

Estimating the Unobservable Components of Electricity Demand Response with Inverse Optimization

Adrian Esteban-Perez

Department of Technology and Operations Management, Rotterdam School of Management, Erasmus University, Rotterdam 3062 PA, The Netherlands estebanperez@rsm.nl

Derek Bunn

London Business School, London NW1 4SA, United Kingdom dbunn@london.edu

Yashar Ghiassi-Farrokhfal

Department of Technology and Operations Management, Rotterdam School of Management, Erasmus University, Rotterdam 3062 PA, The Netherlands y.ghiassi@rsm.nl

Understanding and predicting the electricity demand responses to prices are critical activities for system operators, retailers, and regulators. While conventional machine learning and time series analyses have been adequate for the routine demand patterns that have adapted only slowly over many years, the emergence of active consumers with flexible assets such as solar-plus-storage systems, and electric vehicles, introduces new challenges. These active consumers exhibit more complex consumption patterns, the drivers of which are often unobservable to the retailers and system operators. In practice, system operators and retailers can only monitor the net demand (metered at grid connection points), which reflects the overall energy consumption or production exchanged with the grid. As a result, all "behind-the-meter" activities—such as the use of flexibility—remain hidden from these entities. Such behind-the-meter behavior may be controlled by third party agents or incentivized by tariffs; in either case, the retailer's revenue and the system loads would be impacted by these activities behind the meter, but their details can only be inferred. We define the main components of net demand, as baseload, flexible, and self-generation, each having nonlinear responses to market price signals. As flexible demand response and self generation are increasing, this raises a pressing question of whether existing methods still perform well and, if not, whether there is an alternative way to understand and project the unobserved components of behavior. The data-driven inverse optimization methodology characterizes decomposed consumption patterns without requiring direct observation of behind-the-meter behavior or device-level metering. By analyzing net demand as revealed at the grid connection point, we estimate parameters for a latent optimization model, enabling predictions that infer the unobservable components. We validate the approach using real-world data, demonstrating its superior performance in both point and probabilistic forecasting compared to state-of-the-art time-series analysis and machine learning benchmarks.

Key words: Unobservable Behavior, Electricity, Forecasting, Behind-the-Meter, Inverse Optimization, Demand Forecasting

1. Introduction

Understanding and predicting the electricity demand responses to prices are critical activities for system operators (SO), retailers and regulators. For system operators, accurate demand forecasts

are essential to ensure that energy supply meets demand within the grid's capacity limits and other physical constraints (Shahidehpour et al. 2002). For electricity retailers, these predictions are vital for revenue management, cost control, and tariff design to provide better pricing for consumers (Karthikeyan et al. 2013), whilst for regulators and policy-makers, promoting efficient consumer engagement has become an important element in their net zero ambitions (Pinson et al. 2014).

Electricity demand forecasting at the consumer level has been widely researched and, in practice, had until recently reached acceptable levels of accuracy, largely due to the relatively stable nature of consumer behavior (Hatton et al. 2015). However, recent developments associated with the energy transition have radically altered all elements of the electricity supply and value chain. These changes have been driven by substantial subsidies for decarbonization and the emergence of end-user generation. As a result, system operators have recognized an urgent need for energy *flexibility* to help balance supply and demand (Mathieu et al. 2024, Pinson et al. 2014). Flexibility in this context mainly refers to the ability to adapt energy consumption patterns, such as shedding demand during periods of high energy prices or grid stress, or shifting demand to times when energy is cheaper or the grid is less stressed. To promote consumer flexibility, retailers have introduced new incentive mechanisms, including Time-of-Use (TOU) pricing and dynamic pricing models, which align flexibility with grid status and/or temporal energy prices (Vardakas et al. 2014). In response, cost-conscious consumers have increasingly embraced flexibility e.g., by adopting solar-plus-storage systems and opportunistically scheduling their electric vehicle charging, thereby enhancing their ability to manage energy costs and contribute to grid stability (Danti and Magnani 2017).

As a result of these changes, accurate estimation of the consumer demand response in the evolving retail electricity markets has become both more critical and more complex. This importance has increased because energy procurement is now more challenging due to the added volume and uncertainty, pushing the traditional grid closer to its limits and raising the risk of congestion (Koliou et al. 2015). The complexity has grown because, unlike in the past where demand was generally non-price responsive with a stable profile over time, today's consumer demand, as observed by SOs and retailers at the metering point with the grid, is confounded by new unobservable components. The essential baseload demand (unresponsive to price) is supplemented with price-responsive flexible demand (including batteries and electric vehicles) and reduced by any end-user generation (e.g., solar). Accurately estimating these components would require retailers and SOs to know the types and capacities of flexible devices owned by consumers and to have access to real-time device-level data on how consumers operate these assets to manage costs. However, this information, known

as behind-the-meter (BTM) behavior, is typically hidden from SOs and retailers, as the meters only reveal "net demand". In fact, not only are they often unaware of the devices owned by consumers—sometimes even including solar panels—but they also may be unable to access or utilize real-time device-level measurements. This data is typically classified as personal information and protected by privacy legislation, for example by the General Data Protection Regulation (GDPR) in the EU (European Parliament and Council 2016), and sometimes by considerations of cybersecurity in smart systems (Liu et al. 2012). The difficulty and importance of accurately characterizing BTM flexibility behavior has prompted a call to action among practitioners and policymakers. For instance, an IEEE Task Force is now dedicated to the estimation, optimization, and control of BTM behavior (Srivastava et al. 2024).

Consequently, SOs and retailers need to estimate the unobservable behind-the-meter information in the absence of direct measurements. But existing methodologies, such as time-series analysis (Choi et al. 2020, Dong et al. 2017) and machine learning algorithms have generally been applicable only to the observable demand, as metered at the grid connection point ("net demand") rather than the specific latent components of consumer demand which may be flexing (Nghiem and Jones 2017, Wen et al. 2020). To address this methodology gap, we question if the emerging technique of inverse optimization (IO) could offer a promising solution for estimating the unobserved information on these components, thereby potentially enhancing the prediction of observable net demand. IO operates by modeling a consumer's behind-the-meter price response as a latent optimization problem under different pricing schemes and then integrating this into an optimal prediction problem. Intuitively, this method reverses the typical optimisation methodology: rather than determining optimal actions according to a constrained objective function, it seeks to infer the parameters of the model on the basis of revealed actions, assuming a unobservable latent optimisation approach by the consumers. This approach therefore has the potential to capture nonlinear consumer responses by distinguishing between their flexible, inflexible, and self-generation components of demand, to the extent that they are specified in a latent model of optimal demand response. Whether this decomposition then leads to improved net demand forecasts is the open research question that we address.

This work is positioned as a contribution from the Green Information Systems (Green-IS) body of research to the net zero electricity transition. Green-IS research advances the vital role of information systems and clean technologies in promoting environmentally sustainable practices (Sunar and Swaminathan 2022), aligned with the principles of eco-efficiency and eco-effectiveness

(Loock et al. 2013). Specifically, (Watson et al. 2010, Fig. 1) encourages the Green-IS research themes to include “energy informatics” to help reduce energy consumption and achieve sustainable energy solutions. In this paper, we contribute to this effort by helping retailers and SOs estimate unobservable information in demand behavior.

The key contributions of this paper are therefore summarized as follows:

- **Domain Contributions:** This paper makes a substantial contribution to characterizing the unobserved components of electricity demand response, being, as far as we are aware, the first to do so in the absence of device-level, real-time measurements. We accomplish this by developing an application of inverse optimization (IO) to this critical issue, a method previously used in other fields but not yet in retail electricity. Unlike conventional approaches that depend on direct observations of device-level data, this new method leverages net demand data to estimate unobservable flexibility parameters within a latent optimization model.

- **Experimental Results:** To validate the performance of the proposed methodology, we conduct a two-stage analysis using real-world data from previous studies. In the first stage, we utilize a dataset from Kaggle, which provides access to detailed, device-level, behavior of a US household with the aim of supporting multiple replication studies. Applying inverse optimization (IO) to this data, we find that the estimated flexible and inflexible components of demand align well with actual device behavior, confirming the plausibility of our approach. However, this dataset lacks Time-of-Use (TOU) pricing. To address this, we proceed to the second stage using a dataset from a study in Japan, which includes TOU prices but lacks observable device behavior. Building on the validation of our IO approach in the first stage, we then compare its fitting and forecasting capabilities on this data against benchmark models. The IO approach not only demonstrates superior empirical performance but also provides explicit estimates of the unobservable flexible and baseload demand components.

- **Managerial Implications:** This methodology allows retailers and system operators to gain a deeper understanding of behind-the-meter flexibility without compromising data privacy or requiring real-time device-level measurements. It can help them forecast better and enhance their incentive mechanisms, resulting in benefits for themselves and consumers through more efficient green energy management.

The rest of the paper is organized as follows. Section 2 reviews the related literature. In Section 3, we introduce the problem. In Section 4, we provide the proposed IO model and its tractable reformulation. Section 5 reports the empirical results and Section 6 concludes the paper.

2. Background research

Two research streams relate to our work. One is about the domain focus and the other is about the methodology. We discuss these two research streams, separately.

2.1. Estimating and predicting behind-the-meter demand

Various attempts to decompose "behind-the-meter" consumer demand into its inflexible baseload and flexible, price responsive, components have ranged from statistical and optimization techniques to machine learning approaches (Dobakhshari and Gupta 2018, Hatton et al. 2015, Lei et al. 2020, Sun et al. 2019, Wen et al. 2020, Zhang et al. 2015). The simplest method of estimating baseload demand using the average demand from periods without flexibility incentives, whilst appealing, has led to significant inaccuracies. More advanced statistical methods using linear or non-linear regression models with explanatory variables have resulted in weakly significant baseload estimates, with higher precision requiring large datasets. Regarding the flexible demand component, whilst its price response is confounded with other behavioral and weather-dependent factors, probabilistic methods including Bayesian inference and machine learning have provided useful insights into these uncertainties (Lin et al. 2021, Vallés et al. 2018, Mahdavi et al. 2022, Zhang et al. 2022). Nevertheless, all these methods encounter significant challenges, particularly in capturing temporal dependencies, handling nonlinear characteristics, and in context generalizations.

A fundamental issue with the regression methods, both statistical and machine learning, is that the training (estimation) relies on regressors rather than being driven by actual end-user's price responses. This has been emphasized by Chen et al. (2024), Hekmat et al. (2023). Moreover, various empirical studies have identified limitations associated with nonlinear effects, delayed responses and a general tendency for over-estimated price responses (An et al. 2015, 2016, Cappers et al. 2010, Kirschen et al. 2000). The methods prevalent in existing research struggle with these and other limitations, primarily due to rigid assumptions in the elasticity models (Ruan et al. 2021), the lack of interpretability in machine learning (Ponoćko and Milanović 2018), and the misalignment between demand baseline estimation and flexibility components (Vallés et al. 2018). For example, Ruan et al. (2021) introduce Siamese LSTM networks to estimate time-varying elasticities through a two-stage process. Similarly, Ponoćko and Milanović (2018) propose a two-step approach that first decomposes load using device-level data, followed by load forecasting with an artificial neural network. In another study, Vallés et al. (2018) utilize a quantile regression method to capture demand flexibility by incorporating exogenous features such as requests to adjust demand, incentives, ambient temperature, and time of day. However, this approach requires (i) prior knowledge of the

demand baseline (often predicted in the absence of price incentives) and (ii) the upward/downward consumer's response, which is typically unobservable. Furthermore, the two-step process—baseline estimation followed by quantile regression-based forecasting—can introduce biases in both net load forecasting and the estimation of baseline demand and flexibility estimates.

In consideration of the limitations in the existing research, an important feature of the IO approach, as developed in this paper, is that it integrates the forecasting and demand decomposition elements endogenously within the overall model, rather than as a sequence of separate tasks. This can thereby capture the complex interactions more robustly to provide more accurate predictions whilst maintaining both interpretability and uncertainty quantification. This appears to be beneficial compared to the conventional approach of modeling the elements separately as a sequence of demand estimation, flexibility quantification, forecasting and uncertainty assessment (Srivastava et al. 2024, Table II).

2.2. Inverse optimization

As its name implies, Inverse Optimization takes the results of a presumed optimization process to infer the parameters in the formulation of the decision-maker's optimal actions. Chan et al. (2023) provide a comprehensive review of its theory and applications. In practice, real data does not usually permit an exact derivation of the parameters, with noise in the observed results being due to measurement errors, bounded rationality, or model miss-specifications. Thus, the IO methodology in this case seeks to minimize the fitting errors between the observed data and the prescribed model results to make them “approximately” optimal (Aswani et al. 2018, Mohajerin Esfahani et al. 2018). Nevertheless, empirical evidence has shown that where a latent agent optimization model is plausible, the IO approach has greater predictive accuracy than machine learning and time-series models with limited data (Bian et al. 2023, Fernández-Blanco et al. 2021b, Saez-Gallego and Morales 2017). Power systems have been particularly fertile for the successful applications of IO (Saez-Gallego and Morales 2017).

Related to the demand response context of interest here, Bian et al. (2022, 2023) develop an IO approach using prior model knowledge and a gradient descent algorithm to determine the best-fitting model parameters based on the historical price and response data for forecasting price-responsive behaviors. Shi and Xu (2023) consider a deep learning method embedded into a bi-level model and solved by a gradient-descent method. The research by Fernández-Blanco et al. (2021a,b), Saez-Gallego et al. (2016), Saez-Gallego and Morales (2017) proposes a data-driven IO model for

forecasting purposes by reducing the inherent bi-level optimization in these applications to a single-level one. Then, by relaxing the complementary slackness conditions due to the computational hardness, they solved the problem by employing a heuristic method. Kovács (2021) apply IO to extract parameters from electricity consumer models by a quadratically constrained quadratic program which was solved using successive linear programming (SLP), but expressed concerns about the convergence properties of the SLP approach. Vatandoust et al. (2023) consider a data-driven IO approach to infer the parameters of a flexibility curve through aggregate forecasts. This enabled effective strategies in the Belgian imbalance market, outperforming XGBoost in capturing demand response behavior and thereby generating balancing energy profits.

