

# Clasificación: sesgo y varianza

Matías Carrasco

2 de octubre de 2023

## Índice

<b>1. El problema de clasificación</b>	<b>1</b>
1.1. El clasificador de Bayes . . . . .	1
1.2. Ejemplo: clasificación en un tablero . . . . .	2
<b>2. K vecinos más cercanos (KNN)</b>	<b>3</b>
<b>3. Sesgo y varianza en clasificación</b>	<b>4</b>
<b>4. Descomposición del error</b>	<b>6</b>

## 1. El problema de clasificación

### 1.1. El clasificador de Bayes

Consideremos el siguiente setting:

- Espacio de atributos  $\mathcal{X} \subset \mathbb{R}^D$ .
- Espacio de etiquetas  $\mathcal{Y} = \{1, \dots, m\}$  un conjunto finito.
- Espacio de observaciones  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ .
- Distribución desconocida  $\mathcal{D}$  en  $\mathcal{Z}$  (i.e. una distribución conjunta  $p(\mathbf{x}, y)$ ).
- Conjunto de datos  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^N \sim \mathcal{D}^N$ .
- Función de pérdida:  $\text{Loss}(y, y') = \mathbb{1}_{y \neq y'}$ .

El último punto sobre la función de pérdida requiere una aclaración. La función de pérdida  $\text{Loss}(y, y') = \mathbb{1}_{y \neq y'}$  se llama **0-1 loss**, y equivale a penalizar los errores

de clasificación. Veremos que en casi todos los casos se suele usar una función de pérdida subrogada a la 0-1 loss, como será el caso por ejemplo con máxima verosimilitud.

Así como en un problema de regresión tenemos la *función de regresión* que representa la predicción óptima, en clasificación también tenemos una predicción óptima.

**Definición** (Clasificador de Bayes). Fijado un  $\mathbf{x} \in \mathcal{X}$ , la predicción **óptima** en  $\mathbf{x}$  es

$$c(\mathbf{x}) = \arg \min_{\hat{y} \in \mathcal{Y}} \left\{ \mathbb{E} [\text{Loss}(\hat{y}, y) \mid \mathbf{x}] \right\}$$

Como la pérdida es la 0-1 loss, se ve inmediatamente que

$$c(\mathbf{x}) = \arg \max_{\hat{y} \in \mathcal{Y}} \left\{ \text{Prob} [\hat{y} \mid \mathbf{x}] \right\}$$

La función  $c : \mathcal{X} \rightarrow \mathcal{Y}$  se llama **clasificador de Bayes**.

**Definición** (Error de Bayes). Fijamos un atributo  $\mathbf{x} \in \mathcal{X}$ . El error esperado incurrido por el clasificador óptimo en  $\mathbf{x}$ :

$$N(\mathbf{x}) = \mathbb{E} [\text{Loss}(c(\mathbf{x}), y) \mid \mathbf{x}] = \text{Prob} [c(\mathbf{x}) \neq y \mid \mathbf{x}]$$

se llama **error de Bayes**. También se lo suele llamar *ruido* o *error irreducible*.

## 1.2. Ejemplo: clasificación en un tablero

Consideremos como ejemplo el siguiente problema de predicción del color de un casillero en un tablero:

- Atributos:  $\mathcal{X}$ ,  $\mathbf{x} = (i, j)$  la posición en el tablero.
- Etiquetas:  $\mathcal{Y}$ ,  $y = \text{rojo}(0)/\text{azul}(1)$ .
- Datos:  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ .
- Distribución:  $\mathcal{D}$  está dada por
  - $\mathbf{x}$  es uniforme en el tablero;
  - La curva verde (ver Fig. 1) divide el tablero en dos regiones, superior e inferior, y

$$\text{Prob}(\text{azul} \mid \mathbf{x}) = \begin{cases} 3/4 & \text{si } \mathbf{x} \text{ superior} \\ 1/4 & \text{si } \mathbf{x} \text{ inferior} \end{cases}$$

