

Empirical Findings on Spectral Properties of Return Covariance Matrices

Dayi (Darwin) Yao, Jingyuan Chen

Special thanks to: Prof. Goldberg, Prof. Shkolnik, Rahul, Rahul Pothi Vinoth, Harrison Selwitz, and Jacob Lan

University of California, Berkeley

December 6, 2024; Updated February 21, 2025

Contents

- 1 Spectral decomposition and covariance matrix estimation
- 2 Data and Methodology
- 3 How many factors?
- 4 How do Eigenvalues grow with respect to the number of securities
- 5 We need more factors
- 6 How are the entries of spiked eigenvectors distributed

Big Idea #1: Since they conform to empirically observed properties of financial data and reduce dimension, factor models are used almost universally to generate inputs to mean-variance optimization

The return generating process

$$r = \beta f + \epsilon$$

implies the expected returns:

$$E[r] = \beta E[f] + E[\epsilon]$$

and covariance matrix:

$$\Sigma = \beta F \beta^T + \Delta$$

Returns or excess returns r are the sum of factor returns f scaled by exposures β and specific returns ϵ , which are pairwise uncorrelated and uncorrelated with factor returns. Returns are observable but the factor and specific components are not. The factor and (diagonal) specific return matrices are denoted by F and Δ .

Big Ideas

The Challenge

- Sample covariance matrices are unreliable for portfolio optimization
- High dimension (many securities) and finite observations

The Solution: Factor Models

$$r = \beta f + \epsilon$$

$$\Sigma = \beta F \beta^T + \Delta$$

- Reduces dimensionality while preserving essential structure
- Captures systematic risk ($\beta F \beta^T$) and idiosyncratic risk (Δ)

Extension: RMT

- Provides us tools to estimate factor model parameters in high dimension when data is scarce.

Spectral decomposition and covariance matrix estimation

Eigenvalues and eigenvectors of noisy sample return covariance matrices for large cap equities provide the components of factor-based covariance matrices used in Markowitz optimization

We identify salient characteristics of sample return covariance matrices, and provide some answers to these questions:

- How many spiked eigenvectors (factors) do we typically see, and how does that number vary over time.
- How do spiked eigenvalues depend on the number of securities in the estimation universe?
- How are the entries of spiked eigenvectors distributed?

Data and Methodology

Data

- Our dataset consists of an approximation of the Russell 3000 constituents based on the BlackRock iShares ETF IWW (Oct 09, 2024)
- Using WRDS Center for Research in Security Prices data, we retrieved the Daily Total Return (DlyRet) and Daily Market Capitalization (DlyCap) from January of 2003 to December of 2023 for each Ticker obtained from the above ETF.
- Securities without complete history are dropped, we have an effective maximum of 2340 securities.

Methodology Overview

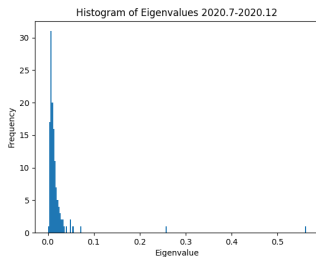
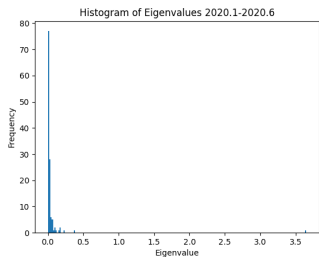
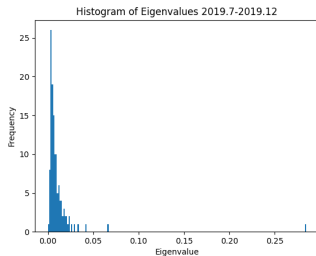
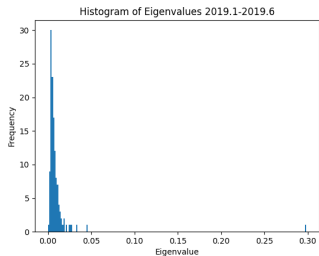
- Covariance Matrix Estimation:
 - Look-back window of 126 trading days
 - Sample covariance matrix computed for each window
- Spectral Analysis:
 - Spectral decomposition of sample covariance matrices
 - Focus on leading eigenvalues and corresponding eigenvectors
 - Comparison between market-cap sorted and "randomly" selected stocks
- Time Period Analysis:
 - Eight distinct 126-day periods (2019-2022)
 - Attention paid to market stress periods (e.g., 2020 pandemic)
- Removing Outliers
 - Remove any stock with a return over 150% or below -100%
 - See Appendix A for a detailed list

How many factors?

