

# A Novel Approach for Phase Identification in Smart Grids Using Graph Theory and Principal Component Analysis

Satya Jayadev P, Aravind Rajeswaran,  
Nirav P Bhatt, Ramkrishna Pasumarthy \*

June 8, 2016

## Abstract

Consumers with low demand, like households, are generally supplied single-phase power by connecting their service mains to one of the phases of a distribution transformer. The distribution companies face the problem of keeping a record of consumer connectivity to a phase due to uninformed changes that happen. The exact phase connectivity information is important for the efficient operation and control of distribution system. We propose a new data driven approach to the problem based on Principal Component Analysis (PCA) and its Graph Theoretic interpretations, using energy measurements in equally timed short intervals, generated from smart meters. We propose an algorithm for inferring phase connectivity from noisy measurements. The algorithm is demonstrated using simulated data for phase connectivities in distribution networks.

**Keywords:** Power Distribution Networks, Smart Meters, Phase Identification, Principle Component Analysis

## 1 Introduction

Electrical power is generally transmitted in three phases due to technical and economic advantages over single-phase transmission [1]. At the distribution level, 3-phase power is distributed to the consumers either in one phase or three phases depending on the customer demand. For single-phase consumers, the service mains is connected to one of the phases of the distribution transformer, A phase, B phase or C phase.

---

\*Satya Jayadev P and Ramkrishna Pasumarthy are with Department of Electrical Engineering, and Nirav P Bhatt and Aravind Rajeswaran are with Department of Chemical Engineering, Indian Institute of Technology Madras, India. [ee15s059@ee.iitm.ac.in](mailto:ee15s059@ee.iitm.ac.in), [ramkrishna@ee.iitm.ac.in](mailto:ramkrishna@ee.iitm.ac.in), [niravbhatt@iitm.ac.in](mailto:niravbhatt@iitm.ac.in) [aravindr@mail.iitm.ac.in](mailto:aravindr@mail.iitm.ac.in)

Accurate phase-load connectivity information is required for balancing loads on all the phases. Balanced loads mitigate the problem of overloading on a phase. Also technical losses can be reduced leading to better efficiency in distribution. Unbalanced loads on phases result in voltage imbalance affecting many consumers, especially with rotating electric machines. From the control point of view, the controllability limits of the power system are affected by voltage imbalance. Balanced loads on phases ensures voltage balance in the three phases.

Accurate phase connectivity information also helps in detection and localization of non-technical losses and state estimation in power distribution system. Hence, reliable phase connectivity information is essential for efficient monitoring and optimization of distribution networks [2].

The problem faced by distribution companies is in maintaining an accurate record of the phase-load connectivity. This information might not be accurately available at all times because of changes that take place due to repairs and maintenance. Also the consumers might have the facility to switch between phases and they do so when phase tripping occurs. The distribution utility is uninformed of such changes. There are techniques such as signal injection and manual verification to determine phase connectivity but utilities refrain from them due to high costs and possible inaccuracies [3].

With the advent of smart grid technologies, distribution companies are installing smart meters at important nodal points in the network including feeders, transformers and consumer service mains. These meters can communicate readings to a central data centre with greater frequency and sometimes in real time.

In this paper, we will deal with identifying the phase connectivity of loads using the smart meter data. We will propose an algorithm by combining graph theory and principal component analysis based on the principle of energy conservation. We also take into account the noise in the data arising due to technical losses, errors in smart meter readings and errors due to imperfect time synchronization of smart meter clocks.

## 2 Related Work

The problem of phase identification is gaining recognition in recent times with the large scale penetration of smart grids. Smart grids have lead to the search for new methods of inferring the connectivity. Chen et al. [4] proposed a phase identification device for phase measurement of underground transformers. Zhiyu designed a signal injector device that can be used for phase identification [5]. However the additional hardware and staff required for these devices to work, makes them costly alternatives. Dilek et al. [6] proposed a search algorithm to determine phase information using power flow analysis and load data but it ignores the noise and uncertainty in data. Arya et al. gives an approach to infer phase connectivity from time series of power measurements using mixed integer programming (MIP) [7] . The MIP solver is time intensive in arriving at the

solution. Pezeshki and Wolfs presented a technique to identify the phases based on cross-correlation method using the time series of voltage measurements [8]. A. Tom proposed a linear regression based algorithm which correlates between consumer voltage and substation voltage [9]. It requires the Geographical Information System (GIS) model which may not always be available. A method using data obtained from micro synchrophasors (uPMU) apart from the voltage magnitudes data is developed [10]. It proposed a brute-force search algorithm based on linear programming optimization structure, to determine phase connectivity with certain constraints on voltage magnitudes and phase angles.

