

Samenvatting

Bescherming van persoonsgegevens in het justitiële domein: Richtlijnen voor Statistical Disclosure Control Project Privacy Utility Tools 2.0

Achtergrond, toepassingsgebied en bijdragen

Overheden proberen hun transparantie, verantwoording en efficiëntie te verbeteren door hun met overheidsgeld gefinancierde datasets proactief open te stellen voor het publiek. Op die manier willen zij participatief bestuur door burgers ondersteunen, innovaties en economische groei voor publieke en/of particuliere ondernemingen bevorderen en het voor burgers en organisaties gemakkelijker maken om onderbouwde beslissingen te nemen. Overheidsinstanties delen ook gegevens met anderen om verschillende redenen, bijvoorbeeld om hun operationele activiteiten te faciliteren, statistische inzichten in de status van hun operationele activiteiten en strategische doelstellingen te verkrijgen, en wetenschappelijk onderzoek uit te voeren dat relevant is voor hun missie, bijvoorbeeld naar de impact van hun beleid. Bescherming van persoonsgegevens is voor overheidsinstanties een belangrijke voorwaarde om hun gegevens op verantwoorde wijze openbaar te kunnen maken en te kunnen delen. Een van de belangrijkste beginselen van de bescherming van persoonsgegevens is het beperken van de hoeveelheid persoonsgegevens die gedeeld worden tot het niveau dat nodig is voor het beoogde gebruik van de gegevens. Als het gaat om open data, waarbij de opengestelde gegevens voor iedereen toegankelijk zijn, ook voor potentiële kwaadwillenden, komt de bescherming van persoonsgegevens voornamelijk neer op gegevensminimalisering.

Er bestaan verschillende technologieën om persoonsgegevens in een dataset te beschermen. SDC-technologieën (Statistical Disclosure Control) zijn daar een onderdeel van. Deze technologieën zijn ontwikkeld om de hoeveelheid persoonsgegevens in een dataset te beperken en tegelijk de bruikbaarheid van de gegevens te behouden. SDC-technologieën kunnen worden toegepast op zowel microdatasets als geaggregeerde datasets. Microdatasets, die (zeer) omvangrijk kunnen zijn, bestaan uit gestructureerde tabellen met een aantal rijen, die staan voor personen, en een aantal kolommen, die staan voor de kenmerken (attributen) van die personen (zoals hun leeftijd, geslacht en beroep). Geaggregeerde datasets zijn opgebouwd uit microdata. Een geaggregeerde dataset bevat een of meer tabellen met een aantal rijen en kolommen, die corresponderen met een aantal groepskenmerken, die weer een subset zijn van de kenmerken in de betreffende microdata. Eerder bestudeerden we SDC-technieken voor de bescherming van microdatasets (Bargh et al., 2018) en voor geaggregeerde datasets (Bargh et al., 2020).

Het hoofddoel van ons onderzoeksproject is het vergroten van de kennis binnen de Nederlandse overheid, en meer in het bijzonder het ministerie van Justitie en Veiligheid, over SDC-technologie, en de mogelijkheden, de beperkingen en het gebruik daarvan. In navolging van eerdere publicaties (Bargh et al., 2018, 2020) zetten we in dit rapport een volgende stap naar de ingebruikname van SDC-technologie binnen overheidsinstanties door enkele initiële richtlijnen te ontwikkelen

voor het gebruik van SDC-technologie in de praktijk. We verwachten dat SDC-technologie op die manier toegankelijker wordt voor data stewards die verantwoordelijk zijn voor het op verantwoorde wijze delen of openstellen van datasets.

De toepassing van SDC-technologie is een multidisciplinaire taak waarvoor onder meer juridische en technische expertise nodig is. De initiële richtlijnen in dit rapport zijn bedoeld om de technische SDC-kennis en de gebruiksvaardigheden van data stewards die nog niet eerder SDC-technologie hebben gebruikt te verbeteren. De belangrijkste bijdrage van dit rapport is het bieden van een eerste reeks richtlijnen voor:

- het werkproces voor het gebruik van SDC-technologie voor de bescherming van microdata en geaggregeerde data;
- de belangrijkste activiteiten in elke stap van het proces; en
- de configuratie van specifieke stappen in de praktijk.

