

A Two-Layer Approach for Estimating Behind-the-Meter PV Generation Using Smart Meter Data

Fankun Bu¹, Graduate Student Member, IEEE, Rui Cheng¹, Graduate Student Member, IEEE, and Zhaoyu Wang¹, Senior Member, IEEE

Abstract—As the cost of the residential solar system decreases, rooftop photovoltaic (PV) has been widely integrated into distribution systems. Most rooftop PV systems are installed behind-the-meter (BTM), i.e., only the net demand is metered, while the native demand and PV generation are not separately recorded. Under this condition, the PV generation and native demand are invisible to utilities, which brings challenges for optimal distribution system operation and expansion. In this paper, we have come up with a novel two-layer approach to disaggregate the unknown PV generation and native demand from the known hourly net demand data recorded by smart meters: 1) At the aggregate level, the proposed approach separates the aggregate PV generation time series from the aggregate net demand time series for customers with PVs. 2) At the customer level, the separated aggregate-level PV generation is allocated to individual PVs. These two layers leverage the spatial correlations of native demand and PV generation, respectively. One primary advantage of our proposed approach is that it is more independent and practical compared to previous works because it does not require PV array parameters, meteorological data and previously recorded solar power exemplars. This paper has verified our proposed approach using real native demand and PV generation data.

Index Terms—Behind-the-meter, distribution system, PV generation estimation, rooftop photovoltaic, smart meter.

I. INTRODUCTION

IN THE last decade, residential rooftop photovoltaic (PV) has been proliferating in distribution systems. In most cases, utilities only install a bi-directional smart meter to record the net demand of customers with PVs. This type of installation is referred to as behind-the-meter (BTM), in which case the net demand equals native demand minus PV generation. Therefore, the PV generation produced by solar array and the native demand consumed by appliances are unknown to utilities. Only metering the net demand can reduce the financial cost for utilities;

Manuscript received October 30, 2021; revised February 27, 2022; accepted March 29, 2022. This work was supported in part by the National Science Foundation under Grant EPCN 2042314, and in part by the Grid Modernization Initiative of the U.S. Department of Energy (DOE) under Grant GMLC project 2.1.1 – FASTDERMS. Paper no. TPWRS-01692-2021. (Corresponding author: Zhaoyu Wang.)

The authors are with the Department of Electrical and Computer Engineering, Iowa State University, Ames, IA 50011 USA (e-mail: fbu@iastate.edu; ruicheng@iastate.edu; wzy@iastate.edu).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TPWRS.2022.3164010>.

Digital Object Identifier 10.1109/TPWRS.2022.3164010

however, as the penetration level of PV increases, the unobservability of notable PV generation and native demand brings significant challenges to distribution systems. We focus on three specific applications to elaborate the necessity of estimating the unknown BTM PV generation and native demand: *First*, the unavailability of native load and PV generation might cause unacceptable forecasting errors because some forecasters require reconstituting the generation and native demand time series [1], [2]. In contrast, knowing BTM PV generation and native load can help utilities forecast generation and load separately, thus provide utilities useful information regarding load/generation growth. *Second*, the invisibility of PV generation and native load can hinder designing optimal service restoration plans [3], [4]. During the restoration stage after an outage, the native demand might be several times higher than the pre-outage demand due to the simultaneous restarting of a large number of air-conditioning appliances. This anomalous demand should be estimated for optimal restoration plans because it can damage electric devices when simultaneously restoring a large number of customers. In practice, utilities usually multiply the normal native demand before outage by a ratio to estimate the anomalous demand during restoration. Also, utilities typically do not consider PVs as reliable restoration sources [3]. Therefore, separating normal native demand and generation is needed for designing optimal restoration plans. *Third*, the unobservability of native demand and solar generation might cause inaccurate reliability analysis. When evaluating a transmission system's reliability, each distribution system is generally simplified as a bus whose native load duration curve is constructed [5], [6]. For those utilities with a high-penetration PV integration, directly using the net demand to construct the load duration curve can significantly underestimate the actual native load [7]. This is because the net demand is typically smaller than the native demand due to the existence of PV generation. In contrast, using the native demand separated from the net demand can help construct more accurate load duration curves. In summary, disaggregating BTM PV generation and native demand from the recorded net demand can enhance distribution system observability and awareness and can also provide more accurate information for transmission system reliability analysis.

Previous works on BTM PV generation disaggregation can be categorized into two types: *Type I - Model-based approaches*: PV array performance model is employed to represent physical

83 PV arrays. In [8], a PV model is combined with a clear sky model
 84 to estimate customer-level solar generation. In [9], a virtual
 85 equivalent PV station model is utilized to represent the aggregate
 86 generation of BTM PVs within a region. In [10] and [11],
 87 a physical PV model and a statistical model are utilized to
 88 estimate BTM solar generation and native demand, respectively.
 89 One primary disadvantage of these model-based approaches is
 90 that detailed PV array parameters or accurate meteorological
 91 data are required. However, in practice, these parameters are
 92 typically unavailable to utilities. Also, acquiring meteorological
 93 data might cause additional costs to utilities. *Type II - Model-free*
 94 *approaches:* In [12] and [13], net demands under heterogeneous
 95 weather conditions are employed to estimate BTM PV capac-
 96 ity, which is then multiplied by a standard solar power time
 97 series to infer BTM PV generations. In [14], native demand
 98 and PV generation are estimated using 1-second net demand
 99 data by identifying appliances' states, which are then leveraged
 100 to estimate appliance demands and solar power. Based on the
 101 variation difference between load and solar power, in [15], an
 102 approach is proposed for estimating service transformer-level
 103 PV generation. In [16], regional-level generation is estimated
 104 by installing additional sensors to record typical PV generation
 105 profiles. In [17], feeder-level solar generation is estimated by
 106 utilizing net load measurements and a nearby PV farm's gener-
 107 ation readings. Using known native loads for customers without
 108 PVs and the generations for a limited number of observable PVs,
 109 in [18], the authors formulate an optimization process to estimate
 110 the aggregated native load and PV generation. In [19], a feder-
 111 ated learning-based framework is proposed to probabilistically
 112 estimate community-level BTM solar generation. In [20], an
 113 approach is developed to estimate the reactive power by taking
 114 advantage of the correlation between the weekly nighttime and
 115 daytime native reactive power demands. Furthermore, previ-
 116 ously in [21] and [22], we have proposed two approaches for
 117 estimating the unknown BTM generation using measured solar
 118 power exemplars. One primary shortcoming of the model-free
 119 approaches is that they rely on contextual information, i.e.,
 120 recorded solar power exemplars or meteorological data, which
 121 might bring additional costs to utilities.

122 Considering the shortcomings of previous approaches, this
 123 paper proposes a novel BTM PV generation and native dem-
 124 and estimation framework which does not require *previously*
 125 recorded solar power and meteorological measurements. Our
 126 approach is based on two findings from real data. The first finding
 127 is the spatial correlation of native load, i.e., the native demands of
 128 two sizeable residential customer groups are strongly correlated
 129 and have highly homogeneous shapes. The second finding is the
 130 spatial correlation of solar power generation, i.e., the generations
 131 for two PVs in a distribution system are significantly correlated
 132 and have highly similar profiles.

133 Our proposed approach contains **two** layers: (1) At the aggre-
 134 gate level, the total generation of all BTM PVs is estimated by
 135 leveraging our first finding. (2) At the customer level, utilizing
 136 our second finding, the estimated aggregate BTM PV generation
 137 is allocated to individual customers. Utilizing the two findings
 138 improves our approach's robustness against the customer-level
 139 load uncertainty [23]. The second layer contains three steps:

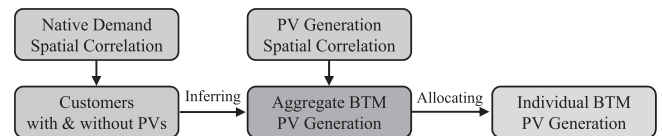


Fig. 1. Overall structure of the proposed BTM PV generation estimation approach.

140 first, our approach trains a model to produce multiple candidate
 141 generation time series, using solar power data generated by a
 142 publicly available tool. Second, our approach determines the
 143 peak generation for each PV. Finally, the allocating procedure
 144 is formulated as an optimization problem. The overall structure
 145 of our proposed approach is shown in Fig. 1. This paper has
 146 verified our proposed approach using real hourly native demand
 147 and PV generation data [24].

148 Smart meters can record individual customers' demands at
 149 an interval of one hour or shorter. Such fine-grained temporal
 150 and spatial granularity can give us more details than traditional
 151 monthly bills. Many researchers have developed advanced ap-
 152 proaches to mine useful information from smart meter data. For
 153 example, [25] utilizes smart meter measurement to perform state
 154 estimation for enhancing distribution system observability, [26]
 155 employs water consumption data recorded by smart water met-
 156 ers to train aggregate water demand forecasters, [27] utilizes
 157 high-resolution phasor measurement units' data to conduct false
 158 data detection, and data redundancy strengthening, [28] converts
 159 smart meter data into manageable load profiles via linearizing
 160 load patterns. Our proposed approach takes advantage of smart
 161 meter data's temporal and spatial granularity to perform BTM
 162 generation estimation.

