

Escala mientro de longitud unitaria

$$X = \begin{bmatrix} 1 & x_1 & \dots & x_n \\ \vdots & \vdots & & \vdots \end{bmatrix}_{n \times p}$$

Nota normalizar es dividir entre la norma (e.g. $0 = \frac{x}{\|x\|}$)

Se coge X_j y se normaliza luego de restarle la media de la columna

$$X_j^* = \frac{1}{\sqrt{\sum_{k=1}^n (x_{kj} - \bar{x}_j)^2}} \begin{bmatrix} x_{1j} - \bar{x}_j \\ x_{2j} - \bar{x}_j \\ \vdots \\ x_{nj} - \bar{x}_j \end{bmatrix}$$

Se llama de longitud unitaria porque las columnas resultantes son de longitud 1.

- i) Y (la respuesta) tambien se normaliza (Y^*)
- ii) Se ajusta el modelo sin intercepto

$$Y_i^* = \sum_{j=1}^K \beta_j^* X_{ij}^*$$

¿Por qué no incluir β_0 ? Una superficie de regresión pasa por el punto (\bar{Y}, \bar{X})

Supongamos que se agrega β_0

$$Y_i^* = \beta_0 + \sum_{j=1}^K \beta_j^* X_{ij}^* \quad \text{Cambiando } X_{ij}^* \text{ por la media muestral se } j=\hat{n} \text{ tiene que}$$

$$Y_i^* = \beta_0 + \sum_{j=\hat{n}}^K \beta_j^* \cdot 0 = \beta_0, \text{ Pero la regresión evaluada en } \bar{X} \text{ es igual a } \bar{Y} \Rightarrow 0 = \beta_0$$

1. Responda las siguientes preguntas.

- a) Suponga que se realiza escalamiento de longitud unitaria en las predictoras pero no en la variable respuesta, ¿qué unidades tienen los coeficientes de la regresión una vez esta es ajustada?
- b) ¿Por qué hay problemas de multicolinealidad cuando se tienen más covariables que observaciones en los datos?
- c) Si la traza de la matriz $\mathbf{X}'\mathbf{X}$ es muy grande, ¿mayor es la distancia entre el vector de parámetros estimados y el verdadero vector de parámetros?
- d) Si la correlación entre las variables X_j y X_k es pequeña, ¿se puede descartar la presencia de multicolinealidad?

a) $Y = \beta_0 + \beta_1 X_1^* + \dots + \beta_k X_k^* + \varepsilon, \varepsilon \sim N(0, \sigma^2 I_n)$

R// Como X_j^* es adimensional para todo j , se tiene que los coeficientes de regresión cargan con las mismas unidades de la respuesta

b) $0 \leq \text{Rango}(\mathbf{X}) \leq \min(n, p) = h$

$$0 \geq -\text{Rango}(\mathbf{X}) \geq -h$$

$$p \geq p - \text{Rango}(\mathbf{X}) \geq p - h > 0$$

$$(p > h)$$

¿?
H variables
libres

H variables libres > 0

Si uno considera $\text{Rango}(\mathbf{X}'\mathbf{X}) > 0 \Rightarrow |\mathbf{X}'\mathbf{X}| \neq 0$

$\mathbf{X}'\mathbf{X} \beta = \mathbf{X}'\mathbf{Y}$ este sistema $\begin{cases} \text{no tiene sol} \\ \text{tiene infinitas} \\ \text{soluciones} \end{cases}$

c) falso, es $(\mathbf{X}'\mathbf{X})^{-1}$

$$\text{TR}([X^T X]^{-1}) = \sum_{i=1}^n \frac{1}{\lambda_i}$$

Si $X^T X$ no es invertible $\Rightarrow \lambda_i = 0$ para algún i

$\Rightarrow \sum_{i=1}^n \frac{1}{\lambda_i}$ explota

1) No, porque la correlación solo mide la dependencia lineal entre dos variables

e.g. Si $X_1 = 2X_2 + 3X_3$ eso probablemente no sea detectado en $\text{Cor}(X_1, X_2)$ y $\text{Cor}(X_1, X_3)$

2) Para la VIF (factor de inflación de varianzas)

VIF_j = $\begin{cases} \leq 5 \Rightarrow \text{no hay problemas} \\ 5 < VIF \leq 10 \text{ hay problemas moderados} \\ > 10 \text{ hay problemas graves} \end{cases}$

```
> car::vif(mod1)
      x1      x2
294997.4 294997.4
> car::vif(mod2)
      x1      x2
172810.7 172810.7
```

} Hay problemas graves

2)

Haciendo análisis espectral

Collinearity Diagnostics						
	Eigen_Value	Condition_Index	Variance Decomposition Proportions			
			Intercept	x1	x2	
1	2.8349e+00	1.000000	0.024457	0.000000	0.000000	x_3
2	1.6508e-01	4.143992	0.923034	0.000001	0.000001	x_4
3	5.4051e-07	2290.167014	0.052509	0.999999	0.999999	0.6 0.1

x_i
está en orden
descendente

$\sqrt{K_j}$
índice
de condición

π_{ij} proporciones de
descomposición de
varianza

- # Condición $\sqrt{K} = \sqrt{\lambda_{\max} / \lambda_{\min}} = 2290.167$ (Se saca de Condition-Index) (es el mayor)

Si $\sqrt{K} \leq 10$, todo ok (no hay problemas)
 $10 < \sqrt{K} \leq 31.62$, problemas moderados
 $\sqrt{K} > 31.62$, problemas graves

- índice de condición $\sqrt{K_j} = \sqrt{\lambda_{\max} / \lambda_j}$

$\sqrt{K} \leq 10$, todo ok (no hay problemas) Para todo j
 $10 < \sqrt{K} \leq 31.62$, problemas moderados Para algún j
 $\sqrt{K_j} > 31.62$, problemas graves

hace referencia a los coeficientes de
regresión

Proporción de descomposición de Varianza: para $i=3$
 π_{13} y π_{23} son 70.5 lo que indica evidencia de
 multicolinealidad entre X_1 y X_2

3)

k	R_sq	adj_R_sq	SSE	Cp	Variables_in_model		
1	0.017	0.016	29042245	0.101	latitude		
1	0.002	0.001	29475559	32.538	longitud		
1	0.000	0.000	29527416	36.420	magnitud		
2	0.017	0.016	29040930	2.002	latitude	longitud	
2	0.017	0.016	29042210	2.098	latitude	magnitud	
2	0.002	0.001	29472811	34.333	longitud	magnitud	
3	0.017	0.015	29040902	4.000	latitude	longitud	magnitud

se mira
 $|C_p - p|$
 pequeño

R^2_{sq} : es la columna del R^2_p y se desea el
 modelo con R^2 más grande
 nos quedamos con el 1 por parsimonia

adj- R^2_{sq} : R^2_p como busca el más alto, nos quedamos
 con el modelo 1

C_p : lo que se busca es que el C_p sea bajito y que a
 su vez $|C_p - p|$ sea mínimo también (posible)

El mejor modelo según este criterio es el 3