

# Differentially Private Anomaly Detection for Energy Market Data

Ana Paixão<sup>1</sup>, Breno da Silva<sup>2</sup>, Filipe H. Cardoso<sup>2</sup>, Alexandre Braga<sup>2</sup>

<sup>1</sup> Institute of Computing, University of Campinas, SP, Brazil.

<sup>2</sup>CPqD - Centro de Pesquisa e Desenvolvimento, Campinas, SP, Brazil.

a272447@dac.unicamp.br, {brenos, filipehc, ambraga}@cpqd.com.br

**Abstract.** *Differential privacy (DP) is often used to protect consumer privacy in controlled data disclosures. This paper applies DP technology in anomaly detection to safeguard the secrecy of business-sensitive data in Brazilian open energy markets, where full metering datasets are publicly disclosed. Anomaly detection on actual metering data can be found wrong when random noise preserves secrecy, but reduces data utility. This paper contributes to solve this real-world issue by demonstrating DP's utility for anomaly detection in energy consumption data. We evaluated this approach across various values for the privacy parameter ( $\epsilon$ ), analyzing Precision, Recall, and F1-Score metrics. Our findings facilitate fine-tuning the trade-off between anomaly detection and business secrecy, minimizing the risk of inaccurate insights from noisy data while ensuring robust privacy.*

## 1. Introduction

Historically, security concerns at energy sector have been related to availability of infrastructures for generation and distribution. Security requirements addressing these concerns were implemented by mechanisms for detection, prevention and reaction to disruption events. As energy sector evolves and incorporates Information and Communication Technologies (ICTs) into its operations, innovation in business models flourishes, but cyber-threats and privacy violations became serious issues. An open energy market is a market where consumers can choose their electricity supplier from a range of different providers. Brazilian open energy market is relatively new and is looking for transparency and reliability in operations. To promote public audit and accountability, whole metering datasets have been published<sup>1</sup>, possibly disclosing business sensitive information. Meanwhile, there is a growing need for providing business secrecy to traders regarding how their clients consume energy and change suppliers over time.

A problem may arise because, in this publicly accessible dataset, highly sensitive information such as client portfolios and market shares can be easily inferred by competitors, as actual metering data for business clients is disclosed in its entirety, without any anonymization, to the public. This exposure allows malicious third-parties or competitors to analyze detailed consumption patterns and track supplier changes, gaining undue insights into rival business strategies and market positions.

Differential privacy is a privacy preserving technology designed to protect an entity's privacy in large datasets [Dwork 2006]. It ensures, with a certain degree of confidence, that adding or removing a single entity's data from a dataset has minimal impact on the overall result of a statistical query or analytical output [Dwork 2006]. In this paper, we use differential privacy technology in an innovative way to protect trader's business sensitive information within

---

<sup>1</sup><https://dadosabertos.ccee.org.br>

metering data for open energy markets, preserving business secrecy. We use publicly available data from Brazilian energy retail market. In this dataset, client portfolios and market shares can be easily inferred by competitors, because actual metering data for business clients is disclosed in the whole, without anonymization, to the public.

The main contribution of this paper lies in investigate a possible solution to the real-world issue of apply differential privacy to protect the secrecy of business-sensitive information in open energy markets while still enabling valuable data analysis. We explore the application of differential privacy in complex analytical tasks crucial for market operations and oversight, specifically the actual scenario of anomaly detection in consumption metering. This work investigates the viability and performance of Differentially Private KMeans (DPKMeans) for identifying anomalous energy consumption patterns within these sensitive datasets. We provide a detailed evaluation of its effectiveness across various privacy levels ( $\epsilon$ ) by examining standard performance metrics (Precision, Recall, and F1-Score) and the behavior of true/false positives/negatives derived from the confusion matrix. Our findings aim to guide the practical implementation of such techniques, allowing for useful, privacy-preserving market analytics that can benefit competing traders, retailers, and third-party analysts by enabling data-driven insights without compromising confidential business intelligence.

