

# Statistical Theory and Related Fields

## Factor Analysis of Correlation Matrices When the Number of Random Variables Exceeds the Sample Size --Manuscript Draft--

<b>Full Title:</b>	Factor Analysis of Correlation Matrices When the Number of Random Variables Exceeds the Sample Size
<b>Manuscript Number:</b>	TSTF-2017-0025R1
<b>Article Type:</b>	Original Article
<b>Keywords:</b>	Alpha spending function; BIC; Cancer surveillance; Eigenvalues; Sparse factor loadings
<b>Abstract:</b>	<p>Factor analysis which studies correlation matrices is an effective means of data reduction whose inference on the correlation matrix typically requires the number of random variables, <math>p</math>, to be relatively small and the sample size, <math>n</math>, to be approaching infinity. In contemporary data collection for biomedical studies, disease surveillance, and genetics, <math>p &gt; n</math> limits the use of existing factor analysis methods to study the correlation matrix. The motivation for the research here comes from studying the correlation matrix of log annual cancer mortality rate change for <math>p=59</math> cancer types from 1969 to 2008 (<math>n=39</math>) in the US. We formalize a test statistic to perform inference on the structure of the correlation matrix when <math>p &gt; n</math>. We develop an approach based on group sequential theory to estimate the number of relevant factors to be extracted. To facilitate interpretation of the extracted factors, we propose a BIC-type criterion to produce a sparse factor loading representation. The proposed methodology outperforms competing ad-hoc methodologies in simulation analyses, and identifies three significant underlying factors responsible for the observed correlation between cancer mortality rate changes.</p>
<b>Order of Authors:</b>	Miguel Marino, Ph.D. Yi Li, Ph.D.
<b>Response to Reviewers:</b>	<p>Prof. Shao Editor-in-Chief Statistical Theory and Related Fields</p> <p>Dear Prof.Shao:</p> <p>We are pleased to resubmit an edited version of our manuscript entitled 'Factor Analysis of Correlation Matrices When the Number of Random Variables Exceeds the Sample Size' which was originally submitted as an Original Article. We thank the editors and reviewers for their thoughtful and constructive feedback, which contributed greatly to many significant improvements in this manuscript. Thank you for your consideration of this very timely piece. We respond to each comment below.</p> <p>COMMENT 1</p> <p>The motivation to develop structure-inferring methods for correlation matrix need to be explained. There already are similar methods developed for covariance matrix. Since usually in practice research obtain covariance matrix first and then obtain correlation matrix. What is the use of the proposed approach if structure can be readily learned from covariance matrix? Is there additional insight from structure learned from correlation matrix?</p> <p>Response: We thank the reviewers for this great point. In some instances, starting with the correlation matrix instead of the covariance matrix is scientifically just. In our proposed data analysis and in the field of cancer mortality change pattern trends, the use of the correlation matrix helps standardize the different types of cancers which will not let factor analysis be driven by cancers with large variances as covariance matrices do. We have included language in the introduction to help motivate this further. Specifically, we added: "In many studies where the random variables of interest are</p>

highly variable (e.g. cancer mortality rates), it is common to standardize the random variables and analyze the correlation matrix. Standardization ensures that results from factor analysis will not be driven by random variables with large variances, which is a challenge when performing factor analysis on covariance matrices." We also added a section in the discussion that calls out further research to identify the similarities and differences between starting with the covariance or correlation matrix. See below.

#### COMMENT 2

In what sense are Models (1) and (2) equivalent? Is  $LF = \sum_{l=1}^m \sqrt{\lambda_l} F_{le_l}$ ?

Response: The reviewer is correct. Using an eigenvalue decomposition,  $L = \lambda_1 e_1 e_1^T + \dots + \lambda_m e_m e_m^T$  with  $m$  orthonormal eigenvectors  $e_l$  for  $l=1, \dots, m$  such that  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$  and  $e_l^T e_k = \delta_{lk}$ , which equals 1 if  $l=k$  and 0 otherwise. Isolating  $L$  and plugging that in model (1) results in  $LF = \sum_{l=1}^m \sqrt{\lambda_l} F_{le_l}$  that is reported in model (2). We have updated the manuscript to make this more clear.

#### COMMENT 3

Page 3, Line 19, provide reference for the theoretical distribution for the largest eigenvalue.

Response: We thank the reviewer for this comment. We have included the appropriate reference.

#### COMMENT 4

Page 3, Line 34, change M to X to make the notations consistent.

Response: The have updated the proposed change.

#### COMMENT 5

Page 5, Line 24, the Tracy-Widom test was used for the largest eigenvalue. How was it used to identify  $\lambda_1, \dots, \lambda_{k-1}$ ?

Response: Given that the Tracy-Widom test can only be used for the largest eigenvalue, we needed to develop a sequential method that removes the information from the first factor and produces new data matrix  $X^{(k)}$  from which we can construct a new correlation matrix and estimate a new 'largest eigenvalue' and test again using the Tracy-Widom test. This approach helps identify the number of factors that are relevant and as such, the  $\lambda_1, \dots, \lambda_{k-1}$ . We have updated the manuscript to clarify this. Specifically, we moved the previous paragraph before section 3.1 into section 3.1 and added the following sentence: "Given that the Tracy-Widom test is used to test the largest eigenvalue, we propose a sequential method that removes the effect of the first factor (if significant) and produces a new data matrix from which we can construct a new correlation matrix and apply once again the Tracy-Widom test on the new largest eigenvalue."

#### COMMENT 6

In data analysis, it will be interesting to compare the results obtained on covariance matrix and those obtained on correlation matrix.

Response: We thank the reviewer for this suggestion. The proposed comparison will answer: 'How do the covariance and correlation matrix results differ?'. We believe that this is important but it is peripheral to the question of interest: 'How do you do factor analysis on correlation matrices?'. Given that this question is outside the scope of this

manuscript, the limited word limit constraints and the short time window for a manuscript revision, we did not produce this analysis. We believe that this is an important question, but it would require more time and space to do it justice. If the editors wish, we would be happy to produce this. This is a great suggestion and we encourage future research to be done to answer the question. We have added a section in the discussion that calls out further research to identify the similarities and differences between starting with the covariance or correlation matrix. Specifically we included: "Lastly, although our focus was on the study of the correlation matrix when  $p > n$ , future studies should compare the performance of the Tracy-Widom test and sequential-rescaling procedure on covariance matrices to compare the performance and generalizability of these methods."

1  
2  
3  
4  
5 ARTICLE

6  
7 **Factor Analysis of Correlation Matrices When the Number of**  
8 **Random Variables Exceeds the Sample Size**  
9

10 Miguel Marino<sup>a</sup> and Yi Li<sup>b</sup>

11  
12 <sup>a</sup>Oregon Health & Science University, 3181 S.W. Sam Jackson Park Rd., Portland, OR,  
13 USA; <sup>b</sup> University of Michigan, 1415 Washington Heights, Ann Arbor, Michigan, USA  
14

15  
16 **ARTICLE HISTORY**

17 Compiled October 27, 2017  
18

19 **ABSTRACT**

20 Factor analysis which studies correlation matrices is an effective means of data re-  
21 duction whose inference on the correlation matrix typically requires the number of  
22 random variables,  $p$ , to be relatively small and the sample size,  $n$ , to be approaching  
23 infinity. In contemporary data collection for biomedical studies, disease surveillance,  
24 and genetics,  $p > n$  limits the use of existing factor analysis methods to study the  
25 correlation matrix. The motivation for the research here comes from studying the  
26 correlation matrix of log annual cancer mortality rate change for  $p = 59$  cancer  
27 types from 1969 to 2008 ( $n = 39$ ) in the US. We formalize a test statistic to perform  
28 inference on the structure of the correlation matrix when  $p > n$ . We develop an ap-  
29 proach based on group sequential theory to estimate the number of relevant factors  
30 to be extracted. To facilitate interpretation of the extracted factors, we propose a  
31 BIC-type criterion to produce a sparse factor loading representation. The proposed  
32 methodology outperforms competing *ad hoc* methodologies in simulation analyses,  
33 and identifies three significant underlying factors responsible for the observed cor-  
34 relation between cancer mortality rate changes.

