Partial Knowledge Data-driven Event Detection for Power Distribution Networks

Yuxun Zhou¹, Reza Arghandeh², Member, IEEE, Costas J. Spanos³ Fellow, IEEE

Abstract—The power system has been incorporating increasing amount of unconventional generations and loads such as distributed renewable resources, electric vehicles, and controllable loads. The induced dynamic and stochastic power flow require high-resolution monitoring technology and agile decision support techniques for system diagnosis and control. This paper discusses the application of micro-phasor measurement unit (μ PMU) data for power distribution network event detection. A novel datadriven event detection method, namely Hidden Structure Semi-Supervised Machine (HS³M), is established. HS³M only requires partial expert knowledge: it combines unlabeled data and partly labeled data in a large margin learning objective to bridge the gap between supervised learning, semi-supervised learning, and learning with hidden structures. To optimize the non-convex learning objective, a novel global optimization algorithm, namely Parametric Dual Optimization Procedure (PDOP), is established through its equivalence to a concave programming. Finally, the proposed method is validated on an actual distribution feeder with installed μ PMUs, and the result justifies the effectiveness of the learning-based event detection framework, as well as its potential to serve as one of the core algorithms for power system security and reliability.

Index Terms—Event Detection, Phasor Measurement Unit, Distribution, Machine Learning.

Nomenclature

 α, β, γ Lagrangian multipliers for the three types of constraints in the primal; Decision variables in the Dual

 η Indicator variables for partly labeled data

 λ Unified hidden variable in the Dual

 θ Unified decision variable in the Dual

Indicator variables for unlabeled data

 $C^{\theta}, C^{\lambda}, C^{0}$ Matrices encoding the coefficients of θ, λ and the constant term in the constraints.

e Column vector of all 1s

Q Kernel matrix in the Dual

T,H,P Composite matrices defined in Theorem 1

w Parameter of the classifier

 x_i The i^{th} data sample (system measurement)

 $\kappa(\cdot,\cdot)$ Kernel function

 Λ,Θ Feasible set for λ and θ , respectively

 \mathcal{A} Index set of active constraints

 \mathcal{H} Hilbert function space

 $\mathcal{J}(\boldsymbol{\theta})$ Value of the dual objective as a function of $\boldsymbol{\theta}$

¹Yuxun Zhou is with Electrical Engineering Computer Sciences, University of California, Berkeley, 94720 CA, USA. yxzhou@berkeley.edu

²Reza Arghandeh is an Assistant Professor in the Electrical and Computer Engineering Department, Florida State University, FL 32306, USA. arghandehr@gmail.com

³Costas J. Spanos is the Andrew S. Grove Distinguished Professor of Electrical Engineering Computer Sciences, University of California, Berkeley, 94720 CA, USA. spanos@berkeley.edu

 $\begin{array}{ll} \mathcal{L}^{+} & \text{Index set of all samples having } y_i = +1 \\ \mathcal{L}^{-} & \text{Index set of completely labeled samples having } y_i = \\ & -1 \\ \mathcal{L}^{+}_{H} & \text{Index set of partly labeled samples having } y_i = +1 \\ \mathcal{L}^{-}_{H} & \text{Index set of partly labeled samples having } y_i = -1 \\ \mathcal{U} & \text{Index of unlabeled data samples} \\ \widehat{y}_i & \text{Tentative label for sample } i \in \mathcal{U} \\ \widehat{y} & \text{Augmented label vector, including tentative and dummy labels} \\ c \ , c_1 - c_3 & \text{Classification loss penalties} \\ K & \text{Number of event classes} \\ k & \text{Index for the } k^{th} \text{ base classifier} \\ \end{array}$

I. INTRODUCTION

Indicator for nominal (0) or event (1)

Event type indicator

A. Background

HISTORICALLY, power distribution networks are behind transmission networks regarding transparency, measurement frequency, and sophisticated monitoring systems. However, the growth of distributed renewable energy resources, electric vehicles, power electronic based grid edge controllers, and controllable loads introduces more short-term and unpredictable disturbances in the electricity flow [1]. It suggests a need for more accurate measurement devices with higher resolution. As utilities have expanded their monitoring systems and sensor networks, the data management, analysis, and inference have become the most significant challenges with the shift from a conventional data acquisition system depended on domain experts (supervised) to an automated intelligent analysis system (unsupervised). This paper specifically discusses high-precision synchrophasors, or micro-phasor measurement units (µPMUs) for high-fidelity measurement of voltage and current waveforms [2]. All μ -PMUs measurements are GPS time stamped to provide timesynchronized observability for three-phase voltage and current magnitude and phase angle with a 0.05% Total Vector Error [3]. Figure 1 shows the actual μ -PMU device and the overall system architecture. The μ -Pnet is the equivalent of the Phasor Data Concentrator (PDC) in a distribution network. More technical information about the μ -PMU device is available in [4]. Topology detection[5], phase labeling [6] and linear state estimation[7] are among applications of μ -PMU data that are explored so far. Moreover, the accuracy and resolution available from such μ -PMU monitoring network enable operators to detect short-time events that would otherwise be unobservable in distribution systems.



Fig. 1: The integration of the proposed event detection tool with existing μ -PMU monitoring and SCADA system.

Events of interest in distribution networks are classified into different categories based on their duration, cause, and location. Two sample events, captured by a μ -PMU, are shown in Figure 2. Events in power systems present themselves as oscillations, transients, short duration variations, long-duration variations, and waveform distortion in voltage, current and frequency values [8]. The causes are usually faults, switching actions, topology changes, controllers, load behavior and source dynamics. Related standards for characterizing events in distribution network including but are not limited to [9], [10], [11]. Moreover, the initial fault or event occurring at one or more components of the grid usually triggers the faults or events of others [12], [13], [14], [15]. For the sake of power systems reliability, stability, security, and resilience, it is crucial to monitor the operating states in real time and detect anomalies quickly to avert disturbances and disruptions[16]. Besides μ -PMU based event detection system in distribution networks, the proposed data-driven tool in this paper is compatible with different types of measurement such as traditional measurement devices, power quality recorders, protection relays, smart meters, and synchrophasors. Within a larger scope, this paper aims to leverage advances in machine learning techniques, global optimization methods and timeseries data analysis for data-driven event detection in power distribution networks.

B. Learning Based Event Detection Methods

A novel framework for event detection is proposed in this paper using μ -PMU data streams. A large body of available literature has been addressing the problem of event detection, in particular, its simplified case of novelty/fault detection with various methods. Available works on event detection can be

divided into two main categories and their combinations, i.e., model-based method and model-less or data-driven method.

The basic idea of model-based approaches is to compare the system behavior, estimated by a dynamic model, to the expected behavior when the system is in a certain state [17]. To list a few, model-based applications range from cyber attacks identification in power system [18], fault diagnosis for switching converters [19], and fault-detection of electric machine [20]. Recently, model-based approaches have also been used in combination with time series analysis to establish semi-model-based algorithms [21]. This type of methods is largely based on the correctness of the dynamic model of the system, as well as system analytic tools such as realtime state estimators and parity equations. Their limitations are evident as dynamics of a system may be (1) hard to specify and (2) nonlinear/coupled in structure. Also, as more and more applications are dealing with a complex system with randomness, the high dimensionality and inherent uncertainty significantly deteriorate the reliability of dynamic models.

On the other hand, the data-driven approaches use methods of machine learning to conduct statistical inference or decision making on available system measurements. As massive amount of data is provided by the advancement of sensor network and information technology, this approach is receiving increasing attention in both application and research domains. The problem of distinguishing abnormalities from the normal states has been extensively studied in literature and is referred to as novelty (or abnormal) detection. Indeed, many classic machine learning tools, such as kernel Principle Component Analysis (kPCA), Partial Least Squares (PLS), and Fisher Discriminant Analysis (FDA) have been widely applied in various fields. Readers are referred to [22] and the references therein for a comprehensive survey.

The event detection problem becomes more challenging as the objective is extended from making a binary decision to distinguishing the type of events. Recently several attempts were made for this purpose, including cost sensitive Support Vector Machine (SVM) for fault type identification in semiconductor manufacturing process [23] and the hierarchical kernel method for building cooling system fault diagnosis [24]. These methods belong to the supervised multi-class classification which require properly labeled data with detailed event types as the training set. However, the occurrence of events might be rare in many real-world applications, and the availability of fully labeled data set is often limited. In the situation dealing with μ -PMU data, complete expert knowledge for power system diagnosis and event labeling can be expensive and scant. Lacking labeled data may cause insufficient supervised learning and eventually lead to degraded performance.

