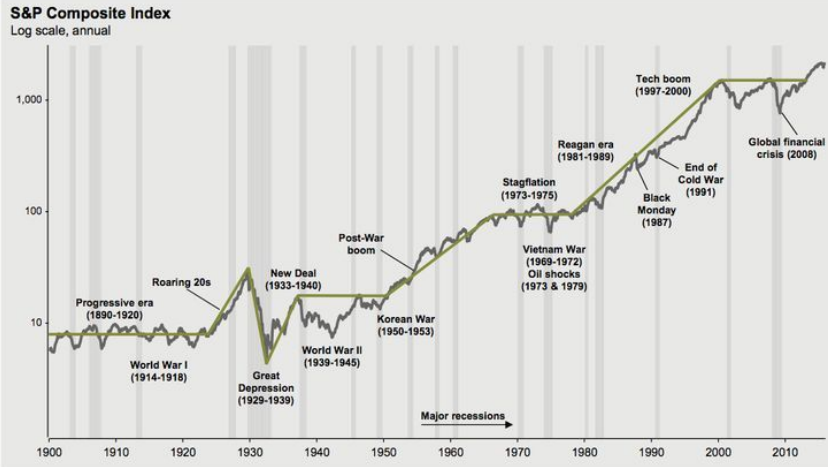# Clustering Fundamental Indicators of S&p 500 Listed Companies

Andrew Koller
Regis MSDS
03.06.2024

# Introduction



Stock market since 1900 — GTM – U.S. | 15

S&P Composite Index
Log scale, annual

Tech boom (1997-2000)

Reagan era (1981-1989)

Global financial crisis (2008)

Stagflation (1973-1975)

End of Cold War (1991)

Black Monday (1987)

Post-War boom

Vietnam War (1969-1972) Oil shocks (1973 & 1979)

Roaring 20s

New Deal (1933-1940)

Progressive era (1890-1920)

Korean War (1950-1953)

World War I (1914-1918)

Great Depression (1929-1939)

World War II (1939-1945)

Major recessions

Source: FactSet, NBER, Robert Shiller, J.P. Morgan Asset Management.
Data shown in log scale to best illustrate long-term index patterns.

The purpose of investing in the US equity market is to make money.

The best way to do this is by 'buying low and selling high'..

Timing the market like this is exceedingly difficult.

An alternative approach is to not time the market, but simply invest in well run companies.

# Hypothesis

In the long run, well-run companies do well & you can use a companies fundamental data to identify well-run  companies. **If well run companies do better in the long run, clustering on fundamental indicators should yield clusters of companies that have better or worse returns.**
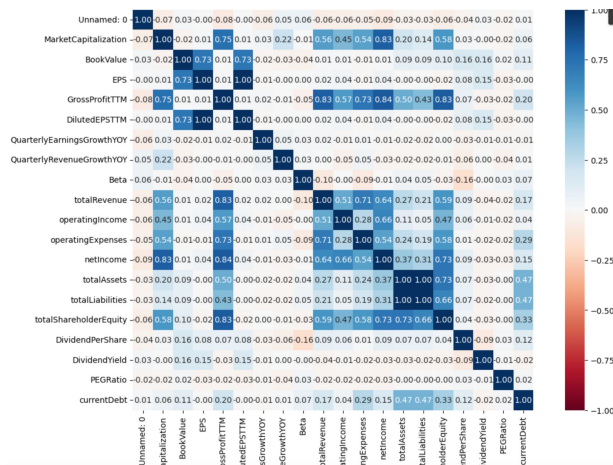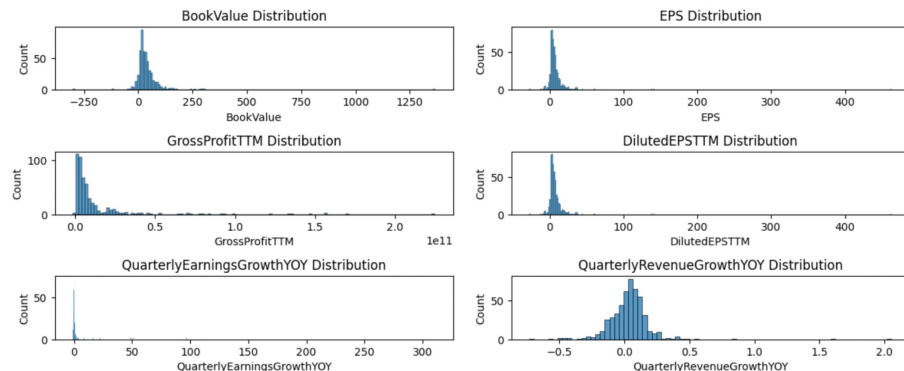
# Data