35 **KEYWORDS**

36 Alpha spending function; BIC; Cancer surveillance; Eigenvalues; Sparse factor  
37 loadings  
38

39 **1. Introduction**  
40

41 Due to its flexibility in characterizing multivariate data, high-dimensional factor analysis  
42 is becoming popular in many scientific disciplines including genetic (Zhou, Wang,  
43 Wang, Zhu, & Song, 2017), biomedical (Shimizu et al., 2016) and economic studies  
44 (Fan, Lv, & Qi, 2011). The objectives of exploratory factor analysis are two-fold: 1)  
45 identify the number of factors that influence a set of random variables; 2) measure the  
46 strength of the relationship between the extracted factors and each random variable.  
47

48 In many studies where the random variables of interest are highly variable (e.g.  
49 cancer mortality rates), it is common to standardize the random variables and analyze  
50 the correlation matrix. Standardization ensures that results from factor analysis will  
51 not be driven by random variables with large variances, which is a challenge when per-  
52 forming factor analysis on covariance matrices. Additionally, in the cases where the  
53 number of random variables exceeds the sample size, a couple statistical challenges  
54  
55

---

56 CONTACT Miguel Marino. Email: marinom@ohsu.edu  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 arise in the analysis of correlation matrices via factor models. First, existing inference  
2 methods rely on the number of random variables,  $p$ , to be relatively small and fixed,  
3 and the sample size,  $n$ , to be approaching infinity (Anderson, 1963; Johnson & Wich-  
4 ern, 1998). Another complication is that factor analysis is not invariant to change on  
5 the scale of variables. Methods that infer structure from covariance matrices (Bickel  
6 & Levina, 2008; Carvalho et al., 2008; Fan, Fan, & Lv, 2008; Ghosh & Dunson, 2008;  
7 Huang, Liu, Pourahmadi, & Liu, 2006; Patterson, Price, & Reich, 2006; West, 2003;  
8 Wong, Carter, & Kohn, 2003) will not always perform similarly on correlation matrices.  
9 There appears to be a lack of methodology for performing inference on correlation  
10 matrices using factor analysis when  $p > n$ . Furthermore, traditional methods for esti-  
11 mating the number of factors to be extracted and their interpretation are insufficient  
12 and need further development for correlation matrices when  $p > n$ .

13 We make several contributions with this paper. First, we formalize a test statistic to  
14 perform inference of the structure of the correlation matrix using the limiting distri-  
15 bution of eigenvalues. This test statistic from Johnstone (2001) but was not delineated  
16 as fully as we do in this paper. Secondly, we extend the work of Johnstone (2001) to  
17 identify the true number of underlying factors present in a factor model, while con-  
18 trolling the type I error. Finally, we propose a BIC-type criterion to produce sparse  
19 factor loadings to ease interpretation of extracted factors.

20 The format of this paper is as follows. In Section 2, we present a test for inference on  
21 the structure of the population correlation matrix, which we term the Tracy-Widom  
22 test. In Section 3, we develop a sequential-rescaling procedure to test for the number  
23 of significant factors in a given factor model. Section 4 describes a sparse factor model  
24 that aids in interpreting the factors detected from the proposed test. Section 5 presents  
25 some designed simulation studies based on the proposed methodology. Section 6 applies  
26 the developed methodology to study the correlation matrix of cancer mortality annual  
27 rate changes data, followed by our concluding remarks in Section 7.

## 2. Methods

### 2.1. Factor model formulation

28 Consider a random vector  $\mathbf{X} = (X_1, \dots, X_p)^T$  where each component  $X_j$  follows a stan-  
29 dard normal distribution. Because different cancers have varying degrees of volatility,  
30 normalization will ensure that the analysis will not be dominated by a few cancer  
31 types. The primary aim of this project is to study the correlation matrix of  $\mathbf{X}$ .

32 A factor model postulates that  $\mathbf{X}$  is linearly dependent on a few underlying, but  
33 unobservable, random quantities  $F_1, \dots, F_m$  called common factors and  $p$  additional  
34 sources of variation  $\epsilon_1, \dots, \epsilon_p$  called white noise or specific factors, such that

$$35 \mathbf{X} = \mathbf{L}\mathbf{F} + \boldsymbol{\epsilon} \tag{1}$$

36 where  $\mathbf{F} = (F_1, \dots, F_m)^T \sim MVN(\mathbf{0}, \mathbf{I}_m)$  is a vector of  $m$  common factors,  $\mathbf{L} =$   
37  $(\ell_1, \dots, \ell_m)$  is a  $p \times m$  matrix of factor loadings with  $\ell_l = (\ell_{l1}, \dots, \ell_{lp})^T$  for  $l =$   
38  $1, \dots, m$  and  $\mathbf{I}_m$  is an identity matrix of dimension  $m$ . We denote the residual as  
39  $\boldsymbol{\epsilon} \sim N(0, \boldsymbol{\Psi})$  where  $\boldsymbol{\Psi}$  is a  $p \times p$  diagonal matrix with the  $l$ th diagonal element being  
40  $\psi_l = 1 - \ell_{1l}^2 - \dots - \ell_{ml}^2$  to ensure that  $Var(X_j) = 1$ .

41 If we assume that  $\mathbf{F}$  and  $\boldsymbol{\epsilon}$  are independent in (1), then it follows that the correlation  
42 matrix for  $\mathbf{X}$  is  $\mathbf{R} = \mathbf{L}\mathbf{L}^T + \boldsymbol{\Psi}$ . Using an eigenvalue decomposition,  $\mathbf{L}\mathbf{L}^T = \lambda_1 \mathbf{e}_1 \mathbf{e}_1^T +$

1  $\dots + \lambda_m \mathbf{e}_m \mathbf{e}_m^T$  with  $m$  orthonormal eigenvectors  $\mathbf{e}_l$  for  $l = 1, \dots, m$  such that  $\lambda_1 \geq$   
2  $\lambda_2 \geq \dots \geq \lambda_m \geq 0$  and  $\mathbf{e}_l^T \mathbf{e}_k = \delta_{lk}$ , which equals 1 if  $l = k$  and 0 otherwise. Hence,  
3  $\mathbf{L}\mathbf{F} = \sum_{l=1}^m \sqrt{\lambda_l} \mathbf{F}_l \mathbf{e}_l$  and results in

$$4 \quad 5 \quad 6 \quad 7 \quad 8 \quad \mathbf{X} = \sum_{l=1}^m \sqrt{\lambda_l} \mathbf{F}_l \mathbf{e}_l + \boldsymbol{\epsilon} \quad (2)$$

9 where  $\lambda_1, \lambda_2, \dots, \lambda_m$  correspond to the  $m$  largest eigenvalues of  $\mathbf{R}$ .

## 12 2.2. Testing Complete Independence of the Correlation Matrix

13 One of the first objectives of studying the correlation matrix of a set of random vari-  
14 ables is to determine if factor analysis is a reasonable method of analysis. This is  
15 equivalent to performing inference on the structure of the correlation matrix with test  
16 of  $H_0 : \mathbf{R} = \mathbf{I}$  versus the alternative  $H_a : \mathbf{R} \neq \mathbf{I}$ . We base our test for  $H_0 : \mathbf{R} = \mathbf{I}$   
17 on the largest eigenvalue of the sample correlation matrix of  $\mathbf{X}$ . A result of random  
18 matrix theory (RMT) suggests that we can build a theoretical distribution for the  
19 largest eigenvalue of random matrices under the null hypothesis of complete inde-  
20 pendence (Johnstone, 2001). A test of complete independence about the  $p$  random  
21 variables compares the observed sample eigenvalue  $\hat{\lambda}_1$  to the theoretical distribution  
22 of  $\lambda_1$  under RMT prediction. This test will reveal one of two possibilities: the first be-  
23 ing that  $\hat{\lambda}_1$  will be determined to not significantly differentiate from RMT prediction.  
24 This suggests that  $H_0 : \mathbf{R} = \mathbf{I}$  cannot be rejected and therefore that factor analysis  
25 will not prove to be useful because specific noise factors play a more dominant role in  
26 the observed correlation than common underlying factors. The second possibility for a  
27 test of the largest eigenvalue is that it will determine  $\hat{\lambda}_1$  to significantly deviate from  
28 RMT prediction (i.e.  $H_0 : \mathbf{R} = \mathbf{I}$  is rejected in favor of the alternative). This scenario  
29 suggests that one (or possibly more) underlying factor(s) could be responsible for the  
30 observed correlation between the random variables.

31 To proceed, we describe the test statistic for testing  $H_0 : \mathbf{R} = \mathbf{I}$ . Suppose that data  
32 matrix  $\mathbf{X} = (X_{ij})_{n \times p}$  has entries that are independent and identically distributed as  
33 standard normal. Let  $\hat{\xi}_1 \geq \hat{\xi}_2 \geq \dots \geq \hat{\xi}_p$  denote the sample eigenvalues of a Wishart  
34 Matrix,  $\mathbf{X}^T \mathbf{X}$ . We can test the significance of  $\hat{\xi}_1$ , the largest eigenvalue of  $\mathbf{X}^T \mathbf{X}$ , with  
35 test statistic

$$36 \quad 37 \quad 38 \quad 39 \quad 40 \quad 41 \quad 42 \quad 43 \quad 44 \quad T_{np} = \frac{(\hat{\xi}_1 - \mu_{np})}{\sigma_{np}} \quad (3)$$

45 where

$$46 \quad 47 \quad 48 \quad 49 \quad \mu_{np} = \begin{cases} (\sqrt{n-1} + \sqrt{p})^2, & \text{when } n \geq p \\ (\sqrt{p-1} + \sqrt{n})^2, & \text{when } p > n \end{cases}$$

