

# Intraday trading volume and non-negative matrix factorization

Hellinton H. Takada and Julio M. Stern

Citation: [AIP Conference Proceedings](#) **1757**, 060006 (2016); doi: 10.1063/1.4959065

View online: <http://dx.doi.org/10.1063/1.4959065>

View Table of Contents: <http://aip.scitation.org/toc/apc/1757/1>

Published by the [American Institute of Physics](#)

---

---

# Intraday Trading Volume and Non-Negative Matrix Factorization

Hellinton H. Takada<sup>\*,†</sup> and Julio M. Stern<sup>†</sup>

<sup>\*</sup>*Quantitative Research, Itaú Asset Management, São Paulo, Brazil*

<sup>†</sup>*Institute of Mathematics and Statistics, University of São Paulo, Brazil*

**Abstract.** The intraday trading volume of a security is the total amount of traded contracts distributed over the day. Consequently, the intraday trading volume captures part of the intraday trading activity and represents a proxy for the intraday liquidity of a market. When executing orders in the market, it is important to avoid impacting the market and, consequently, the prices. Usually, the market impact causes adverse price movements implying in an implicit cost of the execution process. Clearly, the intraday trading volume is important when developing execution strategies. In the literature, the intraday trading volume for equities has been reported to possess an intraday U-shaped pattern. In this paper, we investigate for the first time the statistical factors behind the intraday trading volume using non-negative matrix factorization (NNMF). The obtained factors are directly applicable to the design of execution strategies. Additionally, we observe that the factors obtained using NNMF are more intuitively interpretable than the factors obtained using principal component analysis. Our empirical conclusion helps to corroborate several other applications of NNMF from the literature where the same behavior is observed.

**Keywords:** Information theory, Entropy, Financial markets

**PACS:** 89.70.-a, 89.70.Cf, 89.65.Gh

## INTRODUCTION

The total amount of traded contracts of a security over the trading period of a day is called intraday trading volume. For traders, the intraday trading volume is very important because of its use in technical analysis [1]. Obviously, the intraday trading volume captures part of the intraday trading activity and represents a proxy for the intraday liquidity of a security. The intraday liquidity is the source of contracts from where the execution of an order is made possible. Evidently, the lack of liquidity causes a problem when the amount of contracts to be executed is very large. In such a case, the execution of the order becomes impossible or causes adverse price distortion.

The price distortion is called market impact cost. In other words, the market impact cost measures the adverse change in the market price due to the execution of an order and is rapidly becoming the dominant transaction cost [2]. Discernibly, the intraday trading volume is important when developing execution strategies. The execution strategies are part of the service provided by traders and brokers to their clients. Additionally, the increase of investments in trading technologies enabled the automation of such strategies. An important benchmark for execution strategies is the volume weighted average price (VWAP) which requires models for the intraday trading volume.

In the literature, the intraday trading volume for equities has been reported to possess an intraday U-shaped pattern, i.e. heavy trading volume at the beginning and at the end

of the trading day and the relatively light trading volume at the middle of the trading day (see for example [3]). As a consequence, several approaches were developed to model the intraday trading volume (e.g. it is used a beta density function to fit the U-shaped pattern in [4]). In this paper, we investigate the statistical factors behind the intraday trading volume. Statistical factors are unobserved variables used to describe observed data.

Non-negative matrix factorization (NNMF) is a multivariate data analysis technique aimed to estimate non-negative factors and factor loadings from non-negative data. NNMF was invented by Paatero and Tapper in 1994 under the name positive matrix factorization (PMF) [5] and the name NNMF was established by Lee and Seung in 1999 [6]. There are several applications of NNMF such as text mining [7], image processing [8], sound processing [9], identification of concentrations in chemistry [5], recognition of underlying trends in stock prices [10], modeling of the term structure of interest rates [11], and so on.

