

Simulations for Data of Stock Market in Low-Dimensions

Zupeng Zhao, Weidong Fang

School of Mathematics, South China University of Technology, Guangzhou Guangdong
Email: zhaozupeng950419@163.com

Received: Jan. 27th, 2019; accepted: Feb. 11th, 2019; published: Feb. 19th, 2019

Abstract

The situations of stock market are important references of a country's development in economy, and we often have to deal with the high-dimensional data when we analyze those data from stock market. It's a hard work to analyze those high-dimensional data directly, so we use the simulation methods by reducing the dimensions of variables to decrease the difficulty of analysis. SVD (Singular Value Decomposition), FA (Factor Analysis) and PCR (Principal Component Regression) are three most common simulation methods which are considered to reduce the dimensions of variables. In order to compare simulational effectiveness of the three methods, we used the method of theoretical deduction and demonstration, and got the result that the three methods had the same simulational effectiveness in some case. Hence it was able to draw the conclusion that those three methods of reducing the variable dimensions were coincident in some conditions.

Keywords

Dimensionality Reduction, SVD, FA, PCR, Simulation

股票市场数据的低维模拟

赵祖鹏, 方卫东

华南理工大学数学学院, 广东 广州
Email: zhaozupeng950419@163.com

收稿日期: 2019年1月27日; 录用日期: 2019年2月11日; 发布日期: 2019年2月19日

摘要

股票市场的情况是一个国家在经济上发展水平的重要参考, 在分析股票市场数据时往往需要处理高维度

的变量数据。直接分析这些高维度的变量数据是一件困难的工作, 因此在处理这些数据时会使用降低变量维度的模拟方法来减少分析的难度。奇异值分解、因子分析和主成分回归是三种最常见的被考虑用来降低变量维度的模拟方法。为了比较这三种方法的模拟效果, 本文中使用了理论推导和证明的方法, 得到三种方法在一定条件下有相同模拟效果的结果。于是可以得出在一定的条件下这三种降低变量维度的模拟方法具有一致性的结论。

关键词

降维, 奇异值分解, 因子分解, 主成分回归, 模拟

Copyright © 2019 by authors and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

1. 引言

股票市场中包含各种各样的时间序列数据, 例如股票收盘价、股票价格指数等。在处理多个时间序列的数据时, 需要分析的变量维度有时会很很高, 这给分析和处理带来了一定的难度。尽管每个时间序列都有各自的变化特点, 但它们之间有一定的相关性。利用这种相关性, 就能够减少变量的数量, 从而减少分析的难度和成本。

本文对奇异值分解[1]、因子分析[2] [3] [4]和主成分回归[5]三种常见的降维方法进行讨论。其中, 奇异值分解方法通过提取的奇异值来构造模拟矩阵; 因子分析通过提取因子来构造模拟矩阵; 主成分回归通过提取主成分和线性回归的方法来构造模拟矩阵。最后通过理论推导和证明的方式说明在一定的条件下, 三种方法得到的模拟结果是一致的。

2. 模拟方法

2.1. 基本假设

对于时间序列, 这里有一定的条件限制。第一点, 时间序列为平稳时间序列; 第二点, 时间序列的期望值为 0。本文中以行业价格指数作为例子, 不再对假设条件做过多的讨论。设 $\{X_j(t), t \in \mathbb{N}\}$ 是第 j ($j=1, 2, \dots, r$) 个行业股票指数每日收益率的时间序列, 其中

$$\text{每日收益率} = (\text{当日收盘价} - \text{昨日收盘价}) / \text{昨日收盘价}$$

为了简化操作, 将时间序列 $\{X_j(t), t \in \mathbb{N}\}$ 简化为随机变量 X_j 产生的多个独立同分布的样本。设向量

$$\mathbf{x}_j = [x_{1j}, x_{2j}, \dots, x_{nj}]^T$$

为第 j 个行业的时间序列对应的数值。那么对于所有的 r 个行业, 有数据矩阵(观测值矩阵)

$$\underset{(n \times r)}{\mathbf{A}} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r] \quad (1)$$

对应的样本协方差矩阵为

$$\underset{(r \times r)}{\mathbf{S}} = \begin{bmatrix} s_{11}^2 & \cdots & s_{1r}^2 \\ \vdots & \ddots & \vdots \\ s_{r1}^2 & \cdots & s_{rr}^2 \end{bmatrix}$$

假设 $n > r$, $\frac{1}{n} \sum_{i=1}^n x_{ij} = 0 (j=1, 2, \dots, r)$ 且 $\text{rank}(\mathbf{A}) = r$ 。

2.2. 奇异值分解模型

对式(1)定义的矩阵 \mathbf{A} , 由奇异值分解定理, 存在正交矩阵

$$\mathbf{U}_{(n \times n)} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n] \quad (2)$$

$$\mathbf{V}_{(r \times r)} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r] \quad (3)$$

