

Neural network interpretability for forecasting of aggregated renewable generation

Yucun Lu, Ilgiz Murzakhanov, Spyros Chatzivasileiadis
Center for Electric Power and Energy
Technical University of Denmark
Kgs. Lyngby, Denmark
luyucun003@gmail.com, {ilgmu, spchatz}@elektro.dtu.dk

Abstract—With the rapid growth of renewable energy, lots of small photovoltaic (PV) prosumers emerge. Due to the uncertainty of solar power generation, there is a need for aggregated prosumers to predict solar power generation and whether solar power generation will be larger than load. This paper presents two interpretable neural networks to solve the problem: one binary classification neural network and one regression neural network. The neural networks are built using TensorFlow. The global feature importance and local feature contributions are examined by three gradient-based methods: Integrated Gradients, Expected Gradients, and DeepLIFT. Moreover, we detect abnormal cases when predictions might fail by estimating the prediction uncertainty using Bayesian neural networks. Neural networks, which are interpreted by the gradient-based methods and complemented with uncertainty estimation, provide robust and explainable forecasting for decision-makers.

Index Terms—Interpretable neural networks, solar power output forecasting, feature importance, transparency, uncertainty estimation

I. INTRODUCTION

Small renewable energy producers, such as photovoltaic (PV) solar panels, are relatively new at the distribution grid level. For example, in Australia, the government used to incentivize small households to install solar panels by offering generous feed-in tariff schemes [1]. On top of covering part of their own consumption, the owners of small PVs were paid for injecting the power surplus to the grid. However, with greater penetration of renewable generation at the distribution grid level, PV inverters started to turn off at peak generation hours due to voltage rising in the network [1]. While there can be several solutions to tackle this problem, the one requiring the lowest investments is the forecasting of aggregated renewable generation and aggregated consumption. From the conducted forecast, the decision-makers may conclude about generation and consumption balance and decide to curtail part of solar production. In this work, we propose methods that can (i) interpret the forecasting techniques, and (ii) can assess how certain a predicted value is.

A. Literature Review

Several research efforts have been dedicated to building models to predict solar power generation and load with the focus on either reducing forecasting errors or improving computation efficiency. The traditional methods for load and solar

power predictions are statistical algorithms, such as polynomial regression and auto-regressive and moving average model (ARMA). In [2], the authors present a linear regression model to predict solar radiation using temperature data. In [3], the authors propose an hourly ARMA model to predict power output from a PV panel. Recently, a great number of researchers work on the development and improvement of several machine learning algorithms in solar power forecasting. The algorithms include support vector regression (SVR), gradient boosted regression tree (GBRT), and artificial neural network (ANN). In [4], the researchers present an SVR model for solar power forecasts on a rolling basis for 24 hours ahead. In [5], a GBRT model is proposed using historical power generation data and relevant meteorological variables. The study in [6] suggests ANN for producing solar power forecast. The study in [7] and [8] applies long short-term memory (LSTM) and convolutional neural networks with LSTM (CNNs-LSTM) respectively to provide accurate forecast results.

While machine learning algorithms provide promising results in load and PV forecasting, they lack interpretability. As a result, explainable artificial intelligence (XAI) has become an emerging research direction, which addresses this problem and helps understanding why and how these models make predictions. Currently, only a few works exist in this field. In [9], the authors propose a unified clustering-based prediction framework with two tree-based algorithms and use the SHapley Additive exPlanations (SHAP) XAI tool to interpret the model. In [10], the authors suggest an interpretable forecasting model using the XGBoost algorithm and ELI5 XAI tool. In [11], the authors apply several XAI tools: LIME, SHAP, and ELI to interpret a random forest model.

B. Main contributions

The main contributions of this paper are the following:

- We conduct local and global feature importance analysis for the designed classification and regression neural networks. From the feature analysis, we discover several interesting dependencies between the features and the outputs, which match with the physical phenomena behind the PV generation process.
- We implement three gradient-based methods for interpreting the designed models and find the main reasons why

the neural networks fail. We also give recommendations on building robust neural networks.

- We design a Bayesian neural network to estimate the aleatoric and epistemic uncertainties of the solar power forecast models, which was not done before. As a result, we detect the abnormal PV generation hours, when our models might fail, and consider these hours more carefully.

C. Outline

The remainder of this paper is organized as follows. First, the description of the proposed methodology is given in Section II. Numerical results with the following analysis are provided in Section III. Finally, Section IV concludes the paper and proposes future directions.

