

# GroupProject

Group 16

4/26/2022

## Introduction

How to value a company is one of the greatest debates in the history of finance. One of the most common ways that investors value companies is using revenue multiples. The appropriate revenue multiple to apply to a subject company is obtained from comparable public companies or precedent transaction multiples. Some industries have higher multiples than others. One industry that enjoys some of the highest revenue multiples is software and in particular SaaS which stands for Software as a Service. Some notable SaaS companies include Zoom, Slack, and Dropbox, to name a few. We set out to predict what the revenue multiple of a SaaS company is given several financial metrics. Additionally, we want to find out what gives some SaaS companies higher revenue multiples than others. We are particularly interested in these data questions because we all plan to work as software engineers as graduating. We will build a linear regression model and a random forest regression model to explore these data questions.

## About the data

We define a revenue multiple as the enterprise value of a company divided by the trailing 12 months' revenue of the company. If a company has a 4 billion dollar enterprise value and did 200 million in revenue over the last 12 months, then its multiple would be 20x. All of the companies that we analyzed are publicly traded and we were able to obtain their financial data through lawfully required SEC filings via Ycharts. We put together a dataset of 90 companies on our own but is based on the SEG SaaS Index 2022 Annual Report. SEG / Software Equity Group is a financial analyst firm that focuses on analyzing SaaS companies. The index includes about 100 companies but we removed a handful that had incomplete data. Our response variable is ev\_ttm\_multiple which stands for enterprise value / trailing twelve months multiple. Our first predictor is revenue\_growth which is revenue growth % year over year. Our second predictor is ebitda\_margin % which is a way to measure the profitability of a company. Our third predictor is sales\_efficiency which measures how effective a company is at turning sales and marketing spend into revenue. Our last predictor is gross\_margin which measures the amount of profit made for every dollar in revenue.

```
StockData = read.csv("C:\\\\Users\\\\tobio\\\\Documents\\\\Data Science 4322\\\\Stock_Data.csv")
View(StockData)
#View(StockData)
```

## Linear Regression Model

First we wanted to see the correlation of all the predictors just so we could get a good grasp and what variables would and wouldn't work. Next, we wanted to see if all or only some variables correlated in predicting our response variable. So we used the regsubsets() function that will essentially perform the best subset selection by identifying the best model with a given number of predictors. So in this case we have 4 predictors in total so the function tests for the best 1 variable model all the one to the best 4 variable model.

We saw that our best model was the one with all 4 predictors. To confirm this we then created a table to show the rsq, adjusted r squared, Cp, and BIC for all of the best models from 1 predictor to 4. And we then decided that we would use the model with all four predictors for linear regression.

Here's the equation for the our linear regression model:

$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$ , where  $Y$  repersents the evm\\_ multiple,  $X_1, X_2, X_3, X_4$  represent, revenue growth, sales\_efficiency, ebitda\_margin and gross\_margin respectively.

```
stock_lm4 = lm(ev_ttm_multiple~revenue_growth+sales_efficiency
               +ebitda_margin+gross_margin,data=StockData)
summary(stock_lm4)

##
## Call:
## lm(formula = ev_ttm_multiple ~ revenue_growth + sales_efficiency +
##     ebitda_margin + gross_margin, data = StockData)
##
## Residuals:
##    Min      1Q  Median      3Q      Max
## -37.930 -5.801 -1.200  4.932 55.038
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -23.688    8.769   -2.701  0.00834 **
## revenue_growth 40.315    8.925   4.517 2.01e-05 ***
## sales_efficiency 11.526    4.949   2.329  0.02224 *
## ebitda_margin -16.316    5.517   -2.957  0.00402 **
## gross_margin  32.586   12.229   2.665  0.00922 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 13.35 on 85 degrees of freedom
## Multiple R-squared:  0.5443, Adjusted R-squared:  0.5228
## F-statistic: 25.38 on 4 and 85 DF,  p-value: 7.534e-14
```