In this work, a novel data-driven event detection method, namely the Hidden Structure Semi-Supervised Machine (HS³M), is proposed for the purpose of event type identification. To resolve the scarcity of labeled training data, HS³M also incorporates information from *partly labeled* data and *unlabeled* data. The inclusion of this partial information breaks the convexity of traditional large margin learning formulation. As such a new global optimization algorithm is developed to solve the non-convex learning problem. The proposed method

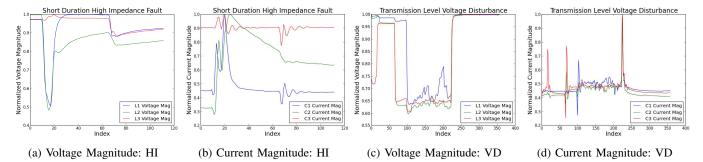


Fig. 2: Real time μ -PMU measurements that contains events. (a)&(b): Voltage/Current measurements when a short duration high impedance fault occurs. (c)&(d)Voltage/Current measurements when a transmission level voltage disturbance occurs.

is applied to a concrete case of power distribution network event detection with μ -PMU data. The integration of the event detection tool as a complement of the existing monitoring and control system is shown in Figure 1.

The contributions are highlighted as follows:

- HS³M provides a unified framework to incorporate data
 of different natures for event detection. The proposed
 method bridges the gap between semi-supervised learning
 and learning with hidden structures.
- The newly developed PDOP algorithm constitutes a promising substitute for previous methods such as Concave-Convex Procedure and alternating optimization.
- HS³M yields significant improvement compared to stateof-the-art methods, showing promising results on distribution network event detection with μ -PMU data.

The paper is organized as follows, in section II we introduce the formulation of HS³M learning objective and the intuition behind it. Section III is devoted to transforming the training problem into a parametric quadratic program. In Section IV, the PDOP algorithm is established for HS³M learning. The case study is given in section V.

II. HS³M Event detection: Motivation and Formulation

Traditionally, there have been two fundamentally different paradigms of machine learning (ML). The first one is the supervised learning, with the goal of learning a mapping from some input \boldsymbol{x} to output y. Usually the observations $(\boldsymbol{x}_i,y_i), i=1,\cdots,n$ are called samples, \boldsymbol{x}_i are referred to as features of sample i, and $y_i \in \mathcal{Y}$ are called labels or targets. To find the "optimal" mapping f, the learning task is commonly formulated as "regularized empirical risk minimization" [25]

$$\min_{f \in \mathcal{H}} \frac{1}{2} ||f||_{\mathcal{H}}^2 + c \sum_i L(y_i, f(\boldsymbol{x}_i))$$
 (1)

in which the second term measures the "goodness of fit" of the classifier f with some loss function L, and the first regularization term controls the complexity of the mapping f to avoid over-fitting¹[26]. By tuning the "hyper-parameter" c (called "model selection" in the jargon of ML), one is able to balance training fitness and model complexity, hence finding the optimal classifier that generalizes well to unseen data set.

The second task of ML is unsupervised learning [27]. Under this setting, only the unlabeled observations $X = \{x_1, \dots, x_n\}$ are given. Typically, the goal of unsupervised learning is to identify interesting structures in the data X, such as clusters, quantiles, support, low-dimensional embedding, or more generally the patterns related with the distribution of the data.

The presence of both labeled and unlabeled data motivates the so-called semi-supervised learning [28]. The hope is that, by combining both types of available data sets, semisupervised learning could find better models/classifiers, and reduce the cost of expert engagement [29]. In the context of event detection using μ -PMU measurement, data with detailed event labels is precious but scant - power system experts are needed to inspect the measurement and mark events. On the other hand partly labeled data with incomplete labels, e.g., nominal/fault, may be obtained less costly by using decision support systems. Moreover, unlabeled data can be acquired in large quantity simply by collecting μ -PMU measurement. Given the characteristics of the different types of data available for power system event detection, a unified ML framework is proposed to leverage all information sources. The following subsection explains the construction of the proposed HS³M step-by-step.

A. Motivation and Intuition of HS³M

To formalize the information availability in different scenarios, consider data of the following three formats:

- 1 Completely labeled data samples, denoted as $\{x_i, y_i, z_i\}$, where i is the sample index, x_i is the system measurement, y_i is a indicator for "nominal/stable state" (y=+1) or an "event" (y=-1). If y=-1, an event type indicator $z_i \in \{1, \cdots, K\}$ is also associated.
- 2 Partly labeled data samples, denoted as $\{x_i, y_i, \cdot\}$, where y_i is still the indicator for "event", but event type is unavailable as differentiating event types is costly.
- 3 Unlabeled data samples, denoted as $\{x_i, \cdot, \cdot\}$, where only system measurement x_i is accessible.

An illustration of different situations is given in Fig. 3. Intuitively, partly labeled data should be helpful: at least it provides discriminating information between stable state and events. The role of unlabeled data might be ambiguous since it does not carry any expert knowledge. However, it does contain distributional information of measurement, which

¹Known as the Occam's razor principle. In some context the regularization can also help alleviate ill-posed problem and induce sparsity.

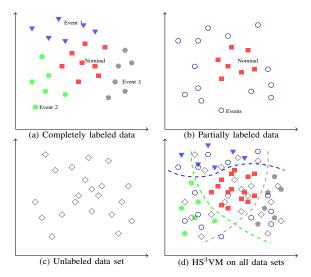


Fig. 3: Different data format and the intuition of HS³M

could be exploited with a proper formulation. In the sequel, an unified learning framework is designed to combine all three information sources to improve detection performance.

To begin with, the regularized empirical risk minimization framework is substantiated by hing loss and L_2 regularization:

$$\min_{\mathbf{w}} \frac{1}{2} ||\mathbf{w}||_{2}^{2} + c \sum_{i} [1 - y_{i}g(\mathbf{w}, \mathbf{x}_{i})]_{+}$$
 (2)

where c is the loss penalties, $[1-y_ig(\boldsymbol{w},\boldsymbol{x}_i)]_+$ is the hinge loss for the miss classification, and \boldsymbol{w} is the parameters of the classifier $g(\boldsymbol{w},\boldsymbol{x}_i)$. The regularization term $||\boldsymbol{w}||_{\mathcal{H}}^2$ also favors the large margin separating hyperplanes. The HS³M formulation resembles this "regularization+hinge loss" minimization, but both partly labeled data and unlabeled data are incorporated as additional information, by constructing a new classifier having a geometric property suitable for power system event detection purposes.

B. A Geometrically Constrained Classifier

Intuitively, data generated from the nominal (stable) state are often "concentrated" in some distribution, whereas data of diverse events are "scattered around" since they deviated from the nominal in different ways. This distributional asymmetric is illustrated in the top-left panel of Figure 3. To encode this particular "geometric" information of event detection, it is natural to describe the stable state by the intersection of the acceptance region of multiple nonlinear decision rules, while the events class by the union of their complements, as is shown in Figure 3(d). Mathematically, a composed classifier $g(\cdot)$ can be written in terms of multiple base classifiers $f_k(\cdot)$ as follows:

$$g(\boldsymbol{w}, \boldsymbol{x}) > 0 \Leftrightarrow \min\{f_1(\boldsymbol{w}, \boldsymbol{x}), \cdots, f_K(\boldsymbol{w}, \boldsymbol{x})\} > 0$$

where K is the number of all possible event types. The construction of the classifier inherently emphasizes the sensitivity to event class because an event is detected by any one of the base classifiers. The classifier also maintains the specificity to stable state class, as all base classifiers have to "agree" for a positive prediction. Given a feature mapping $\phi: \mathbb{R}^d \to \mathcal{H}$, a

hyperplane classifier in the new Hilbert space can be written as $f(\boldsymbol{w}, \boldsymbol{x}) = \langle \boldsymbol{w}, \phi(\boldsymbol{x}) \rangle_{\mathcal{H}} + b \triangleq \boldsymbol{w} \cdot \phi(\boldsymbol{x})$ for short hand notation. Then the proposed classifier is

$$g(\boldsymbol{w}, \boldsymbol{x}) = \min_{k} \{ \boldsymbol{w} \cdot \phi_{k}(\boldsymbol{x}) \}. \tag{3}$$

- 1) Incorporating Labeled Data: Based on the composed classifier, the hinge loss $[1-y_ig(\boldsymbol{w},\boldsymbol{x}_i)]_+$ is computed for each format of data to incorporate their "contribution" to the classifier. With detailed event type being available, the hinge loss for labeled data set can be directly obtained as $[1+\boldsymbol{w}\cdot\phi_{z_i}(\boldsymbol{x}_i)]_+$ when $y_i=-1$ and $[1-\min_k\{\boldsymbol{w}\cdot\phi_k(\boldsymbol{x}_i)\}]_+$ when $y_i=+1$.
- 2) Incorporating Partly Labeled Data: Comparing Figure 3(a) and 3(b), we see that partly labeled data could be viewed as data with "missing detailed labels" (event types). But the hidden clusters of each event class still exist. From a ML perspective, this is the problem of "learning with hidden structures". Composed from multiple base classifiers, the unified classifier g performs implicit "clustering" for the event class, so as to capture hidden subgroups in partly labeled data sets. Thus the hinge loss for a partly labeled data sample $\{x_i, y_i, \cdot\}$ is just $[1 y_i \min_k \{w \cdot \phi_k(x_i)\}]_+$ for $y_i = \pm 1$.

