

Rank - Orange

Scoring Methods

- ☐ Information Gain
- ☒ Information Gain Ratio
- ☐ Gini Decrease
- ☐ ANOVA
- ☐ χ^2
- ☐ ReliefF
- ☐ FCBF

Select Attributes

- ☐ None
- ☐ All
- ☒ Manual
- ☐ Best ranked: 5

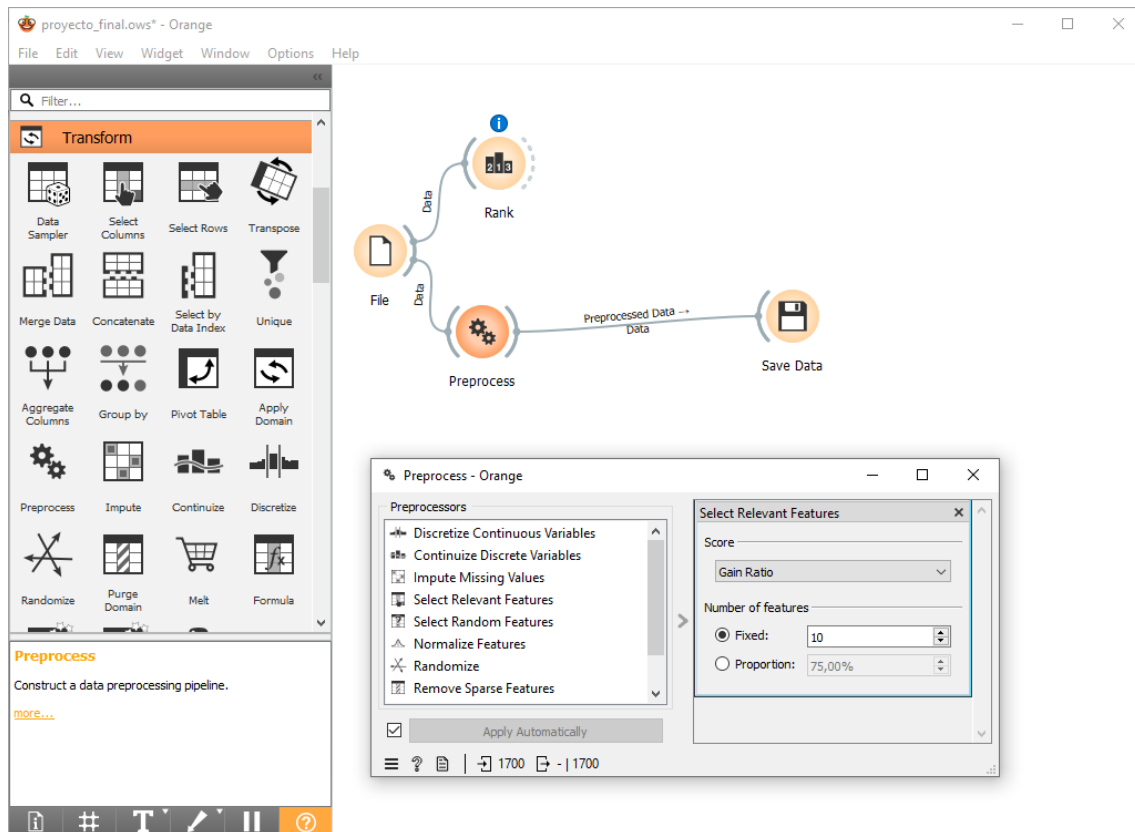
☒ Send Automatically

		#	Gain ratio
1	C n_p_ecg_p_05	2	0.032
2	N GB		0.030
3	C endocr_02	2	0.023
4	N DLIT_AG		0.012
5	C zab_leg_04	2	0.009
6	C fibr_ter_07	2	0.009
7	C RAZRIV	2	0.009
8	C PREDS_TAH	2	0.008
9	C n_r_ecg_p_05	2	0.008
10	C Na_BLOOD	2	0.007
11	C nr04	2	0.007
12	C n_p_ecg_p_11	2	0.007
13	C endocr_01	2	0.007
14	C nr02	2	0.006

1700 | - | 123 | -

Missing values will be imputed as needed.

De aquí cogemos hasta Na_BLOOD (incluido).



Y preprocesamos el dataset cortado con las variables más influyentes. (Datos ausentes)

```
import pandas as pd

df = pd.read_csv('hipertension_cortada.csv')

for col in df.columns:
    if df[col].isnull().any():
        moda = df[col].mode()[0]
        df[col].fillna(modas, inplace=True)

df.to_csv('hipertension_preprocesado.csv', index=False)
```

El dataset quedaría de la siguiente forma:

```
SIM_GIPERT,n_p_ecg_p_05,GB,endocr_02,DLIT_AG,zab_leg_04,fibr_ter_07,RAZRIV,PREDs_TAH,n_r_ecg_p_05,Na_BLOOD
0.0,0.0,3.0,0.0,7.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,1.0,0.0
0.0,0.0,2.0,0.0,2.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,2.0,0.0,3.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,3.0,0.0,7.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,2.0,0.0,7.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,2.0,0.0,7.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,2.0,0.0,6.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,3.0,0.0,6.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,2.0,0.0,3.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,2.0,0.0,2.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,2.0,0.0,6.0,0.0,0.0,0.0,0.0,0.0
0.0,0.0,2.0,1.0,1.0,0.0,0.0,0.0,1.0,0.0
0.0,0.0,2.0,0.0,6.0,0.0,0.0,0.0,0.0,1.0
```

Lo metemos en Weka.

hipertension_preprocesado.csv												
Relation: hipertension_preprocesado												
No.	1: SIM_GIPERT Numeric	2: n_p_ecg_p_05 Numeric	3: GB Numeric	4: endocr_02 Numeric	5: DLIT_AG Numeric	6: zab_leg_04 Numeric	7: fibr_ter_07 Numeric	8: RAZRIV Numeric	9: PRED5_TAH Numeric	10: n_r_ecg_p_05 Numeric	11: Na_BLOOD Numeric	
1	0.0	0.0	3.0	0.0	7.0	0.0	0.0	0.0	0.0	0.0	0.0	
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	
3	0.0	0.0	2.0	0.0	2.0	0.0	0.0	0.0	0.0	0.0	0.0	
4	0.0	0.0	2.0	0.0	3.0	0.0	0.0	0.0	0.0	0.0	0.0	
5	0.0	0.0	3.0	0.0	7.0	0.0	0.0	0.0	0.0	0.0	0.0	
6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
7	0.0	0.0	2.0	0.0	7.0	0.0	0.0	0.0	0.0	0.0	0.0	
8	0.0	0.0	2.0	0.0	7.0	0.0	0.0	0.0	0.0	0.0	0.0	
9	0.0	0.0	2.0	0.0	6.0	0.0	0.0	0.0	0.0	0.0	0.0	
10	0.0	0.0	3.0	0.0	6.0	0.0	0.0	0.0	0.0	0.0	0.0	
11	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
12	0.0	0.0	2.0	0.0	3.0	0.0	0.0	0.0	0.0	0.0	0.0	
13	0.0	0.0	2.0	0.0	2.0	0.0	0.0	0.0	0.0	0.0	0.0	
14	0.0	0.0	2.0	0.0	6.0	0.0	0.0	0.0	0.0	0.0	0.0	
15	0.0	0.0	2.0	1.0	1.0	0.0	0.0	0.0	1.0	0.0	0.0	
16	0.0	0.0	2.0	0.0	6.0	0.0	0.0	0.0	0.0	0.0	1.0	
17	0.0	0.0	2.0	0.0	7.0	0.0	0.0	0.0	0.0	0.0	0.0	
18	0.0	0.0	2.0	0.0	7.0	0.0	0.0	0.0	0.0	0.0	0.0	
19	0.0	0.0	2.0	0.0	2.0	0.0	0.0	0.0	0.0	0.0	0.0	
20	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	

Sacamos los outliers y extreme values.

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter: Choose InterquartileRange -R first-last -O 3.0 -E 6.0 Apply Stop

Current relation: Relation: hipertension_preprocesado-weka.filters.unsupervise... Attributes: 13 Sum of weights: 1700 Instances: 1700

Attributes: All None Invert Pattern

No.	Name
1	<input type="checkbox"/> n_p_ecg_p_05
2	<input type="checkbox"/> GB
3	<input type="checkbox"/> endocr_02
4	<input type="checkbox"/> DLIT_AG
5	<input type="checkbox"/> zab_leg_04
6	<input type="checkbox"/> fibr_ter_07
7	<input type="checkbox"/> RAZRIV
8	<input type="checkbox"/> PRED5_TAH
9	<input type="checkbox"/> n_r_ecg_p_05
10	<input type="checkbox"/> Na_BLOOD
11	<input type="checkbox"/> SIM_GIPERT
12	<input type="checkbox"/> Outlier
13	<input checked="" type="checkbox"/> ExtremeValue

Remove

Selected attribute: Name: ExtremeValue Missing: 0 (0%) Distinct: 2 Type: Nominal Unique: 0 (0%)

No.	Label	Count	Weight
1	no	1474	1474
2	yes	226	226

Class: ExtremeValue (Nom) Visualize All

Status: OK Log x 0

Los eliminamos.