和矩阵

$$\mathbf{M}_{(n \times r)} = \left[\begin{array}{ccc|c} \mu_1 & & & \mathbf{0} \\ & \mu_2 & & \\ & & \ddots & \\ & & & \mu_r \\ & & & \mathbf{0} \end{array} \right]^T \quad (4)$$

使

$$\mathbf{A} = \mathbf{U}\mathbf{M}\mathbf{V}^T$$

其中 $\mu_i (i=1, 2, \dots, r)$ 为 \mathbf{A} 的奇异值(默认 $\mu_1 \geq \mu_2 \geq \dots \geq \mu_r \geq 0$, 下文不再提及)。则

$$\mathbf{B}_{\text{SVD}} = \sum_{i=1}^s \mu_i \mathbf{u}_i \mathbf{v}_i^T \quad (5)$$

是矩阵 \mathbf{A} 的一个秩为 $s (s < r)$ 的同阶模拟矩阵, 均方误差

$$MSE_{\text{SVD}} = \frac{1}{(n-1)r} \text{tr}[(\mathbf{A} - \mathbf{B})(\mathbf{A} - \mathbf{B})^T] = \frac{1}{(n-1)r} \sum_{i=s+1}^p \mu_i^2$$

2.3. 因子分析模拟

假设 X_j 由 $s (s < r)$ 个公共因子 F_i 组成, 即

$$\begin{aligned} X_1 &= l_{11}F_1 + l_{12}F_2 + \dots + l_{1s}F_s + \varepsilon_1 \\ X_2 &= l_{21}F_1 + l_{22}F_2 + \dots + l_{2s}F_s + \varepsilon_2 \\ &\vdots \\ X_r &= l_{r1}F_1 + l_{r2}F_2 + \dots + l_{rs}F_s + \varepsilon_r \end{aligned}$$

或者写为矩阵形式

$$\mathbf{X}_{(r \times 1)} = \mathbf{L}_{(r \times s)} \mathbf{F}_{(s \times 1)} + \boldsymbol{\varepsilon}_{(r \times 1)}$$

设 $(\hat{\lambda}_k, \hat{\mathbf{e}}_k), k=1, 2, \dots, r$ 为样本协方差矩阵 \mathbf{S} 的特征值 - 特征向量对(默认 $\hat{\mathbf{e}}_k$ 为单位向量且 $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_r \geq 0$, 下文不再提及), 则

载荷矩阵 \mathbf{L} 的估计值为

$$\hat{\mathbf{L}}_{(r \times s)} = \left[\sqrt{\hat{\lambda}_1} \hat{\mathbf{e}}_1 \mid \sqrt{\hat{\lambda}_2} \hat{\mathbf{e}}_2 \mid \dots \mid \sqrt{\hat{\lambda}_s} \hat{\mathbf{e}}_s \right] = \{\hat{l}_{ij}\}_{r \times s}$$

矩阵 \mathbf{A} 的近似估计

$$\mathbf{B}_{\text{FA}} = \mathbf{A}\mathbf{S}^{-1} \hat{\mathbf{L}} \hat{\mathbf{L}}^T \quad (6)$$

均方误差

$$MSE_{FA} = \frac{1}{(n-1)r} \text{tr}[(A-B)(A-B)^T] = \frac{1}{r} \sum_{i=1}^r \hat{\psi}_i$$

其中 $\hat{\psi}_i = s_{ii}^2 - \hat{l}_i^2 - \dots - \hat{l}_{is}^2 = \text{Var}(\hat{\epsilon}_i)$ 。

(这里载荷矩阵的估计使用的是主成分法, 因子得分使用的是回归法。)

2.4. 主成分回归模拟

对于 $X_j (j=1, 2, \dots, p)$, 通过样本协方差矩阵 \mathbf{S} 提取其前 s 个主成分的估计值

$$\hat{\mathbf{Y}}_k = [\hat{y}_{k1}, \hat{y}_{k2}, \dots, \hat{y}_{kn}]^T = \mathbf{A}\hat{\mathbf{e}}_k, \quad k=1, 2, \dots, s$$

其中 $(\hat{\lambda}_k, \hat{\mathbf{e}}_k)$ 为 \mathbf{S} 的特征值-特征向量对, 于是有回归函数

$$\mathbf{x}_j = \beta_{j1}\hat{\mathbf{Y}}_1 + \beta_{j2}\hat{\mathbf{Y}}_2 + \dots + \beta_{js}\hat{\mathbf{Y}}_s + \boldsymbol{\epsilon}_j, \quad j=1, 2, \dots, r$$

(这里令常数项为 0。)

