

Capítulo

3

Privacidade de Dados de Localização: Modelos, Técnicas e Mecanismos

Javam C. Machado, Eduardo R. Duarte Neto

Abstract

The growing development of mobile devices has promoted the increasing popularity of location services. However, the preservation of the users' privacy, especially the location data, has been questioned. This short-course describes the problem of the violation of the privacy of individuals' location and presents an in-depth analysis of the main techniques for their preservation, including differentially private mechanisms. Initially, fundamental concepts of data privacy will be presented, as well as vulnerabilities and threats to the privacy of individuals when exposing their location data when using applications on mobile devices. Then the state of the art in preserving location data privacy will be presented and discussed. Finally, we will point out research opportunities in the area and present relevant conclusions on the topic.

Resumo

O desenvolvimento crescente de dispositivos móveis tem promovido uma crescente popularidade dos serviços de localização. Entretanto, a preservação de privacidade dos usuários destes serviços, em especial dos dados de localização, tem sido bastante questionada. Este minicurso descreve o problema da violação da privacidade de localização de indivíduos e apresenta um aprofundamento das principais técnicas para sua preservação, incluindo mecanismos diferencialmente privados. Inicialmente serão apresentados conceitos fundamentais de privacidade de dados, bem como vulnerabilidades e ameaças à privacidade de indivíduos ao expor seus dados de localização quando do uso de aplicações em dispositivos móveis. Em seguida será apresentado e discutido o estado da arte em preservação de privacidade de dados de localização. Por fim iremos apontar oportunidades de pesquisas na área e apresentar conclusões relevantes sobre o tema.

3.1. Introdução

Com o desenvolvimento dos dispositivos móveis, a quantidade de dados coletados por aplicativos a fim de prover os mais diversos serviços tem crescido bastante. Estes dados têm se mostrado bastante valiosos, sendo utilizados nas mais diversas áreas. Por exemplo, muitas empresas têm se utilizado da análise destes dados para traçar o perfil de seus consumidores e assim, adotarem estratégias que venham a potencializar seus lucros. Já na área de saúde, uma análise sobre dados de saúde combinado a dados de localização pode ajudar a identificar que áreas estão mais sujeitas a certos patógenos, e assim, adotar medidas que venham a melhorar o atendimento de seus cidadãos. Por exemplo, um estudo sobre a taxa de infecção de novos casos de COVID-19 por região poderia ajudar em conter o avanço da doença, tratando de forma mais eficiente as áreas mais afetadas. Esse tipo de análise requer acesso a dados privados, o que levanta questionamentos quanto à privacidade dos indivíduos a quem os dados pertencem. Logo, encontrar uma forma de permitir esta análise sem que haja riscos a exposição dos mesmos tem sido objeto de estudo na área de privacidade de dados.

O crescimento da popularidade dos serviços baseado em localização (LBS) tem contribuído bastante para o aumento da quantidade de dados gerados, em especial, dados referentes à localização dos seus usuários. Através dos sensores destes dispositivos, as coordenadas de latitude e longitude são obtidas e utilizadas por estes serviços. Schiller [43] define os serviços de localização como serviços que integram a localização ou posição de um dispositivo móvel a outras informações, de modo a fornecer valor agregado a um usuário. Estes numerosos serviços, tais como navegação, redes sociais, serviços de recomendação, jogos de realidade aumentada, entre outros, têm sido desenvolvidos e integrados às atividades diárias das pessoas, provendo informações úteis sobre seus arredores e sendo capazes de responder perguntas do dia a dia como: qual a melhor rota a ser percorrida para um determinado endereço? Quais os pontos turísticos mais próximos da minha localização atual? Em quanto tempo o táxi que eu solicitei irá demorar para chegar em meu apartamento?

Estas informações de localização geradas por estes serviços podem potencializar vários outros serviços. Empresas e agências governamentais têm utilizado as informações de localização, tais como atividades praticadas nas localizações, para melhorar o serviço prestado, para o lançamento de um novo produto, ou até mesmo para gerar uma nova política pela empresa. Entretanto, acessar dados de localizações de usuários desses serviços, mesmo que com permissão, levanta severas preocupações de privacidade para a maioria dos usuários. Dessa forma, a utilização de serviços baseados em localização pode levar a sérios riscos de violação de privacidade devido a provedores de serviços não confiáveis [28], que podem expor os dados de localização de seus usuários ou até mesmo vender suas informações de localizações a terceiros [54]. De posse dessas informações, os dados obtidos por terceiros são utilizados para descoberta de dados sensíveis dos usuários, *i.e.*, dados de saúde, crenças religiosas, ideologias políticas, questões raciais, preferências sexuais, dentre várias outras.

Para exemplificar o risco da exposição dos dados de localização, podemos observar a Figura 3.1. O usuário Bob realiza em vários momentos requisições a um serviço de localização qualquer. A cada requisição ele envia sua localização corrente. No tempo

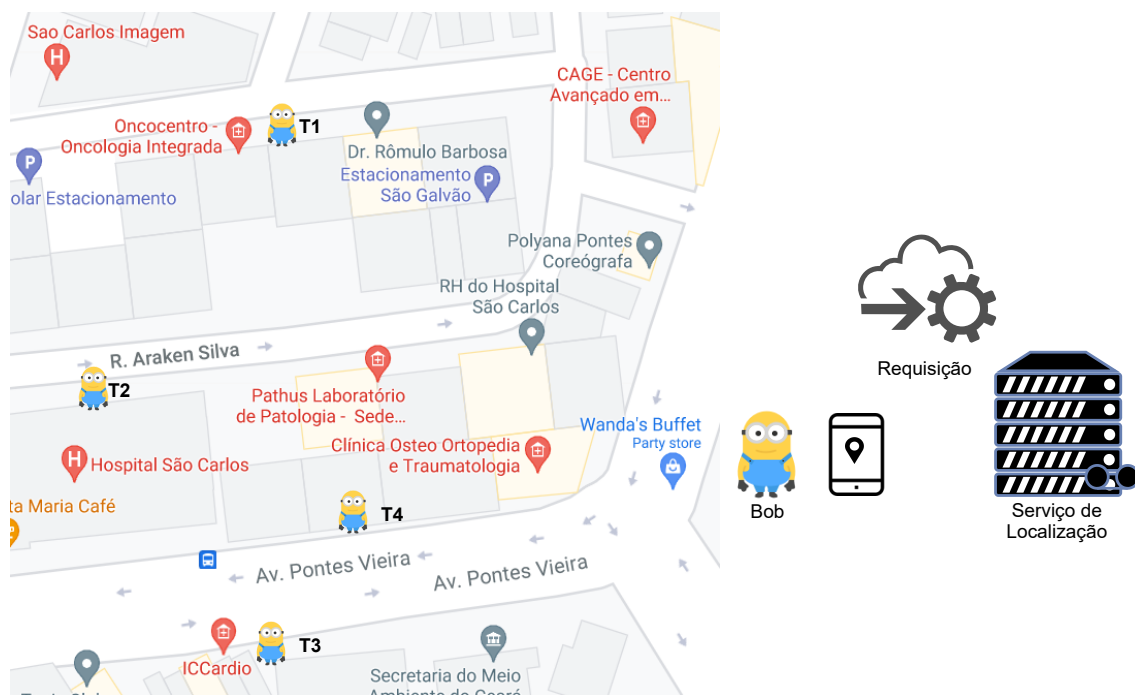


Figura 3.1. Exemplo de requisições realizadas próximas a hospitais e clínicas, permitindo inferências de dados sensíveis do usuário.

t_1 , Bob estava próximo a uma clínica de oncologia. No tempo t_2 , Bob realiza uma nova consulta próxima ao pronto socorro de um hospital. Em outros dois momentos sua localização também está próxima a localizações associadas à área de saúde. Considerando que é de conhecimento do provedor do serviço as requisições feitas pelos usuários, o próprio provedor como um possível agente malicioso pode inferir, com alta probabilidade, que Bob possui algum tipo de doença, ou que ele é da área de saúde, ou está acompanhando alguém enfermo. Estas informações, juntamente com outras informações de contexto, podem aumentar ainda mais o sucesso da inferência sobre dados sensíveis de Bob. Desta forma, o risco de uma violação de privacidade é bastante alto, deixando o usuário exposto.

A aplicação de modelos de privacidade sobre requisições de usuários é imprescindível para evitar que as localizações dos indivíduos não sejam identificadas pelos provedores no uso destes serviços. Todavia, em geral os modelos de privacidade acabam provocando mudanças nos dados, afetando diretamente a sua utilidade, com impacto direto na qualidade do serviço. Portanto, gerenciar essa solução de compromisso (*trade-off*, Figura 3.2) entre privacidade dos indivíduos e utilidade dos seus dados se torna um outro grande desafio. Desta forma, vários modelos de privacidade de dados têm sido propostos por pesquisadores com o objetivo de resolver esta questão.

Este capítulo tem por objetivo introduzir os fundamentos e técnicas para preservação da privacidade de dados dos indivíduos, procurando apresentar os riscos mais comuns e as técnicas mais populares na solução do problema. Em seguida, apresentaremos um aprofundamento sobre o tema privacidade em serviços de localização, apontando os conceitos básicos sobre dados de localização, os tipos de ataques a que estão sujeitos, e os principais modelos de preservação de privacidade de dados de localização na atualidade.

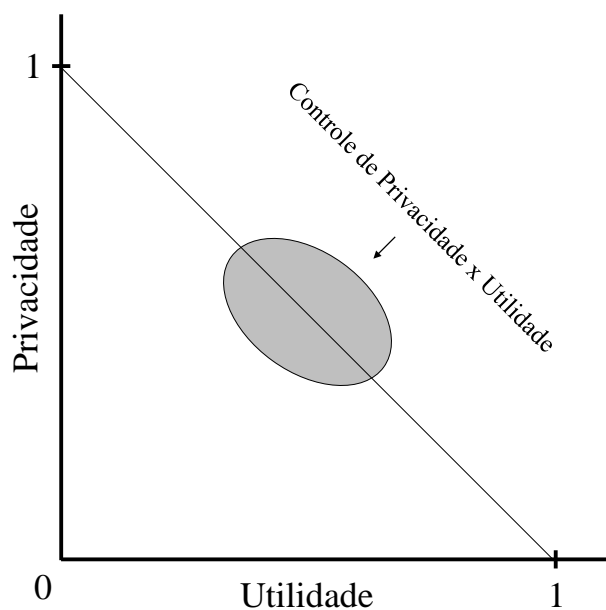


Figura 3.2. Trade-off entre privacidade e utilidade

Na Seção 3.2 apresentaremos os princípios básicos sobre o tema, que tipos de dados estão sujeitos a violação, os principais tipos de ataques, e como a preservação de privacidade pode ser alcançada. Os modelos sintáticos mais populares para preservação de privacidade são descritos na Seção 3.3. A Seção 3.4 apresenta o modelo de privacidade diferencial, estado da arte em preservação de privacidade. Na Seção 3.5 abordaremos o tema privacidade em Serviços de localização, onde descreveremos os serviços de localização e sua arquitetura. Iremos também apontar o problema da exposição temporal sem dados de localização. A Seção 3.5.4.3 descreve os principais tipos de ataque que dados de localização estão sujeitos. Os modelos de privacidade em dados de localização são descritos na Seção 3.6. Apresentaremos desafios de pesquisa na Seção 3.7. E por fim, a Seção 3.8 apresenta as considerações finais do capítulo.

3.2. Fundamentos da Privacidade de Dados

A privacidade é o direito que um indivíduo tem de manter seus assuntos pessoais e relacionamentos secretos [10]. O debate em torno do conceito de privacidade é uma matéria de extrema complexidade que muitas vezes é confundida com o conceito de segurança. Embora privacidade e segurança sejam temas relacionados, elas tratam de pontos bem distintos. No contexto de dados, a segurança define o controle de acesso durante o ciclo de vida do dado. Este controle se refere a regras específicas de quem está autorizado a acessar (ou não) determinados recursos. A forma como o acesso é realizado é papel da privacidade. Normalmente regida por leis e políticas de privacidade que definem o controle de acesso, permitindo a revelação da informação apenas por usuários autorizados. Entretanto, este controle de acesso não é suficiente para garantir a privacidade dos indivíduos, visto que os usuários com acesso àquelas informações podem ser maliciosos, e assim capazes de divulgar informações sensíveis acerca dos donos dos dados.

A quantidade de dados coletados tem aumentado bastante com o desenvolvimento e popularidade dos dispositivos móveis, tais como celulares, relógios inteligentes, dispositivos veiculares, dentre outros. Diversos serviços têm sido ofertados. Muitos dos quais exigem que o usuário abra mão da sua privacidade em favor da prestação destes serviços. Dessa forma, é fundamental identificar quais tipos de dados não devem ser divulgados, e portanto, aplicar técnicas que permitam a proteção destes dados garantindo a privacidade. Todavia, a qualidade de certos serviços pode depender diretamente da precisão destes dados privados. Tornando essencial que mesmo após a aplicação destas técnicas, os dados mantenham uma certa utilidade.

3.2.1. Privacidade em Microdados

De uma forma geral, os dados são representados por tabelas, onde cada linha da tabela corresponde a um registro no conjunto de dados, e as colunas contém os atributos dos registros. A esta representação dá-se o nome de microdados [20]. Os indivíduos estão associados a registros nestas tabelas. Os atributos são características ou propriedades dos indivíduos. No contexto de privacidade de dados, os atributos podem ser classificados em [9]:

1. **Identificadores explícitos:** são aqueles atributos que identificam de maneira única os indivíduos, como "CPF", "nome", etc., e devem ser removidos antes da publicação dos dados;
2. **Semi-identificadores:** são aqueles que não são identificadores explícitos, mas podem identificar o usuário, quando relacionados. "Data de nascimento" e "CEP" são exemplos de atributos semi-identificadores;
3. **Atributos sensíveis:** possuem informações sensíveis a cerca dos indivíduos, como "doença", "salário", etc.;
4. **Atributos não sensíveis:** são aqueles que não se enquadram em nenhuma das categorias citadas anteriormente.

Em privacidade de dados, os atributos sensíveis são aqueles de maior interesse porque apresentam potenciais danos ao seus donos em caso de divulgação. Por esse motivo, tais atributos necessitam ser protegidos. A Tabela 3.1 ilustra um exemplo de registros de indivíduos contendo atributos identificadores explícitos e semi-identificadores, que precisam ser protegidos.

3.2.2. Proteção e ataques à Privacidade

Como explanado na Seção 3.1, a análise dos dados é uma atividade fundamental, seja para melhorar a eficiência de um serviço, seja para ajudar na adoção de estratégias governamentais, ou para auxiliar na atividade econômica. Dessa forma estabelecer a confiança dos indivíduos e assim obter o consentimento para a utilização dos seus dados é um desejo de seus curadores. Portanto, é necessário garantir a proteção dos dados pessoais coletados. A anonimização é uma abordagem promissora para solucionar o problema de preservação de privacidade. Através de transformações do dados antes de sua publicação

Tabela 3.1. Exemplos de identificadores explícitos e semi-identificadores em dados tabulados de indivíduos.

Identificadores Explícitos		Semi-identificadores			
ID	Nome	Idade	Gênero	Endereço	CEP
1	Carla	24	Feminino	Av. I	60127002
2	João	21	Masculino	Av. K	60128001
3	Marcos	27	Masculino	Av. K	60128002
4	Ana	41	Feminino	Rua J	60127001

[21] procura-se impedir a exposição dos dados sensíveis dos indivíduos. Neste processo, um conjunto de dados D é transformado em um conjunto de dados D' , por meio de modificações sobre os dados. Técnicas de generalização, supressão e perturbação, são aplicadas sobre os dados para garantir esta transformação do conjunto de dados.

A generalização modifica os atributos semi-identificadores dos registros no conjunto de dados por valores mais gerais, aumentando a incerteza de um adversário associar um indivíduo a seus dados, ou em especial, a seus atributos sensíveis. Na abordagem mais comum de generalização, o valor de um atributo semi-identificador que se deseja proteger nos diferentes registros é substituído por um valor generalizado. Para exemplificar a aplicação da generalização sobre um conjunto de dados apresentamos a Figura 3.3. Ela ilustra a aplicação do processo de generalização sobre o atributo CEP (Código de Endereçamento Postal) presente nos registros da Tabela 3.1. Podemos observar que nas folhas da árvore têm-se os valores originais para o atributo de 4 registros. No segundo nível agrupam-se os CEPs cujos 5 primeiros dígitos correspondem, enquanto que no nível seguinte em direção ao topo da hierarquia são agrupados os CEPs cujos 4 primeiros dígitos correspondem.