In terms of factorization techniques, the most popular approach is the principal component analysis (PCA) which was introduced by Pearson [12] and developed by Hotelling [13]. In this paper, we compare the factors and factor loadings obtained from intraday trading volume using PCA to those using NNMF. Since the trading volume is represented by non-negative quantities, it is a natural choice the use of NNMF. The NNMF factors and factor loadings are intuitively interpretable and directly applicable to the design of execution strategies with VWAP as their benchmark.

## **INTRADAY TRADING VOLUME PATTERN**

Since 1980s, several works in the related literature report the U-shaped pattern of intraday trading volume [14, 15], i.e. heavy trading volume at the beginning and at the end of the trading day and relatively light trading volume at the middle of the trading day. Typically, markets with defined daily openings and closures (e.g. equity and bond markets) present a distorted U-shaped trading activity over the trading day. Differently, markets with round-the-clock trading (e.g. foreign exchange interbank market) produce more complex patterns.

Several rationalizations were made to explain the U-shape from the interaction of distinct customer groups and market makers. For example, [16] provides a partial explanation of the empirical findings concerning the pattern of volume in intraday transaction data showing that concentrated-trading patterns arises endogenously as a result of the strategic behavior of liquidity traders and informed traders. Alternatively, [17] associates the U-shaped curves to market closure, the power of dealers, and portfolio rebalancing.

Assuming the U-shaped pattern of intraday trading volume, there are approaches available to identify and estimate some parametric specifications [4, 18]. In practice, practitioners obtain the intraday trading volume pattern using the average of historical executed volumes of the last 21 days [19]. As it was mentioned, a model for intraday trading volume is important because of VWAP benchmark. Actually, the VWAP is the target of several execution strategies and the design of such strategies also depends on models for intraday trading volume. In the next section, we explore the U-shaped pattern using factor models such as PCA and NNMF.

## INTRADAY TRADING VOLUME FACTORS

Basically, statistical factors are obtained from factor analysis, a statistical procedure to describe observed data in terms of unobserved variables called factors. The objective of factor analysis is to reduce the dimensionality of the original data  $D = [d_{ij}] \in \mathbb{R}^{m \times p}$ ,  $m \wedge p \in \mathbb{N}_+$ , using an approximation  $\Delta = [\delta_{ij}] \in \mathbb{R}^{m \times p}$  such that

$$D \approx \Delta = \Phi\Lambda, \quad (1)$$

where  $\Phi = [\phi_{ij}] \in \mathbb{R}^{m \times k}$  is the matrix of factors or unobserved (latent) variables;  $\Lambda = [\lambda_{ij}] \in \mathbb{R}^{k \times p}$  is the matrix of factor loadings or weights;  $k$  represents the number of factors ( $k \leq \min(m, p)$ ). In the literature, the most popular approach for factor analysis is PCA. On the other hand, NNMF is a more recent approach.

In our case,  $D$  represents the intraday traded volume in number of contracts,  $p$  represents the number of time bins during a trading day (e.g. a time bin is from 10:00 a.m. until 11:00 a.m.) and  $m$  is the amount of different days and/or equity names. The securities selected are from the Brazilian stock exchange (BM&F Bovespa) for the period from April 2013 until September 2013. Additionally, we present the results obtained with time bins equal to 1 hour with a total of  $p = 8$  time bins a day. We investigate factors and factor loadings from intraday traded volume for only one individual equity name or for a set of different equities. Consequently, we focus not only individual estimations but also joint estimations to obtain intraday trading volume patterns.

In particular, the ticker names of the equities chosen are: PETR3, PETR4, VALE3, VALE5, BBDC3 and BBDC4. PETR3 is the ordinary stock of an oil and gas company, PETR4 is the preferential stock of the same company and one of the most liquid shares traded in Brazil; VALE3 is the ordinary stock of an mining company and VALE5 is the preferential stock of the same company; BBDC3 is the ordinary stock of a financial services company and BBDC4 is the preferential stock of the same company. In the following, we describe the procedures to obtain the PCA and the NNMF factorizations.

