# A Correlation Analysis Method for Power Systems Based on Random Matrix Theory

Xinyi Xu, Xing He, Qian Ai, *Member, IEEE,*, Robert C. Qiu, *Fellow, IEEE,*

*Abstract*—The operating status of power systems is influenced by growing varieties of factors, resulting from the developing sizes and complexity of power systems; in this situation, the model-based methods need be revisited. A data-driven method, as the novel alternative, on the other hand, is proposed in this paper: it reveals the correlations between the factors and the system status through statistical properties of data. An augmented matrix, as the data source, is the key trick for this method; it is formulated by two parts: 1) status data as the basic part, and 2) factor data as the augmented part. The random matrix theory (RMT) is applied as the mathematical framework. The linear eigenvalue statistics (LESs), such as the mean spectral radius (MSR), are defined to study data correlations through large random matrices. Compared with model-based methods, the proposed method is inspired by a pure statistical approach, without *a prior* knowledge of operation and interaction mechanism models for power systems and factors. In general, this method is direct in analysis, robust against bad data, universal to various factors, and applicable for real-time analysis. A case study, based on the standard IEEE 118-bus system, validates the proposed method.

*Index Terms*—correlation analysis, power systems, big data analytics, augmented matrix, random matrix theory, linear eigenvalue statistics

## I. INTRODUCTION

**T**HE operating status of power systems is affected by numerous factors. It is fundamental to understand the statistical correlations between those factors and power systems. These correlations reveal the causes to disturbances and faults [1]. Nowadays, power systems, large in sizes and complex in structure, are penetrated by more and more various elements, such as distributed generations, flexible loads, and electric vehicles. All these elements lead to strong interaction, multiple coupling, and high randomness in power systems. On this occasion, model-based methods, establishing mechanism models with assumptions and simplifications as essential preconditions, are questionable.

Data have become a strategic resource in power systems. The 4Vs data [2], with great potential value, are hard to handle by conventional model-based methods. This situation leads to an emerging paradigm—big data analytics—for power systems. Big data analytics aims to work out statistical features measured by the eigenvalue statistics, without establishing mechanism models. The linear eigenvalue statistics (LESs) are of central interest in statistics [3]. When the matrix size is sufficiently large, LESs tend to deterministic limiting values (expected values). Various forms of the central limit theorems are also established in recent statistical papers. The statistical

error between the limiting expected value and the eigenvalue statistics will be decreasing as a function of data size. For more details on the convergence rate, we refer to [4].

Many studies on big data analytics have been achieved in power systems [5–7]. In our previous work, a universal architecture with big data analytics is proposed [8], and applied successfully in many fields, such as the anomaly detection [9], and the 3D power map for situational awareness [10]. This paper is built upon our previous paper [8] in that the augmented matrix contains a remarkable rich statistical information between two (or more) block matrices. The trick is the observation that a matrix of block random submatrices are also a (large) random matrix. This simple observation delivers many interesting results that are useful in large power systems.

### A. Contribution

This paper, based on random matrix theory (RMT), proposes a data-driven method to reveal the correlations between factors and the power system status. An augmented matrix, formed in a certain manner, is presented as the data source. For each factor, the augmented matrix combines status data (as the basic part), on one hand, and factor data (as the augmented part), on the other hand. According to specific researching purposes, status data can be voltages, frequencies, currents and power flows, while factor data can be loads, distributed generation, wind speed, temperature, electricity price, etc. Then, using the big data architecture proposed previously [8], we conduct real-time analysis based on the augmented matrix, and compare findings with the RMT theoretical predictions (i.e. Ring Law and Marchenko-Pastur Law). During this procedure, the mean spectral radius (MSR), a special case of LESs, is used to indicate data correlations; the kernel density estimation (KDE) is used as an assisted indicator.

In general, the proposed method extracts the correlations in the form of the eigenvalue statistics of measured data. The method involves no knowledge of topologies and parameters of power systems, and is universal to various factors. Besides, the method is robust against random fluctuations in power systems and measuring errors in data. Furthermore, the proposed method is practical for both real-time analysis and off-line analysis, depending on the split-window.

### B. Related Work

Current researches on correlation analysis are mainly model-based methods, for which the mechanism models are essential preconditions. These mechanism models are established based on assumptions and simplifications, and used for

specific power systems and factors. Lian studied the effect of dynamic load characteristics on the voltage stability and sensitivity in power systems, using the P–V and the Q–V curves [11]. In Lian's method, the power system is equivalent to an decentralized system; dynamic loads are approximated as differential equations. These processes increase the complexion and inaccuracy of the analysis. Parinya proposed a stochastic stability index to investigate the small signal stability of power systems incorporating wind power [12]. The status space equations and energy functions need to be rewritten when the grid changes, and the test system is too small in scale to convince.

Also, some data-driven methods for correlation analysis are proposed recently, such as the principal components analysis, the artificial neural networks, the support vector machine [13]. Eltigani utilizes artificial neural networks (ANNs) in assessing the transient stability [14]. In his approach, the power system is described by an equivalent single machine infinite bus system, which cannot reflect accurately the actual state of the system. Moreover, with the scale-up of the system and increase of training samples, the training speed of ANNs progressively slows down.

## II. RANDOM MATRIX THEORY AND DATA PROCESSING

The frequently used notations in this section are given in Table. I.

TABLE I: Notations for RMT and data processing

| Notations | Means |
|---|---|
| $\mathbf{X}$, $\mathbf{x}$, $x$, $x_{i,j}$ | a matrix, a vector, a single value, an entry of a matrix |
| $\mathbf{X}^H$, $\mathbf{x}^H$ | transpose of a matrix and a vector |
| $\mu(\mathbf{x})$, $\sigma^2(\mathbf{x})$ | mean, variance for $\mathbf{x}$ |
| $\mathbf{\Omega}$ | the data source |
| $\mathbb{C}^{N \times T}$ | $N \times T$ dimensional complex space |
| $N$, $T$ | the row size and the column size of the split window |
| $\hat{\mathbf{X}}$ | a raw data matrix |
| $\widetilde{\mathbf{X}}$ | a standard non-Hermitian matrix |
| $\mathbf{X}_u$ | the singular value equivalent of $\mathbf{X}$ |
| $\hat{\mathbf{Z}}$ | the matrix product |
| $\widetilde{\mathbf{Z}}$ | the standard matrix product |
| $\lambda_{\mathbf{Z}}$ | eigenvalues of $\mathbf{Z}$ |
| $|\lambda|$ | radius of eigenvalue $\lambda$ on the complex plane |
| $\kappa_{\text{MSR}}$ | the mean spectral radius |
| $\mathbf{S}$ | the sample covariance matrix |

### A. Random Matrix Theory

The random matrix theory (RMT), developed from several different sources in the early 20th century, is one of the statistical foundations for big data analytics. It is used as an important mathematical tool in various fields, namely, physics, finance, wireless communication engineering, etc.

Massive data can be naturally represented by large random matrices [15]. The random matrix model is the most general: rectangular and complex. In our formulation, we view the $N$ variables as space samples of a random network (or graph). For each variable, we make $T$ observations. As a result, a random matrix of $N \times T$ is obtained as our data matrix, which is the starting part for our analysis.

According to RMT, when the dimensions of a random matrix are sufficiently large, the empirical spectral distribution (ESD) of its eigenvalues converges to some theoretical limits, such as Ring Law and Marchenko-Pastur Law (M-P Law) [16]. For details on the Ring Law and M-P Law, please refer to Appendix A. It is noted that although the asymptotic convergence in RMT is considered under infinite dimensions, the asymptotic results are remarkably accurate for relatively moderate matrix sizes such as tens. This is the very reason why RMT can be used in practical world.