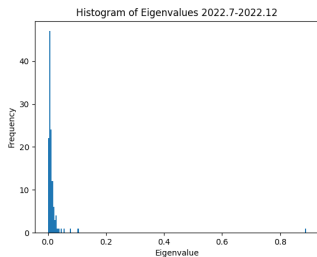
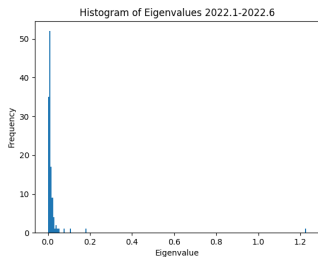
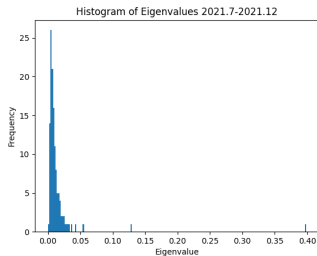
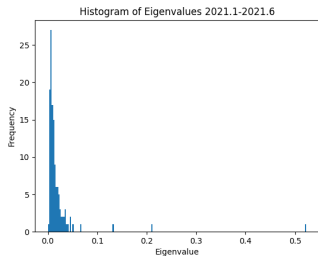
How many factors

- How many spiked eigenvectors (factors) do we typically see
- How does that number vary over time
- Eight time periods of 126 days each are chosen (2019 - 2022)
- The spectrum of eigenvalues is used to identify the number of factors

Observation over 8 time periods (Part I)



Observation over 8 time periods (Part II)



Observation over 8 time periods

Period	Factor Number
2019.1 - 2019.6	3
2019.7 - 2019.12	5
2020.1 - 2020.6	1
2020.7 - 2020.12	3

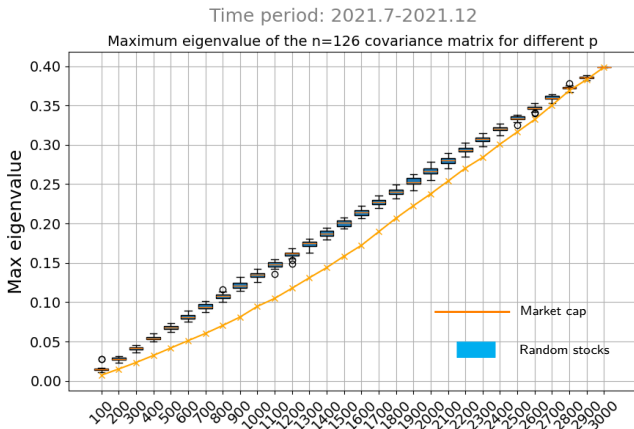
Period	Factor Number
2021.1 - 2021.6	5
2021.7 - 2021.12	2
2022.1 - 2022.6	3
2022.7 - 2022.12	2

- In normal cases, the number of outstanding factors is around 4
- In a financial crisis, the number of factor concentrates to one
- e.g. the pandemic (2020.1 - 2020.6)

Question: Is this (Elbow Method) a good way to count factors? Can we have a more systematic approach (Biometrika Paper)?

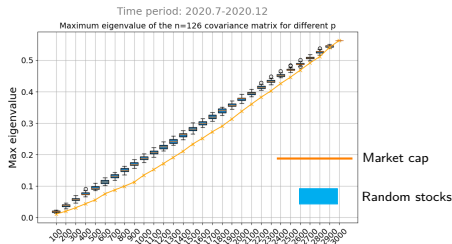
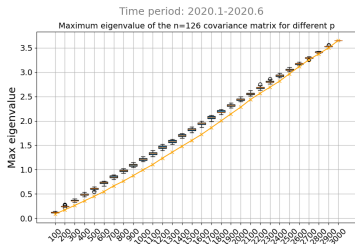
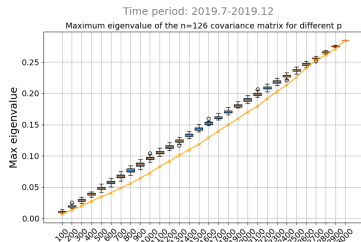
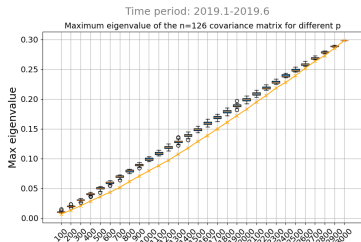
How do Eigenvalues grow with respect to the number of securities