II. PROPOSED METHODOLOGY

Next, we describe two models, which predict whether solar power will exceed load in the next 15-60 minutes. Depending on the models' output, a decision on curtailing part of the PV generation is made.

The first model is the regression neural network, which predicts solar power generation. By comparing it with the predicted load, we can conclude whether solar power will be larger than load. Here we use the previous hour's load as the predicted load because the load does not vary a lot in one hour compared to the solar power generation. As a result, our prediction is largely dependent on the accuracy of solar power prediction. The advantage of the regression model is that we can predict the exact amount of PV generation for curtailment.

The second model is a binary classification neural network that directly predicts whether solar power generation will be larger than load. If the classification model outputs 1, it means that solar power generation will exceed load; and vice versa if the model outputs 0. This model considers both effects of solar power generation and load. The classification model is more robust than the regression one and it can effectively avoid unwanted false-positive cases.

Implementation of both classification and regression models allows us to assess how the interpretability methods perform on the different types of neural networks. We deploy several attribution methods to explain both models by interpreting feature importance globally and locally. Thus, through neural network interpretability, not only do we identify the most influential features in making predictions, but also show how and why neural networks make certain predictions in individual cases. This gives unique insights into the models and can support informed decisions for the human operator.

Additionally, for the regression model of solar power prediction, aleatoric uncertainty and epistemic uncertainty are estimated with Bayesian neural networks. The interpretation of neural networks can tell us why the model predicts certain results, while the uncertainty can tell us how sure the model is about its predictions. Combining interpretation with prediction uncertainty results in more transparent, convincing, and robust decisions.

A. Interpretability

To interpret the neural networks and find the most influential features in predictions, we apply attribution methods. Attribution methods assign an attribution value, sometimes called "contribution" or "relevance", to each input feature of a neural network, which can help to determine how different features contribute to the model's output [12]. By looking into the attribution values of different features, we can provide reasonable explanations of how and why the prediction is made. There are five main gradient-based attribution methods: Gradient, Gradient*Input, Integrated Gradients, Expected Gradients, and Deep Learning Important Features (DeepLIFT). Each of these methods has advantages and shortcomings, which we discuss further.

Historically, one of the first attribution methods adapted in the deep learning domain is Gradient [13]. Since the gradient of the network output measures the instantaneous rate of change of model output with respect to one specific feature, it can naturally be a candidate for attribution values. Another straightforward method is Gradient*Input which is more commonly used when we are concerned with the marginal effect of a feature on the output. This attribution can be computed by multiplying the gradient of the model output with respect to the input itself [14]. Unfortunately, Gradient and Gradient*Input methods have many unsatisfying shortcomings. One of the worst problems is called saturation: the gradients of input features might only have small magnitudes around a sample even if the network depends heavily on those features [15]. The problem is common when the model output flattens after some features get to a certain magnitude. The saturation problem is significant because it might lead us to overlook some influential features but focus on less relevant features, thus yielding the wrong interpretation of networks.

More advanced methods as DeepLIFT, Integrated Gradients, and Expected Gradients, are free of the saturation problem. As a result, we implement these three attribution methods in our work. DeepLIFT is a method for decomposing the output of the network on input features by back-propagating the contributions of all neurons in the network to every input feature [16]. To be more specific, DeepLIFT compares the activation of each neuron x to its "reference" (also called baseline) activation \bar{x} and assigns contribution scores according to the difference. The baseline is defined by the user and often chosen to be zero. The method of Integrated Gradients is similar to the Gradient method in that it also computes the partial derivatives of the output with respect to each input feature [17]. However, instead of computing one single gradient, Integrated Gradients computes several instantaneous gradients on the straight-line path from the baseline to the observation being explained and averages them. The method of Expected Gradients is an extension of Integrated Gradients, designed to remove baseline ambiguity [18]. To make the baseline uninformative, this method uses different baselines and takes an average over them. All three methods satisfy one desirable property called completeness that the attributions

sum up to the target output minus the target output evaluated at the baseline [13]. An overview of implemented gradient-based attribution methods is given in Table I.

B. Uncertainty quantification

The knowledge of uncertainty is fundamental to the development of robust and interpretable machine learning techniques. The uncertainty associated with machine learning models can be broadly classified into two types: aleatoric uncertainty and epistemic uncertainty. The overall uncertainty of any model is a combination of the above two types of uncertainty.

