

A Truncated Gaussian Mixture Model for Distributions of Wind Power Ramping Features

Mingjian Cui, Zhenke Wang, Cong Feng, Jie Zhang
University of Texas at Dallas
Richardson, TX, 75080 USA

Abstract—Wind power ramps (WPRs) are significantly impacting the power balance of the system operations. Better understanding the statistical characteristics of ramping features would help power system operators better manage these extreme events. Toward this end, this paper develops an analytical truncated Gaussian mixture model (TGMM) to fit the probability distributions of different ramping features. The non-linear least square method with the Trust-Region algorithm is adopted to optimize the tunable parameters of mixture components; the optimal number of mixture components is adaptively solved by minimizing the Euclidean distance to the actual probability distribution. A sign function is utilized to truncate the original GMM distribution and obtain the final TGMM. The cumulative distribution function (CDF) of TGMM is analytically derived. Numerical simulations on publically available wind power data show that the parametric TGMM can accurately characterize the irregular and multimodal distributions of each ramping feature.

Index Terms—Probability distribution, statistical analysis, truncated Gaussian mixture model, wind power ramps.

I. INTRODUCTION

With the increase of wind power penetration in the power grid, the intermittence and fluctuation of wind power have drawn more and more attention, especially under extreme weather conditions. As a type of extreme events, wind power ramps (WPRs) have been investigated in recent studies [1], [2]. WPRs have a serious impact on the power balance of the system, and may lead to an instability of the power system frequency, load shedding, and other reliable operations. Statistical analysis of WPRs would help power system operators better understand the characteristics of ramping features, thereby assisting them to manage these extreme ramping events.

However, currently there are few studies in the literature focusing on accurately characterizing the parametric distributions of wind power ramping features, which are apt to be practically integrated into the power system scheduling models like the chance-constrained economic dispatch and unit commitment [3]. Sevlian *et al.* [4] characterized and analyzed ramping magnitude, duration, and rate by empirical distributions. But the empirical distribution is a step function with discrete (rather than continuous) probability values, which cannot be analytically expressed. Cui *et al.* [5] depicted the ramping feature statistics by using the kernel smoothing probability density (ksdensity) estimate. However, the ksdensity distribution was still a nonparametric model. Ganger *et al.* [6] utilized the Fréchet distribution (a generalized extreme value distribution) to fit the empirical wind power

ramping magnitude. However, the Fréchet distribution is a unimodal distribution that cannot accurately fit the multimodal distribution.

The Gaussian mixture model (GMM) has been widely used in the statistics community, and recently been applied in the renewable energy areas [7]–[9]. The GMM specializes in characterizing the multimodal and irregular probability distribution. Ke *et al.* [7] customized the GMM by three Gaussian functions and utilized the GMM to approximate the PDF of wind power generation with triple probability peaks. Singh *et al.* [8] represented all irregular probability distribution functions of load using GMM in various distribution system applications. Valverde *et al.* [9] proposed the use of GMM to represent non-Gaussian correlated wind power output and aggregated load demands for modeling the probabilistic load flow. Wind power ramps are high nonlinear and uncertain, and likely present multi-mode in the distribution of ramping features. Thus, this paper develops a GMM model to fit the PDF and CDF of different ramping features.

The developed GMM model is expected to accurately model the distribution of different ramping features, therefore being used in a variety of power system operations. The main contributions of this paper include: (i) developing a truncated GMM (TGMM) to fit the irregular and multimodal distributions of wind power ramping features; and (ii) deducing the analytical CDF expression of the TGMM.

The organization of this paper is as follows. In Section II, a wind power ramp extraction method using an optimized swinging door algorithm is briefly introduced. Section III presents the analytical expressions of the probability and cumulative distributions for the developed truncated Gaussian mixture model. Case studies and result analysis performed on publically available wind power data are discussed in Section IV. Concluding remarks and future work are discussed in Section V.