A distinctive aspect of this investigation is its foundation on actual, publicly available metering data from the Brazilian open energy market. This grounding in a real-world operational context allows for an assessment of DP-KMeans' practical utility in detecting anomalies within authentic energy consumption patterns. The paper is organized as follows. Section 2 contains related work. Section 3 overviews Brazilian open energy market. Section 4 discusses dataset disclosure in Brazilian open energy market and illustrates the risks of inaccurate insights with differentially private histograms. Section 5 details our methodology using differentially private clustering for anomaly detection on metering datasets and presents its performance evaluation. Section 6 concludes this text.

## **2. Related Work**

Differential privacy was proposed in 2006 in a series of three papers [Dwork 2006, Dwork et al. 2006b, Dwork et al. 2006a]. First, [Dwork 2006] shows that semantic security cannot be achieved with absolute privacy and proposes differential privacy to capture the risk of data leaks for someone present in a database subject to queries. [Dwork et al. 2006b] explains that privacy is protected when the true response from a database query is perturbed by adding random noise generated according to a carefully chosen distribution, and this response (with added noise), is returned to the user. This way, privacy can be preserved by calibrating noise's standard deviation according to information's desired sensitivity. Last, [Dwork et al. 2006a] explains that privacy can also be achieved by perturbing the true response of query by adding a small amount of exponentially distributed noise.

Differential privacy has been used in various application scenarios. At energy sector, proprietary schemes for differential privacy have been used in smart grids targeting residential consumers [Zhao et al. 2014, Peralta-Peterson and Kotevska 2021, Marks et al. 2021, Janghyun et al. 2022] and time series [Leukam Lako et al. 2021, Roman et al. 2021, Roman 2023, McElroy et al. 2023, Shaham et al. 2024] to preserve privacy, and recent work has also focused on evaluating open-source libraries for these purposes (e.g., [Paixão et al. 2025a, Paixão et al. 2025b]). Other application domains include, for instance,

anomaly detection with a focus on public health surveillance, such as in epidemic outbreak detection [Fan et al. 2013]. Also, government agencies have proposed guidelines for evaluating privacy guarantees [Near et al. 2023], as well as risk analysis and hardening guidelines for forecasting demand on electricity grids [ENISA 2023].

All these works, in way or another, contribute to set the stage where government agencies are more concerned than ever about the effectiveness of privacy preserving technology, giving appropriate conditions to evaluated differential privacy guarantees in specific scenarios.

### 3. Brazilian Open Energy Market in Electrical Smart Grids

The Advanced Metering Infrastructure (AMI) is a subsystem of modern power supply networks (a.k.a Electrical Smart Grids) that transmits energy consumption metering and other data, commands, and controls, to network operators. Other participants in this smart grid are *consumers units* (residences, buildings, offices, stores and factories where a smart meter is installed) and *service providers* (energy generators, transmitters and distributors). Open market players are relatively new to Brazilian smart grids and can act as aggregators, brokers, traders, retailers and even data analysts.

In open energy markets, consumers can choose their electricity supplier from a range of different providers. This allows consumers to compare prices and find the best deal for their needs. In Brazil, CCEE (*Câmara de Comercialização de Energia Elétrica*) is a non-profit organization responsible for facilitating commercialization of electricity within the National Interconnected Energy System [CCEE 2024] and its open energy market. CCEE grants that consumption metering data transacted by traders, retailers and distributors are disclosed in the whole to the public, without anonymization, since march 2024.

In Brazilian open energy market, metering aggregators, traders, retailers and outsourced analysts may have different access needs, which may not include viewing other player's data. For instance, trustful aggregators usually grant traders and retailers access to fined-grained metering on energy consumption. Despite perfectly justifiable from a business point of view, open access to detailed metering can potentially jeopardize business secrecy by disclosing not only clients consumption patterns, but also business sensitive information. In principle, traders and retailers need access to their own data to perform various analysis in decision making processes, such as volume analysis, dynamic pricing, energy type aggregations, billing, tax charging, load forecasting, and usage variance, among other needs. However, a malicious third-party or competitor may use open access to perform queries on other's metering data or tracking of retailer membership transactions and, through union of results and external data, infer consumption patterns, customer portfolios and market shares for competitors.

