

Pontificia Universidad Javeriana
Proyecto de Inteligencia artificial
Predicción-de-Stroke
Realizado por: Laura Sofía Gómez Lizarazo

El siguiente proyecto tiene como objetivo la predicción de Stroke en diferentes pacientes, a partir del dataset mostrado a continuación:

<https://www.kaggle.com/datasets/fedesoriano/stroke-prediction-dataset>

En dicho Dataset aparecen las siguientes características:

- 1) "id: unique identifier"
- 2) "gender": "Male", "Female" or "Other"
- 3) "age": age of the patient
- 4) "hypertension": 0 if the patient doesn't have hypertension, 1 if the patient has hypertension
- 5) "heart_disease": 0 if the patient doesn't have any heart diseases, 1 if the patient has a heart disease
- 6) "ever_married": "No" or "Yes"
- 7) "work_type": "children", "Govt_jov", "Never_worked", "Private" or "Self-employed"
- 8) "Residence_type": "Rural" or "Urban"
- 9) "avg_glucose_level": average glucose level in blood
- 10) "bmi": body mass index
- 11) "smoking_status": "formerly smoked", "never smoked", "smokes" or "Unknown"
- 12) "stroke": 1 if the patient had a stroke or 0 if not

Dado lo anterior se procedió a realizar el procesamiento de los datos, que en este caso incluía la limpieza de los mismos y la creación de un nuevo archivo con los nuevos datos. Así mismo se normalizó y se realizaron los siguientes modelos de predicción con su respectiva validación.

- * SVM (Máquinas de soporte vectorial)
- * Regresión logística
- * KNN (K vecinos más cercanos)

A continuación se muestran los resultados de los métodos utilizados, a partir de dos métricas, f1 y AUC_ROC:

Para KNN:

```
[ ] G_S_k.best_params_
```

```
{'algorithm': 'auto', 'n_neighbors': 50, 'weights': 'uniform'}
```

```
[ ] G_S_k.best_score_
```

```
0.8091997710175098
```

```
[ ] G_S_k.cv_results_['mean_test_AUC']
```

```
array([0.52760562, 0.52760562, 0.56523318, 0.56507457, 0.62600452,
       0.6255411 , 0.71344893, 0.70863992, 0.79276095, 0.78619002,
       0.80919977, 0.80481029, 0.80864477, 0.80764015, 0.52760562,
       0.52760562, 0.56523318, 0.56507457, 0.62600452, 0.6255411 ,
       0.71344893, 0.70863992, 0.79276095, 0.78619002, 0.80919977,
       0.80481029, 0.80864477, 0.80764015, 0.52760562, 0.52760562,
       0.56523318, 0.56507457, 0.62600452, 0.6255411 , 0.71344893,
       0.70863992, 0.79276095, 0.78619002, 0.80919977, 0.80481029,
       0.80864477, 0.80764015, 0.52760562, 0.52760562, 0.56523318,
       0.56507457, 0.62600452, 0.6255411 , 0.71344893, 0.70863992,
       0.79276095, 0.78619002, 0.80919977, 0.80481029, 0.80864477,
       0.80764015])
```

```
[ ] G_S_k.cv_results_['mean_test_f1']
```

```
array([0.0977839 , 0.0977839 , 0.01864499, 0.0977839 , 0.01882614,
       0.04708114, 0.          , 0.01          , 0.          , 0.          ,
       0.          , 0.          , 0.          , 0.          , 0.0977839 ,
       0.0977839 , 0.01864499, 0.0977839 , 0.01882614, 0.04708114,
       0.          , 0.01          , 0.          , 0.          , 0.          ,
       0.          , 0.          , 0.          , 0.0977839 , 0.0977839 ,
       0.01864499, 0.0977839 , 0.01882614, 0.04708114, 0.          ,
       0.01          , 0.          , 0.          , 0.          , 0.          ,
       0.          , 0.          , 0.0977839 , 0.0977839 , 0.01864499,
       0.0977839 , 0.01882614, 0.04708114, 0.          , 0.01          ,
       0.          , 0.          , 0.          , 0.          , 0.          ,
       0.          ])
```

Conclusiones

A partir de las diferentes métricas y modelos de predicción utilizados, es posible evidenciar que el mejor método para predecir si un paciente presenta Stroke es el de regresión logística con una métrica de AUC_ROC.

Algo a resaltar es que la métrica que daba más alto en todos los casos fue la de AUC_ROC. Por otra parte el resultado más bajo fue el obtenido con máquinas de soporte vectorial.

En cuanto a los tiempos de ejecución es muy importante destacar que el método más demorado fue el de SVM, tardándose aproximadamente 1 hora.

Otro aspecto importante que vale la pena mencionar es que a lo largo del proyecto se utilizó Gridserch cv, evaluando cada método con dos métricas f1 y AUC-ROC, por lo que es importante que se especifique el refit dentro de los parámetros del Gridserch-cv, pues en caso contrario saldrá error.