II. WIND POWER RAMPING FEATURES EXTRACTION

An optimized swinging door algorithm (OpSDA) [5] is used to detect all the wind power ramps in historical wind power data. A brief example of wind power ramps detection in one day is illustrated in Fig. 1. In the OpSDA, the swinging door algorithm with a predefined parameter, ε , is first applied to segregate the wind power data into multiple discrete segments. Then dynamic programming is used to merge adjacent segments with the same ramping direction and relatively high

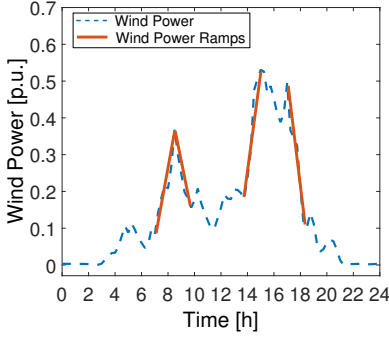


Fig. 1. An example of wind power ramps in one day.

ramping rate. A brief description of OpSDA is introduced here and more details can be found in [5]. Subintervals that satisfy the ramping rules are rewarded by a score function; otherwise, their score is set to zero. The current subinterval is retested as above after being combined with the next subinterval. This process is performed recursively to the end of dataset. Finally, significant wind power ramps with the maximum score are successfully extracted.

III. A TRUNCATED GAUSSIAN MIXTURE MODEL

A. Distribution of Ramping Features

Generally, wind power ramping features consist of ramping magnitude, duration, and change-rate. The statistic distribution of each ramping feature is significantly irregular and asymmetrical with multiple peaks. In addition, due to the definition of WRPs, the distributions of ramping features are truncated. For instance, if WPRs are defined by the threshold of the ramping magnitude without constraining the ramping duration and rate, both probability distributions of ramping magnitude and rate will be truncated. If WPRs are defined by the thresholds of both ramping magnitude and duration, the probability distributions of all three features will be truncated.

B. Truncated Gaussian Mixture Model

The ramping features are detected by the OpSDA based on a large wind power data set. To model the irregular and asymmetric distribution of ramping features, the Gaussian mixture model (GMM) is used and developed in this paper. The GMM model is a probabilistic model that assumes all the data points are generated from a mixture of a finite number of Gaussian distributions with multiple parameters. GMM performs better in irregular distribution models [7], [8], especially for the probability distribution with multiple peaks. The generalized GMM model is formulated as:

$$f_G(x; N_G; \omega_i, \mu_i, \sigma_i) = \sum_{i=1}^{N_G} \omega_i g_i(x; \mu_i, \sigma_i) = \sum_{i=1}^{N_G} \omega_i e^{-\left[\frac{x-\mu_i}{\sigma_i}\right]^2}, \forall x \in \mathcal{X}, \forall i \in \mathcal{I} \quad (1)$$

where \mathcal{X} is the data set of a ramping feature, x , with a total number of N_x . \mathcal{I} is the set of Gaussian mixture models (GMM) with a total number of N_G . A two-stage optimization model is constructed to calculate all the parameters of f_G , i.e., N_G , ω_i , μ_i , and σ_i . The first stage aims to determine the expected value (or mean value) vector M ($\mu_i \in M$), the standard deviation vector Σ ($\sigma_i \in \Sigma$), and the weight coefficient vector Ω ($\omega_i \in \Omega$), i.e., $f_G(x; N_G; \omega_i, \mu_i, \sigma_i) \rightarrow f_G(x; n_G)$. The non-linear least square method with the Trust-Region algorithm [10] is adopted in the first stage to obtain the parameters of the mixture components. The second stage aims to determine the optimal number of components, $N_{G,opt}$, by minimizing the Euclidean distance between the actual PDF, PDF_A , and the primary PDF of GMM, f_G , i.e., $f_G(x; N_G) \xrightarrow{N_{G,opt}} f_G(x)$. Thus, the objective function is formulated as:

