

XGB Stock Forecast

Progetto 'Tree Based Models'

Studente: Davide Zicca

Professoressa: Anna Gottard



UNIVERSITÀ DELLA CALABRIA

DIPARTIMENTO DI ECONOMIA,
STATISTICA E FINANZA

“Giovanni Anania” - DESF

Dipartimento di Economia, Statistica e Finanza
Università della Calabria (UNICAL)

15 Maggio, 2021

L'obiettivo di questo progetto è quello di utilizzare *Tree Based Models* per predire l'andamento (Up-tick vs Down-tick) del prezzo giornaliero delle stock:

- ▶ Amazon ("AMZN")
- ▶ Apple ("AAPL")¹

L'idea è quindi quella di assegnare valore 1 in caso di Up-movement e valore 0 per il Down-movement. Essendo un problema di classificazione binaria utilizzo la libreria, presente in *R Studio*, **xgboost**.

¹Nelle prossime slide verrà mostrato il codice che fa riferimento allo script contenente la serie storica del prezzo delle azioni di Apple. Tuttavia, lo stesso procedimento è stato seguito per studiare anche l'andamento del prezzo azionario di Amazon.

Di importanza cruciale per questo progetto è una breve introduzione al modello oggetto di studio. L'eXtreme Gradient Boosting (**XGBoost**) model (Chen and Guestrin, 2016), è un approccio attraverso il quale nuovi modelli vengono creati (con la capacità di prevedere gli errori di modelli precedenti) e poi vengono uniti al fine di raggiungere la previsione finale.

Il Boosting è una tecnica ensemble nei quali nuovi modelli vengono aggregati per correggere gli errori dei modelli precedenti. I modelli vengono aggiunti sequenzialmente fino a che nessun miglioramento può essere raggiunto. Esso usa il *tree ensemble model* che è un insieme di alberi di classificazione e regressione (*CART models*). L'approccio ensemble viene usato perchè, in alcuni casi, un singolo CART non ha forte potere predittivo. Usando un insieme di *CART* (un *tree ensemble model*) è possibile considerare la somma delle previsioni di più alberi predittivi.

La *objective function* dell'eXtreme Gradient Boosting (**XGBoost**) model è:

$$Obj = L + \Omega$$

Dove, L è la *loss function* che controlla il potere predittivo, e Ω è la componente di *regolarizzazione* che controlla la complessità e l'overfitting del modello. La *loss function* (L), che deve essere ottimizzata, è rappresentata dal *Root Mean Squared Error* per la regressione, *Log-loss* per la classificazione binaria, o *mlogloss* per la classificazione multi-classe. La componente di *regolarizzazione* Ω dipende dal numero di foglie e dallo score previsionale assegnato alle foglie nel *tree ensemble model*.

Viene chiamato *gradient boosting* perchè usa un algoritmo *gradient descent* per minimizzare la perdita quando si aggiungo nuovi modelli. L'algoritmo di *gradient boosting* (M. et al., 2020) supporta sia la regressione che la classificazione per i problemi di modellizzazione previsionali (Mottaghi and Farhangdoost, 2021).

Sono stati collezionati i prezzi azionari delle stock di Apple dal 1° Gennaio 2015. Da questi valori sono stati calcolati alcuni tra i più utilizzati indici azionari:

- ▶ Il Relative Strength Index (RSI) è un *momentum indicator* largamente utilizzato nell'analisi tecnica e che misura la magnitudo dei più cambiamenti di prezzo per valutare fenomeni di eccessiva vendita/acquisti di azioni o altri asset.
- ▶ L'Average Directional Index (ADX) serve per determinare quando una stock è in fase di forte trading sul mercato azionario. Si basa sul calcolo del Moving Average ² del range del prezzo, fissato un intervallo di tempo di interesse.
- ▶ Il SAR parabolico è un indicatore tecnico che determina la direzione del prezzo di un determinato asset: è usato anche come alert dell'imminente cambiamento di direzione del prezzo dell'asset.

²Si veda anche il *MACD* (Moving Average Convergence Divergence) che è un altro indicatore fortemente utilizzato nell'analisi tecnica.

