

Descomposición Sesgo - Varianza

M.Sc. William Caicedo Torres

Universidad Tecnológica de Bolívar

caicedo77@gmail.com

21 de octubre de 2016

Material inspirado en el curso de Caltech “Learning from Data”,
por Yaser S. Abu Mustafá.

- La cota VC:

$$P[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq 4m_{\mathcal{H}}(2N)e^{-\frac{1}{8}\epsilon^2 N} \text{ (Vapnik-Chervonenkis)}$$

- Entre más complejo el conjunto de hipótesis utilizado, menor E_{in} .
- Sin embargo, en virtud de la cota VC, menor complejidad significa una mejor oportunidad de generalizar: Menor E_{out} .
- Lo ideal sería tener la complejidad justa, es decir que $\mathcal{H} = \{f\}$.

Cuantificando el Quid Pro Quo

- La aproximación del análisis VC descompone el error fuera del entrenamiento como: $E_{out} \leq E_{in} + \Omega$
- Hay otra forma de descomponer E_{out} : El análisis Sesgo-Varianza.
- Usando esta aproximación, E_{out} se compone de:
 - 1 Que tan bien aproxima \mathcal{H} a f .
 - 2 Que tan fácil es escoger una buena hipótesis h dentro de \mathcal{H} .
- Lo interesante del análisis Sesgo-Varianza es que lo podemos utilizar fácilmente en algoritmos de **Regresión** usando el **Error Cuadrado Medio**.

- Comenzemos con E_{out} :

$$E_{out}(g^{(\mathcal{D})}) = \mathbb{E}_{\mathbf{x}} \left[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2 \right]$$

$$\mathbb{E}_{(\mathcal{D})} \left[E_{out}(g^{(\mathcal{D})}) \right] = \mathbb{E}_{(\mathcal{D})} \left[\mathbb{E}_{\mathbf{x}} \left[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2 \right] \right]$$

$$\mathbb{E}_{(\mathcal{D})} \left[E_{out}(g^{(\mathcal{D})}) \right] = \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{(\mathcal{D})} \left[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2 \right] \right]$$

- Ahora, concentrémonos en:

$$\mathbb{E}_{(\mathcal{D})} \left[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2 \right]$$

La Hipótesis Promedio

- Para evaluar $\mathbb{E}_{(\mathcal{D})} [(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2]$ debemos definir la hipótesis “promedio”:

$$\bar{g}(\mathbf{x}) = \mathbb{E}_{(\mathcal{D})} [g^{(\mathcal{D})}(\mathbf{x})]$$

- Qué es la hipótesis promedio? Imagine muchos datasets $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$
- Entonces:

$$\bar{g}(\mathbf{x}) \approx \frac{1}{K} \sum_{k=1}^K g^{(\mathcal{D}_k)}(\mathbf{x})$$

- La hipótesis promedio es algo así como la “mejor” hipótesis que puedo encontrar en \mathcal{H} , a partir multiples datasets \mathcal{D}_k

Usando $\bar{g}(\mathbf{x})$ en el análisis

$$\begin{aligned}\mathbb{E}_{(\mathcal{D})} \left[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2 \right] &= \mathbb{E}_{(\mathcal{D})} \left[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}) + \bar{g}(\mathbf{x}) - f(\mathbf{x}))^2 \right] \\&= \mathbb{E}_{(\mathcal{D})} \left[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 + (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2 \right. \\&\quad \left. + 2(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))(\bar{g}(\mathbf{x}) - f(\mathbf{x})) \right] \\&= \mathbb{E}_{(\mathcal{D})} \left[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 \right] + (\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2\end{aligned}$$

Notas: 1. El valor esperado se puede distribuir con respecto a la suma. 2. $\bar{g}(\mathbf{x}) = \mathbb{E}_{(\mathcal{D})} [g^{(\mathcal{D})}(\mathbf{x})]$, por lo que la resta del primer factor del doble producto equivale a cero. 3. $\bar{g}(\mathbf{x}) - f(\mathbf{x})$ es constante con respecto a \mathcal{D} , por lo que se puede factorizar fuera del valor esperado.

