

k-Vecinos más cercanos

Carlos Malanche

12 de abril de 2018

El método de *machine learning* de hoy es uno de los más intuitivos, sólo que lo trataré de presentar de manera formal.

1. Clasificación múltiple

Ahora nos enfrentamos a un problema más duro. Tenemos una serie de pares de datos $S = \{y_i, x_i\}_{i=1}^N$, y sabemos que $y_i \in C = \{\text{Clase1}, \text{Clase2}, \dots, \text{ClaseH}\}$, es decir, hay H clases a las que cada y_i puede pertenecer.

2. k-Vecinos más cercanos (*k-Nearest Neighbours*, *k-NN*)

El único requisito para que este método funcione es que las variables x_i vivan en algún espacio U al que se le pueda asignar una métrica, de tal modo que la cantidad $D(x_i, x_j)$, la distancia entre dos elementos del espacio, esté bien definida para cualesquiera dos elementos $x_i, x_j \in U$. Definimos por $\dot{S}_{\underline{x}}$ la permutación de la serie S en donde se cumple que

$$D((S_{\underline{x}})_i, \underline{x}) \leq D((S_{\underline{x}})_{i+1}, \underline{x}), \text{ para } i = 1, \dots, N-1$$

En donde $(\dot{S}_{\underline{x}})_i$ es el elemento número i de la lista. En el caso de que todas las componentes del vector \underline{x} estén en los reales, $D(x_i, x_j) := \|x_i - x_j\|$. Por último, denotamos por $\dot{S}_{\underline{x}}^k$ al conjunto que contiene los primeros k elementos de $\dot{S}_{\underline{x}}$. Un estimador de k-NN queda escrito entonces como:

$$\hat{f}(\underline{x}) = \frac{1}{k} \sum_{i: x_i \in \dot{S}_{\underline{x}}^k} y_i \quad (1)$$

Es decir, nuestra estimación es el promedio de los k elementos más cercanos a \underline{x} , según la métrica establecida. Cuando se trata de clasificaciones, basta con tomar el elemento de voto mayoritario.

