



Aplicando Modelos de Machine Learning

Luis Zúñiga



Agenda

1. Evaluación de Modelos de ML
2. Bias-Variance Tradeoff
3. Curvas de Aprendizaje
4. Juntando todo...



Evaluación de Modelos de ML

¿Qué modelo elegir?

Redes
Neuronales

Regresión
Logística

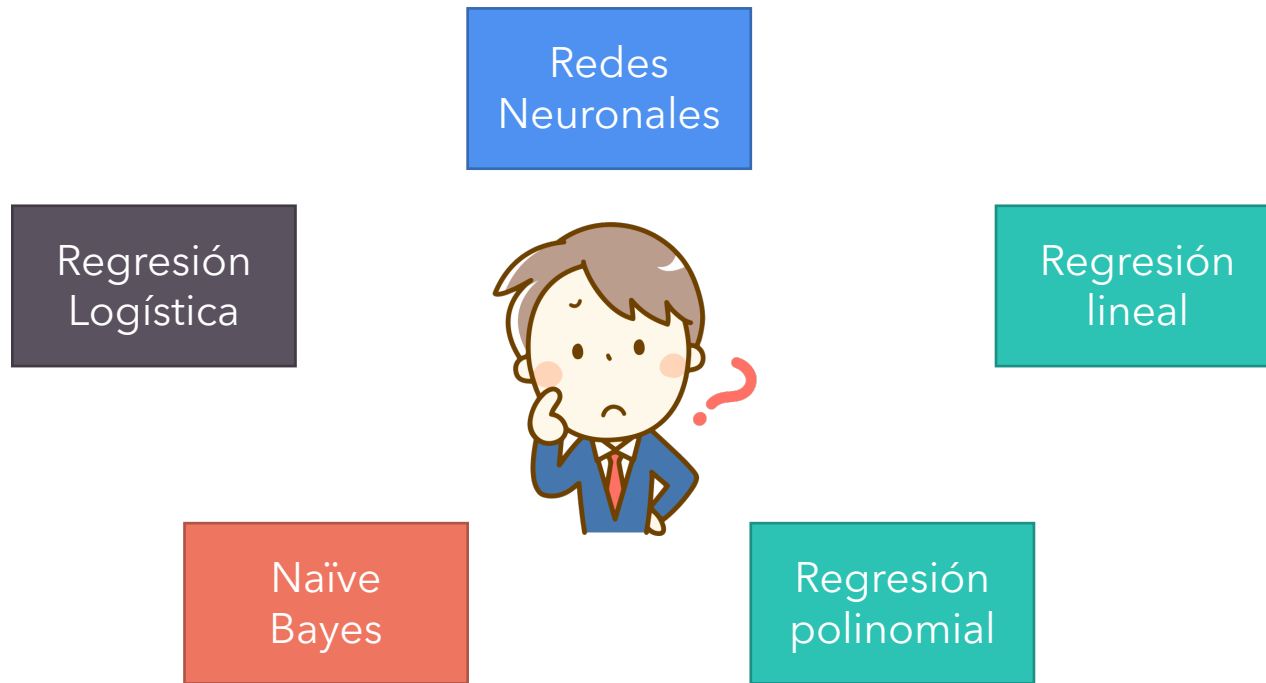
Regresión
lineal



Naïve
Bayes

Regresión
polinomial

¿Qué modelo elegir?



- ¿Es buena idea probar todos los modelos?
- ¿Cómo elegir un buen modelo candidato?
- Si obtengo malos resultados, ¿qué puedo hacer sin perder mucho tiempo?

¿Qué modelo elegir?

Supongamos que entrenamos un modelo para predecir **el precio de una casa con un modelo de regresión lineal con regularización:**

$$J(\theta) = \frac{1}{2n} \left[\sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^p \theta_j^2 \right]$$

Desafortunadamente, al probar nuestro modelo entrenado (la hipótesis) con datos que jamás llegó a ver, nos damos cuenta que **las predicciones realizadas son malas.**

¿Qué se puede hacer?

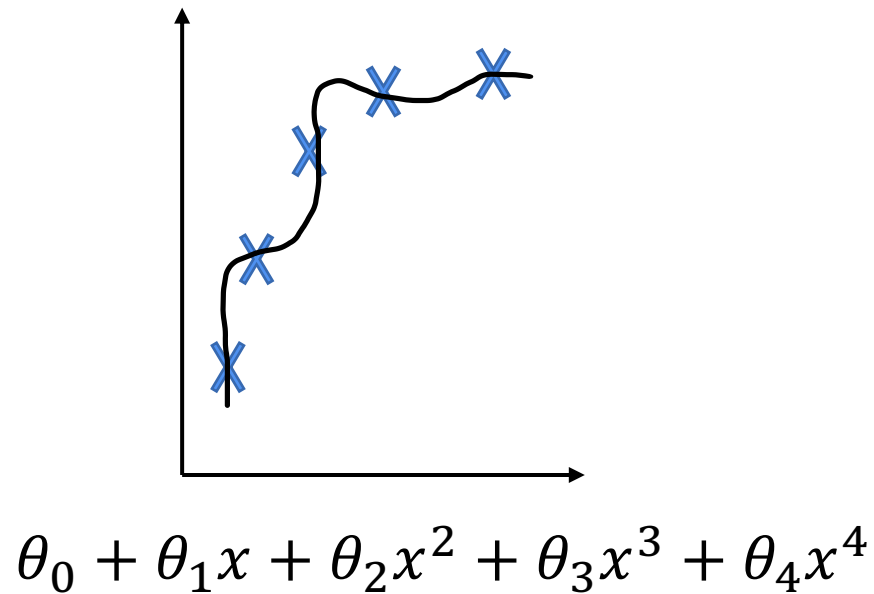
¿Qué modelo elegir?

Después de mucho pensar, llegan a las siguientes opciones:

- Tener más datos para el entrenamiento.
- Considerar menos características.
- Obtener más características.
- Considerar combinaciones de características más complejas.
- Aumentar o reducir el valor de λ .



Evaluación del modelo (hipótesis)



Uno de los problemas que puede surgir al momento de entrenar los modelos es que exista sobreajuste o desajuste.

¿Cómo se detectaba?

Al probar el modelo con datos que jamás vio en el entrenamiento, tiene un pésimo rendimiento.

Evaluación del modelo (hipótesis)

Tamaño (m ²)	Precio (m ²)
150	2500
60	3230
60	3150
80	2120
60	1500
75	1670
300	4350
250	3750
150	2550
60	4500

Conjunto de
entrenamiento
(70%-80%)
al azar

Se entrena con la
partición que
corresponde al conjunto
de entrenamiento.

