

Universidade Federal de Campina Grande
Centro de Engenharia Elétrica e Informática
Coordenação de Pós-Graduação em Ciência da Computação

Preservando a Privacidade em Smart Grids Através
de Adição de Ruído

Pedro Yóssis Silva Barbosa

Dissertação submetida à Coordenação do Curso de Pós-Graduação em
Ciência da Computação da Universidade Federal de Campina Grande -
Campus I como parte dos requisitos necessários para obtenção do grau
de Mestre em Ciência da Computação.

Área de Concentração: Ciência da Computação

Linha de Pesquisa: Redes de Computadores e Sistemas Distribuídos

Andrey Elísio Monteiro Brito e Hyggo Oliveira de Almeida

(Orientadores)

Campina Grande, Paraíba, Brasil

©Pedro Yóssis Silva Barbosa, 27/02/2014

DIGITALIZAÇÃO:
SISTEMOTECA - UFCG

FICHA CATALOGRÁFICA ELABORADA PELA BIBLIOTECA CENTRAL DA UFCG

B238p

Barbosa, Pedro Yóssis Silva.

Preservando a privacidade em Smart Grids através de adição de ruído /
Pedro Yóssis Silva Barbosa. – Campina Grande, 2014.

78 f. : il. Color.

Dissertação (Mestrado em Ciência da Computação) - Universidade
Federal de Campina Grande, Centro de Engenharia Elétrica e Informática.

"Orientação: Prof. Andrey Elísio Monteiro Brito, Prof. Hyggo Oliveira
de Almeida".

Referências.

1. Privacidade. 2. Smart Grids. I. Brito, Andrey Elísio Monteiro. II.
Almeida, Hyggo Oliveira de. III. Título.

CDU 004.056(043)

**"PRESERVANDO A PRIVACIDADE EM SMART GRIDS ATRAVÉS DE ADIÇÃO DE
RÚIDO"**

Pedro Y.S. Barbosa
PEDRO YÓSSIS SILVA BARBOSA

DISSERTAÇÃO APROVADA EM 27/02/2014

Andrey Elísio Monteiro Brito
ANDREY ELÍSIO MONTEIRO BRITO, Dr., UFCG
Orientador(a)

Hyggo Oliveira de Almeida
HYGGO OLIVEIRA DE ALMEIDA, D.Sc, UFCG
Orientador(a)

Raquel Vigolino Lopes
RAQUEL VIGOLVINO LOPES, D.Sc, UFCG
Examinador(a)

Marcos Ricardo Alcântara Moraes
MARCOS RICARDO ALCÂNTARA MORAIS, D.Sc, UFCG
Examinador(a)

CAMPINA GRANDE - PB

Resumo

Companhias de energia começaram a substituir os medidores de energia tradicionais pelos Smart Meters, que podem transmitir valores de consumo para as companhias em curtos intervalos de tempo. Com uma infraestrutura de Smart Meters, existem muitas motivações para as concessionárias de energia coletarem dados de consumo em alta resolução. Entretanto, isto implica em informações bastante detalhadas sobre os consumidores sendo monitoradas. Consequentemente, um problema sério precisa ser resolvido: como preservar a privacidade dos consumidores sem afetar a prestação de certos serviços pelas concessionárias? Claramente, este é um *tradeoff* entre privacidade e utilidade. Existem diversas abordagens para preservar a privacidade, porém muitas delas afetam a utilidade dos dados ou possuem um alto custo computacional. Neste trabalho, nós propomos e avaliamos uma abordagem computacionalmente barata que preserva a privacidade e utilidade dos dados através de adição de ruído. Para validar a privacidade, nós avaliamos possíveis ataques (tal como Monitoramento Não-Intrusivo de Carga de Eletrodomésticos – NIALM, do inglês *Non-Intrusive Appliance Load Monitoring*) utilizando dados reais de consumidores. Para validar a utilidade, nós avaliamos a influência da abordagem em vários benefícios que podem ser providos com o uso de Smart Meters.

Abstract

Power providers have started replacing traditional electricity meters for Smart Meters, which can transmit power consumption levels to the provider within short intervals. With a Smart Metering infrastructure, there are many motivations for power providers to collect high-resolution data of electricity usage from consumers. However, this implies in very detailed information about the consumers being monitored. Consequently, a serious issue needs to be addressed: how to preserve the privacy of consumers but making the provision of certain services still possible? Clearly, this is a tradeoff between privacy and utility. There are several approaches for privacy preserving, but many of them affect the data usefulness or are computationally expensive. In this work, we propose and evaluate a lightweight approach for privacy and utility based on the addition of noise. To validate the privacy, we evaluate possible attacks (such as a NIALM – *Non-Intrusive Appliance Load Monitoring*) using real consumers' data. To validate the utility, we analyze the influence of the approach in various benefits that can be provided through the use of Smart Meters.

Agradecimentos

Agradeço primeiramente a Deus, que me permitiu vencer todas as dificuldades e realizar este mestrado.

Aos meus pais, Antônio Eduardo e Doralice, grandes incentivadores por todo o apoio, pela educação e infraestrutura que me proporcionaram e pelos ensinamentos que me guiaram nas escolhas da vida.

Aos meus irmãos André, Clara e Paulo por me influenciarem, encorajarem e torcerem por mim. Por suas boas companhias nos momentos de descontração.

Agradeço a Larissa, por todo o companheirismo, pela compreensão das minhas ausências e apoio em todos os momentos.

Meus sinceros agradecimentos aos meus amigos e orientadores Andrey Brito e Hyggo Almeida, pela dedicação, empenho e disponibilidade. Sem suas ideias e conhecimentos, certamente eu não teria chegado até aqui.

Aos professores Jacques Sauvé, Marcos Morais e Raquel Lopes da Universidade Federal de Campina Grande, Keiko Fonseca da Universidade Tecnológica Federal do Paraná e Sebastian Clauß da Technische Universität Dresden, pelas sugestões e contribuições neste trabalho.

Aos professores e funcionários da COPIN e do DSC e finalmente, ao Governo Brasileiro, por meio da CAPES, pelo apoio financeiro fornecido para execução das atividades deste mestrado aqui no Brasil e na Alemanha.

Conteúdo

1	Introdução	1
1.1	Smart Grid e Smart Meters	1
1.2	Definição do Problema e Relevância	2
1.3	Objetivo	4
1.4	Metodologia	4
1.5	Trabalhos Relacionados	5
1.6	Organização da Dissertação	7
2	Uma Abordagem de Ofuscamento de Dados	8
2.1	Medições de Consumo	8
2.2	Modelo Empírico	11
2.2.1	Planejamento de Experimentos	11
2.2.2	Analisando o Erro em Função de X	11
2.2.3	Analisando o Erro em Função de N	12
2.2.4	Analisando o Erro em Função de X e N	13
2.3	Modelo Analítico	15
2.4	Métricas	16
3	Demonstrações da Abordagem	19
3.1	Faturamento para um Consumidor Residencial	19
3.2	Faturamento para um Consumidor Industrial	21
3.3	Faturamento com Política de Horário	24
3.4	Monitoramento de uma Região	25

4	Otimizações da Abordagem	28
4.1	Calculando o Erro Permitido (e_p)	28
4.1.1	Janela Saltitante Mensal	29
4.1.2	Janela Deslizante Mensal	32
4.1.3	Janela Saltitante Diária	33
4.1.4	Conclusões Sobre o Cálculo de e_p	35
4.2	Análise de Outras Métricas	37
4.2.1	Correlação	37
4.2.2	Relação Sinal-ruído	38
4.2.3	Erro Quadrático Médio	39
4.2.4	Informação Mútua	40
4.2.5	Investigação de Problemas	42
4.3	Comparando Distribuições de Probabilidade	43
5	Validação da Abordagem	47
5.1	Validação da Privacidade	47
5.1.1	Ataque do NIALM	47
5.1.2	Ataque do Dia da Semana	51
5.1.3	Ataque do Filtro	52
5.1.4	Outros Ataques	56
5.2	Validação da Utilidade	60
5.2.1	Otimizações de Faturamento	60
5.2.2	Monitoramento e Gerenciamento de Carga	61
5.2.3	Detecção de Vazamentos e Roubos de Energia	62
5.2.4	Previsão de Carga para Grupos e Regiões	63
5.2.5	Previsão de Carga para Indivíduos	63
5.2.6	Faturamento com Política de Horário	63
5.2.7	Faturamento com Política de Níveis de Demanda	63
5.2.8	Análise Individual dos Dados	64
5.2.9	Ferramentas para Feedbacks em Casa	65
5.2.10	Outras Funcionalidades	65

<i>CONTEÚDO</i>	vi
6 Considerações Finais	67
6.1 Conclusões e Contribuições	67
6.2 Trabalhos Futuros	68
A Gerador de Carga	74
B Artigos Aceitos para Publicação	78

Lista de Figuras

1.1	Perfil residencial alinhado com assinaturas de eletrodomésticos [21].	3
1.2	Metodologia utilizada.	4
2.1	Como um conjunto de medições pode ser usado.	9
2.2	Analisando o erro variando-se X	13
2.3	Analisando o erro variando-se N	14
2.4	Distribuição e intervalo de confiança dos erros obtidos (e_o) para $X = 0,0217$, $N = 4464$ e $e_p = 2kWh$	15
3.1	Perfil residencial mensal com medições a cada 30 minutos.	20
3.2	Perfil residencial mensal ofuscado com medições a cada 30 minutos.	20
3.3	Perfil residencial diário real (azul) versus perfil residencial diário ofuscado (vermelho) com medições a cada 30 minutos.	21
3.4	Perfil industrial mensal com medições a cada 10 minutos.	22
3.5	Perfil industrial mensal ofuscado com medições a cada 10 minutos.	23
3.6	Perfil industrial diário real (azul) versus perfil industrial diário ofuscado (vermelho) com medições a cada 10 minutos.	23
3.7	Exemplos de tipos de tarifas estabelecidas pela ANEEL [9].	25
3.8	Perfil regional usando dados reais (azul) versus usando dados ofuscados (vermelho) com medições a cada 30 minutos.	26
3.9	Erros obtidos ao longo do tempo para o exemplo da Figura 3.8.	27
4.1	Perfil com alto consumo no mês anterior e baixo consumo no mês atual.	29
4.2	Caso ruim 1 ofuscado com JSM . O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.	30

4.3	Perfil com baixo consumo no mês anterior e alto consumo no mês atual. . .	31
4.4	Caso ruim 2 ofuscado com <i>JSM</i> . O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.	31
4.5	Caso ruim 1 ofuscado com <i>JDM</i> . O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.	32
4.6	Caso ruim 2 ofuscado com <i>JDM</i> . O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.	33
4.7	Caso ruim 1 ofuscado com <i>JSD</i> . O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.	34
4.8	Caso ruim 2 ofuscado com <i>JSD</i> . O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.	35
4.9	Correlações obtidas para diferentes valores de X . Os níveis de significância são de 95%.	39
4.10	SNR obtidos para diferentes valores de X . Os níveis de significância são de 95%.	40
4.11	MSE obtidos para diferentes valores de X . Os níveis de significância são de 95%.	41
4.12	MI obtidos para diferentes valores de X . Os níveis de significância são de 95%.	42
4.13	Perfil diário original (azul) versus perfil diário ofuscado com um gerador de números aleatórios quebrado (vermelho).	43
4.14	Densidades de probabilidade. Arcoseno está representada em verde, Laplace em preto, Normal em vermelho, Uniforme em azul e U-quadrática em roxo.	45
4.15	Erros obtidos (e_o) com cada distribuição de probabilidade.	45
5.1	Fluxo de tarefas para validação da privacidade através de ataques de NIALM.	48
5.2	Exemplo de uma semana obtida do conjunto de dados REDD. Medições são de 1 minuto.	48
5.3	Perfil real do microondas usado no perfil da Figura 5.2.	49
5.4	Microondas inferido pelo INDIC a partir do perfil da Figura 5.2.	49
5.5	Perfil da Figura 5.2 ofuscado com <i>JSD</i> e e_p de 5%.	50

5.6	Microondas inferido pelo INDIC a partir do perfil ofuscado da Figura 5.5.	51
5.7	Perfil diário gerado pelo Load Generator.	53
5.8	Perfil diário ofuscado.	54
5.9	Perfil diário ofuscado filtrado com $P = 2$	54
5.10	Perfil diário ofuscado filtrado com $P = 8$	55
5.11	Perfil diário ofuscado filtrado com $P = 22$	55
5.12	Perfil diário ofuscado filtrado com $P = 200$	56
5.13	Detecção de picos no perfil da Figura 5.7 usando-se S_1 . Os pontos em vermelho são os picos encontrados.	59
A.1	Exemplo de funcionamento do Gerador de Carga.	75
A.2	Exemplo de perfil gerado pelo Gerador de Carga. Em preto está representado o perfil agregado, em verde a máquina de lavar, em vermelho a máquina de lavar louça e em azul a geladeira.	76
A.3	Perfil de consumo representado no PlotWatt [31].	76
A.4	Análises de consumo e identificação de eletrodomésticos pelo PlotWatt [31].	77

Lista de Tabelas

3.1	Resultados de ofuscamento com três tipos de tarifas.	24
4.1	Resultados das estratégias para os casos extremos. Valores em vermelho significam que estão fora do esperado. Valores em verde significam que estão dentro do esperado.	36
4.2	Intervalos de confiança dos erros absolutos entre a média dos e_p e o e_{pr} para 1000 perfis.	36
4.3	Intervalos de confiança das métricas de privacidade e utilidade para as diferentes estratégias ao ofuscarmos 1000 perfis.	37
4.4	Modelos analíticos obtidos para diferentes distribuições de probabilidade.	44
4.5	Níveis de privacidade obtidos com cada métrica para cada distribuição de probabilidade. Os níveis de significância são de 97,5%.	46
5.1	Valores RMS e MNE do INDIC usando diferentes configurações de ofuscamento (sem ofuscamento, ofuscamento com erro permitido de 1%, 2% e 5%.	51
5.2	Efeito do ataque do dia da semana pra um consumidor residencial. Os intervalos de confiança são com níveis de significância de 95%.	52
5.3	Efeito do ataque do filtro para um perfil diário residencial usando diferentes valores de P	57
5.4	Funcionalidades que utilizam dados de medições e o impacto do ofuscamento.	61
5.5	Outras funcionalidades/benefícios que não estão relacionadas com medições de consumo.	66

Capítulo 1

Introdução

Nos dias de hoje, existe um grande interesse pelos dados gerados por diversos dispositivos dos consumidores. Dados que revelam informações sobre consumos de energia, balanços financeiros, estados de saúde e localizações geográficas, são apenas alguns exemplos. Devido à evolução dos algoritmos para extrair informações, grandes volumes de dados agora podem ser analisados pelas prestadoras de serviços para que novas percepções sobre os consumidores sejam obtidas de maneira rápida e eficiente. Alguns especialistas declaram o surgimento da ciência dos dados [11] e cada vez mais empresas prestadoras de serviços investem em tecnologias para o processamento de dados. Atualmente, processamento de dados em larga escala é um dos temas mais discutidos na indústria dos Smart Grids [19].

1.1 Smart Grid e Smart Meters

Um Smart Grid é uma rede elétrica complexa que utiliza sensores e recursos computacionais [36] para melhorar a eficiência, confiabilidade e sustentabilidade da produção e distribuição de eletricidade. Em termos gerais, é a aplicação de tecnologia da informação para o sistema elétrico de potência, integrada aos sistemas de comunicação e infraestrutura de rede automatizada.

A implantação de Smart Meters é a porta de entrada para o processo de transição de uma rede elétrica convencional para um Smart Grid. Um Smart Meter é um medidor inteligente que pode transmitir dados de medições com propósitos de faturamento e de monitoramento de carga para a concessionária de energia, provendo leituras precisas em determinados inter-

valos de tempo de maneira automática. Tal frequência de leituras ainda está para ser definida pelas concessionárias; entretanto, foi especulado que isto pode ser tão granular quanto a cada 1 ou 5 minutos, o que pode ocasionar problemas para a privacidade dos consumidores devido à disponibilidade e processamento de tais dados de medição [10].

Com os Smart Meters, informações detalhadas sobre o uso de energia elétrica pelos consumidores podem ser monitoradas, mudando assim as relações comerciais entre consumidores e concessionárias de energia. Por um lado, tal tecnologia pode criar oportunidades para a prestação de novos serviços pelas concessionárias, por outro lado, levanta preocupações com respeito à privacidade dos consumidores [6]. Rajagopalan et al. [36] descrevem este problema como um *tradeoff* entre utilidade e privacidade dos dados.

1.2 Definição do Problema e Relevância

A partir dos dados enviados pelos Smart Meters, é possível inferir o comportamento dos consumidores. São exemplos de tais comportamentos: quais eletrodomésticos são utilizados; se uma casa está vazia em um certo horário; quando os habitantes acordam, tomam banho ou desligam a televisão; ou mais ainda, se a geladeira ou a máquina de lavar já não estão operando de maneira eficiente.

O processo de analisar perfis de consumo com propósito de inferir quais eletrodomésticos estão sendo utilizados é conhecido como NIALM – *Non-Intrusive Appliance Load Monitoring* (Monitoramento não-intrusivo de carga de eletrodomésticos). Existem muitos algoritmos de NIALM propostos na literatura. Kelly et al. [21] projetam, implementam e avaliam algumas metodologias para identificar o uso de eletrodomésticos a partir de perfis de carga. Se o algoritmo de NIALM é executado remotamente, os proprietários das casas talvez não saibam que os seus dados de comportamento estão sendo monitorados e armazenados. A Figura 1.1 mostra a identificação de assinaturas de alguns eletrodomésticos a partir de um perfil de consumidor residencial [21].

Mesmo que tais informações de comportamento não sejam em princípio úteis para as concessionárias de energia, existe um grande mercado comercial interessado por tais informações de hábitos de consumo. Certamente, este tipo de informação é importante para muitas empresas que desejam identificar potenciais consumidores de seus produtos e servi-

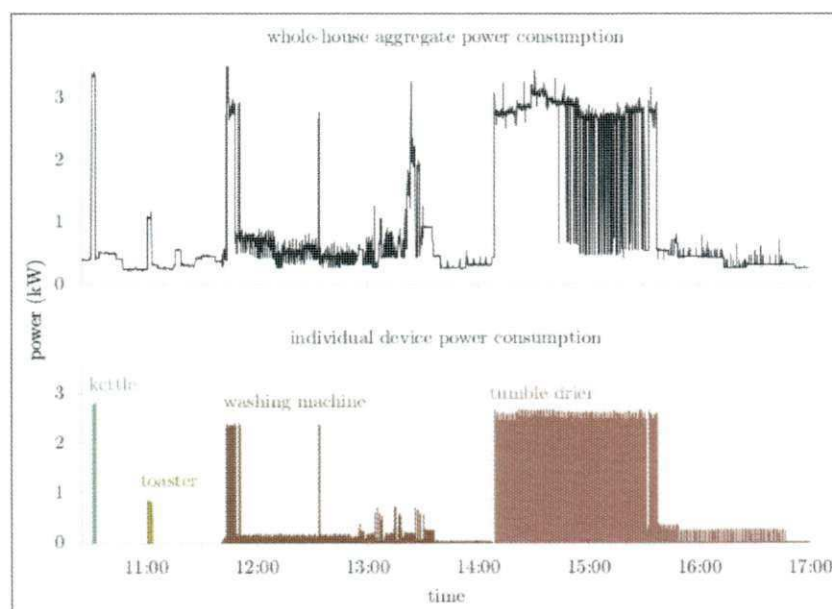


Figura 1.1: Perfil residencial alinhado com assinaturas de eletrodomésticos [21].

ços.

Desta forma, a exposição dos perfis e hábitos dos consumidores levanta questões de privacidade. Regras claras são necessárias para proteger os consumidores do mau uso dos seus dados comportamentais e para evitar que os Smart Grids se tornem uma nova modalidade de Big Brother [6]. Infelizmente, leis de proteção podem levar décadas para serem implantadas, ao passo que Smart Meters já estão operando.

Com respeito à utilidade dos dados para as concessionárias de energia, algumas motivações para a coleta de dados em alta resolução podem ser: identificação de perdas não-técnicas (e.g., roubos de eletricidade); otimização em previsão de carga; otimização em faturamento e em serviços de tarifação; monitoramento dos índices de qualidade de energia. Mais ainda, tais dados podem ser utilizados para facilitar e aprimorar o gerenciamento da rede, reduzir picos de carga, calibrar a distribuição de carga e muitas outras utilidades [36]. Assim, a coleta e a disseminação de informações de consumo são críticas para o Smart Grid. Entretanto, com respeito à privacidade, existe a possibilidade de tais dados serem utilizados para propósitos que não estão relacionados com o gerenciamento de energia, fazendo com que isto seja potencialmente perigoso para a privacidade dos consumidores individuais. De fato, esta é uma das principais razões para a não implantação em massa de Smart Meters em muitos

países [22].

Este trabalho é relevante para resolver o seguinte problema de negócio: como preservar a privacidade dos dados dos consumidores sem afetar a funcionalidade de alguns serviços por parte da concessionária de energia? É necessária uma solução que atenda às necessidades tanto para privacidade quanto para utilidade.

1.3 Objetivo

O objetivo deste trabalho é propor uma solução de privacidade que atenda às necessidades de privacidade e que possui um efeito mínimo na utilidade dos dados para muitos serviços providos pelas companhias de energia. A validação da solução foi realizada considerando tanto aspectos de privacidade quanto de utilidade.

1.4 Metodologia

A Figura 1.2 apresenta a metodologia utilizada para o desenvolvimento deste trabalho. Para a execução das tarefas, foram realizados diversos experimentos utilizando dados reais de consumidores [13], [12], [23] ou utilizando dados simulados gerados por um software (apresentado no Apêndice A).

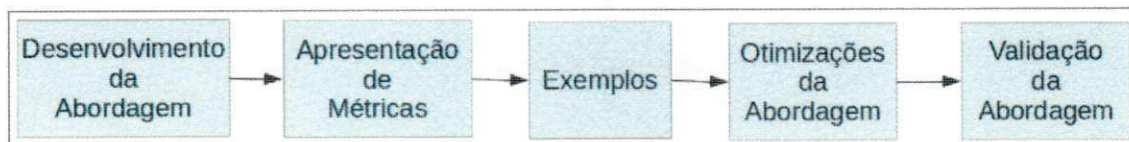


Figura 1.2: Metodologia utilizada.

Na primeira tarefa, desenvolveremos os modelos a serem utilizados pela abordagem. De uma maneira geral, nós propomos que cada consumidor (consumidor aqui se refere a Smart Meter) envie valores ofuscados de consumo em vez de valores reais. Para uma medição individual, o consumidor faz a leitura do consumo real e adiciona um valor aleatório do intervalo $[-X; +X]$ (por exemplo, de acordo com uma distribuição Uniforme). Assim, quando a concessionária de energia somar todos os dados recebidos, ela irá obter uma aproximação

do consumo total real, pois a adição dos números aleatórios tende a ser minimizada após a operação de soma.

Após desenvolver a solução proposta, apresentaremos métricas a serem utilizadas para mensurar tanto a privacidade quanto a utilidade, conforme representada pela segunda tarefa. Para mensurar a privacidade, serão utilizadas métricas como a correlação entre o perfil ofuscado e o perfil original. Para mensurar a utilidade, o erro resultante entre a soma dos valores ofuscados e a soma dos valores originais será utilizado como métrica. Mais ainda, nos modelos desenvolvidos, o valor máximo permitido para X é calculado com base em um erro máximo permitido (pela concessionária ou por uma entidade regulamentadora).

Em seguida demonstraremos o uso da abordagem utilizando dados reais de consumidores, conforme representado pela terceira tarefa.

Na quarta tarefa, realizaremos estudos para otimizar a técnica proposta e as métricas utilizadas. Finalmente, na última tarefa, validaremos a proposta considerando tanto aspectos de privacidade quanto de utilidade. Para validar a privacidade, vários ataques são aplicados e discutidos. Para validar a utilidade, mostramos que mesmo com valores ofuscados, os dados de medições ainda são úteis e não impactam negativamente em muitas aplicações do Smart Grid. No fim, concluímos que a solução proposta atende às necessidades dos consumidores (privacidade) e das concessionárias de energia (utilidade).

1.5 Trabalhos Relacionados

Existem muitos trabalhos propondo mecanismos de criptografia e de chave pública para serem usados entre os Smart Meters e as concessionárias de energia [4]. Criptografia comum pode ser útil para garantir a confidencialidade ao longo do canal e garantir a total utilidade dos dados (porque após a descriptografia, os dados são totalmente legíveis e confiáveis). Entretanto, a privacidade do consumidor dentro da concessionária de energia (ou fora dela caso algum empregado malicioso exporte os dados) ainda é ameaçada.

Efthymiou et al. [10] descrevem um método para que os Smart Meters enviem os dados de medição de maneira segura e anônima. Entretanto, no Smart Grid, a identificação dos consumidores é essencial, uma vez que os dados de medição são usados para propósitos de faturamento.

Mecanismos promissores para resolver o problema da privacidade são os esquemas homomórficos [27], [38]. Tais soluções buscam suportar utilidade e privacidade em diferentes níveis, porém, não possuem uma base teórica e robusta para ambas [36], i.e., são dependentes de aplicação. Além disso, apesar de soluções para esquemas de criptografia homomórfica completa terem sido propostas e aperfeiçoadas, é difícil ignorar questões relacionadas à eficiência. Atualmente, todos os tipos de esquemas de criptografia homomórfica completa propostos ainda possuem um longo caminho de evolução antes que sejam utilizados na prática [26].

O uso de baterias recarregáveis entre os eletrodomésticos e o Smart Meter pode ajudar a reduzir os problemas de privacidade [29], [20], pois as assinaturas dos dispositivos deixam de ser legíveis e apenas as assinaturas das baterias são expostas. Porém, além das próprias assinaturas das baterias ainda serem expostas, nem sempre é viável ter baterias em uma residência.

Abordagens de privacidade parecidas com a nossa é o trabalho de Bohli et. al. [7] e o de Wang et. al. [42]. Bohli et. al. [7] discutem uma abordagem para proteger a privacidade ofuscando os valores com um ruído proveniente de uma distribuição normal. A abordagem é avaliada com a necessidade de se calcular o consumo total para um grupo de consumidores e conclui-se que é necessária uma grande quantidade de consumidores para se obter um ofuscamento considerável. Entretanto, Bohli et. al. não discutem o impacto de tal ofuscamento através de métricas de privacidade e utilidade e não avaliam a abordagem utilizando perfis reais de consumidores (simulações são realizadas com a suposição de que os Smart Meters geram valores seguindo uma distribuição normal). Além disso, a abordagem não foi avaliada com a necessidade de se calcular o consumo total de um consumidor ao fim de um período de faturamento.

Wang et. al. [42] propõem uma abordagem em que os valores são ofuscados com um ruído proveniente de uma distribuição GMM – Gaussian Mixture Model (Modelo Gaussiano de Misturas). Como métrica de privacidade, sugerem o uso de testes estatísticos como o teste F pareado para comparar o perfil real do perfil ofuscado e o teste de Kolmogorov-Smirnov para comparar o perfil de um consumidor com os demais de uma região. Assim, os *p* – valores desses testes servem como métricas para o ofuscamento. Entretanto, tais métricas não são adequadas. Por exemplo, o teste F é normalmente usado para comparar

as variâncias de dois conjuntos de dados, mas não o nível de similaridade entre eles. Dessa forma, para dois perfis que sejam completamente diferentes mas que possuam variâncias iguais, o teste acusará que não existe diferença. O ofuscamento atingido com a abordagem de Wang et. al. [42] é menor do que o atingido com a abordagem de Bohli et. al. [7] e a abordagem não é avaliada utilizando perfis reais de consumidores e nem em aplicações reais de Smart Grid.

1.6 Organização da Dissertação

Este trabalho se organiza da seguinte forma; no Capítulo 2 apresenta-se a solução proposta para, no Capítulo 3, apresentar exemplos sobre o funcionamento geral da solução. No Capítulo 4 possíveis otimizações são avaliadas e discutidas. No Capítulo 5 apresenta-se a validação da solução proposta. O trabalho é finalizado com um capítulo dedicado às conclusões e trabalhos futuros.

Capítulo 2

Uma Abordagem de Ofuscamento de Dados

2.1 Medições de Consumo

Considerando medições de consumo, a abordagem proposta tem como foco:

- Possibilitar o cálculo do consumo total de um consumidor durante um período de tempo (e.g., um mês para faturamento);
- Possibilitar o cálculo do consumo total de todos os consumidores de uma região em um certo instante de tempo;
- Inibir a obtenção do consumo atual de um consumidor individual em um certo instante de tempo.

Desta forma, se cada consumidor enviar uma medição de consumo periodicamente, a concessionária de energia poderá organizar os dados como uma matriz, onde a soma de uma linha se refere ao consumo total de um consumidor ao longo de um período de tempo e a soma de uma coluna se refere ao consumo total de todos os consumidores do grupo em um instante de tempo, como mostrado na Figura 2.1. As linhas da matriz podem ser usadas para propósitos de faturamento (mesmo que os consumidores possuam diferentes períodos de faturamento) e as colunas podem ser usadas para propósitos de monitoramento de carga.

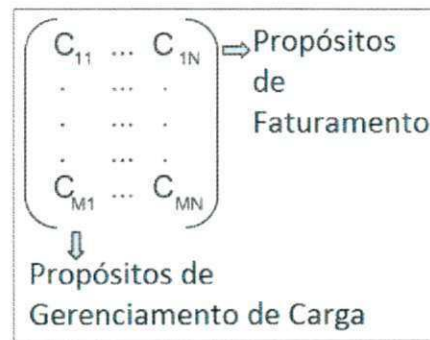


Figura 2.1: Como um conjunto de medições pode ser usado.

Para esconder os consumos instantâneos, nós propomos que os Smart Meters enviem medições ofuscadas de forma que não afetem os resultados das operações de agregação.

Para desenvolver a solução, primeiramente nós consideramos faturamento (soma de uma linha) como aplicação base. Em seguida apresentaremos como a solução funciona para propósitos de gerenciamento de carga (soma de uma coluna).

Para propósitos de faturamento, se a cada medição individual o Smart Meter ler o valor de consumo e adicionar um número aleatório do intervalo $[-X; +X]$, no fim do período de faturamento o resultado será:

$$\sum_{i=1}^N c_i \approx \sum_{i=1}^N (c_i + x_i) \quad (2.1)$$

onde N é o número de medições, x_i é um número aleatório do intervalo $[-X; +X]$ e c_i é uma medição de consumo.

