

Inteligencia de Negocio

Práctica 3

Guillermo Gómez Trenado | 77820354-S
guillermogotre@correo.ugr.es

10 de enero de 2019

Índice

1. Descripción del problema	3
2. Análisis de datos	3
2.1. Valores nulos	4
2.2. Correlación	5
2.3. Relación con la etiqueta de la salida	7
2.4. Métrica y balanceo	7
3. Resultados obtenidos	8
4. Preprocesado	9
4.1. Código original (Base)	9
4.2. Transformaciones	9
4.2.1. Fechas	9
4.2.2. Distancia a un punto	10
4.2.3. Word2Vector	12
4.2.4. Orden por clasificación de las variables categóricas	13
4.3. Reducción de la dimensionalidad	14
4.3.1. Eliminación de características	14
4.3.2. Reducción de características	14
4.3.3. Eliminación de instancias (Ruido)	15
4.4. Imputación	16
4.4.1. Imputación de longitud y latitud	16
4.4.2. Imputación de altura	19
4.4.3. Otras imputaciones	20
5. Clasificadores	20
5.1. Uso de dos clasificadores (Dual)	20
5.2. Uso de tres clasificadores (Triple)	22

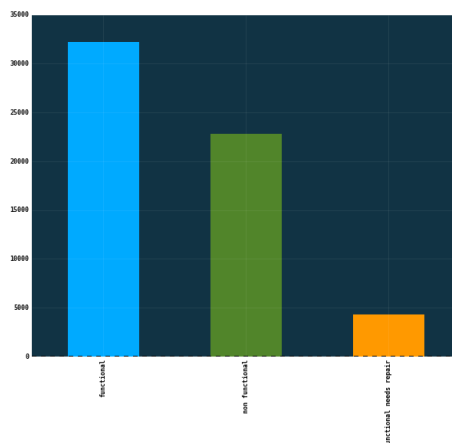
6. Ajuste de hiperparámetros	23
6.1. Algoritmo genético para los hiperparámetros	23
6.2. Evitar sobreajuste	24
6.3. Extra. Differential evolution para ponderación de clasificadores individuales	25
7. Reflexiones finales	28

1. Descripción del problema

El problema que se nos plantea es el etiquetamiento del estado de pozos de agua en Tanzania a partir de diferentes características que se nos presentan, binarias, numéricas y categóricas. Tenemos un conjunto de entrenamiento con sus etiquetas correspondientes —sin separación entre train y validación— y un conjunto de test que se evalúa en la página web, para no hacer sobreajuste con los datos de test sólo se permiten tres subidas al día. La métrica utilizada en la página web es la porción de instancias correctamente etiquetadas.

2. Análisis de datos

Lo primero que hice fue leer la descripción del problema¹. Algunas características están mejor comentadas que otras, pero aun así nos permiten hacernos una idea general del problema, es un problema desbalanceado, hay tres etiquetas desigualmente distribuidas



Después de ver la descripción busqué información sobre el problema en google para un acercamiento preliminar a su análisis, ya que no era la primera vez que se realiza el concurso, y por lo tanto hay muy buenos análisis en internet, especialmente los citados abajo.^{2 3 4 5}

¹DrivenData. Pump it Up: Data Mining the Water Table. Problem Description. <https://www.drivendata.org/competitions/7/pump-it-up-data-mining-the-water-table/page/25/>

²Pump it Up: Data Mining the Water Table. Análisis https://rstudio-pubs-static.s3.amazonaws.com/187095_fa3184e3f1244ba9b98f351570699732.html

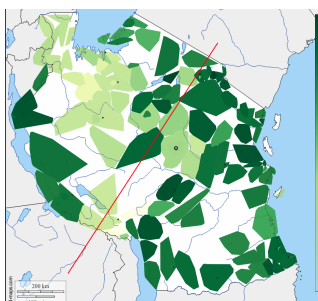
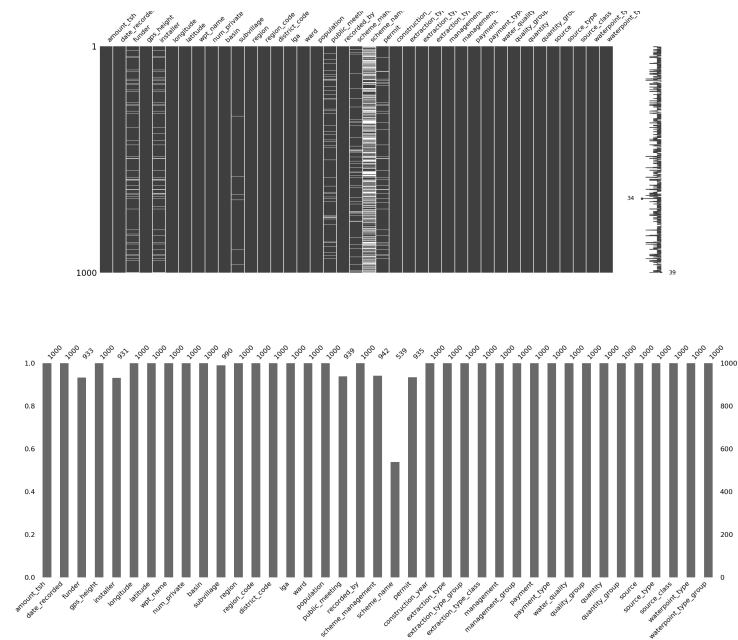
³Vaibhav Shukla. Medium. Pump it Up: Data Mining the Water Table <https://medium.com/@vaibhavshukla182/pump-it-up-data-mining-the-water-table-f903d4cfc7a8>

⁴David Currie. Pump it Up: Data Mining the Water Table. https://currie32.github.io/water_pumps.html

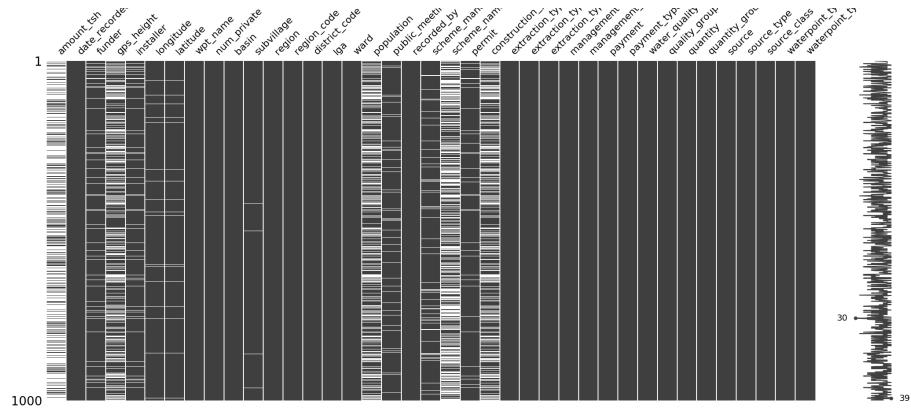
⁵Calin Uioreanu. Community Driven Data. <https://community.drivendata.org/t/pumps-visualization-dashboard-with-tableau/186>. Tableau. https://public.tableau.com/profile/calin.uioreanu#!/vizhome/DataMiningtheWaterTableDrivenData_com/Bubblestatusquantity

2.1. Valores nulos

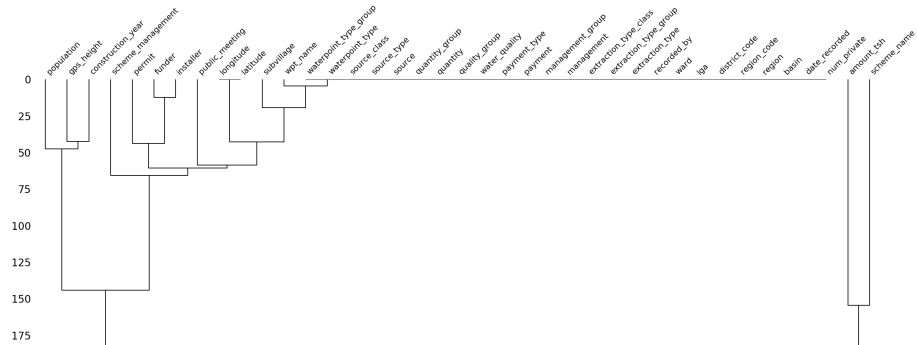
Tras este primer acercamiento al problema, el primer asunto que quiero tratar es la presencia de valores nulos. Para las variables categóricas es inmediato porque vienen etiquetadas con *NaN*. Veamos tres representaciones de la proporción de valores nulos, el primero una representación de una muestra de mil instancias, la segunda imagen un histograma, y en tercer lugar la distribución por código de región —vemos cómo no sólo cómo las variables están irregularmente rellenas entre ellas sino con distinta distribución en distintas zonas del mapa.