Because our p-value is less than  $2.2 \times 10^{-36}$ . We can the reject the null hypothesis which states that  $\beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$ . This further tells us that at least one of the predictors can be used to predict the model as shown above. An average  $R^2$  value of 52.28% was gotten which telss us that about 52% of the data can be explained bu our model. Also, we it is evident from the model that each predictor has some sort of significance. Howevwe, beccause we know that increasing the number of predictors increases the  $R^2$  value, we tried to get the best subset of predictors and all signs point to using the 4 variables as our predictors

```
library(leaps)
stock_fit = regsubsets(ev_ttm_multiple~revenue_growth+sales_efficiency
                      +ebitda_margin+gross_margin,data=StockData)

stock_res = summary(stock_fit)
stock_stat = cbind(stock_res$rsq,
                   stock_res$adjr2,
                   stock_res$cp,
                   stock_res$bic)
colnames(stock_stat) = c("rsq", "Adjr2", "Cp", "BIC")
stock_stat
```

```

##          rsq      AdjR2       Cp       BIC
## [1,] 0.4392949 0.4329232 18.583109 -43.07079
## [2,] 0.4892489 0.4775074 11.265656 -46.96912
## [3,] 0.5152085 0.4982972  8.423636 -47.16404
## [4,] 0.5442865 0.5228411  5.000000 -48.23113

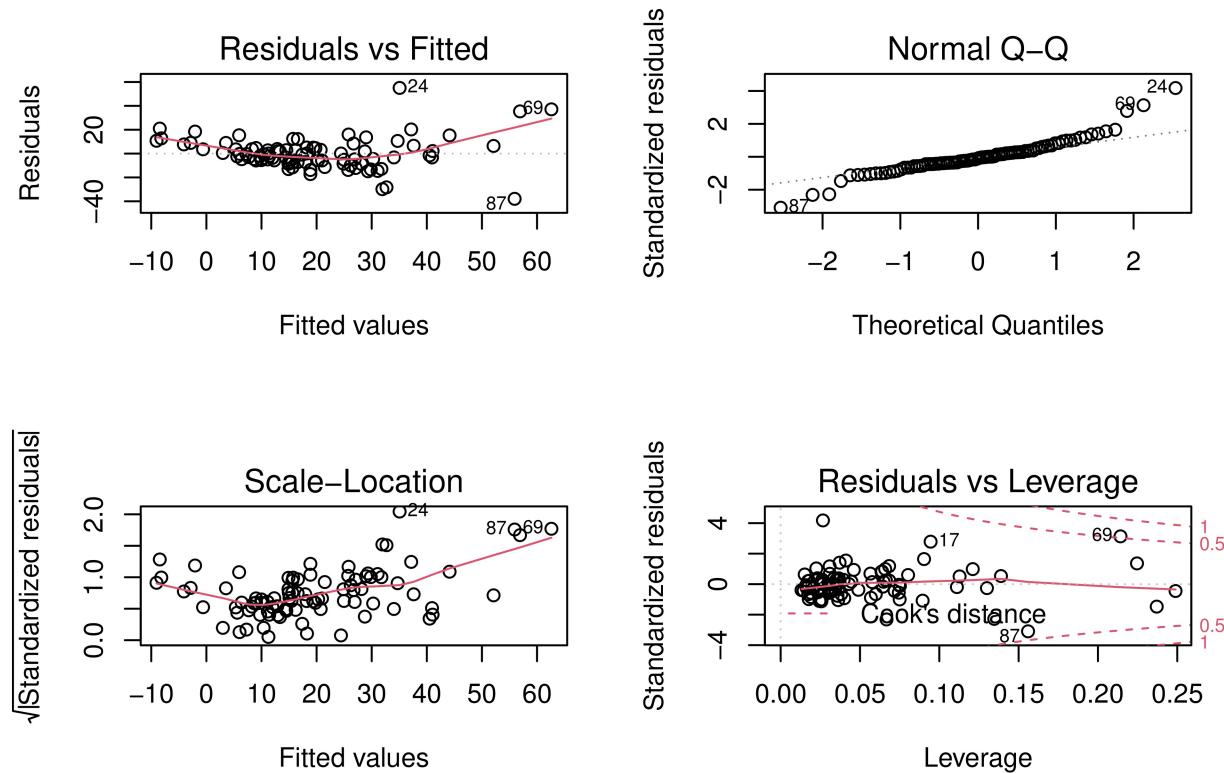
```

To select the best model, we need to look at these 3 statistics. The  $C_p$  value must be close to the number of predictors, The BIC, should be a low value and the Adjusted  $R^2$  should be as high as possible. Based on the following statistics, it is easy to conclude that the model with 4 predictors is the best model to fit our data.

```

par(mfrow = c(2,2))
plot(stock_lm4)

```



Hence, this is the equation for the model

`ev_ttm_multiple = -23.888 + 40.315Xrevenue_growth + 11.526Xsales_efficiency - 16.316xebitda_margin + 32.586Xgross_margin.`

From this equation, ebitda margin will negatively affect the ev\_ttm\_multiple while all others will affect it positively. To be more precise, for every %increase in revenue\_growth, sales\_efficiency and gross\_margin , the ev\_ttm\_multiple increases by 40.315, 11.526 and 32.586 respectively WHILE the ev\_multiple decreases by 16.316 with a percent increase in the ebitda\_margin.