Com essa abordagem, a mudança no procedimento de comunicação entre o consumidor e a concessionária é apenas a geração de um número aleatório e a adição deste número com o valor de consumo. Desta forma, a complexidade de se utilizar tal abordagem é igual à complexidade de se gerar um número aleatório (e.g., o algoritmo de Mersenne-Twister pode ser utilizado com uma complexidade de $O(p^2)$, onde p é o grau do polinômio que indica o período de repetição do algoritmo [28]). Com isto, a abordagem proposta é simples e de baixa complexidade computacional, possibilitando assim a implantação em dispositivos embarcados com recursos limitados (como é o caso de muitos Smart Meters).

Podemos ainda reescrever a formalização acima da seguinte forma:

$$\sum_{i=1}^N c_i = \sum_{i=1}^N (c_i + x_i) - e_o \quad (2.2)$$

onde e_o é um erro obtido devido aos números aleatórios inseridos. Assim, e_o é a soma de todos os números aleatórios inseridos:

$$e_o = \sum_{i=1}^N x_i \quad (2.3)$$

Como observado, tal abordagem de ofuscamento pode inserir um erro no valor total resultante. Entretanto, mesmo com atuais padrões de medição, a existência de erros em aplicações da rede elétrica é algo inerente. Por exemplo, para faturamento, no Brasil o INMETRO (Instituto Nacional de Metrologia, Qualidade e Tecnologia) estabelece limites de erros percentuais para medições com índices de classe. Nesta classificação, as medições de energia ativa para o setor residencial podem ter erros entre +/- 2,00%, enquanto que para o setor industrial as medições podem ter erros entre +/- 3,00%. Para mais informações consulte a portaria de número 375, de 27 de setembro de 2011 [18].

Claramente, se o valor de X aumenta, o nível de ofuscamento também aumenta, ao passo que a precisão diminui. Sendo assim, os seguintes problemas técnicos são evidentes:

- Para que a concessionária obtenha um valor preciso do consumo total de um consumidor em um período de faturamento (i.e., e_o abaixo de um limite aceitável), quão pequeno deve ser X ?
- Para evitar que a concessionária (ou um atacante) infira informações a partir dos dados ofuscados, quão grande deve ser X ?

Foram desenvolvidos um modelo empírico e um analítico de maneira que, dado um erro máximo aceitável, calculam o valor de X que representa o equilíbrio entre utilidade (precisão) e privacidade (ofuscamento). Ambos os modelos apresentaram resultados bastante similares, o que pode servir como validação para ambos.

Para ambos os modelos, as variáveis independentes são X e N . A variável de resposta é e_o . Além disso, a variável e_p é usada para representar o erro máximo permitido (ou seja, e_o deve ser menor que e_p).

2.2 Modelo Empírico

2.2.1 Planejamento de Experimentos

Para encontrar o modelo empírico, as seguintes hipóteses foram definidas:

1. H_{0-0} : A precisão dos dados não é alterada de maneira estatisticamente significativa quando mudamos o valor de X para qualquer outro valor possível.
 - $\frac{Var(e_{oX})}{Var(e_{oX})} = 1$, onde e_{oX} é o erro obtido usando um valor menor de X e e_{oX} é o erro obtido usando um valor maior de X .
2. H_{1-0} : A precisão dos dados não é alterada de maneira estatisticamente significativa quando mudamos o valor de N para qualquer outro valor possível.
 - $\frac{Var(e_{oN})}{Var(e_{oN})} = 1$, onde e_{oN} é o erro obtido usando um valor menor de N e e_{oN} é o erro obtido usando um valor maior de N .
3. H_{2-0} : Não podemos extrair uma correlação entre X e N para a precisão dos dados.
 - Não existe uma função $f(X, N) = e_p$ e conseqüentemente $f(e_p, N) = X$ que obtenha valor de X que representa o equilíbrio entre utilidade e privacidade.

As hipóteses H_{0-0} e H_{1-0} foram projetadas para serem testadas usando um teste de análise de variância (teste F ou ANOVA) com o design de experimentos "One variation at a time" (uma variação por vez), enquanto que a hipótese H_{2-0} foi projetada para ser testada usando uma regressão múltipla com o design de experimentos "K-factorial" (K-fatorial). De acordo com Wohlin et. al. [43], é importante analisar os efeitos individuais de cada fator na variável de resposta (hipóteses H_{0-0} e H_{1-0}) antes de analisar o efeito da interação dos fatores (hipótese H_{2-0}).

2.2.2 Analisando o Erro em Função de X

Para simular o valor de um erro, de acordo com a Equação 2.3, apenas os parâmetros N e X são necessários. Portanto, valores de consumo não são necessários neste estudo. Utilizando o design de experimentos "One variation at a time" para testar apenas o efeito de X no

erro, nós fixamos N (e.g. sempre 1000, que é um valor factível) e variamos apenas X (e.g. de 0.1 até 100, que são valores factíveis). Após isso, nós obtemos várias amostras para cada configuração e armazenamos os piores casos (quando o erro obtido é o máximo na configuração específica). Como os erros seguem uma distribuição normal (com média igual a zero e variância dependendo do valor de X), nós podemos realizar um teste F entre duas populações, uma com o menor valor testado de X e outra com o maior valor testado de X . Usando este teste com 95% de confiança ($\alpha = 0,05$), nós encontramos uma diferença estatística, pois o teste retornou um p -valor muito pequeno ($< 2,2 \cdot e^{-16}$) e a condição para que a hipótese nula seja verdadeira é obter um p -valor maior que α . Portanto, encontramos nossa primeira conclusão:

- Evidência 1: a hipótese nula H_{0-0} é falsa e a hipótese alternativa H_{0-1} é verdadeira:
 - A precisão dos dados foi alterada de maneira significativa quando mudamos o valor de X (i.e. $\frac{Var(e_{oX})}{Var(e_{oX})} \neq 1$).

Para analisar como o valor de X afeta a precisão dos dados, nós podemos representar o gráfico como ilustrado na Figura 2.2. Esta relação é exponencial, conforme iremos apresentar posteriormente.

2.2.3 Analisando o Erro em Função de N

Utilizando o design de experimentos "One variation at a time" para testar apenas o efeito de N no erro, nós fixamos X (e.g. sempre 100, que é um valor factível) e variamos apenas N (e.g. de 1 até 1000, que são valores factíveis). Após isso, obtemos várias amostras para cada configuração e armazenamos os piores casos (quando o erro é o máximo na configuração específica). Nós podemos realizar um teste F entre duas populações, uma com o menor valor testado de N e outra com o maior valor testado de N . Usando este teste com 95% de confiança ($\alpha = 0,05$), nós encontramos uma diferença estatística, pois o teste retornou um p -valor muito pequeno ($< 2,2 \cdot e^{-16}$). Portanto, encontramos nossa segunda conclusão:

- Evidência 2: a hipótese nula H_{1-0} é falsa e a hipótese alternativa H_{1-1} é verdadeira:
 - A precisão dos dados foi alterada de maneira significativa quando mudamos o valor de N (i.e. $\frac{Var(e_{oN})}{Var(e_{oN})} \neq 1$).

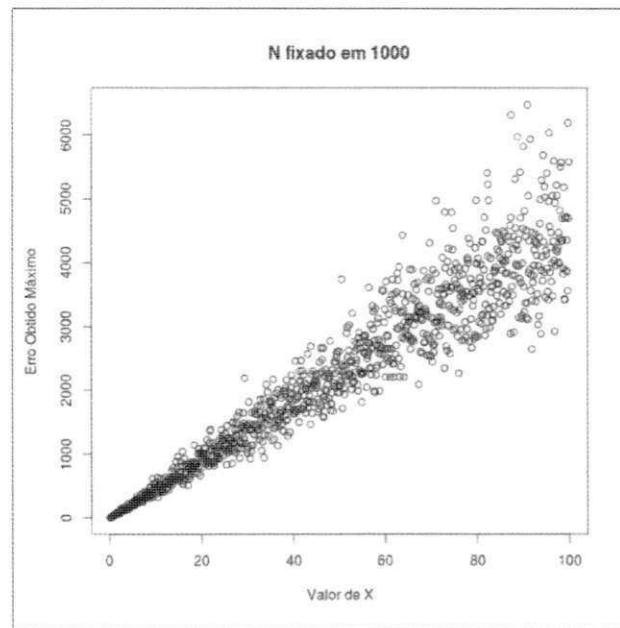


Figura 2.2: Analisando o erro variando-se X .

Para analisar como o valor de N afeta a precisão dos dados, nós podemos representar o gráfico como ilustrado na Figura 2.3. Esta relação é exponencial, conforme iremos apresentar posteriormente.

2.2.4 Analisando o Erro em Função de X e N

Usando o design de experimentos "*K-factorial*" para testar o efeito de X e N no erro, nós variamos ambos (X de 0.1 até 100 e N de 1 até 1000). Após isso, obtemos várias amostras para cada configuração e armazenamos os piores casos (quando o erro é o máximo na configuração específica). Calculando uma regressão curvilínea com logaritmos, nós encontramos o seguinte modelo com um R^2 de 97.5%:

$$\ln(e_p) = 0.32 + \ln(X) + \frac{\ln(N)}{2}$$

$$X = \frac{0,726 \cdot e_p}{\sqrt{N}} \quad (2.4)$$

Com isso, temos a nossa terceira conclusão:

- A hipótese nula H_{2-0} é falsa e a hipótese alternativa H_{2-1} é verdadeira:

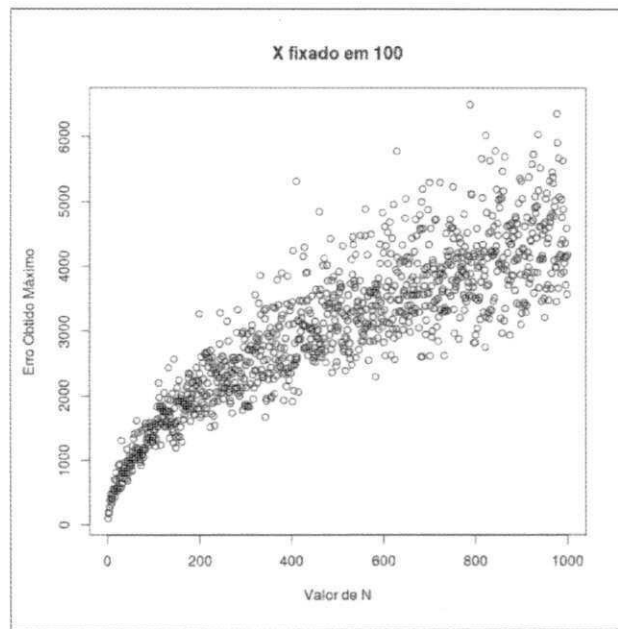


Figura 2.3: Analisando o erro variando-se N .

- Existe uma função $f(X, N) = e_p$ e conseqüentemente $f(e_p, N) = X$. Essa função é $X = \frac{0,726 \cdot e_p}{\sqrt{N}}$

Para mostrar como o modelo empírico pode ser usado, suponha que a concessionária de energia quer calcular o consumo total de um consumidor no fim de um mês de 31 dias. As medições são coletadas a cada 10 minutos, i.e., tem-se um total de $N = 4464$ medições. Se o erro máximo permitido (e_p) pela concessionária (ou por uma agência regulamentadora ou pelo próprio consumidor) é um valor percentual que corresponde a, por exemplo, 2 kWh, o valor obtido de X é: $0,726 \cdot 2 / \sqrt{4464} = 0,0217$.

Foram feitos alguns experimentos (com 1000 amostras) para gerar erros obtidos usando a Equação 2.3 e foi observado que os erros ficam de fato entre -2 e 2 kWh, como apresentado na esquerda da Figura 2.4. Os valores que estão fora do intervalo de -2 a 2 correspondem a aproximadamente 2,5% e representam os valores que a regressão não dá cobertura, uma vez que o R^2 obtido foi de 97,5%. O intervalo de confiança na direita da Figura 2.4 foi obtido usando um nível de significância de 95% e como podemos observar, a média dos erros obtidos é próxima de zero.

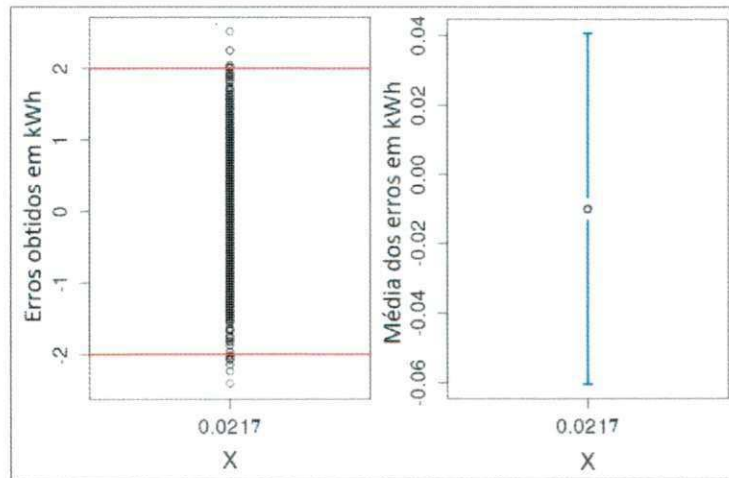


Figura 2.4: Distribuição e intervalo de confiança dos erros obtidos (e_o) para $X = 0,0217$, $N = 4464$ e $e_p = 2kWh$.

2.3 Modelo Analítico

Utilizando teoria da probabilidade, também desenvolvemos um modelo analítico que é bem próximo e correspondente do modelo obtido empiricamente.

Seja x_i uma variável aleatória uniformemente distribuída entre $-X$ e X . A sua variância é $\sigma_x^2 = \frac{(X - (-X))^2}{12} = \frac{X^2}{3}$ [25]. Para N relativamente grande, o teorema do limite central assegura que o erro obtido irá seguir uma distribuição normal com média $\mu_{e_o} = 0$ e variância $\sigma_{e_o}^2 = N^2 \left(\frac{\sigma_x^2}{N} \right) = N \cdot \sigma_x^2 = \frac{N \cdot X^2}{3}$ (da Equação 2.3 perceba que e_o é N vezes a média dos x_i).

Em outras palavras, nós podemos encontrar a probabilidade do erro obtido estar entre dois valores usando a distribuição normal descrita na Equação 2.5.

$$e_o \sim N \left(0, \frac{N \cdot X^2}{3} \right) \quad (2.5)$$

Utilizando o mesmo exemplo anterior, suponha que a concessionária de energia quer obter o consumo total de um consumidor no fim de um mês de 31 dias. Com medições a cada 10 minutos, tem-se um total de $N = 4464$ medições. Se o erro máximo permitido é de $e_p = 2kWh$, temos que encontrar a variância $\sigma_{e_o}^2$ da distribuição normal tal que a probabilidade do erro estar entre -2 kWh e 2 kWh seja alta, e.g., 0.98 ($P(-2 \leq e_o \leq 2) = 0.98$). Essa variância é $\sigma_{e_o}^2 = 0,739113$. Portanto, $X = \sqrt{\frac{3 \cdot \sigma_{e_o}^2}{N}} = \sqrt{\frac{3 \cdot 0,739113}{4464}} = 0,0222$. Este

resultado é próximo do resultado obtido com o modelo empírico. Várias outras configurações foram testadas e os resultados foram próximos, servindo isto como uma forma de validar ambos os modelos.

2.4 Métricas

Como métrica de utilidade para verificar o quão é preciso o valor final obtido pela concessionária de energia, estamos utilizando o erro percentual entre a soma dos dados ofuscados e a soma dos dados reais (i.e., e_o em porcentagem). Um erro percentual perto de 0% significa um alto nível de utilidade, enquanto que um erro percentual distante de 0% significa um baixo nível de utilidade. Como pode ser observado, se no fim do período de faturamento o erro obtido for positivo, o consumidor estará pagando um pouco mais do que realmente consumiu. Porém, talvez no próximo período de faturamento o erro obtido seja negativo e o consumidor tenha que pagar menos. Isto é aceitável porque como apresentado anteriormente, o erro segue uma distribuição normal com média zero. De fato, esta imprecisão pode ser tratada como uma penalidade que o consumidor terá que cumprir para poder obter alguma privacidade em seus dados.

Uma observação importante é que para propósitos de faturamento, que podem ser mais sensíveis aos erros, o Smart Meter pode acumular a soma de todos os números aleatórios adicionados até então e enviar isto junto com a última medição do período de faturamento, garantindo que o erro de faturamento seja zero e ainda sem disponibilizar informação sobre o perfil de consumo detalhado.

Uma funcionalidade chave da abordagem é a possibilidade de permitir que o consumidor escolha o seu próprio nível de privacidade. Se para entrar no Smart Grid um consumidor requer um alto nível de privacidade, ele pode divulgar os seus dados ofuscados usando o valor máximo possível de X . Se após um intervalo de tempo este mesmo consumidor está mais convencido em contribuir com o Smart Grid através da divulgação dos seus dados (e.g., em troca de uma redução na tarifa proposta pela concessionária de energia), ele pode divulgar dados parcialmente ofuscados e que revelam mais informações sobre o seu perfil. Se após outro intervalo de tempo este consumidor estiver totalmente convencido dos benefícios em divulgar os seus dados (e.g., vendedores detectando eletrodomésticos que não estão funci-

onando corretamente e sugerindo novos eletrodomésticos para este consumidor), ele pode divulgar os seus dados reais, i.e., usando um X igual à zero.

Uma vez que a abordagem de ofuscamento considera apenas os dados que são divulgados para a concessionária, isto não afeta a possibilidade do consumidor analisar o seu próprio perfil real de consumo dentro de sua casa. De fato, esta é uma boa maneira para gerenciar o consumo (e.g., transferência de picos) e economizar dinheiro [15].

Para monitoramento de carga, um erro obtido pode não ser considerado crítico, uma vez que a concessionária também possui outras alternativas para obter dados precisos (e.g., o fluxo de carga transmitido para uma região). A solução proposta provê mais informação (que pode ser utilizada para muitos propósitos, como detecção de roubo e vazamento) para a concessionária ao passo que preserva a privacidade dos consumidores.

Com adição de ruído, a informação de quão similar são os dados perturbados e os dados originais é crucial [30]. Existem várias métricas propostas na literatura, mas a princípio estamos usando a correlação entre o perfil ofuscado de consumo e o perfil real como métrica de privacidade. A medida de dependência entre duas quantias mais familiar é a correlação de Pearson. O coeficiente de correlação entre dois perfis A e B com valores esperados μ_A e μ_B e desvios padrões σ_A e σ_B , é definido como:

$$\text{corr}(A, B) = \frac{E[(A - \mu_A) \cdot (B - \mu_B)]}{\sigma_A \cdot \sigma_B} \quad (2.6)$$

onde E é o operador de valor esperando e corr uma notação para correlação.

Em experimentos, quando o valor de X aumenta, a correlação tende a zero. Portanto, para uma correlação próxima de zero tem-se um alto nível de privacidade, enquanto que para uma correlação próxima de 1 tem-se um baixo nível de privacidade. Outras métricas como informação mútua, erro quadrático médio e relação sinal ruído também foram analisadas, conforme apresentamos no Capítulo 4.

O uso de correlação como métrica de privacidade talvez não detecte se o comportamento do consumidor está realmente sendo escondido. Com o nosso modelo matemático, quando o número de medições em um período de tempo é grande, obtém-se um melhor valor de X para ofuscar cada medição individual. Como pode ser visto na Equação 2.4, X é oposto à N , mas realizar mais medições também implica em valores de consumo menores para cada medição individual. Por exemplo, considere o valor calculado de X igual a 0,2 para uma medição de

0,2 kWh em 30 minutos; se essa medição for dividida em duas medições de 15 minutos, os valores podem ser duas medições de 0,1 kWh e o novo valor de X é 0,1414. É melhor ofuscar um valor de 0,1kWh com 0,1414 do que ofuscar um valor de 0,2kWh com 0,2. Entretanto, sabemos que realizar menos medições também implica em maior privacidade. Considerando um exemplo extremo, claramente obtém-se maior privacidade com uma única medição do consumo total no fim do período de faturamento do que com medições a cada 15 minutos. Nossa métrica de correlação não está considerando estes casos.

Quando a frequência de medições aumenta, é fato que o nível de privacidade diminui, mas o poder de ofuscamento da abordagem também aumenta. Com isso, uma vez que uma métrica geral como correlação não é suficiente, outros estudos devem ser feitos para mensurar a privacidade do consumidor. O Capítulo 4 apresenta uma análise de outras métricas e o Capítulo 5 valida a privacidade através de possíveis ataques à solução.

Capítulo 3

Demonstrações da Abordagem

3.1 Faturamento para um Consumidor Residencial

Os dados utilizados no exemplo a seguir são medições coletadas a cada 30 minutos de um consumidor residencial real (anonimizado) da Irlanda [13].

Como um exemplo inicial, iremos assumir que o período de faturamento é de um mês e que o preço da energia é constante (sem nenhuma política de tarifação e diferenciada por horário). A Figura 3.1 mostra o perfil completo de um consumidor residencial ao longo de um mês (Março de 2010).

O consumo deste consumidor durante Março é de 165,04 kWh. Considerando um erro máximo de 5% (8,352 kWh) e usando o modelo obtido empiricamente, o seguinte valor de X foi obtido (para medições a cada 30 minutos, $N = 1488$).

$$X = \frac{0,726 \cdot 8,352}{\sqrt{1488}} = 0,1572$$

A Figura 3.2 apresenta o perfil mensal ofuscado usando o valor obtido de X . Um gráfico com a comparação entre um perfil diário real (10 de Março) e o perfil diário ofuscado também foi traçado, conforme apresentado na Figura 3.3. O coeficiente de correlação entre o perfil mensal real e o perfil mensal ofuscado é de 0,489 (nossa métrica de privacidade). Mesmo que este nível de privacidade não seja considerado muito alto, as medições foram feitas com uma baixa resolução (30 minutos é considerado um período longo), o que contribui para o nível de privacidade, como discutido anteriormente.

Foi considerado que, no fim do mês, a concessionária de energia computa o consumo

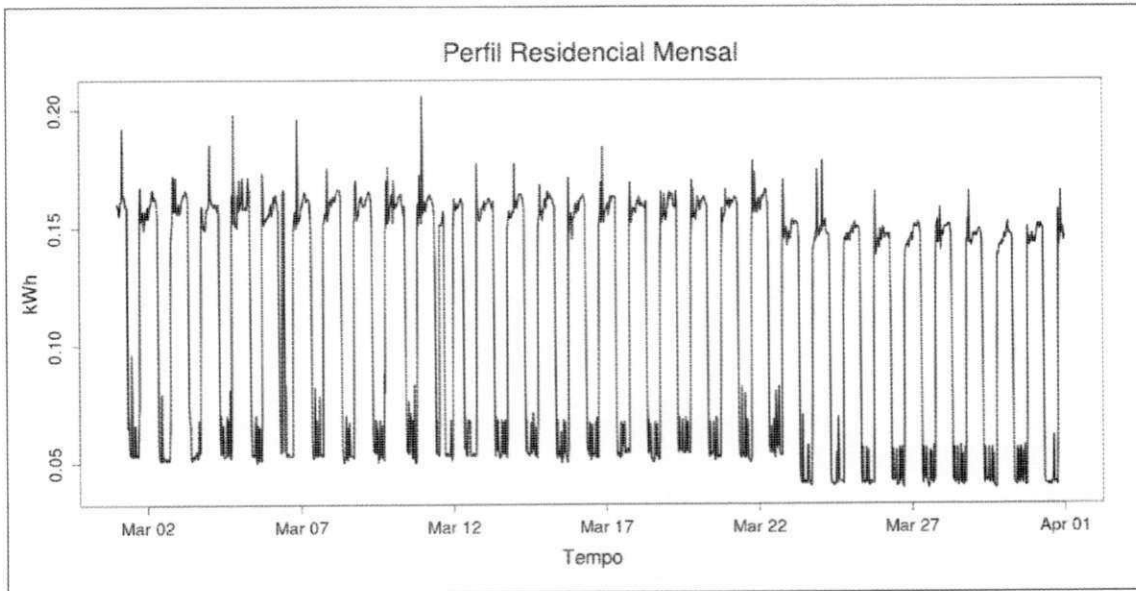


Figura 3.1: Perfil residencial mensal com medições a cada 30 minutos.

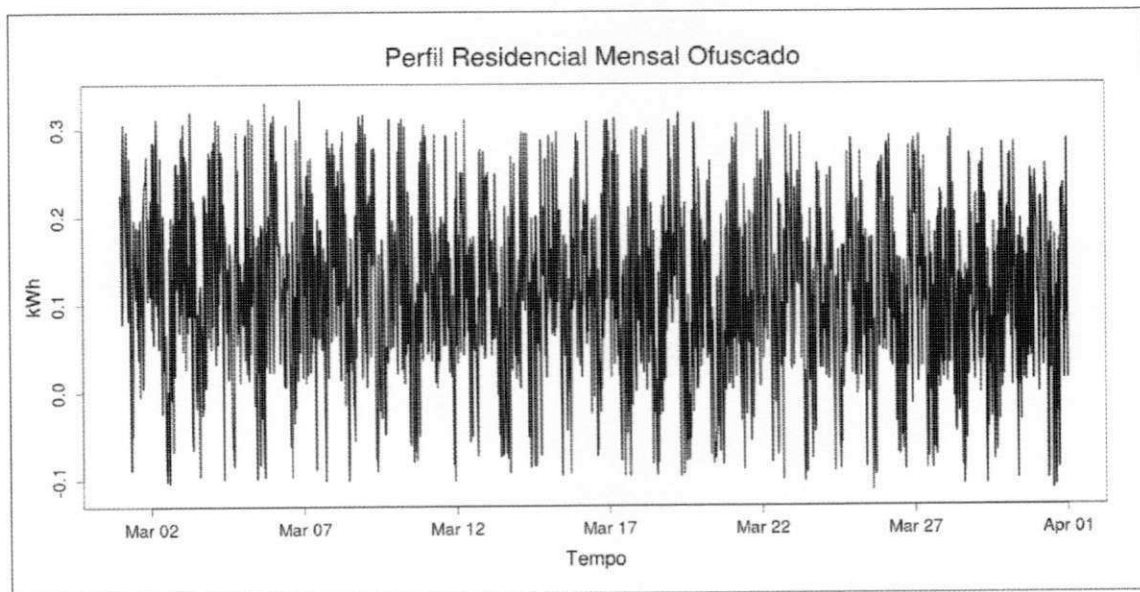


Figura 3.2: Perfil residencial mensal ofuscado com medições a cada 30 minutos.

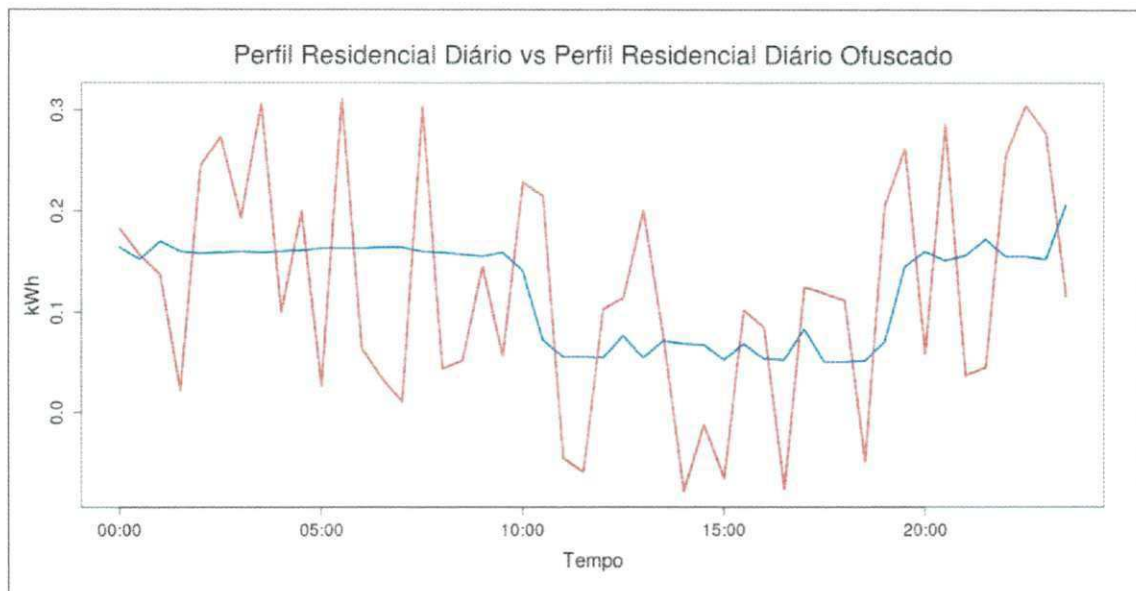


Figura 3.3: Perfil residencial diário real (azul) versus perfil residencial diário ofuscado (vermelho) com medições a cada 30 minutos.

total deste consumidor para propósitos de faturamento. Somando os valores ofuscados informados pelo consumidor neste exemplo específico, a concessionária de energia obteve um valor de 165,346 kWh (o valor real é de 167,04 kWh). A diferença entre estes dois valores é um erro de -1,01% (métrica de utilidade), o que é menor do que o erro máximo permitido (5%). Entretanto, se o firmware do medidor acumular a soma dos números aleatórios adicionados e enviar isto com a última medição do mês, o erro é zero.

3.2 Faturamento para um Consumidor Industrial

Os dados usados no exemplo a seguir são medições coletadas a cada minuto de um consumidor industrial real (anonimizado) dos Estados Unidos [12]. Para uma melhor visualização do perfil em um gráfico, esses dados foram transformados em medições a cada 10 minutos (mas deve-se observar que alta resolução implica em um melhor ofuscamento).

Foi assumido que o período de faturamento é de um mês e que o preço da energia é constante (sem nenhuma política de tarifação diferenciada por horário). A Figura 3.4 mostra o perfil completo de um consumidor industrial ao longo de um mês (Março de 2012).

