

Estudo de Casos sobre Privacidade e Transparência na Publicação de Dados

Gabriel Henrique Nunes

EVCOMP 2020

Departamento de Ciência da Computação
Universidade Federal de Minas Gerais

17 de Fevereiro de 2020



Laboratory of Information Security,
Cryptography, Privacy, and Transparency

Introdução

Transparência

Privacidade

Anonimização

Desidentificação e Pseudonimização

Métodos Determinísticos

Métodos Probabilísticos

Estado da Arte

Google, Microsoft, e Apple

US Census Bureau

Resumo



Laboratory of Information Security,
Cryptography, Privacy, and Transparency

Introdução

Transparência



A Importância da Publicação de Dados



The actual Enumeration shall be made within three Years after the first Meeting of the Congress of the United States, and within every subsequent Term of ten Years, in such Manner as they shall by Law direct.



A Importância da Publicação de Dados

Planejamento

| Políticas públicas



A Importância da Publicação de Dados

Planejamento

- | Políticas públicas
 - | Representação política



A Importância da Publicação de Dados

Planejamento

- | Políticas públicas
 - | Representação política
 - | Distribuição de fundos



A Importância da Publicação de Dados

Planejamento

- | Políticas públicas
 - | Representação política
 - | Distribuição de fundos
- | Investimentos privados



A Importância da Publicação de Dados

Planejamento

- | Políticas públicas
 - | Representação política
 - | Distribuição de fundos
- | Investimentos privados

Pesquisas científicas

- | Demografia
- | Economia

A Importância da Publicação de Dados

Planejamento

- | Políticas públicas
 - | Representação política
 - | Distribuição de fundos
- | Investimentos privados

Pesquisas científicas

- | Demografia
- | Economia

US Code Titles 13 & 14

A Importância da Publicação de Dados

Planejamento

- | Políticas públicas
 - | Representação política
 - | Distribuição de fundos
- | Investimentos privados

Pesquisas científicas

- | Demografia
- | Economia

Legislação específica
(1968)

Lei de Acesso à Informação
(2011)

A Importância da Publicação de Dados

Planejamento

- | Políticas públicas
 - | Representação política
 - | Distribuição de fundos
- | Investimentos privados

Pesquisas científicas

- | Demografia
- | Economia

Princípios Fundamentais das
Estatísticas Oficiais
(2014)

Introdução

Privacidade

O Problema da Publicação de Dados

Princípio 6. Os dados individuais coletados pelos órgãos estatísticos para a elaboração de estatísticas, sejam eles referentes a pessoas físicas ou jurídicas, devem ser estritamente confidenciais e utilizados exclusivamente para fins estatísticos.

O Problema da Publicação de Dados

Princípio 6. Os dados individuais coletados pelos órgãos estatísticos para a elaboração de estatísticas, sejam eles referentes a pessoas físicas ou jurídicas, devem ser estritamente confidenciais e utilizados exclusivamente para fins estatísticos.

O Problema da Publicação de Dados

A Importância da Con dencialidade

- | Precisão dos dados coletados

O Problema da Publicação de Dados

A Importância da Con dencialidade

- | Precisão dos dados coletados
- | Justiça (Fairness)

O Problema da Publicação de Dados

A Importância da Con dencialidade

- | Precisão dos dados coletados
- | Justiça (Fairness)
- | Direitos fundamentais

O Problema da Publicação de Dados

A Importância da Con dencialidade

- | Precisão dos dados coletados
- | Justiça (Fairness)
- | Direitos fundamentais
- | Democracia

O Problema da Publicação de Dados

A Importância da Con dencialidade

- | Precisão dos dados coletados
- | Justiça (Fairness)
- | Direitos fundamentais
- | Democracia

Entretanto, garanti-la não é simples

O Problema da Publicação de Dados

A Importância da Con dencialidade

- | Precisão dos dados coletados
- | Justiça (Fairness)
- | Direitos fundamentais
- | Democracia

Entretanto, garanti-la não é simples

- | Dados sensíveis

O Problema da Publicação de Dados

A Importância da Con dencialidade

- | Precisão dos dados coletados
- | Justiça (Fairness)
- | Direitos fundamentais
- | Democracia

Entretanto, garanti-la não é simples

- | Dados sensíveis
- | Quasi-identi cadores

Con dencialidade e Privacidade

Con dencialidade

Envolve um conjunto de regras ou uma promessa que limitam o acesso ou impõem restrições a certos tipos de informações

Con dencialidade e Privacidade

Con dencialidade

Envolve um conjunto de regras ou uma promessa que limitam o acesso ou impõem restrições a certos tipos de informações

- I Usualmente utilizado no relacionamento com médicos, advogados, ou instituições nanceiras

Con dencialidade e Privacidade

Con dencialidade

Envolve um conjunto de regras ou uma promessa que limitam o acesso ou impõem restrições a certos tipos de informações

- I Usualmente utilizado no relacionamento com médicos, advogados, ou instituições nanceiras

Privacidade

Capacidade de isolar a si mesmo ou de isolar informações sobre si mesmo e, assim, expressar-se de maneira seletiva

Con dencialidade e Privacidade

Con dencialidade

Envolve um conjunto de regras ou uma promessa que limitam o acesso ou impõem restrições a certos tipos de informações

- | Usualmente utilizado no relacionamento com médicos, advogados, ou instituições nanceiras

Privacidade

Capacidade de isolar a si mesmo ou de isolar informações sobre si mesmo e, assim, expressar-se de maneira seletiva

- | Utilizado na Legislação mais recente sobre dados pessoais

Lei Geral de Proteção de Dados Pessoais (LGPD)

Dados pessoais são coletados:

- | por governos

Lei Geral de Proteção de Dados Pessoais (LGPD)

Dados pessoais são coletados:

- | por governos
- | por empresas privadas

Lei Geral de Proteção de Dados Pessoais (LGPD)

Dados pessoais são coletados:

- | por governos
- | por empresas privadas
- | para publicação

Lei Geral de Proteção de Dados Pessoais (LGPD)

Dados pessoais são coletados:

- | por governos
- | por empresas privadas
- | para publicação
- | para uso interno

Lei Geral de Proteção de Dados Pessoais (LGPD)

Dados pessoais são coletados:

- | por governos
- | por empresas privadas
- | para publicação
- | para uso interno

LGPD

- | Entra em vigor em Agosto de 2020

Lei Geral de Proteção de Dados Pessoais (LGPD)

Dados pessoais são coletados:

- | por governos
- | por empresas privadas
- | para publicação
- | para uso interno

LGPD

- | Entra em vigor em Agosto de 2020
- | O que são dados pessoais?