### B. Real-time Data Processing for RMT

Currently, there exists no general standardized definition for big data. In this paper, we use a mathematical definition proposed in our previous work [8]. In a power system, assume that there are $n$ kinds of measurable status variables. At the sampling time $t_i$, measured data of these variables are formed as a column vector $\hat{\mathbf{x}}(t_i) = (\hat{x}_1, \hat{x}_2, \ldots, \hat{x}_n)^H$ [17]. For a series of time, we can arrange these vectors $\hat{\mathbf{x}}$ in chronological order to form a matrix as a data source $\mathbf{\Omega}$ for further analysis ($\hat{\mathbf{x}} \in \mathbf{\Omega}$).

Within $\mathbf{\Omega}$, we can obtain a raw data matrix $\hat{\mathbf{X}} \in \mathbb{C}^{N \times T}$ arbitrarily by using a split-window. Then, we convert it into a standard non-Hermitian matrix $\widetilde{\mathbf{X}}$ with following algorithms.

$$\widetilde{x}_{i,j} = (\hat{x}_{i,j} - \mu(\hat{\mathbf{x}}_i)) \times \frac{\sigma(\widetilde{\mathbf{x}}_i)}{\sigma(\hat{\mathbf{x}}_i)} + \mu(\widetilde{\mathbf{x}}_i) \tag{1}$$

where $\hat{\mathbf{x}}_i = (\hat{x}_{i,1}, \hat{x}_{i,2}, \ldots, \hat{x}_{i,T})$, $\mu(\widetilde{\mathbf{x}}_i) = 0$, $\sigma(\widetilde{\mathbf{x}}_i) = 1$ for $i = 1, 2, \ldots, N$ and $j = 1, 2, \ldots, T$. The matrix $\widetilde{\mathbf{X}}_u \in \mathbb{C}^{N \times N}$ is introduced as the singular value equivalent of $\widetilde{\mathbf{X}}$ by

$$\widetilde{\mathbf{X}}_u = \sqrt{\widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^H} \mathbf{U} \tag{2}$$

where $\mathbf{U} \in \mathbb{C}^{N \times N}$ is a Haar unitary matrix and $\widetilde{\mathbf{X}}_u \widetilde{\mathbf{X}}_u^H = \widetilde{\mathbf{X}}\widetilde{\mathbf{X}}^H$.

For multiple arbitrarily assigned standard non-Hermitian matrices $\widetilde{\mathbf{X}}_i$ ($i = 1, 2, \ldots, L$), the matrix product is obtained by $\hat{\mathbf{Z}} = \prod_{i=1}^{L} \widetilde{\mathbf{X}}_{u,i}$. Then, $\hat{\mathbf{Z}}$ is converted to the standard matrix product $\widetilde{\mathbf{Z}}$ by

$$\widetilde{\mathbf{z}}_i = \frac{\hat{\mathbf{z}}_i}{\sqrt{N}\sigma(\hat{\mathbf{z}}_i)} \tag{3}$$

where $\widetilde{\mathbf{z}}_i = (\widetilde{z}_{i,1}, \widetilde{z}_{i,2}, \ldots, \widetilde{z}_{i,N})$ and $\hat{\mathbf{z}}_i = (\hat{z}_{i,1}, \hat{z}_{i,2}, \ldots, \hat{z}_{i,N})$.

For the standard matrix product $\widetilde{\mathbf{Z}}$, the sample covariance matrix is obtained by

$$\mathbf{S} = \frac{1}{N}\mathbf{Y}\mathbf{Y}^H = \widetilde{\mathbf{Z}}\widetilde{\mathbf{Z}}^H \tag{4}$$

where $\mathbf{Y} = \sqrt{N}\widetilde{\mathbf{Z}}$, and $\sigma^2(\mathbf{y}_i) = \sigma^2(\sqrt{N}\widetilde{\mathbf{z}}_i) = 1$.

In order to conduct real-time analysis, we use a specific split-window to obtain the raw data matrix $\hat{\mathbf{X}}$ from $\mathbf{\Omega}$, namely, the real-time split-window. The real-time split-window truncates measured data at continuous sampling times, where the last sampling time is the current time. In other words, at the sampling time $t_i$, the raw data matrix $\hat{\mathbf{X}}_{t_i}$ is formed by

$$\hat{\mathbf{X}}(t_i) = (\hat{\mathbf{x}}(t_{i-T+1}), \hat{\mathbf{x}}(t_{i-T+2}), \ldots, \hat{\mathbf{x}}(t_i)) \tag{5}$$

where $\hat{\mathbf{x}}(t_j) = (\hat{x}_1, \hat{x}_2, \ldots, \hat{x}_N)^H$ is measured data at the sampling time $t_j$.

The data processing procedure above is organized as following steps. The standard matrix product $\widetilde{\mathbf{Z}}$ is calculated for Ring Law; the sample covariance matrix $\mathbf{S}$ is calculated for M-P Law. For simplicity, we set $L=1$ in the matrix product $\hat{\mathbf{Z}}$.

| Steps of Data Processing for Ring Law |
| --- |
| 1) At the sampling time $t_i$, obtain the raw data matrix $\hat{\mathbf{X}}(t_i) \in \mathbb{C}^{N \times T}$ from the data source $\mathbf{\Omega}$. |
| 2) Convert $\hat{\mathbf{X}}(t_i)$ into the standard non-Hermitian matrix $\widetilde{\mathbf{X}}(t_i)$. |
| 3) Calculate the singular value equivalent $\widetilde{\mathbf{X}}_u(t_i)$ of $\widetilde{\mathbf{X}}(t_i)$. |
| 4) Form the matrix product $\hat{\mathbf{Z}}(t_i) = \widetilde{\mathbf{X}}_u(t_i)$. |
| 5) Convert $\hat{\mathbf{Z}}(t_i)$ into the standard matrix product $\widetilde{\mathbf{Z}}(t_i)$. |
| 6) Calculate the sample covariance matrix $\mathbf{S}(t_i) = \widetilde{\mathbf{Z}}(t_i)\widetilde{\mathbf{Z}}(t_i)^H$ |

### C. Linear Eigenvalue Statistics and Kernel Density Estimation

#### 1) Linear Eigenvalue Statistics:

The linear eigenvalue statistics (LESs), major focuses in this paper, indicate the statistical characteristics of large random matrices. The linear eigenvalue statistic of a random matrix $\mathbf{X}$ is defined as [18]

$$\mathcal{N}_n(\varphi) = \sum_{i=1}^{n} \varphi(\lambda_i) \qquad (6)$$

where $\lambda_i$ $(i = 1, 2, \ldots, n)$ are eigenvalues of $\mathbf{X}$, and $\varphi(\cdot)$ is a test function.

