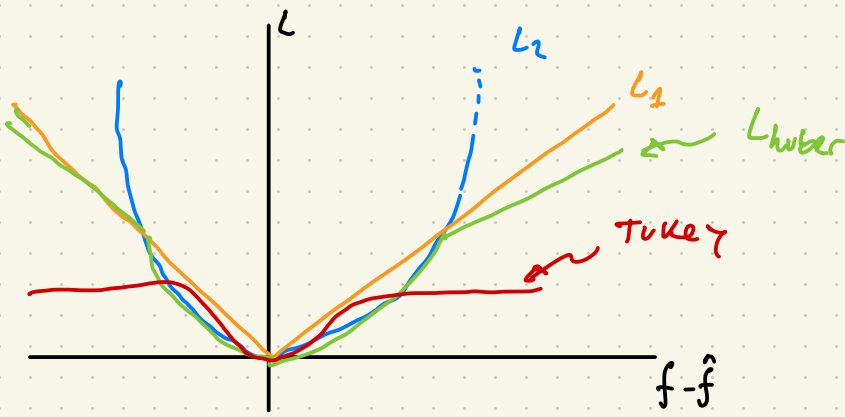
¿cómo puedo mejorar mi tendencia obtenida?

+ Cambiar L :

$$- L_2 = (f - \hat{f})^2 \rightarrow \mathbb{E}(L_2) = \text{ECM}$$

$$- L_1 = |f - \hat{f}| \rightarrow \mathbb{E}(L_1) = \text{EAM}$$

$$- L_I = \mathbb{I}_{|f - \hat{f}| > c} = \begin{cases} 0 & |f - \hat{f}| < c \\ 1 & \text{ow} \end{cases}$$



Cuando queremos robustear el Riesgo empírico vamos a elegir un L adecuado.

\Rightarrow si utilizo modelos robustos no en general tengo más parámetros para optimizar

siempre buscamos minimizar R , si no se puede (o no se justifica) usar regresores robustos \Rightarrow podemos regularizar.

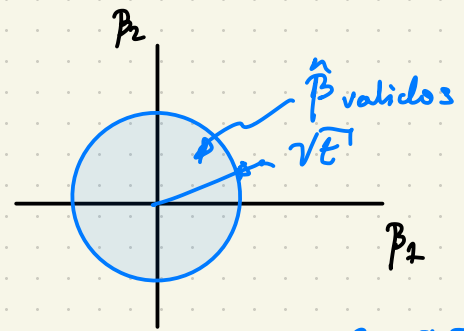
Regularización:

Con reg. buscamos min. el riesgo emp R al mismo tiempo que restringimos (o limitamos) el comportamiento de los parámetros (parameter shrinkage). no buscamos disminuir el error de representación

partimos de $\hat{y} = \beta_0 + \sum_i \beta_i x_i$ no $\hat{\beta} = \arg\min_{\beta} \sum_i (y_i - \sum_j \beta_j x_{ij})^2$

Vamos a condicionar β :

$$\begin{cases} \arg\min_{\beta} \sum_i (\dots)^2 \\ \sum_j \beta_j^2 \leq t \quad (\|\beta\| \leq \sqrt{t}) \end{cases}$$



con esto: $\hat{\beta} = \arg\min_{\beta} \left(\sum_i y_i - \sum_j \beta_j x_j \right)^2 - \lambda \sum_j \beta_j^2$

en g'ral es q_2

término de regularización

parámetro de complejidad

término de regularización = weight decay

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

Observado - Predicción

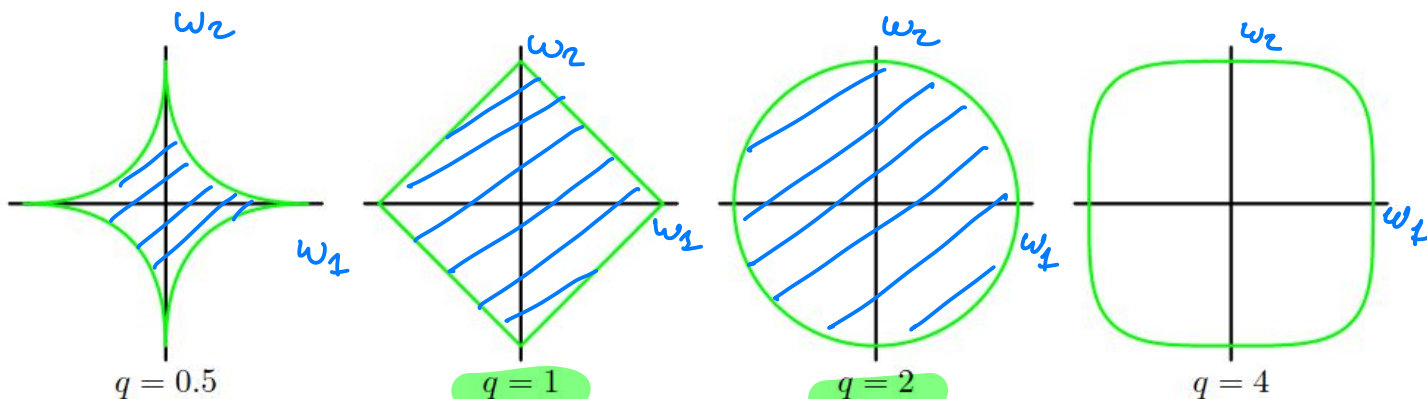


\mathbf{w} está “libre”

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

Término de
regularización
“weight decay”

→ w afecta la pérdida



▣ \hat{w} válidos

Lasso

Ridge

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

Maximum A Posteriori como regularización

Distribución “a priori” de los parámetros \longrightarrow Observar data \longrightarrow Actualizar distribución (Posterior)