Aleatoric uncertainty refers to the irreducible error of the uncertainty which means the error cannot be reduced by choosing a better model because the uncertainty originates from the non-deterministic nature of the sought input/output dependency. Aleatoric uncertainty captures the uncertainty concerning information that our data cannot explain [19]. For example, in lab experiments, even when all input values are similar, the values measured after multiple trials are never the same. This type of uncertainty is a typical aleatoric uncertainty, while uncertainty due to a lack of knowledge about the perfect predictor, for example, caused by uncertainty about the parameters of a model, is called epistemic uncertainty. Epistemic uncertainty captures our limited understanding of the real-world process for which we are building the model, as we are not able to capture all the input features that affect the target variable [19]. In principle, epistemic uncertainty is a reducible error which means this uncertainty can be reduced with more knowledge about the process.

In this paper, we apply Bayesian neural networks to estimate the two types of uncertainties. Bayesian neural networks are similar to normal neural networks except that instead of having a model parameterized by its point weights, each weight of the Bayesian neural network now has a probability distribution with a mean and variance which is tuned during the training process. For each point, epistemic uncertainty is modeled by Monte Carlo sampling. In Bayesian neural networks, there is a prior distribution over the model's weight. As a result, each time when we apply the network to make predictions, we sample the model weight from distributions. In other words, every time we run the model, we would have different model weights and thus a different predicted output. After we run the model hundreds of times, the standard deviation of model predictions can be estimated as epistemic uncertainty. The estimation of aleatoric uncertainty is much easier. Here, apart from the prior distribution put over the model's weights, we also place a distribution over the output of the model [20]. Thus, the Bayesian neural network can directly output a distribution with mean and variance. The mean value is the predicted value while the standard deviation is our estimated aleatoric uncertainty.

III. NUMERICAL RESULTS

In this section, we describe the used dataset, introduce results for the designed neural networks, and provide analysis

on interpretability and uncertainty quantification. All code used for the analysis in this section is available at the GitHub repository [21].

A. Dataset

For numerical simulations, we use the data set from Presumed Open Data: Data Science Challenge, which is provided by Western Power Distribution energy data hub [22]. The load data is extracted from the Stentaway Primary substation near Plymouth, on the south coast of the UK. The solar PV generation data is from a solar farm in Devon, UK, which is not too far from the Stentaway substation. The original load and solar power data consist of half-hourly average power values from November 2017 to July 2020. Besides, hourly irradiance forecast and hourly surface temperature forecast data from January 2015 to July 2020 have been extracted from six different sites which are close to Devon but in different directions.

We follow standard data preprocessing procedures: data cleaning, outlier detection, and data imputation. During this process, we only retain samples from 7:00 - 18:00 in each day because other hours have zero or close to zero solar power generation throughout the whole year. The final data set consists of 10623 hourly instances from November 2017 to July 2020 with 7 different attributes.

A description of all data attributes and their type is presented in Table II. As mentioned before, we have hourly irradiance forecast data and hourly surface temperature forecast data of six sites. However, in the neural network training, we choose forecast data of only one site as input features due to the found high co-linearity in the data. Note that we consider load and generation values for an hour, as a result, use kWh units further.

Note that during model training, the hour index has been encoded using a one-hot (also called 'one-of-K' or 'dummy') encoding scheme which creates new binary columns to indicate the presence of each hour from the original data. One-hot encoding can make the representation of categorical data more expressive and ensure that the neural network does not assume that higher numbers (by usual integer encoding) are more important.

B. Neural network models

Conventionally, the data set is divided into training and test sets. 10% of the data are randomly shuffled for testing and the rest are used for training. The motivation behind random shuffling is that the performance of the neural network depends on how many previously unseen data samples are contained in the test set. By randomly shuffling the data we guarantee that the model trains on as many data samples as possible while it avoids overfitting. When randomly shuffling the data set, a common random seed is used so that we can compare and validate the model with a consistent train/test split. Note that the training and testing sets are the same for the regression and classification models.