The mean spectral radius (MSR), a special case of LESs, is used to reflect the eigenvalue distribution of the standard matrix product $\widetilde{\mathbf{Z}}$; it is defined as the mean distribution radius of eigenvalues, formulated by

$$\kappa_{\mathrm{MSR}} = \frac{1}{N} \sum_{i=1}^{N} |\lambda_{\widetilde{\mathbf{Z}},i}| \qquad (7)$$

where $\lambda_{\widetilde{\mathbf{Z}},i}$ $(i = 1, 2, \ldots, N)$ are eigenvalues of $\widetilde{\mathbf{Z}}$, and $|\lambda_{\widetilde{\mathbf{Z}},i}|$ is the radius of the eigenvalue $\lambda_{\widetilde{\mathbf{Z}},i}$ on the complex plane. Based on MSR, we define the variance of spectral radius (VSR) as a further reflection of eigenvalue distribution; it reflects the dispersion degree of the eigenvalues of $\widetilde{\mathbf{Z}}$, formulated by

$$\kappa_{\mathrm{VSR}} = \frac{1}{N-1} \sum_{i=1}^{N} (|\lambda_{\widetilde{\mathbf{Z}},i}| - \kappa_{\mathrm{MSR}})^2 \qquad (8)$$

Note that since these complex eigenvalues $\lambda_{\widetilde{\mathbf{Z}},i}$ $(i = 1, 2, \ldots, N)$ are highly correlated random variables (each is a complicated matrix function of random matrices $\widetilde{\mathbf{X}}_i$ $(i = 1, 2, \ldots, L)$, $\kappa_{\mathrm{MSR}}$ and $\kappa_{\mathrm{VSR}}$ are both random variables.

#### 2) Kernel Density Estimation:

The kernel density estimation (KDE) depicts the ESD of the sample covariance matrix $\mathbf{S}$ by following PDF

$$f_{\mathrm{KDE}}(\lambda_{\mathbf{S}}) = \frac{1}{Nh} \sum_{i=1}^{N} \mathrm{K}(\frac{\lambda_{\mathbf{S}} - \lambda_{\mathbf{S},i}}{h}) \qquad (9)$$

where $\lambda_{\mathbf{S},i}$ $(i = 1, 2, \ldots N)$ are eigenvalues of $\mathbf{S}$, and $\mathrm{K}(\cdot)$ is the kernel function for the bandwidth parameter $h$.

TABLE II: Notations for correlation analysis

| Notations | Means |
| --- | --- |
| $n$ | the number of status variables of power systems |
| $m$ | the number of influential factors in power systems |
| $t$ | the study duration |
| $t_s$ | the sampling time |
| $\mathbf{B}$ | the status matrix |
| $\mathbf{c}$ | a factor vector |
| $\mathbf{D}$ | the matrix duplicating $\mathbf{c}$ for $k$ times |
| $\mathbf{E}$ | the noise matrix |
| $m_e$ | the magnitude of white noise |
| $\mathbf{C}$ | the factor matrix |
| $\rho$ | the signal-to noise ratio |
| $\mathbf{A}$ | the augmented matrix |

## III. A CORRELATION ANALYSIS METHOD FOR POWER SYSTEMS BASED ON AUGMENTED MATRICES

The frequently used notations are given in Table.II.

### A. Augmented Matrix Method for Power Systems

Different factors have different effects on power systems in their own ways. According to big data analytics, the relationships between factors and power systems can be revealed by data correlations. Based on RMT, it is feasible to extract correlations from a data source including both status data and factors data, namely, augmented matrix.

In a power system, assume that there are $n$ types of status variables and $m$ types of influential factors, both of which are measurable. In a study period $t_i$ $(i = 1, 2, \ldots, t)$, measured data of each factor are formed as a row vector $\mathbf{c}_j \in \mathbb{C}^{1 \times t}$ $(j = 1, 2, \ldots, m)$ (i.e. factor vector), and measured data of status variables for the power system are formed as a matrix $\mathbf{B} \in \mathbb{C}^{n \times t}$ (i.e. status matrix).

In order to balance the proportion of factors data and status data in the data source, we form a factor matrix for each factor vector. First, for the factor $\mathbf{c}_j$, we duplicate it for $k$ times to construct a matrix $\mathbf{D}_j$, formulated by

$$\mathbf{D}_j = \begin{bmatrix} \mathbf{c}_j \\ \mathbf{c}_j \\ \vdots \\ \mathbf{c}_j \end{bmatrix}_{k \times t} \quad (j = 1, 2, \ldots, m) \qquad (10)$$

where $k$ has the similar size with $n$. Then, we introduce white noise into $\mathbf{D}_j$ to reduce the correlations among its own rows. The noise matrix is $\mathbf{E} = \{e_{i,j}\}_{k \times t}$, where $e_{i,j}$ is random variable according with the normal distribution. Finally, the factor matrix is formulated by

$$\mathbf{C}_j = \mathbf{D}_j + m_{e,j}\mathbf{E} \quad (j = 1, 2, \ldots, m) \qquad (11)$$

where $m_{e,j}$ is the magnitude of white noise for the factor matrix $\mathbf{C}_j$. The signal-to-noise ratio (SNR) of the factor matrix $\mathbf{C}_j$ is defined as

$$\rho_j = \frac{\mathrm{Tr}(\mathbf{D}_j\mathbf{D}_j^H)}{\mathrm{Tr}(\mathbf{E}\mathbf{E}^H) \times m_{e,j}^2} \quad (j = 1, 2, \ldots, m) \qquad (12)$$

where $\mathrm{Tr}(\cdot)$ is the trace of the matrix. The value of $\rho_j$, requiring careful selections, will affect the performance of correlation analysis. In this paper, we set the same SNR for

all factors, to ensure the consistency of correlation analysis. Therefore, for each $\mathbf{D}_j (j = 1, 2, \ldots, m)$, the value of $m_{e.j}$ is calculated by

$$m_{e,j} = \sqrt{\frac{\text{Tr}(\mathbf{D}_j \mathbf{D}_j^H)}{\text{Tr}(\mathbf{E}\mathbf{E}^H) \times \rho}} \quad (j = 1, 2, \ldots, m) \quad (13)$$

For each factor $\mathbf{c}_j$, we can construct an augmented matrix for parallel correlation analysis, formulated by

$$\mathbf{A}_j = \begin{bmatrix} \mathbf{B} \\ \mathbf{C}_j \end{bmatrix} \quad (j = 1, 2, \ldots, m) \quad (14)$$

### B. Status Data and Factor Data

#### 1) Status Variables of Power System:

The operating status of power systems can be estimated by various types of status variables, such as frequencies, voltages, currents and power flows. One or more types of these status variables can form the status matrix $\mathbf{B}$. In this paper, we use magnitudes of bus voltages as status data for the following considerations:

a) The voltage magnitude is one of the most basic parameters in power systems. It is measurable and available on every bus with common measuring devices. Therefore, magnitudes of bus voltages have considerable redundancy and accuracy for correlation analysis.

b) The voltage magnitude does not involve the topology of power systems. Therefore, we can conduct analysis without the *a prior* knowledge of network structures and parameters.

#### 2) Factors in Power Systems:

The operating status of power systems is mainly affected by electrical factors, climatic factors and economic factors.

a) Electrical factors–nodal loads and distributed generations, etc.

b) Climatic factors—temperature, wind speed, light intensity, etc.

c) Economic factors—electricity price, gross domestic product, etc.

Since data normalization is used during data processing, factor data and status data in the augmented matrix can have different units and magnitudes. In consideration of different sampling frequencies for status data and factor data, it can be assumed that data with lower sampling frequency are the same values in a sampling interval.

### C. Correlation Analysis Method and Its Advantages

Based on RMT and the augmented matrix method, a correlation analysis method for power systems is designed as following steps.

Step 2) is conducted for the anomaly detection, where the $\kappa_{\text{MSR}} - t$ curve of the status matrix discovers the signals in power systems [9]. Step 3)–4), analyzing the correlations between the system status and each factor, aim to determine the causes to the signals. During the analysis procedure, $\kappa_{\text{MSR}}$ and $f_{\text{KDE}}$ are calculated as correlation indicators.

