

Power Output Reconstruction of Photovoltaic Curtailment

Roman Kohút

*Institute of Information Engineering,
Automation and Mathematics
Slovak University of Technology in Bratislava
Bratislava, Slovakia
roman.kohut@stuba.sk*

Michal Kvasnica

*Institute of Information Engineering,
Automation and Mathematics
Slovak University of Technology in Bratislava
Bratislava, Slovakia
michal.kvasnica@stuba.sk*

Abstract—The use of renewable energy sources in the grid’s energy mix has recently gained popularity. Especially as solar photovoltaic (PV) generation production has almost zero emissions during its operations, they are preferred over fuel-based electricity production. However, expanding PV generation in grid capacity increases the chance of PV curtailment occurrence. Not to mention that microgrids supported with large-scale PV generation almost certainly create PV curtailment regularly. As the forecast of PV production is one of the electricity grid operation cornerstones, the prediction model should be as accurate as possible. The latest trend is utilizing machine learning (ML) models to predict PV output, thanks to their excellent learning and regression capabilities. However, its performance can be highly influenced by measurements used during the model design. Unfortunately, only some of the research on this topic deals with the PV curtailment problem resulting in underperforming ML models. This paper proposes a novel approach to identify and replace curtailed PV measurements. The methodology includes the physical model as a baseline of truly producible energy, which is then investigated and corrected as a piecewise linear system using Pearson correlation and weather measurements. Through real-life comparative scenarios, the suggested data reconstruction method provides increased accuracy of supervised ML-based solar power prediction.

Index Terms—Photovoltaic curtailment, Solar energy prediction, Physical model, Pearson correlation, Machine learning

I. INTRODUCTION

In recent years the world has changed its leaning towards renewable sources, especially solar energy emerges as one of the leading clean and cheap power production. This change in the energy mix is essential as it establishes a sustainable and ecological electrical system [1]. However, the operating grid can not rely on pure renewable production as its generation is volatile and can not be directly shifted towards the electricity load demand [2]. This burden requires power generation maintainers to predict their production most accurately. Based on this generation’s forecasts, microgrid customers can schedule load demand, buy electricity on the short-term market or utilize their model predictive control to optimize the operation of the microgrid [3]. On the other hand, large grid operators can prepare other resources to provide electricity when needed [4]. Hence, the quality of the provided photovoltaic (PV) production predictions is crucial, as many decision-making actions depend on it.

Many approaches exist to designing PV forecasters, with different prediction horizons and steps considering their use. The bases in microgrid control or real-time grid scheduling are day-ahead forecasts with quarter-hourly or hourly steps. The direct solar energy predictions can be categorized as physical, statistical, and machine learning (ML) [5], [6] approaches. The most widely used are ML-based prediction especially deep-learning techniques, thanks to its excellent regression performance through its capability of learning hidden patterns. The bottleneck of this method are data used to train and generate forecasts. To accurately predict future PV production and its fluctuation, it is necessary to correctly identify and predict weather conditions acting on the PV system. In this field, it is standard to use various data preprocessing techniques, as shown in [7], to improve prediction capabilities [8]. However, most of them focus on removing outliers, filtering signals, replacing missing data, and finding the correct feature inputs to the ML models. A few of them solve the problem of the PV curtailment [9], which directly affects solar power measurements see [10], [11]. The critical step to finding the most accurate ML model is to consider all possible influences that can act on the PV system. If the data used contain curtailed PV production model will underperform as it is trained to provide false (lower) PV prediction for specific inputs.

This paper’s main contribution and purpose is to improve any ML-based PV power forecasts of the systems where PV curtailment occurs. Proposed data reconstruction algorithm utilizes the physical model of the PV system, the Pearson correlation coefficient, and a wide range of weather measurements. The capability of used data reconstruction method is tested on a real large-scale system with five different scenarios. To demonstrate the benefits of this approach, multiple ML-based models are trained (using the reconstructed and curtailed data) and tested to produce a day-ahead prediction.

II. PRELIMINARIES

A. System and Data Description

In general solar plants consist of multiple assets. The main two are photovoltaic generators that directly produce electricity from sunlight and power converters that transform direct current (DC) to alternating current (AC).

