

El compromiso entre bias y varianza

Veamos un problema de regresión ordinaria, donde tenemos N puntos

$$(x_i, y_i)$$

y suponemos que hay una relación causal entre ellos

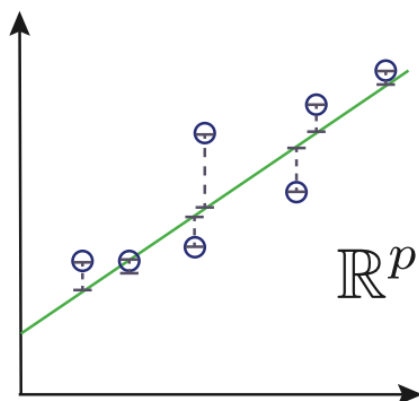
$$y_i = f(x_i)$$

Pero como y_i fue "medido" y no calculado, sabemos que la medición siempre involucra un error

$$y_i = f(x_i) + \epsilon_i$$

Donde ϵ_i es una variable aleatoria, independiente del error de otras mediciones, con distribución gaussiana de media 0 y desviación estándar $\sigma = 1$.

Minimizar el caso en que podamos aproximar la función $f(x)$, como pasa en las ciencias



7 puntos

FIG. 10 Geometric interpretation of least squares regression. The regression function g defines a hyperplane in \mathbb{R}^p (green solid line, here we have $p = 2$) while the residual of data point x_i (hollow circles) is its projection onto this hyperplane (bar-ended dashed line).

En este ejemplo sabemos que

$$f(x) = a + bx$$

O sea, tenemos 2 parámetros. Los elegimos de manera que minimicemos el error cuadrático medio

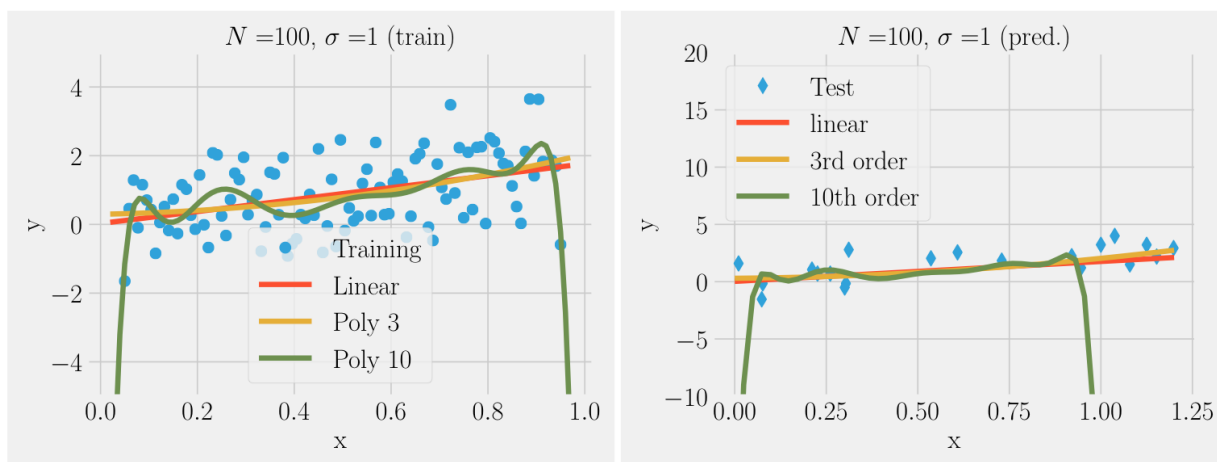
$$ECM = \frac{1}{7} \sum_{i=1}^7 (\hat{y}_i - y_i)^2$$

donde \hat{y}_i es el valor medido (con error) e y es el estimado (sin error).

Aquí conocemos f y solo necesitamos calcular los mejores parámetros

Para otros 7 puntos tendremos otros parámetros por culpa del error.

Cuando mayor sea el # de puntos, mayor será el ECM. Pero mejor interpolareé (estimareé puntos no medidos).

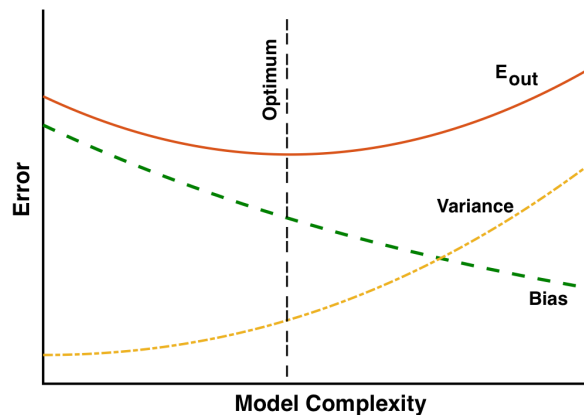
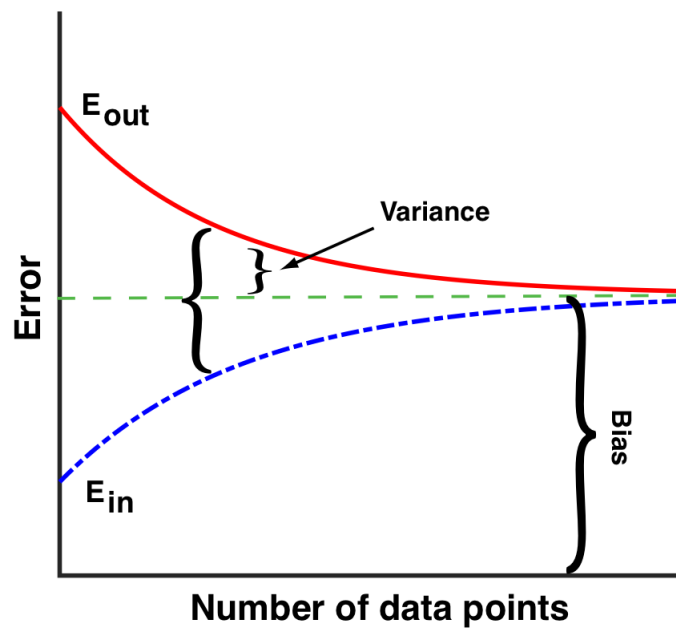