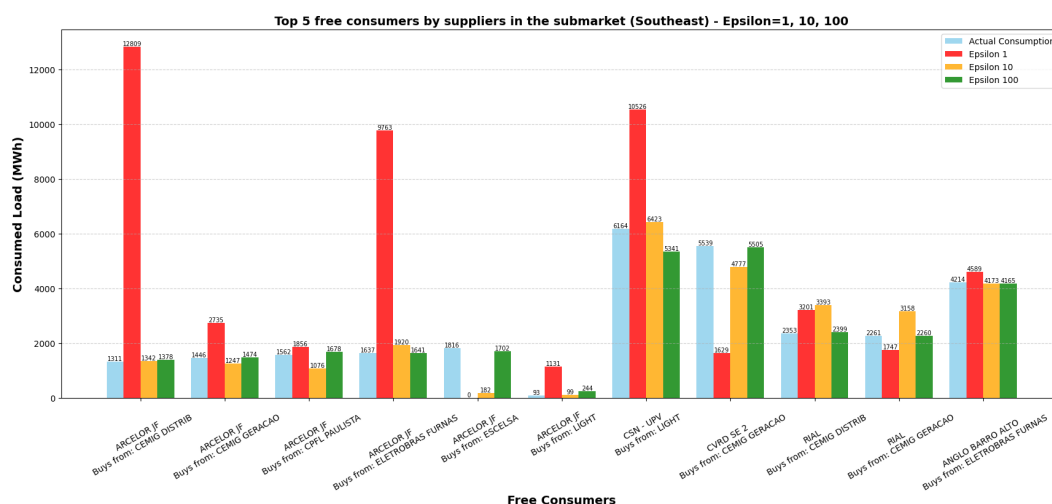
### 4. Differentially Private Histograms and the Risk of Inaccurate Insights

This section uses two examples to highlight the importance of carefully selecting an appropriate epsilon value that minimizes the risk of inaccurate insights while still providing adequate privacy protection. The first uses differentially private histograms and the second uses private average.

First, for illustrative purpose only, Figure 1 shows the impact of differential privacy on energy consumption data for the top 5 free consumers in the Southeast submarket. It compares actual consumption against data protected using three different privacy levels (three  $\epsilon$  values: 1, 10, and 100). In Figure 1, x-axis identifies the consumers and their energy suppliers, while

y-axis shows the consumed load in MWh. In legend, with "Actual Consumption" representing the baseline and the other three representing progressively weaker privacy levels (and thus, less noise added to the data), allowing comparison of consumption levels across the different privacy settings for each consumer. Visually, we can observe how increasing  $\epsilon$  values (from 1 to 100) results in data points that more closely approximate the actual consumption, thereby directly illustrating how a strong privacy budget (low epsilon) can lead to significantly distorted data, posing a substantial risk of inaccurate insights in market analysis.