TABLE I
IMPLEMENTED GRADIENT-BASED ATTRIBUTION METHODS

Method	Attribution
Integrated Gradients	$(x_i - \bar{x}_i) \cdot \int_{\alpha=0}^1 \frac{\partial y_c(\bar{x})}{\partial \tilde{x}_i} \Big _{\bar{x}=\bar{x}+\alpha(x-\bar{x})} d\alpha$
Expected Gradients	$\int_{\bar{x}} \left((x_i - \bar{x}_i) \cdot \int_{\alpha=0}^1 \frac{\partial y_c(\bar{x})}{\partial \tilde{x}_i} \Big _{\bar{x}=\bar{x}+\alpha(x-\bar{x})} d\alpha \right) p_D(\bar{x}) d\bar{x}$
DeepLIFT	$(x - \bar{x}_i) \cdot \frac{\partial^g y_c(x)}{\partial x_i}, g = \frac{\sigma(z) - \sigma(\bar{z})}{z - \bar{z}}$

TABLE II
DESCRIPTION OF ATTRIBUTES IN DATA SET

Feature	Description	Type
Index	Hour index	Categorical
STemp	Surface temperature forecast ($^{\circ}C$)	Float
Irra	Irradiance forecast (W/m^2)	Float
PTemp	Solar panel temperature in last hour ($^{\circ}C$)	Float
HPow	Previous hour's solar power (KWh)	Float
DPow	Solar power at the same hour at previous day (KWh)	Float
HLoad	Previous hour's load (KWh)	Float

1) *Classification model*: We build a neural network of three hidden layers with 50, 30, 10 neurons, and the ReLu activation function. The output layer has one neuron and the activation function is sigmoid so that when the output is larger than 0.5, it predicts that solar power is larger than load. We choose binary cross-entropy as the loss function and Adam as the optimizer. For the classification problem, we use the accuracy metric to evaluate the model. The test accuracy of the model is 91.6%.

2) *Regression model*: For the solar power forecast regression problem, only the first 6 input features listed in Table II are used. A neural network with the same structure as for the classification problem is built for the regression problem. The only difference is that now the output layer is one neuron with the sigmoid activation function combined with a lambda layer whose function is to do simple mathematical operations on the previous layer without adding more trainable weights. Here, the lambda layer is used to multiply the sigmoid output from the last layer with the solar power capacity. Since the sigmoid activation function will output from 0 to 1, when it is multiplied by the capacity of solar PV, the predictions are limited within the range from 0 to solar power capacity of 2000 kW. Thus, the neural network will not output some implausible and unpractical values like negative power generation. We select the mean squared error (MSE) as the loss function and Adam as the optimizer. The root mean square error (RMSE) metric is used to evaluate the model, which can provide a global error measure during the entire forecasting period [23]. The RMSE of the solar power forecast model is around 144 kWh.

C. Interpretability

Neural networks are usually referred to as “black-box” models due to lack of transparency. In this section, we focus on the interpretability aspects of the neural networks: why and how the models make the prediction and which features

have more or fewer contributions to the final predicted results. Next, three attribution methods are implemented to explain the predictions of two models: Integrated Gradient, Expected Gradient, and DeepLIFT.

First, we obtain a global view of model behavior and find out which features are the most influential for the models’ predictions. We accomplish this goal by calculating the average magnitude of feature attributions across all data set points. We conclude that the features with a larger value have more deterministic effects on the prediction results. The most important global features for the classification and regression models are shown in Table III and Table IV, respectively. We use the extreme baseline for Integrated Gradients and DeepLIFT methods, which means that the baseline prediction value is obtained by setting all features to zero.

TABLE III
GLOBAL FEATURE IMPORTANCE OF THE CLASSIFICATION MODEL

Attribution method	Feature importance					
	DPow	HPow	HLoad	Irra	PTemp	STemp
Integrated Gradients	0.185	0.462	0.758	0.336	0.234	0.034
Expected Gradients	0.046	0.228	0.079	0.186	0.116	0.011
DeepLIFT	0.188	0.487	0.776	0.360	0.254	0.037

TABLE IV
GLOBAL FEATURE IMPORTANCE OF THE REGRESSION MODEL

Attribution method	Feature importance				
	DPow	HPow	Irra	PTemp	STemp
Integrated Gradients	37.20	450.39	170.18	146.64	40.36
Expected Gradients	15.81	336.10	114.64	77.50	13.70
DeepLIFT	44.22	288.62	206.80	47.74	31.26

For the classification model in Table III, the most influential features are the previous hour’s load (HLoad), the previous hour’s solar power (HPow), and irradiance forecast (Irra). For Integrated Gradients and DeepLIFT methods, HLoad is more influential than HPow. As for Expected Gradients, HLoad is a less important feature. The reason is that Expected Gradient uses different baselines and takes an average over them, thus the baseline of Expected Gradients is close to the average values of each feature. As discovered from the data set, solar power generation is more unstable than load. As a result, the variance of the load is much smaller than the variance of solar power generation. Thus, the influence of load is reduced during the selection of average values as the baseline. Additionally, we find that global feature importance interpreted by Integrated gradients and DeepLIFT have similar patterns:

their average magnitudes of feature contributions are very close to each other. Our observations match with the previous research, where DeepLIFT is characterized as a good and fast approximation of Integrated Gradient [12].

