

Reporte 2 (Anexos)

Taller Multidisciplinario

Joshué Helí Ricalde Guerrero

2 de julio de 2020

1. Regresión Penalizada

Recordando que en el modelo lineal

$$\underset{n \times 1}{Y} = \underset{n \times p}{X} \underset{p \times 1}{\beta} + \underset{n \times 1}{e},$$

con $\mathbb{E}[e] = 0$ y $\text{Var}(e) = \sigma^2 I$, el *estimador de mínimos cuadrados* cumple que

$$\hat{Y} = X \hat{\beta} = \underset{\text{col}(X)}{\text{proy}} Y.$$

Explícitamente,

$$\hat{\beta} = \underbrace{(X^T X)^{-1}}_{\substack{\text{Inversa} \\ \text{generalizada}}} X^T Y.$$

De ser invertible, la inversa generalizada es única y $(X^T X)^{-1}$.

Algunas de las dificultades que se tienen en éste método son las siguientes:

1. Si $(X^T X)$ no es invertible, $(X^T X)^{-1}$ no es única y $\hat{\beta}$ no está definida unívocamente.
2. Si $(X^T X)$ es invertible pero mal condicionada, el error cuadrático medio es grande:

$$\text{tr}(\text{Var}(\hat{\beta})) = \text{tr}(\sigma^2 (X^T X)^{-1}) = \sigma^2 \sum_{i=1}^p \frac{1}{\lambda_i},$$

con λ_i 's autovalores de $(X^T X)$; luego,

$$\begin{aligned} \text{ecm}(\hat{\beta}) &:= \mathbb{E}[\|\hat{\beta} - \beta\|^2] \\ &= \|\mathbb{E}[\hat{\beta}] - \beta\|^2 + \text{tr}(\text{Var}(\hat{\beta})) \\ &= 0 + \text{tr}(\text{Var}(\hat{\beta})), \end{aligned}$$

donde la última igualdad se da gracias a que el estimador es insesgado ($\mathbb{E}[\hat{\beta}] = \beta$).

Para solucionar estas problemáticas, se sigue la idea de “penalizar” aquellos estimadores que se alejan del origen.

Definición 1.1. Para $\lambda > 0$, se define el estimador *Ridge* o *regularizado* como

$$\tilde{\beta}(\lambda) := \min_{\beta \in \mathbb{R}^p} \left\{ \|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right\},$$

con $\|\cdot\|_2$ la norma euclidiana.

Definición 1.2. Se define el estimador *Lasso* $\hat{\beta}(\lambda)$ del vector de coeficientes $\beta \in \mathbb{R}^p$ como cualquier solución del problema de optimización

$$\hat{\beta}(\lambda) := \min_{\beta \in \mathbb{R}^p} \left\{ \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\},$$

para $\lambda \geq 0$ fijo.

Definición 1.3. Se define el estimador de *Red elástica* al vector de coeficientes $\beta \in \mathbb{R}^p$ que soluciona el problema de optimización

$$\min_{\beta \in \mathbb{R}^p} \left\{ \|Y - X\beta\|_2^2 + \lambda \left[\alpha \|\beta\|_1 + (1 - \alpha) \|\beta\|_2^2 \right] \right\},$$

para $\lambda \geq 0$ y $\alpha \in [0, 1]$ fijos.

2. Validación Cruzada

Como se ha visto hasta ahora, para regresión se ha tomado $\lambda > 0$ definido por el investigador. Sin embargo, queda la cuestión de ¿cómo escoger este λ ? En general, este problema no sólo se reduce a modelos de regresión penalizada, sino a otros contextos que hacen uso de parámetros de “aprendizaje” o *hiperparámetros*, y comúnmente es conocido *Determinación de (hiper)parámetros* ó *Selección de modelos*.

Aunque hay varias metodologías para abordar el tema, por lo pronto únicamente se revisa *Validación cruzada* (Cross-Validation o CV), aunque su uso puede llevarse en paralelo a desigualdades de concentración. También se señala que se revisa CV en general, y no GCV (Generalized Cross-Validation), que aunque tiene “aplicabilidad en Ridge, puede ser limitante.”

La idea principal de CV es básicamente tomar la muestra de datos y segmentarla, de manera que se utilice una parte de la muestra para encontrar el estimador $\hat{\beta}(\lambda)$, y el resto de la muestra para evaluar que tan bueno resultó el estimador. Como puede verse, una de las ventajas de CV es que no se requiere de cálculo teórico adicional; en cambio, la desventaja es que requiere de más cálculo computacional.

Aquí aparece la pregunta de “¿porqué es necesario tomar otros modos de selección?” La respuesta: porque $\|Y - X\hat{\beta}(\lambda)\|$ no es bueno para estimar λ , pues simplemente se elegiría $\lambda = 0$.