$$\min \sqrt{\sum_{x \in \mathcal{X}} [f_G(x; N_G) - PDF_A]^2} \quad (2)$$

A sign function, $sign(x)$, is utilized to truncate the original PDF function of GMM, $f_G(x)$, given by:

$$sign(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases} \quad (3)$$

Then the final PDF of TGMM, $f_{TG}(x)$, can be analytically formulated as:

$$f_{TG}(x) = f_G(x) \times sign(x - Tr) \quad (4)$$

where Tr is the threshold for defining wind power ramps. In this paper, a wind power ramp is defined as the change in wind power output larger than 20% of the rated capacity without constraining the ramping duration and rate. The threshold of ramping magnitude, Tr_M , equals 0.2. The threshold value of ramping duration, Tr_D , equals 0. The threshold of ramping rate, Tr_R , is calculated by $Tr_M / (max(Dr))$, where $max(Dr)$ represents the maximum value of ramping duration.

C. Analytical Expression of the CDF of TGMM

The cumulative distribution function (CDF), F_G , is another essential statistic metric to analyze ramping features due to its monotonicity, which can be analytically expressed as:

$$\begin{aligned} F_G(x; N_G; \omega_i, \mu_i, \sigma_i) &= \int_{-\infty}^x f_G(t; N_G; \omega_i, \mu_i, \sigma_i) dt \\ &= \int_{-\infty}^x \sum_{i=1}^{N_G} \omega_i e^{-\left[\frac{t-\mu_i}{\sigma_i}\right]^2} dt \\ &= \sum_{i=1}^{N_G} \left[\frac{\sqrt{\pi}}{2} \omega_i \sigma_i \operatorname{erf}\left(\frac{\mu_i - x}{\sigma_i}\right) \right] + C \end{aligned} \quad (5)$$

where erf is the Gaussian error function and defined as:

$$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt \quad (6)$$

Equation (5) is an indefinite integral with a constant C , which can be solved by (7). Since the detected ramping

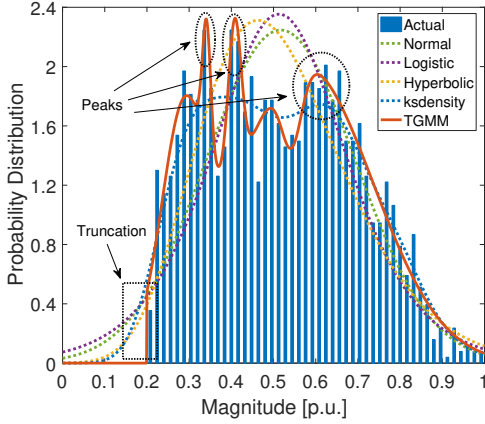


Fig. 2. PDFs of ramping magnitude of five distribution models.

magnitudes are normalized into the range $[0, 1]$, it can be derived that $F_G(x < 0) = 0$ and $F_G(x > 1) = 1$. Hence, we use $x = -0.1$ ($F_G(-0.1) = 0$) to obtain the constant, C , in (5), shown as:

$$C = F_G(-1.1) - \sum_{i=1}^{N_G} \left[\frac{\sqrt{\pi}}{2} \omega_i \sigma_i \operatorname{erf}\left(\frac{-1.1 - \mu_i}{\sigma_i}\right) \right] \quad (7)$$

Considering the sign function in (3), the final CDF of TGMM, $F_{TG}(x)$, can also be analytically formulated as:

$$F_{TG}(x) = F_G(x) \times \operatorname{sign}(x - \operatorname{Tr}) \quad (8)$$

IV. CASE STUDIES

A. Test Case and Benchmarks

The proposed truncated Gaussian mixture model is evaluated and analyzed based on the Wind Integration National Dataset (WIND) Toolkit [11]. The WIND Toolkit data represent wind power generation and forecasts spanning from January 1st 2007 to December 31st 2012, sampled every 5 minutes. 711 wind sites located around Dallas, Texas are selected for this case study. The rated capacity is 9,987 MW. The total number of samples is 631,296, which is sufficiently large for statistical analysis of ramping features. There are 1,586 wind power ramps detected by OpSDA in total. The door width of OpSDA is set as 5% of the rated capacity.