O consumo deste consumidor durante Março foi de 283.959 kWh. Considerando que o

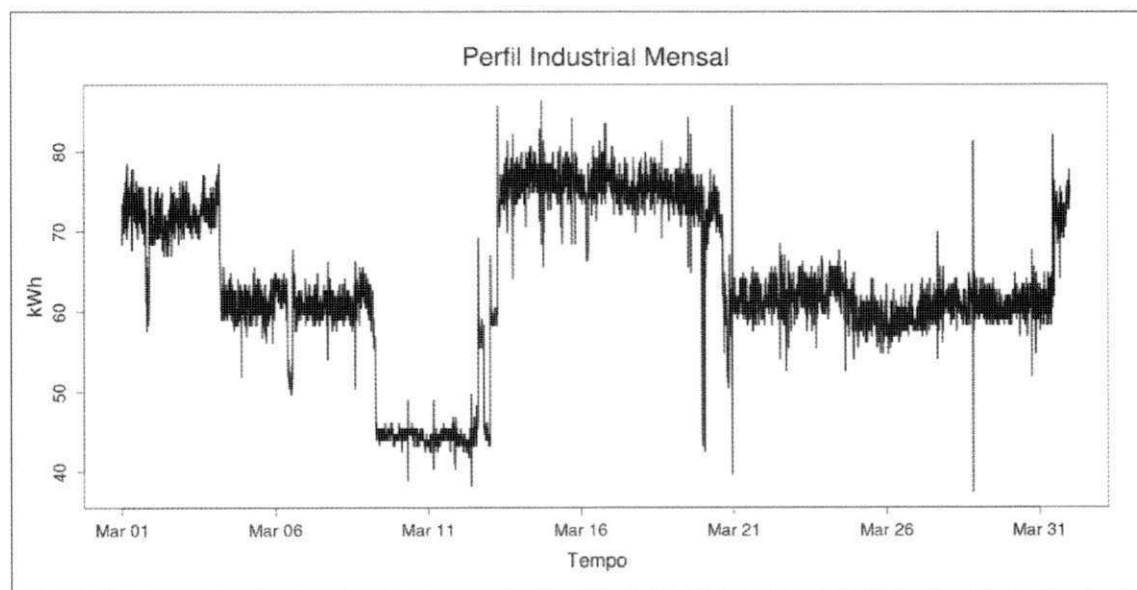


Figura 3.4: Perfil industrial mensal com medições a cada 10 minutos.

erro máximo permitido é de 5% (14.197,95 kWh), o valor obtido de X (para medições de 10 minutos, $N = 4464$) é 154,2766.

A Figura 3.5 apresenta o perfil mensal ofuscado usando o valor obtido de X . Um gráfico com a comparação entre um perfil diário real (12 de Março) e o perfil diário ofuscado também foi traçado, conforme apresentado na Figura 3.6. Como se pode observar, o nível de privacidade obtido foi tão alto que o perfil real (a linha azul) é quase uma reta em comparação ao perfil ofuscado (a linha vermelha). De fato, para este exemplo, o coeficiente de correlação entre o perfil mensal real e o perfil mensal ofuscado é de 0,109 (métrica de privacidade).

Foi considerado que, no fim do mês, a concessionária de energia computa o consumo total deste consumidor para propósitos de faturamento. Somando os valores ofuscados informados pelo consumidor neste exemplo específico, a concessionária de energia obteve um valor de 292.536,6 kWh, mas o valor real é de 283.959 kWh. A diferença entre estes dois valores é um erro de 3,02% (métrica de utilidade), o que é menor do que o erro máximo permitido (5%).

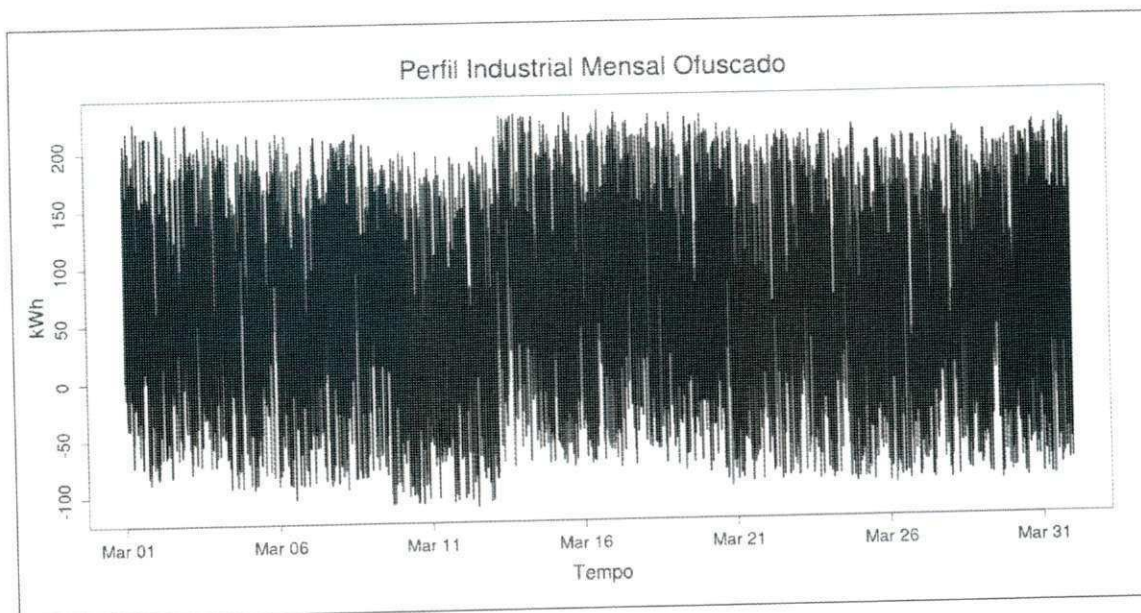


Figura 3.5: Perfil industrial mensal ofuscado com medições a cada 10 minutos.

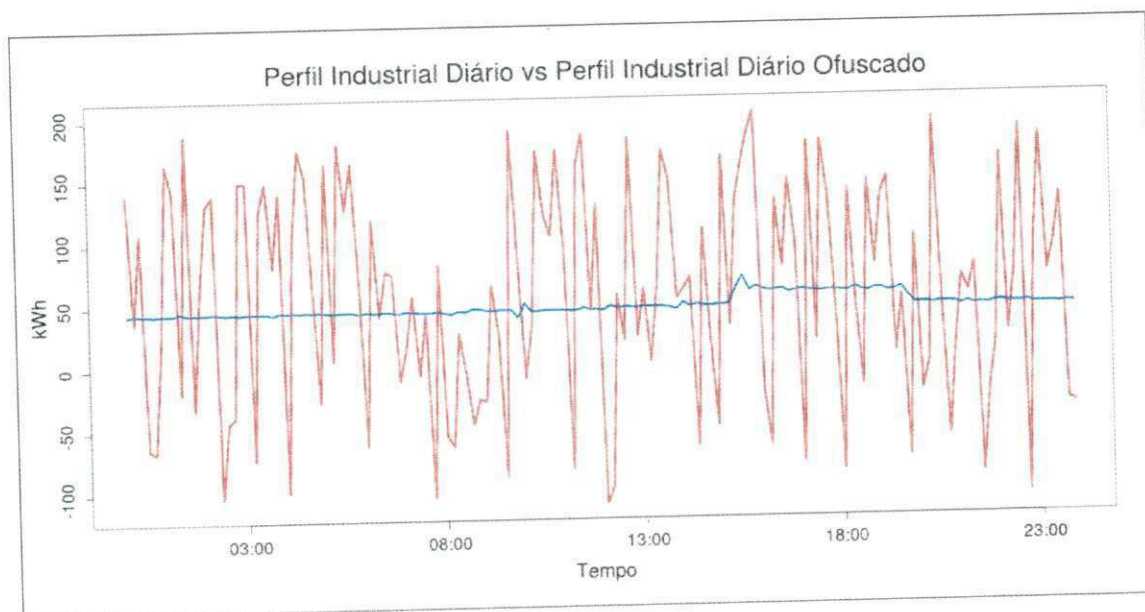


Figura 3.6: Perfil industrial diário real (azul) versus perfil industrial diário ofuscado (vermelho) com medições a cada 10 minutos.

3.3 Faturamento com Política de Horário

Também foi considerado um exemplo de faturamento com política de horário. No Brasil, a ANEEL (Agência Nacional de Energia Elétrica) estabelece um regulamento (PRORET – Procedimentos de Regulação Tarifária [9]) que fixa três tipos de tarifas de acordo com o período do dia:

- Ponta: Três horas consecutivas definidas pela concessionária de energia de acordo com a curva de carga da sua rede elétrica;
- Intermediário: Duas horas, sendo uma hora imediatamente antes e outra hora imediatamente depois do período de ponta;
- Fora de Ponta: As horas complementares (i.e., excluindo os períodos de ponta e intermediário).

A Figura 3.7 mostra um exemplo dos tipos de tarifas estabelecidas pela ANEEL. Os três períodos podem ser determinados pelas somas das colunas da matriz da Figura 2.1. Neste exemplo, nós supomos que a concessionária de energia estabeleceu as seguintes políticas de tarifa baseadas no horário: Ponta: 16:00 ~ 19:00; Intermediário: 15:00 ~ 16:00 e 19:00 ~ 20:00; Fora de ponta: 00:00 ~ 15:00 e 20:00 ~ 00:00.

Para ofuscar um perfil com diferentes períodos de tarifação, valores de X devem ser computados separadamente para cada período. Para o mesmo consumidor industrial da Figura 3.4, a Tabela 3.1 apresenta os consumos totais reais, os valores de X , os consumos totais obtidos (do perfil ofuscado) e os erros obtidos para cada período de tarifação. O perfil ofuscado obtido foi parecido com o apresentado na Figura 3.5, mas o coeficiente de correlação obtido foi igual a 0,131.

Período	Total real (kWh)	X	Total obtido (kWh)	Erro obtido (max. 5%)
Ponta	35.619,12	54,735	35.585,56	-0,094%
Intermediário	23.631,84	44,476	22.931,55	-2,963%
Fora de Ponta	224.708,1	137,21	224.946,6	0,106%

Tabela 3.1: Resultados de ofuscamento com três tipos de tarifas.

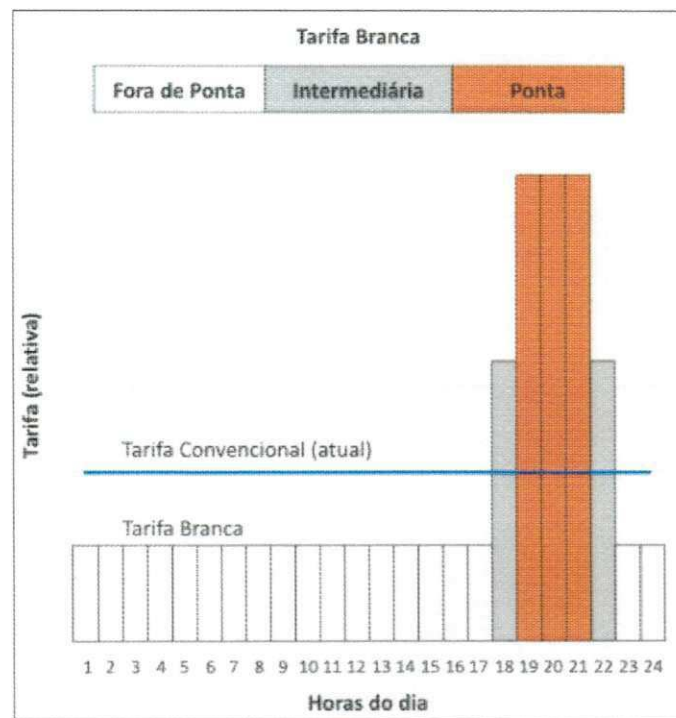


Figura 3.7: Exemplos de tipos de tarifas estabelecidas pela ANEEL [9].

3.4 Monitoramento de uma Região

Se cada consumidor ofuscar os seus dados baseando-se no período de faturamento (cada consumidor calculando o seu próprio X), a concessionária de energia obterá valores precisos para as linhas da matriz. Porém, para obter valores precisos para as colunas, o número de consumidores deve ser o maior possível. Um valor mais alto de M (Figura 2.1) implica em uma maior precisão, pois os consumidores ofuscam os seus dados baseados no período de faturamento. Deseja-se que o ofuscamento seja despercebido nos dados resultantes usados para gerenciamento de carga.

Os dados usados no exemplo a seguir são medições coletadas a cada 30 minutos de consumidores residenciais reais (anonimizados) da Irlanda [13].

Foi considerado que a concessionária de energia quer saber o consumo total de uma região ao longo do tempo para propósitos de monitoramento de carga (e.g., encontrar picos, detecção de perdas não técnicas, previsão de carga e muitas outras aplicações). Para obter valores precisos usando os dados ofuscados enviados pelos consumidores, o número de consumidores precisa ser o maior possível. Por exemplo, usando um período de faturamento de

1 mês e medições a cada 30 minutos, tem-se $N = 1488$ medições durante Março. Com isso, para nosso exemplo, foi considerado $M = 1488$ consumidores (uma matriz quadrada). A Figura 3.8 apresenta o perfil regional durante Março de 2010 obtido através de perfis reais de consumidores versus o perfil regional obtido através de perfis ofuscados de consumidores neste exemplo específico.

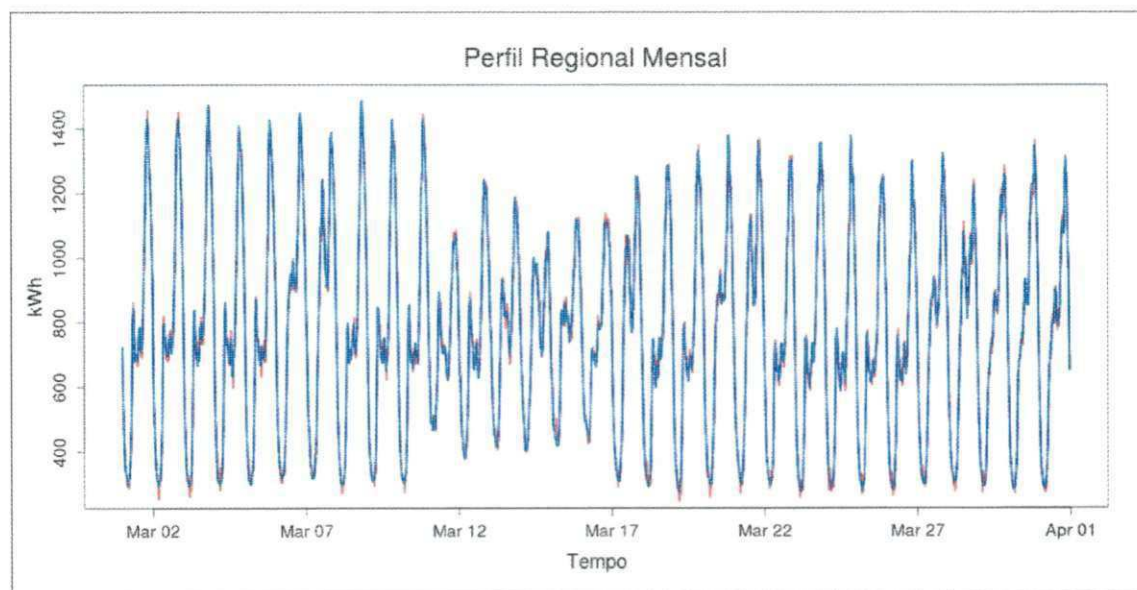


Figura 3.8: Perfil regional usando dados reais (azul) versus usando dados ofuscados (vermelho) com medições a cada 30 minutos.

Como pode ser observado na Figura 3.8, os dados parecem tão precisos de forma que ambas as linhas são quase similares (sobrepostas). Porém, isto é dependente do comportamento da população. Por exemplo, um erro obtido em kWh para um período de alto consumo tem um significado diferente do erro obtido em um período de baixo consumo. Os maiores erros ilustrados na Figura 3.8 foram obtidos em períodos de baixo consumo (e.g., durante a noite), porque enquanto consumindo menos, consumidores ainda estão ofuscando os seus dados usando um valor de X baseado no período de faturamento. A Figura 3.9 apresenta os erros obtidos ao longo do tempo para o exemplo da Figura 3.8. Para uma maior precisão, mais consumidores precisam ser incluídos (um valor maior de M).

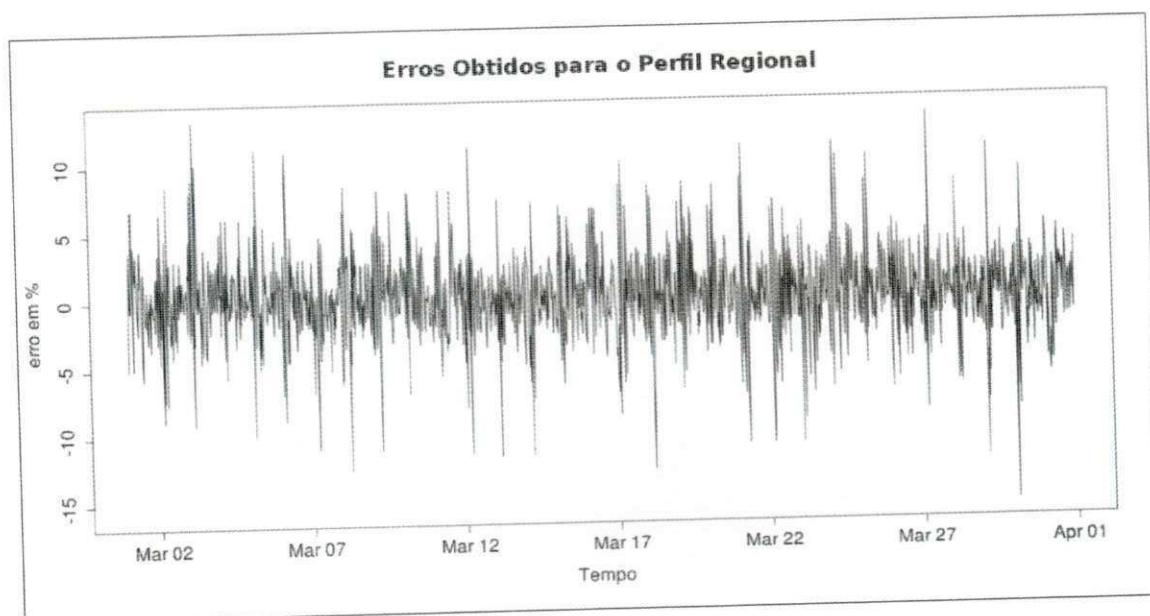


Figura 3.9: Erros obtidos ao longo do tempo para o exemplo da Figura 3.8.

Capítulo 4

Otimizações da Abordagem

4.1 Calculando o Erro Permitido (e_p)

Um dos principais desafios na utilização da nossa abordagem é o cálculo do erro máximo permitido (e_p). Nos exemplos apresentados até o momento, atribuímos a este parâmetro uma porcentagem (normalmente 5%) em relação ao valor do consumo total durante o próprio período de faturamento que queremos ofuscar. A partir de agora, chamaremos este valor de erro permitido real (e_{pr}). Entretanto, sabemos que para ofuscamento em tempo real, tal valor não é possível de ser obtido, pois o Smart Meter não possui uma forma de prever exatamente qual será o consumo total ao fim do período de faturamento corrente. Para contornar este problema, propomos algumas estratégias e através de experimentos as avaliamos em diferentes cenários.

Descrevemos e avaliamos três estratégias para calcular e_p e as denominamos de Janela Saltitante Mensal (*JSM*), Janela Saltitante Diária (*JSD*) e Janela Deslizante Mensal (*JDM*). A estratégia *JSM* calcula e_p com base no mês anterior, a estratégia *JSD* calcula com base no dia anterior e a estratégia *JDM* calcula com base nas últimas N medições (fator de deslizamento 1). De uma maneira geral, a pergunta que queremos responder através de experimentos é: qual estratégia é mais satisfatória? A resposta é aquela que provê, estatisticamente, valores de e_p mais próximos de e_{pr} . Mais ainda, tal estratégia deve se comportar bem quando aplicada em casos extremos, como apresentaremos nas próximas seções.

4.1.1 Janela Saltitante Mensal

Com a estratégia de *JSM*, calcula-se um único valor de e_p para ser usado durante todo o período de faturamento e consequentemente um único valor de X . Tal valor de e_p é calculado com base no consumo total do período de faturamento anterior (e.g., mês anterior). No fim de cada período de faturamento, o valor de e_p é atualizado.

Um dos principais problemas dessa estratégia é a imprecisão na situação em que o consumo total do período de faturamento anterior é muito maior do que o consumo total do período de faturamento atual (e.g., o consumidor viajou). Um outro problema é o baixo ofuscamento na situação em que o consumo do período anterior é muito menor do que o consumo do período atual (e.g., o consumidor chegou de viagem). Considere os casos extremos apresentados nas Figuras 4.1 e 4.3. São perfis obtidos do conjunto de dados da CER [13] em 2010.

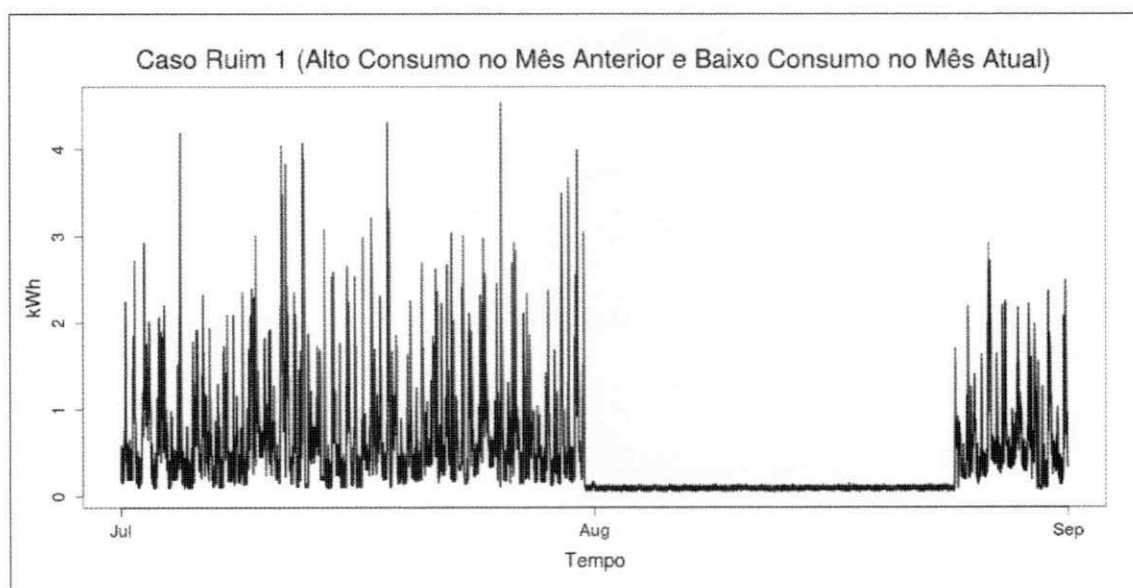


Figura 4.1: Perfil com alto consumo no mês anterior e baixo consumo no mês atual.

Ao ofuscarmos o mês de Agosto do perfil da Figura 4.1 utilizando a estratégia *JSM* e considerando um valor de e_p de 5%, obtemos um perfil como o da Figura 4.2.

Fazendo-se uma análise visual da Figura 4.2, podemos perceber um alto nível de ofuscamento e uma baixa acurácia devido ao consumo total do mês de Julho ter sido maior do que o consumo total do mês de Agosto e consequentemente o e_p calculado com essa estratégia foi

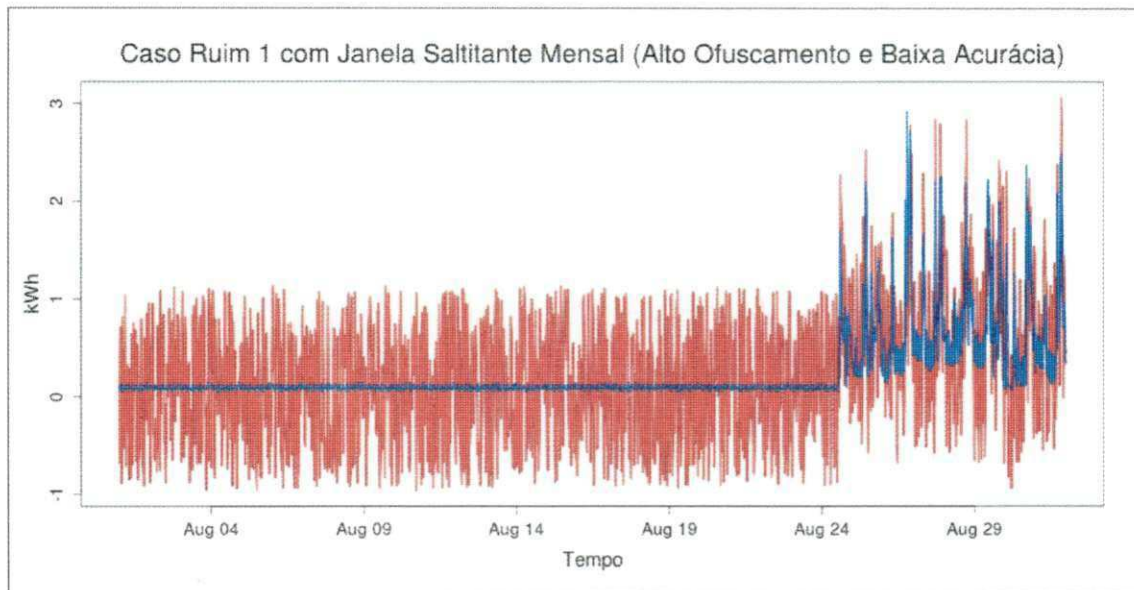


Figura 4.2: Caso ruim 1 ofuscado com *JSM*. O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.

maior do que o e_{pr} . De fato, o consumo total real durante o mês de Agosto é de $345,775kWh$ e o obtido com o perfil ofuscado foi de $313,368kWh$. A diferença entre esses valores é um erro obtido (e_o) de $9,372\%$, que está fora do e_{pr} de 5% . A correlação entre ambos os perfis é de $0,555$.

Em um outro extremo, quando ofuscamos o mês de Agosto do perfil da Figura 4.3 utilizando a estratégia *JSM* e considerando um e_p de 5% , obtemos um perfil como o da Figura 4.4.

Na Figura 4.4, podemos perceber um baixo nível de ofuscamento e uma alta precisão devido ao consumo total do mês de Julho ter sido menor do que o consumo total do mês de Agosto. O consumo total real durante o mês de agosto é de $561,812kWh$ e o obtido com o perfil ofuscado foi de $559,109kWh$. A diferença entre esses valores é um e_o de $0,481\%$, que está dentro do e_{pr} de 5% . A correlação entre ambos os perfis é de $0,987$.

Utilizando *JSM*, uma grande quantidade de informação é liberada quando ocorre uma mudança brusca no perfil devido à transação entre períodos de faturamento com consumos totais muito diferentes. Um atacante pode usar tais informações a seu favor.

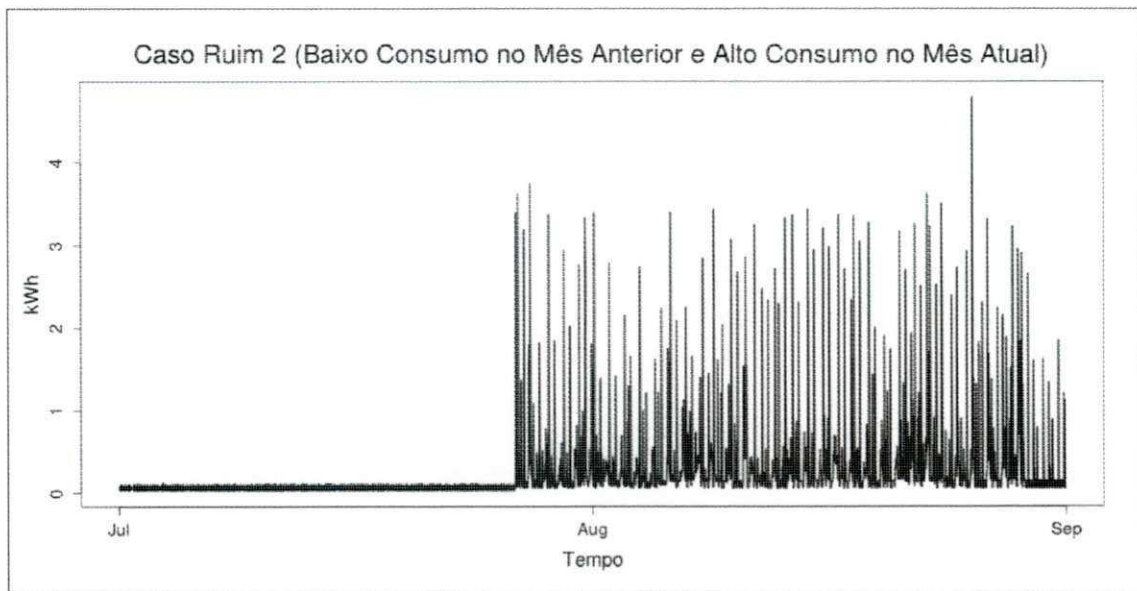


Figura 4.3: Perfil com baixo consumo no mês anterior e alto consumo no mês atual.

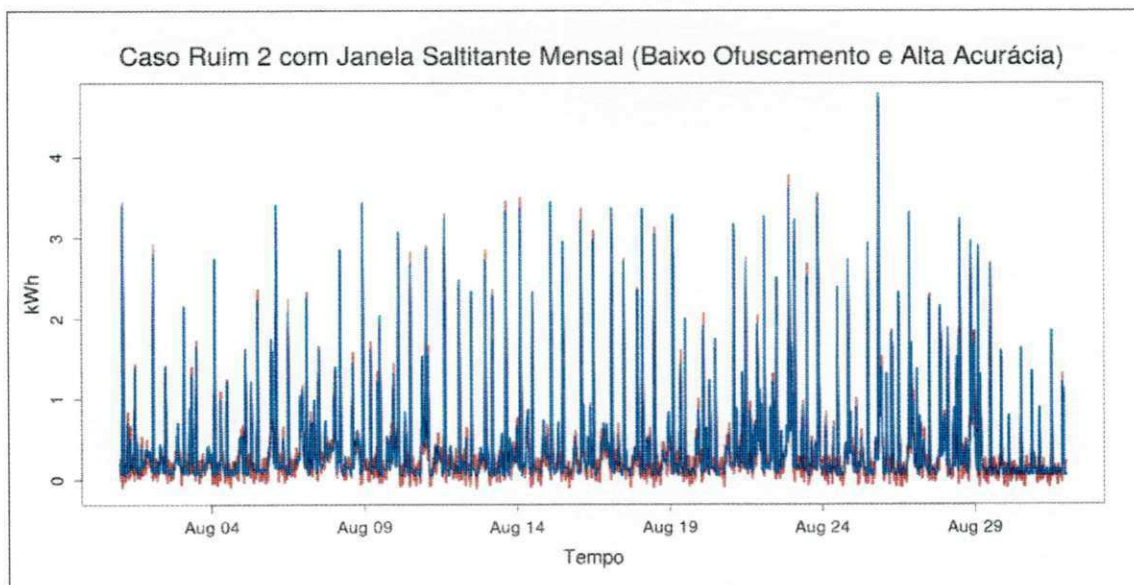


Figura 4.4: Caso ruim 2 ofuscado com *JSM*. O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.

