

Análisis y Reporte sobre el desempeño del modelo (septiembre 2022)

Gerardo Peña Pérez A01701474

En este reporte se evalúa el desempeño de un modelo de regresión logística usado para predecir si un tumor es maligno o benigno. Se hizo uso del framework scikit-learn

RESUMEN Se tomó un modelo de regresión lineal

CONCEPTOS CLAVE Machine Learning, regresión lineal, framework, scikit-learn

I. INTRODUCCIÓN

La regresión logística es un algoritmo que resulta útil para resolver problemas de clasificación binaria, es decir, en el que se tiene que predecir entre 2 diferentes categorías, o bien predecir la probabilidad de que algo pertenezca a una clase, esto basándose en diversos datos de entrada. En este caso se usará dicho algoritmo para predecir si un tumor es maligno o benigno, tomando como entrada características de la masa a analizar, tales como el radio, la textura, la concavidad, entre otras.

Se usa el data set de Wisconsin Breast Cancer Database (January 8, 1991), el cual cuenta con 569 instancias y 32 atributos. Uno de dichos atributos es la variable a pronosticar, es decir la clase de tumor, teniendo un valor de 2 al ser benigno y 4 al ser maligno.

II. SEPARACIÓN DE LOS DATOS

Se realizó una etapa de entrenamiento y una de pruebas, para esto, el set de datos fue separado en 2 subsets (train y test), esto para poder verificar que el modelo esté aprendiendo correctamente y que no esté aprendiendo los valores específicos en lugar de arrojar predicciones reales, a esto se le llama estado de “over fitting”, en el que, a pesar de que se llegara a lograr un error de entrenamiento bajo, el error de validación terminaría siendo bastante alto, debido a que el modelo no reaccionaría bien al enfrentarse a datos que no haya visto durante su entrenamiento, pero de esto hablaré más adelante.

Para dividir los datos se usó la función “train_test_split” del modulo de model_selection de scikit-learn. A dicha función le doy los conjuntos de datos que quiero que divida y una semilla que inicializa el generador de números pseudoaleatorios que básicamente me ayuda a que la división sea aleatoria. Al no

definir el tamaño de cada subset, automáticamente forma el set de entrenamiento con un 75% del total y el de pruebas con 25%

III. RESULTADOS

Una vez entrenado el modelo, se procede a crear un arreglo con las predicciones generadas por el modelo para así poder evaluar el desempeño de este probando con datos que no conoció en la etapa de entrenamiento. Haciendo uso de la función “accuracy_score” podemos conocer que la precisión con la que el modelo genera predicciones es de 0.9714285714285714, es decir que ronda cerca del 97% de precisión.

Para poder describir el desempeño del modelo de clasificación generado se usó una matriz de confusión. Esta es la matriz obtenida con los datos de prueba:

$$\begin{bmatrix} 117 & 1 \\ 4 & 53 \end{bmatrix}$$

Podemos notar que las predicciones acertadas suman 170. Mientras que tuvimos 4 casos de falsos positivos y tan solo 1 caso de falso negativo.

IV. NIVEL DE AJUSTE DEL MODELO

Para saber si el modelo estaba en un estado de overfitting se realizó un cross validation haciendo 5 divisiones del dataset y se obtuvo un accuracy promedio de 0.9513874614594039, lo cual nos dice que el modelo en efecto entrenó bien y no solo para los valores de entrenamiento, evitando así el estado de overfitting

V. DETECCIÓN DE BIAS Y VARIANZA

U

V. REGULARIZACIÓN

En este caso se obtuvo un nivel aceptable de error y accuracy, pero se puede calibrar el modelo para mejorar su desempeño, para esto se utiliza algún método de regularización. En este caso se usó la regularización Ridge (L2).