Conjunto de
prueba
(20%-30%)
al azar

Se evalúa con la
partición que
corresponde al conjunto
de prueba.

Train/Test - Procedimiento

Para el caso de regresión:

- Obtener los parámetros θ mediante el entrenamiento con el conjunto de entrenamiento (minimizando el error mediante gradiente descendiente $J(\theta)$).
- Con los parámetros θ encontrados, determinar el error con el conjunto de prueba:

$$J_{test}(\theta) = \frac{1}{2n_{test}} \sum_{i=1}^{n_{test}} \left(h_{\theta} \left(x_{test}^{(i)} \right) - y_{test}^{(i)} \right)^2$$

Train/Test - Procedimiento

Para el caso de clasificación:

- Obtener los parámetros θ mediante el entrenamiento con el conjunto de entrenamiento

$$J(\theta) = -\frac{1}{n_{test}} \sum_{i=1}^{n_{test}} [y_{test}^{(i)} \log(h_{\theta}(\mathbf{x}_{test}^{(i)})) + (1 - y_{test}^{(i)}) \log(1 - h_{\theta}(\mathbf{x}_{test}^{(i)}))]$$

- Con los parámetros θ encontrados, determinar el error con el conjunto de prueba:

$$err(h_{\theta}(\mathbf{x}), y) = \begin{cases} 1 & \text{si } h_{\theta}(\mathbf{x}) \geq 0.5, y = 0 \\ 1 & \text{si } h_{\theta}(\mathbf{x}) < 0.5, y = 1 \\ 0 & \text{en cualquier otro caso} \end{cases}$$

$$Test\ error = \frac{1}{n_{test}} \sum_{i=1}^{n_{test}} err(h_{\theta}(\mathbf{x}_{test}^{(i)}), y_{test}^{(i)})$$

Accuracy, Recall, Precision o F1

Selección de Modelos – Conjunto de Validación

Supongamos que tenemos los siguientes modelos:

1. $h_{\theta}(x) = \theta_0 + \theta_1 x$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$
- \vdots
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10}$

Existe un nuevo parámetro d que rige el grado del polinomio de regresión lineal.

¿Cómo elegir el mejor?

Selección de Modelos – Conjunto de Validación

Supongamos que tenemos los siguientes modelos:

	Entrenamiento	Prueba
1. $h_{\theta}(x) = \theta_0 + \theta_1 x$	θ_1	
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x$	θ_2	
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x + \theta_3 x$	θ_3	θ_5
\vdots	\vdots	
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x$	θ_{10}	

¿Ven algo mal?

Selección de Modelos – Conjunto de Validación

Supongamos que tenemos los siguientes modelos:

1. $h_{\theta}(x) = \theta_0 + \theta_1 x$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x + \theta_3 x$
- \vdots
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x$

Al elegir θ_5 usamos **el conjunto de prueba para comparar** los modelos entre sí y elegir el mejor.

Es una **estimación optimista del error de generalización** por el parámetro adicional d .

Selección de Modelos – Conjunto de Validación

Tamaño (m ²)	Precio (m ²)
150	2500
60	3230
60	3150
80	2120
60	1500
75	1670
300	4350
250	3750
150	2550
60	4500

Conjunto de
entrenamiento
(60%)

Conjunto de
prueba
(20%)

Conjunto de
validación
(20%)

Se entrena con la
partición que
corresponde al conjunto
de entrenamiento.

Se evalúa con la partición
que corresponde al
conjunto de prueba.

Se elige la que menor
error tenga en el conjunto
de validación como
medida de
generalización.

Selección de Modelos – Conjunto de Validación

Tamaño (m ²)	Precio (m ²)
150	2500
60	3230
60	3150
80	2120
60	1500
75	1670
300	4350
250	3750
150	2550
60	4500

Conjunto de
entrenamiento
(60%)

Conjunto de
prueba
(20%)

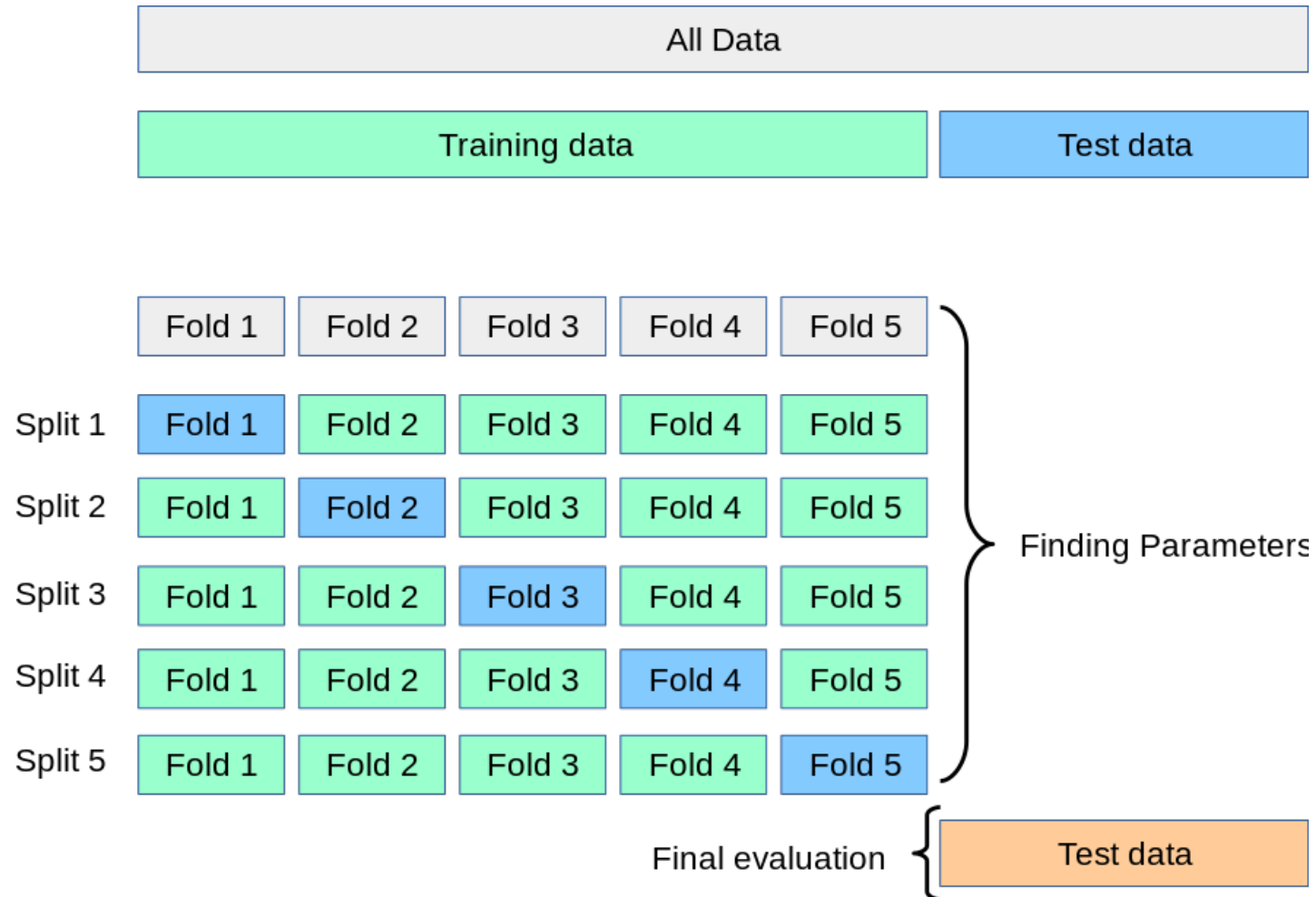
Conjunto de
validación
(20%)

