

Tema 2 – Optimización y Regularización (Parte 4)

Aprendizaje Automático II - Grado en Inteligencia Artificial
Universidad Rey Juan Carlos

Iván Ramírez Díaz
ivan.ramirez@urjc.es

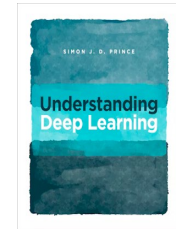
José Miguel Buenaposada Biencinto
josemiguel.buenaposada@urjc.es

2.5 Estimación del rendimiento

- Ruido, sesgo y varianza
- Reducir la varianza
- Reducir el sesgo y el relación sesgo-varianza
- Doble descenso
- Elegir los hiperparámetros

Bibliografía

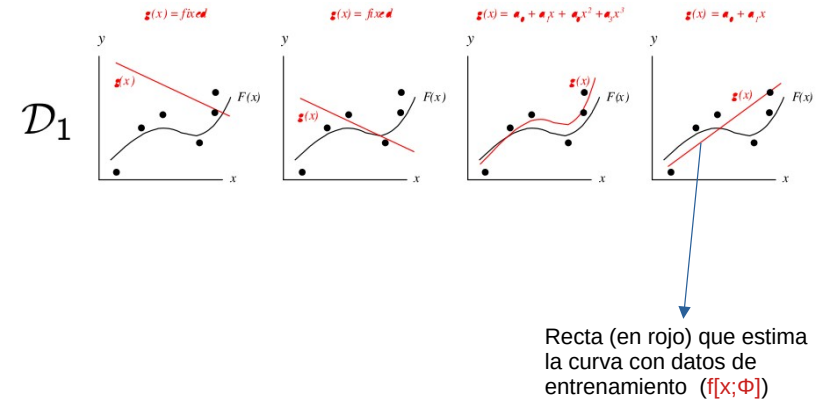
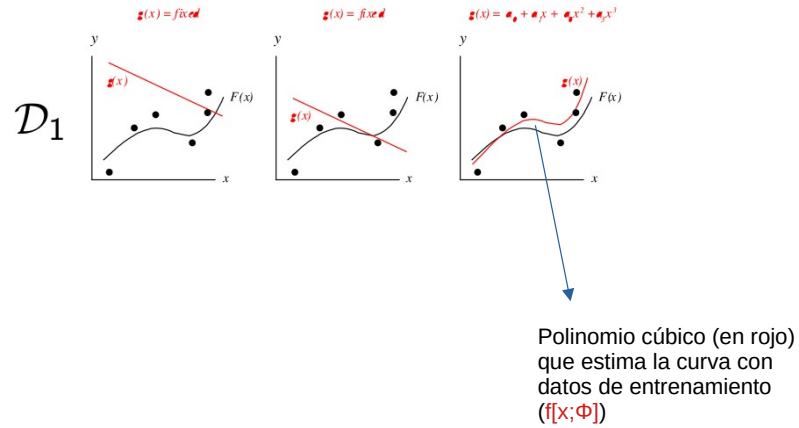
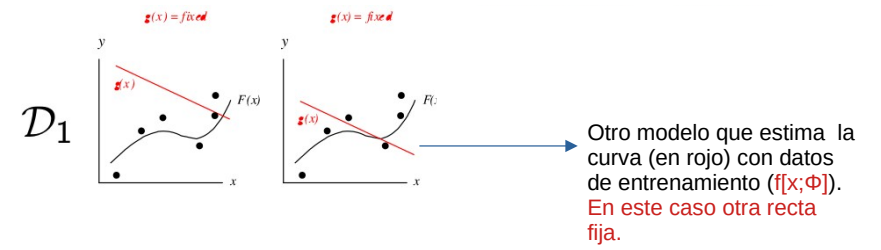
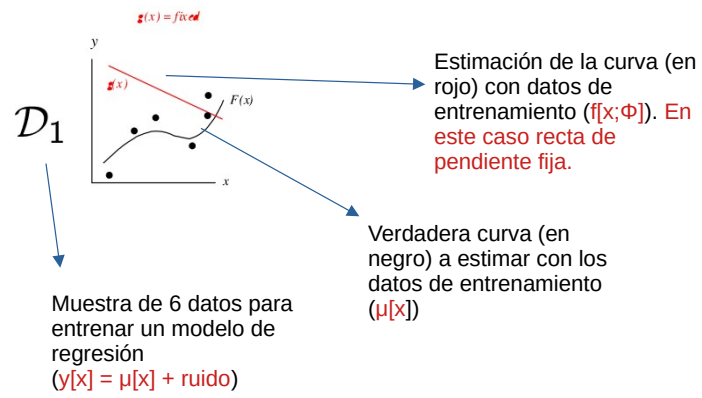
- **Understanding Deep Learning**. Capítulo 8.

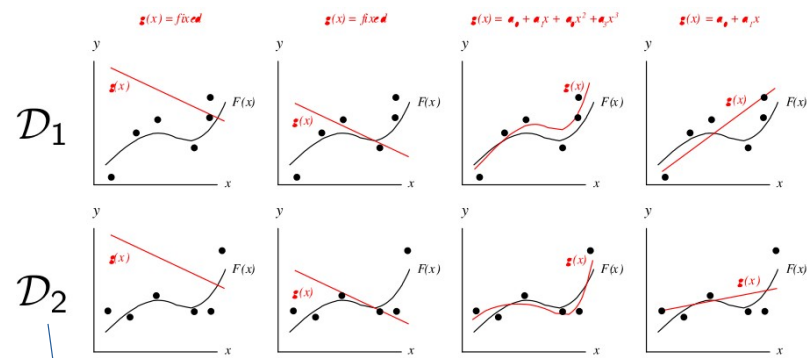


- **Deep Learning: CS 182 2021**. Lecture 3, Part 3.
Sergey Levine. UC Berkeley.
Curso en youtube.