The power plant is capable of functioning in three distinct modes. The first mode is referred to as the off-the-grid system, which operates independently of any national or local electricity distribution network. The remaining two setups, known as the hybrid and grid-connected, are connected to the primary power grid. Both the hybrid and off-the-grid modes allow the utilization of excess energy through battery storage, although the battery's capacity is typically undersized. In addition, the grid connection for the hybrid and grid-connected setups enables the exportation of excess electricity to the primary grid or importing of electricity during any shortages.

Measurement of this solar power plant generation can be done through inverters or smart electricity meters. The negative of this system from the data perspective is that the measurement collection is limited only to the exact energy production. When the current possible power generation on the solar panels is greater than the load's consumption, naturally, the load will consume only the necessary power. The extra energy is lost (curtailed), as the information about it. If the power plant includes batteries, this behavior is only shifted in time as the capacity of the batteries reaches its maximum and can no longer store energy, so the power production is curtailed. Such behavior occurs mainly in an off-grid system, as others can export excessive electric energy. However, grid contestants must meet agreed deals as it can destabilize the whole grid operation. This concluded in the same situation that PV production is curtailed as the load consumption decreases.

B. Physical Model

The PV system's output relies on multiple factors, categorized as weather conditions and mechanical properties. The more significant factor is the corresponding meteorological state, which behaves as a multi-variable non-linear system, including solar irradiance, temperature, humidity, wind speed, pressure, visibility, and many more. Not all the mentioned parameters are essential to developing a good estimator of accessible power production. The most important is solar irradiance, which directly transforms into electrical energy and temperature, influencing PV panel efficiency.

In general, total solar irradiance G_t can be modeled as a sum of three components, surface absorbed irradiance G_s , diffused irradiance G_d , and ground reflected irradiance G_g

$$G_t(k) = G_s(k) + G_d(k) + G_g(k). \quad (1)$$

By direct measurement of $G_{DHI}(k)$ diffuse horizontal irradiance and $G_{DNI}(k)$ direct normal irradiance at each time step k is possible to determine each component separately, concerning α_{STA} and α_{SAA} PV surface tiled and azimuth angle respectively. The final output power production can be calculated as energy produced from total solar irradiance affected by current PV panel efficiency

$$P_t(k) = P_{\max} \frac{G_t(k)}{G_{STC}} N_p \left(1 + \kappa (T_c(k) - T_{STC}) \right), \quad (2)$$

where the left part of the brackets represents total irradiation (1) transformed at defined time step k at standard test conditions (STC) $G_{STC} = 1000 \text{ W/m}^2$, $T_{STC} = 25 \text{ }^\circ\text{C}$, P_{\max} is

a maximal (peak) power production for specific PV panel at STC and N_p is a number of installed panels. The right side of the brackets defines the effect of PV panel efficiency linked to the cell temperature $T_c(k)$, and κ is the temperature coefficient of P_{\max} . As the PV cell temperature is hard to measure in real-time, it can be calculated as follows

$$T_c(k) = T(k) + \frac{G_d(k)}{G_{TC}} (T_{NOCT} - T_{TC}), \quad (3)$$

where T_{NOCT} is nominal operating cell temperature (NOCT) obtained at test condition (TC) $G_{TC} = 800 \text{ W/m}^2$, $T_{STC} = 20 \text{ }^\circ\text{C}$ and $T(k)$ is measured ambient temperature. It follows from the (2) that if the $T_c(k)$ is higher, then T_{STC} efficiency of the PV panel decrease by the coefficient κ .

Defined STC, TC, NOCT, P_{\max} , κ can be obtained from PV manufacturer. The number of PV panels, surface tiled, and azimuth angles are unique for each solar power plant installation. This section provides a brief overview of the model used. For more details, see [12].

