

2023 Privacy Report

개인정보보호 월간동향분석

3월호



2023 Privacy Report

개인정보보호 월간동향분석

3월호

1. 웹 3.0 시대 도래로 부상하는 개인정보보호 이슈 분석
2. 사우디아라비아의 데이터 · 프라이버시 규제 샌드박스 추진현황
3. 개인정보보호와 활용성 강화 기술로 주목받는 재현데이터
기술의 특성과 활용 과제

KISA

웹 3.0 시대 도래로 부상하는 개인정보보호 이슈 분석

[목 차]

1. 웹 3.0의 개념과 특징

2. 웹 3.0 주요 특징별 개인정보 이슈

- (1) 탈중앙화
- (2) 지능화
- (3) 사용자 중심성

3. 요약 및 결론

1. 웹 3.0의 개념과 특징

- ▶ 웹 3.0(Web 3.0)은 지능적이며 사용자 중심의 웹 경험을 보다 탈중앙화된 형태로 제공하는 것을 목표로 하는 차세대 월드와이드웹(WWW)을 의미
 - 인공지능(AI), 시맨틱웹(Semantic Web)¹⁾ 등의 기술을 기반으로 성장 중인 웹 3.0은 블록체인 기술과 접목되어 탈중앙화적인 특징이 강조되고 있으며, 웹에서 새로운 형태의 상호 작용, 협업, 가치 창출을 가능하게 함
 - 특히 최근 마이크로소프트의 검색엔진 Bing과 브라우저 Edge에 대형 언어 모델(LLM) 기반 생성 AI인 ChatGPT가 탑재되는 등 웹브라우저가 고도화되는 가운데 웹 3.0과 지능형 웹 기술에 대한 기대감이 더욱 증폭
- ▶ 웹 2.0까지 '인터넷 브라우저'는 검색과 데이터 저장을 위한 역할에 그쳤으나, 웹 3.0 시대에서는 블록체인, AI 등 최신 IT 기술들이 결합되면서 '초개인화된 웹 생태계'를 형성 중

1) 컴퓨터가 사람을 대신하여 정보를 읽고 이해하고 가공하여 새로운 정보를 만들어 낼 수 있도록, 이해하기 쉬운 의미를 가진 차세대 지능형 웹. XML에 기반한 시맨틱 마크업 언어를 기반으로 함

표 _ 웹의 진화 양상

구분	웹 1.0	웹 2.0	웹 3.0
정의	정적 정보와 제한된 양방향성을 제공하는 읽기 중심의 웹	사용자 생성 콘텐츠와 협업을 촉진하는 참여형 소셜 웹	분산 및 개인화 서비스 창출을 위해 AI와 블록체인 등의 기술을 활용한 읽고, 쓰고, 실행할 수 있는 웹
콘텐츠	하이퍼링크와 이미지를 활용한 텍스트 기반 페이지	비디오, 오디오, 애니메이션 및 양방향 요소를 갖춘 멀티미디어 페이지	메타데이터, 자연어처리 및 스마트계약 등을 지원하는 시맨틱 페이지
사용자	제한된 소스의 수동적 정보 소비자	다양한 소스에서의 정보의 생산자와 소비자	P2P 네트워크 내 정보의 참여자와 소유자
예시	Yahoo!, Britanica Online, Geocities	구글, 위키피디아, YouTube, Facebook, Twitter	이더리움, IPFS ²⁾ , Brave Browser ³⁾ , Steemit ⁴⁾

출처: simpilearn.com, techgart.com 외

▶ 웹 3.0은 탈중앙화(decentralization), 지능화(intelligence), 사용자 중심성(user-centricity) 등을 주요 특징으로 함

• **(탈중앙화)** 웹 3.0은 웹에서 데이터, 콘텐츠 및 서비스를 제어하거나 영향을 미치는 중앙 집중식 플랫폼 또는 중개자에 대한 의존도를 줄이고 최소화시키는 특징을 지님