Sin embargo, para las numéricas no queda claro cuando el cero significa cero y cuándo es una característica no cumplimentada. A la luz de la descripción de la página web, las tres variables a 0 que con seguridad suponen valores nulos son *amount_tsh* —cantidad de agua disponible al pozo—, *longitud* y *latitud* —la coordenada 0,0 está fuera de tanzania— y *gps_height* —altura 0—.



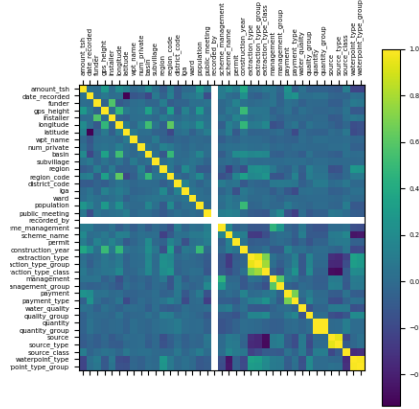
Lo primero que observo es que las variables con mayor porcentaje de valores nulos son las relativas a las coordenadas y altura geográfica, la cantidad de agua disponible, la cantidad de población que utiliza el pozo, el año de construcción y quién opera el pozo. Sobre las geográficas tendremos oportunidad de centrarnos con mayor detenimiento en el preprocesado.



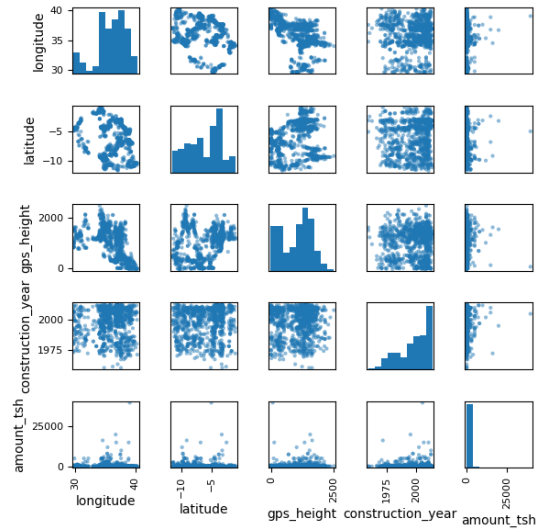
Vemos como las variables agrupadas por su tasa de valores nulos estás por lo general semánticamente relacionadas, esta intuición nos será útil para la imputación.

2.2. Correlación

Ahora me interesa ver la naturaleza de las características y su relación, aplico *LabelEncoder* para poder hacer una matriz de correlación y la imprimo



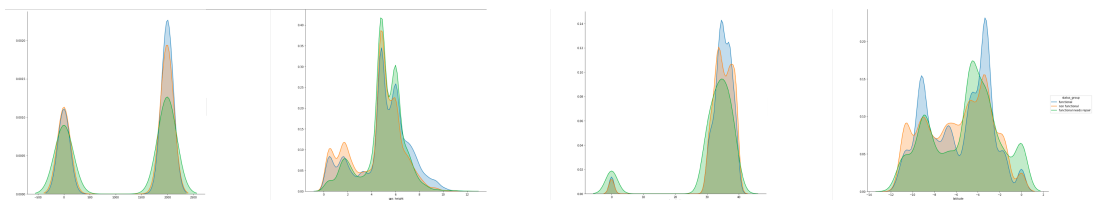
Lo que analizamos es que hay variables altamente correlacionadas, al analizarlas en detalle observo que explican el mismo fenómeno, como sucede con *extraction_type*, *extraction_type_group* y *extraction_type_class*, o aquellas que diferencian entre *type* y *class* o entre *quantity* y *quantity_group*, este tipo de variables serán las primeras que eliminemos en el preprocesado.



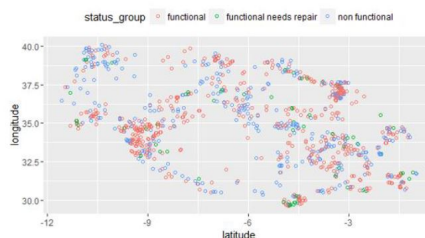
Arriba las gráficas de correlación e histogramas de las variables geográficas — después de eliminar los ceros — porque serán sobre las que centre principalmente el preprocesado por ser altamente inteligibles y fáciles de imputar con relativa confianza.

2.3. Relación con la etiqueta de la salida

Vaibhav Shukla en Medium ⁶ hace un análisis interesantísimo en su descripción del problema de las distribuciones superpuestas según la clase, adjunto algunas como referencia



Analizo las cuatro variables numéricas que veo más fácil imputar, vemos que el año de construcción no parece tener una capacidad predictiva relevante, pero la longitud, latitud y altura parecen ser unos predictores de relativa calidad como podemos confirmar en la siguiente gráfica de Jithin Paul⁷. Por eso nos centraremos en estos valores para la imputación.



2.4. Métrica y balanceo

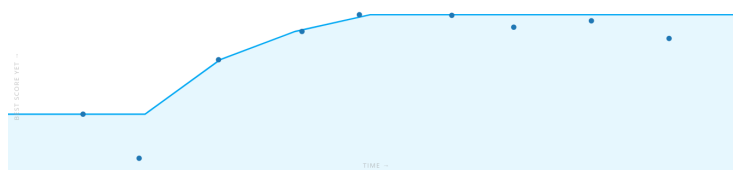
La métrica con la que se nos evalúa en DrivenData es la tasa de acierto, que como ya hemos visto no es la métrica más robusta ni más descriptiva de la calidad del clasificador, el clasificador constante arroja una tasa de acierto de $\sim 0,56$ —con la salida constante *functional*— aun siendo un problema con tres clases, debido al desbalanceo. Sin embargo, debido a la métrica usada no voy a centrar mis esfuerzos en balancear las clases, pues si bien pudiera esto generar un clasificador más robusto y que sirviera mejor a los propósitos últimos de esta competición, lo más posible es que penalizara el rendimiento en DrivenData pues ellos no atajan en su métrica el problema de la falta de balanceo, de hecho si hicieramos un clasificador binario, ignorando la clase *functional need repair* podríamos optar a clasificar correctamente el 93 % de la muestra.

⁶Vaibhav Shukla. Medium. Pump it Up: Data Mining the Water Table <https://medium.com/@vaibhavshukla182/pump-it-up-data-mining-the-water-table-f903d4cfc7a8>

⁷Jithin Paul. PUMP IT UP – Data Mining the Water Table. <https://jitpaul.blog/2017/07/12/pump-it-up/>

3. Resultados obtenidos

Descripción	CV 5fold	All train	DrivenData	DD Posicion	Preprocesado	Algoritmo	Parámetros
Base	0.7954	0.825			LE	LGBM	Estimators=200
Fechas	0.7966	0.8262			LE+Date	LGBM	Estimators=200
Geodist	0.7971	0.8262			LE+Date+GeoDist	LGBM	Estimators=200
KNN GPS	0.7969	0.8267	0.8034	833	LE+Date+GeoDist+GPS	LGBM	Estimators=1000
W2V	0.7963	0.8364	0.7937	833	W2V+Date+GeoDist+GPS	LGBM	Estimators=1000
Height KNN	0.7978	0.826	0.8155	641	LE+Date+GeoDist+GPS+Height	LGBM	Estimators=1000
Ordenadas	0.8099	0.8999	0.8218	295	Sorted LE+Date+GeoDist+GPS+Height	LGBM	Estimators=1000
Dual	0.8146	0.9059			Sorted LE+Date+GeoDist+GPS+Height	LGBM+RF	Es=1000+Es=500,Maxdepth=20
Triple	0.816	0.9503	0.8255	57	Sorted LE+Date+GeoDist+GPS+Height	XGBOOST+LGBM+RF	Es=500, NClass=4, MaxDepth=14 Es=1000 Es=500,Maxdepth=20
Genetico Hiperparámetros	0.8165	0.9411	0.8253	57	Sorted LE+Date+GeoDist+GPS+Height	XGBOOST+LGBM+RF	Descripción en el apartado respectivo
Ruido	0.8191	0.9152	0.8241	57	Sorted LE+Date+GeoDist+GPS+Height+Noise	XGBOOST+LGBM+RF	
No overfit	0.8172	0.8841	0.8202	57	Sorted LE+Date+GeoDist+GPS+Height	XGBOOST+LGBM+RF	
Differential Evolution Pesos*	0.8179	0.9312	0.8258	48	Sorted LE+Date+GeoDist+GPS+Height	XGBOOST+LGBM+RF	