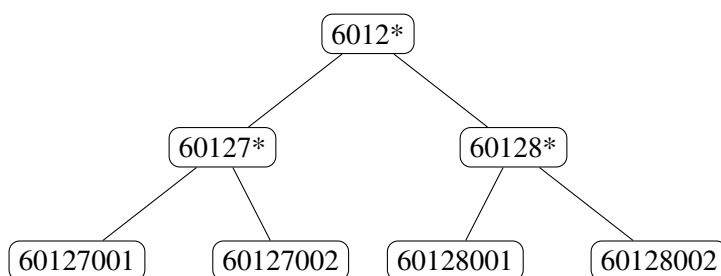


Figura 3.3. Exemplos de generalização do atributo semi-identificador CEP.

Como resultado da generalização sobre o atributo CEP da Tabela 3.1 temos a Tabela 3.2. Todos os CEPs dos registros foram substituídos pelo CEP generalizado da raiz da árvore, tornando este atributo indistinguível dos demais registros. Desta forma, este processo dificulta a re-identificação de um dos indivíduos caso seja de conhecimento externo o CEP de seu endereço. Por outro lado, a generalização diminui a precisão da localização, afetando diretamente a utilidade dos dados. Dessa forma, uma análise sobre os dados de CEP da 3.2 pode não ser útil, por exemplo, para se identificar uma distribuição

de gênero dentro das micro-regiões de uma cidade no caso dessa generalização diminuir bastante a precisão das localizações.

Tabela 3.2. Exemplo de semi-identificador anonimizado por generalização.

Identificadores Explícitos		Semi-identificadores			
ID	Nome	Idade	Gênero	Endereço	CEP
1	Carla	24	Feminino	Av. I	6012****
2	João	21	Masculino	Av. K	6012****
3	Marcos	27	Masculino	Av. K	6012****
4	Ana	41	Feminino	Rua J	6012****

A supressão de dados ocorre pela remoção dos valores de atributos ou pela sua substituição. Neste último caso, um ou mais valores do conjunto de dados é substituído por algum valor especial que dificulte a tarefa de descoberta de semi-identificadores por adversários. Alguns dos principais tipos de supressão:

- **Supressão de registro:** a supressão de registro remove um registro inteiro do conjunto de dados, conseqüentemente nenhum valor de atributo é disponibilizado para uso [6, 27].
- **Supressão de valor:** a supressão de valor remove ou substitui todas as ocorrências de um valor de um atributo semi-identificador por um valor especial, como “*”. Por exemplo, em uma tabela de funcionários de uma empresa, os valores de atributo salário abaixo de R\$ 30.000,00, podem ser removidos ou substituídos por “*”, enquanto os demais valores não sofrem distorções [50, 51].
- **Supressão de células:** nessa técnica, apenas algumas instâncias de valores de um atributo são removidas ou substituídas por um valor especial, caracterizando uma *supressão local* [36]. Por exemplo, pode-se remover apenas metade dos valores de atributo salário abaixo de R\$ 30.000,00, em uma tabela de empregados. Assim, instâncias de salário podem conter valores abaixo ou acima de R\$ 30.000,00, além de valores suprimidos. Entretanto, essa estratégia pode levar a inconsistências em eventuais análises de dados.

Por último, mas, não menos importante, a perturbação substitui os valores dos atributos semi-identificadores originais por valores fictícios, de modo que informações estatísticas calculadas a partir dos dados originais não se diferenciem significativamente de informações estatísticas calculadas sobre os dados perturbados. As técnicas mais comuns de perturbação de dados são:

- **Adição de ruído:** essa técnica é aplicada comumente sobre atributos numéricos. O valor original “ v ” de um atributo será substituído por “ $v + r$ ”, onde r corresponde ao ruído adicionado ao valor original. Para alcançar a proteção desejada, este ruído em geral é escolhido de forma aleatória seguindo alguma distribuição. Outra forma

de se aplicar a adição de ruído é através de um fator multiplicativo, onde o valor original é substituído por " $v \times r$ ". Os valores dos atributos são portanto, perturbados com um determinado nível de ruído, que pode ser adicionado ou multiplicado pelo valor original de cada atributo [47];

- **Permutação de dados:** nesta abordagem os valores de um mesmo atributo de dois registros diferentes são permutados. Isso mantém algumas características estatísticas dos dados, como frequência dos atributos e contagem [15]. Apesar desta técnica não alterar o domínio dos atributos, as possíveis permutações de valores diferentes podem levar a valores nos registros sem sentido, e com isso, informações equivocadas, tendo um impacto indesejável na utilidade dos dados;
- **Geração de dados sintéticos:** nesta técnica, um modelo estático é inicialmente gerado a partir do conjunto de dados e, após isso, são gerados dados sintéticos que seguem o modelo gerado [1]. Esses dados sintéticos são os que devem ser disponibilizados para o uso final. A vantagem desta técnica é que todas as propriedades estatísticas dos dados são mantidas. Entretanto, assim como na permutação de dados, nesta técnica pode-se gerar alguns valores sem sentido e que não são condizentes com o mundo real, embora as propriedades estatísticas mantenham-se fiéis às dos dados originais [33].

Esta transformação dos dados através da generalização, supressão ou perturbação protege os dados, permitindo, portanto, o seu compartilhamento a outras entidades. O nível de utilidade alcançado poderá garantir o uso das informações contidas nos dados sem que haja uma exposição dos indivíduos presentes no conjunto de dados. Todavia, assim como estas técnicas buscam garantir a privacidade dos indivíduos, pessoas mal intencionadas, comumente chamadas de atacantes, ou adversários, buscam se opor a esta proteção utilizando todos os recursos a sua disposição para retirar o máximo de informação dos registros. Esta informação a disposição do atacante utilizada no auxílio de inferências sobre o conjunto de dados é identificada por conhecimento prévio. Um exemplo de conhecimento prévio é o do adversário que trabalha no mesmo local da vítima, e conhece algumas informações privadas da vítima, tais como, endereço residencial, cargo na empresa, além de outras, permitindo a inferência de informações sensíveis, como localização, opção sexual, etc. Em se tratando de publicação de dados, o atacante pode ter acesso a outros conjuntos de dados previamente publicados, e assim, cruzar referências para descobrir novas informações sensíveis da vítima. Portanto, podemos concluir que o conhecimento do adversário é imensurável e imprevisível e deve ser levado em consideração nas soluções de preservação de privacidade, apesar de suas características.

Um adversário é capaz de violar a privacidade dos usuários por meio de diversos ataques que citaremos a seguir:

- **Ataque de Ligação ao Registro:** este ataque tem por objetivo re-identificar o registro de um usuário, cujas informações pertencem ao conjunto de dados publicado. Por exemplo, o atacante tem conhecimento do indivíduo João, mas não tem qualquer conhecimento sobre alguns de seus atributos, como seu endereço. Dessa forma, o atacante busca identificar que o registro de ID 2 da Tabela 3.1 pertence ao indivíduo João;

- **Ataque de Ligação ao Atributo:** o objetivo do adversário é ser capaz de inferir atributos sensíveis do usuário mesmo sem re-identificar seu registro, com base nos valores sensíveis relacionados ao grupo que o usuário pertence. Nesse tipo de ataque o adversário sabe que um usuário pertence a um certo grupo de registros, como por exemplo, que todos os usuário possuem CEP igual a 6012*** como na Tabela 3.2. Através desse conhecimento o adversário busca identificar informações sensíveis comuns a todos do grupo.
- **Ataque de Ligação à Tabela:** este tipo de ataque assume que o adversário sabe que o registro do usuário foi publicado. Neste ataque o intuito é inferir se a vítima está presente ou ausente nos dados publicados.
- **Ataque Probabilístico:** este ataque tem o foco de destacar como o adversário mudaria seu pensamento probabilístico sobre um usuário depois de ter acesso ao conjunto de dados disponível. Por exemplo, após analisar um conjunto de dados publicado, inferir a probabilidade de o usuário ser do sexo masculino.

3.3. Modelos de Privacidade Sintáticos

Os modelos de privacidade sintáticos exigem que o conjunto de dados anonimizados possuam uma forma definida que ajuda a reduzir o risco de quebra de privacidade [14]. Os modelos sintáticos, em geral, aplicam transformações nos registros por meio de técnicas de supressão e/ou generalização até que esta forma seja alcançada. Iremos apresentar alguns dos modelos de privacidade sintáticos mais utilizados em preservação de dados.

3.3.1. k -anonimato

Modelo de privacidade mais conhecido no campo da anonimização de dados [46], o k -anonimato assegura que, para cada combinação de valores de semi-identificadores, existem pelo menos k registros no conjunto de dados, formando uma classe de equivalência. O k -anonimato atua sobre o princípio da indistinguibilidade, isto é, cada registro em um conjunto de dados k -anônimo é indistinguível de pelo menos outros $k - 1$ registros em relação ao conjunto de semi-identificadores. Desta forma, a probabilidade de se ligar qualquer indivíduo a um registro no conjunto de dados é de no máximo $\frac{1}{k}$.

O nível de privacidade é ajustado em função do parâmetro k , afetando diretamente o equilíbrio entre utilidade e privacidade. Assim, um valor de k grande implica em uma maior proteção dos dados, entretanto, diminui a utilidade dos mesmos, por ser necessário adicionar grande volume de ruído a fim de se alcançar classes de equivalência com pelo menos k registros. É importante ressaltar que não existem abordagens analíticas para determinar um valor ótimo para o parâmetro k [12], sendo este um problema NP-difícil [36]. Dessa forma, cabe aos *dataholders* esta complexa tarefa de se definir o nível desejado de privacidade que seja adequado para garantir um equilíbrio adequado entre privacidade e utilidade.

A Tabela 3.3 será tomada como base para aplicar alguns modelos de privacidade sintáticos, entre eles o k -anonimato. São atributos identificadores explícitos: Placa, Motorista e CPF. São atributos sensíveis: Tipo de Multa e Valor da Multa. Os demais atributos são semi-identificadores.

Tabela 3.3. Dados sobre infrações de trânsito [31].

	Placa	Motorista	CPF	Data de Nascimento	Data da Infração	Tipo de Multa	Valor da Multa (R\$)
1	UVW-1840	Gigi	223.512.956	14/03/1980	03/01/2013	1	170
2	AXO-2064	André Luis	523.512.511	04/03/1980	03/01/2013	2	250
3	AUG-1046	Juçara Silva	123.998.687	24/05/1980	03/01/2013	1	170
4	FBI-1001	Bruno Lima	230.320.523	20/04/1982	04/01/2013	1	170
5	ACO-6241	Abu Ali	221.320.876	20/05/1982	04/01/2013	2	250
6	ABA-5012	Pedro Ramires	210.329.890	13/05/1982	05/01/2013	2	250
7	HBV-2002	Eduardo Neto	538.687.045	15/05/1982	05/01/2013	1	170

A Tabela 3.4 foi gerada após a aplicação de técnicas de supressão nos identificadores explícitos e de generalização nos atributos semi-identificadores. Nesta tabela podemos perceber quatro classes de equivalência para os semi-identificadores: Classe A = “03/1980, 01/2013” nas linhas 1 e 2; Classe B = “05/1980, 01/2013” registro 3; Classe C = “04/1982, 01/2013” com o registro 4 e Classe D = “05/1982, 01/2013” nas linhas 5, 6 e 7. Observe, que após aplicar o processo de anonimização, k -anonimato ainda não foi alcançado para um $k = 2$, já que as classes B e C, não possuem uma quantidade mínima requerida de 2 registros, sendo, portanto, necessário, algum novo processo de transformação. Uma estratégia válida seria remover os registros 3 e 4, como podemos observar na Tabela 3.5, onde a tabela agora contém apenas 2 classes de equivalência, e alcançando o k -anonimato para $k = 2$.

Tabela 3.4. Dados sobre infrações de trânsito anonimizados [31].

	Placa	Motorista	CPF	Data de Nascimento	Data da Infração	Tipo de Multa	Valor da Multa (R\$)
1	*	*	*	03/1980	01/2013	1	170
2	*	*	*	03/1980	01/2013	2	250
3	*	*	*	05/1980	01/2013	1	170
4	*	*	*	04/1982	01/2013	1	170
5	*	*	*	05/1982	01/2013	2	250
6	*	*	*	05/1982	01/2013	2	250
7	*	*	*	05/1982	01/2013	1	170

Tabela 3.5. Tabela no modelo 2-anonimato [31].

	Placa	Motorista	CPF	Data de Nascimento	Data da Infração	Tipo de Multa	Valor da Multa (R\$)
1	*	*	*	03/1980	01/2013	1	170
2	*	*	*	03/1980	01/2013	2	250
3	*	*	*	*	*	*	*
4	*	*	*	*	*	*	*
5	*	*	*	05/1982	01/2013	2	250
6	*	*	*	05/1982	01/2013	2	250
7	*	*	*	05/1982	01/2013	1	170

3.3.2. *l*-diversidade

Assim como o *k*-anonimato, o *l*-diversidade age sobre o princípio da indistinguibilidade. Entretanto, apesar de sua popularidade, simplicidade e eficácia contra ataques de ligação ao registro, o *k*-anonimato não se mostra adequado contra ataques de ligação ao atributo, *i.e.*, ataques em que um adversário procura inferir informações sensíveis sobre registros mesmo sem identificá-los. Observando a Tabela 3.5 que garante o *k*-anonimato, para $k = 2$, podemos identificar pelo menos dois registros em cada uma das classes de equivalência. Entretanto, imagine o cenário hipotético, em que o atacante tem conhecimento de um outro conjunto de dados não anonimizado, em que é possível identificar um indivíduo cujo nome é Pedro que nasceu em maio de 1982, e sofreu uma infração em janeiro de 2013. Ele poderá inferir com uma probabilidade de $\frac{2}{3}$ que a multa recebida por Pedro foi do tipo 2, e seu valor foi de 250 reais. Superior à $\frac{1}{2}$, desejada pelo modelo *k*-anonimato.

O *l*-diversidade busca prover proteção a ataques de ligação ao atributo, garantindo que para cada classe de equivalência, exista pelo menos *l* valores distintos para cada atributo sensível. Assim, o que se pretende é que um atacante, mesmo com conhecimento prévio sobre a classe de equivalência de um registro, não seja capaz de inferir o atributo sensível do mesmo com probabilidade maior que $\frac{1}{l}$. Na Tabela 3.6, a probabilidade de se identificar que o indivíduo tem asma, valor do atributo sensível "Doença", caso o atacante tenha conhecimento de que o CEP do indivíduo é 540040, é de 100%, superior a $\frac{1}{4}$ exigido por um modelo 4-anonimato. Convertendo a Tabela 3.6 para o modelo 3-diversidade, não é preciso fazer nenhuma alteração nos registros da classe A (linhas 1 a 4), pois esta já possui no mínimo 3 valores distintos para o atributo sensível. Entretanto, a classe B (linhas 5 a 8) possui todos os valores de atributos sensíveis iguais. Uma solução simples seria suprimir os registros das linhas 5 a 8. Outra solução seria modificar os valores do atributo sensível destas linhas por valores diferentes que garantam a diversidade, conforme Tabela 3.7 que atende, portanto, o modelo 4-anonimato e 3-diversidade.

Tabela 3.6. 4-anonimato [31].

	Idade	CEP	Cidade	Doença
1	<70	560001	*	Sinusite
2	<70	560001	*	Gripe
3	<70	560001	*	Zika
4	<70	560001	*	Hérnia
5	<35	540040	*	Asma
6	<35	540040	*	Asma
7	<35	540040	*	Asma
8	<35	540040	*	Asma

Tabela 3.7. 4-anonimato e 3-diversidade [31].

	Idade	CEP	Cidade	Doença
1	<70	560001	*	Sinusite
2	<70	560001	*	Gripe
3	<70	560001	*	Zika
4	<70	560001	*	Hérnia
5	<35	540040	*	Sinusite
6	<35	540040	*	Zika
7	<35	540040	*	Asma
8	<35	540040	*	Asma

Variações do *k*-anonimato e *l*-diversidade também foram propostas como uma extensão desses modelos com a finalidade de prover uma maior garantia de preservação de privacidade tanto contra ataques de ligação ao registro, como ao atributo [29, 49].