 3) Incorporating Unlabeled Data: Given a classifier
- 3) Incorporating Unlabeled Data: Given a classifer $g(\boldsymbol{w}, \boldsymbol{x})$, the predicted label is its sign, i.e., $\widehat{y}_i = \operatorname{sign}(g(\boldsymbol{w}, \boldsymbol{x}))$. This inspires a tentative labeling strategy to include information provided by unlabeled data. More specifically, with $\widehat{y}_i = \operatorname{sign}(\min_k \{\boldsymbol{w} \cdot \phi_k(\boldsymbol{x}_i)\})$, the corresponding hinge loss has the form

$$\left[1 - \widehat{y}_i \min_{k} \{\boldsymbol{w} \cdot \phi_k(\boldsymbol{x}_i)\}\right]_{+} = \left[1 - \left|\min_{k} \{\boldsymbol{w} \cdot \phi_k(\boldsymbol{x}_i)\}\right|\right]_{+}$$

C. Overall Learning Objective

Putting things together, the following regularized hinge loss minimization is formulated for event detection that incorporates all explicit and partial expert knowledge:

$$\min_{\boldsymbol{w}} \quad \frac{1}{2} ||\boldsymbol{w}||_{\mathcal{H}}^{2} + c_{1} \sum_{i \in \mathcal{L}^{+}} \left[1 - \min_{k} \{\boldsymbol{w} \cdot \phi_{k}(\boldsymbol{x}_{i})\} \right]_{+} \\
+ c_{21} \sum_{i \in \mathcal{L}^{-}} \left[1 + \boldsymbol{w} \cdot \phi_{z_{i}}(\boldsymbol{x}_{i}) \right]_{+} \\
+ c_{22} \sum_{i \in \mathcal{L}_{H}^{-}} \left[1 + \min_{k} \{\boldsymbol{w} \cdot \phi_{k}(\boldsymbol{x}_{i})\} \right]_{+} \\
+ c_{3} \sum_{i \in \mathcal{U}} \left[1 - \left| \min_{k} \{\boldsymbol{w} \cdot \phi_{k}(\boldsymbol{x}_{i})\} \right| \right]_{+}$$
(ORIG)

where \mathcal{L}^+ denotes the index set of all data samples that has $y_i = +1$, including both completely and partly labeled samples, \mathcal{L}^- as the index set of completely labeled samples with $y_i = -1$ and event type z_i (hence the hinge loss only involves the corresponding individual classifier f_{z_i}). The index set \mathcal{L}_H^- contains partly labeled samples that are "events", and \mathcal{U} is the index of all unlabeled data samples. The loss penalty hyper-parameters c_1 - c_3 weight each loss term differently, and should be chosen by taking into account the imbalanced cost for false positive and false negative error, sample size in each category, as well as for model selection considerations. Since the above formulation deals with both hidden structures and unlabeled samples in the available data, we call it Hidden Structured Semi-Supervised Machine (HS³M).

III. LEARNING OBJECTIVE REFORMULATION

The first three terms in the learning objective **ORIG** are convex in decision variables, the last two terms, however are not convex. Existing heuristics for this type of problems use either group alternating optimization or Concave-Convex Procedure (CCCP), which only converge to local minimals and may lead to deteriorated solution. The degradation especially happens as we are in a machine learning context [30].

In this work, a novel optimization algorithm is derived for HS³M as well as for a class of "hidden structure" problems. A remarkable property of the algorithm is that it can approach global optimum by iteratively improving local solutions. The new method is based on an equivalent optimization problem and two important parametric property of its dual. As the first step, we transform **ORIG** by introducing additional "hidden" variables, and write the learning objective in the following joint optimization form:

Proposition 1. ORIG is equivalent to (OPT1)

$$\begin{split} \min_{\boldsymbol{\eta}, \boldsymbol{\zeta}} \min_{\boldsymbol{w}} & & \frac{1}{2} ||\boldsymbol{w}||_{\mathcal{H}}^2 + c_1 \sum_{i \in \mathcal{L}^+} \left[1 - \min_{\boldsymbol{k}} \{ \boldsymbol{w} \cdot \phi_{\boldsymbol{k}}(\boldsymbol{x_i}) \} \right]_+ \\ & & + c_{21} \sum_{i \in \mathcal{L}^-} \left[1 + \boldsymbol{w} \cdot \phi_{z_i}(\boldsymbol{x_i}) \right]_+ \\ & & + c_{22} \sum_{i \in \mathcal{L}^-_H} \sum_{k=1}^K \eta_{ik} \left[1 + \boldsymbol{w} \cdot \phi_{k}(\boldsymbol{x_i}) \right]_+ \\ & & + c_3 \sum_{i \in \mathcal{U}} \sum_{k=1}^K \zeta_{ik} \left[1 + \boldsymbol{w} \cdot \phi_{k}(\boldsymbol{x_i}) \right]_+ \\ & & + c_3 \sum_{i \in \mathcal{U}} \zeta_{i(K+1)} \max_{j} \left\{ 0, 1 - \boldsymbol{w} \cdot \phi_{j}(\boldsymbol{x_i}) \right\} \end{split}$$

subject to $\boldsymbol{\eta}_i \in \mathbb{S}^K$, $\forall i \in \mathcal{L}_H^-$; $\boldsymbol{\zeta}_i \in \mathbb{S}^{K+1}$, $\forall i \in \mathcal{U}$

In addition, the two minimization is interchangeable.

The introduced variables η and ζ can be thought of as hidden state indicators for partially labeled data and unlabeled data, respectively. The corresponding dual of the **inner optimization** of OPT1 is

$$\max_{\boldsymbol{\alpha},\boldsymbol{\beta},\boldsymbol{\gamma}} - \frac{1}{2} \left\| \sum_{k,i \in \mathcal{I}} \alpha_{ik} y_i \phi_k(\boldsymbol{x}_i) + \sum_{i \in \mathcal{L}^-} \beta_i \phi_{z_i}(\boldsymbol{x}_i) + \sum_{k,i \in \mathcal{U}^+} \gamma_{ik} \phi_k(\boldsymbol{x}_i) \right\|_{\mathcal{H}}^2$$

$$+ \sum_{k,i \in \mathcal{I}} \alpha_{ik} + \sum_{i \in \mathcal{L}^-} \beta_i + \sum_{k,i \in \mathcal{U}^+} \gamma_{ik}$$
subject to
$$\begin{cases} \alpha_{ik} \geq 0; & \sum_{k} \alpha_{ik} \leq c_1 \quad \forall i \in \mathcal{L}^+ \\ 0 \leq \beta_i \leq c_{21} \quad \forall i \in \mathcal{L}^- \\ 0 \leq \alpha_{ik} \leq c_{22} \eta_{ik} \quad \forall i \in \mathcal{L}^- \\ 0 \leq \alpha_{ik} \leq c_{32} \eta_{ik} \quad \forall i \in \mathcal{U}^- \\ \gamma_{ik} \geq 0; & \sum_{k} \gamma_{ik} \leq c_{3} \zeta_{i(K+1)} \quad \forall i \in \mathcal{U}^+ \\ \sum_{k,i \in \mathcal{I}} y_i \alpha_{ik} + \sum_{k,i \in \mathcal{L}^-} y_i \beta_i + \sum_{k,i \in \mathcal{U}^+} y_i \gamma_{ik} = 0 \end{cases}$$

where the Lagrangian multipliers α, β, γ are now decision variables. Note that the unlabeled data set \mathcal{U} is used as two dummy copies with tentative labels $y_i = +1$ for $i \in \mathcal{U}^+$ and $y_i = -1$ for $i \in \mathcal{U}^-$, respectively. Also for short hand notation, we denote the index set of all samples by \mathcal{I} , and a unified decision variable in the Dual by

$$oldsymbol{ heta} = \left[oldsymbol{lpha}^T, oldsymbol{eta}^T, oldsymbol{\gamma}^T
ight]^T$$

where $\alpha \triangleq \left[\alpha_{11}, \cdots, \alpha_{|\mathcal{I}|1}, \cdots, \alpha_{|\mathcal{I}|K}\right]^T$, $\beta \triangleq \left[\beta_1, \cdots, \beta_{|\mathcal{L}^-|}\right]^T$ and $\gamma \triangleq \left[\gamma_{11}, \cdots, \gamma_{|\mathcal{I}|1}, \gamma_{12} \cdots, \gamma_{|\mathcal{I}|K}\right]^T$. It is immediate that the norm in the Hilbert space reduces to inner products. Hence the objective of the dual (pulling out a minus sign) can be equivalently written as

$$-\min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{i,j} \theta_i \langle \phi(\boldsymbol{x}_i, y_i, i), \phi(\boldsymbol{x}_j, y_j, j) \rangle \theta_j - \sum_i \theta_i$$
 (4)

in which the so called kernel trick could be used for direct computation of the inner product without the need to compute the explicit feature mapping $\phi(\cdot)$, i.e.,

$$\langle \phi(\boldsymbol{x}, y, i), \phi(\boldsymbol{x}', y', j) \rangle = \kappa_{(iy)(jy')}(\boldsymbol{x}, \boldsymbol{x}') \tag{5}$$