Submissions

BEST	CURRENT RANK	# COMPETITORS	SUBS. TODAY
0.8255	58	6184	0 / 3

SUBMISSIONS			
Score	Submitted by	Timestamp	
0.8034	Guillermo_Gomez_UGR	2019-01-02 21:12:52 UTC	
0.7937	Guillermo_Gomez_UGR	2019-01-02 22:13:48 UTC	
0.8155	Guillermo_Gomez_UGR	2019-01-03 22:30:24 UTC	
0.8218	Guillermo_Gomez_UGR	2019-01-03 23:48:25 UTC	
0.8255	Guillermo_Gomez_UGR	2019-01-04 00:39:57 UTC	
0.8253	Guillermo_Gomez_UGR	2019-01-04 21:30:45 UTC	
0.8227	Guillermo_Gomez_UGR	2019-01-05 03:03:18 UTC	
0.8241	Guillermo_Gomez_UGR	2019-01-05 15:51:18 UTC	
0.8202	Guillermo_Gomez_UGR	2019-01-05 22:43:52 UTC	

Al final del plazo para participar en la competición, el resultado obtenido en DrivenData fue de 82.55, estando en la posición 57. Sin embargo, un día después, fuera de plazo, con otro usuario que creé al finalizar éste y con una estrategia que seguía desarrollando la línea de las anteriores conseguí colocarme en el puesto 48, desplazándome a mí mismo, en la cuenta creada *ex profeso* para la UGR hasta el 58, adjunto también la información relativa a esta última subida como extra.

Además de la documentación adjunto una carpeta con el código y los csv con el mismo título que en la tabla de arriba. Como el proceso que seguí fue progresivo, primero centrándome en el preprocesado y posteriormente en los clasificadores, vamos a seguir esta misma sucesión en la descripción.

Por otro lado, no tengo el resultado en DrivenData para las primeras modificaciones porque como le comenté por correo no recibí el correo de confirmación hasta 3 días después de que comenzara el desarrollo de la práctica, y como el número de subidas estaba limitado decidí invertirlo en mis progresos.

4. Preprocesado

Mi estrategia se centró en ir transformando todas las variables que tuvieran un significado fuertemente ordinal —categóricas y evidentemente numéricas— en variables numéricas con transformaciones personalizadas a sus características y por último atajar el problema de las categóricas puras. Tanto el preprocesado como la imputación, por no basarse en la etiqueta de salida y debido al enorme tiempo que consumen algunas lo he realizado fuera de la validación cruzada, como tenemos una tercera etapa de test, aunque sea incorrecto, podemos comprobar la calidad última de estas transformaciones.

4.1. Código original (Base)

La primera y única modificación que realicé sobre el código que se nos facilitó fue modificar *LabelEncoder* ya que de la forma en la que estaba definido, aunque gracias al orden alfabético estarían próximas las variables entre ellas, no teníamos garantía de que una instancia con la misma etiqueta en train y test tuvieran el mismo valor numérico tras el preprocesado. Sencillamente uní las dos tablas, apliqué *LabelEncoder* sobre las variables y las volví a separar, como el orden es alfabético sin observar la etiqueta de salida, no hay por qué entrenar la operación de transformación sobre el train y aplicarla sobre el test —reservando una etiqueta para las variables desconocidas— porque no se puedo producir *data leaking*.

4.2. Transformaciones

4.2.1. Fechas

La primera variable que modifiqué fue la fecha en que se registró la instancia, como el formato de fecha no es inteligible por los algoritmos disponibles en las librerías de clasificadores con las que trabajamos —scikit-learn, xgboost y lgbm— la dividí en dos nuevas características, el tiempo sucedido desde que se registró la primera de la que tenemos constancia en el conjunto de train en días por un lado, y una segunda categoría, el mes en el que se registró, por si tuviera alguna relevancia que desconocemos —luego descubriremos que no la tiene y la eliminaré—, la idea era no perder información ni del tiempo transcurrido ni de la periodicidad inherente en el tiempo.

```

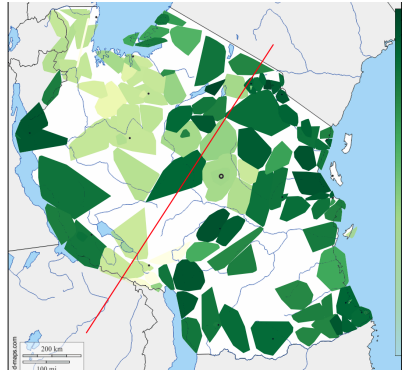
1 # String to Python datetime
2 def parsed(x):
3     TIME_STR = '%Y-%m-%d'
4     return datetime.datetime.strptime(x, TIME_STR)
5
6 # Days from 2002-10-14 (Precomputed)
7 def getDays(date_c,min_date=parsed('2002-10-14')):
8     # Si queremos calcular la fecha en vivo
9     if min_date is None:
10         min_date = parsed(date_c.min())
11     days = [(parsed(d) - min_date).days for d in date_c]
12     return days
13
14 # Month of the year (1-12)
15 def getMonth(date_c):
16     months = [parsed(d).month for d in date_c]
17     return months
18
19 # Add new columns
20 def addDatesPrep(data):
21     # Obtenemos la columna de fecha
22     days = getDays(data['date_recorded'])
23     # Aniadimos días desde la referencia
24     data['days_ellapsed'] = days
25     # Aniadimos mes del año
26     data['month_n'] = getMonth(data['date_recorded'])

```

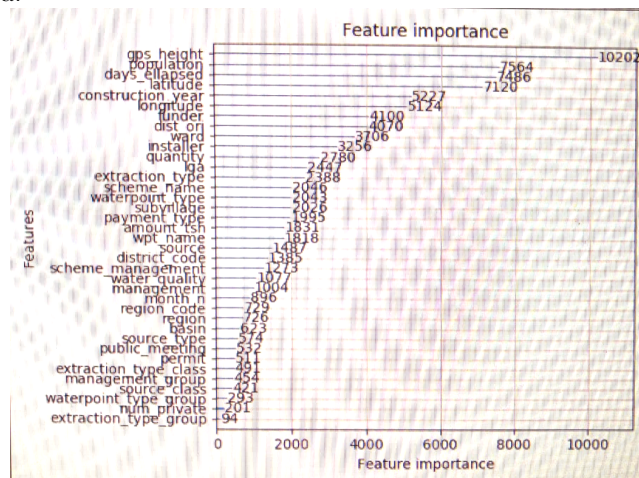
Normalmente el criterio que seguí fue evaluar sobre 5-fold CV y si no empeoraba al aplicar el cambio conservaba la transformación, a fin de ir eliminando las variables no inteligibles por los algoritmos. En este caso pasamos de 0.7954 a 0.7966, que si bien no tenemos confianza para indicar que sea estadísticamente mejor, parece que al menos no es peor.

4.2.2. Distancia a un punto

Esta transformación se volverá a realizar después de la imputación, pero la describimos aquí porque la desarrollé en primer lugar. La idea era introducir una tercera coordenada no ortogonal a las dos ya presentes, probé con varios orígenes de referencia y no encontré una mejora significativa de ninguno respecto al otro, así que elegí la coordenada 0,0. El resultado fue una nueva coordenada relativamente redundante perpendicular la línea roja dibujada, además añadí una segunda variable booleana que reflejaba si la coordenada que imputaría en un paso posterior era la original o una imputada.