50 and

$$51 \quad 52 \quad 53 \quad 54 \quad 55 \quad \sigma_{np} = \begin{cases} (\sqrt{n-1} + \sqrt{p}) \left( \frac{1}{\sqrt{n-1}} + \frac{1}{\sqrt{p}} \right)^{1/3}, & \text{when } n \geq p \\ (\sqrt{p-1} + \sqrt{n}) \left( \frac{1}{\sqrt{p-1}} + \frac{1}{\sqrt{n}} \right)^{1/3}, & \text{when } p > n. \end{cases}$$

56 Johnstone (2001) has shown that under  $H_0$ , and  $n, p \rightarrow \infty$  such that  $n/p \rightarrow \gamma$  for  $\gamma$

1 some constant and the test statistic  $T_{np} \xrightarrow{d} W_1$ , where  $W_1$  is called the Tracy-Widom  
2 distribution (Tracy & Widom, 2000). We term (3) the Tracy-Widom test and will  
3 reject the null hypothesis of  $H_0 : \mathbf{R} = \mathbf{I}$  when  $T_{np} > W_{1,1-\alpha}$  where  $W_{1,1-\alpha}$  is the  
4  $(1 - \alpha) \times 100$  percentile of the Tracy-Widom distribution. One of the strengths of  
5 this test is that it can be applied in the classical setting where  $n > p$  as well as in  
6 high-dimensional settings where  $p > n$ .  
7  
8

### 9 **2.3. Correlation correction of Tracy-Widom test**

11 A technical note suggests that the Tracy-Widom test applies to the study of covariance  
12 matrices and does not directly apply to correlation matrices, which is problematic for  
13 distribution theory (Anderson, 1963). To be able to apply the Tracy-Widom test to  
14 study correlation matrices, we expand on the procedure that was briefly mentioned  
15 in Johnstone (2001) but has not been fully studied. To this end, suppose we draw  
16  $n$  *i.i.d.* row vector samples from  $\mathbf{N}(\mathbf{0}, \mathbf{\Sigma})$  to produce data matrix  $\mathbf{X}_{n \times p}$ . Under the  
17 null hypothesis, the column vectors  $\mathbf{X}_j$  are *i.i.d.* on the unit sphere  $S^{n-1}$ . As a result,  
18 we can multiply each  $\mathbf{X}_j$  by an independent chi-distributed length to synthesize a  
19 Gaussian matrix, call it  $\tilde{\mathbf{X}}$  such that  $\tilde{\mathbf{X}} = [\tilde{\mathbf{X}}_1 \dots \tilde{\mathbf{X}}_p]$  where  $\tilde{\mathbf{X}}_j = \psi_j \mathbf{X}_j$  and  $\psi_j^2 \sim$   
20  $\chi_{(n-1)}^2$ . We can then construct a sample pseudo-covariance matrix  $\tilde{\mathbf{S}} = \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  which  
21 approximately follows a Wishart distribution with  $n - 1$  degrees of freedom. Under the  
22 null, this data augmentation allows us to apply the Tracy-Widom test on the largest  
23 eigenvalue of  $\tilde{\mathbf{S}}$  to test  $H_0 : \mathbf{R} = \mathbf{I}$ .  
24  
25  
26  
27  
28

### 29 **3. Identifying Additional Factors**

31 If  $H_0 : \mathbf{R} = \mathbf{I}$  is rejected, then at least one latent factor is useful in describing the  
32 observed correlation among the  $p$  random variables. One of the most crucial steps of  
33 factor analysis is to estimate the true number of underlying factors,  $m$ , as misspec-  
34 ification of the number of factors retained can lead to poor factor-loading pattern  
35 reproduction and interpretation (Hayton, Allen, & Scarpello, 2004). Furthermore, es-  
36 timation of the number of factors can affect the factor model results more than other  
37 decisions, such as the factor rotation method used (Zwick & Velicer, 1986). In this  
38 section, we extend the work of Johnstone (2001) to identify the number of relevant  
39 factors to be used in a factor model.  
40  
41

42 Previous work on estimating the number of factors have focused on factor analy-  
43 sis for covariance matrices (Bai, 2003; Bai & Ng, 2002; Leek, 2011; Onatski, 2009).  
44 Johnstone (2001) and Baik and Silverstein (2006) have considered the asymptotic be-  
45 havior of  $\hat{\xi}_{r+1}$ , the  $(r + 1) - th$  largest eigenvalue of a covariance matrix, when the  
46 true population covariance follows a spiked model with  $\mathbf{\Sigma} = \text{diag}(\tau_1, \dots, \tau_r, 1, \dots, 1)$ ,  
47 where  $\tau_1 \geq \dots \geq \tau_r > 1$ . As factor analysis is not invariant to changes in the scale of  
48 the variables, it is often recommended that factor analysis be performed for standard-  
49 ized variables. Standardization converts a covariance matrix problem into a correlation  
50 problem and it is unclear how these methods would be applied to the study of sample  
51 correlation matrices.  
52

53 Common *ad hoc* methods of determining the number factors to extract from corre-  
54 lation matrices include the scree plots, Guttman-Kaiser criterion and parallel analysis.  
55 The number of extracted factors based on the scree plot are highly subjective as the  
56 estimate is visually selected as point that resembles an elbow. The Guttman-Kaiser  
57  
58

criterion (Guttman, 1954; Kaiser, 1960) selects the number of factors to be equal to the number of sample eigenvalues of the correlation matrix that are greater than one. Parallel analysis (Horn, 1965) is a simulation-based approach that compares the eigenvalues of the sample correlation matrix to eigenvalues from a matrix of random values of the same dimensionality. The estimated number of factors retained are the number of observed sample eigenvalues greater than the 95th percentile of the distribution of eigenvalues derived from the random data.

### 3.1. Sequential-Rescaling Testing Procedure

We propose to view the testing procedure of extracting relevant underlying factors as a sequential procedure. Given that the Tracy-Widom test is used to test the largest eigenvalue, we propose a sequential method that removes the effect of the first factor (if significant) and produces a new data matrix from which we can construct a new correlation matrix and apply once again the Tracy-Widom test on the new largest eigenvalue. In general, we will test for the significance of  $\lambda_k$  only after verifying that  $\lambda_{k-1}$  are significantly different than RMT prediction and after eliminating the effect of the first  $k - 1$  factors. We remove the effect of the first  $k - 1$  factors because of the phenomenon where the largest eigenvalue has the potential to pull other sample eigenvalues away from unity. The resulting procedure is termed a sequential-rescaling procedure. The advantage of the procedure that follows is that it controls the type I error through the use of an alpha spending function, and it is not a conservative technique based on what has been proposed in Patterson et al. (2006).

Suppose we have declared the first  $k - 1$  eigenvalues to be significantly different than RMT prediction. The following procedure tests the subsequent eigenvalue  $\lambda_k$ . The procedure assumes the Tracy-Widom test has already identified  $\lambda_1, \lambda_2, \dots, \lambda_{k-1}$  to be significant. Associated with eigenvalue  $\lambda_\ell$  is its corresponding eigenvector  $\mathbf{e}_\ell = (e_{\ell 1}, \dots, e_{\ell p})$ . We proceed to test  $\lambda_k$  through the following two-step procedure:

**Step 1.** Construct a data matrix  $\mathbf{X}^{(k)}$  such that

$$\mathbf{X}^{(k)} = \mathbf{D}_{k-1}^{-1/2}(\mathbf{X}^{(k-1)} - \sqrt{\lambda_{k-1}}\mathbf{e}_{k-1}F_{k-1}) \quad (4)$$

where  $\mathbf{D}_{k-1}$  is the rescaling diagonal matrix with its  $i$ -th diagonal element being  $D_{k-1,ii} = 1 - \lambda_{k-1}e_{k-1,i}^2$ . The rescaling matrix,  $\mathbf{D}_{k-1}^{-1/2}$  will assure the desirable property that  $\text{var}(\mathbf{X}^{(k)}) = 1$ . Note that we have removed the effect of the first  $k - 1$  factors in (4) through the  $(\mathbf{X}^{(k-1)} - \sqrt{\lambda_{k-1}}\mathbf{e}_{k-1}F_{k-1})$  term.

**Step 2.** It can be shown that the sample correlation matrix for the rescaled  $\mathbf{X}^{(k)}$ , on which we will test the significance of  $\lambda_k$  using the Tracy-Widom test proposed in (3) is

$$\mathbf{R}_k \stackrel{def}{=} \mathbf{D}_{k-1}^{-1/2}(\mathbf{R}_{k-1} - \lambda_{k-1}\mathbf{e}_{k-1}\mathbf{e}'_{k-1})\mathbf{D}_{k-1}^{-1/2} \quad (5)$$

We perform this two-step procedure applying the Tracy-Widom test on each subsequent eigenvalue until an eigenvalue is no longer significant.

We note that caution should be taken when testing subsequent sample eigenvalues. To circumvent the multiple testing issues that are present in this procedure, we apply methodology from the group sequential analysis literature to control the type I error. Lan and DeMets (1983) proposed an alpha spending technique in which the nominal



significance level needed to reject the null hypothesis at each analysis is less than  $\alpha$  and increases as the study progresses. If an overall type I error ( $\alpha$ ) is desired, we propose to use the following alpha spending function

$$\alpha^*(k) = \alpha/2^k \tag{6}$$

where  $\alpha^*(k)$  is the significance level for the  $k$ th hypothesis test. This is opposed to Lan and DeMets (1983), as alpha spending function (6) does not depend on the overall number of tests being conducted. Therefore one need not specify the maximum number of eigenvalues being tested, which is ideal for unsupervised learning.