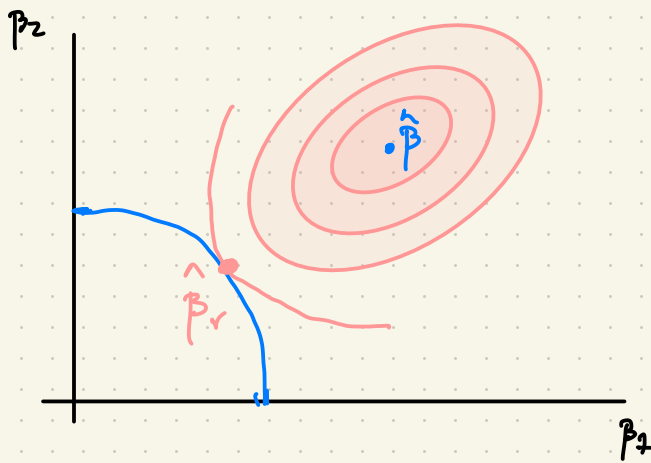
$$p(w) \sim D(\theta)$$

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

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

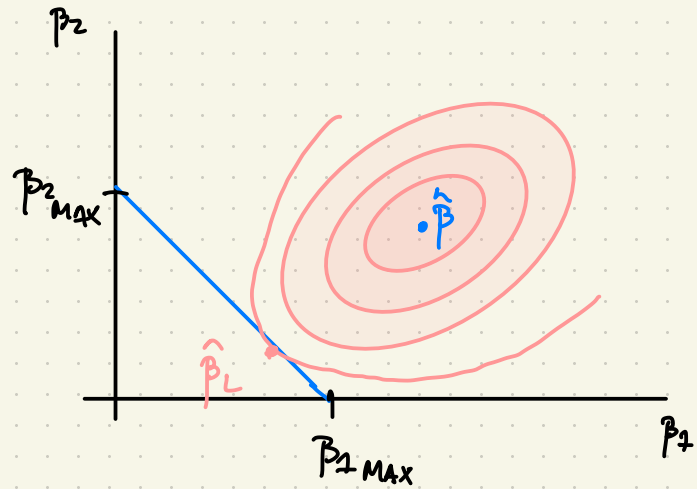
$$w_{map} = (\Phi^T \Phi + \frac{\sigma^2}{b^2} I)^{-1} \Phi^T y$$

Gaussian prior con varianza b^2



Ridge

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



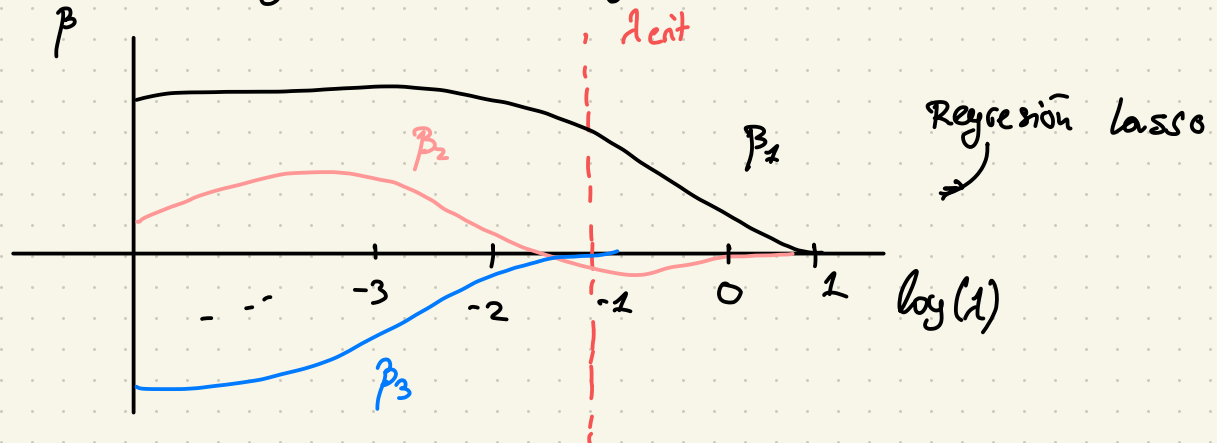
Lasso

Recordemos que $\hat{\beta}$ en gual tiene curvas de nivel elípticas, esto es por construcción del estimador.

El costo de regularizar es que nos alejamos del óptimo

¿Cómo funciona?

1. plantear (elegir) q (lasso 1, ridge 2)
2. elegimos un vector de λ 's 'apropiado' ($\lambda \uparrow \Rightarrow \uparrow$ penalidad)
3. optimizar con un $\lambda_i \Rightarrow \text{RSS}(\lambda; q) = \|y - X\beta\|^2 + \lambda \|\beta\|_q$
4. Calculamos las métricas (Error de representación, bono del de ajuste, etc...)
5. Comparamos y elegimos el mejor.



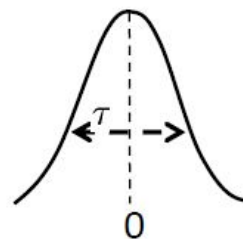
Maximum A Posteriori como regularización - Ridge (L2)

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

I) Gaussian Prior

$$\beta \sim \mathcal{N}(0, \tau^2 \mathbf{I})$$

$$p(\beta) \propto e^{-\beta^T \beta / 2\tau^2}$$



$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2 \quad \text{Ridge Regression}$$

↓
constant(σ^2, τ^2)

$$\hat{\beta}_{\text{MAP}} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

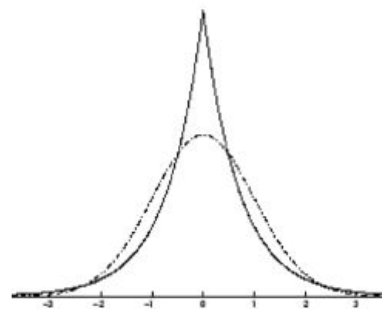
Maximum A Posteriori como regularización - LASSO (L1)

