

# **IECD-325: Métodos de selección automática**

**Felipe Osorio**

[felipe.osorio@uv.cl](mailto:felipe.osorio@uv.cl)

## Métodos de selección automática

Suponga el modelo de regresión<sup>1</sup>

$$Y_i = \beta_1 x_{i1} + \cdots + \beta_K x_{iK} + \epsilon_i, \quad i = 1, \dots, n,$$

que puede ser escrito como  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$  con los supuestos habituales, y  $\mathbf{X} \in \mathbb{R}^{n \times K}$ ,  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_K)^\top$ .

Suponga que particionamos

$$\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2),$$

donde  $\mathbf{X}_1 \in \mathbb{R}^{n \times p}$  y  $\mathbf{X}_2 \in \mathbb{R}^{n \times (K-p)}$  y análogamente  $\boldsymbol{\beta} = (\boldsymbol{\beta}_1^\top, \boldsymbol{\beta}_2^\top)^\top$ . El modelo adopta la forma,

$$\mathbf{Y} = \mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{X}_2 \boldsymbol{\beta}_2 + \boldsymbol{\epsilon}.$$

### Objetivo:

Deseamos identificar las variables “significativas” con coeficientes no nulos.

---

<sup>1</sup>Si  $x_{i1} = 1$ , para  $i = 1, \dots, n$ , tenemos un intercepto en el modelo.

## Métodos de selección automática

Suponga que sospechamos que los coeficientes asociados a las  $(K-p)$ -variables en  $\mathbf{X}_2$  son cero. Es decir, deseamos probar  $H_0 : \beta_2 = \mathbf{0}$ .

En otras palabras, debemos discriminar entre 2 modelos, uno con  $K$  variables y el otro con  $p$  variables. Para esto, podemos considerar el estadístico  $F$ , dado por:

$$F = \left( \frac{n - K}{K - p} \right) \frac{\text{RSS}_p - \text{RSS}_K}{\text{RSS}_K}.$$

Por otro lado, suponga que  $K = p + 1$ , así

$$F = (n - p - 1) \frac{\text{RSS}_p - \text{RSS}_{p+1}}{\text{RSS}_{p+1}}. \quad (1)$$

## Métodos de selección automática

Consideraremos los siguientes procedimientos:

- ▶ Selección forward.
- ▶ Eliminación backward.
- ▶ Método stepwise.

Adicionalmente, un procedimiento que ha ganado popularidad es [regresión lasso](#), el cual resuelve el problema

$$\min_{\beta} Q_1(\beta), \quad Q_1(\beta) = \|\mathbf{Y} - \mathbf{X}\beta\|^2 + \lambda \sum_{j=1}^p |\beta_j|,$$

que corresponde a un método de selección de regresores.

# Métodos de selección automática

## Selección forward (SF):

Este método añade variables regresoras una en cada vez, partiendo desde el modelo más simple, y una vez que una variable ingresa al modelo, esta **no** es retirada.

La elección que cual variable ingresa al inicio es arbitraria. Sin embargo, es usual inicializar el procedimiento con el modelo que solo tiene intercepto.

Calculamos (1) con  $p = 1$  para todas las  $m = K - 1$  variables regresoras restantes, y escogemos aquella tal que (1) sea la mayor.

Luego de eso, repetimos el procedimiento para  $p = 2, 3, \dots$  seleccionando en cada etapa una variable que no se había seleccionado previamente.

Es decir, añadimos la variable  $x_j$  al modelo con  $p$ -regresores, si:

$$F_j = \frac{\text{RSS}_p - \text{RSS}_{p+1}}{\text{RSS}_{p+1} / (n - p - 1)} > F_{\text{IN}},$$

donde  $F_{\text{IN}} = F_{1-\alpha}(1, n - p - 1)$  para algún valor de  $\alpha$ ,<sup>2</sup> o alternativamente  $F_{\text{IN}} = 2$ .

---

<sup>2</sup>Se recomienda usar  $0.15 \leq \alpha \leq 0.25$ .