C. Pearson Correlation Coefficient

Purpose of the correlation coefficients is to find and interpret how strong a relationship is between the investigated variables [13]. Pearson correlation coefficient is one of the most popular tools in the field of feature extraction for machine learning inputs with a large number of available variables. This coefficient represents a linear correlation between two continuous variables (x, y), and it is formulated as follows

$$r_{xy} = \frac{\sum_{i=1}^M (x_i - \bar{x}) \sum_{i=1}^M (y_i - \bar{y})}{\sqrt{\sum_{i=1}^M (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^M (y_i - \bar{y})^2}} \quad (4)$$

where r_{xy} represents correlation coefficient, \bar{x} denotes mean of x and \bar{y} denotes mean of variable y across all samples M . Value form (4) can come from a variance of $r_{xy} \in [-1, 1]$. The closer the coefficient gets to the 1, the higher the positive correlation is between x and y . Conversely, suppose the coefficient gets closer to the -1 . In that case, variables have a higher negative correlation, and if the r_{xy} is close to 0, there is no direct linear correlation between x and y .

III. PROBLEM STATEMENT

As described in Section II-A, power production can be curtailed by many factors despite solar power plant operation modes. The consequence of such a behavior is that the measured information from the smart devices provides incomplete information as other grid parts influence it. This creates unwanted data corruption, which is hard to identify and process before syntheses of the forecasting model. If data are not processed correctly model trained with such data provides underestimated PV production. It affects the utilization of solar power generation as its forecasts are used in model predictive control (MPC), which benefits from knowing the most accurate

predictions. Likewise, various design-making processes are dependent on the best PV forecasts.

This work suggests a novel curtailed data reconstruction approach for supervised machine learning modeling of PV panels production. Utilizing the solar panel physical model and Pearson correlation coefficient to identify and replace curtailed power production with estimated values of a maximal possible generation.

IV. PREDICTION OF ACHIEVABLE ENERGY

In general, the designed PV energy forecasting model includes weather forecasts as input producing the one-step prediction in the following form

$$y(k) = f_P(x(k)), \quad k = 0, \dots, N, \quad (5)$$

where y is prediction of generated power, x represents input features to the predictive model f_P for define time step k .

In the field of ML-based solar power prediction, the adopted methodology consists of standard procedures, including:

- 1) Data collection and pre-processing.
- 2) Identification and feature extraction.
- 3) Model selection (type, structure).
- 4) Model training and validation using weather measurements and historical PV production.
- 5) Testing procedure using weather forecasts and historical PV production to select best performance model.

The following section extends the pre-processing as it aims to avoid biased predictions and model underperforming. The training data are reconstructed to replace curtailed PV production with estimated achievable production.

A. Data Reconstruction

Assume we have acquired a representative sample of historical weather measurements $W \in \mathbb{R}^{M \times n_w}$, where n_w represents the number of unique variables $w \in \mathbb{R}^M$ and power production measurements $Y \in \mathbb{R}^M$ which includes corrupted information.

The first step in PV data reconstruction is to find the physical model (2) and identify its parameters. Using historical weather measurements and the designed model, we can construct dataset $Y_P \in \mathbb{R}^M$, which corresponds to the measured one. This data represents a naive guideline of achievable power generation, which is unaware of grid limitations. The curtailment of power production indicates that only interesting information from the modeled dataset Y_P is that which is greater than measured information from Y . Using this fact created dataset is modified as follows

$$Y_P(m) = y_p, \quad y_p = \begin{cases} Y_P(m), & \text{if } Y(m) < Y_P(m) \\ Y(m), & \text{otherwise} \end{cases}, \quad (6)$$

where $m = 1, \dots, M$ represents an index of the single value of the vector. In this way, it is possible to keep original measurements intact and separate potentially corrupted ones.

To correctly identify curtailed power production, it is necessary to investigate modeled data. Whether the mismatching between Y_P and Y is caused by effects of the connected grid

or exogenous influences of the weather conditions, which are not included in the model (2). Calculating the difference

$$\Delta Y = Y_P - Y, \quad (7)$$

and using it in Pearson guided correlation from Section II-C it is possible to determine if the weather measurements w from W significantly impact modelled and measured data mismatch (7). If the premise given is correct, it should be reflected in reconstructed data. Otherwise, we can assume that there was a curtailment in PV-produced energy.