3.4. Privacidade Diferencial

Nos modelos de privacidade sintáticos uma violação de privacidade ocorre quando o atacante utilizando de quaisquer meios de conhecimento consegue re-identificar indivíduos

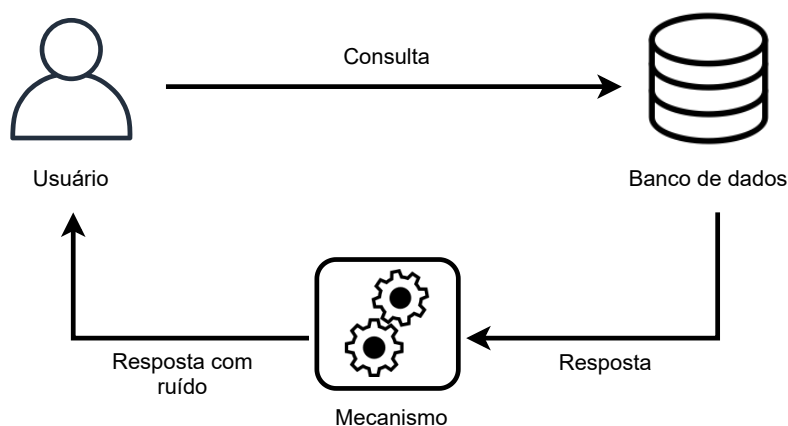


Figura 3.4. Ambiente interativo no modelo de Privacidade Diferencial.

dentro do conjunto de dados publicado em formato tabulado. Sob outra perspectiva, a Privacidade Diferencial investiga a ideia de publicar resultados de consultas, ao invés de dados tabulados, de tal forma que são adicionados ruídos a esses resultados. Dessa forma, a privacidade é garantida ao se anonimizar a resposta de consultas ao conjunto de dados. Assim, um atacante não será capaz de concluir algo com 100% de confiança. A sua principal convicção é de que as conclusões obtidas sobre um indivíduo são referentes aos dados de toda a tabela, e não apenas a um registro em particular que possa ser ligado a um indivíduo. Por esse motivo, o modelo de privacidade em questão propõe evitar ataques probabilísticos.

3.4.1. Conceitos Básicos

Proposto em [17], a Privacidade Diferencial fornece sólidas garantias de privacidade. Seu objetivo é solucionar o paradoxo entre garantir o aprendizado de informações úteis sobre a população de um conjunto de dados sem permitir obter qualquer informação específica sobre indivíduos desta população [19]. Um mecanismo diferencialmente privado garante a privacidade das consultas através da adição de um ruído aleatório controlado. Este modelo foi projetado em um ambiente interativo, onde os usuários submetem consultas a um conjunto de dados e este, por sua vez, responde por através do mecanismo. A Figura 3.4 apresenta este modelo iterativo, onde o mecanismo garantirá a privacidade ao introduzir “aleatoriedade” na geração do ruído, e portanto, protegendo os resultados das consultas realizadas sobre um conjunto de dados em seu formato original.

A Privacidade Diferencial assegura que qualquer sequência de resultados (isto é, resposta de consultas) é igualmente possível de acontecer independente da presença de qualquer indivíduo no conjunto de dados [19]. Assim, a adição ou remoção de um indivíduo não afetará consideravelmente o resultado de qualquer análise estatística realizada no conjunto de dados [13]. Portanto, o conhecimento adquirido por um atacante sobre qualquer indivíduo presente no conjunto de dados após realizar consultas a este conjunto não deve ser maior ao que ele já possuía antes de realizar qualquer consulta a este mesmo conjunto.

3.4.2. Definição Formal

Dado um algoritmo aleatório (mecanismo) M , este mecanismo garante ϵ -Privacidade Diferencial se para todos os conjuntos de dados vizinhos D_1 e D_2 no conjunto de dados, que diferem de no máximo um elemento, e para todo S contido na variação de resultados de M , isto é, para todo $S \subseteq \text{Range}(M)$,

$$Pr[M(D_1) \in S] \leq \exp(\epsilon) \times Pr[M(D_2) \in S],$$

onde Pr é a probabilidade dada a partir da “aleatoriedade” de M . Em outras palavras, a definição formal afirma que a diferença máxima entre as distribuições de probabilidade de uma consulta retornar o mesmo resultado para dois conjuntos de dados vizinhos é limitada pelo parâmetro ϵ . Portanto, para qualquer par de entradas que diferem de apenas um registro, para cada saída, um adversário não será capaz de distinguir entre os conjuntos de dados D_1 e D_2 baseado apenas na resposta fornecida pelo mecanismo.

A Figura 3.5 mostra um exemplo de probabilidades de saída de um algoritmo M , nos conjuntos de dados vizinhos D_1 e D_2 , a partir de um valor de ϵ . O algoritmo M fornece garantias de privacidade adicionando ruído aleatório no seu retorno, i.e., $M(D) = f(D) + \text{ruído}$, onde f é a resposta de uma consulta realizada por um usuário. No exemplo apresentado, foi utilizado o mecanismo de Laplace, seguindo a distribuição de probabilidades de mesmo nome, o que resulta em uma distribuição com pico mais acentuado do que se utilizasse uma distribuição em função da normal. O conceito de mecanismo é definido na Seção 3.4.3.

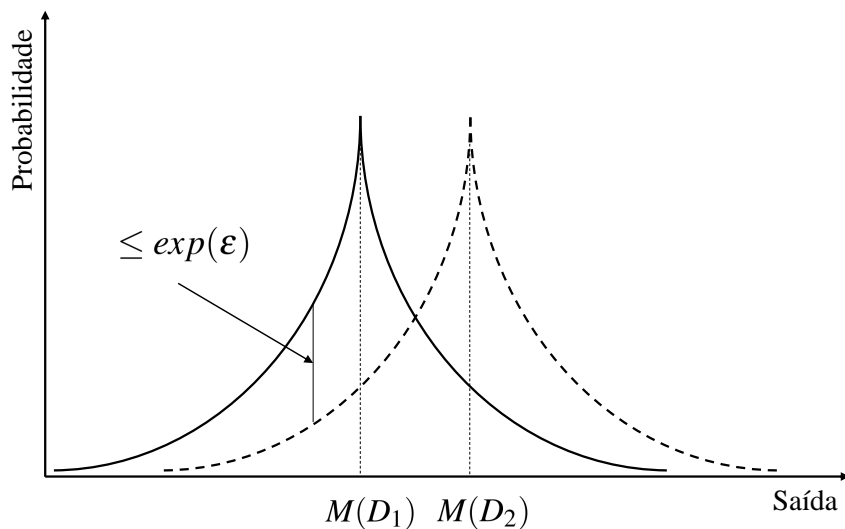


Figura 3.5. Probabilidades de saída de um algoritmo aleatório M sobre os conjuntos de dados vizinhos D_1 e D_2 .

Dessa forma, se um indivíduo pertence a um conjunto de dados D onde análises estatísticas serão feitas através de um mecanismo que é ϵ -Diferencialmente Privado, esse mecanismo irá garantir que a probabilidade de violação de privacidade não irá aumentar caso este indivíduo não pertencesse ao conjunto de dados. Dessa forma, podemos concluir que, como a Privacidade Diferencial é uma propriedade estatística restritiva a que o

mecanismo está sujeito, as garantias que ela oferece são altas, inclusive essas garantias não dependem de poder computacional ou informações que um atacante possa ter obtido.

O parâmetro ϵ que controla o nível de ruído adicionado pelo mecanismo não possui uma correlação explícita com a privacidade dos indivíduos como em outras técnicas vistas. Esse parâmetro depende da consulta que está sendo feita e dos próprios dados que estão no conjunto de dados, devendo, portanto, ser escolhido por um especialista com o objetivo de se obter o melhor equilíbrio entre privacidade e utilidade das respostas. A literatura em geral concorda que o valor de ϵ deva ser pequeno, como por exemplo 0,01, 0,1 ou até logaritmo natural $\ln 2$ ou $\ln 3$ [18]. Quanto menor o valor de ϵ , maior a privacidade. Tradicionalmente, a escolha do valor de ϵ se dá de forma empírica, portanto, para cada mecanismo, deve ser feita uma análise para escolher o parâmetro adequado utilizando métricas [39] para avaliar a precisão da resposta do mecanismo com diversos valores de ϵ [26].

3.4.3. Mecanismo e Sensibilidade

Como dito anteriormente nesta seção, a Privacidade Diferencial é idealizada em um modelo interativo, onde o usuário submete consultas a uma base de dados D , e um determinado mecanismo fornece uma resposta ϵ -Diferencialmente Privada. Porém, existem diversas formas de se atingir a Privacidade Diferencial através de um mecanismo. O objetivo das técnicas que utilizam esse modelo de privacidade é criar um mecanismo M que irá adicionar um ruído adequado para produzir uma resposta a uma consulta f feita pelo indivíduo, de forma que esse ruído seja independente do conjunto de dados D .

A quantidade de ruído necessária depende do tipo de consulta f aplicada sobre um conjunto de dados. Dessa forma precisamos definir o que é a sensibilidade de um conjunto de dados D . Antes disso, porém, precisamos definir de forma prática o que são conjuntos de dados vizinhos.

Definição 1 *Dado um conjunto de dados D , todos os conjuntos de dados D_i decorrentes da remoção de um indivíduo i do conjunto de dados original D são definidos como vizinhos [10].*

Por exemplo, considere o conjunto de dados D na Tabela 3.8. Um possível conjunto de dados vizinhos pode ser obtido através da remoção do registro de $ID = 6$, resultando na Tabela 3.9.

A sensibilidade por sua vez procura quantificar a diferença que um usuário faz ao ser removido do conjunto de dados na resposta da função de consulta. Isso é fundamental para o cálculo adequado do ruído a ser adicionado pelo mecanismo, uma vez que quanto maior o valor de Δf , mais ruído terá de ser adicionado à resposta do mecanismo para mascarar a remoção de um indivíduo, de forma a assegurar a privacidade do mesmo [13].

Tabela 3.8. Conjunto de dados D .

	Idade	CEP	Cidade
1	35	560001	Fortaleza
2	40	560001	Aquiraz
3	55	560001	Messejana
4	21	560001	Eusébio
5	35	540040	Aracati
6	35	540040	Caucaia
7	45	540040	Sobral
8	22	540040	Fortaleza

Tabela 3.9. Conjunto de Dados D' .

	Idade	CEP	Cidade
1	35	560001	Fortaleza
2	40	560001	Aquiraz
3	55	560001	Messejana
4	21	560001	Eusébio
5	35	540040	Aracati
7	45	540040	Sobral
8	22	540040	Fortaleza

Definição 2 *Seja D o domínio de todos os conjuntos de dados. Seja f uma função de consulta que mapeia conjuntos de dados a vetores de números reais. A sensibilidade global da função f é:*

$$\Delta f = \max_{x,y \in D} \|f(x) - f(y)\|_1$$

para todo x, y diferindo de no máximo um elemento, ou seja, vizinhos [18].

O mecanismo de Laplace é o algoritmo de adição de ruído controlado mais comum e simples para alcançar a Privacidade Diferencial. A adição de ruído é baseada na geração de uma variável aleatória da distribuição de Laplace com média μ e escala b de forma que

$$Laplace_{\mu,b}(x) = \frac{1}{2b} \exp\left(-\frac{|x|}{b}\right).$$

Podemos então definir formalmente o mecanismo de Laplace.

Definição 3 *Dada uma função de consulta $f : D \rightarrow \mathfrak{R}$, o mecanismo de Laplace M :*

$$M_f(D) = f(D) + Laplace(0, \Delta f / \epsilon)$$

fornece ϵ -Privacidade Diferencial. Onde $Laplace(0, \Delta f / \epsilon)$ retorna uma variável aleatória da distribuição de Laplace com média zero e escala $\Delta f / \epsilon$.

Utilizaremos a Tabela 3.10 para demonstrar a aplicação do mecanismo de Laplace de forma simplificada para melhor compreensão. A figura contém hipoteticamente um conjunto de dados da Receita Federal, que contém o número de imóveis que um determinado indivíduo declarou.

Suponha a consulta f que retorne o total de imóveis de todos os indivíduos da base de dados. Primeiramente, é necessário calcular a sensibilidade da função f sobre o conjunto de dados. Para isso calculamos f para cada conjunto de dados vizinho. A resposta real da consulta é 14. A Tabela 3.11 mostra os conjuntos de dados vizinhos gerados a partir da base original e suas respectivas respostas da consulta f .

Tabela 3.10. Exemplo de conjunto de dados original contendo o número de imóveis de cada indivíduo.

ID	Contribuinte	Nº de Imóveis
1	João	4
2	Bruno	2
3	Iago	7
4	Malu	1

Tabela 3.11. Conjuntos de dados vizinhos gerados a partir da base original e suas respectivas respostas da consulta f (soma).

ID	Contribuinte	Nº de Imóveis
2	Bruno	2
3	Iago	7
4	Malu	1

$$f(D_1) = 2 + 7 + 1 = 10$$

ID	Contribuinte	Nº de Imóveis
1	João	4
2	Bruno	2
4	Malu	1

$$f(D_3) = 4 + 2 + 1 = 7$$

ID	Contribuinte	Nº de Imóveis
1	João	4
3	Iago	7
4	Malu	1

$$f(D_2) = 4 + 7 + 1 = 12$$

ID	Contribuinte	Nº de Imóveis
1	João	4
2	Bruno	2
3	Iago	7

$$f(D_4) = 4 + 2 + 7 = 13$$

Portanto, a sensibilidade é dada pela variação máxima que a ausência de um indivíduo provoca no resultado da consulta. Essa variação é obtida quando da remoção do registro de $ID = 4$, cuja diferença máxima é de 7. Por fim, o ruído a ser adicionado para atender ao modelo de Privacidade Diferencial, utilizando o mecanismo de Laplace, deve ser igual a $Laplace(0, \frac{7}{\epsilon})$.

Como dito anteriormente, o parâmetro ϵ é definido por um especialista, o detentor dos dados. A Tabela 3.12 apresenta cinco exemplos de ruído, respostas e probabilidade de ocorrência após a aplicação da Privacidade Diferencial sobre o conjunto de dados original da Tabela 3.10, considerando $\epsilon = 1$. Assim, após utilizar o mecanismo de Laplace, o valor de ruído de $-4,58$ possui probabilidade de ocorrência de $3,7\%$ sobre o valor original da consulta f (cujas soma do número de imóveis original é igual a 14), resultando em um valor anonimizado de $9,42$ imóveis. De forma análoga, o valor de ruído de $-0,15$ possui uma probabilidade um pouco maior de ocorrência ($6,98\%$), caso a mesma consulta seja realizada nesse conjunto de dados, conforme mostra a Tabela 3.12.

Apesar do mecanismo de Laplace apresentar bons resultados, em se tratando de consultas do tipo categórica, que não permitem uma aproximação de seus valores categóricos, o mecanismo exponencial emergiu como uma alternativa eficaz satisfazendo ϵ -privacidade diferencial. [35] propôs o mecanismo exponencial justamente para garantir a privacidade em situações onde se deseja a melhor resposta à consulta. Diferente do mecanismo de Laplace, que adiciona ruído aleatório à resposta, o mecanismo exponencial para uma consulta qualquer busca selecionar de forma aleatória uma resposta do conjunto de possíveis respostas. Portanto, este mecanismo é ideal em anonimizar consultas onde necessita medir a utilidade da resposta enquanto se preserva a privacidade diferencial.

Tabela 3.12. Cinco possíveis valores de ruído, resposta e probabilidade de ocorrência após a aplicação da Privacidade Diferencial.

Ruído	$f(D) + \text{ruído}$	$Pr(f(D) + \text{ruído})\%$
-4,58	9,42	3,70
-0,15	13,85	6,98
12,15	26,15	1,25
-6,43	7,57	2,85
2,89	16,89	4,72

Tabela 3.13. Tabela de número de imóveis dos contribuintes e a pontuação retornada pela função de utilidade.