For a more compact form of the dual, let us further define a matrix Q with elements $Q_{(iy)(jy')} = \kappa_{(iy)(jy')}(x,x')$, a column vector of all hidden variables $\boldsymbol{\lambda} \triangleq \{\boldsymbol{\eta}_1; \cdots; \boldsymbol{\eta}_{|\mathcal{L}_H^-|}, \boldsymbol{\zeta}_1, \cdots, \boldsymbol{\zeta}_{|\mathcal{U}|}\}$, and an augmented label vector (with dummy copies of unlabeled set) as

$$\widetilde{\boldsymbol{y}} = [\underbrace{1,\cdots,1}_{\mathcal{L}^+},\underbrace{-1,\cdots,-1}_{\mathcal{L}^-_H},\underbrace{-1,\cdots,-1}_{\mathcal{U}^-},\underbrace{-1,\cdots,-1}_{\mathcal{L}^-},\underbrace{1,\cdots,1}_{\mathcal{U}^+}]^T$$

together with a matrix encapsulated inequality constraints, the (negative) inner optimization becomes

$$\min_{\boldsymbol{\theta}} \mathcal{J}(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^T \boldsymbol{Q} \boldsymbol{\theta} - \boldsymbol{e}^T \boldsymbol{\theta}$$
subject to
$$\begin{cases}
\boldsymbol{C}^{\boldsymbol{\theta}} \boldsymbol{\theta} \leq \boldsymbol{C}^{\lambda} \boldsymbol{\lambda} + \boldsymbol{C}^0 \\
\widetilde{\boldsymbol{y}}^T \boldsymbol{\theta} = 0.
\end{cases}$$
(OPT_D)

where C^{θ} , C^{λ} , C^{0} are constant matrices with $K|\mathcal{L}_{H}^{-}| + (K+1)|\mathcal{U}|$ rows. Similar to other kernel methods in machine learning, HS³M is restricted to Mercer kernels, thence Q is positive definite, and $(\mathbf{OPT}_{-}\mathbf{D})$ is in a convex quadratic program. The learning objective **ORIG** proposed in section II now becomes

$$\max_{\boldsymbol{\lambda} \in \Lambda} \min_{\boldsymbol{\theta} \in \Theta(\boldsymbol{\lambda})} \mathcal{J}(\boldsymbol{\theta}) \tag{6}$$

Next the inner optimization of (6) is analyzed in some depth, and establish the theoretical foundations of the novel global optimization algorithm.

IV. STRUCTURES OF THE DUAL OPTIMALITY AND A NOVEL GLOBAL OPTIMIZATION METHOD

First of all, it is helpful to consider the inner optimization (**OPT_D**) as if "parameterized" by the outer optimization variable λ . Thus solving it with some fixed λ , one can get the optimal solution as a function of the "parameters". Denote the optima as $\theta^*(\lambda)$ to emphasize this dependence, then the original learning objective becomes $\min_{\lambda \in \Lambda} \mathcal{J}(\theta^*(\lambda))$.

Two important properties of $\mathcal{J}(\boldsymbol{\theta}^*(\boldsymbol{\lambda}))$ will be characterized: (1) Locally (in a well defined neighborhood called critical regions) $\boldsymbol{\theta}^*(\boldsymbol{\lambda})$ has an explicit form. (2) Globally the inner optimal objective $\mathcal{J}(\boldsymbol{\theta}^*(\boldsymbol{\lambda}))$ is convex piece-wise quadratic in $\boldsymbol{\lambda}$. These two observations serve as the underpinning for the new optimization algorithm.

The problem of analyzing the dependence between optimal solution and involved parameters have been previously addressed in operational research and our communities, with

the terminology Parametric Programming (PP) or Sensitivity Analysis (SA). In particular, the study of Parametric Quadratic Programming (PQP) can be dated back to [31], and following the pioneer work [32], it has been widely applied for model predictive control problems [33][34]. However, to the best of our knowledge, PQP has not yet been explored for machine learning problems. More importantly, technical difficulties arise because existing results usually rely on the so called Linear Independence Constraint Qualification (LICQ) assumption, which in the current case of (OPT_D) cannot be satisfied due to the existence of the equality constraint, $\tilde{\boldsymbol{y}}^T\boldsymbol{\theta}=0$. In fact, in the jargon of PP or SA, the problem at hand corresponds to a degenerate case for which existing solution is still lacking.

In subsequent parts, we will bridge the gap by firstly replacing the LICQ condition with a sample partition property, and then show that the local explicit form of $\theta^*(\lambda)$ can still be obtained under a mild condition. But before that, it is useful to define the following terms.

Definition 1. (Active Constraint) Assume that an optimal solution of (OPT_D) has been obtained as $\theta^*(\lambda)$. Then the i^{th} row of the constraints is said to be active at λ , if $C_i^{\theta}\theta^*(\lambda) = C_i^{\lambda}\lambda + C_i^{0}$, and inactive if $C_i^{\theta}\theta^*(\lambda) < C_i^{\lambda}\lambda + C_i^{0}$. Let index set of all active inequality constraints i be denoted by A, and all inactive inequality constraints by A^C . C_A denotes the C matrix with only rows that correspond to the active constraints, and C_{A^C} as the matrix with only rows that correspond to inactive constraints.

Definition 2. Samples Partition: Based on the value of θ^* , i.e., $\alpha^*, \beta^*, \gamma^*$ at optimal,

- A j^{th} sample in $\mathcal{L}^+ \cup \mathcal{L}^- \cup \mathcal{U}^+$ is called: Support vector if $\alpha_{jk} > 0$, $\beta_j > 0$ or $\gamma_{jk} > 0$ respectively. Unbounded support vector if in addition $\sum_k \alpha_{jk} < c_1$, $\beta_j < c_{21}$, or $\sum_k \gamma_{jk} < c_3$, and bounded support vector if the upper bound (equality) is reached.
- A j'^{th} sample in $\mathcal{L}_H^+ \cup \mathcal{U}^-$ is called: Support vector of subgroup k if $\alpha_{j'k} > 0$. Unbounded support vector of subgroup k if in addition $\alpha_{j'k} < c_{22}$ for $j' \in \mathcal{L}_H^-$ or $\alpha_{j'k} < c_3$ for $j' \in \mathcal{U}^-$, and bounded support vector if the upper bound (equality) is reached.

This definition is extended from classic binary Support Vector Mathine (SVM) and has very similar geometric interpretations. Next, we define a term based on the characteristics of the sample partition, which will serve as a sufficient condition for the existence of the parametric solution of (**OPT D**).

Definition 3. (Qualified Solution) We say that the solution of the dual problem is Qualified if its corresponding sample partition contains at least one unbounded support vector in both $\mathcal{L}^+ \cup \mathcal{L}^- \cup \mathcal{U}^+$ and $\mathcal{L}^+_H \cup \mathcal{U}^-$

It is worth noting that requiring qualified solution is a very mild condition. Indeed since the unbounded support vectors are essentially the sample points that lie on the decision boundaries that construct the classifiers (and its interception). In order to have meaningful classification in practice this condition is necessary and is expected to be satisfied with even a few samples. The following theorem characterizes the solution structure and provides explicit forms.

Theorem 1. Assume that the solution of (OPT_D) is qualified and induces a set of active and inactive constraints A and A^C , respectively. Denote the composed matrices

$$oldsymbol{H} riangleq rac{oldsymbol{Q}^{-1} \widetilde{oldsymbol{y}} \widetilde{oldsymbol{y}}^T oldsymbol{Q}^{-1}}{\widetilde{oldsymbol{y}}^T oldsymbol{Q}^{-1} \widetilde{oldsymbol{y}}} - oldsymbol{Q}^{-1}$$

$$T \triangleq H(C_A^{\theta})^T; \quad P \triangleq C_A^{\theta} H(C_A^{\theta})^T.$$
 (7)

and the vector $\widetilde{e} \triangleq C_{\mathcal{A}}^{\theta} H e$. Then we have

- 1 The matrix **H** is symmetric negative semi-definite, and **P** is symmetric strictly negative definite hence is invertible.
- 2 The optimal solution is a continuous piecewise affine function of λ . And in the critical region defined by

$$\begin{cases} P^{-1}(C_{\mathcal{A}}^{\lambda}\lambda + C_{\mathcal{A}}^{0} + \widetilde{e}) \ge 0 \\ C_{\mathcal{A}^{C}}^{\lambda}\lambda + C_{\mathcal{A}^{C}}^{0} - C_{\mathcal{A}^{C}}^{\theta}TP^{-1}(C_{\mathcal{A}}^{\lambda}\lambda + C_{\mathcal{A}}^{0} + \widetilde{e}) \ge 0 \end{cases}$$
(8)

the optimal solution $oldsymbol{ heta}^*$ of $(\emph{OPT_D})$ admits a closed form

$$\boldsymbol{\theta}^*(\boldsymbol{\lambda}) = \boldsymbol{T} \boldsymbol{P}^{-1} (\boldsymbol{C}_{\mathcal{A}}^{\lambda} \boldsymbol{\lambda} + \boldsymbol{C}_{\mathcal{A}}^0 + \widetilde{\boldsymbol{e}})$$
 (9)

3 The optimal objective is a continuous piece-wise quadratic (**PWQ**) function of λ .

In essence the theorem indicates that each time the inner optimization (**OPT_D**) is solved, full information (closed form solution) in a well-defined neighborhood (critical region) can be retrieved as a function of outer optimization variables, i.e., those newly introduced "hidden variables".