We follow a method similar to that proposed in [11], in which a water distribution network is reconstructed from flow measurements.

### 3 Preliminaries

#### 3.1 Graph Theory

The incidence matrix ( $\mathbf{A}$ ) of a graph describes the incidence of edges on nodes and is defined as follows for a directed graph:

$$A_{ij} = \begin{cases} +1 & \text{if edge } j \text{ enters node } i \\ -1 & \text{if edge } j \text{ leaves node } i \\ 0 & \text{if edge } j \text{ is not incident on } i \end{cases}$$

**Proposition 1 ( Refer Theorem 8 of Chapter 3 in [12])**

*A directed graph (or a directed forest) can be uniquely constructed from an incidence matrix, provided there are no self loops.*

We make use of this proposition to determine the connectivity from incidence matrix.

#### 3.2 Principal Component Analysis (PCA)

PCA is one of the most widely used techniques in multivariate statistical analysis. PCA provides the best approximation of linear model between a set of variables, of which some are dependent on the other. A set of samples corresponding to the variables, measured at different time instances can be used to apply PCA. The linear model can be estimated even in the presence of Gaussian noise [13].

Let  $\mathbf{Z}$  be the  $(n \times N)$ -dimensional matrix obtained by stacking  $n$  variables of  $N$  samples each. Let  $\mathbf{z}_j$  be the vector of values of  $n$  variables in the  $j^{th}$  sample. These variables are linearly related and described by the following model:

$$\mathbf{C} \mathbf{z}_j = \mathbf{0} \tag{1}$$

where  $\mathbf{C}$  can be referred to as the Constraint matrix of  $(n_d \times n)$  dimension where  $n_d$  is the number of dependent variables. We call it the constraint matrix because it gives the physical constraints of the system.

Note that the data lies in the subspace orthogonal to the space spanned by the rows of the constraint matrix. So to represent the constraint matrix, we need a set of basis vectors orthogonal to the subspace in which the data lies. This is obtained from the eigenvectors of the Covariance matrix corresponding to the least  $m$  eigenvalues, where  $m$  is the number of dependent variables [11]. The Covariance matrix  $\mathbf{S}_z$  is

$$\mathbf{S}_z = \mathbf{Z}\mathbf{Z}^T \quad (2)$$

The eigenvectors of the covariance matrix can be determined using Singular Value Decomposition (SVD) of the data matrix. SVD of the data matrix can be written as:

$$svd(\mathbf{Z}) = \mathbf{U}_1\mathbf{S}_1\mathbf{V}_1^T + \mathbf{U}_2\mathbf{S}_2\mathbf{V}_2^T \quad (3)$$

where  $\mathbf{U}_1$  are the set of orthonormal eigenvectors corresponding to the  $(n - m)$  largest eigenvalues of  $\mathbf{S}_z$  while  $\mathbf{U}_2$  are the orthogonal eigenvectors corresponding to the smallest  $m$  eigenvalues of  $\mathbf{S}_z$ . It has been shown that  $S_R(\mathbf{U}_2^T) \sim S_R(\mathbf{C})$  where  $S_R(\cdot)$  indicates the subspace spanned by the rows of  $(\cdot)$  matrix [14]. Then  $\mathbf{U}_2^T$  satisfies the following relationship:

$$\mathbf{U}_2^T \mathbf{z} = 0 \quad (4)$$

where  $\mathbf{z} = [z_1, z_2, \dots, z_n]^T$  is the vector of  $n$  variables. It is to be observed that the constraint matrix suffers from rotational ambiguity.

