

## PROYECTO FINAL:

Fecha límite de entrega: 6 de junio 2018

Valoración: 25 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 gráficos 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.
- ENTREGAR SOLO EL CODIGO FUENTE, NUNCA LOS DATOS.
- **Forma de entrega:** Subir los ficheros .R y .pdf ... según requerido en la definición de la entrega en el Tablón docente de CCIA. Y en un zip todos ficheros en DECSAI.

## 1. AJUSTE DEL MEJOR MODELO

Este ejercicio se centra en el ajustar el mejor predictor (lineal o no-lineal) a un conjunto de datos. Debemos mostrar que los distintos algoritmos proponen soluciones para los datos pero que unas soluciones son mejores que otras para unos datos dados. El criterio que usaremos en la comparación será el error medio cuadrático para regresión, la curva ROC en clasificación binaria y el número de errores en clasificación multiclase. Además de un modelo lineal se deberán presentar resultados con al menos dos modelos de entre los propuestos

Los posibles modelos no-lineales a usar son:

- **Redes Neuronales.** Considerar tres clases de funciones definidas por arquitecturas con 1,2 y 3 capas de unidades ocultas y número de unidades por capa en el rango 0-50. Definir un conjunto de modelos(arquitecturas) y elegir el mejor por validación cruzada. Recordar que a igualdad de  $E_{out}$  siempre es preferible la arquitectura más pequeña.
- **Máquina de Soporte de Vectores (SVM):** usar solo el núcleo RBF-Gaussiano o el polinomial. Encontrar el mejor valor para el parámetro libre hasta una precisión de 2 cifras (enteras o decimales)
- **Boosting:** Para clasificación usar AdaBoost con funciones “stamp”. Para regresión usar árboles como regresores simples.
- **Random Forest:** Usar los valores que por defecto se dan en la teoría y experimentar para obtener el número de árboles adecuado.

Se habrá de buscar el mejor modelo posible para la base de datos seleccionada y se habrá de justificar cada uno de los pasos dados para conseguirlo. Los puntos de discusión señalados en el trabajo.3 deben de servir como guía.

Se usará para ello una de las siguientes BBDD del repositorio de la UCI (<https://archive.ics.uci.edu/ml/>).

**Bases de datos elegibles:**

1. Pen-Based Recognition of Handwritten digits (clasificación)
2. Page Blocks Classification (clasificación)
3. Amazon Commerce reviews set (clasificación)
4. Breast Cancer Wisconsin (Diagnostic) (clasificación)
5. Communities and Crime (regresión)
6. Parkinson Telemonitoring (regresión)
7. Housing (regresión)
8. Cardiotocography (clasificación)
9. Thyroid Disease (clasificación)
10. Occupancy detection (clasificación)
11. Default of Credit Card Clients (clasificación)
12. Internet Advertisements (clasificación)
13. Human Activity Recognition Using Smartphones (clasificación)
14. Image Segmentation (clasificación)
15. Mushroom (clasificación)

16. Student Performance Data Set
17. Tennis Major Tournament Match Statistics Data Set
18. Arcene (clasificación)
19. APS Failure at Scania Trucks Data Set (clasificación)
20. Bank Marketing Data Set (clasificación)

Se aceptará el uso de bases de otras bases de datos si previamente han sido autorizadas por el profesor.