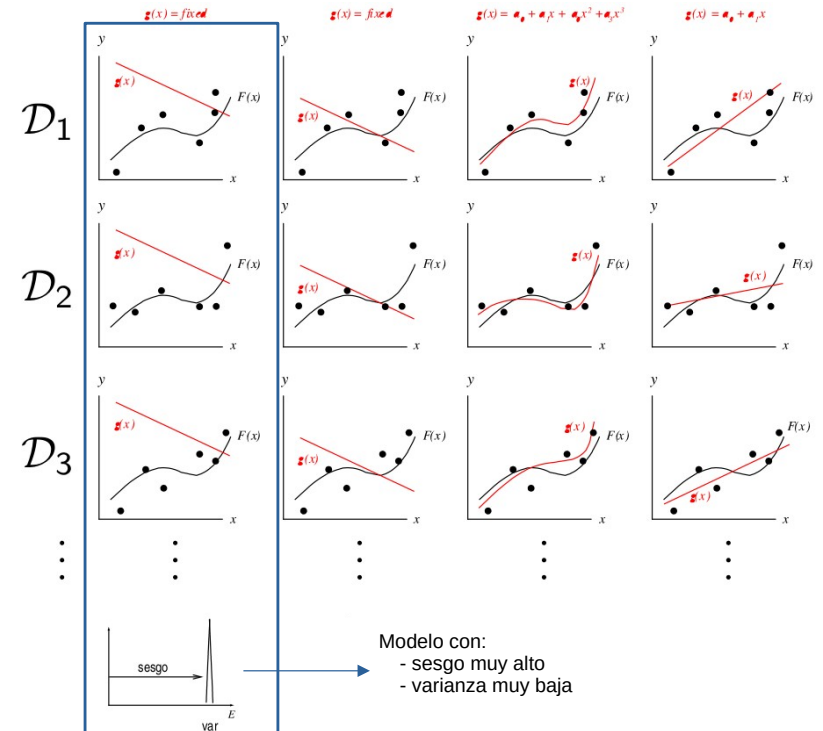
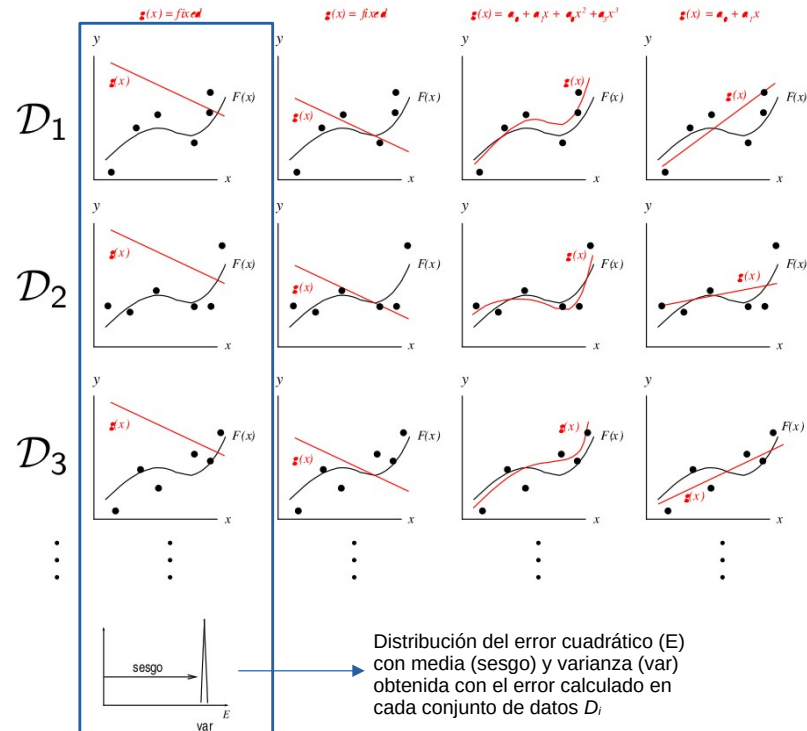
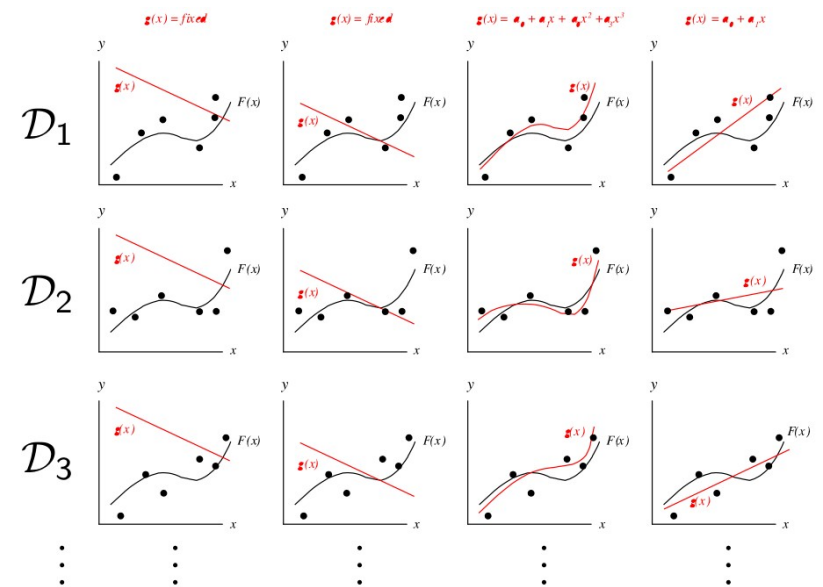
2.5 Estimación del rendimiento