De toepassing van SDC-technologie in de praktijk (d.w.z. de ingebruikname ervan binnen organisaties) is vanwege de complexiteit, het multidisciplinaire karakter en de afhankelijkheid van de context geen eenmalige onderneming. Daarom bieden we ook een raamwerk waarmee SDC-technologie stapsgewijs en geleidelijk in een organisatie kan worden ingevoerd en geïntegreerd. De initiële SDC-richtlijnen uit dit rapport dienen als een van de eerste stappen van dit raamwerk.

NB: in dit rapport duiden we met de term 'anonimisering van gegevens' op het proces waarbij directe en indirecte identificerende gegevens uit een dataset worden verwijderd en SDC-technieken worden toegepast om de dataset te kunnen delen of openstellen.

Methodologie

We hebben verschillende methoden gebruikt om de SDC-richtlijnen en het raamwerk te ontwikkelen, zoals literatuuronderzoek, casusstudies, interviews met deskundigen, experimenten, prototypes en simulaties.

Ons literatuuronderzoek heeft geleid tot twee technische rapporten (Bargh et al., 2018, 2020), die dienen als aanvullende leesstof waarin de theoretische onderbouwing van de initiële richtlijnen wordt toegelicht en enkele voorbeelden worden gegeven. We hebben vier casestudies uitgevoerd om meer te weten te komen over de manier waarop gegevens op dit moment in het justitiële domein in de praktijk geanonimiseerd worden. Daarnaast hebben we vier interviews met deskundigen gehouden om enkele keuzes die we voor de initiële richtlijnen hebben gemaakt te toetsen en onderbouwen. We hebben de eerste versie van de initiële richtlijnen toegepast op vier datasets uit het justitiële domein om de toepasbaarheid van de richtlijnen in de praktijk te evalueren. Sommige gebruikersinterface-aspecten van de initiële richtlijnen zijn geïmplementeerd in een mock-up-prototype van een SDC-softwaretool. Tot slot hebben we een aantal simulaties uitgevoerd met een open-source-SDC-softwaretool om enkele van onze ontwerpbeslissingen voor de richtlijnen te onderbouwen.

Al deze methoden hebben ons in staat gesteld ontwerpkeuzes te maken voor de initiële richtlijnen en de reikwijdte daarvan tot een hanteerbaar niveau te brengen. Bovendien hebben ze ons geholpen de beperkingen van het werk in kaart te brengen en een aantal richtingen voor toekomstig onderzoek te bepalen.

Belangrijkste resultaten

Hieronder beschrijven we kort de belangrijkste resultaten van het onderzoek.

Over (het belang van) anonimisering van gegevens

We hebben aangetoond dat er een aantal drijvende krachten is achter de noodzaak van dataminimalisatie. De noodzaak om de verwerking van persoonsgegevens te beperken in lijn met het beoogde doel wordt benadrukt in allerlei (recente) privacywet- en -regelgeving, zoals de AVG.

Het minimaliseren van persoonsgegevens begint vaak met het verwijderen van informatie uit een dataset die tot directe identificatie kan leiden (zoals namen en burgerservicenummers). Wanneer ook de informatie die tot indirecte identificatie kan leiden, moet worden beperkt, wordt het proces van gegevensbescherming ingewikkelder. In dergelijke gevallen wordt vaak SDC-technologie gebruikt om de hoeveelheid indirect identificeerbare informatie over personen in een dataset aan te passen naar een gewenst, vereist of toegestaan niveau, afhankelijk van het doel waarvoor de gegevens worden gebruikt. SDC-technologie kan aan dit proces bijdragen door inzicht te geven in en mechanismen te creëren voor:

- a het transformeren van ruwe gegevens;
- b het beoordelen van de bruikbaarheid van de oorspronkelijke en de getransformeerde gegevens;
- c het inschatten van de onthullingsrisico's van de oorspronkelijke en de getransformeerde gegevens; en
- d het maken van afwegingen tussen het nut van gegevens en de onthullingsrisico's.