163 The main contributions of our paper are summarized as
 164 follows: (1) This paper proposes an approach that does not
 165 rely on PV array parameters, historical meteorological data,
 166 and pre-recorded generation exemplars. This independence can
 167 significantly improve the viability of our approach because ac-
 168 quiring the above three types of information can bring challenges
 169 or additional costs for utilities. (2) Our approach only relies on
 170 the net demands of customers with PVs and the native demands
 171 of customers without PVs for estimating the aggregate-level
 172 PV generation. These two types of demands - net and native
 173 - are typically available to utilities, making our approach sig-
 174 nificantly practical. (3) Our approach innovatively estimates
 175 individual PV-installed customers' peak generations by mining
 176 net demand data. The peak generations are then utilized to
 177 estimate individual PV-installed customers' BTM generation
 178 time series.

179 Throughout the paper, vectors are denoted using bold italic
 180 letters, and matrices are represented as bold non-italic letters. In
 181 addition, we adopt the sign convention that the native demand
 182 consumed by customers and the power output from PVs are both
 183 positive.

184 The rest of the paper is organized as follows: Section II
 185 introduces our first and second findings regarding spatial cor-
 186 relation of native demand/generation. Section III presents how
 187 we estimate the aggregate generation for customers with PVs.

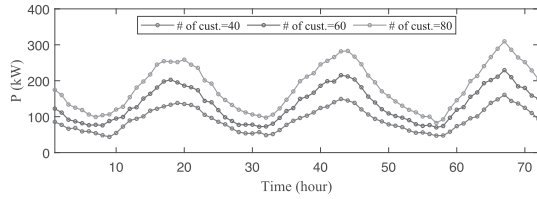


Fig. 2. Three-day actual native demand curves for three example groups with different customer numbers.

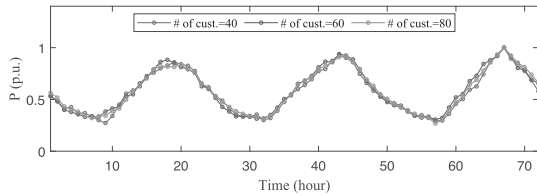


Fig. 3. Three-day normalized native demand curves for three example groups with different customer numbers.

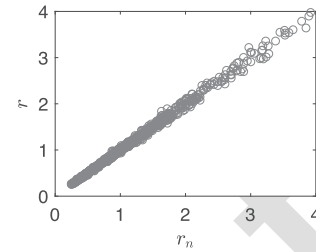


Fig. 4. The relationship between native demand ratio and the nocturnal native demand ratio between two example customer groups.

188 Section IV presents the procedure of formulating and solving
 189 an optimization problem to allocate the estimated aggregate
 190 generation to individual PVs. In Section V, case studies are
 191 analyzed. Section VI concludes the paper.

192 II. SPATIAL CORRELATION OF NATIVE DEMAND/PV 193 GENERATION

194 A. Finding 1: Native Demand Spatial Correlation Between 195 Two Sizeable Groups

196 By examining real residential native demand data, we find
 197 that once the customer numbers for two groups reach a certain
 198 level, their native demands are highly correlated. This finding
 199 is leveraged for estimating the *aggregate* native demand time
 200 series for customers with PVs.

201 Specifically, we use native demand curves to illustrate the
 202 observed spatial correlation. Fig. 2 presents real native demand
 203 curves for three example groups with different customer num-
 204 bers, i.e., 40, 60, and 80, respectively. We can observe that these
 205 three curves demonstrate almost identical shapes, although they
 206 have different magnitudes. The high shape similarity can also
 207 be corroborated by Fig. 3, which presents normalized native
 208 demand curves corresponding to the curves in Fig. 2. Note that
 209 the normalized curves are obtained by dividing the real curves
 210 in Fig. 2 by their peaks, respectively.

211 To stress the importance of Fig. 3, we first define **two** types
 212 of customer groups: the residential customers **with** and **without**
 213 PVs. These two customer groups are denoted as C_w and C_o ,
 214 respectively. For C_o , its native demand is recorded by smart
 215 meters. For C_w , we only know its net demand, and we do not
 216 know its native demand. Our goal is to estimate C_w 's unknown
 217 native demand and thus to estimate its PV generation. Therefore,
 218 Fig. 3 inspires us that given the known native demand curve of
 219 C_o , we can infer the unknown native demand curve of C_w by
 220 multiplying the native demand curve of C_o by a *ratio*, r .

221 Since the native demands for the customers in C_o are directly
 222 recorded by smart meters, the native demand curve of C_o can be
 223 obtained by aggregating the native demand time series over the
 224 customers in C_o . The challenge for inferring the unknown native
 225 demand curve of C_w is that the ratio, r , is unknown and needs to
 226 be estimated. The unknown of r is caused by the unavailability of
 227 the native demand during the daytime for the customers in C_w .
 228 This is because PV generates power during the daytime, which
 229 masks the native demand in the case of net metering. Thus, we
 230 cannot use daytime native demand to compute r . Instead, we use
 231 the nocturnal native demand to estimate r because PV does not
 232 generate power during nighttime, and thus the nocturnal native
 233 demand for C_w is known. Based on the above inference, we
 234 propose first utilizing the nocturnal native demand to compute
 235 a nocturnal native demand ratio, r_n , and then approximating r
 236 as r_n .

237 One *pre-condition* for approximating r as r_n is that r should
 238 be close to r_n . To verify this condition, we randomly select two
 239 groups with different customer numbers ranging from 20 to 80.
 240 Then, for each group, the native demand time series are spatially
 241 aggregated over customers to obtain an aggregate native demand
 242 time series. After that, we compute r using the two groups' native
 243 demand time series throughout a certain period, and compute
 244 r_n using the two groups' native demand time series only during
 245 *nighttime* within that period. Finally, we plot r against r_n , as
 246 shown in Fig. 4. We can see that r is almost identical with r_n .
 247 Therefore, we can accurately estimate r by directly letting it
 248 equal r_n .

249 Once we obtain the estimate of r , we can compute the un-
 250 known native demand of C_w by multiplying the known native
 251 demand of C_o by the estimate of r . After that, estimating
 252 the unknown PV generation of C_w is straightforward, i.e., by
 253 subtracting the recorded net demand measurements from the
 254 estimated native demand.

255 B. Finding 2: Generation Spatial Correlation Between Two 256 PVs

257 There are two primary factors that determine the generation
 258 spatial correlation: (1) In most cases, a distribution system is
 259 geographically bounded in a small district. (2) The most widely
 260 available sampling resolution for smart meters is 1-hour. Under
 261 these two conditions, different PV arrays are subject to nearly
 262 identical meteorological inputs. Thus, the identical inputs can
 263 result in highly similar shapes among PV generation curves.

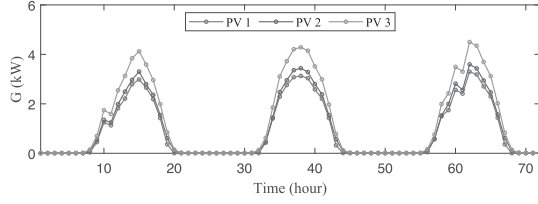


Fig. 5. Three-day real generation curves for three example PVs with different capacities.

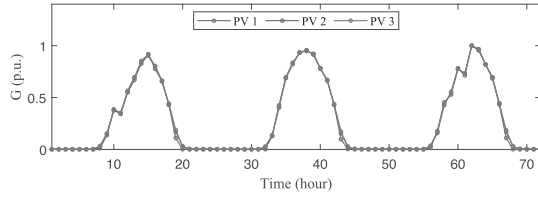


Fig. 6. Three-day normalized generation curves for three example PVs with different capacities.

264 Fig. 5 presents three example PV generation curves correspond-
 265 ing to different PV array capacities. Similar to the native demand
 266 curves for sizeable customer groups, these three generation
 267 curves also demonstrate significant spatial correlation, i.e., they
 268 possess highly similar shapes. This high similarity can also be
 269 corroborated by Fig. 6, where the normalized generation curves
 270 corresponding to the three curves in Fig. 5 overlap with each
 271 other. Most importantly, Figs. 5 and 6 inspire us that estimating
 272 a BTM PV generation curve comes down to two steps: first,
 273 determine the generation curve's shape, and then determine
 274 its magnitude. This two-step method can notably simplify the
 275 estimation of unknown BTM PV generation time series. This
 276 is because compared to model-based methods, our approach is
 277 developed on the foundation of high similarity among generation
 278 curves; therefore, it requires significantly less information.

279 III. ESTIMATING AGGREGATE BTM PV GENERATION FOR 280 CUSTOMERS WITH PVs

281 As elaborated in Section II-A, the native demands of two size-
 282 able customer groups are highly correlated. This high correlation
 283 inspires us that we can infer the unknown native demand of C_w
 284 by multiplying the known native demand of C_o by a ratio:

$$\hat{\mathbf{P}}_w = \mathbf{P}_o r, \quad (1)$$

285 where, $\hat{\mathbf{P}}_w = \{\hat{P}_w(t)\}$ and $\mathbf{P}_o = \{P_o(t)\}$, $t = 1, \dots, T$, denote
 286 the estimated native demand time series for C_w and the actual
 287 native demand time series for C_o , respectively. T is the total
 288 number of native demands in a selected window (e.g., one
 289 month). $P_o(t)$ is computed by aggregating the measured native
 290 demands over customers without PVs:

$$P_o(t) = \sum_{i=1}^{N_o} P_{o,i}(t), \quad t = 1, \dots, T, \quad (2)$$

where, N_o represents the total number of customers in C_o , i.e.,
 customers without PVs. $P_{o,i}(t)$ denotes the measured *native*
 demand at time t for the i 'th customer in C_o .