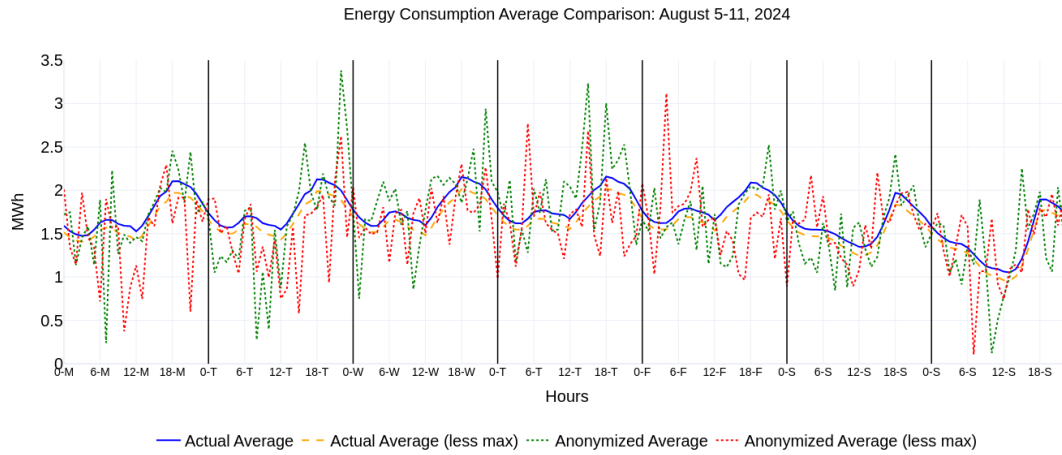
The chart clearly demonstrates the core challenge of differential privacy: balancing data utility with individual privacy. While stronger privacy (lower epsilon) adds more noise and thus protects individual consumer data more effectively, it also distorts the consumption figures, potentially impacting the accuracy of market analysis and decision-making. Conversely, weaker privacy (higher epsilon) provides data closer to the actual values but compromises privacy to a greater degree. The varying impact of noise across different consumers also raises questions about fairness and the potential for skewed market perceptions. This highlights the critical need to identify an optimal privacy budget that minimizes the risk of generating misleading insights while ensuring sufficient protection for sensitive business information.



**Figure 1. Top consumers by suppliers in Brazilian southeast submarket. Actual histogram is compared with anonymised histograms for three different  $\epsilon$  values.**

In a second example, Figure 2 shows a week-long time series of consumption averages: Actual averages (blue line) and actual averages minus one consumer unit (yellow dotted line). The gap between these two lines is the energy consumption of a missing consumer unit. Once consumption patterns are revealed, a customer behavior could be tracked even when she changes energy supplier. Differential privacy can be used to anonymise single consumers in general statistics, making then indistinguishable from each other. As an example, Figure 2 also shows randomized versions of actual averages (green line) and actual averages minus one consumer (red dotted line). The presence or absence of one consumer unit is easily distinguishable from actual values (green and red lines), but should be indistinguishable in differentially private averages. The challenge facing differential privacy in energy metering is to balance secrecy and utility by finding the right amount of random noise to be added to

a time series that preserves business secrecy, while allowing useful analytics.



**Figure 2. Time series of average energy consumption by hour.**

## 5. Differentially Private Anomaly Detection on Metering Datasets

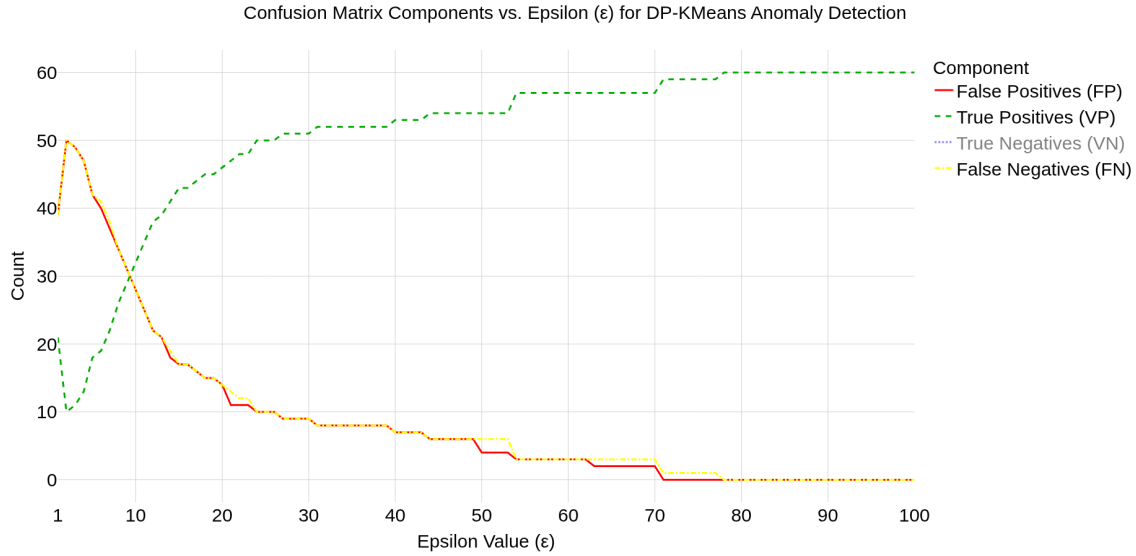
Many practical applications within open energy markets need sophisticated analytical capabilities. A critical example is anomaly detection, aimed at identifying unusual consumption patterns or outliers. Such anomalies can point to various important events, including equipment malfunctions, potential energy theft, significant deviations from contractual energy usage, data integrity problems, or novel market behaviors of interest to both market participants and system operators. Addressing the challenge of performing these advanced analyses without compromising the business-sensitive nature of the data, as highlighted in Section 4, is paramount.