Besides the local explicitness result, the next theorem describes the overall geometric structure of the optimality, showing that globally the optimal objective is convex in λ , which inspires the proposed learning algorithm for HS³M.

Theorem 2. Still assuming qualified solution,

- 1. The dual optimization has finite number of polyhedron critical regions CR_1, \dots, CR_{N_r} which constitute a partition of the feasible set of λ , i.e., each feasible λ belongs to one and only one critical region.
- 2. The optimal objective $\mathcal{J}(\theta^*(\lambda))$ is a **convex** PWQ function of λ , and is almost everywhere differentiable.

With the local explicit solution and global convex PWQ structure of $\mathcal{J}(\theta^*(\lambda))$ revealed by the theoretical analysis, the learning problem is reduced to maximizing (minimizing) a non-smooth but convex (concave) function in the space $\lambda \in \Lambda$. Concave minimization is well-known to be NP-hard. Despite the hardness, there exist several global optimality conditions. In this work one of these conditions [35] and a level set idea are adopted to establish a an optimization algorithm that is able to approach global optima. The derivation and convergence analysis of the algorithm is long and technical, hence is postponed to Appendix for interested readers. It's worth pointing out that the proposed optimization strategy can be extended to a much broader class of non-convex machine learning problems, and the parametric analysis done in the section complements prior studies on PQP. These two points constitutes the theoretical contribution of this paper.

V. EVALUATION

The authors are collaborating in the "Micro-Synchrophasor for Power Distribution Networks" project [2] to develop a network of μ -PMU devices (μ -Pnet) at the distribution level. In this paper, data from μ -PMUs are used to validate the proposed method. Each μ -PMU provides 120 samples per second for three-phase voltage and current magnitude and phase angle [3],[4].

A. Data Collection and Feature Selection

The raw data sets collected from the μ -PMU measurement are multi-stream time series, which are transmitted and stored with the sMAP protocol [36]. Notation-wise, raw data is written as $\{X_1,\cdots,X_T\}$, where each X_t is a $M\times C$ dimensional vector. In the current experiment, M=5 is the number of μ -PMUs and C=12 is the number of channels of each μ -PMU. Because the raw data is in millisecond's resolution and almost all practical events happen at a larger time scale, one can safely use a sliding window to extract useful information. The window size L is set to 12 according to the time scale of the event of interest. We denote $w_t^i\triangleq\{x_t^i,\cdots,x_{t+L}^i\}$ as the t^{th} window of stream i. For the purpose of detecting different types of events, miscellaneous single stream and inter-stream feature extraction are performed. A summary of computed feature are given in Table I.

With the presented feature extraction procedure, a total number of 312 features have been obtained. However, some of them may be redundant as there are significant similarities among extracted features, for example, when the three phases are balanced, their single stream mean, variation, etc., are almost the same. From a machine learning point of view, adding redundant features does not help event detection, but instead introduces extra noise and cause computational difficulties. In this work, we adopt a method developed in [37], called Minimum-redundancy-maximumrelevance (mRMR). The procedure uses mutual information as the metric of goodness of a feature set, and resolve the tradeoff between relevancy and redundancy. For each event, mRMR is conducted to choose 20 most informative features [38][39]. Also note that all numerical experiments in the following are conducted on a workstation having dual Xeon X5687 CPUs and 72GB memory.

B. Performance of HS³M

The overall task is not only identifying the occurrence of events versus stable state, i.e., the binary classification, but also distinguishing 4 types of events including Voltage Disturbance (VD) and Voltage Sag (VS), Motor Start (MS), High Impedance fault (HI), i.e., the multi-class classification. The training set contains about 40000 μ PMU records with detailed labels (completely labeled data). The testing data set contains the similar events and has around 30000 data points, but is collected at a different time. For the training of HS³M which enables the inclusion of partial knowledge, another 36000 partially labeled data and 108000 unlabeled data are also used (the effect of the size of these data sets will be discussed later).

TABLE I: Extracted Candidate Features

Single Stream	Statistics	$\begin{aligned} & \operatorname{mean}(w_t^i), \operatorname{var}(w_t^i), \operatorname{range}(w_t^i) \\ & \operatorname{median}(w_t^i), \operatorname{entropy}(w_t^i), \operatorname{hist}(w_t^i) \end{aligned}$
	Difference	$u_t^i = \text{Diff}(x_t^i)$; Statistics
	Transformation	$fft(w_t^i)$, wavelet (w_t^i)
Inter Stream	Deviation	$x^i - x^j \forall i, \forall j \in \mathcal{N}(i)$
	Correlation	$corr(x^i, x^j) \forall i, \forall j \in \mathcal{N}(i)$

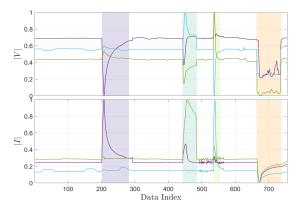


Fig. 4: Detection results over Time. Note that some periods of stable state are shrunk for visualization purpose.

The performance of HS³M is compared with other popular multi-class event detection methods, including cost sensitive versions of Ada Boost, Decision Tree, and Gaussian Process Classification. The hyper-parameters of those models, e.g., the cost balance coefficient, are chosen with 10-folds cross validation (CV). A brief introduction of those methods and their implementation details are listed below:

- Decision tree method: It create a tree-like model that predicts the value of a target variable by learning simple decision rules inferred from the data features. The conditional inference tree algorithm is implemented in this work, which uses multiple significance tests to grow the tree. For cross validation, the maximum tree depth is varied from 10 to 50 by factors of 10, and the splitting threshold is varied from 0.1 to 0.9 with 0.1 intervals.
- Ada Boosting Method: Boosting generates a prediction model by combining many weak classifiers into a stronger classification committee. The AdaBoost procedure is implemented to combine basic tree classifiers for ensemble learning. We vary the maximum tree depth from 10 to 50 by factors of ten. The number of boosting iterations range from 100 to 500 by a step size of 50.
- Gaussian Process Classification (GPC): Instead of directly parameterizing a latent function for classification, GPC models it with a generic Gaussian process. The posterior of the process is updated with training data set, and is squashed through a logistic function for classification. The GPC is implemented with the help of kernlab package ², which includes several approximation algorithms for acceleration. The radial basis kernel is used with kernel width tuned from 2⁻⁵ to 2⁴.

The binary detection results are shown in Figure 4 and its performance are listed in Table II together with other

²http://www.jstatsoft.org/v11/i09/

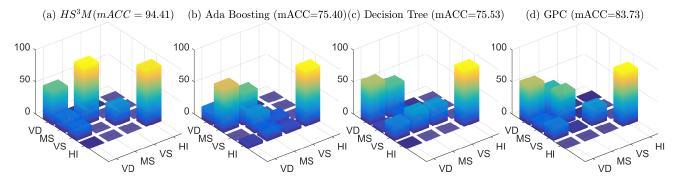


Fig. 5: Confusion Matrix for different methods. Diagonal terms are correct identifications and off-diagonal ones are misclassifications. mACC for multi-class detection accuracy. Note that the class of stable state is not included for better visualization.

methods under comparison. We observe that HS³M performs extremely well in distinguishing stable state and events, with 99.13% binary accuracy (bACC), only 2.33% false positive (false alarm) rate and 0% false negative (miss detection) rate. Compared to other methods, the proposed HS³M outperforms the rest, by at least 9% of bACC and 10% of false positive with respect to the runner-up. The computational cost in terms of training and testing (or prediction) time/memory usage are also listed in Table II. It appears that HS³M requires longer time and about median memory usage in the training phase. This is expected since HS³M is a more comprehensive method incorporating partial information. On the other hand, the testing time/memory usage of HS³M are one of the shortest/smallest. This is due to the solution sparsity (Definition 2) of the HS³M classifier. In practice, testing cost is of major concern because event prediction should be done in real time on distributed systems, while training can be performed "offline" on powerful computers. In this regards, the proposed HS³M is promising also when computational cost is a concern.

To compare the performance on distinguishing event types, the confusion matrices (contingency table) for all methods are shown in Figure 5. Each row of the sub-figure represents the samples in predicted class while each column represents the samples in actual (true) class. The overall multi-class accuracy (mACC) is summarized in the title of each sub-figure. We see from the confusion matrices that significant improvement is achieved: HS³M provides 94.41% mACC, outperforming the best of the other methods by around 11%, while the classic tree based method only yields 75% classification accuracy. Moreover, HS³M gives improvements in differentiating all event types, especially VS, MS, and HI with an accuracy at least 90%. The only issue is that it tends to confuse VD with VS, which is somewhat expected as the criteria for distinguishing VD and VS events are thresholding on the voltage magnitude. In short, the results justified the effectiveness of the proposed HS³M, as well as the idea of incorporating partial information for event detection.

