

Similarity Analysis in Financial Time Series



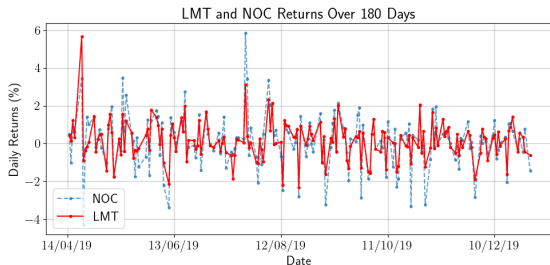
Zella Baig

`zella.baig@maths.ox.ac.uk`

Supervisors :

Mohsin Javed (BlackRock)

Yuji Nakatsukasa (University of Oxford)



- ▶ Problem: We wish to extract some measure of “similarity” between stocks, but they are noisy. Can we separate **noise** and **signal**?
- ▶ A possible solution: **Singular Spectrum Analysis (SSA)**.

Consider¹ a time series of observations $Z_T = (z_1, \dots, z_T)$. With fixed *window length* L and with $K := T - L + 1$:

1. Construct the (Hankel) trajectory matrix:

$$\mathbf{X} := \begin{bmatrix} z_1 & z_2 & z_3 & \dots & z_K \\ z_2 & z_3 & z_4 & \dots & z_{K+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ z_L & z_{L+1} & z_{L+2} & \dots & z_T \end{bmatrix} \quad (1)$$

¹Hassani, Mahmoudvand, et al. 2011.

2. Compute the *singular value decomposition* (SVD) of \mathbf{X} :

$$\mathbf{X} = \sum_{i=1}^n u_i v_i^T \sigma_i$$

3. Truncate the SVD to r rank-1 matrices, with *rank* r chosen s.t. $r \leq n$:

$$\mathbf{X} \approx \mathcal{X} = \sum_{i=1}^r u_i v_i^T \sigma_i$$

4. \mathcal{X} is not necessarily Hankel, so re-diagonalise on the off-diagonals to reconstruct a de-noised series $\bar{Z}_T = (\bar{z}_1, \dots, \bar{z}_T)$ from the averaged Hankel matrix

$$\bar{\mathbf{X}} := \begin{bmatrix} \bar{z}_1 & \bar{z}_2 & \bar{z}_3 & \dots & \bar{z}_K \\ \bar{z}_2 & \bar{z}_3 & \bar{z}_4 & \dots & \bar{z}_{K+1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \bar{z}_L & \bar{z}_{L+1} & \bar{z}_{L+2} & \dots & \bar{z}_T \end{bmatrix} \quad (2)$$

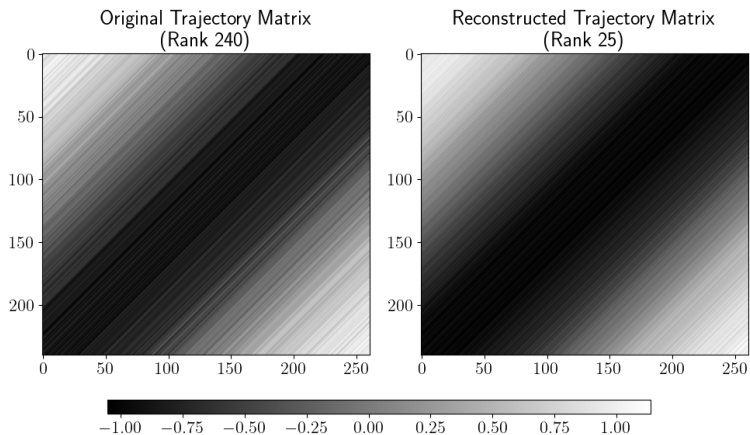


Figure 1: Noisy sinusoidal signal and denoised signals' trajectory matrices.

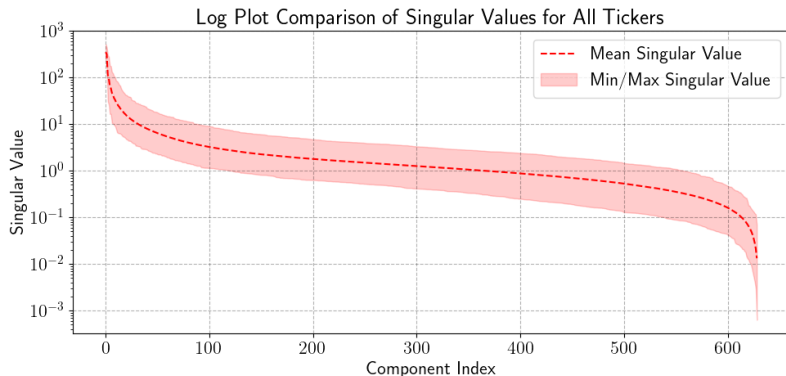


Figure 2: Choose r via examination of the scree plot, with knee at $r \approx 25$.