We extend this body of work methodologically by developing a computationally efficient reformulation of the IO program to solve a conic mixed integer program. To our knowledge, this is the first study that simultaneously provides accurate demand forecasting as well as a decomposition into its unobservable components, without using direct device measurements. The proposed reformulation relies on the underlying conic-based structure, binary variables (to model the intrinsic combinatorial nature of load shifting) and kernel regression techniques, as advocated by Fernández-Blanco et al. (2021b), to minimize the out-of-sample forecasting errors. We avoid the need to apply a gradient descent method (Bian et al. 2022, 2023) being aware, in particular, that gradient-based algorithms have problems handling bi-level formulations with binary variables (Zhou et al. 2024). We also avoid the need to relax the complementary slackness conditions and apply a heuristic, as used by Fernández-Blanco et al. (2021a,b), Saez-Gallego et al. (2016), Saez-Gallego and Morales (2017).

Thus, the contribution of this work is to develop and apply a novel IO formulation that (i) uses net demand data to estimate for the first time the unobservable load components in an interpretable way and (ii) specifies the demand-response behavior of the consumers with a set of predictive functions to better understand potential flexibilities. Our proposed method is also able to provide the shiftability and sheddability estimates going beyond the approaches proposed elsewhere (Chen et al. 2024, Hekmat et al. 2023).

3. Model

In this section, we present the model. Throughout the paper, bold lower-case letters denote vectors, while standard lower-case letters are reserved for scalars. We use the shorthand $[T] = \{1, \dots, T\}$ to represent the set of all integers up to T . For reference, we also include all notations in Table 1.

Notation	Abbreviations
SO	System operator
TOU	Time-Of-Use
IO	Inverse optimization
FOP	Forward optimization program
BTM	Behind-the-meter
IS	Information Systems
GDPR	General Data Protection Regulation
SLP	Successive Linear Programming
Notation	Indices and Sets
$[T]$	Review period set, defined as $[T] = \{1, \dots, T\}$ made by T hours
t	Index of time period (hours) ($t \in [T]$)
$[S]$	Set of S days defined as $[S] = \{1, \dots, S\}$
s	Index of the daily samples ($s \in [S]$)
Θ	Flexibility feasibility set
Notation	Exogenous variables and parameters
$\hat{\mathbf{d}}_s$	Vector T -dimensional of load measurements on day s (kWh)
$\hat{\mathbf{g}}_s$	Vector of T -dimensional renewable energy generation on day s (kWh)
p_t	Price/flat tariff at time t (c/kWh, JPY/kWh)
$p_{s,t}^{sf,+}$	Per unit price incentive of energy of shifting by increasing at time t on day s (c/kWh, JPY/kWh)
$p_{s,t}^{sf,-}$	Per unit price incentive of energy of shifting by decreasing at time t on day s (c/kWh, JPY/kWh)
$p_{s,t}^{sd}$	Per unit price incentive of energy of shedding at time t on day s (c/kWh, JPY/kWh)
\mathbf{z}_s	Vector of T -dimensional exogenous observations build up by price on day s
$\hat{\mathbf{z}}_{s,t}$	Vector of external features at time t on day s (kWh)
α	Forgetting factor parameter
ω_s	Normalized forgetting weight on day s
$c_{s,t}^{sd}$	Comfort cost incurred by load shedding at time t on day s (c/kWh ² , JPY/kWh ²)
$c_{s,t}^{sf,+}$	Comfort cost incurred by load shifting by increasing demand at time t on day s (c/kWh ² , JPY/kWh ²)
$c_{s,t}^{sf,-}$	Comfort cost incurred by load shifting by decreasing demand at time t on day s (c/kWh ² , JPY/kWh ²)
\mathbf{c}_s	Vector of T -dimensional exogenous observations build up by cost on day s
\mathbf{K}_t	Vector of physical upper bounds of the flexibility envelopes at time t defined by (K_t^{flex}, K_t^{sd}) (kWh)
T^{\max}	Maximum number of shifted hours
ℓ	Utility function
Notation	Endogenous variables
d_t^{bl}	Baseload demand at time t (kWh)
$d_{s,t}^{sf,-}$	Downward shifted demand at time t on day s (kWh)
$d_{s,t}^{sf,+}$	Upward shifted demand at time t on day s (kWh)
$d_{s,t}^{sf}$	Shifted demand at time t on day s (kWh)
$d_{s,t}^{sd}$	Sheddable demand part of the net demand at time t on day s (kWh)
$\bar{d}_{s,t}^{sd}$	Sheddability envelope function at time t on day s (kWh)
$\bar{d}_{s,t}^{sf,+}, \bar{d}_{s,t}^{sf,-}$	Shiftability envelope functions at time t on day s (kWh)
$d_{s,t}^{sd,-}$	Sheddable demand at time t on day s (kWh)
d_t^{flex}	Flexible demand at time t on day s (kWh)
$\delta_{s,t}^{sf,+}, \delta_{s,t}^{sf,-}$	Auxiliary binary variables at time t on day s (kWh)
\mathbf{u}_s	Array of demand attributes $(d_t^{bl}, \bar{d}_{s,t}^{sf,+}, \bar{d}_{s,t}^{sf,-}, \bar{d}_{s,t}^{sd})$ on day s
θ_s	Array of decision variables of the s -th FOP on day s
θ_t	Array of decision variables $(d_t^{sf,+}, d_t^{sf,-}, d_t^{sd,-})$ of the consumer's problem at time t

Table 1 Notation



Figure 1 Decomposition of the observable load (by the SO) into its components: baseload demand d^{bl} , downward shifted demand $d^{sf,-}$, upward shifted demand $d^{sf,+}$ and sheddable demand d^{sd} .

The model considers the relationship between two parties: the consumer and the service provider. The service provider may be a retailer, aggregator, or network system operator, the distinction not being important for the general formulation. Therefore, for simplicity, we collectively refer to them as the system operator (SO). The SO is responsible for managing the consumer's transactions on the basis of net demand information, taken from the consumer's meter connection to the grid. We assume consumers aim to minimize their energy costs by adjusting consumption behind the meter, based on the prices they face, and they may be aided in this respect by algorithms or energy service companies. The activities behind the meter are unobserved by the SO. The specific details are outlined below:

3.1. Consumer

In its most general form, we assume that the consumer has baseload (inflexible) demand, some self-generation (such as rooftop solar PV panels), and some flexible demand, including storage or electric vehicle charging. The demand components are illustrated in Figure 1 and detailed below:

Net demand: The self-generation of the consumer at any time t is exogenous and random, denoted by \hat{g}_t . The electricity demand of the consumer is first served by this self-generation, if available. The rest of the demand, after taking the self-consumption and including both baseload and flexible components, is the net demand. The net demand is what the SO can observe and measure at the meter point.

Baseload demand: This component of the net demand refers to the electricity demand from the consumer which is unresponsive to price incentives. We denote this at time t , by d_t^{bl} . This has the connotation of being the necessary demand that the consumer requires at each point in time at any price below the retail market price cap, although in practice it is an empirical estimate of inflexibility across the range of prices as observed in the data.

Shiftable demand: The consumer can shift (advance or delay) some of its consumption when there are price incentives to do so. This can be done through adjusting the use of appliances such as washing machines, or interrupting the periodic requirements for heating and cooling, or making opportunistic use of battery storage and EV charging/discharging. We denote by d_t^{sf} , the shifted part of the net demand at any time t . Note that the shifted demand at any time t could be a positive or negative, respectively, denoted by $d_t^{sf,+}$ and $d_t^{sf,-}$, where $d_t^{sf} = d_t^{sf,+} - d_t^{sf,-}$, and $d_t^{sf,+} \geq 0$, $d_t^{sf,-} \geq 0$. The consumer receives linear price benefits $p_t^{sf,+}$, $p_t^{sf,-}$ per unit of energy shifted, in the positive and negative directions, respectively. However, the consumer also incurs a total cost of $c^{sf,+}(d_t^{sf,+})^2$ and $c^{sf,-}(d_t^{sf,-})^2$, respectively, for shifting in the increasing and decreasing directions at any time t . This cost would, for example, relate to the comfort or amenity loss to the consumer through this shift, and the quadratic assumption is intended to reflect that it is likely to be an increasing function of scale. Finally, we assume that the load shifting in the positive and negative directions are both upper bounded by maximum values, respectively, denoted by $\bar{d}_t^{sf,+}$ and $\bar{d}_t^{sf,-}$, and referred to as ‘shiftability envelopes’.

Sheddable demand: This part of the net demand, denoted by d_t^{sd} , refers to discretionary consumption that can be entirely or partially avoided when there is a price incentive. We denote by \bar{d}^{sd} , the maximum amount of sheddable demand at any time t and by $d_t^{sd,-}$ the amount of sheddable demand that is taken. Thus, denoting by d_t^{sd} the amount of sheddable demand at time t that is still required, we have by definition $d_t^{sd} = \bar{d}^{sd} - d_t^{sd,-}$. As with shifting, the consumer receives remuneration of p_t^{sd} per unit of energy shedding but also incurs a total cost of $c^{sd}(d_t^{sd,-})^2$ for shedding demand.

3.2. Service provider (SO)

The consumer needs to interact with an organization for energy delivery and billing. This organization, in different jurisdictions or applications, could take a different form such as an energy retailer, an aggregator, or a network operator. We abstract away from the specific nature of this organization for the purpose of this study and we call it simply a service provide (SO). What is shared among all of these distinct forms of ‘SOs’ is that they need to have an accurate understanding and characterization of the demand components of the consumer. To be more precise, the SO can observe the net demand, but not necessarily the building blocks of the demand (i.e., sheddable, shifted, and baseload demands). Figure 1 illustrates the ‘observable’ net demand by the SO, as well as the building blocks of the net demand that may be hidden from the SO. Gaining a deeper understanding of these building blocks of the net demand, without direct measurements, would benefit the SO in many ways as indicated in the Introduction. The open question is whether and how

	Perfectly given	Stats given	Objective	Decision variables
Consumer	c, z, u	$\hat{\mathbf{g}}, \hat{\mathbf{d}}$	Cost minimization	θ
SO	\mathbf{z}, \mathbf{K}_t	$\hat{\mathbf{g}}, \hat{\mathbf{d}}, \hat{\xi}_t$	Forecast error minimization	θ, \mathbf{u}

Table 2 Describing SO's and consumer's problems

the SO can hypothesize the building blocks of net consumption and subsequently use it for better prediction. In this work, we show how to achieve this by exploiting some exogenous information in the grid (weather information, contextual grid data, voltage, etc) in addition to the historical values of the observable net demand.

3.3. Consumer and SO problem formulations

Using the above definitions and setup, we proceed with the consumer's and SO's problem formulations. To do so, we also introduce three further notations. We define the array of 'demand attributes' as $\mathbf{u} := (d_t^{bl}, \bar{d}_t^{sf,+}, \bar{d}_t^{sf,-}, \bar{d}_t^{sd})_{t \in [T]}$ and the array of 'costs' as $\mathbf{c} = (c_t^{sf,+}, c_t^{sf,-}, c_t^{sd})_{t \in [T]}$ and the array of 'prices' as $\mathbf{z} = (p_t, p_t^{sf,+}, p_t^{sf,-}, p_t^{sd})_{t \in [T]}$.

3.3.1. Consumer's cost minimization We assume that the consumer aims to minimize its overall cost for the entire review period $[T]$. By review period, we mean the forward horizon over which it is seeking to optimize consumption. In this cost minimization, the consumer is assumed to know its cost parameters (\mathbf{c}), prices (\mathbf{z}), its flexibility potentials and baseload demand (\mathbf{u}), the historical values of its net demand ($\hat{\mathbf{d}}$), and its self-generation ($\hat{\mathbf{g}}$). With this information, the consumer decides how to activate its flexibility (i.e., shifting and shedding) to minimize the cost. We define $\theta := (\theta_t)_{t \in [T]} = (d_t^{sf,+}, d_t^{sf,-}, d_t^{sd,-})_{t \in [T]}$ as the vector that represents the flexibility decision variables of the consumer. With this notation, we can express the consumer problem as

$$\max_{\theta \in \Theta} \left(\ell(\theta) := \sum_{t \in [T]} [Q_t(\theta_t) + L_t(\theta_t)] \right) \quad (1)$$

where ℓ is consumer's utility function and Θ is the flexibility feasibility set both explained below.

The utility function ℓ defined in Eq. (1) is the cost comprising of a quadratic and a linear term, respectively, denoted by Q_t and L_t . The quadratic part represents the total comfort cost of demand shifting and shedding, $Q_t(\theta_t) := -c^{sf,+}(d_t^{sf,+})^2 - c^{sf,-}(d_t^{sf,-})^2 - c^{sd}(d_t^{sd,-})^2$. The linear part represents the financial exchanges with the SO, which is the difference between the demand

flexibility gain and the consumption cost, $L_t(\theta_t) := p_t^{sf,+} d_t^{sf,+} + p_t^{sf,-} d_t^{sf,-} + p_t^{sd} d_t^{sd,-} - p_t(d_t^{bl} + d_t^{sf,+} - d_t^{sf,-} + \bar{d}_t^{sd} - d_t^{sd,-} - \hat{g}_t)$.

To understand the feasible region in which the consumer chooses its optimal flexibility decision, we define the *flexibility feasibility set*, Θ , which comprises the set of possible actions of the consumer for its flexible (shifted and sheddable) behavior. Essentially, Θ incorporates the physical upper bounds on the negative shiftability, positive shiftability, and sheddability, respectively, denoted by $\bar{d}_t^{sf,+}$, $\bar{d}_t^{sf,-}$, \bar{d}_t^{sd} , which we collectively refer to generally as *flexibility envelopes* (or resp. as *shiftability and sheddability envelopes*). These envelopes represent the physical restrictions on the maximum level of their respective flexibilities at time t . We, therefore, represent the flexibility feasibility set Θ as $\Theta(\bar{d}_t^{sf,+}, \bar{d}_t^{sf,-}, \bar{d}_t^{sd})$ to make explicit the dependence with respect to the shiftability and sheddability envelopes. Given the envelopes, the feasibility set is expressed as follows:

$$\Theta := \begin{cases} 0 \leq d_t^{sd,-} \leq \bar{d}_t^{sd} & \forall t \in [T] \\ 0 \leq d_t^{sf,+} \leq \bar{d}_t^{sf,+} \delta_t^{sf,+} & \forall t \in [T] \\ 0 \leq d_t^{sf,-} \leq \bar{d}_t^{sf,-} \delta_t^{sf,-} & \forall t \in [T] \\ \sum_{t \in [T]} d_t^{sf,+} = \sum_{t \in [T]} d_t^{sf,-} & \\ \delta_t^{sf,+} + \delta_t^{sf,-} \leq 1, & \forall t \in [T] \\ \sum_{t \in [T]} (\delta_t^{sf,+} + \delta_t^{sf,-}) \leq T^{\max} & \\ \delta_t^{sf,+}, \delta_t^{sf,-} \in \{0, 1\}, & \forall t \in [T] \end{cases} \quad (2)$$

The first three lines in the feasibility set described above ensure that the shedding, positive shifting, and negative shifting, all satisfy their respective upper envelopes. In those lines, we introduce auxiliary binary variables $\delta_t^{sf,+} \in \{0, 1\}$ and $\delta_t^{sf,-} \in \{0, 1\}$, $t \in [T]$ to reflect practical restrictions on demand shifting. The fourth line ensures that the overall demand shifting is energy neutral and does not carry over to the next review cycle; i.e., any shifting in the positive/negative direction must be reversed before the review period is over. The fifth and sixth lines, respectively, ensure that positive and negative shifting does not occur simultaneously and that the total number of shifted hours is less than or equal to T^{\max} . Note that $T^{\max} \leq T$ represents the time flexibility of demand shifting. Large values of T^{\max} imply a larger time that the shifted demand can be spread over and hence, more time flexibility.