### Factors using PCA

The singular value decomposition (SVD) is a technique from linear algebra used to obtain the factors and factor loadings from PCA [20] and results in the following factorization

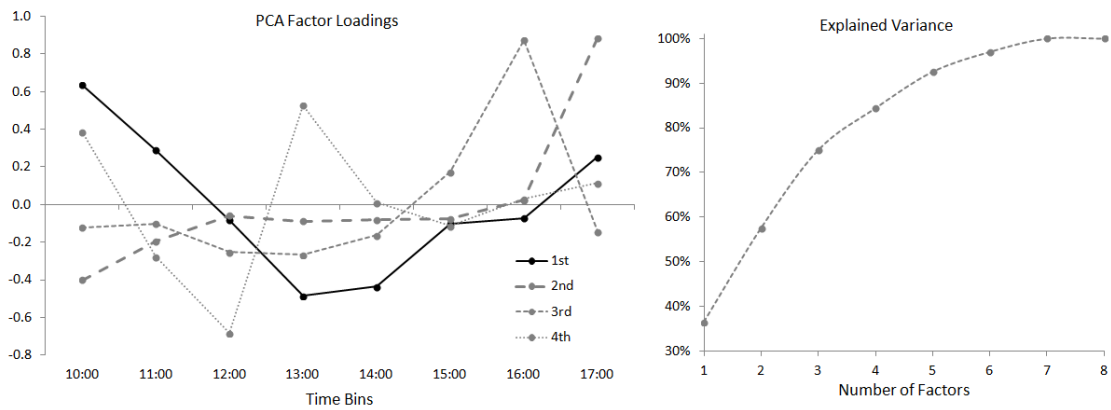
$$\dot{D} = USV', \quad (2)$$

where  $\dot{D} = [\dot{d}_{ij}] \in \mathbb{R}^{m \times p}$  is obtained mean centering the data matrix  $D$ ;  $U = [u_{ij}] \in \mathbb{R}^{m \times p}$ ;  $S = [s_{ij}] \in \mathbb{R}^{p \times p}$  is a diagonal matrix such that  $s_{11} \geq s_{22} \geq \dots \geq s_{pp}$ ;  $V = [v_{ij}] \in \mathbb{R}^{p \times p}$ ;  $VV' = I_p$ . Given the number of factors  $k$ , the PCA  $k$ -factor model is

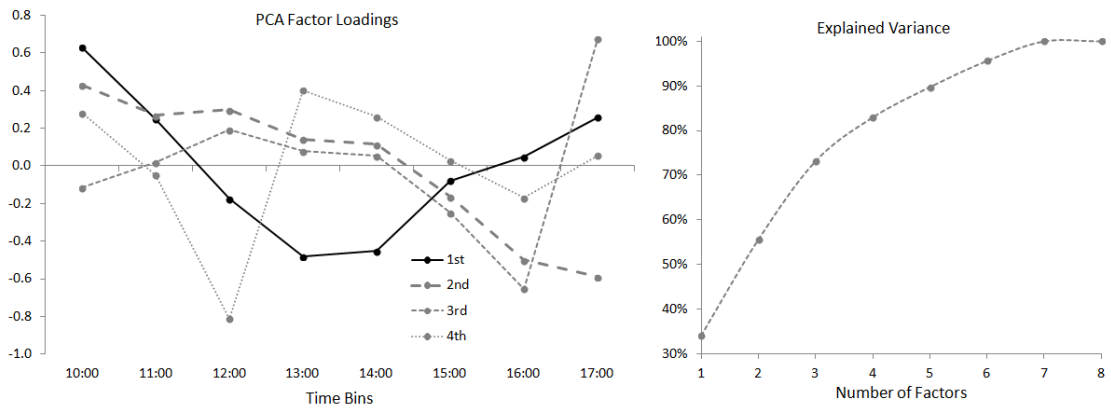
$$\dot{D} \approx \hat{D} = U\hat{S}\hat{V}', \quad (3)$$

where  $\hat{S} = [s_{ij}] \in \mathbb{R}^{p \times k}$  and  $\hat{V} = [v_{ij}] \in \mathbb{R}^{p \times k}$ . The columns of  $\hat{F} = U\hat{S} = [\hat{f}_{ij}] \in \mathbb{R}^{m \times k}$  are the factors and the rows of  $\hat{L} = \hat{V}' = [\hat{l}_{ij}] \in \mathbb{R}^{k \times p}$  are the corresponding factor loadings ( $\hat{D} = \hat{F}\hat{L}$ ).