$$\mathbf{Q}\mathbf{U}_2^T \mathbf{z} = \mathbf{0} \quad (5)$$

where  $\mathbf{Q}$  is a non-singular matrix. So the estimated constraint matrix may not represent the physical relationships even though the correct subspace has been extracted. The estimated constraint matrix,  $\mathbf{C}$ , can at best be a basis for the row space of true constraint matrix. It is shown that by partitioning  $\mathbf{C}$  into columns corresponding to dependent variables and those of independent variables, a matrix  $\mathbf{R}$  (Regression Matrix), which is unique to the system, can be obtained [11].  $\mathbf{R}$  can be computed as:

$$\mathbf{R} = -\mathbf{C}_d^{-1}\mathbf{C}_i \quad (6)$$

$$\mathbf{z}_d = \mathbf{R}\mathbf{z}_i \quad (7)$$

where  $\mathbf{C}_d$  are the columns of  $\mathbf{C}$  corresponding to dependent variables,  $\mathbf{C}_i$  are the columns of  $\mathbf{C}$  corresponding to independent variables,  $\mathbf{z}_i$  is the vector of independent variables and  $\mathbf{z}_d$  is the vector of dependent variables.

Many variants to PCA have been developed which pertain to different cases of error covariance matrix [14]. Hence, this approach can be applied to noisy data.

### 3.3 Integrating Graph Theory & PCA

In our problem, we deal with a forest of directed trees with only one parent node each and many child nodes. The challenge is to determine which child nodes are connected to which parent node. It can be observed that the sub-matrix of the incidence matrix of the forest, with rows corresponding to the parent nodes, is sufficient to infer the connectivity. This connectivity is unique as described in Proposition 1. The incidence matrix can be written as:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_d \\ \mathbf{A}_i \end{bmatrix} \quad (8)$$

where  $\mathbf{A}_d$  are the rows corresponding to parent nodes and  $\mathbf{A}_i$  are the rows corresponding to child nodes. It is pertinent to note that  $\mathbf{A}_d$  only comprises of -1 and 0 as its elements. Further, it is to be noted that each column of  $\mathbf{A}_d$  contains only one -1 and rest zeros.

Mathematically, the parent nodes can be considered dependent variables and child nodes to be independent variables. Hence it can be easily verified that the regression matrix  $\mathbf{R}$  which regresses the dependent variables on the independent variables is in fact the matrix  $-\mathbf{A}_d$  and the uniqueness of  $\mathbf{R}$  makes it comparable element-wise to  $-\mathbf{A}_d$ .

Now we formulate the problem and present an algorithm which determines the regression matrix from the data pertaining to our network.

## 4 Problem Formulation

The phase connectivity network can be represented as a forest with three trees as its components, each tree corresponding to a phase. The child nodes of a tree represent the single-phase meters connected to the corresponding phase.<sup>1</sup> The 3-phase transformer meter is considered equivalent to three single-phase meters by taking the readings corresponding to each phase, separately. Therefore, the root of each tree corresponds to a phase of the transformer meter as shown in Fig. 1.

In our approach, energy measurements in watt-hour (Wh), in equal time intervals, generally 15 minutes or 30 minutes, are collected from all the consumer meters and transformer meter to form the data matrix. Let  $\mathbf{Z}$  be the data matrix,  $n$  be the number of meters (or nodes),  $n_d$  be the number of parent nodes and  $n_i$  be the number of child nodes,

$$\mathbf{Z} = [ z_{ij} ]_{(n \times N)} \quad (9)$$

where  $z_{ij}$  is the energy measurement corresponding to the  $i^{th}$  meter in  $j^{th}$  time interval. By definition,

$$n_d = 3 \quad (10)$$

$$n_i = n - 3 \quad (11)$$

---

<sup>1</sup>This notation can be extended to 3-phase loads by considering them as three separate loads and taking the per phase readings.

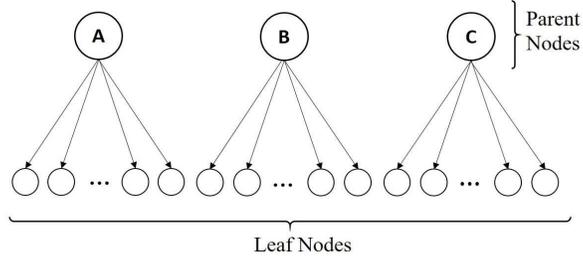


Figure 1: Graph Representation of Phase Connectivity

The principle of conservation of energy implies that the energy supplied by each phase of transformer is equal to sum of energies consumed by all the consumers connected to that phase of the transformer, in a given time interval. In graph theory view point, this principle implies that the parent node (meter) reading is equal to the sum of child nodes (meters) readings. So there exists a linear relationship between the nodal measurements which PCA exploits to determine the constraint matrix and consequently the regression matrix.