- P2P 네트워크와 분산 시스템을 통해 중개자 없이도 사용자가 자신의 데이터와 자산에 대한 소유·제어·공유가 가능

• **(지능화)** 웹 3.0은 AI와 머신러닝을 활용하여 사용자에게 초개인화(hyper-personalized)되고 연관성이 높은 적응형(adaptive) 웹 경험을 제공

- 메타데이터, 온톨로지(Ontology), 자연어 처리, 스마트 계약(Smart contract) 등과 같은 시맨틱웹 기술을 통해 기계와 인간이 웹 데이터를 보다 쉽게 이해하고 실행하도록 지원

• **(사용자 중심성)** 웹 3.0의 사용자는 웹 콘텐츠와 서비스를 직접 생성하고 소비하는 데 보다 적극적으로 참여

2) InterPlanetary File System. 웹 콘텐츠 접근을 위한 기존 HTTP 프로토콜을 대체한 분산형 파일 공유 프로토콜

3) 블록체인 기반 디지털 ID, 암호화화폐 지갑 등의 웹 3.0과 IPFS를 통합한 프라이버시를 강조한 웹브라우저

4) 블록체인 기반 분산 네트워크상에서 운영되는 소셜 미디어 플랫폼으로 콘텐츠 생성 및 큐레이팅에 기여하는 참여자에 대한 보상 지급

- 사용자의 관심도, 참여도, 창의성 등에 대하여 토큰으로 보상을 얻을 수 있는 토큰 경제를 기반으로 사용자는 자신의 기여도에 상응하는 혜택을 누릴 수 있음

▶ 이러한 웹 3.0의 핵심적인 특징은 개인정보보호와 관련하여 여러 과제와 이슈를 제기

2. 웹 3.0 주요 특징별 개인정보보호 이슈

(1) 탈중앙화⁵⁾

① 개인정보 유출

- ▶ 탈중앙화 시스템은 분산 네트워크의 여러 노드 또는 컴퓨터에 데이터를 저장하므로 적절하게 암호화하거나 보호하지 않을 경우 데이터가 노출되거나 유출될 위험이 존재
 - 탈중앙화 아키텍처는 데이터의 기밀성을 유지할 수 있지만, 메타데이터 분석 시 데이터에 대한 안전한 보호가 어려움
 - 따라서 개인정보보호 및 자율성을 증진하기 위한 탈중앙화 인프라가 제대로 설계되지 않을 경우, 오히려 중앙집중형 인프라보다 정부나 기업의 감시에 훨씬 더 취약할 수 있음
 - 일례로 IPFS는 누구나 고유 해시 식별자로 콘텐츠에 액세스할 수 있는 분산형 파일 공유 프로토콜이지만, IPFS에 접속하는 기기는 노드의 로컬 캐시가 삭제되기 전까지 한시적으로 콘텐츠 호스트 역할을 하게 되는데, 이때 네트워크에서 수행된 모든 트랜잭션의 기록을 검색하고 민감한 정보 등에 대한 분석이 가능

② 데이터 보안 및 무결성

- ▶ 탈중앙화된 분산 네트워크에 저장된 개인정보는 허가되지 않은 접근, 수정, 삭제 등으로부터 보호되어야 함
 - 이러한 데이터 보안 및 무결성 관련 문제는 암호화, 해싱(hashing), 디지털 서명, 합의(consensus) 알고리즘을 활용하여 해결을 도모
 - (암호화) 암호 해독 키를 가진 권한 있는 당사자에게만 데이터 접근을 허용
 - (해싱) 데이터의 고유한 지문을 비교하여 데이터의 진위 여부와 무결성을 확인
 - (디지털 서명) 공개 키 암호화를 사용하여 데이터 전송자를 통해 데이터 무결성을 검증

5) Primavera De Filippi, The Interplay between Decentralization and Privacy: The Case of Blockchain Technologies, 2016.9.14.

