Web Scraping and Linear Regression Project

Identify relationship between stock prices and financial statistics

Abstract

The goal of this project is to find out how much a public company's stock price can be explained by its financial numbers. The project uses a few financials as features and stock closed price on Oct 4th as dependent variable. The algorithm used includes Ordinary least square regression, regulations like Lasso, Ridge, and ElasticNet. I also use a dummy variable to improve the model. The research will help potential investors, traders, and whoever insisted to make better informed decisions based on the contents the project offers.

Design

The data used in the project are web-scraped from website https://stockanalysis.com/stocks/ Financial information provided in "Summary" and "Statistics" section will be collected. After data cleaning and wrangling, I did OLS linear regression model, checked assumptions and evaluation. Lasso, Ridge and ElasticNet regression model with cross validation are also performed to seek a better model with less variances. The project later uses a categorical variable "Industry" as dummy variable to make improvement.

Data

I scrape ticker list info from https://stockanalysis.com/stocks/ and for each ticker go to https://stockanalysis.com/s

Tools

- request, selenium for data web scraping;
- BeautifulSoup for HTML syntax parsing;
- pandas and numpy for data manipulation;
- Matplotlit and Seaborn for data visualization;
- Pickle for data object serialization and de-serialization
- sklearn and statsmodels for regression

Algorithms

- 1. Clean data: check/transfer data types, deal with missing value.
- 2. Create a baseline OLS linear regression model on all features.
- 3. Check assumptions and evaluation: get rid of outliers, check VIF and move away features of collinearity, transform dependent feature with log function to correct non_normality.
- 4. Compare R squared values for regressions with / without historical price feature
- 5. Add Polynomial features and try different regulation methods with cross validation: Lasso, Ridge and ElasticNet, and choose the model with best mean R squared value.
- 6. Add available dummy variable "Industry" to improve model.

Result

1. Linear regression with historical price related feature: "52 Week High 3"

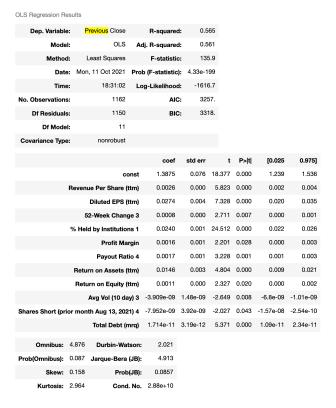
OLS Regression Results				
Dep. Variable:	Previous Close	R-squared:	0.956	
Model:	OLS	Adj. R-squared:	0.955	
Method:	Least Squares	F-statistic:	2073.	
Date:	Mon, 11 Oct 2021	Prob (F-statistic):	0.00	
Time:	18:30:58	Log-Likelihood:	-287.62	
No. Observations:	1162	AIC:	601.2	
Df Residuals:	1149	BIC:	667.0	
Df Model:	12			
Covariance Type:	nonrobust			

Insights:

The R-squared and Adj. R-squared are both terrifically high, it is caused by feature "52 Week High 3". "52 Week High 3" is defined the highest price at which a security, such as a stock, has traded during previous 52 weeks. It is an important factor in the analysis of a stock's current value, so it is very reasonable to have a high R squared value with "52 Week High 3" as one of our feature.

2. Linear regression without historical price related feature.

Insight: Generally speaking, the financial numbers of a public company explains 56.5% of its stock price.



3. Regression with Regulations:

```
print("Linear Reg Mean Score: %3.5f"%(np.mean(lm_scores)))
print( "Lasso CV best mean score:%3.5f"%(grid_Lasso.best_score_), ", param is",grid_Lasso.best_params_)
print( "Ridge CV best mean score:%3.5f"%(grid_ridge.best_score_), ", param is",grid_ridge.best_params_)
print( "Elastic CV best mean score:%3.5f"%(grid_elastic.best_score_), ", param is",grid_elastic.best_params_)

Linear Reg Mean Score: 0.52626
Lasso CV best mean score:0.53012 , param is {'Lasso_alpha': 0.005179474679231213}
Ridge CV best mean score:0.53084 , param is {'Ridge_alpha': 94.75205302806543}
Elastic CV best mean score:0.53121 , param is {'Elastic_alpha': 0.06723357536499334, 'Elastic_l1_ratio': 0.02564102
564102564}
```

Choose model with best mean score: ElasticNet model:

```
intercept, coef df
(3.2421367076548813,
                                       Feature coefficent
                       Revenue Per Share (ttm)
                                                 0.189678
                             Diluted EPS (ttm)
                              52-Week Change 3
                                                 0.066463
                      % Held by Institutions 1
                                                  0.71223
                                 Profit Margin
                                                 0.085392
                                Pavout Ratio 4
                                                 0.091668
                                                 0.194229
                        Return on Assets (ttm)
                        Return on Equity (ttm)
                            Avg Vol (10 day) 3
                                                 -0.076562
     Shares Short (prior month Aug 13, 2021) 4
                                                -0.049132
                              Total Debt (mrq)
                                                 0.162104)
```

Insight:

1. 'Revenue Per Share (ttm)', 'Diluted EPS (ttm)', '52-Week Change 3', '% Held by Institutions 1', 'Profit Margin', 'Payout Ratio 4', 'Return on Assets (ttm)', 'Return on Equity (ttm)''Total Debt (mrq)' have a

Positive relationship with stock price. For example: For every dollar increase in Diluted EPS, the stock price will increase $e^{(0.229)} = 1.26$ times

2. 'Avg Vol (10 day) 3', 'Shares Short (prior month Aug 13, 2021) 4' have negative relationship with stock price. For example, for every share Short by investors, the stock price decrease by (1-np.e**-0.049132) = 4%

Definitions:

Revenue per share: Amount of revenue over common shares outstanding.

