

# Teoria de la decision

Matias Bajac

2024-09-21

## Recapitulando...

Donde estamoss?!!!

$$Y = f(X) + \epsilon$$

- **Objetivo:** estimar  $f()$

-tipos de problemas: regresion vs clasificacion

- Porque estimar  $f$ , prediccion vs inferencia
- Compromiso flexibilidad en ajuste vs interpretabilidad de resultados

## Problema de regresión

**Objetivo:** construir un modelo que permita **predecir** las ventas en nuevos mercados, y entender en que medios es conveniente invertir

2 aproximaciones

- Vecino mas cercano
- Modelo lineal

##Modelo vecino mas cercano##

Comenzamos usando solamente TV como variable explicativa o independiente

$$Sales_i = f(TV_i) + \epsilon = E(Y_i | X_i = x_0) + \epsilon$$

Estimamos  $f(x)$  con el **promedio local** de los  $K$  puntos mas cercanos a  $x$  ( $k$  es el numero de vecinos)

$$\hat{f}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$$

$$N_x = \{x_i \in D : d_i \leq d_k\}$$

## Vecino mas cercano

Donde:

- $d_i^2 = d_2^2 = \sum_j (x = x_{ij})^2$  distancia Euclidiana

## Cantidad de vecinos

El numero de vecinos  $k$ , tiene un efecto importante en el resultado.

- Pocos vecinos: resultados muy pegados a los datos, **overfitting**
- Muchos vecinos: resultado poco sensible, promedio global.

## Modelo con todas las explicativas

- El metodo es el mismo usando las 3 variables

-Antes de comenzar hay que elegir el valor de  $k$ , cantidad de vecinos

- Para evaluar el modelo hay que calcular algun tipo de error

## Evaluacion de modelos

- No hay un metodo estadistico mejores a otros para todos los posibles problemas.
- Algunos metodos funcionan bien en algunos problemas y no en otros.
- Hay que decidir para cada problema cual metodo produce mejor resultado
- Seleccionar la mejor solucion es uno de los puntos mas desafiantes del aprendizaje estadistico

## Trade off - sesgo varianza

Como vimos antes.  $\hat{f}$  se estima a traves de los datos, ya que no estamos suponiendo un modelo para el mismo.

Para evaluar la performance de un metodo en un conjunto de datos se necesita contar con formas de medir que tan cercanas se encuentran las predicciones de los datos observados.

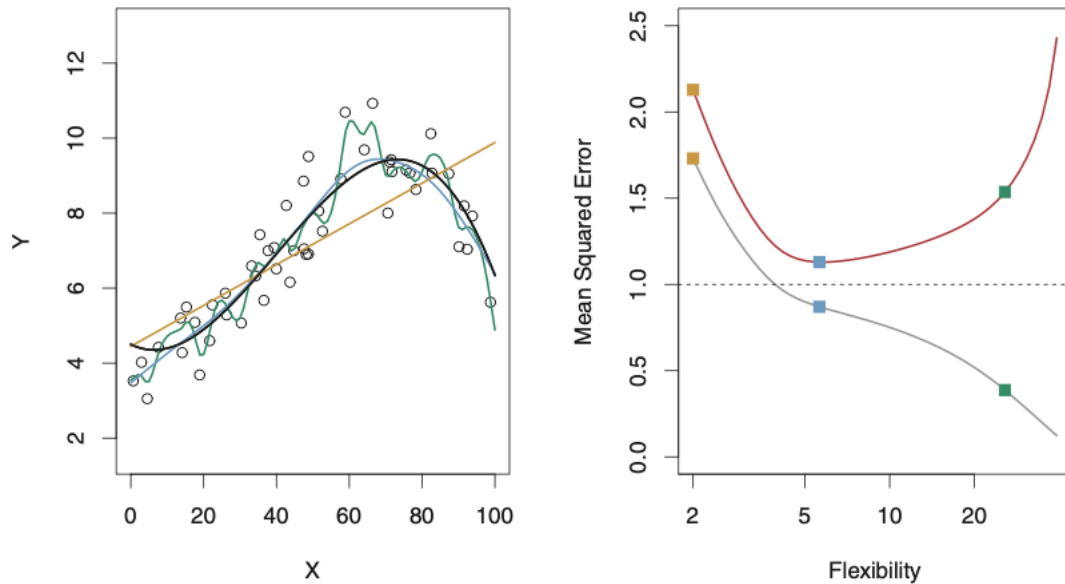
- Conjunto de entrenamiento (train)