The leading eigenvalue shows roughly affine dependence on the number of securities p

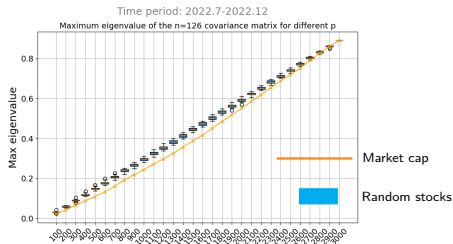
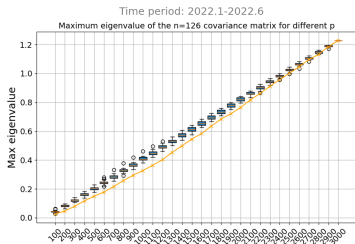
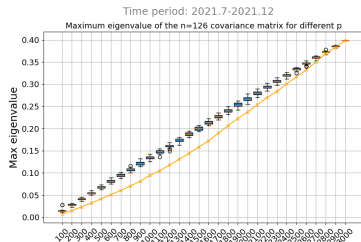
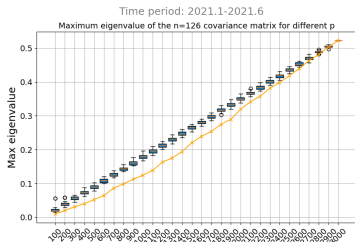


Covariance matrix estimated EOY 2021 using trailing 126 days of data. Stocks are sorted by market capitalization (orange line) or randomly drawn for each p (blue box plots).

Over 8 time periods (Part I)



Over 8 time periods (Part II)



Fit into a one-factor model

- Suppose returns follow a one-factor, homogeneous specific risk model:

$$r = \beta f + \epsilon$$

- β is a p -vector of exposures
- f is the factor return
- ϵ is a p -vector of mean 0 specific returns
- The population covariance matrix of r is given by:

$$\Sigma = \sigma^2 \beta \beta^\top + \delta^2 I$$

- σ and δ are factor and specific volatility and I is the $p \times p$ identity matrix

Fit into a one-factor model

- We draw β s from a normal distribution with mean 1 and standard deviation τ
- β is the leading eigenvector of Σ and the eigenvalue is given by:

$$\lambda^2 = \sigma^2 |\beta|^2 + \delta^2 \approx \sigma^2 p(1 + \tau^2) + \delta^2$$

- The approximation should improve as p grows
- If we fit this to the previous plot of EOY 2021 we'll have:

$$\sigma^2(1 + \tau^2) = \text{Slope} * 252 = 0.033$$

$$\delta^2 = \text{Intercept} * 252 = 0.403 \quad \delta = 0.636$$

- τ can be calculated from the first eigenvector (normalized to mean 1)

$$\tau^2 = 0.385 \quad \text{Therefore,} \quad \sigma^2 = 0.024 \quad \sigma = 0.156$$

How do σ and δ change over time

Period	σ	δ
2019.1 - 2019.6	0.134	nan
2019.7 - 2019.12	0.124	0.518
2020.1 - 2020.6	0.499	nan
2020.7 - 2020.12	0.167	0.337
2021.1 - 2021.6	0.166	1.182
2021.7 - 2021.12	0.155	0.636
2022.1 - 2022.6	0.260	nan
2022.7 - 2022.12	0.239	nan

- Fitted σ s are in a reasonable range and are close to the ones used in reality
- Fitted δ s are larger than expected (due to the hidden factors)
- σ and δ explodes in financial crisis (also true for 2008)