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter: Choose **RemoveWithValues** -S 0.0 -C last -L last Apply Stop

Current relation
Relation: hipertension_preprocesado-weka.filters.unsupervise... Attributes: 13
Instances: 1474 Sum of weights: 1474

Attributes

All None Invert Pattern

No.	Name
1	<input type="checkbox"/> n_p_ecg_p_05
2	<input type="checkbox"/> GB
3	<input type="checkbox"/> endocr_02
4	<input type="checkbox"/> DLIT_AG
5	<input type="checkbox"/> zab_leg_04
6	<input type="checkbox"/> fibr_ter_07
7	<input type="checkbox"/> RAZRIV
8	<input type="checkbox"/> PREDS_TAH
9	<input type="checkbox"/> n_r_ecg_p_05
10	<input type="checkbox"/> Na_BLOOD
11	<input type="checkbox"/> SIM_GIPERT
12	<input type="checkbox"/> Outlier
13	<input checked="" type="checkbox"/> ExtremeValue

Remove

Selected attribute
Name: ExtremeValue
Missing: 0 (0%) Distinct: 1 Type: Nominal
Unique: 0 (0%)

No.	Label	Count	Weight
1	no	1474	1474
2	yes	0	0

Class: ExtremeValue (Nom) Visualize All

Status OK Log x 0

Normalizamos los datos.

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter: Choose **Normalize** -S 1.0 -T 0.0 Apply Stop

Current relation
Relation: hipertension_preprocesado-weka.filters.unsupervise... Attributes: 11
Instances: 1474 Sum of weights: 1474

Attributes

All None Invert Pattern

No.	Name
1	<input type="checkbox"/> n_p_ecg_p_05
2	<input type="checkbox"/> GB
3	<input type="checkbox"/> endocr_02
4	<input type="checkbox"/> DLIT_AG
5	<input type="checkbox"/> zab_leg_04
6	<input type="checkbox"/> fibr_ter_07
7	<input type="checkbox"/> RAZRIV
8	<input type="checkbox"/> PREDS_TAH
9	<input type="checkbox"/> n_r_ecg_p_05
10	<input type="checkbox"/> Na_BLOOD
11	<input checked="" type="checkbox"/> SIM_GIPERT

Remove

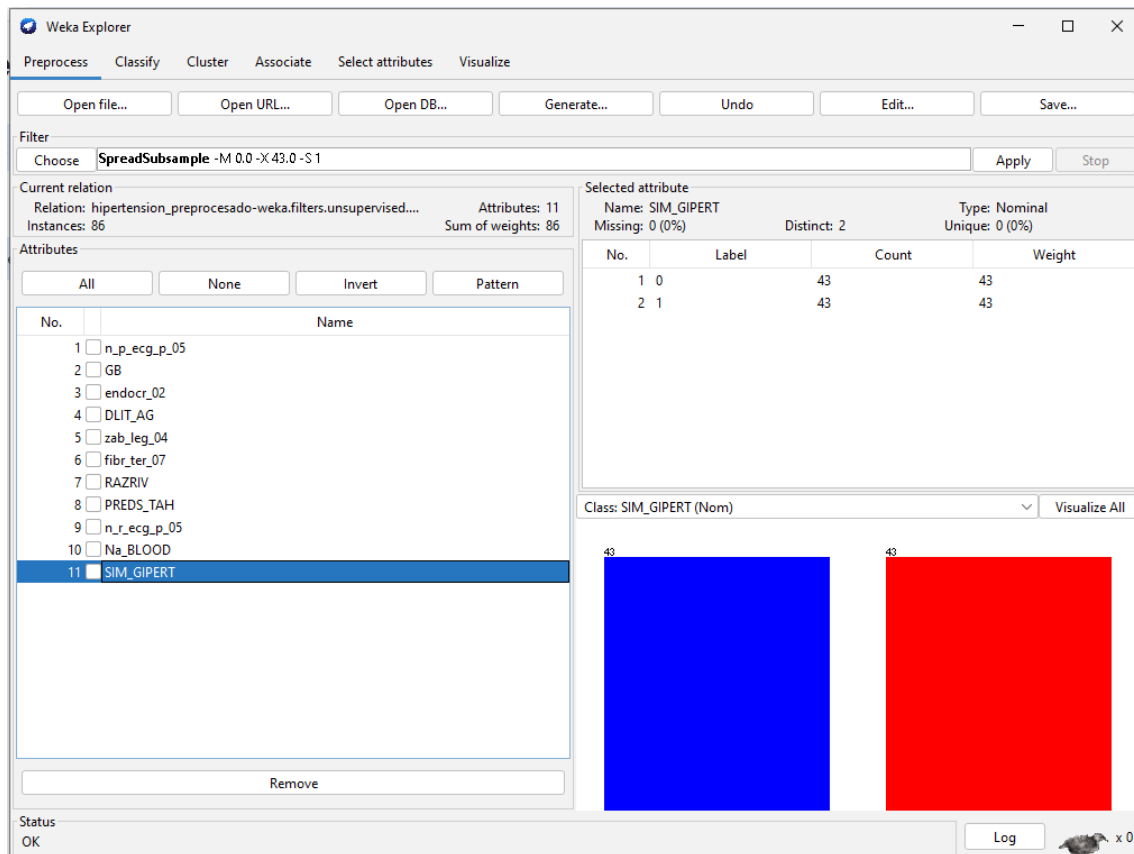
Selected attribute
Name: SIM_GIPERT
Missing: 0 (0%) Distinct: 2 Type: Numeric
Unique: 0 (0%)

Statistic	Value
Minimum	0
Maximum	1
Mean	0.029
StdDev	0.168

Class: SIM_GIPERT (Num) Visualize All

Status OK Log x 0

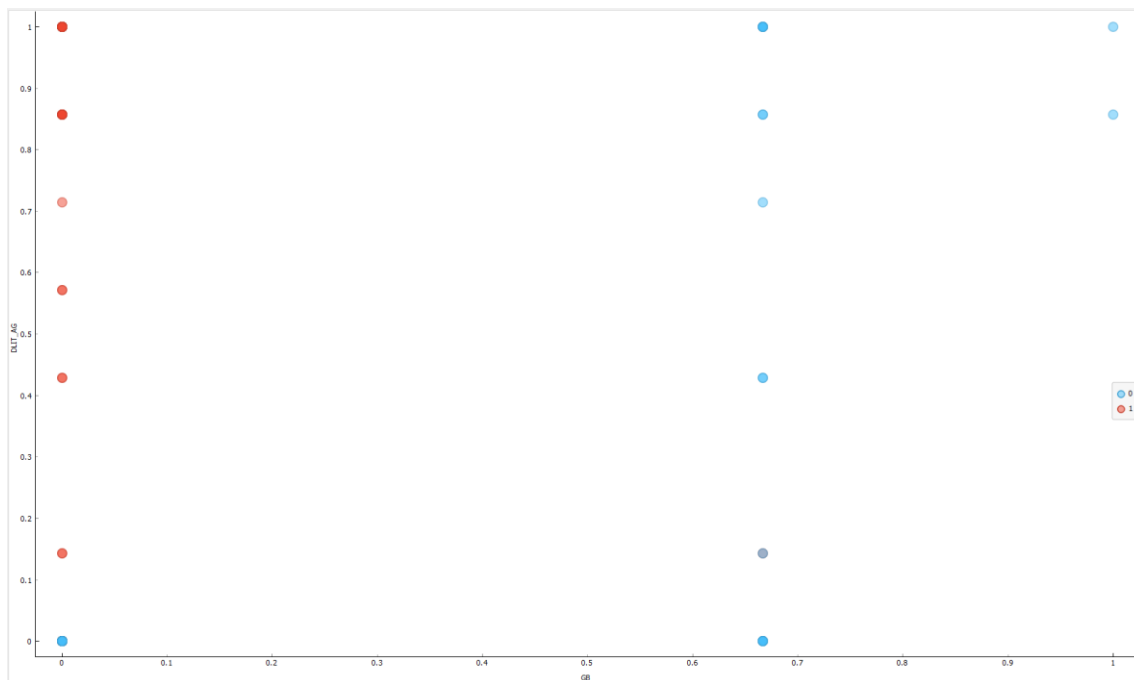
Convertimos la clase target a nominal para poder balancear los datos y los balanceamos.



Visualización de datos (diagramas dispersión, histogramas, barras) de todas las variables.