Lei Geral de Proteção de Dados Pessoais (LGPD)

Dados pessoais são coletados:

- | por governos
- | por empresas privadas
- | para publicação
- | para uso interno

LGPD

- | Entra em vigor em Agosto de 2020
- | O que são dados pessoais?



Possíveis Soluções

Possíveis Soluções

O problema é forma de acesso aos dados. Eles deveriam ser anônimos. A cada pessoa se atribui um código, um ID, e assim as pessoas que lidam com o banco de dados não têm acesso às informações em nível pessoal. Criptografado. Assim se evita invasão de privacidade, stalker, etc

- Usuário do Twitter



Possíveis Soluções

Possíveis Soluções

Controle de acesso

Possíveis Soluções

Controle de acesso

- | Evita acesso não autorizado

Possíveis Soluções

Controle de acesso

- | Evita acesso não autorizado
- | Evita vazamento de informações

Possíveis Soluções

- Controle de acesso e Criptografia
 - | Evitam acesso não autorizado
 - | Evitam vazamento de informações

Possíveis Soluções

- Controle de acesso e Criptogra a
- | Evitam acesso não autorizado
 - | Evitam vazamento de informações
 - | Não protegem contra inferência quando os dados são acessíveis

Possíveis Soluções

Controle de acesso e Criptografia

- | Evitam acesso não autorizado
- | Evitam vazamento de informações
- | Não protegem contra inferência quando os dados são acessíveis

Anonimização

- | Diversos métodos propostos

Possíveis Soluções

Controle de acesso e Criptografia

- | Evitam acesso não autorizado
- | Evitam vazamento de informações
- | Não protegem contra inferência quando os dados são acessíveis

Anonimização

- | Diversos métodos propostos
- | Diversos vetores de ataque já demonstrados

Anonimização

Desidentificação e Pseudonimização

Desidentificação ou Anonimização "inocente"

Definição

Uma base de dados é dita anônima se dela foram retirados os identificadores diretos dos indivíduos

Desidenti cação ou Anonimização "inocente"

Definição

Uma base de dados é dita anônima se dela foram retirados os identificadores diretos dos indivíduos

Exemplos

Nome, números únicos de identificação (CPF, RG), endereço, etc

Desidentificação ou Anonimização "inocente"

Definição

Uma base de dados é dita anônima se dela foram retirados os identificadores diretos dos indivíduos

Exemplos

Nome, números únicos de identificação (CPF, RG), endereço, etc
Entretanto, a desidentificação de identificadores diretos é muito exíbil

Exemplo de Desidentificação

	Nome	Idade	Condição
1	Jon Snow	30	Resfriado
2	Jamie Lannister	39	Mão amputada
3	Arya Stark	16	Dor de estômago
4	Bran Stark	14	Paraplegia
5	Eddad Stark	32	Dor de cabeça
6	Ramsay Bolton	32	Psicopatia
7	Daenerys Targaryen	25	Mania de grandeza

Catuscia Palamidessi, Kostas Chatzikokolakis

Exemplo de Desidenti cação

	Nome	Idade	Condição
1	*	30	Resfriado
2	*	39	Mão amputada
3	*	16	Dor de estômago
4	*	14	Paraplegia
5	*	32	Dor de cabeça
6	*	32	Psicopatia
7	*	25	Mania de grandeza

Quasi-identificadores

Definição

Atributos que podem ser vinculados a informações externas para identificar indivíduos unicamente

Quasi-identificadores

Definição

Atributos que podem ser vinculados a informações externas para identificar indivíduos unicamente

Observação

Definir o conjunto de quasi-identificadores para uma dada publicação ainda é uma questão em aberto

Sweeney's Linkage Attack (1998)

DB1: Contém
dados sensíveis
(anonimizada)

Sweeney's Linkage Attack (1998)

DB1: Contém
dados sensíveis
(anonimizada)

DB2: Coleção
pública de dados
não sensíveis

Sweeney's Linkage Attack (1998)

DB1: Contém
dados sensíveis
(anonimizada)

DB2: Coleção
pública de dados
não sensíveis

Informações
auxiliares

Sweeney's Linkage Attack (1998)

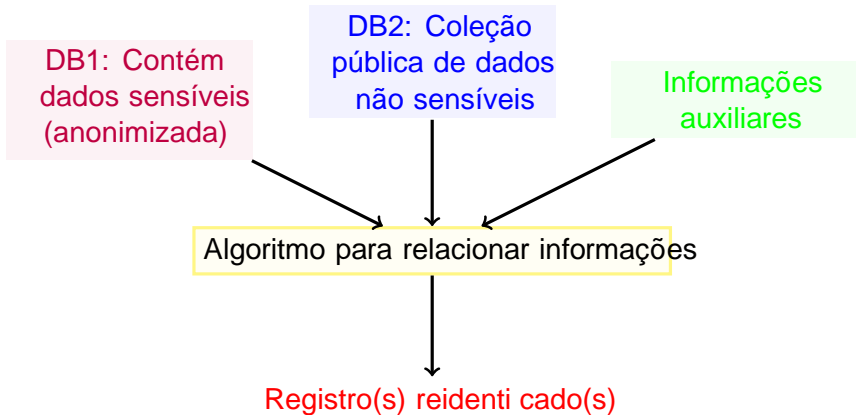
DB1: Contém
dados sensíveis
(anonimizada)

DB2: Coleção
pública de dados
não sensíveis

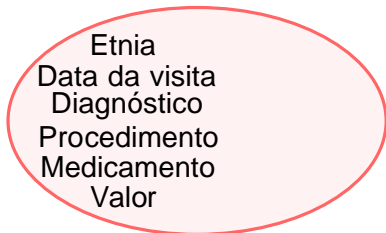
Informações
auxiliares

Algoritmo para relacionar informações

Sweeney's Linkage Attack (1998)