Contribuinte	Nº de Imóveis	Pontuação
João	4	16
Bruno	2	10
Iago	7	28
Malu	1	4

A distribuição de probabilidades de respostas a uma consulta irá considerar uma função de utilidade que irá mapear todos os possíveis conjuntos de dados D e todas as possíveis respostas O em uma pontuação de utilidade $u : (D \times O) \rightarrow \mathfrak{R}$ para um ε . Δu é a sensibilidade da consulta em termos da função de utilidade u :

$$\Delta u = \max_{o \in O} \max_{D_1, D_2: \|D_1 - D_2\|_1 \leq 1} |u(D_1, o) - u(D_2, o)|,$$

para todo D_1, D_2 diferindo de no máximo 1 elemento.

Definição 4 Para qualquer função de utilidade $u : (D \times O) \rightarrow \mathfrak{R}$, e um orçamento de privacidade ε , o mecanismo exponencial $M_u^\varepsilon(D)$ produz o como resposta com uma probabilidade proporcional a $\exp(\frac{\varepsilon u(D, o)}{2\Delta u})$, onde Δu é a sensibilidade da função de utilidade:

$$Pr[M_u^\varepsilon(D) = o] = \frac{\exp(\frac{\varepsilon u(D, o)}{2\Delta u})}{\sum_{o' \in O} \exp(\frac{\varepsilon u(D, o')}{2\Delta u})}$$

Utilizando os dados da Tabela 3.10 podemos aplicar o mecanismo exponencial para responder uma consulta que busque o nome do indivíduo que possui o maior número de imóveis. Para aplicar o mecanismo, primeiramente precisamos definir uma função de utilidade que indique o quão adequado é cada possível resposta. Dessa forma, nossa função de utilidade precisa estar relacionada ao número de imóveis. Suponha que nossa função de utilidade seja o nome do contribuinte vezes o número de imóveis que ele possui, ou seja, $u(D, o) = \text{len}(o.\text{contribuinte}) \times o.\text{imoveis}$, onde D é o conjunto de dados e o é um registro em D .

A Tabela 3.13 apresenta as possíveis respostas para a consulta e a pontuação obtida pela função de utilidade u . Observe que a resposta correta para a consulta é *Iago*, justamente a resposta com a maior pontuação retornada pela função de utilidade. O próximo passo é medir a sensibilidade de f . A Tabela 3.14 apresenta o impacto na função de utilidade ao se remover um elemento do conjunto de dados original. Como Δ_u é a máxima diferença na função de utilidade em conjuntos de dados vizinhos, a sensibilidade da função para os conjuntos de dados vizinhos da Tabela 3.14 é de $\Delta_u = 28$, quando se remove o contribuinte *Iago* do conjunto de dados.

Tabela 3.14. $\Delta_u = 28$, quando se remove o contribuinte *Iago* do conjunto de dados.

ID	Contr.	Nº. Imóveis	Pont.
2	Bruno	2	10
3	Iago	7	28
4	Malu	1	4

$$|u(D, o) - u(D_1, o)| = 16$$

ID	Contr.	Nº. Imóveis	Pont.
1	João	4	16
2	Bruno	2	10
4	Malu	1	4

$$|u(D, o) - u(D_3, o)| = 28$$

ID	Contr.	Nº. Imóveis	Pont.
1	João	4	16
3	Iago	7	28
4	Malu	1	4

$$|u(D, o) - u(D_2, o)| = 10$$

ID	Contr.	Nº. Imóveis	Pont.
1	João	4	16
2	Bruno	2	10
3	Iago	7	28

$$|u(D, o) - u(D_4, o)| = 4$$

Uma vez calculada a sensibilidade da consulta em termos da função de utilidade, a Tabela 3.15 contém as probabilidades de cada uma das possíveis respostas serem retornadas pelo mecanismo exponencial, onde a probabilidade da resposta correta ser retornada pelo mecanismo é de 31,45 %.

3.4.4. Aplicações

Leis de proteção a dados pessoais quanto à coleta e uso têm sido cada vez mais frequentes em diversos países, tais quais o Brasil. Dessa forma, muitas organizações detentoras de dados têm desenvolvido aplicações para garantir a privacidade sobre os dados coletados de seus clientes. A Microsoft, por exemplo, desenvolveu o PINQ [34], uma plataforma de análise de dados projetada para fornecer garantias de privacidade para os registros contidos em suas bases de dados. Agindo como uma camada intermediária entre a base de dados e o dono dos dados, o PINQ utiliza um sistema de consultas próprio, denominado LINQ, permitindo-o realizar consultas sem comprometer a privacidade dos indivíduos pertencentes a base.

Tabela 3.15. Cinco possíveis respostas e probabilidade de ocorrência após a aplicação da Privacidade Diferencial através do Mecanismo Exponencial.

Resposta	$Pr(f(D))\%$
João	25,43
Bruno	22,75
Iago	31,35
Malu	20,45

A Google através dos seus diversos aplicativos coleta uma quantidade enorme de dados de seus usuários, utilizando para diversos fins, incluindo melhorar a qualidade de seus serviços. Visando proteger estas informações ela lançou uma ferramenta denominada *Rappor* (Randomized Aggregatable Privacy Preserving Ordinal Responses). Essa ferramenta utiliza a perturbação na coleta de dados, mantendo as informações estatísticas necessárias para realizar suas análises e preservando a privacidade dos usuários que utilizam seu navegador.

A Apple em 2016 anunciou que aplica nos dispositivos de sua marca a privacidade diferencial na coleta de dados dos usuários. De acordo com a companhia, os dados são coletados de maneira privada, a fim de proteger a privacidade do usuário mas permite fazer análise de dados agregados e implementar melhorias no serviço. Dessa forma, recursos como Siri e até o QuickType poderão prever melhor as palavras que, por exemplo, um determinado conjunto de usuários mais utilizam.

Nesta subseção apresentamos vários modelos e técnicas utilizados por grandes organizações para garantir a privacidade. Além delas, é possível identificar muitos exemplos de aplicações no mundo real que também utilizam as técnicas vistas neste minicurso, como [37], [8] e [32].

3.5. Serviços de Localização

Como já destacado na Seção 3.1, o desenvolvimento dos dispositivos móveis equipados com GPS, juntamente com a disponibilidade de redes sem fio, tem contribuído com um aumento significativo da popularidade dos Serviços de localização [3]. O serviço prestado utiliza da localização do usuário, muitas vezes em tempo real, para prover alguma valia ao solicitante do serviço. São alguns exemplos comuns de serviços de localização:

- **Navegação:** provê ao usuário direcionamentos a um ponto de interesse geograficamente localizado. Os dados de localização do usuário em tempo real são utilizados como referência às instruções de direção. São algumas aplicações: Google Maps e Waze.
- **Aplicações de tempo (clima):** estes serviços apresentam condições climáticas em tempo real, bem como previsões. A localização da solicitação é usada para obter informações relevantes sobre o clima local.
- **Jogos:** utilizam a localização do usuário no contexto do ambiente virtual do jogo. Os mais recentes usam tecnologia de realidade aumentada, onde a movimentação do usuário em tempo real se reflete no jogo. Exemplo desse tipo de jogo é Pokemon GO.
- **Serviços de Recomendação:** estes serviços utilizam da localização do usuário para enviar recomendações de locais de interesse próximos. São exemplos: Foursquare e Yelp.

Nesta seção, iremos destacar a arquitetura dos serviços de localização, a natureza dos dados de localização, e o porquê destes dados serem de grande valia na análise de dados, ao mesmo tempo que são objeto de ataques à privacidade.

3.5.1. Arquitetura

Em um sistema tradicional de serviço de localização as informações de localização são obtidas por meio de sistemas de posicionamento global (GPS), presente na maioria dos dispositivos móveis da atualidade. O usuário, através de seu dispositivo móvel, realiza uma requisição ao provedor do serviço tendo como referência sua localização, como por exemplo, a previsão do tempo. O serviço então atenderá esta requisição.

A Figura 3.6 ilustra um típico serviço baseado em localização com preservação de privacidade [30]. São alguns componentes básicos de um LBS:

- **GPS:** permite determinar a localização dos objetos envolvidos, *i.e.*, usuários, ou outra entidade qualquer. O GPS é o mais popular sistema de posicionamento. Ele é um mecanismo de posicionamento por satélite que fornece a um aparelho receptor a sua posição.
- **Usuários:** são participantes que irão usufruir do serviço baseado em localização prestado. Através de dispositivos como *smartphones*, *notebooks*, *wearables*, os usuários se conectam ao meio de comunicação e enviam requisições ao provedor do serviço
- **Rede de comunicação:** é o meio através do qual acontece o tráfego de informações entre os participantes. Normalmente o meio utilizado é a rede de banda larga móvel, como a 4G.
- **Servidor do LBS:** é o responsável por receber as requisições dos usuários e prestar o serviço baseado em localização de acordo com sua natureza, seja para encontrar uma localização, seja para auxiliar na navegação do usuário, ou um outro tipo de serviço qualquer que utiliza a informação de localização enviada na requisição.
- **Provedor de Conteúdo/Dados:** provedores de Conteúdo/Dados fornecem dados e conteúdo ao servidor LBS. Alguns provedores de LBS possuem seus próprios dados e conteúdo, enquanto outros usam um terceiro para fornecer esse serviço.
- **Servidor de Privacidade:** o servidor de privacidade de localização executa os algoritmos de preservação de privacidade, como anonimização e criptografia e pode ser de propriedade e operado pelo provedor LBS ou por terceiros.

3.5.2. Dados de localização

Os dados de localização, em geral, possuem informações agregadas. Principalmente quando estamos falando de Pontos de Interesse (*POI*), tais como restaurantes, hospitais, fábricas, dentre outros. Um hospital por exemplo, está associado a doenças, horário de atendimento, exames. Uma fábrica está associada a trabalhadores, produtos, horário de expediente. Esta informação agregada garante a existência de uma relação entre as localizações, podendo ser medida através da correlação entre elas [53].

A informação de localização é uma componente chave em uma requisição a um serviço de localização. As informações de localização são constituídas de três partes principais: identidade, tempo e posição [30], Figura 3.7.

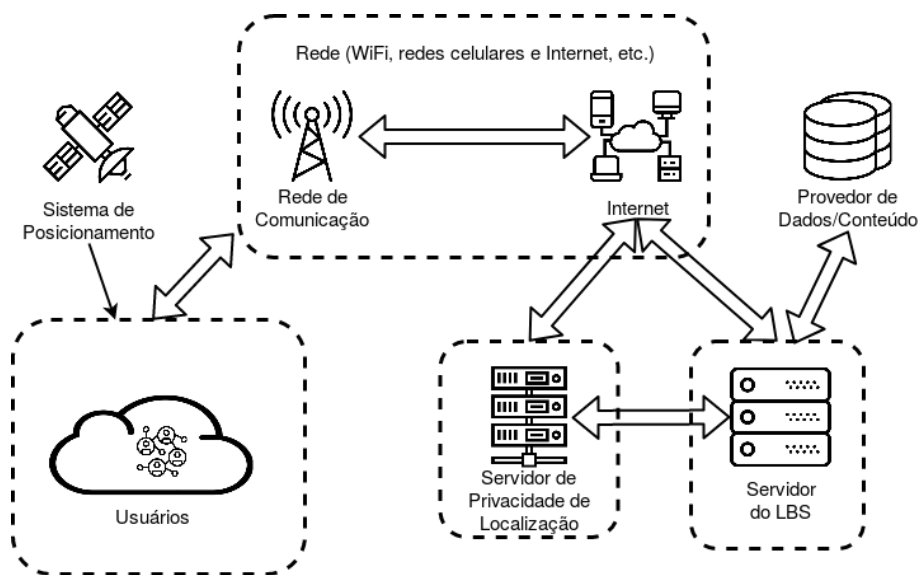


Figura 3.6. Modelo de sistema de serviços baseados em localização



Figura 3.7. Três atributos da informação de localização

A *identidade* está relacionada ao usuário do serviço de localização. Pode ser um endereço de e-mail, nome ou qualquer outra informação que torne um indivíduo distinguível dos outros. Esta identidade pode ser: (i) consistente, que é aquela requisitada obrigatoriamente para o acesso a um serviço, como um nome de usuário; (ii) inconsistente, através do uso de pseudônimos; (iii) anônima, onde há a ausência de uma identificação.

O *tempo* é referente ao momento ao qual as localizações estão associadas. Em geral, os serviços de localização associam as localizações a marcos de tempo. Esta informação temporal pode ser acumulada ou corrente. As aplicações em tempo acumulado não publicam as informações de localização em tempo real, mas em um tempo posterior ao atual. Um exemplo destas aplicações é o sistema de rastreamento do *Fitbit*, que coleta as informações percorrida pelos usuários, mas só publica a trajetória computada depois de encerrada a atividade [30]. As aplicações em tempo real publicam as informações de localizações associadas ao marco de tempo atual, de forma imediata.

A *informação espacial (posição)* é a forma de se identificar no plano espacial a localização do usuário. Em geral, a posição é determinada em função das coordenadas

geo-espaciais do indivíduo, através de sua latitude e longitude, ou através de alguma outra representação desta localização. Ela pode ser única, também chamada de simples, quando as localizações, em geral pela natureza do serviço, não são correlacionadas entre si, ou podem formar uma trajetória no caso em que são fortemente correlacionadas, o que acaba por gerar maiores riscos de exposição.

Quanto à representação, as localizações podem ser classificadas em diretas ou indiretas. As localizações diretas representam a posição precisa da localização, através de suas coordenadas de latitude e longitude. Os serviços de localização tradicionais usam localizações diretas. As localizações indiretas são aquelas estabelecidas com base na proximidade física, substituindo a localização exata pelo POI mais próximo, dotado de informação complementar sobre o mesmo, como o nome do estabelecimento, horário de funcionamento, endereço, entre outras.

Como exposto na Seção 3.1, a análise da informação de localização pode ser extremamente útil, inclusive na otimização de performance e de qualidade do mesmo. Entretanto, essa exposição também traz riscos à privacidade, permitindo a inferência de dados sensíveis do usuário, como por exemplo, seu estado de saúde, suas crenças religiosas, posicionamentos políticos, dentre outros. A privacidade de localização pode ser alcançada através da proteção dos três atributos que formam a informação de localização de uma pessoa, identidade, posição e tempo. Entretanto, é importante destacar que a garantia de privacidade de localização de indivíduos não é absoluta, e não é esperado que seja. Uma privacidade de localização sem qualquer vazamento de informação, inviabilizaria completamente um serviço. Por exemplo, se o indivíduo quer saber a previsão do tempo de sua cidade. É esperado que a localização enviada seja da cidade em que ele quer saber a previsão. Caso contrário, o serviço pode não atingir um nível adequado de qualidade. Desta forma, pode-se identificar dois requisitos principais da privacidade de localização dos indivíduos: a expectativa de privacidade dos indivíduos de "circunstâncias normais", ou seja, o que o indivíduo espera em termos de exposição da sua localização, e a maneira como as informações são coletadas e usadas. A expectativa de privacidade de uma pessoa pode mudar com o tempo, assim como a forma como as informações de localização são coletadas e usadas também mudam. Logo, para avaliar a privacidade de localização do indivíduo, seus principais requisitos devem ser definidos do ponto de vista dos usuários.

3.5.3. Exposição Temporal

A exposição temporal em serviços de localização constitui igualmente informação sensível dos usuários desses serviços. Informações referentes ao momento no tempo em que certas requisições foram realizadas, ou mais especificamente, em que momento no tempo se esteve em determinadas localizações, pode contribuir para a identificação de indivíduos, bem como, inferir informações sensíveis [52]. No contexto de trajetórias, além da exposição temporal, é possível identificar uma correlação forte entre as localizações que a compõem, permitindo a inferência de características semi-identificadoras [41].

Esta exposição é potencializada em função do período de observação. Quanto maior o período, maior é a exposição temporal gerada. O trabalho [40] define como calcular a probabilidade de se identificar a localização do usuário em curto período de observação e em longo período.

