

# Tema 5

## Técnicas de regresión avanzadas (I)

### Técnicas multivariantes

Dr. Antoni Ferragut



- Podemos generalizar
- Tenemos

$$SST = SSE + SSR,$$

donde SST es la suma de cuadrados total, SSE la suma de cuadrados del error y SSR la suma de cuadrados de la regresión.

- SST mide la variabilidad de la muestra  $Y$  respecto a la media muestral.
- SSE mide la variabilidad de los datos de  $Y$  respecto a los valores ajustados.
- SSR es la parte de la varianza que explican todas las variables predictoras  $X_j$  del modelo.

- Podemos definir el coeficiente de determinación general:

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}.$$

- En este caso general, no resulta útil para poder elegir entre diferentes modelos con un número de variables predictoras distinto.
- A medida que añadimos variables, por definición  $R^2$  aumenta.
- Podemos definir un coeficiente de determinación ajustado:

$$R_{\text{adj}}^2 = 1 - \frac{(n-1)SSE}{(n-m-1)SST}.$$

- Con este nuevo coeficiente, al aumentar  $m$  el valor de  $R_{\text{adj}}^2$  decrece.
- Este coeficiente penaliza modelos complejos.

- Una hipótesis extra a comprobar es la colinealidad.
- Colinealidad: relación lineal entre dos o más variables.
- Problemas de identificabilidad: ¿a qué se deben los cambios en la variable respuesta?
- Pequeños cambios en la muestra pueden provocar grandes cambios en la estimación de los parámetros.
- El VIF (factor de inflación de la varianza) detecta esta colinealidad:

$$\text{VIF} = \frac{1}{1 - R^2}.$$

- A mayor VIF, mayor colinealidad. Para  $\text{VIF} > 10$  puede haber un problema de colinealidad.
- El análisis puede hacerse antes de la regresión o después con las variables que se hayan usado para ajustar el modelo.

### **Situaciones:**

- Algunas variables predictoras no aportan información.
- Algunas variables están muy relacionadas entre sí (colinealidad).
- Si hay dos modelos con capacidad predictiva similar, escogemos el más sencillo (parsimonia/Ockham).

### **Best subset:**

- Fijado el número de variables se analizan todos los modelos posibles y se obtiene el mejor según capacidades predictivas (p.ej.  $R^2$ ,  $R^2_{\text{adj}} \dots$ ).
- Para  $m$  grande tiene un coste de computación elevado (crecimiento exponencial:  $2^m$  modelos).

### Stepwise selection:

- Computacionalmente más económica.
- Se busca el mejor modelo de manera secuencial.
- Selección hacia adelante: partiendo del modelo nulo se añade en cada paso la variable que añade más información al modelo. El proceso termina cuando añadir una variable no mejora el modelo.
- Selección hacia atrás: partiendo del modelo con todas las variables, se elimina en cada paso la variable tal que el modelo mejora más. El proceso termina cuando al eliminar variables no mejora el modelo.
- Si  $m > n$  solamente se puede usar la selección hacia adelante y terminar en el paso  $n$ .

- Datos perdidos, medidas no tomadas o tomadas mal, etc.
- Una primera opción es eliminar las observaciones con algún dato faltante en alguna de las variables.
- La usaremos con muestras grandes y pocos datos faltantes, o si estimarlos es difícil.
- Perdemos información que teníamos.
- Como alternativa, podemos usar una regresión para poder imputar el dato faltante o imputarlo por la mediana o la moda.



### **Pasos:**

- Centrar las variables (restar la media).
- Estandarizar (restar la media y dividir por la desviación estándar).
- Escalar la variable a  $[0, 1]$ .