$$z_j^k = \sum_{i=1}^{n_k} z_{ij}, \forall k = A, B, C \quad (12)$$

where  $z_j^k$  is the energy measured in phase  $k$  in the  $j^{th}$  time interval,  $n_k$  is the number of single-phase meters connected to phase  $k$  and  $z_{ij}$  is the energy measured in  $i^{th}$  meter in the  $j^{th}$  time interval.

In practice, due to losses and other errors, this relationship is only approximate. We will try to infer the exact connectivity using the distributions of the noise components. The inclusion of these components in the formulation is discussed later.

Apart from the basic assumptions for PCA to work [13], we make the following assumptions:

1. There is no theft of electricity and un-metered loads are generally estimated and so are separated from the readings before applying the algorithm.
2. The clocks of all the meters in the network are time synchronized.
3. The loads do not change phases while the  $N$  measurements are recorded.

## 5 An Approach for Phase Identification

In this section, two cases for phase identification using data are considered: (A) Noiseless data set and (B) Noisy data set.

## 5.1 Noiseless Case

With no noise, Eq. (12) holds exactly, and we get the exact incidence sub-matrix, except for some numerical residues. The connectivity is determined by the following steps:

1. PCA is applied and the data matrix is decomposed using Eq. (3)
2. As we know that there are only three dependent variables, we take the eigenvectors corresponding to least three eigenvalues of  $\mathbf{S}_z$  to get the constraint matrix  $\mathbf{C}_{(3 \times n)}$ .

$$\mathbf{C} = \mathbf{U}_2^T \quad (13)$$

3. The columns corresponding to dependent and independent variables are separated and the regression matrix  $\mathbf{R}_{(3 \times n_i)}$  is calculated using Eq. (6).
4. The regression matrix is rounded off to truncate any possible numerical residues and the resultant is the incidence sub-matrix from which the phase connectivity can be inferred.

Now the question is how many readings are required to infer the connectivity uniquely. We answer this question in Proposition 2 by establishing the lower bound of  $N$  for any potential data set.

### Proposition 2

*Let  $N$  be the number of readings (or samples). In the noiseless case, the minimum number of linearly independent readings required to infer the graph uniquely is equal to  $n_i$ , the number of independent nodes (or consumers meters) in the network.*

$$N \geq n_i \quad (14)$$

*Proof:* With no noise, all the samples clearly lie in the subspace represented by the constraint matrix and PCA gives a linear relationship between the variables. The linear independence of readings ensures that they span the target subspace. In such a case, the regression matrix can be written as:

$$\mathbf{R} \mathbf{Z}^i = \mathbf{Z}^d \quad (15)$$

where  $\mathbf{Z}^d$  are rows of  $\mathbf{Z}$  corresponding to the dependent variables and  $\mathbf{Z}^i$  are the rows of  $\mathbf{Z}$  corresponding to the independent variables. To obtain  $\mathbf{R}$  uniquely,  $\mathbf{Z}^i$  should be a full row rank matrix. By definition of rank,

$$\text{rank}(\mathbf{Z}^i) \leq \min(n_i, N) \quad (16)$$

For  $\mathbf{Z}^i$  to be full rank, Eq. (14) should be satisfied. Hence Proposition 2 is proved. ■

### Corollary 1

*In the noiseless case,  $n_i$  number of linearly independent energy measurements in different time intervals are sufficient to determine the phase connectivity of the loads, uniquely.*

*Proof:* This follows from propositions 1 and 2. ■

## 5.2 Noisy Case

In practice, measurements are noisy and hence it is imperative to account for the various sources of noise in a distribution network. In the problem of phase identification, we account for the technical losses, random errors in smart meter readings, and smart meter clock synchronization errors.