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

II) Laplace Prior

$$\beta_i \stackrel{iid}{\sim} \text{Laplace}(0, t)$$

$$p(\beta_i) \propto e^{-|\beta_i|/t}$$

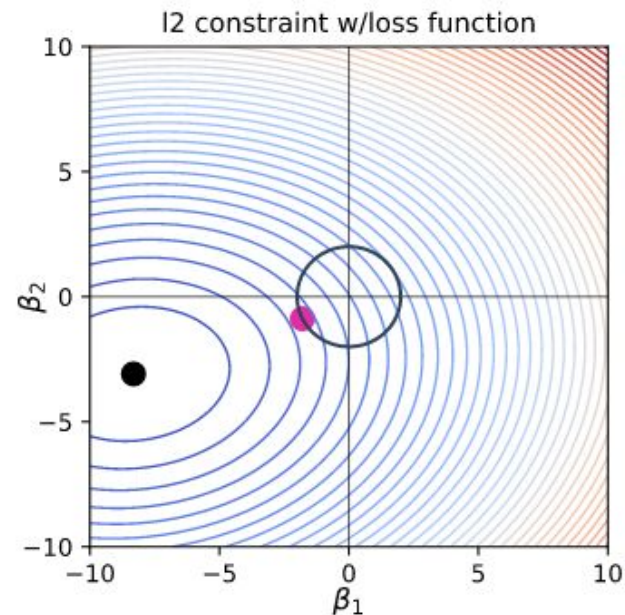
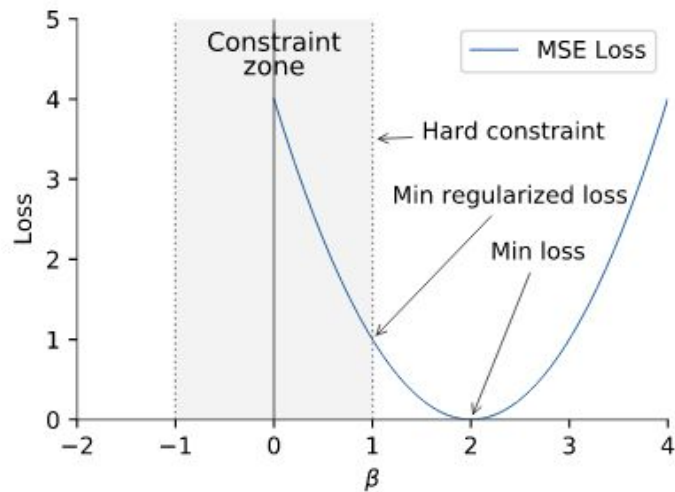


$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_1$$

\downarrow
constant(σ^2, t)

Lasso

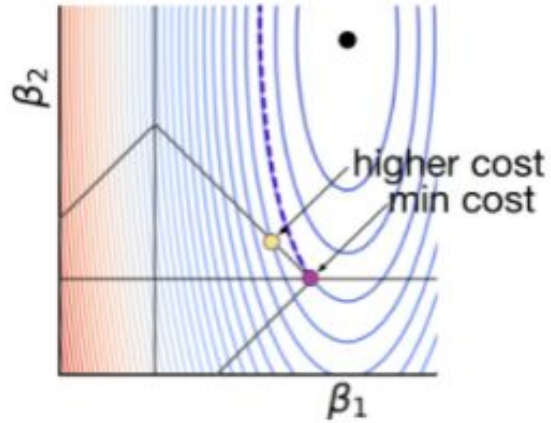
Regularización



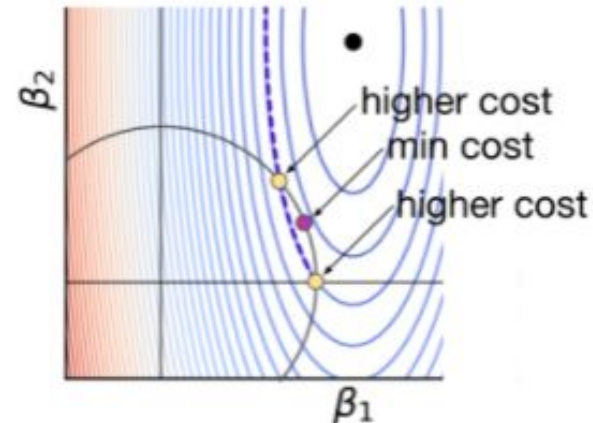
Ejemplo interactivo:
<https://shinylive.io/py/app/#regularization>

Regularización

(a) L1 Constraint Diamond



(b) L2 Constraint Circle

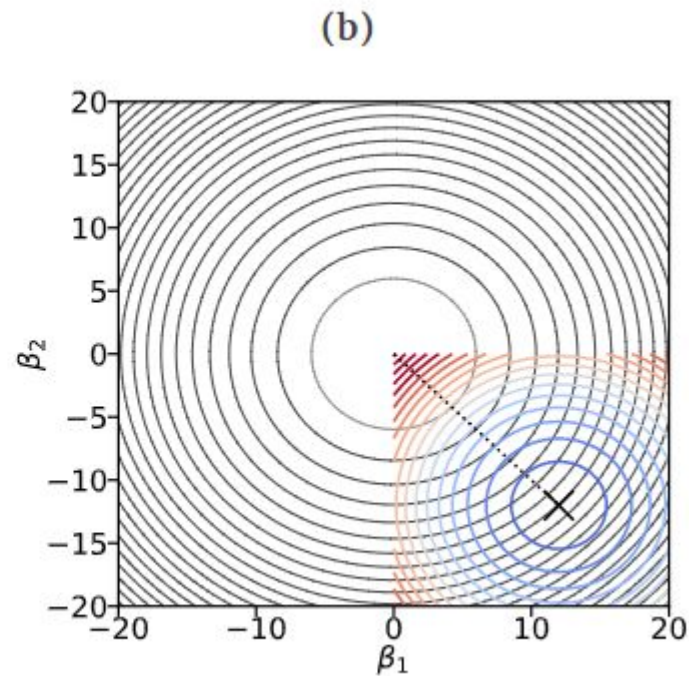
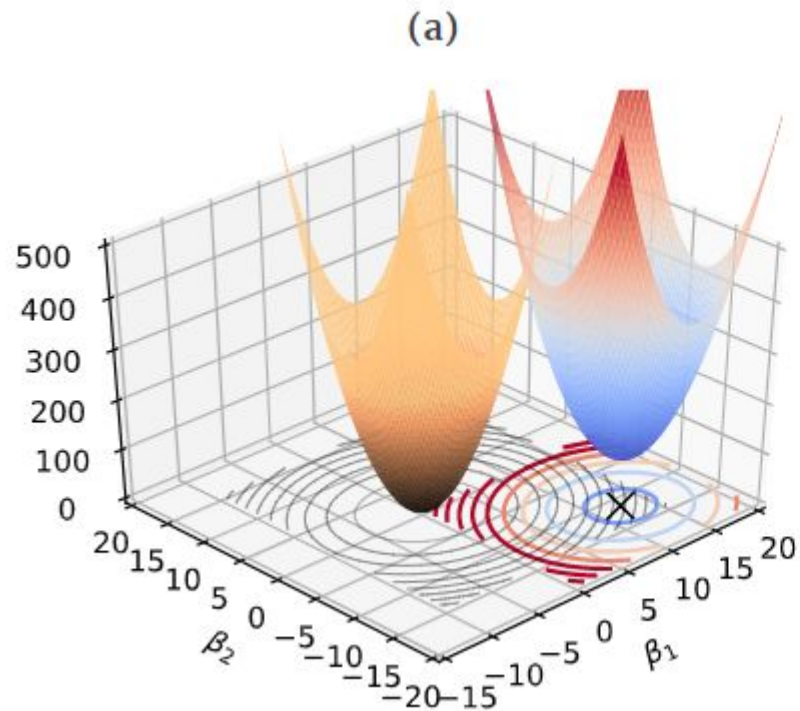