In Figures 1 and 2, we show the first four factor loadings  $\{[\hat{l}_{1j}], [\hat{l}_{2j}], [\hat{l}_{3j}], [\hat{l}_{4j}]\}$  and the percentage of total explained variance of  $D$  according to the number of factors for PETR4 and for the set of equities, respectively. It is possible to notice that the first factor loading (with higher percentage of total explained variance) for both cases has the U-shaped pattern. However, the second, third and fourth factor loadings do not possess a direct interpretation. Additionally, the percentage of total explained variance for the first factors are very low indicating the need of more factors (not only one) to be used in the PCA based factor model for intraday trading volume (for example, in both cases illustrated in Figures 1 and 2, it is necessary to include at least three factors to explain more than 70% of the total data variance).



**FIGURE 1.** The first four factor loadings  $\{[\hat{l}_{1j}], [\hat{l}_{2j}], [\hat{l}_{3j}], [\hat{l}_{4j}]\}$  (left figure) and the percentage of total explained variance of  $D$  according to the number of factors (right figure) for PETR4.



**FIGURE 2.** The first four factor loadings  $\{[\hat{l}_{1j}], [\hat{l}_{2j}], [\hat{l}_{3j}], [\hat{l}_{4j}]\}$  (left figure) and the percentage of total explained variance of  $D$  according to the number of factors (right figure) for the set of equities.

## Factors using NNMF

Observing that the data matrix  $D$  containing the traded volume is non-negative  $D = [d_{ij}] \in \mathbb{R}_{\geq 0}^{m \times p}$  and given the number of factors  $k$ , the NNMF approach aims to find the following approximation

$$D \approx \tilde{D} = \tilde{F}\tilde{L}, \quad (4)$$

where  $\tilde{D} = [\tilde{d}_{ij}] \in \mathbb{R}_{\geq 0}^{m \times p}$ ;  $\tilde{F} = [\tilde{f}_{ij}] \in \mathbb{R}_{\geq 0}^{m \times k}$ ;  $\tilde{L} = [\tilde{l}_{ij}] \in \mathbb{R}_{\geq 0}^{k \times p}$ . It is important to state that the columns of  $\tilde{F}$  are the factors and the rows of  $\tilde{L}$  are the factor loadings.

The NNMF optimization procedures minimize the approximation error between  $D$  and  $\tilde{D}$ . In a generalized way, the Bregman divergence  $D_\varphi(D||\tilde{D})$  is used as the objective function to be minimized [21, 22]. Considering only separable Bregman divergences,

$$D_\varphi(D||\tilde{D}) = \sum_{ij} D_\varphi(d_{ij}||\tilde{d}_{ij}) = \sum_{ij} \{\varphi(d_{ij}) - \varphi(\tilde{d}_{ij}) - \nabla\varphi(\tilde{d}_{ij})[\varphi(d_{ij}) - \varphi(\tilde{d}_{ij})]\}, \quad (5)$$

where  $\varphi(\cdot)$  is a strictly convex function with a continuous first derivative. Formally, the resulting optimization problems are

$$\min_{\tilde{F}, \tilde{L} \geq 0} \{D_\varphi(D||\tilde{F}\tilde{L}) + J(\tilde{F}) + G(\tilde{L})\} \quad (6)$$