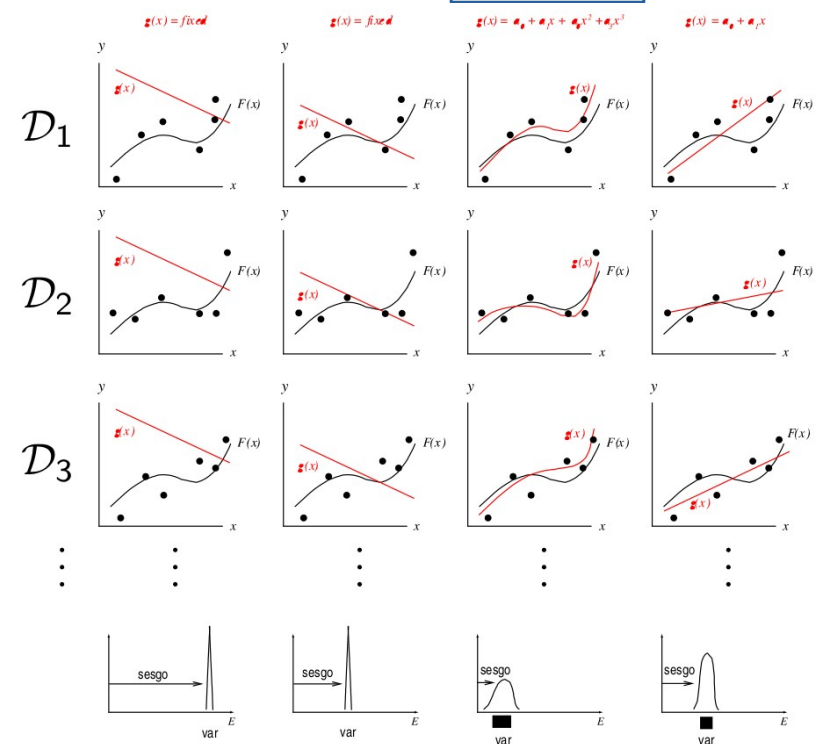
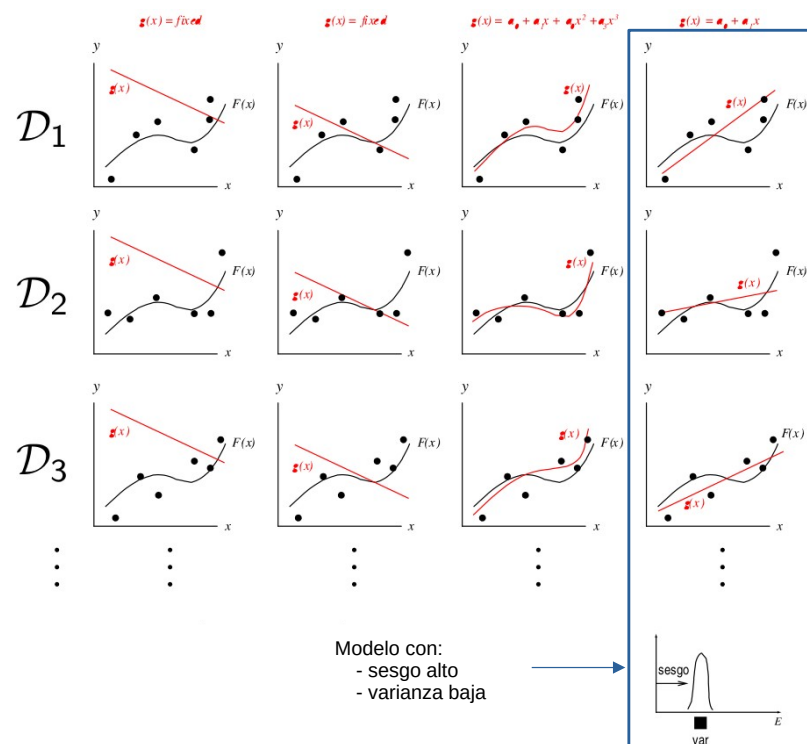
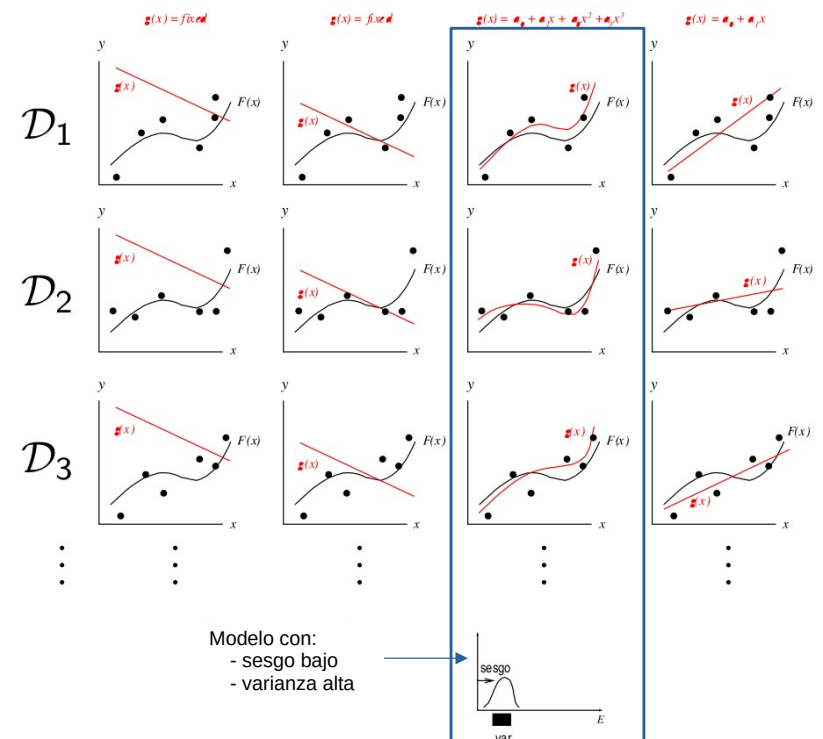
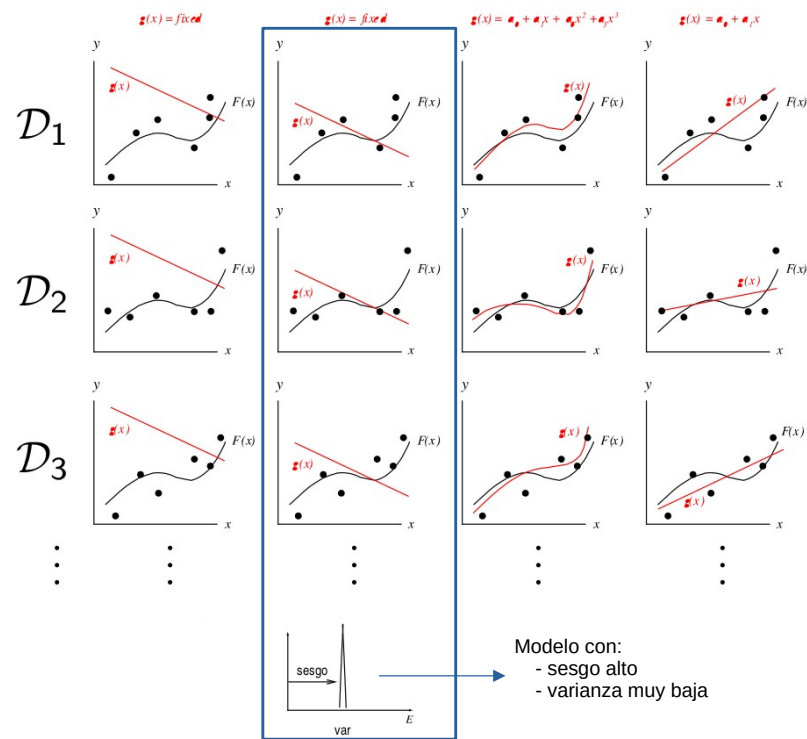
- **Ruido, sesgo y varianza**
- Reducir la varianza
- Reducir el sesgo y el relación sesgo-varianza
- Doble descenso
- Elegir los hiperparámetros