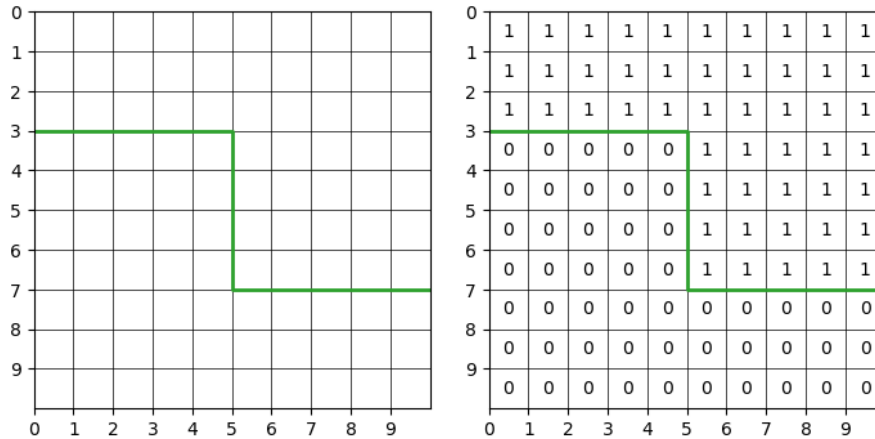


Figura 1: Ejemplo - clasificador de Bayes en el tablero.

Es fácil ver que el clasificador de Bayes es

$$c^*(\mathbf{x}) = \begin{cases} 1 & \text{si } \mathbf{x} \text{ superior} \\ 0 & \text{si } \mathbf{x} \text{ inferior} \end{cases}$$

y el error de Bayes es

$$\text{Prob}_{(\mathbf{x},y) \sim \mathcal{D}} \{y \neq c^*(\mathbf{x})\} = 1/4$$

## 2. K vecinos más cercanos (KNN)

Es un clasificador clásico y simple. El principio básico detrás de su funcionamiento es que instancias con atributos cercanos tienen etiquetas similares. Para ello es necesario elegir una forma de medir la distancia o cercanía entre puntos en el espacio de atributos  $\mathcal{X}$ .

Las distancias mas comunes son:

$$L_2 \text{ o Euclidia: } d_2(\mathbf{x}, \mathbf{x}') = \left( \sum_{j=1}^D |x_j - x'_j|^2 \right)^{1/2}$$

$$L_1 \text{ o Manhattan: } d_1(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^D |x_j - x'_j|$$

$$L_p \text{ o Minkowski: } d_p(\mathbf{x}, \mathbf{x}') = \left( \sum_{j=1}^D |x_j - x'_j|^p \right)^{1/p}$$

Como todo algoritmo de Machine Learning su input son los datos de entrenamiento  $S$  y su salida es un modelo (en este caso un clasificador)  $c(\mathbf{x})$ .

El modelo resultante  $c(\mathbf{x})$  se calcula de la siguiente manera:

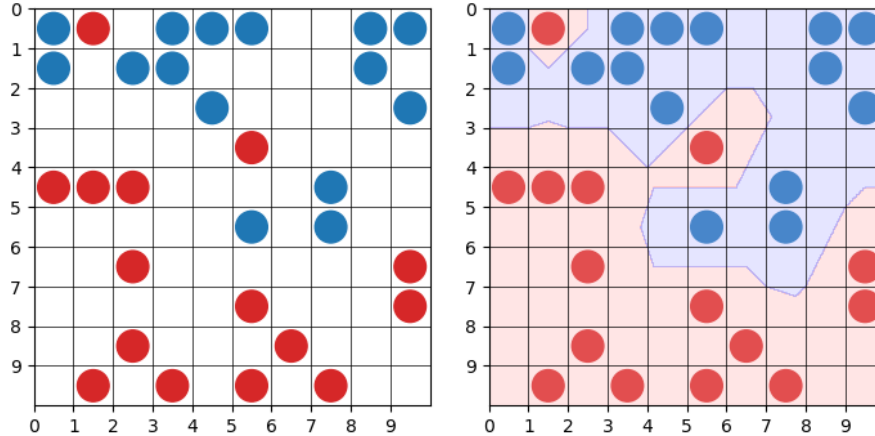


Figura 2: A la izquierda una muestra de datos  $S$  de ejemplo en el problema de predicción en el tablero. A la derecha el clasificador de  $K$  vecinos más cercanos con  $K = 1$ .

1. Calcula la distancia entre  $\mathbf{x}$  y todos los puntos en  $S$ .
2. Ordena las distancias de menor a mayor.
3. Selecciona los  $K$  puntos de  $S$  más cercanos a  $\mathbf{x}$ :  $V_K(\mathbf{x})$ .
4. Voto mayoritario:  $c(\mathbf{x}) = \arg \max_c \text{Prob}_K(y = c | \mathbf{x})$  donde

$$\text{Prob}_K(y = c | \mathbf{x}) = \frac{1}{K} \# \{ \mathbf{x}_i \in V_K(\mathbf{x}) : y_i = c \}$$

Sus hiperparámetros más comunes son: el valor de  $K$  y la distancia.