The proposed correlation analysis method is driven by measured data, and based on statistical theories. The procedure involves no mechanism models for power systems and factors; it eliminates the interference brought by assumptions and

---

**Steps of Correlation Analysis for Power Systems**

1) In a study period, form the status matrix $\mathbf{B} \in \mathbb{C}^{n \times t}$ and factor vectors $\mathbf{c}_j \in \mathbb{C}^{1 \times t} \ (j = 1, 2, \ldots, m)$.
2) At each sampling time,
   2a) acquire the standard matrix product $\widetilde{\mathbf{Z}}$ from $\mathbf{B}$.
   2b) calculate $\kappa_{\text{MSR}}$ and $\kappa_{\text{VSR}}$ of $\widetilde{\mathbf{Z}}$.
3) Meanwhile, for each factor $\mathbf{c}_j$, construct the augmented matrix $\mathbf{A}_j$.
4) At each sampling time,
   4a) acquire the standard matrix product $\widetilde{\mathbf{Z}}$ from $\mathbf{A}_j$.
   4b) calculate $\kappa_{\text{MSR}}$ and $\kappa_{\text{VSR}}$ of $\widetilde{\mathbf{Z}}$.
   4c) acquire the sample covariance matrix $\mathbf{S}$ from $\widetilde{\mathbf{Z}}$.
   4d) calculate $f_{\text{KDE}}$ of $\mathbf{S}$ and draw the $f_{\text{KDE}} - \lambda$ curve.
5) Draw $\kappa_{\text{MSR}} - t$ curves of the status matrix and the augmented matrix for each factor.

---

simplifications. Compared with model-based methods, the proposed method is data-driven, and universal for various factors. Meanwhile, the proposed method has strong robustness against random fluctuations in systems and measuring errors in data. Besides, the method is practical for real-time analysis by using a specific split-window. The advantages above will be verified in case studies.

## IV. CASE STUDIES

The proposed method is tested with simulated data in the standard IEEE 118-bus system. Detailed information of the system is referred to the case118.m in Matpower package and Matpower 4.1 User's Manual [19]. In the simulations, we regard the active load of each bus as a factor; each change of a factor is considered as a signal.

Four cases are designed in different scenarios to validate the effectiveness of the proposed method. In case 1, case 2 and case 3, three kinds of signals are added into single factor to affect the operating status of the test system. In case 4, multiple factors with different kinds of signal produce effects on the system status.

In order to determine the causes to the signals, simulated data of each factor are used to conduct correlation analysis. It is a parallel analysis procedure, and we can only pay our attention to potential factors. Here, we just illustrate the results with the load of bus 117 and bus 54 to show the performance of the proposed method.

Assume that status data are sampled at each time interval, while factor data are sampled every 50 time intervals. Besides, white noise is introduced into both status data and factor data to represent random fluctuations and measuring errors.

For all cases, let $n = 118, t = 1000, N = 118, T = 240, k = \frac{1}{2}N = 59, \rho = 500$.

### A. Correlation Analysis for Single Factor

#### 1) Case 1 - Step signal in the load of bus 117:

In case 1, assumed signals for each factor are shown in Tab. III. The correlation analysis results are shown in Fig. 6. It is noted that the correlation analysis begins at $t_s = 240$, because the real-time split-window needs $T = 240$ times of sampling data, including the present sampling and 239 times of historical sampling.

TABLE III: Assumed Signals for Each Factor in Case 1

| Bus | Sampling Time | Active Load(MW) |
|---|---|---|
| 117 | $t_s = 1 \sim 500$ | 20.0 |
| | $t_s = 501 \sim 1000$ | 120.0 |
| Others | $t_s = 1 \sim 1000$ | Unchanged |

In Fig. 6a, based on the $\kappa_{\text{MSR}} - t$ curve, we can detect signals based on analyses below:

I. During the sampling time $t_s = 240 \sim 500$, $\kappa_{\text{MSR}}$ of status matrix remains constant around 0.86, between the outer and inner circle; it means the system status is normal without signals.

II. At $t_s = 501$, $\kappa_{\text{MSR}}$ starts to decline (from 0.8638 to 0.7002), and deviates from the predicted ring (the inner radius is 0.7130); it means there exist signals that change the system status.

III. $\kappa_{\text{MSR}}$ increases back to 0.8662 at $t_s = 740$ and remains constant inside the ring afterwards.

First, we can tell that the signals occur right at $t_s = 501$. Moreover, in our method, the impact of a signal to MSR will delay for a duration of $T$, due to historical sampling data in the real-time split-window. In Fig. 6a, the signal area (U-shaped curve) is $t_s = 501 \sim 740$, so we can calculate the actual duration of signals as $740 - 501 + 1 - T = 0$. Therefore, we can speculate that there are instantaneous signals occurring at $t_s = 501$.

In order to find out the causes to above signals, we conduct correlation analysis for each factor. In Fig. 6b, when we augment load data of bus 117, $\kappa_{\text{MSR}}$ of the signal area ($t_s = 501 \sim 740$) decreases dramatically (from 0.7812 to 0.2998), below the inner circle (0.5123); it indicates strong correlations between the load of bus 117 and the system status. On the other hand, in Fig. 6c, $\kappa_{\text{MSR}}$ remains inside the ring throughout the signal area; it indicates poor correlations between the load of bus 54 and the system status.

Besides, we can also determine the correlations by $f_{\text{KDE}} - \lambda$ curves in Fig. 4 and Fig. 5. In Fig. 4, at $t_s = 620$ (inside the signal area), when we augment load data of bus 117, $f_{\text{KDE}}$ deviates from $f_{\text{MP}}$; it indicates strong correlations between the load of bus 117 and the system status. However, in Fig. 5, $f_{\text{KDE}}$ with the load of bus 54 accords with $f_{\text{MP}}$ at $t_s = 620$; it indicates poor correlations between the load of bus 54 and the system status.

As a result, we deduce that the load of bus 117, but not bus 54, is the cause for instantaneous signals at $t_s = 501$. This analysis result accords with assumed signals in Tab. III. In fact, we only add signals to the active load of bus 117. Specifically, the active load of bus 117 increases from 20 MW to 120 MW right at $t_s = 501$.

*2) Case 2 - Peak and dip signals in the load of bus 117:*

In case 2, assumed signals for each factor are shown in Tab. IV. The correlation analysis results are shown in Fig. 7.