❖ Diagramas de dispersión

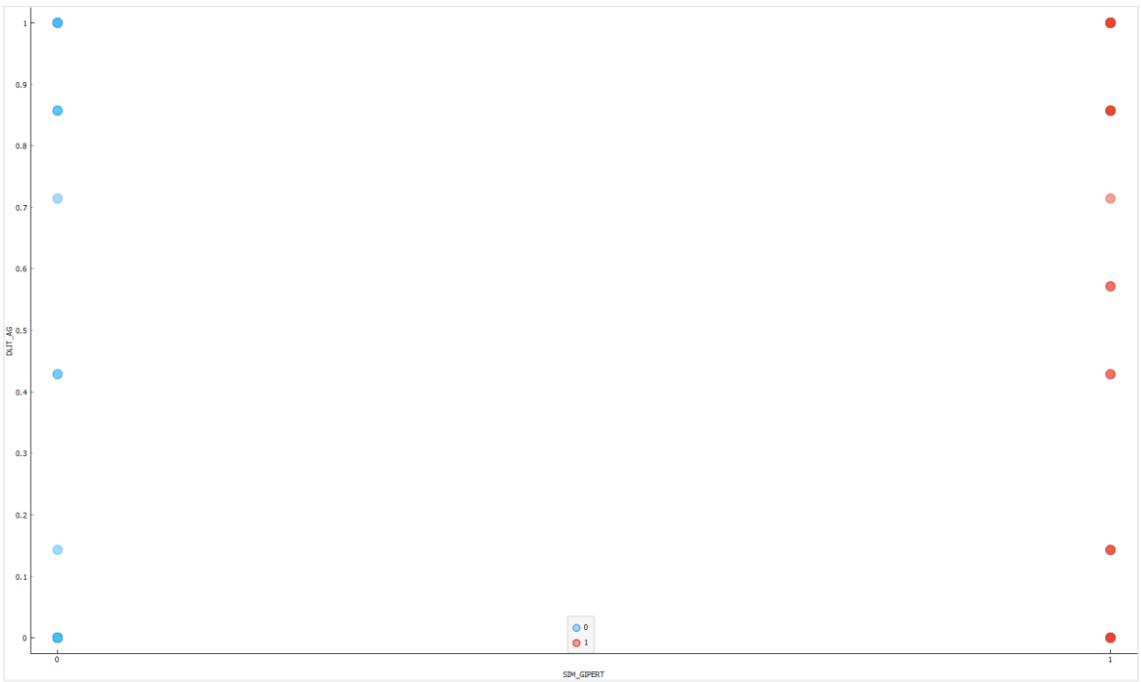
x: GB, y: DLIT_AG



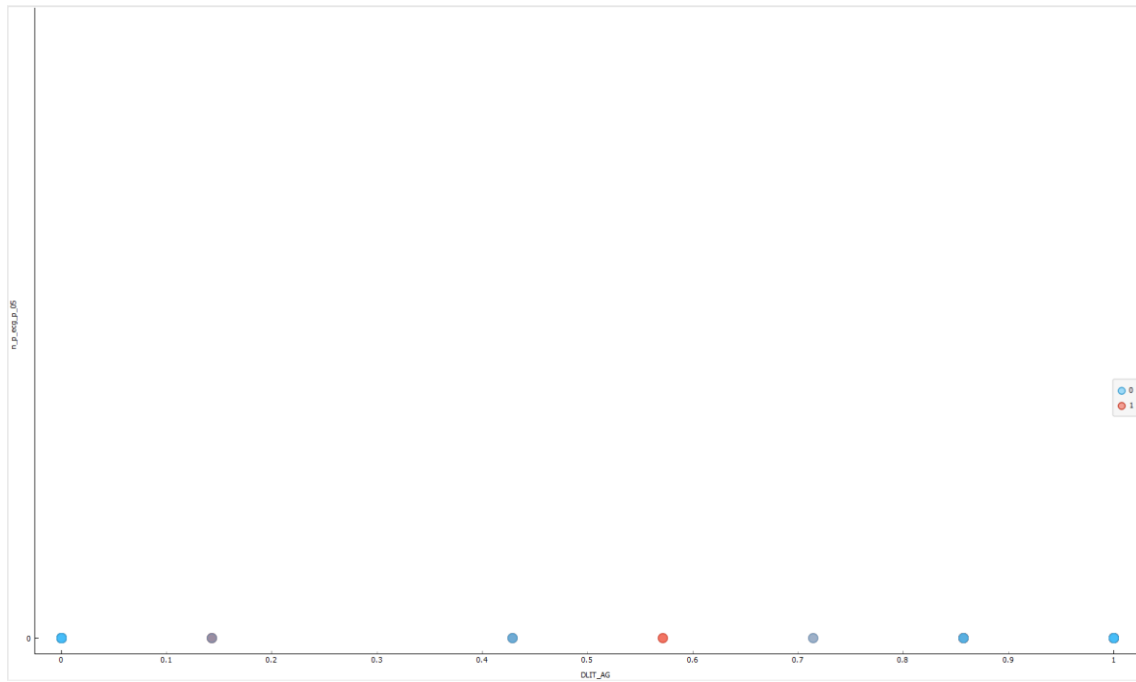
x: n_p_ecg_p_05, y: GB



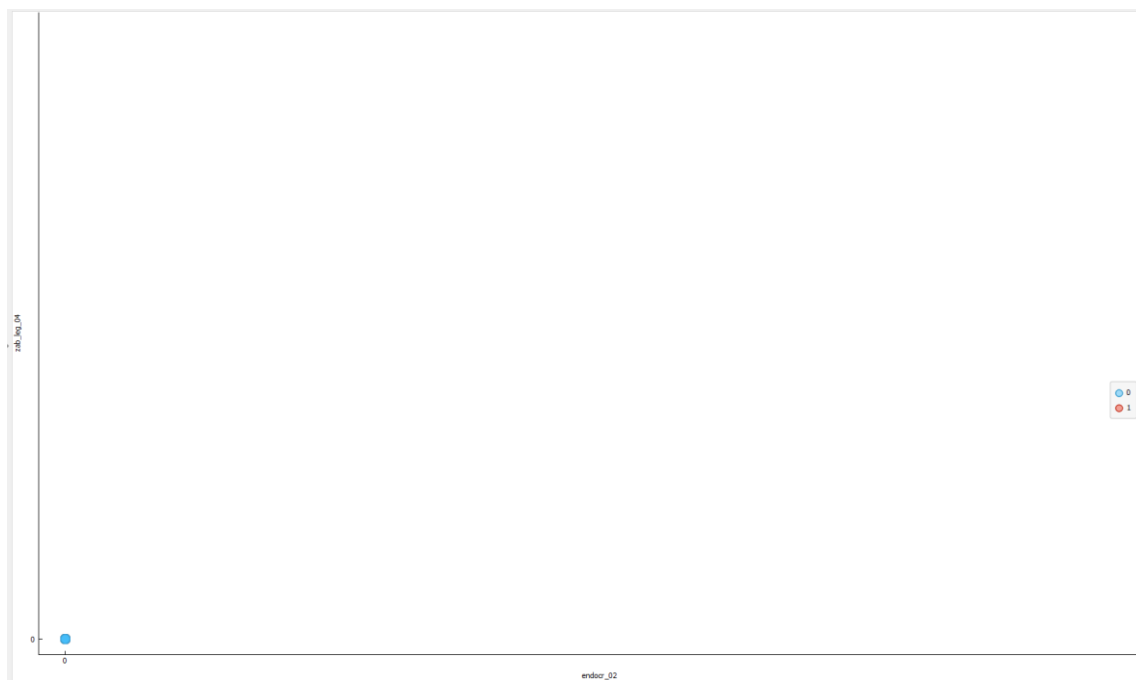
x: SIM_GIMPERT, y: DLIT_AG