The data for this project includes income, cashflow, and balance sheet information for all S&P 500 companies for the most recent fiscal quarter. The data for this project was downloaded using a paid tier of the Alpha Vantage Python API. There was an attempt to impute any missing data with values from yfinance, and features missing from both sources were dropped from the analysis.

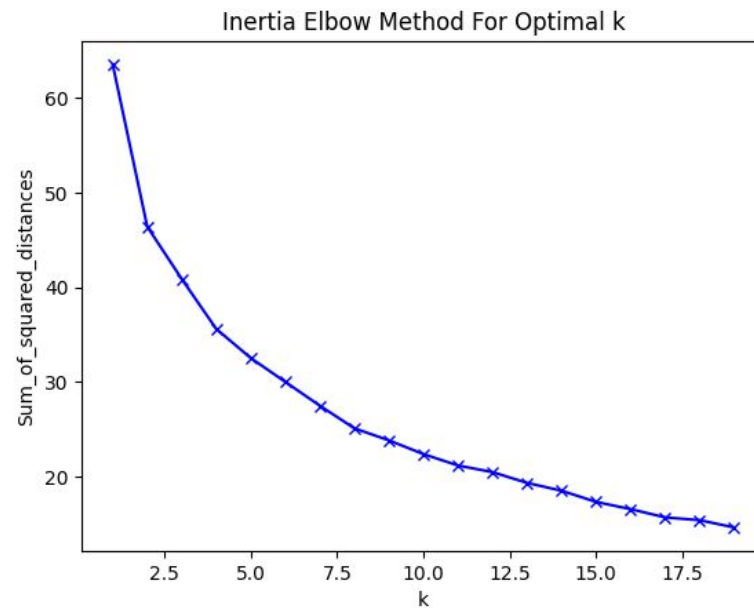| | feature | na_count |
|---|---|---|
| 7 | DividendYield | 97 |
| 6 | DividendPerShare | 97 |
| 25 | currentDebt | 95 |
| 30 | dividendPayout | 81 |
| 3 | PERatio | 30 |
| 2 | EBITDA | 28 |
| 28 | capitalExpenditures | 12 |
| 4 | PEGRatio | 5 |
| 14 | Beta | 3 |
| 29 | profitLoss | 3 |
| 21 | ebit | 2 |
| 16 | grossProfit | 1 |
| 5 | BookValue | 1 |
| 18 | costOfRevenue | 1 |
| 20 | operatingExpenses | 1 |
| 27 | operatingCashflow | 1 |
| 17 | totalRevenue | 1 |

# EDA & Feature Engineering

- The distributions of each feature were checked to ensure enough variation to draw clusters