ElasticNet

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

¿Qué β se reduce más?

Regularización



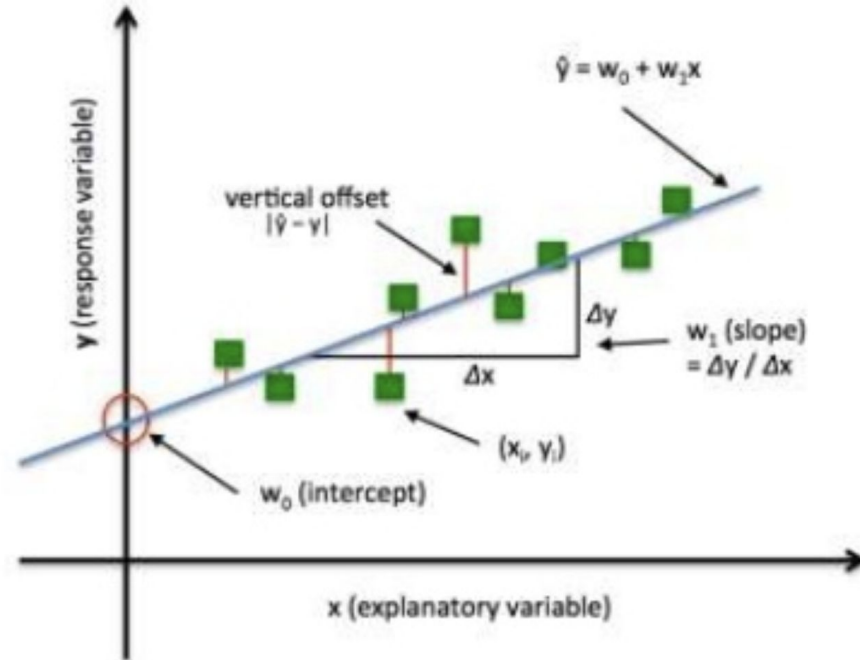
Gradiente Descendente

Implementación de Gradiente Descendente

Solucion analitica

$$\min_W \|Y - XW\|_2^2$$

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



Implementación de Gradiente Descendente

Solución numérica

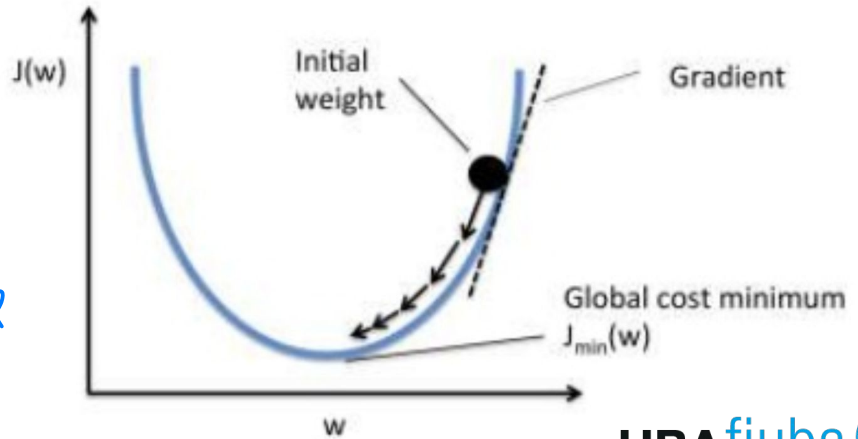
$$\min_W \|Y - XW\|_2^2 \implies \min_W \sum_i (y_i - X_i \cdot W)^2$$