### i. Technical Losses

Major component of technical losses in a distribution network is the copper loss and it is proportional to the square of current in the line. As the voltage is almost constant, the current varies with the load. So in a given time interval, higher the load, higher is the loss and hence the phase meter reading is always greater than the sum of consumer meter readings. This leads to

$$z_{i(m)}^k = z_{i(t)}^k + \sum_{j=1}^{n_k} l_{ij}, \quad \forall i = 1 \text{ to } N \quad (17)$$

where  $z_{i(m)}^k$  is the energy measured with loss included, in phase  $k$  in  $i^{th}$  time interval,  $z_{i(t)}^k$  is true value of energy consumed in phase  $k$  and  $l_{ij}$  is the energy loss in line connecting transformer phase  $k$  and consumer  $j$ , in  $i^{th}$  time interval. Due to this loss, the noise in data is not normally distributed with zero mean. So we need to pre-process the data before applying PCA.

### ii. Random Errors in meter readings

ANSI C12.20-2010, the latest standard for electricity meters, stipulates that electricity meters must be of 0.2 or 0.5 accuracy class [15]. That means the meter reading can be in the range of  $\pm 0.2\%$  of true value for 0.2 accuracy class meter and in the range of  $\pm 0.5\%$  of true value for 0.5 accuracy class meter. Let us consider that all the smart meters in our network are of 0.5 accuracy class. This error can be approximately modelled to be Gaussian with each variable having a different error variance.

$$\mathbf{z}_{j(m)} = \mathbf{z}_{j(t)} + \boldsymbol{\epsilon}_j \quad (18)$$

where  $\mathbf{z}_{j(m)}$  is the vector of measured values of  $n$  variables in  $j^{th}$  time interval,  $\mathbf{z}_{j(t)}$  is the vector of true values of  $n$  variables in  $j^{th}$  time interval and  $\boldsymbol{\epsilon}_j$  is the error in the reading in  $j^{th}$  time interval

$$\boldsymbol{\epsilon}_j \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon) \quad (19)$$

where  $\boldsymbol{\Sigma}_\epsilon$  is the error covariance matrix. As the errors are not correlated,  $\boldsymbol{\Sigma}_\epsilon$  is a diagonal matrix.

### iii. Clock Synchronization errors

Clock synchronization error can also be modelled as Gaussian. The clocks of all the meters may not be synchronized perfectly leading to time intervals of measurements to be varying. The variation is generally in the order of milliseconds.

$$t_{k+1} = t_k + \mathcal{N}(\Delta t, \sigma_t^2) \quad (20)$$

where  $t_k$  is the  $k^{th}$  time instance,  $\Delta t$  is the duration of the time interval and  $\sigma_t^2$  is the variance in time interval. For example, in a 15 minutes time interval, a variance of one second will lead to an error of 0.1% ( $= 1/(15 \times 60)$ ) and even with a drastic change (say 5 times) of load during that second, will lead to an error of only 0.5%.

This error can be formulated similar to the previous one and superimposing them, we get

$$\mathbf{z}_{j(m)} = \mathbf{z}_{j(t)} + \boldsymbol{\epsilon}_j + \boldsymbol{\delta}_j \quad (21)$$

$$\boldsymbol{\epsilon}_j + \boldsymbol{\delta}_j \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\epsilon + \boldsymbol{\Sigma}_\delta) \quad (22)$$

$$\boldsymbol{\Sigma}_e = \boldsymbol{\Sigma}_\epsilon + \boldsymbol{\Sigma}_\delta \quad (23)$$

In the noisy case, we propose some additional steps in the algorithm to handle the noise.

1. The total technical losses can be approximately calculated as:

$$loss_i = \sum_{j=1}^{n_d} z_{ij}^d - \sum_{k=1}^{n_i} z_{ik}^i, \forall i=1 \text{ to } N \quad (24)$$

where  $loss_i$  is the approximate technical loss component in  $i^{th}$  time interval,  $z_{ij}^d$  is the value of the  $j^{th}$  dependent variable in  $i^{th}$  time interval and  $z_{ik}^i$  is the value of the  $k^{th}$  independent variable in  $i^{th}$  time interval.

