

Regularización de modelos lineales y aditivos

Ridge, Lasso y Redes Elásticas

Gabriel Martos Venturini
gmartos@utdt.edu

Universidad Torcuato Di Tella

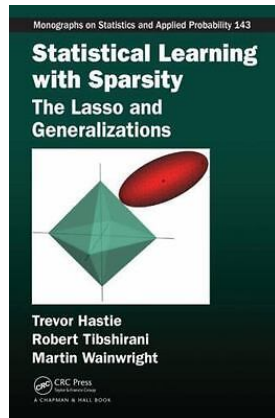
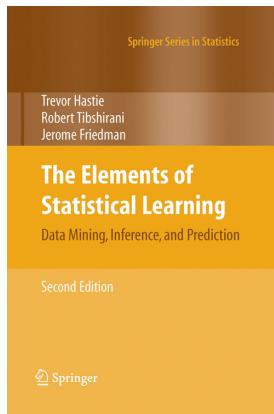
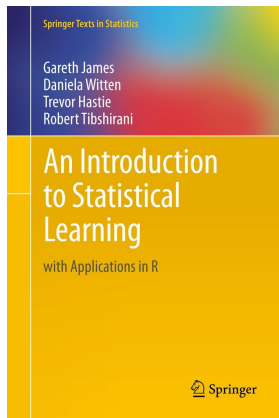


Agenda

Regularización de modelos lineales y aditivos

Apéndice

Bibliografía recomendada



ISL: 6 y 7.1–7.5.

ESL: 3.1–3.4 y 5.1–5.5.

SLS: 1 y 2.

Agenda

Regularización de modelos lineales y aditivos

Apéndice

- ▶ Cuando $p \gg 0$ (típico contexto con modelos aditivos) los modelos de regresión suelen presentar mucha *variabilidad*¹.
 - ▶ Esto se traduce en que ante perturbaciones en los datos de entrenamiento se producen cambios considerables en los $p \gg 0$ parámetros estimados (aprendidos) del modelo en cuestión.
 - ▶ Error Esperado Modelo = $\text{Bias}^2 + \text{Variance} + \text{cte.}$
- ▶ Para intentar controlar este fenómeno vamos a introducir una *restricción de presupuesto* sobre los parámetros del modelo de forma tal de controlar la variabilidad de forma explícita.
 - ▶ Reformular el problema de minimización del riesgo empírico.
- ▶ Será necesario estimar (aprender) un hiperparámetro adicional (λ) que controla el trade-off entre sesgo y variabilidad.
- ▶ Omitimos la referencia a los modelos aditivos para simplificar.

¹ *The curse of dimensionality* ▶ ver apéndice en § 26 .

Notación

- ▶ La discusión vale para modelos aditivos en general:

$$Y = \beta_0 + \sum_{b=1}^{B_1} \beta_{1b} \phi_b(X_1) + \sum_{b=1}^{B_2} \beta_{2b} \phi_b(X_2) + \cdots + \sum_{b=1}^{B_p} \beta_{pb} \phi_b(X_p) + \varepsilon$$

- ▶ Pero para simplificar la notación, siempre se plantea que:

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \varepsilon.$$

- ▶ Notar que se trata simplemente de renombrar los features:

- ▶ $X_1 \equiv \phi_1(X_1)$ y $\beta_1 \equiv \beta_{11}$.

- ▶ $X_2 \equiv \phi_2(X_1)$ y $\beta_2 \equiv \beta_{12}$.

- ▶ ...

- ▶ $X_p \equiv \phi_{B_p}(X_p)$ y $\beta_p \equiv \beta_{pB_p}$ (notar que generalmente $p \gg 0$).

Ridge

- ▶ Asumimos que $Y = \underbrace{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}_{\text{Modelo para } f(X)} + \varepsilon$.

- ▶ Restricción que involucra a las pendientes: $\sum_{i=1}^p \beta_i^2 \leq \tau^2$.

- ▶ Pedimos que $(\beta_1, \dots, \beta_p) \in B(\mathbf{0}, \tau)$.

