

**Analyzing Company Financials to Classify Executive Titles**

Kevin W.S. Baum and Stephen F. Reagin

Shiley-Marcos School of Engineering, University of San Diego

## Introduction

Top executives at publicly-traded companies are among the most highly compensated workers in American society, with total compensation packages often reaching 7 figures or more on an annual basis. The subject of executive pay is also contentious in popular media, where it is often addressed in the context of pay or wealth inequality.

Calls for greater transparency of executive pay have led to disclosure requirements, including both the amounts and performance conditions tied to the pay packages for top executive officers. Public companies must disclose, on an annual basis, the amounts paid to their top executives and the justifications for doing so. And as recently as August 2022 the SEC has adopted new rules requiring disclosure of the connection between company performance and executive pay (United States SEC, 2022).

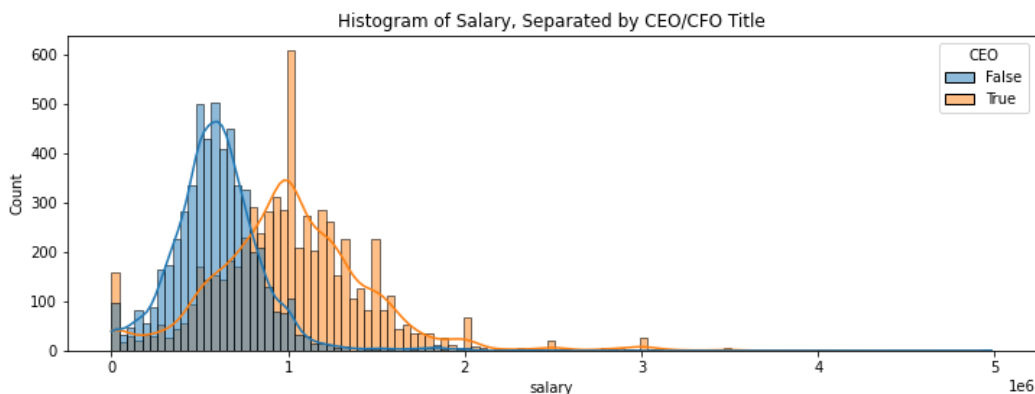
Looking through the *proxy statements* of public companies, which is the document in which this legally-required information is disclosed, one of the most common justifications for large executive pay packages is that higher pay should be tied to higher performance. However, performance can be defined differently for semiconductor stalwarts like Intel and AMD versus beverage companies like Keurig Dr Pepper and Coca Cola. That said, generally speaking, the most commonly agreed-upon metrics for overall company performance tend to focus on revenue and income growth.

A review of the academic and business literature in this space reveals common research themes: (a) the growth of executive pay over time, (Bebchuk & Grinstein, 2005), (b) the effect of increased transparency (Bruce & Skovoroda, 2015), and (c) consultancy services advising companies on how to tie executive pay to company performance (Abbott et al., 2019).

We would like to investigate the relationship, if any, between executive pay and company performance using a data-driven approach. However, it's not clear at the outset that the techniques of machine learning and statistical mining are the appropriate tools for such an investigation. For example, we know CEOs tend to earn more than their CFO counterparts at most public companies with rare exceptions, and this is evidenced by the distribution of salaries (among other pay components) as seen in Figure 1.

**Figure 1**

*Histogram of Salary, Separated by CEO/CFO Title*



If we have a list of pay packages given to CEOs and CFOs for the S&P 500 companies, the classification algorithm is relatively simple for a human observer: (1) look at the two people listed for each company, (2) identify which person is paid more, and (3) label the person with higher pay as the CEO. We want to know whether computers can perform comparably using the standard Python tools of machine learning.

Therefore, the subject of our report is to investigate the feasibility of using mathematical tools to answer such questions. In particular, we ask: “can machines tell the difference between CEOs and CFOs on the basis of their pay packages?”

### Overview

With the course's focus on classification and the encouragement in the instructions to do a project involving classification, we decided to create a classification model that would predict a CEO (Chief Executive Officer) versus a CFO (Chief Financial Officer), using pay components and financial data as predictor variables. This is not a dataset that we found on Kaggle or a site that specifically curated it for the purposes of using and discovering classification techniques. Rather, it is a completely unedited data scrub using SEC data. The process to discover our question was iterative, and we relied on exploratory data analysis and combining our intuitions in order to carve our path. In doing so, our overall goal was to demonstrate knowledge gained in ADS502 by using key learnings from the course to apply to "wild data."

With this in mind, we started with two datasets, one for executive pay, and the other for company financial data, extracting 10 years of quarterly "trailing 12 month" data from current S&P 500 companies into a Pandas data frame. We then grouped the data into average values for every year based on its ticker. For example, Apple will have average values specific to its financials for each year over the last 10 years. Next, we moved to the executive pay dataset, and extracted the executive information for those individuals within the S&P 500 ticker list over the last 10 years. After this, we combined the two separate data frames, left-joining them via their tickers so that financial information was duplicated for multiple executives by company and year. We then performed cleaning, pre-processing, and exploratory data analysis (EDA) before sub-setting the data frames by CEO and CFO, and finally concatenating them together.

Once this data frame was ready, we partitioned the data into training and testing datasets. Next, we next scaled the data because of the large numbers involved when summing over ten years. We proceeded to run CART, C5.0, random forest, logistic regression, naïve bayes, K Neighbors, and support vector machine models and scored them with evaluation metrics.

### **Initial Data**

The ultimate source of all executive compensation information is the United States Securities and Exchange Commission (SEC). Every publicly-traded company must file an annual proxy statement with the SEC for shareholders to review, among other things, the total compensation packages paid to the CEO, CFO, and next three most highly-paid officers, a group is often referred to as the “Top 5 executives.” Of particular importance is the “Summary Compensation Table” which is also required by the SEC, and for all Top 5 officers the company must disclose three years of:

- Salary
- Performance bonus (“Non-Equity Incentive Compensation”)
- Stock and Option awards
- Off-cycle or one-time bonuses
- Other services like personal security, air travel, etc.

However, this data is inconvenient to source directly from .html files on the SEC because of the manual labor involved. We instead pulled data from sec-api.io, a third-party source, using API calls under a paid academic license. This raw dataset includes more than 170,000 observations of Summary Compensation Table data rows for companies in the Russell 3000,

dating from 2006 to 2022. We also feature engineered Boolean title columns (CEO, CFO, Interim) to flag specific roles. We further narrowed our universe to include only the 485 members of the S&P 500 for which we have approximately 28,000 reliable data observations.

We sourced company financial information from stockrow.com using their free Excel-based worksheets of Income Statement information, specifically targeting Revenue and Net Income as the most common metrics of financial performance. Table 1 shows the descriptive statistics of the S&P 500 subset.

**Table 1**

*Descriptive Statistics of Raw Dataset*

Feature	Observations	Median	Mean
Salary	28440	\$618K	\$708K
Bonus	28440	\$0	\$247K
stockAwards	28440	\$1,425K	\$2,788K
optionAwards	28440	\$93K	\$5,133K
nonEquityIncentiveCompensation	28440	\$612K	\$1,054K
otherCompensation	28440	\$61K	\$370K
Total	28440	\$4,083K	\$6,460K
Revenue	28391	\$8.79B	\$22.0B
Net Income Common	28391	\$787M	\$1.99B

### **Data Preparation and Exploratory Analysis**

The first step was to get the data we needed from the SEC API for executive pay, and from Stockrow.com for the company financials. After gathering the data from both sources, we decided we wanted to use only S&P 500 data and so we also loaded a data source containing the 503 tickers. We matched the data from the company financials to only return data for each ticker in the list of 503. For each ticker in the list, we created CSV files and then combined all the CSV files into a Pandas data frame to perform manipulations on. Now that our data frame

was ready, we narrowed the fields into just ones we wanted: Ticker, Date, Revenue, Gross Profit, Operating Income, Income after Tax, and Net Income Common. Additionally, we created a datetime object from the “date” feature and extracted the year out of it. This let us ultimately break down the information by yearly basis so that we could establish yearly trends.

Next, we grouped the data by ticker (representative of company) and year and got the mean values for each ticker per year. To make sure this was done correctly, we manually checked “AAPL,” Apple, and found that the numbers matched when we did an independent search online. From here, we were ready to join in the executive pay data set and did so based on the ticker; only executive pay information for the S&P 500 companies were included. Note that not all companies were in the executive pay dataset, and after merging them together only 485 companies remained with 18 not present in the executive pay from the SEC API.