We define type I error as the probability of incorrectly choosing a model that has extracted more factors than the true model. Compared to Lan and DeMets (1983) who suggest that the alpha spending function should be nondecreasing, our spending function is nonincreasing ( $\alpha^*(1) > \alpha^*(2) > \dots > \alpha^*(K)$ ); because finding a parsimonious model is preferred, we need strong evidence for choosing a more complicated model with more significant eigenvalues over a simpler one. We have shown in supplementary material that the overall type I error rate using the proposed spending function (6) will not exceed  $\alpha$ .

#### 4. Interpretation of Factors

After an estimation is made for the number of factors to be used, the next objective in factor analysis is to provide an interpretation for each underlying factor. In principle, the factor loadings provide the basis for interpreting the factors underlying the data. The size and direction of the extracted factor loadings denote the strength and direction of the correlation between the random variables and the extracted factors. Traditionally, the task of interpreting factors has been subjective and unsatisfactory.

Because the original factor loadings may not be easily interpretable, it has become common to rotate the loadings (e.g. varimax, oblique, etc.) to increase or decrease the size of factor loadings to ease of interpretation. Unfortunately, regardless of the factor rotation used, it is rare for factor loadings to be set exactly to zero which would ease in the interpretation of the underlying factor.

With the recent developments of regularized regression in mind, we propose to implement a regularization technique to detect a set of sparse factor loadings for easier interpretation of identified factors. The resulting sparse factor loading vector sets the loadings of negligible random variables to zero, assuring that they will not contribute to the interpretation of the underlying factor, making the interpretation of the factors more straightforward. Additionally, because negligible random variables are removed, the variance explained by the sparse factor loadings will not suffer much from their removal.

Eigenvalue decomposition (2) provides the factoring of the correlation matrix of  $\mathbf{R}$ . The factor loading matrix  $\mathbf{L}$  is given by  $\mathbf{L} = (\sqrt{\lambda_1}\mathbf{e}_1, \dots, \sqrt{\lambda_m}\mathbf{e}_m)$  where  $(\lambda_l, \mathbf{e}_l)$  are the eigenvalue-eigenvector pairs of  $\mathbf{R}$ . Producing sparse factor loadings is equivalent to setting components of  $\mathbf{e}_l$  to zero. It can be shown that apart from the scale value  $\sqrt{\lambda_l}$ , the factor loading column  $\mathbf{e}_l$  are the coefficients of the principal components of the population. This observation allows us to implement well studied sparse principal components methods to produce sparse factor loadings.

We propose to regularize  $\mathbf{e}_l$  for  $l = 1, \dots, m$  using the sparse principal components analysis (SPCA) method proposed by Zou, Hastie, and Tibshirani (2006). SPCA es-

1 sentially takes the problem of setting PCA loadings to zero and transforms it into  
2 a regression-type problem that uses an elastic net regularization technique to detect  
3 sparse loadings even when  $p > n$ . Details of SPCA methodology can be found in Zou  
4 et al. (2006), but we provide a brief description below.

#### 7 4.1. SPCA for Sparse Factor Loadings

8 We consider the problem of producing sparse factor loadings for the  $m$  estimated  
9 factors. Let  $\mathbf{A}_{p \times m} = (\alpha_1, \dots, \alpha_m)$ ,  $\mathbf{B}_{p \times m} = (\beta_1, \dots, \beta_m)$  and  $\mathbf{X}$  be the  $n \times p$  data  
10 matrix as before and  $\mathbf{X}_i$  denote the  $i$ th row vector of  $\mathbf{X}$ . The problem of producing  
11 sparse factor loadings can be transformed into the following regression-type criterion  
12 with an elastic net penalty  
13

$$14 \quad (\hat{\mathbf{A}}, \hat{\mathbf{B}}) = \underset{\mathbf{A}, \mathbf{B}}{\operatorname{argmin}} \sum_{i=1}^n \|\mathbf{X}_i - \mathbf{A}\mathbf{B}^T \mathbf{X}_i\|^2 + \gamma \sum_{l=1}^m \|\beta_l\|^2 + \sum_{l=1}^m \gamma_{1,l} \|\beta_l\|_1 \quad (7)$$

15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

subject to  $\mathbf{A}^T \mathbf{A} = \mathbf{I}_{m \times m}$

for any  $\gamma > 0$ . The last term in (7) uses the  $L_1$  penalty to produce sparse factor loadings  
because the estimated sparse factor loadings, defined as  $\hat{\mathbf{e}}_l^s = \frac{\beta_l}{\|\beta_l\|}$  for  $l = 1, \dots, m$ ,  
are a function of the sparse  $\beta_l$  vector.

#### 4.2. Selection of Tuning Parameter

The optimization problem in (7) contains two tuning parameters that must be selected.  
The first tuning parameter,  $\gamma$ , is the same for all the  $m$  factors. It has been shown  
(Zou et al., 2006) that when  $p > n$ , a positive  $\gamma$  is required to produce exact loadings  
when the second tuning parameter is set to zero. The tuning parameter  $\gamma$  has been  
studied and is well understood. Empirical evidence has shown that for the case when  
 $n > p$ ,  $\gamma$  can be set to zero. When  $p > n$ ,  $\gamma$  can be set to a small positive number to  
overcome collinearity between the columns of  $\mathbf{X}$ .

The second tuning parameter  $\gamma_{1,l}$  is a factor-specific tuning parameter and requires  
more development. Zou et al. (2006) did not provide clear guidance on selecting  $\gamma_{1,j}$ ,  
other than choosing  $\gamma_{1,j}$  such that it provides a good compromise between explained  
variance and sparsity. Other methods exist for selecting the tuning parameters, such  
as cross validation (Shen & Huang, 2008) which could be computationally extensive  
and requires a large sample size. We add to the current literature on producing sparse  
factor loadings by proposing a BIC-type criterion for selecting the factor specific tuning  
parameters  $(\gamma_{1,1}, \dots, \gamma_{1,m})$ .

For a fixed  $\gamma$ , we propose to use the following BIC-type criterion for selection of  
tuning parameters  $(\gamma_{1,1}, \dots, \gamma_{1,m})$

$$54 \quad BIC = \log \left[ \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p (X_{ij} - \hat{\ell}_j \hat{F}_i)^2 \right] + df(\gamma_{1,l}, \hat{\mathbf{L}}) \frac{\log(np)}{np} \quad (8)$$

55 where  $\hat{\ell}_j = (\hat{\ell}_{1j}, \hat{\ell}_{2j}, \dots, \hat{\ell}_{mj}) = (\sqrt{\hat{\lambda}_1} \hat{\mathbf{e}}_{1j}^s, \dots, \sqrt{\hat{\lambda}_m} \hat{\mathbf{e}}_{mj}^s)$ ,  $\hat{\mathbf{L}} = (\ell_1, \dots, \ell_p)^T$  is the factor  
56 loading matrix and  $\hat{F}_i = (\hat{\lambda}_1^{-1/2} \hat{\mathbf{e}}_1^s \mathbf{X}_i, \dots, \hat{\lambda}_m^{-1/2} \hat{\mathbf{e}}_m^s \mathbf{X}_i)$  where  $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^T$ .

We define the degrees of freedom,  $df(\gamma_{1,l}, \hat{\mathbf{L}})$ , to be the number of nonzero loadings in the loading matrix  $\hat{\mathbf{L}}$ . Zou, Hastie, and Tibshirani (2007) showed that the number of nonzero coefficients in lasso regression provides an unbiased estimate for the degrees of freedom and suggests that BIC can be used to determine the optimal number of nonzero factor loadings.