For comparison, one nonparametric and three parametric distributions widely used in statistical analysis are selected in the case study. The normal distribution has been widely used to design the random number generator and generate load forecasting errors [12]. The logistic distribution is used for growth models in logistic regression and has longer tails and a higher kurtosis than the normal distribution [13]. The hyperbolic distribution has been used to accurately analyze and characterize wind and load forecasting errors [14]. The nonparametric kernel smoothing density (ksdensity) distribution has been widely used in the wind speed distribution characterization and renewable energy forecasting [15].

B. Metrics for Evaluating the Fitting Performance

To compare the performance of TGMM with different distribution models, a suit of widely used metrics in the wind power forecasting community are adopted to assess the distribution accuracy [15]. These metrics include correlation coefficient, normalized root mean square error (NRMSE), maximum absolute error (MaxAE), mean absolute error (MAE), Kolmogorov-Smirnov test integral (KSIPer), standard deviation, and fourth root mean quartic error (4RMQE). The correlation coefficient is a measure of the correlation between the actual PDF and the PDF of fitting distribution models. NRMSE is suitable for evaluating the overall accuracy of the fit while penalizing large fitting errors in a square order. MaxAE is suitable for evaluating the largest fitting error. MAE is suitable for evaluating uniform fitting errors. KSIPer evaluates the statistical similarity between the actual PDF and the PDF of fitting distribution models. Standard deviation quantifies the uncertainty of the fit. 4RMQE is suitable for evaluating the overall accuracy of the fit while penalizing large fitting errors in a quartic order. A smaller value indicates a better forecast for most of the metrics, only except for the correlation coefficient. Detailed information about the metrics can be found in [15].

C. Statistical Comparisons of Ramping Magnitudes

In this case, seven Gaussian components are found to accurately fit the distribution of ramping magnitude. There are totally 21 parameters (3×7) in the TGMM distribution model, which are optimized and listed in Table I. Fig. 2 compares the probability of the actual histogram distribution and five distribution models. For the actual distribution, there are three peaks located around 0.32 p.u., 0.41 p.u., and 0.65 p.u.. It means that the wind power output changes will occur at these three values in a high probability, which is informative and could be used in power system operations. For example, ramping reserve requirements could be designed by considering these three peak values instead of only one peak. Since the normal, logistic, and hyperbolic models conform to the unimodal distribution, only one single peak is depicted with the highest probability. For both the normal and logistic distributions, the peak values are 0.5 p.u.. For the hyperbolic distribution, the peak value is about 0.45 p.u.. This ill-information may mislead power system operators to mainly focus on coping with the WPRs with magnitudes spanning from 0.45 p.u. to 0.5 p.u., and neglecting other significant wind power ramps around 0.32 p.u., 0.41 p.u., and 0.65 p.u.. Though the nonparametric model, ksdensity, is well-known in fitting the irregular distributions, it presents a worse performance than the TGMM from visual inspection. Besides, the nonparametric nature of ksdensity restricts its application in practice. This is because the analytical expressions of both PDF and CDF of the distribution model are generally required in stochastic power system operations, such as chance-constrained constraints in economic dispatch or unit commitment.

Another interesting finding in fitting the distribution is the truncation part in the left tail area in Fig. 2. Due to

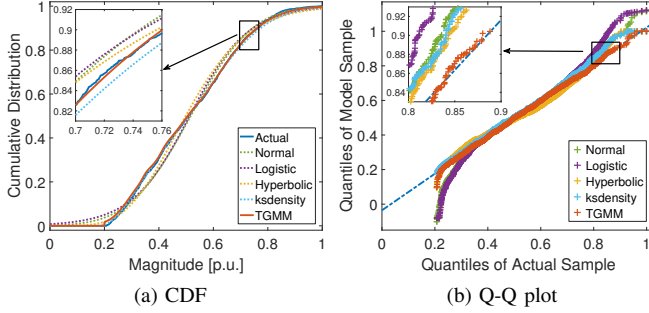


Fig. 3. Cumulative distributions and Q-Q plots of ramping magnitude using five distribution models.