For data cleaning, we noted some duplicate records and performed a standard `drop_duplicates()` method to eliminate these duplicates. The company Agilent had duplicate records that were duplicates in reality, but had only one slight difference that caused them to get past the `drop_duplicates` method. We eliminated these duplicates by deleting the column that was supposed to contain full names but instead contained the words “human resources” and “group.” We noted about 200 outliers in the executive pay data and dropped these as well. One such outlier was a salary of 18 billion US dollars in one year.

Additionally, there were also some missing values in the company financials. For instance, 73 of the 485 companies did not have revenue data in 2012 and 14 did not have revenue data in 2021. We decided not to impute values or drop these companies because our

classification model was classifying executives as either CEOs or CFOs. The more important metrics, in our eyes, were individual pay statistics.

For exploratory data analysis (EDA), we aimed to find patterns in the pay of CEOs over the last 10 years.

**Figure 2**

*CEO Salary Percentiles 2012-2021*

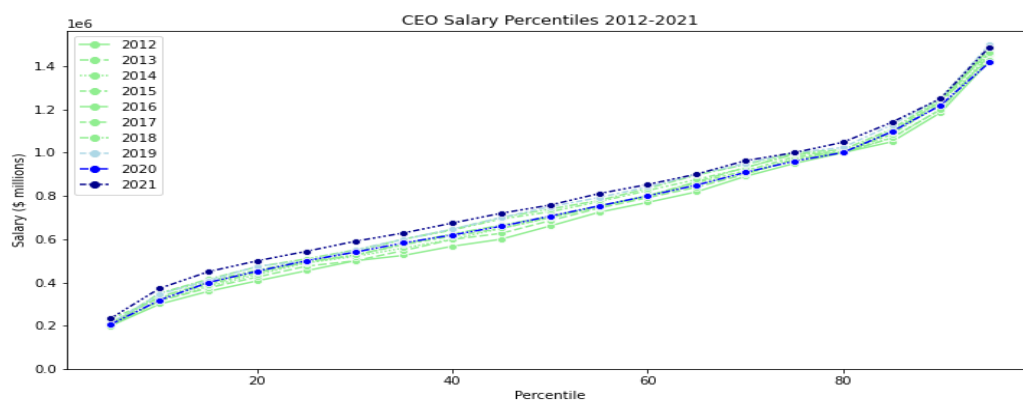


Figure 2 is constructed with the percentiles of CEO salary on the x-axis with the actual salary number on the y-axis. This means if a salary is the highest in the dataset, it will be in the 99<sup>th</sup> percentile, and if a salary is the lowest in the dataset, it will be in the 1<sup>st</sup> percentile. From this visualization, we can see that there is a consistent slope until roughly the 80<sup>th</sup> percentile mark is reached. This means that the increase in income from say the 40<sup>th</sup> percentile to the 60<sup>th</sup> percentile is reliable and predictable. This is not the case after about the 80<sup>th</sup> percentile however, where we see a steepening of the slope of the line for all years. This indicates that those CEOs above the 80<sup>th</sup> percentile earn much more than those below them. A comparable increase in earner percentile, say 20-percentile increase, from the 40<sup>th</sup> to 60<sup>th</sup> is much less in real dollar value than a 20-percentile increase from the 70<sup>th</sup> percentile earner to the 90<sup>th</sup>



percentile earner. Additionally, we highlight the lines for 2020 and 2021 to show the difference between these two years and note that 2021 is much higher for salary than 2020 for any given percentile chosen.

We have also created boxplots, violin plots, kernel density estimate (KDE) plots, a color mapped kernel density estimate (KDE) plot, and bar plots to investigate a multitude of factors. Our focus was on the changing CEO pay landscape, and we used visualizations to investigate changes in salary, stock awards, non-equity incentive compensation, and total compensation. We found a trend of CEO salaries rising in the last 16 years (since 2006), and the most visually telling representation of this trend was the KDE plot in Figure 3.

**Figure 3**

*KDE Plot for CEO Salary Under \$3M*

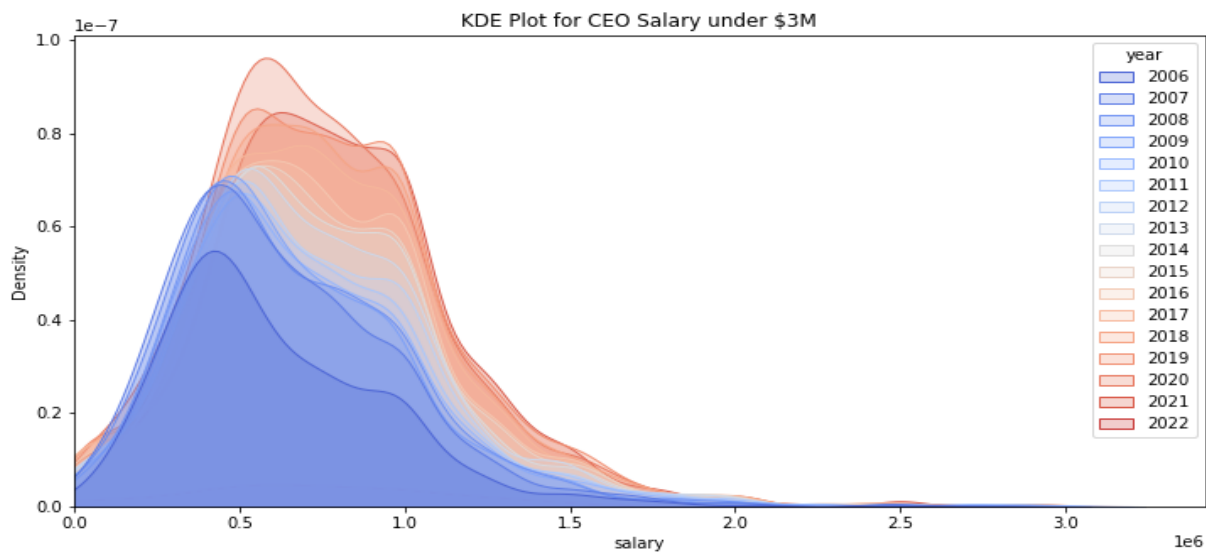


Figure 3 displays the density of incomes from 2006-2022 with blue marking the earlier years and red marking the more recent years, showing a clear trend of rising salaries over time.

When analyzing stock awards via box plot in Figure 4, we see a surprising widening of the interquartile range (IQR) spread in 2021, following a trend of increased IQR over time.