TABLE IV: Assumed Signals for Each Factor in Case 2

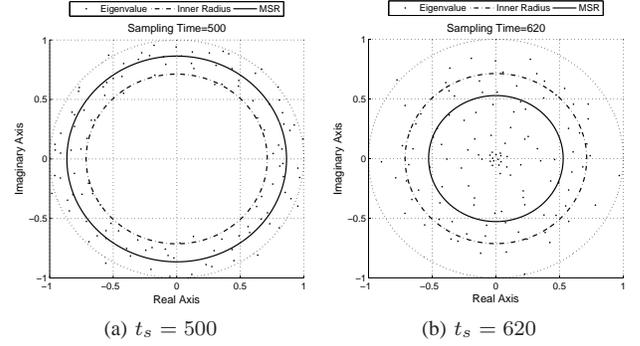| Bus | Sampling Time | Active Load(MW) |
|---|---|---|
| | $t_s = 1 \sim 300$ | 60.0 |
| | $t_s = 301 \sim 350$ | 120.0 |
| 117 | $t_s = 351 \sim 650$ | 60.0 |
| | $t_s = 651 \sim 700$ | 20.0 |
| | $t_s = 701 \sim 1000$ | 60.0 |
| Others | $t_s = 1 \sim 1000$ | Unchanged |



Fig. 1: Eigenvalue distributions of standard matrix products in case 1: the data source is the status matrix.
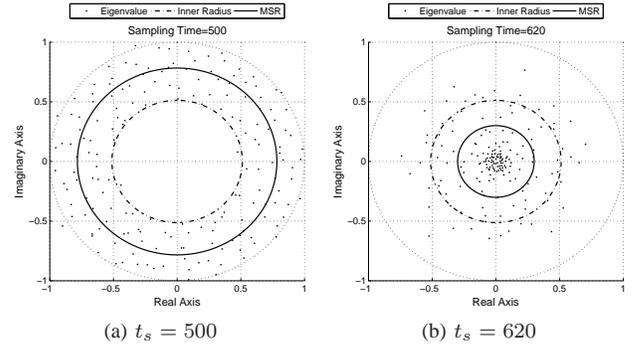


Fig. 2: Eigenvalue distributions of standard matrix products in case 1: the data source is the augmented matrix, including load data of bus 117.
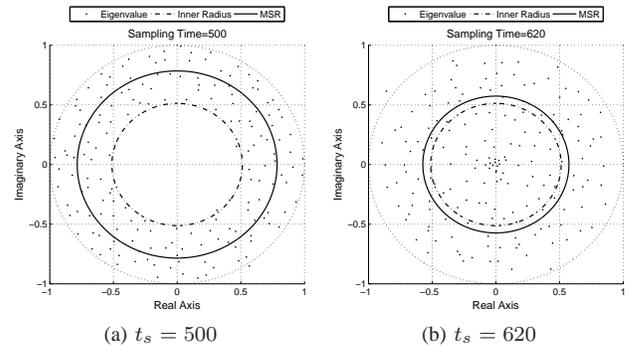


Fig. 3: Eigenvalue distributions of standard matrix products in case 1: the data source is the augmented matrix, including load data of bus 54.

In Fig. 7a, based on the $\kappa_{\text{MSR}} - t$ curve, we can detect signals below:

I. From $t_s = 301$ to $t_s = 590$, the $\kappa_{\text{MSR}} - t$ curve is U-shaped beyond the predict ring. It indicates that the signals occur at $t_s = 301$, and the signal area is $t_s = 301 \sim 590$.

II. From $t_s = 650$ to $t_s = 940$, the $\kappa_{\text{MSR}} - t$ curve is U-shaped beyond the predict ring. It indicates that the signals occur at $t_s = 650$, and the signal area is $t_s = 650 \sim 940$.

In consideration of the delayed effect of a signal to MSR, we can calculate actual durations of above signals as $590 - 301 + 1 - T = 50$ and $940 - 651 + 1 - T = 50$. Therefore, we
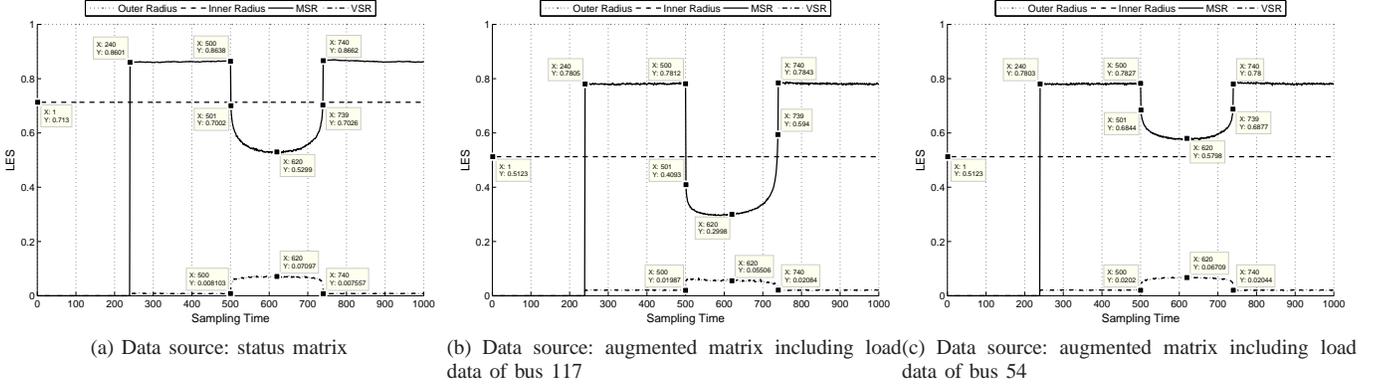
(a) Data source: status matrix     (b) Data source: augmented matrix including load data of bus 117     (c) Data source: augmented matrix including load data of bus 54

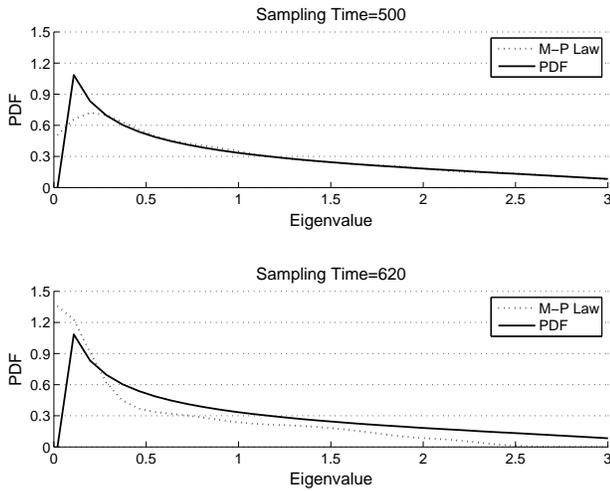Fig. 6: $\kappa_{\mathrm{MSR}} - t$ curves of standard matrix products in Case 1



Fig. 4: $\kappa_{\mathrm{KDE}} - \lambda$ curves of standard matrix products in case 1: the data source is the augmented matrix, including load data of bus 117
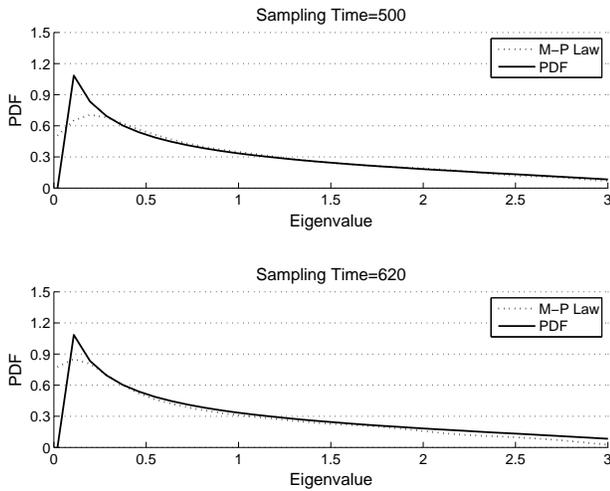


Fig. 5: $\kappa_{\mathrm{KDE}} - \lambda$ curves of standard matrix products in case 1: the data source is the augmented matrix, including load data of bus 54

can speculate that there are continuous signals occurring at

$t_s = 301 \sim 351$ and $t_s = 651 \sim 701$.