3.5.3.1. Curto período de observação Temporal

A localização do usuário é inferida em função de apenas um momento de tempo da localização ofuscada. Para isso, é necessário o conhecimento *a priori* sobre a localização ofuscada. Este conhecimento pode ser calculado de diversas formas. Uma forma bastante utilizada é o histórico do usuário de acessos ao serviço de localização. Desta forma, sendo L o conjunto de todas as localizações possíveis, a distribuição de probabilidade a posteriori $Pr(l|l')$, para $l, l' \in L$, pode ser calculada por:

$$Pr(l, l') = \frac{f(l)M(l'|l)}{\sum_{l \in L} f(l)M(l'|l)},$$

onde l' é a localização ofuscada, $f(l)$ denota o conhecimento a priori sobre a localização atual e $M(l'|l)$ representa a distribuição do mecanismo de ofuscação. Calculada a distribuição a posteriori, a localização atual l^* pode ser estimada de duas formas:

$$l^* = \arg \max_{l \in L} Pr(l|l'),$$
$$l^* = \arg \min_{l^* \in L} \sum_{l \in L} Pr(l|l') deuc(l^*, l).$$

3.5.3.2. Longo período de observação Temporal

As localizações ofuscadas reportadas pelo usuário são observadas por um período específico de tempo. Considerando uma situação em que o usuário aplica a ofuscação sobre a mesma localização várias vezes neste período de tempo, por exemplo, seu local de trabalho, o conjunto $O = \{o_1, \dots, o_n\}$ é o conjunto das frequências com que as localizações são reportadas neste período. A estimativa da localização atual l^* pode ser dada por:

$$l^* = \arg \max_{l \in L} o_l,$$

3.5.4. Tipos e Métodos de Ataques

Um atacante, também chamado de adversário, é qualquer entidade que possa ter acesso aos dados de localização de um ou de vários indivíduos visando se beneficiar [30]. Em geral, modelos de privacidade consideram o provedor de serviço, ou o curador dos dados, honesto mas curioso [23]. Desta forma, um adversário pode ser desde o próprio provedor do serviço de localização, ou até mesmo um cientista de dados que tenha acesso a uma publicação dos dados [42]. O conhecimento adversário é justamente essa informação prévia acessível ao atacante, e que pode ser usada por exemplo para identificar um usuário, uma localização, ou uma sequência de lugares de um objeto móvel [41].

Os dados de localização, como explicado na Seção 3.5.2, por serem carregados de informação, permitem que informações de contexto sejam adicionadas a este conhecimento adversário gerando potenciais riscos de violação de privacidade. São exemplos de conhecimento de contexto: o número de usuários em uma área em uma determinada hora do dia; a relação entre diferentes usuários; as restrições de localização de uma determinada área, como rede de ruas, área de preservação; a distribuição e a probabilidade estatística associada às localizações.

Em particular, assume-se que o atacante possui qualquer base de dados que contém conhecimentos adicionais sobre a semântica das informações de localização dos usuários. Além disso, o provedor do LBS pode identificar que o usuário está utilizando alguma técnica de preservação de privacidade de localização a fim de garantir a utilização do serviço sem expor sua localização real. Ataques a privacidade de localização podem ser aplicados em função da identidade, ou em função da localização do usuário [30].

3.5.4.1. Ataque de Identidade

Os ataques de identidade procuram cruzar conhecimentos adversários de diversas fontes a fim de determinar a identidade do alvo. São alguns exemplos deste tipo de ataque:

- **Ataque de identificação pessoal:** através do conhecimento prévio pessoal de um indivíduo, busca-se identificar o indivíduo dentro do conjunto de dados, a fim de se obter toda a informação a ele associada no conjunto de dados. Considere o exemplo: o atacante tem o conhecimento sobre o endereço residencial de um indivíduo. Através dele, mesmo em um conjunto de dados anonimizado, se o atributo endereço não tiver sido protegido, o atacante poderá identificar o dono do registro em função do endereço residencial exposto.
- **Ataque de presença agregada:** identificar a identidade com base na relação entre dois indivíduos ou através de uma propriedade agregadora, por exemplo pessoas agrupadas próximas a um evento, uma estação de Pokemon Go, ou uma loja com ofertas, dentre outros eventos.

3.5.4.2. Ataque de Localização

Os ataques de localização consistem em identificar as informações espaciais e temporais referentes a um indivíduo. São alguns exemplos de ataques de localização:

- **Ataque a localizações sensíveis:** procura identificar localizações importantes, como residência ou local de trabalho.
- **Ataque de revelação de presença ou ausência:** determina se um usuário está presente ou ausente em determinadas localizações em um determinado horário do dia.
- **Ataque de rastreamento:** identifica uma sequência de eventos para rastrear um usuário.

3.5.4.3. Métodos de Ataque

Os métodos de ataque dizem respeito à forma como o ataque é realizado. São alguns destes métodos:

- **Ataques de vinculação de contexto:** é a forma mais comum em ataques de localização. O conhecimento de contexto é combinado com a informação de localização obtida para se chegar à localização precisa da vítima em um ataque de localização. Por exemplo, um indivíduo ao realizar um *check-in* em um hospital, preenche seus dados informando seu endereço residencial. Se um atacante tiver conhecimento do endereço residencial do indivíduo, ele poderá usá-lo para identificar este indivíduo na lista de *check-in* do hospital.
- **Ataques probabilísticos:** este tipo de ataque é baseado na coleta de informações estatísticas sobre o ambiente [44]. Pode ser aplicado tanto para ataques de identidade, como de localização. Sendo assim, a localização do usuário pode ser inferida em razão da probabilidade de o usuário estar em uma determinada localização em um horário preciso.
- **Ataque de conluio de usuários maliciosos:** é realizado por usuários que usam o mesmo provedor de serviços baseado em localização, que colidem para realizar vários ataques. Por exemplo, usuários em conluio utilizam sua posição para obter, do serviço, a distância da vítima, e baseado nisso calculam a exata localização da vítima.

3.6. Privacidade de Dados de Localização

A privacidade de dados de localização busca proteger a privacidade dos usuários no uso dos serviços de localização. O que se pretende é entregar uma margem de segurança para os indivíduos em serviços de localização. Esta margem de segurança pode ser medida em função do erro de estimativa de um adversário, uma métrica muito utilizada para medir o nível de privacidade em termos da habilidade de um atacante em estimar a localização real de um usuário [2]. A Figura 3.8 exemplifica o erro de estimativa como sendo a distância euclidiana entre a localização real de um usuário e a localização estimada pelo adversário em preto.

Outro ponto importante em serviços de localização se refere a utilidade da informação de localização reportada. Quanto mais preciso for a informação reportada, maior é sua utilidade, e conseqüentemente, melhor é a qualidade do serviço. Esta utilidade pode ser medida em função do erro de utilidade. A Figura 3.9 exemplifica o erro de utilidade como sendo a distância euclidiana entre a localização reportada e a real localização do usuário.

O serviço de localização tem um impacto direto no nível de privacidade e utilidade desejado, portanto, é importante medir o nível de privacidade e utilidade em função destas métricas citadas. Por exemplo, um serviço de requisição de Pontos de Interesse, exige uma precisão na localização reportada bem maior do que um serviço de meteorologia, que certamente possui um relaxamento na precisão da localização na casa de quilômetros. Dessa forma, diversas técnicas de privacidade foram propostas a fim de preservar a privacidade dos indivíduos em serviços de localização. Esta seção apresenta uma visão geral das técnicas mais utilizadas para a preservação de dados de localização.

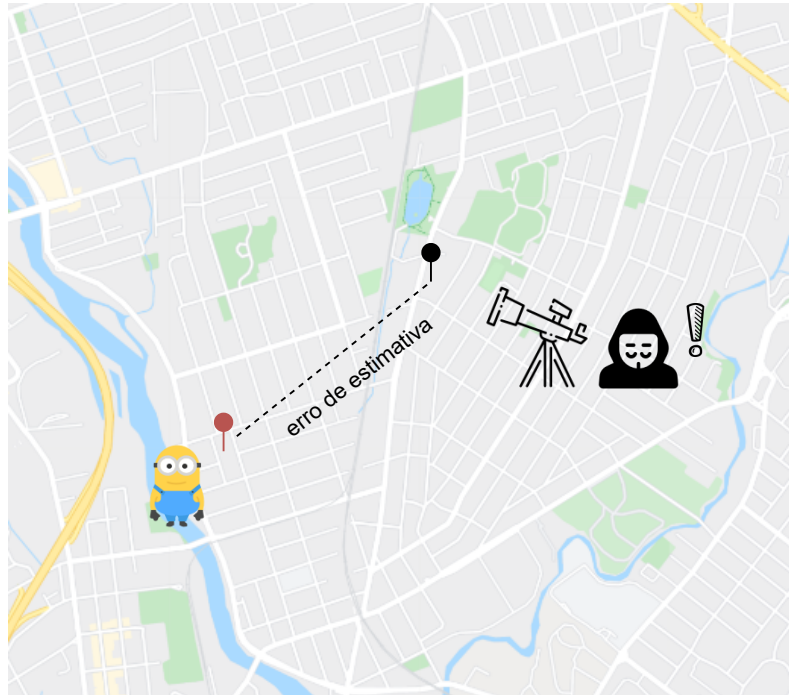


Figura 3.8. Erro de estimativa de um adversário.

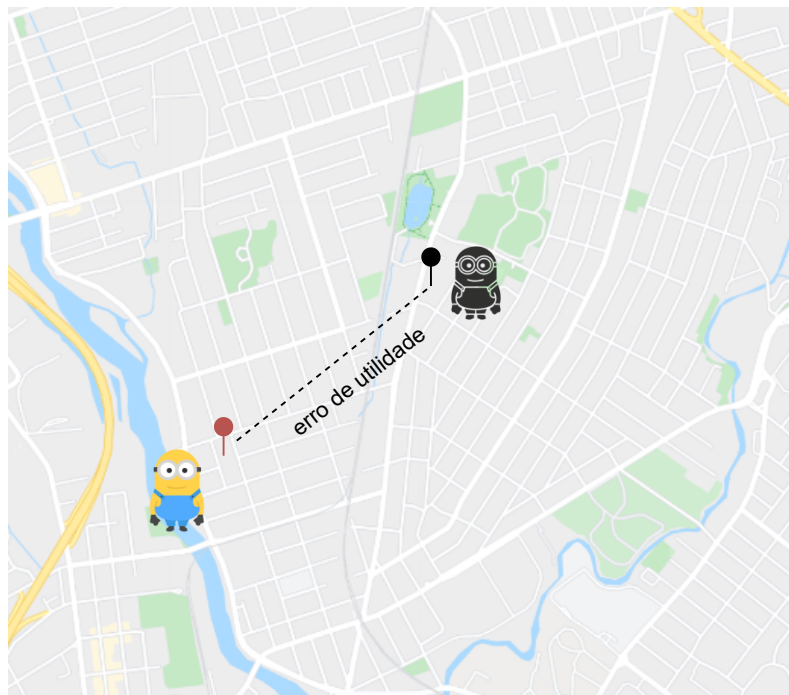


Figura 3.9. Erro de utilidade.

3.6.1. *k*-anonimato de localizações

Conforme discutido na Seção 3.3.1, o *k*-anonimato é um modelo de privacidade proposto por Sweeney et al. [46] com o objetivo de prevenir ataques de ligação ao registro. Através de técnicas de generalização ou supressão de dados, busca dificultar a re-identificação do indivíduo apagando os valores dos semi-identificadores ou diminuindo a sua precisão de tal sorte que o indivíduo tenha suas propriedades assemelhadas às propriedades de outros indivíduos. O *k*-anonimato busca garantir que para cada combinação de *k* atributos semi-identificadores, existem pelo menos *k* registros distintos no conjunto de dados publicado, formando uma classe de equivalência.

No contexto de privacidade de localização, um sujeito é tido *k*-anônimo se sua localização é indistinguível da localização de outros $k - 1$ usuários [22]. Portanto, a probabilidade de um usuário malicioso violar a privacidade de um indivíduo através de um ataque não será maior do que $\frac{1}{k}$.

O parâmetro *k* do modelo define o nível de privacidade requisitada. Assim, quanto maior o valor de *k*, maior será a privacidade dos dados e, conseqüentemente, menor é a probabilidade de se identificar o indivíduo. Entretanto, um valor alto de *k* tem um impacto no desempenho do algoritmo de anonimização, afetando a qualidade do serviço. Além disso, o processo de aproximar coordenadas geográficas pode deslocar virtualmente o solicitante para uma localização semelhante a outros *k* indivíduos, diminuindo a precisão do serviço. Desta forma, encontrar um equilíbrio entre privacidade e qualidade se faz ainda mais importante no contexto de dados de localização. Contudo, encontrar um valor ótimo para o parâmetro *k* é um problema NP-difícil, como citado na Seção 3.3.1. Desta forma, os responsáveis pela anonimização devem especificar o grau de privacidade desejada em função desse parâmetro.

O modelo tradicional da aplicação do *k*-anonimato em dados de localização requer a participação de uma entidade confiável responsável pelo processo de anonimização, o anonimizador. Dessa forma, quando um usuário necessita realizar uma requisição, enviando sua localização, o anonimizador calcula um conjunto de *k* usuários e reporta uma área de ofuscação contendo *k* localizações, incluindo a localização do usuário.

A Figura 3.10 ilustra uma abordagem da utilização dessa técnica, para um $k = 3$, onde o usuário solicitante deseja enviar uma requisição ao serviço de localização informando sua posição em azul. Aplicando o *k*-anonimato para $k = 3$, um terceiro confiável responsável por anonimizar a sua localização, agrupa a localização do usuário a outras $k - 1$ localizações, enviando uma requisição ao LBS contendo *k* localizações no total. Este, por sua vez, irá responder a requisição em função de cada uma das localizações enviadas. Como o terceiro confiável tem conhecimento das localizações dos usuários, ele irá filtrar a resposta referente a localização real presente na requisição e enviar o resultado para o usuário solicitante. Assim, para um atacante, a localização do usuário pode ser qualquer uma das *k* localizações que fazem parte da requisição, garantindo que a probabilidade de se identificar a localização do usuário não seja superior a $\frac{1}{k}$.

Uma desvantagem desse modelo é justamente a necessidade do anonimizador para realizar o processo de agrupamento de usuários. Dessa forma, se o atacante estiver agindo entre o anonimizador e o serviço, interceptando as requisições enviadas, estas estão ano-

nimizadas segundo o modelo k -anonimato. Entretanto, se o atacante estiver agindo entre o usuário e o anonimizador, a localização do usuário se encontra desprotegida.

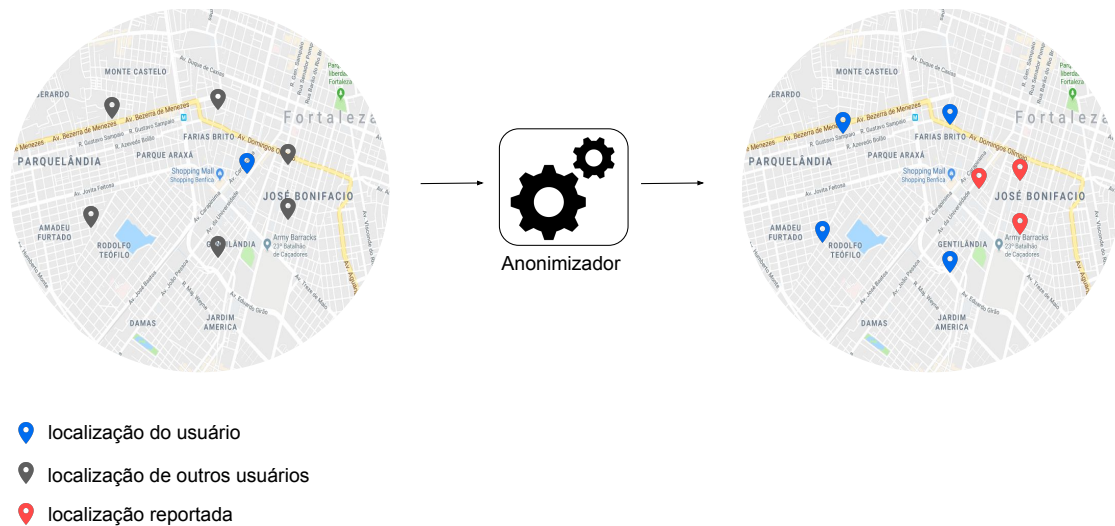


Figura 3.10. Processo de anonimização utilizando k -anonimato.

3.6.2. Zonas de mixagem

As zonas de mixagem procuram proteger a privacidade do usuário de ataques de ligação, ao evitar que seja possível vincular a identidade dos usuários à sua localização. Entretanto, diferentemente do k -anonimato, a zona de mixagem pode ser aplicada sem qualquer informação de identidade do usuário. O conceito de zona de mixagem foi proposto por Beresford et al. [7]. Ele propõe um *framework*, onde os usuários utilizam pseudo-ids que são modificados constantemente garantindo que estes não sejam identificados no uso de serviços de localização. Sendo assim, a real identidade do usuário é protegida através do uso de pseudônimos.

