| Title | |
|---|---|
| Author(s) | Nower, Naushin |
| Citation | |
| Issue Date | 2015-03 |
| Type | Thesis or Dissertation |
| Text version | ETD |
| URL | http://hdl.handle.net/10119/12747 |
| Rights | |
| Description | Supervisor: , , |

JAIST
JAPAN
ADVANCED INSTITUTE OF
SCIENCE AND TECHNOLOGY

Japan Advanced Institute of Science and Technology

# Data Recovery Schemes for Cyber-Physical Systems with Incomplete Feedback

by

Naushin Nower

submitted to
Japan Advanced Institute of Science and Technology
in partial fulfillment of the requirements
for the degree of
Doctor of Philosophy

*Supervisor:* Associate Professor Yuto Lim

*School of Information Science*
*Japan Advanced Institute of Science and Technology*

March, 2015

# Abstract

Now-a-days real-time systems have been intensively explored by the research community around the world due to many future technologies require real-time processing. Enormous efforts have been made on the upcoming technologies such as Internet of Things (IoT), Machine-to-machine (M2M), Cyber-Physical systems (CPS), Big data etc. These new technologies rely on wireless sensor and actuator networks (WSAN) as a communication media to perform real-time control and communication. However, by using WSAN, the point-to-multipoint mode of communication cannot guarantee reliable and real-time communication. Thus unreliable communication severely degrades the overall system performance and as well as it can affect the control and computation of the real-time system. To ensure real-time and guaranteed communication for point-to-multipoint configuration, data recovery scheme is needed. One of the examples of real-time point-to-multipoint systems is CPS, which enable orchestrating networked computational resources with physical systems. Moreover, CPS have many benefits over conventional network control system in terms of network integration and scalability point of view and also attract attention in a variety of different areas such as smart grid, health care, intelligent transportation, etc.

CPS enable the virtual world to interact with the physical world in order to monitor and control the intended parameter in real-time basis through the feedback control loop. Thus, the proper timing and accuracy of feedback data is very important for the interaction between the cyber and the physical world. Therefore a data recovery scheme is designed to ensure uninterrupted control in CPS.

This dissertation concerns research of technological issues for analysis of data, design and evaluation of a data recovery algorithm and error minimization from the recovered data. The overall objective of this dissertation is to develop a data recovery scheme, which provides quality of result in terms of efficiency and real-time.

In the data analysis part, the data patterns of various physical systems are investigated and a general classification is made according to the property such as data series with small variation, or large variation and/or repetition exist on the data series. To recover

various patterns it is important to know the nature of their underlying property. To do this, a data pattern analyzer is proposed which is able to classify various data patterns, as a data pre-processing step. iHouse data and Intel Berkeley Research lab data are examined using the analyzer.

Some data series remain stable with small change time and some time it is highly correlated with space. Thus to recover this, a data recovery scheme, called Efficient spatial data recovery (ESDR) scheme is proposed. In this scheme, a recovery algorithm is presented with Pearson correlation coefficient (PCC) to efficiently solve the long consecutive missing data. The proposed scheme is evaluated on iHouse data. On the other hand some data patterns have a randomness and variation in its nature, which make a great challenge to maintain the real-time control whenever the data is lost. To handle these kind of data, an Efficient Temporal and Spatial Data Recovery (ETSDR) scheme is proposed. The proposed scheme consists of two phases. In the first phase, which is pre-processing step, the temporal model is identified for large variation data and determined the spatial effects of neighbors. Auto Regressive Integrated Moving Average (ARIMA) model is a very powerful model to identify the auto-correlated nature or trend of a data series. In the next phase, which is real-time/online, temporal model and spatial effect is utilized to recover missing data.

Moreover to improve the recovered data, Kalman filter is used to reduce the error from the model estimated data. The temporal model, generated from ARIMA has internal errors and the model parameters may not remain constant. Thus, to improve the accuracy of the estimated data, a Kalman filter is incorporated to reduce the error. Before that, the window for Kalman filter is fixed to determine the proper process noise co-variance in real-time. Numerical results reveal that the proposed ETSDR/EM are very promising regardless of the increment percentage of missing data in terms quality of result (QoR).

This proposed research can help the development of CPS applications by ensuring uninterrupted control.

Keywords: Data recovery, correlation, cyber-physical systems, real-time, quality of result.

# Acknowledgments

I would like to first express my greatest gratitude to my principal advisor, Associate Professor Dr. Yuto Lim, who encouraged and advised the idea from the beginning to the end of my dissertation. I am delighted and grateful to be able to work under his excellent supervision.

I gratefully acknowledge the generous support and cooperation of Professor Dr. Yasuo Tan, who is my sub supervisor. He gave a lot of valuable comments and intellectual effort to my research. He helped me to realize how to do the realistic and scientific research.

I am deeply grateful to my minor research supervisor Associate Professor Dr. Masashi Unoki, who has provided me the knowledge on prediction during my minor research work. Since I am very newbie in the signal processing field, he gave a helpful guidelines, discussions and suggestions on my research.

I also thank Professor Dr. Mineo Kaneko, Japan Advanced Institute of Science and Technology, and Dr. Bing Zhang, NICT, for serving on my dissertation committee.

I am thankful to all members of Tan and Lim Lab for their support and cooperation. Special thanks to my lab member Wai Wai Shein and Kho Lee Chin for their helpful discussions and mental support.

Finally, I would like to thank my lovely family members for their never-ending love, support and prayer.

# Contents

# List of Figures

# List of Figures

# List of Tables

# Chapter 1

# Introduction

In many future technologies, Internet of things (IoT), Big data, Cyber-Physical Systems (CPS), Machine-to-Machine (M2M) require real-time processing. The purpose of these systems is to connect as much as devices and to ensure anywhere anytime communication in real-time. These new technologies rely on WSAN as a communication media to perform real-time. Communication media transports packet in different forms of communications, like unicast, multi-cast, broadcast and any-cast. Among them, multi-cast and broadcast are unreliable. To ensure data are transmitting in real-time, line configuration is used. Line configuration can be point-to-point or point-to-multi point. To deal with real-time, point-to-point is not a problem, but point-to-multi-point can not guarantee reliable data transmission. Thus, to ensure real-time and reliable communication for point-to-multi-point, data recovery scheme is needed for any real-time system. One of the example of emerging real-time point-to-multi point system could be CPS.

Recently, CPS have emerged as a prominent direction because of the communications of physical and virtual worlds. CPS perform sophisticated interaction between cyber and the physical entities through closed feedback control loops. This interaction, that is feedback control loop plays an important role in any networked control system specially for CPS which have a huge number of potential applications. Thus, feedback data loss can severely degrade the overall system performance and as well as it can affect the control and computation of the CPS. CPS have received a great deal of attention recently in a wide varieties of emerging and time-critical applications with different data patterns [1]. In wide varieties of applications, feedback control loop from WSAN must be present on

time to ensure real-time control. To maintain uninterrupted control, it is always needed to ensure the continuous presence of feedback data, which is frequently lost, corrupted or delayed due to nature of WSAN. This dissertation deals with design and implementation of an data recovery scheme to ensure uninterrupted feedback control for CPS.

## 1.1 Overview of Cyber-physical Systems (CPS)

CPS enable the virtual world to interact with the physical world in order to monitor and control the intended parameter in real-time basis. In CPS, technologies such as communication, control, computation, cognition and sensing converge to create new technologies for a smarter society [2]. The area of CPS represents the intersection of several system trends, such as real-time embedded systems, distributed systems, control systems and networked wireless systems.

To facilitate communications between the cyber and the physical world, WSAN is an essential component of CPS. This is because, the traditional wireless sensor network (WSN) is limited in its ability to monitor the physical world [2]. However, CPS achieve this requirement by the combination of sensing, interaction and changing the physical world in real-time by using feedback control loop. In a typical CPS application, sensor nodes collect information from the physical world as a source of CPS input. Upon receiving the input information, a controller makes a corresponding decision by computing and actuators perform a corresponding action in the physical world through the closed-loop feedback. Thus, the proper timing and accuracy of feedback data is very important for the interaction between the cyber and the physical world. Figure 1.1 shows the shows the basic architecture of CPS, where cyber and physical world make interaction through feedback control loop.

The special characteristics of CPS in [3] are: CPS model must represent for physical world, sensors and actuators, hardware platform, software, network and control system. Obviously, CPS are different form desktop computing, traditional embedded/real-time systems, and WSN. However, they have some different characteristics as defined in [3], [4] and [5]:

- ***Cyber capability in every physical component and resource constraint:***

Figure 1.1: Architecture of CPS.

CPS emphasis on computational elements, and link between the computational and physical elements. But in the embedded system, emphasis is provided on the information processing unit only and the system resources are usually limited.

- **Closely integrated**: CPS deeply integrate communications and computation with physical processes.

- **Networked at multiple and extreme scales**: In CPS, networks such as, wired/wireless network, Wi-Fi, Bluetooth and GSM, and etc are included at a multiple and extreme scale to improve the scalability.

- **Complex multiple temporal and spatial scale**: In CPS, the different components likely have unequal granularity of time and spatiality. CPS are strictly constrained by spatiality and real-time capacity.

- **Dynamically recognizing/re-configuring**: CPS, as very complicated and large-scale systems, must have adaptive capabilities.

- **Closed-loop control and high degrees of automation**: CPS provide convenient human-machine interaction, and advanced feedback control technologies are widely applied to these systems.

- **Operation must be dependable and certified in some cases**: Reliability and security are necessary for CPS because of its extreme scales and complexities.

## 1.2 Application Domains of CPS

In recent years, CPS become a very active research field for engineers and researchers because of its potential to realize the vision of smarter world.

The applicability of CPS is found in numerous time-critical applications including smart house to smart grid. Emerging applications of CPS include, medical devices and systems, aerospace systems, transportation vehicles and intelligent highways, defense systems, robotic systems, process control, factory automation, building and environmental control, smart spaces, intelligent home and so on [2]. These systems are equipped with a large network of sensors distributed across different components, which leads to a huge amount of measured data available to the system controller. For example, CPS can be used in the medical health care applications, where various types of sensors are used to monitor patient's condition and then controller communicates with doctor using the feedback closed control loop system. Thus, the doctor can remotely monitor the patient's physical condition, give suggestions or prescriptions and also do remotely guided robotic microsurgery. Moreover, CPS are planning to use in more complex situation, in particular robot-assisted MRI guided interventions on aortic valve implantation, cardiac surgery, etc [7]. In these time-critical applications, the accuracy and real-time presence of feedback data is very essential for the controller to make a real-time and highly reliable decision.

## 1.3 Current Research on CPS

CPS are integration of computation, networking, and physical dynamics, in which embedded devices such as sensors and actuators are networked to sense, monitor and control the physical world. It is expected that in both the academic and industrial communities, CPS will have great technical, economic and societal impacts in the near future. The CPS of tomorrow will far exceed those of today in terms of both performance and efficiency. The realm of CPS is opening up unprecedented opportunities for research and development in numerous disciplines, e.g. computing, communications, and control. In recent years, CPS have been attracting attention from a rapidly-increasing number of researchers and engineers. According to the US National Science Foundation (NSF), CPS is identified as a key area of research [8]. Starting from the 2006, the NSF and other United

States federal agencies organized several workshops, conferences and provided grants on CPS. However, to fully exploit the potential of CPS, many research challenges must be overcome. There are considerable challenges, particularly because the physical components of such systems requires safety, efficiency, real-time and reliability requirements qualitatively different from those in general-purpose computing. Moreover, physical components are qualitatively different from object-oriented software components also. As a result, CPS need advanced approaches for building abstractions and architectures to enable control, communication and computing integration. The architectures and abstraction should allow the integration and interoperability of different heterogeneous systems. Besides these, current research on CPS also concretes on distributed computations and networked control area. To meet the high reliability and security requirements for CPS, new frameworks, algorithms, methods, and tools, software, variable time delays, failures, reconfiguration, and distributed decision support systems are potential area of research to interact with the physical environment. Research on software components, operating systems and middleware are also going on for CPS. Software for CPS must be highly dependable, re-configurable, and, where required, certifiable, from components to fully integrated systems. Some of the current research on specific CPS applications are described below.

CPS play a major role in the design and implementation of intelligent transportation system (ITS). The ITS are designed to address a range of problems including congestion, fuel consumption, and thus improving cost and safety of our roadways [9]. According to [10], there exits numerous research challenges for ITS due to (a) physical/environmental factors, such as mobility and speed of vehicles, density of vehicles, characteristics of the wireless radio channel, and power and bit rate of radio transceivers, and (b) cyber issues, such as MAC layer access point associations and address resolutions (ARP), network layer addressing, routing and hand-offs, and transport layer re-transmissions lead to unpredictability in the timely and reliable dissemination of information to drivers and so on.

The research on CPS in health care is another promising area of research. In CPS, the combination of active user input such as, feedback system, digital records of patient data, and passive user input such as bio-sensors and/or smart devices in health care environ-

ments can support the data acquisition for efficient decision making [11]. This combination of data acquisition and decision making system is yet to be rigorously explored in health care applications and, therefore, such combination is a matter of high research interest. In health-care applications, bio-medical sensors are responsible for collecting important physiological data and these data are fed to the processing and communication system for further use. Opportunities of utilizing CPS in health care include the introduction of coordinated interoperation of autonomous and adaptive devices, as well as new concepts for managing and operating medical physical systems using computation and control, miniaturized implantable smart devices, body area networks, programmable materials, and new fabrication approaches [12], [13].

Smart grid and renewable energy research and development has been in the forefront of research interest and is therefore a high priority for policy makers. In a smart grid, a variety of communication networks are interconnected to the electric grid for the purpose of sensing, protection, monitoring, and control. Most recently, these networks include connections between suppliers, consumers, stakeholders in economic markets, and independent system operators [14]. In order to implement the smart grid promise we have to utilize CPS that can able to monitor, communicate and control information and actions on the real world. CPS is viewed as an integral part of the smart grid, however, several challenges is needed to be effectively addressed.

Besides these research challenges, CPS control imposes considerable quality of service (QoS) requirements on WSAN. Depending on the type of application, QoS in WSAN can be characterized by timeliness, reliability, robustness, availability, and security, among others. The unique requirements of CPS make it quite difficult to provide QoS support in control systems over WSAN. Some major challenges to control in CPS over WSAN are resource constraints, platform heterogeneity, unpredictable wireless channel characteristics, dynamic network topology and heterogeneous data etc.

## 1.4   Incomplete Feedback of CPS

CPS monitor and control the physical processes by using the feedback loops where physical processes affect computations and vice versa [1]. Based on this feedback, the computation and actuation is performed on the physical world. Thus, feedback makes the interaction

Figure 1.2: (a) Conventional CPS and (b) CPS with proposed data recovery scheme

between the physical and the cyber world. The concept of CPS is not valid without this interaction through feedback. Feedback loop deals with the regulation of the characteristics of a CPS. The main idea of feedback control is to exploit measurements of the system, to determine the control commands that yield the desired system behavior. A controller, together with some sensors and actuators, is usually used to sense the operation of the physical system, compare it against the desired behavior, compute control commands, and perform actions onto the system to effect the desired change. This feedback architecture of a cyber-physical control system is regarded as a closed loop, implying that the cyber space and the physical system are able to affect each other. The proper timing and accuracy of feedback data is very important for interaction between the cyber and the physical world. However, the controlling decision is hampered, when the measurement from the sensor missed or lost which forms incomplete feedback. This incomplete feedback affects the control and performance of CPS.

## 1.5 Research Problem and Motivations

Extensive research efforts have been made to develop a variety of time-critical applications in CPS, including smart house to smart grid such as intelligent transport system, environmental monitoring, energy management, heath-care, security control, etc. In all of these applications, CPS exploit the physical information collected by WSAN, thus it also inherit the wireless contention problem of WSAN. This is a challenging issue for

control in real-time. Wireless channels have many adverse properties like path loss, fading, adjacent channel interference, node/link failure, etc. Besides these, wireless signal can be easily affected by noise, physical obstacles, node movement, environmental change and so on [15]-[16]. Because of this unpredictable and dynamic nature, sensing data loss is a common phenomenon, which makes hamper in controlling decision. In particular, for time-critical applications, feedback data must have to arrive on time, to make decision. In many cases, re-transmission cannot provide appropriate solution because of the unpredictable network behavior, which can cause high delay. To maintain uninterrupted control, we always need to ensure the continuous presence of feedback data. The Fig. 1.2(a) shows the conventional CPS and Fig. 1.2(b) depicts proposed data recovery scheme for point-multi point CPS.

On the other hand, in a wide spectrum of CPS applications, different data properties are observed in different applications. In these applications, systems use a large network of sensors distributed across different areas, which leads to a huge amount of measured data available to the system controller. These measurements are collected continuously along the time, they can be regarded as a time series data. These time series data also have different patterns in terms of their shape, trend, variation and periodicity. Some series maintain stable stage, some show stochastic behavior and others exhibit repetition in their evolution. By considering different CPS applications, data patterns are classified: data with small variation, data with large variation

. Therefore, to handle any uncertainty it is better to know the behavior or trend of the data.

Data recovery is a part of most research and there exist several methods to handle this. Even-though, there exist several methods, the recovery of data loss for CPS still poses an open problem because of its unique requirement. The whole recovery process for CPS must be held in real-time and invisible to the outside world.

No universal method seems to be superior for every data set. Even if one methodology works well with one type of data set, the results often cannot be repeated on other data sets. This is due to the underlying distributions in the data sets, temporal and spatial correlations between them, the amount of missing values and the sample size. Data recovery methodologies can create biases with imputed values if the correct underlying

behavior of the data set is not known and applied. Thus understanding the relationships needed to create a superior imputation method is not a luxury when missing values are present. However, it is also stated that, missing rates of less than 1% are generally considered trivial to deal with and 15% is manageable [17]. As missing values increase to $5-15\%$, methods that are more sophisticated are required to handle the downfalls of single imputation methods. Due to these problems outlined, it is clear that more work is needed to advance all fields of scientific research [18].

In addition, in the existing literature, there is no direction of data recovery based on data pattern for CPS. Thus, the recovery process without considering the nature can not provide a solution for all. To recover data accurately, it is important to understand the nature of the data and their spatial relationship with others. For time series data, a general tool is needed that can analyze and determine the pattern from the data. Thus, it is important to build an effective data pattern analyzer to analyze the data, for better understanding the underlying properties of collected time series data that control the system operation. Based on the data properties, it is easier to design an effective data recovery algorithm to provide uninterrupted control. Thus, successful determination of data pattern ensure efficient data recovery to maintain continuous control. To achieve our motivation, it is proposed a data pre-processing stage, where the data analyzer is used to classify the data pattern and based on that property, a model is built for real-time recovery process.

For those reasons, the research problems of this dissertation are defined as followed:

- design and implementation of data pattern analyzer to classify the data patterns of CPS applications

- develop algorithms to recover the missing data to ensure the uninterrupted control

- applying Kalman prediction to minimize the error in recovery

Then, this dissertation is motivated with three main parts:

      i. design and implementation of data pattern analyzer

     ii. Development of pattern based data recovery algorithm

    iii. Error minimization using Kalman filter

In the first part, investigation is made on potential CPS applications. Based on investigation of different physical systems, a classification is made based on specific property. Then, a data pattern analyzer is designed to classified the data stream for CPS. The designed analyzer is used to classify by using fast Fourier transform (FFT), auto-correlation coefficient function (ACF) and cumulative sum model.

In the second part, data recovery algorithms are proposed for different data patterns. Before designing consideration is made on the spatial and temporal correlation among the data pattern. The some pattern almost remains stable with time. Using this stability property with the spatial correlation data is recovered. Some data series is normally highly auto-correlated and have a large temporal variation . ARIMA [19] model is a very powerful model to identify the auto-correlated nature or trend of data series with large variation . The time series data that has inexplicable changes in direction, is analyzed and build a temporal model by modeling it in ARIMA model in pre-processing step. In real-time recovery algorithm, the model and spatial correlation is used to recover data. Then, evaluation of the data recovery algorithm is shown in terms of root means square error (RMSE), mean absolute error (MAE) and integral of absolute error (IAE).

Third part is regarded with error minimization of the system using Kalman filter. Kalman filter is applied on the recovered data to improve the accuracy. To get the better result using Kalman filter, a training is made on pre-processing step to get the proper error co-variance. Then the proposed system is evaluated in terms of quality of result (QoR).

## 1.6   Dissertation Purpose and Objectives

The presence of feedback control loop is very important to make the real-time controlling decision for CPS applications. However, the unpredictable nature of WSN cannot guarantee the presence of feedback data every time. For these reasons designing the data recovery scheme is critical by ensuring the accurate and real-time presence of feedback data. In addition, it is needed to ensure certain level of accuracy of recovered data. In this aspect, Kalman filter becomes one of the solutions to reduce the error of the recovered data.

The purpose of this research is to develop an data recovery scheme for the application

of CPS technology with different data patterns with high accuracy while maintaining real-time.

In particular, the research objectives are summarized as follows.

1. To classify different data patterns for CPS applications and develop a data pattern analyzer for CPS data patterns.

2. To design a data recovery scheme for different patterns of CPS by considering the spatial and temporal correlation between them.

3. To reduce the error, apply Kalman filter to the recovered data.

## 1.7 Dissertation Contribution

The contribution of this dissertation fall in three parts concerning respectively data analysis and model generation, data recovery algorithms and error minimization using Kalman filter. This dissertation can help the development of CPS application by ensuring reliable and real-time feedback data. In this dissertation, the following specific contributions are made to advancing the state of the art in this area.

1. designing a data pattern analyzer for CPS-based applications by studying the properties of different time series pattern.

2. Presenting algorithms for data recovery of different patterns to ensure accuracy and real-time. Simulation results show that the proposed algorithms give better performance in accuracy and time than the conventional approach.

3. Presenting a error reduction method from the recovered data using a Kalman filter. Then, the system is evaluated in terms of QoR that is efficiency and time.