C. Effect of Partial Information

Last but not least, the benefit of including additional partially labeled data and unlabeled data is investigated. To do this, HS³M is evaluated with 0, 7200, 21600 and 36000

TABLE II: Comparison of (Binary) Detection Performance and Computational Cost. bACC for binary Accuracy.

Method	HS ³ M	Ada Boost.	Decision Tree	GPC
bACC (%)	99.13	90.11	89.20	92.73
False Positive (%)	2.33	12.3	16.87	10.1
False Negative (%)	0.0	5.32	0.01	2.31
Time Train (min)	35.7	17.4	11.8	28.6
Time Test (sec)	53.7	206.6	41.2	221.9
Mem. Train (MB)	193	166	59	299
Mem. Test (MB)	0.79	2.91	0.52	292

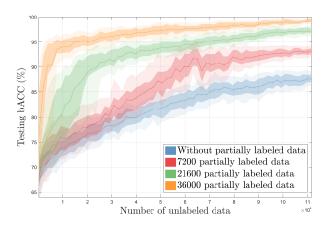


Fig. 6: The incorporation of partially and unlabeled data

partially labeled data samples and 0 - 108000 unlabeled data points. The testing accuracy (bACC) of each test, averaged over 50 random sampled experiments, is plotted in Figure 6 (diamond line), together with the 0.75 and 0.90 confidence intervals (shaded area). Note that when no partially labeled data is included (blue line), HS³M reduces to the semisupervised (kernelized) support vector machine [40]. In the case where unlabeled data is not incorporated (the beginning of each line), HS³M can be viewed as a multi-class variations of the consensus learning recently proposed in [41]. In general, it is observed that the performance improves as more partially labeled data and unlabeled data are added, while the improvements exhibit a "diminishing return" property, i.e., the marginal benefit of including more and more partially/unlabeled data is decreasing. Besides, it appears clearly that HS³M, by leveraging both source of information, significantly outperforms previous semi-supervised and consensus learning.

VI. CONCLUSION AND FUTURE WORK

In summary, with the help of high resolution μ -PMU measurement, a pure data-driven framework is designed for distribution network event detection by incorporating all types of accessible labeled, partially labeled and unlabeled data. More importantly, we developed a new optimization algorithm, by revealing the locally explicit and globally convex property of the dual solution. The experimental results on a real-world μ -PMU dataset justified the effectiveness of the proposed HS 3 M and the PDOP global optimization method, as well as the benefit of incorporating partial knowledge for event detection in power distribution network.

For future work, we will implement HS^3M methods for a large volume of μ -PMUs data streams, as well as developing online learning HS^3M method. More recorded events will be used to train our algorithm and perform feature selection. More importantly, the spatial-temporal characters of large scale events in power systems will be investigated to detect possible cascading effect.

APPENDIX

A. Derivation of the PDOP algorithm

With a slight abuse of notation, we omit the intermediate variable and denote the non-smooth convex function $\mathcal{J}(\theta^*(\lambda))$ by $\mathcal{J}(\lambda)$. From previous discussion, we consider the following non-smooth convex maximization:

$$\max_{\lambda \in \mathbb{P}} \mathcal{J}(\lambda) \tag{10}$$

We first provide a global optimality condition for maximizing non-smooth convex function, then we propose a corresponding optimization algorithm that approaches global solution with level set augmentation.

The global optimality condition for maximizing smooth convex function, in particular convex quadratic functions, has been studied before [35]. However, none of them consider non-smooth piecewise defined functions. In this work, we extend Strekalovsky's condition to non-smooth cases. First of all, the notion of level set is defined as the set of variables that produces the same function values, i.e.,

Definition 4. The level set of function \mathcal{J} at λ is defined by

$$E_{\mathcal{J}(\boldsymbol{\lambda})} = \{ q \in \mathbb{R}^n \mid \mathcal{J}(\boldsymbol{\lambda}) = \mathcal{J}(\boldsymbol{\lambda}) \}$$

The following result gives sufficient and necessary condition for a point λ^* to be the global maximizer of non-smooth convex function $\mathcal{J}(\lambda)$

Theorem 3. Assume that \mathcal{J} is not constant, then λ^* is a global optimal solution of problem (10) if and only if for all $\lambda \in \mathbb{P}$, $\mathbf{q} \in E_{\mathcal{J}(\mathbf{p}^*)}$, $g(\mathbf{q}) \in \partial \mathcal{J}(\mathbf{q})$

$$(\lambda - q)^T g(q) \le 0 \tag{11}$$

where $\partial \mathcal{J}(q)$ is the set of subgradients of \mathcal{J} at λ .

By virtue of Theorem 3, we can check the optimality of any point λ by solving

$$\Delta(\boldsymbol{\lambda}) \triangleq \max_{\substack{\boldsymbol{q} \in E_{\mathcal{J}(\boldsymbol{\lambda})}, \ \boldsymbol{\lambda}' \in \mathbb{P} \\ g(\boldsymbol{q}) \in \partial \mathcal{J}(\boldsymbol{q})}} (\boldsymbol{\lambda}' - \boldsymbol{q})^T g(\boldsymbol{q})$$
(12)

We call the above maximization auxiliary problem at λ . This seems to be a hard problem because (1) it is bilinear in decision variables and (2) usually the level set $E_{\mathcal{J}(\lambda)}$ cannot be calculated explicitly. Next we study solution method for (12) based on approximating the level set with a collection of representative points. To begin with, let us formally define the notion

Definition 5. Given a user specified approximation degree m, the approximation level set for $E_{\mathcal{J}(\lambda)}$ is defined by

$$A_{\boldsymbol{\lambda}}^m = \left\{ \boldsymbol{q}^1, \boldsymbol{q}^2, \cdots, \boldsymbol{q}^m \mid \boldsymbol{q}^i \in E_{\mathcal{J}(\boldsymbol{\lambda})} \ i = 1, 2, \cdots, m \right\}$$

Now consider solving the auxiliary problem approximately by replacing $E_{\mathcal{J}(\lambda)}$ with A^m_{λ} , we obtain that for each q^i , problem (12) becomes

$$\max_{\boldsymbol{\lambda} \in \mathbb{P}, \ g(\boldsymbol{q}^i) \in \partial \mathcal{J}(\boldsymbol{q}^i)} \ (\boldsymbol{\lambda} - \boldsymbol{q}^i)^T g(\boldsymbol{q}^i)$$
 (13)

Since $\mathcal{J}(\lambda)$ is almost everywhere differentiable, in most cases $g(q^i)$ is unique and equals to the gradient $\nabla \mathcal{J}(q^i)$. Then the auxiliary problem is a simple linear program. In the cases when q^i is on the boundary of critical regions, $\partial \mathcal{J}(q^i)$ becomes a convex set and the auxiliary problem becomes a bilinear program. Generally bilinear program is hard, but fortunately (12) has disjoint feasible sets and one can easily check that it is equivalent to

$$\max_{\boldsymbol{\lambda} \in \mathbb{P}} \left\{ \max_{g(\boldsymbol{q}^i) \in V(\partial \mathcal{J}(\boldsymbol{q}^i))} (\boldsymbol{\lambda} - \boldsymbol{q}^i)^T g(\boldsymbol{q}^i) \right\}$$
(14)

which uses the fact that the optimal solution to (13) must be on the vertex of the feasible polyhedron, i.e., $g^*(q^i)$ must be one of the vertices of $\partial \mathcal{J}(q^i)$. Moreover, since $\mathcal{J}(\lambda)$ can be viewed as a pointwise maximum, the subgradient is the convex hull of one side derivatives of the neighboring critical regions. The above analysis suggests an enumerative method for (12) by solving a set of linear programs each with an element in A_{∞}^{m} and a vertex of $\partial \mathcal{J}(q^i)$.

Having solved the approximate auxiliary problem, we can immediately determine if an improvement can be made with current approximate level set. Let $\{(\boldsymbol{u}^i, \boldsymbol{s}^i), i = 1, \cdots, m\}$ be the solution of (13) on the approximate level set $A^{\infty}_{\boldsymbol{\lambda}}$, i.e.,

$$(\boldsymbol{u}^i - \boldsymbol{q}^i)^T \boldsymbol{s}^i = \max_{\boldsymbol{\lambda} \in \mathbb{P}, \ g(\boldsymbol{q}^i) \in V(\partial \mathcal{J}(\boldsymbol{q}^i))} (\boldsymbol{\lambda} - \boldsymbol{q}^i)^T g(\boldsymbol{q}^i) \quad (15)$$

and define

$$\Delta(A_{\lambda}^{m}) = \max_{i=1,\cdots,m} (\boldsymbol{u}^{i} - \boldsymbol{q}^{i})^{T} \boldsymbol{s}^{i}$$
(16)

Then with the convexity of \mathcal{J} we have

Proposition 2. For any $\lambda \in \mathbb{P}$, if there exist $\mathbf{q}^i \in A^m_{\lambda}$, $g(\mathbf{q}^i) \in V(\partial \mathcal{J}(\mathbf{q}^i))$, and \mathbf{u}^i defined in (15), such that $(\mathbf{u}^i - \mathbf{q}^i)^T g(\mathbf{q}^i) > 0$, then $\mathcal{J}(\mathbf{u}^i) > \mathcal{J}(\lambda)$.

