

# Progetto Intelligenza Artificiale

## Decision Tree

Lorenzo Italiano - 6134469

20/03/2019

### 1 Descrizione

Il codice utilizzato permette di fare uso del *DecisionTree* di scikit-learn per determinare se un nuovo impiegato avrà uno stipendio maggiore rispetto alla mediana degli stipendi dei vari impiegati nel dataset.

Il dataset è reperibile su [Stack Overflow Developer Survey, 2017](#).

Sono state usate le librerie:

- **numpy**, **pandas** per operare sul dataset
- **matplotlib** per disegnare i vari grafici
- **sklearn** per definire il decision tree, per creare la roc curve e per utilizzare il 10-fold cross validation
- **pydotplus**, **graphviz**, **collections** per disegnare il grafo dell'albero

Nella prima parte del codice viene riadattato il dataset per poter essere letto da *DecisionTree*.

In seguito è stato definito il target ed è stato diviso il dataset in train-set e test-set per determinare la precisione della previsione dell'albero.

È stato inoltre verificato quale fosse il parametro *max\_depth* adatto a seconda del train-set, seguendo le indicazioni dell'articolo "[InDepth: Parameter tuning for Decision Tree](#)" di *Mohtadi Ben Fraj*, illustrandone il grafico rispetto alla sua variazione.

Infine, dopo aver creato l'albero decisionale ed aver determinato la precisione della previsione mediante il 10-fold cross validation ( $\approx 0.83$ ), è stato disegnato il grafo dell'albero, seguendo il codice ottenuto dall'articolo "[Creating and Visualizing Decision Trees with Python](#)" di *Russel*.

## 1.1 Setup Dataset

Per poter definire il decision tree è risultato essere necessario rimodellare il dataset.

Dopo aver eliminato le righe con valore *NaN* nella colonna Salario, averla estratta dal dataset e aver eliminato le colonne '*ExpectedSalary*' e '*NonDeveloperType*' (poiché una irrilevante e l'altra contenente solo valori *NaN*), sono stati sostituiti i valori *NaN* delle altre colonne: quelle contenenti stringhe con una stringa '*NaN*' e quelle contenenti numeri con la mediana degli altri numeri della colonna.

Dopodiché tutte le colonne sono state divise in un numero di colonne pari al numero di risposte date, tutte con valori booleani (**1** se l'impiegato aveva dato quella risposta e **0** altrimenti), inoltre, poiché il dataset contiene risposte multiple, ogni risposta è stata divisa grazie al metodo *get\_dummies()* della libreria **pandas**.

## 2 Ricerca del Miglior 'max\_depth'

Seguendo l'articolo, citato nella descrizione, di *Mohtadi Ben Fraj* è stato possibile creare un grafico della precisione dell'albero, nella metrica **AUC** (*=Area Under Curve*), al variare del parametro *max\_depth*, il quale indica l'altezza massima dell'albero.

Come mostra la **Figure 1**, aumentando l'altezza dell'albero, si ha un incremento sia nel train-set che nel test-set (in quanto riesce a catturare più informazioni), al quale segue un calo nel test-set a causa del fenomeno dell'*overfitting*, che si osserva quando si ha un numero eccessivo di parametri rispetto al numero di osservazioni.

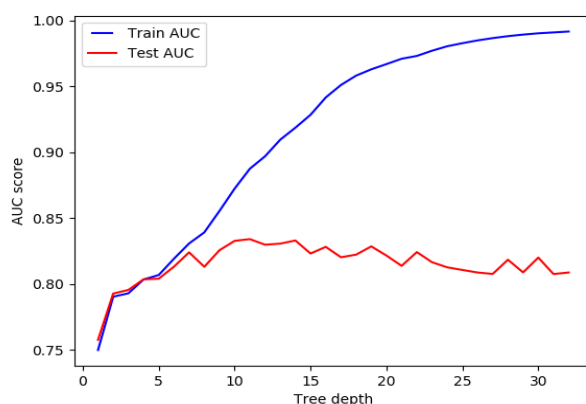


Figure 1: score al variare del max\_depth

Dopodiché è trovata la profondità ideale ad avere una precisione migliore a seconda della divisione effettuata del dataset in train-set e test-set.

La precisione viene calcolata mediante la funzione `score()` del Decision-TreeClassifier preso dalla libreria **sklearn**, la metrica AUC di **matplotlib** ed infine con il metodo del *10-fold cross validation* fornito sempre da **sklearn**.

Utilizzando questa serie di azioni la *precisione* della predizione con l'albero decisionale risulta essere di circa l'**83%**.

### 3 Grafo del Decision Tree

Nell'ultima parte del codice vi è una parte commentata, in quanto serve a disegnare e a memorizzare il grafo dell'albero decisionale ottenuto.

Il metodo utilizzato per disegnare è quello dell'articolo, citato nella descrizione, di *Russel*, che consiste di creare un *dot.data* dell'albero per trasformarlo poi in un grafo. In seguito vengono colorati i nodi in maniera alternata ed infine viene memorizzato il grafo.

A causa della larghezza spropositata del grafo, per poterlo visualizzare correttamente è necessario cliccare [qui](#).