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

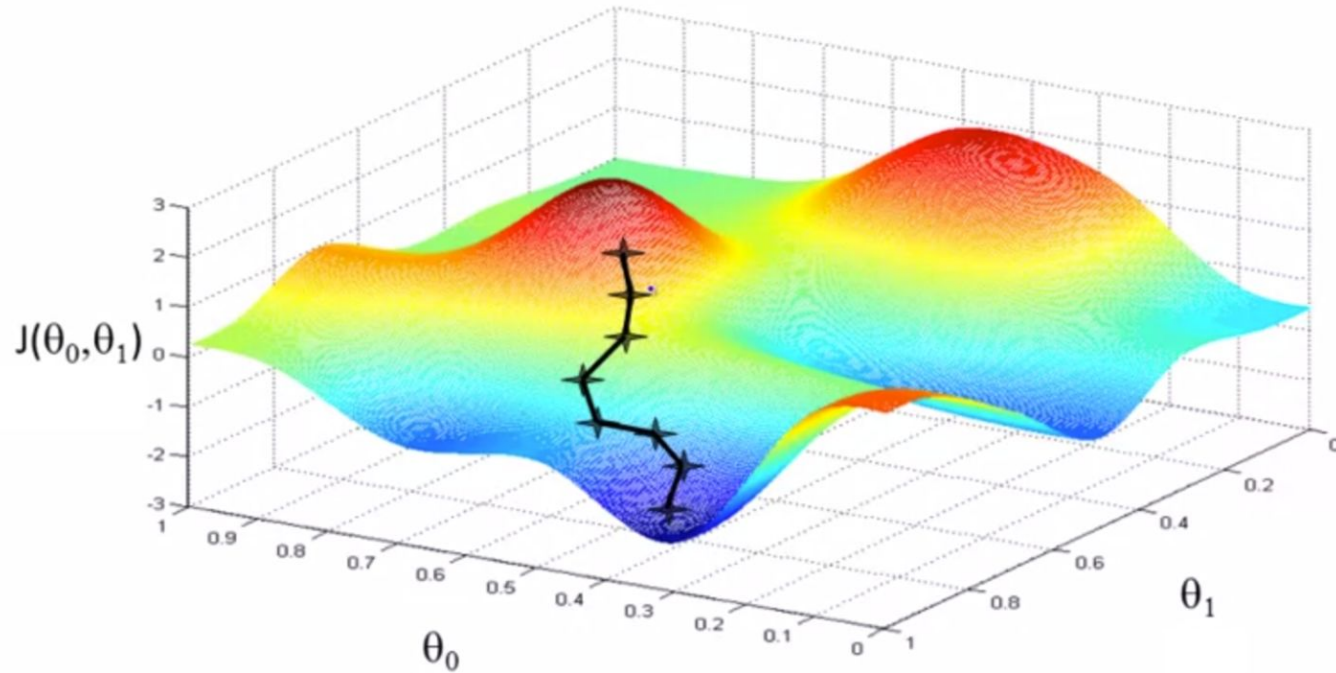
$$w_{t+1} = w_t - \alpha \nabla R$$

↖ gradiente del funcional

↗ learning rate

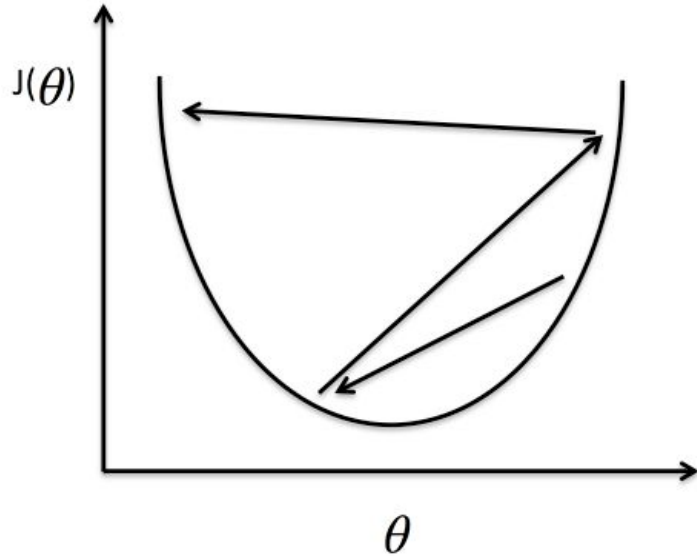


Gradiente Descendente

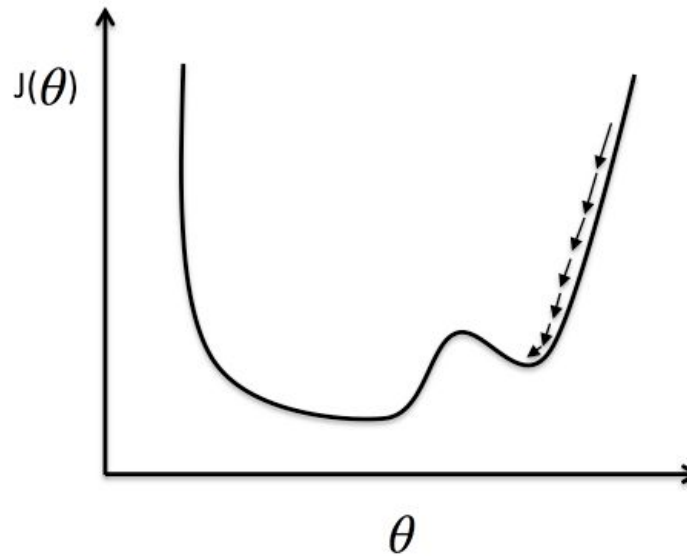


Andrew Ng

Gradiente Descendente



Large learning rate: Overshooting.



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

Implementación de Gradiente Descendente

Solución numérica

$$\begin{aligned}\nabla_w J(w) &= \nabla_w \left(\sum_i (y_i - X_i W)^2 \right) \\ &= \sum_i \left(\nabla_w (y_i - X_i W)^2 \right) \\ &= \sum_i \left(\nabla_w (y_i - (x_{i1}w_1 + x_{i2}w_2 + \dots + x_{im}w_m))^2 \right) \\ &= \sum_i \left(-2(y_i - \hat{y}_i)x_{ij} \right) \quad \forall j \in (1 \dots m)\end{aligned}$$

en sklearn Linear_model, SGDRegressor

Implementación de Gradiente Descendente

Solución numérica

$$w_{t+1} = w_t - \alpha \nabla_w J(w)$$

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

Gradient Descent algorithm (datos w_{inicial})