In (1), r denotes the native demand ratio between C_w and C_o ,
 and is defined as follows:

$$r = \frac{\sum_{t=1}^T P_w(t)}{\sum_{t=1}^T P_o(t)}. \quad (3)$$

However, as presented in Section II-B, since the *diurnal* native
 demand for C_w is masked by PV generation and unavailable to
 utilities, we need to estimate r using *nocturnal* native demand
 measurements. This approximation method is based on the obser-
 vation that PV does not generate power during nighttime and the
 verification that r and r_n are almost identical. Specifically,
 we use r_n to approximate r :

$$\hat{r} = r_n = \frac{\sum_{t \in I_n} P_w(t)}{\sum_{t \in I_n} P_o(t)}, \quad (4)$$

where, I_n denotes the set of nighttime hours. In our paper, I_n
 refers to the hours between 9:00 P.M. and 5:00 A.M. Note that
 for the hours in I_n , since PV does not generate power, $P_w(t)$
 equals the known aggregate *net* demand, $P'_w(t)$. Therefore,

$$\hat{r} = \frac{\sum_{t \in I_n} P'_w(t)}{\sum_{t \in I_n} P_o(t)}, \quad (5)$$

where, $P'_w(t)$ is computed by aggregating the measured net
 demands over customers in C_w :

$$P'_w(t) = \sum_{i=1}^{N_w} P'_{w,i}(t), \quad t = 1, \dots, T, \quad (6)$$

where, N_w represents the total number of customers in C_w .
 $P'_{w,i}(t)$ denotes the measured *net* demand at time t for the i 'th
 customer in C_w .

Then, using the estimate of r and the known native demand
 time series for C_o , we can apply (1) to compute the estimated
 native demand time series for C_w . Finally, inferring the PV
 generation time series for C_w , $\hat{\mathbf{G}}_w = \{\hat{G}_w(t)\}$, $t = 1, \dots, T$,
 is straightforward:

$$\hat{\mathbf{G}}_w = \hat{\mathbf{P}}_w - \mathbf{P}'_w, \quad (7)$$

where, $\mathbf{P}'_w = \{P'_w(t)\}$, $t = 1, \dots, T$, denotes the known net
 demand time series for C_w .

The above procedure for estimating the aggregate-level PV
 generation and native demand for C_w are illustrated in Fig. 7.

IV. ESTIMATING BTM PV GENERATION FOR EACH INDIVIDUAL PV

Knowing the aggregate BTM PV generation and native demand
 might not be sufficient for some applications [29], [30].
 For example, some demand response schemes require known
 customer-level native demand [12]. Therefore, estimating indi-
 vidual customers' BTM native demand and PV generation is of
 significance.

To achieve this goal, we propose an approach to allocate the
 estimated aggregate PV generation/native demand time series
 to individual customers with PVs. As discussed in Section II-B,

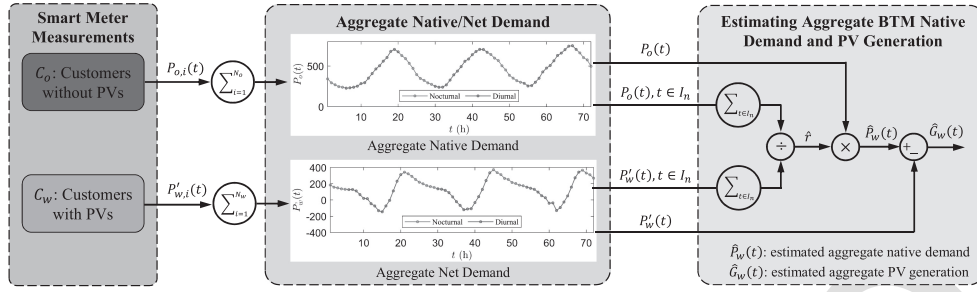


Fig. 7. Detailed structure of the proposed aggregate-level BTM PV generation/native demand estimation.

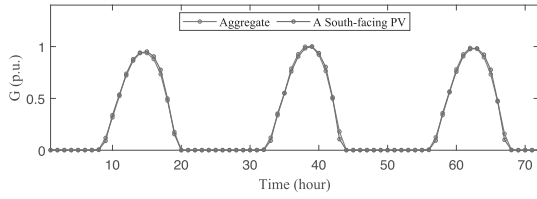


Fig. 8. Three-day normalized aggregate generation curve for all PVs and normalized generation curve for an individual PV facing south.

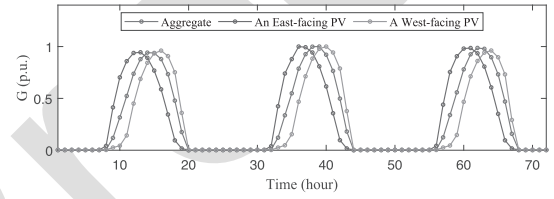


Fig. 9. Three-day normalized aggregate generation curve of all PVs and normalized generation curves for two example PVs facing east and west, respectively.

332 estimating an individual PV's generation curve boils down to
 333 determining the generation curve's shape and its magnitude. In
 334 this section, our approach has three steps to perform allocating:
 335 (Step-I): generate candidate generation curves for individual
 336 PVs; (Step-II): estimate the peak generation for each PV; and
 337 (Step-III): allocate the estimated aggregate PV generation time
 338 series to individual PVs by solving an optimization problem.

339 A. Generating Diverse Candidate Generation Curves for 340 Individual PVs

341 As discussed earlier, in a geographically bounded distribution
 342 system, two primary factors determining a generation curve
 343 are the magnitude and shape. This subsection aims to generate
 344 candidate generation curves for those non-south-facing PVs.
 345 First, we train a regression model using the data generated by
 346 PVWatts Calculator. Then, we feed the estimated generation
 347 curve of a south-facing PV into the trained model to infer the
 348 targeted candidate generation curves for those non-south-facing
 349 PVs.

350 In Section III, we have obtained the estimated time series for
 351 the aggregate generation of all PVs. One question is whether we
 352 can use that shape to represent the unknown shapes of individual
 353 PVs. To answer this question, we have conducted a numerical
 354 experiment. First, we normalized the aggregate generation curve
 355 of all PVs by dividing the aggregate generation time series by
 356 its peak. Then, in the same way, we normalized the generation
 357 curve of an example PV facing *south*. The two normalized
 358 curves are plotted in Fig. 8. It can be seen that the normalized
 359 curve corresponding to the aggregate generation for all PVs
 360 is highly similar to the normalized curve for a south-facing
 361 PV. One primary reason for this similarity is that the majority
 362 of residential PVs face south because a south-facing PV can

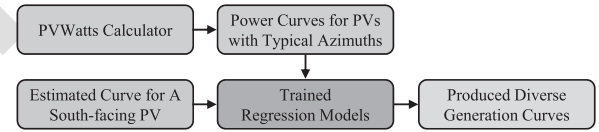


Fig. 10. Overall structure for producing diverse candidate PV generation curves using power output data generated by PVWatts Calculator.

363 typically generate more power than PVs in other directions. Most
 364 importantly, Fig. 8 tells us that a south-facing PV's generation
 365 curve can be accurately represented by the normalized aggregate
 366 generation curve of all PVs.

367 Note that in distribution systems, in addition to the majority
 368 of south-facing PVs, there exist some residential PVs with other
 369 azimuths, such as east or west. These non-south-facing PVs'
 370 generation curves cannot be fully represented by the normalized
 371 aggregate PV generation curve in Fig. 8. Specifically, compared to
 372 the normalized aggregate PV generation curve, the normalized
 373 generation curves for an east-facing PV and a west-facing PV
 374 are somewhat "left-skewed" and "right-skewed," respectively,
 375 as shown in Fig. 9. Therefore, it is necessary to obtain candidate
 376 shapes for those non-south-facing PVs' generation curves. To
 377 achieve this goal, our basic idea is first to feed PV power data
 378 generated by PVWatts Calculator into a regression model to
 379 capture the relationship between the generations for a south-
 380 facing PV and a non-south-facing PV. Then, the aggregate
 381 generation curve estimated in Section III, which can accurately
 382 represent a south-facing PV's generation curve, is fed into the
 383 trained regression model to produce diverse generation curves
 384 corresponding to non-south azimuths. The overall structure is
 385 shown in Fig. 10:

1) *Training a Gaussian Process Regression Model:* Since the shape of a south-facing PV's generation curve can be approximated as the shape of the aggregate generation curve of all PVs, one intuitive way for inferring *non-south-facing* PVs' candidate shapes is to produce diverse shapes based on the *south-facing* PV's estimated generation curve. This idea is based on our observation that there exists a mapping between the generation curves for PVs with different azimuths. Therefore, one critical step for producing diverse candidate generation curves is to identify the relationship between a non-south-facing PV's generation curve and a south-facing PV's generation curve. To capture the relationship, first, we use PVWatts Calculator [31], an online application developed by the National Renewable Energy Laboratory (NREL), to generate power output data for PVs with typical azimuths, e.g., east, south, and west. Then, using the generated PV output power data, we train a Gaussian Process Regression (GPR) model to capture the relationship between the generation curve corresponding to a typical azimuth except for south (e.g., east) and the generation curve corresponding to the azimuth of the south. The primary reason for selecting GPR is that after running numerical tests, GPR demonstrated a relatively better performance when applied to our dataset than some other state-of-the-art nonlinear regression models, such as the Support Vector Machine model and the Polynomial regression model.