For the regression model in Table IV, the previous hour’s solar power (HPow) and irradiance forecast (Irra) are the most important features for prediction. However, we observe that panel temperature (PTemp) is also an informative feature. We explain it by the physical phenomena when the high temperature of solar panels decreases their efficiency, as a result, leading to lower generation of PV panels.

The aforementioned global feature analysis allows defining the most influential features for predictions. However, this analysis does not explain how exactly the features affect the predicted output. To gain more insights on that, we apply the DeepLIFT method and plot the summary plots for the classification and regression models in Fig. 1 and Fig. 2, respectively. The summary plots show the global feature impacts on the models’ output, through which we can reveal the inner workings of neural networks. On the plots, the blue dot means a point with a low feature value, while the red dot means a point with a high feature value.

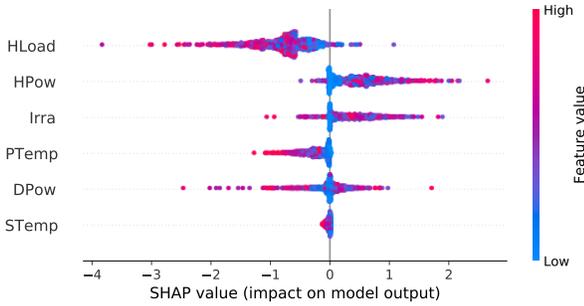


Fig. 1. Summary plot of the classification model

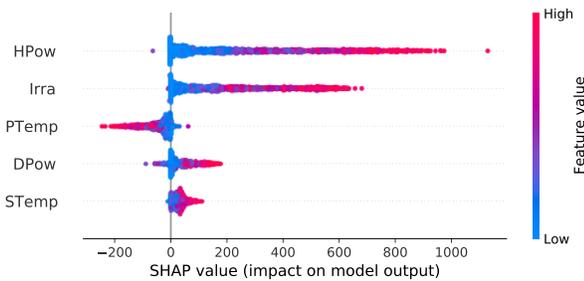


Fig. 2. Summary plot of the regression model

For the classification model in Fig. 1, HLoad generally has a negative contribution while HPow and Irra have positive impacts on model output in most cases. When HPow and Irra have high values, they tend to have large positive impacts on model output, which means the model will be more likely to predict 1 (solar power is greater than load). On the contrary, when HPow and Irra have small contributions, the model’s output is more likely to be 0. For HLoad, it has negative

impacts on the model output, that is, when HLoad has high values, it will promote model output to be 0, and vice versa. Both of these observations perfectly match with common-sense expectations.

For the regression model in Fig. 2, the baseline features are set to be zero. As a result, Irra and HPow always have positive contributions compared to zero irradiance or zero generation in the last hour. When the regression model predicts a large solar power generation, HPow and Irra tend to have larger contributions. We also discover that PTemp has negative effects on solar power predictions. When PTemp value is high, it decreases the model output. As explained earlier, it happens due to high temperature decreases the efficiency of the solar panel.

To get a better understanding of the role of different features contributing to the output and compare the roles those features play in each case, we examine four individual cases of local interpretability: True Positive (TP), True Negative (TN), False Negative (FN), False Positive (FP). The attribution values of the four individual cases interpreted by DeepLIFT for the classification and regression models are presented in Table V and Table VI, respectively. As observations for the models are similar, we analyze Table V and Table VI together.