## 1.8 Dissertation Outline

This dissertation is progressed by the following steps:

Step 1. Introduction of the dissertation background, research problem and motivation, dissertation objectives and contributions are in Chapter 1.

Step 2. Chapter 2 deals with related knowledge of this dissertation such as background study, types of missing data, WSAN, correlation and the proposed framework of CPS based data recovery scheme.

Step 3. Design and implementation and the evaluation of the data recovery algorithm for data pattern with small variation are presented in Chapter 3.

Step 4. Chapter 4 describes design and implementation and the evaluation of the data recovery algorithm for pattern with large temporal variation.

Step 5. Data recovery with error minimization is presented in Chapter 5.

Step 6. Summary of the dissertation and discussion of future research directions are depicted in Chapter 6.

# Chapter 2

# Background, Classification and Framework

## 2.1 Introduction

In recent years, CPS has been attracting attention from a rapidly-increasing number of researchers and engineers. To fully exploit the potential of CPS, however, many challenges must be overcome. Wireless sensor and actuator networks play an essential role in cyber-physical control systems, since they are the bridge between the cyber and physical worlds. In comparison with the filed of general WSN in which significant progress has been made over the years, WSAN is a relatively new research area yet to be explored. In particular, there is only limited work in the WSAN area targeting cyber-physical control applications. To ensure to the continuous control for CPS, there is an urgent need to maintain accurate and real-time feedback control data from WSAN.

In the control community, significant effort has been made for data loss compensation. Despite their differences, most of existing data loss compensation methods are computationally intensive in terms of memory and time and requires iterative steps. For these reasons, they are impractical for real-time CPS control. In addition, they are usually not desirable solutions for resource-constrained WSAN because of overly-large computational overheads [20].

## 2.2 Background Research

Missing data recovery is a part of most research and there exist several methods to handle this. Although there exists several methods, but the recovery of data loss for CPS still poses an open problem because of its unique requirement. The whole recovery process for CPS must be held in real-time and invisible to the outside world.

In statistics, it has been made an extensive study on missing data. Little and Rubin discuss an overview to statistical missing data imputation techniques, such as least squares estimates, Bartletts ANCOVA and likelihood-based approaches in [21]. Maximum likelihood (ML), multiple imputation (MI) and expectation maximization (EM) are widely used method for missing data imputation. ML [22] calculates the likelihood function for given set of data, which is a hypothetical probability that uses past event with known outcome. Then, by using iterative steps, ML makes the likelihood function maximum. EM [23] also uses iterative step to maximize the likelihood function but here, model depends on unobserved or latent variables. Based on mean and covariance matrix of multivariate normal distribution, expectation (E) step initializes the expected values for latent variables. Maximization (M) step plugs the expected values into the log-likelihood function and maximizes the log-likelihood function by repeating the E and M steps. However initialization step directly impact the performance of EM based imputation. On the other hand, in MI [24], missing data are filled by m different times to generate m complete data sets. Generated $m$ data sets are analyzed by standard procedure and then combined for inference. But these well known techniques for missing data imputation are not suitable for WSNs, due to their high space and/or time complexities.

Machine learning based imputation methods require sophisticated procedures that use a predictive model to estimate values. These approaches model the missing data estimation by relying information available in the data set. If the observed data contain useful information then, imputation procedure maintains high precision [25]. Multi-layer perceptron (MLP), self organizing map (SOM), k-nearest neighbors (k-NN) are examples of imputation techniques based on learning. MLP is multi-layer computational unit which is connected by feed-forward way. It estimates the missing data by training an MLP to learn incomplete data by using complete data [26]. On the other hand, in SOM, a set of nodes is organized in 2D grid, where each node has a specific position and weight. The

14

weight is initialized by iterative training steps, and then it is used to estimate missing data [27]. Both of this methods require all data to trained and estimate the missing value. But in k-NN [28], to impute missing data, only $k$ nearest neighbor's data is considered. These techniques are used in WSN to impute data but for real time CPS, these are not suitable.

Compressed sensing (CS) [29] is widely used scheme for signal processing to acquire and reconstruct a signal, based on underdetermined linear systems. This takes advantage of the signals sparseness or compressibility in some domain, allowing the entire signal to be determined from relatively few measurements. However, the main difference between the missing data recovery problem and the conventional CS is that, in the conventional CS, the missing sampling sequence is fixed/ set by the users, and usually random linear projections are preferred, on the other hand, in the missing data recovery problem, the sampling sequence cannot be controlled by the user because it is completely determined by the missing events, e.g., locations or nature of missing nodes in the network which is completely uncertain [46].

Besides these, many researchers combine genetic algorithm (GA) with artificial neural network (ANN) [31], GA with Bayes algorithm [25] and many more to estimate the missing value. Xia, et al. [32] first propose a solution for CPS over WSANs to cope with packet loss. They illustrate three prediction algorithms and show a comparison between them. First algorithm based on the assumption that the state of the physical system does not change during the last sampling period. So, previous sample is used to replace the missing value. The second algorithm computes a moving average of the previous m samples to restore the lost data. Thus it treats every previous measurement equally. In third algorithm weighted average of all previous samples is taken to replace the missing one. Simulation result shows that third algorithm works well compared with others.

Choi, et al. [33] exploit an exponentially weighted moving average (EWMA) based value estimation algorithm to reduce the impact of packet. When some packets are randomly dropped in wireless network environment, the EWMA algorithm filters an abrupt increase or decrease by exponentially smoothing commands or data based on the past value profile.

In [34], the authors proposed a data analysis technique to extract meaningful infor-

mation from the large volume of noisy data. Their designed analyzer named Tru-Alarm, is used to recognize trustworthy alarms from the noisy and false alarms. Tru-Alarm estimates the positions of objects causing alarms, and from that constructs an object-alarm graph and carries out trustworthiness inference based on the graph links. Their studies also reveal that the alarm trustworthiness and sensor reliability could be mutually enhanced. This property is used to ignore the alarms generated by unreliable sensors. Moreover, in [35], the authors proposed a method called IntruMine to detect and verify intruders from the untrustworthy data by modeling the relationships between sensor and intruders. The authors discovered the trajectories of intruders from the untrustworthy data by constructing watching network in [36].

In [37] authors discussed about retrieving the atypical events from massive sensor data and analyzing them with spatial, temporal, and other multidimensional information. Whenever a abnormal event happens such as a congestion is detected in traffic system, the sensor will send out as atypical records. They fixed a threshold for normal event and based on that a atypical event is detected and cluster is formed. The basic cluster is designed to summarize an individual event, and the macro-cluster is used to integrate the information from multiple events. The atypical cluster is then used to effective query execution. Each of the existing analyzer is designed for different purposes and objectives. None of this can be used for data traffic pattern analysis for data recovery.

In the literature, there exists some model based data aggregation scheme. In [38], authors proposed an ARIMA based data aggregation method to reduce the energy consumption and number of communication. In their scheme, both sensor node and aggregator have the same model for data generation. Sensor node checks whether the data predicted from the model and measure data is same or not. Whenever, the original value and predicted value remain within the threshold, then the sensor node will refrain to transmit the data to the aggregator. Otherwise, sensor will send the new to the aggregator.

## 2.3   Types of Missing Data

According to Little and Rubin, there are three types of missing data [21]; missing not at random (MNAR), missing at random (MAR), and missing completely at random (MCAR). In MNAR, the data are missing because of its own observation data. That

is, when the probability of an instance having a missing value for an attribute could depend on the value of that attribute. As an example a person with overweight does not want to reveal his weight, thus the value of weight is missed, because of of its own attribute value. In MAR, the data are missing because of the data is depending on other variables. In this case the probability of an instance having a missing value for an attribute may depend on the value of other attribute, but not on the value of the missing data itself. As an example, women less likely reveal their weight. That is probability of missing (weight) does not depend on data weight, depends on gender. In MCAR, the data are missing because of unpredictable circumstances, e.g., the sending packet of a sensor is loss due to the radio link quality is poor. That is, the probability that an observation $(X_i)$ is missing has no relation to the value of $X_i$ or to the value of any other variables. The focus of this research is to handle MCAR.

Figure 2.1: Sensor data conversion

## 2.4  Data Patterns

In CPS, sensors are used for interaction with the physical world. These sensor can be analog and/or digital sensors. In the case of analog sensor, the raw data is encoded and modulation is performed before sending. In the receiver side, the data is decoded by using digital-to-analog converter as shown in Fig. 2.1 (a). These encoded data is

used to recognize patterns. In this thesis, the analog sensor is considered, however, the same procedure can be applicable for digital sensor also by applying the digital to analog conversion as shown in Fig. 2.1 (b).

From the investigation of typical data characteristics in various physical systems, it is identified that the collected data from those systems have a wide range of varieties from one another. To extract the specific properties from the sensors raw data, we investigate temporal (amplitude and frequency) characteristics of sensor data series in different physical systems. This specific property is called patterns. Based on the investigation, the data patterns in physical systems are classified into the following two types

1. Data patterns with small variation: In many physical system the collected data series does not change much for a long time, or the variation of data is very small. For example, the temperature of specific room is controlled at a specific degree during system operation. As a result, collected temperature measurements are remained constant or vary within a small range.

2. Data patterns with large variation: The observed data from many physical systems does not stable and the data changes with time. The data pattern with large variation represents a wide range, among them in this thesis, periodic, non-periodic, stationary and non-stationary are considered. These data patterns are difficult to handle and thus observed data is pre-processed to know their nature. For example some observed time series data contain indeterminacy and randomness in their evaluation. These data series can be represented by statistical terms or probabilistic forms and may have aperiodicity. These can be further categorized into stationary and non-stationary data regardless to their moments stability. Some measured time series data from different systems show strong periodic patterns due to the regular behavior of physical processes.

## 2.5   Wireless Sensor and Actuator Network (WSAN)

WSAN is a distributed network system of sensors and actuators. Sensors gather information about the physical world and actuators perform actions to change the behavior of the physical world. For real deployment, stand-alone WSAN is insufficient. A gateway

device is essential to enable end-to-end connectivity between the sensors and/or actuators of WSAN and the Internet devices. The sensor nodes used to detect/estimate event features from the environment, are highly correlated with time and space. The correlation among sensor nodes bring significant advantages which can drastically enhance the overall network performance [15] . The characteristics of the correlation in the WSAN is summarized as follows:

## 2.5.1   Correlation in WSAN

There exist spatial and temporal correlation in WSAN. Spatial correlation means that adjacent observations of the same phenomenon are correlated. However, temporal correlation is about proximity in time, while spatial correlation is about proximity in space.

Spatial Correlation: Sensor measurement made at different locations may not be independent. If an environment is highly correlated in space, then the spatial information can be used to estimate missing data and the estimation function can achieve a high accuracy. As an example, measurements made at nearby locations may be similar in value than measurements made at locations farther apart [39]. This phenomenon is called spatial correlation. Spatial correlation measures the correlation of a variable with others through space. Spatial correlation can be positive or negative. Positive spatial correlation occurs when similar values occur near one another. Negative spatial correlation occurs when dissimilar values occur near one another.

Measurement of Spatial Correlation: Most of the spatial correlation regression measures the linear correlation between the nearest neighbors. There are several way to calculate spatial correlation between the sensor nodes. Some are discussed as follows:

1. Moran's I is used to measure the global spatial correlation among the sensor nodes [40]. It is based on cross-products of the deviations from the mean and is calculated for $n$ observations on a variable $x$ at locations $j$ and $k$, as follows:

$$I = \frac{n}{S_0} \frac{\sum_j \sum_k w_{jk}(x_j - \bar{x})(x_k - \bar{x})}{\sqrt{(\sum_j (x_j - \bar{x})^2)}} \tag{2.1}$$

where $\bar{x}_j$ is the mean of the $x$ variable, $w_{jk}$ are the elements of the weight matrix, and $S_0$ is the sum of the elements of the weight matrix: $S_0 = \sum_j \sum_k w_{jk}$.

2. Gearys C statistic is another way to calculate the spatial correlation among the sensors. It is based on the deviations in responses of each observation with one another [41]. The value of Geary's C lies between 0 and 2. 1 means no spatial correlation. Values lower than 1 demonstrate increasing positive spatial correlation, and the values higher than 1 illustrate increasing negative spatial correlation. Geary's C is inversely related to Moran's I, but it is not identical. Moran's I is a measure of global spatial correlation, while Geary's C is more sensitive to local spatial correlation.

$$GC = \frac{n-1}{2S_0} \frac{\sum_j \sum_k w_{jk}(x_j - x_k)^2}{\sum_j (x_j - \bar{x})^2} \tag{2.2}$$

3. Pearson Correlation Coefficient (PCC) is a common measure of the linear correlation between two random variables $x$ and $y$. It reflects the degree of association between two variables [42]. Therefore, the coefficient correlation degree of PCC ($\rho_{xy}$) in between two random variables $x$ and $y$ in specified window size ($W$) can be computed as follows

$$\rho_{x,y} = \frac{\sum_{w=1}^{W} (x(w) - \bar{x})(y(w) - \bar{y})}{\sqrt{\left(\sum_{w=1}^{W} (x(w) - \bar{x})^2\right)} \times \sqrt{\left(\sum_{w=1}^{W} (y(w) - \bar{y})^2\right)}} \tag{2.3}$$

PCC is used in this research because, it a local correlation measurement and using this, the most correlated sensor can be determined.

## 2.6 Framework of CPS based Data Recovery Scheme

In this section, the proposed framework of data recovery scheme for CPS is presented. The designed data recovery framework contains two phases: i) Pre-processing: data analysis and model construction and ii) Real-time processing: data recovery with error reduction. The proposed framework is depicted in Fig. 2.2. This chapter deals with data analyzer only and ESDR is presented on Chapter 3 and in Chapter 4, discussion on temporal model and ETSDR are presented. Error minimization with ETSDR is presented on Chapter 5. The chapter wise overview of the proposed data recovery scheme is depicted on Fig. 2.3.

The aim of the first phase of the proposed framework is to analyze the measured data from the sensor and classify according to the property exits on them. This analysis,

Figure 2.2: Framework of CPS based data recovery scheme

classification and temporal model construction is done in data pre-processing step. In the real-time step, the data recovery algorithm based on the pattern is used to recover the data.

## 2.6.1 Data Pattern Analyzer

The aim of this step is to classify the data using the analyzer, based on the property present in the data. Initially, there are three property checkers in the data analyzer: CUSUM [43] is used for checking whether data has very small variation or not, Auto-correlation coefficient function [44] is used to identify the stochastic nature of a series and periodogram [45] is used for periodicity detection of a pattern. The block diagram of the proposed data pattern analyzer is shown in Fig. 2.4. The following assumptions have been considered. First, $n$ observed sensor data is available for analysis. Second, the group of time series data for a applications follow a specific data property.

To identify the pattern with small variation property of a data series CUSUM is calculated. The main feature of stable pattern is that, they have almost constant value or have a very small variation. CUSUM is a widely used sequential analysis technique in process control, to model series with almost constant values and small deviations. It uses

21

Figure 2.3: Overview of the proposed data recovery scheme

two counters $C+$ and $C-$ for each time series $y_t$, which accumulate the deviation of $y_t$ above the mean, i.e., $y_t + \mu$ and below the mean $y_t - \mu$, respectively. The values stored in counters $C_0{}^+$ and $C_0{}^-$ known as upper CUSUM and lower CUSUM, can be regarded as two time series.

$$C_t{}^+ = max[0, y_t - (\mu + t_v) + C_{(t-1)}^+] \qquad (2.4)$$

$$C_t{}^- = max[0, y_t - (\mu - t_v) + C_{(t-1)}^-] \qquad (2.5)$$

where $t_v$ represents the tolerance range of $y_t$s normal behavior. In this analyzer it is determined as half of $y_t$s standard deviation. The deviation error sequence is defined as a $r_t = max(C_t{}^+, C_t{}^-)$. If the maximum of deviation errors is smaller than the threshold value $\delta$, i.e., $maxr_t < \delta$, the model includes $y_t$ and keeps its profile. The threshold value $\delta$ set as two times of the standard deviation of $y_t$.

To determine data pattern with large variation, auto-correlation coefficient function (ACF) is used. It is assumed that pattern has aperiodic and stochastic nature. These data series can be further categorized into two types: stationary and non-stationary. Stationarity, is defined as a quality of a time series data, in which the statistical parameters (mean and variance) of the series do not change with time. The stationary time series

Figure 2.4: Block diagram of data pattern analyzer

data can be determined by examining the auto-correlation coefficient function (ACF) and partial correlation coefficient function (PACF). The ACF is a set of correlation coefficients between the series and lags of itself over time [44]. Auto-correlation finds the correlation of a series against different versions of itself time-shifted by various amounts. Each time-shift amount is called a lag time. The output of an autocorrelation is the correlation amount as a function of lag time. The maximum value will always be at a lag of zero, since a data is always perfectly correlated with an exact copy of itself. The $k$-order auto-correlation coefficient of a data series $y_1, y_2, .., y_n$ is defined as

$$r_k = \frac{\sum\limits_{j=1}^{n-k} (y_j - \bar{y_j})(y_{(j+k)} - \bar{y_j})}{\sum\limits_{j=1}^{n} (y_j - \bar{y_j})^2} \qquad (2.6)$$

where, $r_k$ is the $k$-lag sample auto-correlation and $\bar{y_j}$ is the average of $n$ observations. The PACF is the partial correlation coefficients between the series and lags of itself over time. The $k$-order partial auto-correlation coefficient of a data series is defined as

23

$$\phi_{11} = r_1 \tag{2.7}$$

$$\phi_{22} = (r_2 - r_1{}^2)(1 - r_1{}^2) \tag{2.8}$$

$$\phi_{kj} = \phi_{(k-1)j} - \phi_{kk}\phi_{(k-1)(k-j)} \tag{2.9}$$

$$\phi_{kk} = \left. r_k - \sum_{j=1}^{k-1}\phi_{(k-1)}r_{k-j} \middle/ 1 - \sum_{j=1}^{k-1}\phi_{(k-1)}r_j \right. \tag{2.10}$$

For the stationary time series, the ACF and PACF trend to zero gradually (die out). On the other hand, for non-stationary data series, the value of ACF and PACF remain for a long time. The analyzer uses this property to determine different the types data. For data series with small variation the ACF remains almost constant and for periodic pattern, ACF shows periodicity.

A periodogram is used to identify the dominant periods (or frequencies) of a time series data. This is very a helpful tool for identifying the dominant cyclical behavior in a series, particularly when the cycles are not related to the commonly encountered seasonality. The periodogram of repeated data series contains a dominant spike in their evaluation. Periodogram is calculated as follows

$$Periodogram = \frac{abs(fft(y_j))^2}{n} \tag{2.11}$$

where $n$ is the number of sample in a series $y$ The checker integrator combines the result from the all three property and makes decision based on the percentage of data follows that property.

## 2.7 Summary

In this chapter, the background study, basic definitions and overall framework of CPS based data recovery scheme are presented. Inside the framework the first step to analyze the data for classification. Initially, three pattern checker is used and in future more can be added to improve it.

# Chapter 3

# Data Recovery Scheme with Spatial Correlations

## 3.1 Introduction

In many physical systems, the sensor data maintain high spatial correlation and their evaluation maintain change with small deviation normally. However, in many applications high accuracy of data is recommended and environment can change suddenly because of fire, earthquake etc. To deal with this kind of scenario, we need a data recovery scheme that can handle insufficient feedback control information. In this chapter, a highly Efficient Spatial Data Recovery (ESDR) scheme is proposed that deals with CPS. To do this, a framework structure for the CPS is designed with the proposed data recovery scheme. The designed framework incorporates the proposed ESDR scheme, which is based on the spatial correlation of neighboring sensors by using the Pearson correlation coefficient (PCC). Since sensor data is highly correlated with space and time, the spatial relationship is utilized to recover the lost data.

One of the contributions is that the proposed ESDR scheme ensures timely data recovery because of minimum computation. Second, the proposed ESDR scheme is used to examine the smart home environment with CPS approach in order to maintain desired room temperature at different locations. Thus, the feedback measured room temperature is very important to keep the desired room temperature steadily at all the times. Another advantage is that, the proposed scheme ensures scalability. Since it uses only one-hop

Figure 3.1: CPS with proposed ESDR data recovery scheme

neighbors, thus the scheme can be applicable in both small and large network.

The rest of this chapter is organized as follows. In Section 3.2, the existing research works are presented. The proposed ESDR scheme is presented in section 3.3. In section 3.4 performance metrics are described. Simulation scenario and result are discussed in section 3.5 and section 3.5 summarizes this chapter.

## 3.2    Related Works

The following section reviews the existing research on spatial correlation based data estimation on wireless sensor network.

Guo, et al. [46] design an algorithm considering spatial-temporal correlations of sensor nodes, which is more suitable with WSNs due to nature of WSNs. Their algorithm first checks if a neighbor sensor node is within the missing sensors sensing range. Then the observation from the neighbor is used for filling in the missing values. This generates a spatially correlated replacement. If there are multiple neighbors within the sensors range and they do not have the same readings, the majority reading is chosen. But in real life, there is no guarantee that all the sensors within one-hop neighbor are spatially and temporally correlated.

In the existing literature, there are other two ways to investigate the spatial correlation for missing data recovery, which is inverse distance weighted averaging (IDWA) [47] and Kriging [48]. The IDWA, which is relatively fast and easy to compute, is one of the most widely used methods for computing spatial interpolation [47]. Assuming the spatial correlation in adjacent sensors is uniform, IDWA tries to estimate the values of missing data in the form of some linear combination of neighboring sensors data. The weights for

the linear combination only depend on the distance between the sensors. The weight is higher for the sensor which is situated in large distance compare to the close one. Thus, IDWA will work well if the values of missing sensors are expected to be similar to values of the neighboring sensors. However, this assumption affects the estimation accuracy in many practical situations, where a physical phenomenon varies rather than uniformly increasing or decreasing in magnitude. The averaging process in IDWA has the tendency to smoother the data, which is not suitable for the situation when data change fast in the area of interest.