Se entrena con la
partición que
corresponde al conjunto
de entrenamiento.

Se evalúa con la partición
que corresponde al
conjunto de prueba.

Se elige la que menor
error tenga en el conjunto
de validación como
medida de
generalización.

Selección de Modelos – Validación Cruzada



Selección de Modelos – Validación Cruzada

"No Free Lunch" :(

D. H. Wolpert. The supervised learning no-free-lunch theorems. In *Soft Computing and Industry*, pages 25–42. Springer, 2002.

Our model is a simplification of reality



Simplification is based on assumptions (model bias)



Assumptions fail in certain situations

Roughly speaking:

"No one model works best for all possible situations."



Bias - Variance

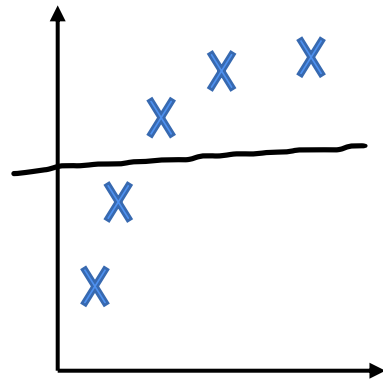
Bias - Variance

Como habíamos comentado hace unas semanas, el problema del *bias - variance tradeoff* es importante determinarlo. Si nuestro modelo no es bueno, puede:

- Tener un alto sesgo (bias).
- Tener alta varianza (variance).

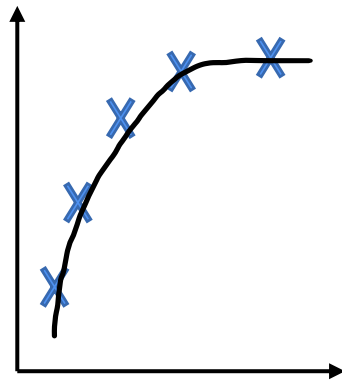
¿Cómo identificarlo?

Bias - Variance



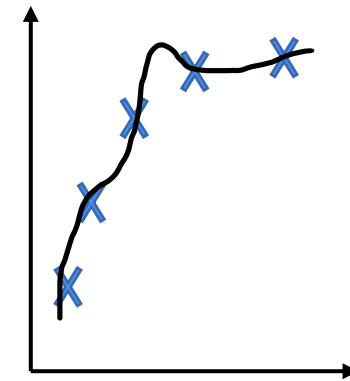
$$\theta_0 + \theta_1 x$$

Underfitting
(high bias)



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

Just Right



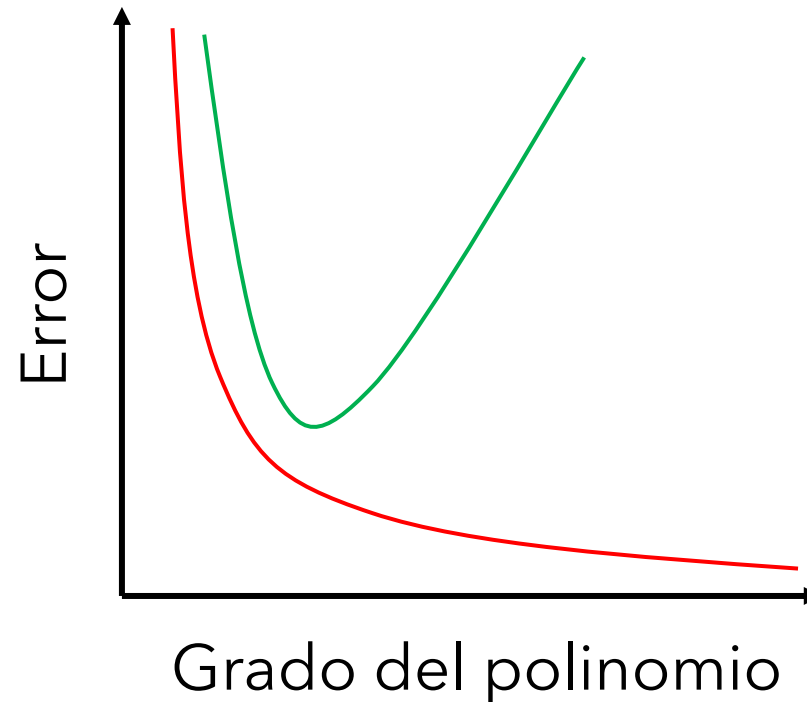
$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Overfitting
(high variance)