3.3.2. SO's forecast error minimization The ultimate goal of the SO is to minimize the (out-of-sample) forecast error for net demand and, in the process, to estimate the consumer's unobserved

demand components, including the flexibility responses/decisions, θ , and baseload demand and flexibility potentials, \mathbf{u} . We assume that the SO knows the array of price signals \mathbf{z} and physical upper bounds of the flexibility potentials, \mathbf{K}_t . It is generally the case that the SO has information about self generation $\hat{\mathbf{g}}$ (as these renewable facilities are usually metered separately for subsidies and/or regulatory requirements), past observed net demand measurements $\hat{\mathbf{d}}$ and various relevant weather/grid data $\hat{\xi}_t$ for the entire review period $[T]$. Table 2 summarizes the problem setting.

4. Solution using inverse optimization (IO)

In this section, we first introduce the IO methodology and then explain how it can be further developed for unobservable demand response.

4.1. Theoretical background on inverse optimization

In the terminology of Inverse Optimization (IO), a *forward* optimization program (FOP) is proposed whose solutions are observed. This FOP is parametrized on some exogenous observations, \mathbf{w} , and some endogenous parameters, \mathbf{u} . Hence the formulation in a deterministic setting can be expressed as:

$$\text{FOP}(\mathbf{w}|\mathbf{u}) := \begin{cases} \max & \ell(\theta; (\mathbf{w}|\mathbf{u})) \\ \text{s.t.} & \theta \in \Theta(\mathbf{u}) \end{cases}, \quad (3)$$

where for any value of \mathbf{u} , we assume that $\Theta(\mathbf{u})$ is a closed convex set and ℓ a differentiable and concave utility function with respect to θ .

In a non-deterministic setting (such as our problem), the IO solution must consider the observation noise, such that the IO model becomes a supervised learning problem with multivariate output. More formally, therefore, the training dataset, \mathcal{D} , is built up by S pairs of observations $\mathcal{D}_S := \{(\mathbf{w}_s, \theta_s)\}_{s \in [S]}$ and the goal is to find $\hat{\mathbf{u}}$ that would make θ_s “optimal” for the s -th FOP. That is, $\text{FOP}(\mathbf{z}_s|\hat{\mathbf{u}})$. Thus, having estimated the parameter \mathbf{u} by $\hat{\mathbf{u}}$, and given a new observation \mathbf{w}' , $\text{FOP}(\mathbf{w}'|\hat{\mathbf{u}})$ is used to estimate θ' (Saez-Gallego and Morales 2017).

4.2. Applying IO to unobservable demand response

The FOP model is naturally delivered by the consumer’s problem as defined in Eq. (1). Adapting the basic approach as described above, we proceed in detail, as follows:

4.2.1. The forward or reconstruction program Given the training dataset $\mathcal{D}_S := \{((\mathbf{c}_s, \mathbf{z}_s, \hat{\mathbf{g}}_s), \hat{\mathbf{d}}_s)\}_{s \in [S]}$ made by S training samples (which may correspond to a set of days), built on pairs of T -dimensional exogenous observations (cost, price and self generation data), $(\mathbf{c}_s, \mathbf{z}_s, \hat{\mathbf{g}}_s)$, and observed load measurements $\hat{\mathbf{d}}_s$. In this setting, $(\mathbf{c}_s, \mathbf{z}_s, \hat{\mathbf{g}}_s)$ plays the role of \mathbf{w} in program (3)

.Thus, the data-driven IO model tries to find estimates of demand attributes $\hat{\mathbf{u}}_s$ of the parameters

$\mathbf{u}_s = (d_t^{bl}, \bar{d}_{s,t}^{sf,+}, \bar{d}_{s,t}^{sf,-}, \bar{d}_{s,t}^{sd})$ which would make the observed s -th demand $\hat{\mathbf{d}}_s$ “optimal” for the s -th

FOP, $\text{FOP}((\mathbf{c}_s, \mathbf{z}_s, \hat{\mathbf{g}}_s) | \mathbf{u}_s)$, defined as follows:

$$\begin{aligned}
 & \max \quad \ell(\boldsymbol{\theta}_s; ((\mathbf{c}_s, \mathbf{z}_s, \hat{\mathbf{g}}_s) | \mathbf{u}_s)) \\
 & \text{s.t. } \boldsymbol{\theta}_s \in \Theta(\bar{d}_{s,t}^{sf,+}, \bar{d}_{s,t}^{sf,-}, \bar{d}_{s,t}^{sd}) \\
 & \quad \boldsymbol{\theta}_s = (\mathbf{d}_s^{sf,+}, \mathbf{d}_s^{sf,-}, \mathbf{d}_s^{sd,-}) \\
 & \quad d_{s,t}^{sf} = d_{s,t}^{sf,+} - d_{s,t}^{sf,-} \quad \forall t \in [T] \\
 & \quad d_{s,t}^{sd} = \bar{d}_{s,t}^{sd} - d_{s,t}^{sd,-} \quad \forall t \in [T]
 \end{aligned} \tag{4}$$

Having introduced the end-users’ FOP, next we explain how the proposed framework works:

(i) a *learning step* where the data-driven IO model is fitted by using a training dataset, and (ii) a *forecasting step* where the time-series forecast output is computed.

4.2.2. Learning Having observed a time series given by set of $[S]$ days of the demand per time

period t , $\{(\hat{d}_{s,t})_{t \in [T]}\}_{s \in [S]}$, which form the training dataset, the learning step is given by the solution

of the following data-driven inverse optimization program:

$$\min \sum_{s \in [S]} \omega_s \left\| \left([\mathbf{d}^{bl} + \mathbf{d}_s^{sf} + \mathbf{d}_s^{sd} - \hat{\mathbf{g}}_s] - \hat{\mathbf{d}}_s \right)_s \right\|_p^p \quad (5a)$$

$$\text{s.t. } d_{s,t}^{sf} = d_{s,t}^{sf,+} - d_{s,t}^{sf,-}, \forall s \in [S], \forall t \in [T] \quad (5b)$$

$$d_{s,t}^{sd} = \bar{d}_{s,t}^{sd} - d_{s,t}^{sd,-}, \forall s \in [S], \forall t \in [T] \quad (5c)$$

$$0 \leq d_t^{bl}, \forall t \in [T] \quad (5d)$$

$$0 \leq \bar{d}_{s,t}^{sf,+}, \forall s \in [S], \forall t \in [T] \quad (5e)$$

$$0 \leq \bar{d}_{s,t}^{sf,-}, \forall s \in [S], \forall t \in [T] \quad (5f)$$

$$\bar{d}_{s,t}^{sf,+} + \bar{d}_{s,t}^{sf,-} + \bar{d}_{s,t}^{sd} \leq K_{s,t}^{flex}, \forall s \in [S], \forall t \in [T] \quad (5g)$$

$$0 \leq \bar{d}_{s,t}^{sd} \leq K_{s,t}^{sd}, \forall s \in [S], \forall t \in [T] \quad (5h)$$

$$\delta_{s,t}^{sf,+}, \delta_{s,t}^{sf,-} \in \{0, 1\}, \forall s \in [S], \forall t \in [T] \quad (5i)$$

$$\delta_{s,t}^{sf,+} + \delta_{s,t}^{sf,-} \leq 1, \forall s \in [S], \forall t \in [T] \quad (5j)$$

$$\sum_{t \in [T]} (\delta_{s,t}^{sf,+} + \delta_{s,t}^{sf,-}) \leq T^{\max}, \forall s \in [S] \quad (5k)$$

$$\boldsymbol{\theta}_s := (\mathbf{d}^{sf,+}, \mathbf{d}^{sf,-}, \mathbf{d}^{sd,-})_s, \forall s \in [S] \quad (5l)$$

$$\mathbf{u}_s = (d_t^{bl}, \bar{d}_{s,t}^{sf,+}, \bar{d}_{s,t}^{sf,-}, \bar{d}_{s,t}^{sd})_s, \forall s \in [S] \quad (5m)$$

$$\boldsymbol{\theta}_s \in \arg \max \ell(\boldsymbol{\theta}_s; ((\mathbf{c}_s, \mathbf{z}_s, \hat{\mathbf{g}}_s) | \mathbf{u}_s)) \quad (5n)$$

$$\text{s.t. } \sum_{t \in [T]} d_{s,t}^{sf,+} = \sum_{t \in [T]} d_{s,t}^{sf,-} : (\kappa_s) \quad (5o)$$

$$d_{s,t}^{sf,+} \leq \bar{d}_{s,t}^{sf,+} \delta_{s,t}^{sf,+} : (\mu_{s,t}^+ \geq 0), \forall t \in [T] \quad (5p)$$

$$d_{s,t}^{sf,-} \leq \bar{d}_{s,t}^{sf,-} \delta_{s,t}^{sf,-} : (\mu_{s,t}^- \geq 0), \forall t \in [T] \quad (5q)$$

$$d_{s,t}^{sd,-} \leq \bar{d}_{s,t}^{sd} : (\mu_{s,t}^0 \geq 0), \forall t \in [T] \quad (5r)$$

$$0 \leq d_{s,t}^{sf,+} : (\nu_{s,t}^+ \geq 0), \forall t \in [T] \quad (5s)$$

$$0 \leq d_{s,t}^{sf,-} : (\nu_{s,t}^- \geq 0), \forall t \in [T] \quad (5t)$$

$$0 \leq d_{s,t}^{sd,-} : (\nu_{s,t}^0 \geq 0), \forall t \in [T] \quad (5u)$$

where $\|\cdot\|_p$, with $p > 1$, denotes the p -norm, ω_s is the weight of the forecast error on day $s \in [S]$ (Saez-Gallego et al. 2016) and the variables in parenthesis in the constraints of the s -th lower level

problem denote the respective dual variables. We set $\omega_s := \left(\frac{s}{S}\right)^\alpha / \left(\sum_{s' \in [S]} \left(\frac{s'}{S}\right)^\alpha\right)$ the normalized weights of the training days, with $\alpha \geq 0$ the forgetting factor parameter to represent discounting older data. If $\alpha = 0$, then all observations are equally weighted and when α increases, more weight is given to recent observations, as suggested by Saez-Gallego et al. (2016), Lu et al. (2018).

The flexibility envelopes $\bar{d}_{s,t}^{sf,+}$, $\bar{d}_{s,t}^{sf,-}$ and $\bar{d}_{s,t}^{sd}$, which are upper bounded by the input parameters $K_{s,t}^{sf,+}$, $K_{s,t}^{sf,-}$ and $K_{s,t}^{sd}$, respectively (in the numerical experiments are equal to the absolute value of the observed demand, $|\widehat{d}_{s,t}|$), are random functions themselves because they may evolve over time from day to day, and hence, should be estimated from past data. To do so, we follow the approach introduced by Fernández-Blanco et al. (2021b), which proposes a decision rule based on a Gaussian Kernel function to capture the non-linear dependence between regressors and power:

$$\bar{d}_{s,t}^{sf,+} = \beta_0^{sf,+} + \sum_{s' \in [S], t' \in [T]} \beta_{s',t'}^{sf,+} e^{-\gamma^{sf,+} \|\hat{\xi}_{s,t} - \hat{\xi}_{s',t'}\|} \quad (6)$$

$$\bar{d}_{s,t}^{sf,-} = \beta_0^{sf,-} + \sum_{s' \in [S], t' \in [T]} \beta_{s',t'}^{sf,-} e^{-\gamma^{sf,-} \|\hat{\xi}_{s,t} - \hat{\xi}_{s',t'}\|} \quad (7)$$

$$\bar{d}_{s,t}^{sd} = \beta_0^{sd} + \sum_{s' \in [S], t' \in [T]} \beta_{s',t'}^{sd} e^{-\gamma^{sd} \|\hat{\xi}_{s,t} - \hat{\xi}_{s',t'}\|} \quad (8)$$

where $\gamma^{sf,+}, \gamma^{sf,-}, \gamma^{sd} > 0$ are the Gaussian kernel bandwidth parameter associated to $\bar{d}_{s,t}^{sf,+}$, $\bar{d}_{s,t}^{sf,-}$ and $\bar{d}_{s,t}^{sd}$ respectively. Moreover, $\hat{\xi}_s$ is the vector of external features (for example weather information or contextual grid data, etc) corresponding to day s at hour t of the inverse optimization program used to estimate the maximum availability of flexible demand at time t on day s . In order to mitigate the over-fitting effect of the flexibility envelopes, we add the term $\lambda \|\beta^{sf,+}, \beta^{sf,-}, \beta^{sd}\|_p^p$ in (5a) by using a regularization parameter $\lambda > 0$. Note that the kernel regression parameters are determined by the IO program, as they are treated as decision variables within the upper-level model.