In Fig. 7b, when we augment load data of bus 117, $\kappa_{\mathrm{MSR}}$ inside two signal areas declines remarkably and deviates from the predicted ring (from 0.7858 to 0.3538 and from 0.7830 to 0.3846); it indicates strong correlations between the load of bus 117 and the system status. On the other hand, in Fig. 7c, when we augment load data of bus 54, $\kappa_{\mathrm{MSR}}$ inside both signal areas remains in the predicted ring throughout the signal area; it indicates poor correlations between the load of bus 54 and the system status. As a result, we deduce that the load of bus 117, but not bus 54, is the cause for continuous signals during $t_s = 301 \sim 351$ and $t_s = 651 \sim 701$.

Above analyses accord with assumed signals in Tab. IV. In this case, we only add signals to the active load of bus 117. To be specific, the active load of bus 117 has a peak during $t_s = 301 \sim 351$, and a dip during $t_s = 651 \sim 701$.

*3) Case 3 - gradual signals in the load of bus 54:*

In case 3, assumed signals for each factor are shown in Tab. V. The correlation analysis results are shown in Fig. 8.

TABLE V: Assumed Signals for Each Factor in Case 3

| Bus | Sampling Time | Active Load(MW) |
|---|---|---|
| | $t_s = 1 \sim 300$ | 113.0 |
| | $t_s = 301 \sim 350$ | 135.6 |
| | $t_s = 351 \sim 400$ | 158.2 |
| | $t_s = 401 \sim 450$ | 180.8 |
| | $t_s = 451 \sim 500$ | 203.4 |
| 54 | $t_s = 501 \sim 550$ | 226.0 |
| | $t_s = 551 \sim 600$ | 248.6 |
| | $t_s = 601 \sim 650$ | 271.2 |
| | $t_s = 651 \sim 700$ | 293.8 |
| | $t_s = 701 \sim 1000$ | 316.4 |
| Others | $t_s = 1 \sim 1000$ | Unchanged |

In Fig. 8a, the U-shaped curve is from $t_s = 301$ to $t_s = 940$. It indicates signals occurring at $t_s = 301$. In consideration of the delayed effect of a signal to MSR, we can calculate actual durations of above signals as $940 - 301 + 1 - T = 400$. Therefore, we can speculate that there are continuous signals occurring at $t_s = 301 \sim 701$.

In Fig. 8c, when we augment load data of bus 54, $\kappa_{\mathrm{MSR}}$ inside the signal area declines remarkably and deviates from the predicted ring (from 0.7803 to 0.4227); it indicates strong correlations between the load of bus 54 and the system status.

On the other hand, in Fig. 8b, when we augment load data of bus 117, $\kappa_{\mathrm{MSR}}$ remains in the predicted ring throughout the signal area; it indicates poor correlations between the load of bus 117 and the system status. As a result, we deduce that the load of bus 54, but not bus 117, is the cause for continuous signals during $t_s = 301 \sim 701$.

Above analyses is in accordance with assumed signals in Tab. V. In this case, we only add signals to the active load of bus 54. In detail, the active load of bus 54 increase gradually during $t_s = 301 \sim 701$.

### B. Correlation Analysis for Multiple Factors

In case 4, assumed signals for each factor are shown in Tab. VI. The correlation analysis results are shown in Fig. 9.

TABLE VI: Assumed Signals for Each Factor in Case 4

| Bus | Sampling Time | Active Load(MW) |
|---|---|---|
| 117 | $t_s = 1 \sim 300$ | 60.0 |
| | $t_s = 301 \sim 350$ | 120.0 |
| | $t_s = 351 \sim 650$ | 60.0 |
| | $t_s = 651 \sim 700$ | 20.0 |
| | $t_s = 701 \sim 1000$ | 60.0 |
| 54 | $t_s = 1 \sim 300$ | 113.0 |
| | $t_s = 301 \sim 350$ | 135.6 |
| | $t_s = 351 \sim 400$ | 158.2 |
| | $t_s = 401 \sim 450$ | 180.8 |
| | $t_s = 451 \sim 500$ | 203.4 |
| | $t_s = 501 \sim 550$ | 226.0 |
| | $t_s = 551 \sim 600$ | 248.6 |
| | $t_s = 601 \sim 650$ | 271.2 |
| | $t_s = 651 \sim 700$ | 293.8 |
| | $t_s = 701 \sim 1000$ | 316.4 |
| Others | $t_s = 1 \sim 1000$ | Unchanged |

Based on the $\kappa_{\mathrm{MSR}} - t$ curve in Fig. 9a, two kinds of signals are detected:

I. Two U-shaped curves are found during $t_s = 301 \sim 590$ and $t_s = 651 \sim 940$. Referring to the analysis in case 2, it indicates two continous signals during $t_s = 301 \sim 351$ and $t_s = 651 \sim 701$.

I. The Third curve are found during $t_s = 301 \sim 940$. Referring to the analysis in case 3, it indicates continuous signals during $t_s = 301 \sim 700$.

During the first two signal areas ($t_s = 301 \sim 590$ and $t_s = 651 \sim 940$), when we augment the load of bus 117, $\kappa_{\mathrm{MSR}}$ deviates from Ring Law, shown in Fig. 9b. During the third signal area ($t_s = 301 \sim 940$), when we augment the load of bus 54, the deviation of $\kappa_{\mathrm{MSR}}$ is shown in Fig. 9c. As a result, we can achieve following speculations:

I. The load of bus 117 affects the system status during $t_s = 301 \sim 351$ and $t_s = 651 \sim 701$.

II. The load of bus 54 affects the system status during $t_s = 301 \sim 700$.

These speculations are validated by Table. VI.

### C. More Discussions in Details

Through above four cases, the effectiveness and performance of the correlation analysis method is verified under different scenarios. However, there are still some interesting details left.