The compensation of weather variables that are not modeled is based on the simple linear model in which coefficients are found as follows

$$\min_{a,b} J = \sum_{m=1}^M \left(\Delta Y(m) - (S(m)a^\top + b) \right)^2, \quad (8a)$$

$$\text{s.t.} \quad S(m)a^\top + b \leq \Delta Y(m), \quad (8b)$$

where $a \in \mathbb{R}^{n_s}$, $b \in \mathbb{R}$ represents coefficient of linear equation and $S \in \mathbb{R}^{M \times n_s}$ includes only dependent variables from W ($n_s \leq n_w$). The weather measurements are selected based on a simple rule. If the investigated weather variable $w^{(i)}$ where $i = 1, \dots, n_w$, has a greater correlation coefficient (4), then the variables from the model (2) (G_{DHI}, G_{DNI}, T) it is included in S . After the linear model is found, the final formulation of the resulting preprocess data is provided as

$$Y_f(m) = Y_P(m) - (S(m)a^\top + b), \quad m = 1, \dots, M \quad (9)$$

where Y_f represents reconstructed data, Y_P represents modeled data, and $Sa^\top + b$ is a compensation of exogenous weather influences acting on the PV system.

However, as mentioned in Section II-B, the meteorological impact on PV power production is highly non-linear. This means that a simple linear model provides unsatisfying results for a large number of historical samples M containing the different weather conditions (rainy, warm, cold days or changing seasons, etc.). Also, Pearson correlation is a measure of linear dependency between two datasets creating the same problem.

The key step is to investigate a small portion of the time series data $M \Rightarrow \{M_1, \dots, M_n\}$ at the time. This solution helps to decompose time series data as a piecewise linear process providing the desired results.

B. Implementation Details

The condition on lines (13, 16) can be interpreted as follows. Suppose the correlation of the exogenous weather variable $w^{(i)}$ with ΔY for the defined portion of data is smaller than the modeled variables (G_{DHI}, G_{DNI}, T). In that case, it is possible to assume that external weather influences did not create the difference ΔY . Otherwise, the linear model (9) will compensate their influence. This condition can be modified by tuning coefficient $p_r \geq 0$, which for non-zero values allows passing less dependant weather variables.

The range of the investigated data Δt should be selected based on the sampling time of the measured data W, Y and weather conditions of the PV panel's location. Chosen time

Algorithm: Data Reconstruction

Input: W, Y
Output: Y_f
Initialization:

1: Set $\Delta t \leftarrow \mathbb{R} \in (1, M), p_r \leftarrow 0$.

Initial Calculation:

2: Set $Y_{P_i} \leftarrow P_i(W)$.

Power Saturation :

3: **for** $m = 1 : M$ **do**

4: **if** $(Y_{P_i}(m) < Y(m))$ **then**

5: Set $Y_{P_i}(m) \leftarrow Y(m)$.

6: **end if**

7: **end for**

8: Set $\Delta Y \leftarrow Y - Y_{P_i}$.

Main Loop:

9: **for** $j = 1 : \Delta t : M - \Delta t$ **do**

10: Set $\delta \leftarrow \Delta Y(j : j + \Delta t), \phi \leftarrow Y_{P_i}(j : j + \Delta t)$ and $\omega \leftarrow W(j : j + \Delta t)$.

Calculate Correlation for Modeled Variables :

11: Set $r_{xy}^{DHI} \leftarrow f(G_{DHI}, \delta)$, $r_{xy}^{DNI} \leftarrow f(G_{DNI}, \delta)$ and $r_{xy}^T \leftarrow f(T, \delta)$.

12: **for** $i = 1 : m_w$ **do**

13: **if** $(\omega^{(i)} \neq \{G_{DHI}, G_{DNI}, T\})$ **then**

14: Set $r_{xy}^\omega \leftarrow f(\omega^{(i)}, \delta)$.

15: **end if**

16: **if** $(|r_{xy}^\omega| < \left\{ |r_{xy}^{DHI}|, |r_{xy}^{DNI}|, |r_{xy}^T| \right\} - p_r)$ **then**

17: Append $S \leftarrow \omega^{(i)}$.

18: **end if**

19: **end for**
Calculate Reconstructed Data with Compensation :

20: Find $a, b \leftarrow \operatorname{argmin} J(\delta, S)$.

