

Profesor: Dr. Oldemar Rodríguez Rojas
Minería de Datos 1
Fecha de Entrega: Jueves 8 de junio - 8am
Instrucciones:

Instrucciones:

- Las tareas deben ser subida la Aula Virtual antes de las 12 media noche. Luego de esta hora pierde 20 puntos y cada día de retraso adicional perderá 20 puntos más.
- Las tareas son estrictamente individuales.
- Tareas idénticas se les asignará cero puntos.
- Todas las tareas tienen el mismo valor en la nota final del curso.
- Cada día de entrega tardía tendrá un rebajo de 20 puntos.

TAREA NÚMERO 12

- **Ejercicio 1:** [30 puntos] Para esta pregunta también usaremos los datos `tumores.csv`.
 1. El objetivo de este ejercicio es calibrar el método de **ADA** para esta Tabla de Datos. Aquí interesa predecir en la variable **tipo**. Usando los paquetes **snow** y **trainR** programe en **paralelo** 5 Validaciones Cruzadas con 10 grupos calibrando el modelo de acuerdo con los tres tipos de algoritmos que permite, **discrete**, **real** y **gentle**. Para medir la calidad de método sume la cantidad de 1's detectados en los diferentes grupos. Luego grafique las 5 iteraciones para los tres algoritmos en el mismo gráfico. ¿Se puede determinar con claridad cuál algoritmo es el mejor? Para generar los modelos predictivos use las siguientes instrucciones:

```
modelo<-train.ada(tipo~.,data=taprendizaje,iter=80,nu=1,type="discrete")
modelo<-train.ada(tipo~.,data=taprendizaje,iter=80,nu=1,type="real")
modelo<-train.ada(tipo~.,data=taprendizaje,iter=80,nu=1,type="gentle")
```
 2. Repita el ejercicio anterior, pero esta vez en lugar de sumar la cantidad de 1's, promedie los errores globales cometidos en los diferentes grupos (folds). Luego grafique las 5 iteraciones para los tres algoritmos en el mismo gráfico. ¿Se puede determinar con claridad cuál algoritmo es el mejor?
 3. Para estar realmente seguros de cuál de los tipos algoritmos es mejor, modifique (una en uno solo código) los códigos de los dos ejercicios anteriores de manera que, en lugar de sumar la cantidad de 1's detectados y de promediar los errores globales, guarde en cada iteración (en listas) las matrices de confusión de cada uno de los 3 algoritmos. Luego al final de los ciclos ejecutados para cada algoritmo calcule una **Matriz de Confusión Promedio** de todas las matrices de confusión guardadas en la lista del respectivo algoritmo; así con estas **Matrices de Confusión Promedio**, mediante gráficos de barras, determine el método que en promedio detecta mayor porcentaje de 1's, determine el método que en promedio detecta mayor porcentaje de 0's y determine el método que en promedio tiene menor error global.

4. ¿Cuál algoritmo usaría con base en la información obtenida en los dos ejercicios anteriores?

■ **Ejercicio 2:** [30 puntos] Para esta pregunta usaremos nuevamente los datos `tumores.csv`.

1. El objetivo de este ejercicio es calibrar el método de `knn` para esta Tabla de Datos. Aquí interesa predecir en la variable `tipo`. Usando los paquetes `snow` y `trainR` programe **en paralelo** 5 Validaciones Cruzadas con 10 grupos calibrando el modelo de acuerdo con todos los tipos de algoritmos que permite `train.knn` en el parámetro `kernel`, estos algoritmos son: `rectangular`, `triangular`, `epanechnikov`, `biweight`, `triweight`, `cos`, `inv`, `gaussian` y `optimal`. Para medir la calidad de método sume la cantidad de 1's detectados en los diferentes grupos. Luego grafique las 5 iteraciones para todos algoritmos en el mismo gráfico. ¿Se puede determinar con claridad cuál algoritmo es el mejor?
2. Repita el ejercicio anterior, pero esta vez en lugar de sumar la cantidad de 1's, promedie los errores globales cometidos en los diferentes grupos (folds). Luego grafique las 5 iteraciones para todos los algoritmos en el mismo gráfico. ¿Se puede determinar con claridad cuál algoritmo es el mejor?
3. Para estar realmente seguros de cuál de los tipos algoritmos es mejor, modifique (una en uno solo código) los códigos de los dos ejercicios anteriores de manera que, en lugar de sumar la cantidad de 1's detectados y de promediar los errores globales, guarde en cada iteración (en listas) las matrices de confusión de cada uno de los 3 algoritmos. Luego al final de los ciclos ejecutados para cada algoritmo calcule una **Matriz de Confusión Promedio** de todas las matrices de confusión guardadas en la lista del respectivo algoritmo; así con estas **Matrices de Confusión Promedio**, mediante gráficos de barras, determine el método que en promedio detecta mayor porcentaje de 1's, determine el método que en promedio detecta mayor porcentaje de 0's y determine el método que en promedio tiene menor error global.
4. ¿Cuál algoritmo usaría con base en la información obtenida en los dos ejercicios anteriores?

■ **Ejercicio 3:** [40 puntos] Esta pregunta también utilizan nuevamente los datos `tumores.csv`.

1. El objetivo de este ejercicio es comparar todos los métodos predictivos vistos en el curso con esta tabla de datos. Aquí interesa predecir en la variable `tipo`, Usando los paquetes `snow` y `trainR` programe **en paralelo** 5 Validaciones Cruzadas con 10 grupos para los métodos SVM, KNN, Árboles, Bosques, Potenciación, eXtreme Gradient Boosting, LDA, Bayes, Regresión Logística, ConsensoPropio y Redes Neuronales, para KNN y Potenciación use los parámetros obtenidos en las calibraciones realizadas en los ejercicios anteriores. Luego grafique las 5 iteraciones para todos los métodos en el mismo gráfico. ¿Se puede determinar con claridad cuál métodos es el mejor?
2. Repita el ejercicio anterior, pero en lugar de sumar la cantidad de 1's, promedie los errores globales cometidos en los diferentes grupos (folds). Luego grafique las 5 iteraciones para todos los métodos vistos en el curso en el mismo gráfico. ¿Se puede determinar con claridad cuál algoritmo es el mejor?
3. Para estar realmente seguros de cuál de los tipos algoritmos es mejor, modifique (una en uno solo código) los códigos de los dos ejercicios anteriores de manera que, en lugar de sumar la cantidad de 1's detectados y de promediar los errores globales, guarde en cada iteración (en listas) las matrices de confusión de cada uno de los 3 algoritmos. Luego al

final de los ciclos ejecutados para cada algoritmo calcule una **Matriz de Confusión Promedio** de todas las matrices de confusión guardadas en la lista del respectivo algoritmo; así con estas **Matrices de Confusión Promedio**, mediante gráficos de barras, determine el método que en promedio detecta mayor porcentaje de 1's, determine el método que en promedio detecta mayor porcentaje de 0's y determine el método que en promedio tiene menor error global.

4. ¿Cuál método usaría con base en la información obtenida en los dos ejercicios anteriores?