Kriging is another way to estimate the missing samples using the combination of available measurements. It defines a semi-variogram by calculating the spatial correlation between sensors. From the semi variogram, the weight for the linear combination is determined. As a result, these weights vary spatially and depend on the correlation [48]. It is assumed that the historical variogram is known and can approximately represent the current variogram. Missing samples are estimated based on the historical variogram function. However, the spatial interpolation may not be right if the semi-variogram varies a lot in the temporal dimension .

## 3.3 Efficient Spatial Data Recovery Scheme

In this section, the proposed data recovery scheme for CPS, called efficient spatial data recovery (ESDR) scheme is presented. In this research, the ESDR scheme is designed to mitigate the problem of MCAR.

To deploy the proposed ESDR scheme, a flowchart with the ESDR scheme for CPS is depicted in Fig. 3.2. The following assumptions have been considered. First, the historical data set for one-hop neighbor is available up to window size to perform the ESDR scheme. Second, the error offset ($e_0$) of the measured data and estimated data is initially computed and known. Third, the maximum number of consecutive missing data counter ($MC$) is fixed at initialization stage. The parameter ($MC$) is also used for terminating the entire system to indicate the estimated data cannot be produced anymore because of the long consecutive missing data. In the flowchart, the ESDR scheme will compute the estimated data when, there is an input measured data from the sensors. If there is no missing data, then the measured data is used as a feedback data. At the same time the difference

27

Figure 3.2: Proposed flowchart with ESDR scheme for CPS.

between the measured and estimated data is computed and if the difference is greater then error offset, ESDR scheme is refined to reduce the error. When there is a missing data, the consecutive missing data is evaluated and the estimated data is used as a feedback data.

Most of the spatial correlation for data recovery scheme is focusing on the data correlation that based on the difference between the nearest neighbors. In proposed ESDR scheme, it is considered the most spatial correlation among the one-hop neighboring sensors based on the Pearson correlation coefficient (PCC) [42]. PCC is a common measure of the linear correlation between a variables of two locations $j$ and $k$. It reflects the degree of association between two variables. From the value of PCC, the nature of correlation can be determined. Whenever the value of $\rho$ is 1, it indicates the redundant information produces by that corresponding sensors. Therefore, the coefficient correlation degree of PCC ($\rho_{jk}$) in between two random variables $j$ and $k$ in specified window size ($W$) can be computed as follows

$$\rho_{j,k} = \frac{\sum_{w=1}^{W} (j(w) - \bar{j})(k(w) - \bar{k})}{\sqrt{\left(\sum_{w=1}^{W} (j(w) - \bar{j})^2\right)} \times \sqrt{\left(\sum_{w=1}^{W} (k(w) - \bar{k})^2\right)}} \tag{3.1}$$

That is this scheme first compute the most correlated sensor using PCC and then

28

Table 3.1: Correlation Degree of Pearson Correlation Coefficient

| Degree of Co-relationship | |
|---|---|
| No Correlation | $0 \cdot 1 > \rho > -0 \cdot 1$ $1 \cdot 0 > \rho$ and $\rho < -1 \cdot 0$ |
| Small | $0 \cdot 1 \leq \rho < 0 \cdot 3$ and $-0 \cdot 1 \geq \rho > -0 \cdot 3$ |
| Medium | $0 \cdot 3 \leq \rho \leq 0 \cdot 5$ and $-0 \cdot 3 \geq \rho \geq -0 \cdot 5$ |
| Large | $0 \cdot 5 < \rho \leq 1 \cdot 0$ and $-0 \cdot 5 > \rho \geq -1 \cdot 0$ |

use that sensor measurement to estimate the missing sensor data. PCC measures local correlation between two sensor rather then global correlation measurement like Morgan's I. In PCC, if an environment is highly correlated in space, then the spatial information can be used to estimate missing data and the estimation function can achieve a high accuracy. Table 3.1 shows the association degree of the $\rho$. The range from $-1.0$ to $1.0$ shows that the $\rho$ has a degree of correlation. The negative value of $\rho$ indicates the negative linear relationship, whereas the positive value of $\rho$ indicates the positive linear relationship.

Fig. 3.3 describes the ESDR algorithm, which is used to produce an estimated data from time to time. In this algorithm, it is assumed that the threshold value of estimation counter $(c_{th})$ is used to optimize the estimation function of the algorithm. Once the ESDR algorithm cannot use the PCC, it is recommended that the estimated data is produced based on the nearest neighbor data. When the number of estimation counter $(c_l)$ for the corresponding of sensor $l$ is above the threshold value, the new corresponding of sensor will be computed again. To maintain high accuracy in estimation, the value of $\rho$ is in between 0.5 to 1.0 is selected.

## 3.4 Performance Metrics

To evaluate the performance of the said algorithms, the root mean square error (RMSE), the mean absolute error (MAE) and the integral of absolute error (IAE) are computed .

**Algorithm: Efficient Spatial Data Recovery(ESDR)**

1: **if** $c_l=0$ **then**

2:   **for** each input data $d_j$ of sensor location $j$ **do**

3:     **for** all $d_k$ within one-hop neighbor of location $j$ **do**

4:       Compute $\rho_{jk}$ with specified window size, $W$

5:       **if** $0.5 <| \rho_{jk} |\leq 1.0$ **then**

6:         $l \leftarrow argmax| \rho_{jk} |; c_l \leftarrow 1$

7:       **end if**

8:     **end for**

9:   **end for**

10:   **else if** $c_l > c_{th}$ **then**

11:     $c_l \leftarrow 0$

12:     $d_{est.}(t) \leftarrow d_j(t) = d_j(t-1)$

13:   **else**

14:     Compute $\rho_{jl}$ with specified window size, $W$

15:     **if** $0.5 <| \rho_{jl} |\leq 1.0$ then

16:       $d_{est.}(t) \leftarrow d_j(t) = d_l(t) + [d_j(t-1) - d_l(t-1)]$

17:     **else**

18:       $l \leftarrow argmin\{distance_{jk}\}$

19:       $d_{est.}(t) \leftarrow d_j(t) = d_l(t) + [d_j(t-1) - d_l(t-1)]$

20:       $c_l \leftarrow c_l + 1$

21:     **end if**

22:   **end if**

23: **end if**

Figure 3.3: Pseudo code for ESDR algorithm

The RMSE is a frequently used measure of the difference between values estimated by an algorithm and the values actually measured from the real environment. The RMSE of an algorithm estimation with respect to the estimated value, $d_{est.}$ is defined as the square root of the mean squared error as written as

$$RMSE = \sqrt{\frac{\sum_{n=1}^{N} \left(d_{org.}(n) - d_{est.}(n)\right)^2}{N}} \qquad (3.2)$$

where $d_{org.}$ is original measured value.

The MAE is another statistical measurement that used to measure how close the estimated values are to the measured values. The MAE is given by

$$MAE = \frac{1}{N} \sum_{n=1}^{N} |d_{est.}(n) - d_{org.}(n)| \qquad (3.3)$$

The MAE measures the average magnitude of the errors in a data set, without considering their direction. It is also an average of the absolute error, $e = |d_{est.} - d_{org.}|$. In other words, it measures the accuracy for the continuous variables. The MAE and the RMSE can be used together to analyze the variation in the errors of the data set. The value of RMSE will always be greater or equal to the MAE. The higher the difference between them, the greater the variance in the individual errors in the sample [49]. If the RMSE is equal to the MAE, then all the errors are the same magnitude. In [49], Wilmott, et al. indicate that the MAE is the most natural and unambiguous measure of average error magnitude.

On the other hand, the IAE is a widely used performance metric in control community, which is recorded to measure the performance of the control application. The IAE is calculated as follows

$$IAE = \int_{0}^{t} |d_{est.}(t) - d_{org.}(t)| \, dt \qquad (3.4)$$

where, $t$ denotes total simulation time. In general, the larger the IAE values imply the worse the performance of the control algorithm.

Figure 3.4: (a) iHouse facilities (b) Layout of 2nd floor

## 3.5 Numerical Simulation

### 3.5.1 Simulation Scenario

In this section, conducted simulation studies to evaluate the proposed ESDR scheme is presented. The comparison is made among the proposed scheme to the WP algorithm [32] and the STI approach [46]. For simulation, an experiment is made to measure the inside temperature of the master bedroom in the iHouse facility which is in situated Nomi city in Japan as shown in Fig. 3.4. The room is equipped with eight sensors at eight corners. The measurements are taken in every two minutes. All the sensors forward their data to reach the base station in single radio hop through the simplest spanning tree topology routing protocol.

Based on the collected information from the experiment, the performance of the proposed scheme is investigated using a MATLAB. In this simulation, it is assume that the single sensor produces a missing sensed data when it transmits its packet to the base station. The simulation is done on the data collected on 1st January 2012-5th January 2012. From the collected original data set, data is randomly deleted according to the percentage of missing data (30% to 60% in steps of 10%) using two different random seeds. Then, recover them using the aforementioned data recovery algorithms. The root mean square error (RMSE), mean absolute error (MAE) and integral of absolute error (IAE) are used to evaluate the performance of the said algorithms.

32

Figure 3.5: The comparison of RMSE of all the data recovery algorithms as the percentage of missing data changes from 30% to 60%.

## 3.5.2 Simulation Result and Discussion

In this section, simulation results are presented and discussions are made on the performance of algorithms. The aim of this simulation is to examine the potential of the proposed ESDR scheme in coping with the data missing for the CPS applications. In the proposed ESDR scheme, the PCC is measured in between the sensors from time to time by specified the window size $(W)$, which is ten data samples. The most correlated value $\rho$ of the one-hop neighbor is used to recover the missing data of sensor. It is realized that not all the sensors within one-hop neighbor are spatially correlated with each other. In the simulation, investigation is made on the impact of increasing percentage of missing data on the data recovery algorithm performance. The percentage of missing data is varied from 30% to 60% in steps of 10%.

Fig. 3.5 depicts the average RMSE comparison among three data recovery algorithms. As the percentage of data missing increases, the proposed algorithm always shows better performance that is compared to the existing two algorithms. At the 40% data missing, the proposed ESDR scheme performs slightly better than the WP algorithm. At the 60% data missing, the proposed ESDR scheme reduces almost half of the RMSE than the WP algorithm. The reason for this dramatic improvement is because the WP algorithm cannot cope with the long consecutive missing data. Through this simulation, it is observed that this problem also can be found at the STI algorithm. Both WP and STI

algorithm use the combination of previous measurements only. Thus, they unable to cope with long consecutive missing and frequent changes in the environment of the conducted experiments.



Figure 3.6: The comparison of MAE of all the data recovery algorithms as the percentage of missing data changes from 30% to 60%.

The MAE comparison among three data recovery algorithms is shown in Fig. 3.6. It is shown that the proposed ESDR scheme outperforms the WP algorithm and the STI algorithm. Besides that, the proposed ESDR scheme can steadily maintain a small value of MAE regardless of the increment of missing data. This also means that the distance between the real measured data and estimated data of the proposed ESDR scheme is always stable. As a result, the proposed ESDR scheme can estimate a better value to recover the missing data. Simulation results reveal that the average MAE of the proposed ESDR scheme is about two times smaller than the WP algorithm and the STI algorithm.

In Fig. 3.7, the accumulated IAE comparison of all the data recovery algorithms is plotted. The simulation results demonstrate that the proposed ESDR scheme outperforms the WP algorithm and the STI algorithm. This is because of the error of the estimation function in the proposed ESDR scheme is minimized by using the PCC approach. The IAE values of the proposed ESDR scheme are 2.8 and 4.3 at the 30% and the 60% data missing, respectively. The average IAE of the WP and STI algorithms is about 3.0 times larger than the proposed ESDR scheme.
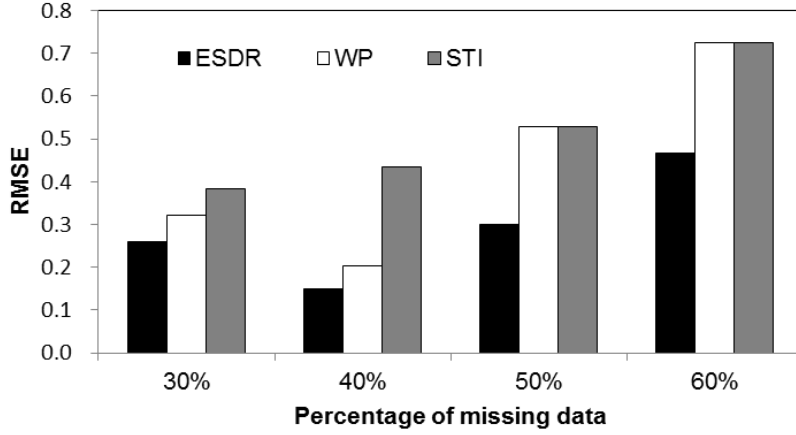
Figure 3.7: The comparison of IAE of all the data recovery algorithms as the percentage of missing data changes from 30% to 60%.

## 3.6   Summary

In this chapter, a new data recovery scheme, called ESDR scheme is presented to ensure a very low error in estimating the missing data. In this research, it is also identified that the nearest one-hop neighbor sensor is not always spatially correlated with the missing input sensor. The simulation results reveals that the proposed ESDR scheme is very beneficial and outperforms the WP and STI algorithms regardless of the increment of missing data. This scheme works well for the the data series whose behavior are stable with time and have a high spatial correlation with others. However, to recover the data pattern with large variation, this scheme can not provide a solution. In the next chapter, the data recovery scheme for data pattern with large variation is discussed.

# Chapter 4

# Data Recovery Scheme with Temporal and Spatial Correlations

## 4.1   Introduction

CPS hold enormous potential for a wide range of emerging applications that include different data patterns. These data patterns have wide varieties of diversities. Among them pattern with large and fast variation is more difficult to recover. In this thesis, the pattern with large variation is used to refer those data series that are not stable with time and may contain randomness with stationary and/or non-stationary or repetitions. To recover the theses data series it is necessary to know the nature of their underlying property. In this chapter, a data recovery scheme for patterns with large variation of CPS is discussed, which comprises data pre-processing step. In the proposed scheme called Efficient Temporal and Spatial Data Recovery (ETSDR), a temporal model is built based on the pattern in pre-processing stage. Then, a data recovery algorithm is proposed to recover the incomplete feedback for CPS to maintain real time control. To determine the nature of pattern with large variation, ARIMA model is used. ARIMA [19] model is a very powerful model to identify the auto-correlated nature or trend of data that has large variation. The data pattern with large variation is identified by ARIMA and then a temporal model is constructed to recover the data.The Fig. 4.1 shows the CPS with proposed ETSDR scheme.

The rest of the chapter is organized as follows. In Section 4.2, the proposed temporal

Figure 4.1: Proposed ETSDR data recovery scheme of CPS

model construction steps are presented. The ETSDR algorithm is presented in section 4.3. Section 4.4 describes the numerical simulation with simulation scenario and the evaluation parameters. Conclusion and discussions are presented in section 4.5.

## 4.2   Temporal Model Construction

For the data pattern with large variation, a temporal model construction step is deployed. It is assumed that, the error offset, that is the maximum difference between the model computed data and the measured data is fixed at initialization step. Historical data size $n$ is available for model generation and verification. The data series is analyzed by modeling it into ARIMA series. The Autoregressive Integrated Moving Average (ARIMA) models, or Box-Jenkins methodology, are a class of linear models that is capable of representing stationary as well as non-stationary time series.

ARIMA model is a very powerful tool that uses historical data to predict future data values. Most of the data series with large variation with stochastic nature can be identified by this model [44]. The ARIMA model, also called Box-Jenkins model, can be divided into three components: auto-regressive (AR), moving-average (MA), and one-step differencing. The AR component estimates the current sample as a linear-weighted sum of previous samples; the MA component captures relationship between prediction errors; and the one-step differencing component captures relationship between adjacent samples. In ARIMA, the AR component captures the temporal correlation in the time series by

modeling a future value as a function of a number of past values. The MA component is modeled as a zero-mean, uncorrelated Gaussian random variable [50].

### Auto-regressive model of order $p$

An auto-regressive (AR) model is a simplified version of ARIMA model which describes random time-varying process. The AR model specifies that the output variable depends linearly on its own previous values [19]. The AR model of sensor $s$ data series $d_{s1}, d_{s2}, .., d_{sn}$ with order $p$ is defined as follows

$$d_{sn} = c + \sum_{j=1}^{p} \varphi_j d_{s(n-1)} + \varepsilon_n \tag{4.1}$$

where $p$ is the order of auto-regressive terms, $\varphi_1, \varphi_2, ..\varphi_p$ are the parameter of the model, $c$ is a constant and $\varepsilon_n$ is white noise. This can be equivalently written using the back-shift operator B as

$$d_{sn} = c + \sum_{j=1}^{p} \varphi_j B^j d_{sn} + \varepsilon_n \tag{4.2}$$

### Moving average model of order $q$

A moving-average (MA) model is a linear regression of the current and previous error of a random series. The MA model of sensor $s$ data series $d_{s1}, d_{s2}, .., d_{sn}$ with order $q$ is defined as follows

$$d_{sn} = \mu + \sum_{j=1}^{q} \theta_j \varepsilon_{n-1} \tag{4.3}$$

where $q$ is the number of moving average terms, $\mu$ is the mean of the series, $\theta_1, \theta_2, ..\theta_q$ are the parameter of the series, and $\varepsilon_n$ is the error. This can be written using back shift operator B as

$$d_{sn} = \mu + \sum_{j=1}^{q} \theta_j B^j \varepsilon_n \tag{4.4}$$

### ARIMA model

An ARIMA model predicts future values of a sensor $s$ data series by a linear combination of its auto-regressive past values, integrated, and moving average of errors. The model

is generally denoted to as an ARIMA$(p, d, q)$ model where, parameters $p$, $d$, and $q$ are non-negative integers used to refer to the order of the auto-regressive, the amount of differencing, and moving average parts of the model respectively. ARIMA is used for non-stationary data time series modeling. If any of $p$, $d$, or $q$ are zero, the corresponding letters are often dropped. For example, if $p$ and $d$ are zero, then model would be denoted MA$(q)$.

$$\theta_p(B)\triangle^d d_{s(t)} = \Theta_q(B)\varepsilon_n \tag{4.5}$$

where $B$ is the backward shift operator, $\triangle$ is the backward difference, $d$ is the order of differencing and $\theta_p$ and $\Theta_q$ are the polynomial of order $p$ and $q$ respectively. In addition, $Bd_{sn} = d_{s(n-1)}$ and $\triangle = 1 - B$. ARIMA$(p, d, q)$ model is the product of an AR part AR$(p)$:

$$\theta_p = 1 - \varphi_1 B - \varphi_2 B^2 - ... - \varphi_p B^p, \tag{4.6}$$

an integrating part:

$$I(d) = \triangle^{-d} \tag{4.7}$$

and a MA part MA(q):

$$\Theta_q = 1 - \theta_1 B - \theta_2 B^2 - ... - \theta_q B^q. \tag{4.8}$$

The flowchart for temporal model identification for data with large variation is shown is Fig. 4.2.

### Step 1: Temporal Model Identification

The aim of this step is to determine whether the series has repetition or not. If the series has no repetition, it is checked for stationary. If the series is not stationary, it is converted to a stationary series by differencing: the original series is replaced by a series of differences and an ARMA model is then specified for the differenced series. Differencing can help stabilize the mean of a time series by removing changes in the level of a time series, and thus eliminates trend and seasonality. The differenced series is the change between consecutive observations in the original series, and can be written as $d_{st}^1 = d_{st} - d_{s(t-1)}$. Whenever the differenced data is not stationary yet then, it is necessary to difference the

Figure 4.2: Flowchart of temporal model construction

data in second times to obtain a stationary series: $d_{st}^2 = (d_{st} - d_{s(t-1)}) - (d_{s(t-1)} - d_{s(t-2)})$. In practice, it is almost never necessary to go beyond second-order differences.

**Step 2: Estimate the Temporal Model**

In this step, the order of $p$ and $q$ is determined from the observed series and the model is identified by comparing the sample ACF and PACF with the theoretical pattern of known model. The properties of ACF and PACF for AR($p$), MA($q$) and ARMA($p,q$) is listed in Table 4.1. From the ACF and PACF, the ARMA model that closely fit to the data can be identified.

**Step 3: Solve the Parameters of Temporal Model**

In this step, the parameters of the identified model is calculated using method of moments and Yule-Walker equations [51].

**Step 4: Verify the Temporal Model**

To verify the model, the model generated data is compared with the observed sensor data. If the verification fails, model estimation is continued until the maximum counter $MC$ is reached. In the case of successful verification, that model is used to estimate the

40

Table 4.1: Properties of ACF and PACF

| | ACF | PACF |
|---|---|---|
| AR($p$) | Tails off (trend to zero gradually) | Cuts off after lag $p$ |
| MA($q$) | Cuts off after lag $q$ (disappear or zero) | Tails off |
| ARMA($p$,$q$) | Tails off after lag $(q-p)$ | Tails off after lag $(q-p)$ |

data in future.

## 4.3   Proposed ETSDR Scheme

To deploy the proposed data recovery scheme, a flowchart of ETSDR scheme is depicted in Fig. 4.3. The proposed ETSDR scheme will compute the model estimated data when there is an input measured data from the sensors. If there is no missing data, then the measured data is used as a feedback data. At the same time the difference between the measured data and model computed data is computed for model verification. If the verification fails, model is updated by computing new parameters.