## 5. Analysis of Simulated Data

To assess the performance of the proposed method, we simulate data from factor models where the  $p$  observable random variables are constructed from zero, one, two or three underlying factors. The zero factor model is given by  $X_j = \epsilon_j$  where  $\epsilon_j \sim N(0, 1)$  for  $j = 1, \dots, p$ . The one factor model is given by

$$\begin{aligned} X_j &= U_1 F_1 + \epsilon_j^1 & \epsilon_j^1 &\sim N(0, 1), & j &= 1, \dots, 30 \\ X_j &= \epsilon_j^0 & \epsilon_j^0 &\sim N(0, 1), & j &= 31, \dots, p \end{aligned}$$

where  $U_1 \sim Unif(0, 1)$  and  $F_1 \sim N(0, 1)$ . The two factor model is simulated from

$$\begin{aligned} X_j &= U_1 F_1 + \epsilon_j^1 & \epsilon_j^1 &\sim N(0, 1), & j &= 1, \dots, 30 \\ X_j &= U_2 F_2 + \epsilon_j^2 & \epsilon_j^2 &\sim N(0, 1), & j &= 31, \dots, 50 \\ X_j &= \epsilon_j^0 & \epsilon_j^0 &\sim N(0, 1), & j &= 51, \dots, p \end{aligned}$$

where  $U_2 \sim Unif(0.5, 1.5)$  and  $(F_1, F_2)' \sim MVN(0, \mathbf{I})$ . Finally, the three factor model is simulated from the following model:

$$\begin{aligned} X_j &= U_1 F_1 + \epsilon_j^1 & \epsilon_j^1 &\sim N(0, 1), & j &= 1, \dots, 30 \\ X_j &= U_2 F_2 + \epsilon_j^2 & \epsilon_j^2 &\sim N(0, 1), & j &= 31, \dots, 50 \\ X_j &= U_3 F_3 + \epsilon_j^3 & \epsilon_j^3 &\sim N(0, 1), & j &= 51, \dots, 75 \\ X_j &= \epsilon_j^0 & \epsilon_j^0 &\sim N(0, 1), & j &= 76, \dots, p \end{aligned}$$

where  $U_3 \sim Unif(1, 1.5)$  and  $(F_1, F_2, F_3)' \sim MVN(0, \mathbf{I})$ .

We consider configurations of the data by taking  $n$  samples from each of the factor models and we vary  $p$  to be less than, equal to, or more than  $n$ . The following parameter configurations are considered:  $(p = 100, n = 500)$ ,  $(p = 500, n = 500)$ ,  $(p = 500, n = 100)$ . We also consider the special case when  $p = 59, n = 39$ , which is the number of distinct cancer types and the sample size of the SEER cancer mortality data. In this special case, the number of random variables loading on  $F_1$  is 25, the number of random variables loading on  $F_2$  is 15 and 10 on  $F_3$ .

### 5.1. Simulation Results for Estimating the Number of Factors

In this section, we use the simulated data sets to demonstrate the behavior of the sequential-rescaling procedure when used to estimate the number of factors in a model with zero, one, two or three underlying factors. We compare the proposed procedure to the Guttman-Kaiser criterion and parallel analysis.

**Table 1.** Simulation results based on 1500 simulated data sets for selecting the true number of factors comparing Guttman criterion (Gu), parallel analysis (Pa) and the proposed methodology (Pr). Presented is the discrete probability of the estimated number of factors ( $\hat{m}$ ) and its corresponding mean and standard deviation for one, two and three factor models. The number of random variables ( $p$ ) and sample size ( $n$ ) are varied.

$(p, n)$	$\hat{m}$	Zero Factor			One Factor			Two Factor			Three Factor		
		Gu	Pa	Pr	Gu	Pa	Pr	Gu	Pa	Pr	Gu	Pa	Pr
(100, 500)	0	0.00	0.95	0.98	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	1	0.00	0.04	0.02	0.00	0.92	0.98	0.00	0.00	0.01	0.00	0.00	0.00
	2	0.00	0.01	0.00	0.00	0.08	0.02	0.00	0.99	0.97	0.00	0.00	0.00
	3	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.03	0.00	1.00	0.97
	4+	1.00	0.00	0.00	1.00	0.00	0.00	1.00	0.00	0.00	1.00	0.00	0.03
	$mean(\hat{m})$		45.32	0.07	0.02	41.68	1.09	1.02	35.81	2.01	2.03	25.26	3.00
$sd(\hat{m})$		0.73	0.34	0.14	0.73	0.33	0.15	0.74	0.10	0.17	0.75	0.00	0.17
(500, 500)	0	0.00	0.95	0.98	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	1	0.00	0.03	0.02	0.00	0.74	0.98	0.00	0.00	0.00	0.00	0.00	0.00
	2	0.00	0.01	0.00	0.00	0.18	0.02	0.00	0.68	0.99	0.00	0.00	0.00
	3	0.00	0.01	0.00	0.00	0.05	0.00	0.00	0.22	0.01	0.00	0.88	0.99
	4+	1.00	0.00	0.00	1.00	0.03	0.00	1.00	0.10	0.00	1.00	0.12	0.01
	$mean(\hat{m})$		195.65	0.07	0.02	194.31	1.39	1.02	191.51	2.45	2.01	185.54	3.14
$sd(\hat{m})$		0.85	0.38	0.15	0.85	0.82	0.13	0.84	0.80	0.10	0.84	0.38	0.06
(500, 100)	0	0.00	0.97	0.97	0.00	0.04	0.01	0.00	0.00	0.00	0.00	0.00	0.00
	1	0.00	0.03	0.03	0.00	0.47	0.98	0.00	0.00	0.01	0.00	0.00	0.00
	2	0.00	0.00	0.00	0.00	0.19	0.01	0.00	0.02	0.98	0.00	0.00	0.01
	3	0.00	0.00	0.00	0.00	0.05	0.00	0.00	0.01	0.01	0.00	0.00	0.98
	4+	1.00	0.00	0.00	1.00	0.25	0.00	1.00	0.97	0.00	1.00	1.00	0.01
	$mean(\hat{m})$		99.00	0.03	0.03	99.00	21.94	1.00	99.00	382.91	2.00	99.00	445.02
$sd(\hat{m})$		0.00	0.18	0.16	0.00	88.39	0.14	0.00	126.46	0.15	0.00	5.50	0.13
(59, 39)	0	0.00	0.96	0.97	0.00	0.01	0.02	0.00	0.00	0.00	0.00	0.00	0.00
	1	0.00	0.03	0.03	0.00	0.85	0.96	0.00	0.00	0.03	0.00	0.00	0.00
	2	0.00	0.01	0.00	0.00	0.11	0.02	0.00	0.91	0.96	0.00	0.00	0.04
	3	0.00	0.00	0.00	0.00	0.03	0.00	0.00	0.08	0.01	0.00	0.97	0.96
	4+	1.00	0.00	0.00	1.00	0.00	0.00	1.00	0.01	0.00	0.00	0.03	0.00
	$mean(\hat{m})$		21.55	0.05	0.03	20.76	1.17	1.00	19.21	2.09	1.98	17.77	3.03
$sd(\hat{m})$		0.67	0.26	0.18	0.70	0.49	0.19	0.77	0.33	0.21	0.80	0.18	0.21

We present simulation results in Table 1 for 1500 simulated data sets derived from zero, one, two or three factor models. The results in Table 1 shows that the proposed method performs well, and in almost all cases outperforms the Guttman-Kaiser criterion and parallel analysis. The Guttman-Kaiser criterion consistently overestimates the number of factors to be retained, compared to other methods (even when  $n > p$ ). We note that when  $p$  is vastly larger than  $n$ , the Guttman-Kaiser criterion always estimates the number of factors to be the rank of  $\hat{\mathbf{R}}$ .

The parallel analysis method of estimating the number of factors is relatively accurate across the range of factor models when  $n > p$ . When  $p$  become comparable to  $n$ , the parallel analysis estimator overestimates  $m$  compared to the proposed methodology. As  $p$  becomes significantly larger than  $n$ , the parallel analysis estimator breaks down and significantly overestimates  $m$ .

## 5.2. Simulation Results of BIC Criterion

For each of the simulated data sets, there are numerous random variables that have zero loadings on the underlying factors. We perform sparse principal components analysis

**Table 2.** Simulation results based on 200 simulated data sets for the proposed BIC-type criterion tuning parameter selection. The  $|\ell_m|$  denotes the true number of random variables that load on each corresponding factor and  $|\hat{\ell}_m|$  denotes the mean number of nonzero factor loadings for each factor across the simulated data sets. FP and FN denote the false positive rate and false negative rate, respectively. The number of random variables ( $p$ ) and sample size ( $n$ ) are varied.

$(p, n)$		One Factor	Two Factors		Three Factors		
		$F_1$	$F_1$	$F_2$	$F_1$	$F_2$	$F_3$
(100, 500)	$ \ell_m $	30.00	30.00	20.00	30.00	20.00	25.00
	$ \hat{\ell}_m $	30.10	29.96	20.00	30.01	20.06	25.76
	FP	0.001	0.000	0.000	0.000	0.001	0.010
	FN	0.000	0.001	0.000	0.000	0.000	0.000
(500, 500)	$ \ell_m $	30.00	30.00	20.00	30.00	20.00	25.00
	$ \hat{\ell}_m $	30.02	28.61	20.00	30.00	20.03	25.75
	FP	0.000	0.000	0.000	0.000	0.000	0.002
	FN	0.000	0.046	0.000	0.000	0.000	0.000
(500, 100)	$ \ell_m $	30.00	30.00	20.00	30.00	20.00	25.00
	$ \hat{\ell}_m $	31.05	33.66	32.40	31.64	29.53	31.77
	FP	0.008	0.012	0.026	0.009	0.020	0.014
	FN	0.085	0.065	0.000	0.083	0.000	0.000
(59, 39)	$ \ell_m $	25.00	25.00	15.00	25.00	15.00	10.00
	$ \hat{\ell}_m $	27.94	23.63	24.36	23.41	22.63	20.00
	FP	0.177	0.122	0.228	0.115	0.310	0.289
	FN	0.123	0.221	0.044	0.220	0.402	0.415

on each of the simulated extracted factor loadings to obtain vectors of factor loadings with zero loadings that can help in interpreting the underlying factors. We choose the factor-specific tuning parameters  $(\gamma_{1,1}, \dots, \gamma_{1,m})$  based on the BIC criterion described in Section 4.2. We consider 200 simulated data sets for one, two and three factors models with varying  $(p, n)$  as describe earlier.