由多元线性回归结果为:

$$\hat{\boldsymbol{\beta}}_j = [\beta_{j1}, \beta_{j2}, \dots, \beta_{js}]^T = (\hat{\mathbf{Y}}^T \hat{\mathbf{Y}})^{-1} \hat{\mathbf{Y}}^T \mathbf{x}_j$$

其中

$$\hat{\mathbf{Y}}_{(n \times s)} = [\hat{\mathbf{Y}}_1 | \hat{\mathbf{Y}}_2 | \dots | \hat{\mathbf{Y}}_s]$$

于是矩阵 \mathbf{A} 有近似估计

$$\mathbf{B}_{\text{PCR}} = \hat{\mathbf{Y}}(\hat{\mathbf{Y}}^T \hat{\mathbf{Y}})^{-1} \hat{\mathbf{Y}}^T \mathbf{A} \quad (7)$$

均方误差

$$MSE_{\text{PCR}} = \frac{1}{(n-1)r} \text{tr}[(A-B)(A-B)^T]$$

3. 三种模拟方法的一致性

引理 3.1 设 \mathbf{A} 是 $n \times r (n > r)$ 阶实矩阵, $s < r = \text{rank}(\mathbf{A})$, 并且有奇异值分解 $\mathbf{U}\mathbf{M}\mathbf{V}^T$, 具体形式见式(2)(3)(4), 则

$$\mathbf{B}^* = \sum_{i=1}^s \mu_i \mathbf{u}_i \mathbf{v}_i^T$$

是 \mathbf{A} 的秩- s 最小二乘逼近, 使得在所有秩小于等于 s 的 $n \times r$ 阶矩阵 \mathbf{B} 中, 平方误差和

$\text{tr}[(A-B)(A-B)^T]$ 最小, 且最小值为 $\sum_{i=s+1}^r \mu_i^2$ (见文献[6])。

#

引理 3.2 设 \mathbf{A} 是 $n \times r (n > r)$ 阶实矩阵, $r = \text{rank}(\mathbf{A})$, 并且有奇异值分解 $\mathbf{U}\mathbf{M}\mathbf{V}^T$, 具体形式见式(2)(3)(4), 则

$$\mathbf{A}^T \mathbf{A} \mathbf{v}_i = \mu_i^2 \mathbf{v}_i, \quad i=1, 2, \dots, r$$

即 $A^T A$ 有特征值-特征向量对 (μ_i^2, \mathbf{v}_i) (见文献[6])。

#

利用之前的三个模型结果和引理 3.1、引理 3.2, 可以证明下面的定理。

定理 3.1 对于在 2.2、2.3 和 2.4 中三种使用同阶的低维度矩阵 B ($\text{rank}(B) = s$) 来模拟原数据矩阵 A ($\text{rank}(A) = r, r > s$) 的方法中(见式(5) (6) (7), 且矩阵 A 满足 2.1.中的假设条件), 并且都使用样本协方差矩阵 S 进行操作时, 三种方法的模拟结果相同, 即模拟矩阵

$$B_{FA} = B_{PCR} = B_{SVD} = AV_s V_s^T$$

其中

$$V_s = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_s] = [\hat{\mathbf{e}}_1, \hat{\mathbf{e}}_2, \dots, \hat{\mathbf{e}}_s]$$

且均方误差

$$MSE_{SVD} = MSE_{FA} = MSE_{PCR} = \frac{1}{r} \sum_{i=s+1}^r \hat{\lambda}_i = \frac{1}{(n-1)r} \sum_{i=s+1}^r \mu_i^2$$

达到最小值。

($(\hat{\lambda}_i, \hat{\mathbf{e}}_i)$ 为矩阵 S 的特征值-特征向量对, \mathbf{v}_i 的定义见式(3), μ_i 的定义见式(4)。)

#

证明: 根据 $\frac{1}{n} \sum_{i=1}^n x_{ij} = 0 (j=1, 2, \dots, r)$ 有

$$\frac{1}{n-1} A^T A = S \quad (8)$$

于是由引理 3.2 有

$$\hat{\lambda}_i = \frac{\mu_i^2}{n-1}, \quad i=1, 2, \dots, r \quad (9)$$

$$\mathbf{v}_i = \hat{\mathbf{e}}_i, \quad i=1, 2, \dots, r \quad (10)$$

(使(10)式成立有时需要做一定的调整, 这里我们不多做考虑。)

令

$$\hat{\Lambda}_{(s \times s)} = \begin{bmatrix} \hat{\lambda}_1 & 0 & \dots & 0 \\ 0 & \hat{\lambda}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \hat{\lambda}_s \end{bmatrix}$$

对于特征值分解方法, 有

$$B_{SVD} = \sum_{i=1}^s \mu_i \mathbf{u}_i \mathbf{v}_i' = U M V^T V_s V_s^T = A V_s V_s^T \quad (11)$$