4.1.2 Janela Deslizante Mensal

Com a estratégia de *JDM*, para cada medição calcula-se um novo valor de e_p com base no consumo total das últimas N medições (como uma janela que desliza ao longo do tempo, em que para uma medição de consumo c_i , X_i é calculado usando um erro permitido e_{pi} que é uma porcentagem do consumo total das medições de c_{i-N} até c_{i-1}).

Utilizando os mesmos casos extremos da seção anterior, ao ofuscarmos o mês de Agosto do perfil da Figura 4.1 utilizando a estratégia de *JDM* e considerando um e_p de 5%, obtemos um perfil como o da Figura 4.5.

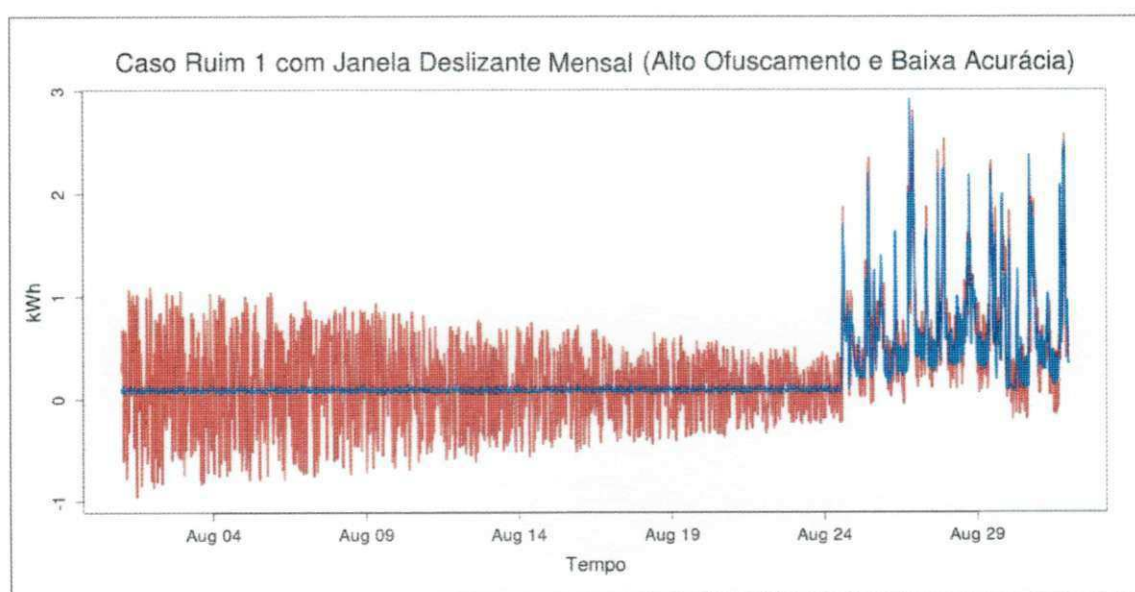


Figura 4.5: Caso ruim 1 ofuscado com *JDM*. O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.

De maneira semelhante à estratégia *JSM*, a estratégia *JDM* para o perfil da Figura 4.1 também apresenta um alto nível de ofuscamento e uma baixa acurácia. De fato, o consumo total real durante o mês de agosto é de $345,775kWh$ e o obtido com o perfil ofuscado foi de $373,96kWh$. A diferença entre esses valores é um e_o de $-8,151\%$, que está fora do e_{pr} de 5%. A correlação entre ambos os perfis é de $0,723$.

Em um outro extremo, quando ofuscamos o mês de Agosto do perfil da Figura 4.3 utilizando a estratégia *JDM* e considerando um e_p de 5%, obtemos um perfil como o da Figura 4.6.

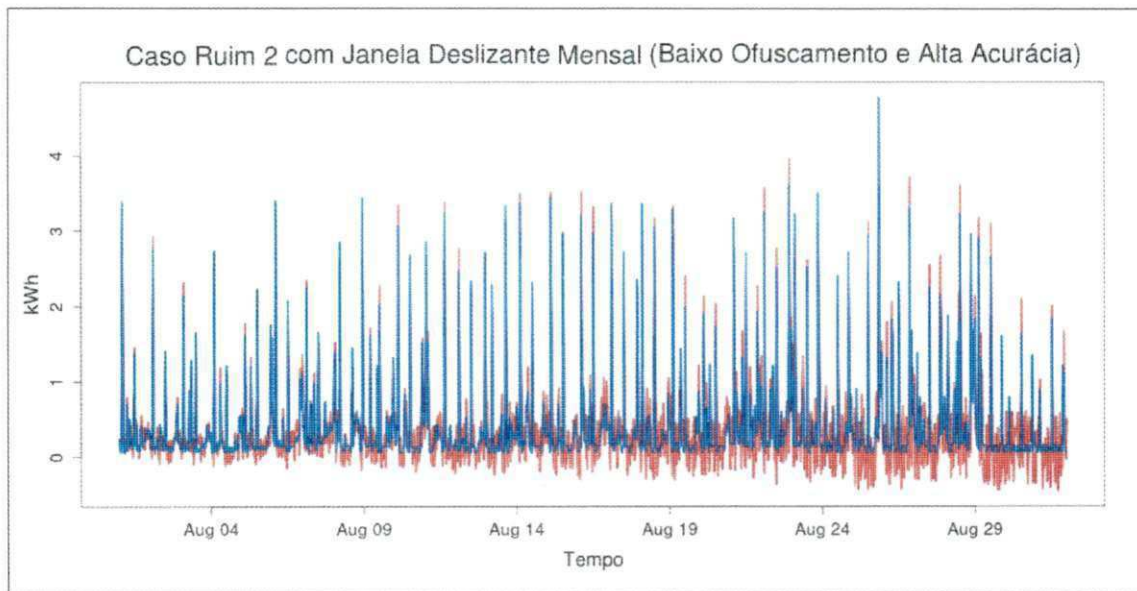


Figura 4.6: Caso ruim 2 ofuscado com *JDM*. O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.

Semelhante à Figura 4.4, na Figura 4.6 também podemos perceber um baixo nível de ofuscamento e uma alta precisão. O consumo total real durante o mês de Agosto é de $561,812kWh$ e o obtido com o perfil ofuscado foi de $566,926kWh$. A diferença entre esses valores é um e_o de $-0,91\%$, que está dentro do e_{pr} de 5% . A correlação entre ambos os perfis é de $0,932$.

Utilizando *JDM*, em relação à estratégia *JSM*, uma quantidade menor de informação é liberada para o atacante porque não existem mais as mudanças bruscas no perfil devido às transações entre períodos de faturamento com consumos totais muito diferentes.

4.1.3 Janela Saltitante Diária

A estratégia *JSD* é parecida com a estratégia *JSM*, a diferença é que ao invés de calcular um único valor de e_p para ser usado durante todo o período de faturamento, calcula-se um e_p para ser usado durante cada dia. Tal valor de e_p é calculado com base no consumo total do dia anterior multiplicado pela quantidade de dias que o período de faturamento possui. No fim de cada dia, o valor de e_p é atualizado.

Formalizando, seja K a quantidade de dias que o período de faturamento possui. Para

cada dia j a estratégia *JSD* ofusca o perfil utilizando um mesmo X_j que é calculado usando um erro permitido e_{pj} . Tal valor de e_{pj} é igual a uma porcentagem de $C_{j-1} \cdot K$, onde C_{j-1} é o consumo total durante o dia $j - 1$.

Essa estratégia resolve o problema da imprecisão na situação em que o consumo do período anterior é muito maior do que o consumo do período atual e também resolve o problema do baixo ofuscamento na situação em que o consumo do período anterior é muito menor do que o consumo do período atual. Ou seja, diferente das estratégias anteriores, a estratégia *JSD* apresenta maior estabilidade nas métricas de utilidade e privacidade mesmo em casos extremos, pois os valores calculados de e_p são mais próximos de e_{pr} .

Ao ofuscarmos o mês de Agosto do perfil da Figura 4.1 utilizando a estratégia de *JSD* e considerando um e_p de 5%, obtemos um perfil como o da Figura 4.7.

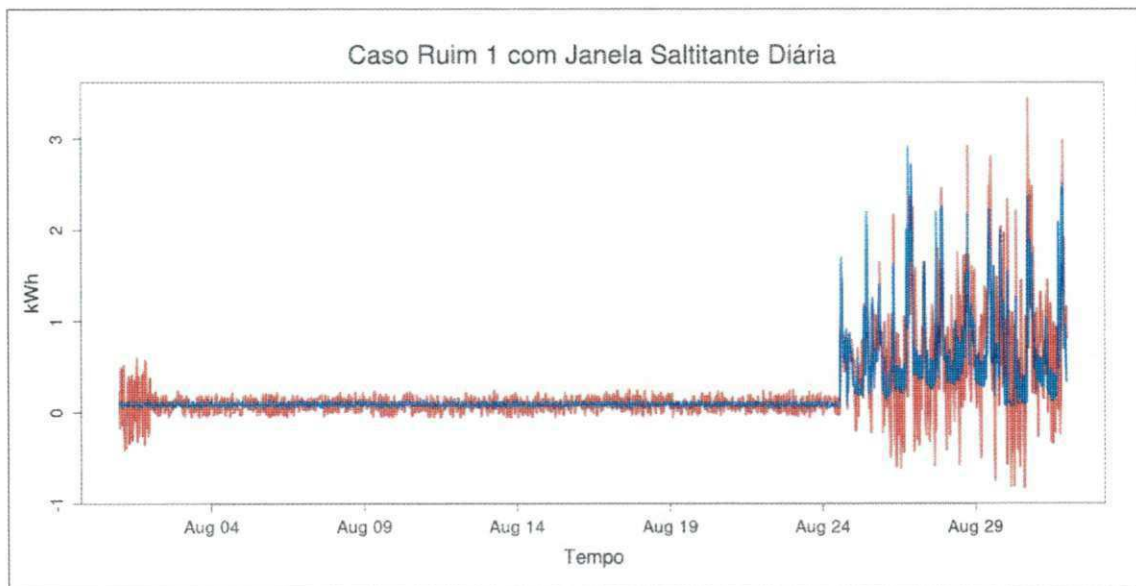


Figura 4.7: Caso ruim 1 ofuscado com *JSD*. O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.

Na Figura 4.7, o consumo total real durante o mês de agosto é de $345,775kWh$ e o obtido com o perfil ofuscado foi de $347,5154kWh$. A diferença entre esses valores é um e_o de $-0,503\%$, que está dentro do e_{pr} de 5% . A correlação entre ambos os perfis é de $0,813$.

Em um outro extremo, quando ofuscamos o mês de Agosto do perfil da Figura 4.3 utilizando a estratégia *JSD* e considerando um e_p entre 5% , obtemos um perfil como o da Figura

4.8.

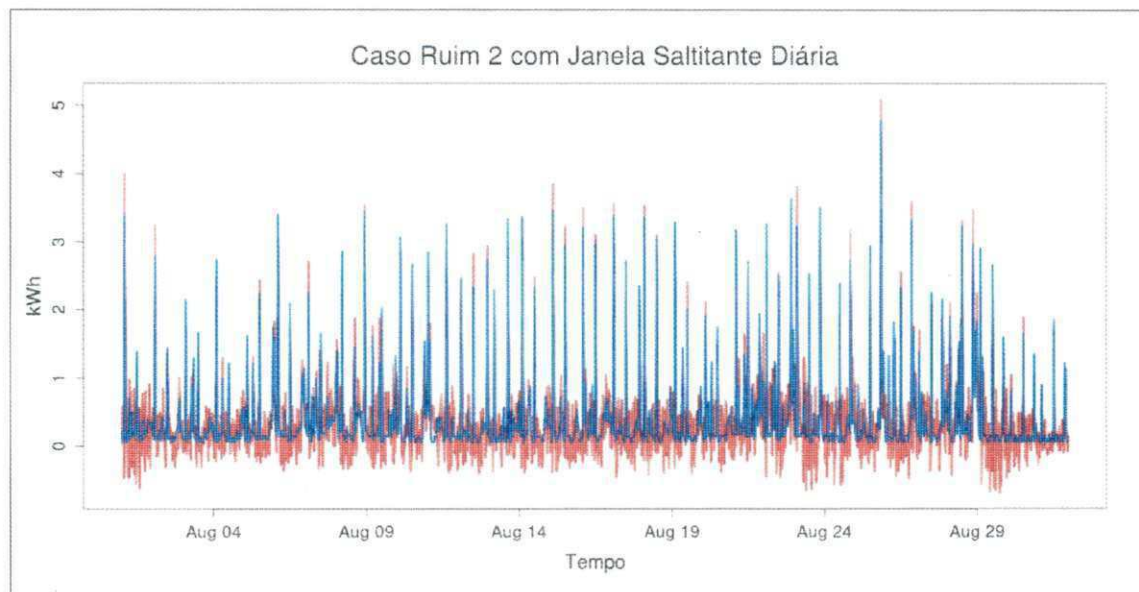


Figura 4.8: Caso ruim 2 ofuscado com *JSD*. O perfil em azul é o perfil real e o em vermelho é o perfil ofuscado.

Na Figura 4.8, o consumo total real durante o mês de agosto é de $561,812kWh$ e o obtido com o perfil ofuscado foi de $540,444kWh$. A diferença entre esses valores é um e_o de $3,803\%$, que está dentro do e_{pr} de 5% . A correlação entre ambos os perfis é de $0,876$.

Utilizando *JSD*, uma quantidade de informação ainda pode ser liberada para o atacante. Isso ocorre quando se tem uma mudança brusca no perfil devido à transação entre dias com consumos muito diferentes (semelhante à *JSM*). Entretanto, diferente das outras estratégias, o balanceamento entre utilidade e privacidade proposto pela nossa abordagem de ofuscamento ainda é mantido, conforme apresentaremos na próxima seção.

4.1.4 Conclusões Sobre o Cálculo de e_p

A Tabela 4.1 resume os resultados obtidos com as diferentes estratégias para os dois casos extremos apresentados anteriormente.

Como se pode observar, a estratégia *JSD* apresentou melhores resultados para os casos extremos. Entretanto, ainda queremos saber qual estratégia calcula valores de erros permitidos e_p mais próximos dos erros permitidos reais e_{pr} (ou seja, qual estratégia se comporta

Caso Extremo 1			Caso Extremo 2		
Estratégia	Correlação	Erro Obtido	Estratégia	Correlação	Erro Obtido
JSM	0,555	9,372%	JSM	0,987	0,481%
JDM	0,723	-8,151%	JDM	0,932	-0,91%
JSD	0,813	-0,503%	JSD	0,876	3,803%

Tabela 4.1: Resultados das estratégias para os casos extremos. Valores em vermelho significam que estão fora do esperado. Valores em verde significam que estão dentro do esperado.

melhor no geral, e não só nos casos extremos apresentados anteriormente).

Para 1000 perfis aleatórios do banco de dados da CER [13], nós comparamos os erros permitidos de cada estratégia com os erros permitidos reais. Para cada perfil, calculou-se uma média μ_{e_p} dos e_p calculados com cada estratégia e a comparou com o e_{pr} do perfil. Para a estratégia *JSM* temos que $\mu_{e_p} = e_p$, visto que utiliza-se um mesmo e_p durante todo o mês. Para a estratégia *JDM* temos que $\mu_{e_p} = \sum_{i=1}^N e_{pi}$, onde N é a quantidade de medições. Para a estratégia *JSD* temos que $\mu_{e_p} = \sum_{j=1}^K e_{pj}$, onde K é a quantidade de dias.

A Tabela 4.2 apresenta os intervalos de confiança (com níveis de confiança de 95%) dos erros absolutos (E_{abs}) entre μ_{e_p} e e_{pr} para os 1000 perfis em cada uma das estratégias, onde $E_{abs} = |e_{pr} - \mu_{e_p}|$.

Estratégia	Int. Conf. (95%) de E_{abs}
JSM	(3,402; 3,954)
JDM	(1,879; 2,184)
JSD	(0,344; 0,39)

Tabela 4.2: Intervalos de confiança dos erros absolutos entre a média dos e_p e o e_{pr} para 1000 perfis.

Conclui-se que a estratégia *JSD* além de apresentar um bom comportamento em casos extremos, ainda calcula erros permitidos mais próximos dos erros permitidos reais. Além disso, realizando experimentos com os mesmos 1000 perfis, concluiu-se que a estratégia *JSD* possui um maior poder de ofuscamento do que as outras estratégias (incluindo a própria estratégia que usa valores de erros permitidos reais e_p). A Tabela 4.3 apresenta os intervalos

de confiança das métricas de privacidade e utilidade quando ofuscamos os 1000 perfis.

Estratégia	Int. Conf. (95%) da Correlação	Int. Conf. (95%) do e_o (%)
Real	(0,8170; 0,8274)	(-0,1379; 0,1236)
JSM	(0,8171; 0,8287)	(-0,0628; 0,2046)
JDM	(0,8171; 0,8280)	(-0,1502; 0,1104)
JSD	(0,8045; 0,8146)	(-0,1193; 0,1509)

Tabela 4.3: Intervalos de confiança das métricas de privacidade e utilidade para as diferentes estratégias ao ofuscarmos 1000 perfis.

Ao realizarmos alguns testes estatísticos (ANOVA junto com Schéffe) para comparar as médias, concluiu-se com 95% confiança que a estratégia *JSD* apresenta, estatisticamente, um maior nível de ofuscamento e mantêm o nível de utilidade.

Sendo assim, considerando todas as estratégias apresentadas para calcular o e_p , concluiu-se que a estratégia de *JSD* é a melhor opção, pois: apresenta bom comportamento em casos extremos, calcula erros permitidos mais próximos dos reais e apresenta um maior poder de ofuscamento. Entretanto, conforme apresentado, o resultado obtido com uma estratégia é bastante dependente do perfil que está sendo ofuscado. Desta forma, soluções dinâmicas ou híbridas também podem ser opções.

4.2 Análise de Outras Métricas

Para a métrica de privacidade, quatro métricas foram consideradas: correlação (*corr*) [2], relação sinal-ruído (*SNR*) [30], erro quadrático médio (*MSF*) e informação mútua (*MI*) [36]. Ofuscamos um perfil qualquer (que está em azul na Figura 4.13) várias vezes usando valores de X entre 0,01 e 0,4 (escolhidos ao acaso, mas factíveis) e analisamos cada uma destas métricas.

4.2.1 Correlação

A medida de dependência entre dois conjuntos de dados mais familiar é o coeficiente de correlação produto-momento de Pearson, ou simplesmente "correlação de Pearson". Esta

medida é obtida dividindo a covariância das duas variáveis pelo produto dos seus desvios padrões.

O coeficiente de correlação $corr(A, B)$ entre duas variáveis aleatórias A e B com valores esperados μ_A e μ_B e desvios padrões σ_A e σ_B , é definido como:

$$corr(A, B) = \frac{cov(A, B)}{\sigma_A \cdot \sigma_B} = \frac{E[(A - \mu_A) \cdot (B - \mu_B)]}{\sigma_A \cdot \sigma_B}$$

onde E é o operador de valor esperando e $corr$ uma notação para correlação.

A correlação de Pearson será +1 quando houver uma perfeita relação linear positiva (correlação), -1 quando houver uma perfeita relação linear negativa (anticorrelação) e um valor entre -1 e +1 nos outros casos, indicando o grau de dependência linear entre as variáveis. A medida que o valor se aproxima de zero existe menos relação entre as variáveis. Se as variáveis são independentes, o coeficiente de correlação de Pearson será zero.

Ofuscando um perfil qualquer (que está em azul na Figura 4.13) várias vezes (10 amostras) com diferentes valores de X (variando-se de 0,01 a 0,4) foram obtidos valores de correlação entre os perfis ofuscados e o perfil original, conforme apresentado na Figura 4.9. Para uma visualização do *tradeoff* entre privacidade e utilidade, os valores absolutos dos erros obtidos também estão sendo representados.

Como se pode observar, quando o valor de X aumenta, a correlação tende a zero. Portanto, se a correlação entre o perfil ofuscado e o perfil real for zero, tem-se um alto nível de privacidade, ao passo que se a correlação for próxima de 1, tem-se um baixo nível de privacidade.

4.2.2 Relação Sinal-ruído

Em telecomunicações, a relação sinal-ruído (SNR) é uma medida que compara o nível do sinal original com o nível do ruído. Tal relação é definida como a razão entre a potência do sinal com a potência do ruído. Ou seja:

$$SNR = \left(\frac{RMS_{sinal}}{RMS_{ruído}} \right)^2$$

onde RMS é a raiz da média dos valores ao quadrado.

Uma razão maior que 1 indica mais sinal do que ruído. Informalmente, em comunicação de dados, esta razão é usada para distinguir a quantidade de informação útil da quantidade

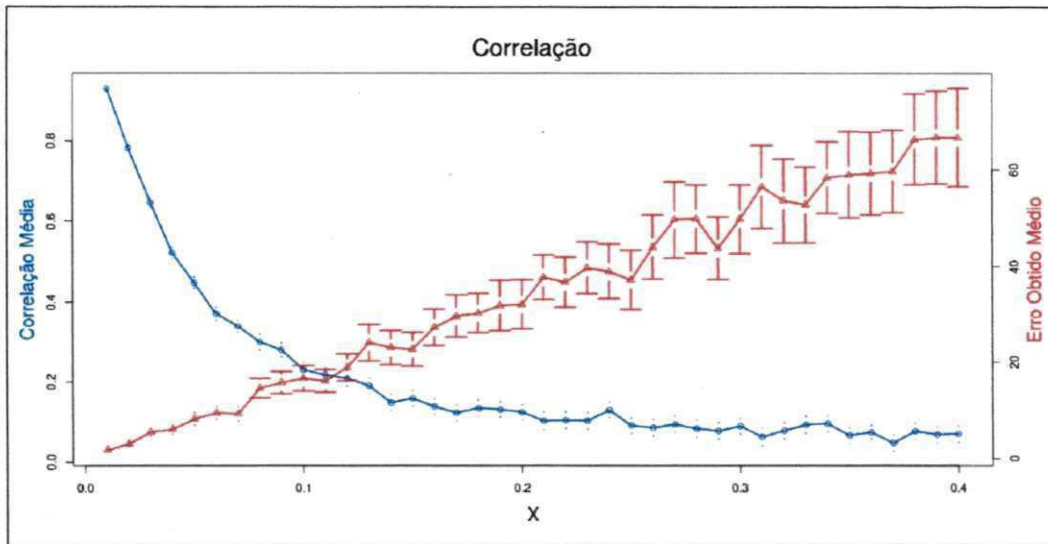


Figura 4.9: Correlações obtidas para diferentes valores de X . Os níveis de significância são de 95%.

falsa ou irrelevante. Em nosso cenário, podemos considerar o sinal como sendo o perfil original e o ruído como sendo o próprio ruído inserido para ofuscar o perfil original.

Ofuscando um perfil qualquer (que está em azul na Figura 4.13) várias vezes (10 amostras) com diferentes valores de X (variando-se de 0,01 a 0,4) obtemos valores de SNR entre o perfil original e os ruídos inseridos, conforme apresentado na Figura 4.10. Para uma visualização do *tradeoff* entre privacidade e utilidade, os valores absolutos dos erros obtidos também estão sendo representados.

Como se pode observar, quando o valor de X aumenta, o SNR tende a zero. Portanto, se o SNR entre o perfil original e o ruído inserido for zero, tem-se um alto nível de privacidade (porque existe mais ruído do que sinal), ao passo que se a relação for alta, tem-se um baixo nível de privacidade.

4.2.3 Erro Quadrático Médio

O erro quadrático médio (MSE) de um estimador é uma das várias formas de quantificar a diferença entre as estimativas e os valores reais. MSE mede a média dos "erros" ao quadrado. O erro é uma quantia que significa o valor em que a estimativa difere do valor real.

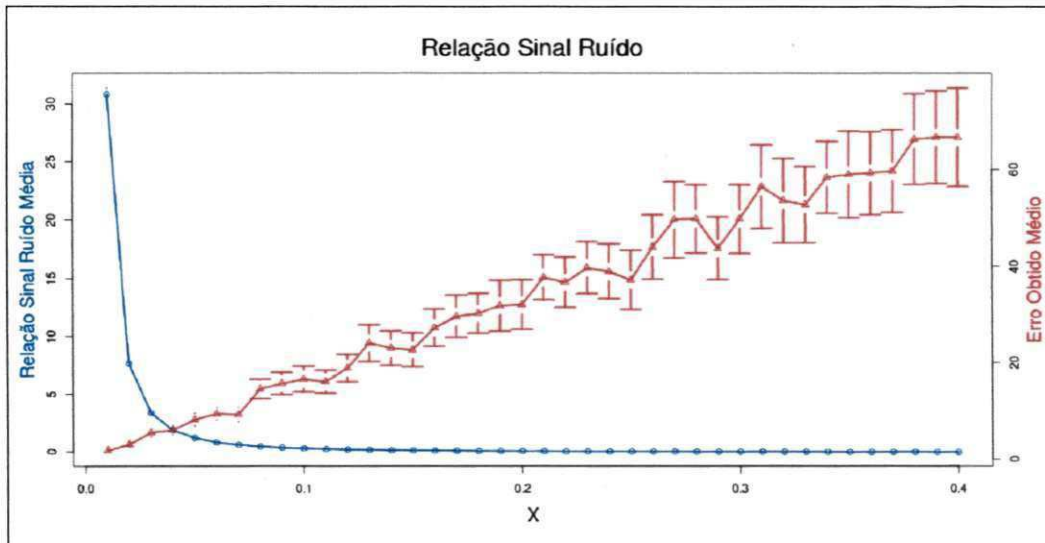


Figura 4.10: SNR obtidos para diferentes valores de X . Os níveis de significância são de 95%.

Se A é um vetor com os valores reais e B é um vetor com as estimativas (ou valores ofuscados), então o valor do MSE é:

$$MSE = \frac{1}{N} \sum_{i=1}^N (a_i - b_i)^2$$

Ofuscando um perfil qualquer (que está em azul na Figura 4.13) várias vezes (10 amostras) com diferentes valores de X (variando-se de 0,01 a 0,4) obtemos valores de MSE entre os perfis ofuscados e o perfil original, conforme apresentado na Figura 4.11. Os valores absolutos dos erros obtidos também estão sendo representados.

Como se pode observar, quando o valor de X aumenta, o MSE também aumenta. Portanto, se o MSE entre o perfil original e o perfil ofuscado for alto, tem-se um alto nível de privacidade, ao passo que se o MSE relação for baixo, tem-se um baixo nível de privacidade. Para uma visualização do *tradeoff*, também analisamos a métrica $\frac{1}{MSE}$. Entretanto, foi obtido um gráfico muito parecido com o da relação sinal-ruído na Figura 4.10.

4.2.4 Informação Mútua

A informação mútua (MI) é uma quantia que mede a dependência mútua entre duas variáveis aleatórias. Formalmente, a informação mútua entre duas variáveis aleatórias A e B pode ser

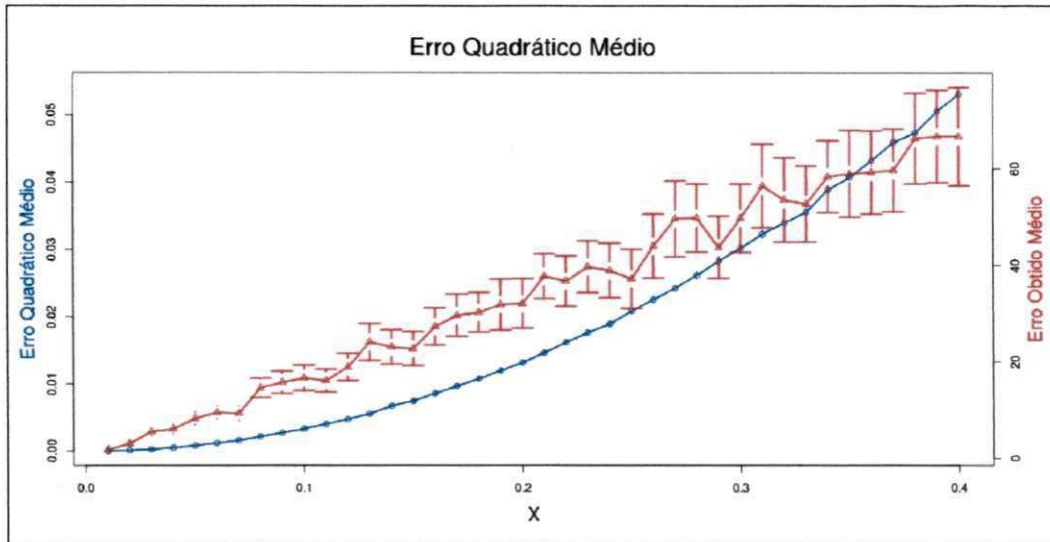


Figura 4.11: MSE obtidos para diferentes valores de X . Os níveis de significância são de 95%.

definida como:

$$MI(A, B) = \int_A \int_B p(a, b) \cdot \log \left(\frac{p(a, b)}{p(a) \cdot p(b)} \right) db da$$

onde $p(a, b)$ é a função de densidade de probabilidade conjunta de A e B , $p(a)$ e $p(b)$ são as funções de densidade de probabilidade marginais de A e B , respectivamente.

Intuitivamente a informação mútua mede a quantidade de informação que A e B compartilham: ela mede o quanto conhecer uma dessas variáveis reduz a incerteza sobre a outra. Por exemplo, se A e B são independentes, então ter o conhecimento de A não dá nenhuma informação sobre B e vice versa, portanto a informação mútua entre essas variáveis é zero. Em outro extremo, se A e B são idênticas, tendo conhecimento de A então é possível determinar B e vice versa.

