

APRENDIZAJE AUTOMÁTICO (2019)
DOBLE GRADO EN INGENIERÍA INFORMÁTICA Y MATEMÁTICAS
UNIVERSIDAD DE GRANADA

Proyecto Final

Ignacio Aguilera Martos, Luis Balderas Ruiz
nacheteam@correo.ugr.es, luisbalderas@correo.ugr.es

5 de junio de 2019

Índice

1	APS Failure Scania Trucks Classification	3
1.1	Introducción	3
2	Preprocesamiento	3
2.1	Análisis explotario y valores perdidos	3
2.1.1	Desbalanceo de clases	4
2.1.2	Selección de características	6
2.1.3	Normalización	6
3	Modelos y clasificación	7
3.1	Gradiente descendente estocástico	7
3.1.1	AdaBoost	8
4	Bibliografía	9

Índice de figuras

2.1.	Representación de las instancias tras la imputación de VP con la mediana	4
2.2.	Matriz de confusión generada por SVM /RF/SGD/Red Neuronal hasta el momento	5
2.3.	Disposición del dataset tras SMOTE	5
3.1.	Matriz de confusión para SGD	7
3.2.	Curva ROC para SGD	8

Índice de tablas

2.1.	Matriz de confusión sin FA	6
2.2.	Matriz de confusión con FA	6
3.1.	Resumen de las medidas obtenidas con SGD	7
3.2.	Resumen de las medidas obtenidas con AdaBoost	8

1. APS Failure Scania Trucks Classification

1.1. Introducción

Nos encontramos ante un problema de clasificación binaria (clase positiva y negativa) en la que la clase positiva consiste en el fallo o mal funcionamiento de un componente específico del sistema APS para la compañía de camiones Scania. Por su parte, la clase negativa está formada por camiones con fallos en componentes que no están relacionado con el sistema APS. En lo que se refiere al dataset en sí, el conjunto de entrenamiento está formado por 60000 instancias y 171 características. El test tiene 16000 muestras. Sabemos que los nombres de los atributos han sido anonimizados por motivos de privacidad. En el presente documento explicamos con detalle cada uno de los pasos que hemos seguido en el estudio y diseño de un sistema inteligente para la clasificación con varios modelos, a saber, un modelo lineal (SGD), una Máquina de Soporte de Vectores, un modelo basado en Boosting (AdaBoost, concretamente) y Random Forest. Previo uso de los modelos, hemos hecho un profundo estudio de los datos que nos han llevado a modificarlos y tratarlos para mejorar los resultados, entre lo que destaca el tratamiento de valores perdidos con la imputación de la mediana, eliminación de estancias poco relevantes, oversampling con SMOTE para compensar el gran desbalanceo y análisis factorial para reducir la dimensionalidad en las características.

2. Preprocesamiento

2.1. Análisis exploratorio y valores perdidos

Nos encontramos ante un problema que tiene una gran cantidad de datos. En principio, este hecho es una ventaja, dado que la capacidad de generalización aumenta enormemente. Sin embargo, uno de los grandes retos ha sido, por un lado, el gran desbalanceo entre las clases y, por otro, la gran cantidad de valores perdidos e inconsistencias en los datos. Nos ocupamos de esto segundo en primer lugar. El análisis exploratorio de los datos nos hizo llegar a la conclusión de que, de entre las pocas instancias que forman parte de la clase positiva (esto es, con un fallo en el sistema APS), los valores de las variables estaban extraviados o, si no era así, más parecían outliers que otra cosa. En seguida nos dimos cuenta de que, desde un punto de vista físico tiene sentido, dado que la avería en un componente puede hacer que el comportamiento en los demás sea errático y, como consecuencia, puede generar observaciones de lo más variopintas. No obstante, nos dedicamos a eliminar los valores perdidos y le dimos validez total a los datos de los que disponemos.