On the other hand, when there is a missing data, the model estimated data is utilized. At the same time, neighbor's model estimated data and neighbor's measured data is compared. Whenever the difference between two data crosses the spatial regressive threshold ($SR_{th}$), the spatial regression is considered. $SR_{th}$ is the maximum tolerable error value as a threshold indicator used to determine the whether spatial regression to be applied or not in the ETSDR scheme. At the initialization step, $SR_{th}$ is a predefined as a constant value in order to cope with the dynamic environmental changes (i.e., the disturbance effects). Since, the temporal model is based only on the property of data series itself, but in real life, the sensor measurement can be effected by the surrounding environment factors. In the case of a missing data of a sensor, the temporal model is utilized to estimate the model computed data and at the same time checking is performed at all the one-hop neighbors' measurements to determine whether consideration of the spatial regression is needed or not. To handle the spatial regression, the comparison is made on the neighbor's measured data and the neighbor's model computed data. This

Observed data from sensor ($d_{sj(t)}$)

$d_{si(t)}$ is measured data from sensor
$d_{msi(t)}$ is model computed data
$e_i$ is the mean error between $d_{sj(t)}$ and $d_{msj(t)}$ for all the $j$ neighbors
$SR_{th}$ is spatial regressive threshold

Is $d_{si(t)}$ missing?

No     Yes

Input $d_{sj(t)}$ as feedback data

Compute $d_{msj(t)}$ from the temporal model

Does data fit with the temporal model?

Does $e_i > SR_{th}$?

Yes

No

No

Yes

Solve the new parameters of temporal model

Adjust $d_{msj(t)}$ with spatial regression

Update the temporal model

Input $d_{msj(t)}$ as feedback data

Figure 4.3: Flowchart of ETSDR scheme

difference value is defined as model generated error. It is defined in this research that $e_i$ is the average error between all the one-hop neighbors' sensor of the measured data and the model computed data. If this $e_i$ is greater the $SR_{th}$, the spatial regression is added to the model computed data. Otherwise, the only the model computed data is used as a feedback data.

Most of the spatial correlation regression measures the linear correlation between the nearest neighbors. If an environment is highly correlated in space, then the spatial information can be used to estimate missing data and the estimation function can achieve a high accuracy. But in real-life environment, the neighbor sensors can be correlated non-linearly with their neighbors also. This phenomenon is considered and the spatial regression is calculated based on the applications.

An example scenario of ETSDR is discussed here. Suppose five sensors are placed randomly with one-hop neighbor to each other. They are denoted as $s_1$, $s_2$, $s_3$, $s_4$, and $s_5$ as illustrated in Fig. 4.4. It is assumed that the measured data of $s_1$ ($d_{s1(t)}$) is lost at time $t$. In the pre-processing step, the temporal model is constructed for each sensor, which is denoted as $m_{s1}$, $m_{s2}$, $m_{s3}$, $m_{s4}$, $m_{s5}$ and their model computed data are as $d_{ms1(t)}$, $d_{ms2(t)}$, $d_{ms3(t)}$, $d_{ms4(t)}$, and $d_{ms5(t)}$, respectively at the time $t$. Whenever the measured data $d_{s1(t)}$

| Time | $s_1$ | | $s_2$ | | $s_3$ | | $s_4$ | | $s_5$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| t-1 | $d_{s1(t-1)}$ | $d_{ms1(t-1)}$ | $d_{s2(t-1)}$ | $d_{ms2(t-1)}$ | $d_{s3(t-1)}$ | $d_{ms3(t-1)}$ | $d_{s4(t-1)}$ | $d_{ms4(t-1)}$ | $d_{s5(t-1)}$ | $d_{ms5(t-1)}$ |
| t | ? | $d_{ms1(t)}$ | $d_{s2(t)}$ | $d_{ms2(t)}$ | $d_{s3(t)}$ | $d_{ms3(t)}$ | $d_{s4(t)}$ | $d_{ms4(t)}$ | $d_{s5(t)}$ | $d_{ms5(t)}$ |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |

Figure 4.4: An example of 5-sensors scenario

is lost, the model computed data, $d_{ms1(t)}$ is used and at the same time between neighbor's sensor model computed data and measured data is checked for spatial regression.

Suppose the measured data from the sensor 1 ($s_1$) at $t$ time is missing, i.e, $d_{s1(t)}$. The temporal model is utilized to compute the model computed data ($d_{ms1(t)}$) of $s_1$. At the same time, the temporal model will also compute and also compare the measured data of other sensors. Then, the average error, $e_1$ is computed and compare with the $SR_{th}$. If the average error is more than the $SR_{th}$, the spatial regression is taken into account. Otherwise the model computed data only being used as a feedback data.

Fig. 4.5 describes the proposed ETSDR algorithm, which is used to produce an estimated data from time to time.

## 4.4 Numerical Simulations

In this section, conducted simulation studies to evaluate the proposed ETSDR scheme and comparison with the WP algorithm [32] and the EWMA algorithm [33] is presented. A small scenario is created for simulation that can resemble to smart grid applications

**Algorithm: Efficient Temporal and Spatial Data Recovery**

1:   **if** $d_{sj(t)}$ = available **then**

2:     **for** each $d_{sj(t)}$ from the sensor $s_i$ **do**

3:       Compute $d_{msj(t)}$ from the temporal model

4:       **if** $abs(d_{sj(t)} - d_{msj(t)}) >$ error offset **then**

5:         Update the model by calculating new parameters

6:       **end if**

7:     **end for**

8:   **else**

9:     **for** all one-hop neighbors, $k$ of sensor $s_j$ **do**

10:       **if** $avg(abs(d_{sk(t)} - d_{msk(t)})) > SR_{th}$ **then**

11:         $d_{est.(t)} \longleftarrow d_{sj(t)} = d_{msj(t)} +$ spatial regression

12:       **else**

13:         $d_{e(t)} \longleftarrow d_{sj(t)} = d_{msj(t)}$

14:       **end if**

15:     **end for**

16:   **end if**

Figure 4.5: Pseudo-code for efficient temporal and spatial data recovery algorithm

for energy consumption control in smart community. It is assumed a community with five houses, where each sensor (e.g., smart meter) in a house measures the energy consumption and communications with the controller that placed in a cloud for computing the energy demand and supply in real-time manner. The energy value that produces by this smart meter has stochastic in nature and depends on its usage profile of consumer on the home appliances in a house that is correlated with its own usage profile. Moreover, the value of created energies (e.g., solar panel, fuel cell, or electric vehicle, wind energy, etc) from different houses may or may not linearly correlate with other houses as a spatial correlation. This kind of scenario is considered for simulation. In the simulation environment, five sensors and one controller are considered. Sunspot data [] series is collected and assigned to the one sensor. For other four sensors, distance based cor-

Table 4.2: Parameter settings of first simulation

| Parameter | Value |
| --- | --- |
| $p$ : order for AR model | 2 |
| $q$ : order for MA model | 0 |
| $n$ : no. of data for model identification | 100 |
| $m$ : no. of data for verification | 80 |
| $C$ : maximum no. of attempts | 6 |
| $SR_{th}$ :Spatial Regressive Threshold | 1.5 |

related data is generated using MATLAB simulator. Moreover, to make the scenario more realistic, some disturbance effects are added at the certain period of time. In the next step, the temporal model is constructed from the generated data by observing the ACF and PACF. The possible value of $p$ and $q$ is identified as $p = 2$ and $q = 0$ for sensor 1. Then, using MATLAB the parameters are solved. Thus, the estimated model is $d_{ms1(t)} = 1.321 \times d_{s1(t-1)} - 0.637 \times d_{s1(t-2)} + 14.9$. Following the same way, the temporal model is constructed for the other four sensors also, where all of them are AR(2) model with different parameters.

To determine the value of $SR_{th}$ for observed data pattern, the history information of all the measured data is needed. Through this information, the value of $SR_{th}$ can be computed before performing the ETSDR scheme. In other words, the value of $SR_{th}$ is predefined at the initialization stage of the ETSDR scheme. To show how the value of $SR_{th}$ is obtained, the errors of four one-hop neighbors of the sensor 1 is plotted without the disturbance effects as shown in Fig. 4.6. The graphs show that the data changes very frequently thus, the model computed data has the model generated error. The model computed data are obtained from one whole day with the interval sensing of 2 minutes. Through this graph, the value of $SR_{th}$ is set as 1.5, which is use for the first simulation. The parameters and values used in first simulation are shown in Table 4.2.

In this part, data series with large variation that contain repetition is considered. As an example data pattern of a ECG is presented. To evaluate the periodic pattern, evaluation is performed on the Electrocardiography(ECG) data collected from [52]. The analyzer determines that ECG data has a periodic pattern since, it matched with one of its stored data. ECG data has a known pattern, which includes P wave, a QRS complex, and

Figure 4.6: Error of the measured data from each sensor and the corresponding model computed data to determine the spatial regressive threshold

T wave. In addition, there is a interval between the waves, such as PR interval indicates the interval between 'P' wave and 'R' wave, RT interval denotes the interval between 'R' wave and 'T' wave, and RR interval indicates the interval between 'R' wave to next 'R' wave. Moreover, 'Q' to 'R' to 'S' wave, formed the most obvious part of ECG, known as QRS complex, which has a fixed duration. These properties are used to generate a model for ECG in the MATLAB simulator. The ECG model parameters for normal adult people is listed in Table 4.3. The experiment is done for 3 lead ECG, where lead I, lead II and lead III initially started with $0\,°$, $60\,°$ and $120\,°$ phase angle respectively. This spatial property is consider for adjustment when lead I's data is lost. It is assumed that lead I sensor data is missing during the transmission.

For repeated data pattern, the measured data has a fixed periodic pattern. In the second simulation, the three lead sensors have a fixed range for normal patient. In this case, the maximum value is considered as the $SR_{th}$ for that repeated data pattern. For

Table 4.3: Parameter settings of second simulation

| Parameter | Value |
|---|---|
| PR interval | 0.12 – 0.20 s |
| QRS duration | 0.12 s |
| QT interval | 0.43 s |
| RR interval | 0.60 – 1.00 s |
| $n$ : no. of data for model identification | 250 |
| $m$ : no. of data for verification | 240 |

example, PR interval has a range $0.12 - 0.20$ seconds for all lead sensor I, sensor II, and sensor III. When the lead sensor I s PR interval crosses 0.20, the PR intervals of the lead sensor II and lead sensor III are also affected. Thus, it is needed to consider the spatial correlation among them.

Based on the generated data, the performance of the proposed scheme is investigated using a MATLAB. In this simulation, it is assumed that the single sensor produces a missing sensed data when it transmits its packet to the base station. The data is randomly deleted according to the percentage of missing data from the original set and recover them using the aforementioned data recovery algorithms. The root mean square error (RMSE), the mean absolute error (MAE) and the integral of absolute error (IAE) are used to evaluate the performance of the said algorithms.

## 4.4.1 Simulation Results and Discussion

In this section, simulation results and some discussions are presented on the performance of algorithms. The aim of this simulation is to examine the potential of the proposed algorithm in coping with the data missing for the CPS application. The percentage of missing data is varied from 10% to 60% in steps of 10%.

Fig. 4.7 depicts the RMSE comparison among data recovery algorithms for patterns with large variation. As the percentage of data missing increases, the proposed algorithm always shows better performance that is compared to the existing two algorithms. The reason for this improvement is because the proposed scheme estimates the data model then uses that model to generate data. On the other hand, other two algorithm always use the

Figure 4.7: The comparison of RMSE of all the data recovery algorithms as the percentage of missing data changes from 10% to 60%

same combinations of previous measurement. In addition, they do not consider the effect from the neighbors. Through this simulation, it is observed that this problem also can be found at the EWMA algorithm. Both WP and EWMA algorithm use the fixed combination of previous measurements only. Thus, they unable to cope with long consecutive missing and frequent changes in the environment of the conducted experiments.

The MAE comparison for data pattern with large variation, among three data recovery algorithms is shown in Fig. 4.8. It is observed that the proposed scheme outperforms the WP algorithm and the EWMA algorithm. Besides that, the proposed scheme can steadily maintain a small value of MAE regardless of the increment of missing data. This also means that the distance between the real measured data and estimated data of the proposed scheme is always stable.

In Fig. 4.9, the accumulated IAE comparison of all the data recovery algorithms is plotted. The simulation results demonstrate that the proposed scheme outperforms the WP algorithm and the EWMA algorithm. In the 30% data missing the proposed algorithm's IAE is 15.73 on the other hand the IAE of WP and EWMA is 22.17 and 31.21 respectively. At 50% data missing, the proposed scheme's IAE is almost two times smaller than the EWMA algorithm.

From Fig. 4.10 to Fig. 4.12, the RMSE, MAE and IAE comparison for ECG data of

Figure 4.8: The comparison of MAE of all the data recovery algorithms as the percentage of missing data changes from 10% to 60%

all the data recovery algorithms is plotted. The simulation results demonstrate that the proposed scheme outperforms the other algorithms dramatically. This is because ECG data has its own pattern and the model is generated based on that pattern. Moreover to handle irregular ECG spatial correlation with lead II and lead III sensors are considered.

## 4.5   Concluding Remarks

A model based ETSDR scheme for data pattern with large variation is proposed in this chapter. Data pattern with large variation is more difficult to estimate than the those of small variation, thus to handle these data, the model is incorporated from that data pattern. The simulation results reveal that the proposed ETSDR scheme is very beneficial and outperforms the WP and the EWMA algorithms regardless of the increment of missing data because of incorporating model before the recovery.

Moreover more analysis is needed on examining the real-time recovery using the proposed ETSDR scheme. In addition, a prediction analysis is needed to reduce the error from the model estimated data that have with large variation.
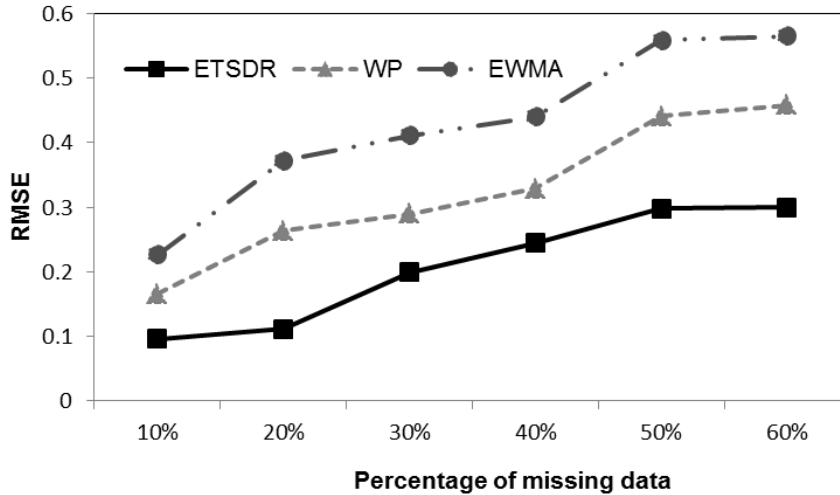
Figure 4.9: The comparison of IAE of all the data recovery algorithms as the percentage of missing data changes from 10% to 60%



Figure 4.10: The comparison of ECG data's RMSE of all the data recovery algorithms as the percentage of missing data changes from 10% to 60%

Figure 4.11: The comparison of ECG data's MAE of all the data recovery algorithms as the percentage of missing data changes from 10% to 60%
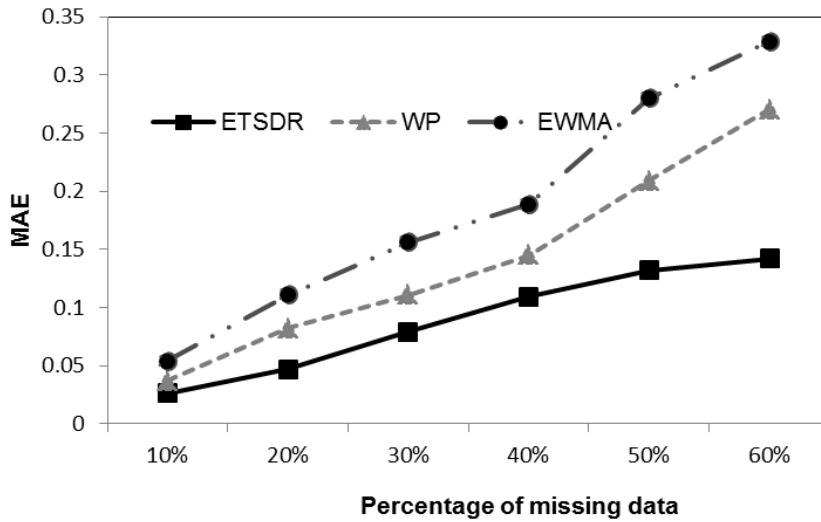


Figure 4.12: The comparison of ECG data's IAE of all the data recovery algorithms as the percentage of missing data changes from 10% to 60%

51

# Chapter 5

# Error Minimization of Data Recovery Scheme

## 5.1 Introduction

Among different traffic patterns, the data pattern with large variation makes a great challenge in efficient and real-time recovery whenever the data is lost. In the chapter 4, the data recovery scheme for data pattern with large variation called ETSDR is proposed. To deal with data pattern with large variation, the temporal model is generated using ARIMA model and consider the spatial effect of the neighbors as a data pre-processing step. However, the performance of the ETSDR algorithm in real-time processing depends on the temporal model identification, more specifically on the parameter estimation and outliers detection of ARIMA model which always has some internal error. Moreover, it is also assumed that the estimated parameters of temporal model are constant throughout the series. But, in real-life CPS applications, the parameters may not remain constant and it is quite impossible to refine the parameter estimation in real-time.

In order to improve the accuracy and ensure real-time computation, in this chapter, a Kalman filter (KF) [54] is incorporated to minimize the error from the estimated data of the temporal model. Before that, a study is performed on Kalman filter to restore noisy speech signal. The proposed scheme for noise reduction and result are put in the appendix. In KF, state model and error covariance act as key role of controlling the performance of KF. The state space model for KF is generated from ARIMA temporal

Figure 5.1: Proposed data recovery scheme for control view of CPS

model. Next, it is needed to determine the correct error covariance to get the best optimal performance of KF. In order to do that, a window is determined to get the proper process noise co-variance. When the error covariance is computed from the actual error of the measurement, satisfactory results are obtained without divergence of Kalman performance. Thus, to get the proper error covariance, the window is fixed for KF whenever the original measurement is available. The Fig. 5.1 shows the architecture of CPS with proposed data recovery scheme.

The rest of this chapter is organized as follows. In section 5.2, the discussion on high-confidence CPS is presented. Section 5.3 presents the proposed real-time data recovery with Kalman filter based scheme is provided. The detail of numerical simulation is presented in Section 5.4. Finally, Section 5.6 summarizes the chapter.

## 5.2   Quality of Result

As it is said before for the interactions between cyber world and the physical world, sensor networks will become a crucial ingredient of CPS due to the need for coupling geographically distributed computing devices with physical elements. In particular, CPS [55] requires the employed sensor networks to support real-time, efficient, dependable, safe, and secure operations. Among them, real-time and efficiency are the most critical criteria for ensuring CPS.

The key issues of CPS are real-time and efficiency. In CPS, the passage of time becomes a central feature to ensure real-time system, in fact, it is one of the important constraint

Figure 5.2: Efficiency vs. execution time in (a) hard and (b)soft real-time systems

distinguishing these systems from distributed computing in general. According to [56], "A real-time system must react to stimuli from the controlled object (or the operator) within time intervals dictated by its environment". Depending on the time constraints, there are two types of real time system: Hard real-time and soft-real time as shown in Fig. 5.2. In a hard real-time system, the system must produces result before the deadline has expired. In a soft real-time system, an answer may still be useful for some time interval after the deadline has expired. The CPS has to be designed to meet the hard real-time, such that the desired outcome is guaranteed within the deadline [57]. Depending on the particular environments and applications the deadline for hard real-time CPS may vary.

QoR [57] is used to evaluate the outcome/result of a scheme or process. It is generally a combination of multiple parameters as a synthetic measure. As an example QoR as a performance indicator of integrated circuits can contains parameters as the area and speed of a chip. As the industry evolved, new chip parameters were considered for coverage by the QoR, illustrating new areas of focus for chip designers (for example power dissipation, power efficiency, routing overhead, etc.). Because of the wide scope of quality assessment, QoR eventually used as a generic representation consisting a number of different parameters, where the value of each parameter was explicitly specified in the QoR analysis document.

In this research, the evaluation merits of accuracy such as root mean square error (RMSE), mean absolute error (MAE) and integral of absolute error (IAE) are used. Among them, MAE provides the unbiased result in terms of accuracy [49]. To define

54

the QoR, efficiency and execution time are used. The efficiency is defined as the improvement of the scheme with respect to the MAE in term of the percentage of missing data. Whereas, the execution time is defined as the elapsed time to produce a loss data of the scheme. In this research, it is defined that the QoR is specified as an acceptable range of efficiency, i.e., above 80% and a deadline of execution time, i.e., below 1 millisecond. If a scheme does not achieve both of the said parameters, then the scheme cannot achieve its QoR.

## 5.3   Proposed Data Recovery with Kalman Filter

To deploy the proposed data recovery with Kalman filter based scheme, a flowchart is proposed which is depicted in Fig. 5.3(a) . Here, the temporal model is used to compute the estimated data and the error is calculated, when there is an input measured data from the sensors. If there is no missing data, then the measured data is used as a feedback data. At the same time, error is computed from the measured data and model computed data to get the actual error for ensuring the better performance from the KF.

On the other hand, when there is a missing data, the model estimated data is utilized and apply KF on the model estimated data to make it more accurate. The KF has been used in a wide range of applications for error minimization. It is an efficient recursive filter that estimates the state of a process in a way that minimizes the mean of the squared error when the process and measurement models are accurate. The details of KF setting is discussed in the following subsection.