We measure similarity of two time (de-noised) time series using the **Time Warped Edit Distance² (TWED)**. Why?

1. Cointegration v.s. correlation.
2. Elasticity
3. (Relatively) cheap

²Marteau 2008.

	1	∞	18	17	16
Y_T	1	∞	13	12	13
	9	∞	4	13	14
	0	∞	∞	∞	∞
			5	-3	1
				X_T	

Figure 3: Populated TWED grid, with $\nu = \lambda = 0.5$. $D_{T,T} = 16$.

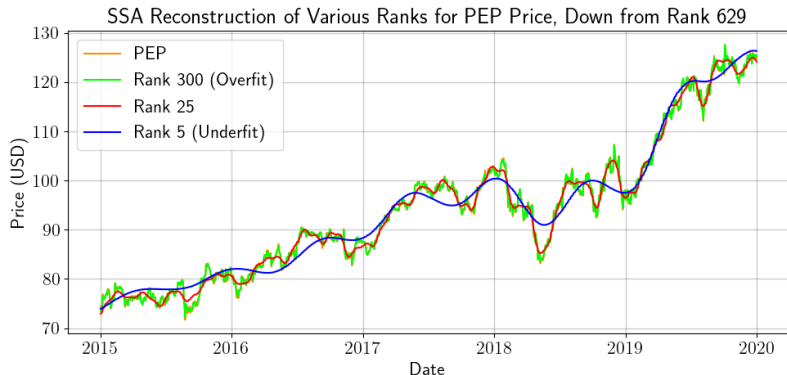


Figure 4: Different rank SSA reconstructions. Note underfitting at $r = 5$, and overfitting at $r = 300$.

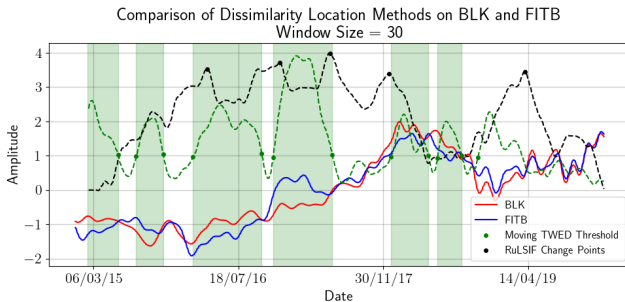
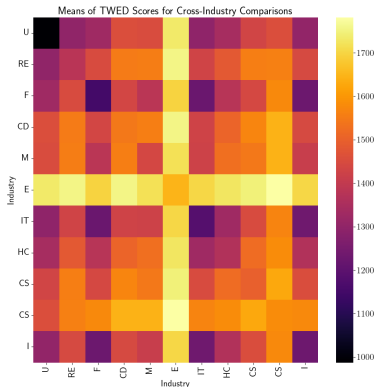


Figure 5: Implementation on our two dissimilarity finding methods on BLK and FITB over five years.

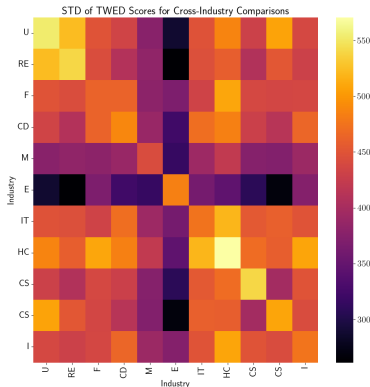
Key takeaways:

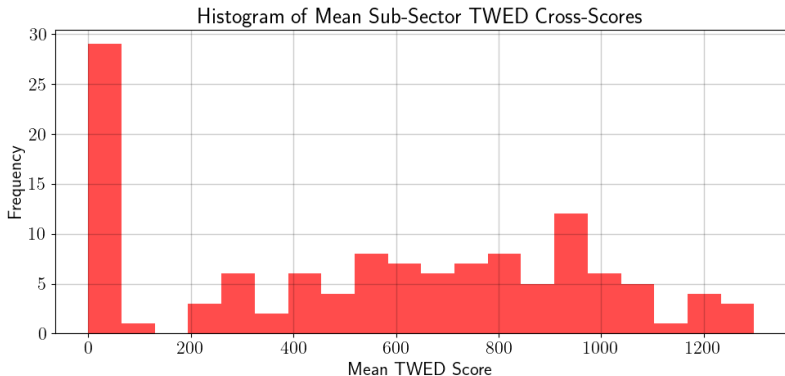
- ▶ Energy, Consumer Staples sector dissimilar to other sectors.
- ▶ Utilities, Finance, IT show strong inter-similarity.