Ofuscando um perfil qualquer (que está em azul na Figura 4.13) várias vezes (10 amostras) com diferentes valores de X (variando-se de 0,01 a 0,4) obtemos valores de MI entre os perfis ofuscados e o perfil original, conforme apresentado na Figura 4.12. Para uma visualização do *tradeoff* entre privacidade e utilidade, os valores absolutos dos erros obtidos também estão sendo representados.

Diferente das outras métricas, a informação mútua apresentou resultados não esperados

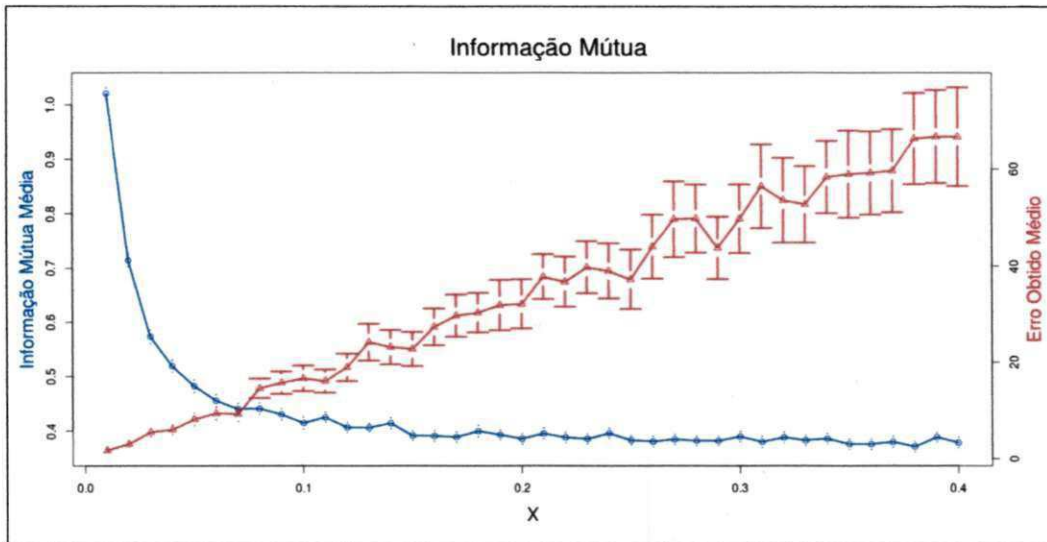


Figura 4.12: MI obtidos para diferentes valores de X . Os níveis de significância são de 95%.

quando aumenta-se o ofuscamento pois ao invés de tender a zero, o valor de MI convergiu para um valor próximo de 0,4. Valores de X maiores que os apresentados na figura 4.12 também foram testados e a convergência persistiu. Sendo assim, por apresentar resultados subjetivos, não consideramos informação mútua como uma boa métrica (é desejável que para ofuscamentos maiores a métrica acuse uma maior privacidade).

4.2.5 Investigação de Problemas

Como apresentado anteriormente, somente as métricas $Corr$, SNR e MSE apresentaram resultados esperados. Destas três métricas, a métrica $Corr$ é a que varia mais devagar, fornecendo assim uma melhor representação do balancemaneto da privacidade com relação à utilidade. Entretanto, outros cenários devem ser investigados para decidir qual a melhor métrica para ser usada.

Supondo que, por uma falta de sorte, o gerador de números aleatórios usado para ofuscar os dados esteja quebrado e sempre gerando um mesmo valor. Com isso, o perfil ofuscado resultante será apenas uma versão do perfil original deslocado na vertical, como exemplificado na Figura 4.13. Uma vez que o perfil divulgado é semelhante ao perfil original, os hábitos de

consumo ainda são expostos e não se tem privacidade neste caso.

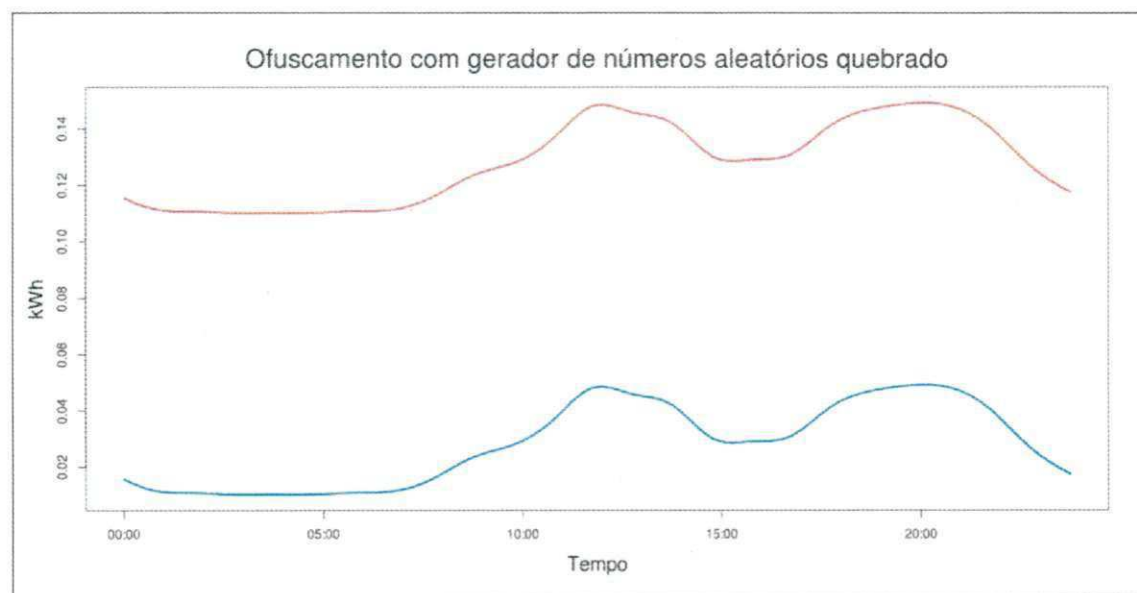


Figura 4.13: Perfil diário original (azul) versus perfil diário ofuscado com um gerador de números aleatórios quebrado (vermelho).

Neste cenário, obtemos os seguintes valores para as métricas: $Corr = 1$, $SNR = 0,09977$, $MSE = 0,01$ e $MI = 1,9366$. De acordo com a semântica destas métricas, a relação sinal-ruído e o erro quadrático médio não apresentaram bons resultados, pois identificaram que existe privacidade.

As métricas de correlação e informação mútua apresentaram bons resultados. De fato, o resultado obtido com a correlação entre o perfil original e o perfil ofuscado é o mesmo obtido com a correlação entre o perfil original e o próprio perfil original, acusando assim que não existe privacidade. O mesmo aconteceu para a informação mútua. Entretanto, como apresentamos anteriormente, existem casos em que a informação mútua não pode ser considerada uma boa métrica. Sendo assim, concluímos que correlação é a melhor métrica de privacidade a ser usada.

4.3 Comparando Distribuições de Probabilidade

Até o momento, em todos os experimentos utilizamos apenas a distribuição de probabilidade uniforme para gerar os números aleatórios que serão utilizados no ofuscamento. Entretanto,

na tentativa de otimizar a técnica proposta, foi feito um estudo com o objetivo de responder a seguinte pergunta: qual a melhor distribuição de probabilidade a ser usada para ofuscar os dados? Com isso, investigamos a seguinte hipótese: H_{0-0} : não existe uma distribuição de probabilidade que ofusca os dados melhor do que as outras distribuições.

Como apresentado anteriormente, o erro obtido segue uma distribuição normal com média zero, i.e., $e_o \sim N\left(0, N \cdot \sigma_{e_o}^2\right)$. Utilizando os mesmos procedimentos apresentados na Seção 2.3, também encontramos modelos analíticos para as seguintes distribuições: Arcoseno [39], Laplace [24], Normal [33] e U-quadrática [41], conforme apresentado na Tabela 4.4.

Distribuição	Modelo Analítico	Comentários
Arcoseno	$e_o \sim N\left(0, \frac{N \cdot X^2}{2}\right)$	X são os extremos da distribuição original
Laplace	$e_o \sim N\left(0, \frac{N \cdot 2}{\lambda^2}\right)$	λ é a taxa da distribuição original
Normal	$e_o \sim N\left(0, N \cdot \sigma_X^2\right)$	σ_X^2 é a variância da distribuição original
Uniforme	$e_o \sim N\left(0, \frac{N \cdot X^2}{3}\right)$	X são os extremos da distribuição original
U-quadrática	$e_o \sim N\left(0, \frac{N \cdot 3 \cdot X^2}{5}\right)$	X são os extremos da distribuição original

Tabela 4.4: Modelos analíticos obtidos para diferentes distribuições de probabilidade.

Para mostrar como estes modelos podem ser usados, suponha que a concessionária de energia quer obter o consumo total de um consumidor no fim de um mês de 31 dias. Com medições a cada 10 minutos, tem-se um total de $N = 4464$ medições. Se o erro máximo permitido é de $e_p = 2 \text{ kWh}$, temos que encontrar a variância $\sigma_{e_o}^2$ da distribuição normal tal que a probabilidade do erro estar entre -2 kWh e 2 kWh seja alta, e.g., 0.98 ($P(-2 \leq e_o \leq 2) = 0.98$). Essa variância é $\sigma_{e_o}^2 = 0,739113$. Substituindo este valor em cada um dos modelos encontrados, obtemos os parâmetros que faltam. A Figura 4.14 apresenta as funções de densidade de probabilidade de cada distribuição original (usada para ofuscar os dados) neste cenário.

Para cada distribuição foram feitos alguns experimentos (com 1000 amostras) para gerar erros obtidos e foi observado que os erros ficam de fato entre -2 e 2 kWh , como apresentado na esquerda da Figura 4.15. Os valores que estão fora do intervalo de -2 a 2 correspondem a aproximadamente 2%, conforme esperado.

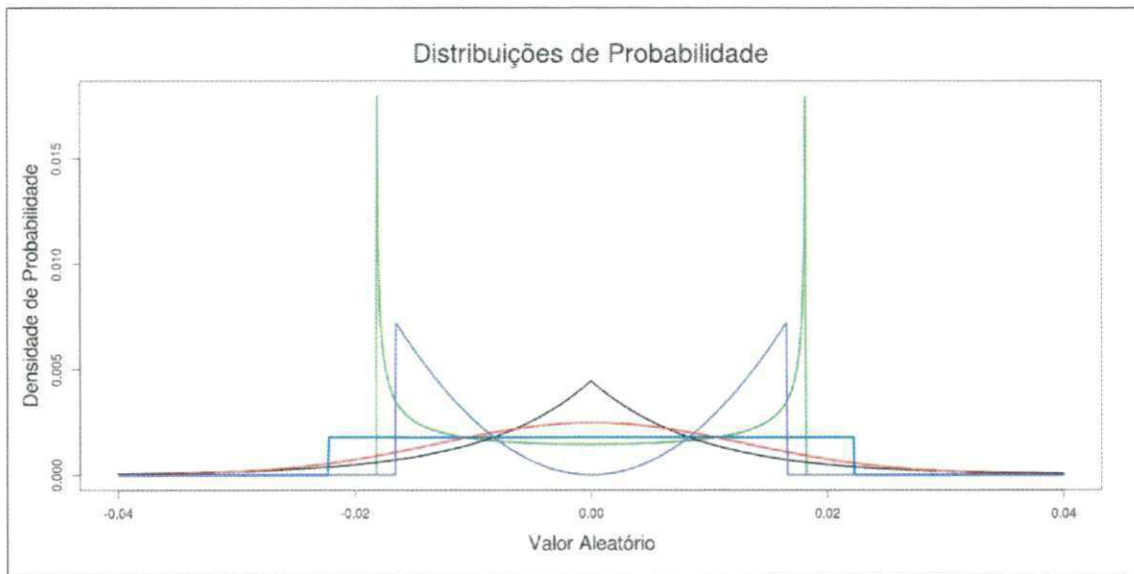


Figura 4.14: Densidades de probabilidade. Arcoseno está representada em verde, Laplace em preto, Normal em vermelho, Uniforme em azul e U-quadrática em roxo.

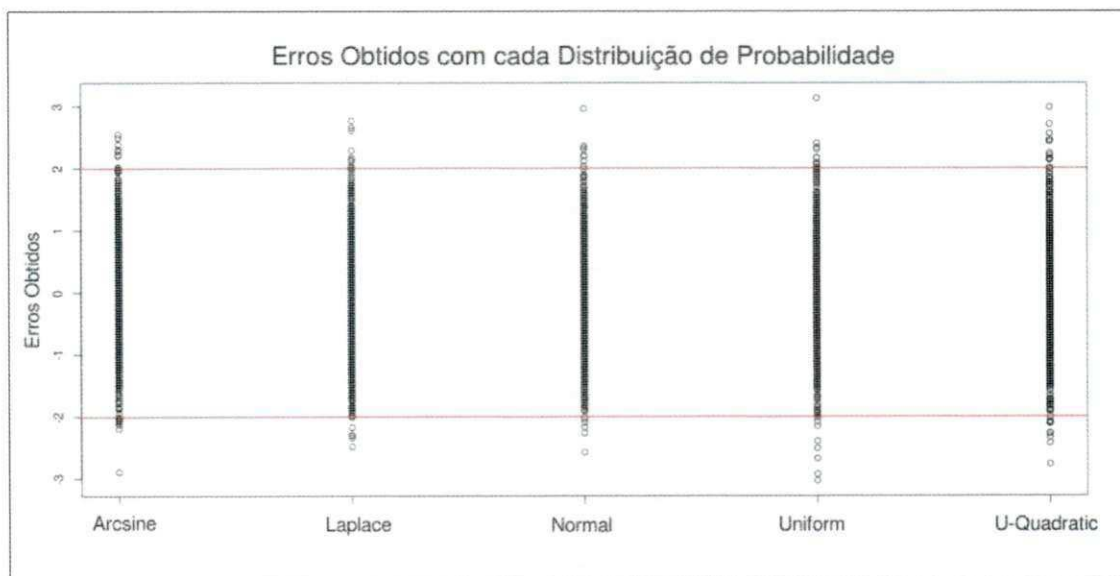


Figura 4.15: Erros obtidos (e_o) com cada distribuição de probabilidade.

Queremos identificar qual distribuição de probabilidade provê uma maior privacidade. Ofuscando o perfil da Figura 3.1 várias vezes (100 amostras) usando cada distribuição, obtemos os resultados apresentados na Tabela 4.5.

Métr. / Distr.	Correlação	SNR	MSE
Arcoseno	(0,475; 0,483)	(1,748; 1,760)	(0,0086; 0,0086)
Laplace	(0,478; 0,486)	(1,725; 1,764)	(0,0086; 0,0087)
Normal	(0,476; 0,483)	(1,740; 1,764)	(0,0085; 0,0087)
Uniforme	(0,472; 0,480)	(1,742; 1,757)	(0,0086; 0,0087)
U-quadrática	(0,477; 0,484)	(1,750; 1,758)	(0,0086; 0,0086)

Tabela 4.5: Níveis de privacidade obtidos com cada métrica para cada distribuição de probabilidade. Os níveis de significância são de 97,5%.

Como se pode observar, independente da métrica, não existe diferença estatística entre os níveis de privacidade alcançados com cada distribuição analisada. Com isso, reforçamos (não conseguimos rejeitar) a hipótese H_{0-0} de que não existe uma distribuição de probabilidade que ofusca os dados melhor do que as outras distribuições.

Capítulo 5

Validação da Abordagem

5.1 Validação da Privacidade

Além de propor uma solução para preservar a privacidade, nós também pesquisamos sobre ataques que podem afetar a solução proposta. Para validar a privacidade, foram realizados ataques e os resultados foram avaliados. Considera-se ataques de sucesso os que de alguma forma podem inferir informações a partir dos dados ofuscados.

5.1.1 Ataque do NIALM

O processo de analisar perfis com propósito de inferir quais eletrodomésticos estão sendo utilizados é conhecido como NIALM – Non-Intrusive Appliance Load Monitoring (Monitoramento não-intrusivo de carga de eletrodomésticos). Como uma forma de validar a privacidade, ataques de NIALM foram realizados e os resultados foram avaliados. Figura 5.1 apresenta o fluxo de tarefas deste procedimento de validação. Nossa abordagem de ofuscamento é aplicada nas tarefas que estão em verde, enquanto que nas tarefas que estão em laranja, nós aplicamos o algoritmo INDIC – Improved NIALM using load DIvision and Ca-libration (NIALM melhorado usando divisão de carga e calibragem) [3].

A Figura 5.2 apresenta um perfil semanal obtido do popular conjunto de dados REDD [23]. Este conjunto de dados possui perfis residenciais e perfis individuais dos eletrodomésticos usados nas residências. Existem muitos eletrodomésticos sendo usados no perfil da Figura 5.2, tais como microondas e geladeira. A Figura 5.3 apresenta o perfil real do mi-

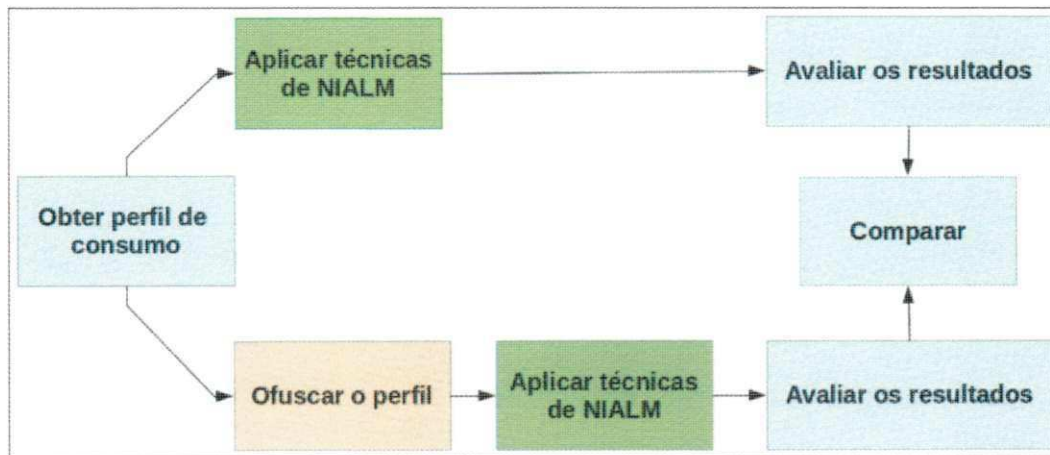


Figura 5.1: Fluxo de tarefas para validação da privacidade através de ataques de NIALM.

croondas. Um microondas foi escolhido como exemplo porque é um eletrodoméstico bem característico e um dos mais difíceis de ofuscar (pois gera muitas variações de pico).

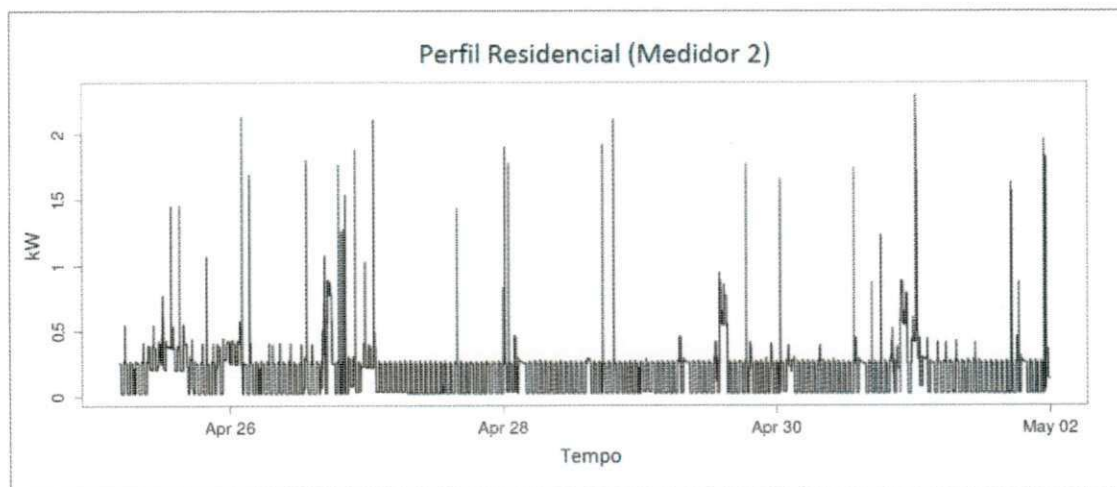


Figura 5.2: Exemplo de uma semana obtida do conjunto de dados REDD. Medições são de 1 minuto.

Após aplicar o algoritmo do INDIC para realizar o NIALM no perfil da Figura 5.2, os eletrodomésticos são desagregados. A Figura 5.4 apresenta o perfil do microondas inferido pelo INDIC. Como observado, este perfil inferido é muito parecido com o perfil real do microondas (pelo menos os momentos e níveis de picos são os mesmos).

Usando a estratégia de JSD e e_p de 5% para ofuscar o perfil da Figura 5.2, obteve-se o

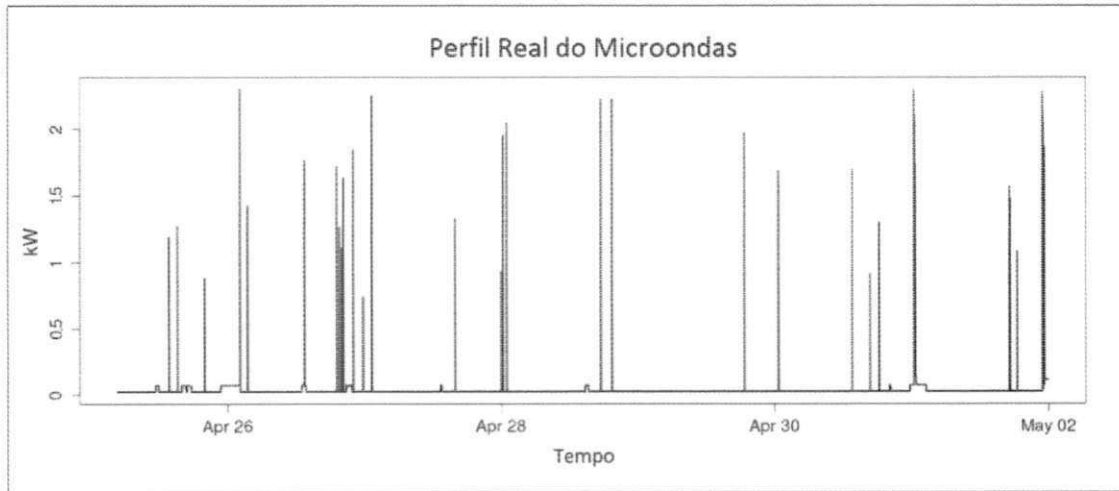


Figura 5.3: Perfil real do microondas usado no perfil da Figura 5.2.

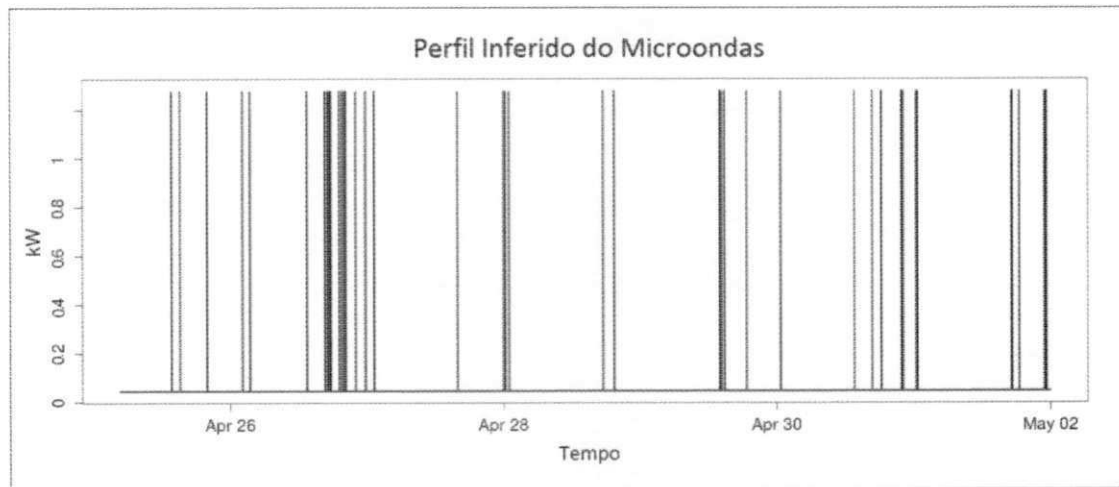


Figura 5.4: Microondas inferido pelo INDIC a partir do perfil da Figura 5.2.

perfil da Figura 5.5. Uma vez que valores negativos de demanda são impossíveis, para ser possível aplicar o NIALM, todos os valores negativos foram substituídos por zero.

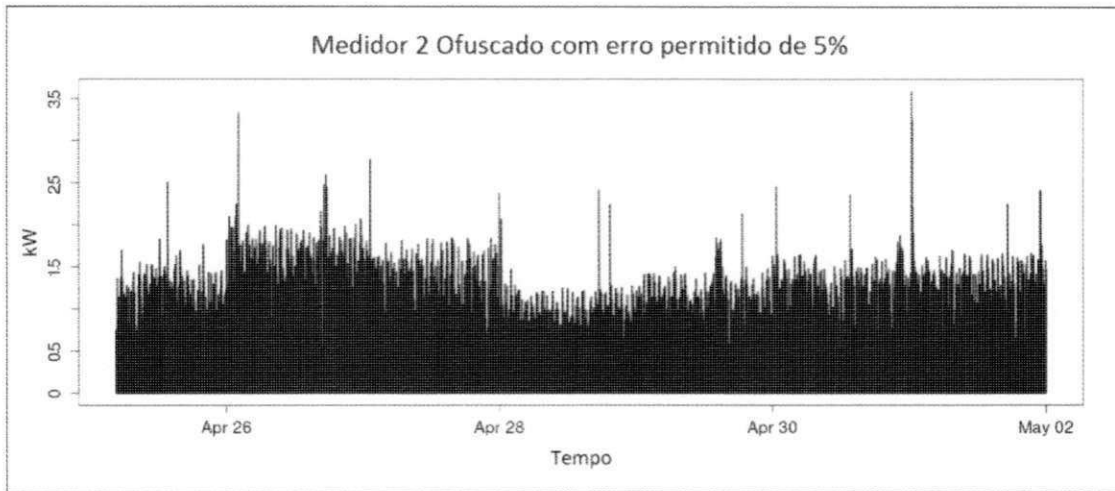


Figura 5.5: Perfil da Figura 5.2 ofuscado com JSD e e_p de 5%.

Após aplicar o mesmo algoritmo (INDIC) para realizar o NIALM, foi observado que a técnica perdeu significativamente o seu poder de detectar o uso de eletrodomésticos. A Figura 5.6 apresenta o perfil do microondas inferido pelo INDIC a partir do perfil ofuscado da Figura 5.5. Inferir um perfil para o microondas não significa que o INDIC detectou que a residência possui um microondas. Na verdade, o algoritmo irá gerar perfis para todos os eletrodomésticos que serviram de entrada para o treinamento do algoritmo de aprendizagem (mesmo que no perfil analisado não exista o uso de tais eletrodomésticos).

Para avaliar o poder das técnicas de NIALM em detectar o uso de dispositivos, Batra et al. [3] utilizam duas métricas: MNE – Mean Normalized Error (Erro Médio Normalizado) e RMS – Root Mean Square Error (Raiz da média dos erros ao quadrado). Valores menores de MNE e RMS implicam em uma melhor precisão da técnica de NIALM. Usando diferentes configurações, nós avaliamos a nossa abordagem de ofuscamento para os dois eletrodomésticos mais significantes do perfil da Figura 5.2, como apresentado na Tabela 5.1. Como observado, mesmo usando um alto nível de utilidade e um baixo ofuscamento (escolhendo um e_p pequeno), a abordagem de ofuscamento ainda diminuiu significativamente o potencial do NIALM em detectar o uso de eletrodomésticos (os valores de MNE e RMS são aumentados significativamente). Desta forma, nós consideramos que este ataque não tem

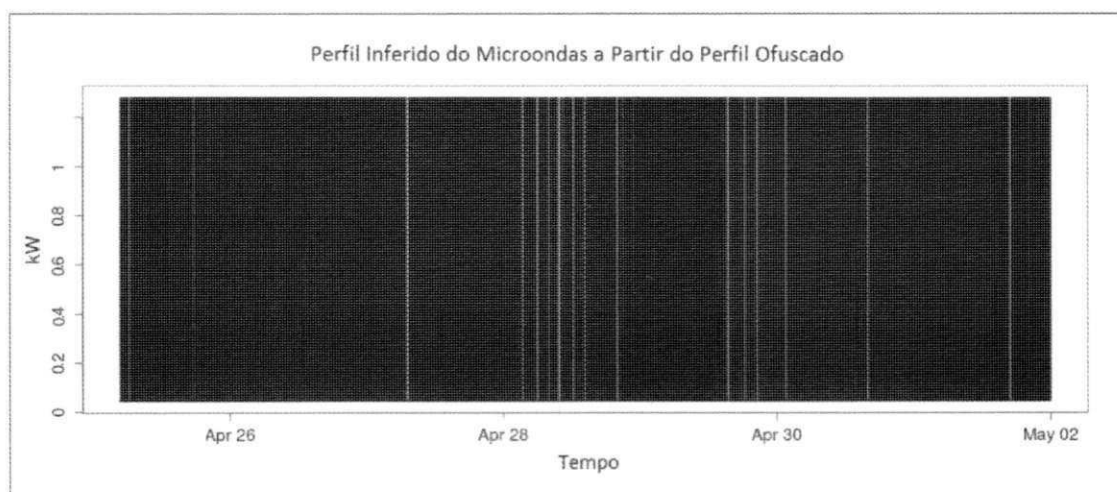


Figura 5.6: Microondas inferido pelo INDIC a partir do perfil ofuscado da Figura 5.5.

nenhum impacto na privacidade.

Eletrodoméstico	Sem ofuscamento		Ofuscamento com $e_p = 1\%$		Ofuscamento com $e_p = 2\%$		Ofuscamento com $e_p = 5\%$	
	RMS	MNE	RMS	MNE	RMS	MNE	RMS	MNE
Microondas	89,93	273,03	123,11	307,23	212,48	458,68	541,64	2119,09
Geladeira	61,02	43,19	85,64	66,51	100,31	83,44	102,87	90,35

Tabela 5.1: Valores RMS e MNE do INDIC usando diferentes configurações de ofuscamento (sem ofuscamento, ofuscamento com erro permitido de 1%, 2% e 5%).