or

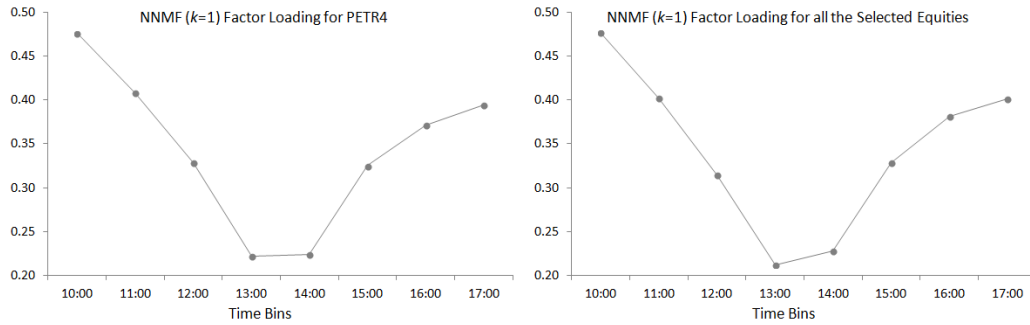
$$\min_{\tilde{F}, \tilde{L} \geq 0} \{D_\varphi(\tilde{F}\tilde{L}||D) + J(\tilde{F}) + G(\tilde{L})\}, \quad (7)$$

where  $J(\cdot)$  and  $G(\cdot)$  are penalty functions to enforce certain application-dependent characteristics of the solution, such as sparsity and /or smoothness. It is also important to remember that the Bregman divergences are not symmetric in general. Consequently, we consider  $D_\varphi(D||\tilde{F}\tilde{L})$ .

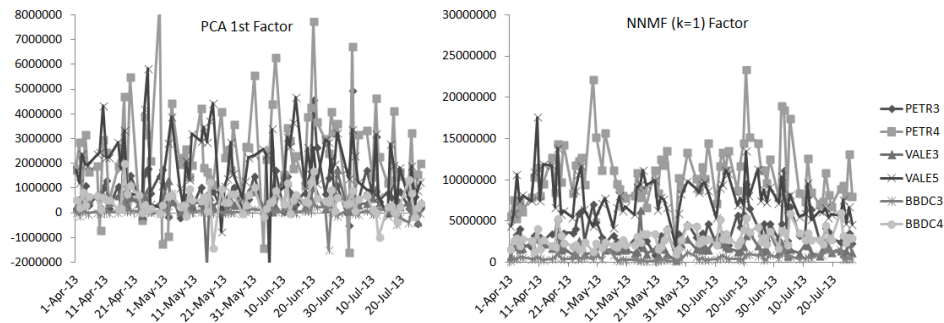
Adopting  $\varphi(x) = x^2/2$  and  $J(\cdot) = G(\cdot) = 0$ , there are some known algorithms to solve the NNMF problem. The algorithms are divided in general classes [23]: gradient descent algorithms, multiplicative update algorithms and alternating least squares algorithms (ALS). Here, the ALS will be adopted (the use of other algorithms does not provide great differences to the empirical applications presented in the following sections).

## One Factor using NNMF

In this section, the objective is to explore the existence of one main factor in the intraday traded volume ( $k = 1$ ). In Figure 3, we present the obtained factor loadings  $[\tilde{l}_{1j}]$  for our data set. As expected, it is possible to identify the U-shaped pattern observing the factor loadings of PETR4 and the set of equities. For the estimation using the set of equities, we present in Figure 4 the first factors obtained from PCA and from NNMF ( $k = 1$ ) for each equity. The NNMF factors can be easily interpreted as the volume level for each equity. Consequently, it is clear that PETR4 has the highest traded volume over the time while BBDC3 has the lowest traded volume over the time. Concerning the PCA factors, since they become negative such an interpretation is not possible.



**FIGURE 3.** The NNMF ( $k = 1$ ) factor loading for PETR4 (left figure) and the NNMF ( $k = 1$ ) factor loading for the set of equities (right figure).