As zonas de mixagens são definidas como áreas circulares de raio r , onde nenhum dos usuários dentro da zona de mixagem possui qualquer registro de chamada de retorno ao serviço, ou seja, estão anônimos em relação ao serviço. Esta técnica procura garantir a indistinguibilidade dos usuários no uso do serviço dentro da zona, através do uso de pseudo-ids e a ausência das informações de localização de seu usuário. Assim, como o k -anonimato, a técnica de zona de mixagem em sua forma tradicional exige a figura de um terceiro confiável, responsável pela anonimização. Este terceiro confiável tem o papel de gerenciar os pseudo-ids dos usuários, garantindo que sempre que um usuário entre na zona de mixagem, ele possua um pseudo-id único que não tenha sido registrado por nenhuma aplicação. Dessa forma, como o serviço não recebe qualquer informação de localização dos usuários, sua identidade está "misturada" com a dos outros usuários dentro da zona. Qualquer informação de localização presente na zona de mixagem diz respeito à zona de uma forma geral, como por exemplo o seu ponto central.

A eficácia desta técnica depende de ajustar um tamanho adequado destas zonas de mixagem em função de seu raio r e também em função da quantidade mínima de indivíduos presentes. Se o valor de r for muito grande, impossibilitando que um indivíduo

percorra a distância que ela cobre e ingresse em outra zona em um único período de tempo, é possível que esta informação seja usada para identificar usuários que saem e entram em outras zonas, tornando incompleta a função de mixagem. Além disso, se a quantidade de usuários presentes em uma zona de mixagem for muito pequena, novamente ela não cumpre seu propósito de proteger a privacidade dos indivíduos presentes dentro da zona.

A Figura 3.11 ilustra a aplicação da técnica de zona de mixagem. A princípio, um usuário qualquer, cuja localização se encontra marcada em azul, ingressa na zona de mixagem 1 e obtém um pseudo-id que será usado na comunicação com o serviço, e só então passa a dispor do serviço coberto pela zona de mixagem. Ao se deslocar para a zona de mixagem 2, o usuário se comunica com o anonimizador e obtém um novo pseudo-id para esta zona. Requisições realizadas ao serviço não conterão qualquer informação de localização do usuário, uma vez que a única informação de localização utilizada é referente à zona de mixagem a qual ele serve.

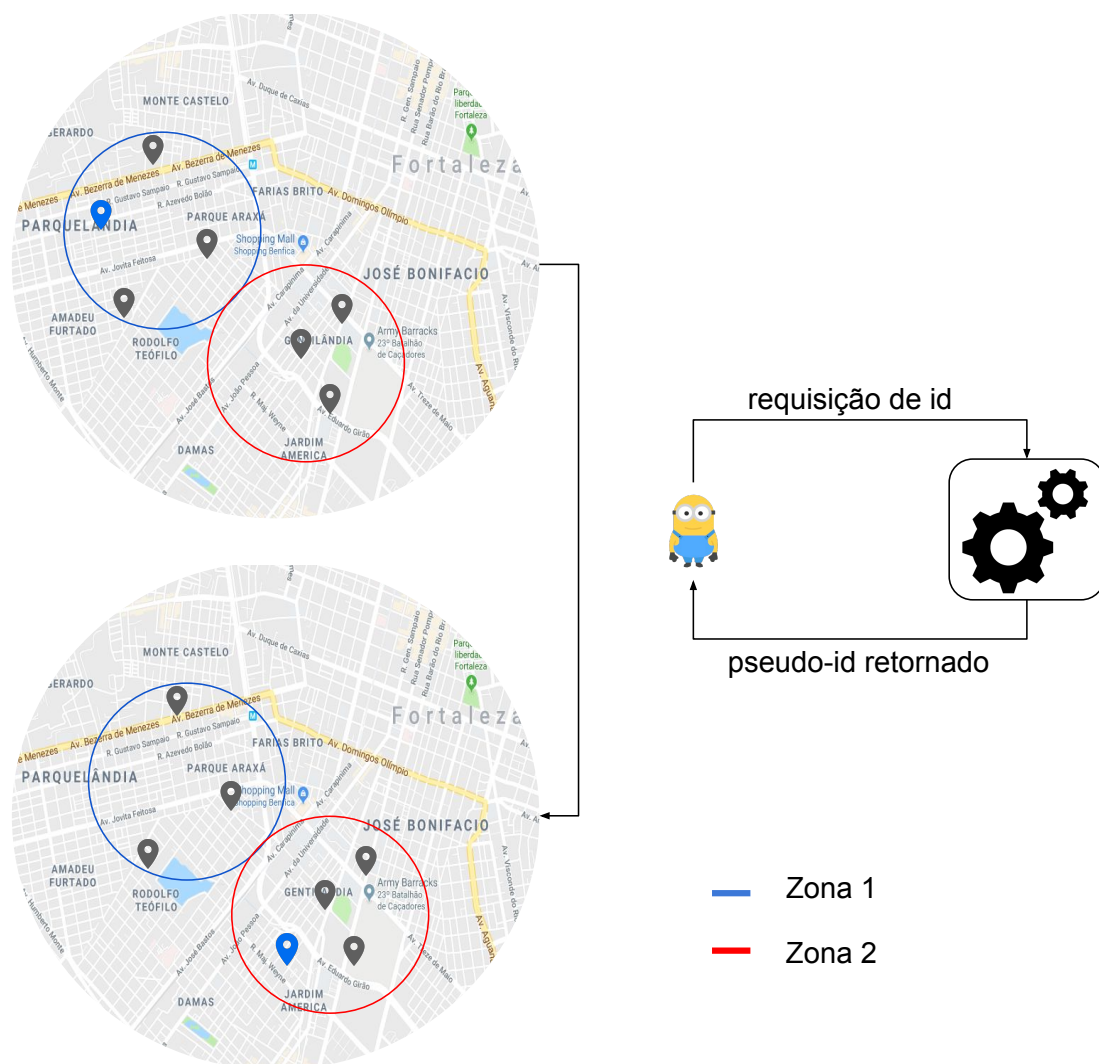


Figura 3.11. Processo de anonimização utilizando zona de mixagem.

3.6.3. Localizações falsas

Consultas realizadas a serviços de localização possuem dois componentes básicos: a informação de localização e o conteúdo da informação. Enquanto o primeiro contém qualquer informação relevante à localização, o segundo contém a requisição do serviço, por exemplo, a solicitação de dados de navegação em função da informação de localização informada [16].

A técnica de localizações falsas [25] procura garantir a privacidade dos dados de localização do usuário através de um processo de anonimização da *informação de localização* presente na requisição do usuário. Dessa forma, a informação de localização do usuário em uma requisição contém as informações da localização atual acrescida de localizações falsas, cujo objetivo é mascarar a localização verdadeira do usuário [16].

Diferentemente dos modelos tradicionais de k -anonimato e zona de mixagem, a técnica de seleção de localizações falsas não necessita da presença de servidores confiáveis para a realização do processo de anonimização, diminuindo o risco de exposição. Esse processo é realizado pelo cliente. A seleção de localizações falsas é realizada nos próprios dispositivos utilizados pelos usuários para consumir o serviço de localização.

O objetivo da presença de localizações falsas é garantir que a probabilidade de se identificar a localização real dentre aquelas presentes na requisição não seja maior que $\frac{1}{k}$, em que k é o grau de privacidade desejado para uma requisição com uma quantidade de k localizações presentes.



Figura 3.12. Técnica de Localizações Falsas.

A Figura 3.12 ilustra a aplicação da técnica de localizações falsas, onde o usuário, por intermédio de seu dispositivo, envia uma requisição anonimizada ao provedor de serviço. O processo de anonimização seleciona $k - 1$ localizações, na Figura 3.12 representado em preto, que serão enviadas na informação de localização da requisição, juntamente com a localização real do usuário, em laranja na figura. O servidor receberá a requisição contendo k localizações, e responderá a solicitação contida no conteúdo de

requisição, tendo como referência cada uma das localizações presentes na informação de localização. O dispositivo então filtra a resposta referente à localização real.

O mecanismo de seleção das localizações falsas é fundamental para garantir a privacidade de localização do usuário, uma vez que esta técnica, assim como as técnicas que abordam o modelo de anonimização estão sujeitas a ataques de conhecimento, que, conforme discutido na seção 3.5.4.3, utilizam de conhecimentos prévios para inferir dados sensíveis dos usuários. Desta forma, várias abordagens foram propostas, a fim de solucionar este problema e garantir uma seleção de localizações falsas que minimize o risco de exposição dos dados de localização do usuário, utilizando para isto, do próprio conhecimento prévio disponível ao provedor de serviço [45, 38].

Embora esta técnica garanta que a qualidade do serviço não terá qualquer impacto pela anonimização das informações de localização, já que não existe qualquer perda de utilidade da mesma, a estratégia de seleção de localizações, bem como o valor de k tem um impacto direto tanto na privacidade garantida, como no esforço computacional desta seleção. É importante que a estratégia de seleção considere fatores que elimine localizações improváveis que aumentariam a probabilidade de se identificar a real localização. Por exemplo, se for selecionado localizações pouco habitadas, ou de circulação proibida, juntamente com localizações bem frequentadas, estas últimas possuem uma probabilidade bem superior de ser a localização real do usuário dentre as outras. Além disso, se o valor de k for muito grande, pode ser bem custoso selecionar localizações que atendem aos requisitos de seleção do algoritmo utilizado.

3.6.4. Ofuscação de localização

As técnicas de ofuscação de localização procuram garantir a preservação de privacidade do usuário através da redução deliberada da precisão da localização do mesmo. Em sua abordagem tradicional apresentada por Ardagna *et al.* [4, 5], a informação de localização não mais apresenta as coordenadas da localização atual do usuário. Em vez disso, é enviado uma área circular, representada pela tupla $Area(r, x_c, y_c)$, centralizada nas coordenadas geográficas (x_c, y_c) e raio r , onde a probabilidade de a localização do usuário estar contida dentro dessa área é igual a 1. A Figura 3.13 ilustra uma requisição anonimizada pela técnica de ofuscação, onde o usuário, com localização em laranja, ao enviar sua requisição, envia como informação de localização uma área circular de raio r , centrada nas coordenadas do ponto central em preto, semelhante a todos os outros usuários presentes na área de ofuscação.

Uma abordagem alternativa é proposta por Gutscher [24], onde operações geométricas (*i.e.*, rotação, translação) são executadas sobre as coordenadas geográficas da localização atual, reduzindo sua precisão. A Figura 3.14 ilustra a aplicação desta técnica através da operação de translação sobre a localização atual do usuário.

Estas técnicas têm a vantagem de proteger a privacidade do usuário através da diminuição da precisão da localização apresentada, entretanto, este ganho de privacidade implica na perda de utilidade desta informação, que dependendo do serviço pode ter um impacto na qualidade do serviço, muitas vezes inviabilizando. Por exemplo, imagine que você realize uma requisição a um aplicativo de transporte e passe sua localização imprecisa com uns duzentos metros de imprecisão. Certamente, o motorista não vai lhe



Figura 3.13. Técnica de Ofuscação de Localização usando área circular.

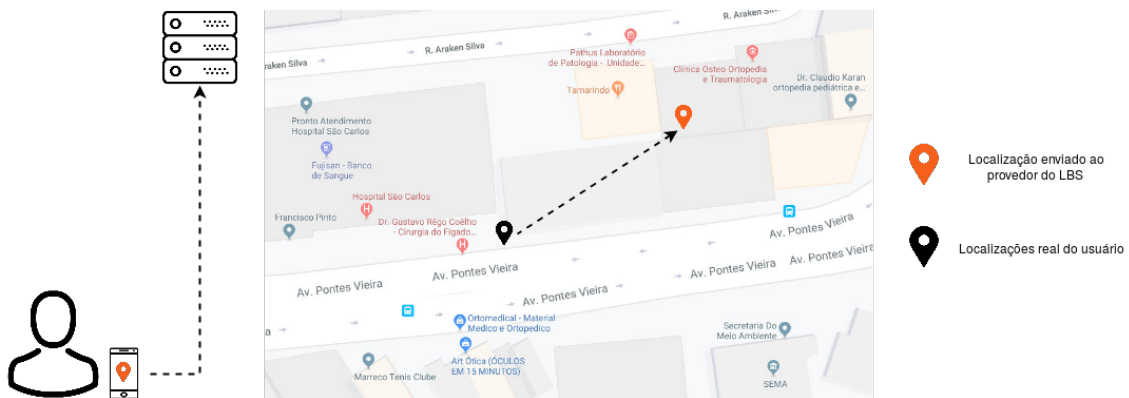


Figura 3.14. Técnica de Ofuscação de Localização usando operações geométricas.

encontrar. Entretanto, se você gostaria de utilizar um serviço buscando os restaurantes próximos dessa mesma localização imprecisa em duzentos metros, já será possível obter informação útil dessa requisição.

3.6.5. Geo-indistinguibilidade

Em se tratando de privacidade de localização, falta uniformidade quanto a conceitos básicos, como por exemplo, como quantificar a privacidade oferecida pelo serviço, ou qual a expectativa de privacidade do usuário no uso do serviço? O erro de estimativa de distância apresentado no início desta seção é uma métrica bastante usada para medir a privacidade de um mecanismo. Entretanto, essa noção de privacidade muitas vezes é definida em função do conhecimento adversário, que em geral é desconhecido.

A Geo-indistinguibilidade [3] é uma generalização da técnica de privacidade diferencial para o contexto de privacidade de localização. Desta forma, ela formaliza uma definição de privacidade de localização com sólidas garantias matemáticas que independe do conhecimento adversário. A ideia é garantir um nível de privacidade dentro de uma região geográfica de raio $r > 0$.

Antes de definir a noção de geo-indistinguibilidade é necessário definir o conceito de l -privacidade. Assim, podemos dizer que um indivíduo goza de l -privacidade dentro

de uma região de raio r , se quaisquer duas localizações a uma distância de no máximo r produzem distribuições similares, onde o nível de similaridade é definido em função de l . O parâmetro l representa o nível de privacidade dentro de r . Quanto maior for l , menor é o nível de privacidade. Para garantir uma maior utilidade do serviço, l precisa ser proporcional a r , $l = \epsilon r$. Assim, quanto maior for o raio, menor é o nível de privacidade garantido pelo mecanismo.

Um mecanismo K garante ϵ -geo-indistinguibilidade se para quaisquer duas localizações x e $x' \in X$:

$$D_p(K(x), K(x')) \leq \epsilon d(x, x'),$$

onde $D_p(K(x), K(x'))$ é a distância entre as distribuições produzidas por duas localizações x e x' pertencentes ao conjunto de todas as possíveis localizações X , e $d(x, x')$ é a distância euclidiana entre as localizações x e x' . Desta forma, se considerarmos a distância máxima r entre x e x' , a definição força que a distância máxima entre as distribuições seja ϵr . Dessa forma, o mecanismo gera uma variável aleatória seguindo uma distribuição, onde para uma localização x , será retornado uma localização x' com certa probabilidade. Falaremos com mais detalhes sobre mecanismo na Seção 3.6.5.1.

Aplicar o mecanismo sobre uma localização é bem direta. Uma vez definido o orçamento de privacidade ϵ , também chamado de *budget*, a região geográfica no entorno do usuário em função do raio r , basta aplicar o mecanismo e ele retornará uma localização que será usada na requisição do usuário. Entretanto, é possível que o usuário queira reportar mais de uma localização, como por exemplo, o conjunto de localizações frequentemente visitadas por ele. Nesse caso, a estratégia adotada de anonimização tem impacto no nível de privacidade garantido. Assim, considere o conjunto $x = \{x_1, x_2, \dots, x_n\}$ contendo n localizações que o usuário pretende anonimizar. Uma estratégia é anonimizar cada localização de forma independente. Dessa forma, para o conjunto x , o mecanismo reportará o conjunto $x' = \{x'_1, x'_2, \dots, x'_n\}$. A Figura 3.15 demonstra a aplicação do mecanismo sobre as localizações do usuário de forma independente. Para cada localização do usuário é aplicado o mecanismo e reportado uma localização com ruído adicionado em função de um mecanismo geo-indistinguível.