We proceeded to train our model based on an 80%-20% division in order to calculate the MSE in order to get its prediction accuracy and perform cross validation. Here are our findings

```

set.seed(10)
errors = rep(0,10)

#TO GET MSE'S
for (i in 1:10){
  sample = sample.int(n = nrow(StockData), size=round(.80*nrow(StockData),0),replace=F)
  training_data = StockData[sample,]
  testing_data = StockData[-sample,]
  my_fit = lm(ev_ttm_multiple~revenue_growth+sales_efficiency
              +ebitda_margin+gross_margin,data=training_data)
  errors[i] = mean((StockData$ev_ttm_multiple - predict(my_fit,StockData))[-sample]^2)
}
errors

## [1] 144.97062 321.90086 75.84935 139.90945 269.98391 221.93285 207.32634
## [8] 242.54906 205.05763 285.90040

mean(errors)

## [1] 211.538

```

Upon doing 10 iterations, we arrived at a wide variety of MSE values. MSE are important because it is a means of showing how far off in prediction are we. Statistically, it tells us the square of the difference between the predicted and the actual observation. Our MSE values range from a low of 75.8 to a high of 285.90 which shows how fairly accurate the model is in predicting the ev\_ttm\_multiple.

## DECISION TREES

Decision trees are easier to interpret than regression models. That is the main reason why we are using a decision tree. However we have to remember that decision trees perform worst in prediction than most regression model. This code makes a tree

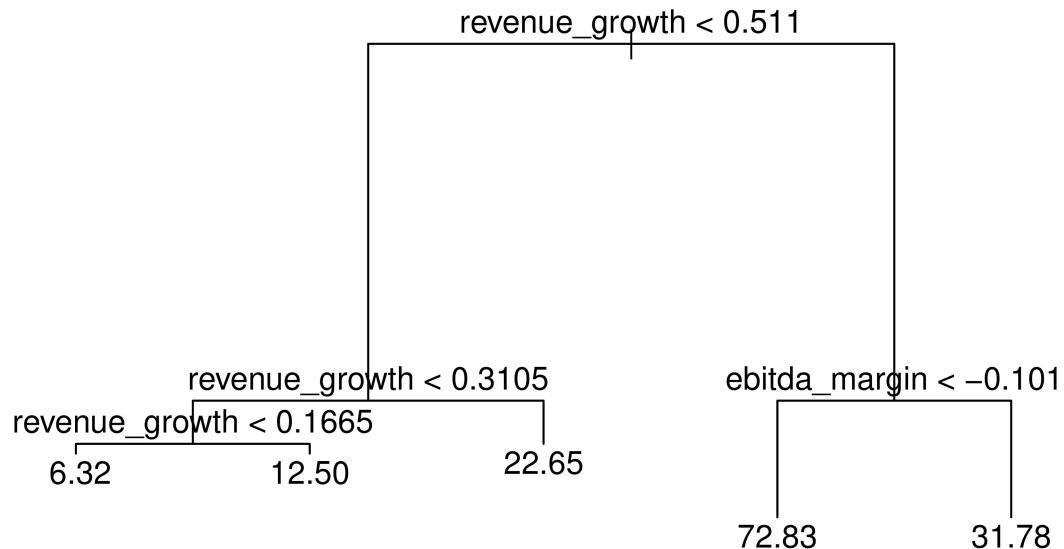
```

set.seed(20)
library(tree)
tree.model <- tree(ev_ttm_multiple~revenue_growth+sales_efficiency+gross_margin+ebitda_margin, data=training_data)
summary(tree.model)

##
## Regression tree:
## tree(formula = ev_ttm_multiple ~ revenue_growth + sales_efficiency +
##       gross_margin + ebitda_margin, data = training_data)
## Variables actually used in tree construction:
## [1] "revenue_growth" "ebitda_margin"
## Number of terminal nodes: 5
## Residual mean deviance: 127.2 = 8525 / 67
## Distribution of residuals:
##    Min. 1st Qu. Median Mean 3rd Qu. Max.
## -33.830 -4.715 -1.210 0.000 4.050 27.620