A pesar de ser redundante vi que los resultado mejoraban ligeramente sobre CV —de 0.7966 a 0.7971— y en la gráfica de importancia sobre las características utilizadas por lgbm estaba entre las primeras. Siempre estamos a tiempo de eliminarla posteriormente para reducir la dimensionalidad. En un exabrupto de genialidad tras pasarle una foto a un compañero de esta gráfica la eliminé , por eso tiene una calidad tan pobre, porque he tenido que reutilizar una fotografía de la gráfica.



```

1 # Librería para la distancia esférica
2 from geopy.distance import geodesic
3
4 def addGeoDist(data):
5     # Obtenemos tupla long, lat
6     coords = list(zip(data['longitude'], data['latitude']))
7
8     # Eliminamos los valores nulos
9     msk = np.abs(np.sum(coords, axis=1)) < 1e-5
10

```

```

11 # Booleano coordenada imputada
12 data['reliable_gps'] = ['T' if m else 'F' for m in msk]
13
14 # Calculo distancia
15 kms = [geodesic(p1,(0,0)).km for p1 in coords]
16 data['dist_ori'] = kms
17
18 return msk

```

4.2.3. Word2Vector

La siguiente estrategia utilizada fue transformar las variables categóricas en numéricas con el mismo algoritmo word2vec que ha demostrado buenos resultados con variables categóricas ⁸, para ello utilicé la librería de *gensim* ⁹.

El código es demasiado extenso para copiarlo aquí, el procedimiento es el siguiente. Elegimos las columnas que queremos transformar, creamos una bolsa de palabras con el conjunto completo —unión de todas las columnas— con la forma “Feature {value} {column_name}”, definimos un tamaño de ventana y un tamaño para los vectores —definido experimentalmente—. Aunque este procedimiento —idealmente— es capaz de separa espacialmente las frases según su parecido dadas sus características semánticas, el aumento de la dimensionalidad es tremendo, número de características categóricas por tamaño de los vectores—.

```

1 CAT_KEYS =
  → ["funder", "installer", "wpt_name", "basin", "subvillage", "region",
2 "lga", "ward", "recorded_by", "scheme_management", "scheme_name",
3 "extraction_type", "extraction_type_group", "extraction_type_class",
4 "management", "management_group", "payment", "payment_type", "source",
5 "source_type", "source_class", "water_quality", "quality_group",
6 "quantity", "quantity_group", "waterpoint_type", "waterpoint_type_group"]
7 BIG_CAT_KEYS = ['funder', 'installer', 'wpt_name', 'subvillage',
  → 'ward', 'scheme_name']

```

Probé en primer lugar con *CAT_KEYS* y la explosión de la dimensionalidad era inasumible, y posteriormente con *BIG_CAT_KEYS* y el rendimiento no mejoraba en train —de 0.7969 a 0.7963—, empeoraba en DrivenData —0.8034 a 0.7937—, y era costosísimo a nivel computacional, así que lo descarté.

⁸MyYellowRoad. Yonatan Hadar. Using categorical data in machine learning with python: from dummy variables to Deep category embedding and Cat2vec -Part 2<https://blog.myyellowroad.com/using-categorical-data-in-machine-learning-with-python-from-dummy-variables-to-deep-category-42fd0a43b009>

⁹Gensim. Word2Vec <https://radimrehurek.com/gensim/models/word2vec.html>

4.2.4. Orden por clasificación de las variables categóricas

Tras descartar la intentona anterior probé otra estrategia también decrita por Yonatan Hadar¹⁰. La idea es que aquellas etiquetas que clasifiquen de una forma similar estén próximas entre ellas, este procedimiento, al realizarlo fuera de la validación cruzada sí que produce *data leaking*, y encontramos información sobre las etiquetas de salida de las instancias de validación en el conjunto de entrenamiento —a través de las estadísticas sobre el conjunto completo—, sin embargo tomé dos precauciones, esta fue la última modificación antes de pasar a modificar los clasificadores y lo validé sobre el test en la página de DrivenData, mal hecho lo anterior pero decidí priorizar el desarrollo rápido y poder realizar muchas pruebas por encima del rigor —y a la luz de los resultados en DataDriven parece que no fue una estrategia equivocada—. Como tenemos tres categorías podía o crear dos columnas de salida por cada columna de entrada para reflejar toda la información o elegir la mayoritaria —opté por esta última opción—.

```
1 def ordenateCategories(tr,tst,y):
2     # Obtenemos las columnas categóricas
3     keys = list(filter(lambda x: tr[x].dtype == np.dtype('O'),
4         ↪ tr.columns))
5     # Por cada columna
6     for k in keys:
7         # Obtenemos los valores distintos sobre el train
8         vals = list(Counter(tr[k].values.reshape(-1)))
9         # Los ordenamos por su porcentaje de clasificacion como
10        ↪ funcionales
11        impV = sorted(map(
12            lambda x: (np.sum((y[tr[k] == x] ==
13                ↪ 'functional').values) / np.sum((tr[k] ==
14                ↪ x).values), x),vals),reverse=True)
15        # Creamos un diccionario asignando un índice según su
16        ↪ posición
17        d = dict(map(lambda x: (x[1][1], x[0]), enumerate(impV)))
18
19        # Modificamos train
20        tr[k] = list(map(lambda x: d.get(x,-1),tr[k].values))
21        # Modificamos test
22        tst[k] = list(map(lambda x: d.get(x,-1),tst[k].values))
```

¹⁰MyYellowRoad. Yonatan Hadar. Using categorical data in machine learning with python: from dummy variables to Deep category embedding and Cat2vec -Part 2<https://blog.myyellowroad.com/using-categorical-data-in-machine-learning-with-python-from-dummy-variables-to-deep-category-42fd0a43b009>

4.3. Reducción de la dimensionalidad

Las siguientes transformaciones no se encuadran temporalmente en este momento concreto sino que forman parte de otras de modificaciones que veremos más adelante, sin embargo, materialmente tiene sentido agruparlas aquí.

4.3.1. Eliminación de características

El siguiente paso fue utilizar el análisis realizado previamente para eliminar las columnas estrictamente redundantes o sin significancia para el problema, estas son *recorded_by*, *payment*, *quality_group*, *quantity_group*, excepto *recorded_by* las otras estaban duplicadas en el dataset, y eran explicadas de la misma manera o de forma más precisa por otra variable. De esta modificación no registré la evaluación porque entendí que sólo podía beneficiar a los algoritmos y evitar la *maldición de la dimensionalidad*.

4.3.2. Reducción de características

Al substituir *LabelEncoder* por la transformación personalizada que veremos a continuación añadí una transformación preliminar para reducir la dimensionalidad del problema, tomé aquellas variables categóricas que restaban por transformar y por cada una de ellas transformé los valores nulos a una nueva categoría llamada “NaN”, conté las apariciones de cada valor y aquellos con menos de $MIN = 50$ apariciones los agrupé en una nueva etiquetada llamada “Others”. El valor de MIN lo definí experimentalmente, a medida que aumentaba el valor el resultado sobre la validación cruzada mejoraba pero llegaba un punto en el que penalizaba el rendimiento, 50 estaba próximo a ese punto dulce.

```
1 def reduceBigCat(tr,tst):
2     MIN = 50
3     gen_label = 'Other'
4     # Unimos train y test
5     X = pd.concat([tr,tst],axis=0)
6
7     # Por cada columna de tipo no numérico
8     for k in list(filter(lambda x: tr[x].dtype == np.dtype('O'),
9         ↪ tr.columns)):
10         # Sustituimos valores nulos por nueva etiqueta nula
11         col = X[k].fillna('NaN')
12         # Contamos ocurrencias
13         c = Counter(col.values.reshape(-1))
14         # Creamos una máscara con instancias a reducir
15         msk = col.isin(list(filter(lambda x: c[x] < MIN, c)))
16         # Cambiamos la etiqueta
17         col[msk] = gen_label
18         X[k] = col
```

```

18
19 # Devolvemos los dos subconjuntos de instancias separados
20 return X.iloc[:tr.shape[0],:], X.iloc[tr.shape[0]:,:]

```

4.3.3. Eliminación de instancias (Ruido)

En este caso, la reducción del ruido la realicé después de la imputación de las coordenadas y altura que veremos a continuación porque esta imputación no se basaba en la etiqueta de salida, y por lo tanto el ruido no se vería amplificado, además, las coordenadas nulas ya las había detectado y eliminado, al comprobar que entraran dentro de valores razonables para la posición geográfica de Tanzania. Sin embargo la analizo aquí por seguir un orden lógico en el tratamiento de los datos.