5.1.2 Ataque do Dia da Semana

Este ataque é baseado na hipótese de que o consumidor em cada semana tende a ter um comportamento similar. O atacante pode coletar os dados e calcular uma semana esperada para este consumidor. Uma semana esperada é composta por sete dias esperados (de domingo a sábado) e um dia esperado é composto pelas médias das medições de cada horário deste dia específico (e.g., o atacante pode calcular o domingo esperado a partir de todos os domingos já coletados anteriormente). Usando a semana esperada, o atacante pode tentar inferir o comportamento do consumidor para semanas futuras. Claramente, o efeito do ataque depende da quantidade de dados disponíveis para o atacante.

Para o mesmo consumidor residencial da Figura 3.1, este ataque foi aplicado e os resultados são apresentados na Tabela 5.2. Para analisar o efeito do ataque, a correlação média entre as semanas ofuscadas e as semanas reais foi comparada com a correlação média entre a semana esperada e as semanas reais.

Número de semanas disponíveis	Correlação média entre semanas ofuscadas e semanas reais	Correlação média entre a semana esperada e semanas reais
2	(0,499; 0,510)	(-0,046; -0,033)
4	(0,335; 0,344)	(-0,005; 0,004)
8	(0,369; 0,374)	(0,154; 0,166)
16	(0,326; 0,331)	(0,171; 0,180)
32	(0,436; 0,439)	(0,182; 0,189)
52	(0,432; 0,434)	(0,316; 0,323)

Tabela 5.2: Efeito do ataque do dia da semana pra um consumidor residencial. Os intervalos de confiança são com níveis de significância de 95%.

Como pode ser observado na Tabela 5.2, quando o número de semanas disponíveis para o atacante aumenta, o efeito do ataque também aumenta, pois as correlações entre a semana esperada e as semanas reais são maiores. Porém, em nossos experimentos, mesmo escolhendo um consumidor que quase sempre repete o seu comportamento (veja Figura 3.1) e usando um longo período de observação (52 duas semanas, um ano completo), estas correlações são menores que as correlações entre as semanas ofuscadas e as semanas reais. Isto significa que para um atacante, é melhor inferir o comportamento do consumidor a partir das próprias semanas ofuscadas do que a partir da semana esperada. Desta forma, nós consideramos que este ataque não tem nenhum impacto na privacidade.

5.1.3 Ataque do Filtro

Este ataque é baseado no cálculo de uma média móvel ao longo do perfil. Abaixo segue o algoritmo de filtragem deste ataque.

- Seja T a série temporal que representa o perfil ofuscado e T_f uma nova série (o perfil filtrado).
- Os primeiros P valores de T_f serão iguais a zero.
- O valor de índice $P + 1$ de T_f será igual à média dos valores de índices 1 até $P + 1$ de T . O valor de índice $P + 2$ de T_f será igual à média dos valores de índices de 2 até $P + 2$ de T , e assim por diante. Este procedimento é uma média móvel para eliminar o ruído em alta frequência.

Para analisar a eficácia do ataque, fizemos experimentos para verificar se o coeficiente de correlação entre um perfil original e um perfil ofuscado é aumentado após a filtragem (ou seja, verificar se o nível de privacidade é diminuído). A Figura 5.7 apresenta um perfil residencial diário com medições realizadas a cada 1 minuto. Este perfil foi obtido através do Gerador de Carga, um software que simula o uso de eletrodomésticos em uma residência e gera o correspondente perfil de consumo (software criado pelo mesmo projeto que deu suporte à esta pesquisa de mestrado), conforme apresentaremos no Apêndice A.

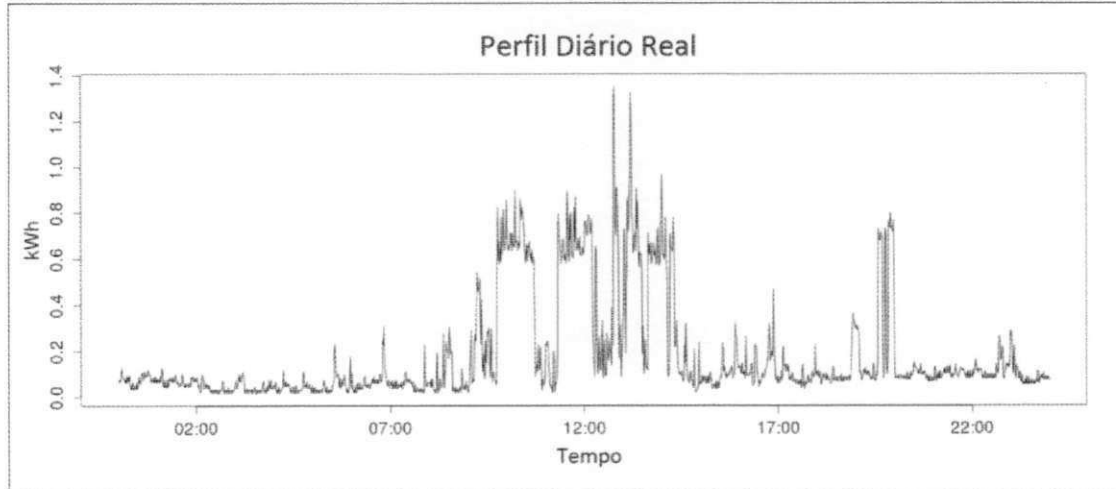


Figura 5.7: Perfil diário gerado pelo Load Generator.

Ao ofuscarmos o perfil da Figura 5.7, obtemos o perfil da Figura 5.8. Queremos agora filtrar este perfil para verificar se o ruído inserido é removido e se a correlação com o perfil original aumenta. Fizemos este procedimento utilizando diferentes valores de P , conforme apresentado nas Figuras 5.9, 5.10 e 5.11.

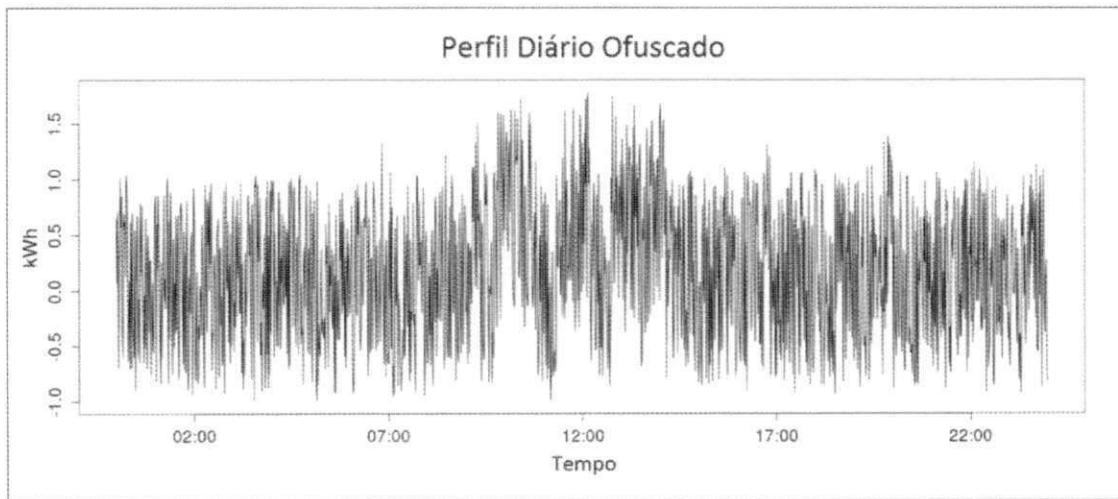
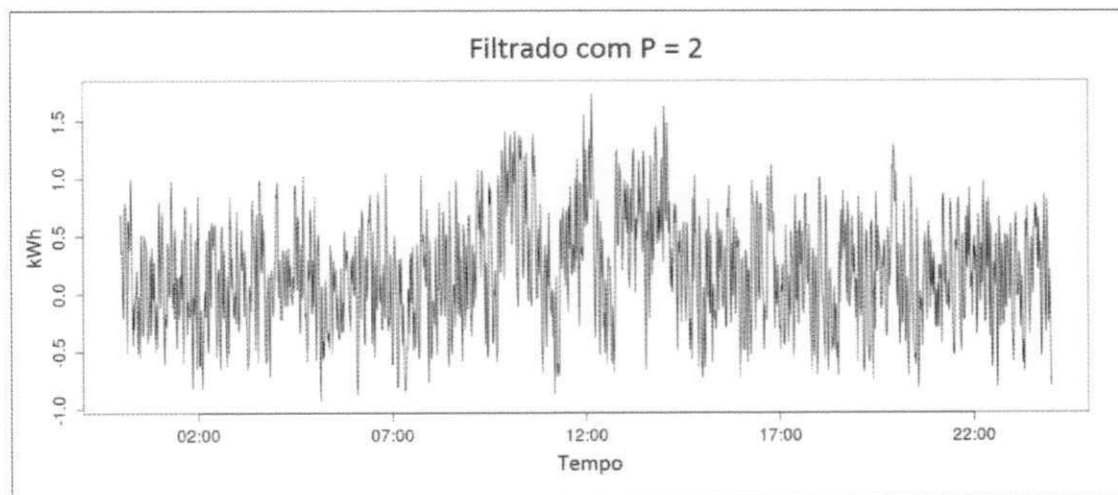


Figura 5.8: Perfil diário ofuscado.

Figura 5.9: Perfil diário ofuscado filtrado com $P = 2$.

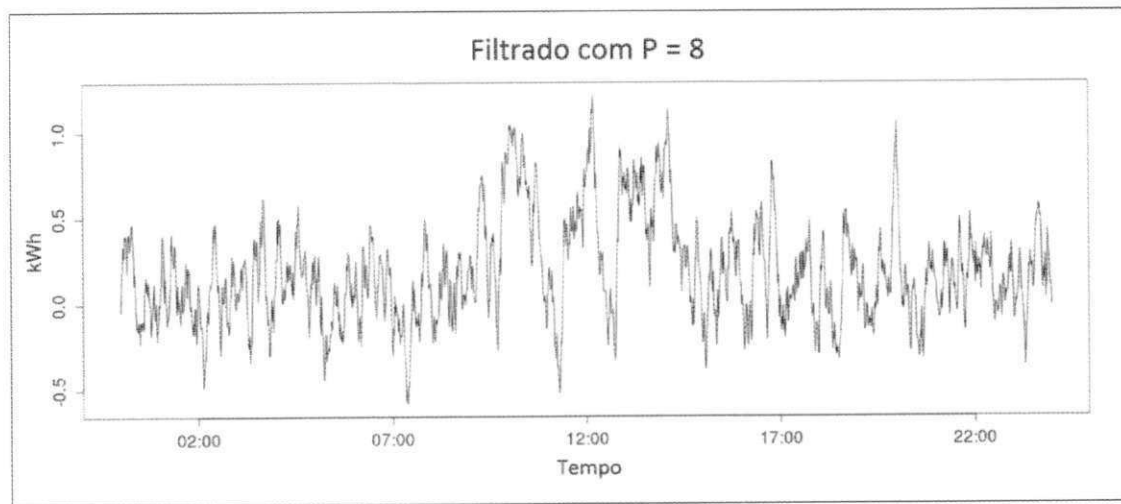


Figura 5.10: Perfil diário ofuscado filtrado com $P = 8$.

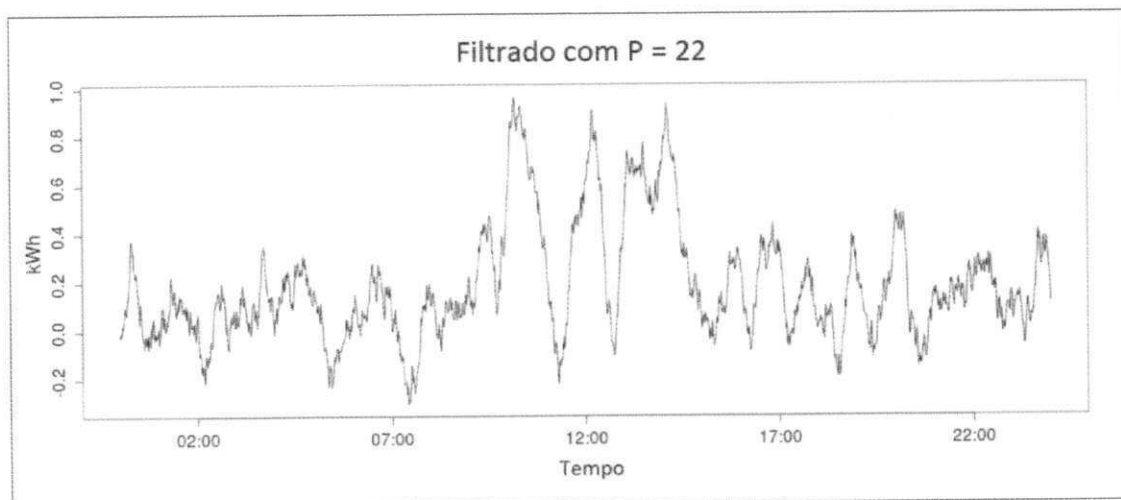


Figura 5.11: Perfil diário ofuscado filtrado com $P = 22$.

Como pode ser observado, o poder de filtragem do ataque é maior quando aumentamos o valor de P , pois o ruído em alta frequência é cada vez mais eliminado. De fato, a correlação entre o perfil ofuscado e o perfil original é de 0,350. Quando aplicamos os filtros com $P = 2$, $P = 8$ e $P = 22$ obtemos correlações de 0,450, 0,595 e 0,6144, respectivamente. Ou seja, podemos considerar que este ataque afeta a privacidade. Entretanto, ao usarmos valores de P maiores que 22, a correlação obtida é cada vez menor e isto se deve à saturação da filtragem. Para este exemplo, usando valores de P maiores que 22, o filtro eliminará não somente o ruído inserido pelo ofuscamento, mas também as próprias características do perfil original (e isto contribui para uma menor correlação). A Figura 5.12 apresenta este ataque de filtragem usando um $P = 200$ (a correlação obtida foi de 0,4619).

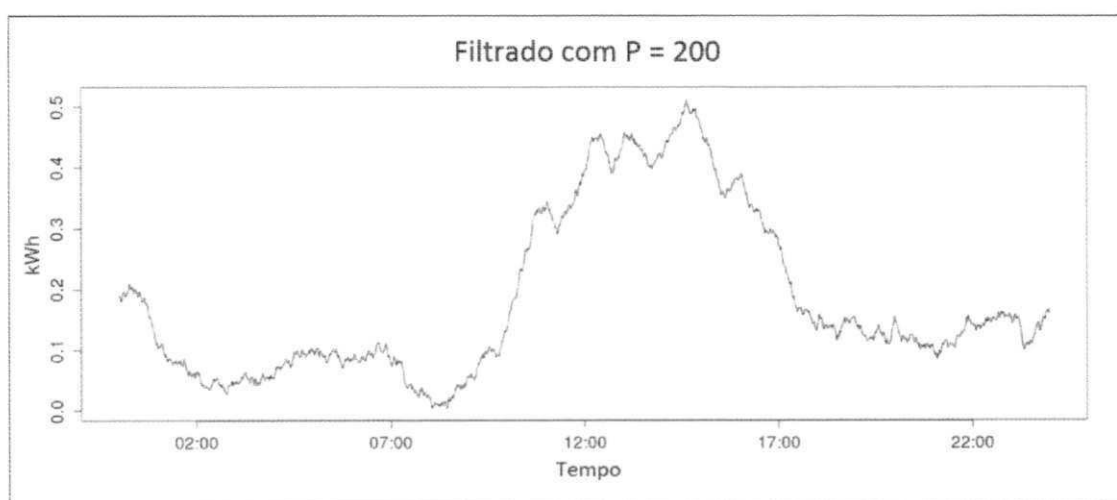


Figura 5.12: Perfil diário ofuscado filtrado com $P = 200$.

A Tabela 5.3 apresenta os resultados obtidos usando diferentes valores de P . Como se pode observar, para este cenário, $P = 22$ é o que apresentou um melhor resultado.

Utilizando ofuscamentos gerados por outras distribuições de probabilidade, também realizamos ataques de filtragem. Concluiu-se que, estatisticamente, nenhuma das distribuições fornece uma maior resistência ao ataque do filtro.

5.1.4 Outros Ataques

Nesta seção, apresentaremos ideias de outros ataques que não foram diretamente analisados neste trabalho mas que podem ser abordados em trabalhos futuros.

Valor de P	Correlação	Valor de P	Correlação
0	0,350	16	0,6132
2	0,450	18	0,6134
4	0,543	20	0,6143
6	0,578	22	0,6144
8	0,595	23	0,6141
10	0,600	24	0,6123
12	0,6057	100	0,5645
14	0,6089	200	0,4619

Tabela 5.3: Efeito do ataque do filtro para um perfil diário residencial usando diferentes valores de P .

Detecção de Picos e Ausências de Consumo

A detecção de picos e ausências de consumo em perfis certamente tem muito a dizer sobre o comportamento dos consumidores.

Com a detecção de ausências de consumo é possível inferir quando as pessoas estão dormindo ou viajando. Técnicas de detecção de ausências de consumo são simples e podem ser resumidas na detecção de consumos iguais a zero (ou próximos à zero) e sem muita variação.

Com a detecção de picos é possível inferir em quais horários existem mais pessoas na residência, quando as pessoas tomam banho ou recarregam seus veículos elétricos ou ainda quais os horários de funcionamento das fábricas. Existem várias técnicas para detectar picos em séries temporais [32]. O algoritmo abaixo apresenta o procedimento genérico.

- Seja $T = c_1, c_2, \dots, c_N$ uma série temporal uniformemente amostrada contendo N valores.
- Seja S uma dada função de pico que associa uma pontuação $S(i, c_i, T)$ ao i -ésimo elemento c_i da série temporal T . Um dado ponto c_i em T é um pico se $S(i, c_i, T) \geq h$, onde h é um valor limiar especificado pelo usuário ou calculado de maneira adequada.
- A questão importante é: como calcular a função S ?

As equações abaixo apresentam algumas formas de se computar a função S . Em todas elas, para determinar se uma medição c_i é um pico, leva-se em consideração as k medições vizinhas à esquerda e as k medições vizinhas à direita.

$$S_1(k, i, T) = \frac{\max(c_i - c_{i-1}, \dots, c_i - c_{i-k}) + \max(c_i - c_{i+1}, \dots, c_i - c_{i+k})}{2}$$

$$S_2(k, i, T) = \frac{\frac{(c_i - c_{i-1} + \dots + c_i - c_{i-k})}{k} + \frac{(c_i - c_{i+1} + \dots + c_i - c_{i+k})}{k}}{2}$$

$$S_3(k, i, T) = \frac{\left(c_i - \frac{(c_{i-1} + \dots + c_{i-k})}{k}\right) + \left(c_i - \frac{(c_{i+1} + \dots + c_{i+k})}{k}\right)}{2}$$

$$S_4(k, w, i, T) = H_w(N^-(k, i, T)) - H_w(N^+(k, i, T))$$

Em S_1, S_2, S_3 e S_4 , k significa a quantidade de medições vizinhas que serão consideradas (à esquerda e à direita), i significa o índice da medição que se quer determinar se é um pico e T é a série temporal em questão (perfil de consumo). Em S_4 , N^- é a sequência dos k vizinhos à esquerda de c_i , N^+ a sequência dos k vizinhos à direita, H a função para cálculo de entropia e w é o parâmetro usado na função kernel para estimar a função de densidade de probabilidade usada em H .

Usando a primeira técnica (S_1), detectamos os picos do perfil da Figura 5.7, conforme apresentado na Figura 5.13. Os parâmetros utilizados foram $k = 5$ e $h = 5 \cdot \sigma^2$, onde σ^2 é a variância do perfil.

Usando detecção de picos e ausências de consumo, possíveis hipóteses para verificar a eficácia dos ataques podem ser:

- H_{0-0} : Ao ofuscarmos um perfil, não é possível detectar os momentos de picos de consumo.
- H_{1-0} : Ao ofuscarmos um perfil, não é possível detectar os momentos de ausências de consumo.

Usando as estratégias apresentadas anteriormente, cabe ao estudo do ataque tentar refutar as hipóteses H_{0-0} e H_{1-0} .

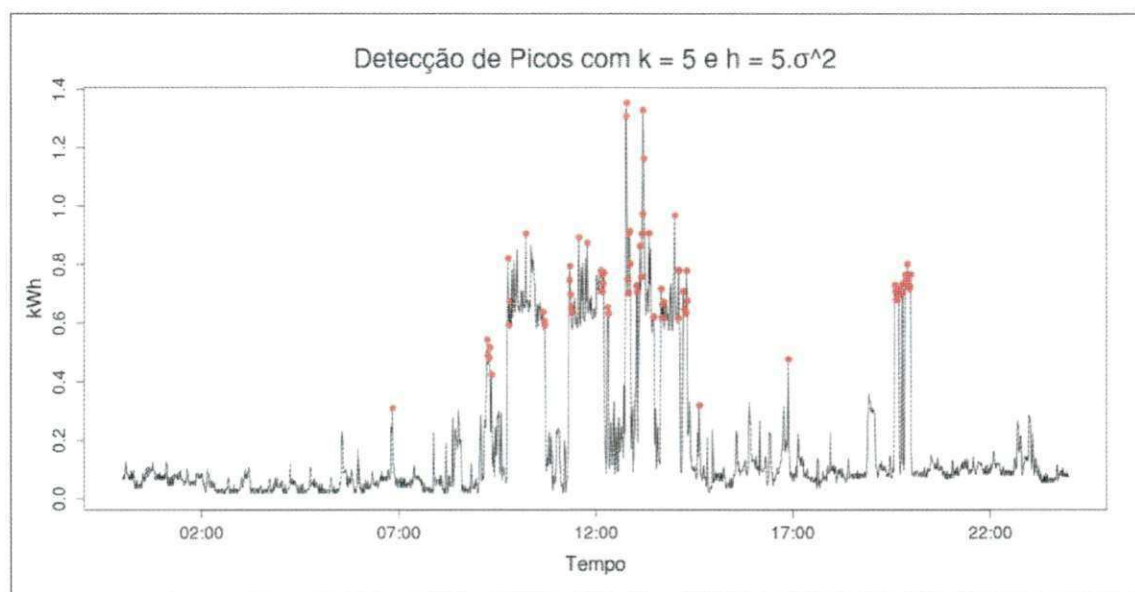


Figura 5.13: Detecção de picos no perfil da Figura 5.7 usando-se S_1 . Os pontos em vermelho são os picos encontrados.

Ataque da Mineração de Dados

Sramka [40] propõe um ataque usando mineração de dados para remover ruídos que foram inseridos em dados originais. Através de experimentos utilizando alguns conjuntos de dados que satisfazem certas condições teóricas de privacidade como ϵ -Differential Privacy (Privacidade Diferencial), mostrou-se que o ataque possui forte impacto nestes dados.

O ataque baseia-se na seguinte ideia: múltiplos algoritmos de mineração de dados são aplicados nos dados ofuscados para obter múltiplas predições dos valores reais. Estas múltiplas predições são fundidas e os valores resultantes representam estimativas dos dados originais.

Utilizando um banco de dados com idades de pessoas, Sramka [40] faz duas suposições fortes: o ruído inserido é diferente para idades distintas e fixo para idades iguais. Assim, um mesmo valor pode aparecer várias vezes no mesmo conjunto de dados. Com isso, o atacante pode fundir predições para valores iguais em uma mesma linha e em uma mesma coluna.

Para a aplicação usada em nosso trabalho, em que os dados são consumos de energia e não idades, tais suposições não são válidas, pois o ruído inserido pode ser o mesmo para valores de consumo distintos e diferente para valores de consumo iguais. Além disso, da-

dos de consumo são quase contínuos e com uma grande quantidade de possibilidades, ao passo que valores de idades são discretos. Entretanto, possíveis ataques para inferir um valor usando mineração de dados podem ser com base no comportamento do consumidor ao longo do tempo (i.e., linha da matriz) e no comportamento de todos os consumidores da região no instante de tempo desejado (i.e., coluna da matriz), pois certamente existe uma correlação forte entre estes dados.

Uma possível hipótese para verificar a eficácia de tal ataque pode ser: H_{0-0} : não existe um algoritmo de mineração que elimine o ruído e afete a privacidade dos consumidores.

5.2 Validação da Utilidade

Para validar a utilidade, algumas funcionalidades ou benefícios que usam dados de medição foram listadas e avaliadas para verificar se ainda são suportadas mesmo usando dados ofuscados. Uma vez que a abordagem proposta preserva os valores agregados (e.g., soma de linhas e soma de colunas da matriz da Figura 2.1), o procedimento para verificar se uma funcionalidade é suportada ou não pode ser mapeada como uma verificação se tal funcionalidade utiliza apenas valores agregados ou se utiliza também valores individuais. A Tabela 5.4 apresenta a lista destas funcionalidades e o impacto do ofuscamento. É importante observar que, como a solução proposta é genérica, as funcionalidades suportadas podem ser providas simultaneamente.

Caso o consumidor queira que todas as funcionalidades sejam suportadas, ele pode divulgar os seus dados sem nenhum ofuscamento. Neste caso ele não terá privacidade.

Nas próximas seções, analisaremos o impacto do ofuscamento em cada uma destas funcionalidades.

5.2.1 Otimizações de Faturamento

Com a implantação dos Smart Meters, faturamentos de consumo serão gerados com base no consumo real, e não mais com base em estimativas, como é feito algumas vezes no Brasil. Além disso, consumidores não terão mais problemas quando eles mudarem de residência ou empresa.

Funcionalidade	Suporte
Otimizações de faturamento [2]	Sim
Monitoramento e gerenciamento de carga para grupos específicos ou regiões [2]	Sim
Deteção de vazamentos e roubos de energia [1]	Sim
Previsão de carga para grupos e regiões [17]	Sim
Previsão de carga para indivíduos [17]	Não
Faturamento com política de horário [2]	Sim
Faturamento com política de níveis de demanda (e.g., diferentes preços para diferentes demandas) [35]	Não
Análise individual dos dados (e.g., NIALM e mercados oferecendo produtos adicionais à seus consumidores) [31]	Não
Ferramentas para feedbacks em casa (e.g., faturamento estimado, gerenciamento do uso de energia, perfis de eletrodomésticos, etc.) [5]	Sim

Tabela 5.4: Funcionalidades que utilizam dados de medições e o impacto do ofuscamento.

Conforme apresentamos nas Seções 3.1 e 3.2, nossa abordagem dá suporte às otimizações de faturamento.

5.2.2 Monitoramento e Gerenciamento de Carga

Uma curva de carga é um gráfico ilustrando a variação da demanda/consumo de carga/energia ao longo do tempo. Concessionárias de energia usam estas informações para planejar quanto de energia devem gerar em um instante de tempo. Uma vez que energia elétrica é uma forma de energia que não pode ser armazenada de maneira eficiente, é necessário gerar, distribuir e consumir imediatamente. Quando a carga em um sistema alcança a capacidade máxima de geração, operadores da rede precisam encontrar meios adicionais para suprir a demanda ou racionar. Se estes processos forem malsucedidos, o sistema irá se tornar instável e apagões podem ocorrer.

Gerenciamento de carga é uma estratégia que as concessionárias de energia podem usar para reduzir a demanda durante ocasiões de pico, tais como em dias quentes de verão. Na

prática, isto significa aumentar os preços de energia nos horários de pico para forçar os consumidores a consumirem menos.

Conforme apresentamos na Seção 3.4, nossa abordagem dá suporte ao monitoramento e gerenciamento de carga.

5.2.3 Detecção de Vazamentos e Roubos de Energia

Perdas podem ser classificadas como técnicas (e.g., devido à resistência nas linhas de transmissão), e não técnicas (e.g., roubos, principalmente). Roubo de energia resultam em altas taxas para os consumidores legítimos. Medidores distribuídos (diferente dos instalados em residências ou empresas) localizados em pontos estratégicos podem identificar a quantidade de eletricidade transmitida para áreas específicas. Combinados com os Smart Meters dos consumidores, concessionárias podem detectar roubos de energia com mais precisão e a localização pode ser identificada rapidamente [1].

Perdas podem ser calculadas através da subtração da quantidade de energia fornecida pela quantidade de energia cobrada (informada pelos Smart Meters dos consumidores). Se quisermos calcular as perdas não técnicas, uma estratégia simples seria:

$$EP = EF - EC \quad (5.1)$$

$$EP = PNT + PT \quad (5.2)$$

Onde EP é o total de energia perdida, EF é o total de energia fornecida, EC é o total de energia cobrada, PNT é o total devido a perdas não técnicas e PT é o total devido a perdas técnicas. Para todos estes valores, o período de análise pode ser o período completo de faturamento (e.g., um mês). Combinando as Equações 5.1 e 5.2, temos:

$$PNT = EF - EC - PT$$

Nossa abordagem dá suporte a esta funcionalidade, pois mesmo com o ofuscamento, o valor total de EC ainda é preservado.

5.2.4 Previsão de Carga para Grupos e Regiões

Previsão de carga é vitalmente importante para a indústria elétrica na sociedade atual. As informações providas pelas previsões de carga podem ajudar as companhias a tomar decisões como comprar e gerar mais energia, mudar demandas, avaliar melhor os contratos e projetar infraestrutura. Para previsão de carga de certos grupos ou regiões, algoritmos de aprendizagem de máquina usam medições passadas como conjunto de treinamento [17]. Como com a nossa abordagem é possível obter um histórico preciso dos valores de consumo de grupos e regiões, esta funcionalidade é suportada.

5.2.5 Previsão de Carga para Indivíduos

Assim como no caso anterior, em previsão de carga para indivíduos, algoritmos de aprendizagem de máquina usam medições passadas como conjunto de treinamento [17]. Entretanto, como com a nossa abordagem não é possível obter um histórico preciso dos valores de consumo de indivíduos, esta funcionalidade não é suportada.

5.2.6 Faturamento com Política de Horário

Políticas de faturamento podem ser de diferentes formas e podem fornecer diferentes preços de energia de acordo com o horário do dia, dia da semana e mês do ano. Sendo assim, as concessionárias de energia podem usar preços mais altos para momentos de muita demanda e mais baixos para momentos de pouca demanda.

