

MODEL CARD

Visió general del model

El model present és bàsicament un Support Vector Machine per classificar l'estat de supervivència dels pacients donats unes característiques clíniques. Concretament, la base de dades és la que proporciona UCI (la Universitat de Califòrnia a Irvine) anomenada "Cirrhosis Patient Survival Prediction". En aquest "Model card" trobaràs informacions tant sobre la variable resposta com les variables predictores, la intenció principal del model, les seves limitacions i les seves puntuacions sobre les mètriques d'avaluació que hi ha.

Versió

name: hdayiuy8e1231dasshi91
data: 28 - 12 - 2023

Autor

- Zhihao Chen, zhihao.chen@estudiantat.upc.edu

Referència

[https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+\(Diagnostic\)](https://archive.ics.uci.edu/ml/datasets/Breast+Cancer+Wisconsin+(Diagnostic))
<https://minds.wisconsin.edu/bitstream/handle/1793/59692/TR1131.pdf>
<https://scikit-learn.org/stable/>

Consideracions

Usuaris destinats

- Estudiants de IAA
- Professors de IAA

Cas d'ús

El conjunt de dades en què es va entrenar aquest model es va crear per donar suport a la comunitat d'aprenentatge de màquines en la realització d'anàlisis empíriques dels algorismes de ML. La base de dades es pot utilitzar en estudis relacionats amb la predicció de la supervivència dels pacients donats les seves mesures clíniques.

Limitacions

El principal problema de la base de dades investigada són el desbalanceig de les classes. Per aquest motiu, el model podria tenir una taxa veritable positiva baixa i una taxa veritable negativa alta. Per resoldre-ho s'ha aplicat el mètode de donació de pesos a les classes, però probablement no ho resoldrà del tot. Per altra banda també hi són l'existència dels valors atípics, la gran quantitat de valors buits i la falta de mostres observades. Tot això ha limitat també la capacitat predictiva del mode.

Consideracions ètiques

Risc: La predicció errònia pot causar decisions inadequades en el tractament del pacient.

Estratègia de mitigació: Només utilitzar el model per a estudis relacionats amb ML.

Detalls del model

Descripció del Model

La idea fonamental darrere de les SVM és trobar un hiperplà òptim que pugui separar de la millor manera possible dos conjunts de dades en un espai de característiques d'alta dimensió. Té dos hiperparàmetres més rellevants:

- El paràmetre de regularització C controla l'equilibri entre dos objectius oposats a l'entrenament d'una SVM: maximitzar el marge i minimitzar la classificació errònia dels punts de dades.
- El tipus de kernel determina la forma del hiperplà que separa el conjunt de dades.

Preprocessament

Per les variables numèriques s'ha fet la transformació dels valors atípics a valors buits, la imputació amb la mediana i l'estandardització. Per les variables categòriques no ordinals s'ha creat una nova modalitat per guardars els valors buits i s'ha fet la recodificació binària. Per les variables ordinals s'ha fet la imputació amb la moda i la recodificació ordinal.

De les particions de dades no s'ha decidit guardar una part per la validació, per tant només hi són train i test, amb una proporció habitual de 8:2. A l'hora de necessitar la partició de validació es farà servir l'algorisme de 10-fold Cross Validation sobre la partició de train.

Entrenament del model

Degut al desbalanceig de classes que hi ha en la base de dades, s'ha decidit aplicar la tècnica de donació de pesos a les classes. Per trobar la millor configuració dels hiperparàmetres s'ha fet servir l'algorisme de GridSearchCV, provant diverses combinacions d'hiperparàmetres disponibles. Finalment els hiperparàmetres són:

{C=0.01, class_weight='balanced', kernel='linear', random_state=1}, la resta de paràmetres són els de defecte.