Neste caso, como cada localização anonimizada x'_i a partir de x_i tem uma probabilidade de ocorrer em função da variável aleatória gerada pelo mecanismo, a probabilidade combinada do conjunto anonimizado é dada pelo produto das probabilidades de cada localização presentes no conjunto reportado, i.e., $Pr[K(x) = x'] = \prod_i Pr[K_o(x_i) = x'_i]$. Como o mecanismo K permite uma combinação de n observações sobre as localizações de x , o nível de privacidade alcançado é de $n\epsilon$ -geo-indistinguibilidade.

Uma alternativa para este problema de escalabilidade quando o n é muito grande seria aplicar o mecanismo sobre alguma função agregadora. Ou seja, caso seja possível, em virtude da natureza do serviço, realizar uma requisição reportando alguma informação agregada sobre x , como por exemplo, o centroide do conjunto. Assim, seja f uma função de agregação que retorna uma localização \hat{x} que represente a informação agregada de x , a localização anonimizada reportada pelo mecanismo K pode ser dada por $K(f(\hat{x}))$. A Figura 3.16 demonstra a aplicação do mecanismo sobre o centroide do conjunto de localizações de um usuário, calculado pela função f . Em vermelho tem-se a localização anonimizada a partir deste centroide. Portanto, se f for Δ -sensível com respeito à d (distância

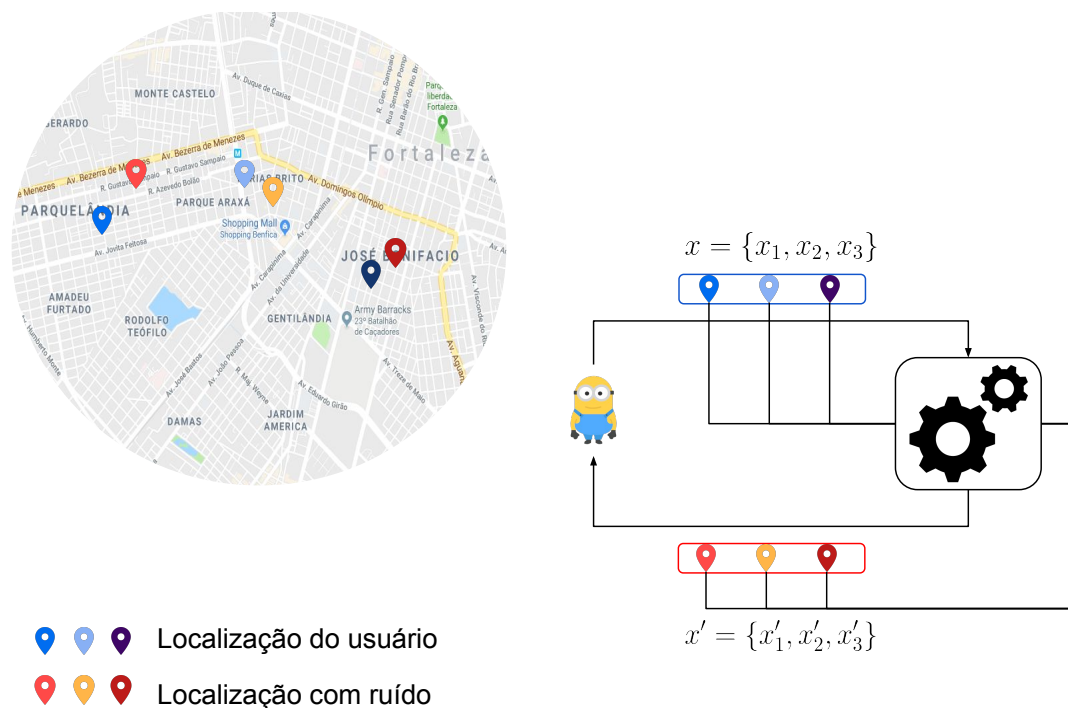


Figura 3.15. Localizações anonimadas usando ϵ -geo indistinguibilidade para cada localização de forma independente.

euclidiana entre dois pontos) e d_∞ (distância máxima entre as localizações dos conjuntos de pontos), o que significa que $d(f(x), f(x')) \leq \Delta d_\infty(x, x')$ para todos x e x' , e K satisfaz ϵ -geo-indistinguibilidade, então a composição $K \circ f$ satisfaz $\Delta\epsilon$ -geo-indistinguibilidade.

3.6.5.1. Mecanismo

O mecanismo é o responsável pela adição do ruído aleatório à localização. A ideia é que para qualquer localização x em um plano contínuo, o mecanismo irá reportar uma localização x' no mesmo plano de acordo com a função de aleatoriedade. Esta função deve garantir que a probabilidade de se reportar um ponto dentro de uma certa área no entorno de x' , quando as localizações atuais são x_0 ou x_1 respectivamente, se difere de no máximo $e^{\epsilon d(x_0, x_1)}$. Intuitivamente, o que se busca é que quanto menor for a distância entre esta localização na área de x' e a localização atual, maior é a probabilidade de esta localização ser reportada.

O mecanismo de Laplace Planar [3] consegue capturar esta distribuição desejada. Para um $\epsilon \in \mathbb{R}^+$, e a localização atual x , a função de densidade de probabilidade (PDF) do mecanismo de Laplace centralizada em x é dada por:

$$Laplace_\epsilon(x)(x') = \frac{\epsilon^2}{2\pi} e^{-\epsilon d(x, x')},$$

onde $\frac{\epsilon^2}{2\pi}$ é o fator de normalização. Portanto, o mecanismo K calcula uma localização x' que será reportado em função da distribuição de Laplace para a localização real x do

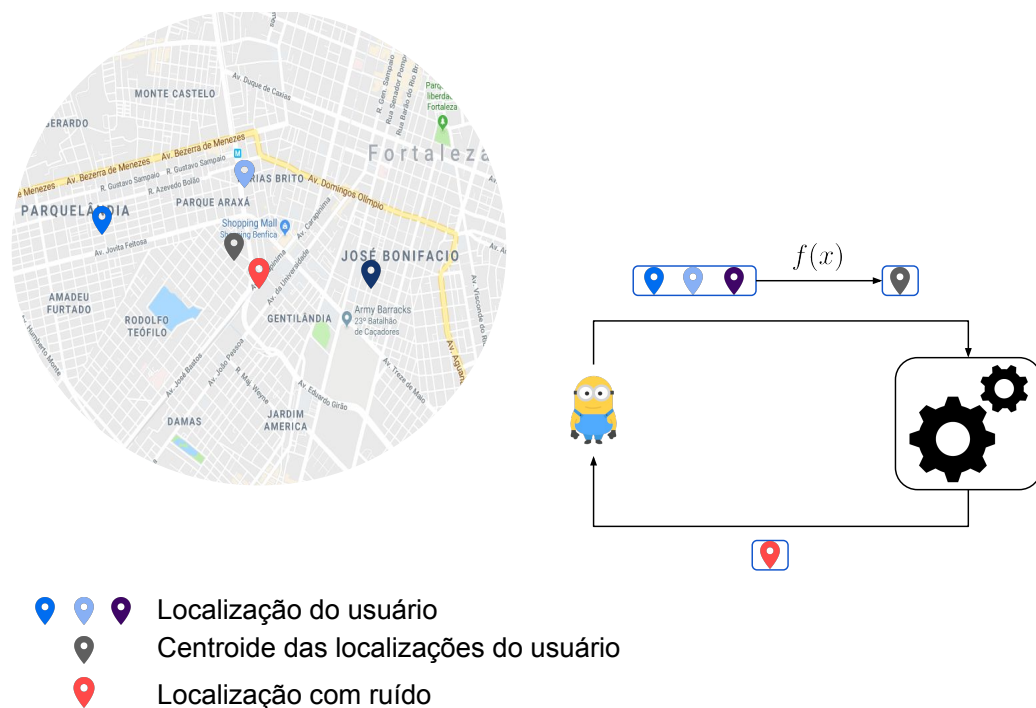


Figura 3.16. Localizações anonimadas usando ϵ -geo indistinguibilidade sobre a localização $f(x)$.

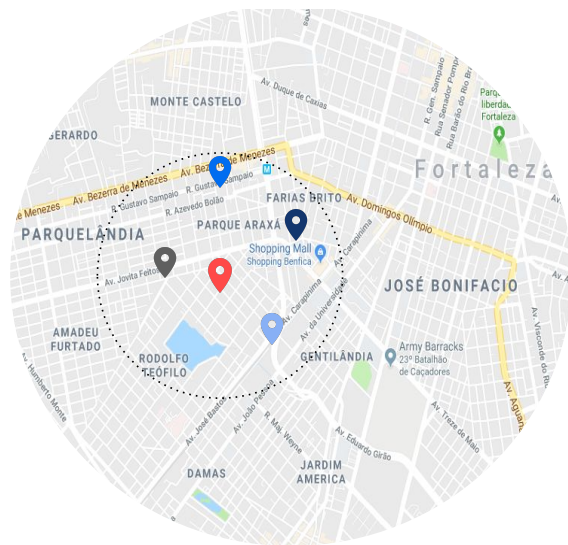
usuário.

A Figura 3.17 exemplifica esta situação. Em função da localização atual do usuário em vermelho na figura, o mecanismo apresenta uma probabilidade de reportar cada uma das localizações presentes. A localização reportada foi a localização em cinza, cuja probabilidade era de 12% segundo a tabela presente na figura. Observe, que as localizações que estão a uma distância mais próxima da localização atual possuem uma probabilidade maior de serem reportadas. Isto é decorrência da distribuição de Laplace utilizada.

3.6.6. Geo-indistinguibilidade Adaptativa

Buscando um melhor equilíbrio entre utilidade e privacidade, vários trabalhos se inspiraram na geo-indistinguibilidade [48, 11]. Entretanto, embora esta técnica tenha se mostrado promissora no controle da informação de localização publicada, ela não considera a potencial correlação existente entre as localizações em consultas contínuas ao serviço de localização. Conseqüentemente, um adversário pode explorar esta vulnerabilidade para reduzir o nível de privacidade da localização do usuário e obter uma melhor estimativa, violando os preceitos de privacidade garantidos pelo mecanismo.

Esta correlação é originada principalmente pela existência de padrões de movimento entre as diversas localizações que compõem as trajetórias dos usuários [2]. Como dito anteriormente, os serviços de localização necessitam que as informações de localização tenham um certo nível de utilidade para manter sua qualidade. Dessa forma, os mecanismos de privacidade precisam revelar parte das informações de localização, o que em geral, permite a existência dessa correlação mesmo que esta informação esteja anoni-



Localização	Prob.
	2,2%
	12%
	5,3%
	3%

Localização do usuário

Figura 3.17. Distribuição de probabilidade do mecanismo de Laplace planar.

mizada.

O estudo da correlação entre as localizações pode ser bastante útil, servindo para aumentar a eficiência de serviços, como por exemplo, na predição das próximas localizações visando uma melhoria no tempo de resposta de requisições. Entretanto, esta correlação pode também ser usada para reduzir o erro de estimativa de um adversário que busque identificar a localização dos usuários do serviço.

A geo-indistinguibilidade adaptativa procura apresentar um modelo adaptativo que seja capaz de preservar a privacidade apesar da correlação existente entre as localizações. Garantindo, que o nível de ruído possa ser ajustável em função da correlação existente entre as localizações reportadas. Para isso, o mecanismo utiliza uma janela de predição, que define o período de tempo de observação das localizações reportadas, a fim de se medir a correlação existente entre elas e a localização atual do usuário.

3.6.6.1. Modelo

O objetivo da geo-indistinguibilidade adaptativa é reportar localizações que sejam resistentes a ataques de análise de correlação, ou seja, ataques que busquem identificar a correlação existente entre as localizações. Estes tipos de ataques utilizam das localizações reportadas anteriormente para prever a próxima localização do usuário. O mecanismo adaptativo busca ajustar o nível de ruído em função dessa capacidade de se prever a próxima localização.

Mensurar o nível adequado de privacidade e utilidade em função da correlação é uma tarefa complexa e precisa ser adequada também ao tipo de serviço. Serviços com maior precisão, como serviços de navegação, apresentam uma expectativa de estimativa de erro baixo, já que as localizações reportadas deverão ter uma maior correlação entre si, já que existe uma necessidade maior de acurácia na localização reportada. Já um serviço com baixa precisão, como serviços de meteorologia, permite uma expectativa de estimativa de erro alta, já que as localizações reportadas deverão apresentar uma correlação baixa. Para melhor ajustar a quantidade de ruído o modelo adaptativo proposto define quatro parâmetros a ser ajustado em função da exigência do serviço:

- Δ_1 : Limite inferior de correlação;
- Δ_2 : Limite superior de correlação;
- α : fator de multiplicação do orçamento de privacidade para baixa correlação;
- β : fator de multiplicação do orçamento de privacidade para alta correlação;

O erro de estimativa pode ser classificado em alto, médio ou baixo. Um erro de estimativa alto indica que o nível de privacidade é alto, portanto, um ruído baixo pode ser utilizado para melhorar a qualidade do serviço sem afetar a privacidade. Um erro de estimativa baixo indica que o nível de privacidade está baixo, necessitando um maior ruído para garantir a privacidade. Um erro de estimativa médio indica que um ruído de nível mediano deve ser utilizado.

Seja $T = \{x_1, \dots, x_n\}$ a trajetória do usuário, e $x_i = (lat, lon, ts)$ uma localização com coordenadas de latitude lat e longitude lon em um momento de tempo ts . $OT = \{z_1, \dots, z_n\}$ é a trajetória ofuscada do usuário ao aplicar um mecanismo geo-indistinguível, como o mecanismo de Laplace planar apresentado na Seção 3.6.5.1 para um orçamento de privacidade ε_i . A localização x'_i representa a localização predita dada a localização x_i ao se aplicar uma simples regressão linear sobre as localizações ofuscadas em duas janelas de tamanho ws , $LatW$ e $LonW$. As janelas de predição $LatW$ e $LonW$ são usadas para prever as coordenadas de latitude e longitude de x'_i . O erro estimado err_i é calculado em função da distância euclidiana entre x_i e x'_i . Assim, a correlação em função do erro estimado err_i pode ser definida como:

- baixa: $err_i < \Delta_1$
- média: $\Delta_2 \geq err_i \geq \Delta_1$
- alta: $\Delta_2 < err_i$

A quantidade de ruído necessário é calculada em função da correlação existente. Se a correlação for baixa o ruído necessário é baixo, portanto, é aplicado um fator multiplicação α sobre o ε , assim, o orçamento utilizado pelo mecanismo será de $\varepsilon_i = \alpha\varepsilon$. Se o nível de correlação for alto, é aplicado um fator multiplicação β , e o orçamento utilizado será de $\varepsilon_i = \beta\varepsilon$. Se o nível for médio, então o orçamento utilizado é o próprio ε .

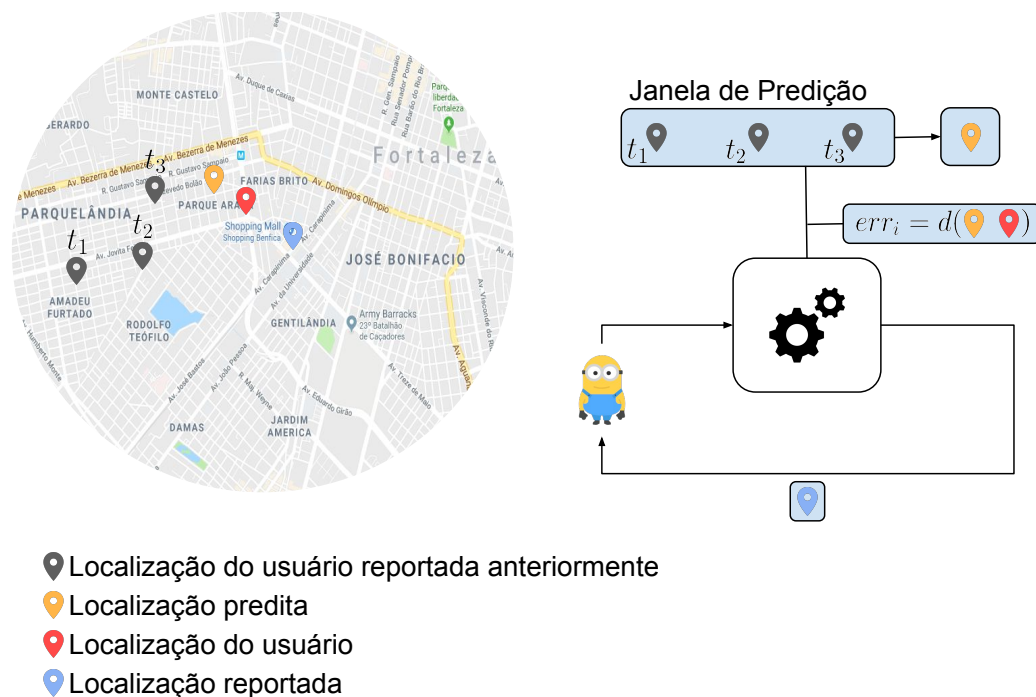


Figura 3.18. Geo-indistinguibilidade adaptativa.