We present the estimated number of nonzero factor loadings, the false positive rate and false negative rate for each factor based on sparse PCA using the proposed BIC criterion in Table 2. Across all factor models, the BIC tuning parameter selection method selects the true nonzero loadings with good consistency when  $n$  is large and when  $n$  is larger or comparable to  $p$ . When  $p > n$ , the BIC tends to select larger models for factor 2 and factor 3 and the false positive rate and false negative rate are no longer negligible.

## 6. Data Analysis

The motivation for the proposed statistical methodology is derived from work on identifying change patterns in cancer mortality trends. Cancer mortality data for the United States comes from the National Cancer Institute’s Surveillance, Epidemiology and End Results (SEER) Program. We analyze age-adjusted cancer mortality change patterns separately for males and females. For the sake of brevity and to avoid redundancy, we only present the results for males.

In the study of cancer mortality change pattern trends, it is common to use the log transformed annual rate change (ARC) instead of the actual mortality rate. The ARC of cancer type  $j$  in year  $i$  denoted by  $ARC_{ij}$  is defined as  $ARC_{ij} = \log r_{ij} -$

**Table 3.** Sequential-rescaling procedure: Largest four estimated eigenvalues of the pseudo covariance matrix ( $\hat{\xi}_i$ ) are denoted in the second column. The p-value for the corresponding Tracy-Widom test and alpha spending function,  $\alpha^*(k)$ , are in the third and fourth column, respectively. The last column denotes the decision to retain or not retain the factor.

Factor	$\hat{\xi}_i$	$p$ -value	$\alpha^*(k)$	Decision
1	268.90	<0.0001	0.0250	Retain
2	234.42	<0.0001	0.0125	Retain
3	219.46	0.0003	0.0063	Retain
4	207.17	0.0033	0.0031	Do not retain

$\log r_{i-1,j}$ , where  $r_{ij}$  denotes the cancer mortality rate of cancer type  $j$  in year  $i$ , and the log transformation is applied to normalize the data and the difference to construct independent components (Kim, Fay, Feuer, & Midthune, 2000). Because the different cancer types have varying levels of volatility, we will center and standardize  $ARC_{ij}$  such that it has mean 0 and variance 1. We denote the standardized rate change as  $X_{ij}$ . We obtain an estimate of the correlation matrix,  $\hat{\mathbf{R}} = \frac{1}{n-1} \mathbf{X}^T \mathbf{X}$  where  $\mathbf{X}_{n \times p}$  is the data matrix with  $X_{ij}$  as its  $(i, j)$ -th entry. Cancer mortality rates were obtained for  $p = 59$  distinct male cancer types over  $n = 39$  years (1969-2008).

### 6.1. Application of proposed methodology to SEER data

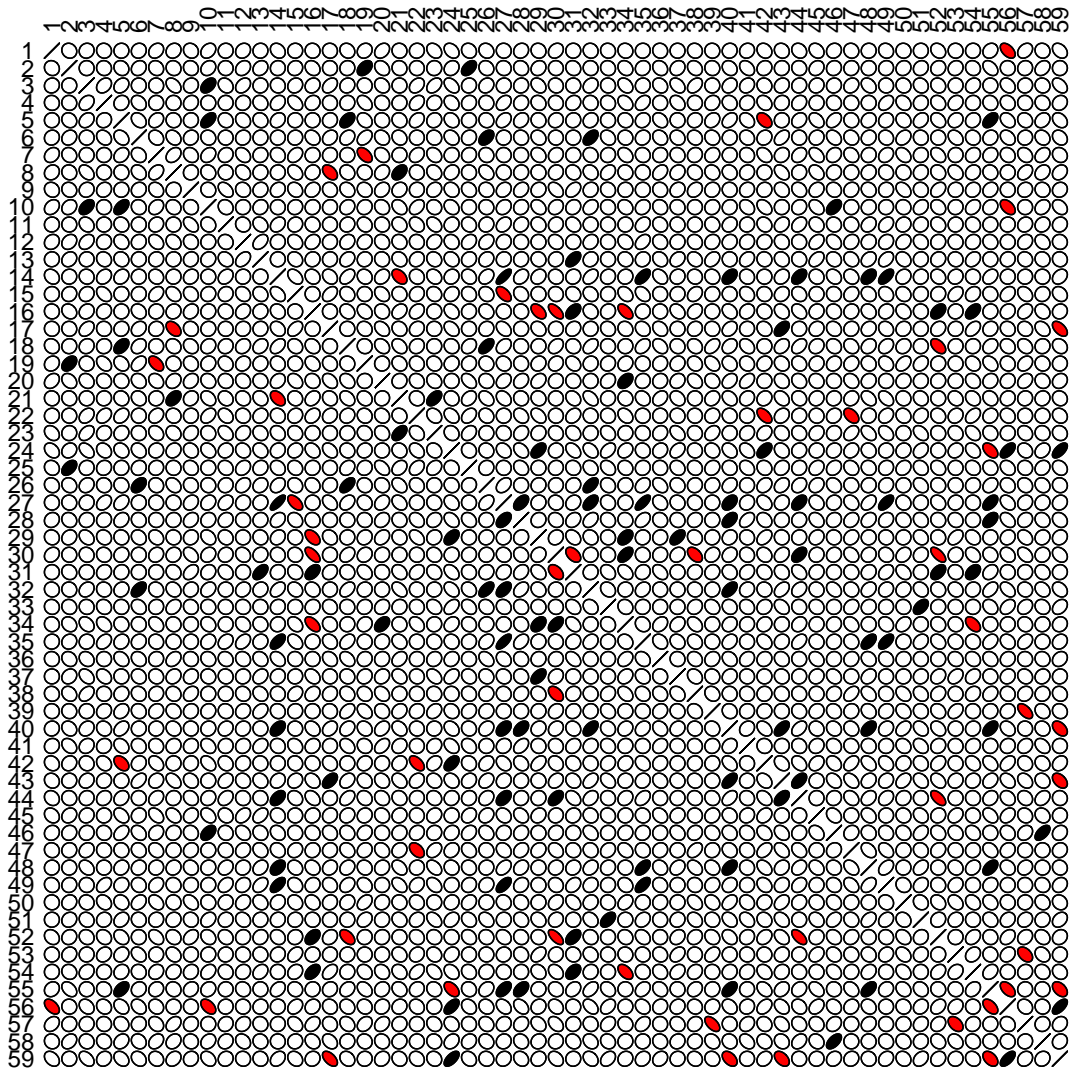
To visualize the correlation matrix of cancer ARC, we construct a correlation matrix using ellipse-shared glyphs for each entry in Figure 1. Overall, Figure 1 displays how the correlation matrix is dominated by low correlations between the cancer types. It is feasible that the population correlation matrix of ARC could be equal to the identity matrix and that the few moderate observed correlations are simply noisy estimates.

We begin our investigation of the correlation matrix of ARC by testing the null hypothesis  $H_0 : \mathbf{R} = \mathbf{I}$  versus the alternative  $H_A : \mathbf{R} \neq \mathbf{I}$ . To test this hypothesis, we study the largest eigenvalue of  $\hat{\mathbf{R}}$  which was estimated to be 7.12. After performing the correlation correction of the Tracy-Widom test described in Section 2.3, the estimated largest eigenvalue of the pseudocovariance matrix is 268.90. Applying the Tracy-Widom test on this value, we calculate the test statistic  $T_{np} = 8.63$  where  $\mu_{np} = (\sqrt{59-1} + \sqrt{38})^2 = 189.89$  and  $\sigma_{np} = (\sqrt{59-1} + \sqrt{38}) \left( \frac{1}{\sqrt{59-1}} + \frac{1}{\sqrt{38}} \right)^{1/3} = 9.16$ . Compared to the Tracy-Widom distribution of order 1, the test statistic results in a p-value  $< 0.0001$ . We reject the null hypothesis of complete independence in ARC between cancer types which suggests that at least one factor is sufficient to describe the observed correlation among the cancer types.

Next, we determine the number of factors to be used in the analysis using the sequential rescaling procedure described in Section 3.1 and present those results in Table 3. Table 3 suggests that three underlying factors are important in characterizing the correlation matrix of cancer mortality ARC.

Next, for the three extracted factors we performed sparse principle components analysis described in Section 4 to regularize the factor loadings. Only cancer types with meaningful associations to each underlying factor will have a nonzero factor loading, and we consider these to be important for the interpretation of the factors. We set  $\gamma = 1.0 \times 10^4$  in our SPCA analysis because the number of cancer types exceeded the number of data points available. To determine the degree of sparsity for each factor, we selected  $(\gamma_{1,1}, \gamma_{1,2}, \gamma_{1,3})$  to be the values that minimized the BIC criterion in (8).