- ▶ Dada una muestra de train $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, los parámetros del modelo $\hat{\beta}_\tau^{\text{ridge}} \equiv (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)_\tau^{\text{ridge}}$ se aprenden resolviendo:

$$\min_{b_0, b_1, \dots, b_p} \underbrace{\sum_{i=1}^n \left(y_i - (b_0 + b_1 x_{1,i} + \dots + b_p x_{p,i}) \right)^2}_{\text{RSS}(\text{Datos}, b_0, b_1, \dots, b_p)} \quad \text{s.a.} \quad \underbrace{\sum_{j=1}^p b_j^2}_{\text{Restricción de presupuesto}} \leq \tau^2$$

- ▶ Veamos una ilustración (cuando $p = 2$) para tener intuición sobre las consecuencias de hacer variar τ sobre el fitting del modelo.

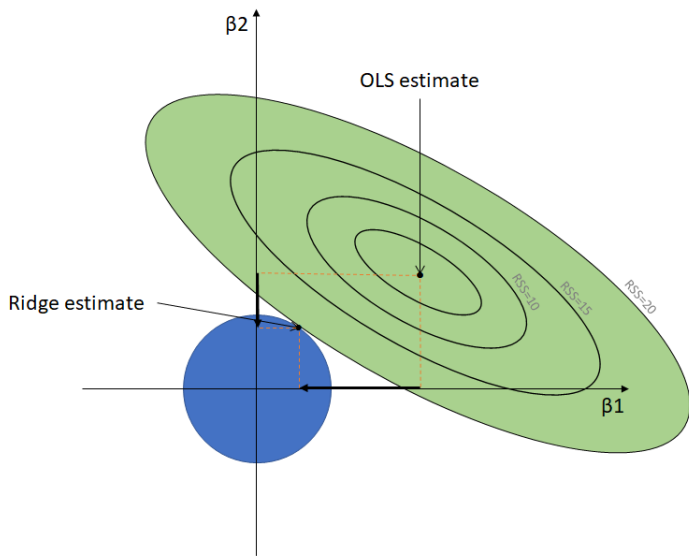


Figure: $\min_{b_1, b_2} \text{RSS}(\text{Datos}, b_1, b_2)$ s.a. $(b_1^2 + b_2^2) \leq \tau^2$.

Recapitulemos

- ▶ $\downarrow \tau \rightarrow$ modelo con menos varianza (y más sesgo).
 - ▶ Parámetros *recortados* (shrink).
- ▶ Trabajamos con features estandarizados.
- ▶ De manera compacta (y equivalente*) definimos:

$$\hat{\boldsymbol{\beta}}_{\lambda}^{\text{ridge}} = \underset{b_1, \dots, b_p}{\operatorname{argmin}} \quad \underbrace{\text{RSS}(b_1, \dots, b_p)}_{\text{Bias}} + \lambda \underbrace{\|\mathbf{b}\|_2^2}_{\text{Variance}}$$

donde $\mathbf{b} = (b_1, \dots, b_p)$ y $\|\mathbf{b}\|_2 = \sqrt{\sum_{i=1}^p b_i^2}$ (norma ℓ_2).

- ▶ $\|\mathbf{b}\|_2$ mide la complejidad del modelo.
- ▶ $\downarrow \tau$ equivale a $\uparrow \lambda$.
- ▶ λ distribuye peso entre términos (**validación–cruzada**).

*(Back Up slide)

- Para todo τ que tiene asociado:

$$\hat{\beta}_{\tau}^{\text{ridge}} = \underset{b_1, \dots, b_p}{\operatorname{argmin}} \operatorname{RSS}(b_1, \dots, b_p), \text{ s.a. } \|\mathbf{b}\|_2^2 \leq \tau^2,$$

- Existe un solo valor de λ que verifica que $\hat{\beta}_{\tau}^{\text{ridge}} = \hat{\beta}_{\lambda}^{\text{ridge}}$, con:

$$\hat{\beta}_{\lambda}^{\text{ridge}} = \underset{b_1, \dots, b_p}{\operatorname{argmin}} \operatorname{RSS}(b_1, \dots, b_p) + \lambda \|\mathbf{b}\|_2^2.$$

- Dualidad de Lagrange

Algunos detalles sobre el estimador Ridge

$$\min_{\mathbf{b}} \underbrace{(\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{y} - \mathbf{X}\mathbf{b})}_{\text{RSS}(\text{Datos}, \mathbf{b})} + \lambda \underbrace{\mathbf{b}^T \mathbf{b}}_{\|\mathbf{b}\|_2^2}$$