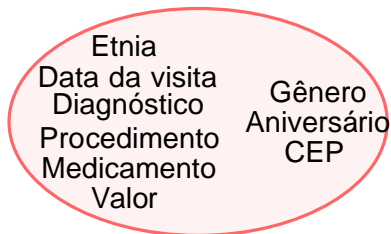


Sweeney's Linkage Attack na Prática



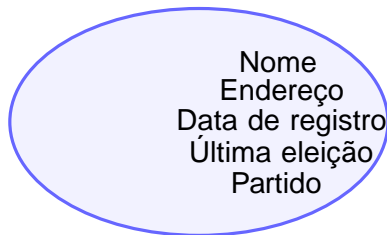
Base 1: Dados médicos

Sweeney's Linkage Attack na Prática



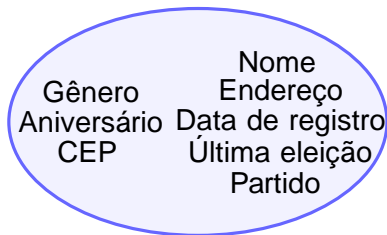
Base 1: Dados médicos

Sweeney's Linkage Attack na Prática



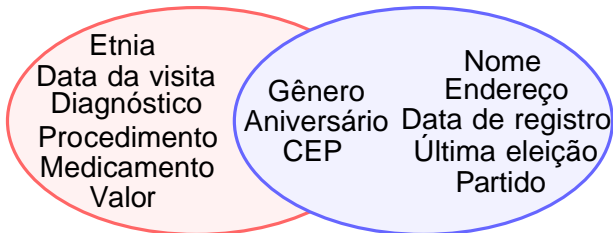
Base 2: Lista de eleitores

Sweeney's Linkage Attack na Prática



Base 2: Lista de eleitores

Sweeney's Linkage Attack na Prática



Base 1: Dados médicos

Base 2: Lista de eleitores

Sweeney's Linkage Attack na Prática

Conclusões

- I Com apenas gênero, data de aniversário, e CEP, Sweeney foi capaz de reidentificar o governador de Massachusetts

Sweeney's Linkage Attack na Prática

Conclusões

- | Com apenas gênero, data de aniversário, e CEP, Sweeney foi capaz de reidentificar o governador de Massachusetts
- | O mesmo foi possível para 67% da população dos EUA, também unicamente identificável pelos mesmos atributos

Pseudonimização

Definição

Uma base de dados é dita pseudonimizada se dela foram retirados os identificadores diretos dos indivíduos em substituição por um código único para cada indivíduo

Pseudonimização

Definição

Uma base de dados é dita pseudonimizada se dela foram retirados os identificadores diretos dos indivíduos em substituição por um código único para cada indivíduo

Entretanto...

- | Apresenta os mesmos problemas que a Desidentificação

Pseudonimização

De nição

Uma base de dados é dita pseudonimizada se dela foram retirados os identi cadores diretos dos indivíduos em substituição por um código único para cada indivíduo

Entretanto...

- | Apresenta os mesmos problemas que a Desidenti cação
- | Publicação de dados longitudinais no tempo

Pseudonimização

De nição

Uma base de dados é dita pseudonimizada se dela foram retirados os identi cadores diretos dos indivíduos em substituição por um código único para cada indivíduo

Entretanto...

- | Apresenta os mesmos problemas que a Desidenti cação
- | Publicação de dados longitudinais no tempo
 - | Pseudonimização é ainda pior do que Desidenti cação

Pseudonimização

De nição

Uma base de dados é dita pseudonimizada se dela foram retirados os identi cadores diretos dos indivíduos em substituição por um código único para cada indivíduo

Entretanto...

- | Apresenta os mesmos problemas que a Desidenti cação
- | Publicação de dados longitudinais no tempo
 - | Pseudonimização é ainda pior do que Desidenti cação
 - | Facilita a reidenti cação ao longo do tempo



Exemplo de Pseudonimização

Exemplo de Pseudonimização

- | Divulgou 20 milhões de pesquisas realizadas por mais de 650 mil usuários ao longo de três meses

Exemplo de Pseudonimização

- | Divulgou 20 milhões de pesquisas realizadas por mais de 650 mil usuários ao longo de três meses
- | Objetivo era fomentar pesquisas científicas

Exemplo de Pseudonimização

- | Divulgou 20 milhões de pesquisas realizadas por mais de 650 mil usuários ao longo de três meses
- | Objetivo era fomentar pesquisas científicas
- | Os dados foram liberados para toda a Internet

Exemplo de Pseudonimização

- | Divulgou 20 milhões de pesquisas realizadas por mais de 650 mil usuários ao longo de três meses
- | Objetivo era fomentar pesquisas científicas
- | Os dados foram liberados para toda a Internet
- | The New York Times cruzou as pesquisas de um usuário com listas telefônicas

Exemplo de Pseudonimização

- | Divulgou 20 milhões de pesquisas realizadas por mais de 650 mil usuários ao longo de três meses
- | Objetivo era fomentar pesquisas científicas
- | Os dados foram liberados para toda a Internet
- | The New York Times cruzou as pesquisas de um usuário com listas telefônicas
 - | 60 anos de idade , homens solteiros , cão que urina em tudo , paisagistas em Lilburn, GA

Exemplo de Pseudonimização

- | Divulgou 20 milhões de pesquisas realizadas por mais de 650 mil usuários ao longo de três meses
- | Objetivo era fomentar pesquisas científicas
- | Os dados foram liberados para toda a Internet
- | The New York Times cruzou as pesquisas de um usuário com listas telefônicas
 - | 60 anos de idade , homens solteiros , cão que urina em tudo , paisagistas em Lilburn, GA
 - | Reidentificou Thelma Arnold, viúva de 62 anos, com três cães, de Lilburn, Geórgia

Anonimização

Métodos Determinísticos

k-Anonimização (Sweeney & Samarati, 1998)

Definição

Uma base de dados é dita Anônima se cada registro for indistinguível de $k - 1$ outros registros, considerando-se os quase-identificadores, através de:

k-Anonimização (Sweeney & Samarati, 1998)

Definição

Uma base de dados é dita Anônima se cada registro for indistinguível de $k - 1$ outros registros, considerando-se os quase-identificadores, através de:

- | generalização de atributos

k-Anonimização (Sweeney & Samarati, 1998)

Definição

Uma base de dados é dita Anônima se cada registro for indistinguível de $k - 1$ outros registros, considerando-se os quase-identificadores, através de:

- | generalização de atributos
- | supressão de atributos

k-Anonimização (Sweeney & Samarati, 1998)

Definição

Uma base de dados é dita Anônima se cada registro for indistinguível de $k - 1$ outros registros, considerando-se os quase-identificadores, através de:

- | generalização de atributos
- | supressão de atributos
- | adição de registros sintéticos

Exemplo de 4-Anonimização

	Não Sensível			Sensível Condição
	CEP	Idade	País	
1	13053	28	Rússia	Cardíaco
2	13068	29	EUA	Cardíaco
3	13068	21	Japão	Virose
4	13053	23	EUA	Virose
5	14853	50	Índia	Câncer
6	14853	55	Rússia	Cardíaco
7	14850	47	EUA	Virose
8	14850	49	EUA	Virose
9	13053	31	EUA	Câncer
10	13053	37	Índia	Câncer
11	13068	36	Japão	Câncer
12	13068	35	EUA	Câncer