Tree Based Models - Progetto Xgb Stock Forecast

Davide Zicca

9/5/2021

Data Preparation

```
setwd("C:/Davide/MASTER IN DATA SCIENCE/Materiale del Master/Tree Based Models/PROGETTO/XGB stock forecast")
library(quantmod)
```

```
## Loading required package: xts
## Loading required package: zoo
##
## Attaching package: 'zoo'
## The following objects are masked from 'package:base':
##
##   as.Date, as.Date.numeric
## Loading required package: TTR
## Registered S3 method overwritten by 'quantmod':
##   method      from
##   as.zoo.data.frame zoo
```

```
library(TTR)
library(xgboost)
```

```
## Warning: package 'xgboost' was built under R version 4.0.4
```

```
library(readr)
```

```
getSymbols("AAPL", from = "2015-01-01")
```

```
## 'getSymbols' currently uses auto.assign=TRUE by default, but will
## use auto.assign=FALSE in 0.5-0. You will still be able to use
## 'loadSymbols' to automatically load data. getOption("getSymbols.env")
## and getOption("getSymbols.auto.assign") will still be checked for
## alternate defaults.
##
## This message is shown once per session and may be disabled by setting
## options("getSymbols.warning4.0"=FALSE). See ?getSymbols for details.
## [1] "AAPL"
```

```
# AAPL
write.csv(AAPL, file = "AAPL.csv")
AAPL <- read_csv("AAPL.csv", col_types = cols(
  X1 = col_number(),
  AAPL.Open = col_number(),
```

```

AAPL.High = col_number(),
AAPL.Low = col_number(),
AAPL.Close = col_number(),
AAPL.Volume = col_number(),
AAPL.Adjusted = col_number()
))

## Warning: Missing column names filled in: 'X1' [1]

symbol = "AAPL"
fileName = paste(getwd(), "/", symbol, ".csv", sep="") ;
df = as.data.frame(AAPL)
colnames(df) = c("Date", "Open", "High", "Low", "Close", "Volume", "Adjusted")
head(df)

##   Date      Open      High      Low      Close      Volume Adjusted
## 1    1 27.8475 27.8600 26.8375 27.3325 212818400 24.81924
## 2    2 27.0725 27.1625 26.3525 26.5625 257142000 24.12005
## 3    3 26.6350 26.8575 26.1575 26.5650 263188400 24.12232
## 4    4 26.8000 27.0500 26.6750 26.9375 160423600 24.46056
## 5    5 27.3075 28.0375 27.1750 27.9725 237458000 25.40040
## 6    6 28.1675 28.3125 27.5525 28.0025 214798000 25.42763

rsi = RSI(df$Close, n=14, maType="WMA")
adx = data.frame(ADX(df[,c("High", "Low", "Close")]))
sar = SAR(df[,c("High", "Low")], accel = c(0.02, 0.2))
trend = df$Close - sar
rsi = c(NA, head(rsi, -1))
adx$ADX = c(NA, head(adx$ADX, -1))
trend = c(NA, head(trend, -1))
# Create the target variable
price = df$Close - df$Open
class = ifelse(price > 0, 1, 0)
# Create a Matrix
model_df = data.frame(class, rsi, adx$ADX, trend)
model = matrix(c(class, rsi, adx$ADX, trend), nrow=length(class))
model = na.omit(model) # Drop NA
colnames(model) = c("class", "rsi", "adx", "trend")

```

Predisposizione del modello: training e CV

```

# Split data into train and test sets
train_size = 2/3
breakpoint = nrow(model) * train_size
training_data = model[1:breakpoint,]
test_data = model[(breakpoint+1):nrow(model),]
# Split data training and test data into X and Y
X_train = training_data[,2:4] ; Y_train = training_data[,1]
class(X_train)[1] ; class(Y_train)

## [1] "matrix"
## [1] "numeric"

```

Prima di procedere oltre, è stato creato un *lag* negli indicatori precedentemente calcolati per evitare il *look-ahead bias* ³.

Le *input features* sono state così create.

Si noti che tale modello è una semplificazione di un modello di *stock forecast*. Altri indicatori e migliori assunzioni sono necessarie per migliorare le performance predittive del modello. Come si vedrà le performance di tale modello si attestano su percentuali non soddisfacenti per basare una solida strategia di trading azionario. Le input features e la variabile target create precedentemente sono unite per creare una singola matrice: al fine di seguire la struttura matriciale imposta dal modello XGBoost. Il dataset viene poi diviso in training e test dataset. La funzione **xgboost** contenuta in R viene usata per addestrare il modello.