To consider the spatial effect, neighbor's model estimated data and neighbor's measured data is compared. Whenever the difference between two data crosses the spatial regressive threshold $(SR_{th})$, the spatial regression is considered. $SR_{th}$ is the maximum tolerable error value as a threshold indicator to determine the spatial regression to be applied or not in the ETSDR algorithm. At the initialization step, $SR_{th}$ is a predefined constant value in order to cope with the dynamic environmental changes (i.e., the disturbance effects). Since the temporal model is based only on the property of data series itself, but in real life, the sensor measurement can be effected by the surrounding environment factors. In the case of a missing data of a sensor, we utilize the temporal model to estimate the model computed data and at the same time we check all the one-hop

55

Figure 5.3: (a) Flowchart of ETSDR/EM algorithm (b)Steps of KF for error reduction.

neighbors' measurements to determine whether we should consider the spatial regression or not. To handle the spatial regression, we compare the neighbor's measured data and the neighbor's model computed data. In this paper, we define that $e_i$ is the average error between all the one-hop neighbors' sensor of the measured data and the model computed data. If this $e_j$ is greater the $SR_{th}$, the spatial regression is added to the model computed data. Otherwise, the only the model computed data is used as a feedback data.

## 5.3.1 Modeling of Temporal Model in Kalman Filter

KF is based on a state-space approach in which a state equation models the dynamics of the data generation process with process error and an observation equation models the generated data with observation error. Thus, it is needed to convert our temporal model into a state-space approach that contains state and observation equations. The performance of KF depends on the proper modeling of these equations and error co-variances. The steps of KF for error reduction is depicted in Fig. 5.3(b).

The temporal pattern of the data is identified by ARIMA model in the data pre-processing phase. An auto-regressive (AR) model is a simplified version of ARIMA model which describes linear stochastic process. The AR model of sensor $s$ data series $d_{s1}, d_{s2}, ..., d_{sn}$ with order $p$ is defined as follows

$$d_s(n) = c + \varphi_1 d_s(n-1) + \varphi_2 d_s(n-2) + ... + \varphi_p d_s(n-p) + W_s(n) \qquad (5.1)$$

where, $p$ is the order of auto-regressive terms, $\varphi_1, \varphi_2, ....\varphi_p$ are the parameter of the model, $c$ is a constant and $W_s n$ is error. The variables $\varphi_1, \varphi_2, ....\varphi_p$ are the state-space model framework. From this, the state equation is formed as follows

$$\begin{bmatrix} d_s(n) \\ d_s(n-1) \end{bmatrix} = \begin{bmatrix} \varphi_1, \varphi_2, ..., \varphi_p \\ 1, 0, ..., 0 \end{bmatrix} \begin{bmatrix} d_s(n-1) \\ d_s(n-2) \\ ... \\ d_s(n-p) \end{bmatrix} + \begin{bmatrix} c \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} W_s(n) \qquad (5.2)$$

where, $\begin{bmatrix} \varphi_1, \varphi_2, ....\varphi_p \\ 1, 0, ..., 0 \end{bmatrix} = A$ is a state transition matrix and $W_s(n-1)$ is the process error.
The observation equation is as follows.

$$O_s(n) = \begin{bmatrix} 1, 0, ..., 0 \end{bmatrix} \begin{bmatrix} d_s(n) \\ d_s(n-1) \\ ... \\ d_s(n-p) \end{bmatrix} + V_s(n) \qquad (5.3)$$

where, $H = \begin{bmatrix} 1, 0, ..., 0 \end{bmatrix}$ is the observation matrix and $V_s$ is the measurement error. Thus, from (1) and (2) we get the state space model as follows

$$d_s(n) = A \sum_{j}^{p} d_s(n-j) + W_s(n) \qquad (5.4)$$

$$O_s(n) = H d_s(n) + V_s(n) \qquad (5.5)$$

Equation (5.4) represents a linear stochastic equation where, $d_s n$ is a linear combination of its previous value and a process error. Equation (5.5) indicates that any measurement value is a linear combination of the data value with the measurement error. From the temporal model, the state transition matrix and observation matrix $A$ and $H$ are derived and then left is to set the co-variances $(Q)$ and $(R)$ of process error $(W_s(n-1))$ and measurement error $V_s(n)$ respectively. By using the same process, the state space model for MA and ARMA temporal model can be derived .

The KF requires that all of the error co-variances to be known exactly. Error co-variances in the KF play a key role in controlling the Kalman gain. At first, the values of $Q$ and $R$ is chosen from the pre-processing phase during the temporal model verification step. The main purpose is to minimize the error from the model, which is the process

Figure 5.4: Determination of window size for stable $Q$

error $V_s n$. Thus, it is needed to set the co-variance $Q$ of process error $V_s n$ properly. Since, it is assumed that the measurement noise is almost zero, $R$ is set close to zero. Initially, the value of $Q$ can be set from the pre-processing step. In the next step, whenever the sensor measured data is available, the actual error is computed and then $Q$ is refined to be more appropriate. In order to do that, a window is defined in which the value of $Q$ will become stable. Before doing that, it is assumed that, there is no missing data within this window length, thus the actual error is used to converge to the stable value of $Q$. As long as the window length is higher, the more accurate $Q$ can be achieved, but at the same time, the assumption become unrealistic. To determine the suitable window length, seven data series with 200 samples without any missing data is analyzed and found that the $Q$ become stable within window length 14.8 (on average) as shown in Fig. 5.4. From this analysis, the window length 15 is fixed to get the stable process noise co-variance in real-time without any missing data. The window utilization is easier, simple and requires less memory compare to training process.

The KF algorithm involves two stages: Time update (prediction) and measurement update (correction). The time update equations are responsible for projecting forward (in time) the current state and error co-variance estimates to obtain the a previous estimates

for the next time step.

$$\hat{d}_s n = A \sum_{j}^{p} \hat{d}_s(n-j) \tag{5.6}$$

$$\hat{P}_k = A\hat{P}_{k-1}A^T + Q \tag{5.7}$$

The measurement update equations are responsible for the feedback of KF, it incorporates a new measurement into the a previous estimate to obtain a next improved estimate.

$$\hat{d}_s(n+1) = \hat{d}_s n + K_k(y_s n - H\hat{d}_s n) \tag{5.8}$$

$$P_k = (1 - K_k H)\hat{P}_k \tag{5.9}$$

where, $K_k$ is a Kalman gain, which is defined as follows

$$K_k = \hat{P}_k H^T (H\hat{P}_k H^T + R)^{-1} \tag{5.10}$$

Figure 6.5 describes the proposed ETSDR/EM algorithm, which is used to produce an estimated data from time to time.

## 5.4   Numerical Simulations

In this section, conducted simulation studies are presented to evaluate the proposed ETSDR/EM scheme compared to the ETSDR algorithm[58], the WP algorithm [32] and the EWMA algorithm [33]. A small scenario is created for simulation that can resemble to smart grid applications for energy consumption control in smart community. It is assumed that a community with five houses, where each sensor (e.g., smart meter) in a house measures the energy consumption and communications with the controller, that placed in a cloud for computing the energy demand and supply in real-time manner. The value of created energies (e.g., solar panel, fuel cell, or electric vehicle, wind energy, etc) from different houses may or may not linearly correlate with other houses as a spatial correlation. In the simulation environment, five sensors and one controller are considered. The sunspot data were verified and recognized as a aperiodic and stationary data series and assigned to the one sensor. The data series for other four sensors are generated using

59

---

**Algorithm: Efficient Temporal and Spatial Data Recovery with Kalman Filter**

---

1:   **if** $d_s(t) =$ available **then**

2:    **for** each $d_s(t)$ from the sensor $s$ **do**

3:     Compute $d_{ms}(t)$ from the temporal model

4:      Apply KF on $d_{ms}(t)$ to reduce error

6:       **end if**

7:    **end for**

8:   **else**

9:    **for**  all one-hop neighbors, $r$ of sensor $s$ **do**

10:      **if** $avg(abs(d_r(t) - d_{mr}(t))) > SR_{th}$ **then**

11:       $d_{est.(t)} \longleftarrow d_s(t) = KF(d_{ms}(t)) +$ spatial regression

12:      **else**

13:       $d_{est.(t)} \longleftarrow d_s(t) = KF(d_{ms}(t))$

14:      **end if**

15:    **end for**

16:   **end for**

---

Figure 5.5: Pseudo code for ETSDR/EM

MATLAB simulator and then assigned it to the four sensors. It is assumed that the one-hop sensors are linearly co-related. Moreover, to make the scenario more realistic some disturbance effects are added at the certain period of time. $SR_{th}$ [58] is set as 1.09 to cope with the spatial effect. The temporal model is constructed from the generated data by following the steps as described in chapter 2. The identified the temporal model is $d_s(n) = 0.11 * d_s(n-1) - 0.96 * d_s(n-2)$ which is a AR(2) model. From this model I get the matrix $A = [.11 - .96]$ and use $H = [1, 0]$ for Kalman filtering. Since the goal is to reduce the error form the model and it is assumed that there is almost no measurement error, thus the value of $R$ is set as $= 1e - 3$. In order to determine the process error, the window length is used as (15) and the value of $Q$ converges from 2.009 and become stable at 0.5.

Based on the generated data, the performance of the proposed scheme is investigated

using a MATLAB. In this simulation, it is assumed that the single sensor produces a missing sensed data when it transmits its packet to the base station. The data is randomly deleted according to the percentage of missing data from the original set and recover them using the aforementioned data recovery algorithms. The root mean square error (RMSE), the mean absolute error (MAE) and the integral of absolute error (IAE) are used to evaluate the performance of the said algorithms.

## 5.5 Simulation Results and Discussion

In this section, we present our simulation results and make some discussions on the performance of algorithms based on QoR that is efficiency and execution time. The aim of this simulation is to examine the potential of the proposed algorithm in coping with the data missing for the CPS application. The percentage of missing data is varied from 10% to 60% in steps of 10%. Figure 5.6(a) depicts the RMSE comparison among data recovery algorithms. As the percentage of data missing increases, the proposed ETSDR/EM always shows better performance that is compared to the ETSDR and other two algorithms. The reason for this improvement over ETSDR is because ETSDR/EM reduces the model generated error using Kalman filter. On the other hand, WP and EWMA algorithm always use the same combinations of previous measurement. In addition, they do not consider the effect from the neighbors. Through this simulation, we can observe that this problem also can be found at the EWMA algorithm. Both WP and EWMA algorithm use the fixed combination of previous measurements only.

The MAE comparison among four data recovery algorithms is shown in Fig. 5.6(b).It is shown that the ETSDR/EM outperforms the ETSDR, the WP algorithm and the EWMA algorithm. Besides that, the proposed scheme with Kalman filter can steadily maintain a small value of MAE regardless of the increment of missing data because of accurate setting of process error co-variance through the window. This also means that the distance between the real measured data and estimated data of the proposed scheme is always stable. In Fig. 5.6(c), the accumulated IAE comparison for all the data recovery algorithms is plotted. The simulation results demonstrate that the proposed scheme with Kalman filter outperforms the others. This is for properly incorporating the Kalman filter with ETSDR algorithm.

Figure 5.6: Comparison of (a) RMSE (b) MAE and (c) IAE of all the data recovery algorithms as the percentage of missing data changes from 10% to 60% and (d) average percentage of QoR of algorithms for 10% to 60% missing data

Table 5.1: Execution time in unit of seconds

| Algorithms | | | |
| --- | --- | --- | --- |
| ETSDR/EM | ETSDR | WP | EWMA |
| 1.5565e-04 | 1.0263e-06 | 2.7368e-06 | 2.0526e-06 |

To measure the execution time of the all the said real-time algorithms, the computer with the Intel Core i7 3.0 GHz processor and the 8 GB memory is used to run each algorithm 10 times. The average execution time of each algorithm is given in Table 5.1, which shows that all of the said algorithms can meet the deadline. To illustrate the QoR of all the algorithms, we depicted one more graph in Fig. 5(d) which shows the average MAE of all said algorithms from 10% to 60% missing data. It is easily observed that only ETSDR/EM can achieve more then 80% in terms of QoR. On the other hand, none of the others can achieve 80% efficiency. Although ETSDR/EM requires higher execution time compare to other three, but it still maintain the deadline. Thus ETSDR/EM maintains

QoR in terms of efficiency and deadline compare to the others.

## 5.6 Concluding Remarks

In this chapter, a data recovery with KF based scheme is proposed for data pattern with large variation of CPS. Since, data pattern with large variation is more difficult to estimate than the others, Kalman filter is incorporated to improve the accuracy in estimation. The simulation results reveal that the proposed ETSDR/EM scheme is very beneficial and outperforms the ETSDR and the WP and the EWMA algorithms regardless of the increment of missing data. Moreover, further research will focus on examining the real communication environment to ensure the performance of real-time execution of ETSDR/EM scheme.

# Chapter 6

# Conclusion and Future Research

## 6.1 Conclusion

This dissertation deals with the research issues for data pattern analysis, designing and evaluation of data recovery scheme and error minimization using Kalman filter of the estimated data for ensuring accuracy and real-time control. Since now-a-days many emerging technologies require real-time processing, thus it is needed to ensure the accurate presence of data on-time.

This thesis illustrates how to ensure continuous presence of data for taking decision on time. In this aspect the analysis of different data pattern and classification is made. Upon this study, for the different data patterns different schemes are proposed. For data pattern with small variation, a data recovery is proposed considering spatial correlation. Since data pattern with almost constant or small variation remain stable for a time, the spatial correlation is considered with its own redundant data. But for the data with large variation, is more complicated thus, a pre-processing is needed where temporal model is construction. The purpose of this pre-processing stage is to know the underlying property of the data and construct a model. The model is constructed to know the temporal correlation with its own data and spatial correlation is applied to handle spatial effect. The constructed model is used in real-time to recover the data. The proposed scheme is evaluated in terms of RMSE, MAE and IAE in terms of accuracy. Then, study on Kalman estimation is performed for error reduction. Then, Kalman filter is used to reduce the error from the estimated data to make it more accurate for CPS. Finally the performance

Figure 6.1: Overall of dissertation.

of ETSRD/EM is evaluated in terms of accuracy and real-time. Figure 6.1 shows the overall of this dissertation.

## 6.2 Future Work

The main objective of this research is to propose a data recovery scheme that can ensure accurate and real-time presence of feedback control data. The proposed data recovery schemes can provide accuracy and timely presence of various pattern of data. However, there are still several issues need to consider for more improvement.

- Considering not only the single sensor's missing data but also multiple sensors missing data

- Need to consider heterogeneous data from different sources

- A general tool is needed to model the periodic data and more analysis is required on other modeling techniques also

- More concentration is needed for finding model for non-linear data

- Need to examine the real-communication environment to ensure the real-time performance

- Need to consider optimization to get the best possible result

- Also need to concentrate a error refinement scheme to analyze the error and refine them

Considering the present research, the future research will focus on specific time-critical application for providing uninterrupted control.

# Appendix A

# Noisy Data Restoration Scheme

In real world scenarios, the desired speech signal is frequently smeared by various kinds of noises. These noises not only degrade the perceptual aspects of speech quality and intelligibility but also reduce the performance of various automated speech systems, such as automatic speech recognition systems, speaker recognition systems, and hearing aids. Therefore, the quality and intelligibility of speech signal in the noisy environments have to be enhanced.

Speech enhancement is concerned with improving the quality and intelligibility of corrupted speech in the presence of noises. During the past two decades, various methods for speech enhancement have already been proposed to remove the effects of noise from the noisy speech to improve its quality. Among them, classical speech enhancement methods such as, spectral subtraction (SS) [60], the Ephraim-Malah algorithm (MMSE-STSA estimator) [61], the Scalart-Filho algorithm (Wiener filtering) [62], have attracted a great deal of attention, because of simplicity and efficient spectral magnitude estimation. SS method [60] performs subtraction of an estimated noise magnitude spectrum from a noisy speech magnitude spectrum, where the noise spectrum can be estimated and updated during periods when the speech is absent. The Wiener filter [62] algorithm performs filtering of a noisy speech signal by using a filter derived based on the minimum mean-square error (MMSE) criterion. These methods employ the short-time Fourier analysis-modification-synthesis (AMS) framework for speech enhancement.

Besides these, there exists various statistical model-based speech enhancement methods in the literature. In the model-based approach, the modeling is done using the statistical properties of the speech signal over multiple frames. This modeling is performed using hidden Markov model (HMM) [63, 64, 65], Gaussian mixture model (GMM), or codebook-based methods [66]. HMM-based speech enhancement is the renowned model-based technique and resolves the common problems of classical speech enhancement methods in dealing with rapid variation of noise characteristics [67]. In [68] the authors combine independent component analysis (ICA) based noise estimator with multi-channel-wise non-linear signal processing to achieve higher noise reduction performance. But all of their improvements are limited due to the consideration of spectral domain only.

Recent research investigate the importance of speech enhancement in modulation domain, such as, modulation-domain Kalman filter (MDKF)[69, 70]. Consequently, Cor-

pus based approach [71], model based speech enhancement with spectral estimation [72], speech enhancement using non-matrix factorization (NMF) [73, 74] are included in modern speech enhancement method. All of the existing methods process the corrupted speech signals by modifying or correcting the speech in either temporal or spectral magnitude only and keeping the phase component unchanged. This is because the phase spectrum conventionally considered is unimportant and has been shown not to contribute much towards speech enhancement. Wang and Lim emphasized this point [75] and this is perhaps the most cited work to justify the unimportance of phase for speech enhancement.

However, recent studies have reported that the use of phase spectrum in the short-time Fourier transform (STFT) can significantly improve speech enhancement [76, 77, 78, 79]. Shannon and Paliwal [76] reported that magnitude only and phase only experiments have been carried out to investigate the effect of phase in speech enhancement. In the magnitude-only experiment, the clean magnitude was used and the phase was set to a random value. In contrast, in the phase-only experiment, the clean phase was used and the magnitude was set to one. Their results showed that the phase spectrum also contains useful and important information. After that, Paliwal and Alsteris [77] investigated whether the shape and length of the window function used in the STFT for phase manipulation are important factors for speech enhancement.

Paliwal *et al.* [78] showed that modifying the phase spectrum can greatly improve speech enhancement. For this, they investigated various cases where the different combinations of noisy, clean (noiseless), and compensated amplitude and phase spectra are considered. This suggested that significant speech enhancement can be possible if the clean phase is known or the compensated phase spectrum is available. They also studied the effect of mismatched or matched windows for both amplitude and phase spectra estimation during AMS in the STFT. The results show that the proper choice of an analysis window and AMS setting on the phase spectrum can significantly improve the speech enhancement. All existing researches on phase spectrum either emphasizes the importance of phase in speech enhancement or investigates the suitable size and shape of the window for phase manipulation. Thus, for better speech enhancement, we obviously need to consider both the amplitude and phase of the noisy signal.

It is well-known that all existing speech enhancement algorithms based on STFT-AMS can improve speech quality but not speech intelligibility [80]. The reasons for that are still unclear so that many researchers have investigated the expected strategy for reducing distortions and enhancing features related to speech intelligibility. On the other hand, from psychoacoustical studies, it is found that temporal envelope (TE) and temporal fine structure (TFS) are important cues for speech perception [81, 82]. It is also revealed that TE and TFS play an important role of improving intelligibility of noise-degraded speech [83, 84]. Therefore, AMS in the filterbank is suitable framework for speech enhancement, rather than AMS in the STFT. Hence, it is expected that temporal amplitude and phase manipulations as the ASM in the filterbank can drastically improve quality as well as intelligibility of noise degraded speech.

Motivated about the effectiveness of phase manipulation from the existing literature, our aim is to propose a speech enhancement scheme as the ASM on the filterbank to enhance both the instantaneous amplitude and phase by using recursive Kalman filter in a Gammatone filterbank (GTFB) [85]. We deal with the instantaneous amplitude and phase on the GTFB because temporal smoothed information (amplitude and phase) are directly related to improve quality and intelligibility of speech. The Kalman filter is of particular interest in smooth prediction method for dealing with the instantaneous amplitude and phase in sub-band. Moreover it can be viewed as a joint estimator for both the magnitude and phase spectrum of speech, under non-stationary condition [69]. We believe, the ability of the Kalman filter to process non-stationary signals as well as estimate both the magnitude and phase spectrum makes it preferable over STFT-based enhancement methods.

In this research, the assumptions of Kalman filter are deeply and carefully analyzed for speech enhancement in noisy environment. There are two issues in Kalman filter on background noises of instantaneous amplitude and phase. That is, the observation and driving noise of both instantaneous amplitude and phase have to be modeled as a white and Gaussian noise for best optimal estimation. In addition, derivation of the accurate transition matrices, is the main core processing technique of Kalman filtering. This is because of enhancement performance of Kalman filter, is dependent on the accuracy and reliability of transition matrices. These transition matrices of the state-equation of both

Figure A.1: Block diagram of proposed scheme (non-blind method) for speech enhancement.

instantaneous amplitude and phase are unknown, thus, it is difficult to set these transition matrices in the Kalman filtering for suitable speech enhancement.

In this research, the speech signal is analyzed as a highly temporal correlated time series signal and the linear prediction is used as a state of the art concept to derive these coefficients as modulation characteristics of amplitude and phase in each sub-band of Kalman filtering. Investigation is made on the characteristics of linear predictive coding (LPC) in modulation domain and the stability property of line spectral frequency (LSF) is utilized to train the LPC in a effective way that can be regarded as a speaker, gender and content independent LPC for Kalman filtering. Thus, the main contribution of this research is to verify the concepts validity and concentrate on the core processing of Kalman filter for providing a state of art concept of the speech enhancement.

The rest of the chapter is organized as follows. In Section 5.2, the proposed scheme for speech enhancement is presented. We describe the details of algorithm implementation in Section 5.3. The evaluation results and discussions are described in Section 5.4. Section 5.5 concludes with conclusion and future works.

# A.1 Proposed Scheme

The proposed method for speech enhancement is intended to improve both instantaneous amplitude and instantaneous phase as the ASM on a Gammatone filterbank. The model consists of three steps: (i) Analysis stage, where instantaneous amplitude and instantaneous phase are extracted from the noisy speech by the GTFB. (ii) Modification stage, where instantaneous amplitude and instantaneous phase are enhanced by a Kalman filter with linear prediction (LP) (with/without training phase). (iii) Re-synthesis by the inverse Gammatone filterbank. The block diagram of the proposed method is shown in Fig. A.1.