This section investigates the application of Differentially Private KMeans (DP-KMeans) for anomaly detection in energy metering data. The anomaly detection methodology employs standard KMeans (from Scikit-learn) and Differentially Private KMeans (DP-KMeans, from ‘diffprivlib’ [Holohan et al. 2019]). Both algorithms were configured with  $k=3$  clusters and a `random_state=42` for reproducibility; standard KMeans also used `n_init=10`. DP-KMeans utilized data-derived bounds and varying  $\epsilon$  values to analyze the privacy-utility trade-off.

To illustrate the approach and evaluate its performance, our analysis utilizes a dataset comprising 31 days of consumption data from August 2024, focusing on a representative segment of commercial consumers – specifically, selected shopping malls within the Southeast submarket of the Brazilian open energy market (described in Section 3). While these data are publicly available, the consumption patterns of commercial entities, including shopping malls, can still reveal sensitive business insights such as operational schedules, peak activity periods, and even strategic changes in energy management, which could be exploited by competitors. We selected shopping malls as our running example and use case because most people have a good understanding of how shopping malls work.

Figure 3 details the evolution of the raw counts of the confusion matrix components — True Positives (VP), False Positives (FP), True Negatives (VN), and False Negatives (FN) — as a function of  $\epsilon$ . Figure 3 reveals that for very low  $\epsilon$  values (e.g.,  $\epsilon = 1$ ), the number of False Positives (FP) (orange line) is significantly high, starting at approximately 58. This large

count of FPs decreases sharply as  $\epsilon$  increases, approaching near-zero values for  $\epsilon > 45$ . This behavior is the primary explanation for the low Precision observed at low  $\epsilon$  values in Figure 4.