x: DLIT_AG, y: n_p_ecg_p_05

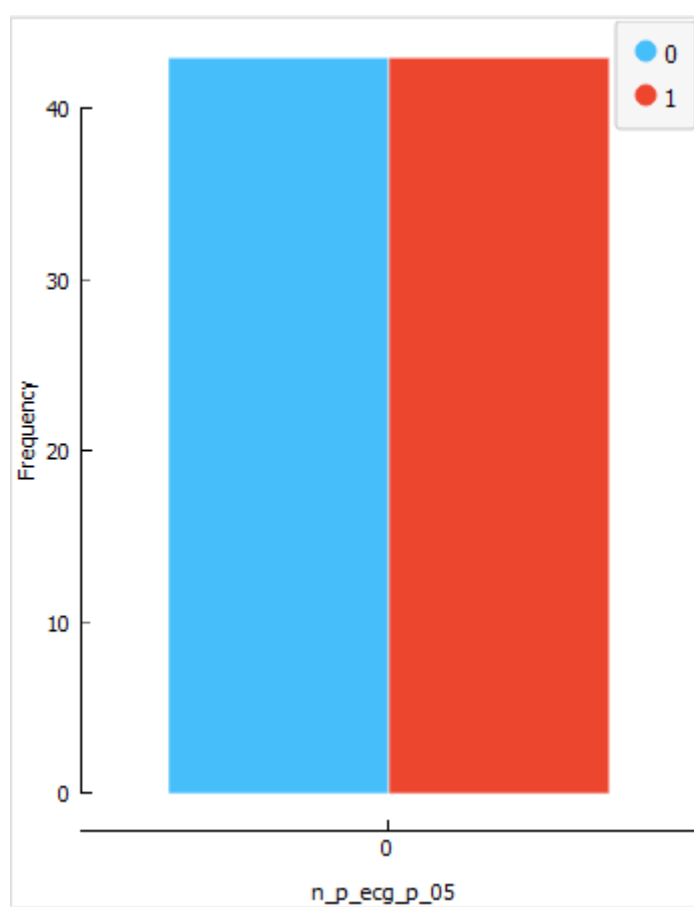
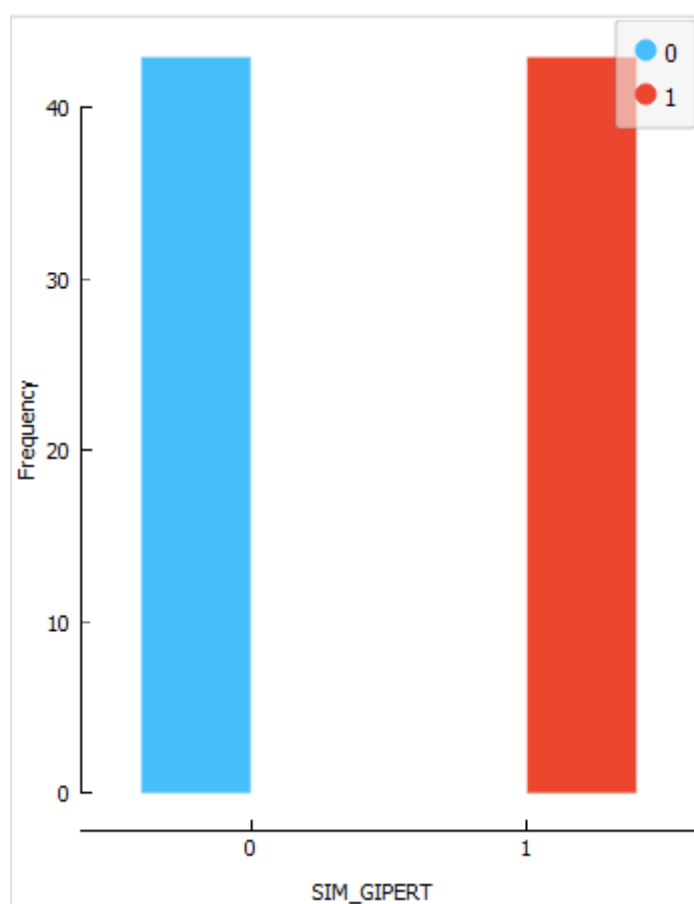


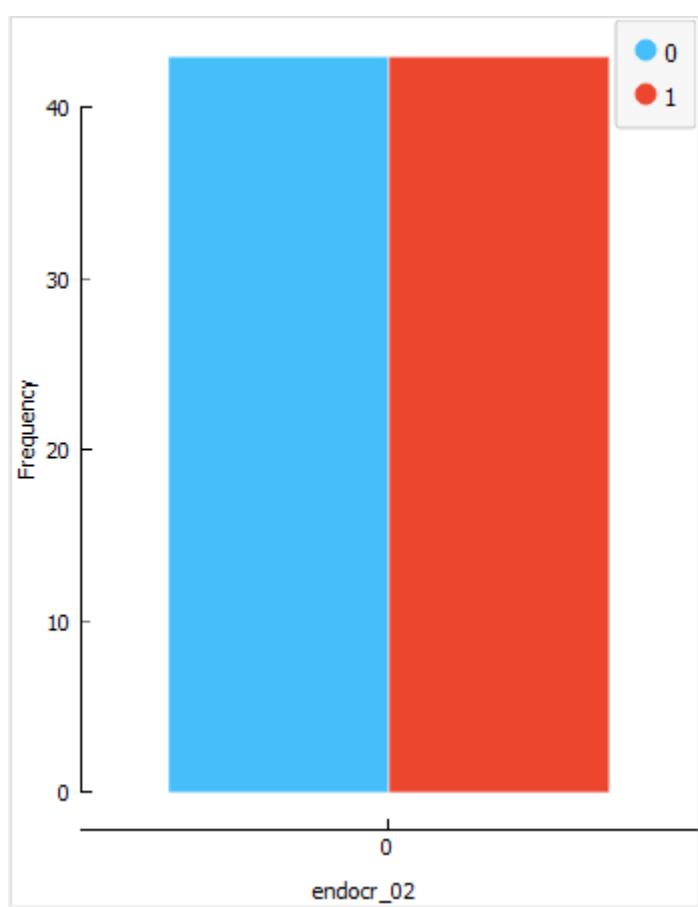
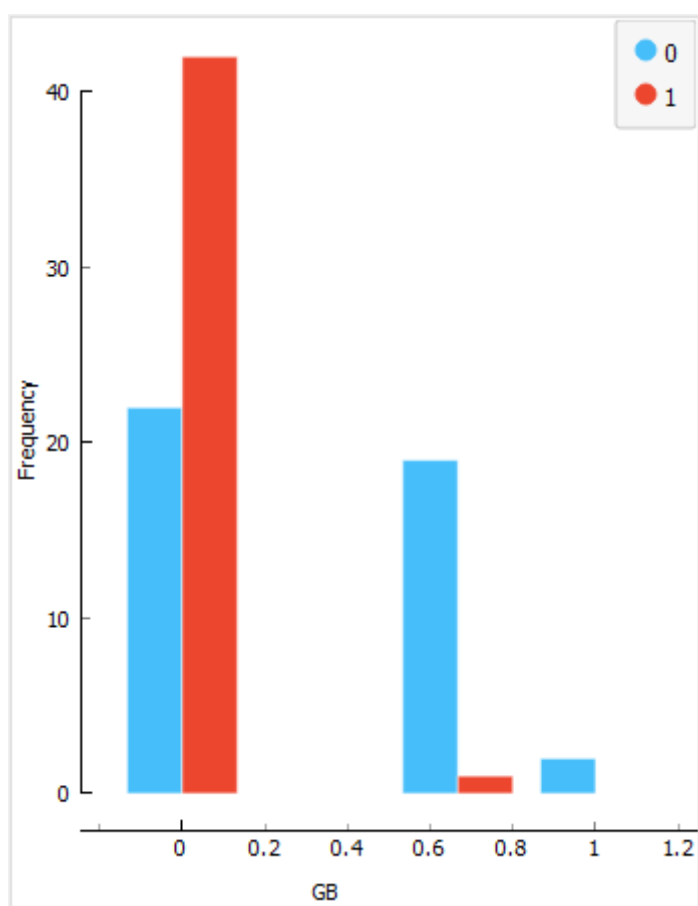
x: endocr_02, y: zab_leg_04