### 3. Sesgo y varianza en clasificación

Para ilustrar los conceptos de sesgo y varianza en clasificación usaremos el ejemplo del tablero con el algoritmo de  $K$  vecinos más cercanos. En la Fig. 2 tenemos una muestra de datos  $S$  de ejemplo con  $N = 30$  observaciones. A su vez se muestra el clasificador KNN para  $K = 1$ .

Al cambiar el conjunto de datos  $S$  cambia el clasificador, y por lo tanto cambia la predicción  $c_S(\mathbf{x})$  en un  $\mathbf{x}$  dado.

**Definición** (Predicción más frecuente). Fijemos un atributo  $\mathbf{x} \in \mathcal{X}$ . Llamamos  $\bar{c}(\mathbf{x})$  la **predicción más frecuente** en  $\mathbf{x}$  al variar el conjunto de datos  $S$ :

$$\bar{c}(\mathbf{x}) = \arg \max_{\hat{y} \in \mathcal{Y}} \{ \text{Prob}_{S \sim D^N} [c_S(\mathbf{x}) = \hat{y}] \}$$

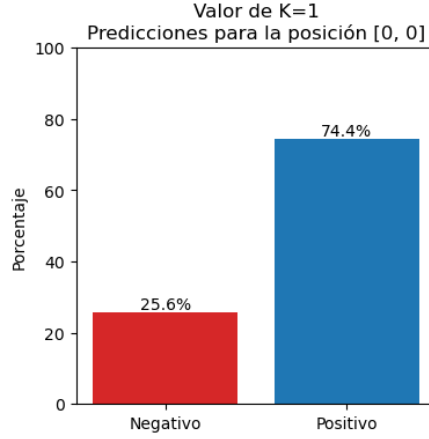


Figura 3: Simulación de 5mil conjuntos de datos  $S$  y cálculo de la predicción (con  $K = 1$ ) en el casillero  $[0, 0]$ .

En el ejemplo del tablero, es muy difícil calcular la predicción más frecuente explícitamente. Sin embargo, podemos hacer una simulación. En la Fig. 3 se muestra una simulación de la distribución de predicciones para el casillero  $[0, 0]$  usando 5mil datasets  $S$  aleatorios:

- Fijamos el casillero  $\mathbf{x} = [0, 0]$ .
- Generamos  $M = 5k$  muestras  $\{S_j\}_{j=1}^M$  todas ellas de tamaño  $N = 30$ .
- Para cada muestra  $S_j$  ajustamos KNN  $c_j$  con  $K = 1$ .
- Calculamos las predicciones  $\{c_j(\mathbf{x})\}_{j=1}^M$ .

Vemos así que en este caso  $\bar{c}([0, 0]) = 1$ .

En esta simulación, el error en  $\mathbf{x} = [0, 0]$  se aproxima de la siguiente manera

$$\text{Prob}_{S \sim D^N, y} \{y \neq c_S(\mathbf{x}) \mid \mathbf{x} = [0, 0]\} \approx \frac{1}{M} \sum_{j=1}^M \text{Prob}_y (y \neq c_j(\mathbf{x}) \mid \mathbf{x} = [0, 0])$$

Descomponiendo en  $c_j(\mathbf{x}) = 0$  o  $1$  tenemos:

$$0,256 \times \underbrace{\text{Prob}_y (y \neq 0 \mid \mathbf{x} = [0, 0])}_{0,75} + 0,744 \times \underbrace{\text{Prob}_y (y \neq 1 \mid \mathbf{x} = [0, 0])}_{0,25} = 0,378$$

A partir de la predicción más frecuente y de la predicción óptima podemos definir el sesgo.

**Definición.** Fijamos un atributo  $\mathbf{x} \in \mathcal{X}$ . El **sesgo** en  $\mathbf{x}$  se define como

$$B(\mathbf{x}) = \text{Loss}[c^*(\mathbf{x}), \bar{c}(\mathbf{x})] = \begin{cases} 1 & \text{si } c^*(\mathbf{x}) \neq \bar{c}(\mathbf{x}) \\ 0 & \text{si } c^*(\mathbf{x}) = \bar{c}(\mathbf{x}) \end{cases}$$

En el ejemplo del tablero de la Fig. 3, como tenemos  $c^*([0, 0]) = 1$  y  $\bar{c}([0, 0]) = 1$  el sesgo es  $B([0, 0]) = 0$ .

Por otro lado tenemos la varianza.