Mètrica d'avaluació

A causa de tenir dades desequilibrades, les mètriques d'avaluació no poden ser les estàndards:

- L'exactitud balancejada (Balanced accuracy)
- La F1 score ponderada

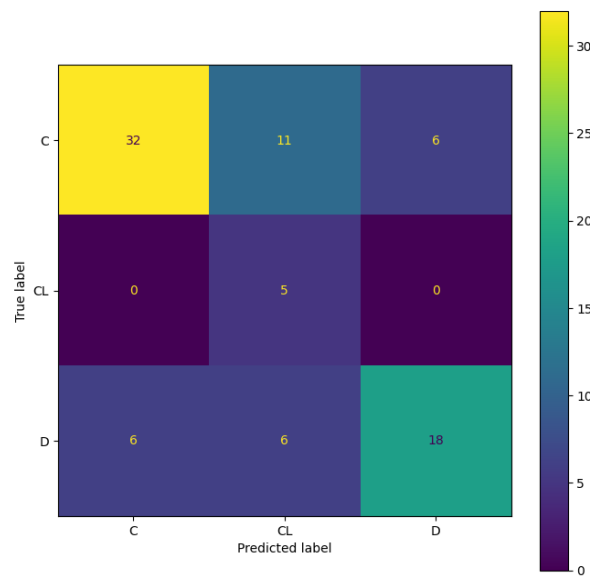
Aquestes dues no són influenciades pel desbalanceig de classes i proporcionen informacions suficients fiables sobre el rendiment del model: tant com la taxa de veritables positius com la de veritables negatius són importants en la present investigació.

Rendiment del model

A continuació són els resultats obtinguts sobre la partició de test:

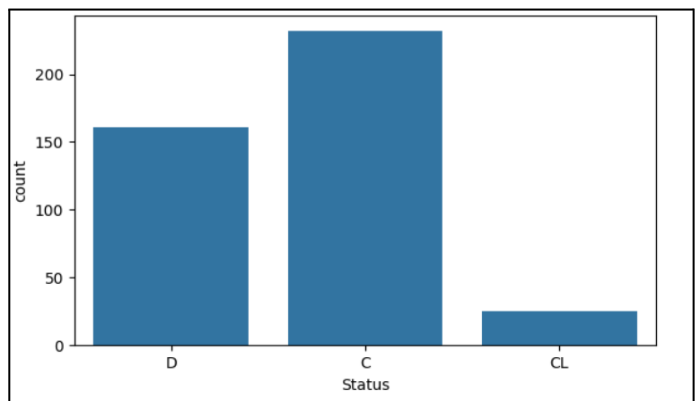
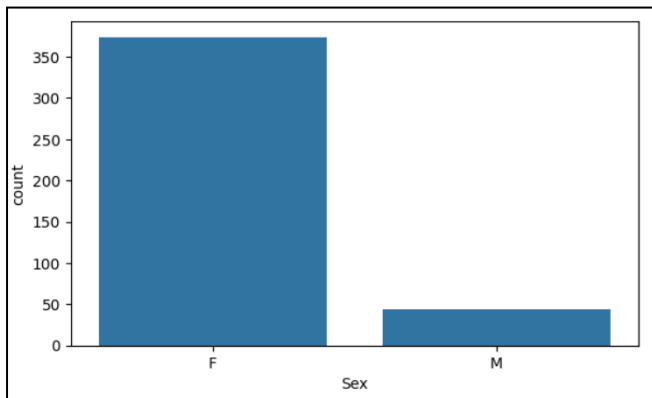
- Balanced accuracy: 0.75102
- F1 score average weighted: 0.68926

En la matriu de confusió de les prediccions sobre la partició de test s'observa que totes les mostres que pertanyen a la classe "CL", que és minoritària, són classificades correctament, la qual cosa vol dir que la donació de pesos a les classes va tenir èxit. A més, per les altres dos classes la predicció tampoc és gens mala: la majoria s'ha classificat correctament.



Variables categòriques

En aquesta secció inclouen els diagrames de barres de la variable "Sex" i la variable objectiu "Status". S'ha seleccionat especialment aquestes dos variables perquè s'ha considerat rellevant que l'usuari vegi el desbalanceig de classes que existeix tant en les variables independents com en la variable resposta.



Variables numèriques

La intenció d'aquesta secció és mostrar a l'usuari la quantitat de valors atípics i valors buits que hi ha en les variables numèriques mitjançant els histogrames i els diagrames de caixes. Només s'ha mostrat una variable perquè ja que les altres tenen el mateix comportament majoritàriament.