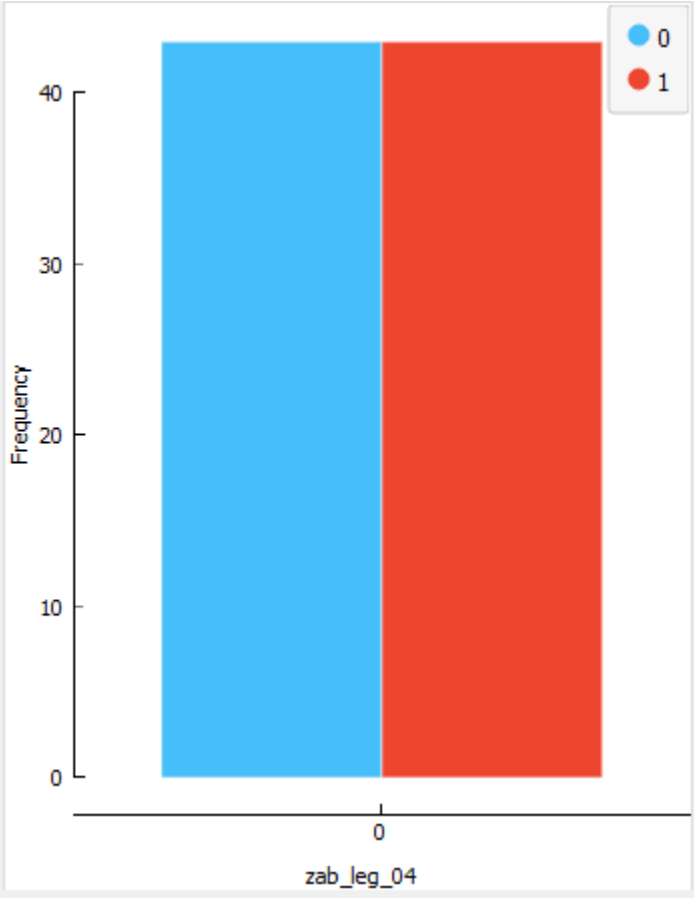
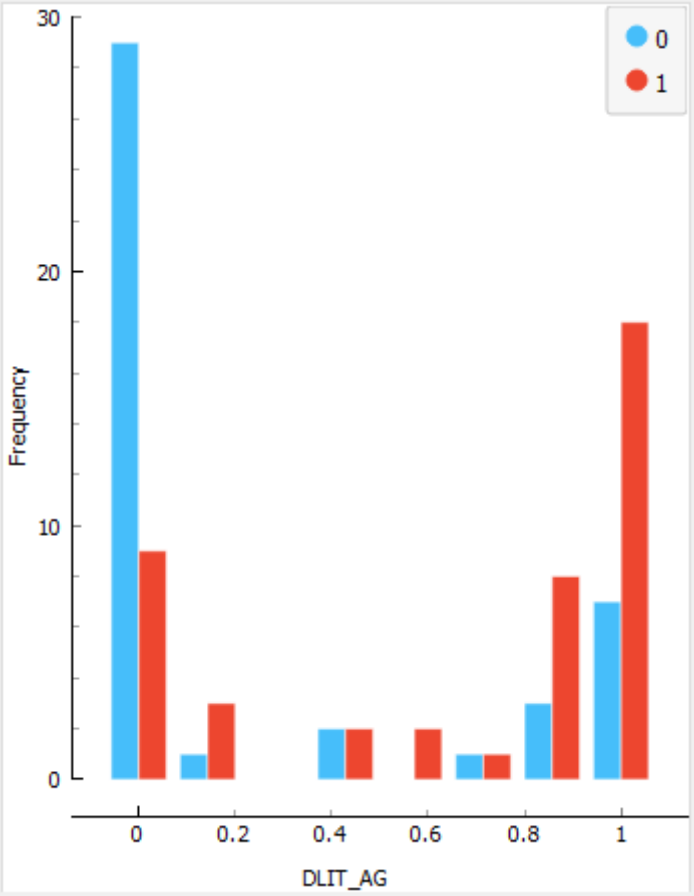


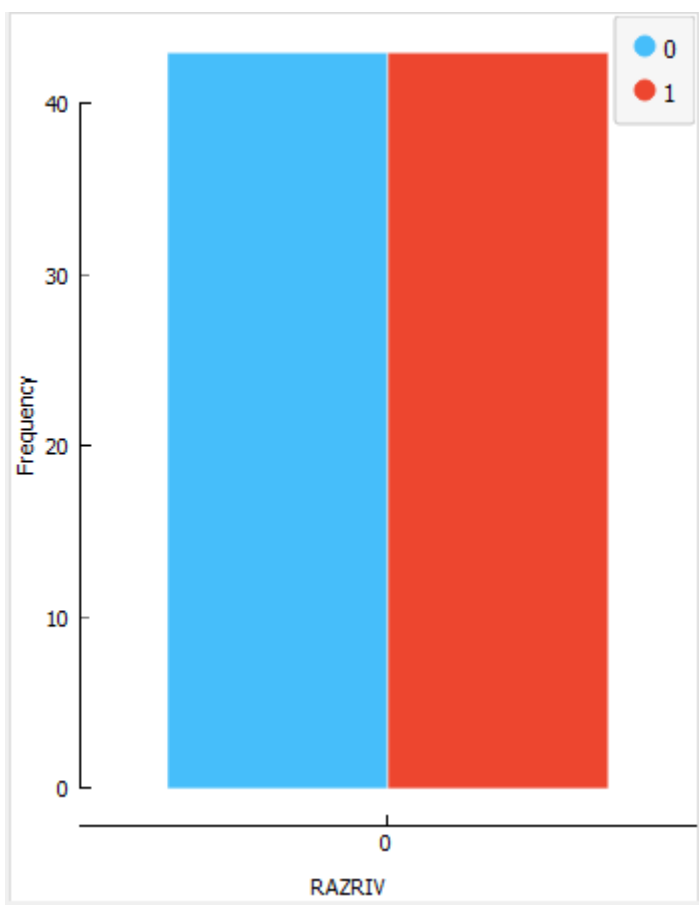
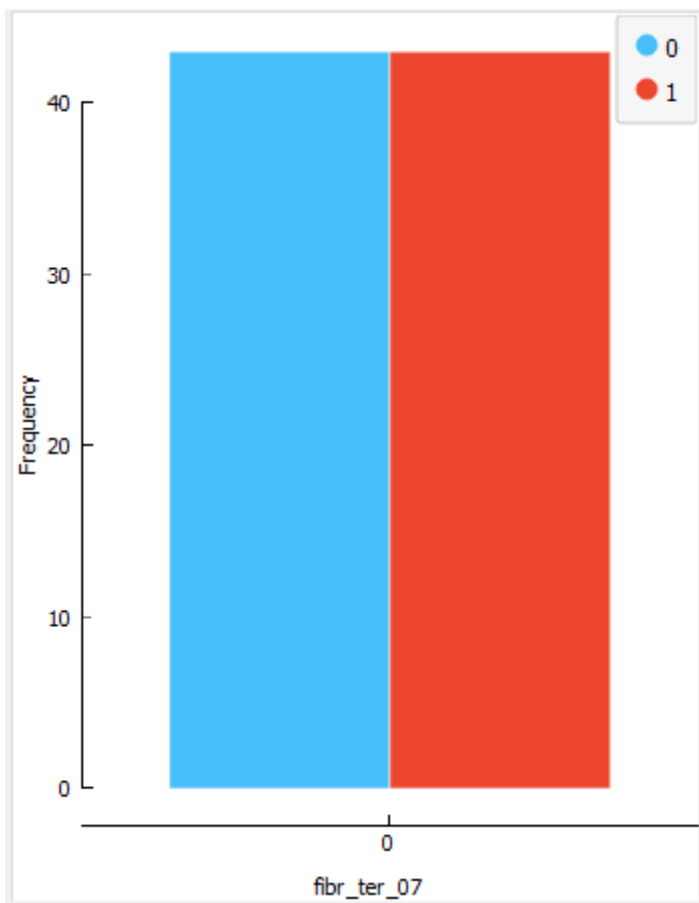
El resto de los gráficos no aportan información útil y son muy parecidos a este último, por lo que hemos decidido no ponerlos.

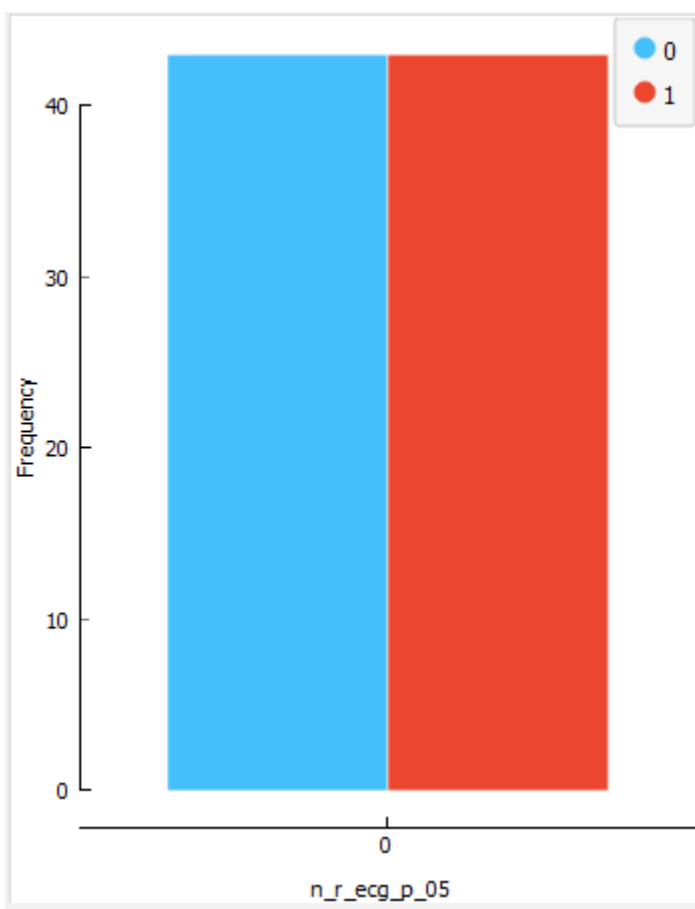
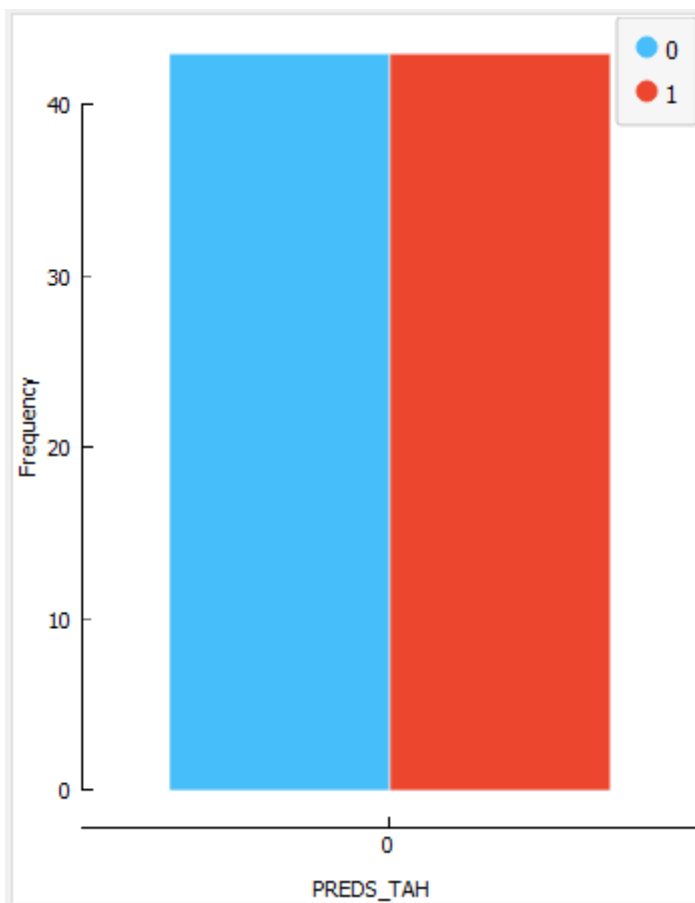
❖ Histogramas