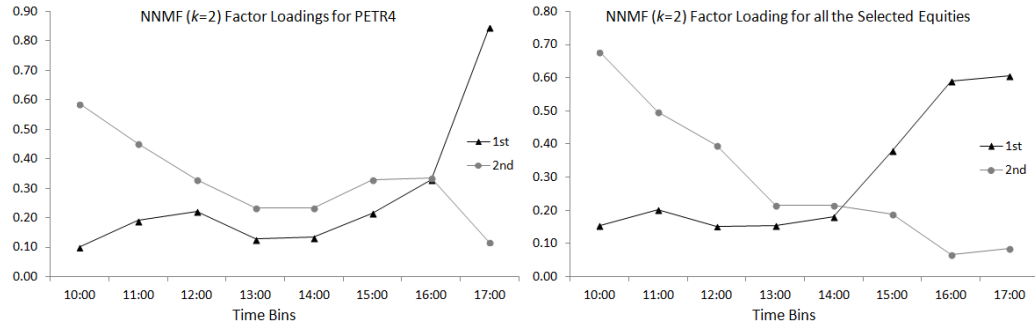
**FIGURE 4.** The PCA first factors (left figure) and the NNMF ( $k = 1$ ) factors (right figure) for each equity in the set of equities.

## Two Factors using NNMF

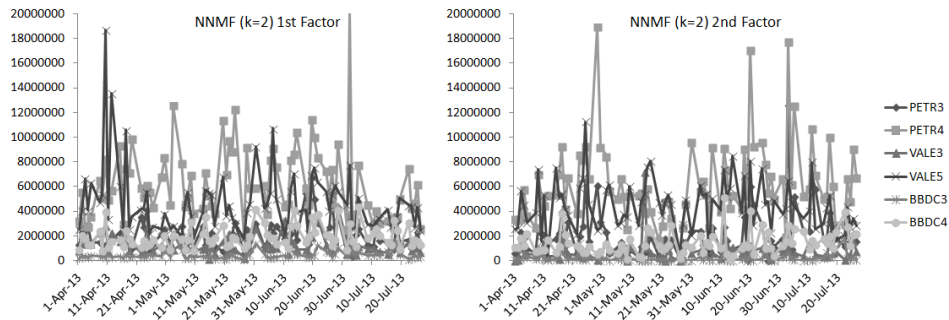
In this section, the objective is to explore the existence of two factors in the intraday trading volume ( $k = 2$ ). In Figure 5, we present the obtained factor loadings  $\{[\tilde{l}_{1j}], [\tilde{l}_{2j}]\}$  for our data set. As expected, it is possible to identify an increasing end of day pattern (first factor loading) and a decreasing start of day pattern (second factor loading) observing the factor loadings for PETR4 and the set of equities.

In Figure 6, we present the first and second factors obtained for each selected equity. The NNMF two factors can be easily interpreted as the volume level at the beginning and at the end of the day for each equity name. Consequently, using the two NNMF factors a financial analysis could study individually the volume levels at the beginning and at the end of the day. Again, concerning the PCA factors, since they become negative such an interpretation is not possible.

Finally, the residual sum of squares (RSS) of the factor models for the joint estimation of the selected equities is given in Table 1. As expected, for a same number of  $k$  the NNMF reduces the RSS compared with the PCA. In Table 2, we present the explained percentage of the total data variance of the factor models for the joint estimation of selected equities. Again, as expected, for a same number of  $k$  the NNMF explains a higher percentage of the data variance compared with PCA.



**FIGURE 5.** The NNMF ( $k = 2$ ) factor loadings for PETR4 (left figure) and the NNMF ( $k = 2$ ) factor loadings for the set of equities (right figure).



**FIGURE 6.** The NNMF ( $k = 2$ ) first factors (left figure) and the NNMF ( $k = 2$ ) second factors (right figure) for each selected equity.