Key takeaways:

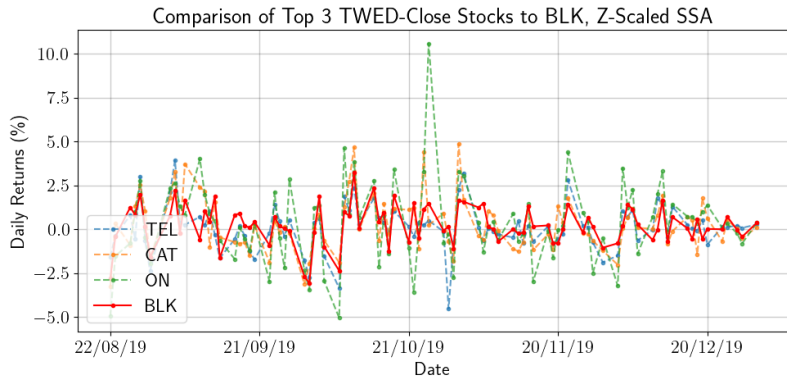
- ▶ Energy sector conclusions strong.
- ▶ Utilities, Health Care conclusion very weak.



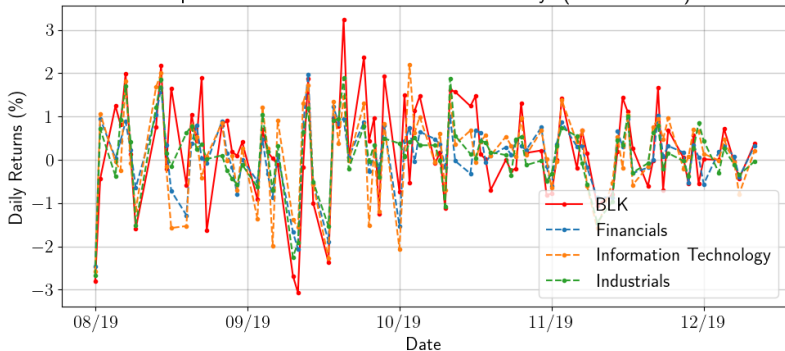


Similar Tickers	Top Five Tickers by Similarity Method			
	d_E	d_{DTW}	d_{TWED}	...
TIC_1	TIC_1	TIC_2	TIC_7	...
TIC_2	TIC_{23}	TIC_{11}	TIC_{99}	...
TIC_3	TIC_2	TIC_8	TIC_3	...
	TIC_{366}	TIC_{223}	TIC_{23}	...
	TIC_3	TIC_3	TIC_7	...

Table 1: Example scoring method used to help choose a similarity measure.

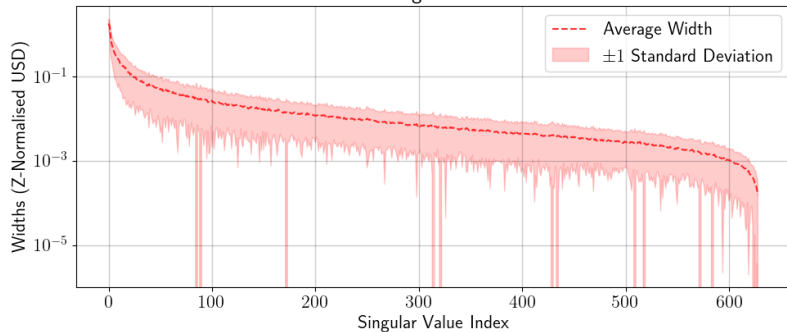


Top 3 Closest Industries to BLK Over 90 Days (Intra Method)



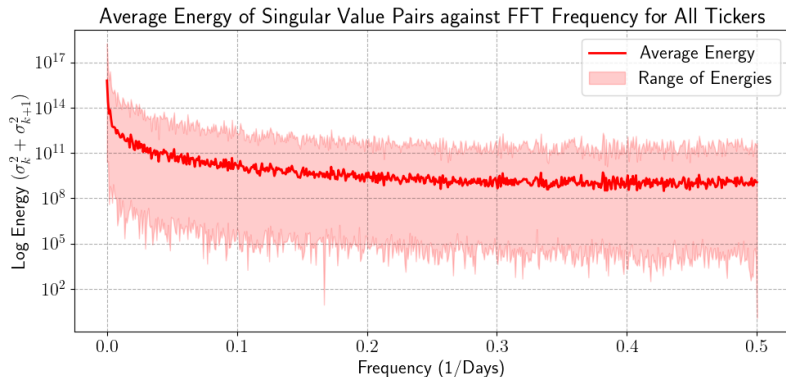
Example pairs of (potentially surprising!) similar tickers:

- ▶ Visa/Microsoft (Financials/IT)
- ▶ Mastercard/Intuit (Financials/IT)
- ▶ MSCI/Cintas (Financials/Industrials)
- ▶ TransDigm/FICO (Industrials/IT)
- ▶ Steris/CoStar (Health Care/Industrials)

Average Z-Normalised "Band Width" of Reconstructed Series
for Each Singular Value Index

Average Frequency of Dominant Frequency Component
for Each Singular Value Index





- ▶ Proper orthogonal decomposition
- ▶ Similarity v.s. distance
- ▶ `series_scorer`: A Python package for multiple time series scoring.
- ▶ Back testing!