$$\mathbb{E}_{(\mathcal{D})} \left[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2 \right] = \underbrace{\mathbb{E}_{(\mathcal{D})} \left[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 \right]}_{\text{Varianza}} + \underbrace{(\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2}_{\text{Sesgo}}$$

Recordemos que,

$$\mathbb{E}_{(\mathcal{D})} \left[E_{out}(g^{(\mathcal{D})}) \right] = \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{(\mathcal{D})} \left[(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x}))^2 \right] \right]$$

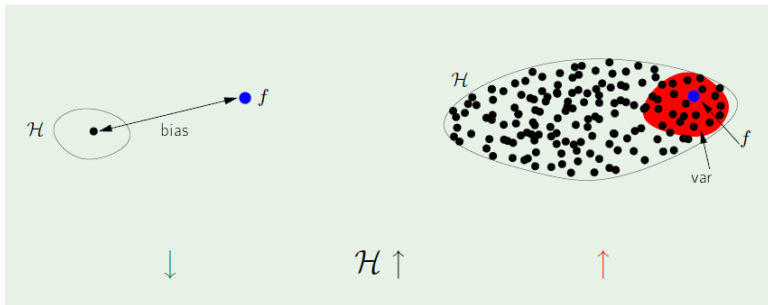
Por lo tanto,

$$\begin{aligned} \mathbb{E}_{(\mathcal{D})} \left[E_{out}(g^{(\mathcal{D})}) \right] &= \mathbb{E}_{\mathbf{x}} [\text{Varianza}(\mathbf{x}) + \text{Sesgo}(\mathbf{x})] \\ &= \text{Varianza} + \text{Sesgo} \end{aligned}$$

El Quid Pro Quo según nuestra descomposición

$$\text{sesgo} = \mathbb{E}_{\mathbf{x}} [(\bar{g}(\mathbf{x}) - f(\mathbf{x}))^2]$$

$$\text{varianza} = \mathbb{E}_{\mathbf{x}} \left[\mathbb{E}_{(\mathcal{D})} \left[(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x}))^2 \right] \right]$$



Ejemplo: Función Seno

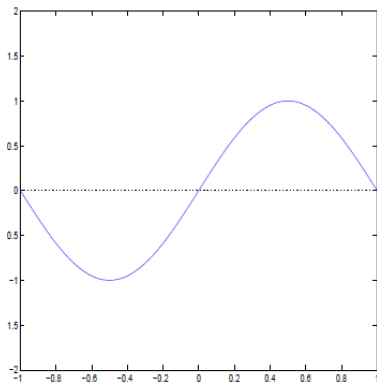
$$f : [-1, 1] \rightarrow \mathbb{R} \quad f(x) = \text{sen}(\pi x)$$

Solo disponemos de 2 ejemplos de entrenamiento Tenemos 2 modelos de aprendizaje:

$$\mathcal{H}_0 : h(x) = b$$

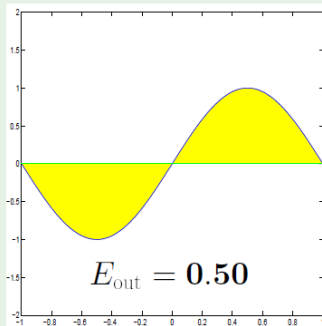
$$\mathcal{H}_1 : h(x) = ax + b$$

Cuál es mejor, \mathcal{H}_0 o \mathcal{H}_1 ?

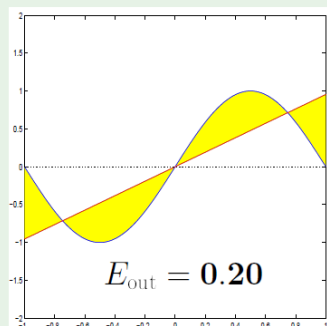


En un mundo ideal...