First, we can observe that $\kappa_{\mathrm{MSR}} - t$ curves in case 4 are approximately superposed by corresponding curves in case 2 and case 3. In fact, the assumed signals in case 4 are combined with those in case 2 and case 3. In conclusion, there should be some relationships between phenomena in power systems, and LES in mathematics.
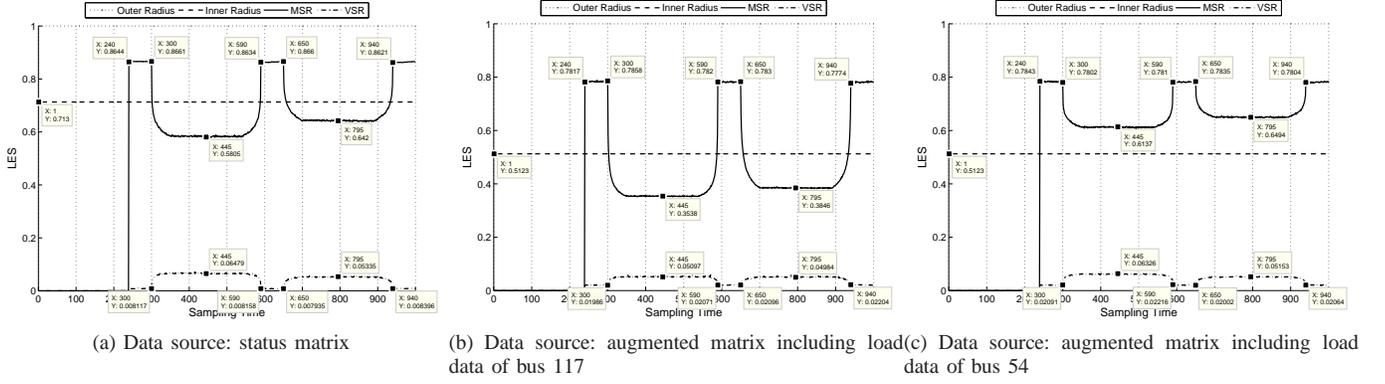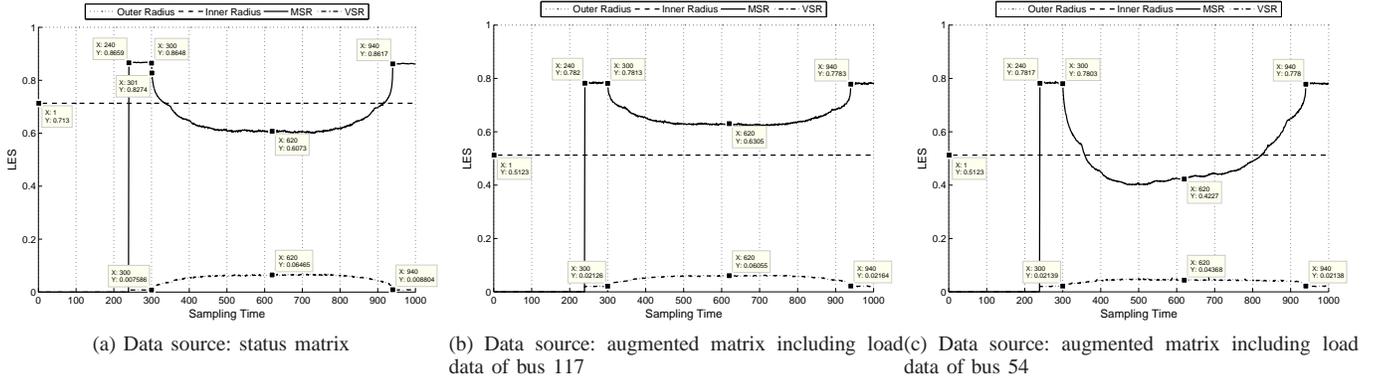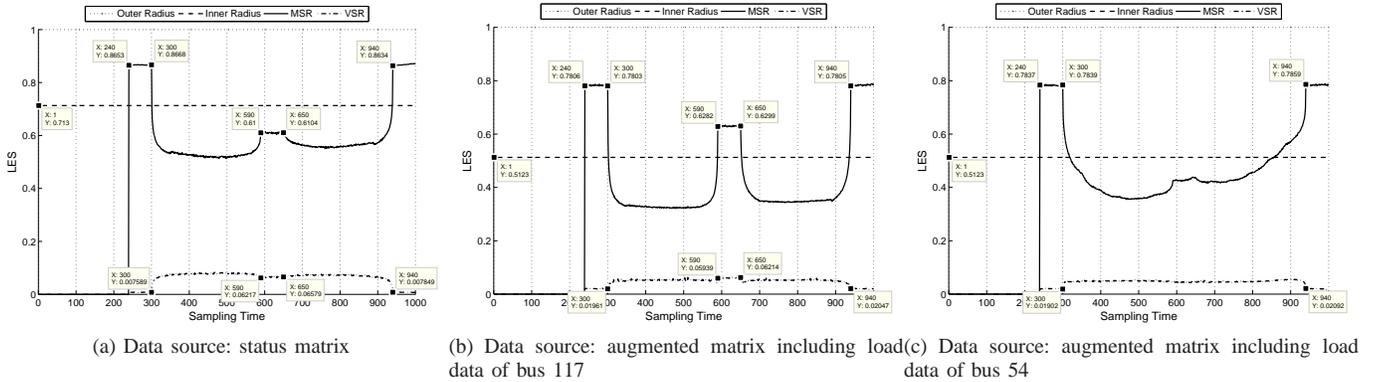
Secondly, when data become more correlated, the VSR increases as well. This phenomenon can be explained based on the eigenvalue distributions of standard matrix products. In Fig. 1a, Fig. 2a and Fig. 3a, at $t_s = 500$, the eigenvalues distribute between the outer circle and the inner circle, conforming to the Ring Law; it indicates poor correlations in data. During the signal area (such as $t_s = 620$), when we augment data of correlated factors, all the eigenvalues gather towards the circle center, shown in Fig. 2b. When we augment data of irrelevant factors, only some of the eigenvalues gather towards the circle center, shown in Fig. 3b.

Thirdly, to represent random fluctuations in the system and measuring errors in sampling data, white noise is introduced into both status data and factor data throughout the simulation. Observations in case studies indicate that $\kappa_{\mathrm{MSR}}$ does not change dramatically in a system dominated by white noise. Therefore, $\kappa_{\mathrm{MSR}}$ is a reliable indicator to identify signals from random fluctuations and measuring errors.

### V. CONCLUSION

This paper proposes a data-driven method to reveal the correlations between factors and the power system status. First, for each factor, we construct an augmented matrix as a data source, combining status data and factors data reasonably. Secondly, to conduct real-time analysis, we use a specific split-window to obtain the raw data matrix at each sampling time. Then, we conduct correlation analysis for the raw data matrix based on random matrix theory (RMT). The mean spectral radius (MSR), a kind of linear eigenvalue statistics (LESs), is calculated as a correlation indicator; the kernel density estimation (KDE) is used as a assisted tool to reveal data correlations. The proposed method is direct, robust against bad data and universal to varieties of factors. Case studies demonstrate the effectiveness of the method for both single factor and multiple factors.

The current work is only a preliminary exploration of correlation analysis based on RMT. Much more researches are needed along this direction. The degree of correlations needs to be further quantified. For example, in a study period, we can use the integration of MSR to quantify the correlations. Besides, we can use data of subarea to construct the basic matrix; in this way we can fix the data missing problems. Furthermore, the proposed method can be used to reveal the correlation between any two types of variables, as long as combining their data reasonably in the data source. Our method can be applied for specific issues in power systems, such as voltage stability, weak buses identification, abnormal and fault diagnosis. Combinations with model-based methods and data processing methods will improve the performance, and uncover more connections between electrical phenomena and statistical characteristics.

(a) Data source: status matrix

(b) Data source: augmented matrix including load data of bus 117

(c) Data source: augmented matrix including load data of bus 54

Fig. 7: $\kappa_{\mathrm{MSR}} - t$ curves of standard matrix products in Case 2



(a) Data source: status matrix

(b) Data source: augmented matrix including load data of bus 117

(c) Data source: augmented matrix including load data of bus 54

Fig. 8: $\kappa_{\mathrm{MSR}} - t$ curves of standard matrix products in Case 3



(a) Data source: status matrix

(b) Data source: augmented matrix including load data of bus 117

(c) Data source: augmented matrix including load data of bus 54

Fig. 9: $\kappa_{\mathrm{MSR}} - t$ curves of standard matrix products in Case 4

# APPENDIX A
# RANDOM MATRIX THEORY

## A. Ring Law

Let $\widetilde{\mathbf{X}} \in \mathbb{C}^{N \times T}$ be a standard non-Hermitian random matrix, whose entries are independent identically distributed (i.i.d.) variables with

$$\mu(\widetilde{\mathbf{x}}_i) = 0, \sigma^2(\widetilde{\mathbf{x}}_i) = 1 \quad (i = 1, 2, \ldots, N) \tag{15}$$

where $\widetilde{\mathbf{x}}_i = (\widetilde{x}_{i,1}, \widetilde{x}_{i,2}, \ldots, \widetilde{x}_{i,T})$. For multiple standard non-Hermitian random matrices $\widetilde{\mathbf{X}}_i$ $(i = 1, 2, \ldots, L)$, we define the