对于因子分解方法, 由特征值和特征向量的定义, 有

$$S V_s = V_s \hat{\Lambda} \quad (12)$$

于是

$$\mathbf{B}_{\text{FA}} = \mathbf{A}\mathbf{S}^{-1}\hat{\mathbf{L}}\hat{\mathbf{L}}^{\text{T}} = \mathbf{A}\mathbf{S}^{-1}\mathbf{V}_s\hat{\mathbf{\Lambda}}^{1/2}\left(\mathbf{V}_s\hat{\mathbf{\Lambda}}^{1/2}\right)^{\text{T}} = \mathbf{A}\mathbf{S}^{-1}\mathbf{V}_s\hat{\mathbf{\Lambda}}\mathbf{V}_s^{\text{T}} = \mathbf{A}\mathbf{V}_s\mathbf{V}_s^{\text{T}} \quad (13)$$

对于主成分回归方法, 有

$$\hat{\mathbf{Y}} = \mathbf{A}\mathbf{V}_s$$

再使用式(8), (12)得到

$$\mathbf{B}_{\text{PCR}} = \hat{\mathbf{Y}}\left(\hat{\mathbf{Y}}^{\text{T}}\hat{\mathbf{Y}}\right)^{-1}\hat{\mathbf{Y}}^{\text{T}}\mathbf{A} = \mathbf{A}\mathbf{V}_s\left(\mathbf{V}_s^{\text{T}}\mathbf{S}\mathbf{V}_s\right)^{-1}\mathbf{V}_s^{\text{T}}\mathbf{S} = \mathbf{A}\mathbf{V}_s\mathbf{V}_s^{\text{T}} \quad (14)$$

综合式(11), (13)和(14)得到

$$\mathbf{B}_{\text{FA}} = \mathbf{B}_{\text{PCR}} = \mathbf{B}_{\text{SVD}} = \mathbf{A}\mathbf{V}_s\mathbf{V}_s^{\text{T}}$$

于是

$$MSE_{\text{SVD}} = MSE_{\text{FA}} = MSE_{\text{PCR}}$$

最后, 根据引理 3.1 和式(9)得到均方误差

$$MSE_{\text{SVD}} = MSE_{\text{FA}} = MSE_{\text{PCR}} = \frac{1}{p} \sum_{i=s+1}^p \hat{\lambda}_i = \frac{1}{(n-1)p} \sum_{i=s+1}^p \mu_i^2$$

达到最小值。

#

4. 结论

当由多个时间序列构成的数据矩阵满足对应时间序列的期望为零, 且特征值和特征向量均由对应的样本协方差矩阵提取时, 奇异值分解、因子分析和主成分回归构造的降维模拟方法具有一致性(这里的一致性仅限于上文提到的构造方法)。其中, 模拟矩阵的结果仅依赖于所提取的特征向量(或奇异值分解的其中一个正交矩阵), 模拟矩阵均方误差的结果由所提取的特征值(或奇异值)完全决定。

参考文献

- [1] 罗小桂. 矩阵奇异值分解(SVD)的应用[J]. 井冈山医学学报, 2005, 12(4): 133-135.
- [2] 范龙振, 余世典. 中国股票市场的三因子模型[J]. 系统工程学报, 2002, 17(6): 537-546.
- [3] 任福匀. 因子分析法在我国股票市场行业投资价值评价中的应用[D]: [硕士学位论文]. 长沙: 中南大学, 2005.
- [4] Zhang, W. (2011) APT Model Based on Factor Analysis and an Empirical Study in China's Growth Enterprise Market. 2011 2nd International Conference on Artificial Intelligence, Management Science and Electronic Commerce (AIMSEC), Dengleng, 8-10 August 2011, 994-997. <https://doi.org/10.1109/AIMSEC.2011.6010620>
- [5] 张亚梅. 基于主成分回归分析科技创新对金融业的影响——以甘肃省为例[J]. 甘肃科技纵横, 2018, 47(9): 81-85.
- [6] 理查德·A.约翰逊, 迪安·W.威克恩. 实用多元统计分析[M]. 第 6 版. 陆璇, 叶俊, 译. 北京: 清华大学出版社, 2008: 76-78.

知网检索的两种方式：

1. 打开知网页面 <http://kns.cnki.net/kns/brief/result.aspx?dbPrefix=WWJD>
下拉列表框选择：[ISSN]，输入期刊 ISSN：2325-2251，即可查询
2. 打开知网首页 <http://cnki.net/>
左侧“国际文献总库”进入，输入文章标题，即可查询

投稿请点击：<http://www.hanspub.org/Submission.aspx>

期刊邮箱：sa@hanspub.org