Diluted EPS: Earnings per share (EPS) if all convertible securities were exercised.

Profit margin: The difference between sales and the cost of goods sold divided by revenue.

Payout Ratio: The percentage of net income that a company pays out as dividends to common shareholders.

Return on Asset: The rate of return (after tax) being earned on all of the firm's assets regardless of financing structure.

Return on Equity: Rate of return on the money invested by common stock owners and retained by the company thanks to previous profitable years.

Avg vol(10days): The average number of shares traded within a day in a given stock.

Shares Short: Number of shares that investors don't own but selling to other investors.

4. What if feed all non-price-related features to ElasticNet model

```
print(grid_elastic_check.best_params_,", best mean score:",grid_elastic_check.best_score_)
print("train score = %3.5f" %(grid_elastic_check.score(X_2,logy)),",test score = %3.5f" %(
```

	Feature	coefficent
0	MARKET_CAP	-0.000644
1	EPS (TTM)	0.025295
2	Short Ratio (Sep 15, 2021) 4	-0.0
3	Revenue Per Share (ttm)	0.147927
4	Diluted EPS (ttm)	0.02288
5	Total Cash Per Share (mrq)	0.12549
6	Current Ratio (mrq)	0.017106
7	Book Value Per Share (mrq)	0.1757
8	52-Week Change 3	0.070312
9	% Held by Insiders 1	0.158532
10	% Held by Institutions 1	0.831912
11	Profit Margin	0.075249
12	Operating Margin (ttm)	0.0
13	Short % of Shares Outstanding (Sep 15, 2021) 4	0.073475
14	Short % of Float (Sep 15, 2021) 4	-0.079532
15	Payout Ratio 4	0.100258

16	Return on Assets (ttm)	0.173257
17	Return on Equity (ttm)	0.113667
18	Avg Vol (3 month) 3	0.125056
19	Avg Vol (10 day) 3	-0.129229
20	Shares Outstanding 5	-0.105899
21	Shares Short (Sep 15, 2021) 4	-0.09158
22	Shares Short (prior month Aug 13, 2021) 4	-0.0
23	Revenue (ttm)	-0.141354
24	Gross Profit (ttm)	0.244014
25	EBITDA	0.018669
26	Net Income Avi to Common (ttm)	-0.0
27	Total Cash (mrq)	0.032115
28	Total Debt (mrq)	0.073748
29	Operating Cash Flow (ttm)	0.060896

Insight: The mean score increased slightly, but our model will become much more complex.

5. linear regression with dummy variables "Industry"

The R squared value will increase to 0.601.

Dep. Variable:	Previous Close	R-squared:	0.601
click to scroll output; do Model:	uble click to hide OLS	Adj. R-squared:	0.592
Method:	Least Squares	F-statistic:	63.30
Date:	Mon, 11 Oct 2021	Prob (F-statistic):	1.43e-204
Time:	18:27:32	Log-Likelihood:	-1566.5
No. Observations:	1162	AIC:	3189.
Df Residuals:	1134	BIC:	3331.
Df Model:	27		
Covariance Type:	nonrobust		

Industry_Chemicals	-0.0130	0.221	-0.059	0.953	-0.446	0.420
Industry_Electronic Equipment, Instruments & C	0.0238	0.197	0.121	0.904	-0.363	0.411
Industry_Energy Equipment & Services	-0.7968	0.226	-3.532	0.000	-1.239	-0.354
Industry_Equity Real Estate Investment Trusts	-0.0743	0.194	-0.383	0.702	-0.455	0.306
Industry_Food Products	0.0814	0.227	0.359	0.720	-0.363	0.526
Industry_Health Care Equipment & Supplies	0.3378	0.179	1.892	0.059	-0.013	0.688
Industry_Health Care Providers & Services	-0.1748	0.209	-0.835	0.404	-0.586	0.236
Industry_Hotels, Restaurants & Leisure	-0.0616	0.214	-0.287	0.774	-0.482	0.359
Industry_IT Services	0.4952	0.194	2.549	0.011	0.114	0.876
Industry_Insurance	0.1214	0.213	0.569	0.570	-0.297	0.540
Industry_Machinery	0.1427	0.208	0.684	0.494	-0.266	0.552
Industry_Oil, Gas & Consumable Fuels	-0.2092	0.201	-1.041	0.298	-0.603	0.185
Industry_Pharmaceuticals	-0.2141	0.201	-1.066	0.286	-0.608	0.180
Industry_Semiconductors & Semiconductor Equipment	0.7565	0.221	3.416	0.001	0.322	1.191
Industry_Software	0.7854	0.169	4.643	0.000	0.454	1.117
Industry_Specialty Retail	-0.3562	0.210	-1.694	0.091	-0.769	0.056
Industry_other	-0.0620	0.132	-0.470	0.638	-0.321	0.197
intercept	1.3713	0.144	9.522	0.000	1.089	1.654