	Não Sensível			Sensível Condição
	CEP	Idade	País	
1	130**	< 30	*	Cardíaco
2	130**	< 30	*	Cardíaco
3	130**	< 30	*	Virose
4	130**	< 30	*	Virose
5	1485*	40	*	Câncer
6	1485*	40	*	Cardíaco
7	1485*	40	*	Virose
8	1485*	40	*	Virose
9	130**	3*	*	Câncer
10	130**	3*	*	Câncer
11	130**	3*	*	Câncer
12	130**	3*	*	Câncer

Exemplo de 4-Anonimização

	Não Sensível			Sensível Condição
	CEP	Idade	País	
1	130**	< 30	*	Cardíaco
2	130**	< 30	*	Cardíaco
3	130**	< 30	*	Virose
4	130**	< 30	*	Virose
5	1485*	40	*	Câncer
6	1485*	40	*	Cardíaco
7	1485*	40	*	Virose
8	1485*	40	*	Virose
9	130**	3*	*	Câncer
10	130**	3*	*	Câncer
11	130**	3*	*	Câncer
12	130**	3*	*	Câncer

Exemplo de 4-Anonimização

	Não Sensível			Sensível Condição
	CEP	Idade	País	
1	130**	< 30	*	Cardíaco
2	130**	< 30	*	Cardíaco
3	130**	< 30	*	Virose
4	130**	< 30	*	Virose
5	1485*	40	*	Câncer
6	1485*	40	*	Cardíaco
7	1485*	40	*	Virose
8	1485*	40	*	Virose
9	130**	3*	*	Câncer
10	130**	3*	*	Câncer
11	130**	3*	*	Câncer
12	130**	3*	*	Câncer

Exemplo de 4-Anonimização

	Não Sensível			Sensível Condição
	CEP	Idade	País	
1	130**	< 30	*	Cardíaco
2	130**	< 30	*	Cardíaco
3	130**	< 30	*	Virose
4	130**	< 30	*	Virose
5	1485*	40	*	Câncer
6	1485*	40	*	Cardíaco
7	1485*	40	*	Virose
8	1485*	40	*	Virose
9	130**	3*	*	Câncer
10	130**	3*	*	Câncer
11	130**	3*	*	Câncer
12	130**	3*	*	Câncer

`-Diversidade (Kifer et al., 2007)

Definição

Uma base de dados é dita Diversa se cada agrupamento de registros apresentar uma diversidade de, ao menos, atributos sensíveis

Exemplo de 3-Diversidade

	Não Sensível			Sensível Condição
	CEP	Idade	Gênero	
1	130**	50	*	Câncer
2	130**	50	*	Câncer
3	130**	50	*	Hemo lia
4	130**	50	*	Virose
5	122**	> 50	*	Hemo lia
6	122**	> 50	*	Câncer
7	122**	> 50	*	Virose
8	122**	> 50	*	Virose
9	130**	50	*	Câncer
10	130**	50	*	Câncer
11	130**	50	*	Hemo lia
12	130**	50	*	Virose

Exemplo de 3-Diversidade

	Não Sensível			Sensível Condição
	CEP	Idade	Gênero	
1	130**	50	*	Câncer
2	130**	50	*	Câncer
3	130**	50	*	Hemo lia
4	130**	50	*	Virose
5	122**	> 50	*	Hemo lia
6	122**	> 50	*	Câncer
7	122**	> 50	*	Virose
8	122**	> 50	*	Virose
9	130**	50	*	Câncer
10	130**	50	*	Câncer
11	130**	50	*	Hemo lia
12	130**	50	*	Virose

Exemplo de 3-Diversidade

	Não Sensível			Sensível Condição
	CEP	Idade	Gênero	
1	130**	50	*	Câncer
2	130**	50	*	Câncer
3	130**	50	*	Hemo lia
4	130**	50	*	Virose
5	122**	> 50	*	Hemo lia
6	122**	> 50	*	Câncer
7	122**	> 50	*	Virose
8	122**	> 50	*	Virose
9	130**	50	*	Câncer
10	130**	50	*	Câncer
11	130**	50	*	Hemo lia
12	130**	50	*	Virose

Exemplo de 3-Diversidade

	Não Sensível			Sensível Condição
	CEP	Idade	Gênero	
1	130**	50	*	Câncer
2	130**	50	*	Câncer
3	130**	50	*	Hemo lia
4	130**	50	*	Virose
5	122**	> 50	*	Hemo lia
6	122**	> 50	*	Câncer
7	122**	> 50	*	Virose
8	122**	> 50	*	Virose
9	130**	50	*	Câncer
10	130**	50	*	Câncer
11	130**	50	*	Hemo lia
12	130**	50	*	Virose



Limitações de Diversidade

- | Pode ser muito rigorosa

Limitações de Diversidade

- | Pode ser muito rigorosa
 - | A sensibilidade dos atributos pode variar

Limitações de Diversidade

- | Pode ser muito rigorosa
 - | A sensibilidade dos atributos pode variar
- | Pode não ser su ciente

Limitações de Diversidade

- | Pode ser muito rigorosa
 - | A sensibilidade dos atributos pode variar
- | Pode não ser su ciente
 - | Permite a inferência de atributos sensíveis, e em alguns casos com elevada probabilidade

Limitações de Diversidade

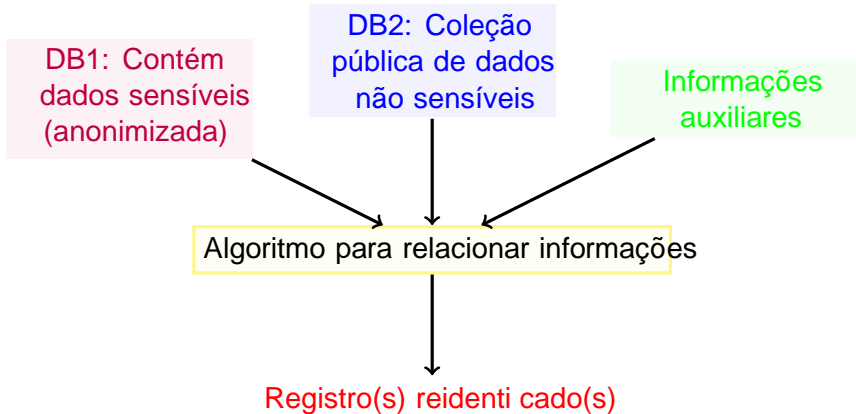
- | Pode ser muito rigorosa
 - | A sensibilidade dos atributos pode variar
- | Pode não ser suficiente
 - | Permite a inferência de atributos sensíveis, e em alguns casos com elevada probabilidade

t-Proximidade (Li et al., 2007)