**Figure 4**

*Boxplot Distribution of Stock Awards Less Than \$30M*

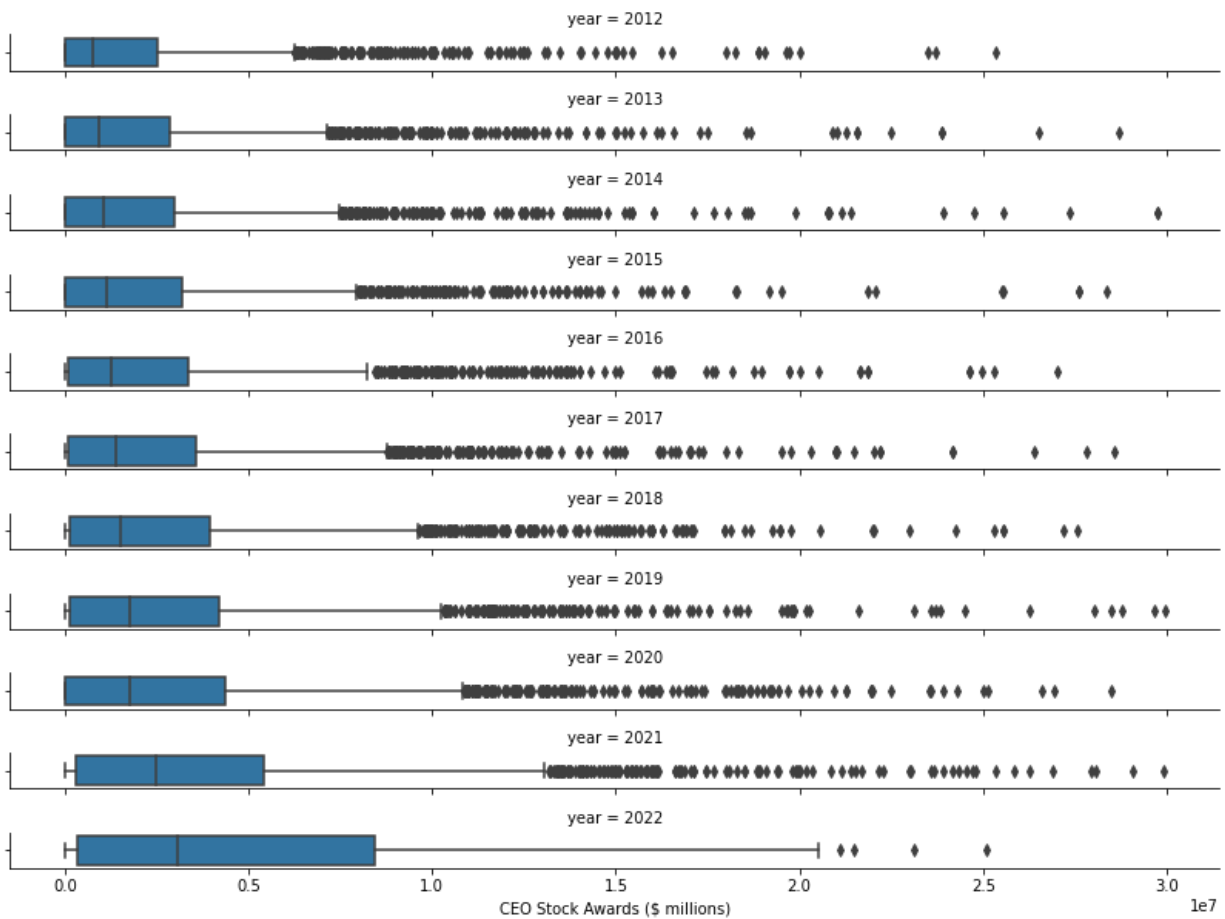


Figure 4 displays the overall trend of widening IQR ranges between 2012 and 2021, which may be due to stock awards being a more common form of compensation now compared with years past. Such a short box with extreme positive skewness in 2012 indicates that a high proportion of CEOs had smaller stock awards as part of their compensation package. In fact, for all years the data is very positively skewed, which indicates that the mean is larger than the median as evidenced in Table 1 and is heavily impacted by outlying very large values.

This is a logical conclusion based on the observation that the top 20% of CEO earners are making much more than the lower 80%. The large stock awards going to these top earners drags up the mean and results in a high positive skewness.

Finally, it is worth noting that companies vary wildly in how they pay their top earners. For instance, Microsoft has a very exaggerated change the year that they elected to have Satya Nadella become the new CEO in 2014. In years prior, the CEO position received relatively low compensation from the company. This is because Bill Gates and Steve Ballmer, the prior CEOs, were original founders and had already earned many tens of billions of dollars, thus any annual compensation was considered unlikely to motivate increased performance. But when it was time to make a change, the chief executive position started to make thousands of times the amount as before. Microsoft is only one example, and it is conceivable that other CEOs that founded their companies likewise have decided to have very low compensation as well.

To address this inherent variability in executive pay and financial performance, our final training dataset consisted of the *sum* of compensation components and financial metrics over the previous 10 years for each company. In this way, even if transitions occurred, we were still capturing the total amount which the company paid to *the CEO role* over the relevant period.

### **Data Mining Preliminary Results**

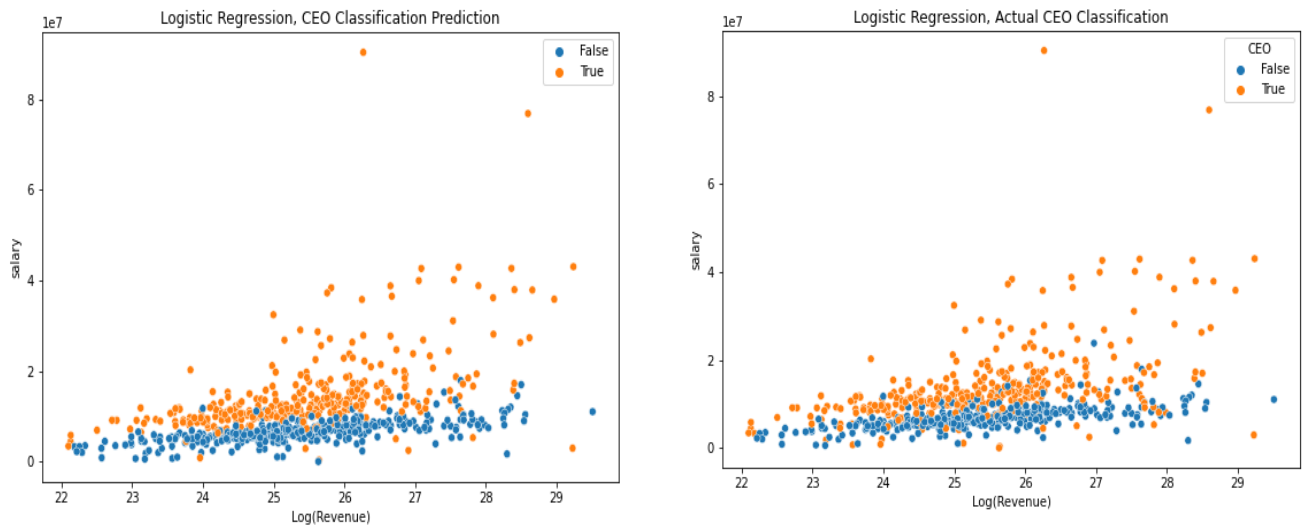
Our basic goal for model selection was to follow the textbook's advice, "to select the model that shows lowest generalization error rate" (Tan et al., 2019). We considered other possibilities for optimal selection, such as minimizing the count of false positives or false negatives. Ultimately, we felt that a generalized accuracy measure would work well for this type of data, as there is not a need to be extra conservative like there would be for, say,

diagnosing whether a patient's health indicators classify them as having a disease or not having it. In such a case, it would be better to have false positives (patient does not really have the disease but tests say they do) rather than false negatives (patient really has the disease but tests say they do not).

To evaluate our models, we partitioned the data to “assess the performance of a learned model on a labeled test set [which] has not been used at any stage of model selection” (Tan et al., 2019). We felt that the data set was sufficiently small to follow the book's advice for smaller or less complex datasets, where “one should retain sufficient records for accurate assessment, so that the training sets would contain only 50–67% of the original data” (Tan et al., 2019). Our main concern was balancing the proportion so that the training model could effectively learn the behavior but not be so large that the testing model's error metric was unreliable because it was not able to compute over enough instances. Ultimately, we decided on a 70% partition into the training set and 30% into the test set.

We tried different models to see which would produce the best results for reducing our error rate, producing the best accuracy for true positives. The model which we felt was best for this purpose was the logistic regression model.

## Figure 5

*Comparison of Classification Prediction Versus Actual*

Many of the models tended to overfit the training data. For instance, random forest produced a higher overall accuracy, but we felt it overfit the training data with an accuracy of 95% on the training data with 84% on the testing data. The CART and C5.0 models had lower accuracy for both training and testing than the logistic regression model, while also being slightly overfit in the training data. Naïve Bayes, KNeighbors, and SVM produced good results, but also had more of a difference in accuracy between the training and testing datasets than the logistic regression model. Ultimately, logistic regression works well for the general purpose for our model because it is a “discriminative model for classification that directly computes the poster[ior] probabilities without making any assumption about the class conditional probabilities” making it “generic” and “applicable in diverse applications” (Tan et al., 2019).

The baseline condition was 52%, representing the proportion of CEOs to CFOs in the dataset. While all of the models in Table 2 represent a significant improvement above the baseline 52% accuracy, logistic regression was our preferred choice.

**Table 2***Comparison of Model Accuracy Results*

Model	Train Data Accuracy	Test Data Accuracy
CART	87%	83%
C5.0	87%	80%
Random Forest	95%	84%
Logistic Regression	90%	89%
Naïve Bayes	87%	83%
KNeighbors	91%	87%
Support Vector Machine	91%	88%

**Further Discussion**

Classifying CEOs and CFOs is only scratching the surface of the type of information this data source is capable of. Our initial goal was to solve the more politically hard-hitting question of, “do executives actually get paid based on their companies do or do they reward themselves no matter what the company does?” We wanted to create a classification model that would predict if a CEOs or CFOs compensation trends would correctly classify if a company was a rapid grower, flatliner, or decliner. Even though the 2010s was a very prosperous time in the American economy, we found that about 70 of the companies in the S&P 500 had declining revenues during the period. One of the beautiful elements of the economy though is that growth is limitless, while the maximum amount of decline possible is just 100%. Thus, due to the fantastic success of companies like Tesla and Amazon, the growth of the overall economy far outweighed the companies bringing down the overall index.