- ▶ Estimador Ridge:

$$\hat{\beta}_{\lambda}^{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1} \mathbf{X}^T \mathbf{y}$$

- ▶ Programación cuadrática (escalable con n y p).
- ▶ En el caso de diseños ortonormales: $\hat{\beta}_{\lambda}^{\text{ridge}} = \hat{\beta}^{\text{ols}} / (1 + \lambda)$.
 - ▶ Por tanto: $\hat{\beta}_{\lambda}^{\text{ridge}} \xrightarrow{\lambda \rightarrow 0} \hat{\beta}^{\text{ols}}$ y $\hat{\beta}_{\lambda}^{\text{ridge}} \xrightarrow{\lambda \rightarrow \infty} \mathbf{0}$.
- ▶ Cuando $\lambda \rightarrow \infty$, la solución Ridge restringe todos los parámetros estimados a cero (salvo la ordenada al origen).

Implementación en R

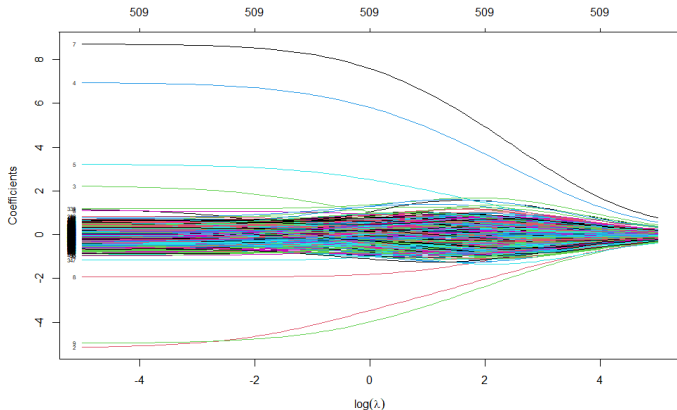


Figure: Simulaciones. Las curvas representan los valores estimaciones de los parámetros del modelo Ridge para distintos valores de $\log(\lambda)$.

Comentarios

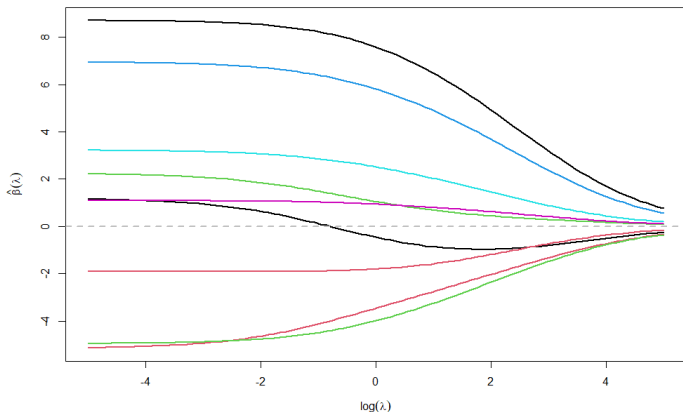


Figure: Sólo las variables relevantes en el modelo. Algunos parámetros estimados cambian de signo al variar λ (inconsistencia).

- ▶ Con $\lambda \gg 0$ algunos $\hat{\beta} \approx 0$ (no hacemos selección de modelo).
- ▶ Los parámetros estimados cambian de signo :(

Lasso (Least Absolute Shrinkage and Selection Operator)

- Dada $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ y $\tau \geq 0$, Lasso resuelve:

$$\begin{aligned} (\hat{\beta}_1, \dots, \hat{\beta}_p)_{\tau}^{\text{lasso}} = \operatorname{argmin}_{b_1, \dots, b_p} \quad & \sum_{i=1}^n \left(y_i - (b_1 x_{1,i} + \dots + b_p x_{p,i}) \right)^2 \\ \text{subject to} \quad & \sum_{j=1}^p |b_j| \leq \tau \end{aligned}$$

- Variables *estandarizadas*.

De manera equivalente y compacta:

$$\hat{\beta}_{\lambda}^{\text{lasso}} = \operatorname{argmin}_{b_1, \dots, b_p} \text{RSS}(b_1, \dots, b_p) + \lambda \|\mathbf{b}\|_1$$

donde $\mathbf{b} = (b_1, \dots, b_p)$ y $\|\mathbf{b}\|_1 = \sum_{i=1}^p |b_i|$ (norma ℓ_1).