Zoals elke andere (gegevensbeschermings)technologie moeten de mogelijkheden van SDC-technologie met enige terughoudendheid worden bekeken. Eén voorbehoud is dat het minimaliseren van persoonsgegevens door middel van SDC-technologie geen gegarandeerde anonimiteit oplevert in de zin van de definitie uit de AVG. Daarom mag SDC-technologie niet worden beschouwd als een wondermiddel, want zij biedt geen volledige, kant-en-klare oplossing voor gegevensbescherming. Toch is het gebruik van op SDC gebaseerde inzichten en SDC-technologie *noodzakelijk* om te voldoen aan de regelgeving die geldt voor het delen of openstellen van persoonsgegevens. Vaak moet SDC-technologie dan worden gecombineerd met aanvullende technische maatregelen en/of niet-technische mechanismen voor gegevensbeheer (zoals contracten, beleid en organisatorische procedures) om de restrisico's te beperken.

Over het beoogde raamwerk voor integratie in organisaties

SDC-technologie is een geavanceerd expertisegebied waarin actief onderzoek wordt gedaan en dat voortdurend verder ontwikkelt. Het gebruik van SDC-technologie vereist een holistische benadering waarbij rekening wordt gehouden met technologische, juridische, ethische, publieke en organisatorische aspecten. Bovendien is de toepassing van SDC in de praktijk afhankelijk van vele omgevingsfactoren, zoals de beschikbaarheid van achtergrondkennis voor indringers en de gevoeligheid van de te delen gegevens. Daarom presenteren wij een raamwerk om SDC-technologie in te bedden in een organisatie.

Het beoogde raamwerk omvat een model om SDC-verantwoordelijkheden binnen een organisatie te verdelen, en een iteratief leerproces om relevante SDC-kennis binnen de organisatie te ontwikkelen op basis van de (toenemende) behoeften van de organisatie.

- Het beoogde model beoogt een optimale verdeling van SDC-kennis en -vaardigheden binnen een organisatie: enerzijds bij de afzonderlijke partijen of afdelingen die privacygevoelige datasets beheren en anderzijds bij een centrale partij die over expertise op het gebied van SDC-technologie beschikt. Een mogelijke vorm is dat routine SDC-taken aan individuele partijen worden gedelegeerd en complexe en geavanceerde taken aan de centrale partij.
- Het iteratieve leerproces binnen ons raamwerk is erop gericht de SDC-taken geleidelijk zoveel mogelijk te distribueren en te delegeren aan de individuele partijen. Op die manier proberen we op termijn SDC-expertise te creëren bij alle individuele partijen, zonder hen onmiddellijk te belasten met en de verantwoordelijkheid te geven voor het leren en toepassen van complexe SDC-taken. Volgens dit leerproces wordt de initiële reeks SDC-richtlijnen geleidelijk uitgebreid (en/of gewijzigd) doordat men leert van de praktijk.

Over de algemene richtlijnen

De initiële SDC-richtlijnen uit dit rapport zijn bedoeld om enerzijds het proces van het gebruik van SDC-technologie te beschrijven en anderzijds enkele aanbevelingen te doen voor de te gebruiken beschermingsmodellen en -methoden en de parameterconfiguraties. We beginnen met een reeks basisrichtlijnen, die zorgvuldig zijn ontwikkeld op basis van onze onderzoeksmethoden. Deze initiële SDC-richtlijnen kunnen worden gebruikt voor onderwijsdoeleinden en bij het uitvoeren van routinematige gegevensanonimiseringstaken. Volgens het beoogde ontwikkelingsproces zouden de initiële SDC-richtlijnen geleidelijk uitgebreid (en/of gewijzigd) moeten worden door van de praktijk te leren.

De initiële SDC-richtlijnen zijn onderverdeeld in algemene en specifieke richtlijnen. De specifieke richtlijnen zijn bedoeld voor de bescherming van microdata of geaggregeerde data en worden in de volgende paragraaf beschreven. In de algemene richtlijnen worden de taken van het SDC-proces beschreven die van toepassing zijn op de bescherming van beide soorten data. Deze algemene taken hebben betrekking op:

- het selecteren van de te delen gegevens;
- het specificeren van de doelstelling(en) van het delen van gegevens;
- het specificeren van de dataomgeving; en
- het transformeren, analyseren en delen van gegevens.