Specifically, first, we use PVWatts Calculator to generate time-series data for a south-facing PV and a PV with other typical azimuth (e.g., east). Then, each time series is normalized so that the peak generation is 1 p.u. The two normalized time series corresponding to the south-facing PV and the non-south-facing PV are denoted as $\mathbf{G}_s^* = \{G_s^*(t)\}$ and $\mathbf{G}_{ns}^* = \{G_{ns}^*(t)\}$, $t = 1, \dots, T$, respectively. $G_s^*(t)$ and $G_{ns}^*(t)$ denote the normalized generation at time t for a south-facing PV and a non-south-facing PV, respectively. Our goal is to use $G_s^*(t)$ to explain $G_{ns}^*(t)$ because PVs in a geographically bounded distribution system typically have highly correlated generations. By conducting numerical experiments, we find that in addition to $G_s^*(t)$, the hour-in-day, $H_d(t)$, and day-in-year, $D_y(t)$, are also related with $G_{ns}^*(t)$. Therefore, we use $G_s^*(t)$, $H_d(t)$, and $D_y(t)$ as the input variables and $G_{ns}^*(t)$ as the output variable, respectively, to train a GPR model. The function of GPR is to capture the relationship between $G_{ns}^*(t)$ and $G_s^*(t)$. The basic idea behind GPR is that if the distance between two explanatory variables is small, the difference between their corresponding dependent variables will also be relatively small. Specifically, the output, $G_{ns}^*(t)$, is denoted as a function of the input vector, $\mathbf{X}^*(t)$:

$$G_{ns}^*(t) = f(\mathbf{X}^*(t)), \quad (8)$$

where, $\mathbf{X}^*(t) = [G_s^*(t), H_d(t), D_y(t)]^T$. For GPR, $f(\mathbf{X}^*(t))$ is assumed to be a random variable reflecting the uncertainty of functions evaluated at $\mathbf{X}^*(t)$. Specifically, the function $f(\mathbf{X}^*(t))$ is distributed as a Gaussian process:

$$f(\mathbf{X}^*(t)) \sim \mathcal{GP}(\mu(\mathbf{X}^*(t)), K(\mathbf{X}^*(t), \mathbf{X}^*(t'))), \quad (9)$$

where, $\mu(\mathbf{X}^*(t))$ represents the expected value of $f(\mathbf{X}^*(t))$, i.e., the value of $G_{ns}^*(t)$. The covariance function, $K(\mathbf{X}^*(t), \mathbf{X}^*(t'))$, represents the dependence between $G_{ns}^*(t)$'s at different times. In our problem, the covariance function,

$K(\cdot, \cdot)$, is specified by the Squared Exponential Kernel function expressed as:

$$K(\mathbf{X}^*(t), \mathbf{X}^*(t')) = \sigma_f^2 \exp\left(-\frac{\|\mathbf{X}^*(t) - \mathbf{X}^*(t')\|_2^2}{2\sigma^2}\right), \quad (10)$$

where, $\|\cdot\|_2$ represents l_2 -norm, σ_f and σ are hyper-parameters, which are determined using cross-validation. Intuitively, (10) measures the distance between $\mathbf{X}^*(t)$ and $\mathbf{X}^*(t')$, which can also reflect the similarity between $G_{ns}^*(t)$ and $G_{ns}^*(t')$.

Note that $G_s^*(t)$ and $G_{ns}^*(t)$ are generated solar powers using PVWatts Calculator; thus, they are known and a T -dimensional joint Gaussian distribution can be constructed as:

$$\begin{bmatrix} f(\mathbf{X}^*(1)) \\ \vdots \\ f(\mathbf{X}^*(T)) \end{bmatrix} \sim \mathcal{N}(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*), \quad (11)$$

where,

$$\boldsymbol{\mu}^* = \begin{bmatrix} \mu(\mathbf{X}^*(1)) \\ \vdots \\ \mu(\mathbf{X}^*(T)) \end{bmatrix}, \quad (12a)$$

$$\boldsymbol{\Sigma}^* = \begin{bmatrix} K(\mathbf{X}^*(1), \mathbf{X}^*(1)) & \cdots & K(\mathbf{X}^*(1), \mathbf{X}^*(T)) \\ \vdots & \ddots & \vdots \\ K(\mathbf{X}^*(T), \mathbf{X}^*(1)) & \cdots & K(\mathbf{X}^*(T), \mathbf{X}^*(T)) \end{bmatrix}. \quad (12b)$$

The joint Gaussian distribution formulated in (11) represents a trained non-parametric model, which captures the relationship between $G_{ns}^*(t)$ and $G_s^*(t)$.

2) *Inferring a Non-South-Facing PV's Generation Curve:* As shown in Fig. 8, the normalized generation curve for a south-facing PV, $\mathbf{G}_s = \{G_s(t)\}$, $t = 1, \dots, T$, can be approximated as the normalized estimated aggregate generation curve for all PVs:

$$\mathbf{G}_s = \frac{\hat{G}_w}{\hat{G}_m}, \quad (13)$$

where, \hat{G}_m denotes the peak of \hat{G}_w . To infer the unknown generation time series for a non-south-facing PV, $\mathbf{G}_{ns} = \{G_{ns}(t)\}$, $t = 1, \dots, T$, we assume $G_{ns}(t)$ is a function of $G_s(t)$, i.e., $G_{ns}(t) = f(G_s(t))$. By appending $f(G_s(t))$ to the end of (11), an $(N + 1)$ -dimensional joint Gaussian distribution can be constructed as:

$$\begin{bmatrix} G_{ns}^*(1) \\ \vdots \\ G_{ns}^*(T) \\ G_{ns}(t) \end{bmatrix} = \begin{bmatrix} f(\mathbf{X}^*(1)) \\ \vdots \\ f(\mathbf{X}^*(T)) \\ f(\mathbf{X}(t)) \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_* \\ \mu_1 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_* & \boldsymbol{\Sigma}_{*1} \\ \boldsymbol{\Sigma}_{*1}^T & \Sigma_{11} \end{bmatrix}\right), \quad (14)$$

where, $\mathbf{X}(t) = [G_s(t), H_d(t), D_y(t)]^T$ is a vector of explanatory variables. $\boldsymbol{\Sigma}_{*1}$ represents the training-test set covariances and Σ_{11} is the test set covariance. Since $G_{ns}^*(t)$, $\mathbf{X}^*(t)$, and $\mathbf{X}(t)$ are known, using the Bayes rule, the distribution of $G_{ns}(t)$

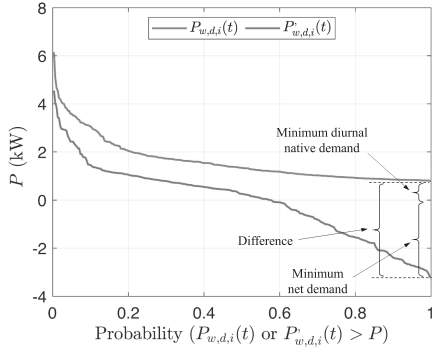


Fig. 11. Load duration curves for an example customer's diurnal native demand and diurnal net demand.

466 conditioned on \mathbf{G}_{ns}^* can be computed as follows:

$$G_{ns}(t) | \mathbf{G}_{ns}^* \sim \mathcal{N}(\mu_1(t), \Sigma_1(t)), \quad (15)$$

467 where, $\mu_1(t) = \Sigma_{*1}^T \Sigma_{*1}^{-1} \mathbf{G}_{ns}^*$ and $\Sigma_1(t) = \Sigma_{11} -$
 468 $\Sigma_{*1}^T \Sigma_{*1}^{-1} \Sigma_{*1}$. Note that $\mu_1(t)$ denotes the most probable value
 469 of the estimated generation at time t for a non-south-facing PV.
 470 By conducting the above inferring procedure for all the t 's, we
 471 can obtain a candidate generation time series corresponding
 472 to a *particular* typical PV azimuth. Since there are multiple
 473 typical azimuths, such as east, and west, we can infer multiple
 474 candidate PV generation time series:

$$\mathbf{G}_{ns}^j = \{G_{ns}^j(t)\}, \quad t = 1, \dots, T, \quad j = 1, \dots, N_{ns}, \quad (16)$$

475 where, $G_{ns}^j(t)$ denotes the inferred PV generation at time t ,
 476 for the j 'th typical non-south-facing azimuth. N_{ns} denotes the
 477 total number of typical non-south-facing PV azimuths and is
 478 determined by conducting numerical experiments.

479 B. Estimating Peak Generation for Each Individual PV

480 Simply knowing the candidate shapes for unknown generation
 481 curves is insufficient for allocating the estimated aggregate
 482 generation to individual PVs. As discussed earlier, we should
 483 also know the magnitudes for the candidate generation curves.
 484 To estimate the peak generation, we employ our observation
 485 from real data that the peak generation is almost identical with
 486 the difference between the minimum diurnal native demand and
 487 the minimum net demand.