TABLE V
FEATURE CONTRIBUTION TO THE PREDICTION OF THE CLASSIFICATION MODEL

Features	Attribution			
	TP	TN	FN	FP
Base value	0.634	0.634	0.634	0.634
Index	-0.001	-0.001	-0.003	-0.001
STemp	-0.034	-0.008	-0.132	-0.031
Irra	0.118	0	1.231	0.625
PTemp	-0.296	-0.04	-0.803	-0.296
HPow	0.799	0.154	0.816	0.425
DPow	0.14	0.018	0.196	-0.047
HLoad	-0.36	-0.755	-1.489	-0.318
Predicted value	1.000	0.002	0.450	0.991

TABLE VI
FEATURE CONTRIBUTION TO THE PREDICTION OF THE REGRESSION MODEL

Features	Attribution			
	TP	TN	FN	FP
Base value	1.32	1.32	1.32	1.32
Index	-68.73	-11.13	-42.53	-59.76
STemp	45.67	7.57	45.3	31.93
Irra	534.43	3.09	226.76	421
PTemp	-89.63	-11.35	-77.31	-74.51
HPow	636.33	17.18	278.1	510.92
DPow	122.83	3.5	45.38	104.21
Predicted value	1182.22	10.18	477.02	935.11
Real value	1180.00	6.00	784.00	470.00
Predicted load	400	810	560	414

For True Positive (TP) case, we see that the large attribution values of Irra and Hpow lead to a high predicted solar power generation. After comparing solar generation with the predicted load, the regression model can give a prediction that solar power will exceed load. In the classification model, due to large positive attributions of Hpow and large negative attributions of HLoad, the model also gives output close to

1, which shows that the model is quite confident that solar power will exceed load. In this case, both models agree on the forecast, which is correct.

For True Negative (TN) case, Irra and HPow have small contributions leading to a low predicted solar power generation in the regression model. Thus, when compared to the previous hour's relatively high load, it predicts that solar power is smaller than load in this hour. For the classification model, because of low Irra and HPow contributions and large negative HLoad contributions, the model also gives a quite confident prediction (close to 0) that solar power will be smaller than load. In this case, both models make perfect predictions too.

Next, let's look into the cases when models fail and find out why the models make wrong predictions. For False Negative (FN) case, we see that solar power prediction is 477.02 KWh, which is slightly lower than the predicted load of 560 KWh. However, real solar generation in this hour turns out to be 784 KWh which is much larger than the predicted value of 477.02 KWh. The classification model also outputs an unsure prediction (close to 0.5) that solar power will be smaller than load. As a result, both models provide the wrong prediction. The reason is that solar power has a dramatic increase of around 42% in one hour, which could not be forecasted by the models.

For False Positive (FP) case, high contribution of HPow and Irra and low contribution of HLoad output a high predicted solar power and a confident prediction (close to 1) that solar power is larger than load. However, it turns out that both models are wrong. The reason is that solar power generation has a dramatic decrease in this hour. Although the solar power is high in the last hour which is 1012 KWh, it drops sharply to 470 KWh. Besides, the failure of irradiance forecast also contributes to the wrong prediction of two models.

We use the SHAP tool [11] to make a force plot that visualizes the contribution of each feature. A force plot for the regression model of TP case is given in Fig. 3. We see that Irra, HPow, and DPow have large positive contributions and promote the predicted solar power generation to a higher value. On the contrary, high solar panel temperature slightly decreases the solar power output.

To sum up, we find that the most important features for solar power forecast are HPow (previous hour's generation) and Irra (Irradiance forecast). These two features, together with HLoad (previous hour's load) have significant impacts on the classification model too. If the solar irradiance forecast is accurate and if the solar power does not have a sudden dramatic change within one hour, our models provide accurate predictions.

D. Uncertainty quantification

The neural networks only provide a point estimation of solar power generation. As a result, the regression neural network in section III-C is "equally confident" in all described TP, TN, FP, FN cases, while hours during a day are not equivalent in terms of solar irradiance. We propose to address this problem

by quantifying the uncertainty of the solar power forecast. Uncertainty estimation of the models allows increasing efficiency of the decision-making process. So, designed neural networks can perform autonomously in hours with low uncertainty. On the contrary, in the hours with high uncertainty, the models can request the intervention of the domain expert. We quantify aleatoric and epistemic uncertainties by designing Bayesian neural networks. Uncertainty estimation using Bayesian neural network is shown in Fig. 4, where red lines depict different predictions made through Monte Carlo sampling, and the standard deviation of these values can be interpreted as epistemic uncertainty. Green lines represent the standard deviation of the model's output distribution which can be seen as aleatoric uncertainty. We find that both uncertainties are high when solar power generation is high, which happens in the middle of the day. Moreover, aleatoric uncertainty is high in some abnormal hours. Summing up aleatoric and epistemic uncertainties, we plot a probabilistic forecast with 95% confidence in Fig. 5. While for some abnormal hours, the predicted solar output is far from the real one, it is still within the confidence interval. Bayesian neural networks allow finding the reason for this difference, which is a sudden change in solar power generation.