```

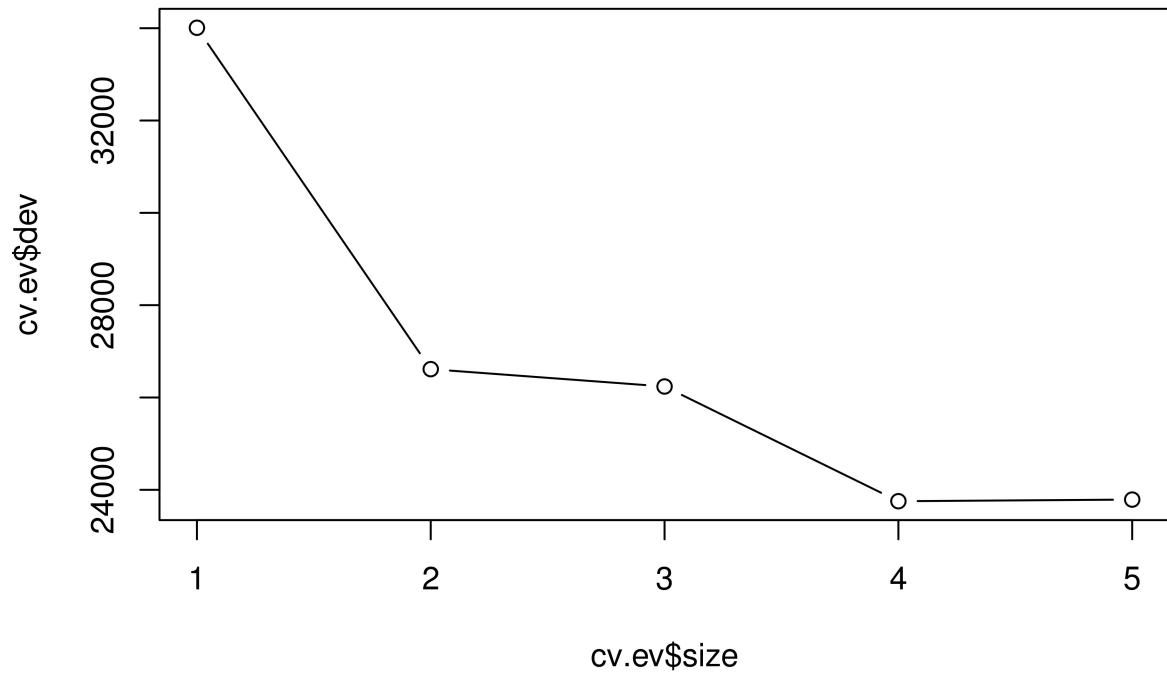
```
plot(tree.model)
text(tree.model)
```



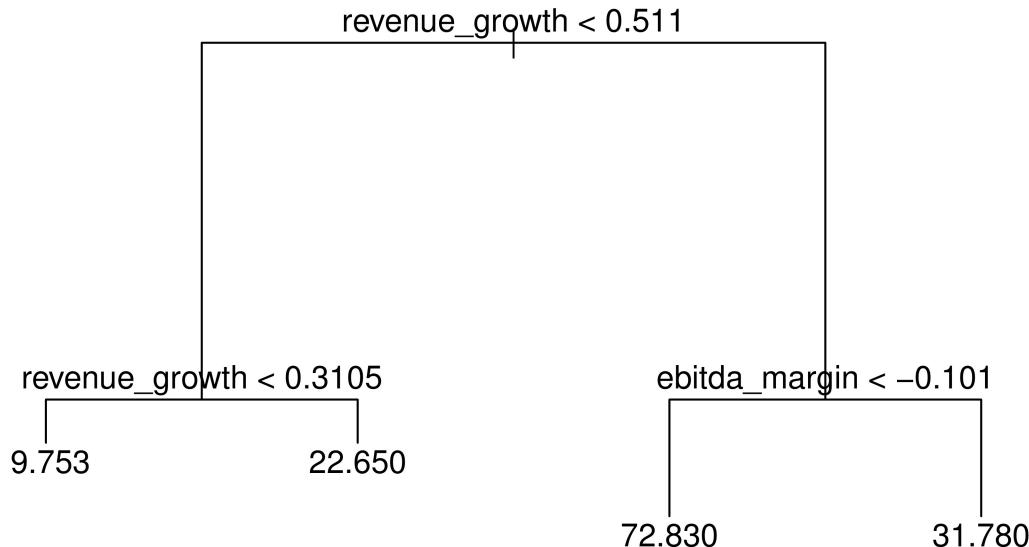
## Prune Tree

Then we proceeded to prune the tree at the best node using this code.

```
cv.ev=cv.tree(tree.model)
plot(cv.ev$size, cv.ev$dev, type="b")
```



```
pruned_tree = prune.tree(tree.model, best = 4)
plot(pruned_tree)
text(pruned_tree)
```



We can see that `revenue growth`, `sales_efficiency` and `ebitda_margin` are the predictor that affect the most our data.

MSE Next we are interested in check how well this model predict our data in a goal to compare it to our linear regression model.