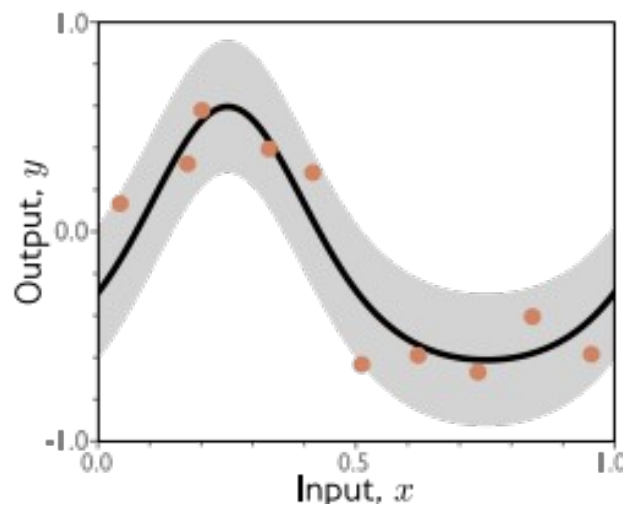


Otra muestra diferente
de 6 datos para entrenar
un modelo de regresión



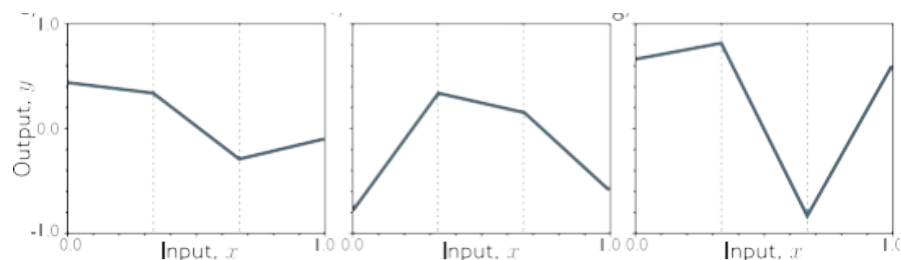


Ejemplo de regresión

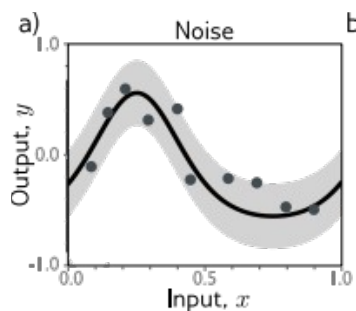


Modelo básico para la estimación

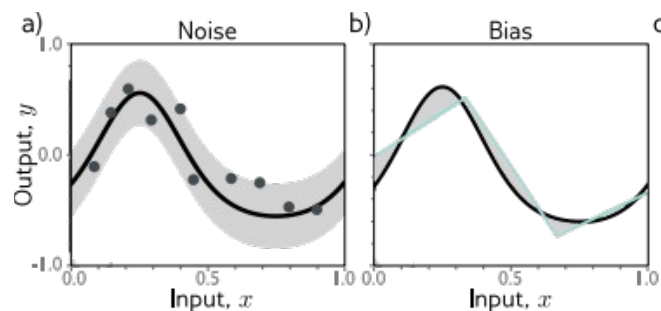
- Modelo que estima tres modelos lineales en el intervalo $[0, 1]$ en tramos de igual longitud.



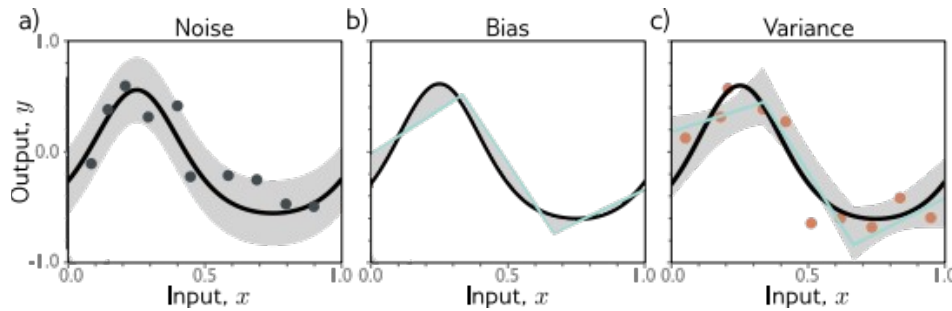
Ruido, sesgo (bias) y varianza (variance)



- Posibles causas:
 - Ruido en las medidas
 - Algunas variables del modelo son ocultas
 - Datos mal etiquetados



Ruido, sesgo (bias) y **varianza** (variance)



Ruido, sesgo (bias) y **varianza** (variance)

- La **varianza** es la incertidumbre en el modelo entrenado debido a la elección del conjunto de datos.
- El **sesgo** es la desviación sistemática de la media de la verdadera función que estamos estimando debida a las limitaciones de nuestro modelo.
- El **ruido** es la incertidumbre inherente en la función verdadera que lleva un dato x a una salida y .

Para un problema de mínimos cuadrados ...

$$L[x] = (f[x, \phi] - y[x])^2$$

Regresor con parámetros ϕ que estima y dado un x .

Verdadera función que genera los datos de entrenamiento con un cierto ruido de dist. Normal ($y[x] = \mu[x] + \text{ruido}$)

Para un problema de mínimos cuadrados ...

$$L[x] = (f[x, \phi] - y[x])^2$$

Se puede demostrar que:

$$\mathbb{E}_{\mathcal{D}} [\mathbb{E}_y [L[x]]] = \underbrace{\mathbb{E}_{\mathcal{D}} [(f[x, \phi[\mathcal{D}]] - f_{\mu}[x])^2]}_{\text{variance}} + \underbrace{(f_{\mu}[x] - \mu[x])^2}_{\text{bias}} + \underbrace{\sigma^2}_{\text{noise}}$$