Mira el ejemplo con 100 puntos. los bolos azules son datos conocidos usados para ajustar. los puntos (a la derecha) son datos para verificar (test).

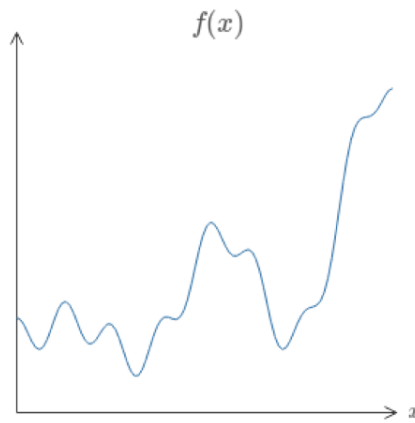
- * El polinomio de grado 10 requiere ajustar 11 parámetros
- * El " " " 3 " " 4 "
- * El " " " 1 " " 2 "

El polinomio de grado 10 ajuste mejor los datos conocidos, pero extrapola mal.

El polinomio de grado 2 ajuste peor los datos conocidos, pero extrapola mejor.

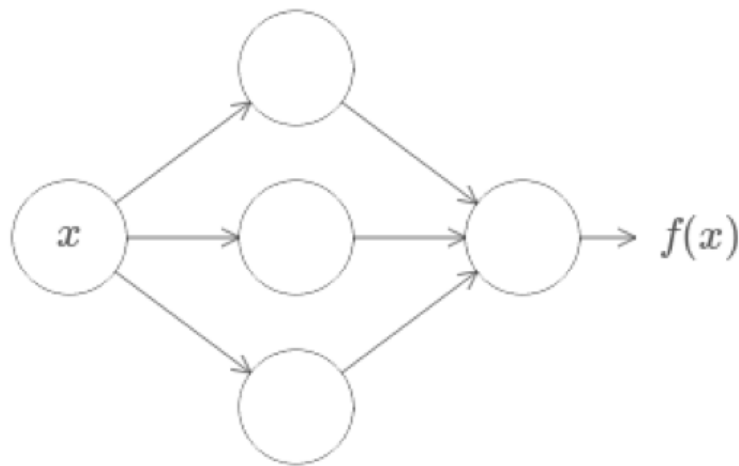


Ahora mismo Machine Learning. Es el mismo problema pero no conocemos $f(x)$, solo los datos.



función
definida
a mano

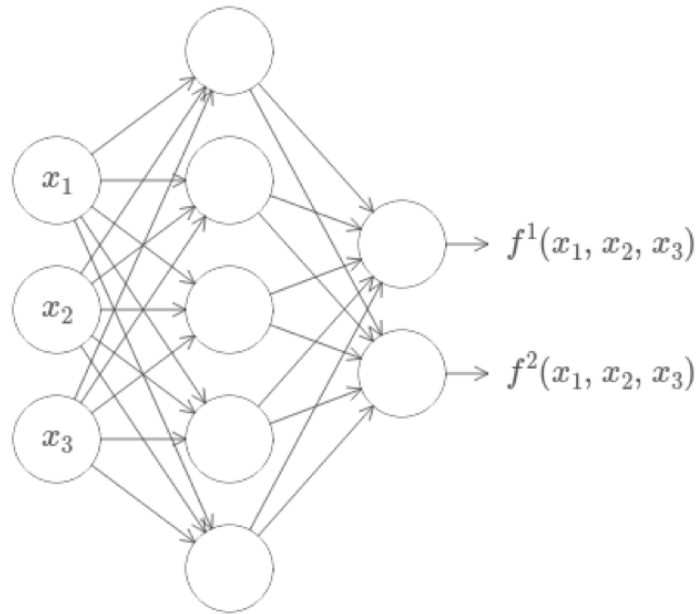
Una de las cosas más curiosas de redes neuronales es que podemos "aproximar" cualquier función por arbitraria que sea.



Esto requiere, por supuesto, conocer "ejemplos".

Esto se puede generalizar a muchos dimensiones

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$



Para esto es suficiente una única capa oculta

Teorema de universalidad

"Approximation by superpositions of sigmoidal functions"

G. Cybenko

Mathematics of Control, Signal y Systems 2 (303)

1989

Lo que nos dice el teorema es que, para cierta tolerancia $\epsilon > 0$ tal que

$$|f(x) - g(x)| < \epsilon$$

\uparrow función \nwarrow output de la red

existe un número N de neuronas ocultas.

La función debe ser continua

Una red neuronal de una capa oculta se puede usar para aproximar cualquier función continua para cualquier precisión deseada