De te delen of publiceren dataset wordt doorgaans geselecteerd op basis van de relevante strategische doelstellingen, beleidslijnen en overwegingen. Dit kan reactief gebeuren, als reactie op een concreet dataverzoek, of proactief, om transparantie te creëren (bv. in het geval van Open Data).

Deze specificatie van het doel waarmee gegevens worden gedeeld, wat gewoonlijk buiten het anonimiseringsproces gebeurt, kan vervolgens in het SDC-proces gebruikt worden om bepaalde aspecten van de dataomgeving te specificeren, passende maatregelen te kiezen om de bruikbaarheid van gegevens en de onthullingsrisico's te beoordelen, en afwegingen te maken tussen risico's en bruikbaarheid.

Het specificeren van de dataomgeving heeft betrekking op het modelleren van de context waarbinnen de gegevens worden gedeeld en gebruikt. Een dergelijk contextmodel is belangrijk om de onthullingsrisico's te kunnen bepalen. Bij het modelleren van de context spelen onder meer de volgende factoren een rol:

- de betrokken actoren (bv. de soorten indringers);

- de aanvullende gegevensbronnen die potentiële indringers kunnen gebruiken als achtergrondkennis om persoonsgegevens uit een gedeelde dataset openbaar te maken;
- de mechanismen voor gegevensbeheer die worden gebruikt om restricties te beperken; en
- de technische middelen die zijn gebruikt om de gedeelde dataset te beschermen of die gebruikt kunnen worden om persoonsgegevens uit de gedeelde dataset af te leiden.

De richtlijnen over het transformeren en analyseren van gegevens zijn afhankelijk van het gegevenstype (d.w.z. microdata of geaggregeerde data) en worden daarom hieronder beschreven. Ten slotte is bij het delen van gegevens een van de belangrijkste taken het documenteren van het anonimiseringsproces voor zowel intern als extern gebruik.

Over de gegevensspecifieke richtlijnen

De gegevenstransformatie- en -analysetaken uit de algemene richtlijnen zijn vergelijkbaar, maar verschillen technisch per gegevenstype. Daarom worden ze in het rapport afzonderlijk besproken.

De SDC-taken die specifiek zijn voor microdata, draaien om het transformeren en analyseren van gegevens. Voor *het transformeren van gegevens* definiëren we vervolgens een aantal taken, namelijk: het bepalen van de privacybenadering, het beoordelen van de gevoeligheid van de attributen, het kiezen van privacymodellen en -methoden, het configureren van parameters, het toepassen van de gekozen modellen en methoden en de configuraties daarvan.

- Er worden twee *privacybenaderingen* onderscheiden voor de bescherming van gegevens, namelijk een syntactische en een op ruis gebaseerde benadering. We hebben die eerste gekozen als basis voor de initiële SDC-richtlijnen, omdat deze waarheidsgetrouw is en gemakkelijk kan worden gekoppeld aan juridische privacyconcepten. In toekomstige edities van de richtlijnen kan ook worden overwogen op ruis gebaseerde benaderingen op te nemen.
- Voor de syntactische benadering worden de attributen in de te delen microdataset onderverdeeld in vier typen, namelijk expliciet identificerende attributen, quasi-identificerende attributen, gevoelige attributen en niet-gevoelige attributen. De expliciet identificerende attributen (zoals namen en persoonsnummers) worden over het algemeen verwijderd. De quasi-identificerende attributen kunnen door indringers potentieel gebruikt worden om personen in de dataset te identificeren.
- In de initiële richtlijnen stellen wij voor twee *privacymodellen* te gebruiken om quasi-identificerende en gevoelige attributen te beschermen: k-anonimiteit en l-diversiteit. Daarvoor moet men zich echter wel bewust zijn van de mogelijkheden en beperkingen van deze modellen.
- Vervolgens moeten de parameters van het gekozen privacymodel of de gekozen privacymodellen en methoden worden *geconfigureerd*. Daarbij kunnen verschillende factoren in aanmerking worden genomen, zoals de gevoeligheid van de gedeelde attributen, de gevoeligheid van de waarden van de gedeelde attributen, het doel van de gegevensdeling, het al dan niet bestaan van aanvullende mechanismen voor gegevensbescherming, het verwachte soort aanvallers, de reputatie van de gegevensverwerkers, de frequentie en het type van de gegevensselectie (ten opzichte van de volledige dataset), om er maar een paar te noemen.
- Met behulp van een softwaretool kunnen de gekozen privacymodellen, -methoden en -parameters op de dataset worden *toegepast*.