Thank you!

`zella.baig@maths.ox.ac.uk`

- (1) Dey, S.; Dutta, A.; Toledo, J. I.; Ghosh, S. K.; Lladós, J.; Pal, U. SigNet: Convolutional Siamese Network for Writer Independent Offline Signature Verification, Number: arXiv:1707.02131, 2017.
- (2) Ghodsi, M.; Hassani, H.; Rahmani, D.; Silva, E. S. *Journal of Applied Statistics* **2018**, *45*, Publisher: Taylor & Francis
_eprint: <https://doi.org/10.1080/02664763.2017.1401050>, 1872–1899.
- (3) Hassani, H.; Kalantari, M.; Yarmohammadi, M. *Comptes Rendus Mathématique* **2017**, *355*, 1026–1036.
- (4) Hassani, H.; Mahmoudvand, R.; Zokaei, M. *Comptes Rendus Mathématique* **2011**, *349*, 987–990.

- (5) Hou, L.; Jin, X.; Zhao, Z. In *2019 12th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI)*, 2019 12th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI), 2019, pp 1–6.
- (6) Marteau, P.-F. Time Warp Edit Distance, Number: arXiv:0802.3522, 2008.
- (7) Serrà, J.; Arcos, J. L. *Knowledge-Based Systems* **2014**, *67*, 305–314.

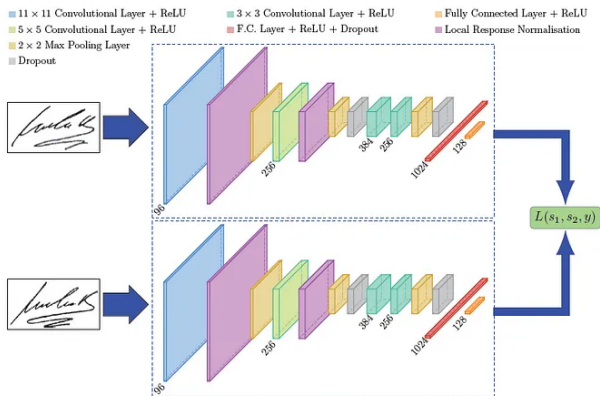


Figure 6: Overview of an SNN, as used in SigNet³.

³Dey et al. 2017.

There exist two different types of SSA forecasting: recurrent, and vector. We go over them in turn:

1. **Recurrent forecasting**⁴: Consider the left singular vectors u_1, u_2, \dots, u_r . Take their L^{th} components, denoted π_i , and define

$$v^2 := \sum_{i=1}^r \pi_i^2. \quad (3)$$

Denote by \hat{u}_i the $L - 1 \times 1$ vector which is u_i with the final component removed.

⁴Ghodsi et al. 2018.

Then define

$$A = (\alpha_{L-1}, \dots, \alpha_1)^T = \frac{1}{1-v^2} \sum_{i=1}^r \pi_i \hat{u}_i,$$

and thus construct

$$z_t = \begin{cases} \bar{z}_t & t = 1, \dots, T, \\ \sum_{i=1}^{L-1} \alpha_i z_{t-i} & t = T+1, \dots, T+h, \end{cases}$$

for a forecast to h steps ahead.

2. **Vector forecasting**⁵: First define

$$\hat{\mathbf{U}} = [\hat{u}_1, \dots, \hat{u}_r],$$

and construct

$$\mathbf{\Pi} = \hat{\mathbf{U}}\hat{\mathbf{U}}^T + (1 - v^2) AA^T.$$

Finally, construct the operator Θ s.t.

$$\Theta V := \begin{bmatrix} \mathbf{\Pi}\hat{V} \\ A^T\hat{V} \end{bmatrix},$$

where \hat{V} denotes the vector V with the last element removed.

⁵Ghodsi et al. 2018.

Define now

$$\Xi_i = \begin{cases} \mathcal{X}_i & i = 1, \dots, K, \\ \Theta \Xi_{i-1} & i = K + 1, \dots, K + h + L - 1, \end{cases}$$

where \mathcal{X}_i are the columns of \mathcal{X} . Next construct

$$\Xi = [\Xi_1, \dots, \Xi_{K+h+L-1}],$$

and hankelise to get the matrix $\bar{\Xi}$ from which we recover an “extended” time series containing forecasted values.