matrix product as

$$\hat{\mathbf{Z}} = \prod_{i=1}^{L} \widetilde{\mathbf{X}}_{u,i} \tag{16}$$

where $\widetilde{\mathbf{X}}_{u,i}$ is the singular value equivalent of $\widetilde{\mathbf{X}}_i$. The matrix product $\hat{\mathbf{Z}}$ can be converted to the standard matrix product $\widetilde{\mathbf{Z}}$, whose $\sigma^2(\widetilde{\mathbf{z}}_i) = \frac{1}{N}$ in each row. According to Ring Law [20, 21], the ESD of $\widetilde{\mathbf{Z}}$ converges almost surely to the limit with a probability density function (PDF) given by

$$f_{\mathrm{RL}}(\lambda_{\widetilde{\mathbf{Z}}}) = \begin{cases} \frac{1}{\pi c L}|\lambda_{\widetilde{\mathbf{Z}}}|^{\frac{2}{L}-2} & (1-c)^{\frac{L}{2}} \leq |\lambda_{\widetilde{\mathbf{Z}}}| \leq 1 \\ 0 & \text{otherwise} \end{cases} \tag{17}$$

as $N, T \to \infty$ with a constant ratio $c = \frac{N}{T} \in (0, 1]$. On the complex plane of eigenvalues, the inner circle radius is $(1-c)^{\frac{L}{2}}$ and the outer circle radius is unity.

### B. Marchenko-Pastur Law

Let $\mathbf{X} = \{x_{i,j}\}_{N \times T}$ be a random matrix, whose entries are i.i.d. variables with

$$\mu(\mathbf{x}_i) = 0, \sigma^2(\mathbf{x}_i) = d \quad (i = 1, 2, \ldots, N) \tag{18}$$

where $d < \infty$ is a constant, and $\mathbf{x}_i = (x_{i,1}, x_{i,2}, \ldots, x_{i,T})$. The sample covariance matrix of $\mathbf{X}$ is defined as

$$\mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^H \tag{19}$$

According to M-P Law [22, 23], the ESD of $\mathbf{S}$ converges to the following PDF

$$f_{\mathrm{MP}}(\lambda_{\mathbf{S}}) = \begin{cases} \frac{1}{2\pi c d \lambda_{\mathbf{S}}} \sqrt{(b - \lambda_{\mathbf{S}})(\lambda_{\mathbf{S}} - a)} & a \le \lambda_{\mathbf{S}} \le b \\ 0 & \text{otherwise} \end{cases} \tag{20}$$

as $N, T \to \infty$ with a constant ratio $c = \frac{N}{T} \in (0, 1]$, where $a = d(1 - \sqrt{c})^2$, $b = d(1 + \sqrt{c})^2$.

## REFERENCES

[1] M. Kezunovic, L. Xie, and S. Grijalva, "The role of big data in improving power system operation and protection," in *Bulk Power System Dynamics and Control-IX Optimization, Security and Control of the Emerging Power Grid (IREP), 2013 IREP Symposium.* IEEE, 2013, pp. 1–9.

[2] IBM, "The four vs of big data," [EB/OL], http://www.ibmbigdatahub.com/infographic/four-vs-big-data.

[3] R. Qiu and P. Antonik, *Smart Grid and Big Data.* John Wiley and Sons, 2015.

[4] S. ORourke and D. Renfrew, "Central limit theorem for linear eigenvalue statistics of elliptic random matrices," *Journal of Theoretical Probability*, pp. 1–71, 2014.

[5] X. Miao and D. Zhang, "The opportunity and challenge of big data's application in distribution grids," in *Electricity Distribution (CICED), 2014 China International Conference on*, Sept 2014, pp. 962–964.

[6] J. Zhang and M. L. Huang, "5ws model for big data analysis and visualization," in *Computational Science and Engineering (CSE), 2013 IEEE 16th International Conference on.* IEEE, 2013, pp. 1021–1028.

[7] M. Mayilvaganan and M. Sabitha, "A cloud-based architecture for big-data analytics in smart grid: A proposal," in *Computational Intelligence and Computing Research (ICCIC), 2013 IEEE International Conference on.* IEEE, 2013, pp. 1–4.

[8] X. He, Q. Ai, R. Qiu, W. Huang, L. Piao, and H. Liu, "A big data architecture design for smart grids based on random matrix theory," *Smart Grid, IEEE Transactions on*, vol. PP, no. 99, pp. 1–1, 2015.

[9] X. He, R. C. Qiu, Q. Ai, and X. Xu, "An unsupervised learning method for early event detection in smart grid with big data," *arXiv preprint arXiv:1502.00060*, 2015.

[10] X. He, Q. Ai, J. Ni, L. Piao, Y. Xu, X. Xu *et al.*, "3d power-map for smart grids—an integration of high-dimensional analysis and visualization," *arXiv preprint arXiv:1503.00463*, 2015.

[11] S. Lian, S. Morii, T. Ishii, and S. Kawamoto, "Voltage stability and sensitivity analysis considering dynamic load for smart grid," in *Innovative Smart Grid Technologies (ISGT), 2010.* IEEE, 2010, pp. 1–6.

[12] P. Parinya, A. Sangswang, K. Kirtikara, D. Chenvidhya, S. Naetiladdanon, and C. Limsakul, "A study of impact of wind power to power system stability using stochastic stability index," in *Circuits and Systems (ISCAS), 2014 IEEE International Symposium on.* IEEE, 2014, pp. 2656–2659.

[13] L. Xie, Y. Chen, and P. R. Kumar, "Dimensionality reduction of synchrophasor data for early event detection: Linearized analysis," *Power Systems, IEEE Transactions on*, vol. 29, no. 6, pp. 2784–2794, 2014.

[14] D. M. Eltigani, K. Ramadan, and E. Zakaria, "Implementation of transient stability assessment using artificial neural networks," in *Computing, Electrical and Electronics Engineering (ICCEEE), 2013 International Conference on.* IEEE, 2013, pp. 659–662.

[15] R. Qiu and M. Wicks, *Cognitive networked sensing and big data.* Springer, 2014.

[16] C. Zhang and R. C. Qiu, "Data modeling with large random matrices in a cognitive radio network testbed: Initial experimental demonstrations with 70 nodes," *arXiv preprint arXiv:1404.3788*, 2014.

[17] D. L. Donoho *et al.*, "High-dimensional data analysis: The curses and blessings of dimensionality," *AMS Math Challenges Lecture*, pp. 1–32, 2000.

[18] I. Jana, K. Saha, and A. Soshnikov, "Fluctuations of linear eigenvalue statistics of random band matrices," *arXiv preprint arXiv:1412.2445*, 2014.

[19] R. D. Zimmerman and C. E. Murillo-Sánchez, "Matpower 4.1 users manual," *Power Systems Engineering Research Center, Cornell University, Ithaca, NY*, 2011.

[20] A. Guionnet, M. Krishnapur, and O. Zeitouni, "The single ring theorem," *arXiv preprint arXiv:0909.2214*, 2009.

[21] F. Benaych-Georges and J. Rochet, "Outliers in the single ring theorem," *Probability Theory and Related Fields*, pp. 1–51, 2013.

[22] R. C. Qiu, Z. Hu, H. Li, and M. C. Wicks, *Cognitive radio communication and networking: Principles and practice.* John Wiley & Sons, 2012.

[23] V. A. Marčenko and L. A. Pastur, "Distribution of eigenvalues for some sets of random matrices," *Sbornik: Mathematics*, vol. 1, no. 4, pp. 457–483, 1967.