\mathcal{H}_0

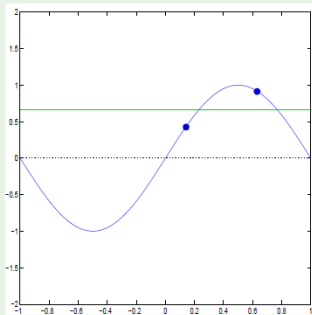


\mathcal{H}_1

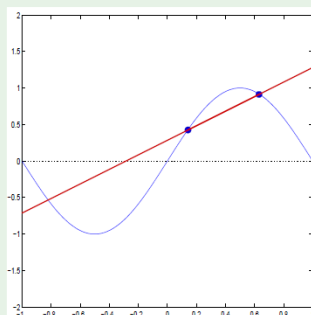


Pero solo tenemos 2 ejemplos de entrenamiento!

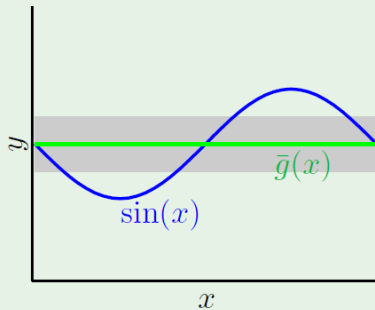
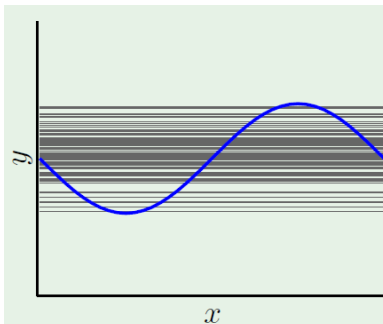
\mathcal{H}_0



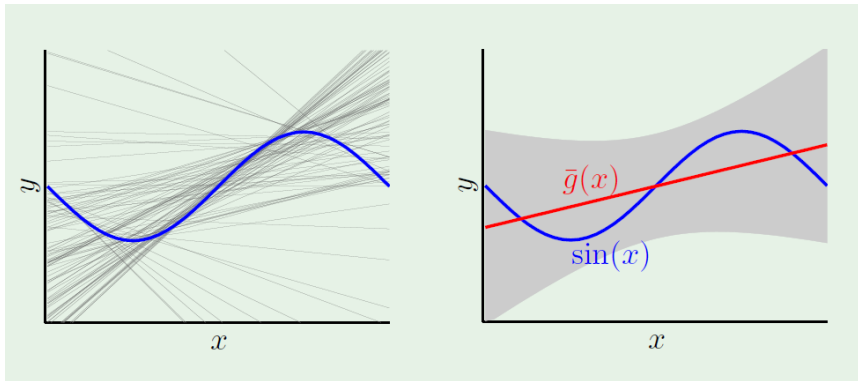
\mathcal{H}_1



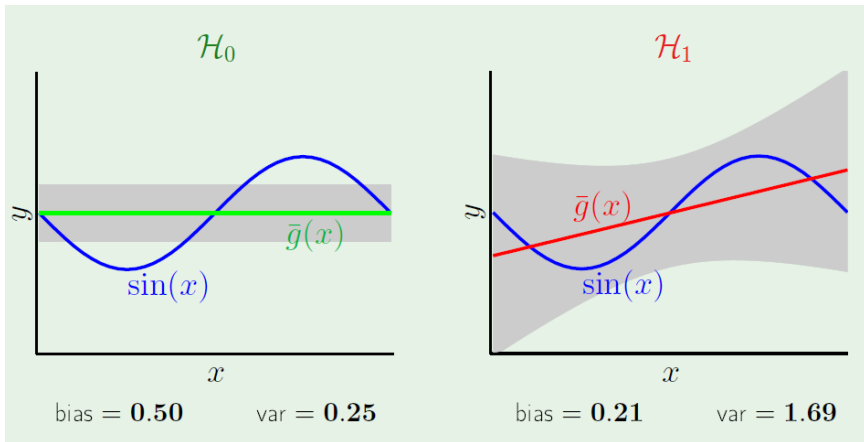
Sesgo y varianza para \mathcal{H}_0



Sesgo y varianza para \mathcal{H}_1



Y el ganador es...



Decida la complejidad del modelo de aprendizaje a utilizar de acuerdo a los datos disponibles, no de acuerdo a la complejidad que ud supone tiene la función blanco!

Muchas gracias!

Preguntas?