It is a good approximation because the magnitude of the Gaussian errors is much smaller compared to these losses. As the other technical losses are very small compared to copper losses and as the copper losses are in proportion to load, the technical losses can be separated from the parent nodes in proportion to the readings.

$$\hat{z}_{ij}^d = z_{ij}^d - \frac{loss_i z_{ij}^d}{\sum_{j=1}^{n_d} z_{ij}^d}, \forall i = 1 \text{ to } N \text{ and } j = 1 \text{ to } n_d \quad (25)$$

where  $\hat{z}_{ij}^d$  are the estimated values of phase meter readings, free from losses. We can apply PCA by substituting the measured values with these estimated values because the expectation of the error in these estimated values is approximately zero.

2. Now we deal with Gaussian error by using MLPCA [16]. It is shown that if the errors in variables are not correlated and error variances are known, the constraint matrix corresponding to error free readings can be obtained by scaling the data matrix by standard deviations of corresponding errors [16],[14]. We take the error variances to be equal to 1% of the mean readings for each variable, which is the maximum possible variance of the Gaussian errors described above.

Now, MLPCA is applied as follows:

$$\boldsymbol{\Sigma}_e = \text{diag}(\bar{z}_1, \bar{z}_2, \dots, \bar{z}_n) \quad (26)$$

where  $\bar{z}_i$  is the mean of  $N$  samples of  $i^{\text{th}}$  variable

$$\bar{z}_i = \sum_{j=1}^N z_{ij} \quad (27)$$

Cholesky decomposition of  $\boldsymbol{\Sigma}_e$  is given by

$$\boldsymbol{\Sigma}_e = \mathbf{L}\mathbf{L}^T \quad (28)$$

where  $\mathbf{L}$  is called Cholesky factor and it is a diagonal matrix in this case. The noisy data matrix is transformed as follows:

$$\mathbf{Z}_s = \mathbf{L}^{-1}\mathbf{Z} = \mathbf{L}^{-1}\mathbf{Z}_t + \mathbf{L}^{-1}\mathbf{E} \quad (29)$$

where  $\mathbf{E}$  is the error matrix and  $\mathbf{Z}_t$  is the data matrix with true values. The covariance matrix of the transformed data matrix is

$$\mathbf{S}_{z_s} = \mathbf{Z}_s\mathbf{Z}_s^T \quad (30)$$

By taking expectation of the  $\mathbf{S}_{z_s}$  [14], we get

$$\mathbf{E}(\mathbf{S}_{z_s}) = \mathbf{M}_z + \alpha^2\mathbf{I} \quad (31)$$

where  $\mathbf{M}_z = \mathbf{L}^{-1}\mathbf{Z}_t\mathbf{Z}_t^T\mathbf{L}^{-T}$ ,  $\mathbf{I}$  is an  $n$ -dimensional Identity Matrix and  $\alpha^2$  is a scalar ( $\alpha^2 \leq 1$ ). It is shown that by applying PCA on  $\mathbf{Z}_s$ , the constraint matrix that we get represents the sub-space in which the data lies [14]. We get the constraint matrix as follows:

$$\text{svd}(\mathbf{Z}_s) = \mathbf{U}_{1s}\mathbf{S}_{1s}\mathbf{V}_{1s}^T + \mathbf{U}_{2s}\mathbf{S}_{2s}\mathbf{V}_{2s}^T \quad (32)$$

$$\mathbf{C} = \mathbf{U}_{2s}^T\mathbf{L}^{-1} \quad (33)$$

3. The columns corresponding to dependent and independent variables are separated and the regression matrix is calculated using Eq. (6).
4. Now, the elements in  $\mathbf{R}$  are rounded off to truncate any deviations due to noise and numerical residues, by taking the elements closest to 1 to be 1 and rest 0, in each column. The resultant matrix is same as the sub-matrix which infers the connectivity between dependent and independent nodes.

## 6 Simulation Results

The proposed approach is demonstrated using simulated data. Since the noiseless case is trivial, we generated noisy data for different networks and the results are presented. The network was built randomly using the random number generators in MATLAB, as follows:

1. Three numbers between 5 and 100 were chosen random (uniformly) to assign the number of consumers connected to each phase.
2. The  $N$  readings for each of the consumer meters were sampled from one of the three uniform distributions, with different ranges, to account for consumers with different ranges of loads. Different values of  $N$  are chosen as multiples of  $n_i$ .
3. Now the  $N$  readings for each of the three phase meters were determined by summation of the respective consumer meter readings, connected to them. The technical losses in proportion to the consumer readings were added to the summation.  
World Bank data indicates that countries with good power infrastructure have average losses in the range of 2% to 10% [17]. Also the losses are projected to come down further in a smart grid set up. In our Simulation, we considered two cases with losses in the range of 2% to 5% and with losses in the range of 5% to 10% of the energy transmitted.
4. To account for the other errors, we added Gaussian noise to all the readings with mean equal to the reading and standard deviation in the range 0.5% to 1% of the reading.