Now the remaining problem is to construct approximate level set given current point λ and the degree m. The following

Algorithm 1 Parametric Dual Optimization Procedure

```
Choose \boldsymbol{\lambda}^{(0)} \in \Lambda; set k=0; compute \boldsymbol{\lambda}_* with subgradient descent. while k \leq \text{iter\_max do}
Starting from \boldsymbol{\lambda}^{(k)}, find local minimizer \boldsymbol{r}^{(k)} \in \Lambda with existing methods. Construct approximate level set A^m_{\boldsymbol{r}^{(k)}} at \boldsymbol{r}^{(k)} for \boldsymbol{q}^i \in A^m_{\boldsymbol{r}^{(k)}} do for \boldsymbol{g}^j \in V(\partial \mathcal{J}(\boldsymbol{q}^i)) do Solve \boldsymbol{u}_{ij} = \operatorname{argmax}_{\boldsymbol{\lambda} \in \Lambda} \ (\boldsymbol{\lambda} - \boldsymbol{q}^i)^T \boldsymbol{g}^j #linear programming end for Let j^* = \max \operatorname{Index}_j \{\boldsymbol{u}_{ij}\}; \ (\boldsymbol{u}^i, \boldsymbol{s}^i) = (\boldsymbol{u}_{ij^*}, \boldsymbol{g}^{j^*}); end for Let i^* = \max \operatorname{Index}_{i=1,\cdots,m} \{(\boldsymbol{u}^i - \boldsymbol{q}^i)^T \boldsymbol{s}^i\}; \Delta(A^m_{\boldsymbol{\lambda}}) = (\boldsymbol{u}^{i^*} - \boldsymbol{q}^{i^*})^T \boldsymbol{s}^{i^*}; \ \boldsymbol{u}^{(k)} = \boldsymbol{u}^{i^*}; if \Delta(A^m_{\boldsymbol{\lambda}}) > 0 then Set \boldsymbol{\lambda}^{(k+1)} = \boldsymbol{u}^{(k)}; \ k = k+1; # improvement found else Terminate and output \boldsymbol{\lambda}^{(k)}; # Global optimality condition end if Collecting explored critical region and explicit forms in Theorem 1.
```

results shows that this is theoretically possible with the help of a global minimizer.

Lemma 1. Let the global minimizer of $\mathcal{J}(\lambda)$ be λ_* , then for any $h \in \mathbb{R}^n$ and $\lambda \neq \lambda_*$, there exist a unique positive scalar γ , such that $\lambda_* + \gamma h \in E_{\mathcal{J}(\lambda)}$.

With this guarantee, we write approximate level set by

$$A_{\lambda}^{m} = \{q^{1}, q^{2}, \cdots, q^{m} \mid q^{i} = \lambda_{*} + \gamma_{i} h^{i} \in E_{\mathcal{J}(\lambda)}\}$$
 (17)

To explore directions for improvement, a natural choice of h is a set of orthogonal basis. Specifically, we could start with a random h^1 and use Gram-Schmidt algorithm to extend it to m orthogonal basis. For each h^i , the corresponding γ_i is found by solving:

$$\Phi(\gamma_i) \triangleq \mathcal{J}(\lambda_* + \gamma_i h^i) - \mathcal{J}(\lambda) = 0$$
 (18)

As stated in Lemma 1, the above function has a unique root, which can be computed efficiently with line searching method such as the Bisection algorithm. While another problem is to solve $\lambda_* = \operatorname{argmin} \mathcal{J}(\lambda)$, which is a convex minimization problem. We adopt a sub-gradient descent method since in each critical region, the gradient can be calculated explicitly with Theorem 1. By further considering the convex PWQ structure of the problem, we have

Lemma 2. Let $\sup_{\lambda} ||\lambda^{(1)} - \lambda|| = B$, and the Lipschitz constant of \mathcal{J} be G, then sub-gradient descent with iteration T and optimal step size $\tau_i = B/G\sqrt{T}$ converges to global minimum within $O\left(n/\sqrt{T}\right)$. To be specific, let \mathcal{J}_* be the global minimum then

$$\mathcal{J}(\boldsymbol{\theta}^*(\boldsymbol{\lambda}_{best}^{(T)})) - \mathcal{J}_* \leq \frac{BG}{\sqrt{T}} \leq O\left(\frac{n}{\sqrt{T}}\right), \quad \text{where}$$

$$\mathcal{J}(\boldsymbol{\theta}^*(\boldsymbol{\lambda}_{best}^{(T)})) \triangleq \min \left\{ \mathcal{J}(\boldsymbol{\alpha}^*(\boldsymbol{\lambda}^{(1)})), \cdots, \mathcal{J}(\boldsymbol{\theta}^*(\boldsymbol{\lambda}^{(T)})) \right\}$$

The overall PDM procedure for $\max_{\pmb{\lambda}\in\mathbb{P}}\mathcal{J}(\pmb{\lambda})$ is summarized in Algorithm 1. Given current solution $\pmb{\lambda}^{(k)}$, the algorithm first tries to improve it with existing methods such as AO, CCCP, and SGD. After finding a local minimizer $\pmb{r}^{(k)}$, approximate level set $A^m_{\pmb{r}^{(k)}}$ is obtained by solving (18) and constructing (17). With $A^m_{\pmb{r}^{(k)}}$ and vertices of current subgradient set, a series of linear programming is solved to pick up vector $\pmb{u}^{(k)}$ and subgradient \pmb{s} that maximize condition (11)

of Theorem 3. If this maximal value, i.e., $\Delta(A_{\lambda}^m)$, is greater than 0, then by Proposition 2, $\boldsymbol{u}^{(k)}$ must be strict improvement of $\boldsymbol{r}^{(k)}$. The algorithm iterate until no improvement could be found at current point with the associated approximate level set. Combining Theorem 3 and Proposition 2, we have

Theorem 4. Algorithm I generates a sequence $\{\lambda^{(1)}, \dots, \lambda^{(k)}, \dots\}$ that converges to the global maximizer of $\mathcal{J}(\lambda)$ in a finite number of steps or finds an approximate maximizer at user specified degree m.

REFERENCES

- Alexandra Von Meier, David Culler, Alex McEachern, and Reza Arghandeh. Micro-synchrophasors for distribution systems. In *Innovative Smart Grid Technologies Conference (ISGT)*, pages 1–5. IEEE, 2014.
- [2] Lawrence Jones. Every moment counts: Synchrophasors for distribution networks with variable resources. In *Renewable Energy Integration*. Academic Press, Boston, 2014.
- [3] Alexandra von Meier and Reza Arghandeh. Diagnostic applications for micro-synchrophasor measurements, 2014.
- [4] http://www.powersensorsltd.com/PQube3.php. Micro synchrophasor catalog, 2015.
- [5] Reza Arghandeh, Martin Gahr, Alexandra von Meier, Guido Cavraro, Monika Ruh, and Goran Andersson. Topology detection in microgrids with micro-synchrophasors. In *Power & Energy Society General Meeting*. IEEE, 2015.
- [6] Miles Wen, Reza Arghandeh, Alexandra von Meier, Kameshwar Poolla, and Victor Li. Phase identification in distribution networks with microsynchrophasors. In *Power & Energy Society General Meeting*. IEEE, 2015
- [7] Luca Schenato, Grazia Barchi, David Macii, Reza Arghandeh, Kameshwar Poolla, and Alexandra Von Meier. Bayesian linear state estimation using smart meters and pmus measurements in distribution grids. In *International Conference on Smart Grid Communications (SmartGrid-Comm)*, pages 572–577. IEEE, 2014.
- [8] Mark F McGranaghan and Surya Santoso. Challenges and trends in analyses of electric power quality measurement data. EURASIP Journal on Advances in Signal Processing, (1), 2007.
- [9] IEEE recommended practice for monitoring electric power quality. IEEE Std 1159-2009 (Revision of IEEE Std 1159-1995), June 2009.
- [10] IEEE draft guide for the use of IEEE Std 1641, standard for signal and test definition. *IEEE P1641.1/D3*, Aug 2012.
- [11] Alexander. Eigeles. Ieee standard 1459: a long overdue document [power quality]. In *IEEE Technical Conference on Industrial and Commercial Power Systems*, May 2003.
- [12] Yihai Zhu, Jun Yan, Yan Lindsay Sun, and Haibo He. Revealing cascading failure vulnerability in power grids using risk-graph. *IEEE Transactions on Parallel and Distributed Systems*, 25(12):3274–3284, 2014.
- [13] Yihai Zhu, Jun Yan, Yufei Tang, Yan Lindsay Sun, and Haibo He. Resilience analysis of power grids under the sequential attack. *IEEE Transactions on Information Forensics and Security*, 9(12):2340–2354, 2014.
- [14] Yihai Zhu, Jun Yan, Yufei Tang, Yan Lindsay Sun, and Haibo He. Joint substation-transmission line vulnerability assessment against the smart grid. *IEEE Transactions on Information Forensics and Security*, 10(5):1010–1024, 2015.
- [15] M Vaiman, Keith Bell, Y Chen, B Chowdhury, I Dobson, P Hines, M Papic, S Miller, and P Zhang. Risk assessment of cascading outages: Methodologies and challenges. *IEEE Transactions on Power Systems*, 27(2):631, 2012.
- [16] Sakis Meliopoulos, George Cokkinides, Renke Huang, Evangelos Farantatos, Sungyun Choi, and Yonghee Lee. Wide area dynamic monitoring and stability controls. In iREP Symposium onBulk Power System Dynamics and Control (iREP)-VIII (iREP). IEEE, 2010.
- [17] Rolf Isermann. Model-based fault-detection and diagnosis-status and applications. Annual Reviews in control, 29(1), 2005.
- [18] Borhan M Sanandaji, Eilyan Bitar, Kameshwar Poolla, and Tyrone L Vincent. An abrupt change detection heuristic with applications to cyber data attacks on power systems. In American Control Conference (ACC). IEEE, 2014.

