

Trabajo.1: Programación

Fecha límite de entrega: 27 de Marzo 2017

Valoración máxima: 13 puntos

NORMAS DE DESARROLLO Y ENTREGA DE TRABAJOS

Para este trabajo como para los demás es obligatorio presentar un informe escrito con sus valoraciones y decisiones adoptadas en el desarrollo de cada uno de los apartados. Incluir en el informe los gráficos generados. También deberá incluirse una valoración sobre la calidad de los resultados encontrados. (obligatorio en pdf). **Sin este informe se considera que el trabajo NO ha sido presentado.**

Normas para el desarrollo de los Trabajos: EL INCUMPLIMIENTO DE ESTAS NORMAS SIGNIFICA PERDIDA DE 2 PUNTOS POR CADA INCUMPLIMIENTO.

- El código se debe estructurar en un único script R con distintas funciones o apartados, uno por cada ejercicio/apartado de la práctica.
- Todos los resultados numéricos o gráficas serán mostrados por pantalla, parando la ejecución después de cada apartado. No escribir nada en el disco.
- El path que se use en la lectura de cualquier fichero auxiliar de datos debe ser siempre "datos/nombre_fichero". Es decir, crear un directorio llamado "datos" dentro del directorio donde se desarrolla y se ejecuta la práctica.
- Un código es apto para ser corregido si se puede ejecutar de principio a fin sin errores.
- NO ES VÁLIDO usar opciones en las entradas. Para ello fijar al comienzo los parámetros por defecto que considere que son los óptimos.
- El código debe estar obligatoriamente comentado explicando lo que realizan los distintos apartados y/o bloques.
- Poner puntos de parada para mostrar imágenes o datos por consola.
- Todos los ficheros (*.R *.pdf) se entregan juntos dentro de un único fichero zip, sin ningún directorio que los contenga.
- ENTREGAR SOLO EL CODIGO FUENTE, NUNCA LOS DATOS.
- **Forma de entrega:** Subir el zip al Tablón docente de CCIA.

1. EJERCICIO SOBRE LA COMPLEJIDAD DE H Y EL RUIDO (3 PUNTOS)

En este ejercicio debemos aprender la dificultad que introduce la aparición de ruido en las etiquetas a la hora de elegir clases de funciones simples. Haremos uso de tres funciones R ya programadas:

- *simula_unif*($N, dim, rango$), que calcula una lista de N vectores de dimensión dim . Cada vector contiene dim números aleatorios uniformes en el intervalo $rango$.
- *simula_gaus*($N, dim, sigma$), que calcula una lista de longitud N de vectores de dimensión dim , donde cada posición del vector contiene un número aleatorio extraído de una distribución Gaussiana de media 0 y varianza dada, para cada dimension, por la posición del vector $sigma$.
- *simula_recta*($intervalo$), que simula de forma aleatoria los parámetros, $v = (a, b)$ de una recta, $y = ax + b$, que corta al cuadrado $[-50, 50] \times [-50, 50]$.

1. Dibujar una gráfica con la nube de puntos de salida correspondiente.
 - a) Considere $N = 50$, $dim = 2$, $rango = [-50, +50]$ con *simula_unif*($N, dim, rango$).
 - b) Considere $N = 50$, $dim = 2$ y $sigma = [5, 7]$ con *simula_gaus*($N, dim, sigma$).
2. Con ayuda de la función *simula_unif*() generar una muestra de puntos 2D a los que vamos añadir una etiqueta usando el signo de la función $f(x, y) = y - ax - b$, es decir el signo de la distancia de cada punto a la recta simulada con *simula_recta*().
 - a) Dibujar una gráfica donde los puntos muestren el resultado de su etiqueta, junto con la recta usada para ello. (Observe que todos los puntos están bien clasificados respecto de la recta)
 - b) Modifique de forma aleatoria un 10 % etiquetas positivas y otro 10 % de negativas. Dibuje de nuevo la gráfica anterior. (Ahora hay puntos mal clasificados respecto de la recta)
3. Supongamos ahora que las siguientes funciones definen la frontera de clasificación de los puntos de la muestra en lugar de una recta
 - $f(x, y) = (x - 10)^2 + (y - 20)^2 - 400$
 - $f(x, y) = 0,5(x + 10)^2 + (y - 20)^2 - 400$
 - $f(x, y) = 0,5(x - 10)^2 - (y + 20)^2 - 400$
 - $f(x, y) = y - 20x^2 - 5x + 3$