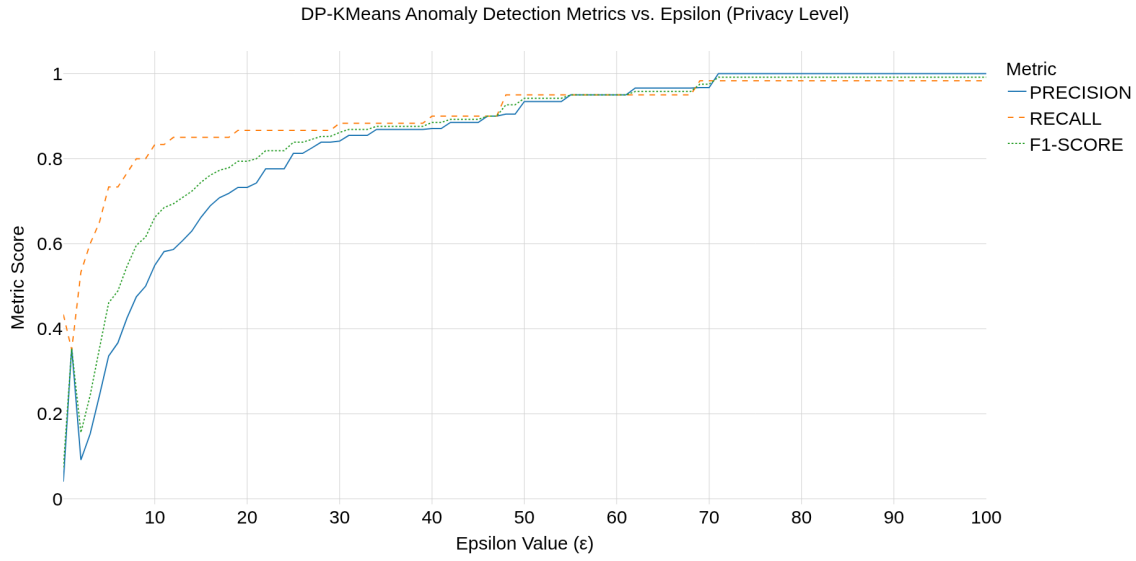
**Figure 3. Variation of confusion matrix components (True Positives, False Positives, True Negatives, and False Negatives) for anomaly detection using DP-KMeans across different Epsilon ( $\epsilon$ ) values. Provides insight into the types of classification errors under varying privacy budgets.**

The True Positives (VP) (blue line) begin at a count of approximately 25 for  $\epsilon = 1$  and increase steadily with  $\epsilon$ , reaching a plateau at approximately 60 (presumably the total number of actual anomalies in the dataset) around  $\epsilon \approx 45$ . Concurrently, the False Negatives (FN) (purple line) start at approximately 35 (Total Anomalies - Initial VP) for  $\epsilon = 1$  and consistently decrease, approaching zero as  $\epsilon$  nears 45. The combined effect of increasing VPs and decreasing FNs explains the rapid rise and subsequent high stabilization of the Recall metric observed in Figure 4.

The True Negatives (VN) (green line) start at a moderate count (around 40-45) and exhibit an upward trend as  $\epsilon$  increases, corresponding to the reduction in False Positives, eventually stabilizing at a high count (around 95-100) for  $\epsilon > 45$ . This component-wise analysis effectively illustrates how the balance of correct and incorrect classifications shifts under varying privacy budgets. This detailed component analysis not only explains the behavior of the aggregated metrics but also directly highlights the practical implications of different privacy budgets on the accuracy of anomaly detection, demonstrating how a balanced  $\epsilon$  value can mitigate classification errors. The pronounced decrease in False Positives alongside the increase in True Positives with increasing  $\epsilon$  are key drivers for the significant improvement in Precision and, consequently, the overall F1-Score for the DP-KMeans anomaly detection methodology.

Figure 4 illustrates the impact of the privacy parameter  $\epsilon$  on key anomaly detection metrics: Precision, Recall, and F1-Score. The x-axis represents the values of  $\epsilon$  ranging from 0 to 100, where increasing  $\epsilon$  indicates decreasing privacy protection (less noise). The y-axis shows the scores for each metric, scaled from 0 to 1.

For very low values of  $\epsilon$  (e.g.,  $0 < \epsilon \leq 10$ ), corresponding to strong privacy guarantees, Precision (blue line) initiates at a notably low value of approximately 0.28 when  $\epsilon = 1$ , suggesting that a large fraction of instances identified as anomalies are, in fact, false positives.



**Figure 4. Impact of the privacy parameter Epsilon ( $\epsilon$ ) on anomaly detection performance metrics (Precision, Recall, and F1-Score) using DP-KMeans. Illustrates the trade-off between the level of privacy and the utility of the detection results.**

In contrast, Recall (orange line) starts at a comparatively higher value of approximately 0.42 at  $\epsilon = 1$  and demonstrates a more rapid initial increase, indicating that a considerable portion of true anomalies are still captured despite the high noise levels.