$$L = \{(x_i, y_i)\}_{i=1}^n = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

-Conjunto de prueba(test):  $T = \{(x_0, y_0)\}$

## Performance del modelo

- **Regresion** : ECM
- con los datos de entrenamiento:  $MSE_{train} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$
- con los datos de testeo:  $MSE_{test} = AVE(y_0 - \hat{f}(x_0))^2$
- **Clasificacion** ERROR DE CLASIFICACION
- con los datos de entrenamiento:  $error_{train} = \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{f}(x_i))$



Datos simulados a partir de una cierta función  $f$  (gráfico en negro). Se ajusta recta de regresión y dos modelos flexibles (splines)

- $\uparrow \text{flexibilidad} \rightarrow \downarrow \text{MSE}_{\text{train}}$  pero no el  $\text{MSE}_{\text{test}}$  (forma de U)

Esto de que ante mayor flexibilidad menor error cuadrático medio de entrenamiento es porque el modelo ajusta mejor a los datos que ya conoce.

El MSE de test empieza a aumentar en modelos flexibles debido a que los datos de testeo no generalizan bien los nuevos datos.

- Si tenemos  $\text{MSE}_{\text{train}}$  bajo pero  $\text{MSE}_{\text{test}}$  alto hay *overfitting*

hay que tener cuidado porque a veces un modelo muy flexible puede generar *overfitting*

## Compromiso sesgo - varianza

La forma de U en las curvas  $\text{MSE}_{\text{test}}$  es el resultado del compromiso sesgo varianza.

$$E(y_0 - \hat{f}(x_0))^2 = V(\hat{f}(x_0)) + \text{bias}(\hat{f}(x_0))^2 + V(\epsilon)$$

Para minimizar el error test esperado necesitamos  $f$  que tenga simultáneamente varianza y sesgo pequeños.

- La varianza de un modelo de AA mide cuánto varía  $\hat{f}$  al modificar el conjunto de entrenamiento, es decir es una medida de la estabilidad técnica (en general los métodos más flexibles tienen mayor varianza)
- El sesgo de un modelo AA cuantifica el error de elección del modelo (en general los métodos más flexibles tienen menos sesgo)

El objetivo consiste en encontrar una función que minimice el riesgo de predecir mal, para ello se define una función de pérdida  $L(X, f(X), Y)$  y se busca  $f^*$  entre todas las funciones de una cierta clase  $C$ , que haga mínimo el valor esperado de  $L$

$$E(L(X, f(X), Y))$$

que llamamos riesgo, en la práctica buscamos aquella función  $\hat{f}$  que minimiza el riesgo empírico

$$\frac{1}{n} \sum_{i=1}^{i=n} L(X_i, f(X_i), Y_i)$$

## Errores

Si  $f^*$  es el mejor entre todos los predictores posibles,  $f^{**}$  el mejor de los predictores posibles en una cierta clase de funciones  $C$  y  $\hat{f}$  el predictor que usamos en la practica tenemos dos tipos de error:

- Error de modelizacion  $f - f^{**}$ : depende de la eleccion de la clase  $C$ , si consideramos como la familia de todas las funciones posibles, tendremos overfitting.
- Error de estimacion  $f^{**} - \hat{f}$ : es un error estadistico, si el tamaño  $n$  de la muestra es grande, bajo cierta hipotesis sobre la clase  $C$ , se cumple que  $\hat{f}$  converge a  $f^{**}$