Bias - Variance

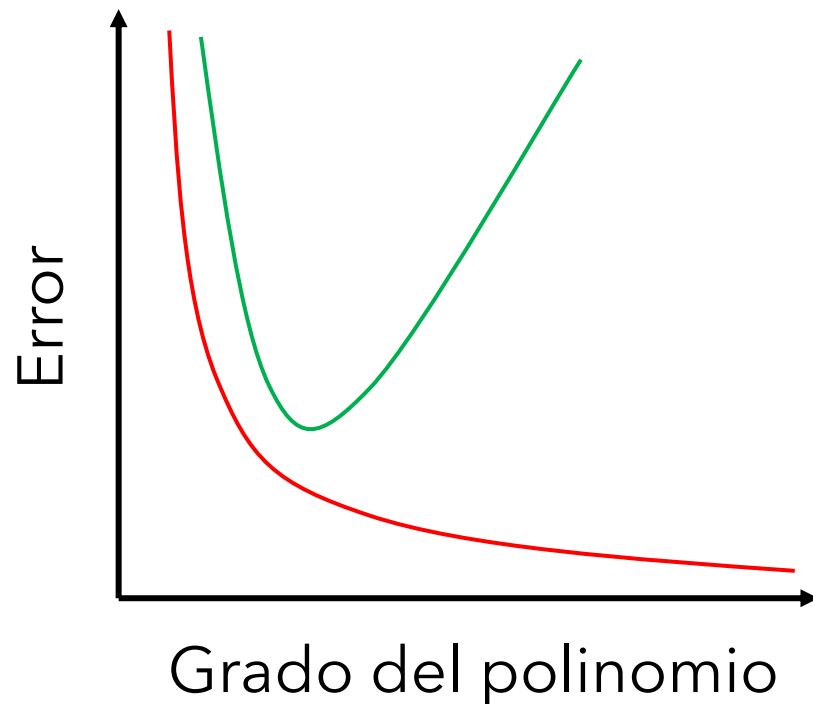
Error de entrenamiento: $J_{train}(\theta) = \frac{1}{2n_{train}} \sum_{i=1}^{n_{train}} \left(h_{\theta}(\mathbf{x}_{train}^{(i)}) - y_{train}^{(i)} \right)^2$

Error de validación: $J_{val}(\theta) = \frac{1}{2n_{val}} \sum_{i=1}^{n_{val}} \left(h_{\theta}(\mathbf{x}_{val}^{(i)}) - y_{val}^{(i)} \right)^2$



Bias - Variance

Considerando lo anterior, el error de prueba/validación $J(\boldsymbol{\theta})$ es alto, lo cual indica que algo anda mal. Pero, ¿qué es? ¿Sesgo o varianza?



Sesgo (*underfit*)

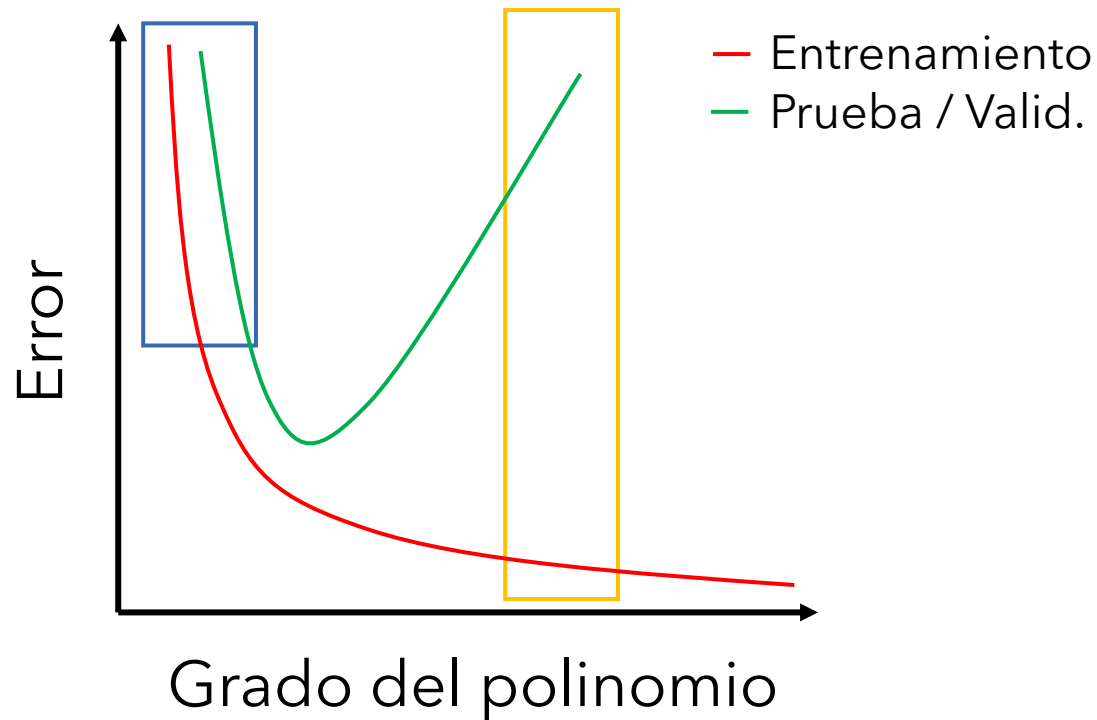
$$J_{train}(\boldsymbol{\theta}) \text{ alto}$$
$$J_{test}(\boldsymbol{\theta}) \approx J_{train}(\boldsymbol{\theta})$$

Varianza (*overfit*)

$$J_{train}(\boldsymbol{\theta}) \text{ bajo}$$
$$J_{test}(\boldsymbol{\theta}) \gg J_{train}(\boldsymbol{\theta})$$

Bias - Variance

Considerando lo anterior, el error de prueba/validación $J(\boldsymbol{\theta})$ es alto, lo cual indica que algo anda mal. Pero, ¿qué es? ¿Sesgo o varianza?



Sesgo (*underfit*)

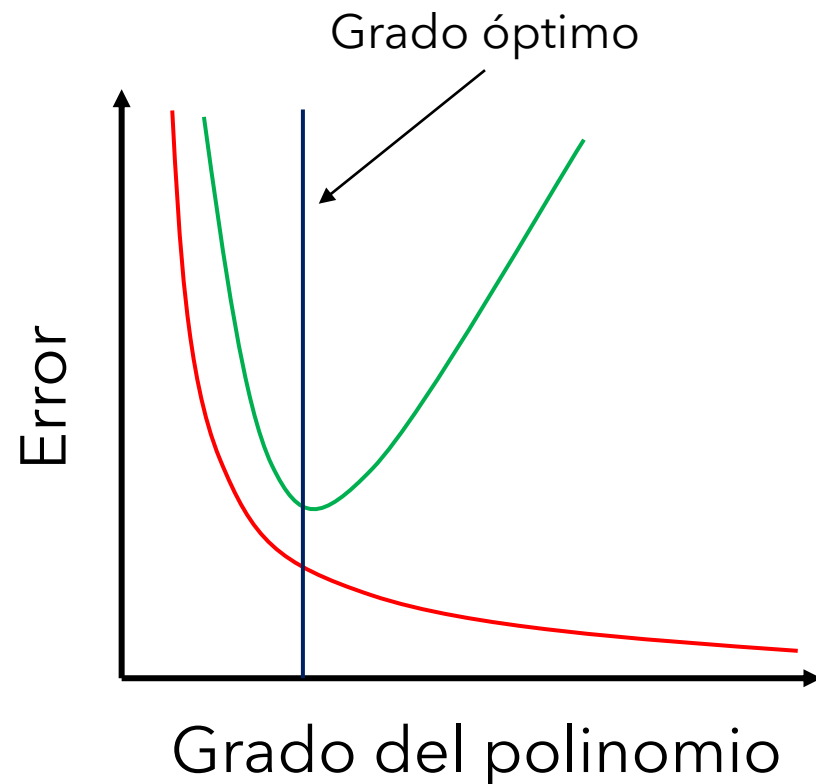
$$J_{train}(\boldsymbol{\theta}) \text{ alto}$$
$$J_{test}(\boldsymbol{\theta}) \approx J_{train}(\boldsymbol{\theta})$$