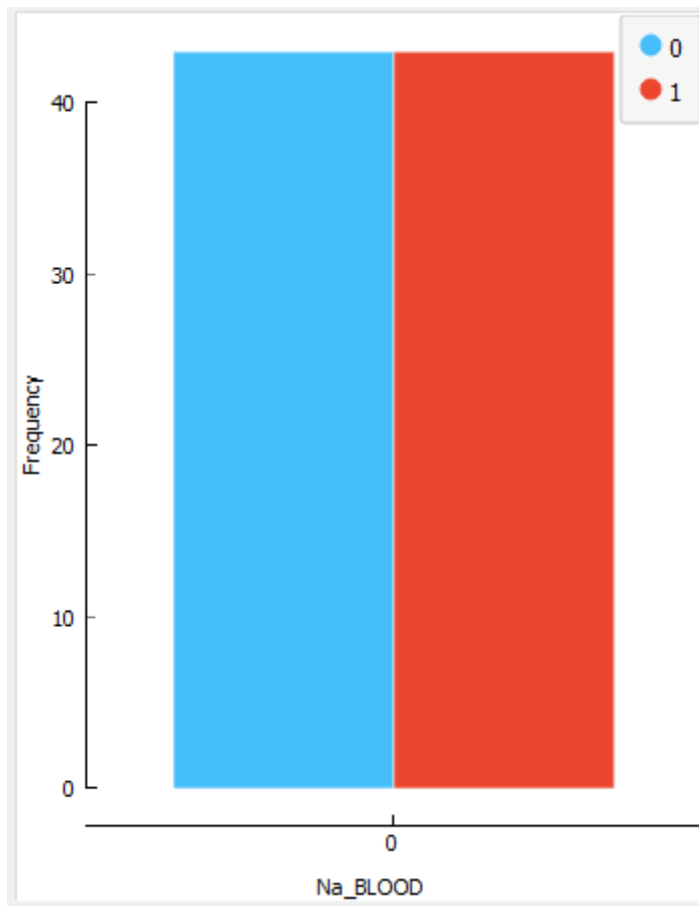




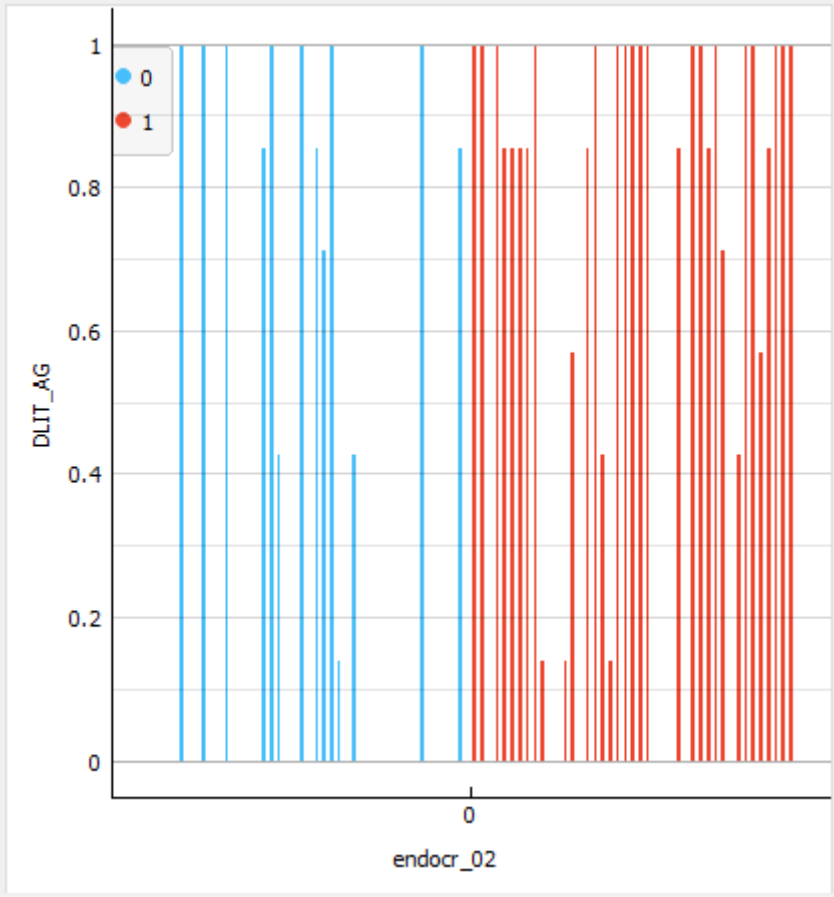
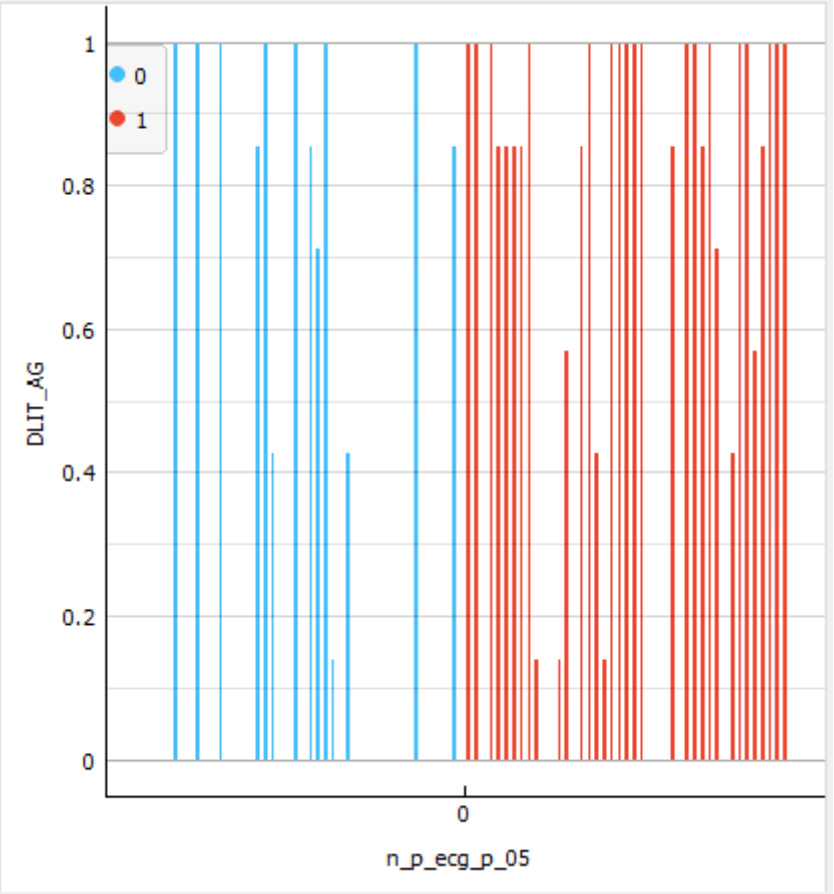