We need more factors

A 4-factor model

- Let r_t be the p -vector of returns at time t , where p are our securities.
- Given centered returns $Y = R - \bar{R}$, we compute sample covariance $S = YY^\top/n$
- Consider its spectral decomposition

$$S = \sum_{(s^2, h)} s^2 h h^\top = HH^\top + N$$

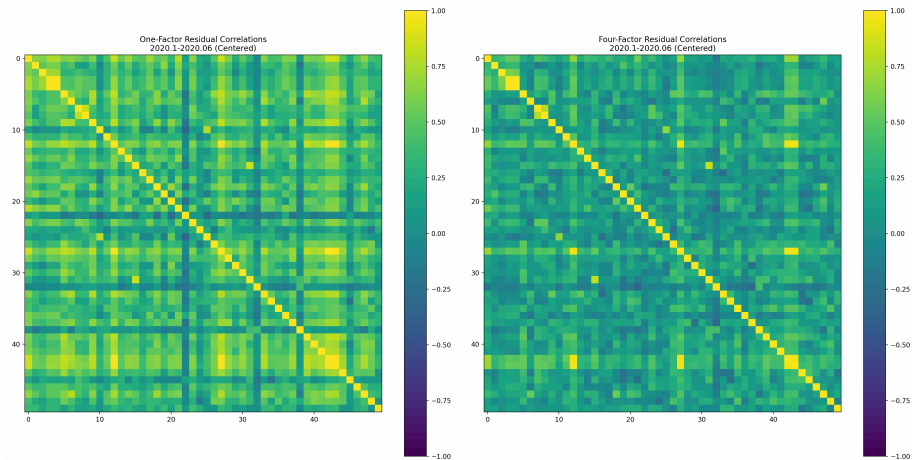
- H is a $p \times k$ matrix with columns of form sh from k largest eigenvalues s^2
- $N = S - HH^\top$ represents the residual matrix
- The PCA covariance matrix is

$$\Sigma_{PCA} = HH^\top + \Delta$$

where $\Delta = \text{diag}(N)$ sets off-diagonal elements of N to zero

- We only center the means of returns and do not standardize variances since:
 - All variables are in the same unit (daily percentage returns)
 - Variance differences between stocks contains information we desire

One vs 4-Factor Residual Correlation Heatmaps

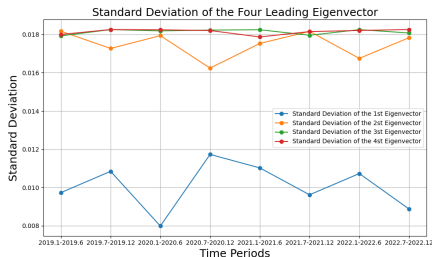
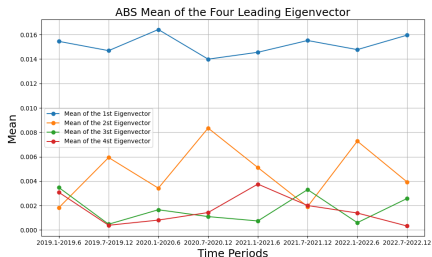


How are the entries of spiked eigenvectors distributed

How are the entries of spiked eigenvectors distributed

- Look at how the mean and variance of spiked eigenvectors change over time (before normalization)
- Normalize the first factor to mean 1
- Z-score normalize the rest factors

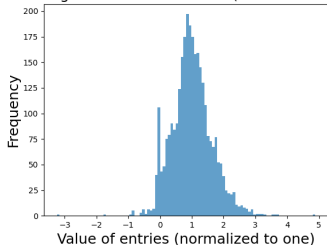
A look at how the mean and variance of spiked eigenvectors change over time (before normalization)



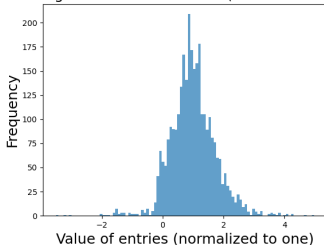
- The first factor has an outstanding mean, while the rest are rather close to 0
- The standard deviation of the first factor is lower and responds more to market changes (e.g. 2020.1 - 2020.6)

Normalize the first factor to mean 1 (Part I)