Dessa forma, o nível de privacidade, definido em função do valor do *budget*, é adaptado de acordo com a capacidade de um atacante prever a próxima localização do usuário ao analisar a correlação existente entre as localizações presentes na janela de predição.

A Figura 3.18 exemplifica a aplicação do mecanismo adaptativo. Em cinza estão as localizações ofuscadas reportadas anteriormente. Aplicando uma regressão linear sobre janelas de tamanho 3, a localização estimada está em laranja. O erro estimado para a localização correta é medido entre a localização predita e a localização atual em vermelho. Supondo que o erro é baixo, indicando uma forte correlação, será necessário um ruído maior a fim de diminuir a correlação na próxima janela de predição. Dessa forma é aplicado o mecanismo de Laplace sobre a localização atual, com o orçamento igual a $\epsilon_i = \beta\epsilon$, reportando a localização ofuscada em azul.

3.7. Desafios de Pesquisa

Em privacidade de dados de forma geral, alcançar um equilíbrio entre privacidade e utilidade dos dados já é uma tarefa bem complexa. Ao longo deste curso falamos bastante neste tema, inclusive, sendo talvez um dos maiores motivadores no surgimento de novos trabalhos que buscam otimizar técnicas já existentes buscando garantir um valor ótimo para os parâmetros de privacidade, justamente para que não seja adicionado mais ruído do que necessário, e nem que seja subestimado esta quantidade.

Em serviços de localização, onde a precisão da localização do demandante do serviço impacta diretamente na qualidade do serviço prestado ao usuário. Esta análise do ruído necessário para que a informação de localização exposta seja o suficiente para

garantir a privacidade e a qualidade do serviço continua sendo uma tarefa de extrema complexidade, e depende diretamente do serviço utilizado e do nível de privacidade desejado. A privacidade diferencial em serviços de localização, dada a sua estratégia de adição de ruído aleatório, embora promissora ainda apresenta limitações quanto à utilidade dos dados fornecido ao provedor do serviço, principalmente pela dificuldade em se tratar os dados correlacionados. Desta forma vemos grande potencial de estudo no estudo da correlação de dados de localização e seu ajuste no uso de mecanismos diferencialmente privados.

3.8. Conclusão

Este capítulo conclui que a preservação da privacidade de dados acerca de indivíduos é um problema desafiador. Técnicas de anonimização têm sido utilizadas para a disponibilização de dados sensíveis, procurando encontrar o melhor balanceamento entre privacidade e utilidade que atenda às diversas partes envolvidas no processo de disponibilização de dados. Diferentes tipos de ataques à privacidade têm sido empregados por usuários maliciosos com a intenção de violar informações sensíveis de bases de dados abertas. Para tal fim, os atacantes utilizam conhecimento que muitas vezes é imensurável, devido aos diversos cenários em que informações podem ser obtidas. No contexto de dados de localização, este risco se potencializa, em virtude das informações agregadas ao dado geográfico buscado quando de uma solicitação a um serviço de localização, que servem de munição para os agentes maliciosos. Este capítulo apresentou as principais técnicas no estado da arte em preservação de privacidade de dados de localização. Os modelos de anonimização buscam proteger de ataques de ligação ao registro, ou seja, prevenir a vinculação entre a identidade do usuário e sua localização, evitando a re-identificação de indivíduos, geralmente utilizando técnicas de supressão e generalização. Os modelos de ofuscação, por sua vez, buscam proteger a localização em si, garantindo que esta não seja revelada, mesmo no uso de serviços de localização. A Privacidade Diferencial se destaca por fornecer soluções de preservação de privacidade, onde um ruído aleatório controlado é adicionado a localização do usuário, garantindo que a localização real do usuário estará protegida independentemente do conhecimento do atacante.

Finalmente, entendemos que o problema da garantia de privacidade de dados de localização dos usuários de serviços de localização continua cientificamente relevante. A busca por um ponto ideal na curva de solução de compromisso entre privacidade do indivíduo e a utilidade do dado fornecido para esse tipo de serviço deve pautar os próximos passos da pesquisa. Este aspecto é particularmente importante no contexto de localizações pois a qualidade do serviço é dependente da precisão do dado de localização, portanto o envio de dado perturbado para o provedor de serviço tende a impactar negativamente na qualidade. Tanto o paradigma de anonimização sintática, quanto o modelo de Privacidade Diferencial apresentam aspectos de revisão que devem ser vistos como oportunidades de pesquisas e desenvolvimento. Avanços em ambos os paradigmas são necessários para garantir que o futuro ofereça cada vez mais proteção à privacidade de indivíduos e ao mesmo tempo haja dados úteis e disponíveis para pesquisadores, testadores e analistas de dados.

Agradecimentos

Esta trabalho foi parcialmente financiada pela Lenovo, como parte do seu investimento em pesquisa e desenvolvimento de acordo com a Lei de Informática, pela CAPES (1836136), CNPq (122201/2018-3) e pelo LSB/D/UFC.

Referências

- [1] Aggarwal, C. C. and Philip, S. Y. (2008). A framework for condensation-based anonymization of string data. *Data Mining and Knowledge Discovery*, 16(3):251–275.
- [2] Al-Dhubhani, R. and Cazalas, J. M. (2018). An adaptive geo-indistinguishability mechanism for continuous lbs queries. *Wireless Networks*, 24(8):3221–3239.
- [3] Andrés, M. E., Bordenabe, N. E., Chatzikokolakis, K., and Palamidessi, C. (2013). Geo-indistinguishability: Differential privacy for location-based systems. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 901–914.
- [4] Ardagna, C. A., Cremonini, M., Damiani, E., De Capitani di Vimercati, S., and Samarati, P. (2007). Location privacy protection through obfuscation-based techniques. In Barker, S. and Ahn, G.-J., editors, *Data and Applications Security XXI*, pages 47–60.
- [5] Ardagna, C. A., Cremonini, M., De Capitani di Vimercati, S., and Samarati, P. (2011). An obfuscation-based approach for protecting location privacy. *IEEE Transactions on Dependable and Secure Computing*, 8(1):13–27.
- [6] Bayardo, R. J. and Agrawal, R. (2005). Data privacy through optimal k-anonymization. In *21st International conference on data engineering (ICDE'05)*, pages 217–228.
- [7] Beresford, A. R. and Stajano, F. (2003). Location privacy in pervasive computing. *IEEE Pervasive computing*, (1):46–55.
- [8] Blocki, J., Datta, A., and Bonneau, J. (2016). Differentially private password frequency lists. *IACR Cryptology ePrint Archive*, 2016:153.
- [9] BRITO, F. T. and MACHADO, J. C. (2017). Preservação de privacidade de dados: Fundamentos, técnicas e aplicações. *Jornadas de atualização em informática*, pages 91–130.
- [10] Brito, F. T. and Machado, J. C. (2017). Preservação de privacidade de dados: Fundamentos, técnicas e aplicações. In Delicato, F. C., Pires, P. F., and Silveira, I. F., editors, *Jornadas de A tualização em Informática 2017*. Sociedade Brasileira de Computação - SBC.
- [11] Chatzikokolakis, K., Palamidessi, C., and Stronati, M. (2015). Constructing elastic distinguishability metrics for location privacy. *arXiv preprint arXiv:1503.00756*.

- [12] Dewri, R., Ray, I., Ray, I., and Whitley, D. (2008). On the optimal selection of k in the k -anonymity problem. In *24th ICDE International Conference on Data Engineering*, pages 1364–1366, Cancun, Mexico.
- [13] Domingo-Ferrer, J., Sánchez, D., and Soria-Comas, J. (2016). Database anonymization: Privacy models, data utility, and microaggregation-based inter-model connections. *Synthesis Lectures on Information Security, Privacy, & Trust*, 8(1):1–136.
- [14] Domingo-Ferrer, J. and Soria-Comas, J. (2015). From t -closeness to differential privacy and vice versa in data anonymization. *Knowledge-Based Systems*, 74:151–158.
- [15] Domingo-Ferrer, J. and Torra, V. (2001). A quantitative comparison of disclosure control methods for microdata. *Confidentiality, disclosure and data access: theory and practical applications for statistical agencies*, pages 111–134.
- [16] Duarte Neto, E. R., Machado, J. C., and Mendonça, A. L. (2019). PrivLBS: Preserving privacy in location based services. *J. Inf. Data Manag.*, 10(2):81–96.
- [17] Dwork, C. (2006). Differential privacy. In *33rd International Colloquium*, pages 1–12.
- [18] Dwork, C. (2008). Differential privacy: A survey of results. In *International Conference on Theory and Applications of Models of Computation*, pages 1–19. Springer.
- [19] Dwork, C., Roth, A., et al. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3-4):211–407.
- [20] Fung, B. C., Wang, K., Fu, A. W.-C., and Yu, P. S. (2010a). *Introduction to Privacy-Preserving Data Publishing: Concepts and Techniques*. Chapman & Hall/CRC, 1st edition. ISBN 978-1-4200-9148-9.
- [21] Fung, B. C. M., Wang, K., Chen, R., and Yu, P. S. (2010b). Privacy-preserving data publishing: A survey of recent developments. *ACM Computing Surveys*, 42(4):1–53.
- [22] Gedik, B. and Liu, L. (2007). Protecting location privacy with personalized k -anonymity: Architecture and algorithms. *IEEE Transactions on Mobile Computing*, 7(1):1–18.
- [23] Goldreich, O. (2003). Cryptography and cryptographic protocols. *Distributed Computing*, 16(2-3):177–199.
- [24] Gutscher, A. (2006). Coordinate transformation - a solution for the privacy problem of location based services? In *Proceedings 20th IEEE International Parallel Distributed Processing Symposium*, pages 7 pp.–.
- [25] Kido, H., Yanagisawa, Y., and Satoh, T. (2005). An anonymous communication technique using dummies for location-based services. In *Proceedings of the Int. Conf. on Pervasive Services, ICPS'05*, pages 88–97. IEEE.

- [26] Lee, J. and Clifton, C. (2011). *How Much Is Enough? Choosing ϵ for Differential Privacy*, pages 325–340. Springer Berlin Heidelberg.
- [27] LeFevre, K., DeWitt, D. J., and Ramakrishnan, R. (2005). Incognito: Efficient full-domain k-anonymity. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*, pages 49–60. ACM.
- [28] Li, H., Sun, L., Zhu, H., Lu, X., and Cheng, X. (2014). Achieving privacy preservation in wifi fingerprint-based localization. In *INFOCOM, 2014 Proceedings IEEE*, pages 2337–2345. IEEE.
- [29] Li, N., Li, T., and Venkatasubramanian, S. (2007). t-closeness: Privacy beyond k-anonymity and l-diversity. In *2007 IEEE 23rd International Conference on Data Engineering*, pages 106–115. IEEE.
- [30] Liu, B., Zhou, W., Zhu, T., Gao, L., and Xiang, Y. (2018). Location privacy and its applications: A systematic study. *IEEE Access*, 6:17606–17624.
- [31] Machado, J., Neto, E. D., and Bento Filho, M. (2019). Técnicas de privacidade de dados de localização. *Sociedade Brasileira de Computação*.
- [32] Machanavajjhala, A., Kifer, D., Abowd, J. M., Gehrke, J., and Vilhuber, L. (2008). Privacy: Theory meets practice on the map. In *Proceedings of the 24th International Conference on Data Engineering, ICDE 2008, April 7-12, 2008, Cancún, México*, pages 277–286.
- [33] Matejka, J. and Fitzmaurice, G. (2017). Same stats, different graphs: generating datasets with varied appearance and identical statistics through simulated annealing. In *Proceedings of the 2017 CHI conference on human factors in computing systems*, pages 1290–1294.
- [34] McSherry, F. (2010). Privacy integrated queries: an extensible platform for privacy-preserving data analysis. *Commun. ACM*, 53(9):89–97.
- [35] McSherry, F. and Talwar, K. (2007). Mechanism design via differential privacy. In *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS'07)*, pages 94–103. IEEE.
- [36] Meyerson, A. and Williams, R. (2004). On the complexity of optimal k-anonymity. In *Proceedings of the twenty-third ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 223–228. ACM.
- [37] Mir, D. J., Isaacman, S., Cáceres, R., Martonosi, M., and Wright, R. N. (2013). DP-WHERE: differentially private modeling of human mobility. In *Proceedings of the 2013 IEEE International Conference on Big Data, 2013, Santa Clara, CA, USA*, pages 580–588.
- [38] Neto, E. R. D., Mendonça, A. L. C., Brito, F. T., and Machado, J. C. (2018). Privlbs: uma abordagem para preservação de privacidade de dados em serviços baseados em localização. In *Brazilian Symposium on Databases SBBD*, Rio de Janeiro, Brazil.

- [39] Nguyen, H. H., Kim, J., and Kim, Y. (2013). Differential privacy in practice. *Journal of Computing Science and Engineering*, 7(3):177–186.
- [40] NIU, B., Chen, Y., Wang, Z., Wang, B., Li, H., et al. (2020). Eclipse: Preserving differential location privacy against long-term observation attacks. *IEEE Transactions on Mobile Computing*.
- [41] Portela, T. T., Vicenzi, F., and Bogorny, V. (2019). Trajectory data privacy: Research challenges and opportunities. In *GEOINFO*, pages 99–110.
- [42] Primault, V., Boutet, A., Mokhtar, S. B., and Brunie, L. (2018). The long road to computational location privacy: A survey. *IEEE Communications Surveys & Tutorials*.
- [43] Schiller, J. and Voisard, A. (2004). *Location-based services*. Elsevier.
- [44] Shokri, R., Theodorakopoulos, G., Le Boudec, J.-Y., and Hubaux, J.-P. (2011). Quantifying location privacy. In *2011 IEEE symposium on security and privacy*, pages 247–262. IEEE.
- [45] Sun, G., Chang, V., Ramachandran, M., Sun, Z., Li, G., Yu, H., and Liao, D. (2017). Efficient location privacy algorithm for internet of things (iot) services and applications. *Journal of Network and Computer Applications*, 89:3–13.
- [46] Sweeney, L. (2002). k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(05):557–570.
- [47] Tan, V. Y. F. and Ng, S.-K. (2007). Generic probability density function reconstruction for randomization in privacy-preserving data mining. In *International Workshop on Machine Learning and Data Mining in Pattern Recognition*, pages 76–90. Springer.
- [48] Theodorakopoulos, G. (2015). The same-origin attack against location privacy. In *Proceedings of the 14th ACM Workshop on Privacy in the Electronic Society*, pages 49–53.
- [49] Truta, T. M., Campan, A., and Meyer, P. (2007). Generating microdata with p-sensitive k-anonymity property. In *Workshop on Secure Data Management*, pages 124–141. Springer.
- [50] Wang, K., Fung, B. C., and Yu, P. S. (2005). Template-based privacy preservation in classification problems. In *Fifth IEEE International Conference on Data Mining (ICDM'05)*, pages 8–pp. IEEE.
- [51] Wong, R. C.-W. and Fu, A. W.-C. (2010). Privacy-preserving data publishing: An overview. *Synthesis Lectures on Data Management*, 2(1):1–138.
- [52] Yang, D., Fang, X., and Xue, G. (2013). Truthful incentive mechanisms for k-anonymity location privacy. In *2013 Proceedings IEEE INFOCOM*, pages 2994–3002. IEEE.

- [53] Zheng, Y., Zhang, L., Xie, X., and Ma, W.-Y. (2009). Mining correlation between locations using human location history. In *Proceedings of the 17th ACM SIGSPATIAL international conference on advances in geographic information systems*, pages 472–475.
- [54] Zhu, X., Chi, H., Niu, B., Zhang, W., Li, Z., and Li, H. (2013). Mobicache: When k-anonymity meets cache. In *Global Communications Conference (GLOBECOM), 2013 IEEE*, pages 820–825, Atlanta, GA, USA. IEEE.