**Definición.** Fijamos un atributo  $\mathbf{x}$  cualquiera. La **varianza** en  $\mathbf{x}$  la definimos como

$$V(\mathbf{x}) = \mathbb{E}_{S \sim D^N} \left\{ \text{Loss}[\bar{c}(\mathbf{x}), c_S(\mathbf{x})] \right\} = \text{Prob}_{S \sim D^N} \left\{ c_S(\mathbf{x}) \neq \bar{c}(\mathbf{x}) \right\}$$

En el ejemplo con  $\mathbf{x} = [0, 0]$  tenemos  $\bar{c}([0, 0]) = 1$  por lo que

$$V([0, 0]) = \text{Prob}_{S \sim D^N} [c_S([0, 0]) \neq 1] = 0,256$$

## 4. Descomposición del error

La descomposición del error, al igual que en regresión, toma la forma

$$\text{Error esperado} = \text{FUNCIÓN}(\text{Sesgo}, \text{Varianza}, \text{Error Irreducible})$$

En el caso de clasificación, y para la 0-1 loss, la descomposición es un poco más compleja.

Fijemos un atributo  $\mathbf{x}$ . En clasificación **binaria** con la 0-1 loss la descomposición es

$$\underbrace{\text{Prob}\{y \neq c_S(\mathbf{x}) \mid \mathbf{x}\}}_{\text{Error}} = \underbrace{B(\mathbf{x})}_{\text{Sesgo}} + [\text{Factor}_1] \times \underbrace{V(\mathbf{x})}_{\text{Varianza}} + [\text{Factor}_2] \times \underbrace{N(\mathbf{x})}_{\text{Error Irreducible}}$$

- El factor de la varianza es

$$\text{Factor}_1 = \begin{cases} 1 & \text{si } B(\mathbf{x}) = 0 \\ -1 & \text{si } B(\mathbf{x}) = 1 \end{cases}$$

- El factor del error irreducible es

$$\text{Factor}_2 = 2 \text{Prob}\{c_S(\mathbf{x}) = c^*(\mathbf{x})\} - 1$$

Veamos cómo queda la descomposición en el ejemplo:

- Fijamos el casillero  $\mathbf{x} = [0, 0]$ .
- Tenemos:
  - Error esperado:  $\text{Prob} \{y \neq c_S(\mathbf{x}) \mid \mathbf{x} = [0, 0]\} = 0,378$
  - Sesgo:  $B([0, 0]) = 0$
  - Varianza:  $V([0, 0]) = 0,256$
  - Factor<sub>1</sub>: 1
  - Error irreducible:  $N([0, 0]) = 0,25$
  - Factor<sub>2</sub>:

$$2 \text{Prob} \{c_S(\mathbf{x}) = c^*(\mathbf{x})\} - 1 \approx 2 \times 0,744 - 1 = 0,488$$

- La descomposición queda:

$$\underbrace{0,378}_{\text{Error}} = \underbrace{0}_{\text{Sesgo}} + [1] \times \underbrace{0,256}_{\text{Varianza}} + [0,488] \times \underbrace{0,25}_{\text{Error Irreducible}}$$

La descomposición de error para todos los casilleros la podemos ver (a partir de una simulación) en la Fig. 4. Vemos que la región cercana a la frontera de decisión del clasificador de Bayes (el óptimo) es en donde el error es mayor.

Veamos la demostración de la descomposición. Recordar que estamos tratando el caso de clasificación **binaria**.

Fijemos un atributo  $\mathbf{x} \in \mathcal{X}$ . Empezamos mostrando que

$$\text{Prob} [y \neq c_S(\mathbf{x}) \mid S, \mathbf{x}] = \text{Loss} [c_S(\mathbf{x}), c^*(\mathbf{x})] + [\text{Factor}_0] N(\mathbf{x})$$

en donde

$$\text{Factor}_0 = \begin{cases} 1 & \text{si } c_S(\mathbf{x}) = c^*(\mathbf{x}) \\ -1 & \text{si } c_S(\mathbf{x}) \neq c^*(\mathbf{x}). \end{cases}$$

Recordar que  $N(\mathbf{x}) = \text{Prob} [y \neq c^*(\mathbf{x}) \mid \mathbf{x}]$ .