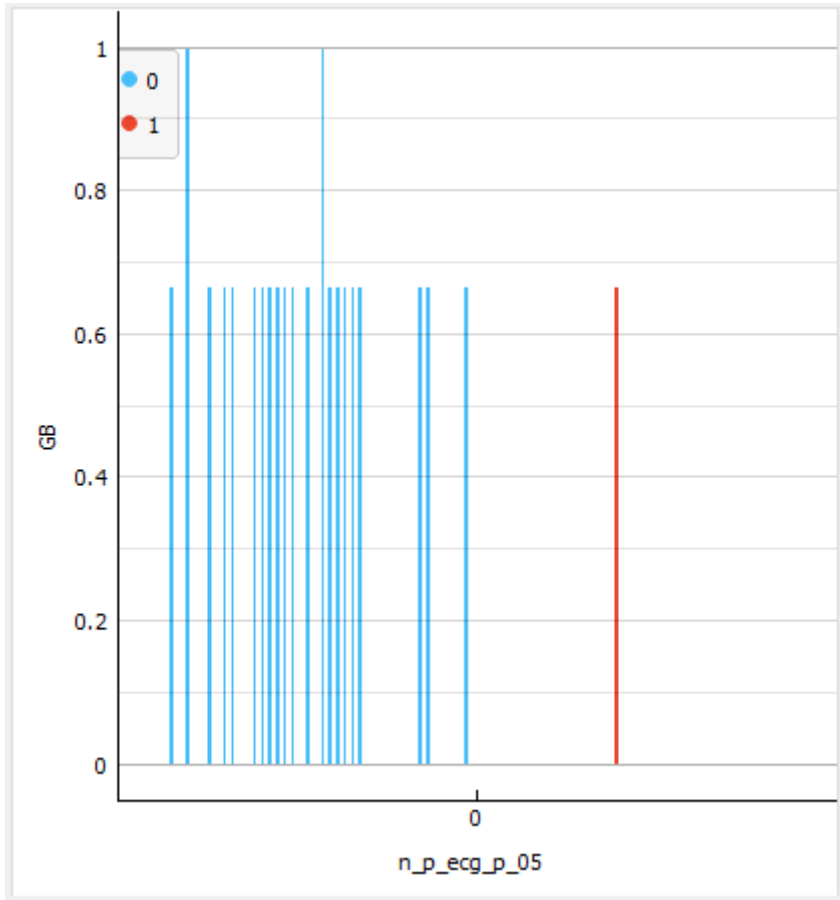
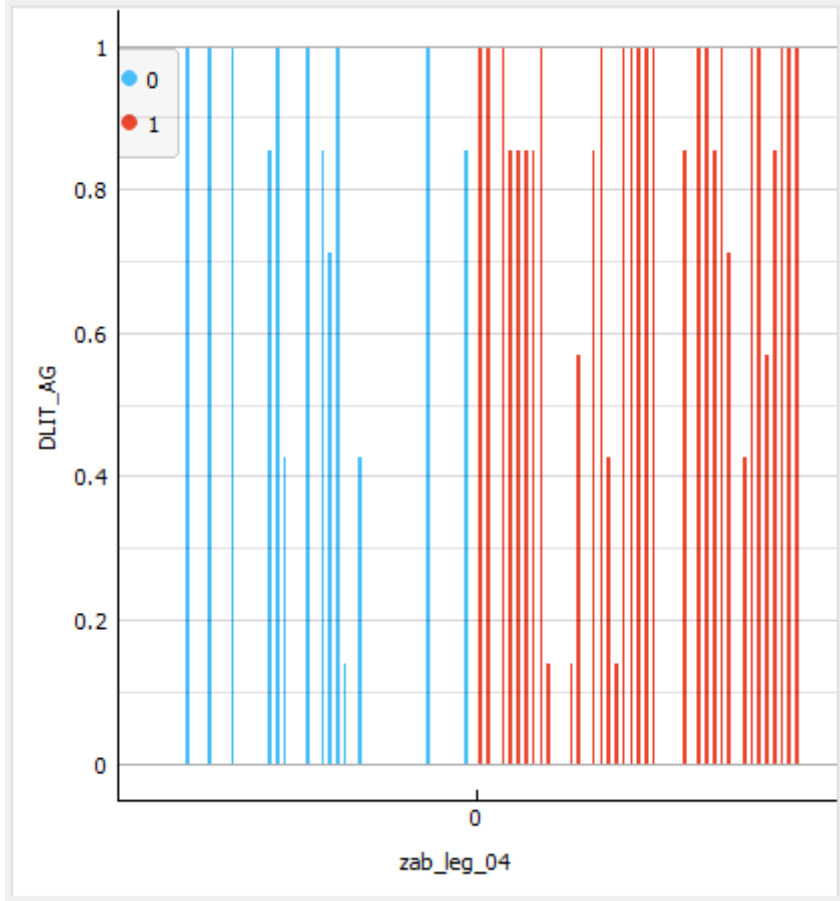


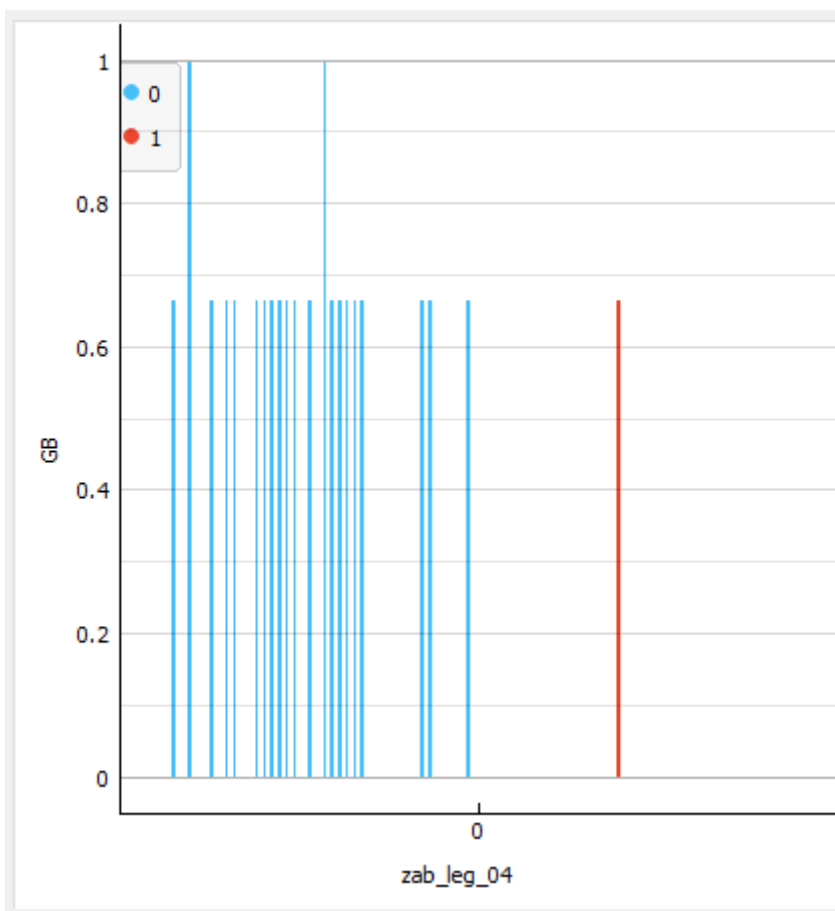
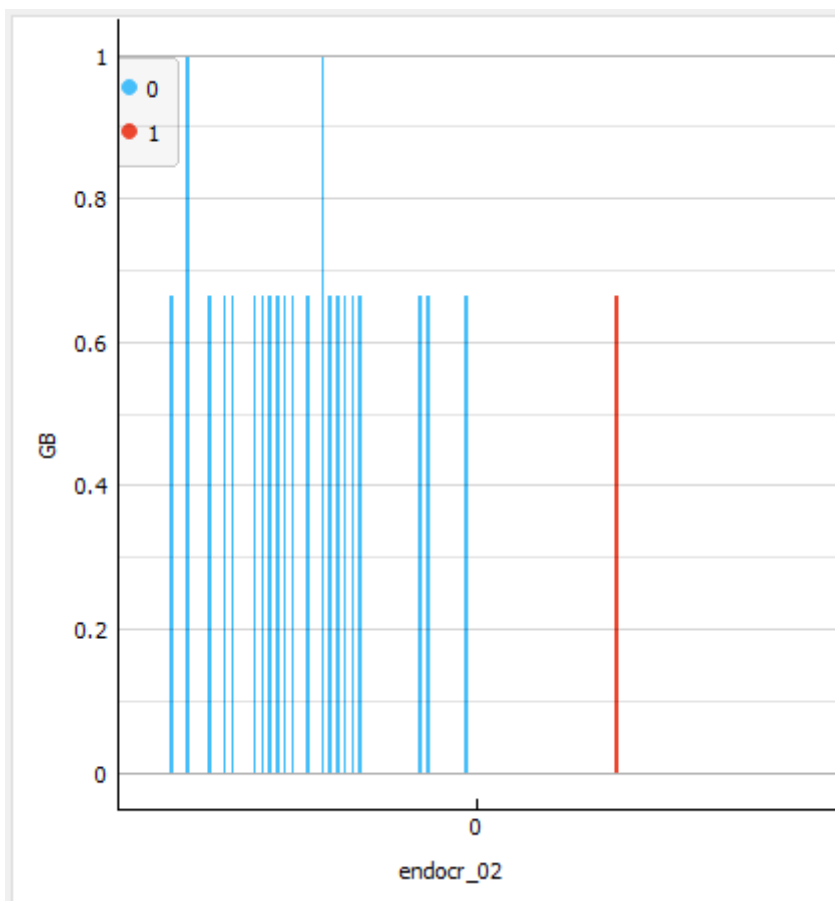




❖ Diagramas de barras







Al igual que ocurre con los diagramas de dispersión, los diagramas de barras restantes no nos aportan información distinta de los que ya tenemos por lo que, de nuevo, no los incluimos.