488 Specifically, to explain our observation regarding the correlation,
 489 we start with Fig. 11, showing the load duration curves for
 490 the i 'th customer's diurnal *native* demand, $P_{w,d,i}(t)$, and diurnal
 491 *net* demand, $P'_{w,d,i}(t)$. Thus, we can compute the difference
 492 between the minimums of $P_{w,d,i}(t)$ and $P'_{w,d,i}(t)$:

$$D_{w,i} = \underline{P}_{w,d,i} - \underline{P}'_{w,d,i}, \quad (17)$$

493 where, $\underline{P}_{w,d,i}$ and $\underline{P}'_{w,d,i}$ denote the minimums of $P_{w,d,i}(t)$
 494 and $P'_{w,d,i}(t)$ during a selected window, respectively. Note that
 495 $\underline{P}_{w,d,i}$ is positive, and $\underline{P}'_{w,d,i}$ is negative. Then, our finding is
 496 that $D_{w,i}$ is highly similar to the peak generation, $G_{w,m,i}$, as
 497 shown in Fig. 12. This relationship inspires us to approximate

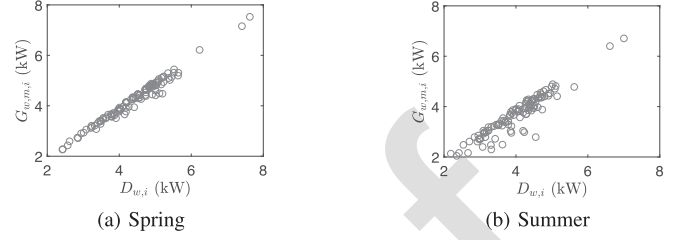


Fig. 12. The relationship between peak generation and the difference between minimum diurnal *native* demand and minimum *net* demand. (a) Spring (b) Summer

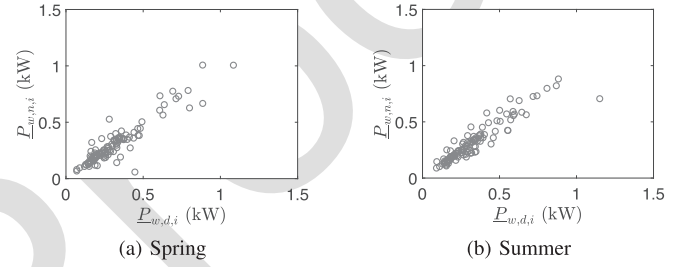


Fig. 13. The relationship between minimum *diurnal* native demand and minimum *nocturnal* native demand. (a) Spring (b) Summer

$G_{w,m,i}$ as $D_{w,i}$:

$$\hat{G}_{w,m,i} = D_{w,i}, \quad i = 1, \dots, N_w, \quad (18)$$

499 where, $\hat{G}_{w,m,i}$ is the estimate of $G_{w,m,i}$. However, one challenge
 500 is that $D_{w,i}$ depends on $\underline{P}_{w,d,i}$, which is unknown due to BTM
 501 PV generation. Therefore, we need to estimate $\underline{P}_{w,d,i}$, which
 502 is involved with another finding from real native demand data.
 503 Specifically, as shown in Fig. 13, the minimum *diurnal* native
 504 demand, $\underline{P}_{w,d,i}$, can be approximated as the minimum *nocturnal*
 505 native demand, $\underline{P}_{w,n,i}$:

$$\underline{P}_{w,d,i} \approx \underline{P}_{w,n,i}, \quad i = 1, \dots, N_w. \quad (19)$$

Note that since PV does not generate power during nighttime,
 506 $\underline{P}_{w,n,i}$ is known to utilities. Finally, using the estimate of $\underline{P}_{w,d,i}$
 507 and the known $\underline{P}'_{w,d,i}$, we can compute $D_{w,i}$ using (17), and
 508 then compute $\hat{G}_{w,m,i}$ using (18).
 509

510 C. Allocating the Estimated Aggregate PV Generation to 511 Individual PVs

512 Sections III, IV-A, and IV-B provide the estimated aggregate
 513 generation time series of all PVs, inferred candidate generation
 514 curves for individual PVs, and estimated generation peaks for
 515 individual PVs, respectively. Therefore, estimating individual
 516 PVs' generation curves comes down to allocating the estimated
 517 aggregate generation time series to individual PVs. This allocat-
 518 ing procedure is formulated as an optimization process:

$$\min_{\mathbf{K}, \gamma} \|\mathbf{G}_e * \mathbf{K} * \mathbf{1} - \hat{\mathbf{G}}_w\|_2^2 + \lambda * \|\gamma\|_2^2 \quad (20a)$$

$$s. t. \quad \mathbf{G}_e * \mathbf{K} \leq \mathbf{1} * (\hat{\mathbf{G}}_{w,m} + \gamma)^T, \quad (20b)$$

$$\mathbf{0} \leq \gamma \leq P_0 * \mathbf{1}, \quad (20c)$$

519 where, $\mathbf{G}_e = [\mathbf{G}_s, \mathbf{G}_{n_s}^1, \dots, \mathbf{G}_{n_s}^{N_s}]$ is a T -by- N_e matrix, which
 520 denotes a collection of candidate generation curves. $N_e =$
 521 $N_s + 1$ denotes the total number of candidate generation curves.
 522 $\mathbf{K} = [\mathbf{K}_1, \dots, \mathbf{K}_{N_w}]$ is an N_e -by- N_w matrix of decision vari-
 523 ables, which denote the weights assigned to candidate generation
 524 curves for individual PVs. $\mathbf{K}_i, i = 1, \dots, N_w$, is an N_e -by-1
 525 vector, which denotes the weights assigned to candidate gener-
 526 ation curves for the i 'th PV. The first $\mathbf{1}$ is an N_w -by-1
 527 vector of ones. $\mathbf{G}_e * \mathbf{K}$ results in a T -by- N_w matrix, which is a collection
 528 of estimated generation time series for individual PVs. The
 529 first term in the objective function (20a) reflects the difference
 530 between the estimated aggregate PV generation, $\hat{\mathbf{G}}_w$, and the
 531 weighted summation of individual PV's estimated generations,
 532 $\mathbf{G}_e * \mathbf{K} * \mathbf{1}$. The second term in the objective function (20a)
 533 considers the estimation errors of peak generations. λ is a
 534 tuning parameter. $\boldsymbol{\gamma}$ is an N_w -by-1 vector with non-negative
 535 elements, which reflect the errors of approximating $G_{w,m,i}$ as
 536 $D_{w,i}$, as shown in (18). The second $\mathbf{1}$ is a T -by-1 vector of ones.
 537 $\hat{\mathbf{G}}_{w,m} = [\hat{G}_{w,m,1}, \dots, \hat{G}_{w,m,N_w}]^T$ denotes an N_w -by-1 vector
 538 of the estimated generation peaks for all PVs. $(\hat{\mathbf{G}}_{w,m} + \boldsymbol{\gamma})$
 539 denotes the corrected generation peaks with consideration of es-
 540 timation errors. $\mathbf{1} * (\hat{\mathbf{G}}_{w,m} + \boldsymbol{\gamma})^T$ produces a T -by- N_w matrix,
 541 in which each column contains the same element. Constraint
 542 (20b) ensures that the estimated generation time series for each
 543 PV is smaller than its estimated peak generation. $\mathbf{0}$ is an N_w -by-1
 544 vector of zeros. P_0 denotes the maximum error of approximating
 545 $G_{w,m,i}$ as $D_{w,i}$ for individual PVs. The third $\mathbf{1}$ is an N_w -by-1
 546 vector of ones. Constraint (20c) ensures that the estimation
 547 errors for individual PVs are non-negative and smaller than an
 548 upper bound. The reason for constraining the elements of $\boldsymbol{\gamma}$ as
 549 non-negative is that $D_{w,i}$ typically under-estimates $G_{w,m,i}$, as
 550 shown in Fig. 13.

551 The optimization process represented in (20) is a convex
 552 quadratic programming problem, thus, we can obtain a unique
 553 solution for \mathbf{K} , i.e., $\mathbf{K}^* = [\mathbf{K}_1^*, \dots, \mathbf{K}_{N_w}^*]$. Then, the estimated
 554 generation time series for the i 'th PV, $\hat{\mathbf{G}}_{w,i} = \{\hat{G}_{w,i}(t)\}, t =$
 555 $1, \dots, T$, can be computed as:

$$\hat{\mathbf{G}}_{w,i} = \mathbf{G}_e * \mathbf{K}_i^*, \quad i = 1, \dots, N_w. \quad (21)$$

556 Then, the estimated native demand time series for the i 'th
 557 customer, $\hat{\mathbf{P}}_{w,i} = \{\hat{P}_{w,i}(t)\}, t = 1, \dots, T$, can be computed as:

$$\hat{\mathbf{P}}_{w,i} = \mathbf{P}'_{w,i} + \hat{\mathbf{G}}_{w,i}, \quad i = 1, \dots, N_w. \quad (22)$$

558 where, $\mathbf{P}'_{w,i} = \{P'_{w,i}(t)\}, t = 1, \dots, T$, denotes the known net
 559 demand time series recorded by smart meter for the i 'th customer
 560 with PVs.

561 Note that (20) can be solved for a selected window. The
 562 window size, T , can impact estimation accuracy and runtime,
 563 which will be examined in the Case Study Section. The detailed
 564 steps for estimating customer-level PV generation are illustrated
 565 in Fig. 14.

566 V. CASE STUDY

567 In this section, the proposed two-layer BTM solar power and
 568 native demand estimation approach is verified using real PV
 569 generation and native demand data.