2.1. Half-Cross-Validation

Consiste en dividir la muestra original M en 2, M_1 , M_2 , $M_1 \cap M_2 = \emptyset$, $M_1 \cup M_2 = M$, una para estimación (a veces conocida también como *entrenamiento*), y otra para validación.

1. Con M_1 obtener $\hat{\beta}^{M_1}(\lambda)$.
2. Con M_2 obtener

$$\hat{C}(\lambda) := \sum_{i \in M_2} (y_i - x_i^T \hat{\beta}^{M_1}(\lambda)).$$

3. Seleccionar $\hat{\lambda}$ tal que

$$\hat{\lambda} = \arg \min_{\lambda \geq 0} \hat{C}(\lambda). \quad (2.1)$$

2.2. Ordinary Cross-Validation

También conocida como *Leave-One-Out Cross-Validation*, consiste básicamente en dividir la muestra original M , de n datos, en n diferentes sub-muestras en vez de sólo dos.

1. Para todo $i = 1, \dots, n$, obtener $M_{(i)} = M \setminus \{(x_i, y_i)\}$.
2. Con $M_{(i)}$ estimar $\hat{\beta}^{(i)}(\lambda)$.
3. Obtener el i -ésimo predictor

$$\hat{y}_i = \sum_{j=1}^p x_{ij} \hat{\beta}_j^{(i)}(\lambda),$$

obtener la función objetivo (a veces llamada *oráculo*)

$$\begin{aligned} \hat{C}(\lambda) &:= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \\ &= \frac{1}{n} \sum_{i=1}^n \left(y_i - \sum_{j=1}^p x_{ij} \hat{\beta}_j^{(i)}(\lambda) \right). \end{aligned}$$

4. Resolver (2.1).

Una de las ventajas de este tipo de CV es que se tiene poco sesgo; sin embargo, por el mismo motivo se tiene mucha varianza. Además, el cálculo numérico es pesado, pues se requiere estimar n parámetros distintos.

2.3. K -fold CV

Se sigue la misma idea que se ha estado llevando hasta ahora, particionando la muestra original M en K partes. ¿Cómo determinar estos pedazos? Al azar, ya que puede haber factores relacionados (por ejemplo, el tiempo, en el sentido que “las muestras tomadas al comienzo no son las mismas que las tomadas al final”).

1. Se especifica K , y con ello se particiona la muestra M en K partes distintas M_k , $k = 1, \dots, K$.
2. Con M_k estimar $\hat{\beta}^{(k)}(\lambda)$.
3. Para el i -ésimo dato (que está en la k -ésima partición), se calcula el predictor

$$\hat{y}_i = \sum_{j=1}^p x_{ij} \hat{\beta}_j^{(k)}(\lambda).$$

4. Sumando sobre todos los índices de todas las particiones, se obtiene la función oráculo

$$\begin{aligned} \hat{C}(\lambda) &:= \frac{1}{K} \sum_{k=1}^K \left[\frac{1}{\#M_k} \sum_{i \in M_k} (y_i - \hat{y}_i) \right] \\ &= \frac{1}{K} \sum_{k=1}^K \left[\frac{1}{\#M_k} \sum_{i \in M_k} \left(y_i - \sum_{j=1}^p x_{ij} \hat{\beta}_j^{(k)}(\lambda) \right) \right], \end{aligned}$$

donde el factor $\frac{1}{K}$ es el que contribuye a la disminución de la varianza.

5. Resolver (2.1).

Comentario: Otro método popular (para selección de modelos) es Bootstrap. Sin embargo, éste no se aborda porque *Lasso* (y en general penalización l_p) de por sí tiene mucha complejidad computacional.

3. Modelo de Mezclas Dirichlet-Multinomial (DMM)

3.1. Muestreo Multinomial

Sea $X = (X_{ij})_{N \times S}$ una matriz de abundancias, donde X_{ij} es la abundancia observada de la taxa j en el individuo i , y se denota por X_i a la fila i -ésima de X , $1 \leq i \leq N$. Se asume que cada uno de éstos renglones es la realización de una variable aleatoria multinomial de parámetro $p_i = (p_{i1}, \dots, p_{iS})$. Entonces, la verosimilitud de la muestra está dada por

$$L(X|p_1, \dots, p_N) = \prod_{j=1}^S L_i(X_i|p_i),$$

donde

$$L_i(X_i|p_i) = J_i! \prod_{j=1}^S \frac{p_{ij}^{X_{ij}}}{X_{ij}}, \quad J_i = \sum_{j=1}^S X_{ij}, \quad i = 1, \dots, N.$$

Puesto que se trata de un enfoque Bayesiano, se supondrá que los parámetros mismos p_i son variables aleatorias.