TABLE I
PARAMETERS OF TGMM WITH SEVEN COMPONENTS ($N_{G,opt} = 7$)

Number of Components	ω	μ	σ
1	1.0427	0.4081	0.0217
2	0.8444	0.3413	0.0138
3	-0.0001	0.6327	7.11×10^{-4}
4	-0.0001	0.5378	2.22×10^{-14}
5	1.3551	0.2851	0.0719
6	-0.5211	0.5427	0.0342
7	2.0049	0.5725	0.2312

the definition of WPRs (20% of the rated capacity), all the ramping magnitudes are greater than or equal to 0.2 p.u.. For the ramping magnitude that is less than 0.2 p.u., the occurrence probability should be zero. Under this circumstance, the TGMM distribution performs much better than any other distributions due to the truncation process, which makes the fitting distribution of the TGMM more realistic.

For quantitative comparison, Table II lists the fitting metrics for different distribution models. Regarding the correlation coefficient metric in green, the TGMM shows the largest value. Regarding other metrics in blue, the TGMM shows the smallest value. This indicates that the TGMM outperforms other parametric distributions, and even performs better than the nonparametric distribution, ksdensity, as a parametric model.

Fig. 3a and Fig. 3b compare the performance of the CDF and Q-Q plot for different distributions. Both the CDF and Q-Q curve of the TGMM distribution fit the actual curve better than other four distributions.

D. Statistical Comparisons of Ramping Duration and Rate

In addition to ramping magnitude, ramping duration and rate are another two important ramping features. In this section, the probability distributions of ramping duration and rate are also characterized and analyzed by using the TGMM.

Fig. 4a compares the probability distribution of ramping duration by using five distribution models. There are five Gaussian components that optimally fit the distribution of ramping duration. Due to the irregular and asymmetric characteristics of the ramping duration distribution, the parametric models of the normal, logistic, and hyperbolic distributions fail to track

TABLE II
METRICS VALUES ESTIMATED FOR RAMPING MAGNITUDE

Metrics	Distribution Models				
	Normal	Logistic	Hyper.	ksdensity	TGMM
Correlation coefficient	0.88	0.84	0.89	0.96	0.98
NRMSE	0.16	0.18	0.16	0.09	0.06
MaxAE	0.41	0.44	0.45	0.22	0.18
MAE	0.12	0.14	0.11	0.07	0.04
KSIPer	33.16	42.98	30.14	23.73	10.28
Standard dev.	0.16	0.18	0.16	0.09	0.06
4RMQE	0.22	0.24	0.21	0.12	0.09

TABLE III
METRICS VALUES ESTIMATED FOR RAMPING DURATION

Metrics	Distribution Models				
	Normal	Logistic	Hyper.	ksdensity	TGMM
Correlation coefficient	0.95	0.94	0.98	0.99	0.99
NRMSE	0.10	0.11	0.07	0.05	0.05
MaxAE	0.29	0.29	0.18	0.14	0.11
MAE	0.07	0.07	0.05	0.03	0.03
KSIPer	11.37	18.13	17.70	12.06	10.24
Standard dev.	0.10	0.11	0.07	0.05	0.05
4RMQE	0.14	0.15	0.10	0.07	0.06