³Tale fenomeno si manifesta quando si usano dati, all'interno di una simulazione o di un'analisi, che non sarebbero altrimenti disponibili durante il periodo di interesse. Questo può portare a risultati inappropriati. Talvolta, può anche portare i risultati della simulazione verso quelli desiderati dallo studio: dando *overconfidence* verso i risultati ottenuti dal modello di forecast. Infine, il *look-ahead bias* può alterare le strategie di trading, quando queste ultime fanno troppo riferimento alle serie storiche.


```

X_test = test_data[,2:4] ; Y_test = test_data[,1]
class(X_test)[1]; class(Y_test)

## [1] "matrix"
## [1] "numeric"

# Train the xgboost model using the "xgboost" function
dtrain = xgb.DMatrix(data = X_train, label = Y_train)
xgModel = xgboost(data = dtrain, nround = 5, objective = "binary:logistic",
                  eval_metric= "error")

## [1] train-error:0.374761
## [2] train-error:0.301147
## [3] train-error:0.267686
## [4] train-error:0.239006
## [5] train-error:0.224665

# Using cross validation
dtrain = xgb.DMatrix(data = X_train, label = Y_train)
cv = xgb.cv(data = dtrain, nround = 10, nfold = 5, objective = "binary:logistic",
            eval_metric= "error")

## [1] train-error:0.331741+0.018158 test-error:0.523864+0.023829
## [2] train-error:0.292782+0.013763 test-error:0.515284+0.012411
## [3] train-error:0.272709+0.014816 test-error:0.519121+0.013701
## [4] train-error:0.254546+0.016846 test-error:0.524858+0.010658
## [5] train-error:0.241640+0.018721 test-error:0.531556+0.009635
## [6] train-error:0.227778+0.023305 test-error:0.523887+0.021587
## [7] train-error:0.220848+0.021945 test-error:0.517188+0.026885
## [8] train-error:0.210809+0.018373 test-error:0.520078+0.028453
## [9] train-error:0.198857+0.019596 test-error:0.504780+0.037977
## [10] train-error:0.195271+0.019460 test-error:0.516250+0.038586

```

Forecast, Performance del Modello e Visualizzazione

```

# Make the predictions on the test data
preds = predict(xgModel, X_test)
# Determine the size of the prediction vector
print(length(preds))

## [1] 523

# Limit display of predictions to the first 6
print(head(preds))

## [1] 0.5860648 0.4273933 0.5050582 0.5050582 0.5050582 0.5050582

prediction = as.numeric(preds > 0.5)
print(head(prediction))

## [1] 1 0 1 1 1 1

# Measuring model performance
error_value = mean(as.numeric(preds > 0.5) != Y_test)
print(paste("test-error=", error_value))

## [1] "test-error= 0.466539196940727"

```

Tale funzione mostra il *gradient* e *second order gradient*. L'output è l'errore di classificazione sul training dataset.

In aggiunta viene usata la funzione di *cross-validation* attraverso il comando: `xgb.cv`. In questo caso, il campione originale è partizionato randomicamente in *nfold* sub-campioni di ugual misura. Tra i *nfold* sub-campioni, un unico sub-campione viene mantenuto come *validation data* per testare il modello, e i rimanenti ($nfold - 1$) sub-campioni vengono usati come *training data*. Il processo di cross-validation viene ripetuto *nrounds* volte, con ognuno dei *nfold* sub-campioni usati uno alla volta come *validation data*. La funzione `xgb.cv` mostra una *data.table* contenente i risultati di cross validation.

Al fine di fare previsioni sui dati non osservati (ovvero sul test dataset), applichiamo il modello fittato sul test dataset. Al fine di leggere i dati ottenuti per la nostra classificazione binaria (0= Down-movement; 1= Up-movement) effettuiamo una semplice trasformazione. Compariamo il risultato ottenuto con una *threshold* fissata pari a 0.5. Come metrica per testare le performance del nostro modello, viene usato l'errore medio. Compariamo lo score ottenuto con la *threshold* pari a 0.5.

Regola decisionale: se lo score è più piccolo di 0.5, allora otteniamo un valore pari a 0. Se tale valore non è uguale al valore effettivo dal test dataset, allora è interpretato come risultato erraneo. Tutte le previsioni sono state comparate con i rispettivi valori del Y_{test} (con il rispettivo errore medio). L'alternativa è quella di creare una matrice di confusione per misurare le performance del modello. Le *important features* del modello sono estratte con la funzione **xgb.importance**.