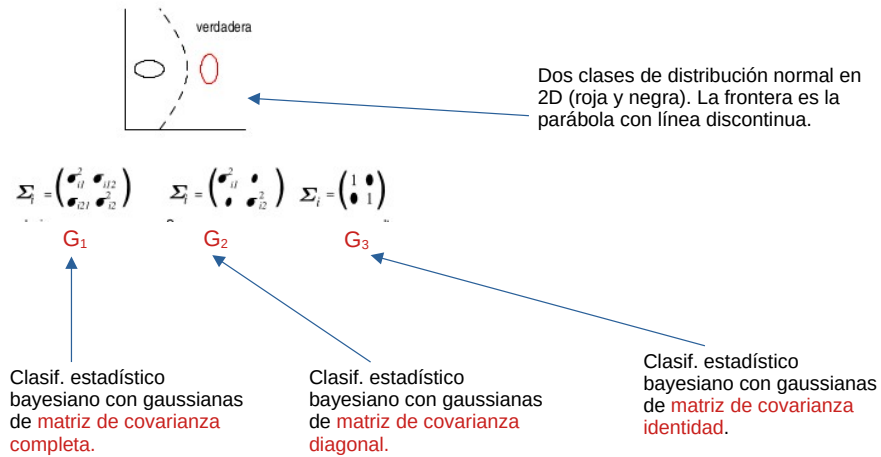
Esperanza sobre las muestras de datos

Esperanza sobre el ruido en el conjunto de test

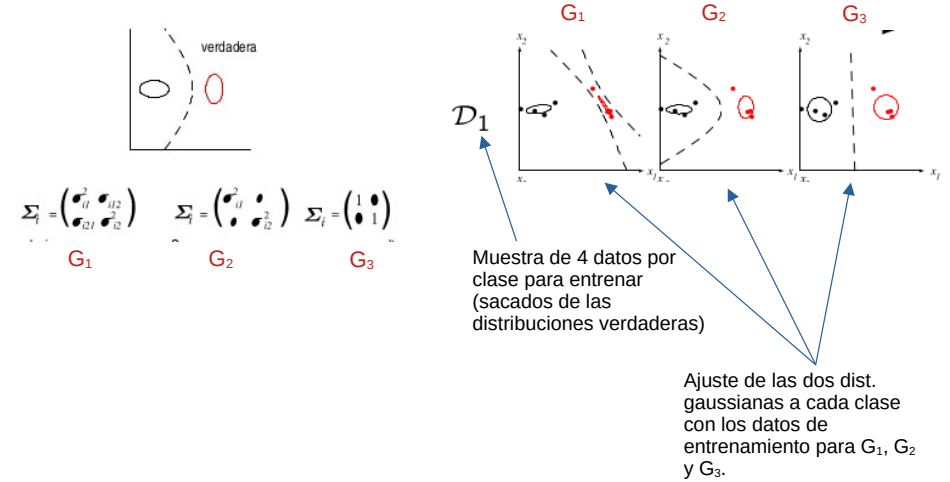
Media de las estimaciones de f para infinitos conjuntos de datos de entrenamiento \mathcal{D}

Verdadera función sin ruido

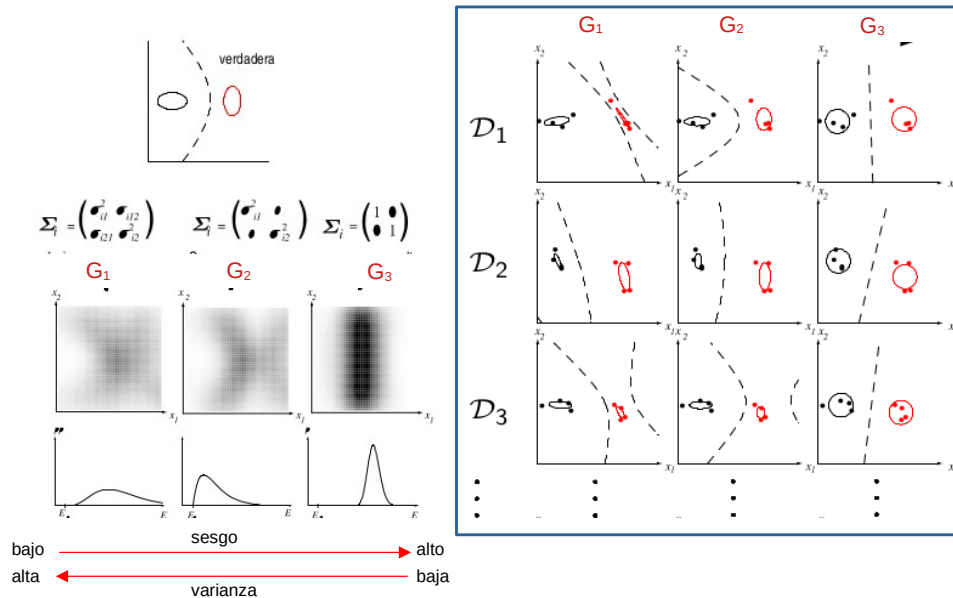
Para un problema de clasificación ...



Para un problema de clasificación ...



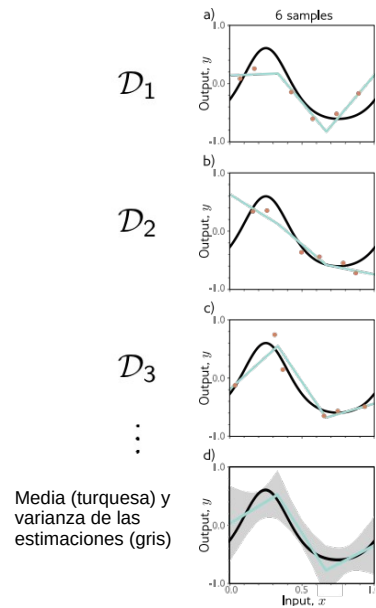
Para un problema de clasificación ...



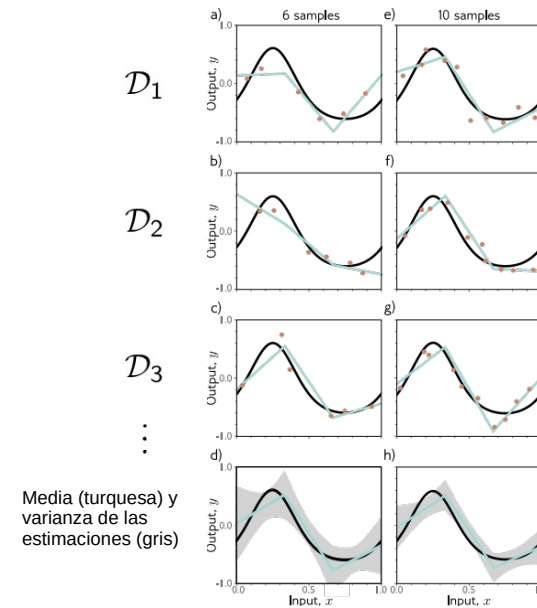
2.5 Estimación del rendimiento

- Ruido, sesgo y varianza
- **Reducir la varianza**
- Reducir el sesgo y el compromiso sesgo-varianza
- Doble descenso
- Elegir los hiperparámetros