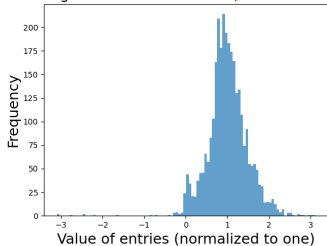
Histogram of the First Factor (2019.1-2019.6)



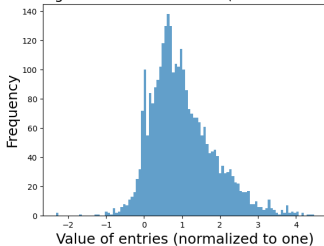
Histogram of the First Factor (2019.7-2019.12)



Histogram of the First Factor (2020.1-2020.6)

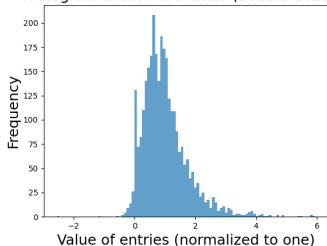


Histogram of the First Factor (2020.7-2020.12)

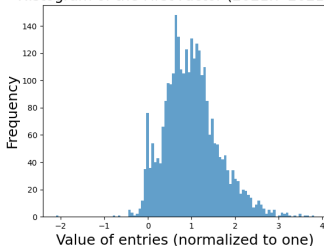


Normalize the first factor to mean 1 (Part II)

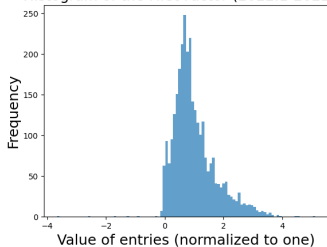
Histogram of the First Factor (2021.1-2021.6)



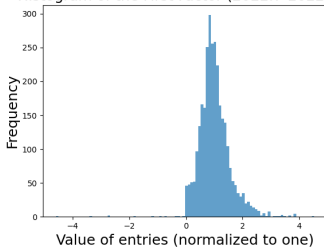
Histogram of the First Factor (2021.7-2021.12)



Histogram of the First Factor (2022.1-2022.6)

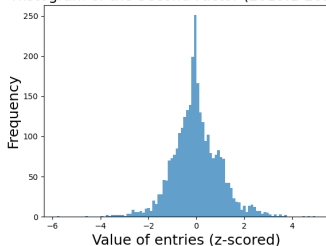


Histogram of the First Factor (2022.7-2022.12)

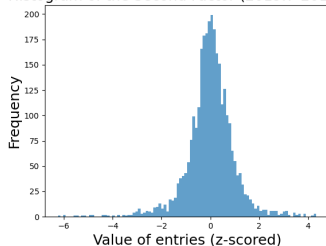


Z-score the second factor (Part I)

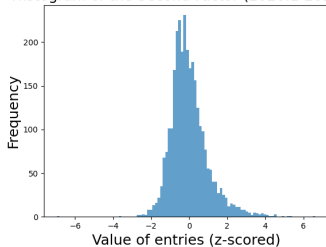
Histogram of the Second Factor (2019.1-2019.6)



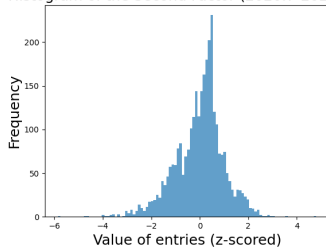
Histogram of the Second Factor (2019.7-2019.12)



Histogram of the Second Factor (2020.1-2020.6)

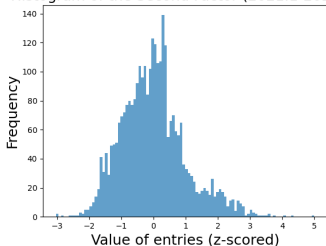


Histogram of the Second Factor (2020.7-2020.12)

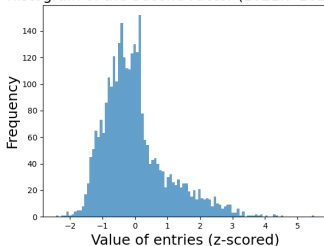


Normalize the first factor to mean 1 (Part II)