21: Set $\sigma \leftarrow \phi - (Sa^T + b)$.

22: Append $Y_f \leftarrow \sigma$.

23: **end for**

range of investigated data is crucial as it is directly linked to linear models. Using a too large portion of the dataset at once (week, month, year) will lead to inaccurate correlations as the modeled variables are often the most dependent and directly affect power production. Consequences are such that algorithm will never or very rarely compensate possible exogenous effects, and if so, it will be inaccurate. On the contrary, selecting a too small portion of the dataset compared to the sampling time of the measurement will lead to too much compensation of modeled data, and we will end up with an underestimation of power production. It is worth mentioning that the range of the reconstructed data Δt may be different in each iteration. The proposed algorithm does not restrict that. If the data indicates it, it is even a recommended step. Another point is that provided physical model from Section II-B can take any complexity, form, or additional variables. The only difference will be that the S will be a smaller subset of W as the number of the exogenous variables will decrease with the modeled variables increase.

V. CASE STUDY

This work presents a real large-scale PV system on which we test the proposed algorithm and provide its reliability for any case of PV setup. Investigated PV installation is located in Romania and includes 11680 individual solar panels of the same efficiency and peak power P_{\max} . The system contains three master slave inverter control, each containing four separate inverters with individual smart meters. The primary inverter controls the other three, which are connecting and disconnecting based on the power production. The system

TABLE I: Solar panel model specification under the standard test condition: $G_{\text{STC}} = 1000 \text{ W/m}^2$, $T_{\text{STC}} = 25 \text{ }^\circ\text{C}$ and test conditions: $G_{\text{TC}} = 800 \text{ W/m}^2$, $T_{\text{TC}} = 20 \text{ }^\circ\text{C}$.

$\alpha_{\text{STA}} [^\circ]$	$\alpha_{\text{SAA}} [^\circ]$	$P_{\max} [\text{W}]$	$\kappa [\%/^\circ\text{C}]$	$T_{\text{NOCT}} [^\circ\text{C}]$
19.5	0	250	-0.44	20

is directly connected to the main grid without local power consumption, so it does not show frequent power curtailment. This makes it a perfect example as a proposed algorithm can be tested against true power production. Excluding power inverters data one by one from the aggregated dataset, we can simulate PV power curtailment, which is then reconstructed based on the Algorithm 1 shown in the previous section.

VI. RESULTS AND DISCUSSION

The data included in this work represents the total PV productions, weather measurements, and forecasts within 230 days with one hour sampling time $T_s = 1 \text{ h}$. Weather measurements and forecasts for defined location are imported from third-party open-source API Tomorrow.io, providing 20 different weather variables, including possible 430h prediction.

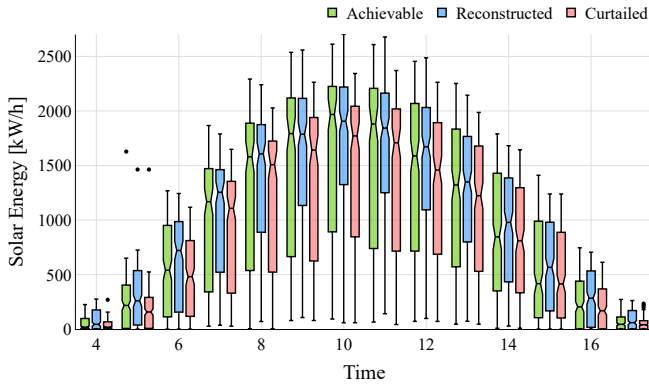
The total PV production is used in five different showcase simulations to provide representative results across a wide range of solar power curtailment. In the presented results, we apply the following scenario scheme:

TABLE II: Utilized data of process grope across all scenarios.

