

Tarea 2

Instrucciones de entrega:

- Entregue un informe en formato **pdf** con sus resultados, pseudo código (si corresponde) y conclusiones. Además incluya su código en R o Python.
- Comente todas sus soluciones, y utilice herramientas gráficas, por ejemplo use el paquete **rgl** (u otro), si es que estima que esto le ayudará a presentar de mejor manera sus resultados.
- **Fecha de entrega:** 20 de Septiembre

Ejercicio 1 (Regresión ridge). Considere datos $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, donde para cada individuo indexado en $i = 1, \dots, n$, tenemos una variable respuesta y_i , y un vector de covariables $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$ de dimensión d . Además, considere el modelo lineal que relaciona la respuesta y_i con el vector de covariables \mathbf{x}_i :

$$\begin{aligned} y_i &= \beta_0 + \beta_1 \cdot x_{i1} + \dots + \beta_d \cdot x_{id} + \epsilon_i \\ &= \beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta} + \epsilon_i, \end{aligned}$$

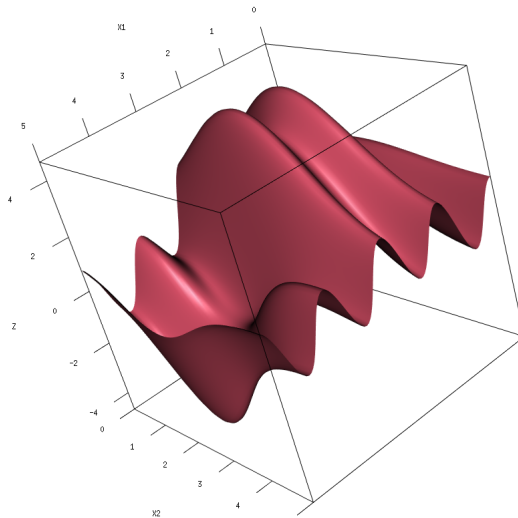
donde $\epsilon_1, \dots, \epsilon_n$ son ruidos independientes con media 0, y varianza constante.

1. Sea $\boldsymbol{\beta} = (\beta_1, \dots, \beta_d)$, y considere la función de pérdida dada por:

$$L_\lambda(\beta_0, \boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + \lambda \|\boldsymbol{\beta}\|^2$$

Encuentre explícitamente (i.e., resuelva el problema de optimización derivando e igualando a cero) estimadores $\hat{\beta}_0$ y $\hat{\boldsymbol{\beta}}$ que minimizan la función de pérdida anterior. (Note que este es el problema estándar resuelto en Ridge Regression, pero donde no se penaliza al intercepto). Discuta que ocurre cuando i) $\lambda = 0$, ii) $\lambda \rightarrow \infty$.

Ejercicio 2 (Kernel Ridge Regression). En este ejercicio trabajaremos con datos simulados que se encuentran en el archivo `sim.txt`. Este archivo contiene información sobre 3 variables: y (variable respuesta), y $X1, X2$ (covariables 2-dimensionales asociadas). La función que generó los datos luce de esta forma:



1. Implemente un modelo de Kernel Ridge Regression basado en la función de pérdida:

$$\mathcal{L}_\lambda(f) = \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda \|f\|_{\mathcal{H}}^2, \quad f \in \mathcal{H}$$

donde \mathcal{H} es el RKHS asociado a la función de kernel dada por el kernel squared-exponential $K : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$:

$$K(\mathbf{x}, \mathbf{x}') = \exp \left\{ -\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\ell^2} \right\}, \quad \text{donde } \mathbf{x}, \mathbf{x}' \in \mathbb{R}^2.$$

En su implementación elija $\ell = 0,5$ y $\lambda = 1$. Evalúe la estimación del modelo en términos del R cuadrado, y del error cuadrático medio predictivo. Para esto último considere $k = 10$ cross-validation. Muestre sus resultados ¿Qué opina de los resultados obtenidos?

2. Implemente el mismo modelo anterior, pero esta vez encuentre parámetros (ℓ, λ) óptimos en el sentido de que minimicen el error cuadrático medio predictivo. Utilizando estos parámetros óptimos vuelva a evaluar su estimación en términos del R cuadrado, y del error cuadrático medio predictivo. Repita este procedimiento usando $k = 10$ cross-validation y compare con el ítem anterior. Comente sobre sus resultados.

Hint 1: Para encontrar los valores óptimos de los parámetros cree una grilla de valores para (ℓ, λ) .

Hint 2: Para una comparación más justa, use los mismos conjuntos para $k = 10$ cross-validation en este, y en el ítem anterior.

3. Note que en el ítem anterior, al usar $k=10$ cross-validation, se obtiene una colección de 10 parámetros óptimos $(\ell_1, \lambda_1), \dots, (\ell_{10}, \lambda_{10})$ (uno por cada set de testeo). ¿Qué puede decir sobre la distribución de estos parámetros? (utilice algún gráfico para mostrar la distribución) ¿Son muy distintos de los ocupados en el ítem 1?
4. Una generalización del squared exponential kernel $K : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ esta dada por la siguiente ecuación:

$$\begin{aligned} K(\mathbf{x}, \mathbf{x}') &= \exp \left\{ -(\mathbf{x} - \mathbf{x}')^\top \Sigma^{-1} (\mathbf{x} - \mathbf{x}') \right\} \\ &= \exp \left\{ -\left(\frac{x_1 - x'_1}{\ell_1} \right)^2 - \left(\frac{x_2 - x'_2}{\ell_2} \right)^2 \right\} \end{aligned}$$

donde $\mathbf{x} = (x_1, x_2)$, $\mathbf{x}' = (x'_1, x'_2)$ y $\Sigma = \begin{pmatrix} \ell_1^2 & 0 \\ 0 & \ell_2^2 \end{pmatrix}$, (note que esta generalización permite 2 parámetros de length-scale distintos). Repita el ítem 2, pero esta vez use la generalización del squared-exponential kernel y busque parámetros óptimos $(\ell_1, \ell_2, \lambda)$ (en esta ocasión necesitará buscar en una grilla 3 dimensional). Comente los resultados obtenidos. Según los resultados que encuentre, comente porque utilizar esta generalización podría ser ventajoso, o desventajoso.