We present in Table 4, the 59 unique cancer types and their corresponding sparse factor loadings for the extracted factors. Of all the 59 cancers, 28 cancer types had zero loadings on all three factors. We note that lung and bronchus, prostate, and colon cancer sites load heavily on the first factor but has exactly zero loadings for factors 2 and 3. Factor 1 might provide more support to the hypothesis that, as for colorectal cancer, early detection through screening and advances in treatment for prostate cancer are important factors that underlie the change in mortality rate. Factor 2 appears to contrast soft tissue cancers and leukemia, however, it is not clearly evident what is driving to their observed correlation. The interpretation of factor 3 appears to be highly related to miscellaneous cancer types (miscellaneous malignant cancer,



**Figure 1.** Visualization of the standardized log-annual cancer mortality rate change correlation matrix of the 59 unique male cancer types. Ellipse-shaped glyphs for each entry represent the level curve of a bivariate normal density with the matching correlation. Ellipses in black denote a positive correlation greater than 0.4 and a red ellipse denotes a negative correlation of more than -0.4. The cancer type can be matched to the number on the figure and in Table 4.

1 other myeloid/monocytic leukemia, other digestive organs, etc.). Appendix Figure 1  
2 provides additional information on each factor and their ARC change over time.  
3  
4

## 5 **7. Discussion**

6

7 We have described a methodology based on random matrix theory that uses factor  
8 analysis to make inference on correlation matrices for settings where  $p > n$ . The meth-  
9 ods described herein are applicable to a wide range of data, because it can be applied  
10 to cases where  $p > n$  as well as to traditional cases where  $n > p$ . We observed that  
11 current methods for selecting the number of factors (Guttman-Kaiser criterion and  
12 parallel analysis) do not perform well when  $p > n$ . Thus, we developed a sequential-  
13 rescaling procedure to determine the number of significant factors in a factor model  
14 using the Tracy-Widom test. This procedure is based on group sequential theory to  
15 control for the overall type I error. We described a practical approach to interpret the  
16 significant factor loadings using sparse principal components analysis and a novel BIC-  
17 type criterion which regularizes the noisy estimates of the factor loadings. Simulation  
18 studies demonstrate great performance for the proposed methodology in selecting the  
19 number of factors to be extracted and for identifying the important random variables  
20 that load on the underlying factors.  
21

22 A number of open problems present themselves. The methods herein were con-  
23 structed under the normality assumption. It is unclear how to determine complete  
24 randomness against any deviation from normality. For future work, it would be ideal  
25 to study the robustness of this methodology and the Tracy Widom test against differ-  
26 ent distributional assumptions. Another limitation is that we have not explored any  
27 methods that test whether the change patterns of any two specific cancer types are  
28 correlated over time. Factor analysis identifies groups of cancers that are linearly de-  
29 pendent upon a few unobservable latent random variables, but cannot make specific  
30 statements about pairwise correlations. Identifying specific pairs of cancers that share  
31 similar change patterns could be extremely useful for cancer researchers. One avenue  
32 to explore related to the identification of significant pairwise change patterns would be  
33 to regularize the elements of the correlations themselves, which have been extensively  
34 studied for covariance matrices (Cai & Liu, 2011; Fan, Liao, & Liu, 2016; Rothman,  
35 Levina, & Zhu, 2009). Lastly, although our focus was on the study of the correlation  
36 matrix when  $p > n$ , future studies should compare the performance of the Tracy-  
37 Widom test and sequential-rescaling procedure on covariance matrices to compare the  
38 performance and generalizability of these methods.  
39  
40  
41  
42  
43  
44

## 45 **Acknowledgement(s)**

46 The authors thank Armin Schwartzman for his insightful comments and Amanda  
47 Delzer Hill for her editing work that greatly improved the manuscript.  
48  
49  
50

## 51 **Disclosure statement**

52

53 The authors have no conflict of interests to declare.  
54  
55  
56  
57  
58



**Table 4.** Specific male cancer types and their corresponding sparse factor loadings. Sparse loadings are estimated by SPCA.

#	Cancer Type	Factor 1	Factor 2	Factor 3
1	Lip			
2	Tongue			
3	Salivary gland			-0.10
4	Floor of mouth			
5	Gum and other mouth			
6	Nasopharynx			
7	Tonsil			
8	Oropharynx			
9	Hypopharynx		0.01	
10	Other oral cavity and pharynx			-0.04
11	Esophagus			
12	Stomach			
13	Small Intestine			
14	Colon excluding rectum	-0.14		
15	Rectum and rectosigmoid junction	0.02		
16	Anus, anal canal and anorectum		0.22	
17	Liver			-0.11
18	Intrahepatic bile duct	-0.04		
19	Gallbladder			
20	Other biliary			
21	Pancreas	0.04		
22	Retroperitoneum			
23	Peritoneum, omentum and mesentery			
24	Other digestive organs			0.30
25	Nose, nasal cavity and middle ear			
26	Larynx	-0.01		
27	Lung and bronchus	-0.16		
28	Pleura	-0.02		
29	Trachea, mediastinum and other respiratory organs			
30	Bones and joints		-0.24	
31	Soft tissue including heart		0.18	
32	Melanoma of the skin	-0.05		-0.05
33	Other non-epithelial skin			0.05
34	Breast		-0.14	
35	Prostate	-0.12		
36	Testis			
37	Penis			
38	Other male genital organs			
39	Urinary bladder			
40	Kidney and renal pelvis	-0.06		-0.08
41	Ureter			
42	Other urinary organs			0.08
43	Eye and orbit			-0.13
44	Brain and other nervous system	-0.05	-0.05	
45	Thyroid			
46	Other endocrine including thymus			
47	Hodgkin Lymphoma			
48	Non-hodgkin lymphoma	-0.08		
49	Myeloma	-0.10		
50	Acute lymphocytic leukemia		-0.03	
51	Chronic lymphocytic leukemia			
52	Other lymphocytic leukemia		0.14	
53	Acute myeloid leukemia			
54	Acute monocytic leukemia		0.10	
55	Chronic myeloid leukemia	-0.05		-0.12
56	Other myeloid/monocytic leukemia			0.24
57	Other acute leukemia			
58	Aleukemic, subleukemic and NOS			
59	Miscellaneous malignant cancer			0.34

## Funding

This work was supported by National Institutes of Health Grants RO1 CA95747 and P01CA134294-01002.

## 8. References

### References

- Anderson, T. (1963). Asymptotic theory for principal component analysis. *The Annals of Mathematical Statistics*, *34*(1), 122–148.
- Bai, J. (2003). Inferential theory for factor models of large dimensions. *Econometrica*, *71*(1), 135–171.
- Bai, J., & Ng, S. (2002). Determining the number of factors in approximate factor models. *Econometrica*, *70*(1), 191–221.
- Baik, J., & Silverstein, J. (2006). Eigenvalues of large sample covariance matrices of spiked population models. *Journal of Multivariate Analysis*, *97*(6), 1382–1408.
- Bickel, P., & Levina, E. (2008). Regularized estimation of large covariance matrices. *The Annals of Statistics*, *36*(1), 199–227.
- Cai, T., & Liu, W. (2011). Adaptive thresholding for sparse covariance matrix estimation. *Journal of the American Statistical Association*, *106*(494), 672–684.
- Carvalho, C., Chang, J., Lucas, J., Nevins, J., Wang, Q., & West, M. (2008). High-dimensional sparse factor modeling: Applications in gene expression genomics. *Journal of the American Statistical Association*, *103*(484), 1438–1456.
- Fan, J., Fan, Y., & Lv, J. (2008). High dimensional covariance matrix estimation using a factor model. *Journal of Econometrics*, *147*(1), 186–197.
- Fan, J., Liao, Y., & Liu, H. (2016). An overview of the estimation of large covariance and precision matrices. *The Econometrics Journal*, *19*(1), C1–C32.
- Fan, J., Lv, J., & Qi, L. (2011). Sparse high-dimensional models in economics. *Annu. Rev. Econ.*, *3*, 291–317.
- Ghosh, J., & Dunson, D. (2008). Bayesian Model Selection in Factor Analytic Models. *Random Effect and Latent Variable Model Selection*, 151.
- Guttman, L. (1954). Some necessary conditions for common-factor analysis. *Psychometrika*, *19*(2), 149–161.
- Hayton, J., Allen, D., & Scarpello, V. (2004). Factor retention decisions in exploratory factor analysis: A tutorial on parallel analysis. *Organizational Research Methods*, *7*(2), 191–205.
- Horn, J. (1965). A rationale and test for the number of factors in factor analysis. *Psychometrika*, *30*(2), 179–185.
- Huang, J., Liu, N., Pourahmadi, M., & Liu, L. (2006). Covariance matrix selection and estimation via penalised normal likelihood. *Biometrika*, *93*(1), 85–98.
- Johnson, R., & Wichern, D. (1998). *Applied multivariate statistical analysis*. Prentice Hall Englewood Cliffs, NJ.
- Johnstone, I. (2001). On the distribution of the largest eigenvalue in principal components analysis. *The Annals of Statistics*, *29*, 295–327.
- Kaiser, H. (1960). The application of electronic computers to factor analysis. *Educational and Psychological Measurement*, *20*(1), 141–151.
- Kim, H., Fay, M., Feuer, E., & Midthune, D. (2000). Permutation tests for joinpoint regression with applications to cancer rates. *Statistics in Medicine*, *19*(3), 335–351.
- Lan, G., & DeMets, D. (1983). Discrete sequential boundaries for clinical trials. *Biometrika*, *70*(3), 659–663.
- Leek, J. (2011). Asymptotic conditional singular value decomposition for high-dimensional genomic data. *Biometrics*, *67*(2), 344–352.