Visualizar el etiquetado generado en 2b junto con cada una de las gráficas de cada una de las funciones. Comparar las formas de las regiones positivas y negativas de estas nuevas funciones con las obtenidas en el caso de la recta ¿Hemos ganado algo en mejora de clasificación al usar funciones más complejas que la dada por una función lineal ? Explicar el razonamiento.

2. EJERCICIO SOBRE EL ALGORITMO PERCEPTRON (2 PUNTOS)

1. Implementar la función *ajusta_PLA*($datos, label, max_iter, vini$) que calcula el hiperplano solución a un problema de clasificación binaria usando el algoritmo PLA. La entrada

datos es una matriz donde cada item con su etiqueta está representado por una fila de la matriz, *label* el vector de etiquetas (cada etiqueta es un valor $+1$ o -1), *max_iter* es el número máximo de iteraciones permitidas y *vini* el valor inicial del vector. La función devuelve los coeficientes del hiperplano.

2. Ejecutar el algoritmo PLA con los datos simulados en los apartados 2a de la sección.1. Inicializar el algoritmo con: a) el vector cero y, b) con vectores de números aleatorios en $[0, 1]$ (10 veces). Anotar el número medio de iteraciones necesarias en ambos para converger. Valorar el resultado relacionando el punto de inicio con el número de iteraciones.
3. Hacer lo mismo que antes usando ahora los datos del apartado 2b de la sección.1. ¿Observa algún comportamiento diferente? En caso afirmativo diga cual y las razones para que ello ocurra.

3. EJERCICIO SOBRE REGRESIÓN LINEAL (5 PUNTOS)

En la web del curso se encuentran disponibles la descripción

1. Leemos datos:
 - Abra el fichero Zip.info disponible en la web del curso y lea la descripción de la representación numérica de la base de datos de números manuscritos que hay en el fichero Zip.train. Lea el fichero Zip.train dentro de su código y visualice las imágenes (usando paraTrabajo1.R). Seleccione solo las instancias de los números 1 y 5. Guardelas como matrices de tamaño 16×16 .
 - También está disponible el fichero Zip.test que deberemos usar más adelante.
2. De cada matriz de números (imagen) vamos a extraer dos características: a) su valor medio; y b) su grado de simetría vertical.
 - Para calcular el grado de simetría haremos lo siguiente: a) calculamos una nueva imagen invirtiendo el orden de las columnas; b) calculamos la diferencia entre la matriz original y la matriz invertida; c) calculamos la media global de los valores absolutos de la matriz. Conforme más alejado de cero sea el valor más asimétrica será la imagen.
 - Representar en los ejes $\{X = \text{Intensidad Promedio}, Y = \text{Simetría}\}$ las instancias seleccionadas de 1's y 5's.
3. Ajustar un modelo de regresión lineal usando la transformación SVD sobre los datos de (Intensidad promedio, Simetría) y pintar la solución obtenida junto con los datos usados en el ajuste. Las etiquetas serán $\{-1, 1\}$. Valorar la bondad del resultado usando E_{in} y E_{out} (usar Zip.test). (usar `Regress_Lin(datos, label)` como llamada para la función).
4. En este apartado exploramos como se transforman los errores E_{in} y E_{out} cuando aumentamos la complejidad del modelo lineal usado. Ahora hacemos uso de la función `simula_unif(N, 2, size)` que nos devuelve N coordenadas 2D de puntos uniformemente muestreados dentro del cuadrado definido por $[-size, size] \times [-size, size]$
 - EXPERIMENTO-1 (1.5 punto):
 - a) Generar una muestra de entrenamiento de $N = 1000$ puntos en el cuadrado $\mathcal{X} = [-1, 1] \times [-1, 1]$. Pintar el mapa de puntos 2D.