Definição

Uma base de dados é dita próxima se cada agrupamento de registros apresentar uma distribuição de atributos sensíveis próxima, a no máximo uma distância, da distribuição geral

Composicionalidade



Composicionalidade

	Não Sensível		País	Sensível Condição
	CEP	Idade		
1	130**	< 30	*	AIDS
2	130**	< 30	*	Cardíaco
3	130**	< 30	*	Virose
4	130**	< 30	*	Virose
5	130**	40	*	Câncer
6	130**	40	*	Cardíaco
7	130**	40	*	Virose
8	130**	40	*	Virose
9	130**	3*	*	Câncer
10	130**	3*	*	Câncer
11	130**	3*	*	Câncer
12	130**	3*	*	Câncer

	Não Sensível		País	Sensível Condição
	CEP	Idade		
1	130**	< 35	*	AIDS
2	130**	< 35	*	Tuberculose
3	130**	< 35	*	Gripe
4	130**	< 35	*	Tuberculose
5	130**	< 35	*	Câncer
6	130**	< 35	*	Câncer
7	130**	35	*	Câncer
8	130**	35	*	Câncer
9	130**	35	*	Câncer
10	130**	35	*	Tuberculose
11	130**	35	*	Virose
12	130**	35	*	Virose

Alice tem 28 anos de idade, mora no CEP 13012, e visita ambos os hospitais

Net ix Prize Dataset (2008)

- | 500.000 registros anônimos de classificações de filmes

Net ix Prize Dataset (2008)

- | 500.000 registros anônimos de classificações de filmes
- | Objetivo era fomentar pesquisas científicas

Net ix Prize Dataset (2008)

- | 500.000 registros anônimos de classificações de filmes
- | Objetivo era fomentar pesquisas científicas
- | Narayanan & Shmatikov cruzaram as informações com perfis públicos no IMBD (The Internet Movie Database)

Net ix Prize Dataset (2008)

- | 500.000 registros anônimos de classificações de filmes
- | Objetivo era fomentar pesquisas científicas
- | Narayanan & Shmatikov cruzaram as informações com perfis públicos no IMBD (The Internet Movie Database)
- | Saber apenas algumas preferências (2-8 filmes, imprecisas) de um assinante foram suficientes para realizar a reidentificação

Net ix Prize Dataset (2008)

- | 500.000 registros anônimos de classificações de filmes
- | Objetivo era fomentar pesquisas científicas
- | Narayanan & Shmatikov cruzaram as informações com perfis públicos no IMBD (The Internet Movie Database)
- | Saber apenas algumas preferências (2-8 filmes, imprecisas) de um assinante foram suficientes para realizar a reidentificação
- | Foi possível inferir posicionamento político e outras informações sensíveis

Anonimização

Métodos Probabilísticos

Ataques Composicionais

- | Problema de métodos determinísticos

Ataques Composicionais

- | Problema de métodos determinísticos
- | Cada **observável** corresponde a um conjunto de **segredos**

Ataques Composicionais

- | Problema de métodos determinísticos
- | Cada **observável** corresponde a um conjunto de **segredos**
- | Combinar observações diferentes permite determinar intersecções no domínio dos segredos

Visão Geral

- I Todo **segredo** pode gerar qualquer **observável** de acordo com uma distribuição de probabilidade

Visão Geral

- I Todo **segredo** pode gerar qualquer **observável** de acordo com uma distribuição de probabilidade

Visão Geral

- | Todo **segredo** pode gerar qualquer **observável** de acordo com uma distribuição de probabilidade
- | Teorema de Bayes:

$$p(s|o) = \frac{p(s)}{p(o)} p(o|s)$$

Visão Geral

- | Todo **segredo** pode gerar qualquer **observável** de acordo com uma distribuição de probabilidade
- | Teorema de Bayes:

$$p(s|o) = \frac{p(s)}{p(o)} p(o|s)$$

Visão Geral

- | Todo **segredo** pode gerar qualquer **observável** de acordo com uma distribuição de probabilidade
- | Teorema de Bayes:

$$p(s|o) = \frac{p(s)}{p(o)} p(o|s)$$

Exemplo

Idade mínima com doença:

Nome	Idade	Doença
Alice	30	não
Bob	30	não
Carl	40	não
Don	40	sim
Ellie	50	não
Frank	50	sim

Exemplo

Nome	Idade	Doença
Alice	30	não
Bob	30	não
Carl	40	não
Don	40	sim
Ellie	50	não
Frank	50	sim

Idade mínima com doença:

- | 30 com probabilidade 1/4

Exemplo

Nome	Idade	Doença
Alice	30	não
Bob	30	não
Carl	40	não
Don	40	sim
Ellie	50	não
Frank	50	sim

Idade mínima com doença:

- | 30 com probabilidade $1/4$
- | 40 com probabilidade $1/2$

Exemplo

Nome	Idade	Doença
Alice	30	não
Bob	30	não
Carl	40	não
Don	40	sim
Ellie	50	não
Frank	50	sim

Idade mínima com doença:

- | 30 com probabilidade $1/4$
- | 40 com probabilidade $1/2$
- | 50 com probabilidade $1/4$

Exemplo

Nome	Idade	Doença
Alice	30	não
Bob	30	não
Carl	40	não
Don	40	sim
Ellie	50	não
Frank	50	sim

Idade mínima com doença:

- | 30 com probabilidade 1/4
- | 40 com probabilidade 1/2
- | 50 com probabilidade 1/4

Alice	Bob
Carl	Don
Ellie	Frank

Exemplo

Peso mínimo com doença:

Nome	Peso	Doença
Alice	60	não
Bob	90	não
Carl	90	não
Don	100	sim
Ellie	60	não
Frank	100	sim

Exemplo

Nome	Peso	Doença
Alice	60	não
Bob	90	não
Carl	90	não
Don	100	sim
Ellie	60	não
Frank	100	sim

Peso mínimo com doença:

- | 60 com probabilidade 1/7

Exemplo

Nome	Peso	Doença
Alice	60	não
Bob	90	não
Carl	90	não
Don	100	sim
Ellie	60	não
Frank	100	sim

Peso mínimo com doença:

- | 60 com probabilidade $1/7$
- | 90 com probabilidade $2/7$

Exemplo

Nome	Peso	Doença
Alice	60	não
Bob	90	não
Carl	90	não
Don	100	sim
Ellie	60	não
Frank	100	sim

Peso mínimo com doença:

- | 60 com probabilidade 1/7
- | 90 com probabilidade 2/7
- | 100 com probabilidade 4/7