TABLE IV
METRICS VALUES ESTIMATED FOR RAMPING RATE

Metrics	Distribution Models				
	Normal	Logistic	Hyper.	ksdensity	TGMM
Correlation coefficient	0.92	0.95	0.98	0.99	0.99
NRMSE	0.11	0.09	0.04	0.03	0.03
MaxAE	0.31	0.28	0.16	0.11	0.07
MAE	0.07	0.06	0.02	0.02	0.02
KSIPer	21.26	9.10	6.54	6.29	5.29
Standard dev.	0.11	0.09	0.04	0.03	0.03
4RMQE	0.16	0.12	0.07	0.05	0.04

the actual probability values very well, especially for the peak values. However, the TGMM and the nonparametric model can fit most probability values of the actual distribution. This is specifically illustrated in Table III with the numerical metrics. Both the TGMM and the nonparametric ksdensity model show better performance than other distribution models, and the TGMM provides equal-to-better performance comparing to the nonparametric ksdensity model. It is noted that the TGMM performs much better than the nonparametric model in terms of the KSIPer indicator. It means the TGMM can show more statistical similarity to the actual histogram distribution, comparing to the nonparametric model.

The probability model of ramping rate is also a truncated distribution due to the truncated ramping magnitude distribution. The truncated PDF of ramping rate is illustrated in Fig. 4b, where the truncation threshold is 3.54 MW/min. Eight Gaussian components are found to optimally fit the actual histogram distribution of ramping rate. It is shown that the

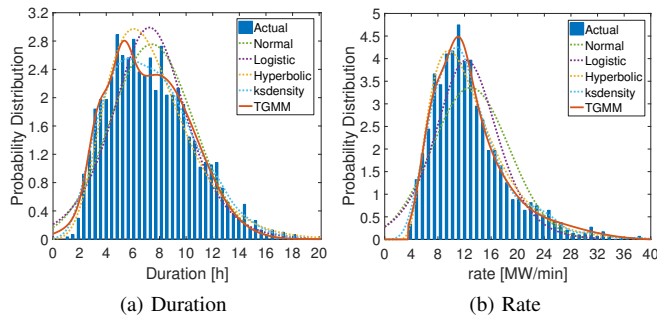


Fig. 4. Probability distributions of ramping duration and rate using five distribution models.

TGMM characterizes the peak of probability better than other distributions. Table IV lists the metrics for the fitting performance of different distribution models. Among all parametric models, the TGMM performs better than the normal, logistic, and hyperbolic models for all metrics. Comparing with the nonparametric model, the TGMM can also provide an equal-to-better performance in an analytical way.

Moreover, comparing to the ramping magnitude distribution simulation in Section IV-C, it is shown that the TGMM performs significantly well at fitting the probability and cumulative distributions of ramping magnitude. For the ramping duration and rate, the TGMM can provide much better performance than unimodal models (i.e., the normal, logistic, and hyperbolic), and equal-to-better performance comparing to the nonparametric model, ksdensity.

V. CONCLUSION

This paper developed a truncated Gaussian mixture model (TGMM) to characterize the probability and cumulative distributions of wind power ramping features. The TGMM was analytically expressed as a parametric form. First, the non-linear least square method with the Trust-Region algorithm was adopted to estimate all the parameters of mixture components. Second, the optimal number of mixture components was adaptively solved by minimizing the Euclidean distance to the actual probability distribution. Finally, the sign function was utilized to truncate the original GMM distribution and obtain the developed TGMM. Moreover, the cumulative distribution function (CDF) of TGMM was analytically deduced. Numerical simulations on the publically available wind power data showed that:

- (i) The TGMM distribution could optimally fit the actual probability and cumulative distributions of ramping features. All the evaluation metrics presented the best performance of the TGMM comparing to both the unimodal parametric models and the nonparametric model.
- (ii) Regarding the multimodal distribution of ramping magnitude, the TGMM performed significantly better than the normal, logistic, hyperbolic, and ksdensity models, especially when multiple peaks present in the distribution.

- (iii) Regarding the unimodal and asymmetric distributions of ramping duration and rate, the TGMM provided an equal-to-better performance comparing to the nonparametric model, and much better performance than the normal, logistic, and hyperbolic models.

In the future, this research can be further improved by: (i) using the analytically developed TGMM in the chance-constrained scheduling of power system operation models; and (ii) applying to multiple wind farms on different geographic locations.