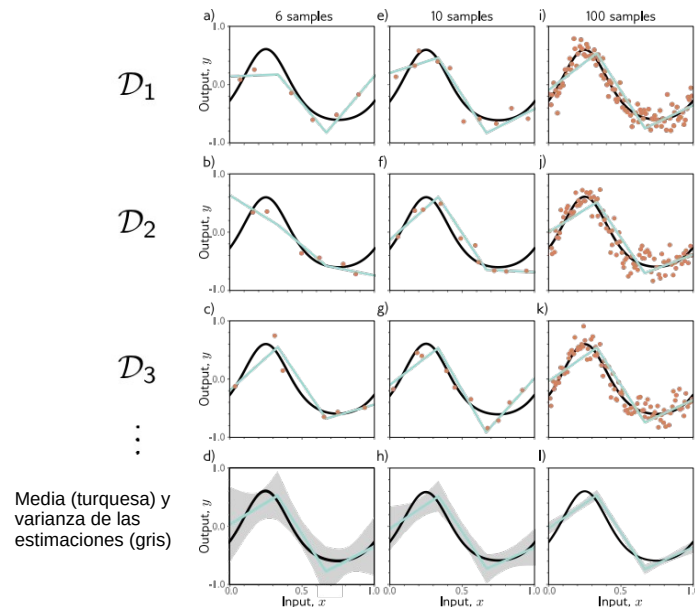
Reducir la varianza



Reducir la varianza



Reducir la varianza

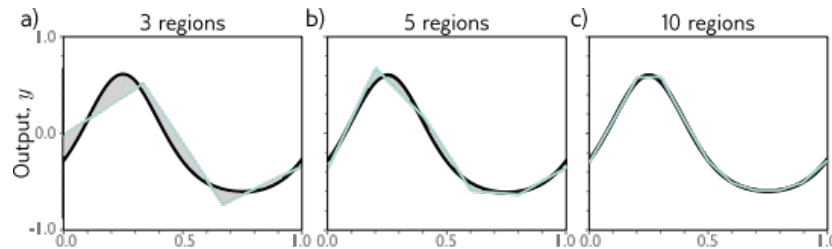


Aumentando el n.º de
datos de entrenamiento
reducimos la varianza.

2.5 Estimación del rendimiento

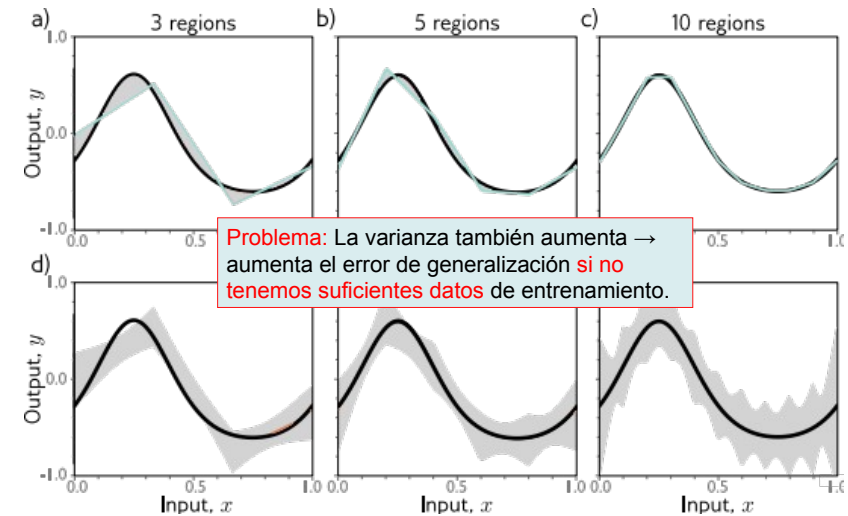
- Ruido, sesgo y varianza
- Reducir la varianza
- Reducir el sesgo y la relación sesgo-varianza
- Doble descenso
- Elegir los hiperparámetros

Reducir el sesgo

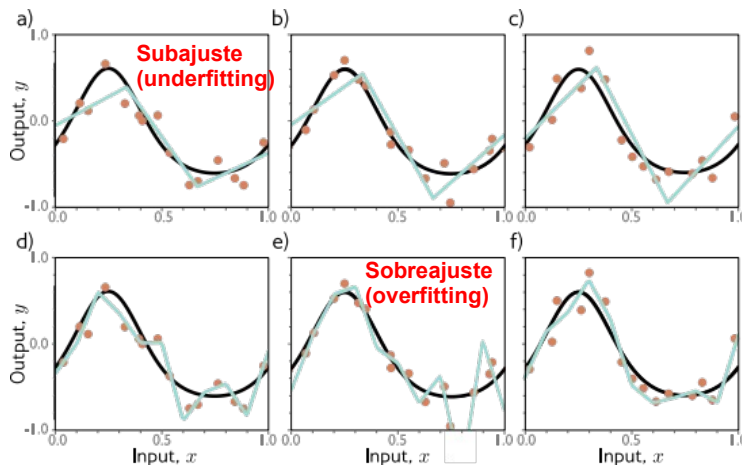


Aumentamos el número de parámetros del estimador
(= flexibilidad para ajustarse a los datos) →
 $f_h[x]$ se acerca más a $\mu[x]$

Reducir el sesgo

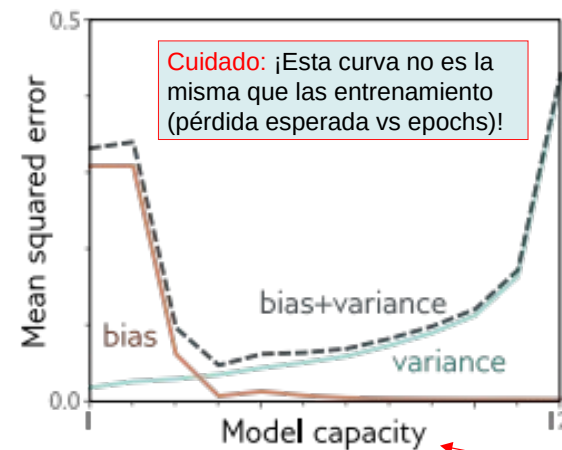


¿Por qué aumenta la varianza?