Exemplo

Nome	Peso	Doença
Alice	60	não
Bob	90	não
Carl	90	não
Don	100	sim
Ellie	60	não
Frank	100	sim

Peso mínimo com doença:

- | 60 com probabilidade 1/7
- | 90 com probabilidade 2/7
- | 100 com probabilidade 4/7

Alice	Bob
Carl	Don
Ellie	Frank

Exemplo

Nome	Idade	Doença
Alice	30	não
Bob	30	não
Carl	40	não
Don	40	sim
Ellie	50	não
Frank	50	sim

Nome	Peso	Doença
Alice	60	não
Bob	90	não
Carl	90	não
Don	100	sim
Ellie	60	não
Frank	100	sim

Alice	Bob
Carl	Don
Ellie	Frank

Observações

- I Deve-se escolher a distribuição de probabilidade com cuidado

Observações

- | Deve-se escolher a distribuição de probabilidade com cuidado
- | O mecanismo deve proporcionar um equilíbrio entre privacidade e utilidade

Definições

Distância de Hamming

A Distância de Hamming entre dois bancos de dados x_1 e x_2 é igual ao número de registros que diferem entre x_1 e x_2

Definições

Distância de Hamming

A Distância de Hamming entre dois bancos de dados x_1 e x_2 é igual ao número de registros que diferem entre x_1 e x_2

Bancos de Dados Adjacentes

Dois bancos de dados x_1 e x_2 são adjacentes se a Distância de Hamming entre eles é igual a um. Denotamos essa propriedade como $x_1 \sim x_2$

Definições

Distância de Hamming

A Distância de Hamming entre dois bancos de dados x_1 e x_2 é igual ao número de registros que diferem entre x_1 e x_2

Bancos de Dados Adjacentes

Dois bancos de dados x_1 e x_2 são adjacentes se a Distância de Hamming entre eles é igual a um. Denotamos essa propriedade como $x_1 \sim x_2$

- | x_1 e x_2 diferem em relação a apenas um indivíduo,

Definições

Distância de Hamming

A Distância de Hamming entre dois bancos de dados x_1 e x_2 é igual ao número de registros que diferem entre x_1 e x_2

Bancos de Dados Adjacentes

Dois bancos de dados x_1 e x_2 são adjacentes se a Distância de Hamming entre eles é igual a um. Denotamos essa propriedade como $x_1 \sim x_2$

- | x_1 e x_2 diferem em relação a apenas um indivíduo,
- | ou o indivíduo foi adicionado a apenas uma das bases

Definições

Distância de Hamming

A Distância de Hamming entre dois bancos de dados x_1 e x_2 é igual ao número de registros que diferem entre x_1 e x_2

Bancos de Dados Adjacentes

Dois bancos de dados x_1 e x_2 são adjacentes se a Distância de Hamming entre eles é igual a um. Denotamos essa propriedade como $x_1 \sim x_2$

- | x_1 e x_2 diferem em relação a apenas um indivíduo,
 - | ou o indivíduo foi adicionado a apenas uma das bases
 - | ou removido de apenas uma das bases

Definições

Distância de Hamming

A Distância de Hamming entre dois bancos de dados x_1 e x_2 é igual ao número de registros que diferem entre x_1 e x_2

Bancos de Dados Adjacentes

Dois bancos de dados x_1 e x_2 são adjacentes se a Distância de Hamming entre eles é igual a um. Denotamos essa propriedade como $x_1 \sim x_2$

- | x_1 e x_2 diferem em relação a apenas um indivíduo,
 - | ou o indivíduo foi adicionado a apenas uma das bases
 - | ou removido de apenas uma das bases
 - | ou teve seu valor alterado em apenas uma das bases

Privacidade-Diferencial (Dwork, 2006)

Definição

Uma base de dados é dita Privacidade-Diferencial se para todos os bancos de dados x, x^0 adjacentes, e para todo $z \in Z$, a resposta reportada, temos que:

$$\frac{p(K = z | X = x)}{p(K = z | X = x^0)} \leq e$$

Privacidade-Diferencial (Dwork, 2006)

Definição

Uma base de dados é dita Privacidade-Diferencial se para todos os bancos de dados x, x^0 adjacentes, e para todo $z \in Z$, a resposta reportada, temos que:

$$\frac{p(K = z | X = x)}{p(K = z | X = x^0)} \leq e$$

Interpretação

A presença ou não de informações sobre um indivíduo na base de dados não altera significativamente os resultados obtidos

Propriedades da Privacidade Diferencial

- I Independência prior , a distribuição de probabilidade sobre os segredos antes de se consultar o banco de dados

Propriedades da Privacidade Diferencial

- | Independência *a priori*, a distribuição de probabilidade sobre os segredos antes de se consultar o banco de dados
 - | Independe do adversário

Propriedades da Privacidade Diferencial

- | Independência d_{prior} , a distribuição de probabilidade sobre os segredos antes de se consultar o banco de dados
 - | Indepe de do adversário
- | Robustez quanto à Composicionalidade

Propriedades da Privacidade Diferencial

- | Independência *a priori*, a distribuição de probabilidade sobre os segredos antes de se consultar o banco de dados
 - | Independe do adversário
- | Robustez quanto à Composicionalidade
 - | Dados dois mecanismos Privados Diferencial e ϵ_2 -Diferencial, a composição é um mecanismo $(\epsilon_1 + \epsilon_2)$ -Diferencial

Propriedades da Privacidade Diferencial

- | Independência *a priori*, a distribuição de probabilidade sobre os segredos antes de se consultar o banco de dados
 - | Independe do adversário
- | Robustez quanto à Composicionalidade
 - | Dados dois mecanismos Privados Diferencial e ϵ_2 -Diferencial, a composição é um mecanismo $(\epsilon_1 + \epsilon_2)$ -Diferencial
 - | A privacidade diminui linearmente com o número de consultas

Propriedades da Privacidade Diferencial

- | Independência *a priori*, a distribuição de probabilidade sobre os segredos antes de se consultar o banco de dados
 - | Independe do adversário
- | Robustez quanto à Composicionalidade
 - | Dados dois mecanismos Privacidade Diferencial e ϵ_2 -Diferencial, a composição é um mecanismo $(\epsilon_1 + \epsilon_2)$ -Diferencial
 - | A privacidade diminui linearmente com o número de consultas
 - | É possível definir um orçamento de privacidade, ou seja, um máximo aceitável de violação da privacidade

Exemplos de Privacidade Diferencial

Qual a altura média?

A distribuição varia entre
50 cm e 250 cm

	Altura (cm)
Alice	140
Bob	180
Carol	160
Daniel	120 (p) 190 ($1 - p$)

Exemplos de Privacidade Diferencial

Qual a altura média?