**TABLE 1.** The RSS of the factor models for the joint estimation of selected equities.

$k$	PCA	NNMF
1	2.04e+16	3.16e+15
2	1.95e+16	2.25e+15

**TABLE 2.** The explained % of total data variance of the factor models for the joint estimation of selected equities.

$k$	PCA	NNMF
1	34.09%	72.35%
2	55.60%	80.30%



## CONCLUSIONS

In this paper, NNMF was for the first time applied to capture the intraday trading volume patterns. Considering NNMF with only one factor, we identified for our selected equities the well-known U-shaped intraday trading volume pattern. The U-shaped pattern is very important for execution strategies based on VWAP. Additionally, we also identified interpretable factors when considering NNMF with two factors. One factor represents the volume level at the start of the trading day and the other factor represents the volume level at the end of the trading day. The two factors enable the individual study of the trading volume level at the start and at the end of the trading day. As expected, our empirical results show that for a given number of factors the NNMF has a higher percentage of explained variance and lower RSS than PCA.

## REFERENCES

1. M. Leibovit, *The Trader's Book of Volume: The Definitive Guide To Volume Trading*, McGraw-Hill, New York, 2011.
2. I. Aldridge, *High-Frequency Trading: A Practical Guide to Algorithmic Strategies and Trading Systems*, John Wiley & Sons, Inc., Hoboken, 2013.
3. P. J. Jain, and G. Joh, *J. Financ. Quant. Anal.* **23** (3), 269–284 (1988).
4. E. Panas, *Appl. Econ.* **37** (2), 191–199 (2005).
5. P. Paatero, and U. Tapper, *Environmetrics* **5** (2), 111–126 (1994).
6. D. Lee, and H. Seung, *Nature* **401** (6755), 788–791 (1999).
7. M. W. Berry (editor), *Computational Information Retrieval*, Philadelphia: Society for Industrial and Applied Mathematics, 2001.
8. D. Lee, and H. Seung, *Advances in Neural Information Processing Systems* **13**, 556–562 (2001).
9. P. Smaragdakis, and J. C. Brown, "Non-negative matrix factorization for polyphonic music transcription," *Proc. IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 177–180 (2003).
10. K. Drakakis, S. Rickard, R. de Fréin, and A. Cichocki, *Int. Math. Forum* **3** (38), 1853–1870 (2008).
11. H. H. Takada, and J. M. Stern, *AIP Conf. Proc.* **1641** (369), 369–377 (2015).
12. K. Pearson, *Philos. Mag.* **2** (11), 559–572 (1901).
13. H. Hotelling, *J. Educational Psychol.* **24** (6), 417–441 (1933).
14. R. A. Wood, T. H. McInish, and J. K. Ord, *J. Financ.* **40** (3), 723–739 (1985).
15. L. Harris, *J. Financ. Econ.* **16** (1), 99–117 (1986).
16. A. R. Admati, and P. Pfleiderer, *Rev. Financ. Stud.* **1** (1), 3–40 (1988).
17. W. A. Brock, and A. Kleidon, *J. Econ. Dyn. Control* **16** (3-4), 451–489 (1992).
18. S. V. Aradhyula, and A. T. Ergün, *Applied Financial Economics* **14** (13), 909–913 (2004).
19. R. Kissell, and M. Glantz, *Optimal Trading Strategies: Quantitative Approaches for Managing Market Impact and Trading Risk*, AMACOM, Inc., New York, 2003.
20. G. H. Golub, and C. F. V. Loan, *Matrix Computations*, The Johns Hopkins Univ. Press, 1996.
21. I. S. Dhillon, and S. Sra, *Adv. Neural Inf. Process. Syst.* **18**, 283–290 (2005).
22. L. Li, G. Lebanon, and H. Park, "Fast Bregman divergence NMF using Taylor expansion and coordinate descent," *Proc. 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, August 12–16, 2012.
23. M. W. Berry, M. Browne, A. N. Langville, V. P. Pauca, R. J. Plemmons, *Comput. Stat. Data An.* **52**, 155–173 (2007).