ACKNOWLEDGMENT

This work was supported by the National Renewable Energy Laboratory under Subcontract No. XGJ-6-62183-01 (under the U.S. Department of Energy Prime Contract No. DE-AC36-08GO28308).

REFERENCES

- [1] Y. Qi and Y. Liu, "Wind power ramping control using competitive game," *IEEE Trans. Sustain. Energy*, vol. 7, no. 4, pp. 1516–1524, Oct. 2016.
- [2] M. Cui, D. Ke, Y. Sun, D. Gan, J. Zhang, and B.-M. Hodge, "Wind power ramp event forecasting using a stochastic scenario generation method," *IEEE Trans. Sustain. Energy*, vol. 6, no. 2, pp. 422–433, Apr. 2015.
- [3] H. Wu, M. Shahidehpour, Z. Li, and W. Tian, "Chance-constrained day-ahead scheduling in stochastic power system operation," *IEEE Trans. Power Syst.*, vol. 29, no. 4, pp. 1583–1591, Jul. 2014.
- [4] R. Sevljan and R. Rajagopal, "Detection and statistics of wind power ramps," *IEEE Trans. Power Syst.*, vol. 28, no. 4, pp. 3610–3620, Nov. 2013.
- [5] M. Cui, J. Zhang, A. R. Florita, B.-M. Hodge, D. Ke, and Y. Sun, "An optimized swinging door algorithm for identifying wind ramping events," *IEEE Trans. Sustain. Energy*, vol. 7, no. 1, pp. 150–162, Jan. 2016.
- [6] D. Ganger, J. Zhang, and V. Vittal, "Statistical characterization of wind power ramps via extreme value analysis," *IEEE Trans. Power Syst.*, vol. 29, no. 6, pp. 3118–3119, Nov. 2014.
- [7] D. Ke, C. Chung, and Y. Sun, "A novel probabilistic optimal power flow model with uncertain wind power generation described by customized gaussian mixture model," *IEEE Trans. Sustain. Energy*, vol. 7, no. 1, pp. 200–212, Jan. 2016.
- [8] R. Singh, B. C. Pal, and R. A. Jabr, "Statistical representation of distribution system loads using gaussian mixture model," *IEEE Trans. Power Syst.*, vol. 25, no. 1, pp. 29–37, Feb. 2010.
- [9] G. Valverde, A. Saric, and V. Terzija, "Probabilistic load flow with non-gaussian correlated random variables using gaussian mixture models," *IET Gener. Transm. Distrib.*, vol. 6, no. 7, pp. 701–709, 2012.
- [10] J. J. Moré and D. C. Sorensen, "Computing a trust region step," *SIAM J. Sci. Statist. Comput.*, vol. 4, no. 3, pp. 553–572, 1983.
- [11] C. Draxl, A. Clifton, B.-M. Hodge, and J. McCaa, "The wind integration national dataset (WIND) Toolkit," *Appl. Energy*, vol. 151, pp. 355–366, 2015.
- [12] Y. V. Makarov, C. Loutan, J. Ma, and P. De Mello, "Operational impacts of wind generation on california power systems," *IEEE Trans. Power Syst.*, vol. 24, no. 2, pp. 1039–1050, May 2009.
- [13] MathWorks. Logistic distribution. [Online]. Available: <https://www.mathworks.com/help/stats/logistic-distribution.html>.
- [14] B.-M. Hodge, D. Lew, and M. Milligan, "Short-term load forecasting error distributions and implications for renewable integration studies," in *Proc. IEEE 5th Green Technol. Conf.*, Denver, CO, USA, 2013, pp. 435–442.
- [15] J. Zhang, A. Florita, B.-M. Hodge, S. Lu, H. F. Hamann, V. Banunaryanan, and A. M. Brockway, "A suite of metrics for assessing the performance of solar power forecasting," *Solar Energy*, vol. 111, pp. 157–175, Jan. 2015.