- (합의 알고리즘) 네트워크의 모든 노드가 특정 규칙에 따라 데이터를 동기화하고 전체의 동의를 이끌어 내도록 보장

③ 개인정보보호 및 기밀성

- ▶ 공공 원장(public ledger)에서 공유되는 개인정보가 원치 않는 당사자에게 노출되거나 의도하지 않은 목적으로 사용되지 않도록 하기 위해서는 개인정보보호 및 기밀성이 담보되어야 함
- 개인정보보호 및 기밀성을 위하여 영지식 증명, 차등 개인정보보호, 동형 암호화, 안전한 다자간 연산 등 여러 개인정보보호 기술(PET, Privacy Enhancing Technologies) 활용이 가능
 - (영지식 증명, Zero-Knowledge Proof) 증명자가 자신이 알고 있는 지식과 정보를 공개하지 않으면서 그 지식을 알고 있다는 사실을 증명하는 방식
 - (차등 개인정보보호, Differential Privacy) 사용자가 검색 결과의 근사치를 수집함으로써 특정 데이터의 역추적을 통한 개인 신상 확인이 불가능하도록 한 기술로, 개별 기록을 공개하지 않고도 데이터에 대한 통계 정보 공유가 가능
 - (동형 암호화, Homomorphic Encryption) 사용자는 암호화된 데이터를 해독하지 않고도 암호화된 데이터에 대한 계산을 수행

④ 규제 준수

- ▶ 탈중앙화 시스템에서는 국경 간 데이터 이동이 활발할 수 있으므로 다양한 지역과 국가의 개인정보보호 법률과 규정 준수 이슈가 부각
- GDPR은 개인정보보호를 위해 접근, 수정, 삭제, 제한 등 정보주체의 권리를 엄격히 보장하고 있으나, 명시적으로 데이터 컨트롤러나 프로세서가 존재하지 않는 탈중앙화 시스템에서는 이러한 권리 보장을 시행하기 어려울 수 있음⁶⁾

표 _ 탈중앙화 시스템에서의 GDPR 이슈

구분	주요 내용
데이터 컨트롤러 및 프로세서 식별	데이터 처리의 여러 단계에 여러 행위자가 관여할 수 있으므로 개인정보 컨트롤러 또는 프로세서 파악이 어려움
동의 획득	누가 자신의 데이터를 처리하는지 또는 어떤 목적으로 처리하는지 파악이 어려우므로 사용자 동의 획득이 어려움

6) Heleen Janssen, et al, Decentralized data processing: personal data stores and the GDPR, 2020.12.28.

구분	주요 내용
권리 행사	사용자가 개인정보 관련 컨트롤러 또는 프로세서와 접촉하거나 요청이 이행되었는지 확인하기 어려우므로 이러한 권리 행사도 제약
보안 보장	데이터의 저장 또는 전송에 사용되는 기반 기술이나 프로토콜과 관련된 취약점이나 위험이 존재

출처: Heleen Janssen(2020) 외

(2) 지능화⁷⁾

① 데이터 수집

- ▶ AI 기반 지능화 시스템은 알고리즘과 모델을 학습하고 개선하기 위해 대량의 데이터를 필요로 하는데, 이로 인해 웹 3.0 생태계의 다양한 주체들이 데이터를 수집, 저장, 처리 및 공유하는 방식에 대한 우려가 제기
- 사용자는 AI 시스템의 데이터 수집 절차를 인지하지 못하거나 그러한 절차에 동의하지 않을 수 있으나, 상업적 활용을 포함하여 개인정보의 활용을 통제하기 어려움

② 데이터 분석

- ▶ AI 시스템은 자연어 처리, 머신러닝, 컴퓨터 비전 등의 첨단 기술을 활용하여 데이터를 분석한 후 이를 바탕으로 예측이나 시사점을 도출하므로 편향적으로 수집된 개인정보에 의한 해석 오류나 조작 등 기타 부정적인 영향이 나타날 수 있음
- 즉, 사용자는 AI 시스템이 수행한 분석의 정확성, 신뢰성 또는 편향성을 검증할 수 없거나 AI 시스템이 내린 결정에 대해 이의를 제기하거나 수정이 불가할 수 있음

