

## Desafío - Mecanismos de Votación

- Para realizar este desafío debes haber revisado la lectura y videos correspondiente a la unidad.
- Crea una carpeta de trabajo y guarda todos los archivos correspondientes (notebook y `csv`).
- Una vez terminado el desafío, comprime la carpeta y sube el `.zip` a la sección correspondiente.

### Descripción

Para esta actividad desarrollaremos un ensamble heterogéneo que prediga la popularidad o inpopularidad de contenido en línea de un sitio web. El archivo `csv` tiene un total de 39797 registros, donde cada uno representa características del artículo publicado en el sitio web. Cada registro tiene un total de 61 atributos que miden características del texto, cantidad de imágenes, keywords, etc. Los artículos provienen de la publicación *K. Fernandes, P. Vinagre and P. Cortez. A Proactive Intelligent Decision Support System for Predicting the Popularity of Online News. Proceedings of the 17th EPIA 2015 - Portuguese Conference on Artificial Intelligence, September, Coimbra, Portugal.* Los artículos registrados son de la página web [www.mashable.com](http://www.mashable.com).

Más información sobre los atributos se puede encontrar en el archivo `OnlineNewsPopularity.names`.

### Ejercicio 1 - Preparación del Ambiente de Trabajo

- A continuación se le presenta un código que agrupa las variables por grupos. Cada uno de éstos hace referencia a alguna característica de los registros en la base de datos. Genere una análisis descriptivo de cada uno de los atributos. Puede utilizar la función `grid_plot_batch` que se encuentra en el archivo `helpers.py`.

```
# las etiquetas de las columnas presentan un espacio extra, con ésto lo
podemos eliminar
df.columns = [i.replace(' ', '') for i in df.columns]
# eliminamos el string de url que no sirve para el análisis
df = df.loc[:, 'n_tokens_title':'shares']
# generamos el conjunto de variables
qty = df.filter(regex='^n_', axis=1)
channel = df.filter(regex='^data_', axis=1)
days = df.filter(regex=re.compile("weekday|weekend"), axis=1)
sentiments = df.filter(regex=re.compile("negative|positive|subjectivity"),
axis=1)
lda = df.filter(regex='^LDA_\d', axis=1)
keywords = df.filter(regex='^kw_', axis=1)
```

- Describa el comportamiento de las variables.
- Dada la naturaleza de los atributos, es probable que algunas mediciones estén correlacionadas entre sí. Para ello, genere un diagnóstico previo de multicolinealidad utilizando la función `identify_high_correlations` que se encuentra en el archivo `helpers.py`. Para todos aquellos atributos que tengan una correlación de .8, reporte sus nombres.
- Antes de generar los conjuntos de entrenamiento y validación, preprocese los datos con los siguientes pasos:
  - Recodifique la variable `shares` en una variable binaria que identifique como 1 todos los registros con más de 1400 "compartir" y 0 de lo contrario. Para evitar multicolinealidad, elimine la variable `shares` posteriormente.
  - Elimine todas las variables que presentaban una correlación mayor a .8. Este paso es para evitar la multicolinealidad de los atributos.
  - Genere un análisis de Componentes Principales para extraer las principales 30 dimensiones. Guarde estas dimensiones en un nuevo objeto.

## Ejercicio 2 - Evaluación de modelos individuales

- A continuación generará una serie de modelos que se incorporarán posteriormente al comité de votación. Para ello, se solicita que:
  - Importe los módulos correctamente.
  - Para cada uno de ellos, genere un reporte en las métricas de desempeño respecto a `Precision`, `Recall`, `F1`. Puede hacer uso de la función `plot_classification_report` disponible en el archivo `helpers.py`.
  - Comente el desempeño general de cada uno
- La lista de modelos es la siguiente. Cabe destacar que la mayoría de éstos corresponden a implementaciones *vanilla*, salvo que se indique lo contrario:
  - Regresión Logística.
  - Algoritmo de KMedias.
  - Árbol de Clasificación con un `max_depth=1`.
  - Árbol de Clasificación con un `max_depth=4`.

### **Ejercicio 3 - Entrenamiento de Comité**

- Entrene el comité de clasificadores sin modificar el esquema de votación.
- Reporte el desempeño a nivel de cada clase para cada métrica.

### **Ejercicio 4 - Calibración de Comité con Ponderadores**

- El base al comportamiento de los clasificadores individuales del ensamble, proponga dos esquemas de ponderación para mejorar el desempeño del modelo.
- Reporte el desempeño del mejor ensamble heterogéneo.