Varianza (*overfit*)

$$J_{train}(\boldsymbol{\theta}) \text{ bajo}$$
$$J_{test}(\boldsymbol{\theta}) \gg J_{train}(\boldsymbol{\theta})$$

Bias - Variance

Considerando lo anterior, el error de prueba/validación $J(\boldsymbol{\theta})$ es alto, lo cual indica que algo anda mal. Pero, ¿qué es? ¿Sesgo o varianza?



Sesgo (*underfit*)

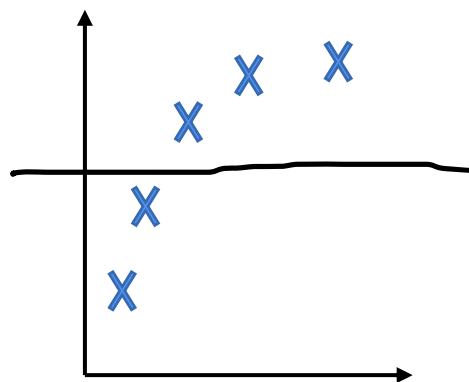
$$J_{train}(\boldsymbol{\theta}) \text{ alto}$$
$$J_{test}(\boldsymbol{\theta}) \approx J_{train}(\boldsymbol{\theta})$$

Varianza (*overfit*)

$$J_{train}(\boldsymbol{\theta}) \text{ bajo}$$
$$J_{test}(\boldsymbol{\theta}) \gg J_{train}(\boldsymbol{\theta})$$

Bias - Variance

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \left[\sum_{i=1}^n (h_{\boldsymbol{\theta}}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^p \theta_j^2 \right] \quad h_{\boldsymbol{\theta}}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

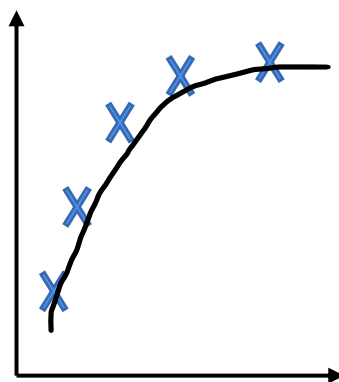


λ grande

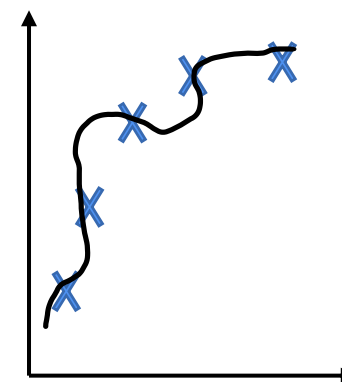
Alto sesgo (Underfitting)

$\lambda = 10000 \rightarrow \theta_1 \approx 0, \theta_2 \approx 0, \dots$

$h_{\boldsymbol{\theta}}(x) \approx \theta_0$



λ intermedia



λ pequena

Alto varianza (Overfitting)

$\lambda = .0001 \rightarrow \theta_1 \approx \theta_1, \theta_2 \approx \theta_2, \dots$

Bias – Variance y Regularización

¡Se determinan los
parámetros con
regularización!

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$J(\theta) = \frac{1}{2n} \left[\sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^p \theta_j^2 \right]$$

$$J_{train}(\theta) = \frac{1}{2n_{train}} \sum_{i=1}^{n_{train}} (h_{\theta}(x_{train}^{(i)}) - y_{train}^{(i)})^2 \quad J_{test}(\theta) = \frac{1}{2n_{test}} \sum_{i=1}^{n_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

$$J_{val}(\theta) = \frac{1}{2n_{val}} \sum_{i=1}^{n_{val}} (h_{\theta}(x_{val}^{(i)}) - y_{val}^{(i)})^2$$

¡Se evalúa el error
sin regularizar!

Bias – Variance y Regularización

Modelo:

$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$J(\theta) = \frac{1}{2n} \left[\sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^p \theta_j^2 \right]$$

1. Probar con $\lambda = 0$
2. Probar con $\lambda = 0.01$
3. Probar con $\lambda = 0.02$
4. Probar con $\lambda = 0.04$
5. Probar con $\lambda = 0.08$
6. Probar con $\lambda = 0.16$

⋮

$$\min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{val}(\theta^{(1)})$$

$$\min_{\theta} J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{val}(\theta^{(2)})$$

⋮

$$\min_{\theta} J(\theta) \rightarrow \theta^{(6)} \rightarrow J_{val}(\theta^{(6)})$$

⋮

Entrenar el modelo con el conjunto de entrenamiento CON regularización y obtener nuestra hipótesis.

Evaluamos con el conjunto de validación/prueba con las definiciones sin regularización.