Beyond classification, the data set also has potential for linear regression projects, as well as unsupervised algorithms like clustering. For example, a linear regression project could investigate which factors are most important to overall compensation for executives. Variables

such as company revenue, company net income, position title, age, location, and many more could provide indicators as to which factors affect executive pay statistics the most. Clustering algorithms could be used to find trends that are worth exploring for supervised algorithms that are not known from an initial view of the data.

### References

- Abbott, S., Aksoy, M., Groysberg, B., & Marino, M.R. (2019, February). Compensation Packages That Actually Drive Performance. Principles for designing executive pay. *Harvard Business Review*. <https://hbr.org/2021/01/compensation-packages-that-actually-drive-performance>
- Bebchuk, L., & Grinstein, Y. (2005). The Growth of Executive Pay. *Oxford Review of Economic Policy*, 21(2), 283–303. <https://doi.org/10.1093/oxrep/gri017>
- Bruce, A., & Skovoroda, R. (2015). *The Empirical Literature of Executive Pay: Context, the Pay-Performance Issue and Future Directions*. Nottingham University Business School. [http://irep.ntu.ac.uk/id/eprint/35402/1/12805\\_Bruce.pdf](http://irep.ntu.ac.uk/id/eprint/35402/1/12805_Bruce.pdf)
- Tan, P. N., Steinbach, M., Karpatne, A., & Kumar, V. (2019). *Introduction to Data Mining* (2nd ed.). Pearson.
- United States Securities and Exchange Commission. (2022, August 25). *SEC Adopts Pay Versus Performance Disclosure Rules*. <https://www.sec.gov/news/press-release/2022-149>







			id	clk	ticker	name	position	year	salary	bonus	stockAwards	optionA
			0	73b3a60ba2037432f0c068330c96b70eb66	1090872	A	Sam Raha	Senior Vice President, Diagnostics	2021	563500	0	1541332
			1	97393f60cd9f32f1650e472673d8aa70c	1090872	A	Michael R. McMullen	Chief Executive Officer	2021	1280000	0	9165390
			2	83ba9cc2bca477b16e23a0cbe56e70c66	1090872	A	Jacob Thaysen	Senior Vice President, Life Sciences	2021	625000	0	1812285
			3	7c60bb804071675a6e15ec930f6dea190	1090872	A	Robert McMahon	Senior Vice President, Chief Financial Officer	2021	663500	0	2291271
			4	259b70a8b9ef67c61a2c3d4ef47344a	1090872	A	Padraig McDonnell	Senior Vice President, Chief Cross-Lab Group	2021	495000	0	1249771
			10	6a9fd7caaf9d2f19409e40f8c289efcd3	1090872	A	Jacob Thaysen	Senior Vice President, Life Sciences	2020	614904	0	2117779
			11	cccf820f85897ef63a1266f4a096840	1090872	A	Michael R. McMullen	Chief Executive Officer	2020	1172853	0	11190749
			12	0a8f57350e9f3b644a4f98024090ca25f	1090872	A	Robert McMahon	Senior Vice President, Chief Financial Officer	2020	638333	0	2614535
			13	5c5621b65ca11f037e274c80b19fcd2	1090872	A	Sam Raha	Senior Vice President, Diagnostics and	2020	538782	0	1673282
			14	4deeb4d5b673a6f53e85a86af1eddd5990	1090872	A	Genomics Group	Michael Tang Senior Vice President, Co...	2020	561026	0	1566688
			15	9637116a014dc78e135af065a0831e78	1090872	A	Mark Doak	Senior Vice President, President, Agilent Cross...	2020	316538	500000	2614535
			22	e22445c4e02f15347b2ba0cbee5789855	1090872	A	Michael R. McMullen	Chief Executive Officer	2019	1220833	0	9560588
			23	ca8648a0544fa160a9e44b6a20d2a1e482	1090872	A	Robert McMahon	Senior Vice President, Chief Financial Officer	2019	610000	1000000	2221067
			24	55ef5c9c8b7f5fe16ed7c7b2fe70659d6	1090872	A	Genomics Group	Michael Tang Senior Vice President, General Co...	2019	545833	0	1351856
			25	d41442da4b3801c123196f9932d80c2bfe	1090872	A	Mark Doak	Senior Vice President, President, Agilent Cross...	2019	597917	0	2240337
			26	93f8d0f60316eaf16044439bc3c141b3d9	1090872	A	Jacob Thaysen	Senior Vice President, Life Sciences	2019	597917	0	1858905
			27	8830ad89737780b9a6e9833a60265300	1090872	A	Michael Tang	Senior Vice President, General Counsel	2019	545833	0	1351856
			34	55f648eec48f4dcedc4ca6a5994f0481	1090872	A	Jacob Thaysen	Senior Vice President, President, Life Sciences...	2018	556250	0	1845713
			35	3d29f1ea2763dfe20f44d6e0f07da9	1090872	A	Robert McMahon	Senior Vice President, Chief Financial Officer	2018	147879	2000000	4137198
			36	9916458088f7667c5d56319a065a07e1	1090872	A	Michael R. McMullen	Chief Executive Officer	2018	1168750	0	8476147
			37	93f6834cfe5a5cd3dc174d07a9b167864	1090872	A	Human Resources	Didier Hirsch Senior Vice President, Ch...	2018	648333	0	2808278
			38	8dafe022270355f6a9421f9c86c78437	1090872	A	Mark Doak	Senior Vice President, President, Cross-Lab Group	2018	570833	0	2088406
			39	a8e9654000c2e2b96c1e9cee384f182e8	1090872	A	Dominique Grau	Senior Vice President, Chief Cross-Lab Group	2018	470825	0	1276453
			46	e048daf0786cf1a0b7e26a2d48f6cb99	1090872	A	Human Resources	Didier Hirsch Former Senior Vice President, Ch...	2017	627500	0	2676250
			47	2a318ae1fb2ac86daf17486a49563820	1090872	A	Michael R. McMullen	Chief Executive Officer	2017	1095833	0	7205347
			48	b846615fa7b6357bd5a5c48b468696f	1090872	A	Patrick Kaltenbach	Senior Vice President, President, Life Sciences...	2017	532083	0	1787948
			49	9d5fef9b9e9e99e0c38903c125f8d0895	1090872	A	Didier Hirsch	Senior Vice President, Chief Financial Officer	2017	627500	0	2676250
			50	8f0626122d95f4a0d7e4da9d3d5181b14	1090872	A	Mark Doak	Senior Vice President, President, Cross-Lab Group	2017	520833	0	1899068
			51	1fec88b3f1ebccdc16e98aa3b3b71d	1090872	A	Jacob Thaysen	Senior Vice President, Chief Diagnostics and	2017	495000	0	1338104
			58	a13b28a1a078d707abcbba143b7abc66	1090872	A	Didier Hirsch	Senior Vice President, Chief Financial Officer	2016	600000	0	2474816
			59	a1e6892081434400721bc17714d9b2d64	1090872	A	Jacob Thaysen	Senior Vice President, President, Diagnostics and	2016	436667	0	845472
			60	2d3f7dacc047db8fd0c3fcacfe73d701a3	1090872	A	Applied Markets Group	Jacob Thaysen Senior Vice President, President, Cross-Lab Group	2016	436667	0	845472
			61	65622904851a0f79530120e2dfb6b57	1090872	A	Human Resources	Didier Hirsch Former Senior Vice President, Ch...	2016	600000	0	2474816
			62	6e25b30af9d2f16e44b1c1d67072833	1090872	A	Michael R. McMullen	Chief Executive Officer	2016	1041667	0	6341255
			63	ece95b49afb267f963e22254db3749d	1090872	A	Mark Doak	Senior Vice President, President, Cross-Lab Group	2016	470833	0	1638114
			64	733eaaac28284cb81dd47f93f5edf955f	1090872	A	Patrick Kaltenbach	Senior Vice President, President, Life Sciences	2016	487976	0	1479571
			72	bd78dd68f1458ba8a44a48be1edf173	1090872	A	Patrick Kaltenbach	Senior Vice President, President, Life Sciences...	2015	407262	0	505138
			73	88ba9193971115c64b4338c80a0c369	1090872	A	Didier Hirsch	Senior Vice President, Chief Financial Officer	2015	600000	0	1049012
			74	73bacfb18a46e148680eb4a9ebcc6c2344	1090872	A	Chain	Patrick Kaltenbach Senior Vice President, Pres...	2015	407262	0	505138
			75	1fa34df5935e9de343892077f0ca48ad	1090872	A	Michael R. McMullen	President, Chief Executive Officer	2015	845833	0	2520205
			76	8338bde68a45b23c5892077f0ca48ad	1090872	A	Mark Doak	Senior Vice President, President, Cross-Lab Group	2015	422917	0	549457
			77	e96905309e246a7c50b160c798988614	1090872	A	Group	William P. Sullivan Former President and Chief...	2015	630000	0	3057116
			84	05a042698e63b1f1fd6c5e326c5791dfd	1090872	A	Didier Hirsch	Senior Vice President, Chief Financial Officer	2014	600000	0	1044826
			85	605f659e14e30f1fec5b6ba2041d28b4	1090872	A	Marie Oh Huber	Senior Vice President, General Counsel and Sec...	2014	525000	0	1007380
			86	74d8234da77e693bf7b0c888c50653825	1090872	A	Michael R. McMullen	President and Chief Operating Officer	2014	606250	0	2332309
			87	d3af951ba9f407439f8a44affe1b6c27f	1090872	A	Ronald S. Nersisyan	Executive Vice President and Keyight President...	2014	798333	0	2109864
			88	56fe7f3a486a52f083e1dded8482b048	1090872	A	Group	William P. Sullivan Former President and Chief...	2014	1050000	0	6632830
			89	a6a0b6e0d3212415b499a2961d5f1589	1090872	A	William P. Sullivan	Chief Executive Officer	2014	1050000	0	6632830
			96	649b860f6a4e4a654140ba6e0aa7c2c25	1090872	A	Group	William P. Sullivan Former President and Chief...	2013	1045000	0	3789936
			97	b4add3c41ed5c86ed22bfc5309856fa	1090872	A	Didier Hirsch	Senior Vice President, Chief Financial Officer	2013	597917	0	869456
			98	5ce6068d033adac55108770021fe1918	1090872	A	Designate, Keysight	Nicolas Rodier Former Senior Vice President, ...	2013	550000	0	668802
			99	9bd549134abbda364e60c272f6daa7	1090872	A	Michael R. McMullen	President and Chief Operating Officer	2013	575000	0	802571
			100	4ca37b05139280eb9d6436a55eb575	1090872	A	William P. Sullivan	President and Chief Operating Officer	2013	1045000	0	3789936
			101	494a3f0983f34c17b6917f520729069	1090872	A	and Diagnostics	Michael R. McMullen Senior Vice President, ...	2013	575000	0	802571
			102	c4d72990c76de80b5d28196a9696b69	1090872	A	Lars Holmkvist	Former Senior Vice President, Life Sciences	2013	938588	0	891721
			103	8f98917e1baf986ed5115cfd841b331	1090872	A	Ronald S. Nersisyan	Executive Vice President and Chie...	2013	741667	0	1560567