- Chen, Tianqi and Carlos Guestrin (2016). "XGBoost: a Scalable Tree Boosting System". In: *CoRR* abs/1603.02754. arXiv: 1603.02754. URL: <http://arxiv.org/abs/1603.02754>.
- M., Nabipour et al. (2020). "Deep Learning for Stock Market Prediction". In: *Entropy* 22. URL: <http://dx.doi.org/10.3390/e22080840>.
- Mottaghi, Navid and Sara Farhangdoost (2021). "Stock Price Forecasting in Presence of Covid-19 Pandemic and Evaluating Performances of Machine Learning Models for Time-Series Forecasting". In: URL: <https://arxiv.org/abs/2105.02785>.

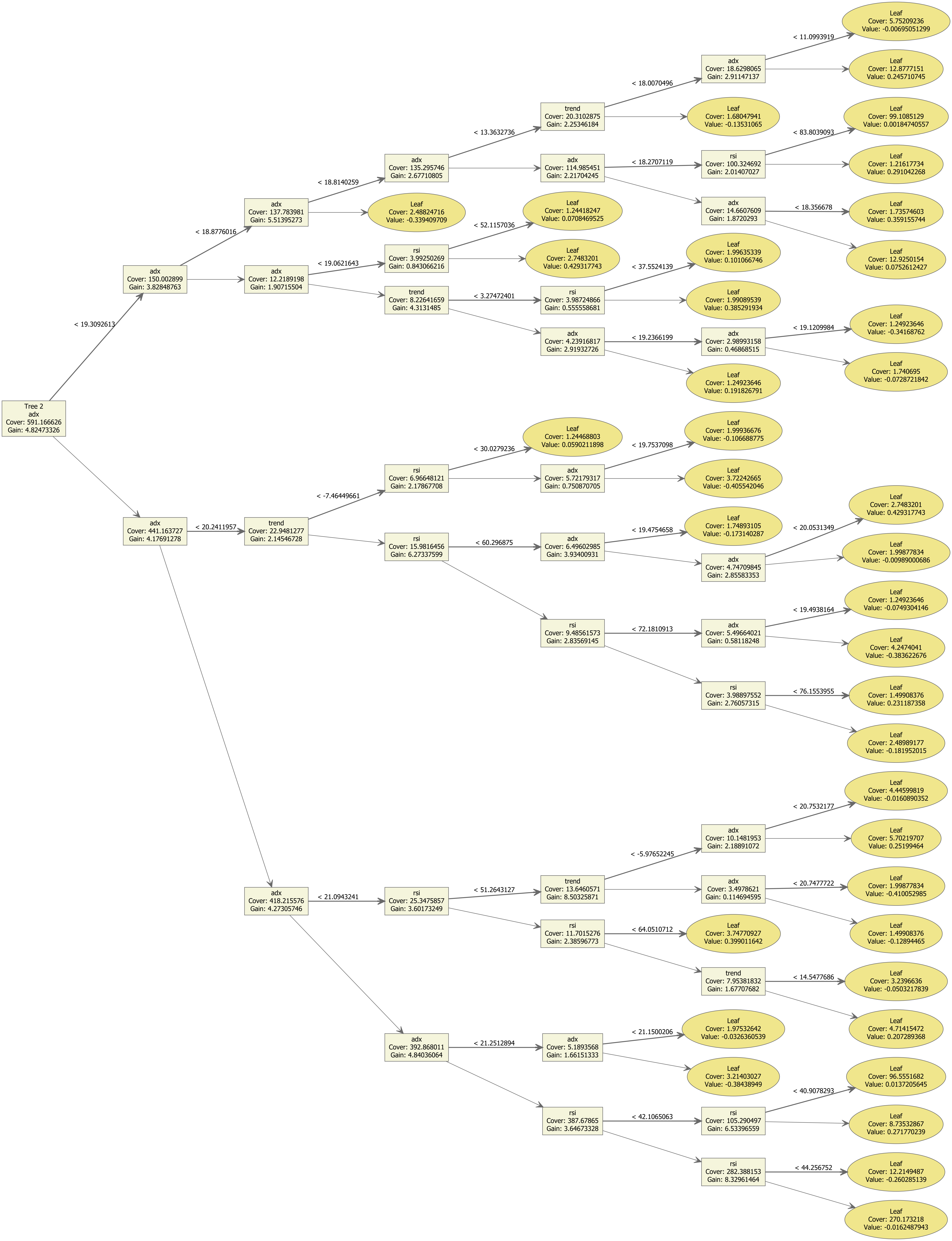
```

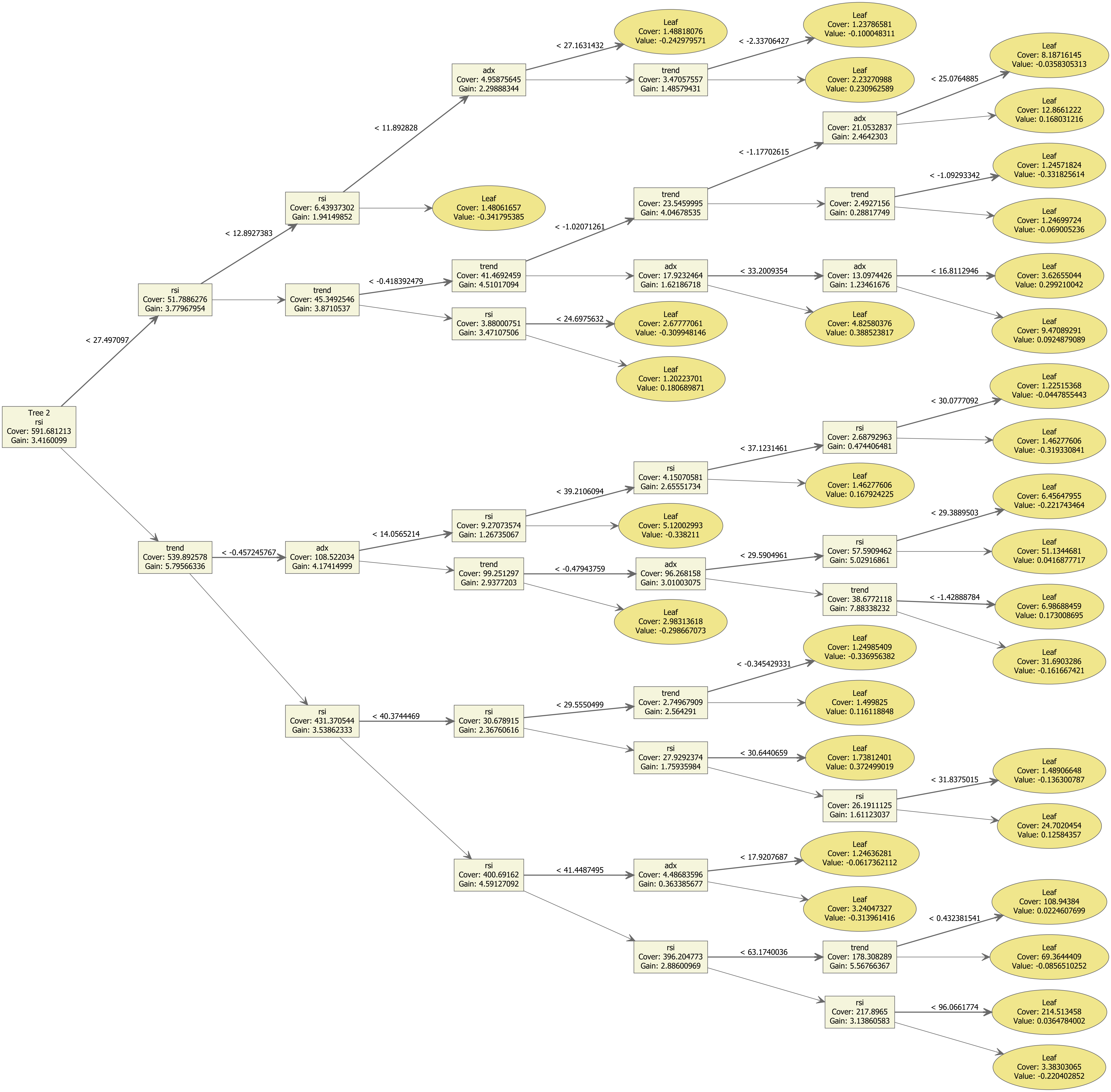
# View feature importance from the learnt model
importance_matrix = xgb.importance(model = xgModel)
print(importance_matrix)

##      Feature      Gain      Cover Frequency
## 1:      rsi 0.3970108 0.3238952 0.3860759
## 2:     trend 0.3545263 0.4393252 0.3481013
## 3:       adx 0.2484628 0.2367796 0.2658228

# View the trees from a model
xgb.plot.tree(model = xgModel)
# View only the first tree in the XGBoost model
xgb.plot.tree(model = xgModel, trees = 1)
xgb.plot.tree(model = xgModel, trees = 3, render = F, plot_height = 2000, plot_width = 1500)

```



GRAZIE PER L'ATTENZIONE


```
table(prediction)
```

```
## prediction
```

```
##      0      1
```

```
## 203 320
```