Para reducir el ruido utilicé la técnica basada en CV donde entreno con un subconjunto pero predigo la totalidad de los datos, tras las cinco particiones, aquellas instancias que no se hayan clasificado correctamente ninguna vez son eliminadas. Para la eliminación creo una máscara para que a la hora de la validación cruzada se entrene con el nuevo subconjunto pero se valida con la totalidad de los datos.

```

1 def cleanNoise(modelo,X,y):
2     # 5fold CV
3     cv = StratifiedKFold(n_splits=5, shuffle=True,
4         ↪ random_state=123456)
5     # Matriz de clasificación
6     corr = np.zeros((X.shape[0],5),dtype=np.bool)
7
8     for i = 0
9     for train, test in cv.split(X, y):
10         # Entrenamos con subconjunto
11         modelo = modelo.fit(X[train], y[train])
12         # Evaluamos todos los datos
13         y_pred = modelo.predict(X)
14         # Añadimos la columna con los aciertos
15         corr[:,i] = y_pred == y
16         i += 1
17     # Reduzco a una sola columna de booleanos
18     # donde los 0 deben ser eliminados
19     m = np.apply_along_axis(np.sum,axis=1,arr=corr).astype(np.bool)
20     return m

```

El resultado parecía prometedor, pues conseguí pasar en train de 0.8165 a 0.8191, sin embargo en DrivenData, bajó de 0.8253 a 0.8241. Esto me hace

pensar que si bien estos casos no pueden ser clasificados y esta estrategia ayuda a la clasificación más robusta del conjunto de entrenamiento, el desplazamiento de los *centroides* —en el sentido más poético posible— de las hojas de los árboles de decisión ayuda a una mejor clasificación sobre el test. Por lo tanto asumo que no es efectivamente ruido sino instancias que representan a un grupo muy reducido y que quedan absorbidas por subconjuntos mayores, sin embargo el incremento del tamaño de los árboles, para crear grupos más especializados repercutía negativamente sobre la validación cruzada, pues sobreajustaba los datos, con los parámetros que definí y veremos luego intenté equilibrar este fenómeno para obtener el mejor balance de esta situación indeseable.

4.4. Imputación

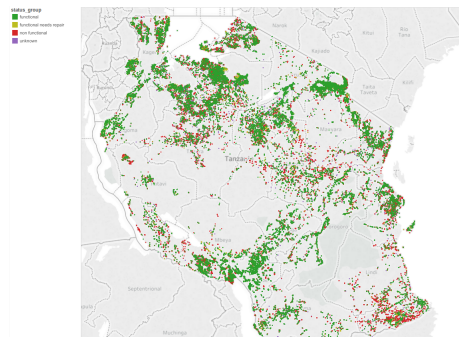
A la luz de la gráfica de importancia del LGBM decidí imputar longitud, latitud y altura, éstas tres por distintos motivos, el primero es que tenemos variables categóricas que hacen referencia a la posición geográfica bien cumplimentadas y en segundo lugar, tenía confianza en que la nulidad de estos datos no refleja ninguna condición del pozo sino de la toma de los datos, sin embargo el resto de variables —a excepción de *amount_tsh* que también intenté imputar pero el resultado no fue bueno— no está claro, por su pobre definición o por su propia naturaleza, que su nulidad sea por una pobre cumplimentación de los datos o por las propias características de los objetos que describen. También probé otras estrategias de imputación menos robustas como la mediana o KNN sobre variables categóricas transformadas pero en primer lugar no daban buenos resultados y en segundo lugar no estaba seguro que semánticamente reflejaran la naturaleza de las instancias.

4.4.1. Imputación de longitud y latitud

Esta era una de las variables que tenía más interés en imputar, por la importancia de éstas en las gráficas antes vistas —importancia de variables y distribución de clases en el mapa—. Podía hacer una imputación con suficiente confianza sobre su robustez basándome en las etiquetas categóricas que describen la posición en el espacio, y además es evidente que todos los puntos de extracción de agua tienen una posición en el espacio.

El mapa creado por Calin Uioreanu¹¹ fue un gran motivador para esta decisión, al ver que las distintas etiquetas se distribuían con distintas densidades dependiendo de la zona.

¹¹Calin Uioreanu. Community Driven Data. <https://community.drivendata.org/t/pumps-visualization-dashboard-with-tableau/186>. Tableau. https://public.tableau.com/profile/calin.uioreanu#!/vizhome/DataMiningtheWaterTableDrivenData_com/Bubblestatusquantity



El problema que enfrentaba era que las variables categóricas no se llevan bien con los algoritmos de regresión disponibles en las librerías que estábamos usando —LGBM y XGBooster sí que lo admiten—, por lo tanto decidí implementar mi propio KNN con mi propia métrica de distancia, ya que el orden alfabético o el orden estadístico en principio no debería reflejar ninguna relación espacial. La función de similitud es el número de etiquetas en las características de interés (*GPS_KEYS*) que comparten.

```

1 GPS_KEYS = ["basin", "subvillage", "region", "region_code",
  ↪ "district_code", "lga", "ward"]
2
3 # En preprocess
4 def preprocess(data_x, data_y, data_x_tst):
5     ...
6     data_x_copy = data_x.copy()
7     # Máscara para las etiquetas de interés de referencia
8     label_mask = [x in GPS_KEYS for x in data_x.columns]
9     # Máscara para longitud y latitud
10    coord_mask = [x in ['longitude', 'latitude'] for x in
  ↪ data_x.columns]
11
12    # Coordenadas predichas
13    coords = replaceGps(
14        data_x_copy.iloc[np.logical_not(msk), label_mask].values,
15        data_x_copy.iloc[np.logical_not(msk), coord_mask].values,
16        data_x_copy.iloc[msk, label_mask].values)
17
18    # Sustituimos en los datos
19    data_x.iloc[msk, coord_mask] = coords
20    ...
21
22 # Función con el KNN
23 def replaceGps(ref, ref_y, dst):
24     # Etiquetas de interés
25     gGps = ref

```

```

26  # Coordenadas de entrenamiento
27  gY = ref_y
28  # Instancias a predecir
29  bGps = dst
30
31  # Función de distancia
32  dist = lambda x1,x2: np.sum([l1==l2 for l1,l2 in zip(x1,x2)])
33
34  Tres vecinos
35  KNN = 3
36
37  res = []
38  # Definición del KNN por columna
39  def knn(row):
40      # Calculamos todas las distancias
41      d = np.apply_along_axis(lambda x: dist(x, row), 1, gGps)
42      # Apilamos distancias y coordenadas
43      m = np.hstack((d.reshape((-1, 1)), gY))
44      # Ordenamos por distancia
45      m = sorted(m, key=lambda x: x[0], reverse=True)
46      # Elegimos los vecinos
47      mnn = np.array(m[:KNN])
48      # Ponderamos por la similitud
49      mnn[:, 1] *= mnn[:, 0]
50      mnn[:, 2] *= mnn[:, 0]
51      # Devolvemos la media para longitud y latitud dividido por
52      ↪ la media de los pesos
53      return np.mean(mnn, axis=0)[1:] / np.mean(mnn, axis=0)[0]
54
55  # Multithreading
56  p = mp.Pool(4)
57  res = p.map(knn, bGps)
58  p.close()
59  p.join()
60  return res

```

El resultado no lo pude contrastar en DrivenData porque fue la primera subida que pude hacer —hasta entonces no tuve acceso— pero sobre el entrenamiento el resultado no empeoraba, 0.7971 a 0.7969, y me permitiría después imputar la altura con mayor seguridad. Cogí tres casos aleatorios y comprobé en Google Maps que las coordenadas coincidían con la descripción de las etiquetas de interés, así que quedé muy satisfecho.