4.3. Component estimation and forecasting

Having introduced the parameterizations in terms of the kernels and the decision variables for the flexibility envelopes, we compute $\widehat{\mathbf{d}}^{bl}$ and $\widehat{\beta}$ being the optimal solutions of \mathbf{d}^{bl} and $\beta :=$

$[\beta_0^{sf,+}, \beta_0^{sf,-}, \beta_0^{sd}, (\beta_s^{sf,+}, \beta_s^{sf,-}, \beta_s^{sd})_{s \in [S]}]$ in the problem defined by (5a)-(5u), respectively. Thus, we can compute $\hat{\mathbf{u}} := (\hat{\mathbf{d}}^{bl}, \bar{d}_{S+1,t}^{sf,+}, \bar{d}_{S+1,t}^{sf,-}, \bar{d}_{S+1,t}^{sd})$ where

$$\bar{d}_{S+1,t}^{sf,+} = \left(\hat{\beta}_0^{sf,+} + \sum_{\substack{s' \in [S] \\ t' \in [T]}} \hat{\beta}_{s',t'}^{sf,+} e^{-\gamma^{sf,+} \|\hat{\xi}_{S+1,t} - \hat{\xi}_{s',t'}\|} \right)_+ \quad (9)$$

$$\bar{d}_{S+1,t}^{sf,-} = \left(\hat{\beta}_0^{sf,-} + \sum_{\substack{s' \in [S] \\ t' \in [T]}} \hat{\beta}_{s',t'}^{sf,-} e^{-\gamma^{sf,-} \|\hat{\xi}_{S+1,t} - \hat{\xi}_{s',t'}\|} \right)_+ \quad (10)$$

$$\bar{d}_{S+1,t}^{sd} = \left(\hat{\beta}_0^{sd} + \sum_{\substack{s' \in [S] \\ t' \in [T]}} \hat{\beta}_{s',t'}^{sd} e^{-\gamma^{sd} \|\hat{\xi}_{S+1,t} - \hat{\xi}_{s',t'}\|} \right)_+ \quad (11)$$

and $(\cdot)_+$ stands for the positive part. Having computed $\hat{\mathbf{u}}$, the forecast output for the day $S+1$ is given by solving the following forward program given the exogenous observations for the day $S+1$, that is cost, price and local generation data on day $S+1$, $(\mathbf{c}_{S+1}, \mathbf{z}_{S+1}, \hat{\mathbf{g}}_{S+1})$:

$$\begin{aligned} & \max \ell(\theta_{S+1}; ((\mathbf{c}_{S+1}, \mathbf{z}_{S+1}, \hat{\mathbf{g}}_{S+1}) | \hat{\mathbf{u}})) \\ & \text{s.t. } d_{S+1,t}^{sf} = d_{S+1,t}^{sf,+} - d_{S+1,t}^{sf,-} \quad \forall t \in [T] \\ & \quad d_{S+1,k,t}^{sd} = \bar{d}_{S+1,t}^{sd} - d_{S+1,t}^{sd,-} \quad \forall t \in [T] \\ & \quad \theta_{S+1} = (\mathbf{d}^{sf,+}, \mathbf{d}^{sf,-}, \mathbf{d}^{sd,-})_{S+1} \\ & \quad \theta_{S+1} \in \Theta(\bar{d}_{S+1,t}^{sf,+}, \bar{d}_{S+1,t}^{sf,-}, \bar{d}_{S+1,t}^{sd}) \end{aligned}$$

After solving the above program and getting $(\mathbf{d}^{sf,+}, \mathbf{d}^{sf,-}, \mathbf{d}^{sd,-})_{S+1}$, we can compute the forecast net load per time t for day $S+1$, $\hat{d}_{S+1,t}$, as the sum of the (net) baseload demand, $\hat{d}_t^{bl} - \hat{g}_{S+1,t}$, plus the flexible demand, $d_{S+1,t}^{flex} := d_{S+1,t}^{sf} + d_{S+1,t}^{sd}$ as follows:

$$\hat{d}_{S+1,t} := (\hat{d}_t^{bl} - \hat{g}_{S+1,t}) + d_{S+1,t}^{flex} \quad (12)$$

The Supplementary material develops a single-level tractable reformulation by exploiting the Karush-Kuhn-Tucker (KKT) conditions for the lower-level and the underlying conic structure of the data-driven IO problem.

5. Empirical research

Using an adaptation of the IO methodology (details and notation in the Methods Supplement), the logical flow of the empirical research questions that we address are as follows:

Q1: Using net demand with structural assumptions and a presumed optimization model for the flexible and baseload components of domestic demand and their price responses, how well is the above IO formulation able to estimate the latent components of flexible and baseload demand? This will establish proof of concept and computational tractability.

Q2: Can the resulting latent component estimates from IO be accurately predicted out of sample, using market-wide information such as weather and periodic effects, but without assuming knowledge of the appliance activities? The predictability of the components needs to be established against benchmark statistical and machine learning methods.

Q3: Are forecasts of net demand obtained by summing the latent components more accurate than forecasts of the aggregate net demand series? This is needed to validate the value of the decomposed approach for forecasting.

Q4: With data that provides explicit knowledge of the domestic appliance activities, do these latent components from IO have a credible relationship to the appliance activities and thereby appear valid? Plausible statistical relationships to the data will establish construct validity.

Q5: Overall, can the price response of domestic customers be more accurately estimated based upon an IO decomposition of flexible and baseload components than by conventional benchmark methods? This is needed to justify the overall proposition that IO can be useful in practice compared to benchmark forecasting methods.

We progress the analysis of Q1-Q5 through two data sets for clarity of focus. For the first part, Q1 to Q4, we use a set of smart household data, complied by Singh (2024) and openly available on the Kaggle platform to facilitate replication studies. It contains the smart meter readings in kW from a US household of all of the main household appliances at minute-by-minute granularity for 350 days of house appliances, together with the weather conditions of that particular region. It does not contain any retail price data. We do however know the typical retail prices for that region (Ohio) in 2018. This example replicates the perspective of a network operator who would not know the retailer-customer tariff details precisely. We use this to test the application of the IO approach and sense-check the identified components of flexible and inflexible components against actual appliance use. The forecasting questions are also evaluated on this data, but for Q5, which relies upon precise TOU responses, we use consumer data from Japan, made available by Kiguchi et al.

(2021). This is intraday net load data from a sample of customers responding to price signals for several months in 2018. This example therefore replicates the other case of a retailer setting TOU tariffs, observing the metered net loads but not knowing exactly how the appliances were being used. The logic of our analysis is, therefore, that having established the construct validity of the component estimates on the first data set, we are justified in applying it in the second data set in order to compare its forecasting performance under TOU prices against conventional benchmarks.

The IO method is benchmarked against three methods: a linear regression model (Linear), a time-series model with exogenous variables (SARIMAX), and also against a non-linear, machine-learning method, XGboost (Chen and Guestrin 2016). The performance metrics used for evaluations of the forecasts are the usual mean absolute error (MAE) and the root mean square error (RMSE) defined on the Methods section.

5.1. Performance effectiveness and construct validity

The smart household data from Singh (2024) is used to test the empirical performance of the proposed IO method against conventional benchmarks. It also allows us to check that the identified flexible and baseload components do in fact represent what they are intended to model, i.e., that they have construct validity. We do that by regressing the estimated flexible and baseload component activities on the known usage of household appliances as well as autoregressive lags. We expect the discretionary use of appliances to feature more strongly in the flexible component but less so in the baseload component, whilst the autoregressive effects would relate to adaptation and habit and should feature more strongly in the baseload component. The data measurements include the kitchen appliances (e.g. fridge, dishwasher, microwave), various rooms (e.g. kitchen, living room, home office), furnace, well, barn, garage door, etc., as well as solar generation and weather conditions. The minute-by-minute data is from a sample of US households in Ohio taken over 2018 (as also analyzed in Tlenshiyeva et al. (2024)). We restrict the analysis to weekdays. Recall that the net load is consumed energy minus the solar generation.

Even though we do not know the precise tariff details for the customer, we expect the IO method to be robust to noisy assumptions and we introduce a typical tariff for that location at that time (based upon inspection of published rates from various local utilities) available at Appendix A.2. The training data uses two months, weekdays only, and the goal is to predict the first five days of the next month. We consider two seasonal sets of results.

We show the summary results in Tables 3-4. Evidently, the IO is more accurate than the benchmark methods for forecasting. Figure 2 shows the net demand (demand is referred to as load in the Kaggle

data source) for the two cases (upper January to March, and lower July to September) under the optimal hyperparameter tuning via grid-search. Similarly, Figures 3 and 4 show the net baseload demand and flexible profiles estimated as part of the modeling.

In all figures, the peak price periods are shown shaded in pink. Whilst the effect of peak pricing is not evident in the net load profiles shown in Figure 2, nor in the baseload components shown in Figure 3, it is very evident in Figure 4. In Figure 4, we can see that only the IO approach is showing a distinct reduction in the flexible load component during the hours of peak pricing. This is, of course, what we would hope to see, as the IO approach is designed to identify this component. Visually, it is a very reassuring display that the IO is indeed providing a new component that the benchmark methods cannot reveal. It is also interesting to see in Figure 4 that during the Summer, the flexible load reduction is compensated by prior increases in load, whereas in the winter it is compensated with load increases afterwards. This seasonal distinction in the pattern of demand shifting is a new insight revealed by the methodology.

Regarding relative accuracy, according to Tables 3, for Winter, IO demonstrates the lowest MAE (0.1937) and RMSE (0.2589) for the “Net load” category, outperforming all the benchmark methods. XGBoost and SARIMAX have higher MAE and RMSE values in comparison to IO but still outperform Linear in the same category. Linear shows the highest errors across most categories. Tables 4 shows that the results for summer exhibit slight shift. While IO maintained the lowest MAE (0.1571) and RMSE (0.2226) for “Net load”, Linear and XGBoost performed comparably but with higher errors than IO and SARIMAX slightly better across all categories. Linear again shows the highest MAE and RMSE for “Net load”.

According to these figures, IO demonstrates the most accurate performance across all seasons, consistently achieving the lowest errors in the “Net load”. SARIMAX and Linear generally follow. Surprisingly, XGBoost, despite having some strong performances in the baseload, d^{bl} , and the flexible load, d^{flex} , often has the highest errors for “Net load” metrics. Most importantly for our purposes, from Figures 3 and 4 it is evident that the baseload component is more regular and predictable than the flexible component, which is intuitive and reassuring. Overall, these results appear to convincingly support research questions, Q1 ad Q2, in establishing the feasibility and comparative value of the IO approach.

In relation to Q3, which focuses upon the inherent modeling element of predicting the components separately, and then summing to get the Net load forecast, Tables 3-4 also report some justification. The rows designated Net Load and Net Load 2 display the accuracies for predicting on the basis only

Table 3 Error metrics for Winter

	IO		XGboost		SARIMAX	
	MAE	RMSE	MAE	RMSE	MAE	RMSE
Net load	0.1937	0.2589	0.3854	0.4632	0.2974	0.3740
d^{bl}	-	-	0.0951	0.1533	0.0746	0.0998
d^{flex}	-	-	0.1298	0.1455	0.0522	0.0695
Net load 2	0.1937	0.2589	0.2196	0.2707	0.2097	0.2793

Table 4 Error metrics for Summer

	IO		XGboost		SARIMAX	
	MAE	RMSE	MAE	RMSE	MAE	RMSE
Net load	0.1571	0.2226	0.3290	0.4090	0.2734	0.3187
d^{bl}	-	-	0.1002	0.1603	0.1115	0.1333
d^{flex}	-	-	0.0882	0.1030	0.1014	0.1123
Net load 2	0.1571	0.2226	0.2226	0.2719	0.2088	0.2574

of the aggregate Net Load data versus the disaggregated approach of predicting the components and summing.

Across the two seasons, the methods show varying degrees of improvement from “Net load” to “Net load 2.” XGBoost verified substantial improvements, with reductions in MAE ranging from approximately 43.0% in Winter to 32.4% in Summer, and reductions in RMSE from approximately 41.6% in Winter to 33.5% in Summer. Linear also shows notable improvements, with MAE reductions of about 24.7% in Winter and 18.9% in Summer, and RMSE reductions of about 27.1% in Winter and 20.2% in Summer. SARIMAX exhibited consistent improvements, with MAE reductions from 29.5% in Winter to 23.7% in Summer, and RMSE reductions from 25.3% in Winter to 19.2% in Summer. Overall, all methods generally performed better with “Net load 2,” with XGBoost and SARIMAX showing the most significant improvements. In other words, the principle of component-based forecasting appears to be justified.

Finally, whilst the error metrics for the point forecasts validate the relative accuracy of the IO approach for forecasting, it is becoming increasingly important in demand forecasting to understand the tail risks in predictions. This involves producing density forecasts and assessing the calibrations

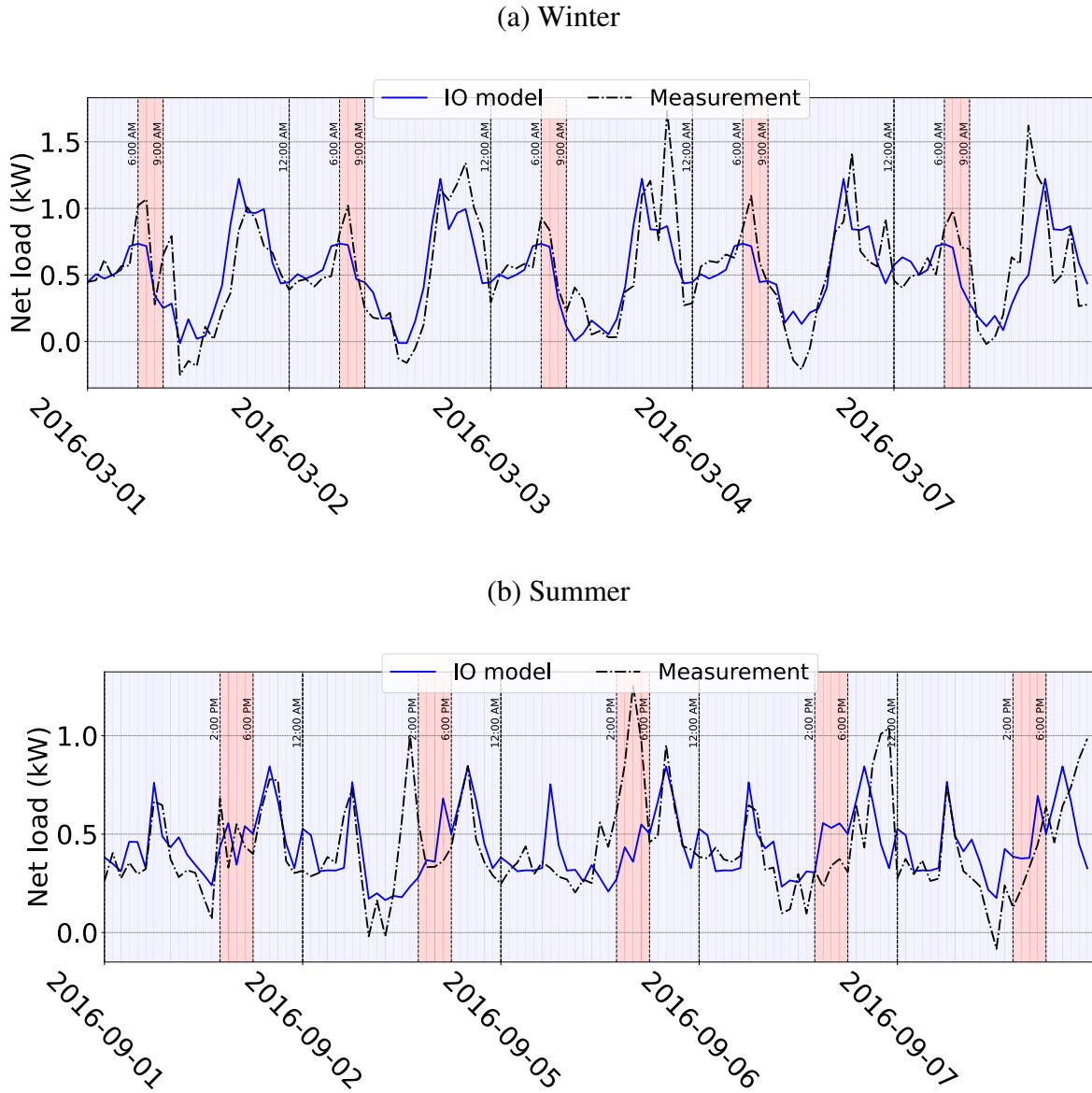


Figure 2 Net load profiles and performance error metrics under best hyperparameter tuning

of the predicted quantiles. For this purpose, we report the conventional continuous rank probability scores (CRPS) averaged over the forecast horizon T^{fo} .