De hecho, si  $c_S(\mathbf{x}) = c^*(\mathbf{x})$  entonces el primer término (del lado derecho de la igualdad) es cero, el factor es uno y la igualdad es clara. Si  $c_S(\mathbf{x}) \neq c^*(\mathbf{x})$  la afirmación es

$$\text{Prob} [y \neq c_S(\mathbf{x}) \mid S, \mathbf{x}] = 1 - \text{Prob} [y \neq c^*(\mathbf{x}) \mid \mathbf{x}]$$

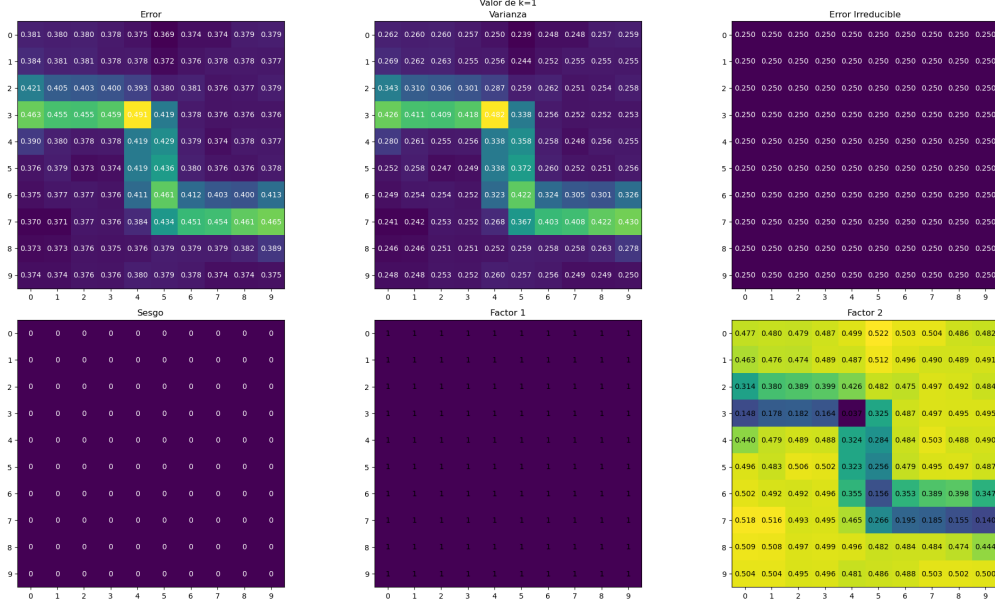


Figura 4: Descomposición del error simulada con  $5k$  datasets en el ejemplo del tablero para el algoritmo KNN con  $K = 1$ .

lo cual es cierto por tratarse de un problema de clasificación binaria.

Veamos ahora que

$$\text{Prob}_{S \sim D^N} [c_S(\mathbf{x}) \neq c^*(\mathbf{x})] = B(\mathbf{x}) + [\text{Factor}_1] V(\mathbf{x})$$

Si  $B(\mathbf{x}) = \text{Loss}[\bar{c}(\mathbf{x}), c^*(\mathbf{x})] = 0$ , entonces  $\bar{c}(\mathbf{x}) = c^*(\mathbf{x})$ , y recordando que

$$V(\mathbf{x}) = \text{Prob}_{S \sim D^N} \{c_S(\mathbf{x}) \neq \bar{c}(\mathbf{x})\}$$

vemos que la igualdad es clara con  $\text{Factor}_1 = 1$ . Por el contrario, si  $B(\mathbf{x}) = 1$ , es decir  $\bar{c}(\mathbf{x}) \neq c^*(\mathbf{x})$ , la igualdad también es clara por ser un problema de clasificación binaria (complemento igual que antes).

Por último, resta observar que el valor promedio de  $\text{Factor}_0$  al variar  $S \sim D^N$  es igual a

$$1 \times \text{Prob}_{S \sim D^N} \{c_S(\mathbf{x}) = c^*(\mathbf{x})\} + (-1) \times \text{Prob}_{S \sim D^N} \{c_S(\mathbf{x}) \neq c^*(\mathbf{x})\}$$

y otra vez, usando que es un problema de clasificación binaria, y por lo tanto

$$\text{Prob}_{S \sim D^N} \{c_S(\mathbf{x}) \neq c^*(\mathbf{x})\} = 1 - \text{Prob}_{S \sim D^N} \{c_S(\mathbf{x}) = c^*(\mathbf{x})\}$$



vemos que

$$\begin{aligned}\mathbb{E}_{S \sim D^N} [\text{Factor}_0] &= 2 \times \text{Prob}_{S \sim D^N} \{c_S(\mathbf{x}) = c^*(\mathbf{x})\} - 1 \\ &= \text{Factor}_2\end{aligned}$$

Esto demuestra la descomposición.