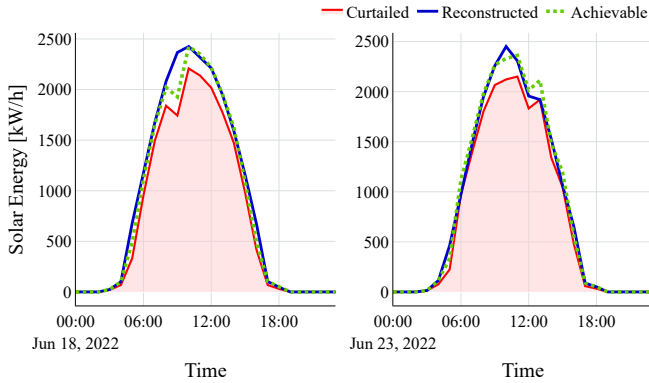
Grope	Case 1	Case 2	Case 3	Case 4	Case 5
1	100%	72.2%	72.2%	53.0%	72.2%
2	84.0%	100%	84.0%	100%	72.5%
3	100%	100%	100%	83.5%	72.1%

From now on, we will refer to the aggregated (the sum of all gropes) information of curtailed production, achievable production (without curtailment), and reconstructed data (curtailed data that have undergone a reconstruction process).

Each case is individually processed, as it contains unique curtailed PV measurements Y , which are then transformed to reconstructed dataset Y_f representing maximal possible power generation for a defined time. The range of the investigating data is chosen as $\Delta t = 24 \text{ h}$. By comparing these datasets Y, Y_f with achievable production Y_R , we generate results shown in Fig. 1 and Tab. III. The presented Fig. 1a shows a detailed



(a) Statistical results of reconstructed data compare to achievable and curtailed PV production for each hour of the daily active generation.



(b) Reconstructed data for 24 h of the day.

Fig. 1: Visualization of data reconstruction procedure (case 2).

TABLE III: Statistic results for reconstructed data Y_f compare to the curtailed Y and achievable Y_R PV power generation [kW/h] across all cases.

Case	P_{miss} [%]	Mean			Variance $\cdot 10^4$		
		Y_R	Y	Y_f	Y_R	Y	Y_f
1	5.3	498.7	472.1	504.8	53.5	46.2	49.5
2	9.2	498.7	453.0	499.8	53.5	44.7	49.8
3	14.5	498.7	427.4	483.9	53.5	38.3	45.8
4	21.1	498.7	395.7	474.1	53.5	32.9	44.3
5	27.8	498.7	363.4	461.5	53.5	28.9	43.0

statistical comparison across all measurements (case 2) for each hour of the active PV production (from 17 pm to 4 am is negligible power generation). As we can see, the mean of the reconstructed data is shifted toward the mean of the achievable PV production. As well as, the variance of reconstructed data is stretched compared to curtailed PV measurements, leading to better correlation with achievable power within every hour of the day. As the results suggest, the reconstruction procedure succeeded in mimicking the real measurements (achievable power), as we can see in the example Fig. 1b. The difference between reconstructed and achievable PV production is caused by a one-time change in the weather conditions, which are not captured as significant by Pearson correlation (4), so in

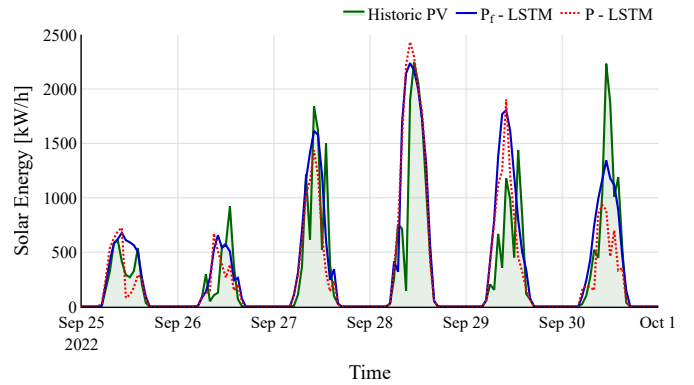


Fig. 2: Day-ahead PV forecast, using LSTM trained with reconstructed P_f and curtailed P PV data compared to achievable historical PV production (case 2).

conclusion, they are ignored by the reconstruction procedure.

In Tab. III, we can see statistic results for the whole dataset across all cases. Naturally, increases in the missing total produced power P_{miss} decrease the overall accuracy of the proposed algorithm. Looking at the mean and variance of the reconstructed data, we can observe a slower decline compared to curtailed PV production, which decreases at a much greater rate. Showing that the reconstruction procedure provides beneficial results as it scales with higher power curtailment.