Table 5 contains the 0.05, 0.95-quantiles, $q_{0.05}, q_{0.95}$, and the CRPS (averaged over the forecast horizon). Note that lower scores are better and we observe that the IO method outperforms very substantially across all seasons and Linear and SARIMAX generally have intermediate performance. Overall, the superior fitting and forecasting accuracy of the IO method for Net load appears to be well established against conventional benchmarks, in terms of both average accuracy and predictive density calibration. This established the feasibility and overall merit of the approach. The next

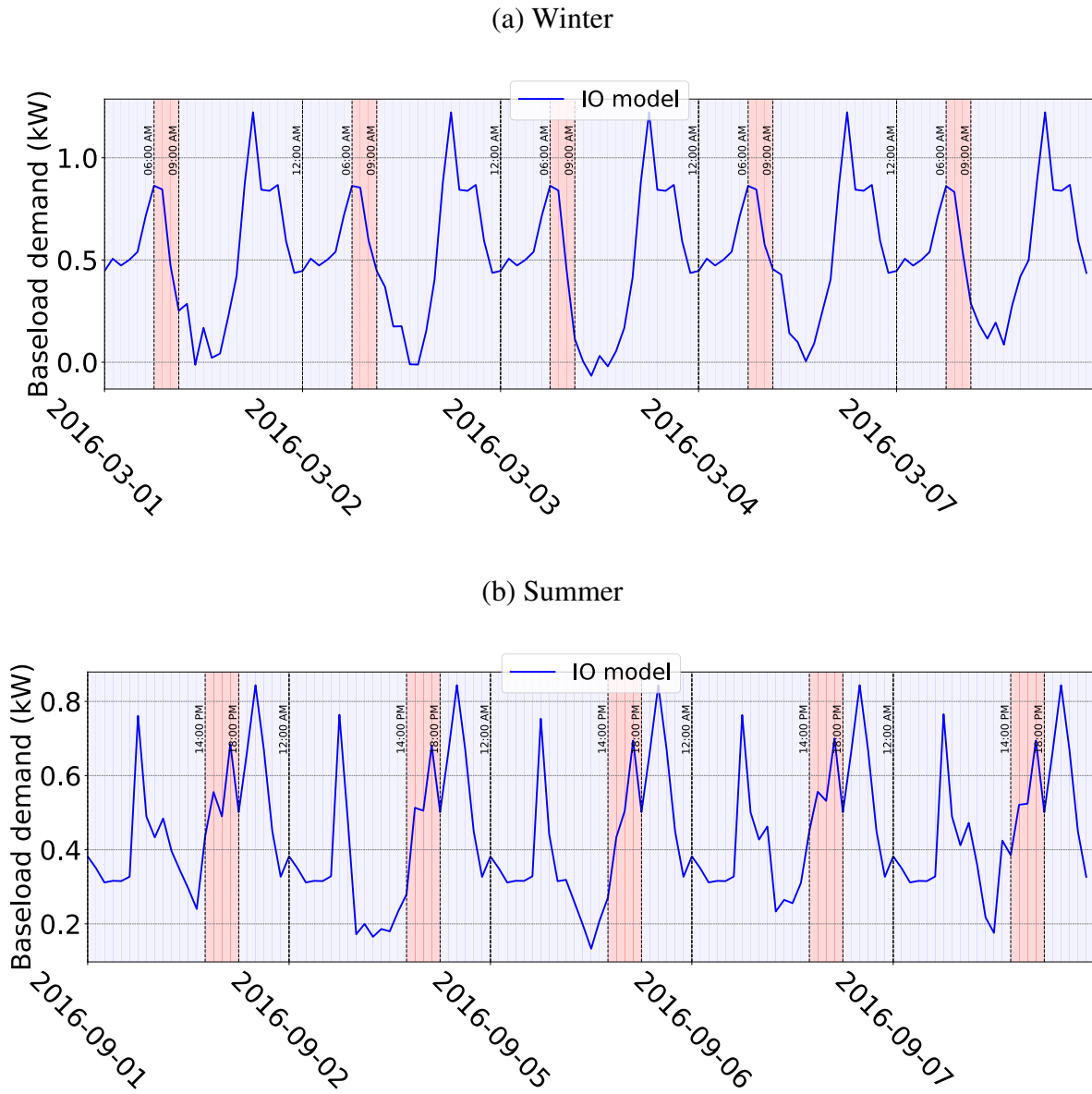


Figure 3 Net baseload demand profiles and performance error metrics under best hyperparameter tuning

Table 5 Error metrics for the net load probabilistic forecasting: Quantiles and CRPS

	IO			XGboost			SARIMAX		
	$q_{0.05}$	$q_{0.95}$	CRPS	$q_{0.05}$	$q_{0.95}$	CRPS	$q_{0.05}$	$q_{0.95}$	CRPS
Winter	0.370	1.012	0.212	0.647	1.114	0.355	0.532	1.105	0.301
Summer	0.263	0.805	0.139	0.388	1.047	0.246	0.379	1.002	0.180

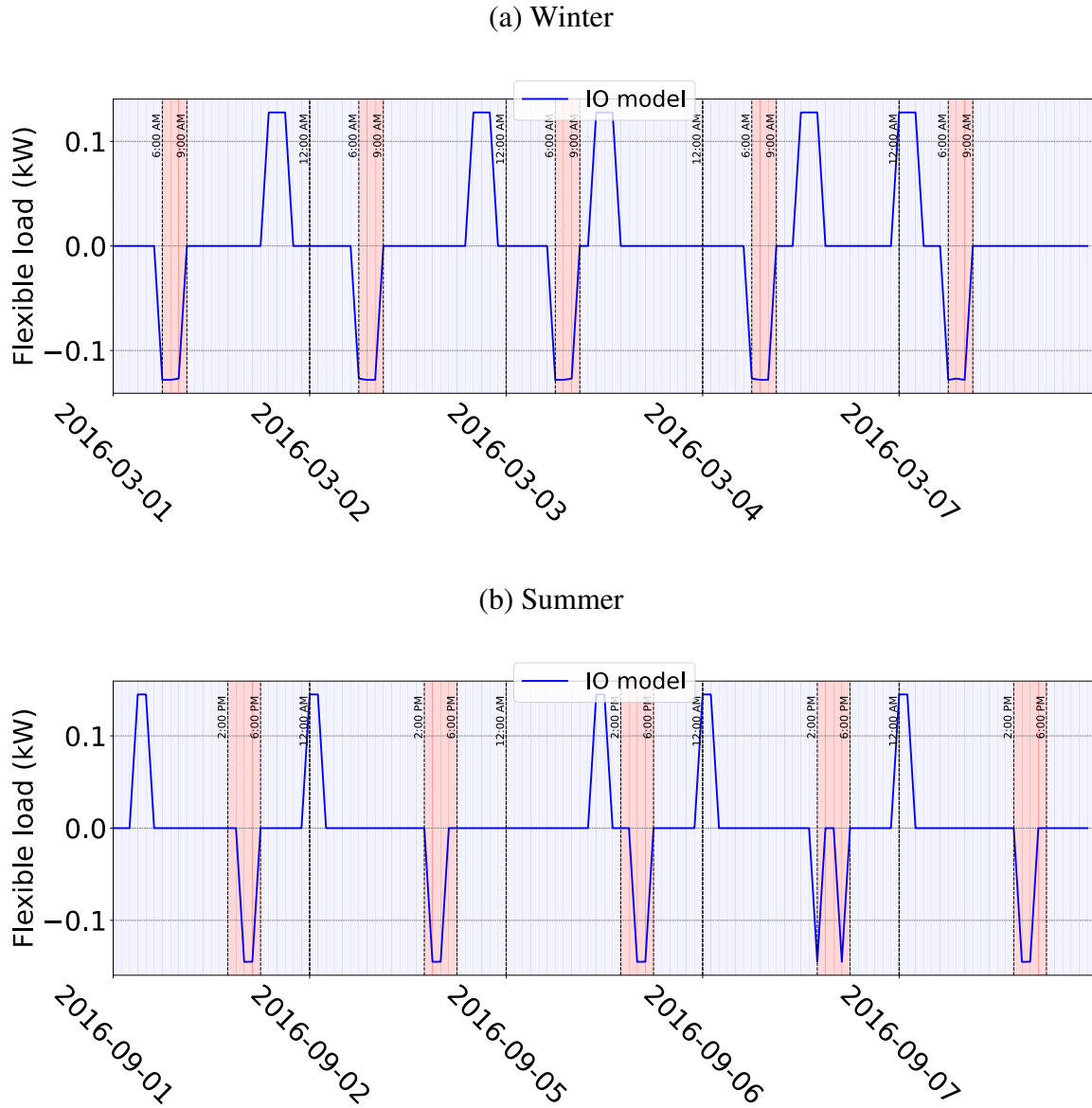


Figure 4 Flexible load profiles and performance error metrics under best hyperparameter tuning

question (Q4) is whether the estimated components really do represent what they are intended to represent, i.e., their construct validity.

To assess the construct validity of the flexible and baseload components delivered by the IO method, we regress each estimated component time series on the known metered appliance activities to assess their relative impact on the component's variation. We use the conventional SHAP (SHapley Additive exPlanations, Lundberg and Lee (2017)) values to represent the relative feature importance of the appliance activities. SHAP is a technique for analyzing the explainability of machine learning models by providing a way to understand how a specific prediction is arrived

at. In this way, we can see if the IO method provides meaningful load decompositions. Figures 5 and 6 display the relative feature strength of each appliance on the baseload and flexible components. The features are ranked by their importance. Figure 5 reveals that the baseload component primarily depends on habitual consumption patterns rather than discretionary appliance usage. This is evidenced by the dominance of autoregressive lags across both seasons, especially lag 1, which exhibits a strong positive relationship with baseload consumption. However, a key seasonal difference emerges: while lag 1 consistently increases baseload consumption in both seasons, the impact of other lagged terms is reversed (lags 6 and 2 in winter and lags 4, 5 and 6 in summer), revealing that higher values contribute to lower baseload consumption. Additionally, during the summer season, the "Living room" feature plays a distinct role, indicating that higher consumption levels in this area lead to increased baseload demand. This suggests a potential correlation with sustained cooling or ventilation use. For the flexible load component, a contrasting pattern is observed. Figure 6 highlights the predominance of specific appliances, such as the furnace and the dishwasher, over autoregressive effects except in summer. This confirms that flexible load is primarily driven by discretionary appliance usage. Unlike the baseload component, autoregressive terms are less influential, further reinforcing the distinction between habitual and flexible consumption. Moreover, Figure 6 systematically demonstrates that in both seasons, variations in appliance usage and certain lagged values translate directly into corresponding changes in flexible consumption. Notably, an autoregressive shift occurs between lags 2 and 3 in winter and summer respectively, suggesting seasonal variations in flexibility dynamics. Moreover, the feature importance of the Barn during winter is remarkable and drops in summer and hence, suggesting a potential correlation of "Barn" activities like holding the "Barn" items during winter season. These findings support the claim that the IO results effectively decompose total load into meaningful baseload and flexible components.

For additional diagnostic metrics, in comparison to the mainly visual insights from the SHAP displays, standard multiple regression results are summarised in a Regression Diagnostics Supplement, revealing the significance levels of the factors and the R-squared fits of 0.855 for winter and 0.847 for summer. Also, we report the respective residual QQ-plots for both seasons showing the normality of the residuals.

Thus, in concluding this section, we observe that, with confidence in the approach, we can now turn to the second set of consumer response data where there is a well defined TOU tariff, but the details of their responses are unobserved.

