

Prueba de Aprendizaje Estadística EMALCA 2025

26 de Junio, 2025

Nombre: Soluciones

Ejercicio 1. (15 puntos)

Sea la función verdadera $f(x) = x^2$, y suponga que las observaciones siguen el modelo

$$Y_i = f(x_i) + \varepsilon_i = x_i^2 + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1).$$

Disponemos siempre de los mismos tres puntos de entrenamiento $x_1 = 1, x_2 = 2, x_3 = 4$, pero cada re-muestreo genera nuevos ruidos ε_i . Deseamos predecir en el punto de consulta

$$x_* = 3 \implies f(x_*) = 3^2 = 9.$$

a) Para cada modelo k -NN $k = 2, 3$, calcule el error cuadrático medio esperado $\text{MSE} = \text{Sesgo}^2 + \text{Varianza}$.

b) Determine qué valor de k minimiza la MSE esperada y explique el resultado en términos del compromiso sesgo-varianza.

a) $k=2$

$$\hat{f}_{(2)} = \frac{Y_2 + Y_3}{2} = \frac{(4 + \varepsilon_2) + (16 + \varepsilon_3)}{2} = 10 + \frac{(\varepsilon_2 + \varepsilon_3)}{2}$$

$$k=3 \quad \hat{f}_{(3)} = \frac{Y_1 + Y_2 + Y_3}{3} = \frac{(1 + \varepsilon_1) + (4 + \varepsilon_2) + (16 + \varepsilon_3)}{3} = 7 + \frac{(\varepsilon_1 + \varepsilon_2 + \varepsilon_3)}{3}$$

Error Cuadrático Medio: (ECM)

$$\begin{aligned} 2\text{-NN} \quad \text{sesgo} &: (9 - 10)^2 = 1 \\ \text{Varianza} &: \text{Var} \left(\frac{\varepsilon_2 + \varepsilon_3}{2} \right) = \frac{1}{2} \\ \text{ECM} &: 1 + \frac{1}{2} = \boxed{1.5} \end{aligned}$$

$$\begin{aligned} 3\text{-NN} \quad \text{sesgo} &: (9 - 7)^2 = 4 \\ \text{Varianza} &: \text{Var} \left(\frac{\varepsilon_1 + \varepsilon_2 + \varepsilon_3}{3} \right) = \frac{1}{3} \\ \text{ECM} &: 4 + \frac{1}{3} = \boxed{\frac{13}{3}} \end{aligned}$$

El sesgo de 2-NN es menor que 3-NN y la varianza es mayor.
En combinación el ECM de 2-NN es mejor

Ejercicio 2. (5 puntos)

Considere una situación en que tiene predecir si un paciente tiene cáncer de pulmón. En la población, solo el 1 % de los pacientes tienen cáncer de pulmón. La red neuronal que usted entreno determina la categoría (cáncer, no cáncer) correctamente 99 % de los casos del conjunto de datos de prueba. Explique de manera concisa si tiene suficiente información para determinar si su método funciona bien basado en esos porcentajes.

No hay suficiente información. Por ejemplo, el método podría clasificar a todos los pacientes "no cancer" y obtener el 99% de los casos correctos pero cometería un error para todos los pacientes que si tienen cancer!