First, only the noisy speech $y(t)$, where $y(t) = x(t) + n(t)$, is observed in the proposed model. Here, $x(t)$ indicates the clean speech and $n(t)$ represents a background noise or the other signal. $t$ is continuous time and $m$ is sampling number ($m = 0, 1, 2, \cdots, M$; $\delta t = m/F_s; t = m\delta t$) where $M$ is the number of time samples and $F_s$ is the sampling frequency. From the observed signal $y(t)$, instantaneous amplitude and instantaneous phase are extracted into the frequency components by the Gammatone filterbank (the number of channels is $K$). The output of the $k$th channel is represented as the analytical form by

$$
\begin{aligned}
Y_k(t) &= Y_{1,k}(t) + Y_{2,k}(t) & \text{(A.1)} \\
&= A_k(t) \exp\left(j\omega_k t + j\phi_k(t)\right) & \text{(A.2)}
\end{aligned}
$$

where $Y_{1,k}(t)$ and $Y_{2,k}(t)$ are the components of $x(t)$ and $n(t)$ that have passed through the filterbank, respectively. In addition, $\omega_k$ is the center frequency of the $k$th channel, $A_k(t)$ is the instantaneous amplitude, and $\phi_k(t)$ is the instantaneous phase of the noisy speech.

## A.1.1 Modeling of Instantaneous Amplitude and Phase in Kalman Filter

The Kalman filter, an efficient computational recursive solution for estimating a signal, is widely used in fields related to statistical processing. It not only exploits the statistical characteristics of signal and noise but also utilizes the speech production model based

on the source-filter model. Therefore, we believe that the Kalman filter can be used to remove noise from both instantaneous amplitude and instantaneous phase.

The state and observation equations are the main equations in the Kalman filter, which are defined as:

$$\boldsymbol{S}[m] = \boldsymbol{F}\boldsymbol{S}[m-1] + \boldsymbol{W}[m], \tag{A.3}$$

$$\boldsymbol{O}[m] = \boldsymbol{H}\boldsymbol{S}[m] + \boldsymbol{V}[m], \tag{A.4}$$

where $\boldsymbol{S}[m]$ is the state in discrete time $m$ and $\boldsymbol{O}[m]$ is the observation in discrete time $m$. $\boldsymbol{W}[m]$ and $\boldsymbol{V}[m]$ are driving noise and observation noise that are assumed to be Gaussian white noise. The state equations of $k$th channel for instantaneous amplitude and instantaneous phase are as follows

$$\boldsymbol{S}_{A,k}[m] = \boldsymbol{F}_A\boldsymbol{S}_{A,k}[m-1] + \boldsymbol{W}_{A,k}[m], \tag{A.5}$$

$$\boldsymbol{S}_{\phi,k}[m] = \boldsymbol{F}_\phi\boldsymbol{S}_{\phi,k}[m-1] + \boldsymbol{W}_{\phi,k}[m], \tag{A.6}$$

where $\boldsymbol{S}_{A,k}[m]$ and $\boldsymbol{S}_{\phi,k}[m]$ are the states of instantaneous amplitude and instantaneous phase of $k$th channel respectively. Since, instantaneous amplitude and phase can be modeled with an auto regressive (AR) process of order $p$, the state vector can be represented as:

$$\boldsymbol{S}_{A,k}[m] = [S_{A,k}[m-p+1], \cdots, S_{A,k}[m]]^T, \tag{A.7}$$

$$\boldsymbol{S}_{\phi,k}[m] = [S_{\phi,k}[m-p+1], \cdots, S_{\phi,k}[m]]^T, \tag{A.8}$$

where $\boldsymbol{F}_A$ and $\boldsymbol{F}_\phi$ are the transition matrices that can be obtained by the linear prediction method. $\boldsymbol{W}_{A,k}[m]$ and $\boldsymbol{W}_{\phi,k}[m]$ are assumed to be Gaussian white noise of $k$th channel, and the variances of $\boldsymbol{W}_{A,k}[m]$ and $\boldsymbol{W}_{\phi,k}[m]$ are $Q_A$ and $Q_\phi$, respectively.

The observation equations for the instantaneous amplitude and instantaneous phase are defined as

$$\boldsymbol{O}_{A,k}[m] = \boldsymbol{H}_A\boldsymbol{S}_{A,k}[m] + \boldsymbol{V}_{A,k}[m], \tag{A.9}$$

$$\boldsymbol{O}_{\phi,k}[m] = \boldsymbol{H}_\phi\boldsymbol{S}_{\phi,k}[m] + \boldsymbol{V}_{\phi,k}[m], \tag{A.10}$$

where $\boldsymbol{O}_{A,k}[m]$ and $\boldsymbol{O}_{\phi,k}[m]$ are the observed instantaneous amplitude and phase of the noisy speech at time $m$ in $k$th channel, respectively. $\boldsymbol{H}_A$ and $\boldsymbol{H}_\phi$ are the observation

matrices, which are $[0, 0, \cdots, 1]$ in this research. $\boldsymbol{V}_{A,k}[m]$ and $\boldsymbol{V}_{\phi,k}[m]$ are observation noise (white Gaussian noise) and the variances of $\boldsymbol{V}_{A,k}[m]$ and $\boldsymbol{V}_{\phi,k}[m]$ are $R_A$ and $R_\phi$, respectively.

Five steps are used to calculate the optimal estimations for both instantaneous amplitude and instantaneous phase.

**Step 1:** The initial state vectors are set as $\boldsymbol{S}_{A,k}[1|1] = [10^{-12} \cdots 10^{-12}]$ and $\boldsymbol{S}_{\phi,k}[1|1] = [10^{-12} \cdots 10^{-12}]$. Then, it can be estimated the instantaneous amplitude and phase of clean speech at next time step from the initial state vector. Repeating this step, it is estimated the instantaneous amplitude and phase of clean speech of time $m$ from the optimal estimation of time $m - 1$.

$$\boldsymbol{S}_{A,k}[m|m-1] = \boldsymbol{F}_A \hat{\boldsymbol{S}}_{A,k}[m-1|m-1], \tag{A.11}$$

$$\boldsymbol{S}_{\phi,k}[m|m-1] = \boldsymbol{F}_\phi \hat{\boldsymbol{S}}_{\phi,k}[m-1|m-1]. \tag{A.12}$$

**Step 2:** The initial error covariance matrix is set as $\boldsymbol{P}_A[1|1] = \mathrm{diag}(R_A \cdots R_A)$ and $\boldsymbol{P}_\phi[1|1] = \mathrm{diag}(R_\phi \cdots R_\phi)$.

$$\boldsymbol{P}_A[m|m-1] = \boldsymbol{F}_A \boldsymbol{P}_A[m-1|m-1]\boldsymbol{F}_A^{\mathrm{T}} + Q_A \tag{A.13}$$

$$\boldsymbol{P}_\phi[m|m-1] = \boldsymbol{F}_\phi \boldsymbol{P}_\phi[m-1|m-1]\boldsymbol{F}_\phi^{\mathrm{T}} + Q_\phi \tag{A.14}$$

**Step 3:** We estimate the current value and smooth the previous value is estimated as follows.

$$\boldsymbol{S}_{A,k}[m|m] = \boldsymbol{S}_{A,k}[m|m-1] + \boldsymbol{e}_A, \tag{A.15}$$

$$\boldsymbol{S}_{\phi,k}[m|m] = \boldsymbol{S}_{\phi,k}[m|m-1] + \boldsymbol{e}_\phi, \tag{A.16}$$

where $\boldsymbol{G}_A[m]$ and $\boldsymbol{G}_\phi[m]$ are the Kalman gains. Here, $\boldsymbol{e}_A = \boldsymbol{G}_A[m] \times (\boldsymbol{O}_{A,k}[m] - \boldsymbol{H}_A \boldsymbol{S}_{A,k}[m|m-1])$ and $\boldsymbol{e}_\phi = \boldsymbol{G}_\phi[m](\boldsymbol{O}_{\phi,k}[m] - \boldsymbol{H}_\phi \boldsymbol{S}_{\phi,k}[m|m-1])$ are called innovation.

**Step 4:** Updating the Kalman gains.

$$\boldsymbol{G}_A[m] = \frac{\boldsymbol{P}_A[m|m-1]\boldsymbol{H}_A^{\mathrm{T}}}{(\boldsymbol{H}_A \boldsymbol{P}[m|m-1]\boldsymbol{H}_A^{\mathrm{T}} + R_A)} \tag{A.17}$$

$$\boldsymbol{G}_\phi[m] = \frac{\boldsymbol{P}_\phi[m|m-1]\boldsymbol{H}_\phi^{\mathrm{T}}}{(\boldsymbol{H}_\phi \boldsymbol{P}_\phi[m|m-1]\boldsymbol{H}_\phi^{\mathrm{T}} + R_\phi)} \tag{A.18}$$

**Step 5:** Updating the error covariances matrix.

$$\boldsymbol{P}_A[m|m] = (\boldsymbol{I} - \boldsymbol{G}_A[m]\boldsymbol{H}_A)\boldsymbol{P}_A[m|m-1] \tag{A.19}$$

$$\boldsymbol{P}_\phi[m|m] = (\boldsymbol{I} - \boldsymbol{G}_\phi[m]\boldsymbol{H}_\phi)\boldsymbol{P}_\phi[m|m-1] \tag{A.20}$$

where $\boldsymbol{I}$ is the unit matrix.

## A.1.2 Assumptions and Noise Modeling in Kalman Filter

In order to verify our proposed concept, we concentrate on the noise assumptions of Kalman filtering. The Kalman filter provides an optimal solution of the estimated parameters whenever, the observation noise and driving noise maintain the white and Gaussian noise properties. In order to ensure this point, we checked the statistical properties of noise in each sub-band. We calculated the normalized power spectrum density (PSD) and distribution of the noise in each sub-band. PSD shows how much power is contained in each of the spectral component. For white Gaussian noise, the PSD plot gives almost fixed power in all the frequencies. In addition, for white noise, the mean values of $\boldsymbol{W}_{A,k}[m]$, $\boldsymbol{W}_{\phi,k}[m]$, $\boldsymbol{V}_{A,k}[m]$ and $\boldsymbol{V}_{\phi,k}[m]$ are zero and variances denoted as $Q_A$, $Q_\phi$, $R_A$ and $R_\phi$ have finite value. We consider ten samples of noise as a statistical sample and verify the PSD and Gaussian distribution in each sub-band. To get the observation noise of amplitude and phase in sub-band, we use the following equations.

$$\boldsymbol{V}_{A,k} = \boldsymbol{S}_{A,k} - \boldsymbol{A}_{A,k} \tag{A.21}$$

$$\boldsymbol{V}_{\phi,k} = \boldsymbol{S}_{\phi,k} - \boldsymbol{A}_{\phi,k} \tag{A.22}$$

where, $\boldsymbol{A}_{A,k}$ and $\boldsymbol{A}_{\phi,k}$ are the clean amplitude and phase in $k$th sub-band. Driving noise is considered as an LP residue. As an example, the result of PSD of a particular sub-band ($k = 28$) for the observation noise of both the instantaneous amplitude and phase are shown in Figs. A.2(a) and A.3(a). The results show that the PSD of observation noise, both in amplitude and phase have flat constant power and we observe the flat PSD property in other sub-bands also. For driving noise we get the similar results also, which is shown in Figs. A.4(a) and A.5(a) as a white noise.

To verify the Gaussian property of the noise, we plot the histogram of the distribution. Figures A.2(b) and A.3(b), show the distribution of observation noise of amplitude and phase in a particular sub-band ($k = 28$) as an example. It is observed that, the observation noise in amplitude follows the Gaussian distribution. Although, the observation noise in phase, shown in Fig. A.3(b) does not perfectly follow the Gaussian distribution, but the Kalman filter can still provide a better solution. The distribution for driving noise of
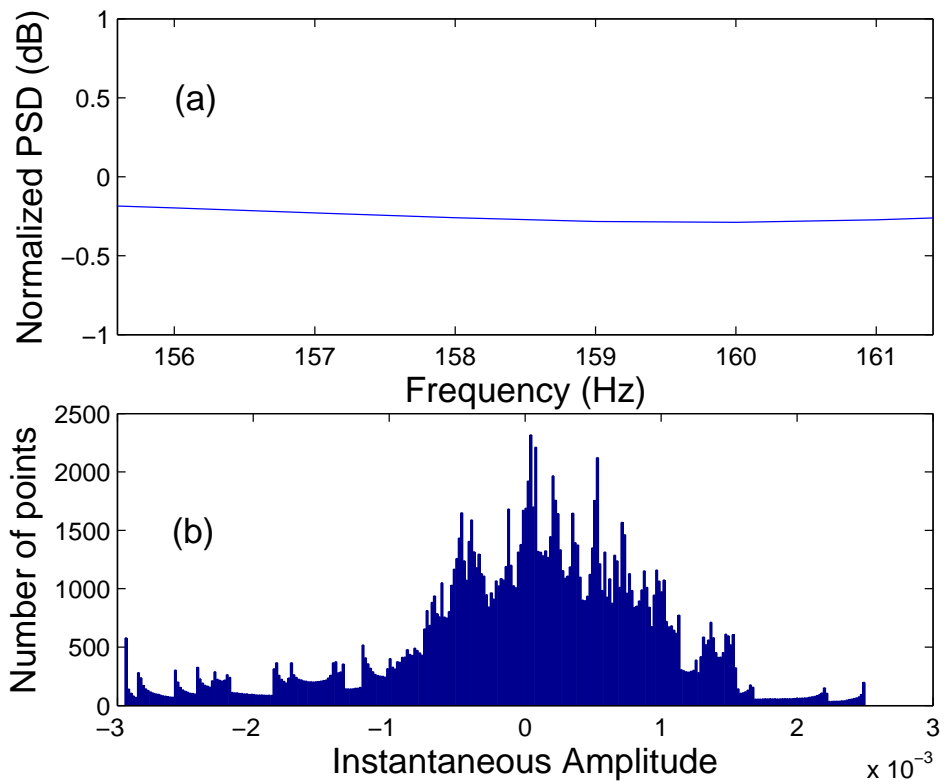
Figure A.2: Analysis results of observation noise in Eq. (A.9): (a) normalized PSD of $\boldsymbol{V}_{A,k}$ and (b) distribution of $\boldsymbol{V}_{A,k}$ in 28th channel.
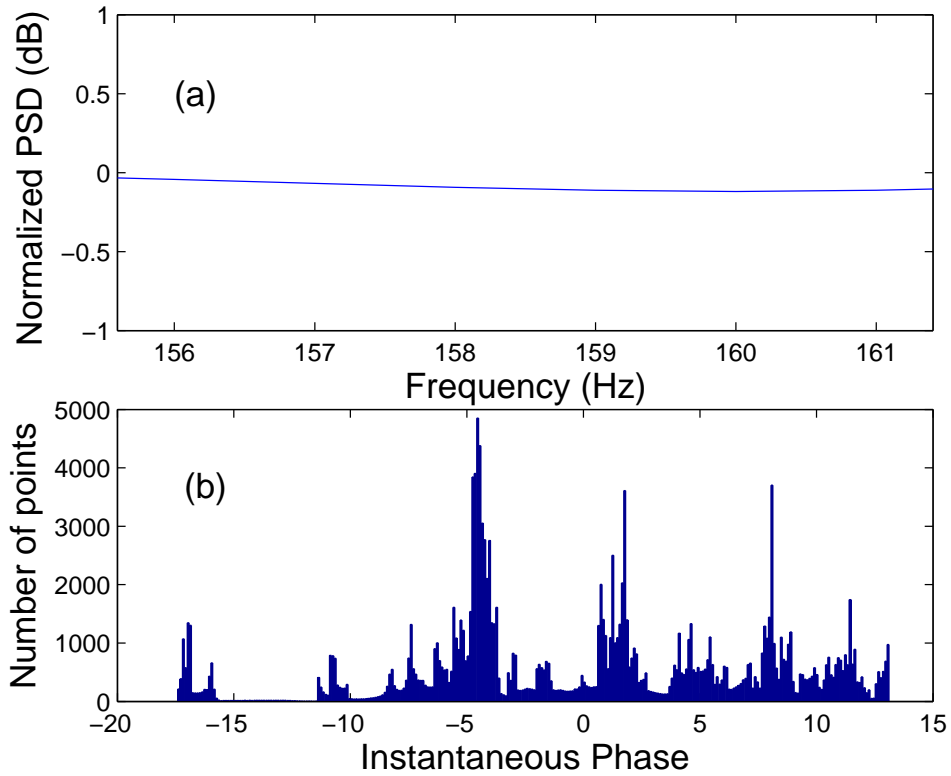


Figure A.3: Analysis results of observation noise in Eq. (A.10): (a) normalized PSD of $\boldsymbol{V}_{\phi,k}$ and (b) distribution of $\boldsymbol{V}_{\phi,k}$ in 28th channel.
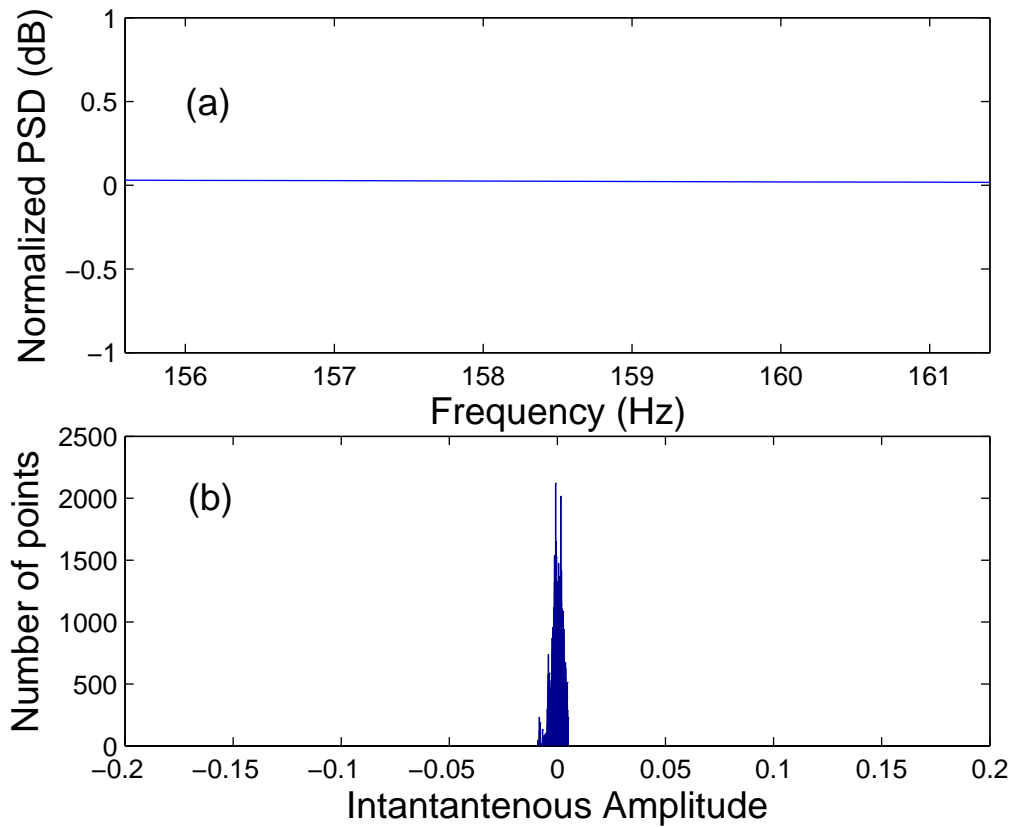
Figure A.4: Analysis results of driving noise in Eq. (A.7): (a) normalized PSD of $\boldsymbol{W}_{A,k}$ and (b) distribution of $\boldsymbol{W}_{A,k}$ in 28th channel.

instantaneous amplitude and phase are shown in Figs. A.4(b) and A.5(b) in a specific sub-band ($k = 28$). The results show that both amplitude and phase perfectly follow Gaussian distribution.

Due to space limitation, we show the result of observation and deriving noise verification in only a specific sub-band. But, we analyze the result on each sub-band and confirm that in all channels the observation and deriving noise in amplitude and phase are white noise since their PSD maintain a flat constant power. In addition, for other sub-bands, the distribution of observation noise in amplitude also show the Gaussian distribution. Although the observation noise in phase does not perfectly Gaussian, we observe that it maintain summation of Gaussian in distribution as similar to Fig. A.3(b). Moreover, the distribution of the deriving noise on amplitude and phase in all channel provide Gaussian distribution.

Hence, in general, both driven noise in Eqs. (A.7) and (A.8) and observation noise in

Figure A.5: Analysis results of driving noise in Eq. (A.8): (a) normalized PSD of $\boldsymbol{W}_{\phi,k}$ and (b) distribution of $\boldsymbol{W}_{\phi,k}$ in 28th channel.

Eqs. (A.9) and (A.10) can be regarded as white Gaussian noise.

# A.2 Algorithm Implementation

## A.2.1 An Auditory-motivated Filterbank

We used an auditory motivated Gammatone filterbank which is designed by considering the properties of auditory system and detection of a discontinuous point dealing with the complex spectrum. For details of filterbank design please see [85].

## A.2.2 Calculation of Instantaneous Amplitude and Phase

We extract the instantaneous amplitude $S_{A,k}[m]$ and instantaneous unwrapped phase $S_{\phi,k}[m]$ of Eqs.(5) and (6) by using Gammatone filterbank. Instantaneous amplitude

Figure A.6: Example of (a) instantaneous amplitude $S_{A,k}[m]$ and (b) instantaneous unwrapped phase $S_{\phi,k}[m]$ extraction in a sub-band (channel $k$=28) using Gammatone filterbank.