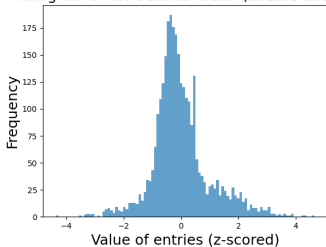
Histogram of the Second Factor (2021.1-2021.6)



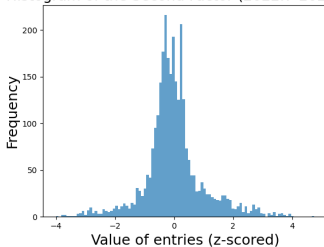
Histogram of the Second Factor (2021.7-2021.12)



Histogram of the Second Factor (2022.1-2022.6)

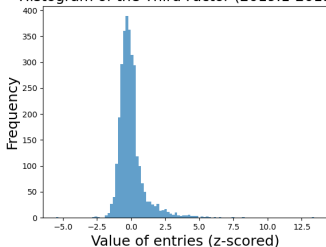


Histogram of the Second Factor (2022.7-2022.12)

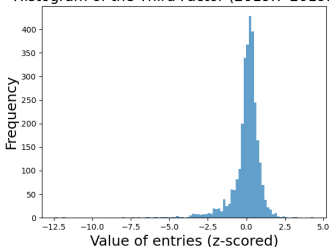


Z-score the third factor (Part I)

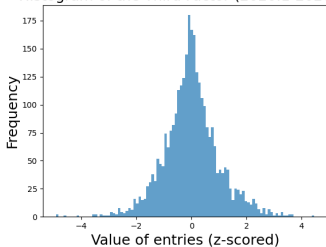
Histogram of the Third Factor (2019.1-2019.6)



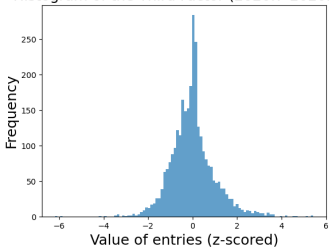
Histogram of the Third Factor (2019.7-2019.12)



Histogram of the Third Factor (2020.1-2020.6)

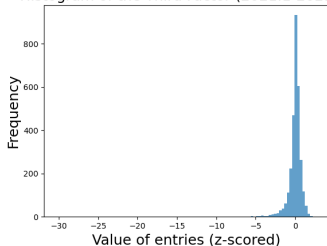


Histogram of the Third Factor (2020.7-2020.12)

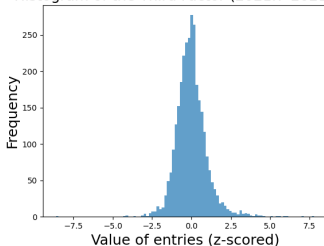


Normalize the first factor to mean 1 (Part II)

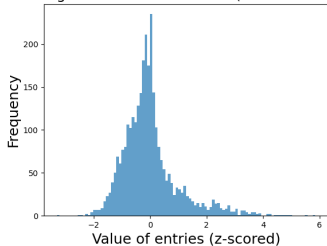
Histogram of the Third Factor (2021.1-2021.6)



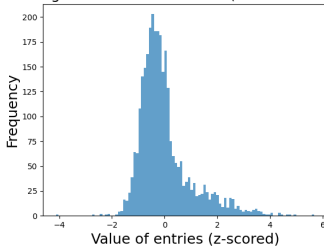
Histogram of the Third Factor (2021.7-2021.12)



Histogram of the Third Factor (2022.1-2022.6)

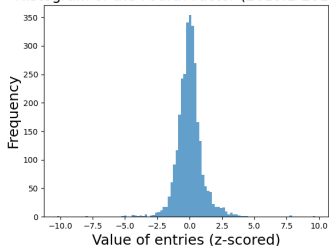


Histogram of the Third Factor (2022.7-2022.12)

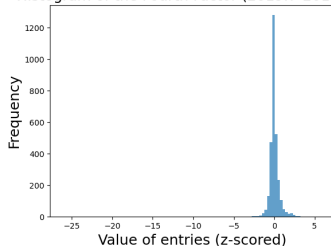


Z-score the fourth factor (Part I)