56 rows x 22 columns

```
In [28]: # Need to drop duplicate Agilent ones with the name "Group" or "human resources"
main_df = main_df.drop(main_df[(main_df.name == "Group") | (main_df.name == "Human Resources")].index)

# Creating the separate dataframes for CEO/CFPO so that we can perform aggregate
df_CEO = main_df[(main_df["CEO"] == True)]
df_CFO = main_df[(main_df["CFO"] == True)]

# Getting sum of values per ticker.
df_CEO = df_CEO.groupby('ticker').sum()
df_CFO = df_CFO.groupby('ticker').sum()

# Subsetting for the fields we want
df_CEO = df_CEO[['CEO','salary','bonus','stockAwards','nonEquityIncentiveCompensation',
'total','Revenue','Net Income Common']]
df_CFO.loc[df_CFO["CEO"] > 0, "CEO"] = True

df_CFO = df_CFO[['CEO','salary','bonus','stockAwards','optionAwards','nonEquityIncentiveCompensation',
'total','Revenue','Net Income Common']]
df_CFO.loc[df_CFO["CEO"] == 0, "CEO"] = False

In [29]: # Preparing to join the data frames back together
frames = [df_CEO, df_CFO]

# Joining the two data frames together.
new_main = pd.concat(frames)
new_main = new_main.sort_values(by=['ticker'])

# NOTE: The following commented out code gets percentages of executive pay as part of revenue.
# It is commented out because it is not used in the final classification model but is
# interesting as a potential variable in future projects with this dataset.
# new_main['salary/revenue'] = new_main['salary'] / new_main['Revenue']
# new_main['bonus/revenue'] = new_main['bonus'] / new_main['Revenue']
# new_main['stockAwards/revenue'] = new_main['stockAwards'] / new_main['Revenue']
# new_main['optionAwards/revenue'] = new_main['optionAwards'] / new_main['Revenue']
# new_main['nonEquityIncentiveCompensation/revenue'] = new_main['nonEquityIncentiveCompensation'] / new_main['R
# new_main['all_Equity_Awards'] = new_main['stockAwards'] + new_main['optionAwards']

## Dropping Where CEO is not equal to True or False
indexNumbers = new_main[ (new_main['CEO'] == 3) |
(new_main['CEO'] == 2) | (new_main['CEO'] == 5)].index
new_main.drop(indexNumbers , inplace=True)

## Converting True/False CEO column from object to boolean because the model
new_main['CEO'] = new_main['CEO'].astype(bool)
print(new_main.dtypes)
new_main

CEO          bool
salary      int64
bonus       int64
stockAwards int64
optionAwards int64
nonEquityIncentiveCompensation int64
total       int64
Revenue     float64
Net Income Common float64
dtype: object

Out[29]:
```

		CEO	salary	bonus	stockAwards	optionAwards	nonEquityIncentiveCompensation	total	Revenue	Net Inct Comr	
	ticker	A	True	11458269	0	68552878	15107596	16133895	115186782	5.834650e+10	7.676250e
	A	False	50895129	3000000	199646031	29996887		4580051	36893668	4.623975e+10	6.070500e
	AAL	True	2441140	0	130227654	0		10203331	159413533	4.842889e+11	1.304125e
	AAL	False	7218184	0	38471191	0		11200278	62707472	4.588179e+11	1.395625e
	AAP	True	9668318	3485000	44487383	14649081		8633582	83061159	9.050150e+10	4.4663842e
	...	...	...	...	...	...	...	...	...	...	...
	ZBRA	True	12679100	0	59775901	12967242		19148913	105531711	4.446959e+10	3.304858e
	ZION	False	6419692	3652500	5839992	985436		1808345	17642731	2.865551e+10	5.218222e
	ZION	True	21152988	11702000	15629354	3350878		7071579	68692342	7.047798e+10	1.245739e
	ZTS	True	10643533	0	51539576	20646680		16059793	101710673	5.388575e+10	9.552260e
	ZTS	False	7704090	26000000	18872778	7109989		6657597	45597108	8.033900e+10	1.366375e