Geometría de los problemas de optimización

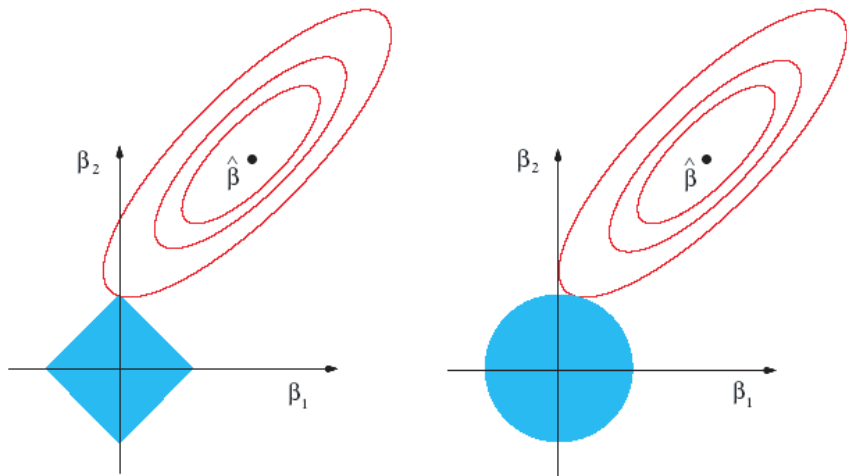


Figure: Lasso ($|\beta_1| + |\beta_2| \leq \tau$) vs Ridge ($\beta_1^2 + \beta_2^2 \leq \tau^2$). Cuando $p \gg 0$, la soluciones de esquina con Lasso son más frecuentes. ((run-giff))

- ▶ Problema de optimización convexo escalable en n y p .

- ▶ Relación con OLS en diseños ortonormales:

$$\hat{\beta}_{\lambda,j}^{\text{lasso}} = \text{sign}(\hat{\beta}_j^{\text{ols}}) \max(|\hat{\beta}_j^{\text{ols}}| - \lambda, 0) \text{ para } j = 1, \dots, p.$$

- ▶ Cuando $\lambda \rightarrow 0$, Lasso se aproxima a OLS.

- ▶ Lasso hace **selección de modelos**: Cuando $|\hat{\beta}_j^{\text{ols}}| \leq \lambda$ desaparece el feature/covariable j -ésimo del modelo.

- ▶ A medida que $\lambda \rightarrow \infty$, la solución de Lasso consiste en descartar todas las features (solo ordenada al origen).

- ▶ Coherencia: Lasso provee estimaciones que no cambian de signo ante cambios en λ (Ridge no cumple siempre esta prop).

Conexiones con los modelos Bayesianos (simplificado)

- ▶ T. Park and G. Casella: *The Bayesian Lasso*. JASA, 2012.

$$\begin{aligned} \mathbf{y}|\mathbf{X}\boldsymbol{\beta}, \lambda, \sigma^2 &\sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n) \\ \boldsymbol{\beta}|\lambda, \sigma^2 &\sim \prod_{j=1}^p \frac{\lambda}{2\sigma} e^{-\frac{\lambda}{\sigma}|\beta_j|} \text{ (iid Laplace).} \end{aligned}$$

- ▶ Llamamos $\pi(\boldsymbol{\beta}|\pi(\boldsymbol{\beta}|\text{Datos}, \sigma, \lambda))$ a la posterior del modelo:

$$-\ln(\pi(\boldsymbol{\beta}|\text{Datos}, \sigma, \lambda)) \propto \sum_{i=1}^n \frac{1}{2\sigma^2} \left(y_i - (\beta_1 X_{1i} + \dots + \beta_p X_{pi}) \right)^2 + \frac{\lambda}{\sigma} \|\boldsymbol{\beta}\|_1.$$

- ▶ $\hat{\boldsymbol{\beta}}_{\frac{\lambda}{\sigma}}^{\text{lasso}}$ es la moda a posteriori del modelo Bayesiano.
- ▶ Selección de modelo e inferencia vía $\pi(\boldsymbol{\beta}|\text{Datos}, \sigma, \lambda)$.
 - ▶ Ver § 3 de Park and G. Casella.