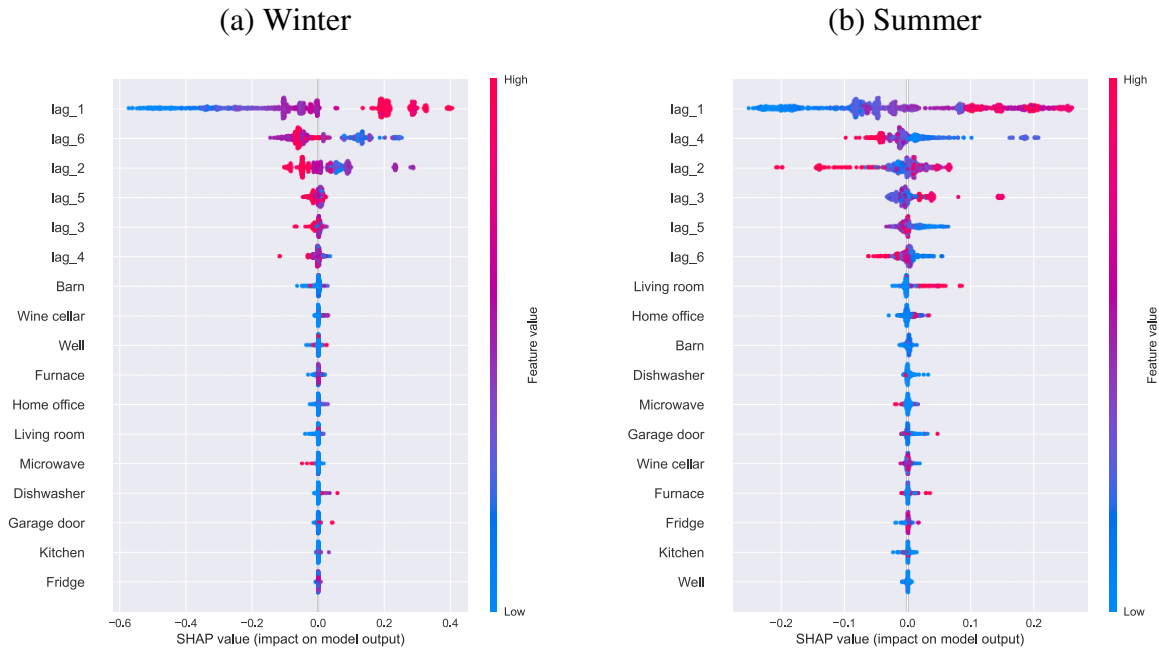


Figure 5 Time series' appliances impact into net baseload demand IO estimates via SHAP values

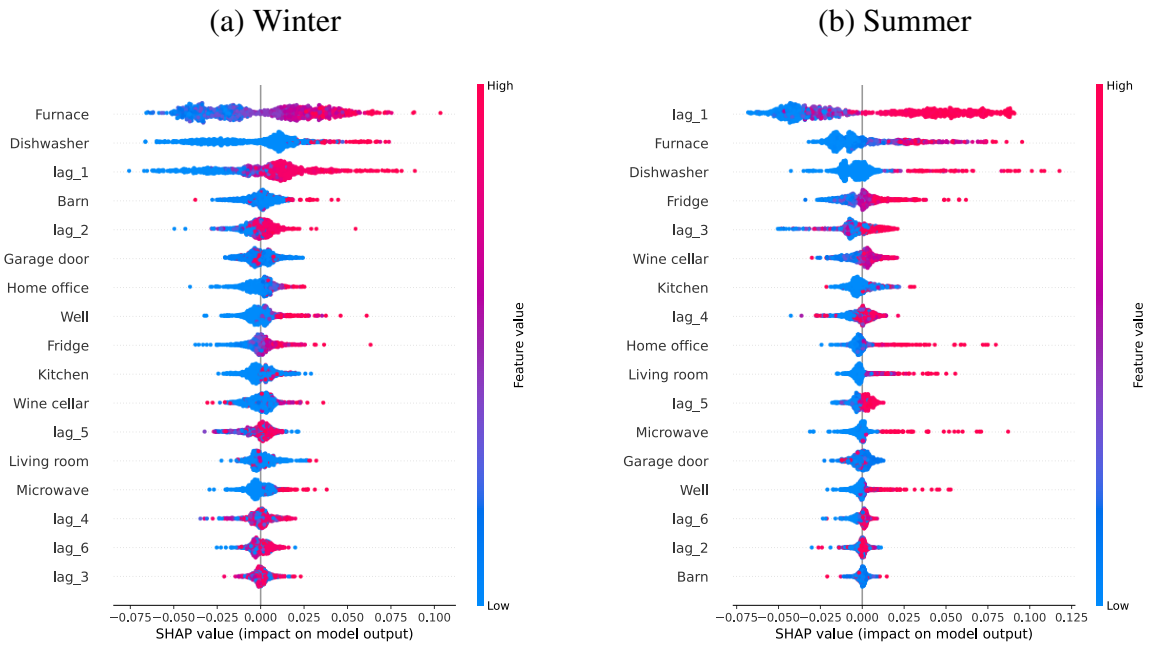


Figure 6 Time series' appliances impact into flexible IO estimates via SHAP values

5.2. Predictability with estimated price responsive components

This example is more common in practice where the TOU electricity tariffs are known and the net loads of the consumer at the meter are observable, but, for privacy, technical or agency reasons, the

Table 6 Error metrics

	IO		XGboost		SARIMAX	
	MAE	RMSE	MAE	RMSE	MAE	RMSE
C137	0.018	0.023	0.022	0.028	0.020	0.024
C162	0.041	0.057	0.053	0.068	0.047	0.062

Table 7 Error metrics for the net load probabilistic forecasting for the conscious demand-response program: Quantiles and CRPS

	IO			XGboost			SARIMAX		
	$q_{0.05}$	$q_{0.95}$	CRPS	$q_{0.05}$	$q_{0.95}$	CRPS	$q_{0.05}$	$q_{0.95}$	CRPS
C137	0.046	0.087	0.017	0.041	0.219	0.020	0.036	0.129	0.017
C162	0.128	0.287	0.039	0.112	0.299	0.039	0.129	0.284	0.042

the actual customer's appliance behavior is not observable. The dataset is taken from a Japanese utility as reported by Kiguchi et al. (2021) and the tariff details are available at Appendix A.2. We consider the weekdays corresponding to the period from July to September 2018 and the goal is to predict the last week of September. Exogenous weather variables are the temperatures and ground solar irradiance values at Tokyo and the day-ahead market price, from the Japan Electric Power Exchange (JPEX). We selected two customers, being the customers 137 and 162 available by Kiguchi et al. (2021).

Figures 8 and 9 show the time series of the load components estimated by the IO method for customers 137 and 162. As before with the Kaggle data, the peak price time periods are highlighted in pink in the figures. Again it is reassuring to see that the demand drops during the peak prices in the flexible load profiles but not in the baseload, and that demand shifting is evident with the flexible loads increasing at the end of those periods. The baseload components are very regular in their periodicities. The ability to identify these series from the unobserved consumer demand responses is the most remarkable aspect of this methodology. From the previous results, we also expect greater forecasting accuracy as a consequence.

Thus, Table 6 shows the performance metrics for the four methods under consideration. It is evident that the IO approach yields better point forecasts than the other methods. To be more precise, for customer C137, XGBoost shows approximately 22.2% higher MAE and 21.7% higher RMSE

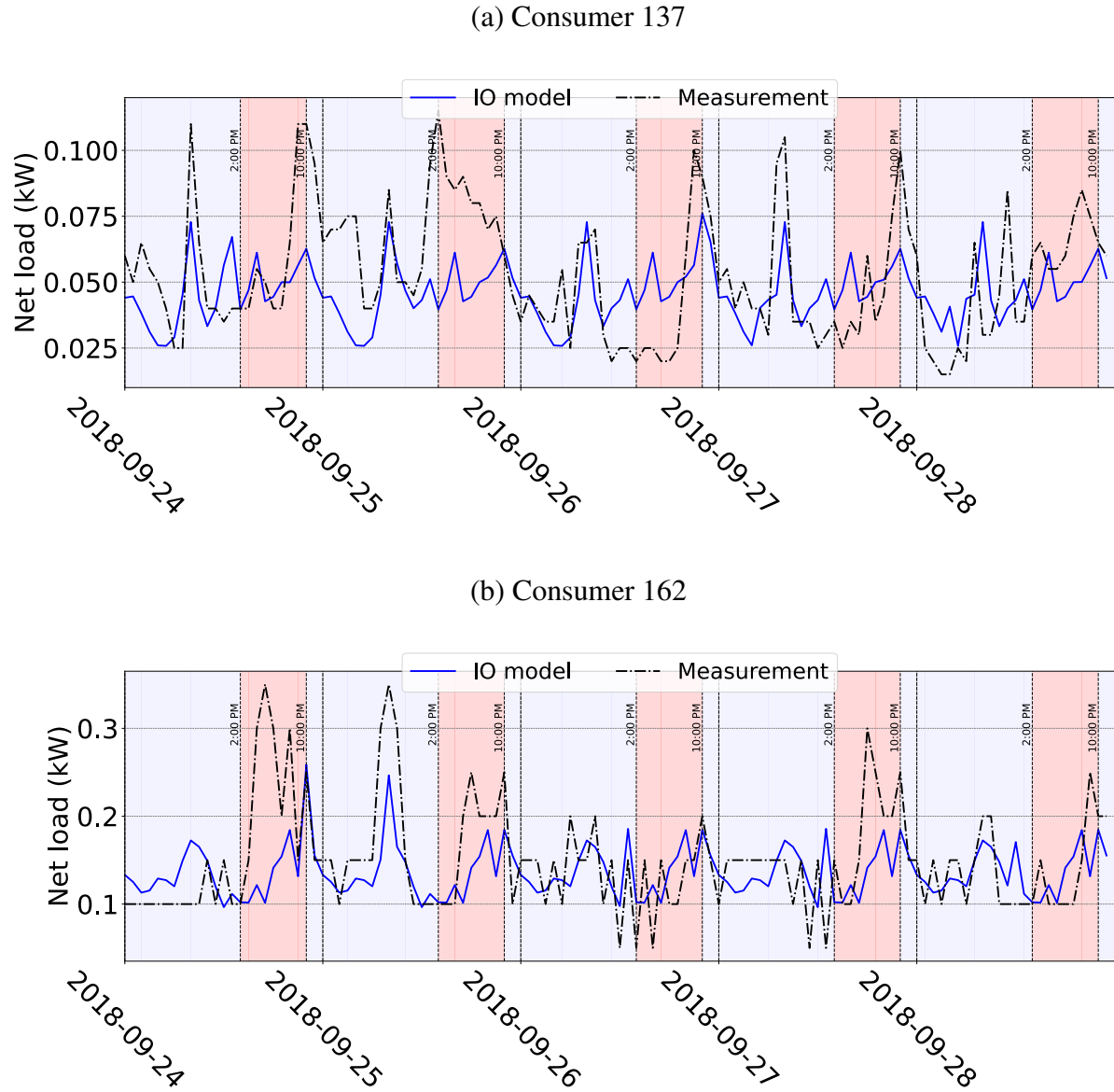


Figure 7 Net Load profiles and performance error metrics under best hyperparameter tuning for the conscious demand-response program case study

compared to IO method. Linear exhibits about 16.7% higher MAE and 21.7% higher RMSE, while SARIMAX shows approximately 11.1% higher MAE and 4.3% higher RMSE than IO method. Similarly, for customer C162, XGBoost demonstrates around 29.3% higher MAE and 19.3% higher RMSE, Linear shows about 24.4% higher MAE and 5.3% higher RMSE, and SARIMAX exhibits approximately 14.6% higher MAE and 8.8% higher RMSE compared to IO method.

As with the previous study, a probabilistic forecasting analysis was also conducted to assess the quality of density forecasting. The data in Table 7 shows the quantiles and CRPS (averaged over the

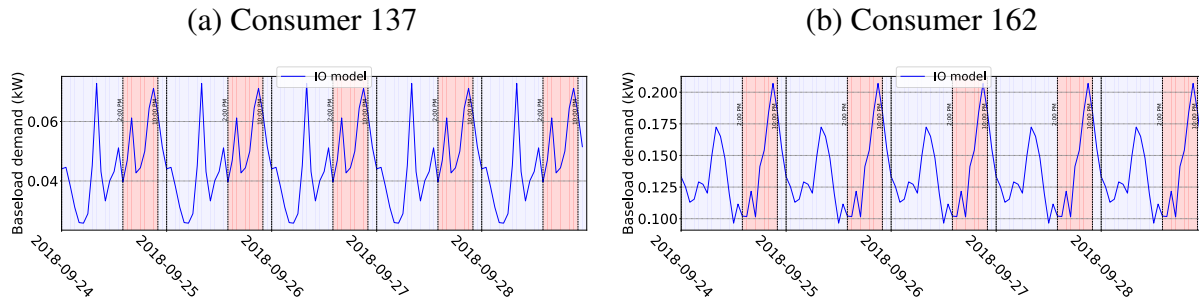


Figure 8 Baseload demand profile estimates

forecasting horizon). XGBoost consistently reports higher errors across metrics for both customers. For Customer C137, XGBoost has quantile errors with $q_{0.05}$ 10.9% lower and $q_{0.95}$ 151.7% higher compared to IO. Linear has $q_{0.05}$ 2.2% higher and $q_{0.95}$ 13.8% higher than IO, while SARIMAX has $q_{0.05}$ 21.7% lower and $q_{0.95}$ 48.3% higher than IO. In terms of CRPS, XGBoost performs 17.6% worse than IO, Linear 11.8% worse, and SARIMAX is on par with IO for this specific customer, C137. Regarding Customer C162, XGBoost exhibits quantile errors 12.5% lower for $q_{0.05}$ and 4.2% higher for $q_{0.95}$ compared to IO. Linear, on the other hand, has $q_{0.05}$ 24.2% higher and $q_{0.95}$ 20.9% lower than IO. SARIMAX has quantile errors 0.8% higher for $q_{0.05}$ and 1.0% lower for $q_{0.95}$. In terms of CRPS, XGBoost matches IO, Linear performs 38.5% worse, and SARIMAX 7.7% worse for this customer, C162.

In summary, while XGBoost consistently delivers higher errors across deterministic and probabilistic metrics, the IO method proves to be more accurate in load forecasting and uncertainty quantification. SARIMAX performs similarly to IO in probabilistic forecasting goals and a slightly better than Linear. However, IO delivers a better interpretability of demand flexibility while SARIMAX just captures some autoregressive and seasonal behavior during the summer period. Thus, in the more realistic setting provided by this dataset, the IO approach appears to be more accurate and more transparent in its revelations of the unobserved components.