IV. CONCLUSION AND FUTURE WORK

In this paper, we consider the problem of designing explainable neural networks, particularly, for solar generation forecasting. Without becoming human interpretable, neural networks cannot be implemented in such critical infrastructure, as power systems. To interpret the designed classification and regression neural networks, we conduct global and local feature analysis. We discover several interesting dependencies between the weather parameters and outputs of the neural networks, which originate from the physical processes during the operation of PV panels. In order to determine the most influential features for forecasting, we implement three state-of-the-art attribution methods: Integrated Gradients, Expected Gradients, and DeepLIFT. We find out that load and solar generation of the previous hour and solar irradiance influence the most to the outcome of the neural networks. In addition, we discover that the main source of the wrong forecasts is the sudden dramatic change of solar irradiance. Such change can be caused, for example, by cloud movements. To resolve this problem we implement Bayesian neural networks, which with the use of historical data, determine hours with high uncertainty of the PV output. As a result, these hours can be known beforehand and considered more carefully by the domain expert.

We see our work as the first step towards the creation of continuously learning neural networks for PV forecasting. Currently, the designed neural networks can perform autonomously in hours with low uncertainty, while in the hours with high uncertainty, the models can request the intervention of the domain expert. The neural network can learn from the decisions made by the human expert, which will further enhance the neural network performance and make

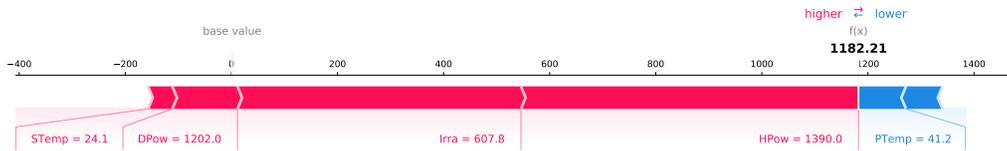


Fig. 3. Force plot of the regression model for TP case

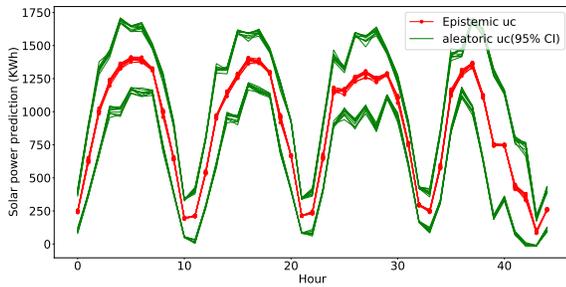


Fig. 4. Uncertainty estimation using Bayesian neural networks

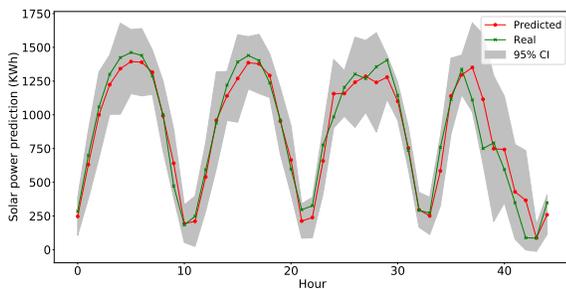


Fig. 5. Probabilistic forecast

a step towards autonomous forecasting and decision-making algorithms performing on the above human level.