Conforme apresentamos na Seção 3.3, nossa abordagem dá suporte ao faturamento com política de horário.

5.2.7 Faturamento com Política de Níveis de Demanda

Concessionárias de energia investem em equipamentos de geração e distribuição de energia para poder prover a demanda que os consumidores requerem em um instante de tempo. Utilizar taxas separadas para consumo e demanda é mais justo devido à distribuição dos custos em prover tais serviços (mais justa principalmente para consumidores que demandam pouca eletricidade). Desta forma, algumas empresas utilizam taxas diferentes para os consumidores

que excedem a demanda contratada.

O faturamento para a demanda contratada excedida é calculada como [35]:

$$FDE = TE \cdot (DM - DC)$$

onde FDE é o valor a ser pago pelo faturamento da demanda excedida, TE é a tarifa de excedência, DM é a demanda medida e DC é a demanda contratada.

Nossa abordagem não suporta faturamento com política de níveis de demanda, pois o perfil divulgado para a concessionária é verticalmente ofuscado e isto dificulta a identificação dos momentos de excedência da demanda contratada.

5.2.8 Análise Individual dos Dados

Como apresentado anteriormente, análise dos dados de medições individuais pode revelar informações detalhadas a respeito dos consumidores e seus comportamentos. Possíveis usos dessas informações podem ser [31]:

- Fabricantes de eletrodomésticos podem usar estas informações para vender os seus produtos e suas garantias;
- Detecção de padrões de comportamento específicos (por exemplo, seguros saúde podem detectar consumidores com sono irregular, que podem indicar problemas de saúde);
- Sistemas de recomendação podem gerar perfis para anúncios direcionados a produtos ou atividades;
- Forças da lei podem identificar atividades suspeitas ou ilegais (e.g., plantações de maconha); investigações; vigilância em tempo real para determinar se os residentes estão presentes e atividades atuais dentro de casa;
- Senhores de terras podem identificar se os inquilinos cumprem as locações;
- A imprensa pode investigar e publicar as atividades dos famosos;

- Criminosos e outros usuários não autorizados podem identificar os melhores momentos para furtos; determinar se os residentes estão presentes; espionagem corporativa; determinar os processos ou dados confidenciais de propriedades;

Com a nossa abordagem, análise individual dos dados não é suportada porque medições individuais são ofuscadas.

5.2.9 Ferramentas para Feedbacks em Casa

Pesquisas mostram que feedbacks em casa podem ajudar os consumidores a conservarem acima de 15% de eletricidade [5]. Por exemplo, uma ferramenta de feedbacks pode fornecer ao consumidor informações sobre o seu consumo atual e o preço da energia que a concessionária está cobrando atualmente. Como isso, o consumidor pode tomar decisões como: escolher usar a máquina de lavar em um horário em que o preço esteja mais baixo ou desligar eletrodomésticos que não estejam sendo utilizados em horários em que o preço esteja alto.

Uma vez que a abordagem de ofuscamento considera apenas os dados que são divulgados para a concessionária, isto não afeta a possibilidade do consumidor analisar o seu próprio perfil real de consumo dentro de sua casa. Portanto esta funcionalidade é suportada por nossa abordagem.

5.2.10 Outras Funcionalidades

Outras funcionalidades ou benefícios que não estão relacionadas com medições de consumo, mas que são providas devido à implantação dos Smart Meters, são listadas na Tabela 5.5.

Funcionalidade	Suporte
Conexão/Desconexão remota	Sim
Fluxo bidirecional (possibilidade de consumidores venderem a energia elétrica que produzem)	Sim
Notificações de falhas	Sim
Diagnósticos de qualidade da energia	Sim
Reduções de custo nas leituras	Sim

Tabela 5.5: Outras funcionalidades/benefícios que não estão relacionadas com medições de consumo.

Capítulo 6

Considerações Finais

6.1 Conclusões e Contribuições

Discutimos várias questões relacionadas à privacidade e utilidade dos dados em uma infraestrutura de Smart Meters e propomos uma abordagem simples e barata para preservar a privacidade dos consumidores sem afetar significativamente a utilidade dos dados para as concessionárias de energia. A modificação no procedimento de comunicação entre um Smart Meter e a concessionária de energia é apenas a geração de um número aleatório e a adição deste número com a medição a ser enviada para a concessionária. Usando exemplos de consumidores em aplicações reais de Smart Grids, a abordagem mostrou-se promissora.

Avaliamos possíveis otimizações na abordagem proposta e concluímos que a estratégia de Janela Saltitante Diária (*JSD*) é a melhor entre as avaliadas (*JSM*, *JDM* e *JSD*) para se calcular o erro permitido (e_p), a correlação é a melhor métrica de privacidade e que não existem diferenças estatísticas (com relação à privacidade e utilidade) entre as distribuições de probabilidade utilizadas para realizar o ofuscamento.

Validamos a privacidade através de possíveis ataques e concluímos que o ataque do NI-ALM e o ataque do dia da semana são ineficazes, enquanto que o ataque do filtro possui uma certa eficácia. Validamos a utilidade através da análise de várias funcionalidades que podem ser alcançadas com o uso de Smart Meters e concluímos que, mesmo com o uso de dados ofuscados, as funcionalidades mais importantes para as concessionárias ainda são suportadas.

6.2 Trabalhos Futuros

Em trabalhos futuros, possíveis otimizações da abordagem e novos ataques podem ser explorados.

Utilizamos conceitos genéricos de privacidade de dados e adição de ruído instanciados para o contexto de Smart Grids. Entretanto, as mesmas ideias podem ser exploradas para outros contextos, como por exemplo, aplicações de finanças, saúde e serviços de localização.

Privacidade através de ofuscamento de dados tem se tornado um método bastante popular em mineração de dados. Na literatura existem várias abordagens e com diferentes definições de privacidade [14]. Um trabalho futuro pode ser comparar a nossa abordagem com as abordagens existentes para mineração de dados e classificar a nossa abordagem de acordo com outras definições de privacidade.

Bibliografía

- [1] M. Anas, N. Javaid, A. Mahmood, S. M. Raza, U. Qasim, and Z. A. Khan. Minimizing electricity theft using smart meters in ami. In *Proc. of the IEEE Seventh International Conference on P2P, Parallel, Grid, Cloud and Internet Computing*, Victoria, Canada, August 2012.
- [2] P. Barbosa, A. Brito, H. Almeida, and S. Clauß. Lightweight privacy for smart metering data by adding noise. In *ACM 29th Symposium on Applied Computing*, Gyeongju, South Korea, March 2014.
- [3] N. Batra, H. Dutta, and A. Singh. Indic: Improved non-intrusive load monitoring using load division and calibration. In *IEEE 12th International Conference on Machine Learning and Applications*, MIAMI, USA, December 2013.
- [4] T. Baumeister. Literature review on smart grid cyber security. Master's thesis, Collaborative Software Development Laboratory at the University of Hawaii, 2010.
- [5] BCHydro. Smart metering and infrastructure program business case. <https://www.bchydro.com/content/dam/BCHydro/customer-portal/documents/projects/smart-metering/smi-program-business-case.pdf>, 2011.
- [6] C. Boccuzzi. Smart grid e o big brother energético. *Metering International América Latina*, 3:82–83, 2010.
- [7] J. Bohli, C. Sorge, and O. Ugus. A privacy model for smart metering. In *Proc. IEEE International Conference Communications Workshops (ICC)*, pages 1–5, Cape Town, South Africa, May 2010.
- [8] Green Button. Green button data. <http://www.greenbuttondata.org>, 2014.

- [9] Agência Nacional de Energia Elétrica do Brasil (ANEEL). Tarifa branca. <http://www.aneel.gov.br/area.cfm?idArea=781&idPerfil=4>, 2011.
- [10] C. Efthymiou and G. Kalogridis. Smart grid privacy via anonymization of smart metering data. In *IEEE 1st International Conference on Smart Grid Communications*, pages 238–243, Gaithersburg, USA, October 2010.
- [11] K. E. Emam. *Guide to the de-identification of personal health information*. CRC Press, 2013.
- [12] EnerNOC. Boston cleanweb hackathon and challenge 2012. http://boston.cleanwebhack.com/wp/?page_id=11, 2012.
- [13] Commission for Energy Regulation (CER). Cer smart metering project. <http://www.ucd.ie/issda/data/commissionforenergyregulation>, 2012.
- [14] B. C. M. Fung, K. Wang, A. W-C. Fu, and P. S. Yu. *Introduction to Privacy-Preserving Data Publishing: Concepts and Techniques*. Chapman and Hall/CRC, 2010.
- [15] V. Giordano, I. Onyeji, G. Fulli, M. S. Jimnez, and C. Filiou. Guidelines for cost benefit analysis of smart metering deployment. *JRC Scientific and Tech. Research*, 2012.
- [16] IBGE. Pesquisa nacional por amostra de domicílios. <http://www.ibge.gov.br/home/presidencia/noticias/imprensa/ppts/00000010135709212012572220530659.pdf>, 2012.
- [17] D. Ilić, P. G. Silva, S. Karnouskos, and M. Jacobi. Impact assessment of smart meter grouping on the accuracy of forecasting algorithms. In *Proc. of the 28th Annual ACM Symposium on Applied Computing*, pages 673–679, Coimbra, Portugal, March 2013.
- [18] INMETRO. Portaria número 375, de 27 de setembro de 2011. <http://www.inmetro.gov.br/legislacao/rtac/pdf/RTAC001738.pdf>, 2011.
- [19] J. S. John. Big data on the smart grid: 2013 in review and 2014 outlook. <http://www.greentechmedia.com/articles/read/Big-Datas-5-Big-Steps-to-Smart-Grid-Growth-in-2014>, 2013.

- [20] G. Kalogridis, C. Efthymiou, S. Z. Denic, T. A. Lewis, and R. Cepeda. Privacy for smart meters: towards undetectable appliance load signatures. In *IEEE 1st International Conference on Smart Grid Communications*, pages 232–237, Gaithersburg, USA, October 2010.
- [21] J. Kelly and W. Knottenbelt. Disaggregating smart meter readings using device signatures. Master's thesis, Imperial Computing Science, London, UK, September 2011.
- [22] O. Koehle. *Just say no to big brother's smart meters. The latest in bio-hazard technology*. ARC Reproductions, 2012.
- [23] J. Z. Kolter and M. J. Johnson. Redd: A public data set for energy disaggregation research. In *Proc. of the ACM workshop on Data Mining Applications in Sustainability*, pages 1–6, USA, August 2011.
- [24] S. Kotz, J. T. Kozubowski, and K. Podgórski. *The Laplace distribution and generalizations*. Birkhauser, 2001.
- [25] L. Kuipers and H. Niederreiter. *Uniform Distribution of Sequences*. Dover, 2006.
- [26] K. Lauter, M. Naehrig, and V. Vaikuntanathan. Can homomorphic encryption be practical? In *Proc. of 3rd ACM workshop on Cloud computing security*, pages 113–124, Illinois, USA, October 2011.
- [27] F. Li, B. Luo, and P. Liu. Secure information aggregation for smart grids using homomorphic encryption. In *IEEE 1st International Conference on Smart Grid Communications*, pages 327–332, Gaithersburg, USA, October 2010.
- [28] M. Matsumoto and T. Nishimura. Mersenne twister: A 623-dimensionally equidistributed uniform pseudorandom number generator. *ACM Transactions on Modeling and Computer Simulations: Special Issue on Uniform Random Number Generation*, pages 3–30, 1998.
- [29] S. McLaughlin, P. McDaniel, and W. Aiello. Protecting consumer privacy from electric load monitoring. In *Proc. of the 18th ACM Conference on Computer and Communications Security*, pages 87–98, Illinois, USA, October 2011.

- [30] K. Mivule. Utilizing noise addition for data privacy an overview. In *International Conference on Information and Knowledge Engineering*, pages 65–71, Las Vegas, USA, July 2012.
- [31] NIST. Guidelines for smart grid cybersecurity: Vol. 2 - privacy and the smart grid. http://csrc.nist.gov/publications/nistir/ir7628/nistir-7628_vol2.pdf, 2010.
- [32] G.K. Palshikar. Simple algorithms for peak detection in time-series. In *Proc. 1st Int. Conf. Advanced Data Analysis, Business Analytics and Intelligence*, Ahmedabad, India, June 2009.
- [33] J. K. Patel and C. B. Read. *Handbook of the Normal Distribution*. Statistics: A Series of Textbooks and Monographs, 1996.
- [34] PlotWatt. Plotwatt. <https://plotwatt.com>, 2014.
- [35] Procel. Manual de tarifação da energia elétrica. programa nacional de conservação de energia. <http://www.eletrobras.com/elb/services/DocumentManagement/FileDownload.EZTSvc.asp?DocumentID=%7B35211210-AF30-4A65-B798-B2FC47107AC8%7D&ServiceInstUID=%7BAEBE43DA-69AD-4278-B9FC-41031DD07B52%7D>, 2014.
- [36] S. R. Rajagopalan, L. Sankar, S. Mohr, and H. V. Poor. Smart meter privacy: a utility-privacy tradeoff framework. In *IEEE 2nd International Conference on Smart Grid Communications*, pages 150–155, Brussels, BE, October 2011.
- [37] A. Reinhardt, P. Baumann, D. Burgstahler, M. Hollick, H. Chonov, M. Werner, and R. Steinmetz. On the accuracy of appliance identification based on distributed load metering data. In *Proceedings of the 2nd IFIP Conference on Sustainable Internet and ICT for Sustainability*, pages 1–9, 2012.
- [38] A. Rial and G. Danezis. Privacy-preserving smart metering. In *Proc. of the 18th ACM Conference on Computer and Communications Security*, pages 49–60, Illinois, USA, October 2011.
- [39] K. M. Schmidt and A. Zhigljavsky. A characterization of the arcsine distribution. *Statistics and Probability Letters*, 79:2451–2455, December 2009.

-
- [40] M. Sramka. A privacy attack that removes the majority of the noise from perturbed data. In *Proc. of the IEEE World Congress on Computational Intelligence*, pages 1–8, Barcelona, Spain, July 2010.
- [41] L. M. Surhone, M. T. Timpledon, and S. F. Marseken. *U-Quadratic Distribution*. VDM publisher, 2010.
- [42] S. Wang, L. Cui, J. Que, D.-H. Choi, X. Jiang, and L. Xie. A randomized response model for privacy preserving smart metering. *IEEE Transactions on Smart Grid*, pages 1317–1324, September 2012.
- [43] C. Wohlin, P. Runeson, M. Höst, M. C. Ohlsson, B. Regnell, and A. Wesslén. *Experimentation in Software Engineering*. Kluwer Academic Publishers, 2000.

Apêndice A

Gerador de Carga

Um problema para a realização de pesquisas em Smart Metering (e.g., pesquisas de privacidade, utilidade e NIALM) é a baixa disponibilidade de dados de medições. Devido a isto, nós promovemos o Gerador de Carga, um software desenvolvido pelo mesmo projeto que deu suporte a esta pesquisa de mestrado, na Universidade Federal de Campina Grande. O software tem como objetivo gerar perfis residenciais de consumo (sintéticos, porém realísticos) para que possam ser utilizados por outros pesquisadores ou por outras ferramentas.

Para gerar os perfis residenciais, perfis individuais de eletrodomésticos são agregados, simulando assim o comportamento dos habitantes. Tais perfis de eletrodomésticos foram obtidos do conjunto de dados TraceBase [37]. O TraceBase possui perfis de diversos eletrodomésticos com medições a cada segundo. Os eletrodomésticos usados em cada perfil são escolhidos aleatoriamente e baseados na probabilidade de uma residência ter ou não um dispositivo. Algumas probabilidades foram coletadas do IBGE (PNAD 2011) [16].

A Figura A.1 apresenta um exemplo de funcionamento do Gerador de Carga para gerar um perfil. As seguintes opções são disponíveis: gerar perfil, escolher base de dados de eletrodomésticos e diretório para salvar os resultados, configurar data inicial e final do perfil, configurar intervalo entre medições e o número de consumidores (perfis a serem gerados).

A Figura A.2 apresenta um perfil diário gerado pelo nosso software. Como a saída do software é o perfil agregado e os vários perfis dos eletrodomésticos individuais, é possível visualizar o uso de cada um deles.

Uma característica importante do nosso gerador de carga, é a possibilidade de gerar dados em diferentes formatos. Por exemplo, o software pode gerar dados de medições em tempo


```

-----X-----
LoadGenerator
Developed by Joëffison Silvério de Andrade
Laboratório de Sistemas Distribuídos
-----X-----
[Main Menu]
Choose an option:
-----X-----
1- Generate Load
2- Set Dataset's Path
3- Set Output's Path
4- Set Default Begin
5- Set Default End
6- Set Default Seconds Between Measurements
7- Set Default Number of Consumers
8- Set Default Profile Type
-----X-----
Your option is: 1
-----X-----
[Generate Load]
To Generate Load requires some parameters (the order's not important).
Missing parameters will be replaced by the default values.
-----X-----
Legend:
-b The initial date of the load;
-e The final date of the load;
-s The number of seconds between the measurements;
-n The number of consumers;
-t The type of the consumer profiles;
-f The path to the output.
-----X-----
Example:
-b 05/10/1992 10:30:00 -e 10/10/1992 10:30:00 -s 60 -n 1 -t Default -f /home/user/path/to/output
-----X-----
Your option is: -b 05/10/1992 10:30:00 -e 10/10/1992 10:30:00 -s 60 -n 1 -t Default -f /home/lg_output

```

Figura A.1: Exemplo de funcionamento do Gerador de Carga.

real e disponibilizar para outros componentes que funcionam na nuvem ou gerar dados no formato Green Button [8] para prover compatibilidade com outros aplicativos da web ou de dispositivos móveis.

O PlotWatt [34] é um exemplo de aplicativo que funciona na web e que aceita dados de medições em tempo real ou no formato Green Button. Algumas funcionalidades do PlotWatt são: monitoramento de eletrodomésticos, recomendações para economia, alertas de atividades, detecção de picos e otimizações de taxa. A Figura A.3 apresenta um perfil gerado pelo nosso Gerador de Carga e que foi submetido ao PlotWatt. Estes dados foram gerados no formato Green Button (medições a cada 15 minutos).

A Figura A.4 mostra um exemplo do uso de eletrodomésticos visualizado no aplicativo PlotWatt. Este é o exemplo padrão disponibilizado pelo próprio PlotWatt.

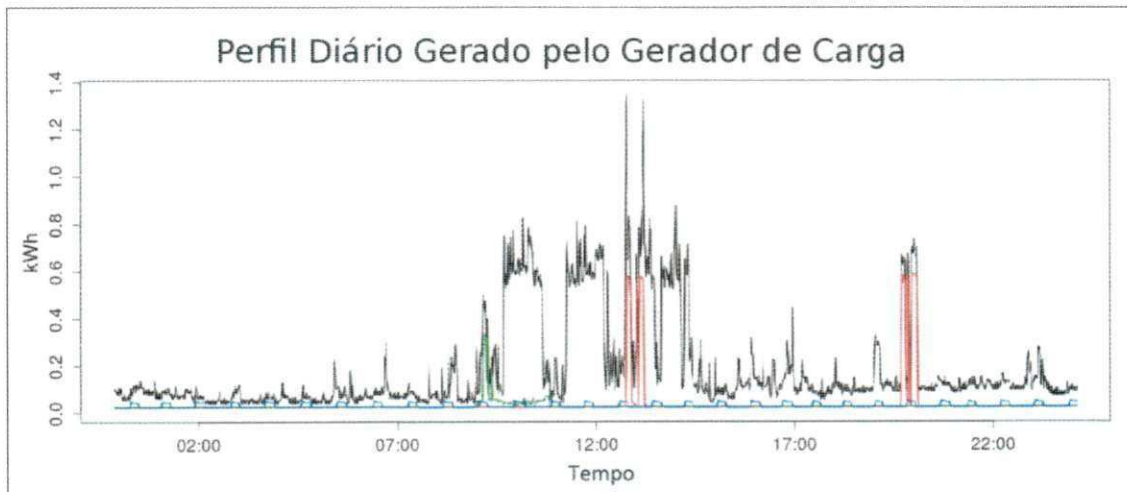


Figura A.2: Exemplo de perfil gerado pelo Gerador de Carga. Em preto está representado o perfil agregado, em verde a máquina de lavar, em vermelho a máquina de lavar louça e em azul a geladeira.



Figura A.3: Perfil de consumo representado no PlotWatt [31].



Figura A.4: Análises de consumo e identificação de eletrodomésticos pelo PlotWatt [31].

Apêndice B

Artigos Aceitos para Publicação

Lightweight Privacy for Smart Metering Data by Adding Noise

Pedro Barbosa

Federal University of Campina Grande - UFCG, Brazil
pedroyossis@copin.ufcg.edu.br

Andrey Brito

Federal University of Campina Grande - UFCG, Brazil
andrey@dsc.ufcg.edu.br

Hyggo Almeida

Federal University of Campina Grande - UFCG, Brazil
hyggo@dsc.ufcg.edu.br

Sebastian Clauß

Dresden University of Technology - TUD, Germany
sebastian.clauss@tu-dresden.de

ABSTRACT

With a Smart Metering infrastructure, there are many motivations for power providers to collect high-resolution data of electricity usage from consumers. However, this collection implies very detailed information about the energy consumption of consumers being monitored. Consequently, a serious issue needs to be addressed: how to preserve the privacy of consumers but making the provision of certain services still possible? Clearly, this is a tradeoff between privacy and utility. There are approaches for preserving privacy in various ways, but many of them affect the data usefulness or are computationally expensive. In this paper, we propose and evaluate a lightweight approach for privacy and utility based on the addition of noise. Furthermore, using real consumers' data, we discuss the influence of the technique in various Smart Grid scenarios. Finally, we also design and evaluate possible attacks to our solution.

Categories and Subject Descriptors

K.4.1 [Computers and Society]: Public Policy Issues – Privacy.
G.3 [Probability and Statistics]: Correlation and regression analysis, Random number generation.

Keywords

Smart metering, smart grid systems, data masking, noise addition.

1. INTRODUCTION

Smart Meters may be the entering gate of the fully deployment of a Smart Grid. These devices can transmit information to the power provider for load monitoring and billing purposes, providing accurate readings automatically at requested time intervals. The expected frequency of such readings is yet to be defined; it has been speculated that this could be as high as every few (1-5) minutes, raising important privacy issues

regarding the availability and processing of such data [7].

With Smart Meters, detailed information about the energy usage by consumers can be monitored, changing the commercial relations between consumers and power providers. On the one hand, it creates opportunities for new services by power providers, on the other, it raises concerns about the consumer's privacy [2]. Rajagopalan *et al.* [21] describe this problem as a tradeoff between utility and privacy.

The solutions offered thus far have been tied to specific technologies (e.g., using batteries [12, 18, 25]), affect the provision of some services by a power provider (e.g., with anonymization [7] it is not possible to provide differentiated tariffs), or are very expensive (e.g., using homomorphic schemes [6, 9, 16, 22]). Moreover, existing solutions have not quantified the loss of benefit (utility) that results from any such privacy-preserving approach [21], and have also not exemplified the use of the solutions in real Smart Grid scenarios.

In this paper, we propose a privacy solution that meets the needs of consumers (privacy) and still have a minor effect on the data usefulness (utility) to many services provided by power companies. After validation with real Smart Grid scenarios and real consumers, we believe the approach is promising.

1.1 Problem Statement

From data sent by Smart Meters, it is possible to infer behavior of a consumer unity. Examples of such behavior are the following: used appliances; if the house is empty at a certain time; when the inhabitants wake up, take a shower, or shut down the television; or even if the refrigerator and washer machine are not operating at a desired level of efficiency anymore.

The process of analyzing consumption profiles in order to deduce which appliances are being used is known as Non-intrusive Appliance Load Monitoring (NIALM). There are many NIALM algorithms proposed in literature [18, 25]. Kelly *et al.* [13] design, implement and evaluate some methodologies to identify the use of appliances through load profiles. If the algorithm is running remotely, the homeowners may not know that their behavior is being monitored and recorded. Figure 1 shows the identification of appliance signatures from a residential consumer [13].

Even if such behavior information is not in principle useful for a power provider, there is a large commercial market eager for

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from Permissions@acm.org.
SAC 2014, March 24 - 28 2014, Gyeongju, Republic of Korea
Copyright 2014 ACM 978-1-4503-2469-4/14/03...\$15.00.
<http://dx.doi.org/10.1145/2554850.2554982>.

such information about consumer habits. Certainly, this type of information is of interest for many businesses that want to identify the profile of a potential consumer of their products and services.

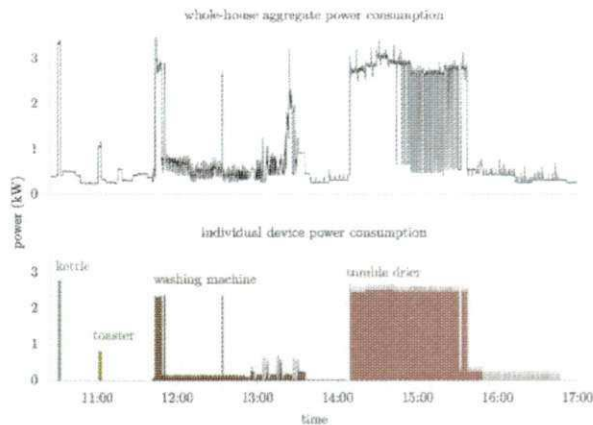


Figure 1. Profile aligned with appliance signatures [13].

The exposure of such profiles and habits of grid users evokes, therefore, issues about privacy. Clear rules are needed to protect consumers from misuse of their behavioral data and to avoid that Smart Grid becoming a new type of Big Brother [2]. Unfortunately, regulatory protection may take decades to be in place while Smart Meters are already operational.

Regarding the usefulness of the data for power providers, some motivations to collect high-resolution data are the following: identification of non-technical losses (e.g., power thefts), optimizations in load forecasting [11], billing and pricing services and monitoring of energy quality scores. Further, such data may be used to facilitate and improve network management, reduce peak load, calibrate the load distribution, and many other uses [21]. Therefore, collection and dissemination of energy information are critical to the Smart Grid. However, regarding privacy, information of power consumption, that are collected and harnessed for a more efficient and multi-faceted grid, may be used for purposes that are not related to energy management, thereby making it potentially dangerous to individual privacy of consumers. In fact, this is one of the main reasons for the lack of mass roll-out of Smart Meters in many countries [14].

Given this background, the problem statement of this paper is the following: how to preserve the privacy of energy consumption data of consumers but making the provision of certain services by the power provider still possible? This paper proposes a lightweight approach to mask each individual data (e.g., energy consumption measurements given by a Smart Meter) without affecting the aggregate data (e.g., the total consumption in a specific region during an interval of time, or the total consumption of a specific consumer throughout a billing period).

1.2 Solution Overview

To preserve the individual privacy, we propose a specification where each consumer (here consumer refers to Smart Meter) sends masked consumption values to the power provider instead of real values. For an individual measurement, the consumer reads the real energy consumption and adds a random number from an interval of $[-X; +X]$ (for instance, according to a

random uniform distribution). Thus, when the power provider sums all received data, it may obtain an approximation of the real total consumption, because the addition of the random numbers tends to be minimized after the sum operation.

This approach may insert error in the resulting total value. However, even with standard metering approaches, the existence of errors in applications of electrical networks is often found nowadays (since these errors should be less than established limits). For example, in Brazil, the INMETRO (National Institute of Metrology, Standardization and Industrial Quality) establishes percentage error limits to measurements with billing purposes. For residential consumers, the relative percentage error for active energy must stay between $\pm 2\%$, whereas for industrial consumers, this error must stay between $\pm 3\%$.¹

With the proposed approach, the change in the procedure between consumer and power provider is only the generation of a random number and the addition of that with the consumption information. This complexity is equals to the complexity of generating a random number (e.g., Mersenne-Twister algorithm can be used with a complexity of $O(p^2)$, where p is the polynomial degree that indicates the repetition period of the generator [17]). In this way, the proposed solution is simple and lightweight, making possible its deployment in devices with limited resources.

We claim that the solution meets the needs of consumers and power providers. In order to validate the consumer side (privacy), the correlation between the masked and the real profile as privacy metric is being used. To validate the power provider side (utility), we show that the masked values are still useful to many Smart Grid applications, and the resulting error between the sum of the masked data and the sum of the real data is being considered as utility metric. Moreover, the maximum allowed value of X is calculated based on the maximum error allowed.

Note that for billing purposes, which may be more sensitive to errors, the meter firmware could accumulate the sum of the random numbers added, and sends it with the last measurement of the billing period, ensuring that the billing error is zero and still providing no information about the detailed profile.

A key feature of the approach is the possibility to allow a consumer to choose his own privacy level. If to join into the Smart Grid a consumer requires a high privacy level, he can disclose a masked data using the maximum possible value of X . If after a period of time this consumer is more convinced to contribute with the Smart Grid by disclosing his data (e.g., in exchange of a tariff reduction proposed by the power provider), he can disclose a partially masked data, which reveal more information about his profile. If after another period this consumer is totally convinced about some benefits of disclosing his data (e.g., vendors detecting appliances that are not working properly and suggesting to this consumer efficient appliances), he can disclose his real data, i.e., using X equals to 0.

Since this masking approach considers only the data that is disclosed to the power provider, it does not affect the customer's possibility of analyze his own real consumption profile inside his

¹ For more information, see ordinance 375 of Sept. 27, 2011 [20].

home (this is a good approach to manage the consumption and reduce costs, e.g., peak load transfer [10]).

1.3 Literature Review and Related Work

There are many works proposing encryption and public key mechanisms to be used between Smart Meters and power providers [1]. Common encryption can be useful to ensure data confidentiality along the channel and ensure the total data utility to the power provider, since, after the decryption, the data is totally clean and trustworthy. However, the consumer's privacy inside the power provider (or outside if some malicious employee exports it) is still threatened.

Efthymiou *et al.* [7] describe a method for securely anonymizing frequent electrical metering data sent by a Smart Meter. However, the consumer identification is essential, since this data may be used for billing purposes.

A mechanism that may solve the privacy problem is homomorphic scheme [6, 9, 16, 22]. These solutions seek to support utility and privacy in different ways; however, they do not have a robust general theoretical basis for privacy and utility [21], i.e., many of them are application dependent. Moreover, in the last years, solutions for fully homomorphic schemes have been proposed, but it is hard to ignore efficiency concerns. Today, all known fully homomorphic encryption schemes have a long way to go before they can be used in practice [15].