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

Implementación de Gradiente Descendente Estocástico

Solución numérica

$$\nabla ((y_i - f_W(X_i))^2)$$

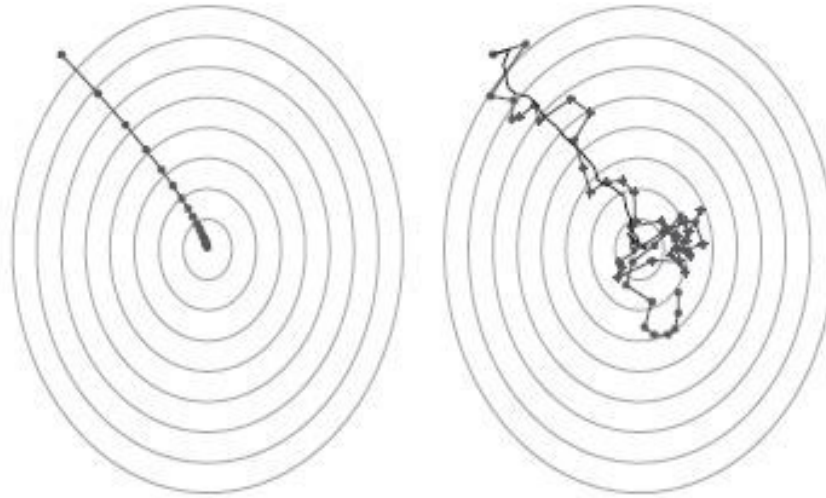
Stochastic Gradient Descent algorithm

for epoch in n_epochs: (Samples = 10^6)

- shuffle the samples $n_samples = 10^2$
- **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

Implementación de Gradiente Descendente Estocástico

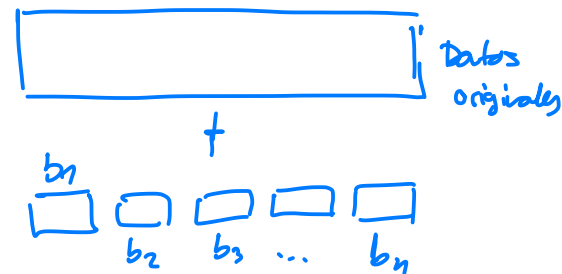
Solución numérica



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:

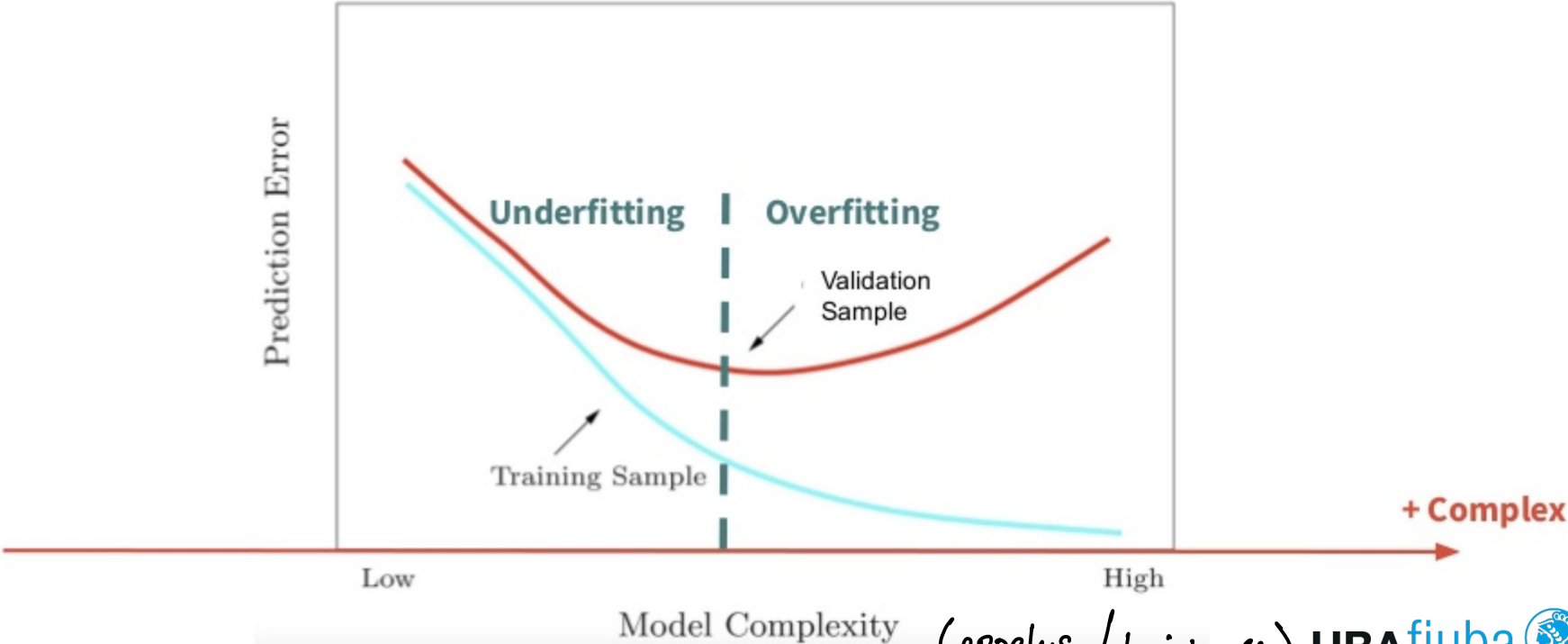
- shuffle the batches
- **for batch in n_batches:**
 - compute the predictions for **the batch**
 - compute the error for the batch
 - compute the gradient for **the batch**
 - update the parameters of the model