## Métodos de selección automática

### Eliminación backward (EB):

Este método se inicia con aquél que tiene  $K$  variables (i.e., el modelo más complejo) y elimina una variable en cada ocasión.

Es decir, se debe calcular (1) y debemos retirar la variable que ocasiona la menor contribución.

Es decir, eliminamos la variable  $x_j$  del modelo con  $p$ -regresores, si:

$$F_j = \frac{\text{RSS}_{p-1} - \text{RSS}_p}{\text{RSS}_p / (n - p)} < F_{\text{OUT}},$$

donde  $F_{\text{OUT}} = F_{1-\alpha}(1, n - p)$  para algún valor de  $\alpha$ , o alternativamente  $F_{\text{OUT}} = 2$ .

# Métodos de selección automática

## Método stepwise:

Este método combina selección forward y eliminación backward el cual corresponde a desarrollar un paso SF seguido por un paso EB en cada etapa.

Este algoritmo inicia con el modelo más simple y añade variables si el estadístico  $F$  asociado<sup>3</sup> es mayor que  $F_{IN}$  o elimina variables si es menor que  $F_{OUT}$ , siempre que  $F_{OUT} \leq F_{IN}$ .

---

<sup>3</sup>ver Ecuación (1).

## Métodos de selección automática

### *Observaciones:*

- ▶ El conjunto de regresores obtenidos mediante cada uno de los métodos puede ser diferente.
- ▶ Estos métodos pueden fallar en detectar el mejor subconjunto de regresores.
- ▶ El método EB asume que  $X$  tiene rango completo, en caso contrario (o si el número de potenciales regresores es muy alto), las únicas opciones factibles son SF o stepwise.
- ▶ Criterios de selección de modelos, como  $C_p$ ,  $s_p^2$  y  $AIC_p$  pueden ser usados como mecanismos de búsqueda.<sup>4</sup>

---

<sup>4</sup>En R están disponibles las funciones `step` o `stepAIC`.

## usando stepAIC

```
1 # ajustando el modelo de regresión
2 fm <- lm(y ~ ., data = portland)
3
4 # cargando la biblioteca MASS
5 library(MASS)
6
7 # ejecutando eliminación backward
8 > EB <- stepAIC(fm, trace = TRUE, direction = "backward")
9 Start: AIC=26.94
10 y ~ x1 + x2 + x3 + x4
11
12          Df Sum of Sq    RSS     AIC
13 - x3      1   0.1091  47.973  24.974
14 - x4      1   0.2470  48.111  25.011
15 - x2      1   2.9725  50.836  25.728
16 <none>             47.864  26.944
17 - x1      1  25.9509  73.815  30.576
18
19 Step: AIC=24.97
20 y ~ x1 + x2 + x4
21
22          Df Sum of Sq    RSS     AIC
23 <none>                 47.97  24.974
24 - x4      1     9.93  57.90  25.420
25 - x2      1    26.79  74.76  28.742
26 - x1      1  820.91 868.88 60.629
```

## usando stepAIC

```
1 # resumen del proceso
2 > EB$anova
3 Stepwise Model Path
4 Analysis of Deviance Table
5
6 Initial Model:
7 y ~ x1 + x2 + x3 + x4
8
9 Final Model:
10 y ~ x1 + x2 + x4
11
12
13   Step Df Deviance Resid. Df Resid. Dev      AIC
14 1                   8    47.86364 26.94429
15 2 - x3  1    0.10909          9    47.97273 24.97388
```

## usando stepAIC

```
1 # modelo inicial y 'full'
2 > empty <- lm(y ~ 1, data = portland)
3 > full <- formula(y ~ x1 + x2 + x3 + x4)
4
5 # ejecutando selección forward
6 > SF <- stepAIC(empty, direction = "forward", scope = full)
7 Start: AIC=71.44
8 y ~ 1
9
10          Df Sum of Sq      RSS      AIC
11 + x4      1    1831.90  883.87  58.852
12 + x2      1    1809.43  906.34  59.178
13 + x1      1    1450.08 1265.69  63.519
14 + x3      1     776.36 1939.40  69.067
15 <none>            2715.76 71.444
16
17 ...
```