En primer, vimos que había características que tenían más de un 70 % de sus valores como NA. Decidimos eliminarlas, ya que iban a entorpecer enormemente el trabajo de los modelos (en alguno de ellos, incluso son incompatibles). Eliminamos por tanto las características 2,75,76,77,78,79 y 113, por lo que pasamos de 171 a 164. En lo que se refiere a instancias, observamos que gran cantidad de ellas también estaban pobladas

de valores perdidos. Teniendo en cuenta el gran desbalanceo de las clases, eliminamos aquellas instancias con etiqueta negativa que tuvieran más de un 15 % de sus entradas como NA, de forma que pasamos de 60000 instancias a 55964.

A partir de ahí, nuestro conjunto de datos es más tratable. Sin embargo, aunque más dispersos, aún hay muchos valores perdidos con los que los modelos no pueden lidiar. Para solucionarlo, consultando la bibliografía, utilizamos dos tipos de imputación de valores perdidos: imputación de la mediana y de la media. La mediana nos ha dado mejores resultados en general, por lo que continuamos nuestro desarrollo con la mediana.

2.1.1. Desbalanceo de clases

Sin duda alguna, uno de los grandes problemas a los que nos enfrentamos cuando tratamos este conjunto de datos es el gran desbalanceo entre las clases. Hicimos una primera visualización de los datos observando el desbalanceo pero, a su vez, leve solapamiento (lo que a priori es una ventaja).

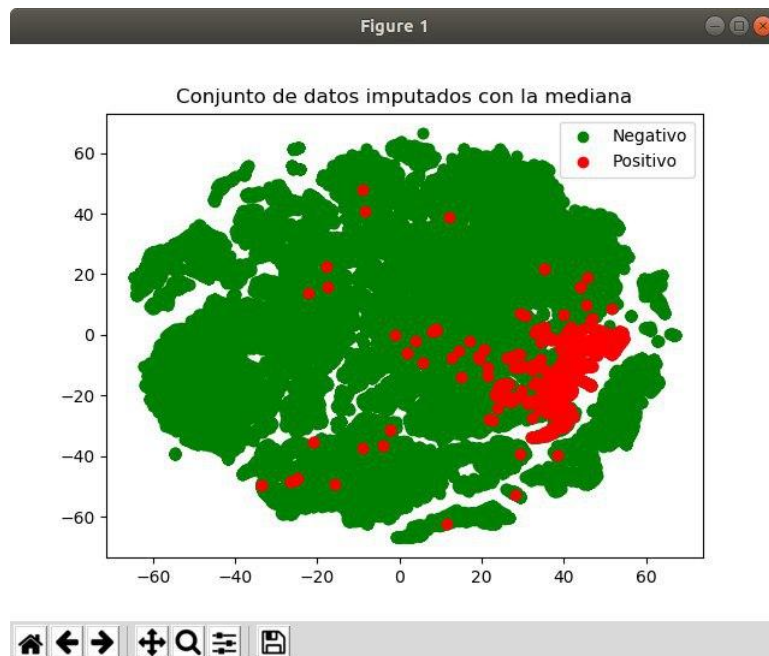


Figura 2.1: Representación de las instancias tras la imputación de VP con la mediana

Sin embargo, a pesar de la separación de las clases éramos incapaces de generar buenos resultados con los modelos y los clasificadores siempre etiquetaban las instancias de test como pertenecientes a la clase negativa. Esto daba lugar a un buen accuracy, pero las demás medidas que tenemos en cuenta, como son recall o f1-score eran nefastas.

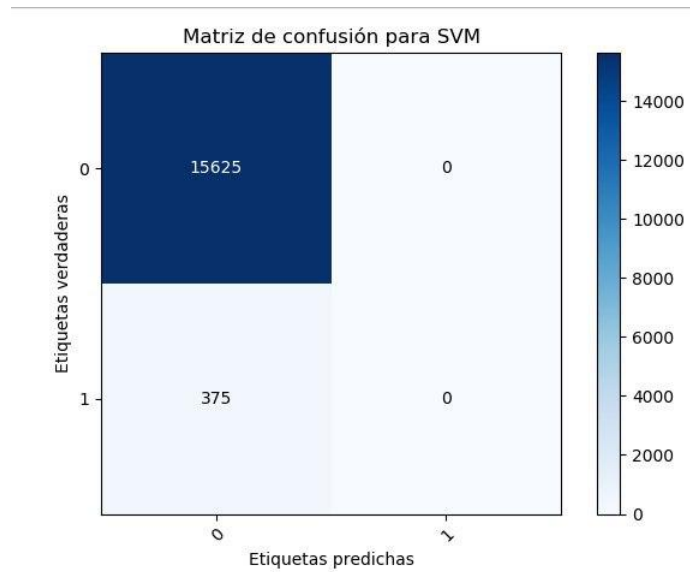


Figura 2.2: Matriz de confusión generada por SVM /RF/SGD/Red Neuronal hasta el momento