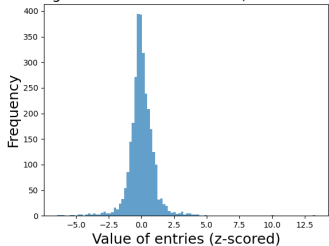
Histogram of the Fourth Factor (2019.1-2019.6)



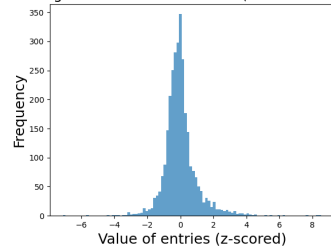
Histogram of the Fourth Factor (2019.7-2019.12)



Histogram of the Fourth Factor (2020.1-2020.6)

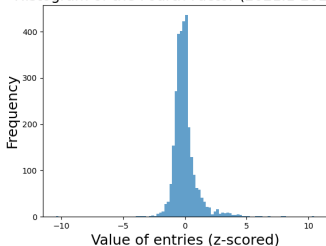


Histogram of the Fourth Factor (2020.7-2020.12)

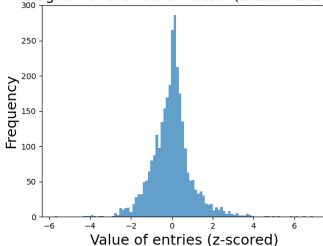


Normalize the first factor to mean 1 (Part II)

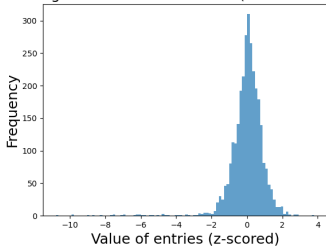
Histogram of the Fourth Factor (2021.1-2021.6)



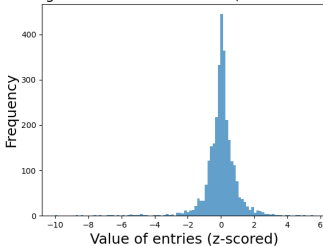
Histogram of the Fourth Factor (2021.7-2021.12)



Histogram of the Fourth Factor (2022.1-2022.6)



Histogram of the Fourth Factor (2022.7-2022.12)



Limitations and work in progress

- Further interpreting the histograms of eigenvectors
- Siamak idea: Instead of only looking at how the mean and variance of spiked eigenvectors change over time, can we use ML (say random forest) to make predictions from the histograms? I guess this may entail computing metrics from the histograms and training the ML alg on a bunch of histogram metrics. I'm not sure how much people have done this in this context.
- Siamak idea: Can we compute kurtosis/peakedness of the presented histograms? I would think this would be a telling comparative point among the different cases.
- Include or exclude outliers in the dataset
- Look at the plots on the same horizontal scales and look at the four factor panels for one date at a time
- And much more to do...

Appendix A: List of Removed Outliers

- 2019.1 - 2019.6
 - AXSOME THERAPEUTICS INC
- 2019.7 - 2019.12
 - CHEMOCENTRYX INC
 - SYNTHORX INC
 - KARUNA THERAPEUTICS INC
 - INTRA CELLULAR THERAPIES INC
 - NEXTCURE INC
- 2020.1 - 2020.6
 - ENABLE MIDSTREAM PARTNERS LP
 - SORRENTO THERAPEUTICS INC
 - MACROGENICS INC
 - ADAPTIMMUNE THERAPEUTICS PLC
 - M F A FINANCIAL INC
 - WINS FINANCE HOLDINGS INC
- 2020.7 - 2020.12
 - SERES THERAPEUTICS INC

List of Removed Outliers (cont.)

- 2021.1 - 2021.7
 - A M C ENTERTAINMENT HOLDINGS INC
 - HISTOGENICS CORP
- 2021.7 - 2021.12
 - STATE AUTO FINANCIAL CORP
- 2022.1 - 2022.7
 - None
- 2022.7 - 2022.12
 - PROMETHEUS BIOSCIENCES INC
 - MADRIGAL PHARMACEUTICALS INC THERAPEUTICS PLC
 - PLIANT THERAPEUTICS INC
 - SUMMIT THERAPEUTICS INC