Comparativa de gradientes

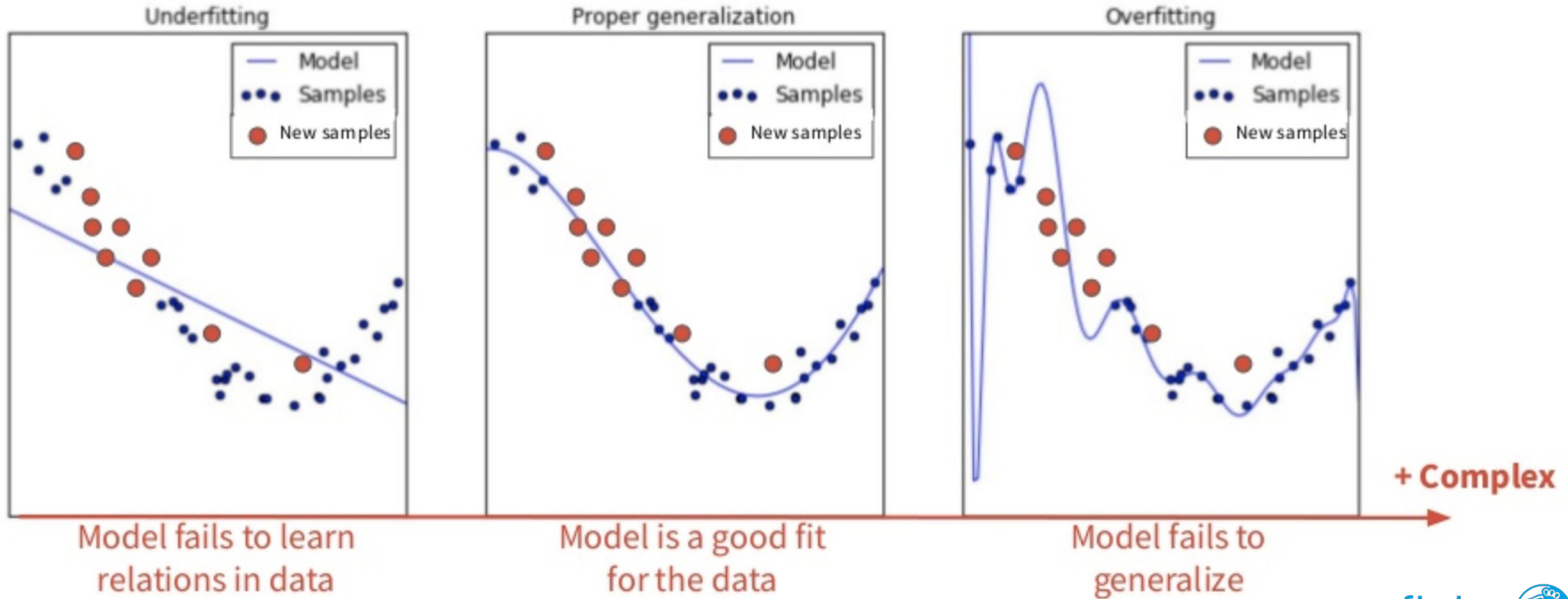
	Gradient Descent	Stochastic Gradient Descent	Mini-Batch Gradient Descent
Gradient	$\nabla \left(\sum_{\text{all samples}} (y_i - f_W(X_i))^2 \right)$	$\nabla \left((y_i - f_W(X_i))^2 \right)$	$\nabla \left(\sum_{\text{batch samples}} (y_i - f_W(X_i))^2 \right)$
Speed	Very Fast (vectorized)	Slow (compute 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