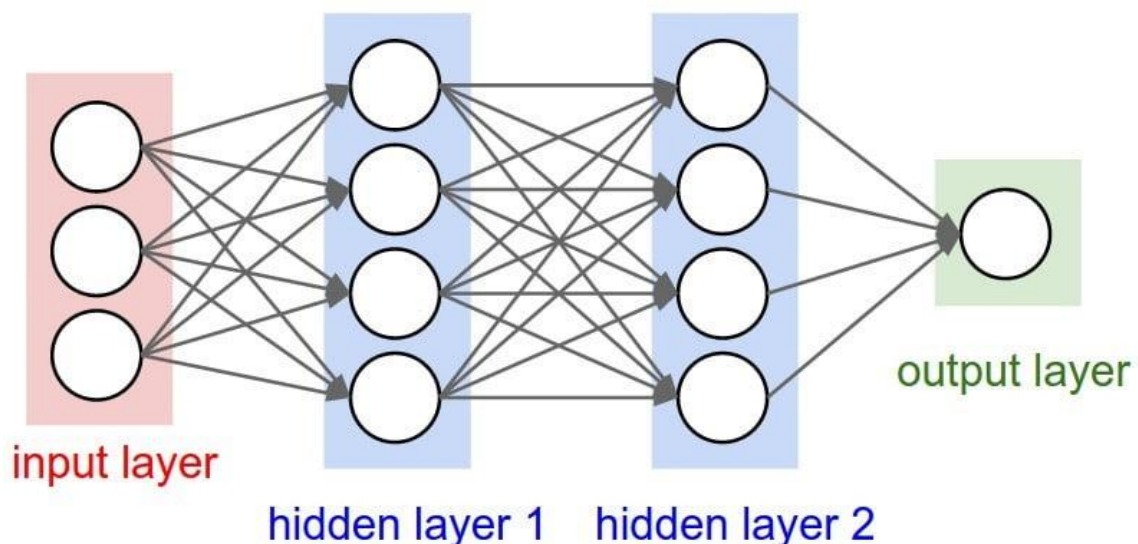
Subir el CSV preprocesado

- **Métodos (Explicar todas las técnicas utilizadas como ANOVA, Chi cuadrado, redes neuronales perceptrón multi capa con referencias en inglés)**

Para seleccionar las 10 variables más influyentes hemos utilizado el Scoring Method: *Gain Ratio*.

Para entrenar a nuestro conjunto de datos hemos utilizado un *Basic Multi_Layer Perceptron*.

- **Estructura de la red neuronal**



Nº capas de entrada: 1 capa de entrada con 10 neuronas

Nº capas ocultas: 2 capas ocultas con 15 neuronas cada una

Nº capas de salida: 1 capa de salida con 1 neurona

- **Experimentación (explicar la experimentación paso a paso)**

La experimentación en nuestro caso se ha basado en ir cambiando 4 parámetros de nuestra red neuronal y de nuestro modelo: el porcentaje de nuestro dataset que iba destinado a entrenamiento y a prueba (train-test), el número de capas ocultas, el número de neuronas que tenían nuestras capas ocultas y el número de iteraciones que realizaba nuestro modelo. Variando los 4 parámetros teníamos que encontrar la combinación de estos que más precisión diera a nuestro modelo.

Inicialmente, íbamos cambiando estos parámetros pero nos salían resultados muy dispares, incluso cuando ejecutábamos el mismo código sin cambiar nada nos daba un resultado distinto, descubrimos que se debía a la aleatoriedad con la que se dividía el dataset. Para arreglar esto introdujimos `random_state=45` en la parte del código en la que dividíamos el dataset, controlando así la aleatoriedad de la división de los datos.

Una vez hecho esto, ya podíamos ir cambiando los parámetros para comprobar que le venía mejor a nuestro modelo. Debido a que es un dataset con pocos datos, nos dimos cuenta de que tanto el número de capas ocultas como el número de iteraciones no tenía que ser muy alto. Finalmente, fuimos probando que número de neuronas por capa nos venía mejor y vimos que porcentaje de train-test nos daba más precisión. El resultado final es: 60% train – 40% test, 2 capas ocultas, 15 neuronas por capa y 2000 iteraciones.

Resultados

Matriz de confusión:

```
Confusion matrix:
[[20  0]
 [ 1 14]]
```

En la matriz de confusión resultante vemos cómo de los 35 datos destinados al test, 34 están bien clasificados (20 Verdaderos Negativos y 14 Verdaderos Positivos) y tan solo 1 dato está erróneamente clasificado, 1 Falso Negativo (El modelo predijo 0 cuando en realidad era 1). Esto indica que la precisión de nuestro modelo es realmente alta.

Precisión:

	precision	recall	f1-score	support
0	0.95	1.00	0.98	20
1	1.00	0.93	0.97	15
accuracy			0.97	35
macro avg	0.98	0.97	0.97	35
weighted avg	0.97	0.97	0.97	35

La precisión de nuestro modelo es del 97% como podemos observar, sin embargo, vista la matriz de confusión se podría esperar que fuera más alta. Esto ocurre porque nuestro dataset tiene muy pocos datos y un fallo supone mucho porcentaje de la precisión del modelo.

• Conclusiones

En este trabajo se ha desarrollado una red neuronal para la predicción de hipertensión sintomática (SIM_GIPERT) utilizando técnicas de minería de datos y aprendizaje automático. El preprocesamiento de los datos se realizó empleando herramientas como Orange y Weka, donde se seleccionaron las variables más relevantes mediante el método

de *Information Gain Ratio*, optimizando así el conjunto de características utilizadas en el modelo.

Posteriormente, se construyó una red neuronal tipo multitask utilizando *MLPClassifier* de *Scikit-learn*, configurada con dos capas ocultas de 15 neuronas cada una. El modelo fue entrenado sobre el dataset preprocesado, alcanzando un accuracy del 97% en el conjunto de prueba, lo que demuestra una alta capacidad de predicción.

Estos resultados reflejan la eficacia de combinar un adecuado preprocesamiento de los datos, una correcta selección de atributos y un diseño cuidado del modelo de red neuronal. Además, resaltan el potencial del uso de inteligencia artificial en el ámbito clínico como apoyo en el diagnóstico temprano de condiciones como la hipertensión.

- **Referencias (IEEE y mínimo 10 referencias)**

- World Health Organization, "Hypertension," 2021. [Online]. Available: <https://www.who.int/news-room/fact-sheets/detail/hypertension>.
- UCI Machine Learning Repository, "Dataset: Myocardial infarction complications" University of California, Irvine, Available: <https://archive.ics.uci.edu/dataset/579/myocardial+infarction+complications>.
- M. A. Al-Bakri, M. M. Al-Kadi, and N. A. El-Sayed, "Artificial neural networks for the prediction of hypertension," *Biomedical Research*
- J. Demšar et al., "Orange: Data mining toolbox in Python," *Journal of Machine Learning Research*
- M. Hall et al., "The WEKA data mining software: an update," *ACM SIGKDD Explorations Newsletter*
- J. Schmidhuber, "Deep learning in neural networks: An overview," *Neural Networks*
- S. Vellido, J. D. Martín-Guerrero, and P. J. G. Lisboa, "Making machine learning interpretable for clinical decision support: a brief state-of-the-art survey," *International Journal of Computational Intelligence Systems*
- D. W. Hosmer, S. Lemeshow, and R. X. Sturdivant, *Applied Logistic Regression*
- G. van Rossum and F. L. Drake, *Python 3 Reference Manual*, Scotts Valley, CA: CreateSpace, 2009
- C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006