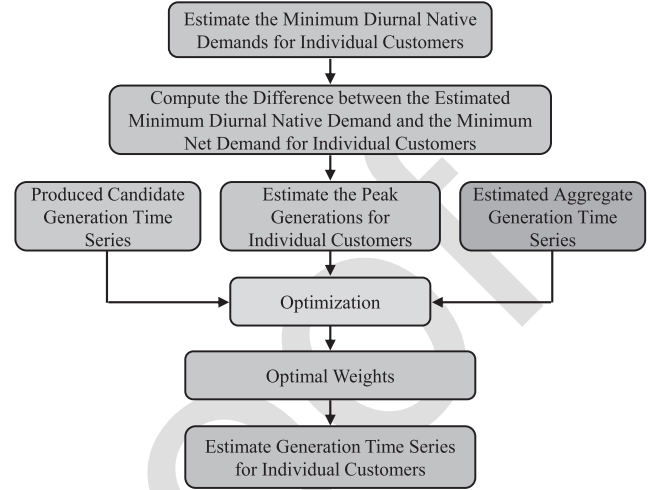


Fig. 14. Detailed steps of the individual customer-level BTM PV generation estimation.

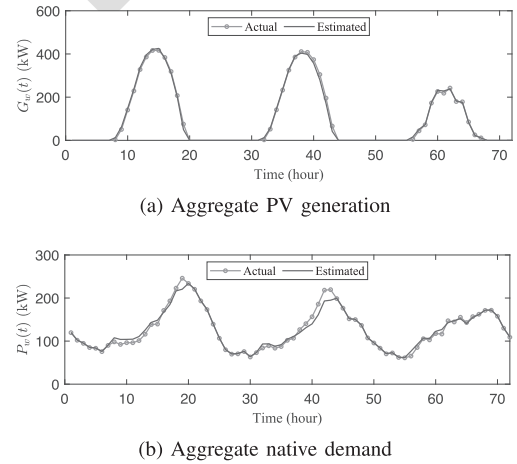


Fig. 15. Three-day actual and estimated aggregate PV generation and native demand curves. (a) Aggregate PV generation (b) Aggregate native demand

A. Dataset Description

570 The hourly native demand and PV generation data used in this
 571 paper are from a public dataset [24]. The time range of native
 572 demand and solar power is one year. This dataset contains a total
 573 number of 100 customers with PVs and 115 customers without
 574 PVs. For the customers with PVs, the net demand is obtained by
 575 subtracting PV generation from native demand.
 576

B. Aggregate-Level BTM PV Generation Estimation Validation

577 Fig. 15 shows three-day actual and estimated aggregate PV
 578 generation/native demand curves. It can be seen that the esti-
 579 mated curves can accurately follow the actual curves. To quanti-
 580 tatively evaluate the estimation accuracy, we compute the mean
 581
 582

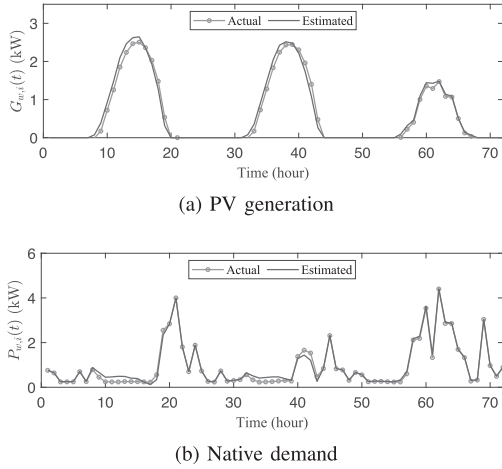


Fig. 16. Three-day actual and estimated PV generation and native demand curves for an example customer with PV. (a) PV generation (b) Native demand

absolute percentage error (MAPE) as follows:

$$MAPE = \frac{100\%}{N_d} \sum_{t \in I_d} \left| \frac{\hat{Y}_w(t) - Y_w(t)}{Y_{w,m}} \right|, \quad (23)$$

where, $\hat{Y}_w(t)$ represents $\hat{G}_w(t)$ or $\hat{P}_w(t)$. $Y_w(t)$ represents $G_w(t)$ or $P_w(t)$. $Y_{w,m}$ represents $G_{w,m}$ or $P_{w,m}$, where $G_{w,m}$ and $P_{w,m}$ denote the actual peaks of PV generation and native demand, respectively. I_d denotes the set of daytime hours. N_d denotes the total number of hours in I_d .

To comprehensively evaluate the performance of our approach, we also compute the mean squared error (MSE) and coefficient of variation (CV):

$$MSE = \frac{1}{N_d} \sum_{t \in I_d} (\hat{Y}_w(t) - Y_w(t))^2, \quad (24)$$

$$CV = \frac{\sigma}{\mu}, \quad (25)$$

where,

$$\mu = \frac{1}{N_d} \sum_{t \in I_d} (\hat{Y}_w(t) - Y_w(t)), \quad (26a)$$

$$\sigma = \sqrt{\frac{1}{N_d - 1} \sum_{t \in I_d} ((\hat{Y}_w(t) - Y_w(t)) - \mu)^2}. \quad (26b)$$

The computed $MAPE$'s for PV generation and native demand are 1.21% and 1.28%, respectively. The computed MSE 's for PV generation and native demand are about 58.09. Note that the actual peaks for the PV generation and native demand are 462.5 and 437.1 kW, respectively. The computed CV 's for PV generation and native demand are about -3.48. The above error metrics reflect the high accuracy of our proposed approach.

C. Customer-Level BTM PV Generation Estimation Validation

1) *Estimation Performance*: Fig. 16 shows three-day actual and estimated PV generation and native demand curves for an example customer with PV. We can see that the estimated curves

TABLE I
EMPIRICAL CDF OF ESTIMATION ERROR METRICS

Empirical CDF	0.1	0.2	0.5	0.7	0.9
$MAPE$ of \hat{G} (%)	2.84	4.05	4.96	6.38	8.80
$MAPE$ of \hat{P} (%)	1.63	2.15	2.80	3.67	4.92
MSE of \hat{G}	0.04	0.06	0.10	0.19	0.33
MSE of \hat{P}	0.03	0.05	0.09	0.18	0.29
CV of \hat{G}	-11.80	-5.13	-2.60	2.37	16.12
CV of \hat{P}	-11.30	-4.65	-2.59	1.77	10.90

can accurately fit the actual curves. To comprehensively examine the performance of our approach, we compute the $MAPE$ for all customers with PVs. Specifically, the $MAPE$'s for the i 'th customer are computed as follows:

$$MAPE_i = \frac{100\%}{N_d} \sum_{t \in I_d} \left| \frac{\hat{Y}_{w,i}(t) - Y_{w,i}(t)}{Y_{w,m,i}} \right| \quad (27)$$

where $Y_{w,i}(t)$ represent $G_{w,i}(t)$ or $P_{w,i}(t)$, $\hat{Y}_{w,i}(t)$ represent $\hat{G}_{w,i}(t)$ or $\hat{P}_{w,i}(t)$, and $Y_{w,m,i}$ represent $G_{w,m,i}$ or $P_{w,m,i}$. $G_{w,m,i}$ and $P_{w,m,i}$ denote the actual generation and native demand peaks for the i 'th customer, respectively. We also compute the MSE and CV for each PV-installed customer:

$$MSE_i = \frac{1}{N_d} \sum_{t \in I_d} (\hat{Y}_{w,i}(t) - Y_{w,i}(t))^2, \quad (28)$$

$$CV_i = \frac{\sigma_i}{\mu_i}, \quad (29)$$

where,

$$\mu_i = \frac{1}{N_d} \sum_{t \in I_d} (\hat{Y}_{w,i}(t) - Y_{w,i}(t)), \quad (30a)$$

$$\sigma_i = \sqrt{\frac{1}{N_d} \sum_{t \in I_d} ((\hat{Y}_{w,i}(t) - Y_{w,i}(t)) - \mu_i)^2}. \quad (30b)$$

Table I summarises the empirical cumulative distribution functions (CDFs) for the estimation $MAPE$, MSE , and CV , which are constructed using all the computed $MAPE$'s, MSE 's, and CV 's, respectively. As can be seen, for the estimated hourly PV generation, 70% of the $MAPE$'s are less than 6.38%. Regarding the estimated hourly native demand, 70% of the $MAPE$'s are less than 3.67%. This effectively verifies the estimation accuracy of our proposed approach. We also provide the percentiles of MSE and CV based on all the PV-installed customers' generation and native demand estimates, which can more comprehensively evaluate the performance of our approach.

Note that the above results are obtained under the conditions that (1) five produced candidate generation curves are employed ($N_e = 5$), (2) the tuning parameter in (20a) is 100 ($\lambda = 100$), and (3) the optimization process specified in (20) is executed for individual windows with a time length of one month ($T = 720$ hours, the entire year is divided into 12 windows).

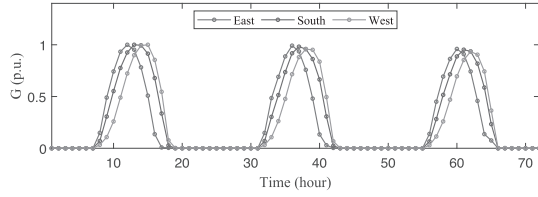


Fig. 17. Three-day produced candidate generation curves corresponding to three typical azimuths, i.e., east, south, and west.