The algorithm is then applied on 100 such generated data sets, with different number of readings and loss components and the results were noted. These simulations were carried out in MATLAB version R2014a. The time taken for the algorithm to give the solution was also noted in all the cases (Windows 10, Intel i5-4200U 1.64 Ghz processor, 6 GB RAM). Now we plot our results as follows:

1. Fig. 2 shows the time taken to arrive at the solution against the number of nodes for different number of readings, with losses in the range 2% to 5%.
2. Fig. 3 shows the time taken to arrive at the solution against the number of nodes for different number of readings, with losses in the range 5% to 10%.
3. Fig. 4 shows the success rate in % for 100 test cases against the number of readings expressed as ratio  $N/n_i$ .

Plots Fig: 2 and Fig: 3 indicate that the time taken increases with the number of nodes and readings but it is not dependent on the errors. This is apparent

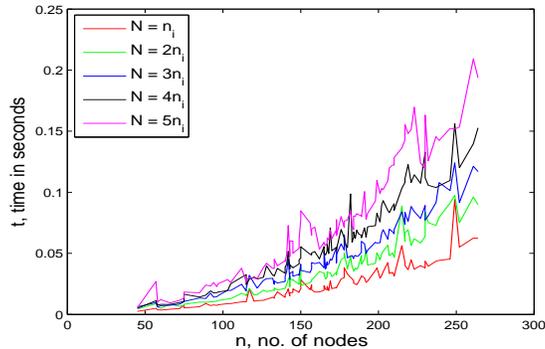


Figure 2: No. of nodes ( $n$ ) Vs Simulation Time ( $t$  in seconds) for losses in the range 2% to 5%

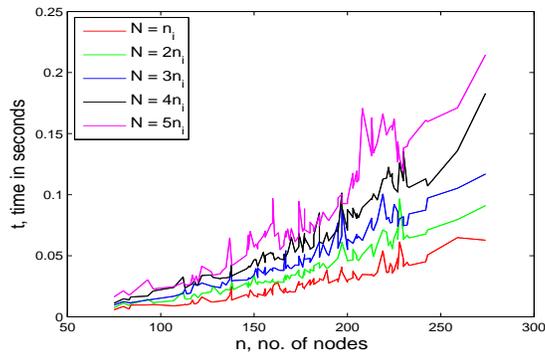


Figure 3: No. of nodes ( $n$ ) Vs Simulation Time ( $t$  in seconds) for losses in the range 5% to 10%

because the cost of computations in the algorithm depends only on the number of nodes and readings and is independent of the error. The important point to be observed is that the time taken is in the order of milliseconds. It shows that our algorithm is time efficient compared to search algorithms like the one proposed in [7], which takes time in the order of seconds to determine connectivity.

Plot Fig: 4 indicates that the success rate is 100% in all the cases when  $N \geq 3n_i$ . Therefore  $3n_i$  readings, satisfying the assumptions, are sufficient to infer the phase connectivity exactly.

## VII. Conclusions and Future Work

In this paper, we show that the phase connectivity problem can be solved through a novel data-driven approach. Our formulation enables the use of PCA

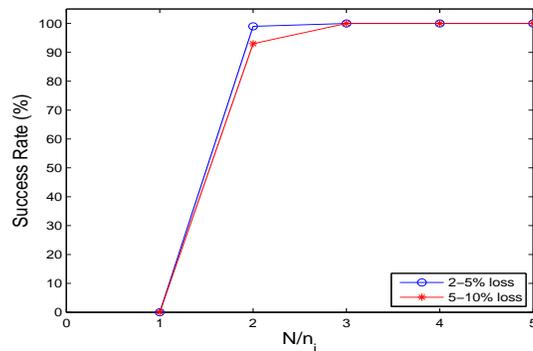


Figure 4:  $N/n_i$  Vs Success Rate (in %)

and its graph theoretic interpretation to infer the connectivity directly.