Lasso vs Ridge: Ejercicio realizado en clase con R

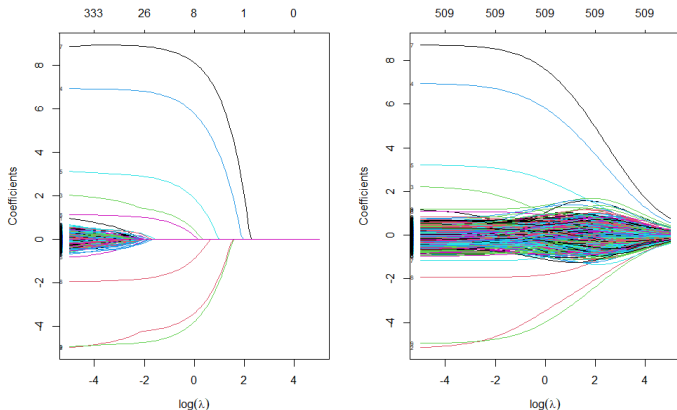


Figure: Comparativa Lasso (izquierda) vs Ridge (derecha).

- ▶ Lasso selecciona variables (descarta lo irrelevante).
- ▶ Lasso evita que los parámetros estimados alternen signo.

Comparativa

- ▶ LASSO \succ RIDGE: Muchas pendientes son cero.
- ▶ RIDGE \succ LASSO: Covariables altamente correladas.

Simulación en R^2 :

| Método | α | λ^* | Parámetros | VC Train (se) | Test |
|--------|----------|-------------|------------|---------------|-------------|
| MCO | 0 | 0 | 509 | 0.62 (–) | 7.59 |
| RIDGE | 0 | 0.12 | 509 | 14.67 (0.74) | 7.44 |
| LASSO | 1 | 0.13 | 17 | 2.71 (0.09) | 2.87 |

Table: Estimaciones del ECM (SE entre paréntesis)

- ▶ Con lasso estimamos modelos parsimónicos.
- ▶ En escenarios raros, lasso tiene propiedades asintóticas (ver simulación en R para discutir *sparsistencia*).

²El escenario de sparcidad hace a lasso el método óptimo.

Generalización del método

$$\hat{\beta}_{\lambda,q} = \operatorname{argmin}_{b_1, \dots, b_p} \sum_{i=1}^n \left(y_i - (b_1 X_{1,i} + \dots + b_p X_{p,i}) \right)^2 + \lambda \|\mathbf{b}\|_q^2,$$

donde $\mathbf{b} = (b_1, \dots, b_p)$ y $\|\mathbf{b}\|_q = (\sum_{i=1}^p |b_i|^q)^{1/q}$ (norma ℓ_q).

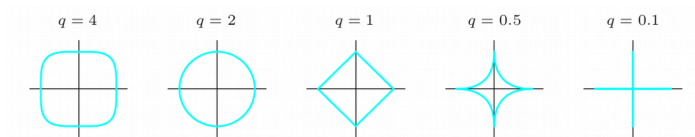


FIGURE 3.12. Contours of constant value of $\sum_j |\beta_j|^q$ for given values of q .

- ▶ Lasso como una relajación convexa cuando $q < 1$.
- ▶ No se gana demasiado aprendiendo “el mejor q ”.
- ▶ Redes elásticas.

Elastic Nets (Redes Elásticas)

$$\hat{\beta}_{\lambda,\alpha}^{\text{EN}} = \underset{b_1,\dots,b_p}{\operatorname{argmin}} \sum_{i=1}^n \left(y_i - (b_1 X_{1,i} + \dots + b_p X_{p,i}) \right)^2 + \lambda \left((1-\alpha) \|\mathbf{b}\|_2^2 + \alpha \|\mathbf{b}\|_1 \right)$$

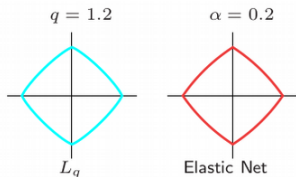


FIGURE 3.13. Contours of constant value of $\sum_j |\beta_j|^q$ for $q = 1.2$ (left plot), and the elastic-net penalty $\sum_j (\alpha \beta_j^2 + (1-\alpha) |\beta_j|)$ for $\alpha = 0.2$ (right plot). Although visually very similar, the elastic-net has sharp (non-differentiable) corners, while the $q = 1.2$ penalty does not.