TABLE II
IMPACT OF CANDIDATE GENERATION CURVES

Case	I	II	III
Average $MAPE$ of \hat{G} (%)	5.677	5.474	5.473
Average $MAPE$ of \hat{P} (%)	3.924	3.086	3.086
Runtime (s)	40	125	194

2) *Testing the Candidate Generation Curves:* As elaborated in Section IV-A, diverse candidate generation curves are produced for representing the unknown BTM generation. Thus, it is of interest to examine the effectiveness of producing candidate curves. Fig. 17 shows three produced candidate generation curves corresponding to three typical azimuths, i.e., east, south, and west, respectively. We can observe that compared to the generation curve corresponding to the south, the produced curve corresponding to the east is “left-skewed,” and the produced curve corresponding to the west is “right-skewed”. Therefore, the produced curves demonstrate diversity, which is consistent with our observation on real PV generation curves shown in Fig. 9.

In addition, we have also quantitatively examined the effectiveness of producing diverse candidate generation curves. Specifically, we test the impact of the number of candidate generation curves, i.e., we solve (20) separately for three cases with different numbers of candidate curves: (I) one candidate generation curve corresponding to the azimuth of south; (II) three candidate generation curves corresponding to the east, south, and west, respectively; and (III) five candidate generation curves corresponding to the east, southeast, south, southwest, and west, respectively. The other conditions for the three cases are the same: $\lambda = 100$ and $T = 720$ hours. To evaluate the impact of candidate number, we compute the average $MAPE$ over all PVs’ $MAPE$ ’s obtained from (27). The results are summarized in Table II. We can see that as the candidate number increases, the estimation error decreases, and the execution time increases. In addition, the $MAPE$ for Case I is relatively greater than Case II and III, and Case II and Case III provide nearly identical $MAPE$ ’s. This is because three candidate curves - corresponding to the east, south, and west - can comprehensively represent the unknown BTM generation curve; adding extra candidate curves simply result in a slight accuracy improvement.

3) *Testing the Tuning Parameter λ :* As discussed in Section IV-C, λ in (20) reflects the confidence of estimating peak generations for individual PVs. One general principle for determining λ is that the largest element in γ is a couple

TABLE III
IMPACT OF WINDOW SIZE T

T (month)	1	2	3	4
Average $MAPE$ of \hat{G} (%)	5.47	5.30	5.18	5.08
Average $MAPE$ of \hat{P} (%)	3.09	2.99	2.92	2.87

of kilo-watts. In addition, the solutions for (20) should not be sensitive to λ , i.e., (20) should be robust to λ . To verify the robustness of our proposed approach, we solve (20) based on different values of λ , and then compute the corresponding average $MAPE$ ’s for the estimated PV generation and native demand. Other conditions are that $T = 720$ hours and five candidate generation curves - corresponding to the south, southeast, south, southwest, and west - are employed. The results show that for the λ ’s ranging from 100 to 500 with an interval of 100, the average $MAPE$ ’s for PV generation and native demand do not change (5.47% and 3.09%). The invariant average $MAPE$ ’s demonstrate the robustness of our proposed approach.

4) *Testing the Window Size T :* Since our proposed approach can be conducted for each divided window, it is of importance to examine the impact of window size on estimation accuracy. To do this, we perform our approach for windows with different lengths and then compute the estimation $MAPE$. In Table III, it can be seen that the average $MAPE$ decreases as T increases. This is because for a wider window, the probability for the minimum diurnal native demand, $\underline{P}_{w,d,i}$, equaling the minimum nocturnal native demand, $\underline{P}_{w,n,i}$, is larger. Thus, we have a smaller estimation error for $\underline{P}_{w,d,i}$, as seen in (19). Then, based on (17) and (18), it can be seen that the smaller estimation error for $\underline{P}_{w,d,i}$ results in a more accurate $D_{w,i}$, which then brings a more accurate estimate for $G_{w,m,i}$. Finally, more accurate peak generation estimates result in smaller estimation errors for the PV generation and native demand time series.

D. Performance Comparison

This paper compares our proposed approach with previous works from two perspectives, qualitatively and quantitatively.

1) *Qualitative Analysis:* From a qualitative point of view, one primary advantage of our approach is that it does not require meteorological data and solar generation exemplars. For the aggregate level, our approach can perform PV generation estimation by only using recorded net demand data. For the customer level, our approach can also work by only relying on recorded smart meter data, although leveraging PVWatts Calculator’s generated data can improve the estimation accuracy.

2) *Quantitative Comparison:* For the customer level, we have also compared our approach with previous works. Specifically, we focus on comparing our approach with the method presented in [22] and [11], which demonstrate better performance compared to previous works. Table IV summarizes the computed $MAPE$ ’s for our approach and the compared approach. Note that the average $MAPE$ ’s for our approach have lower and upper bounds because the considered window size, T , ranges from one month to four months. As can be seen, the approach in [22]

TABLE IV
AVERAGE *MAPE* (%) COMPARISON

Approaches	Our Approach	Approach in [11]	Approach in [22]
\hat{G}	[5.08, 5.47]	7.38	5.24
\hat{P}	[2.87, 3.09]	9.94	2.95

TABLE V
AGGREGATE-LEVEL ESTIMATION *MAPE* (%)

	W/O noise	Case 1	Case 2	Case 3	Case 4	Case 5
\hat{G}	1.21	1.17	1.22	1.38	1.53	1.73
\hat{P}	1.28	1.28	1.33	1.43	1.58	1.76

TABLE VI
AVERAGE CUSTOMER-LEVEL ESTIMATION *MAPE* (%)

	W/O noise	Case 1	Case 2	Case 3	Case 4	Case 5
\hat{G}	5.47	5.84	5.86	5.64	5.54	5.62
\hat{P}	3.09	3.53	3.68	3.62	3.63	3.80

demonstrates a similar estimation accuracy as our approach does. However, our approach does not require solar exemplars, which makes it more independent and practical. The approach in [11] employs a statistical model and a physical model to represent the native load and the PV generation, respectively. Table IV shows that our approach has a better performance than the approach in [11] in terms of the average *MAPE*.

E. Robustness Against Measurement and Communication Noises

To test the robustness of our proposed approach, we add measurement and communication noises to the net demand measurements of customers with PVs and the native demand measurements of customers without PVs. For the measurement noise, we consider the Class 0.5 (having $\pm 0.5\%$ error) specified by ANSI C12.20. For the communication noise, we test five different packet loss rates considering that the packet loss rate depends on the communication bandwidth and data volume. For example, we purposely change 1% of the measurements to zero to achieve a 1% packet loss rate. To comprehensively evaluate our approach's performance, we set up five cases: Case 1 – 1% measurement lost + 0.5% random noise, Case 2 – 2% measurement lost + 0.5% random noise, Case 3 – 3% measurement lost + 0.5% random noise, Case 4 – 4% measurement lost + 0.5% random noise, and Case 5 – 5% measurement lost + 0.5% random noise. Then, we apply our approach to the above five cases and compute the *MAPE* for evaluating the robustness. The results are summarized in Tables V and VI. We can observe that the *MAPE*'s slowly increase while the noise level increases, demonstrating the robustness of our approach.

F. Limitations of the Proposed Approach

Every method has its limitations, and there is no omnipotent method that can apply to all cases. The limitation of our proposed approach is that it requires time-series smart meter data with a temporal granularity that can distinguish daytime and nighttime. This is because our approach innovatively utilizes the temporal correlation between the aggregate *nocturnal* native demand and the aggregate *diurnal* native demand. Under this condition, only having access to the monthly demands of those PV-installed customers brings challenges to our approach because it cannot split the monthly demand into two parts, the diurnal and nocturnal demands, for computing the nocturnal native demand ratio. We intend to address this challenge in our future work.

VI. CONCLUSION

This paper is dedicated to proposing an independent and practical BTM solar power/native demand estimation approach. Our proposed approach contains two interconnected layers. The aggregate level leverages the spatial correlation of native demand to perform the aggregate PV generation/native demand estimation. The customer level utilizes the spatial correlation of PV generation to allocate the estimated aggregate PV generation/native demand to individual customers. The Case Study verifies that our approach can accurately estimate BTM PV generation/native demand, significantly enhancing distribution system observability and situation awareness. The numerical experiments also demonstrate that our approach does not require meteorological data and measured solar power exemplars. Therefore, our approach is more independent and thus is practical for utilities to implement.

REFERENCES

- [1] J. Black and V. Rojo, "Long-term load forecast methodology overview," Sep. 2019. [Online]. Available: https://www.iso-ne.com/static-assets/documents/2019/09/p1_load_forecast_methodology.pdf
- [2] F. Wang, Z. Xuan, Z. Zhen, K. Li, T. Wang, and M. Shi, "A day-ahead PV power forecasting method based on LSTM-RNN model and time correlation modification under partial daily pattern prediction framework," *Energy Convers. Manage.*, vol. 212, no. 112766, pp. 1–14, May 2020.
- [3] R. Seguin, J. Woyak, D. Costyk, J. Hambrick, and B. Mather, "High penetration PV integration handbook for distribution engineers," Nat. Renew. Energy Lab., Golden, CO, USA, Tech. Rep. NREL/TP-5D00-63114, 2016.
- [4] T. A. Short, *Electric Power Distribution Handbook*. Boca Raton FL, USA: CRC Press, 2014.
- [5] V. Krishnan and J. D. McCalley, "Building foresight in long-term infrastructure planning using end-effect mitigation models," *IEEE Syst. J.*, vol. 11, no. 4, pp. 2040–2051, Dec. 2017.
- [6] W. Buehring, C. Huber, and J. Marques, "Expansion planning for electrical generating systems," Int. Atomic Energy Agency, Vienna, Austria, Tech. Rep. STI/DOC/10/241, 1984.
- [7] W. Liu, D. Guo, Y. Xu, R. Cheng, Z. Wang, and Y. Li, "Reliability assessment of power systems with photovoltaic power stations based on intelligent state space reduction and pseudo-sequential Monte Carlo simulation," *Energies*, vol. 11, no. 6, 2018, Art. no. 1431.
- [8] D. Chen and D. Irwin, "SunDance: Black-box behind-the-meter solar disaggregation," in *Proc. 8th Int. Conf. Future Energy Syst.*, 2017, pp. 16–19.
- [9] Y. Wang, N. Zhang, Q. Chen, D. S. Kirschen, P. Li, and Q. Xia, "Data-driven probabilistic net load forecasting with high penetration of behind-the-meter PV," *IEEE Trans. Power Syst.*, vol. 33, no. 3, pp. 3255–3264, May 2018.