The following code does a cross validation on the tree model

```

library(tree)
mse.tree = rep(0, 10)
for (i in 1:10){
  set.seed(i)
  sample<- sample.int(n=nrow(StockData), size=round(.80*nrow(StockData), 0))
  train<-StockData[sample, ]
  test<-StockData[-sample, ]
  tree.model <- tree(ev_ttm_multiple~revenue_growth+sales_efficiency+gross_margin+ebitda_margin, data=train)
  yhat.tree=predict(tree.model, newdata=test)
  mse.tree[i]=mean((test$ev_ttm_multiple - yhat.tree)^2)

}
mean(mse.tree)

## [1] 252.2586
  
```

The obtained MSE is 252.2586 which is much higher than the one we have for the linear regression model.

## Random Forest Model

Although three-based methods are simple and easier to interpret, they lack the accuracy provided by regression models. For that reason, we decided to use a random forest model to add more accuracy to our prediction. A random forest is a good choice because it reduces the variance introduced in decision trees and tree bagging. However, using random forest comes at a cost of longer training periods thus there is a higher computation time and we cannot see the threes so we loose the simplicity of single decision trees.

The following code does a random forest using ev\_ttm\_multiple~revenue\_growth+sales\_efficiency+gross\_margin+ebitda\_margin

```
p=4
B=100
library(randomForest)

## randomForest 4.7-1

## Type rfNews() to see new features/changes/bug fixes.

rf.model=randomForest(ev_ttm_multiple~revenue_growth+sales_efficiency+gross_margin+ebitda_margin, data=StockData)
rf.model

##
## Call:
##   randomForest(formula = ev_ttm_multiple ~ revenue_growth + sales_efficiency + gross_margin + ebitda_margin, data = StockData)
##   Type of random forest: regression
##   Number of trees: 100
##   No. of variables tried at each split: 1
##
##       Mean of squared residuals: 228.528
##       % Var explained: 38.15
```

From this preliminary model, we get an MSE of 228.22 and explain about 38.23% of the variance in the data set. Seeing that this value is lower than the MSE for the decision tree, we go on making a training and a testing set.

```
library(randomForest)
mse.randforest = rep(0, 10)
p=4
B=100
for (i in 1:10) {
  set.seed(i)
  sample<- sample.int(n=nrow(StockData), size=round(.80*nrow(StockData)))
  train<-StockData[sample, ]
  test<-StockData[-sample, ]
  bag.model=randomForest(ev_ttm_multiple~revenue_growth+sales_efficiency+gross_margin+ebitda_margin, data=StockData)
  yhat.randforest=predict(bag.model, newdata=test)
  mse.randforest[i] = mean((test$ev_ttm_multiple - yhat.randforest)^2)
}

mean(mse.randforest)

## [1] 217.2678
```

The MSE obtained is 217.2678 which is still lower than decision tree MSE, proving that random forest is better in prediction than decision trees.

## Boosting

Note: This is not part of our project but we were curious to see how a boosting model would perform compared to our other models so proceeded on making a model. The following does a boosting and show the mse

## Conclusion

The multiple linear model was our best model as it had the lowest MSE. The linear regression was able to give us more insight into which predictors are more associated with a higher revenue multiple. Revenue\_Growth was our best predictor and sales efficiency was the second best. Decision Tree was the best model in terms of inference. We were surprised that the Random Forest has performed worse than the Linear Regression. It

is unreasonable to expect any model to predict the value of a company with high precision, however that are clear patterns that our models bring to light. It is logical that revenue growth plays a large role in the value of a SaaS company, considering that revenue is a company's oxygen. One of our most surprising finding was to see that ebitda\_margin % was somewhat negatively correlated with a higher multiple. A lower ebitda\_margin indicates that a company is "losing money," but it may be doing so to prioritize revenue growth which we found to be the most predictive. There is an expression in SaaS - grow at all costs and that may be true. Though it was not part of our plan, we performed Boosting and we found out that it was the best model for the data.

To tie everything back to why we chose to explore this data, every one of our group members plans to work professionally in software after graduating. Stock is often big part of compensation packages. If any of us are considering a job offer from a SaaS company we will definitely be looking through their financials to see what kind of revenue multiple they should have.

## Endnotes

Analysis for the linear regression model was performed by Oluwatobiloba Oladunjoye, Michael Thomas. Analysis for the random forest and boosting regression model was performed by Mahamadou Dagnoko and Min Shwe Maung Htet. Decision Trees was performed by Rafay Alam and John Jackson. To view our data set and source code, visit our GitHub repository [https://github.com/johnjackson59/saas\\_multiples](https://github.com/johnjackson59/saas_multiples).