To show the advantages of reconstructed data in PV power prediction, we design pair of the long-short term memory (LSTM see [14]) for every case from Tab. II. Each scenario is then represented by LSTM trained using curtailed PV measurement, and LSTM trained with reconstructed data to predict PV power P, P_f notated subsequently. We divide our datasets into three groups. First, the training dataset includes measurements within 138 days, and it is used to estimate the trainable parameters of the LSTMs. The second group represents validation data (23 days) for the model validation. Finally, the last group includes 69 days of measurements on which we perform accuracy analysis presented in Tab. IV and Fig. 2. The structure of the LSTM networks is selected as follows. Weather feature inputs are chosen using provided Algorithm 1 by including all modeled variables and those which were used in model compensation (9). In summary, inputs include (direct normal irradiance, diffuse horizontal irradiance, outside temperature, cloud base, dew point, humidity, precipitation intensity, wind speed, and wind direction). Ideal hyperparameters of the ML model are found using Bayesian optimization, the number of LSTM units is 543, the learning rate is selected as $3.715 \cdot 10^{-4}$. The 2900 epochs of minimizing the sum of squared errors were sufficient in order to find optimal weights and basis of the LSTM units with ADAM optimizer. The inputs to the LSTM are not changing, and outputs are not so much different in provided scenarios, so we fix these parameters for all ten recurrent neural networks.

In Fig. 2, we can see a day-ahead forecast (prediction for 24 h from the midnight of the previous day) comparison of

TABLE IV: Overall accuracy results for day-ahead forecasting of PV power generation [MW/h].

Case	P_{tot}	Maximal		MSE		Surplus		Deficit	
		ΔP	ΔP_f	P	P_f	P	P_f	P	P_f
1	372.5	4.4	2.0	220	86	10.6	19.1	42.6	18.7
2	372.5	4.4	2.1	238	91	9.0	17.3	47.1	21.9
3	372.5	4.6	2.6	329	135	6.3	12.1	60.8	33.5
4	372.5	4.9	3.2	612	201	4.3	11.4	92.3	44.3
5	372.5	5.1	3.4	716	243	3.1	10.6	101.1	49.7

constructed models with historical values for case 2. As shown, LSTM trained with reconstructed data is more reliable on days with higher power generation. For the days with lower production, both LSTMs provide similar results as the training data are not often curtailed in such small production. More detailed results can be found in Tab. IV. Where P_{tot} represents the total power produced without curtailment from our testing dataset, which is the same across all scenarios. From the mean squared error (MSE), we can declare that LSTM trained with reconstructed data provide superior results. The maximal absolute error between aggregated true daily production and the day-ahead forecast is shown in the third column. As we can see, LSTM trained with curtailed power production has almost twice larger maximal error $2\Delta P \approx \Delta P_f$. We chose the sum of daily aggregated surplus and deficit as the last indicator across all daily samples. As it already been pointed out, LSTM trained with reconstructed data naturally provides a larger surplus (overestimated PV prediction) compared to LSTM trained with achievable PV measurements, which on the contrary, provides a more significant deficit (underestimated PV prediction). However, as we can see, the sum of the deficit and surplus for both of the LSTM models compared to P_{tot} are proportionally different on a large scale showing the benefits of the proposed data reconstruction.

VII. CONCLUSION

This paper proposes a novel approach to photovoltaic (PV) power production measurement reconstruction for systems where PV generation occasionally exceeds load consumption. The methodology is based on constructing the physical model of a real PV system providing the baseline of maximal production. Investigating the curtailed PV measurements as a piecewise linear system, we were able to utilize the Pearson correlation coefficient and linear model of weather measurements. Both tools are used to compensate the difference between the modeled and curtailed PV power production. Using this approach, we reconstruct historical PV measurements and preserve the behavior of a non-linear system. The provided algorithm was tested on a real large-scale PV system, including 11680 individual panels. To investigate the performance of our scheme, we have simulated various levels of curtailment by artificially excluding information from a subset of PV inverters. In all scenarios algorithm successfully reconstruct corrupted (curtailed PV) measurements with some degree of accuracy, which is linked to the size of power curtailment.