$S_{A,k}[m]$ and unwrapped phase $S_{\phi,k}[m]$ are calculated as follows

$$S_{A,k}[m] = |\tilde{f}(c,m)| \tag{A.23}$$

$$S_{\phi,k}[m] = \int_0^m \left( \frac{\mathrm{d}}{\mathrm{d}\tau} arg(\tilde{f}(c,m) - \omega_k) \right) \mathrm{d}\tau \tag{A.24}$$

where, $c = \alpha^{k-K/2}$, $\alpha$ is the scale of Gammatone filterbank. $|\tilde{f}(c,m)|$ is the amplitude spectrum defined by the wavelet transform and arg $(\tilde{f}(c,m))$ is the unwrapped phase spectrum defined by the complex wavelet transform. An example of instantaneous amplitude and instantaneous unwrapped phase extraction in a particular sub-band ($k = 28$) is shown in Fig. A.6.

## A.2.3 Linear Prediction

We developed a linear prediction method to extract the LP coefficients for Kalman filtering from the clean speech, which is referred as to non-blind Kalman filtering. However, in real

Figure A.7: Block diagram of proposed scheme (blind method) for speech enhancement.

life clean speech may not be available for LP extraction. Thus, before applying Kalman filtering we incorporate a training phase on closed data set to estimate the LP coefficients, which is known as blind Kalman filtering. These trained LP coefficients are incorporated whenever the clean speech is not available.

**Non-blind Linear Prediction**

In the non-blind linear prediction method, we assume that the sampling sequences of clean speech's amplitude and phase are $S_{A,k}[m]$ and $S_{\phi,k}[m]$, where $m = 0, 1, 2, \cdots,$ $M$. These can be regarded as the output of a $p$-th order AR process. The models of linear prediction can be represented as:

$$S_{A,k}[m] = \sum_{i=1}^{p} a_i S_{A,k}[m - i] \tag{A.25}$$

$$S_{\phi,k}[m] = \sum_{i=1}^{p} b_i S_{\phi,k}[m - i]. \tag{A.26}$$

Here, $S_{A,k}[m]$ and $S_{\phi,k}[m]$ are the optimal estimation of $S_{A,k}[m]$ and $S_{\phi,k}[m]$ under the principle of the minimum mean-square error (MMSE), $\{a_1, a_2, \cdots, a_p\}$ and $\{b_1, b_2, \cdots, b_p\}$ are LP coefficients and $p$ is the prediction order.

There are two types of methods for calculating the LP coefficients: an auto-correlation (AC) method and a covariance method. We chose the AC method to calculate the LP coefficients, and $\{a_i\}$ $(i = 1, 2, \cdots, p)$ and $\{b_i\}$ $(i = 1, 2, \cdots, p)$ could be obtained by solving the Yule-Walker equation as

$$R[q_a] - \sum_{i=1}^{p} a_i R[q_a - i] = 0 \tag{A.27}$$

$$R[q_b] - \sum_{i=1}^{p} b_i R[q_b - i] = 0. \tag{A.28}$$

Here, $R[q_a]$ and $R[q_b]$ are the AC functions of instantaneous amplitude and phase of the clean speech, $S_{A,k}[m]$ and $S_{\phi,k}[m]$, $R[q_a] = E\{S_{A,k}[m]S_{A,k}[m - q_a]\}$ and $R[q_b] = E\{S_{\phi,k}[m]S_{\phi,k}[m - q_b]\}$, where $E\{\cdot\}$ is the expectation. Then, we can obtain the transition matrix $\boldsymbol{F}_A$ and $\boldsymbol{F}_\phi$ for estimating instantaneous amplitude and instantaneous phase by using the Kalman filter.

**Blind Linear Prediction**

In the blind linear prediction, we assume that clean speech is not available for LP estimation. The block diagram of the proposed blind method is shown in Fig. A.7. As stated before, the effectiveness of Kalman filtering depends on the proper estimation of LP coefficients. To estimate the LP coefficients, we study the properties of LP coefficients on modulation domain for amplitude and phase. From this investigation, we found that LP coefficients on modulation domain had some similarities on values. We also checked that this similarities were not sensitive to the specific gender, individual speaker or the contents of the speech. This property is quite different from the LP properties of spectral envelop. The example of LP coefficient similarities on modulation domain for three different speakers is shown in Fig. A.8.

Based on the LP characteristics on modulation domain, we incorporate an offline training phase on closed data set to train the LP coefficient for blind Kalman filtering. We calculate the LP coefficient of each sub-band from the closed clean data and convert it in line spectral frequencies (LSF), which is an an alternative representation to the LP parameters. We average the computed LSF and convert to LP coefficient as a train LP coefficient.

Figure A.8: Example of spectrum analysis of LP coefficient similarities on three different speakers and contents: (a) instantaneous amplitude and (b) phase.

It is known that, the LSF have some useful properties over LP coefficient. LP parameters have a large dynamic range of values, which is not good for quantization [87]. The LSF on the other hand, have a well behaved dynamic range. It is easier to guarantee the stability of the resulting synthesis filter in the case of interpolation in the LSF domain. Whenever, the LP coefficients are encoded as LSFs, we do not need to spend the same number of bits for each LSF. This is because higher LSFs correspond to the high frequency components and high frequency components have less effect in speech perception. So higher LSFs can be quantized using fewer bits than lower LSFs. This reduces the bit rate while keeping the speech quality almost the same. In the LPC analysis of speech, a short segment of speech is assumed to be generated as the output of an all-pole filter $H(z) = 1/A(z)$, where, $A(z)$ is the inverse filter is defined as

$$A(z) = 1 + a_1 z^{-1} + ... + a_p z^{-p} \qquad (A.29)$$

In order to define LSF, the inverse filter polynomial is used to construct two polynomials:

$$P(z) = A(z) + z^{-(M+1)}A(z^{-1}), \tag{A.30}$$

$$Q(z) = A(z) - z^{-(M+1)}A(z^{-1}) \tag{A.31}$$

The roots of the polynomials are called as LSF.

## A.3 Evaluation and Discussion

To evaluate the effectiveness of both proposed methods, we carried out experiments using 9 different English sentences uttered by male and female speakers from the Texas Instruments and Massachusetts Institute of Technology (TIMIT) database. Half of the utterances were selected from male and another half were selected from female speakers from the open data set. Before doing this, we create a closed data set, containing five sentences from two male and three female speakers of TIMIT database to trained the LP coefficient. We define other data set as a open data set. We use white, Babble and pink noise to evaluate the proposed scheme. The signal to noise ratios (SNRs) for white noise between $x(t)$ and $n(t)$ were fixed at from 20 dB to $-10$ dB at intervals of 10 dB. All noisy signals $y(t)$ were generated by adding $x(t)$ with $n(t)$. We used a Gammatone filterbank [85] to divide the signal into 128 channels ($K = 128$). We used the sampling frequency ($F_s$) of 20 kHz. We utilized a 25-ms-long rectangular window. The LP order, $p$, was set to 12.

We have evaluated the improvement of the restored speech by measuring correlation and signal to error ratio (SER). Correlation shows the similarity between the shapes of clean instantaneous amplitude and phase and restored instantaneous amplitude and phase and SER shows the level of the error that we can reduce. Correlation and SER are defined as follows

$$\text{Corr}(x_k, \hat{x}_k) = \frac{\int_0^T \left(x_k(t) - \overline{x_k}\right)\left(\hat{x}_k(t) - \overline{\hat{x}_k}\right) dt}{\sqrt{\left\{\int_0^T (x_k(t) - \overline{x_k})dt\right\}\left\{\int_0^T (\hat{x}_k(t) - \overline{\hat{x}_k})dt\right\}}}, \tag{A.32}$$

$$\text{SER}(x_k, \hat{x}_k) = 10\log_{10}\frac{\int_0^T (x_k(t))^2 dt}{\int_0^T (x_k(t) - \hat{x}_k(t))^2 dt}, \tag{A.33}$$

where $x_k(t)$ is the clean speech of $k$th channel and $\hat{x}_k(t)$ is the restored speech of $k$th channel.

Figure A.9: Improvements in restoration accuracy of the non-blind Kalman filter method: (a) improved Corr. and (b) improved SERs. SNR = 20 dB to −10 dB.

Figure A.9 shows the improvement in correlation and SER in each channel using non-blind method (non-blind Kalman filtering) under the mentioned white noise conditions. In the figure, the height of the bar indicates the mean value of the improvement in SER. All the channels have positive improvement in SER in 20, 10, 0, and −10 dB noise conditions, except with case of higher channels in 20-dB condition. This is because signal components in higher channels in 20-dB condition are almost similar to those of clean signal. Thus, it is easy to see that the proposed method can effectively reduce the noise in both instantaneous amplitude and phase.

The performance of the blind method (blind Kalman filtering) is shown in Fig. A.10. The results prove that blind Kalman filter with the trained LP coefficients also works well, thus we always obtain positive improvements in correlation and SER in all noise conditions, except with the same case mentioned in the above. From the comparison of result between non-blind and blind Kalman filtering, it is observable that, we achieve almost same improvement in correlation and SER. This is because, our trained LP coefficients act as a clean LP coefficients and it can be used as gender and content independent
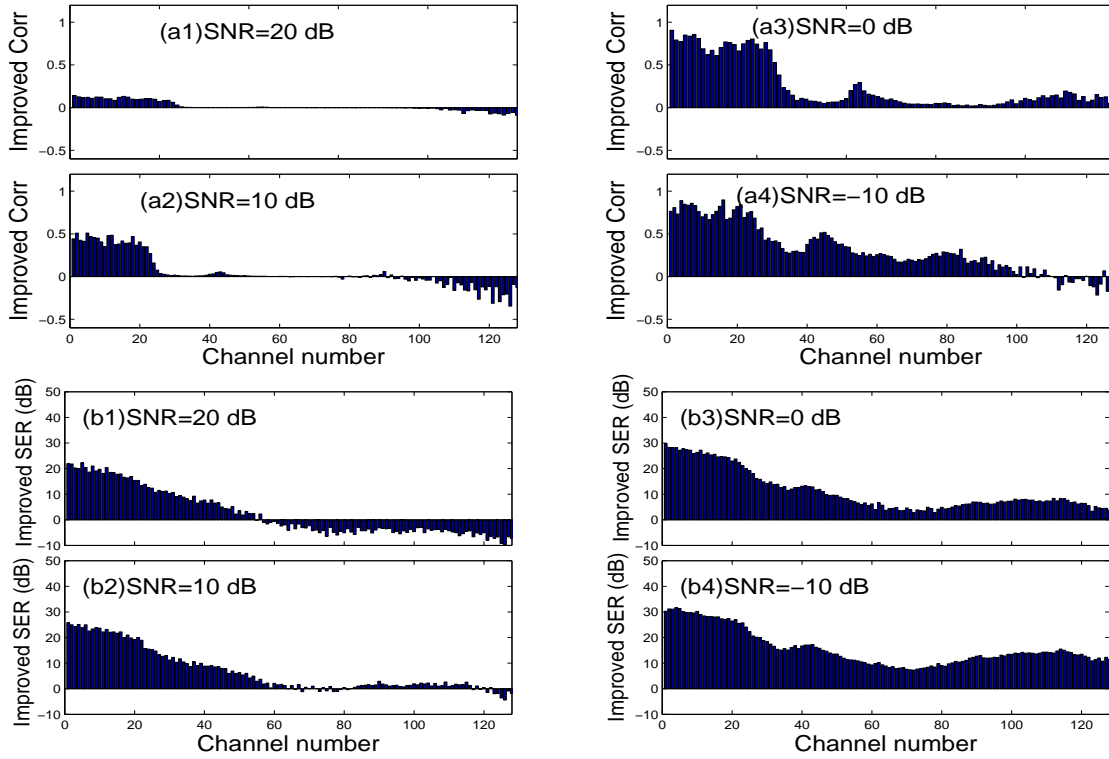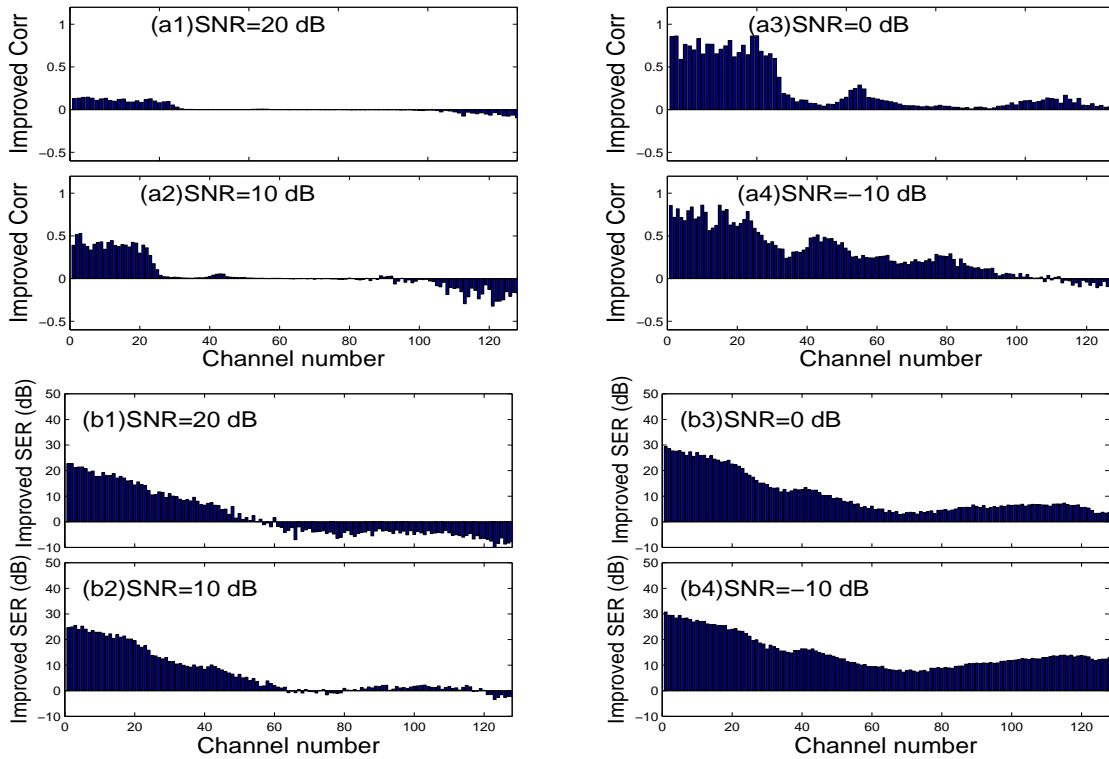
84

Figure A.10: Improvements in restoration accuracy of the blind Kalman filter method: (a) improved Corrs and (b) improved SERs. SNR = 20 dB to −10 dB.

LP coefficients. Figure A.11 shows the example of restored instantaneous amplitude and phase in a particular sub-band by the proposed method (blind Kalman filtering). We can observe that the restored amplitude and phase are matched with the clean amplitude and phase.

Moreover, we also choose the Wiener filtering method (Scalart-Filho algorithm) under the same conditions to compare its effectiveness with that of our proposed method. Based on the results in Fig. A.12, we can see that both of our proposed methods can obviously improve the SER and Corr much more than the Wiener filtering method.

To evaluate the quality and intelligibility of the restored speech, we calculated the perceptual evaluation of sound quality (PESQ) [88] and SNR loss [89] for all stimuli that we used the above evaluations. PESQ in the objective difference grades (ODGs) that covers from −0.5 (very annoying) to 4.5 (imperceptible) was used to evaluate subjective quality. SNR loss that ranges from 0.0 to 1.0 was used to evaluate intelligibility of speech. SNR losses (0 to 1.0) are corresponded to the percent correctness (100% to 0%). The results of objective measures are listed in Table A.2. The results indicate that both of

Figure A.11: Example of comparison among (a) Clean, Restored and Noisy instantaneous amplitude and (b) Clean, Restored and Noisy instantaneous unwrapped phase in a sub-band(channel $k = 28$) by proposed blind Kalman filtering. SNR$= -10$ dB noise (white).

our proposed methods provide better quality and improved intelligibility in the restored speech much more than the existing speech enhancement methods. From the result of evaluations, we can say that the proposed method can effectively reduce the noise from both the amplitude and phase and also improve the quality and speech intelligibility.

## A.3.1 Effectiveness of Phase

To evaluate the effectiveness of phase in speech enhancement, we restore the instantaneous amplitude by applying the blind Kalman filtering and combine with the noisy instantaneous phase to get the enhanced speech. The result of amplitude restoration accuracy in terms of correlation and SER is shown in Fig. A.13. From the comparison of results between without incorporating phase in Fig. A.13 and with restoring phase in Fig. A.10, it is easily observed that, for better speech enhancement we need to incorporate both

Figure A.12: Improvements in restoration accuracy of the Wiener filter method: (a) improved Corrs and (b) improved SERs. SNR= 20 dB to −10 dB.

amplitude and phase. In all SNR conditions, we get better positive improvement by restoring both amplitude and phase as shown in Fig. A.10. The same result is observable for quality and intelligibility improvement also. That is, we get the restored speech with better quality and intelligibility whenever, restoration is applied both on amplitude and phase as listed in Table A.2. We compare the PESQ and SNR loss among the restored speech (both amplitude and phase), amplitude only restored speech and phase only restored speech. From the results of evaluations, we can say that the for better speech enhancement we need to consider both the amplitude and phase to reduce the noise and to improve the quality and intelligibility effectively.

## A.3.2 Evaluation on Pink and Babble Noise

We carried out our simulation in pink and Babble noise condition. The restoration accuracy of the proposed blind Kalman filtering on pink and Babble noise condition are shown in Figs. A.14 and A.15. The proposed method can reduce the error and improve the correlations in the pink and Babble noise in all channels except some higher channels.

Figure A.13: Improvements in restoration accuracy of amplitude only using the blind Kalman filter method: (a) improved Corr. and (b) improved SERs. SNR= 20 dB to −10 dB.

We investigate that, the signal components in higher channels are almost similar to clean speech and have a very high SNR condition. Thus, it is very difficult to reduce noise in higher channel since it is almost similar to the clean speech.

## A.4   Conclusion

We proposed a speech enhancement scheme by using the Kalman filter with/without training phase in the sub-bands as the ASM on the Gammatone filterbank. We presented a greater speech enhancement scheme by verifying the assumptions on noise and provided a concrete enhancement scheme by improving the instantaneous amplitude and instantaneous phase simultaneously. We concentrate on central processing of Kalman filtering by providing an the effective way to train the LP coefficients for blind Kalman filtering. Our simulation results revealed that the proposed scheme with training phase performs almost like non-blind Kalman filtering and also outperforms the existing conventional speech en-

Figure A.14: Improvements in restoration accuracy of the blind Kalman filter method in pink noise condition: (a) improved Corr. and (b) improved SER. SNR= $-2.07$ dB.



Figure A.15: Improvements in restoration accuracy of the blind Kalman filter method in babble noise condition: (a) improved Corr. and (b) improved SER. SNR= $-5.60$ dB.

hancement algorithms in terms of improvements for speech quality and intelligibility. We believe that this is the effect of combining phase enhancement with amplitude enhancement and accurate LP estimation in the sub-band representation. Although the proposed method always shows improvement, but still, future research is required to ensure improvement in all channels for high SNR conditions. We will plan to investigate whether the proposed scheme can be extended to be a speech enhancement in noisy reverberant environments.

Table A.1: Comparison of result of PESQ and SNR loss (averaged values).

| | Methods | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Noisy Speech | | Proposed (Non-blind) | | Proposed (Blind) | | SS [60] | | MMSE [61] | | Wiener filter [62] | |
| SNR | PESQ | SNR loss | PESQ | SNR loss | PESQ | SNR loss | PESQ | SNR loss | PESQ | SNR loss | PESQ | SNR loss |
| 20 dB | 2.98 | 0.57 | 3.62 | 0.53 | 3.50 | 0.54 | 3.01 | 0.82 | 2.74 | 0.84 | 2.79 | 0.882 |
| 10 dB | 2.31 | 0.74 | 3.61 | 0.60 | 3.06 | 0.62 | 2.41 | 0.78 | 2.11 | 0.93 | 1.96 | 0.94 |
| 0 dB | 1.65 | 0.87 | 3.10 | 0.67 | 2.47 | 0.71 | 1.75 | 0.91 | 1.68 | 0.94 | 1.77 | 0.97 |
| −10 dB | 1.12 | 0.95 | 2.84 | 0.73 | 1.62 | 0.82 | 1.09 | 0.90 | 1.01 | 0.94 | 1.34 | 0.96 |

Table A.2: Comparison of restored speech with amplitude only restoration and phase only restoration (averaged values).

| | Methods | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Noisy Speech | | Restored Speech (Blind) | | Amplitude only Restoration (Blind) | | Phase only Restoration (Blind) | |
| SNR | PESQ | SNR loss | PESQ | SNR loss | PESQ | SNR loss | PESQ | SNR loss |
| 20 dB | 2.98 | 0.57 | 3.50 | 0.54 | 3.46 | 0.55 | 3.08 | 0.61 |
| 10 dB | 2.31 | 0.74 | 3.06 | 0.62 | 2.97 | 0.66 | 2.44 | 0.73 |
| 0 dB | 1.65 | 0.87 | 2.47 | 0.71 | 2.39 | 0.77 | 1.83 | 0.85 |
| −10 dB | 1.12 | 0.95 | 1.62 | 0.82 | 1.47 | 0.88 | 1.30 | 0.94 |

# Appendix B

# List of Abbreviations

| | |
|---|---|
| ACF | Auto Correlation Coefficient |
| AR | Auto Regressive |
| ARMA | Auto Regressive Moving Average |
| ARIMA | Auto Regressive Integrated Moving Average |
| CUSUM | Cumulative Sum Model |
| CPS | Cyber-Physical Systems |
| ESDR | Efficient Spatial Data Recovery |
| ETSDR | Efficient Temporal Spatial Data Recovery |
| ETSDR/EM | Efficient Temporal Spatial Data Recovery with Error Minimization |
| EWMA | Exponentially Weighted Moving Average |
| IAE | Integral of Absolute Error |
| IoT | Internet of Things |
| ITS | Intelligent Transportation Systems |
| KF | Kalman Filter |
| NCS | Networked Control System |
| NMAR | Not Missing at Random |
| MA | Moving Average |
| MAE | Mean Absolute Error |
| M2M | Machine-to-Machine |
| MAR | Missing at Random |
| MCAR | Missing Completely at Random |
| PCC | Pearson Correlation Coefficient |
| PID | Proportional Integral Derivative |
| PACF | Partial Auto Correlation Coefficient |
| QoR | Quality of Result |
| RMSE | Root Mean Square Error |
| RTS | Real Time Systems |
| STI | Spatial Temporal Imputation |
| TCP | Transport Control Protocol |
| V2I | Vehicle-to-Infrastructure |
| V2V | Vehicle-to-Vehicle |
| WP | Weighted Prediction |
| WSAN | Wireless Sensor and Actuator Network |

# Appendix C

# List of Symbols

| C+ | Upper Cusums |
|---|---|
| C- | Lower Cusums |
| $\delta$ | Threshold |
| N | number of samples |
| $r_k$ | k-order auto-correlation coefficient |
| $\rho_{xy}$ | PCC coefficient correlation of x and |
| I | Morans I |
| GC | Geary's C |
| $w_{jk}$ | matrix of spatial weights |
| $e_0$ | error offset |
| W | window size |
| MC | Maximum number of consecutive missing data |
| $d_{im}$ | Input measured data |
| $d_{est.}$ | estimated data |
| $SR_{th}$ | maximum tolerable error value |
| $e_j$ | mean error between neighbors measured data and neighbors model generated data for all the k n |
| p | order for AR model |
| q | order for MA model |
| $W_s$ | process error |
| $V_s$ | observation error |
| A | state transition matrix |
| H | observation matrix |

# References

[1] A.L. Edward, "Cyber physical systems: Design challenges," IEEE Symp. on Object Oriented Real-Time Distributed Computing, pp.363–369, (2008).