6. Conclusions

In this paper, we have addressed an emerging information management challenge in retail electricity. Accurately characterizing the components of consumer electricity demand response is critical for system operators and energy retailers, but achieving this requires usually real-time, device-level data. However, such data is often unavailable or restricted due to data protection policies. Instead, all that can be observed is the net demand at the meters connecting customers to the grid, whilst

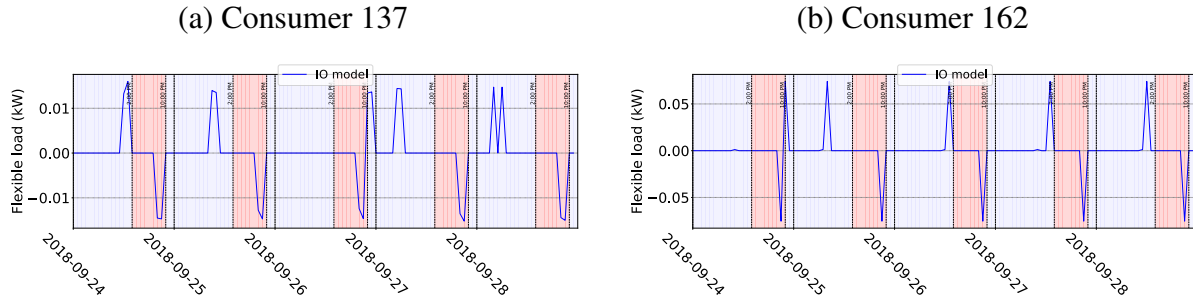


Figure 9 Flexible Load profile estimates

the “behind-the-meter” demand components remain unobservable. As a result, alternative methods which can infer or impute this unobservable information have substantial value. Current approaches fall short in this regard. To tackle this challenge, we developed and tested a data-driven inverse optimization (IO) methodology that effectively addresses this methodological gap.

Unlike conventional methods, the IO approach does not depend on direct observation of device-level data or behind-the-meter activities. Instead, it utilizes net demand data to estimate parameters within a latent optimization model, enabling the decomposition of consumption patterns into both observable and unobservable components. This includes net demand influenced by consumer behavior, as well as underlying baseload and flexible elements, such as shifted and sheddable loads. The inverse optimization (IO) methodology effectively captures the nonlinear relationships between price responses and external factors, providing a more nuanced and accurate understanding of consumer electricity demand response.

Having developed a solution methodology for the specific IO formulation, we conducted a comprehensive experimental study to assess the outcome in two different settings. In the first setting, we evaluated the empirical efficacy of the proposed IO methodology against established benchmarks by assessing of the construct validity of the identified flexible and baseload components. By regressing the estimated activities of these components on the known usage patterns of household appliances and autoregressive lags, we showed that discretionary appliance usage is prevalent in the flexible component, while autoregressive influences are more prominent in the baseload component, reflecting adaptive and habitual behaviors. In the second setting, through a precise TOU responses of consumer data from Japan, we showed that the price response of domestic customers can be more accurately estimated using an IO decomposition of flexible and baseload components compared to conventional benchmark methods. This finding supports the proposition that IO can be practically useful compared to traditional forecasting methods without jeopardizing interpretability and forecast accuracy.

Notes

Appendix A: Computational details

Computations were performed on a Windows-based laptop with 4 cores and 8 logical processors clocking at 2.6 GHz and 16 GB of RAM using the solver MOSEK under Pyomo (Bynum et al. 2021). The IO model depends on five hyperparameters, namely: $param = [T^{\max}, \alpha, \gamma^{sf,+}, \gamma^{sf,-}, \gamma^{sd}, \lambda]$, which were tuned via a grid-search. For the SARIMAX, we choose the parameters that minimize the AIC (the choice (2, 1, 2, 6) seems reasonably well in all cases). XGboost is applied by using the skforecast Python package (Amat Rodrigo and Escobar Ortiz 2024) via default options with 6 lags.

A.1. Performance metrics: MAE, RMSE and CRPS

$$MAE(d^{true}, d^{fo}) = \frac{1}{|T^{fo}|} \sum_{t \in T^{fo}} |d_t^{true} - d_t^{fo}| \quad (13)$$

$$RMSE(d^{true}, d^{fo}) = \sqrt{\frac{1}{|T^{fo}|} \sum_{t \in T^{fo}} (d_t^{true} - d_t^{fo})^2} \quad (14)$$

where d^{true}, d^{fo} are the load measurements and forecast values, and T^{fo} is the planning horizon set (in our case, 5 days of 24 hours). We assume that on every day at hour t the external regressors are perfectly known. This is often assumed in forecasting research in order to focus more directly upon the comparative methodologies so that the results are not confounded by estimation errors in variables.

CRPS is a strictly proper scoring rule which is widely used within density forecasting research. It needs to be approximated in practice, and one approach is based on exploiting its equivalent quantile-based representation Gneiting and Ranjan (2011):

$$CRPS(F, y) = 2 \int_0^1 QL_q(F^{-1}(q), y) dq \quad (15)$$

where F, F^{-1} are the distribution and quantile function, respectively, y is the measurement, and QL is the quantile loss function. Computing the proxy for CRPS requires the q -quantiles, where $q \in \{0.01, 0.05, \dots, 0.95, 0.99\}$. We used the forecasting method with the pinball loss function (also known as quantile score function) available in (Amat Rodrigo and Escobar Ortiz 2024).

A.2. Consumer peak and off-peak pricing in the case studies

Ohio Case Study. We estimated a typical flat tariff of 22 c/kWh and prices of 29 c/kWh and 15 c/kWh for the peak and off-peak periods, with the peak defined as 2 – 6 pm in summer and 6 – 9 am in winter. We also define $p_t^{sf,+} = \max(p_t - TOU_t, 0)$, $p_t^{sf,-} = \max(TOU_t - p_t, 0)$. The parameters for load shedding are defined by $\sum_{t \in [T]} p_t^{sf,+} / 24$ and the cost c^{sd} is defined by p^{sd} . The parameters $c_t^{sf,+}, c_t^{sf,-}$ are set to $|p_t - TOU_t|$ when $p_t^{sf,+}, p_t^{sf,-}$ are equal to zero, respectively, and 0, otherwise. We take as regressors the temperature difference, $Tdiff := T - T^{app}$, the dew point and the exponential function of the solar generation for the upward, downward and sheddable envelopes, respectively.

Tokyo Case Study. We estimated a typical flat tariff of 22 c/kWh and prices of 29 c/kWh and 15 c/kWh for the peak and off-peak periods, with the peak defined as 2 – 6 pm in summer and 6 – 9 am in winter. We also define $p_t^{sf,+} = \max(p_t - \text{TOU}_t, 0)$, $p_t^{sf,-} = \max(\text{TOU}_t - p_t, 0)$. The parameters for load shedding are defined by $\sum_{t \in [T]} p_t^{sf,+} / 24$ and the cost c^{sd} is defined by p^{sd} . The parameters $c_t^{sf,+}$, $c_t^{sf,-}$ are set to $|p_t - \text{TOU}_t|$ when $p_t^{sf,+}$, $p_t^{sf,-}$ are equal to zero, respectively, and 0, otherwise. -For p_t in the IO model the flat tariff was 26 JPY/kWh and we defined $p_t^{sf,+} = \max(p_t - \text{TOU}_t, 0)$, $p_t^{sf,-} = \max(\text{TOU}_t - p_t, 0)$, where $\text{TOU}_t = 35$ JPY/kWh between 2pm and 10pm, and 20 JPY/kWh otherwise. The price incentive for load shedding, p^{sd} , is defined by $\sum_{t \in [T]} p_t^{sf,+} / 24$ and the cost c^{sd} is given by the price incentive p^{sd} . The parameters $c_t^{sf,+}$, $c_t^{sf,-}$ are set to $|p_t - \text{TOU}_t|$ when $p_t^{sf,+}$, $p_t^{sf,-}$ are equal to zero, respectively, and 0, otherwise.

References

- Amat Rodrigo J, Escobar Ortiz J (2024) skforecast. URL <http://dx.doi.org/10.5281/zenodo.8382788>.
- An J, Kumar P, Xie L (2015) On transfer function modeling of price responsive demand: An empirical study. *2015 IEEE power & energy society general meeting*, 1–5 (IEEE).
- An J, Kumar P, Xie L (2016) Dynamic modeling of price responsive demand in real-time electricity market: Empirical analysis. *arXiv preprint arXiv:1612.05021*.
- Aswani A, Shen ZJ, Siddiq A (2018) Inverse optimization with noisy data. *Operations Research* 66(3):870–892.
- Bian Y, Zheng N, Zheng Y, Xu B, Shi Y (2022) Demand response model identification and behavior forecast with optnet: a gradient-based approach. *Proceedings of the Thirteenth ACM International Conference on Future Energy Systems*, 418–429, e-Energy '22 (New York, NY, USA: Association for Computing Machinery), ISBN 9781450393973, URL <http://dx.doi.org/10.1145/3538637.3538871>.
- Bian Y, Zheng N, Zheng Y, Xu B, Shi Y (2023) Predicting strategic energy storage behaviors. *arXiv preprint arXiv:2306.11872*.
- Bynum ML, Hackebeil GA, Hart WE, Laird CD, Nicholson BL, Sirola JD, Watson JP, Woodruff DL (2021) *Pyomo—optimization modeling in python*, volume 67 (Springer Science & Business Media), third edition.
- Cappers P, Goldman C, Kathan D (2010) Demand response in US electricity markets: Empirical evidence. *Energy* 35(4):1526–1535.
- Chan TC, Mahmood R, Zhu IY (2023) Inverse optimization: Theory and applications. *Operations Research*.
- Chen L, Zhu X, Xu B, Ding F (2024) Demand side flexibility envelope quantification under data scarcity. *2024 IEEE Power & Energy Society Innovative Smart Grid Technologies Conference (ISGT)*, 1–5 (IEEE).
- Chen T, Guestrin C (2016) Xgboost: A scalable tree boosting system. *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, 785–794.
- Choi DG, Lim MK, Murali K, Thomas VM (2020) Why have voluntary time-of-use tariffs fallen short in the residential sector? *Production and Operations Management* 29(3):617–642.
- Danti P, Magnani S (2017) Effects of the load forecasts mismatch on the optimized schedule of a real small-size smart prosumer. *Energy Procedia* 126:406–413.

- Dobakhshari DG, Gupta V (2018) A contract design approach for phantom demand response. *IEEE Transactions on Automatic Control* 64(5):1974–1988.
- Dong C, Ng CT, Cheng T (2017) Electricity time-of-use tariff with stochastic demand. *Production and operations management* 26(1):64–79.
- European Parliament and Council (2016) Regulation (eu) 2016/679 of the european parliament and of the council. Official Journal of the European Union, URL <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A02016R0679-20160504>, accessed: 2024-07-29.
- Fernández-Blanco R, Morales JM, Pineda S (2021a) Forecasting the price-response of a pool of buildings via homothetic inverse optimization. *Applied Energy* 290:116791.
- Fernández-Blanco R, Morales JM, Pineda S, Porras Á (2021b) Inverse optimization with Kernel regression: Application to the power forecasting and bidding of a fleet of electric vehicles. *Computers & Operations Research* 134:105405.
- Gneiting T, Ranjan R (2011) Comparing density forecasts using threshold- and quantile-weighted scoring rules. *Journal of Business & Economic Statistics* 29(3):411–422, URL <http://dx.doi.org/10.1198/jbes.2010.08110>.
- Hatton L, Charpentier P, Matzner-Løber E (2015) Statistical estimation of the residential baseline. *IEEE Transactions on Power Systems* 31(3):1752–1759.
- Hekmat N, Cai H, Zufferey T, Hug G, Heer P (2023) Data-driven demand-side flexibility quantification: Prediction and approximation of flexibility envelopes. *2023 IEEE Belgrade PowerTech*, 1–6 (IEEE).
- Karthikeyan SP, Raglend IJ, Kothari DP (2013) A review on market power in deregulated electricity market. *International Journal of Electrical Power & Energy Systems* 48:139–147.
- Kiguchi Y, Weeks M, Arakawa R (2021) Predicting winners and losers under time-of-use tariffs using smart meter data. *Energy* 236:121438.
- Kirschen DS, Strbac G, Cumperayot P, de Paiva Mendes D (2000) Factoring the elasticity of demand in electricity prices. *IEEE Transactions on Power Systems* 15(2):612–617.
- Koliou E, Bartusch C, Picciariello A, Eklund T, Söder L, Hakvoort RA (2015) Quantifying distribution-system operators’ economic incentives to promote residential demand response. *Utilities Policy* 35:28–40.
- Kovács A (2021) Inverse optimization approach to the identification of electricity consumer models. *Central European Journal of Operations Research* 29(2):521–537.
- Lei S, Hong D, Mathieu JL, Hiskens IA (2020) Baseline estimation of commercial building HVAC fan power using tensor completion. *Electric Power Systems Research* 189:106624.
- Lin J, Ma J, Zhu J (2021) A privacy-preserving federated learning method for probabilistic community-level behind-the-meter solar generation disaggregation. *IEEE Transactions on Smart Grid* 13(1):268–279.
- Liu J, Xiao Y, Li S, Liang W, Chen CP (2012) Cyber security and privacy issues in smart grids. *IEEE Communications surveys & tutorials* 14(4):981–997.

- Loock CM, Staake T, Thiesse F (2013) Motivating energy-efficient behavior with green is: an investigation of goal setting and the role of defaults. *MIS quarterly* 1313–1332.
- Lu T, Wang Z, Wang J, Ai Q, Wang C (2018) A data-driven stackelberg market strategy for demand response-enabled distribution systems. *IEEE Transactions on Smart Grid* 10(3):2345–2357.
- Lundberg SM, Lee SI (2017) A unified approach to interpreting model predictions. *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 4768–4777, NIPS’17 (Red Hook, NY, USA: Curran Associates Inc.), ISBN 9781510860964.
- Mahdavi N, Weeraddana D, Guo Y (2022) Probabilistic estimation of pv generation at customer and distribution feeder levels using net-demand data. *IEEE Transactions on Smart Grid* 14(3):1974–1984.
- Mathieu JL, Verbič G, Morstyn T, Almassalkhi M, Baker K, Braslavsky J, Bruninx K, Dvorkin Y, Ledva GS, Mahdavi N, et al. (2024) A new definition of demand response in the distributed energy resource era. *arXiv preprint arXiv:2410.18768*.
- Mohajerin Esfahani P, Shafieezadeh-Abadeh S, Hanasusanto GA, Kuhn D (2018) Data-driven inverse optimization with imperfect information. *Mathematical Programming* 167:191–234.
- Nghiem TX, Jones CN (2017) Data-driven demand response modeling and control of buildings with gaussian processes. *2017 American Control Conference (ACC)*, 2919–2924 (IEEE).
- Pinson P, Madsen H, et al. (2014) Benefits and challenges of electrical demand response: A critical review. *Renewable and Sustainable Energy Reviews* 39:686–699.
- Ponočko J, Milanović JV (2018) Forecasting demand flexibility of aggregated residential load using smart meter data. *IEEE Transactions on Power Systems* 33(5):5446–5455.
- Ruan G, Kirschen DS, Zhong H, Xia Q, Kang C (2021) Estimating demand flexibility using siamese lstm neural networks. *IEEE Transactions on Power Systems* 37(3):2360–2370.
- Saez-Gallego J, Morales JM (2017) Short-term forecasting of price-responsive loads using inverse optimization. *IEEE Transactions on Smart Grid* 9(5):4805–4814.
- Saez-Gallego J, Morales JM, Zugno M, Madsen H (2016) A data-driven bidding model for a cluster of price-responsive consumers of electricity. *IEEE Transactions on Power Systems* 31(6):5001–5011, URL <http://dx.doi.org/10.1109/TPWRS.2016.2530843>.
- Shahidehpour M, Yamin H, Li Z (2002) *Market operations in electric power systems: forecasting, scheduling, and risk management* (John Wiley & Sons).
- Shi Y, Xu B (2023) Demand-side price-responsive flexibility and baseline estimation through end-to-end learning. *IET Renewable Power Generation*.
- Singh T (2024) Smart home dataset with weather information [online] Available at: URL <https://www.kaggle.com/taranvee/smart-home-dataset-with-weather-information>.