959 rows x 9 columns

```
We have now finished preprocessing for our data to run through classification models

In [30]: new_main.to_csv('new_main_df.csv',index=False)
```







```
[27]: plt.figure(figsize=(12,3))
g = sns.FacetGrid(data=ceo_df[ceo_df['total'] < 50000000], row='year', aspect=10,height=1)
g.map(sns.boxplot, 'total')
Out[27]:
<Figure size 864x216 with 0 Axes>

year = 2006

year = 2007

year = 2008

year = 2009

year = 2010

year = 2011

year = 2012

year = 2013

year = 2014

year = 2015

year = 2016

year = 2017

year = 2018

year = 2019

year = 2020

year = 2021

year = 2022

0 1 2 3 4 5
total
1e7

In [28]: plt.figure(figsize=(12,9))
sns.violinplot(data=ceo_df[ceo_df['total'] < 50000000],x='total',y='year')
plt.xlim(0,5)

In [29]: print("Count of CEOs by year:")
print(states_df[starter_df['CEO'] == True] & (starter_df['Interim'] == False).groupby('year')['CEO'].sum())

Count of CEOs by year:
year
2004      1
2005      3
2006     1435
2007     2001
2008     2063
2009     2026
2010     2061
2011     2115
2012     2187
2013     2330
2014     2382
2015     2446
2016     2531
2017     2674
2018     2797
2019     2949
2020     3040
2021     2813
2022      182
Name: CEO, dtype: int64

In [30]: ceo_df['salary'].quantile(0.99)

Out[30]:
2000000.0

In [31]: ticker_series = starter_df.groupby('year')['ticker'].nunique()
ceo_series = starter_df[starter_df['CEO'] == True] & (starter_df['Interim'] == False).groupby('year')['CEO']
pd.DataFrame([ticker_series, ceo_series]).T

Out[31]:
      ticker  CEO
year
2004      1    1
2005     49   53
2006    196  1435
2007    186  2001
2008    196  2063
2009    196  2026
2010    196  2061
2011    196  2115
2012    196  2187
2013    196  2330
2014    196  2382
2015    196  2446
2016    196  2531
2017    196  2674
2018    196  2797
2019    196  2949
2020    196  3040
2021    196  2813
2022     16   182

In [32]: pd.set_option('display.float_format', lambda x: '%.0f' % x)
print(states_df[starter_df['CEO'] == True] & (starter_df['Interim'] == False).groupby('year')['CEO'].sum())

Out[32]:
      ticker  salary      bonus  stockAwards  optionAwards  nonEquityIncentiveCompensation  otherCompensation  to
count  172774  172774      172774      172774      172774      172774      172774      1727
mean    907963  481363      931896      1493095      1425116      966832      441096      32924
std     509158  367194      149921671      55342801      184662312      83971598      91501722      170409
min       1750         0         0         0         0         0         0         0
50%     723125      288000         0         0         0         0         8560      7856
75%     940942      412948         0      337822         0      180000      26206      16569
max     1915657      6000000      52320562950      17980451374      4999991100      18771084337      37381610730      56821463

In [ ]:

In [33]: plt.figure(figsize=(9,3))
sns.barplot(data=pd.melt(ceo_df[ceo_df['ticker'] == 'MSFT'] [['year','salary','nonEquityIncentiveCompensation','stockAwards'],
                    id_vars='year', var_name='component', value_name='amount'),
            x='year',y='amount', hue='component',ci=80000)
plt.ylim(0,600000000)
plt.legend(loc='upper left')
plt.title('Pay Components for Microsoft CEOs 2006-2021')
plt.ylabel('Amount in $100k')

Out[33]:
Text(0, 0.5, 'Amount in $100k')

Pay Components for Microsoft CEOs 2006-2021

Amount in $100k
5
4
3
2
1
0
2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022
year
salary
nonEquityIncentiveCompensation
stockAwards
total

In [34]: pd.melt(ceo_df[ceo_df['ticker'] == 'P'] [['year','salary','nonEquityIncentiveCompensation','stockAwards','total'])
id_vars='year', var_name='component', value_name='amount')

Out[34]:
      year  component  amount
0  2021      salary    1700000
1  2020      salary    1425000
2  2020      salary    1800000
3  2019      salary    1800000
4  2019      salary    1800000
...  ...  ...  ...
91 2008      total    13565378
92 2007      total    21670674
93 2007      total     8677747
94 2006      total    28183476
95 2006      total    2747100

96 rows x 3 columns

In [ ]:
```



```
In [85]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.preprocessing import MinMaxScaler, StandardScaler
from sklearn.model_selection import train_test_split

from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn.ensemble import RandomForestClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import GaussianNB

from sklearn.metrics import confusion_matrix, classification_report
```

```
In [87]: new_main = pd.read_csv('./datasets/new_main_df.csv')

new_main = new_main[['CEO', 'salary', 'stockAwards', 'optionAwards', 'nonEquityIncentiveCompensation',
'Revenue', 'Net Income Common', 'all_Equity_Awards']]

new_main['Revenue'] = np.log(new_main['Revenue'])

new_main.head()
```

Given that we have a small data set here, I am going to go off the guidance of the book in partitioning. "On the other hand, for smaller or less complex data sets, one should retain sufficient records for accurate assessment, so that the training sets would contain only 50-67% of the original data." I think we should sufficient records to train on though, so I posit that we do a 60% training, 40% testing split (if you think we should do otherwise, by all means I am all ears).

```
In [88]: #Partitioning the data
fin_train, fin_test = train_test_split(new_main,
                                       test_size = 0.3, random_state = 7)

# Testing to see if we have partitioned correctly.
print("Shape of entire dataset: ", new_main.shape, "\n")
print("Shape of training set: ", fin_train.shape, "\n")
print("Shape of testing set: ", fin_test.shape, "\n")
```

Given that we have a small data set here, I am going to go off the guidance of the book in partitioning. "On the other hand, for smaller or less complex data sets, one should retain sufficient records for accurate assessment, so that the training sets would contain only 50-67% of the original data." I think we should sufficient records to train on though, so I posit that we do a 60% training, 40% testing split (if you think we should do otherwise, by all means I am all ears).

```
In [88]: #Partitioning the data
fin_train, fin_test = train_test_split(new_main,
                                      test_size = 0.3, random_state = 7)

# Testing to see if we have partitioned correctly.
print("Shape of entire dataset: ",new_main.shape,"\n")
print("Value counts of CEO vs CFO:\n",new_main['CEO'].value_counts(),"\n")
print("Shape of training dataset: ",fin_train.shape,"\n")
print("Shape of testing dataset: ",fin_test.shape)

Shape of entire dataset: (959, 8)

Value counts of CEO vs CFO:
True      497
False     462
Name: CEO, dtype: int64

Shape of training dataset: (671, 8)

Shape of testing dataset: (288, 8)
```

Baseline model: We call the "positive" class "True" for CEO and see that a baseline model would predict correctly with (497/959) = 52% accuracy. If our model beats 52% accuracy, then it beats the baseline and is considered useful. I am not going to choose the X variables as proportions of revenue, although I can see there being multicollinearity concerns with this approach if we are incorporating the Y variable into the X variable in this way. I am curious to hear your thoughts on this. Note that I am not taking the "total comp" as this would definitely be a multicollinearity concern if we are using the other 5 variables.

## Setting up train/test datasets and scaling

```
In [89]: #setting up X and y dataframes
y_train = fin_train[['CEO']]
X_train = fin_train.drop(columns='CEO')

# Setting up the variables from the test dataset.
y_test = fin_test['CEO']
X_test = fin_test.drop(columns='CEO')

#####
#scale the data
sc = StandardScaler()
X_train_sc = sc.fit_transform(X_train)
X_test_sc = sc.transform(X_test)
#####
```

## CART Decision Tree Model

```
In [90]: #instantiate and fit DT classifier
cart01 = DecisionTreeClassifier(criterion = "gini", max_leaf_nodes=5).fit(X_train_sc,y_train)

#make predictions with training dataset
y_predCart_train = cart01.predict(X_train_sc)

print("Confusion matrix for training data:")
print(confusion_matrix(y_true=y_train, y_pred=y_predCart_train))
print(classification_report(y_true=y_train, y_pred=y_predCart_train))

#make predictions with test dataset
y_predCart_test = cart01.predict(X_test_sc)

print("\nConfusion matrix for testing data:")
print(confusion_matrix(y_true=y_test, y_pred=y_predCart_test))
print(classification_report(y_true=y_test, y_pred=y_predCart_test))

Confusion matrix for training data:
[[292  31]
 [ 57 291]]

      precision    recall  f1-score   support

   False      0.84      0.90      0.87      323
    True      0.90      0.84      0.87      348

 accuracy      0.87      0.87      0.87      671
 macro avg      0.87      0.87      0.87      671
weighted avg      0.87      0.87      0.87      671

Confusion matrix for testing data:
[[121  18]
 [ 33 116]]

      precision    recall  f1-score   support

   False      0.79      0.87      0.83      139
    True      0.87      0.78      0.82      149

 accuracy      0.83      0.82      0.82      288
 macro avg      0.83      0.82      0.82      288
weighted avg      0.83      0.82      0.82      288
```