The Simulations results show that as long as the errors are within the limits assumed, which is so in most of the cases, the connectivity can be exactly determined. Also the time taken is observed to be in the order of milliseconds establishing the efficiency of the algorithm.

In future, we will extend this technique to reconstructing the complete distribution network, from sub-station to consumers, given the meter readings at all nodal points. We would also like to handle the problem of localizing non-technical losses, especially theft, and also address the problem of corrupted or missing measurements.

## Acknowledgment

We would like to thank Prof. Shankar Narasimhan of IIT Madras for his valuable inputs. The finance support to Satya Jayadev P. and Aravind Rajeswaran from Data Science Initiative Grant of IIT Madras, and Nirav Bhatt from Department of Science & Technology, India through INSPIRE Faculty Fellowship is acknowledged.

## References

- [1] U. Bakshi and V. Bakshi, *Electrical Technology*, 4th ed. Technical Publications, Pune, 2009.
- [2] G. Giannakis, V. Kekatos, N. Gatsis, K. Seung-Jun, and B. Wollenberg, "Monitoring and optimization for power grids: A signal processing perspective," *IEEE Signal Processing Magazine*, vol. 30, no. 5, pp. 107–128, 2013.

- [3] K. Caird, “Meter phase identification,” *US Patent Application 20100164473*, 2010, patent No.12/345702.
- [4] C. S. Chen, T. T. Ku, and C. H. Lin, “Design of phase identification system to support three-phase loading balance of distribution feeders,” in *Industrial and Commercial Power Systems Technical Conference (I&CPS)*, Baltimore, USA, 2011, pp. 1–8.
- [5] S. Zhiyu, M. Jaksic, P. Mattavelli, D. Boroyevich, J. Verhulst, and M. Belkhat, “Three-phase ac system impedance measurement unit (imu) using chirp signal injection,” in *Applied Power Electronics Conference and Exposition (APEC), 2013 Twenty-Eighth Annual IEEE*, 2013.
- [6] M. Dilek, R. P. Broadwater, and R. Sequin, “Phase prediction in distribution systems,” *IEEE Power Engineering Society Winter Meeting*, 2002.
- [7] V. Arya, D. Seetharam, S. Kalyanaraman, K. Dontasn, C. Pavlovski, S. Hoy, and J. R. Kalagnanam, “Phase identification in smart grids,” in *IEEE International Conference on Smart Grid Communications, Brussels, Belgium*, 2011, pp. 1–6.
- [8] H. Pezeshki and Wolfs, “Consumer phase identification in a three phase unbalanced lv distribution network,” *IEEE PES Innovative Smart Grid Technologies, Europe*, 2012.
- [9] A. Tom, “Advanced metering for phase identification, transformer identification, and secondary modeling,” *IEEE Transactions on Smart Grid*, vol. 4, 2013.
- [10] M. H. Wen, R. Arghandeh, A. von Meier, Poolla, and V. O. Li, “Phase identification in distribution networks with micro-synchrophasors,” *IEEE Power and Energy Society General Meeting, Denver, CO*, 2015.
- [11] A. Rajeswaran and S. Narasimhan, “Network topology identification using PCA and its graph theoretic interpretations,” in *arXiv preprint arXiv:1506.00438*, 2015.
- [12] B. Andrasfai, *Graph Theory: Flows, Matrices*. Akademiai Kiado, Budapest, 1991.
- [13] I. Jolliffe, *Principal Component Analysis*, 2nd ed. Springer-Verlay, New York, 2002.
- [14] S. Narasimhan and N. P. Bhatt, “Deconstructing principal component analysis using a data reconciliation perspective,” *Computers & Chemical Engineering*, 2015.
- [15] “Ansi c12.20-2010,” *American National Standard for Electricity Meters*, pp. 1–11, 2010.

- [16] P. D. Wentzell, D. T. Andrews, D. C. Hamilton, K. Faber, and B. R. Kowalski, "Maximum likelihood principal component analysis," *J. Chemometrics*, vol. 11, pp. 339–366, 1997.
- [17] W. Bank. (1991) World bank data on electric power transmission and distribution losses. [Online]. Available: <http://data.worldbank.org/indicator/EG.ELC.LOSS.ZS>