Por tanto, a pesar de la buena disposición de los datos, no se obtenían buenos resultados. A la vista de los mismos, nos decantamos por utilizar técnicas de oversampling para equilibrar el número de instancias de cada clase. En particular, elegimos SMOTE ([2]) como herramienta. Tras su aplicación, contamos finalmente con 109928 instancias de cada clase, apareciendo la siguiente disposición de datos:

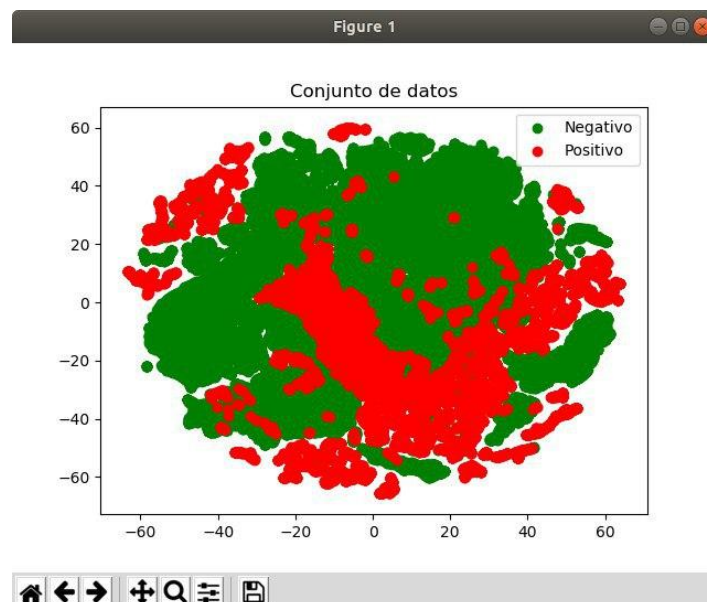


Figura 2.3: Disposición del dataset tras SMOTE

Como se puede apreciar, hemos perdido la localidad de las clases, aunque, como se verá, va a ser muy positivo.

2.1.2. Selección de características

Tras la eliminación de características por poseer demasiados valores perdidos, poseemos 164. No nos parece un número adecuado para utilizar PCA (más recomendado para cuando hay muchas más variables, ya que el sacrificio de la interpretabilidad es menos costoso) así que nos decantamos por otras técnicas de estadística multivariante, como es análisis factorial ([4],[5],[1], [3]). Análisis factorial sí nos ayuda a darle interpretabilidad al modelo y a encontrar qué características son las más determinantes. Reducimos así de 164 a 151. Sin embargo, los resultados en los modelos nos harán comprobar que, a pesar de mejorar el accuracy, se empeora la clasificación de las instancias de la clase positiva, como se puede ver en la siguiente comparación entre matrices de confusión para Random Forest:

Verd \ Pred	0	1
0	56160	2840
1	77	923

Tabla 2.1: Matriz de confusión sin FA

Verd \ Pred	0	1
0	57077	1923
1	113	887

Tabla 2.2: Matriz de confusión con FA

Por tanto, como nuestro objetivo es maximizar todas las medidas, especialmente recall (la clase minoritaria sigue siendo, de forma natural, la clase positiva), optamos por no aplicar de forma definitiva análisis factorial. Después de explorar estas dos técnicas, llegamos a la conclusión de que el tamaño del dataset hace que el número de características no sea muy elevado, por lo que acabamos por no seleccionar ningún subconjunto de ellas.

2.1.3. Normalización

En el análisis exploratorio observamos que el rango de las variables es absolutamente dispar. Todos los clasificadores basados en instancias (que incorporan la distancia) necesitan de un rango de variables equivalente para poder dar el peso adecuada a cada una. Por eso, realizamos una normalización estándar a los datos.

