UT3 - Algoritmos y herramientas para el aprendizaje supervisado

Actividad 3.7 - Predicción de Riesgo de derrumbamiento - Terremotos

El objeto de esta actividad es participar en la competición de ofrecida de la web de DrivenData denominada: Richter's Predictor: Modeling Earthquake Damage.

Título: Richter's Predictor: Modeling Earthquake Damage

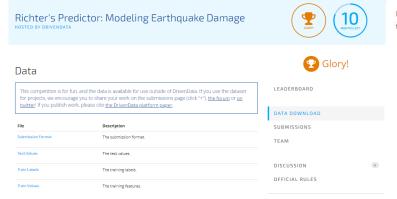
Url: https://www.drivendata.org/competitions/57/nepal-earthquake/data/

La actividad consiste en subir a dicha web un fichero csv con la estimación del nivel de daño provocado en un edificio/vivienda tras un terremoto (1 - low damage, 2 - medium amount of damage, 3 - complete destruction) que hayamos obtenido al aplicar al menos tres modelos de árboles y uno de SVM.

Podrás observar que se pueden realizar hasta un máximo de 3 subidas diarias y la propia web realizará una valoración de tu solución. La valoración se realiza utilizando el **micro averaged F1 score** (<u>variante</u> de <u>F1 score</u> – Ver D18 <u>del UT3 - Algoritmos y herramientas para el aprendizaje(I)_v2.pdf</u>) como criterio de valoración de calidad. Esta técnica es recomendable utilizarla cuando los datos no están balanceados.

Título: sklearn.metrics.f1_score – (Fijarse en average= 'micro')

Url: https://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1 score.html



Based on aspects of building location and construction, your goal is to predict the level of damage to buildings caused by the 2015 Gorkha earthquake in Nepal.



Richter's Predictor: Modeling Earthquake Damage

regression is sometimes described as an problem somewhere in between classification and regression.)

To measure the performance of our algorithms, we'll use the **F1 score** which balances the precision and recall of a classifier. Traditionally, the F1 score is used to evaluate performance on a binary classifier, but since we have three possible labels we will use a variant called the **micro averaged F1 score**.

$$F_{micro} = rac{2 \cdot P_{micro} \cdot R_{micro}}{P_{micro} + R_{micro}}$$

where

$$P_{micro} = rac{\sum_{k=1}^{3} TP_k}{\sum_{k=1}^{3} (TP_k + FP_k)}, \;\; R_{micro} = rac{\sum_{k=1}^{3} TP_k}{\sum_{k=1}^{3} (TP_k + FN_k)}$$

and TP is True Positive, FP is False Positive, FN is False Negative, and k represents each class in 1,2,3.

In Python, you can easily calculate this loss using sklearn.metrics.fl_score with the keyword argument average='micro'. Here are some references that discuss the micro-averaged F1 score further:

Consideraciones a tener en cuenta (Leer la rúbrica):

- Dado que el dataset que nos ofrece el reto contiene un número de filas muy elevado de instancias/filas para realizar el ejercicio, se hace necesario el realizar una selección de un conjunto de ellas que consideres oportuno. Sea cual sea la técnica/herramienta o criterio que propongas lo has de justificar. Este punto se valorará mejor cuanto menos aleatorio sea.
- Aplicar al menos tres modelos de árboles y uno de SVM entre los vistos en clase.
- Realizar pruebas de hiper-parametrización con las dos técnicas explicadas: GridSearch y RandomSearch.
- Es fundamental el utilizar la herramienta de los dendogramas para determinar qué conjunto de características nos conviene seleccionar para realizar la predicción.
- Utilizar la librería (Lazy Predict) indicada en la siguiente publicación con el objeto de determinar si realmente puede resultar interesante a la hora de decidir el modelo a utilizar en la resolución de un problema.

Título: lazypredict 0.2.12

Url: https://pypi.org/project/lazypredict/

Título: Do you need to build a ML classifier but don't know what model to use? Try

"Lazy Predict" to get an idea of the most promising models

Url: https://www.linkedin.com/posts/agostino-calamia_mlops-data-datascience-

activity-7044956022095945730-

mwVB?utm source=share&utm medium=member desktop

La rúbrica a utilizar para evaluar

	Actividad 3.7 - Predicción de Riesgo de derrumbamiento - Terremotos
Peso %	Tareas (Se evalúa entre 0 y 10)
2	Utilizar el <u>drive/github</u> como origen de ficheros para la importación del <u>dataset</u>
5	Importación del dataset: Preparación de los datos: Normaliza, ajusta la calidad de los datos
5	Selección de <u>carcaterísticas</u> : Explora herramientas gráficas o no gráficas, que no sean los <u>dendogramas</u> , para la elección de <u>caractaerísticas</u>
5	Selección de carcaterísticas: Utiliza dendogramas para la elección de las características
5	Además de la división de los datos de train y test, incorpora la utilización de datos de validación.
12	Entrenamiento: Trabaja con tres modelos de árboles - Desarrolla las diversas pruebas propuestas para la selección y justific criterio de calidad para la selección del modelo que mejores resultados ofrece. Utiliza Cross Validation
13	Entrenamiento: Elegir un modelo de regresión o SVM - Desarrolla las diversas pruebas propuestas para la selección y justific criterio de calidad para la selección del modelo. Utiliza Cross Validation y pruebas de hiperparámetros (GridSearch y RandomSearch)
10	Entrenamiento: Utilizar la librería Lazy Predict para determinar otros posibles modelos a utilizar para la resolución del probl elije uno de los modelos propuestos y compara la predicción con el modelo de regresión o SVM elegido
10	Predicción: Utiliza herramientas gráficas para ayudar a entender la precisión de los resultados obtenidos
5	Predicción: Describe con claridad una valoración de los resultados obtenidos.
10	Submit del fichero con la predicción y captura de la valoración /posicionamiento obtenido en la competición
5	Propone soluciones creativas e innovadoras
3	Registra en el pdf final, un cuaderno de bitácora/seguimiento donde se muestra los submitis realizados, explicando los ajus mejoras que han motivado cada subida o un grupo de submits
5	El gdf final tiene una portada., utiliza un índice, apartado de conclusiones y referencias (web). Se hace mención a referencia externas, no recogidas en el material suministrado.
5	Comenta con claridad cada uno de los pasos realizados.

Formato de entrega

- Entregar un fichero en un Archivo PDF con capturas del código y resultados obtenidos, así como la url de GitHub y Google Colab donde has publicado el código.
- Nombrar el archivo siguiendo el siguiente patrón:

 $SNS_ACT3_7_Nombre Apellidos.pdf$