4.4.2. Imputación de altura

El siguiente paso fue imputar la altura, una vez que tenía todas las coordenadas imputadas, podía entrenar un regresor para la altura, en este caso probé con la implementación de KNNRegressor de Scikit-Learn y con LGBMRegressor y el segundo me dio mejor resultado, así que este fue el que conservé

```
1 GPS_KEYS = ["basin", "subvillage", "region", "region_code",
2             ↪ "district_code", "lga", "ward"]
3
4 # En preprocess
5 def preprocess(data_x, data_y, data_x_tst):
6     ...
7     # Imputamos train consigo mismo
8     impZeroHeight(data_x, data_x)
9     # Imputamos el test con el modelo del train
10    impZeroHeight(data_x, data_x_tst)
11    ...
12
13 def impZeroHeight(src, dst):
14     # Obtenemos X(long,lat) e y(height) origen
15     coords_tr = src[['longitude', 'latitude']]
16     height_tr = src[['gps_height']]
17     # Obtenemos X(long,lat) e y(height) destino
18     coords_tst = dst[['longitude', 'latitude']]
19     height_tst = dst[['gps_height']]
20
21     # Máscara para el test
22     tst_msk = (height_tst == 0).values.reshape(-1)
23     coords_tst = coords_tst.values[tst_msk]
24
25     # Máscara para el train
26     msk = (height_tr != 0).values.reshape(-1)
27     # X(long,lat) no nulo
28     nonzero_coords = coords_tr.values[msk]
29     # y(height) no nula
30     nonzero_height = height_tr.values[msk]
31
32     # Creamos el regresor y lo entrenamos
33     rg = lgb.LGBMRegressor(n_estimators=500)
34     rg.fit(nonzero_coords, nonzero_height)
35
36     # Predecimos el destino
37     lbs = rg.predict(coords_tst)
38
```

```

39  # Sustituimos en el destino
40  dst.iloc[tst_msk, list(map(lambda x: x == 'gps_height',
    ↪   dst.columns))] = lbs.reshape((-1,1))

```

Este cambio fue significativo, me hizo pasar sobre CV de 0.7969 a 0.7978 y en DD de 0.8030 a 0.8155, un salto de casi 200 posiciones, ya estaba más cerca de mi objetivo.

4.4.3. Otras imputaciones

Intenté imputar otras variables, como *amount_tsh* pero los modelos arrojaban peor rendimiento, asimismo intenté sustituir los valores nulos en aquellas variables que no pudiera imputar de forma razonable con la mediana y con el knn sobre las variables ordenadas, pero seguía sin mejorar el valor de *accuracy*. Estas funciones pueden encontrarse en el código fuente pero no las comento porque las descarté rápidamente.

5. Clasificadores

Hasta el momento había estado trabajando con una configuración inicial del LGBM con 1000 estimadores, que era el punto dulce de estimadores que obtuve nada más afrontar el problema, más o menos estimadores repercutía negativamente, sin embargo desde entonces no había modificado estos valores. Los clasificadores más populares en el foro del problema eran Random Forest de Scikit-Learn que me daba un resultado ligeramente peor que LGBM y XGBoost, que arrojaba un resultado similar pero era insufriblemente lento y obstaculizaba la dinámica de desarrollo que llevaba hasta el momento.

5.1. Uso de dos clasificadores (Dual)

Tras aplicar la ordenación de etiquetas por sus estadísticas había conseguido colocarme por encima de los 300 mejores, y no conseguía mejorar el rendimiento con LGBM, así que resolví utilizar la decisión combinada de LGBM y Random Forest, que eran los más rápidos. Leí que Random Forest trabajaba mejor con *OneHot Encoding* así que sin eliminar las columnas anteriores, expandí 5 columnas que según la gráfica de importancia parecían ayudar a clasificar el problema y tenían un número reducido de etiquetas posibles por característica, estas son *lga*, *payment_type*, *source*, *waterpoint_type* y *basin*. En este momento tenía unas 140 columnas en total. Tras esto intenté podar características que se presentaban en la gráfica como poco útiles a la clasificación pero empeoraban ligeramente el resultado.

```

1  INTERSTING_SMALL = ['basin', 'lga', 'payment_type', 'source',
    ↪   'waterpoint_type']
2  from sklearn.preprocessing import OneHotEncoder

```

```

3 def oneHot(tr,tst):
4     keys = INTERSTING_SMALL
5     # Generamos el encoder
6     oh = OneHotEncoder()
7     # Lo entrenamos sobre train (no había variables en test que no
      ↪ estuvieran en train)
8     oh.fit(tr)
9
10    # Transformamos train y test y lo devolvemos añadido al resto
      ↪ de columnas
11    mini_tr = oh.transform(tr).toarray()
12    mini_tst = oh.transform(tr).toarray()
13
14    return \
15        pd.concat([tr,pd.DataFrame(mini_tr)],axis=1),\
16        pd.concat([tst,pd.DataFrame(mini_tst)],axis=1)

```

Ahora definí mi clasificador, posteriormente añadí la opción de devolver la matriz tridimensional de probabilidades para optimizar la ponderación de cada uno de ellos. El procedimiento es sencillo, entreno los clasificadores de forma individual y devuelvo la clase con mayor suma.

```

1 class CustomClassifier:
2     clfs = []
3     def __init__(self,*clf):
4         # Creamos la lista de clasificadores
5         self.clfs = clf
6
7     # fit: Entrenamos cada clasificador por separado
8     def fit(self,X,y):
9         for clf in self.clfs:
10             clf.fit(X,y)
11         return self
12
13    # predict: Obtenemos el máximo de la suma de las probabilidades
14    def predict(self,X,return_prob=False):
15        # Obtenemos las predicciones individuales
16        res = [clf.predict_proba(X) for clf in self.clfs]
17        # Obtenemos la lista de etiquetas
18        k = self.clfs[0].classes_
19
20        # Definimos una matriz vacía
21        z = np.zeros(res[0].shape)
22        # Sumamos cada probabilidad
23        for r in res:
24            z += r

```

```

24
25     # Tomamos el índice de la etiqueta más probable
26     pos = np.apply_along_axis(lambda x: np.where(x == x.max()),
    ↪     axis=1, arr=z).reshape(-1)
27
28     # Esta es la opción original
29     if not return_prob:
30         return list(map(lambda i: k[i], pos))
31     # Esto se añadió posteriormente (predicción + volumen de
    ↪     probabilidades)
32     else:
33         return list(map(lambda i: k[i], pos)),
    ↪     np.stack(res,axis=2)
34
35     # predict_proba: Devolvemos el volumen de probabilidades
    ↪     individuales
36     def predict_proba(self, X):
37         # Obtenemos las predicciones individuales
38         res = [clf.predict_proba(X) for clf in self.clfs]
39
40         # Definimos una matriz vacía
41         z = np.zeros(res[0].shape)
42         # Sumamos cada probabilidad
43         for r in res:
44             z += r
45
46         return z

```

Con este cambio, conseguí crear un clasificador más robusto que utilizaba las confianzas individuales para clasificar las instancias, el resultado fue pasar en train de 0.8099 a 0.8146. A la luz de estos resultados decidí añadir un tercer clasificador.

5.2. Uso de tres clasificadores (Triple)

Añadí un tercer clasificador, XGBooster, que estaba muy referenciado en el foro del problema, lo evalué y el resultado fue extraordinario, pasé en train de 0.8146 a 0.816 y en DD desde 0.8218 (Variables categóricas ordenadas) a 0.8255, esto me puso en la posición 57 en el momento de la subida y hasta el final del plazo para subir nuestras predicciones.

Viendo lo prometedor de esta solución vi que XGBooster tenía una opción sobre GPU, lo cual me permitiría acelerar el proceso de entrenamiento —en ese momento tardaba entre 500 y 800 segundos por entrenamiento y predicción debido al alto número de columnas—. Así que instalé CUDA, conseguí bajar cada etapa de fit/predict a unos 200 segundos y seguí avanzando. LGBM tam-

bién tiene opción sobre GPU pero utiliza OpenCL, no CUDA, y el tiempo de ejecución de este algoritmo era el menor de los tres, así que no me preocupé por el pues desde luego no era el cuello de botella.