## usando stepAIC

```
1 # ... continuación
2
3 Step: AIC=58.85
4 y ~ x4
5
6          Df Sum of Sq    RSS    AIC
7 + x1      1     809.10  74.76 28.742
8 + x3      1     708.13 175.74 39.853
9 <none>                    883.87 58.852
10 + x2     1     14.99 868.88 60.629
11
12 Step: AIC=28.74
13 y ~ x4 + x1
14
15          Df Sum of Sq    RSS    AIC
16 + x2      1     26.789 47.973 24.974
17 + x3      1     23.926 50.836 25.728
18 <none>                    74.762 28.742
19
20 Step: AIC=24.97
21 y ~ x4 + x1 + x2
22
23          Df Sum of Sq    RSS    AIC
24 <none>                    47.973 24.974
25 + x3      1     0.10909 47.864 26.944
```

## usando stepAIC

```
1 # resumen del proceso
2 > SF$anova
3 Stepwise Model Path
4 Analysis of Deviance Table
5
6 Initial Model:
7 y ~ 1
8
9 Final Model:
10 y ~ x4 + x1 + x2
11
12
13   Step Df     Deviance Resid. Df Resid. Dev      AIC
14   1                   12 2715.76308 71.44443
15   2 + x4   1 1831.89616    11 883.86692 58.85164
16   3 + x1   1  809.10480    10 74.76211 28.74170
17   4 + x2   1   26.78938      9 47.97273 24.97388
```

## usando stepAIC

```
1 # ejecutando el método stepwise
2 > SW <- stepAIC(empty, direction = "both", scope = full)
3 Start: AIC=71.44
4 y ~ 1
5
6          Df Sum of Sq      RSS      AIC
7 + x4      1   1831.90  883.87  58.852
8 + x2      1   1809.43  906.34  59.178
9 + x1      1   1450.08 1265.69  63.519
10 + x3     1    776.36 1939.40  69.067
11 <none>                    2715.76 71.444
12
13 Step: AIC=58.85
14 y ~ x4
15
16          Df Sum of Sq      RSS      AIC
17 + x1      1    809.10   74.76  28.742
18 + x3     1    708.13  175.74  39.853
19 <none>                    883.87  58.852
20 + x2      1     14.99  868.88  60.629
21 - x4      1   1831.90 2715.76 71.444
22
23 ...
```

## usando stepAIC

```
1 # ... continuación
2
3 Step: AIC=28.74
4 y ~ x4 + x1
5
6          Df Sum of Sq      RSS      AIC
7 + x2     1    26.79  47.97  24.974
8 + x3     1    23.93  50.84  25.728
9 <none>                    74.76  28.742
10 - x1    1   809.10  883.87  58.852
11 - x4    1  1190.92 1265.69  63.519
12
13 Step: AIC=24.97
14 y ~ x4 + x1 + x2
15
16          Df Sum of Sq      RSS      AIC
17 <none>                    47.97  24.974
18 - x4     1     9.93  57.90  25.420
19 + x3     1     0.11  47.86  26.944
20 - x2     1    26.79  74.76  28.742
21 - x1     1   820.91  868.88  60.629
```

## usando stepAIC

```
1 # resumen del proceso
2 > SW$anova
3 Stepwise Model Path
4 Analysis of Deviance Table
5
6 Initial Model:
7 y ~ 1
8
9 Final Model:
10 y ~ x4 + x1 + x2
11
12
13   Step Df     Deviance Resid. Df Resid. Dev      AIC
14   1                   12 2715.76308 71.44443
15   2 + x4   1 1831.89616    11 883.86692 58.85164
16   3 + x1   1  809.10480    10 74.76211 28.74170
17   4 + x2   1   26.78938      9 47.97273 24.97388
```