- The correlation matrix shows that there are not many redundancies in the data

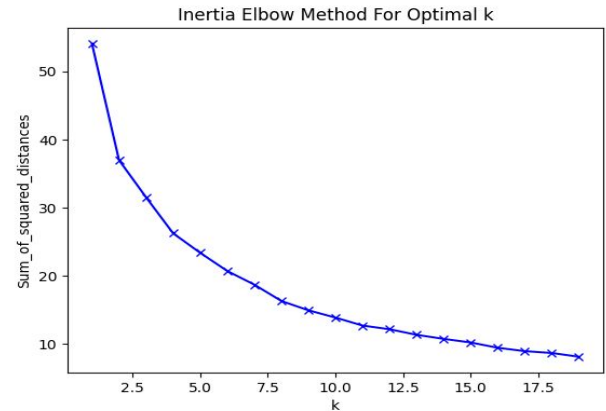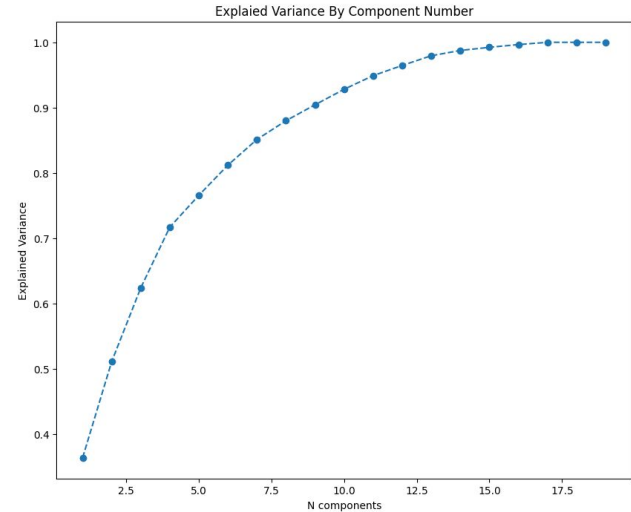- All data was scaled between 0 and 1 using a min/max scaler

# Optimal Number of Clusters

- The 'Elbow Method' was used to choose K

- Borth Distortion and Inertia were considered

- Both charts had ambiguous 'elbows' but 8 seemed like a good choice

- K = 8. Lower cluster numbers make little sense.



Inertia Elbow Method For Optimal k

# PCA + Elbow Method

- Distance measures in higher dimensional spaces become ambiguous

- Use PCA to reduce the dimensionality

- Choose appropriate number of components based on 80% of explained variance
  - Components = 7

- Elbow method with PCA yields 5 clusters



Explaied Variance By Component Number



Inertia Elbow Method For Optimal k

# Clustering Results

| Kmeans | N |
|---|---|
| 0 | 220 |
| 1 | 33 |
| 2 | 8 |
| 3 | 61 |
| 4 | 168 |
| 5 | 4 |
| 6 | 1 |
| 7 | 5 |

| kmeans_pca | | N |
|---|---|---|
| 0 | 192 | 192 |
| 1 | 219 | 219 |
| 2 | 9 | 9 |
| 3 | 6 | 6 |
| 4 | 74 | 74 |

Both models were successfully able to segment different equities off into different clusters. The Silhouette scores for the two models were .32 and .45 respectively. This suggests that there is a bit of overlap, and the clusters might not be very well defined, however doing PCA did yield more consistent clusters.

# Cluster Returns

The returns for each cluster was calculated by finding the percent change in price for each each equity from Jan 2024 - Feb 2024. Each equity is assumed to carry an even weight, and the average percent change was calculated for each cluster.

The overall percent change for the S&P500 was ~5% over this time period

There are clusters which have better and worse returns than the overall index

Hypothesis is validated

| Kmeans Cluster % Return | | Kmeans w/ PCA Cluster % Returns | |
|---|---|---|---|
| 0 | 0.041109 | 0 | 0.065213 |
| 1 | 0.036862 | 1 | 0.042154 |
| 2 | 0.042866 | 2 | 0.040002 |
| 3 | 0.055585 | 3 | 0.055241 |
| 4 | 0.064135 | 4 | 0.038837 |
| 5 | 0.079568 | | |
| 6 | 0.089605 | | |
| 7 | 0.028114 | | |

# Discussion

Kmeans is likely not the best clustering method for this type of analysis. Hierarchical clustering might be a better option.

The clustering DID work as intended and the hypothesis was validated.

The optimal number of clusters is a crap shoot. A combination of heuristics and industry knowledge is needed.