- Srivastava A, Zhao J, Zhu H, Ding F, Lei S, Zografopoulos I, Haider R, Vahedi S, Wang W, Valverde G, et al. (2024) Distribution system behind-the-meter ders: Estimation, uncertainty quantification, and control. *IEEE Transactions on Power Systems* .
- Sun M, Wang Y, Teng F, Ye Y, Strbac G, Kang C (2019) Clustering-based residential baseline estimation: A probabilistic perspective. *IEEE Transactions on Smart Grid* 10(6):6014–6028.
- Sunar N, Swaminathan JM (2022) Socially relevant and inclusive operations management. *Production and Operations Management* 31(12):4379–4392.
- Tlenshiyeva A, Tostado-Véliz M, Hasanien HM, Khosravi N, Jurado F (2024) A data-driven methodology to design user-friendly tariffs in energy communities. *Electric Power Systems Research* 228:110108.
- Vallés M, Bello A, Reneses J, Frías P (2018) Probabilistic characterization of electricity consumer responsiveness to economic incentives. *Applied Energy* 216:296–310.
- Vardakas JS, Zorba N, Verikoukis CV (2014) A survey on demand response programs in smart grids: Pricing methods and optimization algorithms. *IEEE Communications Surveys & Tutorials* 17(1):152–178.
- Vatandoust B, Zad BB, Vallée F, Toubreau JF, Bruninx K (2023) Integrated forecasting and scheduling of implicit demand response in balancing markets using inverse optimization. *2023 19th International Conference on the European Energy Market (EEM)*, 1–6 (IEEE).
- Watson RT, Boudreau MC, Chen AJ (2010) Information systems and environmentally sustainable development: energy informatics and new directions for the is community. *MIS quarterly* 23–38.
- Wen L, Zhou K, Li J, Wang S (2020) Modified deep learning and reinforcement learning for an incentive-based demand response model. *Energy* 205:118019.
- Zhang XY, Watkins C, Kuenzel S (2022) Multi-quantile recurrent neural network for feeder-level probabilistic energy disaggregation considering roof-top solar energy. *Engineering Applications of Artificial Intelligence* 110:104707.
- Zhang Y, Chen W, Xu R, Black J (2015) A cluster-based method for calculating baselines for residential loads. *IEEE Transactions on smart grid* 7(5):2368–2377.
- Zhou B, Jiang R, Shen S (2024) Learning to solve bilevel programs with binary tender. *Conference paper at ICLR 2024* .

Supplementary information. Estimating the Unobservable Components of Electricity Demand Response with Inverse Optimization

Adrian Esteban-Perez^{1*}, Derek Bunn² and Yashar Ghiassi-Farrokhfal¹

^{1*}Department of Technology and Operations Management, Rotterdam School of Management, Erasmus University, Rotterdam, 3062 PA, The Netherlands.

²London Business School, London, NW1 4SA, United Kingdom.

*Corresponding author(s). E-mail(s): estebanperez@rsm.nl;
Contributing authors: dbunn@london.edu; y.ghiassi@rsm.nl;

1 Single-level conic reformulation Supplement

The Karush-Kuhn-Tucker (KKT) conditions for the lower-level problem enables us to reformulate the (regularized) IO model as a single-level optimization program:

$$\min \sum_{s \in [S]} \omega_s u_s + \lambda t \quad (1a)$$

$$\text{s.t. } u_s \geq \left\| \left(\mathbf{d}^{bl} + \mathbf{d}_s^{sf} + \mathbf{d}_s^{sd} - \hat{\mathbf{g}}_s - \hat{\mathbf{d}}_s \right)_s \right\|_p^p \quad \forall s \in [S] \quad (1b)$$

$$t \geq \left\| [\beta^{sf,+}, \beta^{sf,-}, \beta^{sd}] \right\|_p^p \quad (1c)$$

$$(5b) - (5h) \quad (1d)$$

$$(6), (7), (8) \quad \forall s \in [S], \forall t \in [T] \quad (1e)$$

$$-2c_{s,t}^{sf,+} d_{s,t}^{sf,+} + p_{s,t}^{sf,+} - p_{s,t} + \kappa_s - \mu_{s,t}^+ + \nu_{s,t}^+ = 0 \quad \forall s \in [S], \forall t \in [T] \quad (1f)$$

$$-2c_{s,t}^{sf,-} d_{s,t}^{sf,-} + p_{s,t}^{sf,-} + p_{s,t} - \kappa_s - \mu_{s,t}^- + \nu_{s,t}^- = 0 \quad \forall s \in [S], \forall t \in [T] \quad (1g)$$

$$-2c_{s,t}^{sd} d_{s,t}^{sd,-} + p_{s,t}^{sd} + p_{s,t} - \mu_{s,t}^0 + \nu_{s,t}^0 = 0 \quad \forall s \in [S], \forall t \in [T] \quad (1h)$$

$$\mu_{s,t}^+ (\bar{d}_{s,t}^{sf,+} \delta_{s,t}^{sf,+} - d_{s,t}^{sf,+}) = 0 \quad \forall s \in [S], \forall t \in [T] \quad (1i)$$

$$\mu_{s,t}^- (\bar{d}_{s,t}^{sf,-} \delta_{s,t}^{sf,-} - d_{s,t}^{sf,-}) = 0 \quad \forall s \in [S], \forall t \in [T] \quad (1j)$$

$$\begin{aligned}
\mu_{s,t}^0(\bar{d}_{s,t}^{sd} - d_{s,t}^{sd,-}) &= 0 & \forall s \in [S], \forall t \in [T] & \quad (1k) \\
\nu_{s,t}^0 d_{s,t}^{sd,-} &= 0 & \forall s \in [S], \forall t \in [T] & \quad (1l) \\
\nu_{s,t}^+ d_{s,t}^{sf,+} &= 0 & \forall s \in [S], \forall t \in [T] & \quad (1m) \\
\nu_{s,t}^- d_{s,t}^{sf,-} &= 0 & \forall s \in [S], \forall t \in [T] & \quad (1n) \\
\mu_{s,t}^+, \mu_{s,t}^-, \mu_{s,t}^0, \nu_{s,t}^+, \nu_{s,t}^-, \nu_{s,t}^0 &\geq 0 & \forall s \in [S], \forall t \in [T] & \quad (1o) \\
\theta_s &\in \Theta(\bar{d}_{s,t}^{sf,+}, \bar{d}_{s,t}^{sf,-}, \bar{d}_{s,t}^{sd}) & \forall s \in [S] & \quad (1p) \\
\kappa_s \in \mathbb{R}, u_s &\geq 0 & \forall s \in [S] & \quad (1q) \\
t &\geq 0 & & \quad (1r)
\end{aligned}$$

The above program can be recast as a mixed-integer conic-constrained program solvable by off-the-shelf conic optimization software such as MOSEK (MOSEK 2024). The complementary conditions can be rewritten as special ordered sets of type 1 constraints (SOS1), which are a set of non-negative variables that only one of them can be positive, by avoiding the need to select the big-M parameters for the complementary conditions without jeopardizing the computational performance (Kleinert and Schmidt 2023, Siddiqui and Gabriel 2013). Also, the above objective function via the p -norms can be reformulated as $S + 1$ conic constraints amenable for the efficient conic solver of MOSEK. For the applications below, we consider $p = 2$, and hence the above problem defined by (1a)-(1q) can be modeled as a mixed-integer program with linear objective function and $S + 1$ rotated conic quadratic constraints amenable for MOSEK.

2 Regression Diagnostics Supplement

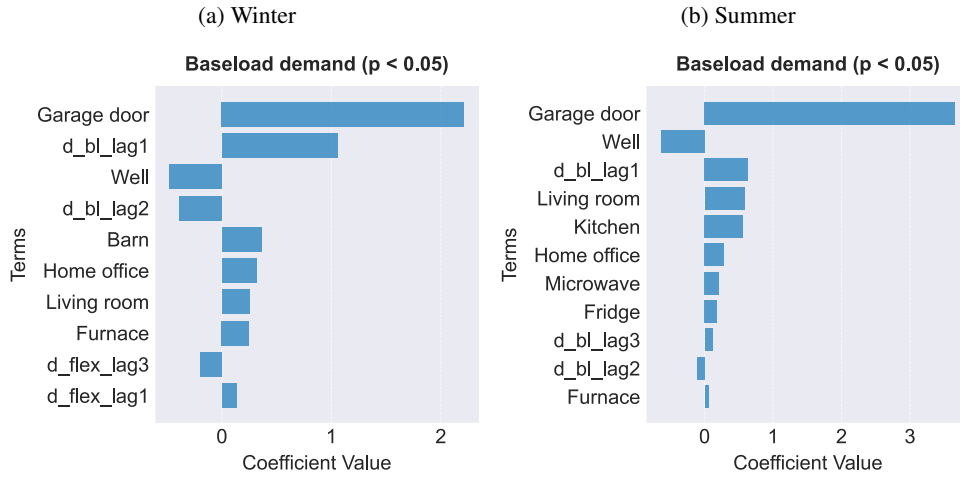


Fig. 1: Significant linear regression coefficients for baseload estimates

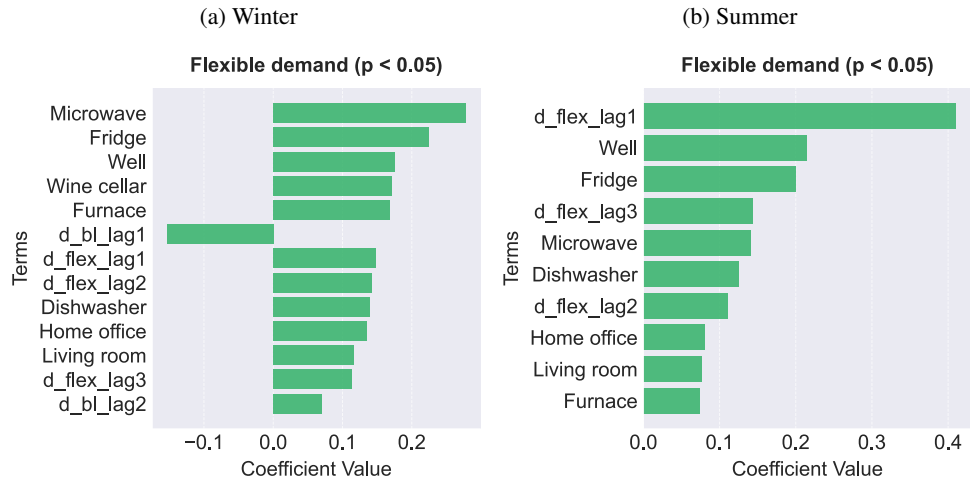


Fig. 2: Significant linear regression coefficients for flexible estimates

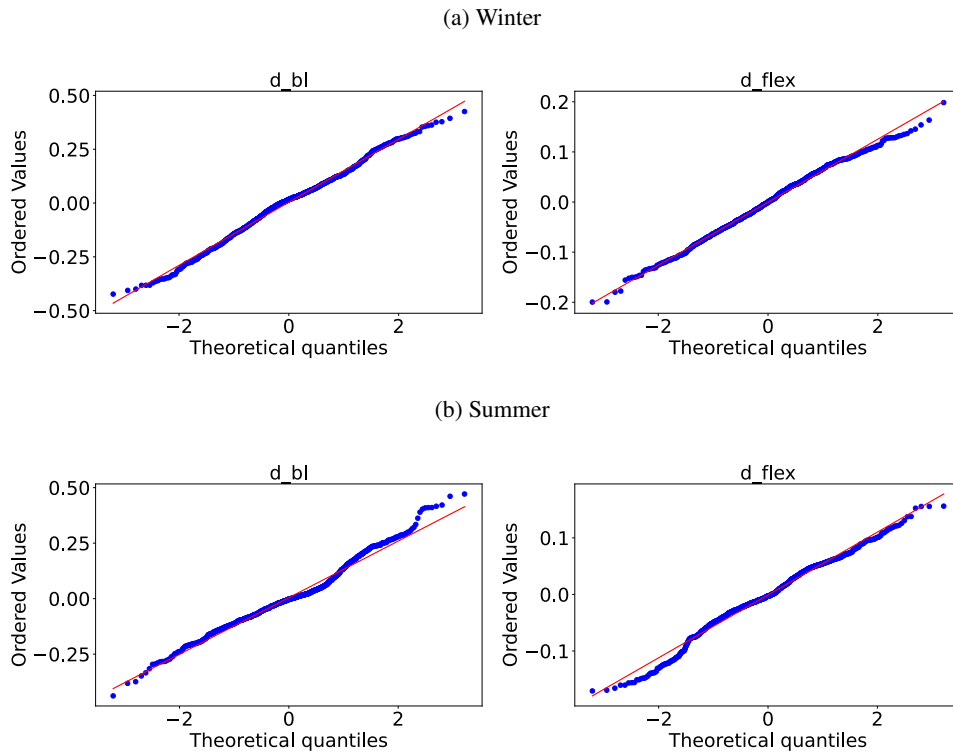


Fig. 3: Residual QQ-plots of load IO estimates

References

- MOSEK: MOSEK ApS. (2024). <https://www.mosek.com>
- Kleinert, T., Schmidt, M.: Why there is no need to use a big-m in linear bilevel optimization: A computational study of two ready-to-use approaches. *Computational Management Science* **20**(1), 3 (2023)
- Siddiqui, S., Gabriel, S.A.: An sos1-based approach for solving mpecs with a natural gas market application. *Networks and Spatial Economics* **13**, 205–227 (2013)