REFERENCES

- [1] J. T. The Guardian, "Under new rules for selling solar power, is it still worth it?" 2019. [Online]. Available: <https://www.theguardian.com/money/2019/jun/30/solar-panels-smart-export-guarantee-is-it-still-worth-it>
- [2] S. Ibrahim, I. Daut, Y. M. Irwan, M. Irwanto, N. Gomesh, and Z. Farhana, "Linear regression model in estimating solar radiation in perlis," *Energy Procedia*, vol. 18, pp. 1402–1412, 2012.
- [3] B. Singh and D. Pozo, "A Guide to Solar Power Forecasting using ARMA Models," *arXiv*, pp. 0–3, 2018.
- [4] M. Abuella and B. Chowdhury, "Solar power forecasting using support vector regression," *arXiv*, 2017.
- [5] C. Persson, P. Bacher, T. Shiga, and H. Madsen, "Multi-site solar power forecasting using gradient boosted regression trees," *Solar Energy*, vol. 150, pp. 423–436, 2017. [Online]. Available: <http://dx.doi.org/10.1016/j.solener.2017.04.066>
- [6] M. Abuella and B. Chowdhury, "Solar power forecasting using artificial neural networks," *2015 North American Power Symposium, NAPS 2015*, 2015.
- [7] H. Sharadga, S. Hajimirza, and R. S. Balog, "Time series forecasting of solar power generation for large-scale photovoltaic plants," *Renewable Energy*, vol. 150, pp. 797–807, 2020. [Online]. Available: <https://doi.org/10.1016/j.renene.2019.12.131>
- [8] V. Suresh, P. Janik, J. Rezmer, and Z. Leonowicz, "Forecasting solar PV output using convolutional neural networks with a sliding window algorithm," *Energies*, vol. 13, no. 3, 2020.
- [9] X. Chang, W. Li, J. Ma, T. Yang, and A. Y. Zomaya, "Interpretable Machine Learning in Sustainable Edge Computing: A Case Study of Short-Term Photovoltaic Power Output Prediction," *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*, vol. 2020-May, pp. 8981–8985, 2020.
- [10] S. Sarp, M. Kuzlu, U. Cali, O. Elma, and O. Guler, "An Interpretable Solar Photovoltaic Power Generation Forecasting Approach Using An Explainable Artificial Intelligence Tool," *2021 IEEE Power & Energy Society Innovative Smart Grid Technologies Conference (ISGT)*, 2021.
- [11] M. Kuzlu, U. Cali, V. Sharma, and Ö. Güler, "Gaining insight into solar photovoltaic power generation forecasting utilizing explainable artificial intelligence tools," *IEEE Access*, vol. 8, pp. 187 814–187 823, 2020.
- [12] M. Ancona, E. Ceolini, C. Öztireli, and M. Gross, "Towards better understanding of gradient-based attribution methods for deep neural networks," *6th International Conference on Learning Representations, ICLR 2018 - Conference Track Proceedings*, pp. 1–16, 2018.
- [13] K. Simonyan, A. Vedaldi, and A. Zisserman, "Deep inside convolutional networks: Visualising image classification models and saliency maps," *2nd International Conference on Learning Representations, ICLR 2014 - Workshop Track Proceedings*, pp. 1–8, 2014.
- [14] M. Ancona, E. Ceolini, C. Öztireli, and M. Gross, "Gradient-Based Attribution Methods," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 11700 LNCS, pp. 169–191, 2019.
- [15] P. Sturmfels, S. Lundberg, and S.-I. Lee, "Visualizing the impact of feature attribution baselines," *Distill*, 2020, <https://distill.pub/2020/attribution-baselines>.
- [16] A. Shrikumar, P. Greenside, and A. Kundaje, "Learning important features through propagating activation differences," *34th International Conference on Machine Learning, ICML 2017*, vol. 7, pp. 4844–4866, 2017.
- [17] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," *arXiv*, 2017.
- [18] G. Erion, J. D. Janizek, P. Sturmfels, S. Lundberg, and S.-I. Lee, "Improving performance of deep learning models with axiomatic attribution priors and expected gradients," 2019. [Online]. Available: <http://arxiv.org/abs/1906.10670>
- [19] E. Hüllermeier and W. Waegeman, "Aleatoric and epistemic uncertainty in machine learning: an introduction to concepts and methods," *Machine Learning*, vol. 110, no. 3, pp. 457–506, 2021.
- [20] A. Kendall and Y. Gal, "What uncertainties do we need in Bayesian deep learning for computer vision?" *Advances in Neural Information Processing Systems*, vol. 2017-December, no. Nips, pp. 5575–5585, 2017.
- [21] Y. Lu, "Interpretability," 2021. [Online]. Available: <https://github.com/yucunlu/Interpretability>
- [22] Western Power Distribution Innovation, "Presumed Open Data: Data Science Challenge," 2021. [Online]. Available: <https://www.westernpower.co.uk/pod-data-science-challenge>
- [23] J. Zhang, B.-M. Hodge, A. Florita, S. Lu, H. F. Hamann, and V. Banunarayanan, "Metrics for Evaluating the Accuracy of Solar Power Forecasting," *3rd International Workshop on Integration of Solar Power into Power Systems*, vol. 17436, no. October, pp. 1–10, 2013.