- ▶ Ridge y Lasso en un solo método.
- ▶ Elastic net selecciona variables.
- ▶ Aprendemos (λ, α) por validación cruzada.

Regularización con glmnet

```
install.packages('glmnet')  
library(glmnet)
```

```
glmnet (x ,y , alpha = 0 , lambda )
```

x,y = en formato MATRIZ.

alpha = 0 (Ridge).

alpha = 1 (LASSO).

alpha en (0,1) (ENet).

lambda = penalización del regularizador.

Output:

Coeficientes estimados para cada valor de lambda.

Preguntas / implementación en R

- ▶ Analiza los resultados del código Simulación Lasso (sparsistencia). Investiga por tu cuenta a que hace referencia esta propiedad de los modelos estimados mediante lasso denominada **sparsistencia**.
- ▶ Siguiendo el ejercicio planteado en clase para ilustrar Ridge y Lasso, completa la Tabla en la slide 18 utilizando ENets.
- ▶ Utiliza los datos `wines` para implementar Elastic Nets (tendrás que construir un pequeño bucle para poder optimizar el valor del hiperparámetro α).
- ▶ Conceptuales: Resuelve los puntos 2,3 y 4 planteados en § 6.8 de ISLR (pp. 259–a–261).

Agenda

Regularización de modelos lineales y aditivos

Apéndice

Propiedades estadísticas de lasso (simplificado)

- ▶ Modelo de regresión: $\mathbf{y} = \mathbf{X}\beta_0 + \varepsilon$, con $\beta_0 \in \mathbb{R}^p$.
 - ▶ Active set: $\mathcal{S} = \{i | \beta_i \neq 0 \text{ con } \beta_i \in \beta_0\}$.
 - ▶ Modelo raro: $|\mathcal{S}| \ll p$ (formalmente $|\mathcal{S}| = o(\sqrt{n/\log(p)})$).
- ▶ A medida que $n \rightarrow \infty$ con $p/n \rightarrow 0$ (y bajo condiciones adecuadas sobre como elegir asintóticamente el valor de λ):
 - ▶ Consistencia: $\hat{\beta}_i^{\text{lasso}} \rightarrow_P \beta_i$.
 - ▶ Identificamos el modelo correcto: $\hat{\mathcal{S}}_n \rightarrow_P \mathcal{S}$.
 - ▶ El estimador lasso $\hat{\beta}^{\text{lasso}}$ es “esparsistente”.
- ▶ La norma $\|\mathbf{X}(\beta_0 - \hat{\beta}^{\text{lasso}})\|_2^2/n$ asociada al estimador lasso decrece a una tasa de al menos \sqrt{n} (cond. gales. en \mathbf{X}).

$$\mathbf{y}_{1 \times n} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{1 \times n}.$$

- ▶ Consideremos \mathbf{X} fija y por lo tanto $\text{Cov}(\mathbf{y}) = \sigma^2 \mathbf{I}_{n \times n}$.
- ▶ El ECM de predicción asociado a $\hat{\boldsymbol{\beta}}^{\text{ols}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ es:

$$\frac{1}{n} E \left\{ \underbrace{\|\mathbf{X} \hat{\boldsymbol{\beta}}^{\text{ols}} - \mathbf{y}\|^2}_{\sum_{i=1}^n (\hat{y}_i - y_i)^2} \right\} = \frac{1}{n} E \left\{ \underbrace{\|\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} - \mathbf{X} \boldsymbol{\beta} - \boldsymbol{\varepsilon}\|^2}_{P_{\mathbf{X}}} \right\}.$$

- ▶ Recordá que $\|P_{\mathbf{X}}\| = \text{tr}(P_{\mathbf{X}}^T P_{\mathbf{X}}) = \text{tr}(P_{\mathbf{X}}) = p$, luego:

$$\frac{1}{n} E \left\{ \|\mathbf{X} \hat{\boldsymbol{\beta}}^{\text{ols}} - \mathbf{y}\|^2 \right\} = \frac{\sigma^2}{n} (\text{tr}(P_{\mathbf{X}}) + \text{tr}(\mathbf{I})) = \frac{\sigma^2}{n} p + \sigma^2.$$

- ▶ Con n fijo, el ECM de predicción crece de forma lineal con p .