As  $\epsilon$  increases, particularly within the range of  $10 < \epsilon < 45$ , both Precision and Recall exhibit significant upward trends. Recall reaches near-optimal values, plateauing very close to 1.0, around  $\epsilon \approx 30$ . Precision continues its steep ascent throughout this range as the influence of noise diminishes, leading to fewer false alarms. Consequently, the F1-Score (green line), which represents the harmonic mean of Precision and Recall, rises sharply, reflecting a substantial improvement in the overall effectiveness of the anomaly detection mechanism.

For  $\epsilon$  values greater than approximately 45, all three metrics tend to stabilize at high performance levels. Recall remains consistently near 1.0, while Precision and F1-Score also achieve and maintain high values, around 0.97 and 0.98 respectively. This stabilization suggests that beyond this threshold, further increasing  $\epsilon$  (i.e., further relaxing privacy) yields only marginal gains, if any, for these specific performance metrics under the tested experimental conditions. The graph vividly demonstrates the inherent trade-off: stringent privacy settings (low  $\epsilon$ ) severely compromise Precision, but as  $\epsilon$  is relaxed, a robust and effective balance for anomaly detection can be attained.

## 6. Conclusion

As the energy sector increasingly incorporates Information and Communication Technologies into the metering infrastructures of Electrical Smart Grids, innovation in open energy markets becomes a reality, alongside heightened risks from cyber-threats and privacy violations. Differential privacy offers a robust framework to protect an individual's privacy in public datasets. In this work, we demonstrated an innovative application of differentially private computations to safeguard the secrecy of business-sensitive data within metering datasets from the Brazilian open energy market. We argued that publishing only anonymized outputs, rather than disclosing whole

metering datasets, is crucial for preserving business secrecy for retailers and traders, protecting consumption patterns, market shares, and client portfolios from malicious competitors.

In this paper we investigate an actual issue in energy smart grids by evaluating the use of Differentially Private KMeans (DP-KMeans) for anomaly detection in energy consumption metering. The findings reveal a clear and quantifiable trade-off between the strength of privacy guarantees, as controlled by the parameter  $\epsilon$ , and the utility of the anomaly detection, measured by Precision, Recall, F1-Score, and an analysis of the confusion matrix components. Specifically, our results (presented in Figure 4 and Figure 3) show that while very strong privacy settings (low  $\epsilon$  values, e.g.,  $\epsilon < 10$ ) significantly impact detection accuracy, primarily by substantially increasing the number of false positives which drastically reduces Precision, performance improves markedly as  $\epsilon$  is relaxed. For instance, Recall achieves near-optimal levels relatively early (around  $\epsilon \approx 30$ ), while Precision and F1-Score reach high utility (e.g.,  $> 0.95$ ) for  $\epsilon > 45$ .

Our results demonstrate that operational ranges for  $\epsilon$  can be identified where a practical and effective balance between anomaly detection capabilities and the stringent secrecy requirements of the open energy market is achievable. This finding is crucial for guiding the practical implementation of privacy-preserving analytics in a real-world context, showcasing the potential for market participants to gain valuable insights while mitigating business risks.

While a semantic interpretation of the identified anomalies was beyond the scope of this current methodological evaluation, requiring energy sector domain expertise, not utilized herein, future work will prioritize this qualitative analysis through direct collaboration with specialists from the energy sector. Such insights, combined with further explorations into adaptive  $\epsilon$  selection and the impact of diverse anomaly types on DP-KMeans performance, are crucial for advancing responsible, data-driven decision-making in the energy sector without compromising sensitive business information, thereby fostering greater trust and security.

Finally, at this exploratory work we aimed at enhancing data protection awareness within the energy sector. Given the current open access to consumption data at the open energy market, it's possible for energy traders to scrutinize consumption patterns of their competitors' clients, possibly inferring portfolios and market shares, which underscores the need for robust data protection. As the sector matures in its approach to security and privacy, the authors anticipate a growing commitment among stakeholders to implement and advance the proposed solution.

