

REGRESIÓN - CLASIFICACIÓN

Ejercicio 1

Justificar la expresión del intervalo de confianza para la recta de regresión.

Ejercicio 2

Considerar la siguiente tabla de datos sobre el rendimiento de cultivos de papas y el registro de lluvias acumuladas en el período de duración del cultivo:

Datos de entrenamiento		Datos de validación	
$x =$ Lluvia (mm)	$y =$ Rendimiento (ton/ha)	$x =$ Lluvia (mm)	$y =$ Rendimiento (ton/ha)
206	29	213	30
188	25	80	16
219	31	391	25
372	25	250	26
345	29	57	9
231	30	303	28
203	26	263	28
170	23	157	25
55	12	72	13
91	15	157	23
292	28	188	26
141	24	216	25
129	23	362	28
170	22	283	33
324	30	308	30

1. Correr una regresión lineal para predecir el rendimiento y en función de la lluvia x .
2. Comparar el MSE en entrenamiento, validación y CV.
3. Determinar el grado óptimo en caso de aplicar una regresión polinomial.
4. Hallar el valor de λ óptimo para la regresión polinomial de grado 5 con regularización.

Ejercicio 3

El objetivo del ejercicio es implementar el algoritmo de descenso de gradiente para estimar los parámetros óptimos de la regresión logística a partir de los siguientes datos:

Paciente	Glucosa	Tiene diabetes (Si/No)
A	90	No
B	160	Si
C	100	No
D	200	Si
E	130	Si

1. Mostrar que la función de pérdida de la regresión logística es $L(\beta) = \frac{1}{N} \sum_{i=1}^N \ln(1 + e^{-y_i \beta' x_i})$ donde $\beta = (\beta_0, \beta_1)$ es el parámetro que se busca.
2. Mostrar que el gradiente de L respecto de β está dado por $\frac{\partial L(\beta)}{\partial \beta_j} = \frac{1}{N} \sum_{i=1}^N -\frac{y_i x_{ij} e^{-y_i \beta' x_i}}{1 + e^{-y_i \beta' x_i}}$
3. Implementando la formula del descenso por gradiente $\beta_{k+1} = \beta_k - \alpha \nabla L(\beta_k)$, calcular el óptimo $\hat{\beta}$ para este ejemplo

Ejercicio 4

Generar 100 observaciones a partir de una distribución normal bivariada $\mathcal{N}(\mu_1, \Sigma_1)$ con $\mu_1 = (3, 1)'$ y $\Sigma_1 = I$ (matriz identidad) y etiquetarlas como 1. Generar otras 100 observaciones a partir de una

distribución gaussiana bivariada gaussiana bivariada $\mathcal{N}(\mu_2, \Sigma_2)$ con $\mu_2 = (1, 3)'$ y $\Sigma_2 = I$ y etiquetarlas como 0. En conjunto, estas 200 observaciones constituyen el conjunto de entrenamiento.

1. Escribir un código R para generar este conjunto de datos. Dibujar estos datos utilizando diferentes colores para las dos clases.
2. Suponiendo que las probabilidades a priori son iguales, encontrar el clasificador de Bayes. Calcular el error de entrenamiento.
3. Entrenar un modelo de regresión lineal, utilizando la función $\text{lm}(y \sim x)$, con el conjunto de entrenamiento.
4. Trazar la frontera del clasificador de Bayes y la recta obtenida por la regresión lineal.
5. Generar un conjunto de prueba de 50 observaciones y calcular el error de prueba del clasificador de Bayes y del modelo lineal.

Ejercicio 5

Considere la siguiente tabla.

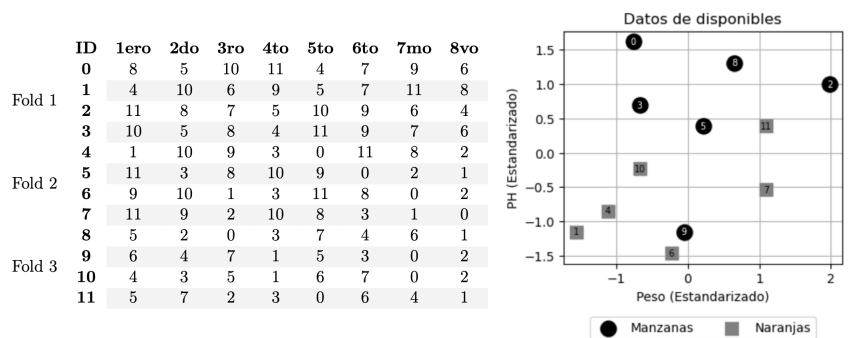
x_1	x_2	x_3	y
0	3	0	Rojo
2	0	0	Rojo
0	1	3	Rojo
0	1	2	Azul
1	0	1	Azul
1	1	1	Rojo

1. Con la distancia euclídea, ¿cuál es la predicción con $k = 1$ y con $k = 3$ para la observación de prueba $(0, 0, 0)$?
2. Si la frontera de decisión de Bayes en este problema es altamente no lineal, entonces ¿esperaríamos que el mejor valor para k fuera grande o pequeño? ¿Por qué?

Ejercicio 6

Se desea implementar el algoritmo de K vecinos más cercanos para clasificar Manzanas y Naranjas en base a su Peso y su PH.

El gráfico a continuación (derecha) muestra los datos disponibles estandarizados. En el lado izquierdo, se presenta el conjunto de datos segmentado en tres folds. Cada fila simboliza una observación, y su ID correspondiente está alineado con el gráfico situado a la derecha. Las columnas, por otro lado, exhiben los IDs de los 8 puntos que son parte de los dos folds que no contienen la observación de la fila en cuestión, todos ellos ordenados por su proximidad a dicha observación. Por ejemplo, en la primera fila se visualizan las observaciones de los folds 2 y 3. Están dispuestas en un orden que va desde la más cercana a la más lejana respecto a la observación con ID=0.



Calcular el error de 3-fold cross-validation para los valores impares de K. ¿Cuál de estos valores de K elegiría?