Selección de modelos

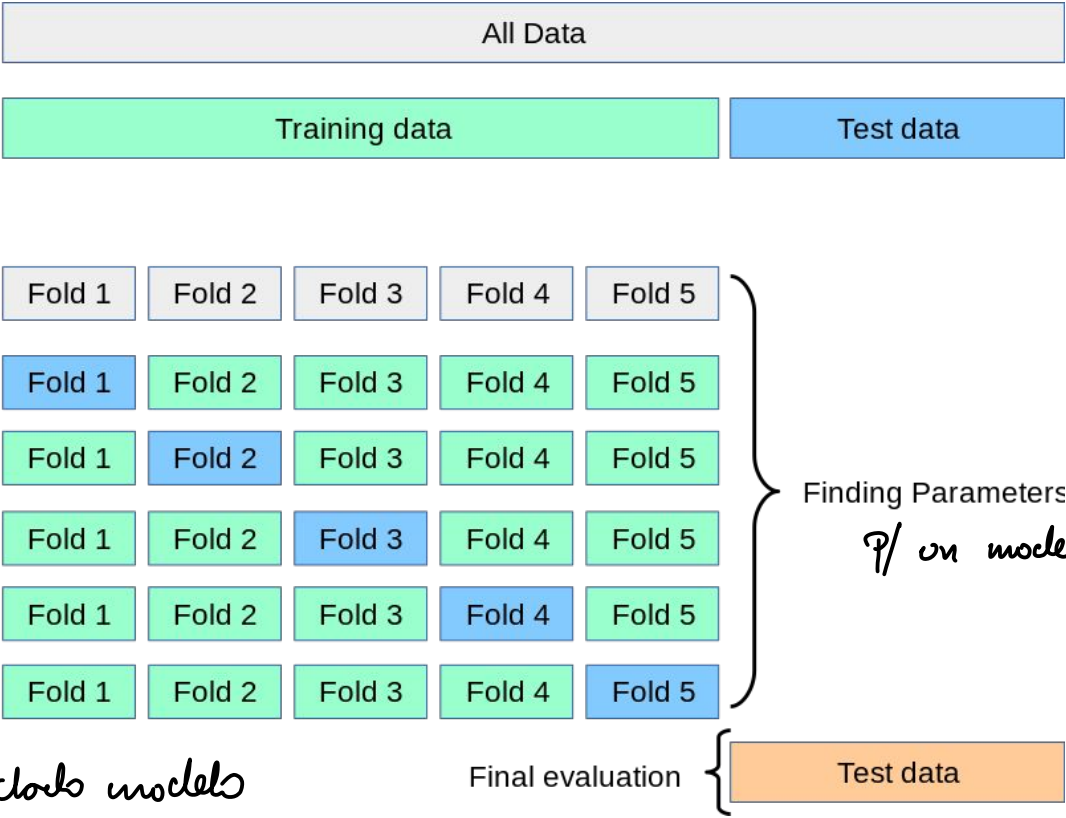
Selección de modelos



Selección de modelos



Cross-Validation



Δ_{cs_1}

→ Split 1

Δ_{cs_2}

→ Split 2

⋮

Split 3

⋮

Split 4

Δ_{cs_5}

→ Split 5

$\overline{\Delta_c}$

para un dato modelo

$models = \{ \dots \}$

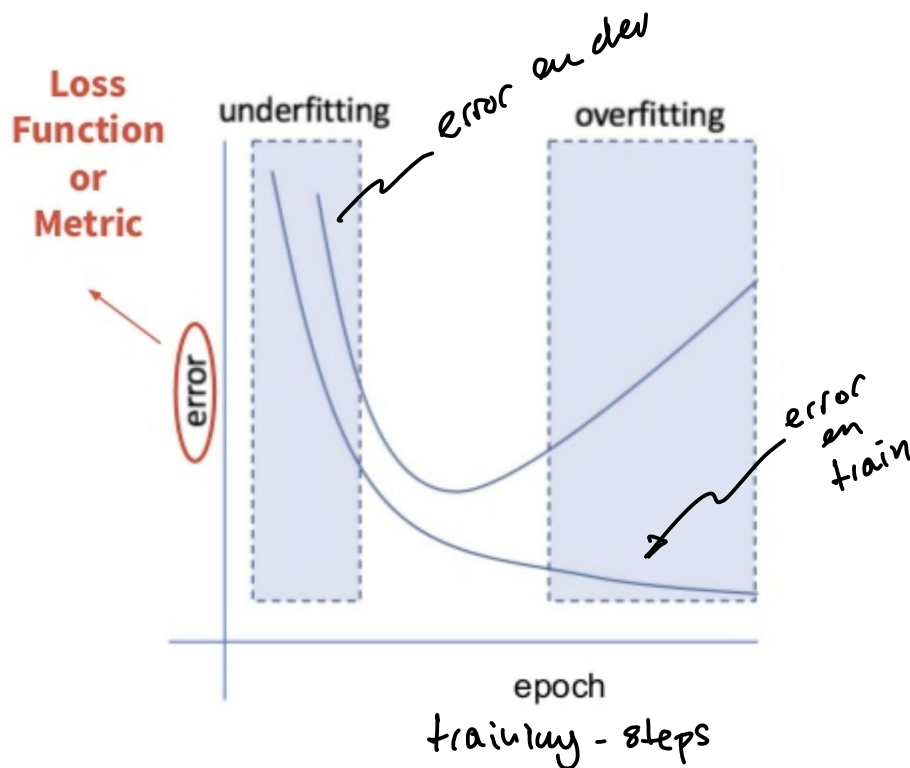
$SVM(k=1, dist=eu)$

$SVM(k=rbf, dist=eu)$

$SVM(k=rbf, dist=m)$

Finding Parameters
P/ on modelo

Entrenamiento numérico del modelo seleccionado - Obtención de parámetros



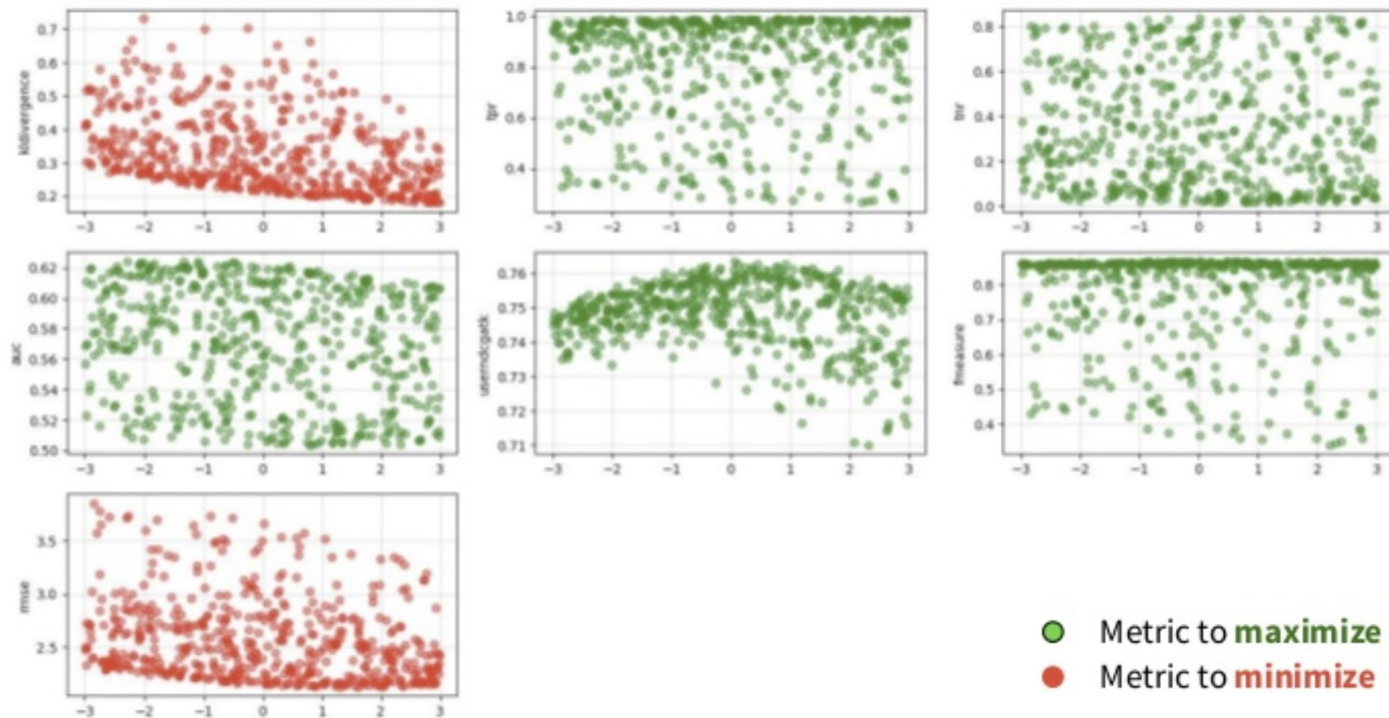
Dataset \rightarrow train 80%
 \hookrightarrow dev 5%
 \hookrightarrow test 15%

Mini-Batch Gradient Descent

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
- plot error vs epoch

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
- Deep Learning con Pytorch, Chintala