- b) Consideremos la función $f(x_1, x_2) = \text{sign}((x_1 + 0,2)^2 + x_2^2 - 0,6)$ que usaremos para asignar una etiqueta a cada punto de la muestra anterior. Introducimos ruido sobre las etiquetas cambiando aleatoriamente el signo de un 10 % de las mismas. Pintar el mapa de etiquetas final.
- c) Usando como vector de características $(1, x_1, x_2)$ ajustar un modelo de regresión lineal al conjunto de datos generado y estimar los pesos w . Estimar el error de ajuste E_{in} .
- d) Ejecutar todo el experimento definido por (a)-(c) 1000 veces (generamos 1000 muestras diferentes) y
 - Calcular el valor medio de los errores E_{in} de las 1000 muestras.
 - Generar 1000 puntos nuevos por cada iteración y calcular con ellos el valor de E_{out} en dicha iteración. Calcular el valor medio de E_{out} en todas las iteraciones.
- e) Valore que tan bueno considera que es el ajuste con este modelo lineal a la vista de los valores medios obtenidos de E_{in} y E_{out}
- EXPERIMENTO-2 (1.5 punto):
 - a) Ahora vamos a repetir el mismo experimento anterior pero usando características no lineales. Ahora usaremos el siguiente vector de características: $\Phi_2(x) = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$. Ajustar el nuevo modelo de regresión lineal y calcular el nuevo vector de pesos \hat{w} . Calcular el error E_{in}
 - b) Al igual que en el experimento anterior repetir el experimento 1000 veces calculando con cada muestra el error dentro y fuera de la muestra, E_{in} e E_{out} respectivamente. Promediar los valores obtenidos para ambos errores a lo largo de las muestras.
 - c) Valore el resultados de este EXPERIMENTO-2 a la vista de los valores medios de los errores E_{in} y E_{out}
- A la vista de los resultados de los errores promedios E_{in} y E_{out} obtenidos en los dos experimentos ¿Que modelo considera que es el más adecuado? Justifique la decisión.

3.1. BONUS

El BONUS solo se tendrá en cuenta si se ha obtenido al menos el 75 % de los puntos de la parte obligatoria.

1. (3 puntos) En este ejercicio exploramos cómo funciona regresión lineal en problemas de clasificación. Para ello generamos datos usando el mismo procedimiento que en ejercicios anteriores. Suponemos $\mathcal{X} = [-10, 10] \times [-10, 10]$ y elegimos muestras aleatorias uniformes dentro de \mathcal{X} . La función f en cada caso será una recta aleatoria que corta a \mathcal{X} y que asigna etiqueta a cada punto de \mathcal{X} con el valor del signo de f en dicho punto. En cada apartado generamos una muestra y le asignamos etiqueta con la función f generada. En cada ejecución generamos una nueva función f
 - a) Fijar el tamaño de muestra $N = 100$. Usar regresión lineal para encontrar una primera solución g y evaluar E_{in} , (el porcentaje de puntos incorrectamente clasificados). Repetir el experimento 1000 veces y promediar los resultados ¿Qué valor obtiene para E_{in} ?
 - b) Fijar el tamaño de muestra $N = 100$. Usar regresión lineal para encontrar g y evaluar E_{out} . Para ello generar 1000 puntos nuevos y usarlos para estimar el error fuera

de la muestra, E_{out} (porcentaje de puntos mal clasificados). De nuevo, ejecutar el experimento 1000 veces y tomar el promedio. ¿Qué valor obtiene de E_{out} ? Valore el resultado.

- c) Ahora fijamos $N = 10$, ajustamos regresión lineal y usamos el vector de pesos encontrado como un vector inicial de pesos para PLA. Ejecutar PLA hasta que converja a un vector de pesos final que separe completamente la muestra de entrenamiento. Anote el número de iteraciones y repita el experimento 1.000 veces ¿Cual es valor promedio de iteraciones que tarda PLA en converger? (En cada iteración de PLA elija un punto aleatorio del conjunto de mal clasificados). Valore los resultados