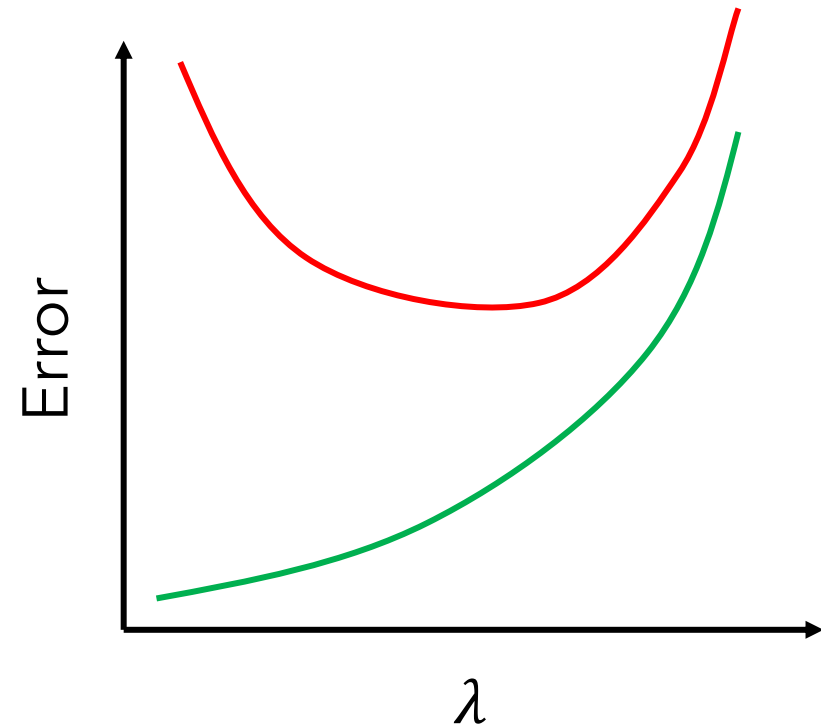
Se elige el que tenga menor error. Supongamos es θ_3 .

Bias – Variance y Regularización

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \left[\sum_{i=1}^n (h_{\boldsymbol{\theta}}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^p \theta_j^2 \right]$$

$$J_{train}(\theta) = \frac{1}{2n_{train}} \sum_{i=1}^{n_{train}} (h_{\theta}(x_{train}^{(i)}) - y_{train}^{(i)})^2$$

$$J_{val}(\theta) = \frac{1}{2n_{val}} \sum_{i=1}^{n_{val}} (h_{\theta}(x_{val}^{(i)}) - y_{val}^{(i)})^2$$

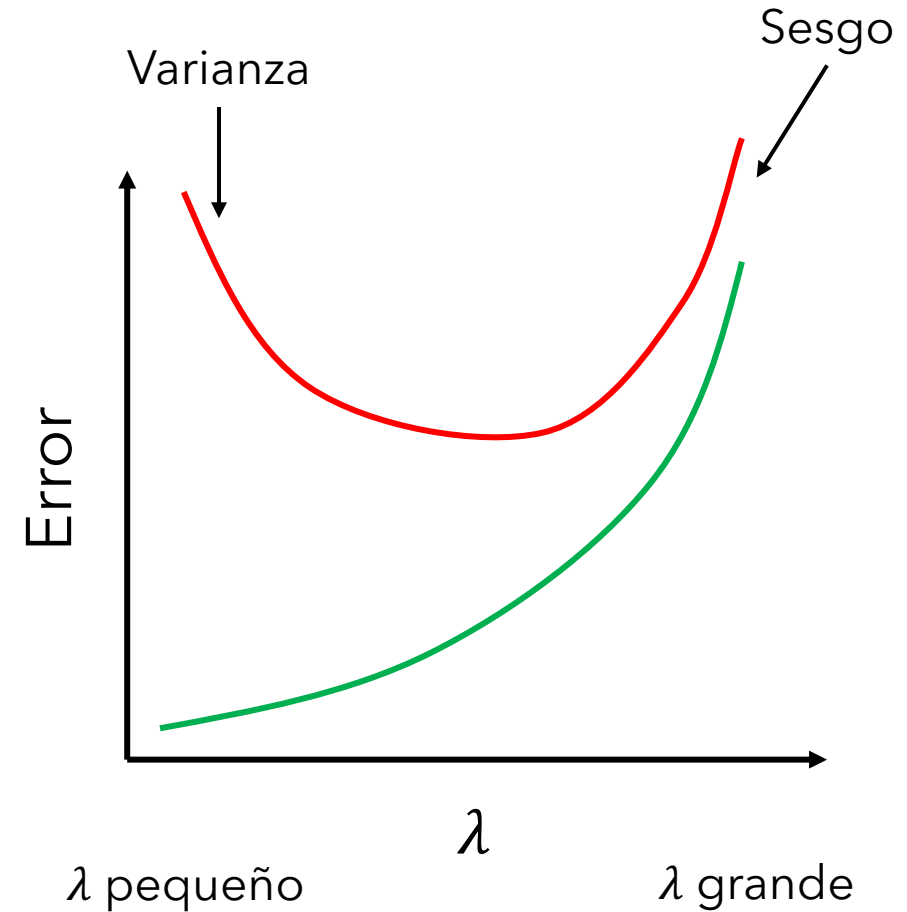


Bias – Variance y Regularización

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \left[\sum_{i=1}^n (h_{\boldsymbol{\theta}}(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^p \theta_j^2 \right]$$

$$J_{train}(\theta) = \frac{1}{2n_{train}} \sum_{i=1}^{n_{train}} (h_{\theta}(x_{train}^{(i)}) - y_{train}^{(i)})^2$$

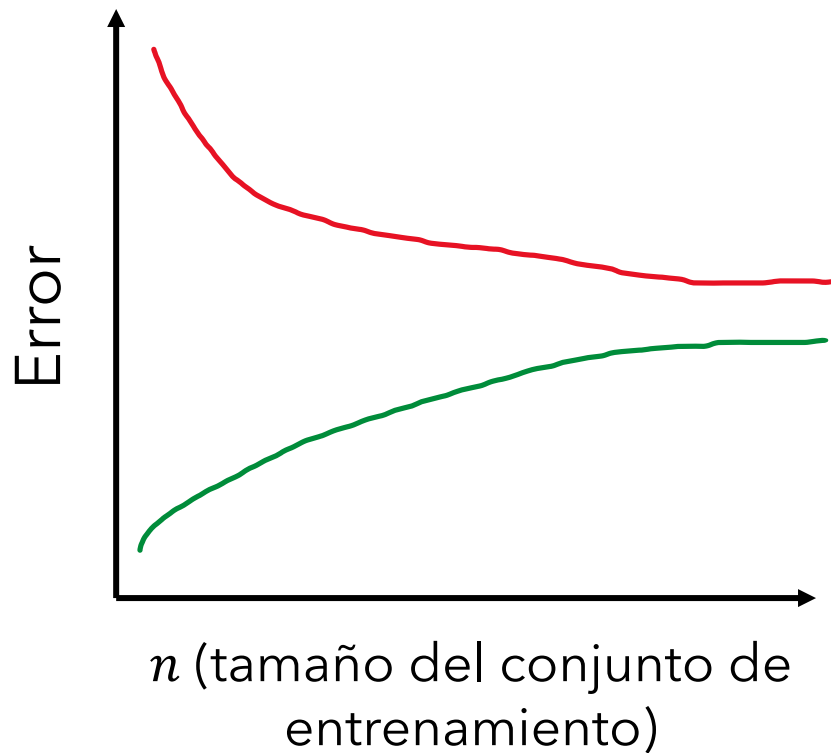
$$J_{val}(\theta) = \frac{1}{2n_{val}} \sum_{i=1}^{n_{val}} (h_{\theta}(x_{val}^{(i)}) - y_{val}^{(i)})^2$$