The direct benefit of the proposed approach is improved machine learning (ML) based photovoltaic power production forecasts. Our study involves generating and comparing day-ahead predictions using ML models for a defined system, both with and without data reconstruction. From the investigation of different scenarios, we conclude that LSTM trained with reconstructed data provide superior results. Not only maximal error between the historical values and its predictions is almost twice lower. We also provide a good trade-off between the prediction deficit and surplus by possible tuning of the data reconstruction algorithm.

ACKNOWLEDGMENT

This research is funded by the European Commission under the grant no. 101079342 (Fostering Opportunities Towards Slovak Excellence in Advanced Control for Smart Industries). The authors gratefully acknowledge the contribution of the Scientific Grant Agency of the Slovak Republic under the grant 1/0490/23 and the Slovak Research and Development Agency under the project APVV-20-0261.

REFERENCES

- [1] S. Sen and S. Ganguly, "Opportunities, barriers and issues with renewable energy development – a discussion," *Renewable and Sustainable Energy Reviews*, vol. 69, pp. 1170–1181, 2017.
- [2] M. Anvari, G. Lohmann, M. Wächter, P. Milan, E. Lorenz, D. Heineemann, M. R. R. Tabar, and J. Peinke, "Short term fluctuations of wind and solar power systems," *New Journal of Physics*, vol. 18, no. 6, p. 063027, 2016.
- [3] G. Gust, T. Brandt, S. Mashayekh, M. Heleno, N. DeForest, M. Stadler, and D. Neumann, "Strategies for microgrid operation under real-world conditions," *European Journal of Operational Research*, vol. 292, no. 1, pp. 339–352, 2021.
- [4] N. Sunar and J. R. Birge, "Strategic commitment to a production schedule with uncertain supply and demand: Renewable energy in day-ahead electricity markets," *Management Science*, vol. 65, no. 2, pp. 714–734, 2019.
- [5] J. Antonanzas, N. Osorio, R. Escobar, R. Urraca, F. J. Martinez-de Pison, and F. Antonanzas-Torres, "Review of photovoltaic power forecasting," *Solar energy*, vol. 136, pp. 78–111, 2016.
- [6] C. Voyant, G. Notton, S. Kalogirou, M.-L. Nivet, C. Paoli, F. Motte, and A. Foulloy, "Machine learning methods for solar radiation forecasting: A review," *Renewable Energy*, vol. 105, pp. 569–582, 2017.
- [7] P. Mishra, A. Biancolillo, J. M. Roger, F. Marini, and D. N. Rutledge, "New data preprocessing trends based on ensemble of multiple preprocessing techniques," *TrAC Trends in Analytical Chemistry*, vol. 132, p. 116045, 2020.
- [8] M. Malvoni, M. G. De Giorgi, and P. M. Congedo, "Forecasting of pv power generation using weather input data-preprocessing techniques," *Energy Procedia*, vol. 126, pp. 651–658, 2017.
- [9] E. O'Shaughnessy, J. R. Cruce, and K. Xu, "Too much of a good thing? global trends in the curtailment of solar pv," *Solar Energy*, vol. 208, pp. 1068–1077, 2020.
- [10] M. Z. Liu, A. T. Procopiou, K. Petrou, L. F. Ochoa, T. Langstaff, J. Harding, and J. Theunissen, "On the fairness of pv curtailment schemes in residential distribution networks," *IEEE Transactions on Smart Grid*, vol. 11, no. 5, pp. 4502–4512, 2020.
- [11] N. Stringer, N. Haghdadi, A. Bruce, and I. MacGill, "Fair consumer outcomes in the balance: Data driven analysis of distributed pv curtailment," *Renewable Energy*, vol. 173, pp. 972–986, 2021.
- [12] S. A. Kalogirou, *Solar energy engineering: processes and systems*. Academic press, 2013.
- [13] H. Xu and Y. Deng, "Dependent evidence combination based on shearmen coefficient and pearson coefficient," *Ieee Access*, vol. 6, pp. 11 634–11 640, 2017.
- [14] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.