A distribuição varia entre
50 cm e 250 cm

	Altura (cm)
Alice	140
Bob	180
Carol	160
Daniel	120 (p) 190 ($(1-p)$)

Um mecanismo que sempre reporta a resposta verdadeira (150 cm) não é privado, qualquer que seja o

$$\frac{p(\text{media} = 150 | \text{Daniel} = 120)}{p(\text{media} = 150 | \text{Daniel} = 190)} = \frac{1}{0} = e^1$$

Exemplos de Privacidade Diferencial

Qual a altura média?
A distribuição varia entre
50 cm e 250 cm

	Altura (cm)
Alice	140
Bob	180
Carol	160
Daniel	120 (p) 190 ($1 - p$)

Um mecanismo que sempre reporta
uma resposta errada (168 cm) é
completamente privado ($\epsilon = 0$), mas
completamente inútil

$$\frac{p(\text{media} = 168 | \text{Daniel} = x)}{p(\text{media} = 168 | \text{Daniel} = x^0)} = \frac{1}{1} = e^0$$

$$\frac{p(\text{media} = z | \text{Daniel} = x)}{p(\text{media} = z | \text{Daniel} = x^0)} = \frac{0}{0}$$

($z \notin \{168\}$)

Exemplos de Privacidade Diferencial

Qual a altura média?

A distribuição varia entre
50 cm e 250 cm

	Altura (cm)
Alice	140
Bob	180
Carol	160
Daniel	120 (p) 190 ($1 - p$)

Um mecanismo que sempre reporta 100 cm se a resposta verdadeira (150 cm) for menor ou igual a 150 cm, ou reporta 200 cm caso contrário, não é privado, qualquer que seja o

$$\frac{p(\text{media} = 100 | \text{Daniel} = 120)}{p(\text{media} = 100 | \text{Daniel} = 190)} = \frac{1}{0} = e^1$$

Exemplos de Privacidade Diferencial

Qual a altura média?

A distribuição varia entre
50 cm e 250 cm

	Altura (cm)
Alice	140
Bob	180
Carol	160
Daniel	120 (p) 190 ($1 - p$)

O mecanismo que reporta a resposta verdadeira com probabilidade $\frac{1}{2}$ e todos os outros inteiros no intervalo $[50; 250]$ com probabilidade $\frac{1}{200}$, é privado -diferencial



Estado da Arte

Google, Microsoft, e Apple

Privacidade Diferencial Local

De nição

A randomização e a adiço de ruido so realizadas pelo software em execuo no dispositivo do usurio

Privacidade Diferencial Local

Exemplos

As seguintes empresas privadas reportam utilizar Privacidade Diferencial Local para coleta de Telemetria

Privacidade Diferencial Local

Exemplos

As seguintes empresas privadas reportam utilizar Privacidade Diferencial Local para coleta de Telemetria

- | Google (Chrome): RAPPOR (Randomized Aggregatable Privacy-Preserving Ordinal Response)

Privacidade Diferencial Local

Exemplos

As seguintes empresas privadas reportam utilizar Privacidade Diferencial Local para coleta de Telemetria

- | Google (Chrome): RAPPOR (Randomized Aggregatable Privacy-Preserving Ordinal Response)
- | Apple: iOS 10 em diante

Privacidade Diferencial Local

Exemplos

As seguintes empresas privadas reportam utilizar Privacidade Diferencial Local para coleta de Telemetria

- | Google (Chrome): RAPPOR (Randomized Aggregatable Privacy-Preserving Ordinal Response)
- | Apple: iOS 10 em diante
- | Microsoft: Windows 10

Privacidade Diferencial Local

Exemplos

As seguintes empresas privadas reportam utilizar Privacidade Diferencial Local para coleta de Telemetria

- | Google (Chrome): RAPPOR (Randomized Aggregatable Privacy-Preserving Ordinal Response)
- | Apple: iOS 10 em diante
- | Microsoft: Windows 10

O objetivo nal é aumentar a aceitação pública da coleta de dados

Privacidade Diferencial Local

Exemplos

As seguintes empresas privadas reportam utilizar Privacidade Diferencial Local para coleta de Telemetria

- | Google (Chrome): RAPPOR (Randomized Aggregatable Privacy-Preserving Ordinal Response)
- | Apple: iOS 10 em diante
- | Microsoft: Windows 10

O objetivo principal é aumentar a aceitação pública da coleta de dados

Entretanto, para garantir a utilidade dos dados coletados, utilizam um orçamento de privacidade (muito grande)



Estado da Arte

US Census Bureau

Visão Geral

- | Responsável pela realização dos Censos decenais

Visão Geral

- | Responsável pela realização dos Censos decenais
- | Os resultados são publicados agregados e servem de base para

Visão Geral

- | Responsável pela realização dos Censos decenais
- | Os resultados são publicados agregados e servem de base para
 - | a repartição dos assentos na Câmara dos Deputados entre os estados

Visão Geral

- | Responsável pela realização dos Censos decenais
- | Os resultados são publicados agregados e servem de base para
 - | a repartição dos assentos na Câmara dos Deputados entre os estados
 - | a distribuição de mais de US\$675 bilhões em fundos federais para estados e organizações locais

Visão Geral

- | Responsável pela realização dos Censos decenais
- | Os resultados são publicados agregados e servem de base para
 - | a repartição dos assentos na Câmara dos Deputados entre os estados
 - | a distribuição de mais de US\$675 bilhões em fundos federais para estados e organizações locais
- | Por Lei, a Confidencialidade deve ser garantida

Visão Geral

- I Até o Censo de 1970, apenas desidenti cava os registros

Visão Geral

- | Até o Censo de 1970, apenas desidenti cava os registros
- | Entre os Censos de 1980 e de 2010, aplicou também métodos determinísticos

Visão Geral

- | Até o Censo de 1970, apenas desidenti cava os registros
- | Entre os Censos de 1980 e de 2010, aplicou também métodos determinísticos
- | Teorema da Reconstrução da Base de Dados (Dinur & Nissim, 2003)

Visão Geral

- | Até o Censo de 1970, apenas desidenti cava os registros
- | Entre os Censos de 1980 e de 2010, aplicou também métodos determinísticos
- | Teorema da Reconstrução da Base de Dados (Dinur & Nissim, 2003)
 - | Muitas estatísticas publicadas com muita precisão a partir de uma base de dados con dencial expõem a base de dados inteira com quase certeza