[2] E. A. Lee, "CPS Foundations", ACM/IEEE Design Automation," Conference (DAC), pp. 737-742, (2010).

[3] R. Rajkumar, I. Lee, L. Sha and J. Stankovic: "Cyber-physical System: The Next Computing Revolution", 47th ACM/IEEE, Design Automation Conf. (DAC), pp. 731-736, (2010).

[4] B. H. Krogh: "Cyber Physical Systems: The Need for New Models and Design Paradigms", Presentation Report, (2008).

[5] J. Shi, J. Wan, H. Yan and H. Suo: "A Survey of Cyber-physical Systems", Int. Conf. on Wireless Communications and Signal Processing, pp. 1-6, (2011).

[6] K. Wan, D. Hughes, K. L. Man, and T. Krilavicius: "Composition Challenges and Approaches for Cyber Physical Systems", IEEE Int. Conf. on Networked Embedded Systems for Enterprise Applications (NESEA), pp. 1-7, (2010).

[7] E. Yeniaras, J. Lamaury, Z. Deng, and N.V. Tsekos, "Towards a new cyber-physical system for MRI-guided and robot-assisted cardiac procedures," IEEE Int. Conf. on Information Technology and Applications in Biomedicine, pp.1–5, November, (2010).

[8] http://www.nsf.gov/pubs/2014/nsf14542/nsf14542.htm.

[9] A. Gokhale, S. Tambe, L. Dowdy, and G. Biswas, "Towards High Confidence Cyberphysical Systems for Intelligent Transportation Systems", Department of EECS, Vanderbilt University, Nashville, TN (2008): 3.

[10] A. Gokhale, MP. McDonald, S. Drager, and W. McKeever, "A Cyber physical Systems perspective on the real-time and reliable dissemination of information in Intelligent Transportation Systems", Air force research lab Rome NY.

[11] S. A. Haque, S.M. Aziz, M. Rahman, "Review of Cyber-Physical System in Healthcare", International Journal of Distributed Sensor Networks, (2014).

[12] D. Min, Medical Cyber Physical Systems and Big data Platforms.

[13] High-Confidence Medical Devices: Cyber-Physical Systems for 21st Century Health Care [Online], Feb. 2009. Available at http://www.nitrd.gov/About/MedDevice-FINAL1-web.pdf

[14] S. Karnouskos, "Cyber-physical systems in the smart grid", 9th IEEE International Conference in Industrial Informatics (INDIN, pp. 20–23, 2011.

[15] F.J. Wu, Y.F. Kao, and Y.C. Tseng,"From wireless sensor networks towards cyber physical systems," J. Pervasive and Mobile Comp., vol.7, no.4, pp.397–413, 2011.

[16] F. Martincic, and L. Schwiebert, Introduction to wireless sensor networking, Handbook of Sensor Networks-Algorithms and Architectures, John Wiley Sons, New York, USA, 2005.

[17] W. Young, G. Weckman, and W. Holland, "A survey of Methodologies for the Treatment of missing values within datasets. Limitations and benefits", Theoretical Issues in Ergonomics Science vol(12)1, pp.15–43, 2011.

[18] O. Harel, and X. Zhou, "Multiple Imputation: Review of Theory, implementation and software", Statistics in Medicine, vol. 26(16), pp. 3057–3077, 2007.

[19] G. Box, G. Jenkins, G. Reinsel, Time Series Analysis: Forecasting and Control, 4th edn., NJ: Wiley, pp. 47–92, 2008.

[20] JP. Hespanha, P. Naghshtabrizi, Y. Xu, Y, " A Survey of recent results in networked Control Systems", Proc. IEEE vol.95(1), pp.138–162 2007.

[21] R.J.A. Little and D.B. Rubin, Statistical Analysis with Missing Data, 2nd Ed., Wiley-Interscience, New York, 2002.

[22] D.C. Howell. University of Vermont, Treatment of missing data 2012 .

[23] J.G. Ibrahim, H. Zhu and N. Tang, "Model selection criteria for missing-data problems using the EM algorithm," J. American Statistical Association. pp.1648-1658, 2008.

[24] H.Y. Chen, H. Xie and Y. Qian, "Multiple imputation for missing values through conditional semi parametric odds ratiomodels", J. Biometrics vol.67, no.3, pp.799-809, 2011.

[25] J.M.I. Molina, P.J. Garcia-Laencina, E. Alba, N. Ribelles, M. Martin and L. Franco, "Missing data imputation using statistical and machine learning methods in a real breast cancer problem", J. Artificial Intelligence in Madicine, Elsevier Science Publishers, vol.50, no.2, pp.105-115, 2010.

[26] C.M. Bishop, Pattern Recognition and Machine Learning (Information Science and Statistics), Springer Science Business Media, 2007.

[27] T. Kohonen, Self-organizing Maps, Springer Series in Information Sciences, Springer-Verlag; 3rd Ed., 2001.

[28] C.C. Huang and H.M. Lee, "A grey-based nearest neighbor approach for missing attribute value prediction", J. Applied Intelligence, vol.20, no.3, pp.239-252, 2004.

[29] W. Bajwa, "Compressive wireless sensing", ACM Conf. on Inf. Processing in Sensor Networks, pp.134-142, 2006.

[30] D. Guo, X. Qu, L. Huang and Y. Yao., "Sparsity-based spatial interpolation in wireless sensor networks", J. Sensors, vol.11, no.3, pp.2385-2407, 2011.

[31] Y.Y. Li and L.E Parker, Classification with missing data in wireless sensor network", IEEE Southeastcon, pp.533-538, April 2008.

[32] F. Xia, X. Kong and Z. Xu, "Cyber-physical control over wireless sensor and actuator networks with packet loss", Wireless Networking Based Control, Springer, pp.85-102, 2011.

[33] R.H. Choi, S.C. Lee, D.H. Lee and J. Yoo, "WiP abstract: Packet loss compensation for cyber-physical control systems", IEEE/ACM Int. Conf. on Cyber-Physical Systems (ICCPS), pp.205, 2012.

[34] L. Tang, X. Yu, S. Kim, Q. Gu, J. Han, A. Leung, T. L. Porta, "Trustworthiness analysis of sensor data in cyber-physical systems", J. of Comp. and System Sciences, vol.79, pp.383–401, 2013.

[35] L. Tang, Q. Gu, X. Yu, J. Han, T.L. Porta , A. Leung, T. Abdelzaher and L. Kaplan, "IntruMine: Mining Intruders in Untrustworthy Data of Cyber-Physical Systems", Int. Conf. on Data Mining (SDM), pp. 600–611, 2012.

[36] L. Tang, X. Yu, Q. Gu, J. Han, A. Leung, and T.L. Porta,, "Mining Lines in the Sand: On Trajectory Discovery From Untrustworthy Data in Cyber-Physical System", ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD), pp. 410–418 , 2013.

[37] L. Tang, Q. Gu, S. Kim, J. Han, W. Peng, Y. Sun, A. Leung, and T. L. Porta, "Multidimensional sensor data analysis in cyber-physical system: An atypical cube approach", Int. J. of Distributed Sensor Network, vol.2012, pp.1–19, 2012.

[38] G. Li, and Y. Wang, "Automatic ARIMA modeling-based data aggregation scheme in wireless sensor networks", EURASIP J. Wireless Comm. and Networking, vol.2013(85), pp.1–13, 2013.

[39] C. Vuran, C. Mehmet, B.A. Ozgur, and I F. Akyildiz, "Spatio-temporal correlation: theory and applications for wireless sensor networks", Computer Networks 45.3 (2004): 245-259.

[40] Paradis, Emmanuel, "Moran's Autocorrelation Coefficient in Comparative Methods", R Foundation for Statistical Computing, Vienna (2009).

[41] http://faculty.salisbury.edu/ ajlembo/419/lecture15.pdf

[42] Y. Ke, J. Cheng and J.X. Yu, " Efficient discovery of frequent correlated subgraph pairs", IEEE Int. Conf. on Data Mining (ICDM), pp.239-248, 2009.

[43] L.C Alwan, and V. H. Roberts, "Time-series modeling for statistical process control", Journal of Business and Economic Statistics 6.1 (1988): 87–95.

[44] B.L. Bowerman, R.T. O' Connell, Forecasting and Time Series: An Applied Approach, China Machine Press, Beijing, 2003.

[45] Cohen, Leon. Time-frequency analysis. Vol. 778. Englewood Cliffs, NJ:: Prentice Hall PTR, 1995.

[46] D. Guo, X. Qu, L. Huang and Y. Yao., "Sparsity-based spatial interpolation in wireless sensor networks", J. Sensors, vol.11, no.3, pp.2385-2407, 2011.

[47] G.Y. Lu and D.W. Wong, "An adaptive inverse-distance weighting spatial interpolation technique", J. Comput. Geosci., vol.34, pp.1044-1055, 2008.

[48] M. Umer, L. Kulik and E. Tanin, "Kriging for localized spatial interpolation in sensor networks , Int. Conf. on Scientific and Statistical Database Management, pp.525-532, 2008.

[49] C.J. Wilmott and K. Matsuura, "Advantages of the mean absolute error (MAE) over root means square (RMSE) in assessing average model performance , Climate Research, vol.30, pp.79-82, 2005.

[50] G.E. Ljung, G.E.P. Box, On a measure of lack of fit in time series models, Biometrika, 1978.

[51] R.J. Hyndman, "Yule-Walker estimators for continuous-time autoregressive models", J. of Time Series Analysis, vol.14(3), pp.281–296, 1993.

[52] G.B.Moody, R.G. Mark, and A.L. Goldberger, "PhysioNet: a Web-based resource for study of physiologic signals", IEEE Trans. on Eng. in Medicine and Biology, vol.20, no.3, pp:70–75,2001.

[53] A. Azadeh, S.M. Asadzadeh, R.J. Marandi, S.N. Shirkouhi, G.B. Khoshjhou and S. Talebi, Optimal estimation of missing values in randomized complete block design by genetic algorithm, Knowledge-Based Systems, Elsevier, vol.37, pp.3747, 2013.

[54] Welch, G., Bishop, G.: An introduction to the Kalman filter. (1995).

[55] Xia, F., Mukherjee, T., Zhang, Y., Song, Y.: Sensor Networks for High-Confidence Cyber-Physical Systems. Int. J. of Distributed Sensor Networks Volume 2011 (2011).

[56] Laplante, P.A.: Real-time systems design and analysis An Engineers Handboook. 2nd ed. IEEE Press, Piscataway, NJ, USA, http://www.ieee.org, 1997.

[57] Kremer, U.: Cyber-Physical Systems: A case for soft real-time.

[58] Nower, N., Yasuo, T. and Lim, A.O.: Traffic Pattern based Data Recovery Scheme for Cyber-physical System, IEICE Trans. on Fundamental E97-A(9), 1926–1936 (2014)

[59] N.Nower, T.Yasuo, A.O. Lim, "Efficient Spatial Data Recovery Scheme for Cyber-physical System", IEEE Int. Conf. on Cyber-Physical Systems, Networks and Applications, pp.92–97, Taipei, Taiwan, August 2013.

[60] Boll, S., 1979. Suppression of Acoustic Noise in Speech Using Spectral Subtraction. IEEE Trans. Acoust. Speech Signal Process., ASSP-27 (2), 113–120.

[61] Ephraim, Y., Mlah, D., 1984. Speech enhancement using a minimum-mean square error short-time spectral amplitude estimator. IEEE Trans. Acoust. Speech Signal Process., ASSP-32 (6), 1109–1211.

[62] Scalart, P., Filho, J. V., 1996. Speech enhancement based on a priori signal to noise estimation. In Proc. IEEE Int. Conf. Acoustics, Speech, and Signal Process., (ICASSP), 629–623.

[63] Ephraim, Y., Malah D., Juang, B-H., 1989. On the application of hidden Markov models for enhancing noisy speech. IEEE Trans. Acoust. Speech Signal Process., ASSP-37 (12), 1846–1856.

[64] Ephraim, Y., 1992. A Bayesian estimation approach for speech enhancement using hidden Markov models. IEEE Trans. Signal Process., 40 (4), 725–735.

[65] Zhao, D.Y., Kleijn, B. W., 2007. HMM-based gain modeling for enhancement of speech in noise. IEEE Trans. Audio, Speech, and Language Process., 15 (3),882–892.

[66] Sriram, S., Jonas, S., Kleijn, W. B., 2007. Codebook-based Bayesian speech enhancement for nonstationary environments. IEEE Trans. Audio, Speech, and Language Process., 15 (2), 441–452.

[67] Veisi, H., Hossein, S., 2013. Speech enhancement using hidden Markov models in Mel-frequency domain. Speech Commun., 55 (2), 205–220.

[68] Nishikawa, T., Saruwatari, H., Shikano.K., 2003. Blind Source Separation of Acoustic Signals Based on Multistage ICA Combining Frequency-Domain ICA and Time-Domain ICA. IEICE Trans. Fundamentals of Electron., Commun. and Comput. Sci., 86 (4), 846–858.

[69] So, S., Paliwal, K. K., 2011. Modulation-domain Kalman filtering for single-channel speech enhancement. Speech Commun., 53 (6), 818–829.

[70] Paliwal, K. K., Schwerin, B., Wójcicki, K., 2012. Speech enhancement using a minimum mean-square error short-time spectral modulation magnitude estimator Speech Commun., 54 (2), 282–305.

[71] Ji, M., Srinivasan, R., Crooks, D., 2011. A corpus-based approach to speech enhancement from nonstationary noise. IEEE Trans. Audio, Speech, and Language Process., 19 (4), 822–836.

[72] Ruofei, C., Cheung-Fat, C., Cheung, H.S., 2012. Model-based speech enhancement with improved spectral envelope estimation via dynamics tracking. IEEE Trans. Audio, Speech, and Language Process., 20 (4), 1324–1336.

[73] Mohammadiha, N., Smaragdis, P., Leijon, A., 2013. Supervised and unsupervised speech enhancement using nonnegative matrix factorization. IEEE Trans. Audio, Speech, and Language Process., 21 (10), 2140–2151.

[74] Sawada, H., Kameoka, H., Araki, S., Ueda, N., 2013. Multichannel Extensions of Non-Negative Matrix Factorization With Complex-Valued Data. IEEE Trans. on Audio, Speech and Language Process., 21 (5), 971–982.

[75] Wang, D., Lim, J. S., 1982. The unimportance of phase in speech enhancement IEEE Trans. Acoust. Speech Signal Process., ASSP-30 (4), 679–681.

[76] Shannon, B. J., Paliwal, K.K, 2006. Role of Phase Estimation in Speech Enhancement. In Proc. IEEE SAPA@ INTERSPEECH, 1427–1430.

[77] Paliwal, K. K., Alsteris, L.D., 2005. On the usefulness of STFT phase spectrum in human listening tests. Speech Commun., 45 (2), 153–170.

[78] Paliwal, K. K., Wójcicki, K., Shannon, B., 2011. The importance of phase in speech enhancement. Speech commun., 53 (4), 465–494.

[79] Roux, J. L., Ono, N., Sagayama, S., 2008. Explicit consistency constraints for STFT spectrograms and their application to phase reconstruction. In Proc. IEEE SAPA@ INTERSPEECH, 23–28.

[80] Loizou, P. C., Kim, G., 2011. Reasons why current speech-enhancement algorithms do not improve speech intelligibility and suggested solutions. IEEE Trans. Audio, Speech, and Language Process., 19 (1), 47–56.

[81] Drullman, R., 1995. Temporal envelope and fine structure cues for speech intelligibility. J. Acoust. Soc. Am., .97 (1), 585–592.

[82] Moore, B.C., 2008. The role of temporal fine structure processing in pitch perception, masking, and speech perception for normal-hearing and hearing-impaired people. J. Assoc. Research Otolaryngology, 9 (4), 399–406.

[83] Swaminathan, J., 2010. The role of envelope and temporal fine structure in the perception of noise degraded speech. Ph.D Thesis, Purdue University.

[84] Swaminathan, J., Heinz, M.G., 2012. Psychophysiological Analyses Demonstrate the Importance of Neural Envelope Coding for Speech Perception in Noise. J. Neuroscience, 32 (5), 1747–1756.

[85] Unoki, M., Masato, M., 1999. A method of signal extraction from noisy signal based on auditory scene analysis. Speech Commun., 27 (3), 261–279.

[86] Nower, N., Liu, Y., Unoki, M., 2014. Restoration of instantaneous amplitude and phase using Kalman filter for speech enhancement. In Proc. IEEE Int. Conf. Acoustics, Speech, and Signal Process. (ICASSP), 4666–4670.

[87] Paliwal, K. K., Atal, B.S., 1993. Efficient vector quantization of LPC parameters at 24 bits/frame. IEEE Trans. Audio, Speech, and Language Process., 1 (1), 3–14.

[88] Hu, Y., Loizou, P.C., 2008. Evaluation of objective quality measures for speech enhancement. IEEE Trans. Audio, Speech, and Language Process., 16 (1), 229–238.

[89] Ma, J., Loizou, P.C., 2011. SNR loss: A new objective measure for predicting the intelligibility of noise-suppressed speech. Speech Commun., 53 (3), 340–354.

# Publications

## International Journal

[1] N. Nower, Y. Tan and A. O. Lim: "Traffic Pattern based Data recovery Scheme for Cyber-physical systems", IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences, vol.E97-A, no.9, pp.1926–1936, Sep. (2014).

[2] N. Nower, Y. Liu and M. Unoki: "Restoration Scheme of Instantaneous Amplitude and Phase using Kalman Filter with Efficient Linear Prediction for Speech Enhancement," Speech Communication, Elsevier 2015. (In press)

## International Conferences

[3] N. Nower, Y. Tan and A. O. Lim, "Efficient Spatial Data Recovery Scheme for Cyber-physical systems", The 1st IEEE International Conference on Cyber-Physical Systems, Networks, and Application, pp.92–97, Taipei, Taiwan, August (2013).

[4] N. Nower, Y. Tan and A. O. Lim, "Efficient Temporal and Spatial Data Recovery Scheme for Stochastic and Incomplete Feedback Data of Cyber-physical systems", The 8th IEEE International Symposium on Service Oriented System Engineering, pp.192–197,Oxford, United Kingdom, April (2014).

[5] N. Nower, Y. Liu and M. Unoki, "Restoration of Instantaneous Amplitude and Phase using Kalman filter for Speech Enhancement", The 39th IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp.4633–4637, Florence, Italy, May (2014).

[6] Y. Liu, N. Nower and M. Unoki, "Restoration of Instantaneous Amplitude and Phase for Speech Signal using Kalman filter in Noisy Reverberant Environment"(To be submitted)

[7] N. Nower, Y. Tan and A. O. Lim, "Stochastic and Incomplete Feedback Data Recovery Scheme with Kalman Filter for High-Confidence Cyber-Physical Systems"(To be submitted)

## Domestic Conferences

[8] N. Nower, Y. Tan and A. O. Lim, "Data Recovery for Cyber-physical System with Incomplete Feedback Sensor Data", IEICE Society Conference, Toyama, Japan, September (2012).

[9] N. Nower, Y. Tan and A. O. Lim, "Efficient Spatial Data Recovery Scheme for Different Traffic Patterns of Cyber-physical Systems", IEICE Society Conference, Fukuoka, Japan, September (2013).

[10] N. Nower, Y. Tan and A. O. Lim, "Data Recovery Scheme for Stochastic and Incomplete Feedback Data of Cyber-physical Systems", IEICE General Conference, Niigata, Japan, March (2014).

[11] N. Nower, Y. Liu and M. Unoki, "Study on Restoration of Instantaneous Amplitude and Phase using Kalman filter for Speech Enhancement", Acoustical Society of Japan (ASJ) Spring Meeting, Tokyo, March (2014).

[12] Y. Liu, N. Nower and M. Unoki, "Instantaneous Amplitude and Phase Restoration using Kalman filter for Speech Enhancement", IEICE Technical Report on Signal Processing (IEICE-SP), vol.114, no.91, pp.2732, Kanazawa, June (2014).