The use of rechargeable batteries between appliances and Smart Meters can help to reduce the privacy issues in Smart Metering [12, 18, 25] as the signatures of appliances are no longer legible. Nevertheless, besides the signatures of battery cycles still being exposed, it is not always feasible to have batteries.

Bohli *et al.* [3] describe an approach to mask the data using noise from a normal distribution. The approach is evaluated to calculate the total consumption from a group of consumers and they conclude that a large number of consumers is necessary to obtain a considerable obfuscation level. Moreover, as the approach targets the computation of the aggregate consumption of the group, it does not solve more granular problems such as billing. Also related is the work from Wang *et al.* [23], which propose an approach to mask the data using GMM (Gaussian Mixture Models). The achieved obfuscation level with this approach is lower than the achieved by Bohli *et al.* and the approach is not evaluated using real profiles (they simulate that a consumption measurement is a random value from a normal distribution).

This paper proposes a masking approach to balance the consumer's privacy without affecting the data usefulness. Moreover, unlike other studies, it is applied the approach in Smart Grid applications and using data from real consumers to conclude that the approach is promising. Design, execution and analysis of possible attacks to the solution are also contributions.

2. A DATA MASKING APPROACH

The approach considers that the power provider is:

- Interested in knowing the total consumption of a consumer throughout a time period (e.g., one month for billing);
- Interested in knowing the total consumption of all consumers in a region at a certain time;

- Not interested in knowing the current consumption of an individual consumer.

Therefore, if each consumer sends a consumption measurement periodically to the power provider, it can organize these values as a matrix, where the sum of a row refers to the total consumption of a consumer throughout the time period, and the sum of a column refers to the total consumption of all consumers from this group at an instant of time, as shown in Figure 2. The rows of the matrix can be used for billing purposes (even if consumers have different billing periods) and the columns can be used for load monitoring.

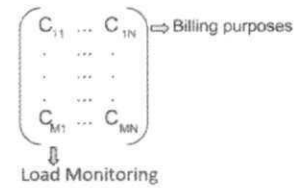


Figure 2. How measurement data can be used.

To hide the instantaneous power consumption, it is proposed that the Smart Meter sends masked measurements to the provider in a way that does not affect the results of the aggregating operations. To develop the solution, first billing (sum of a row) as a base application is considered and after that, how the approach works for load monitoring (sum of a column) is presented.

For billing purpose, if at each individual measurement the Smart Meter reads the consumption and adds a random number from an interval $[-X; +X]$, at the end of a billing period the result will be:

$$\sum_{i=1}^N c_i \approx \sum_{i=1}^N (c_i + x_i) \quad \sum_{i=1}^N c_i = \sum_{i=1}^N (c_i + x_i) - e_o$$

Where N is the number of measurements, x_i is a random number from $[-X; +X]$, c_i is a consumption measurement and e_o is the obtained error by adding random numbers. Therefore, the obtained error is the sum of all random values added:

$$e_o = \sum_{i=1}^N x_i \quad (1)$$

Given this formalization, if the X value increases, the masking level also increases and the accuracy level decreases. Therefore, the following technical problems are evident.

- Enabling the power provider to get a precise value (i.e., e_o less than an acceptable value) of the total consumption of a consumer over a billing period, how small should X be?
- Avoiding the power provider to infer information from the data masked by the consumer, how large should X be?

An empirical and an analytical mathematical model were developed such that, given a maximum acceptable error, finds the value of X that represents the equilibrium between privacy (masking) and utility (accuracy). Both models present very similar results and this may be a way to validate both.

For both models, the independent variables are X and N . The response variable is e_o . Moreover, the variable e_a is used to represent the maximum allowed error (e_o should be less than e_a).

2.1 Empirical Model

To find the empirical model, the following hypothesis was defined:

- H_{0-0} : The correlation between X and N for accuracy cannot be extracted.
 - There is no function $f(X,N)=e_o$ and consequently $f(e_o,N)=X$, to obtain the value of X that represents the equilibrium between utility and privacy.

Using the experimental design "K-factorial" [24] to test the effect of X and N in the error, both factors were varied. After that, we obtained many samples of each configuration and stored the worst cases (when the obtained error e_o is the maximum). Calculating a curvilinear regression with logarithms it was found the following model with an R^2 (coefficient that describes how well a regression line fits a set of data) of 97.5%:

$$\ln(e_o) = 0.32 + \ln(X) + \frac{\ln(N)}{2}$$

$$X = \frac{0.726 \cdot e_o}{\sqrt{N}} \quad (2)$$

Thus, the null hypothesis H_{0-0} is false (the function exists).

The value of X is calculated on the consumer side. To do that, the way that the allowed error e_o (in kWh) is calculated should be considered. Since it is a percentage of the real total value, the consumer does not have a way to predict how much energy he will consume throughout the billing period. According to experiments (including the consumer population used), it does not matter if e_o is a percentage of the real total consumption or a percentage of the total consumption in the previous billing period. Another way is calculating a new value of X for each measurement (instead of the same for all measurements during the billing period) and using allowed error values equals to percentages of previous N measurements (like a sliding window throughout the time, where for a consumption measurement c_i , X_i is calculated using an allowed error e_{oi} that is a percentage of the total consumption from the measurement c_{i-N} to c_{i-1}).

To show how this model can be used, suppose that the power provider wants to compute the total consumption of a consumer at the end of a month with 31 days. The measurements are collected at each 10 minutes, i.e., a total of $N=4,464$ measurements. If the maximum allowed error by the power provider (or by a regulatory agency) is a percentage value that corresponds to, for example, 2 kWh, the X value obtained is: $0.726 \cdot 2 / \sqrt{4,464} = 0.0217$.

Some experiments were made (e.g., 1,000 samples) to generate obtained error values using (1), and it was observed that the errors in fact stayed between -2 and 2 kWh, as presented in the left of Figure 3. The values that are outside of the -2 to 2 interval corresponds to approximately 2.5% and represents the values that the regression does not explain, since the obtained R^2 is 97.5%. The confidence interval in right of Figure 3 was obtained using a significance level of 95% and as we can see, the average of the obtained errors tends to zero.

2.2 Analytical Model

Moreover, an analytical model using probability theory was also developed and the results match with those obtained with the empirical model. Let x_i be a random variable uniformly distributed

between $-X$ and X . Its variance is $\sigma_x^2 = (X - (-X))^2 / 12 = X^2 / 3$. Now, for a large N , the central limit theorem ensures that the obtained error follows a normal distribution with mean $\mu_{e_o} = 0$ and variance $\sigma_{e_o}^2 = N^2 \cdot (\sigma_x^2 / N) = N \cdot \sigma_x^2 = N \cdot X^2 / 3$ (note that from (1), e_o is N times the mean of x_i).

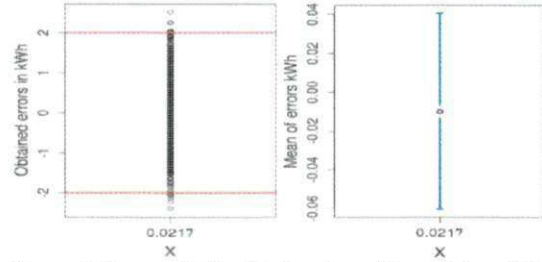


Figure 3. Errors obtained (e_o) and confidence interval for $X=0.0217$, $N=4464$ and $e_o=2$ kWh.

In other words, we can find the probability that the obtained error stays between two values using the normal distribution:

$$e_o \sim N\left(0, \frac{N \cdot X^2}{3}\right) \quad (3)$$

Using the same example from the previous section, suppose that the power provider wants to know the total consumption of a consumer at the end of a month with 31 days. With measurements of 10 minutes, it has a total of $N=4,464$ measurements. If the maximum allowed error e_o is 2 kWh, we have to find the variance $\sigma_{e_o}^2$ of the normal distribution such that the probability of the obtained error stays between -2 and 2 kWh is high, e.g., 0.98 ($P(-2 \leq e_o \leq 2) = 0.98$). This variance is $\sigma_{e_o}^2 = 0.739113$. So, $X = \sqrt{3 \cdot \sigma_{e_o}^2 / N} = \sqrt{3 \cdot 0.739113 / 4,464} = 0.0222$. This result is very close to the result obtained with the empirical model (also with other configurations), providing a validation to both models.

2.3 Utility and Privacy Metrics

As utility metric to see how the obtained final value by the power provider is accurate, the resulting percentage error between the sum of the masked data and the sum of the real data (i.e., e_o in percentage) has been used. A resulting percentage error close to 0% means a high utility level whereas a resulting percentage error distant from 0% means a low utility level. As it can be observed, if at the end of a billing period the obtained error is positive, the consumer will be paying a little more than actually consumed. But maybe in the next billing period the obtained error could be negative and the consumer will pay a little less. It is acceptable because as presented before, the error follows a normal distribution with mean 0. In fact, this imprecision is a penalty that the consumer pays to obtain some data privacy.

For load monitoring, an obtained error is not considered critical, since the power provider also has other ways to obtain accurate data (e.g., the power flow transmitted to a region). The proposed approach provides more information (that can be used for many purposes, such as leak and theft detection) to the power provider whereas providing privacy to consumers.

With noise addition, the measurement of how similar the original data and the perturbed data are is crucial [19]. There are

many metrics in literature, but we used the correlation between the masked and the real consumption profile as privacy metric. The most familiar measure of dependence between two quantities is the Pearson's correlation. The correlation coefficient between two consumption profiles A and B , with expected values μ_A and μ_B and standard deviations σ_A and σ_B , is defined as:

$$\text{corr}(A, B) = \frac{E[(A - \mu_A) \cdot (B - \mu_B)]}{\sigma_A \cdot \sigma_B}$$

E is the expected operator and corr a notation for correlation.

In the experiments, when the value of X increases, the correlation tends to 0. So, if the correlation between the real profile and the masked profile is close to 0, it has a high privacy level, whereas if it is close to 1, it has a low privacy level.

Other metrics were also analyzed for privacy, such as mutual information [21], mean square error and signal-to-noise ratio [19]. For some examples, these metrics were not satisfactory in a way of privacy level information. The correlation presented the best results (no false positives).

The use of correlation as privacy metric maybe does not detect if the consumer behavior is really being hidden. With its mathematical model, when the number of measurements in a time period is large, it obtains a better value of X to hide each individual measurement. As can be seen in (2), X is opposite to N , but more measurements also imply in lower consumption for each single measurement. For example, if the calculated value of X is 0.2 for one measurement of 0.2 kWh through 30 minutes, if it is divided in two measurements of 15 minutes, the values can be two measurements of 0.1 kWh and the new value of X is 0.1414. It is better to mask a 0.1 kWh value with 0.1414 than mask a 0.2 kWh value with 0.2. However, it is known that less measurements also implies in higher privacy. Taking an extreme example, of course that if a consumer sends to the power provider only one measurement with the total consumption at the end of the billing period, the privacy level is much better than sending measurements at each 15 minutes. The correlation metric is not considering this case.

When the frequency of measurements increases, it is fact that the privacy level decreases, but the masking power increases also.

3. BILLING

3.1 Example with Residential Consumer

The data used in the example below are measurements collected at each 30 minutes from a real residential consumer (anonymised) from Ireland [5].

As a starting example, it assumes the billing period as one month and the energy price as constant (time independent and without tariff policy). Figure 4 shows the full profile of a residential consumer throughout a month (March).

The consumption for this consumer during March is 167.04 kWh. Considering the maximal allowed error as 5% (8.352 kWh) and using the model obtained empirically, the following value of X was obtained (for measurements of 30 minutes, $N = 1,488$):

$$X = \frac{0.726 \cdot 8.352}{\sqrt{1,488}} = 0.1572$$

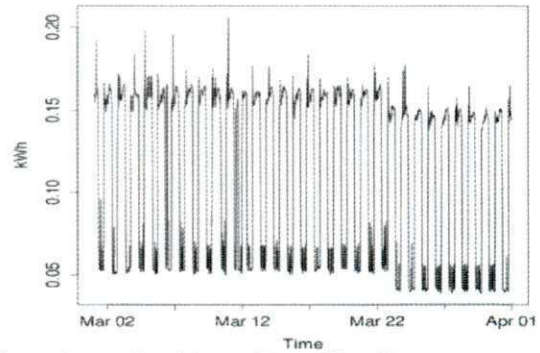


Figure 4. A residential monthly profile with measurements at each 30 min.

Figure 5 presents the monthly masked profile using the obtained X . A daily (March, 10) real profile and the daily masked profile together to see the obfuscation level were also plotted, as presented in Figure 6. The correlation coefficient between the real monthly profile and the masked monthly profile is 0.489 (our privacy metric). Even though this level of privacy is not too high, the measurements were made with a low resolution (30 minutes is considered a long period), which somehow contributes to the privacy level, as discussed before.

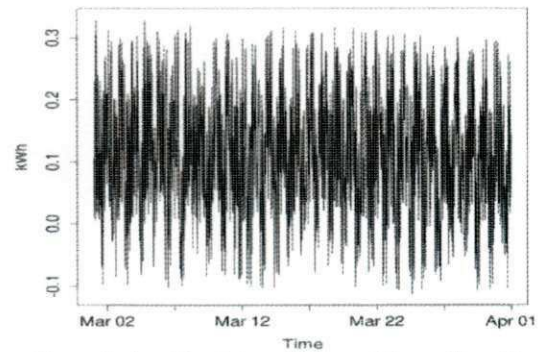


Figure 5. A residential masked monthly profile with measurements at each 30 min.

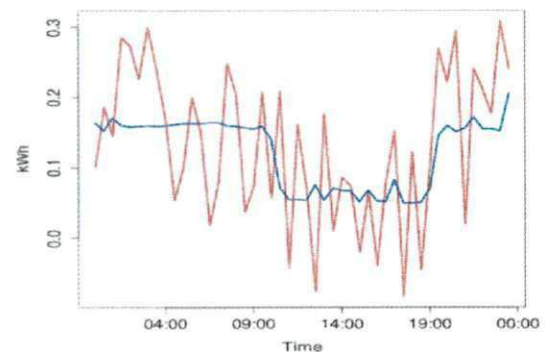


Figure 6. Real residential daily profile (blue) versus masked residential profile (red) with measurements at each 30 min.

It was considered that at the end of the month, the power provider computes the total consumption of this consumer for billing purposes. Summing the informed masked values by the consumer, the power provider obtained a value of 165.346 kWh

(the real value is 167.04 kWh). The difference between these values is an error of -1.01% (utility metric), less than the maximum error allowed (5%). However, if the meter firmware accumulates the sum of the random numbers added and sends that with the last measurement of the month, the error is zero.

3.2 Example with Industrial Consumer

The data used in the example below are measurements collected from a real industrial consumer (anonymised) at each 1 minute [8]. For a better visualization of the profile in a graphic, these data were transformed in measurements of 10 minutes (but remember that higher resolution implies in a better obfuscation).

It was supposed that the billing period is one month and the energy price as constant (time independent and without tariff policy). Figure 7 shows the full profile of an industrial consumer throughout a month (March).

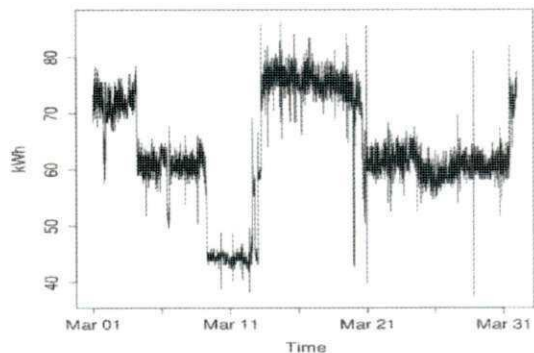


Figure 7. An industrial monthly profile with measurements at each 10 min.

The total consumption for this consumer during the month is 283,959 kWh. Considering that the maximal allowed error is 5% (14,197.95 kWh), the obtained value of X (for measurements of 10 minutes, $N = 4,464$) is 154,2766.

Figure 8 presents the monthly profile masked using the obtained value of X . The daily (March, 12) real profile and the daily masked profile were also plotted together to see the obfuscation level, as presented in Figure 9. As it can be seen, the privacy level is so high in a way that the real profile (blue line) is almost a straight in comparison to the masked profile (red line). In fact, for this example, the correlation coefficient between the real monthly and the masked monthly profile is 0.109 (privacy metric).

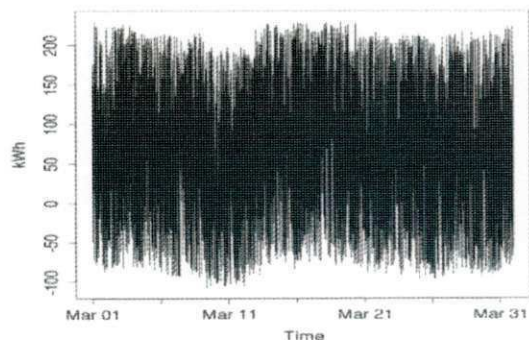


Figure 8. An industrial masked monthly profile with measurements at each 10 min.

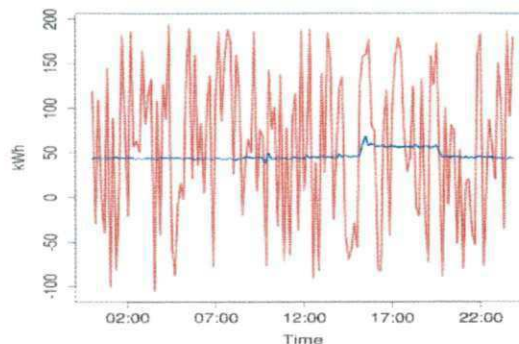


Figure 9. Real industrial profile (blue) versus masked industrial profile (red) with measurements at each 10 min.

It was considered that at the end of the month, the power provider wanted to know the total consumption of this consumer for billing purpose. Summing the informed masked values by the consumer, the power provider obtained a value of 292,536.6 kWh, but the real value is 283,959 kWh. The difference between these values is an error of 3.02% (the utility metric), less than the maximum error allowed by the power provider (5%).

3.3 Example with Tariff Policy

An example with tariff policy using the proposed approach was also considered. In Brazil, ANEEL (National Agency of Electrical Energy) establishes a regulation (PRORET - Procedure of Tariff Regulation [4]) that fixes three types of tariffs according to the period of the day:

- Peak: three consecutive hours defined by the power provider considering the load curve of its electrical grid;
- Intermediate: two hours, being one hour immediately before and another hour immediately after the peak period;
- Off-peak: the complementary hours (i.e., excluding the peak and intermediate periods).

Figure 10 shows an example of the tariff types established by ANEEL. The time periods can be determined by the sums of the columns of the matrix (Figure 2). In our example, we supposed that the power provider has established the following time based tariff policy: Peak: 16:00~19:00; Intermediate: 15:00~16:00 and 19:00~20:00; Off-peak: 00:00~15:00 and 20:00~00:00.

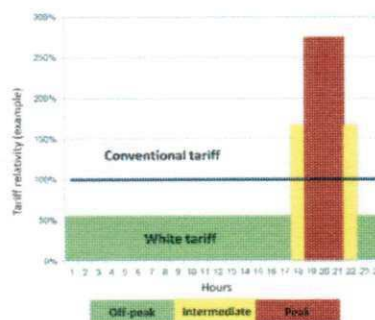


Figure 10. Example of the tariff types established by ANEEL.

In order to know the total electricity consumption per tariff period, X values should be computed separately per period. For

the same industrial consumer used in Section 3.2. Table 1 presents the real total consumptions, the X values, the obtained total consumptions (from the masked profile) and the computed errors for each tariff period. The obtained masked profile is like the presented in Figure 8, but the correlation coefficient obtained is equals to 0.131.

Table 1. Example with three types of tariffs

Tariff period	Total real (kwh)	X	Total obtained (kwh)	Error (max. 5%)
Peak	35,619.12	54.735	35,585.56	-0.094%
Intermediate	23,631.84	44.476	22,931.55	-2.963%
Off-peak	224,708.1	137.21	224,946.6	0.106%

4. LOAD MONITORING

If each consumer mask their data based on a billing period, the power provider may obtain accurate values for the rows of the matrix. But to get accurate values for the columns, the number of consumer must be as big as possible. A higher value of M (Figure 2) implies in higher accuracy, since consumers mask their data based on the billing period. It is wanted that the masking be unnoticed in the resulting data used for load monitoring.

The data used in the example below are measurements collected at each 30 minutes from real residential consumers (anonymised) from Ireland [5].

It was supposed that the power provider wanted to know the total consumption in a region with many consumers throughout the time for load monitoring (e.g., find peak times, leak detection, load forecasting and many other applications). To get accurate aggregate values using masked values sent by consumers, the number of consumers must be as big as possible. A higher value of M implies in higher accuracy and it is dependent of the population behavior. As an example, using a billing period of 1 month and measurements at each 30 minutes, it has $N = 1,488$ measurements during March. So, for experiments, was considered $M = 1,488$ consumers (a square matrix). Figure 11 shows the region profile during March obtained from real consumer profiles versus the region profile obtained from masked consumer profiles.

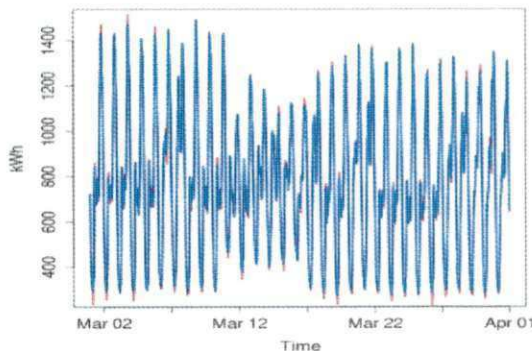


Figure 11. Region profile using real data (blue) versus using masked data (red) with measurements at each 30 minutes.

As it can be seen in Figure 11, the data looks so accurate in a way that both lines are almost similar (overlapped). However, an obtained error in kWh for a period of high consumption has a different meaning from the obtained error in a period of low

consumption. The large errors in Figure 11 were obtained in periods of low consumption (e.g., during the night), because while consuming less, consumers are still masking their data using a value of X based on the billing period. For a better accuracy, more consumers must be included (a higher value of M).

5. AN ATTACK TO THE SOLUTION

Another issue has also been researched: attacks to the proposed solution. It is considered as successful attacks those that can somehow minimize the effects of masking.

5.1 Attack Based on Similar Weekly Behavior

This attack is based on the hypothesis that the consumer tends to have a similar weekly behavior. The attacker can collect the consumer's data and calculate an expected week for this consumer. An expected week is composed by the seven expected days (from Sunday to Saturday), and an expected day is composed by the averages of each time from this specific day (e.g., the attacker calculate the expected Sunday from all available Sundays). Using the expected week, the attacker can try to predict the consumer behavior in future weeks. The attack effect is dependent on the amount of data available to the attacker.

5.2 Attack Evaluation

For the same residential consumer used in Section 3.1, this attack was done and the results are presented in Table 2. The average correlation between masked weeks and real weeks was compared versus the average correlation between the expected week and real weeks to analyze the attack effect. We used confidence intervals with significance levels of 95%.

As it can be seen in Table 2, when the number of weeks available to the attacker increases, the attack effect increases also, because the correlation between the expected week and real weeks is higher. But for our experiments, even choosing a consumer who repeats a behavior almost always (see Figure 4) and using a long period of observation (52 weeks or a full year), these correlations are less than the correlations between masked weeks and real weeks. It means that it is better to guess the consumer behavior from the own masked week than from the expected week. In that way, we considered this attack as unsuccessful.

Table 2. Effect of the attack of the days of week for a residential consumer

Number of available weeks to the attacker	Average correlation between masked weeks and real weeks	Average correlation between the expected and real weeks
2	(0.499, 0.510)	(-0.046, -0.033)
4	(0.335, 0.344)	(-0.005, 0.004)
8	(0.369, 0.374)	(0.154, 0.166)
16	(0.326, 0.331)	(0.171, 0.180)
32	(0.436, 0.439)	(0.182, 0.189)
52	(0.432, 0.434)	(0.316, 0.323)

6. CONCLUSION

We discussed several privacy and utility issues in the Smart Metering Infrastructure and proposed a lightweight mechanism that preserves the privacy of consumers without affecting significantly the data usefulness to the power provider. The modification in the communication procedure between a Smart Meter and the power provider is only the generation of a random

number and the addition of this number to the measurement to be sent to the power provider. Therefore, our approach is simple and lightweight. Using real examples of consumers and Smart Grid applications, the approach can be regarded as promising.

Although in this paper was used the uniform distribution to generate random numbers, further study is required to compare probability distributions for optimizations in the masking process. To do that, a better metric for privacy should be used (e.g., the number of discovered appliances after applying a NIALM algorithm) to really determine with probability distribution provides a better privacy to consumers and which one is more resistant to attacks. Furthermore, in future work, it is wanted to derive additional attacks to the solution. For example, given a masked profile of a consumer, try to filter the random noise using similar mechanisms to those used in signal processing of communication systems.

7. ACKNOWLEDGMENTS

We would like to thank the members of the TrueGrid project (www.truegrid.eu). Special thanks to Keiko V. O. Fonseca (UTFPR) for the feedback on this research. The research leading to these results has received funding from CAPES, DAAD and GIZ through the NoPa program (2011-2013).

8. REFERENCES

- [1] Baumeister, T. Literature review on smart grid cyber security. Collaborative Software Development Laboratory at the University of Hawaii, (2010).
- [2] Boccuzzi, C. Smart grid and the energetic big brother. *Metering International América Latina*, 3, (2010), 82-83.
- [3] Bohli, J., Sorge, C., and Uguo, O. A privacy model for smart metering. *Proc. IEEE Intl. Conf. Commun. Workshops (ICC)*, (2010), 1-5.
- [4] Brazilian Electricity Regulatory Agency (ANEEL). PRORET - Procedure of Tariff Regulation. (2011).
- [5] Commission for Energy Regulation (CER). CER smart metering project. (2012).
- [6] Dimitriou, T., Karame, G. Privacy-Friendly Tasking and Trading of Energy in Smart Grids. *Proc. of the 28th Annual ACM Symp. on Appl. Comp.*, Coimbra, Portugal, (Mar. 18-22, 2013), 652-659.
- [7] Efthymiou, C., and Kalogridis, G. Smart grid privacy via anonymization of smart metering data. *IEEE 1st Intl. Conf. Smart Grid Commun.*, Gaithersburg, MD, (Oct. 4-6, 2010), 238-243.
- [8] EnerNOC. 2012 Boston Cleanweb Hackathon and challenge. (May. 4-6, 2012).
- [9] Garcia, F. D., and Jacobs, B. Privacy-friendly energy-metering via homomorphic encryption. *Security and Trust Management*, 6710, (2011), 226-238.
- [10] Giordano, V., Onyeji, I., Fulli, G., Jimnez, M. S., and Filiou, C. Guidelines for cost benefit analysis of smart metering deployment. *JRC Scientific and Tech. Research*, (2012).
- [11] Ilić, D., Silva, P. G., Karnouskos, S., Jacobi, M. Impact assessment of smart meter grouping on the accuracy of forecasting algorithms. *Proc. of the 28th Annual ACM Symp. on Appl. Comp.*, Coimbra, Portugal, (Mar. 18-22, 2013), 673-679.
- [12] Kalogridis, G., Efthymiou, C., Denic, S. Z., Lewis, T. A., and Cepeda, R. Privacy for smart meters: towards undetectable appliance load signatures. *IEEE 1st Intl. Conf. Smart Grid Commun.*, Gaithersburg, MD, (Oct. 4-6, 2010), 232-237.
- [13] Kelly, J., and Knottenbelt, W. Disaggregating Smart Meter Readings using Device Signatures. *Imperial Computing Science MSc Individual Project*, (Sep. 2011).
- [14] Koehle, O. Just say no to big brother's Smart Meters. The latest in Bio-Hazard technology. *ARC Reproductions*, (2012).
- [15] Lauter, K., Naehrig, M., and Vaikuntanathan, V. Can homomorphic encryption be practical?. *Proc. of the 3rd ACM workshop on Cloud Comp. Sec.*, Chicago, Illinois, USA, (Oct. 17-21, 2011), 113-124.
- [16] Li, F., Luo, B., and Liu, P. Secure information aggregation for smart grids using homomorphic encryption. *IEEE 1st Intl. Conf. Smart Grid Commun.*, Gaithersburg, MD, (Oct. 4-6, 2010), 327-332.
- [17] Matsumoto, M., and Nishimura, T. Mersenne twister: a 623-dimensionally equidistributed uniform pseudorandom number generator. *ACM Trans. on Modeling and Comp. Simul.*, 8, 1, (Jan. 1998), 3-30.
- [18] McLaughlin, S., McDaniel, P., and Aiello, W. Protecting consumer privacy from electric load monitoring. *Proc. of the 18th ACM Conf. on Comp. and Commun. Sec.*, Chicago, Illinois, USA, (Oct. 17-21, 2011).
- [19] Mivule, K. Utilizing noise addition for data privacy, an overview. *Intl. Conf. on Inform. and Know. Engin.*, Las Vegas, USA, (Jul. 16-19, 2012).
- [20] National Institute of Metrology, Standardization and Industrial Quality (INMETRO). Ordinance number 375 of September 27, 2011, (2011).
- [21] Rajagopalan, S. R., Sankar, L., Mohr, S., and Poor, H. V. Smart meter privacy: a utility-privacy tradeoff framework. *IEEE 2nd Intl. Conf. Smart Grid Commun.*, Brussels, (Oct. 17-27, 2011), 150-155.
- [22] Rial, A., and Danezis, G. Privacy-preserving smart metering. *Proc. of the 18th ACM Conf. on Comp. and Commun. Sec.*, Chicago, Illinois, USA, (Oct. 17-21, 2011).
- [23] Wang, S., Cui, L., Que, J., Choi, D.-H., Jiang, X., and Xie, L. A randomized response model for privacy preserving smart metering. *IEEE Trans. on Smart Grid*, Vol. 3, No. 3, (Sep. 2012), 1317-1324.
- [24] Wohlin, C., Runeson, P., Höst, M., Ohlsson, M. C., Regnell, B., and Wessln, A. Experimentation in software engineering. Springer Publisher, (2012).
- [25] Yang, W., Li, N., Qi, Y., Qardaji, W., McLaughlin, S., and McDaniel, P. Minimizing private data disclosures in the smart grid. *Proc. of the 19th ACM Conf. on Comp. and Commun. Sec.*, NC, USA, (Oct. 16-18, 2012).