## Acknowledgments

Authors thank MCTI, Softex, ANEEL, CESP and Auren Energia for financial support (Project CTP - 1283 - *Plataforma Tecnológica para Digitalização da Portabilidade e Agregação da Medição no Ambiente Varejista de Energia*) and grant 03b1629f-e9d9-40bb-bc95-89d7362e28c2. We also thank Way2 for collaboration and CPQD for institutional support.

## References

- CCEE (2024). White paper - ccee. <https://www.ccee.org.br/>. Accessed 2024-11-18.
- Dwork, C. (2006). Differential privacy. pages 1–12.
- Dwork, C., Kenthapadi, K., McSherry, F., Mironov, I., and Naor, M. (2006a). Our data, ourselves: Privacy via distributed noise generation. pages 486–503.



- Dwork, C., McSherry, F., Nissim, K., and Smith, A. (2006b). Calibrating noise to sensitivity in private data analysis. pages 265–284.
- ENISA (2023). Cybersecurity and privacy in ai – forecasting demand on electricity grids. Technical report. ENISA.
- Fan, L., Xiong, L., and Sunderam, V. (2013). Differentially private anomaly detection with a case study on epidemic outbreak detection. In *2013 IEEE 13th International Conference on Data Mining Workshops (ICDMW)*, pages 833–840.
- Holohan, N., Braghin, S., Aonghusa, P. M., and Levacher, K. (2019). Diffprivlib: the ibm differential privacy library. *ArXiv e-prints*, 1907.02444 [cs.CR].
- Janghyun, K., Barry, H., Tianzhen, H., et al. (2022). A review of preserving privacy in data collected from buildings with differential privacy. *Journal of Building Engineering*, 56:104724.
- Leukam Lako, F., Lajoie-Mazenc, P., and Laurent, M. (2021). Privacy-preserving publication of time-series data in smart grid. *Security and Communication Networks*, 2021(1):6643566.
- Marks, J., Montano, B., Chong, J., Raavi, M., Islam, R., Cerny, T., and Shin, D. (2021). Differential privacy applied to smart meters: A mapping study. pages 761–770. Association for Computing Machinery.
- McElroy, T., Roy, A., and Hore, G. (2023). Flip: a utility preserving privacy mechanism for time series. *Journal of Machine Learning Research*, 24(111):1–29.
- Near, J. P., Darais, D., Lefkovitz, N., Howarth, G., et al. (2023). Guidelines for evaluating differential privacy guarantees (nist sp 800-226). Technical report, National Institute of Standards and Technology.
- Paixão, A. C. P., Camargo, G. d. F. L., and Braga, A. M. (2025a). Testing open-source libraries for private counts and averages on energy metering time series. In *Proceedings of the 2025 20th European Dependable Computing Conference (EDCC)*, pages 100–104. IEEE.
- Paixão, A. C. P., da Silva, B. R., Silva, R. L., Cardoso, F. H., and Braga, A. (2025b). Understanding how to use open-source libraries for differentially private statistics on energy metering time series. In *Proceedings of the 10th International Conference on Internet of Things, Big Data and Security (IoTBDSC 2025)*, pages 289–296. SCITEPRESS - Science and Technology Publications.
- Peralta-Peterson, M. and Kotevska, O. (2021). Effectiveness of privacy techniques in smart metering systems. pages 675–678.
- Roman, A.-S. (2023). Evaluating the privacy and utility of time-series data perturbation algorithms. *Mathematics*, 11(5):1260.
- Roman, A.-S., Genge, B., Duka, A.-V., and Haller, P. (2021). Privacy-preserving tampering detection in automotive systems. *Electronics*, 10(24):3161.
- Shaham, S., Ghinita, G., Krishnamachari, B., and Shahabi, C. (2024). Differentially private publication of electricity time series data in smart grids. *arXiv preprint arXiv:2408.16017*.
- Zhao, J., Jung, T., Wang, Y., and Li, X. (2014). Achieving differential privacy of data disclosure in the smart grid. pages 504–512.