- [19] Xuemei Ding, Josiah Poon, Ivan Celanovic, and Alejandro D Dominguez-Garcia. Fault detection and isolation filters for three-phase ac-dc power electronics systems. *IEEE Trans. on Circuits and Systems* 1: Regular Papers, 60(4), 2013.
- [20] Anselm Schwarte, Frank Kimmidi, and Rolf Isermann. Model-based fault detection of a diesel engine with turbo charger-a case study. In Fault Detection, Supervision and Safety of Technical Processes 2003 (SAFEPROCESS 2003): A Proceedings Volume from the 5th IFAC Symposium, Washington, DC, USA, 9-11 June 2003, volume 1. Elsevier, 2004.
- [21] Barry Mather. Quasi-static time-series test feeder for PV integration analysis on distribution systems. In *Power and Energy Society General Meeting*. IEEE, 2012.
- [22] Joe Qin. Survey on data-driven industrial process monitoring and diagnosis. *Annual Reviews in Control*, 36(2), 2012.
- [23] Jae Yeon Baek, Philippe Leray, Anne-Laure Charley, and Costas J Spanos. Real-time inspection system utilizing scatterometry pupil data. *Journal of Micro/Nanolithography, MEMS, and MOEMS*, 13(4), 2014.
- [24] Dan Li, Yuxun Zhou, Guoqiang Hu, and Costas J Spanos. Fault detection and diagnosis for building cooling system with a tree-structured learning method. *Energy and Buildings*, 2016.
- [25] Vladimir Naumovich Vapnik and Vlamimir Vapnik. Statistical learning theory, volume 1. Wiley New York, 1998.
- [26] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [27] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. Unsupervised learning. In *The elements of statistical learning*, pages 485–585. Springer, 2009.
- [28] Olivier Chapelle, Bernhard Scholkopf, and Alexander Zien. Semisupervised learning (chapelle, o. et al., eds.; 2006)[book reviews]. *IEEE Transactions on Neural Networks*, 20(3):542–542, 2009.
- [29] Xiaojin Zhu and Andrew B Goldberg. Introduction to semi-supervised learning. *Synthesis lectures on artificial intelligence and machine learning*, 3(1):1–130, 2009.
- [30] Christodoulos Floudas. *Nonlinear and mixed-integer optimization: fundamentals and applications*. Oxford University Press, 1995.
- [31] John. Boot. On sensitivity analysis in convex quadratic programming problems. *Operations Research*, 11(5), 1963.
- [32] Petter Tndel, Tor Arne Johansen, and Alberto Bemporad. An algorithm for multi-parametric quadratic programming and explicit MPC solutions. *Automatica*, 39(3), 2003.
- [33] Tor Johansen, Warren Jackson, Robert Schreiber, Petter Tondel, et al. Hardware architecture design for explicit model predictive control. In American Control Conference. IEEE, 2006.
- [34] Martin Monnigmann and Matthias Jost. Vertex based calculation of explicit mpc laws. In American Control Conference (ACC). IEEE, 2012.
- [35] Alexander S Strekalovsky. Global optimality conditions for nonconvex optimization. *Journal of Global Optimization*, 12(4), 1998.
- [36] Stephen Dawson-Haggerty, Xiaofan Jiang, Gilman Tolle, Jorge Ortiz, and David Culler. smap: a simple measurement and actuation profile for physical information. In *Proceedings of the 8th ACM Conference* on Embedded Networked Sensor Systems. ACM, 2010.
- [37] Hanchuan Peng, Fuhui Long, and Chris Ding. Feature selection based on mutual information criteria of max-dependency, max-relevance, and minredundancy. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 27(8), 2005.
- [38] Yuxun Zhou, Reza Arghandeh, Ioannis Konstantakopoulos, Shayaan Abdullah, and Costas J Spanos. Data-driven event detection with partial knowledge: A hidden structure semi-supervised learning method. In American Control Conference (ACC), pages 5962–5968. IEEE, 2016.
- [39] Yuxun Zhou, Reza Arghandeh, Ioannis Konstantakopoulos, Shayaan Abdullah, Alexandra von Meier, and Costas J Spanos. Abnormal event detection with high resolution micro-pmu data. In *Power Systems Computation Conference (PSCC)*, pages 1–7. IEEE, 2016.
- [40] Xilan Tian, Gilles Gasso, and Stéphane Canu. A multiple kernel framework for inductive semi-supervised sym learning. *Neurocomputing*, 90, 2012.
- [41] Yuxun Zhou, Ninghang Hu, and Costas J Spanos. Veto-consensus multiple kernel learning. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [42] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. ACM Trans. on Intelligent Systems and Technology, 2, 2011.

ACKNOWLEDGMENT

The authors would like to thank Prof. Alexandra von Meier for her helpful advice and insights. This research is sponsored in part by the U.S. NSF award 1640587 and by the Republic of Singapore's National Research Foundation through a grant to the Berkeley Education Alliance for Research in Singapore (BEARS) for the Singapore-Berkeley Building Efficiency and Sustainability in the Tropics (SinBerBEST) Program.



Yuxun Zhou is currently a Ph.D candidate at Department of EECS, UC Berkeley. He obtained a Diplome d'Ingenieur in Applied Mathematics from Ecole Centrale Paris in 2012, and a B.S. degree in Electrical Engineering from Xian Jiaotong University in 2009. Yuxun's research interest includes statistical learning theory and paradigms for modern information rich, large-scale and human-involved systems.



Reza Arghandeh (SM'01-M'13) is an assistant professor in the ECE Dept and the Center for Advanced Power System. He has been a postdoctoral scholar at the University of California, Berkeley's California Institute for Energy and Environment 2013-2015. He has five years industrial experience in power and energy systems. He completed his Ph.D. in Electrical Engineering with a specialization in power systems at Virginia Tech. He holds Master's degrees in Industrial and System Engineering from Virginia Tech 2013 and in Energy Systems from the University

of Manchester 2008. From 2011 to 2013, he was a power system software designer at Electrical Distribution Design Inc. in Virginia. Dr. Arghandehs research interests include, but are not limited to data analysis and decision support for smart grids and smart cities using statistical inference, machine learning, information theory, and operations research. He is a recipient of the Association of Energy Engineers (AEE) Scholarship 2012, the UC Davis Green Tech Fellowship 2011, and the best paper award from the ASME 2012 Power Conference and IEEE PESGM 2015. He is the chair of the IEEE Task Force on Big Data Application for Power Distribution Network and secretary of the IEEE Working Group on Distribution Power Quality.



Costas J. Spanos (M'77-SM'92-F'98) received the EE Diploma from the National Technical University of Athens, Greece in 1980 and the M.S. and Ph.D. degrees in ECE from Carnegie Mellon University in 1981 and 1985, respectively. In 1988 he joined the Faculty at the department of Electrical Engineering and Computer Sciences of the University of California at Berkeley. He has served as the Director of the Berkeley Microlab, the Associate Dean for Research in the College of Engineering and as the Chair of the Department of EECS. He works in

statistical analysis in the design and fabrication of integrated circuits, and on novel sensors and computer-aided techniques in semiconductor manufacturing. He also works on statistical data mining techniques for energy efficiency applications. He has participated in two successful startup companies, Timbre Tech, (acquired by Tokyo Electron) and OnWafer Technologies (acquired by KLA-Tencor). He is presently the Director of the Center of Information Technology Research in the Interest of Society (CITRIS) and the Chief Technical Officer for the Berkeley Educational Alliance for Research in Singapore (BEARS).