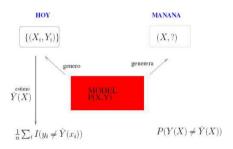


KNN, K VECINOS MÁS CERCANOS

ALAN REYES-FIGUEROA INTELIGENCIA ARTIFICIAL

(AULA 12) 17.FEBRERO.2025



- Datos tienen dos componentes: (\mathbf{x}_i, y_i) , \mathbf{x}_i son las variables "explicativas", y_i variable de respuesta.
- Distinguimos dos casos: $y_i \in \mathbb{R}$ ó $\mathbf{y}_i \in \mathbb{R}^p \Rightarrow$ regresión, $y_i \in \mathbb{N}$ ó y_i discreto \Rightarrow clasificación.

- Los datos (\mathbf{x}_i, y_i) deben ser representativos (muestra suficientemente grande).
- Pragmático. Enfoque geométrico vs. enfoque probabilístico.
- Diferentes tipos de **error**: error empírico (error de entrenamiento), error de generalización (error de validación), error de prueba.
- Cada modelo tiene asociada una complejidad.



• ¿Es cierto que $\frac{1}{n}\sum_{i}\mathbf{1}(y_{i}\neq\widehat{y}_{i})$ converge a $\mathbb{P}(Y(X)\neq\widehat{Y}(X))$?

La ley (débil) de grandes números dice que

$$\lim_{n\to\infty}\frac{1}{n}\sum_{i}Z_{i}=\mathbb{E}(Z_{i})=\mathbb{P}(Z_{i}=1).$$

Respuesta: En general, no.

- Otro problema con esta función de costo empírica es que no es continua, menos diferenciable.
- Si usas tu método de derivación favorito, no funciona.
 Pregunta: ¿cómo optimizar? Respuesta: buscamos mejores métricas de error.

Tratamos de responder la pregunta

$$\frac{1}{n}\sum_{i}\underbrace{\mathbf{1}(y_{i}\neq\widehat{y}(\mathbf{x}_{i}))}_{Z_{i}}\xrightarrow{n\to\infty}\mathbb{P}(Y(X)\neq\widehat{Y}(X))?$$

Mencionamos que en el caso de variables aleatorias Bernoulli $Z_i \sim Ber(p)$, la ley de grandes números establece que

$$\frac{1}{n}\sum_{i}Z_{i}\xrightarrow[n\to\infty]{}\mathbb{E}(Z)=\mathbb{P}(Z=1).$$

¿Vale en este caso?

No. La ley de grandes número requiere independencia de las Z_i .

En este caso, tenemos

$$\frac{1}{n}\sum_{i}\mathbf{1}(y_{i}\neq\widehat{y}(\mathbf{x}_{i}))\xrightarrow[n\to\infty]{}\mathbb{P}(Y(X)\neq\widehat{Y}(X)),$$

donde la función \hat{y} depende de todos los datos (\mathbf{x}_i, y_i) (de modo que no hay independencia de las Z_i). No aplica la ley de grandes números.

Solución ad hoc:

Separamos el conjunto (\mathbf{x}_i, y_i) en dos:

- Conjunto de entrenamiento: lo usamos para construir la función \hat{y} .
- Conjunto de validación: calculamos el error empírico $\frac{1}{n} \sum_i \mathbf{1}(y_i = \widehat{y}(\mathbf{x}_i))$. Ahora sí hay independencia, y este error empírico de validación converge al error de generalización $\mathbb{P}(Y = \widehat{Y})$.

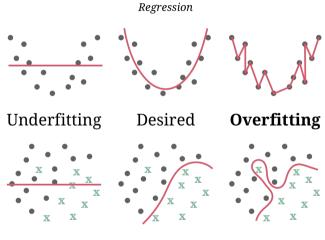
Discutimos el concepto de **complejidad** de un modelo. Este se refiere al número de parámetros involucrados en el modelo.

- en regresión: está claro, relacionado al número de variables
- en clasificación: no es tan evidente.

El concepto es importante por varias razones:

- Esto es lo que directamente va a afectar a los errores (empírico y de generalización).
- Nos va a permitir comparar diferentes modelos (en términos de simplicidad, no de exactitud).
 - Veremos que existen diferentes métricas que miden la complejidad, y nos va a permitir una segunda opinión a la hora de elegir entre diferentes modelos con similar desempeño.





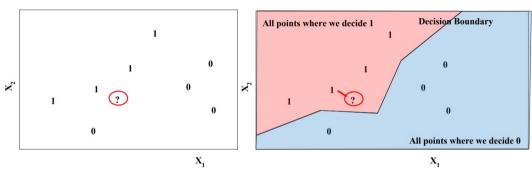


Consideramos el conjunto de datos $\{(\mathbf{x}_i, y_i)\}$, con $\mathbf{x}_i \in \mathbb{R}^d$ (en ocasiones denotamos $\mathbb{X} = (\mathbf{x}_i) \in \mathbb{R}^{n \times d}, \mathbf{Y} = (y_i) \in \mathbb{R}^n$).