$f[x, \Phi]$ (turquesa) describe mejor los datos de
entrenamiento, pero no la verdadera curva
(en negro), $\mu[x] \rightarrow$ **Sobreajuste (overfitting)**

Relación sesgo-varianza (bias-variance trade-off)



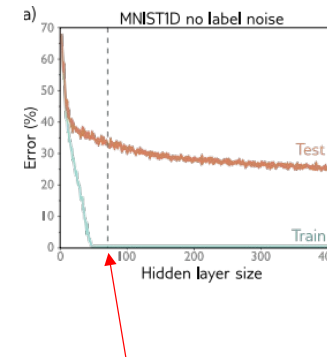
Relacionado con:

- número de datos de entrenamiento
que el modelo puede ajustar sin error
- número de parámetros del modelo

2.5 Estimación del rendimiento

- Ruido, sesgo y varianza
- Reducir la varianza
- Reducir el sesgo y el relación sesgo-varianza
- **Doble descenso**
- Elegir los hiperparámetros

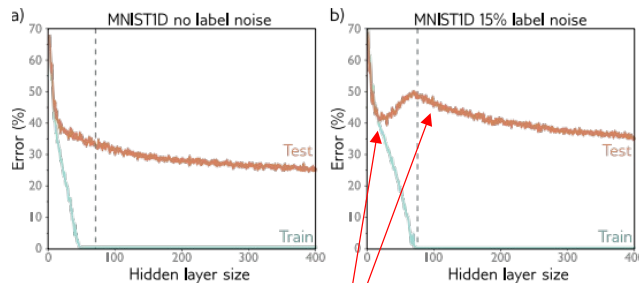
Doble descenso



Cuidado: ¡Esta curva no es la misma que las entrenamiento (pérdida esperada vs epochs)!

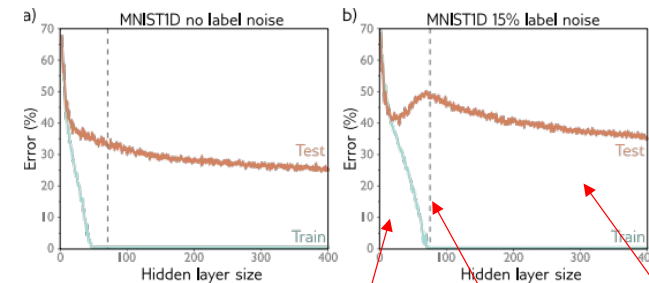
Nº parámetros del modelo = Nº datos entrenamiento

Doble descenso



Aparece el **Doble Descenso**

Doble descenso

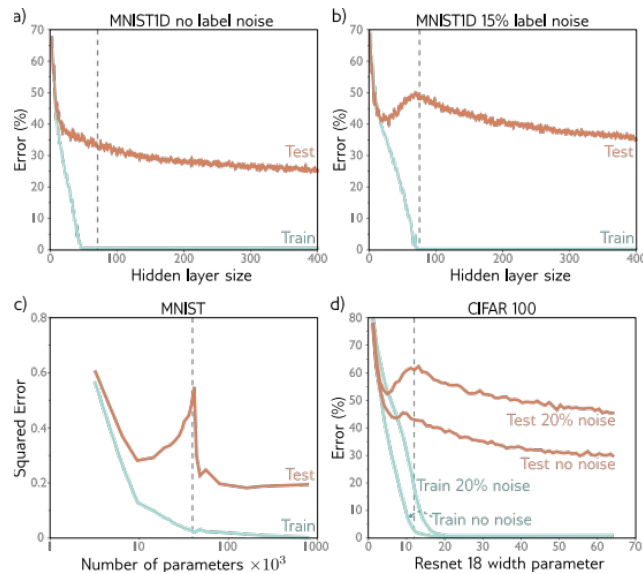


Régimen subparametrizado
(o clásico)

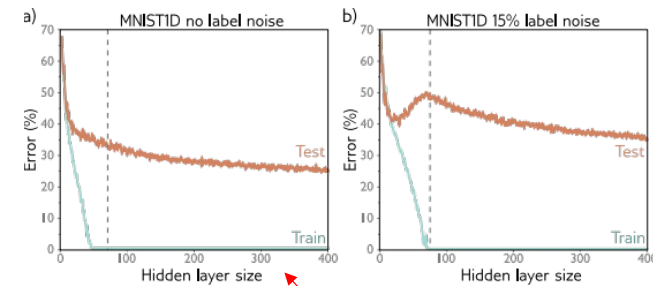
Régimen sobreparametrizado
(o moderno)

Régimen crítico

Doble descenso

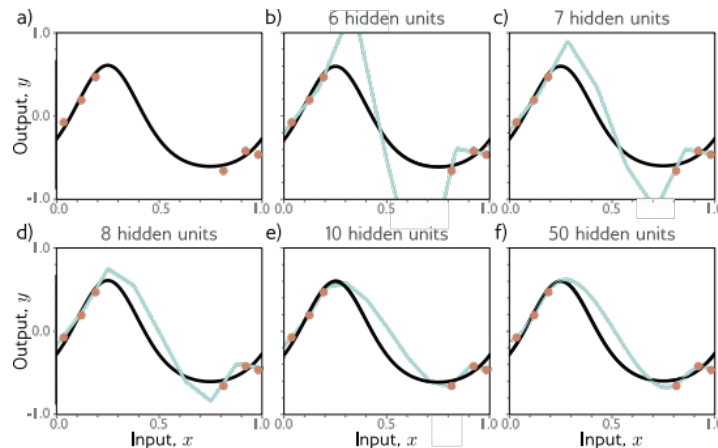


Doble descenso



- El error de entrenamiento es cercano a 0
- Lo que esté ocurriendo no está afectando a los datos de entrenamiento
- ¿Debería de ocurrir fuera de los puntos usados en entrenamiento?

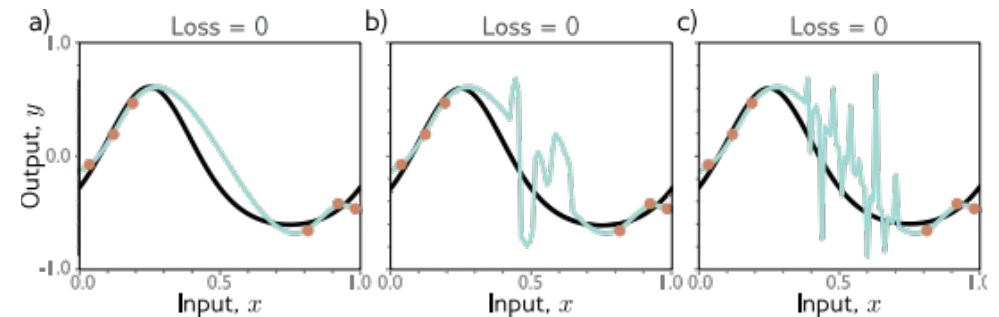
En puntos distintos a los de entrenamiento ...