3.2. Mezclas Dirichlet como distribuciones previas

Recordando que una variable aleatoria Dirichlet tiene como realizaciones vectores con entradas no negativas, cuya suma es igual a 1, se puede considerar a ésta como una distribución de probabilidad sobre distribuciones de probabilidad:

$$Dir(p_i|\bar{\alpha} = \theta \bar{m}) = \Gamma(\theta) \prod_{j=1}^S \frac{p_{ij}^{\theta m_j - 1}}{\Gamma(\theta m_j)} \delta\left(\sum_{j=1}^S p_{ij} - 1\right),$$

donde Γ denota la función Gamma, $\bar{\alpha} = (\alpha_1, \dots, \alpha_S)$, $\alpha_j > 0$, $\theta = \sum_{j=1}^S \alpha_j$, y $m = (m_1, \dots, m_S)$ tal que $\sum_{j=1}^S m_j = 1$.

Ahora, supóngase que se trabajan con K meta comunidades; entonces, se tienen $\bar{\alpha}_k$ parámetros de distribuciones Dirichlet, cada una con un peso π_k , $k = 1, \dots, K$. Para cada observación i se tiene un vector $z_i = (z_{ik})$ con entradas binarias, que indica la metacomunidad a la que pertenece (es decir, es un vector de ceros con un sólo 1 en alguna entrada). Las distribuciones previas de éste vector no son más que la mezcla de los pesos,

$$P(z_i) = \prod_{k=1}^K \pi_k^{z_{ik}};$$

en consecuencia, la distribución competa de la mezcla es

$$P(p_i|Q) = \sum_{k=1}^K Dir(p_i|\bar{\alpha}_k) \pi_k,$$

donde $Q = (K, \bar{\alpha}_1, \dots, \bar{\alpha}_K, \pi_1, \dots, \pi_K)$ es el vector de hiperparámetros.

3.3. Distribución posterior de parámetros multinomiales

Añadiendo la información de la muestra (la verosimilitud multinomial) a la previa, se obtiene directamente que la posterior de los parámetros p_i es

$$P(p_i|X_i) = \frac{\sum_{k=1}^K L_i(X_i|p_i) Dir(p_i|\bar{\alpha}_k) \pi_k}{\sum_{k=1}^K Dir(p_i|\bar{\alpha}_k) \pi_k}.$$

Pero, de las propiedades de ambas distribuciones, tras hacer las manipulaciones necesarias, se tiene que

$$P(p_i|X_i) = \sum_{k=1}^K Dir(p_i|\bar{\alpha}_k + X_i) P(z_{ik} = 1|X_i Q).$$

3.4. Comparación de Modelos por aproximación de Laplace

Sea H_K el modelo DMM con K componentes (metacomunidades). Puesto que se está trabajando en un contexto Bayesiano, denótese por $p(\mathcal{H}_K|X)$ la probabilidad de observar el modelo dada la matriz de datos (de abundancia) X . Del teorema de Bayes,

$$p(\mathbb{H}_K) \propto p(\mathcal{H}_K)p(X_K|\mathcal{H}_K),$$

donde $p(\mathcal{H}_K)$ es la probabilidad previa por un modelo de K clusters y $p(X|\mathcal{H}_K)$ es la evidencia a favor de tal modelo. Dicho de otra manera, la primera es la preferencia del investigador por el modelo dado, mientras que la segunda es la preferencia de los mismos datos por el mismo modelo. Ésta a su vez está dada por

$$p(X|\mathcal{H}_K) = \int p(X|Q, \mathcal{H}_K) p(Q|\mathcal{H}_K) dQ.$$

Entonces, el criterio que se busca para seleccionar el número de componentes adecuado es *escoger el K máximo a partir del cual los incrementos en la preferencia $p(X|\mathcal{H}_K)$ sean despreciables*; i.e., a partir del cual incluir más clusters aporte muy poco. Lamentablemente, la integral no puede ser resuelta analíticamente, pero sí puede ser estimada usando la aproximación de Laplace

$$\log p(X|\mathcal{H}_K) \approx \log p(X|\hat{Q}, \mathcal{H}_K) + \log p(\hat{Q}|\mathcal{H}_K) + \frac{M}{2} \log 2\pi - \frac{1}{2} |H|,$$

donde M es el número de parámetros en Q , \hat{Q} son los parámetros que maximizan la distribución posterior, y H es la matriz Hessiana del negativo de la posterior en \hat{Q} :

$$H = -\nabla^2 \log p(\hat{Q}|X).$$

Tomando el negativo de la \log -aproximación, el criterio resulta comparable con el de Akaike (AIC) y Bayes (BIC).