- 1 Onatski, A. (2009). Testing hypotheses about the number of factors in large factor models.  
2 *Econometrica*, 77(5), 1447–1479.
- 3 Patterson, N., Price, A., & Reich, D. (2006). Population structure and eigenanalysis. *PLoS*  
4 *genetics*, 2(12), e190.
- 5 Rothman, A., Levina, E., & Zhu, J. (2009). Generalized thresholding of large covariance  
6 matrices. *Journal of the American Statistical Association*, 104(485), 177–186.
- 7 Shen, H., & Huang, J. (2008). Sparse principal component analysis via regularized low rank  
8 matrix approximation. *Journal of multivariate analysis*, 99(6), 1015–1034.
- 9 Shimizu, H., Arimura, Y., Onodera, K., Takahashi, H., Okahara, S., Kodaira, J., ... others  
10 (2016). Malignant potential of gastrointestinal cancers assessed by structural equation  
11 modeling. *PloS one*, 11(2), e0149327.
- 12 Tracy, C., & Widom, H. (2000). The distribution of the largest eigenvalue in the Gaussian  
13 ensembles. In *Calogero-Moser-Sutherland Models*, 4, 461–472.
- 14 West, M. (2003). Bayesian factor regression models in the large p, small n paradigm. *Bayesian*  
15 *statistics*, 7, 723–732.
- 16 Wong, F., Carter, C., & Kohn, R. (2003). Efficient estimation of covariance selection models.  
17 *Biometrika*, 90(4), 809–830.
- 18 Zhou, Y., Wang, P., Wang, X., Zhu, J., & Song, P. X.-K. (2017). Sparse multivariate factor  
19 analysis regression models and its applications to integrative genomics analysis. *Genetic*  
20 *Epidemiology*, 41(1), 70–80.
- 21 Zou, H., Hastie, T., & Tibshirani, R. (2006). Sparse principal component analysis. *Journal of*  
22 *computational and graphical statistics*, 15(2), 265–286.
- 23 Zou, H., Hastie, T., & Tibshirani, R. (2007). On the degrees of freedom of the lasso. *The*  
24 *Annals of Statistics*, 35(5), 2173–2192.
- 25 Zwick, W., & Velicer, W. (1986). Comparison of five rules for determining the number of  
26 components to retain. *Psychological Bulletin*, 99(3), 432–442.
- 27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## 9. Appendices

In this section, we show that the proposed alpha spending function

$$\alpha^*(k) = \frac{\alpha}{2^k}$$

to test the number of significant factors will not exceed  $\alpha$ , by calculating three probabilities.

- (1) Probability that a model with one or more factors is chosen given a true zero factor model.

Let  $L_m$  be the event that the true model has  $m$  significant factors and  $\hat{L}_m$  the estimated number of factors. Then it follows that

$$P(\hat{L}_{k \geq 1} | L_0) = 1 - P(\hat{L}_0 | L_0) = 1 - (1 - \alpha) = \alpha$$

- (2) Probability that a model is selected with  $k$  factors given a true zero factor model for any  $k \geq 1$ .

$$P(\hat{L}_k | L_0) = \left(1 - \frac{\alpha}{2^{k+1}}\right) \prod_{q=1}^k \frac{\alpha}{2^q}$$

Because  $\left(\frac{\alpha}{2^{k+1}}\right) < 1$  it follows that

$$P(\hat{L}_k | L_0) < \prod_{q=1}^k \frac{\alpha}{2^q} < \alpha \prod_{q=1}^k \frac{1}{2^q}$$

where the last equality follows as the result of  $\alpha^k < \alpha$  as  $\alpha \in (0, 1)$ . Finally, as  $\prod_{q=1}^k \frac{1}{2^q} < 1$ , we get the result that

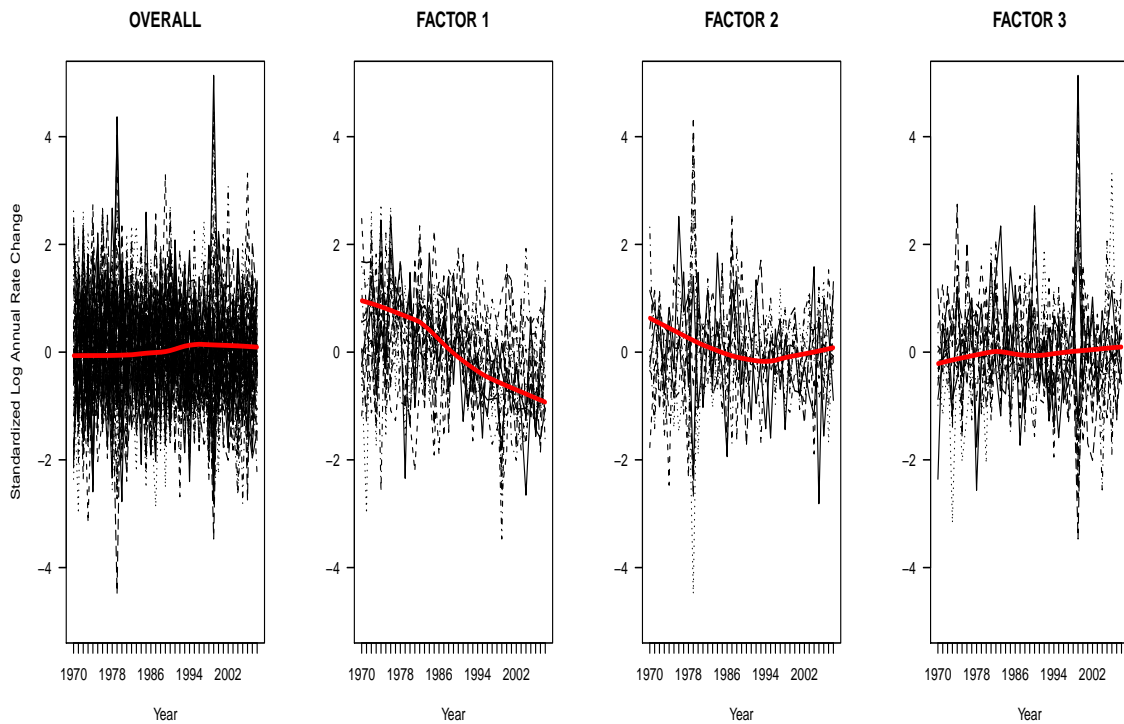
$$P(\hat{L}_k | L_0) < \alpha$$

- (3) Probability that a model is selected with more than  $k$  factors given a true factor model with  $k$  factors.

$$\begin{aligned} P(\hat{L}_{q > k} | L_k) &= \sum_{w=1}^{\infty} \left[ \left(1 - \frac{\alpha}{2^{k+w}}\right) \prod_{q=k+1}^{k+w} \frac{\alpha}{2^q} \right] \\ &< \sum_{w=1}^{\infty} \prod_{q=k+1}^{k+w} \frac{\alpha}{2^q} \\ &< \sum_{w=1}^{\infty} \frac{\alpha}{2^w} = \alpha \end{aligned}$$

Thus, the type I error does not exceed alpha in any of the settings.

1 In Appendix Figure A1, we plot the cancer mortality standardized log ARC over  
2 time for for all 59 cancer types and also three separate plots for the cancer types that  
3 have non-zero loadings for each factor. To visualize the pattern over time, we fit a  
4 Lowess smoothing line across time. Overall, when we consider the change patterns  
5 of ARC for all 59 cancer types simultaneously, we do not observe much change in  
6 ARC over time. The factor analysis performed identifies three distinct cancer mortality  
7 patterns of ARC over time. Factor 1 is a collection of cancer types (primarily influenced  
8 by colon, prostate and lung cancers) that have exhibit a decrease in ARC cancer  
9 mortality across time. The cancer types in factor 2 have decreasing ARC that levels  
10 off after the year 1990. Finally, factor 3 (miscellaneous) cancer types show no change  
11 in ARC over time.  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30



54 **Figure A1.** Line plots of standardized log annual cancer mortality rate change over time. Left panel includes  
55 all cancer types and the last three panels plots the standardized log ARC for the cancer types that have  
56 non-zero loadings on factors 1, 2 and 3, respectively. The solid red line denotes Lowess smoothing curves.  
57  
58