Explicación potencial:

- Con más parámetros estimamos funciones más suaves
- Suavidad fuera de los datos de entrenamiento es algo razonable

¿Explicación del doble descenso?



Todas estas soluciones equivalentes en términos de pérdida esperada.

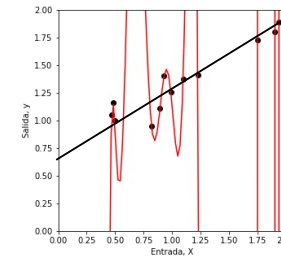
- ¿Por qué el modelo debería elegir la función más suave en a)?
- Explicaciones potenciales:
 - La inicialización elige funciones suaves y la optimización no se sale de ellas
 - El algoritmo de entrenamiento "prefiere" converger a funciones suaves.

2.5 Estimación del rendimiento

- Ruido, sesgo y varianza
- Reducir la varianza
- Reducir el sesgo y el compromiso sesgo-varianza
- Doble descenso
- Elegir los hiperparámetros

Algunas preguntas ...

- ¿Cómo **saber** si estamos sobreajustando o subajustando?
- ¿Cómo elegir el algoritmo de optimización / modelo?
- ¿Cómo elegir los **hiperparámetros**?
- **Idea:** elegir lo que haga la **pérdida esperada** muy baja



$$J(\phi) = 0$$

¡No se puede detectar **sobreajuste** con la pérdida esperada en el entrenamiento!

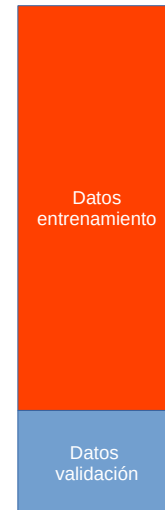
Método de trabajo en Aprendizaje Automático

Datos disponibles



Método de trabajo en Aprendizaje Automático

Datos disponibles



Usar para entrenar

$$J(\phi, D_{train})$$

Reservar para:

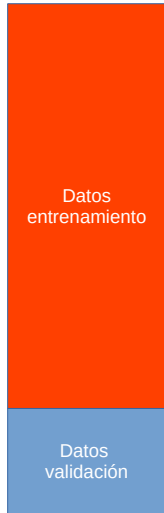
- elegir hiperparámetros
- añadir/eliminar características
- elegir el tipo de modelo

1. Entrenar ϕ con $J(\phi, D_{train})$
Si $J(\phi, D_{train})$ no pequeña
Estás **subajustando**
bajar la regularización
mejorar el optimizador
2. Mirar $J(\phi, D_{val})$
Si $J(\phi, D_{val}) \gg J(\phi, D_{train})$
Estás **sobreajustando**
incrementar regularización

$$J(\phi, D_{val})$$

Método de trabajo adecuado

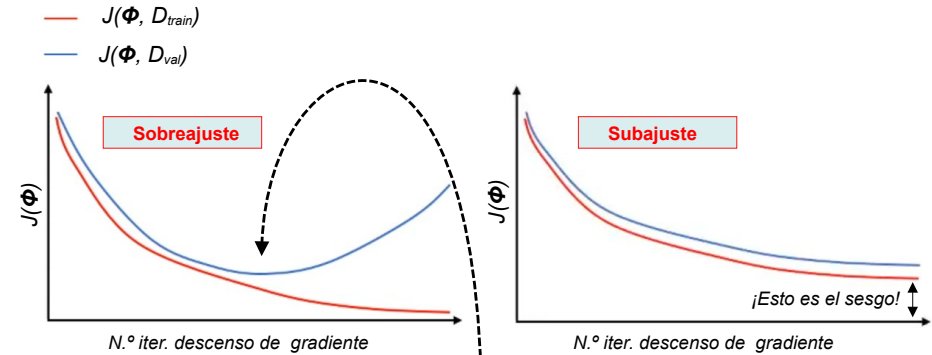
Datos disponibles



➔ Usar para seleccionar:
 Φ (vía optimización)
 Hiperparámetros optimización (p.ej. learning rate)

➔ Usar para seleccionar:
 Tipo de modelo (p.ej. Logistic regresion vs otro)
 Hiperparámetros regularización
 Selección características

Curvas de aprendizaje



Pregunta: ¿Podemos parar aquí?
 Early stopping

El conjunto de test final

Datos disponibles



1. Entrenar Φ con $J(\Phi, D_{\text{train}})$
 Si $J(\Phi, D_{\text{train}})$ no pequeña
 Estás **subajustando**
 bajar la regularización
 mejorar el optimizador
2. Mirar $J(\Phi, D_{\text{val}})$
 Si $J(\Phi, D_{\text{val}}) \gg J(\Phi, D_{\text{train}})$
 Estás **sobreajustando**
 incrementar regularización

¿Hemos terminado?
 ¿Nos quedamos con el clasificador que minimiza $J(\Phi, D_{\text{val}})$?
 No es buena idea - hemos usado el conjunto de validación para elegir hiperparámetros

Método de trabajo en Aprendizaje Automático

Datos disponibles



➔ Usar para seleccionar:
 Φ (vía optimización)
 Hiperparámetros optimización (p.ej. learning rate)

➔ Usar para seleccionar:
 Tipo de modelo (p.ej. Logistic regresion vs otro)
 Hiperparámetros regularización
 Selección características

➔ Usar **únicamente** para reportar rendimiento final

Resumen

- ¿De dónde vienen los errores?
 - **Varianza**: demasiada capacidad, información insuficiente para encontrar los parámetros adecuados
 - **Sesgo**: poca capacidad, no se puede representar la función correcta
 - $\text{Error} = \text{Variance} + \text{Bias}^2$ (en regresión)
 - **Sobreajuste (overfitting)**: demasiada varianza
 - **Subajuste (underfitting)**: demasiado sesgo

Resumen

- ¿Cómo equilibrar sesgo y varianza?
 - Seleccionar el tipo de modelo cuidadosamente
 - Seleccionar las características cuidadosamente
 - **Regularización**: añadida a la función de coste para reducir la varianza

Resumen

- ¿Cómo seleccionar hiperparámetros?
 - Separación entrenamiento/validación
 - Datos de entrenamiento para optimización (aprendizaje)
 - Datos de validación para seleccionar hiperparámetros
 - ¡Datos de test para **obtener el resultado final y nada más!**