6. Ajuste de hiperparámetros

6.1. Algoritmo genético para los hiperparámetros

Como el número de parámetros era de 10, era inasumible utilizar una rejilla de hiperparámetros, así que definí un algoritmo genético personalizado, muy orientado a la intensificación y poca exploración, ya que cada ejecución era larga y sólo necesitaba resultados ligeramente mejores para colocarme entre los cincuenta mejores. Definí un vector de 10 elementos limitado entre 0 y 1, que representaban los parámetros, y funciones personalizadas que los convertían a la horquilla —enteros o coma flotante— que definiera.

Con una población inicial de 6 —5 aleatorios y la configuración utilizada hasta el momento—, en cada iteración se eligen dos padres, se crean dos hijos cruzándolos con un valor comprendido entre los dos padres más un factor multiplicador de la horquilla de 1.15, mutan con una probabilidad del 0.1 —como son diez características siempre muta un valor con un número aleatorio entre 0 y 1— y se evalúan. Se añaden los dos a la población —en vez de añadir sólo el mejor—, se ordenan y se seleccionan de nuevo los seis mejores. Como se puede ver este algoritmo tiene una capacidad exploratoria muy reducida, casi determinada por la población aleatoria inicial, pero como las evaluaciones son muy costosas no quería encontrar el óptimo global sino un óptimo en la localidad de la solución actual. La condición de salida la puse a 200 evaluaciones pero es irrelevante porque nunca llega a evaluar 400 configuraciones, lo interrumpo manualmente.

Probé también con differential evolution, de la librería *scipy*, pero se estabilizaba en la primera evaluación y salía demasiado rápido. No adjunto el código del algoritmo porque es relativamente extenso pero se puede encontrar en el código fuente adjunto. En cada evaluación de una configuración utilizaba como estimador sólo las dos primeras evaluaciones del 5-fold CV, que eran un indicador pesimista —comprobado experimentalmente— del resultado de las cinco evaluaciones. Las horquillas de los parámetros las definí experimentalmente considerando la calidad de los resultados y el tiempo de ejecución. En el algoritmo cacheo las respuestas para cada configuración una vez transformada, pues valores decimales próximos arrojan la misma configuración.

```
1 def intbetween(min,max):
2     return lambda x: int(x*(max-min) + min)
3
4 def floatbetween(min,max):
5     return lambda x: x*(max-min) + min
6
```

```

7 # Dentro del train
8 def train(X,y):
9     ...
10    # XGBooster: max_depth, n_estimators
11    p1 = [intbetween(3,21),intbetween(50,800)]
12    # LGBM: boosting_type, num_leaves, max_depth, n_estimators
13    p2 =
14        ↳ [oneof(['gbdt','dart','goss']),intbetween(20,101),intbetween(5,100),intbetween(200,150)]
15    # Random Forest: n_estimators, max_depth, min_samples_split,
16        ↳ min_samples_leaf=
17    p3 =
18        ↳ [intbetween(100,1200),intbetween(10,31),intbetween(2,20),intbetween(1,20)]
19
20    p_transf = p1+p2+p3
21    ...
22
23 def buildCls(c, c_transf):
24     p = getParams(c, c_transf)
25     clf1 = xgb.XGBClassifier(n_estimators=p[1], max_depth=p[0],
26        ↳ n_jobs=NPROC,
27        ↳ objective='multi:softmax',tree_method='gpu_hist')
28     clf2 = lgb.LGBMClassifier(boosting_type=p[2], num_leaves=p[3],
29        ↳ max_depth=p[4], n_estimators=p[5],
30        ↳ objective='multiclass', n_jobs=NPROC)
31     clf3 = RandomForestClassifier(n_estimators=p[6],
32        ↳ max_depth=p[7], min_samples_split=p[8],
33        ↳ min_samples_leaf=p[9],
34        ↳ n_jobs=NPROC)
35
36     clf = CustomClassifier(clf1, clf2, clf3)
37
38     return clf

```

Sin embargo, pese a todo lo anterior no conseguí mejorar el resultado obtenido, aunque sobre el train sí, que conseguía pasar a 0.817 e incluso 0.818, a la hora de evaluarlo en DrivenData el resultado no era superior al obtenido anteriormente.

6.2. Evitar sobreajuste

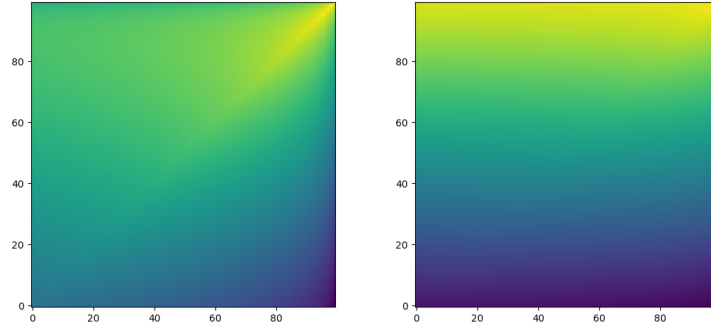
El siguiente paso fue reducir el ruido como ya comenté con anterioridad y el resultado no era prometedor. Con todo esto, al ver que el problema mejoraba sobre el train pero no en DD, sospeché que podía estar sobreajustando algunas características presentes en el train que no aparecían en el test, por eso definí una nueva fórmula para calcular el *fitness* del algoritmo genético, en la que

penalizaba el sobreajuste, pero seguía reflejando la medida de *accuracy* que utilizaban en la web, para ella utilizaba la media sobre la validación cruzada y la media de la evaluación sobre el total de los datos en cada paso de 5-fold —menos costoso que entrenar con el conjunto de datos para el algoritmo genético—. La gráfica de la fórmula tiene el siguiente aspecto, a la izquierda con $\alpha = 0,5$ y a la derecha 0,9 —el valor que mejor reflejaba el comportamiento que esperaba—. Siendo x_v —eje y en la gráfica— el valor sobre la validación y x_a —eje x en la gráfica— sobre el total de los datos.

$$x'_a = 1 - (x_a - x_v)/(1 - x_v)$$

$$x'_a = \frac{1}{x'_a} \quad \text{if } x'_a > 1 \quad \text{else } x'_a$$

$$f(x_v, x_a) = x_v * \alpha + \sqrt[5]{x'_a} * (1 - \alpha)$$



El resultado es una fórmula basada en la tasa de acierto que maximiza el valor en validación y penaliza el sobreajuste, ponderado experimentalmente. El resultado sin embargo, aunque parecía prometedor, pasaba de 0.816/0.9503 en el clasificador triple a 0.08172/0.8841 lo que me hizo estimar que había conseguido el objetivo, sin embargo, pasé en DD de 0.8255 a 0.8202.

Esta fue la última subida con el usuario creado para la competición con el nombre de *Guillermo_Gomez_UGR*.

6.3. Extra. Differential evolution para ponderación de clasificadores individuales

Como veía tan cerca aparecer entre los cincuenta primeros y seguía obsesionado con que el ajuste de hiperparámetros me permitiría subir esas décimas que me separaban del primero —0.0031— a la mañana siguiente me desperté pensando que había ajustado los hiperparámetros de los clasificadores individuales sin atender a las características de las predicciones de cada uno, LGBM era el más robusto, pero era posible que hubiera situaciones que RF y XGBoost se aliaran en su contra o que fueran sobreoptimistas. Así que se me ocurrió otra optimización sencilla para ajustar los clasificadores individuales, basada en dos parámetros —posteriormente probé con tres pero no suponían una mejora sobre dos—.

1. Una potencia que modificara los valores de probabilidad de cada clasificador de forma individual —elimina las probabilidades de suma uno pero sigue conservando la ordenación—, potencias por encima de 1 penalizan las confianzas pobres, y potencias entre 0 y 1 hacen a los clasificadores más optimistas.
2. Una suma ponderada de las distintas probabilidades para cada clase, hasta ahora la ponderación era simétrica ($\frac{1}{3}$ en nuestro caso), sin embargo es posible que no todos los clasificadores sean igual de interesantes para la predicción.

Modifiqué la función de validación cruzada para que devolviera el volumen de probabilidades, como el algoritmo de *Differential Evolution* de la librería *SciPy* permite definir las horquillas no tuve que codificar cada individuo de la población entre 0 y 1 y el código se simplificó bastante —llegados a este punto el código fuente es un campo de minas de modificaciones y comentarios—.