746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805

- [10] F. Kabir, N. Yu, W. Yao, R. Yang, and Y. Zhang, "Estimation of behind-the-meter solar generation by integrating physical with statistical models," in *Proc. IEEE Int. Conf. Commun. Control Comput. Technol. Smart Grids*, 2019, pp. 1–6.
- [11] F. Kabir, N. Yu, W. Yao, R. Yang, and Y. Zhang, "Joint estimation of behind-the-meter solar generation in a community," *IEEE Trans. Sustain. Energy*, vol. 12, no. 1, pp. 682–694, Jan. 2021.
- [12] K. Li, F. Wang, Z. Mi, M. Fotuhi-Firuzabad, N. Duić, and T. Wang, "Capacity and output power estimation approach of individual behind-the-meter distributed photovoltaic system for demand response baseline estimation," *Appl. Energy*, vol. 253, 2019, Art. no. 113595.
- [13] F. Wang *et al.*, "A distributed PV system capacity estimation approach based on support vector machine with customer net load curve features," *Energies*, vol. 11, no. 7, Jul. 2018, Art. no. 1750.
- [14] C. Dinesh, S. Welikala, Y. Liyanage, M. P. B. Ekanayake, R. I. Godaliyadda, and J. Ekanayake, "Non-intrusive load monitoring under residential solar power influx," *Appl. Energy*, vol. 205, pp. 1068–1080, Aug. 2017.
- [15] F. Sossan, L. Nespoli, V. Medici, and M. Paolone, "Unsupervised disaggregation of photovoltaic production from composite power flow measurements of heterogeneous prosumers," *IEEE Trans. Ind. Informat.*, vol. 14, no. 9, pp. 3904–3913, Sep. 2018.
- [16] H. Shaker, H. Zareipour, E. Muljadi, and D. Wood, "A data-driven approach for estimating the power generation of invisible solar sites," *IEEE Trans. Smart Grid*, vol. 7, no. 5, pp. 2466–2476, Sep. 2016.
- [17] E. C. Kara, C. M. Roberts, M. D. Tabone, L. Alvarez, D. S. Callaway, and E. M. Stewart, "Disaggregating solar generation from feeder-level measurements," *Sustain. Energy, Grids Netw.*, vol. 13, pp. 112–121, 2018.
- [18] K. Li, J. Yan, L. Hu, F. Wang, and N. Zhang, "Two-stage decoupled estimation approach of aggregated baseline load under high penetration of behind-the-meter PV system," *IEEE Trans. Smart Grid*, vol. 12, no. 6, pp. 4876–4885, Nov. 2021.
- [19] J. Lin, J. Ma, and J. Zhu, "A privacy-preserving federated learning method for probabilistic community-level behind-the-meter solar generation disaggregation," *IEEE Trans. Smart Grid*, vol. 13, no. 1, pp. 268–279, Jan. 2022.
- [20] S. Talkington, S. Grijalva, M. J. Reno, and J. A. Azzolini, "Solar PV inverter reactive power disaggregation and control setting estimation," *IEEE Trans. Power Syst.*, to be published, doi: 10.1109/TPWRS.2022.3144676.
- [21] F. Bu, K. Dehghanpour, Y. Yuan, Z. Wang, and Y. Zhang, "A data-driven game-theoretic approach for behind-the-meter PV generation disaggregation," *IEEE Trans. Power Syst.*, vol. 35, no. 4, pp. 3133–3144, Jul. 2020.
- [22] F. Bu, K. Dehghanpour, Y. Yuan, Z. Wang, and Y. Guo, "Disaggregating customer-level behind-the-meter PV generation using smart meter data and solar exemplars," *IEEE Trans. Power Syst.*, vol. 36, no. 6, pp. 5417–5427, Nov. 2021.
- [23] F. Bu, K. Dehghanpour, Y. Yuan, and Z. Wang, "Quantifying load uncertainty using real smart meter data," in *Proc. IEEE Smart Grid Commun*, 2020, pp. 1–6.
- [24] K. Nagasawa, C. R. Upshaw, J. D. Rhodes, C. L. Holcomb, D. A. Walling, and M. E. Webber, "Data management for a large-scale smart grid demonstration project in austin, texas," in *Proc. 6th Int. Conf. Energy Sustain.*, 2012, pp. 1027–1031.
- [25] M. A. Khan and B. Hayes, "Smart meter based two-layer distribution system state estimation in unbalanced MV/LV networks," *IEEE Trans. Ind. Informat.*, vol. 18, no. 1, pp. 688–697, Jan. 2022.
- [26] A. A. Nasser, M. Z. Rashad, and S. E. Hussein, "A two-layer water demand prediction system in urban areas based on micro-services and LSTM neural networks," *IEEE Access*, vol. 8, pp. 147647–147661, 2020.
- [27] Q. Wang, W. Tai, Y. Tang, M. Ni, and S. You, "A two-layer game theoretical attack-defense model for a false data injection attack against power systems," *Int. J. Electr. Power Energy Syst.*, vol. 104, pp. 169–177, 2019.
- [28] Z. A. Khan, D. Jayaweera, and M. S. Alvarez-Alvarado, "A novel approach for load profiling in smart power grids using smart meter data," *Electr. Power Syst. Res.*, vol. 165, pp. 191–198, 2018.
- [29] Q. Zhang, Y. Guo, Z. Wang, and F. Bu, "Distributed optimal conservation voltage reduction in integrated primary-secondary distribution systems," *IEEE Trans. Smart Grid*, vol. 12, no. 5, pp. 3889–3900, Sep. 2021.

- [30] R. Cheng, Z. Wang, Y. Guo, and F. Bu, "Analyzing photovoltaic's impact on conservation voltage reduction in distribution networks," Oct. 2021, arXiv:2110.14777.
- [31] A. P. Dobos, "PVWatts Version 5 Manual," Nat. Renew. Energy Lab., Golden, CO, USA, Tech. Rep. NREL/TP-6A20-62641, Sep. 2014.



Fankun Bu (Graduate Student Member, IEEE) received the B.S. and M.S. degrees from North China Electric Power University, Baoding, China, in 2008 and 2013, respectively. He is currently working toward the Ph.D. degree with the Department of Electrical and Computer Engineering, Iowa State University, Ames, IA, USA. From 2008 to 2010, he was a Electrical Commissioning Engineer for NARI Technology Company, Ltd., Nanjing, China. From 2013 to 2017, he was an Electrical Engineer for State Grid Corporation of China, Nanjing, China. His research interests include power system data analytics, behind-the-meter PV generation estimation, load and PV generation forecasting, distribution system modeling, renewable energy integration, and power system relaying.



Rui Cheng (Graduate Student Member, IEEE) received the B.S. degree in electrical engineering from Hangzhou Dianzi University, Hangzhou, China, in 2015 and the M.S. degree in electrical engineering from North China Electric Power University, Beijing, China, in 2018. He is currently working toward the Ph.D. degree with the Department of Electrical & Computer Engineering, Iowa State University, Ames, IA, USA. From 2018 to 2019, he was an Electrical Engineer with State Grid Corporation of China, Hangzhou, China. His research interests include power distribution systems, voltage/var control, transactive energy markets, and applications of optimization and machine learning methods to power systems.



Zhaoyu Wang (Senior Member, IEEE) received the B.S. and M.S. degrees in electrical engineering from Shanghai Jiaotong University, Shanghai, China, and the M.S. and Ph.D. degrees in electrical and computer engineering from the Georgia Institute of Technology, Atlanta, GA, USA. He is currently the Northrop Grumman Endowed Associate Professor with Iowa State University, Ames, IA, USA. His research interests include optimization and data analytics in power distribution systems and microgrids. He was the recipient of the National Science Foundation CAREER

Award, the Society-Level Outstanding Young Engineer Award from IEEE Power and Energy Society (PES), the Northrop Grumman Endowment, College of Engineering's Early Achievement in Research Award, and the Harpole-Pentair Young Faculty Award Endowment. He is the Principal Investigator for a multitude of projects funded by the National Science Foundation, the Department of Energy, National Laboratories, PSERC, and Iowa Economic Development Authority. He is the Chair of IEEE PES PSOPE Award Subcommittee, the Co-Vice Chair of PES Distribution System Operation and Planning Subcommittee, and the Vice Chair of PES Task Force on Advances in Natural Disaster Mitigation Methods. He is an Associate Editor for IEEE TRANSACTIONS ON POWER SYSTEMS, IEEE TRANSACTIONS ON SMART GRID, IEEE OPEN ACCESS JOURNAL OF POWER AND ENERGY, IEEE POWER ENGINEERING LETTERS, and *IET Smart Grid*.