3. Modelos y clasificación

3.1. Gradiente descendente estocástico

El primer modelo que utilizamos es Gradiente Descendente Estocástico. Establecemos un número máximo de iteraciones de 10000 y una tolerancia de 1^{-6} . Se produce la convergencia tras 46 épocas con los siguientes resultados:

Score	0.97105
Precision	0.9883861
Recall	0.97105
F1-Score	0.97734082

Tabla 3.1: Resumen de las medidas obtenidas con SGD

Además, obtenemos la siguiente matriz de confusión:

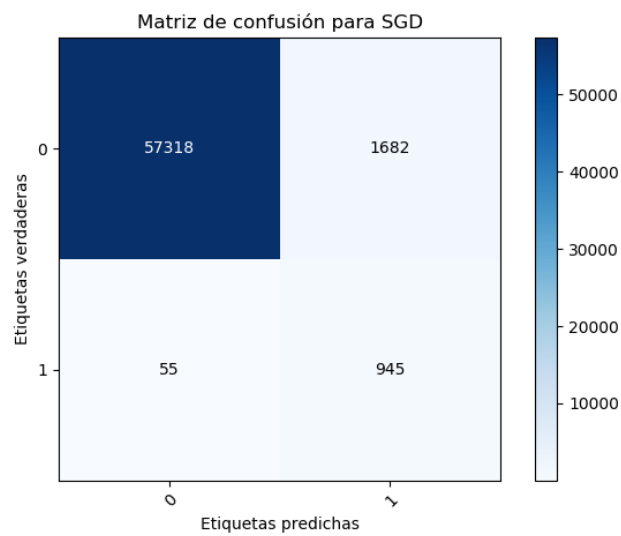


Figura 3.1: Matriz de confusión para SGD

y la siguiente curva ROC:

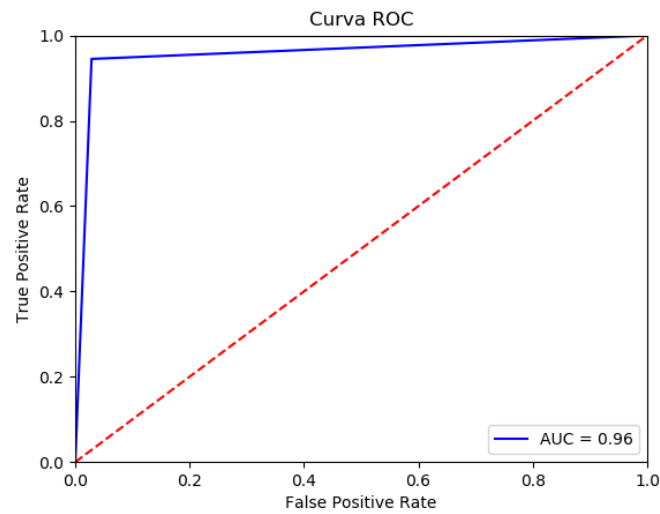


Figura 3.2: Curva ROC para SGD

Como se puede ver, los resultados son altamente satisfactorios. Además, con un problema que posee un volumen de datos tan grande, el tiempo de computación es un factor a tener en cuenta, y, en efecto, SGD apenas tarda un minuto en concluir.

3.1.1. AdaBoost

Score	0.9694
Precision	0.98717
Recall	0.9694
F1-Score	0.976038

Tabla 3.2: Resumen de las medidas obtenidas con AdaBoost

4. Bibliografía

Referencias

- [1] R.B. Cattell. Factor analysis. *New York: Harper*, 1952.
- [2] N.V. Chawla, K.W. Bowyer, L.O. Hall, and W.P. Kegelmeyer. SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 2002.
- [3] B. Fruchter. Introduction to Factor Analysis. *Van Nostrand*, 1954.
- [4] W.K. Härdle and Z. Hlávka. Multivariate Statistics. *Springer*, 2015.
- [5] L. Simar and W.K. Härdle. Applied Multivariate Statistical Analysis. *Springer*, 2015.