Curvas de Aprendizaje

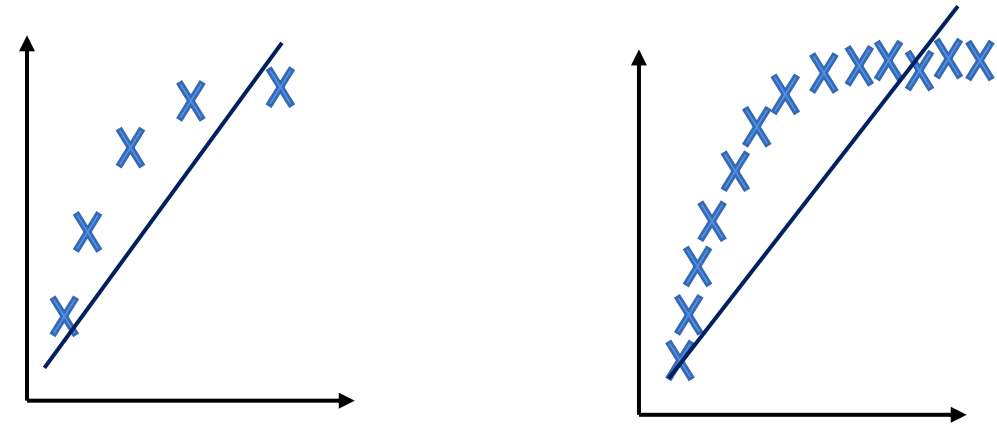
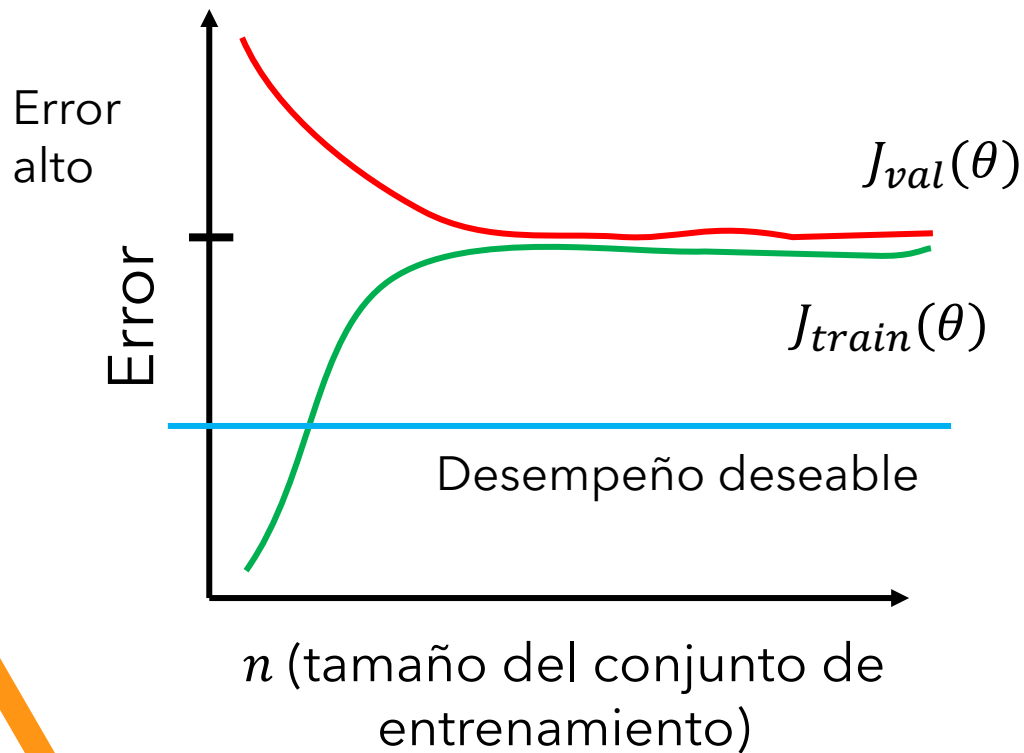
Curvas de Aprendizaje



$$J_{train}(\theta) = \frac{1}{2n_{train}} \sum_{i=1}^{n_{train}} \left(h_{\theta} \left(x_{train}^{(i)} \right) - y_{train}^{(i)} \right)^2$$

$$J_{val}(\theta) = \frac{1}{2n_{val}} \sum_{i=1}^{n_{val}} \left(h_{\theta} \left(x_{val}^{(i)} \right) - y_{val}^{(i)} \right)^2$$