## Setting up a function which creates a confusion matrix based on model and train/test datasets

```
In [91]: def predictor_function(model, X_train, X_test, y_train, y_test):

#make predictions with training dataset
y_pred_train = model.predict(X_train)

print("Confusion matrix for training data:")
print(confusion_matrix(y_true=y_train, y_pred=y_pred_train))
print(classification_report(y_true=y_train, y_pred=y_pred_train))

#make predictions with test dataset
y_pred_test = model.predict(X_test)

print("\nConfusion matrix for testing data:")
print(confusion_matrix(y_true=y_test, y_pred=y_pred_test))
print(classification_report(y_true=y_test, y_pred=y_pred_test))
```

```
In [92]: predictor_function(cart01,X_train_sc, X_test_sc, y_train, y_test)

Confusion matrix for training data:
[[292  31]
 [ 57 291]]

      precision    recall  f1-score   support

   False      0.84      0.90      0.87      323
    True      0.90      0.84      0.87      348

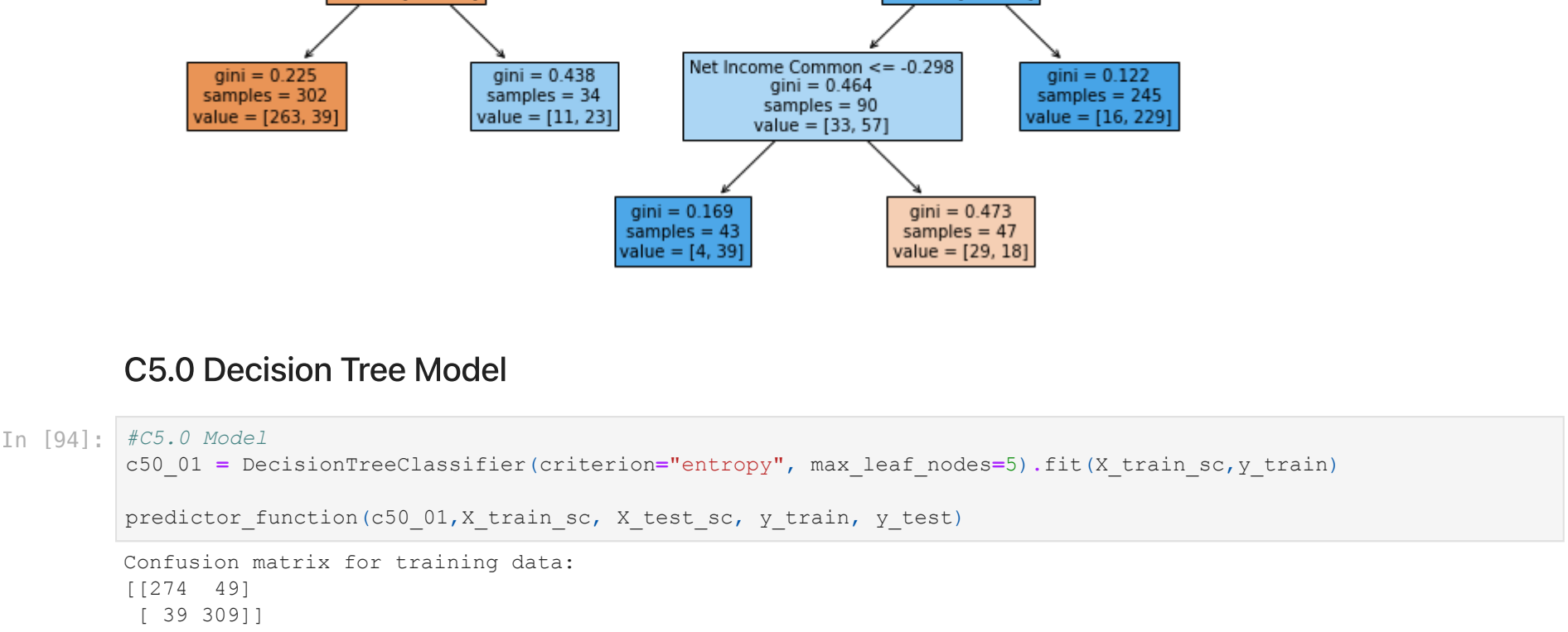
 accuracy      0.87      0.87      0.87      671
 macro avg      0.87      0.87      0.87      671
weighted avg      0.87      0.87      0.87      671

Confusion matrix for testing data:
[[121  18]
 [ 33 116]]

      precision    recall  f1-score   support

   False      0.79      0.87      0.83      139
    True      0.87      0.78      0.82      149

 accuracy      0.83      0.82      0.82      288
 macro avg      0.83      0.82      0.82      288
weighted avg      0.83      0.82      0.82      288
```



## C5.0 Decision Tree Model

```
In [94]: #C5.0 Model
c50_01 = DecisionTreeClassifier(criterion="entropy", max_leaf_nodes=5).fit(X_train_sc,y_train)

predictor_function(c50_01,X_train_sc, X_test_sc, y_train, y_test)

Confusion matrix for training data:
[[274  49]
 [ 39 309]]

      precision    recall  f1-score   support

   False      0.89      0.85      0.86      323
    True      0.86      0.89      0.88      348

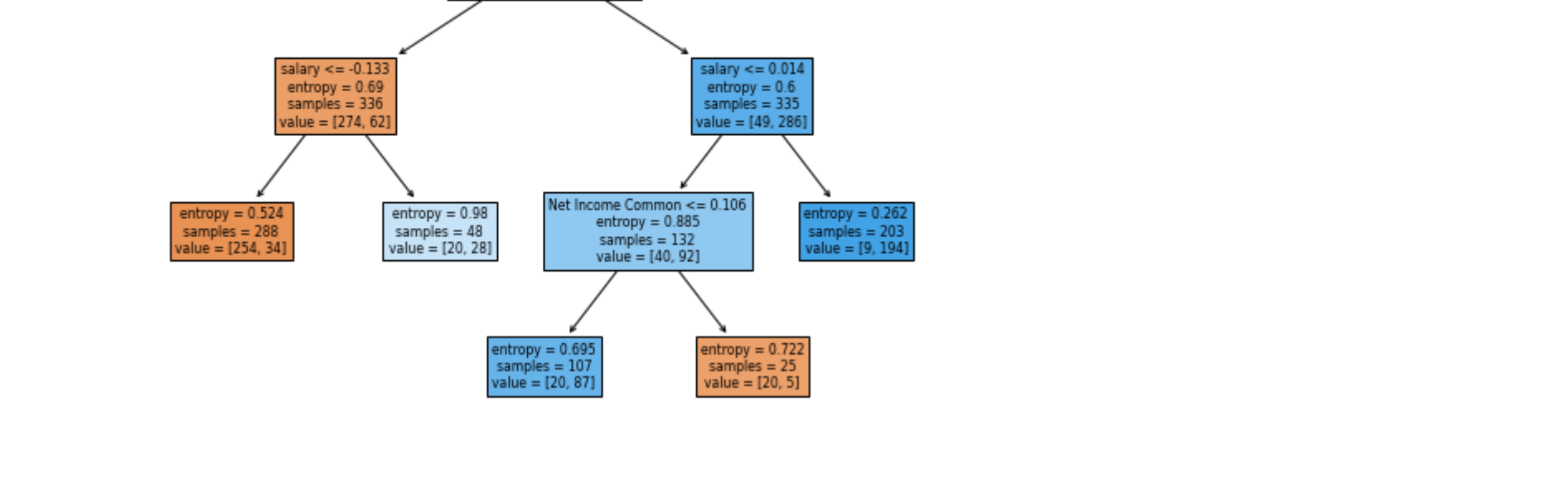
 accuracy      0.87      0.87      0.87      671
 macro avg      0.87      0.87      0.87      671
weighted avg      0.87      0.87      0.87      671

Confusion matrix for testing data:
[[113  26]
 [ 31 118]]

      precision    recall  f1-score   support

   False      0.78      0.81      0.80      139
    True      0.82      0.79      0.81      149

 accuracy      0.80      0.80      0.80      288
 macro avg      0.80      0.80      0.80      288
weighted avg      0.80      0.80      0.80      288
```