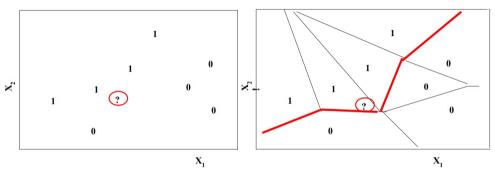
Dado $\mathbf{x} \in \mathbb{R}^d$, para decidir el valor de $\widehat{y}(\mathbf{x})$, construimos $N_k(\mathbf{x})$ el conjunto de las k observaciones más cercanas a \mathbf{x} .

- Para clasificación: decidimos por votación, esto es, asignamos a \mathbf{x} la categoría más frequente en $\{y_i : i \in N_k(\mathbf{x})\}$.
- Para regresión: decisión por promedio, *i.e.* asignamos a \mathbf{x} el promedio de $\{y_i : i \in N_k(\mathbf{x})\}.$

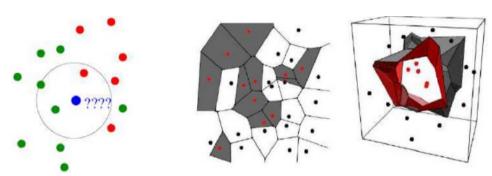
Obs! comentarios sobre cómo romper empates / métodos robustos. El caso k=1 se llama el clasificador de **vecino más cercano**.



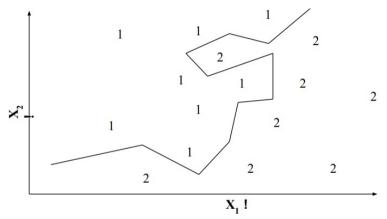
Ejemplo de KNN en el caso de clasificación.



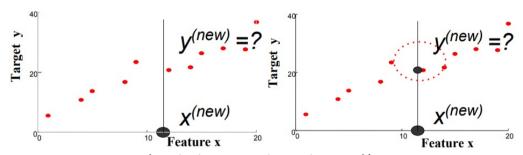
Ejemplo de KNN en el caso de clasificación. Para k=1, la frontera de clasificación coincide con un diagrama de Voronoi.



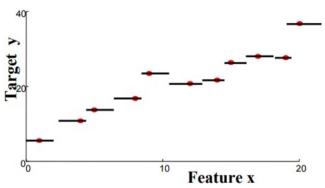
Ejemplo de KNN en el caso de clasificación. Para k = 1, la frontera de clasificación coincide con un diagrama de Voronoi.



Ejemplo de KNN en el caso de clasificación. En el caso general k > 1, la frontera sigue siendo formada por piezas poligonales (o poliedrales).

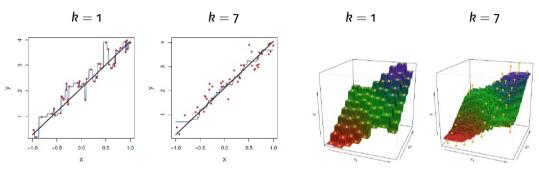


Ejemplo de KNN en el caso de regresión.



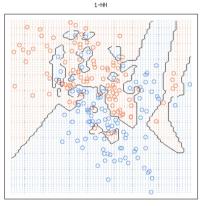
Ejemplo de KNN en el caso de regresión en el caso k = 1. Las discontinuidades ocurren en los puntos medios entre dos observaciones consecutivas.

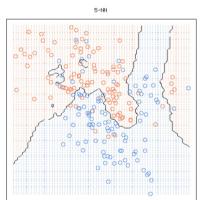
Comportamiento al variar el valor de k:



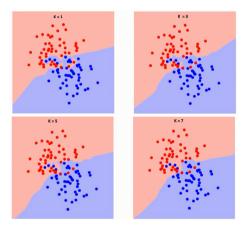
Ejemplo de KNN en el caso de regresión.

Comportamiento al variar el valor de k:

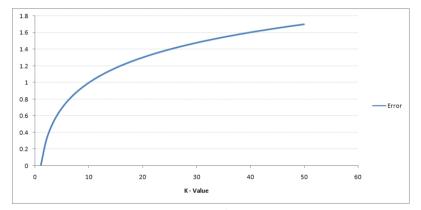




Al aumentar k las fronteras de clasificación se suavizan.

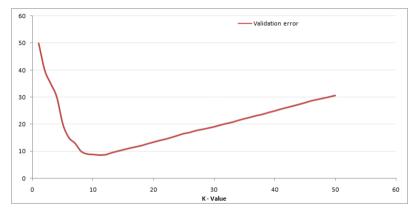


¿Cómo elegir k?

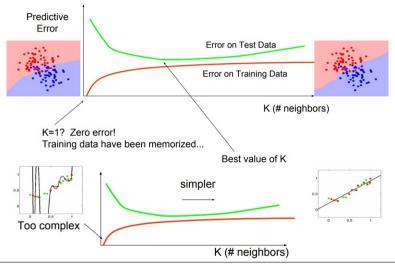


Error de entrenamiento en KNN.

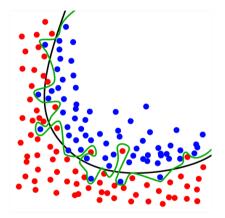
¿Cómo elegir k?



Error de validación en KNN.



Pregunta: ¿Cómo medir la complejidad en KNN?



Está relacionada con el valor $\frac{1}{k}$: Complejidad(KNN) = $\frac{1}{k}$.