Curvas de Aprendizaje – Sesgo Alto

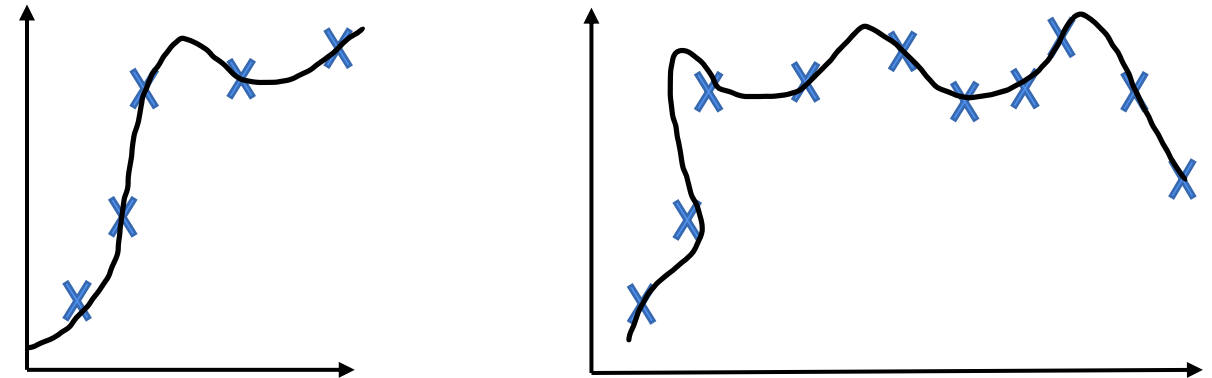
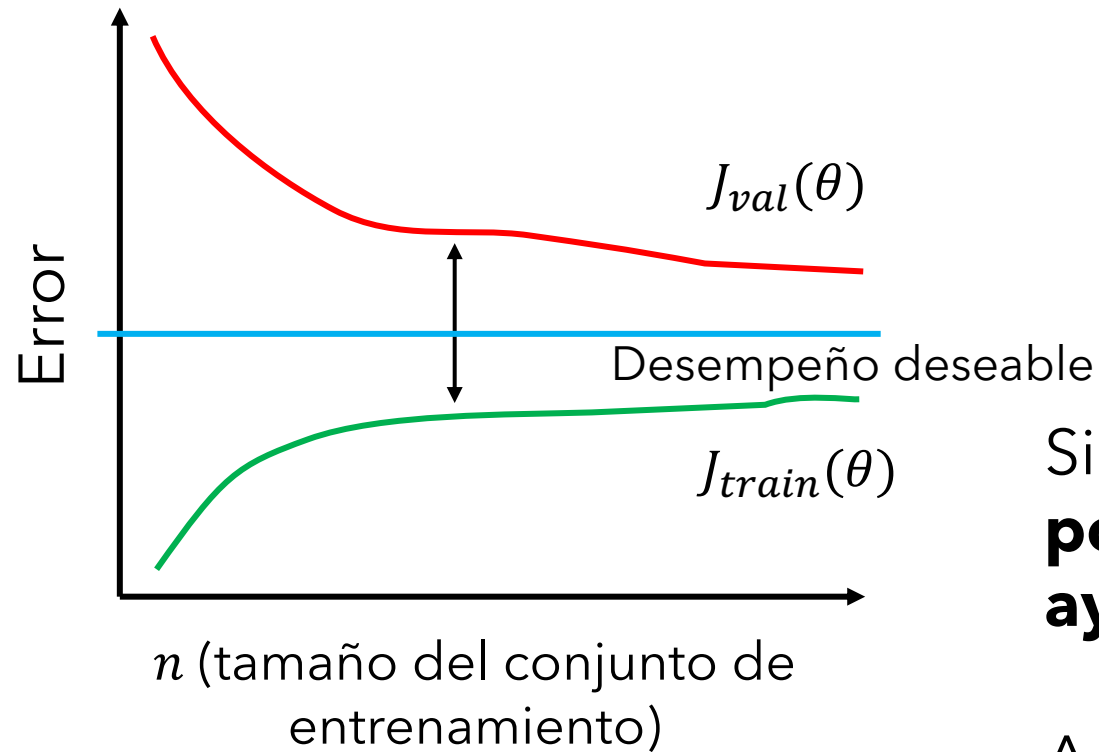


$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Si un algoritmo de aprendizaje sufre de alto sesgo, **obtener más información no ayuda a mejorar su rendimiento.**

Al ser un modelo muy simple, no logrará capturar la complejidad de la información, sin importar cuanta sea.

Curvas de Aprendizaje – Varianza Alta



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \dots + \theta_{100} x^{100}$$

Si un algoritmo tiene alta varianza, **es posible que añadir más información ayude a su desempeño.**

A largo plazo, ambas curvas pueden converger a una asíntota.




Juntemos todo...


¿Qué modelo elegir?

Considerando las opciones iniciales, tenemos que:

- Tener más datos para el entrenamiento. → *resuelve* varianza alta
- Considerar menos características. → *resuelve* varianza alta
- Obtener más características. → *resuelve* sesgo alto
- Considerar combinaciones de características más complejas. → *resuelve* sesgo alto
- Aumentar el valor de λ . → *resuelve* varianza alto
- Reducir el valor de λ → *resuelve* sesgo alto




Conclusiones y Notas Adicionales



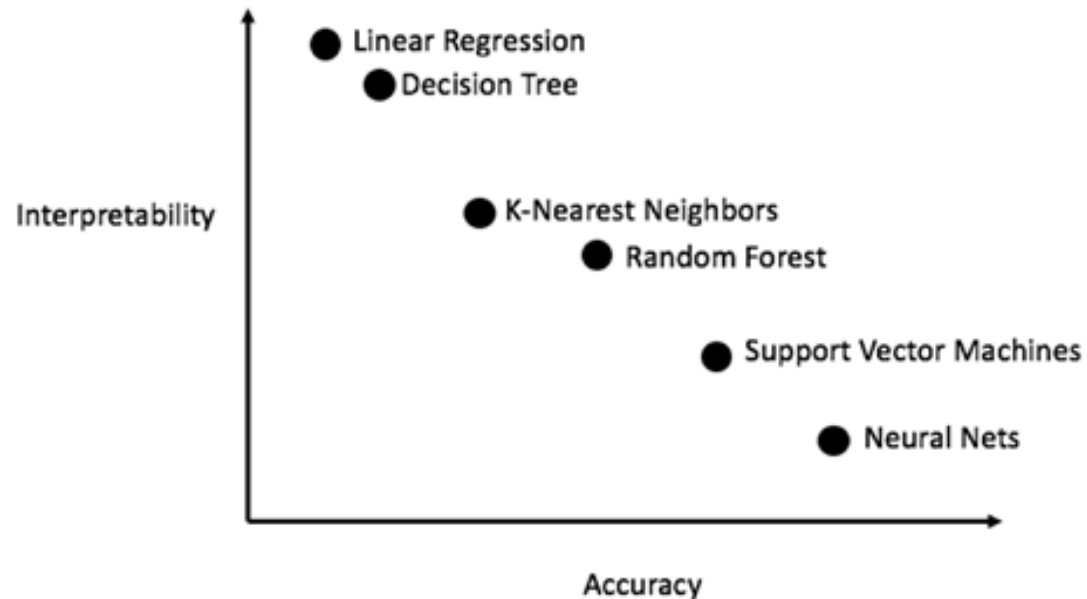


Conclusiones y Notas Adicionales

- Debemos elegir un modelo con un error aceptable para que sea capaz de generalizar.
 - En el caso de regresión, esto implica elegir un polinomio de grado medio.
 - Muy bajo nos lleva al underfitting.
 - Muy alto nos lleva al overfitting.
 - En el caso de clasificación, se debe tener cuidado con las funciones que activan umbrales y los métodos de clasificación.
- 

Conclusiones y Notas Adicionales

- En todos los casos anteriores, alto sesgo implica baja varianza y viceversa.
- También se debe considerar el aspecto de la interpretabilidad de los modelos.





¡Gracias!

Luis Zúñiga

p40887@correo.uia.mx

Website