Visão Geral

- | Até o Censo de 1970, apenas desidenti cava os registros
- | Entre os Censos de 1980 e de 2010, aplicou também métodos determinísticos
- | Teorema da Reconstrução da Base de Dados (Dinur & Nissim, 2003)
 - | Muitas estatísticas publicadas com muita precisão a partir de uma base de dados con dencial expõem a base de dados inteira com quase certeza
 - | Para se ter privacidade, é necessário adicionar uma perturbação de magnitude ao menos $\frac{1}{\sqrt{N}}$, onde N é o tamanho da população

2020 Census

- | Métodos determinísticos foram descartados devido à legislação

2020 Census

- | Métodos determinísticos foram descartados devido à legislação
- | Entretanto, Privacidade Diferencial Local adiciona muito ruído, caso seja utilizada para garantir de fato a privacidade

2020 Census

- | Métodos determinísticos foram descartados devido à legislação
- | Entretanto, Privacidade Diferencial Local adiciona muito ruído, caso seja utilizada para garantir de fato a privacidade
- | Todo o banco de dados deve ser processado de uma só vez para garantir melhor precisão

2020 Census

- | Métodos determinísticos foram descartados devido à legislação
- | Entretanto, Privacidade Diferencial Local adiciona muito ruído, caso seja utilizada para garantir de fato a privacidade
- | Todo o banco de dados deve ser processado de uma só vez para garantir melhor precisão
- | Todos os usos dos dados privados devem ser previamente considerados no orçamento de privacidade, ou seja, antes da publicação dos dados

2020 Census

- | Métodos determinísticos foram descartados devido à legislação
- | Entretanto, Privacidade Diferencial Local adiciona muito ruído, caso seja utilizada para garantir de fato a privacidade
- | Todo o banco de dados deve ser processado de uma só vez para garantir melhor precisão
- | Todos os usos dos dados privados devem ser previamente considerados no orçamento de privacidade, ou seja, antes da publicação dos dados
- | Foi desenvolvido um novo método **Top-Down**, que permite a criação das tabelas com informações agregadas sem grandes perdas de utilidade

2020 Census

- | Métodos determinísticos foram descartados devido à legislação
- | Entretanto, Privacidade Diferencial Local adiciona muito ruído, caso seja utilizada para garantir de fato a privacidade
- | Todo o banco de dados deve ser processado de uma só vez para garantir melhor precisão
- | Todos os usos dos dados privados devem ser previamente considerados no orçamento de privacidade, ou seja, antes da publicação dos dados
- | Foi desenvolvido um novo método Top-Down, que permite a criação das tabelas com informações agregadas sem grandes perdas de utilidade
- | Novos resultados referentes ao Censo de 2010 foram publicados para servirem de exemplo da aplicação do novo método

2020 Census

- | Métodos determinísticos foram descartados devido à legislação
- | Entretanto, Privacidade Diferencial Local adiciona muito ruído, caso seja utilizada para garantir de fato a privacidade
- | Todo o banco de dados deve ser processado de uma só vez para garantir melhor precisão
- | Todos os usos dos dados privados devem ser previamente considerados no orçamento de privacidade, ou seja, antes da publicação dos dados
- | Foi desenvolvido um novo método Top-Down, que permite a criação das tabelas com informações agregadas sem grandes perdas de utilidade
- | Novos resultados referentes ao Censo de 2010 foram publicados para servirem de exemplo da aplicação do novo método
- | US National Census Day: 01 de Abril de 2020

Resumo

Governos e empresas privadas coletam e continuarão coletando dados pessoais

Governos e empresas privadas coletam e continuarão coletando dados pessoais

- | É importante que essas entidades garantam a segurança dessa informações

Governos e empresas privadas coletam e continuarão coletando dados pessoais

- | É importante que essas entidades garantam a segurança dessa informações
 - | Controle de acesso e Criptogra a são muito importantes, mas não resolvem todos os problemas

Governos e empresas privadas coletam e continuarão coletando dados pessoais

- | É importante que essas entidades garantam a segurança dessa informações
 - | Controle de acesso e Criptogra a são muito importantes, mas não resolvem todos os problemas
- | É importante que essas entidades continuem divulgando dados, mas de forma responsável

Governos e empresas privadas coletam e continuarão coletando dados pessoais

- | É importante que essas entidades garantam a segurança dessa informações
 - | Controle de acesso e Criptogra a são muito importantes, mas não resolvem todos os problemas
- | É importante que essas entidades continuem divulgando dados, mas de forma responsável
 - | Muitos métodos de Anonimização já foram propostos, mas diversos problemas permanecem em aberto

Governos e empresas privadas coletam e continuarão coletando dados pessoais

- | É importante que essas entidades garantam a segurança dessa informações
 - | Controle de acesso e Criptogra a são muito importantes, mas não resolvem todos os problemas
- | É importante que essas entidades continuem divulgando dados, mas de forma responsável
 - | Muitos métodos de Anonimização já foram propostos, mas diversos problemas permanecem em aberto
 - | Como de nir o conjunto de quasi-identi cadores?

Governos e empresas privadas coletam e continuarão coletando dados pessoais

- | É importante que essas entidades garantam a segurança dessa informações
 - | Controle de acesso e Criptogra a são muito importantes, mas não resolvem todos os problemas
- | É importante que essas entidades continuem divulgando dados, mas de forma responsável
 - | Muitos métodos de Anonimização já foram propostos, mas diversos problemas permanecem em aberto
 - | Como de nir o conjunto de quasi-identi cadores?
 - | Como garantir um nível aceitável de utilidade?

Governos e empresas privadas coletam e continuarão coletando dados pessoais

- | É importante que essas entidades garantam a segurança dessa informações
 - | Controle de acesso e Criptogra a são muito importantes, mas não resolvem todos os problemas
- | É importante que essas entidades continuem divulgando dados, mas de forma responsável
 - | Muitos métodos de Anonimização já foram propostos, mas diversos problemas permanecem em aberto
 - | Como de nir o conjunto de quasi-identi cadores?
 - | Como garantir um nível aceitável de utilidade?
 - | Como anonimizar dados longitudinais no tempo?



Como equilibrar Transparência e Privacidade?

Como equilibrar Transparência e Privacidade?

Lembrando que a Lei Geral de Proteção de Dados
entra em vigor em Agosto de 2020 no Brasil!

Departamento de Ciência da Computação

Equipe da EVCOMP 2020

Prof. Mário Alvim e Prof. Annabelle McIver

Prof. Catuscia Palamidessi e Prof. Kostas Chatzikelakos



Laboratory of Information Security,
Cryptography, Privacy, and Transparency

Obrigado pela atenção!

ghn@nunesgh.com 
nunesgh.com/sobre 

AGMZYRE
INSCRYPT

Laboratory of Information Security,
Cryptography, Privacy, and Transparency