Bij de *analyse van gegevens*, worden de bruikbaarheid van de gegevens en de onthullingsrisico's van de getransformeerde dataset geanalyseerd. Vervolgens wordt besloten of de afweging tussen de bruikbaarheid en de privacy van de gegevens (d.w.z. de onthullingsrisico's) bevredigend is.

Geaggregeerde data zijn een aggregatie van microdata. De SDC-taken voor de bescherming van geaggregeerde datasets zijn functioneel vergelijkbaar met die voor de bescherming van microdatasets. Vergeleken met microdata hebben geaggregeerd data minder dimensies (d.w.z. minder attributen), maar ze kunnen wel veel meer afhankelijkheden hebben. Deze afhankelijkheden, die bijvoorbeeld ontstaan wanneer meerdere verschillende tabellen uit dezelfde microdataset worden samengesteld, kunnen mogelijk worden gebruikt om privacygevoelige informatie te onthullen. Daarom is de bescherming van geaggregeerde data meer gericht op het identificeren en verhelpen van deze afhankelijkheden. Door deze andere focus verschilt het proces voor de bescherming van geaggregeerde data enigszins van het proces bij microdata.

De SDC-taken specifiek voor geaggregeerde data met betrekking tot gegevenstransformatie en -analyse, zijn: het ontwerpen van de tabel, het kiezen van maatregelen om het onthullingsrisico te verkleinen, het kiezen van beschermingsmethoden, het configureren van parameters en het analyseren van gegevens.

- Tijdens het *tabelontwerp* wordt het gewenste type tabel bepaald (d.w.z. frequentietabel of kwantitatieve tabel), worden de groeperingsattributen en hun waarden geselecteerd en wordt de structuur van de tabel gespecificeerd in termen van de bestaande relaties binnen de tabel en met de andere tabellen die in de dataomgeving aanwezig (zullen) zijn.
- Bij het *kiezen van de maten voor het bepalen van het onthullingsrisico* wordt bepaald welke cellen mogelijk risico lopen op onthulling op basis van (een beperkt aantal) gevoeligheidsregels. Voor de initiële reeks richtlijnen zijn de voorgestelde regels de frequentieregel, de p%-regel, en nulcellen en scheve verdelingen.
- Door het *kiezen van beschermingsmethode(n)* wordt geprobeerd de bedreigingen te verminderen voor de cellen die risico lopen. Daartoe stellen wij een algemene volgorde van stappen voor het kiezen en toepassen van beschermingsmethoden voor, op basis van de gewenste eigenschappen van de beschermde tabel.
- Door *parameters te configureren* kunnen de parameters van de gekozen maatregelen voor het verminderen van het onthullingsrisico en de gekozen beschermingsmethoden nauwkeurig worden afgesteld.
- Tot slot worden, door gegevens te analyseren, de bruikbaarheid en de privacy van gegevens (d.w.z. de onthullingsrisico's) beoordeeld op basis van een aantal maten (d.w.z. gevoeligheidsregels). Vervolgens wordt beslist of de afweging tussen het nut en de privacy van de gegevens bevredigend is.

Deze specifieke SDC-richtlijnen voor de gegevenstransformatie- en -analysetaken worden voor zowel de microdata als de geaggregeerde data herhaald totdat een bevredigende balans is gevonden tussen de bruikbaarheid en de privacy van de gegevens.