En cada evaluación del algoritmo genético, tras hacer la validación cruzada, obtengo la matriz de de probabilidades y busco los pesos óptimos para las características. Como esta parte quedaba fuera de la práctica el código es menos legible, lo adjunto por consistencia con el resto de la documentación. El exponente finalmente lo acoté entre 1 y 10, porque nunca elegía valores negativos y ralentizaba la ejecución; y el factor de ponderación entre 0 y 10, aunque igual habría dado entre 0 y 1 siempre que sea la misma horquilla para las tres.

```

1 # Obtener Accuracy ponderada
2 getAcc = lambda x,y: np.sum(np.where(x ==
    → np.repeat(np.max(x,axis=1),3,axis=0).reshape((-1,3))) [1] ==
    → y)/y.size
3 # Obtener lista de índices
4 getYIdx = lambda x: np.where(x ==
    → np.repeat(np.max(x,axis=1),3,axis=0).reshape((-1,3))) [1]
5 # Ponderar las probabilidades para cada matriz n_clasificadores x
    → n_instancias
6 pond = lambda y,pw,fc: np.power(y, np.repeat(pw,
    → y.shape[0]).reshape((len(pw), -1)).transpose()*np.repeat(fc,
    → y.shape[0]).reshape((len(fc), -1)).transpose())
7
8
9 # Accuracy ponderado con etiquetas codificadas entre 0 y 2
10 def getPondAcc(prob,y,pw,fc):
11     z = np.zeros(prob.shape)
12     # Lo aplico a cada matriz n_clasificadores x n_instancias por
    → no hacer más ilegible el código
13     for i in range(prob.shape[1]):
14         z[:, i, :] = pond(prob[:, i, :], pw, fc)
15     return getAcc(np.sum(z,axis=2),y)

```

```

16
17 # Accuracy ponderado con etiquetas de texto
18 def getPondAcc_text(prob,y,v):
19     pred_y = getPondRes(prob,v)
20     return np.sum(pred_y == y)/len(y)
21
22 # Predicción ponderada para aplicar al test
23 def getPondRes(prob,v):
24     k = ['functional', 'functional needs repair', 'non functional']
25
26     # Obtenemos las probabilidades modificadas
27     z = np.zeros(prob.shape)
28     for i in range(prob.shape[1]):
29         z[:, i, :] = pond(prob[:, i, :], v[:3], v[3:])
30     # Obtenemos los índices de las etiquetas más votadas
31     y_pred = getYIdx(np.sum(z, axis=2))
32     # Transformamos a etiquetas de texto
33     y_pred_lab = [k[i] for i in y_pred]
34
35     return y_pred_lab
36
37
38 from scipy.optimize import differential_evolution
39
40 # Aplica DE para obtener la mejor configuración
41 def bestWeight(prob_all, y):
42     # Transformamos la clase a numérica
43     k = ['functional', 'functional needs repair', 'non functional']
44     y_n = np.zeros(y.shape)
45     for i, ke in enumerate(k):
46         y_n[y == ke] = i
47
48     # Evaluación candidato
49     def eval_cad(v):
50         nonlocal prob_all, y_n
51         #v[:3] = potencias, v[3:] = ponderación
52         return getPondAcc(prob_all, y_n, v[:3], v[3:])
53
54     de = differential_evolution(
55         # Como el algoritmo minimiza una función multiplicamos por
56         ↪ -1 el valor de ajuste
57         lambda x: -eval_cad(x),
58         [(1,10)]*3+[(0,10)]*3,
59         # El que mejor resultado arroja
60         strategy='rand1exp',

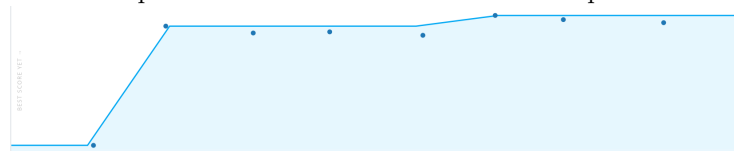
```

```

60     # Defino las iteraciones máximas a 100 aunque he comprobado
    ↪     que no las consume
61     maxiter=100
62 )
63 #Devolvemos el ajuste y la configuración que lo logra
64 return -de.fun, de.x

```

Y finalmente, tras ejecutarlo durante una noche, conseguí una configuración que me permitió obtener 0.8179/0.9312 frente a los 0.816/0.9503 de la mejor solución hasta el momento, y conseguir en DD una precisión de 0.8258, sólo 0.0005 mejor que mi anterior solución pero que me permitió colarme en el puesto 48 con el usuario *guillermogotre* —adjunto una captura para comprobar que no utilicé esta cuenta para realizar subidas mientras duró la práctica—.



Submissions

BEST	CURRENT RANK	# COMPETITORS	SUBS. TODAY
0.8258	48	6188	0 / 3

SUBMISSIONS

Score	Submitted by	Timestamp
0.8029	guillermogotre	2019-01-06 19:51:01 UTC
0.8239	guillermogotre	2019-01-06 23:58:48 UTC
0.8227	guillermogotre	2019-01-06 23:59:17 UTC
0.8229	guillermogotre	2019-01-07 16:56:05 UTC
0.8223	guillermogotre	2019-01-07 17:18:00 UTC
0.8258	guillermogotre	2019-01-08 00:05:14 UTC
0.8251	guillermogotre	2019-01-08 22:34:13 UTC
0.8245	guillermogotre	2019-01-08 23:26:31 UTC

7. Reflexiones finales

La práctica ha sido un ejercicio divertidísimo aunque me ha robado más tiempo del que me podía permitir. El principal elemento de conflicto ha sido el sobreajuste, llegados a este punto sigo sin poder utilizar esta información para estimar la calidad de un modelo, y es evidente que el resultado obtenido en la validación cruzada no es suficiente para explicar el rendimiento sobre el

conjunto de test. Se observa que es interesante cierta cantidad de sobreajuste para explicar las instancias del test —no he comprobado si hubiera instancias repetidas en ambos conjuntos—, aunque de ninguna manera es indicador de la calidad del modelo. Por otro lado, la aleatoriedad juega un papel crítico en los resultados, tanto en la formación de los conjuntos de validación como en el funcionamiento estocástico de los algoritmos de clasificación, y dos ejecuciones con dos semillas distintas con los mismos hiperparámetros pueden suponer una puntuación de 0.8258 o de 0.8230, que no son dramáticamente distintos pero en una competición con tantos candidatos puede suponer varios cientos de puestos de diferencia.

Por otro lado, es posible que mi discurrir durante la práctica no haya sido el correcto, cuando llegué a la evaluación del modelo con tres clasificadores distintos y conseguí llegar a la posición 57 asumí que una mejor configuración de los hiperparámetros me reservaría una plaza entre los mejores resultados, y ha sido así sólo relativamente, conseguí con esto escalar diez posiciones pero sigo estando muy lejos de la primera posición —después de dos días de pruebas sólo conseguir mejorar 0.0005, mientras que el primero está a una distancia de 0.0028, 6 veces más—. Se me presenta ahora evidente que la limitación está en los datos, y que quizás si hubiera invertido ese tiempo en mejor imputación de datos, o en analizar mejor las instancias que eliminaba para discriminar la que efectivamente eran ruido de aquellas que eran casos raros pero de interés predictivo podría haber obtenido una mejor puntuación, por ejemplo, al hacer la documentación y revisar la etapa de visualización, he visto que hay valores raros en *amount_tsh* que ignoré al principio porque no conocía aún el problema, y es posible que esta variable tenga interés predictivo suficiente para tratarla con mayor profundidad pese al primer acercamiento frustrado que tuve.

En definitiva, cuando te enfrentas a un problema a ciegas, donde las referencias no tienen contenido material y sólo vemos los resultados de otros, es difícil tener claro donde focalizar el esfuerzo, más aún cuando el tiempo con el que contamos es limitado, y nos vemos obligados a tomar decisiones intuitivas sobre la marcha, sin embargo estoy muy satisfecho, y no hay mejor manera de entrenar la intuición que con la práctica. Dentro de plazo conseguí quedar en la posición 57 y dos días después en la 48, dentro del podio que era donde me había propuesto llegar, todo esto con una experiencia formativa de un gran interés —y mucho más ameno que KNIME :) —.