## Random Forest Model

```
In [96]: #Random Forest Model
rfy = np.ravel(y_train)
rf01 = RandomForestClassifier(n_estimators = 100, criterion="gini",max_depth=5)
rf01.fit(X_train_sc,rfy)

predictor_function(rf01,X_train_sc, X_test_sc, y_train, y_test)

Confusion matrix for training data:
[[314  9]
 [ 28 320]]

      precision    recall  f1-score   support

   False      0.92      0.97      0.94      323
    True      0.97      0.92      0.95      348

 accuracy      0.95      0.95      0.94      671
 macro avg      0.95      0.94      0.94      671
weighted avg      0.95      0.94      0.94      671

Confusion matrix for testing data:
[[119  20]
 [ 25 124]]

      precision    recall  f1-score   support

   False      0.83      0.86      0.84      139
    True      0.86      0.83      0.85      149

 accuracy      0.84      0.84      0.84      288
 macro avg      0.84      0.84      0.84      288
weighted avg      0.84      0.84      0.84      288
```

```
In [97]: rf01.decision_path(X_train_sc)

Out[97]: <671x4852 sparse matrix of type '<class 'numpy.int64'>'
with 389344 stored elements in Compressed Sparse Row format>,
[[ 0,  51,  98, 153, 202, 251, 302, 349, 398, 447, 498,
  537, 588, 625, 672, 715, 764, 813, 862, 909, 956, 1009,
  1062, 1109, 1162, 1213, 1254, 1293, 1336, 1383, 1420, 1465, 1518,
  1565, 1618, 1667, 1720, 1773, 1824, 1867, 1918, 1977, 2028, 2071,
  2120, 2173, 2218, 2259, 2310, 2365, 2424, 2473, 2526, 2573, 2610,
  2653, 2704, 2753, 2800, 2847, 2898, 2945, 2994, 3037, 3088, 3137,
  3186, 3237, 3298, 3349, 3400, 3445, 3492, 3541, 3592, 3641, 3682,
  3737, 3786, 3839, 3888, 3939, 3988, 4041, 4096, 4139, 4186, 4239,
  4282, 4329, 4372, 4425, 4478, 4529, 4570, 4619, 4664, 4717, 4762,
  4809, 4852]])
```

```
In [98]: rf01.get_params()

Out[98]: {'bootstrap': True,
'ccp_alpha': 0.0,
'class_weight': None,
'criterion': 'gini',
'max_depth': 5,
'max_features': 'auto',
'max_leaf_nodes': None,
'max_samples': None,
'min_impurity_decrease': 0.0,
'min_samples_leaf': 1,
'min_samples_split': 2,
'min_weight_fraction_leaf': 0.0,
'n_estimators': 100,
'n_jobs': None,
'out_of_sample': False,
'random_state': None,
'verbose': 0,
'warm_start': False}
```

## Logistic Regression Model (default settings)

```
In [99]: lry = np.ravel(y_train)
logreg01 = LogisticRegression()
logreg01.fit(X_train_sc,lry)

predictor_function(logreg01,X_train_sc, X_test_sc, y_train, y_test)

Confusion matrix for training data:
[[300  23]
 [ 46 302]]

      precision    recall  f1-score   support

   False      0.87      0.93      0.90      323
    True      0.93      0.87      0.90      348

 accuracy      0.90      0.90      0.90      671
 macro avg      0.90      0.90      0.90      671
weighted avg      0.90      0.90      0.90      671

Confusion matrix for testing data:
[[127  12]
 [ 20 129]]

      precision    recall  f1-score   support

   False      0.86      0.91      0.89      139
    True      0.91      0.87      0.89      149

 accuracy      0.89      0.89      0.89      288
 macro avg      0.89      0.89      0.89      288
weighted avg      0.89      0.89      0.89      288
```

## Naive Bayes Model (default settings)

[https://scikit-learn.org/stable/modules/naive\\_bayes.html](https://scikit-learn.org/stable/modules/naive_bayes.html)

```
In [100... gnb = np.ravel(y_train)
gnb = GaussianNB()
gnb.fit(X_train_sc,gnb)

predictor_function(cart01,X_train_sc, X_test_sc, y_train, y_test)

Confusion matrix for training data:
[[292  31]
 [ 57 291]]

      precision    recall  f1-score   support

   False      0.90      0.84      0.87      323
    True      0.90      0.84      0.87      348

 accuracy      0.87      0.87      0.87      671
 macro avg      0.87      0.87      0.87      671
weighted avg      0.87      0.87      0.87      671

Confusion matrix for testing data:
[[121  18]
 [ 33 116]]

      precision    recall  f1-score   support

   False      0.79      0.87      0.83      139
    True      0.87      0.78      0.82      149

 accuracy      0.83      0.82      0.82      288
 macro avg      0.83      0.82      0.82      288
weighted avg      0.83      0.82      0.82      288
```

## KNeighbors Model

Notice how changing the weights from "uniform" to "distance" results in a very overfit model.

```
In [101... from sklearn.neighbors import KNeighborsClassifier

In [102... knn = KNeighborsClassifier(n_neighbors=5, weights='uniform')

knn = np.ravel(y_train)
knn.fit(X_train_sc,knn)

predictor_function(knn,X_train_sc, X_test_sc, y_train, y_test)

Confusion matrix for training data:
[[306  17]
 [ 44 304]]

      precision    recall  f1-score   support

   False      0.87      0.95      0.91      323
    True      0.95      0.87      0.91      348

 accuracy      0.91      0.91      0.91      671
 macro avg      0.91      0.91      0.91      671
weighted avg      0.91      0.91      0.91      671

Confusion matrix for testing data:
[[125  14]
 [ 24 125]]

      precision    recall  f1-score   support

   False      0.84      0.90      0.87      139
    True      0.90      0.84      0.87      149

 accuracy      0.87      0.87      0.87      288
 macro avg      0.87      0.87      0.87      288
weighted avg      0.87      0.87      0.87      288
```

```
In [103... knn = KNeighborsClassifier(n_neighbors=5, weights='distance')

knn = np.ravel(y_train)
knn.fit(X_train_sc,knn)

predictor_function(knn,X_train_sc, X_test_sc, y_train, y_test)

Confusion matrix for training data:
[[323  0]
 [ 0 348]]

      precision    recall  f1-score   support

   False      1.00      1.00      1.00      323
    True      1.00      1.00      1.00      348

 accuracy      1.00      1.00      1.00      671
 macro avg      1.00      1.00      1.00      671
weighted avg      1.00      1.00      1.00      671

Confusion matrix for testing data:
[[125  14]
 [ 24 125]]

      precision    recall  f1-score   support

   False      0.84      0.90      0.87      139
    True      0.90      0.84      0.87      149

 accuracy      0.87      0.87      0.87      288
 macro avg      0.87      0.87      0.87      288
weighted avg      0.87      0.87      0.87      288
```

```
In [104... sns.scatterplot(data=X_train, x='Revenue',y='salary', hue=knn.predict(X_train_sc))
plt.title('Predicting CEO')


```

## Suppor Vector Machine

```
In [105... from sklearn.svm import SVC

In [106... svcy = np.ravel(y_train)
svc = SVC()
svc.fit(X_train_sc,svcy)

predictor_function(svc,X_train_sc, X_test_sc, y_train, y_test)

Confusion matrix for training data:
[[303  20]
 [ 44 304]]

      precision    recall  f1-score   support

   False      0.87      0.94      0.90      323
    True      0.94      0.87      0.90      348

 accuracy      0.91      0.91      0.90      671
 weighted avg      0.91      0.90      0.90      671

Confusion matrix for testing data:
[[128  11]
 [ 25 124]]

      precision    recall  f1-score   support

   False      0.84      0.92      0.88      139
    True      0.92      0.83      0.87      149

 accuracy      0.88      0.88      0.88      288
 macro avg      0.88      0.88      0.87      288
weighted avg      0.88      0.88      0.87      288
```

```
In [107... svc.get_params()

Out[107]: {'C': 1.0,
'break_ties': False,
'cache_size': 200,
'class_weight': None,
'coef0': 0.0,
'decision_function_shape': 'ovr',
'degree': 3,
'gamma': 'scale',
'kernel': 'rbf',
'kernel_coef': -1,
'probability': False,
'random_state': None,
'shrinking': True,
'tol': 0.001,
'verbose': False}
```

```
In [108... svc.support_vectors_

Out[108]: array([[ 0.50443643,  0.59707948, -0.19249866, ...,  2.26253538,
  3.96678542,  0.13051731],
[-0.53998382, -0.60718945, -0.15549085, ..., -0.69358501,
-0.48241272, -0.43420506],
[-0.32642185, -0.3594091, -0.13628877, ..., -0.22029892,
-0.35876796, -0.29493724],
...,
[ 3.7341026 ,  0.82377441,  0.70739546, ...,  1.22142338,
  0.68566547,  1.01518815],
[ 4.07664129,  2.56042895, -0.19249866, ...,  1.24806762,
  1.74942383,  1.10289119],
[ 0.06781806, -0.72883012, -0.19249866, ..., -0.49593355,
-0.22353882, -0.52622117]])
```