③ 데이터 기밀성

- ▶ AI 시스템은 허가받지 않은 접근에 대응하여 개인정보를 보호하기 위해 익명화, 차등 개인정보보호 등 고도화된 다양한 기술을 활용할 수 있으나 모든 상황에서 완벽한 대안은 존재하지 않음
- 익명화된 데이터라도 다른 정보 소스와 연결하여 개인에 대한 재식별(de-identification)이 가능
 - 익명화 또는 비식별처리는 일부 정보를 제거하고 수정하는 작업을 수반하는 분석이나 연구 목적으로 활용하는 데이터의 유용성과 정확성이 떨어질 수 있음

7) Primavera De Filippi, The Interplay between Decentralization and Privacy: The Case of Blockchain Technologies, 2016.9.14.

- 차등 개인정보보호는 개인에 대한 정보가 노출되지 않도록 데이터 또는 쿼리 결과에 무작위 노이즈를 추가하는 기법을 활용하는데, 개인정보보호와 정보의 효용성 사이에서 적절한 개인정보보호 수준을 설정하기 어렵고 고도의 기술적 전문성과 신중한 설계를 필요로 함⁸⁾
 - 차등 개인정보보호는 기술적인 어려움으로 인해 사용자와 이해관계자에게 해당 기술에 대한 충분한 설명이 필요
- 이외에도 추론, 적대적 학습 등의 공격에 대한 지속적인 대응 기법 개발을 통해 데이터 기밀성과 무결성을 강화해야하는 과제가 존재

(3) 사용자 중심성

① ID 관리⁹⁾

- ▶ 웹 3.0에서 사용자 중심성을 구현하기 위해서는 사용자가 자신의 신원(ID)에 대한 강화된 통제권을 갖는 것이 중요
- 그러나 다수의 플랫폼과 서비스에서 사용자에게 대한 추적이 가능하기 때문에 ID 도용이나 사기, ID 유출에 따른 피해 구제에 제약이 있으며, 개인정보의 국가 간 이동으로 여러 규제에 영향을 받게 되는 등 개인정보보호와 관련한 여러 위험이 존재

표 _ 사용자 ID 관리와 관련된 개인정보보호 이슈

구분	주요 내용
ID 도용 및 사기	<ul style="list-style-type: none"> • 사용자는 개인 키 또는 비밀번호를 분실하거나 해킹 또는 피싱 공격으로 인해 자산의 디지털 ID에 대한 접근 권한을 잃게 될 수 있음 • 또한 악의적인 행위자에 의해 디지털 ID가 사칭되거나 복제될 수 있어 ID 도용 또는 사기 위험에 노출
제한적인 피해 구제책	<ul style="list-style-type: none"> • 사용자의 디지털 ID와 관련된 분쟁, 갈등 또는 손해가 발생할 경우, 탈중앙화된 웹 3.0의 특성상 중재 조직의 부재로 인해 구제 수단이나 보호 조치가 제한적일 가능성 농후
규정 준수	<ul style="list-style-type: none"> • 웹 3.0 환경에서는 사용자 개인정보가 국경을 넘어 다양한 국가에서 이용되므로 개인정보의 보호, 보안 및 활용 동의와 관련하여 다수의 법률 및 규정에 대한 고려가 필요 • 또한 사용자는 자신의 개인정보보호 선호도와 특정 거래 또는 상호 작용에 대한 신원 공개 또는 확인의 법적 의무 사이에서 적절한 균형 유지가 필요

출처: Jessica Groopman(2023) 외

8) NIST, Differential Privacy for Privacy-Preserving Data Analysis: An Introduction to our Blog Series, 2020.7.27.

9) Jessica Groopman, Web 3.0 security risks: What you need to know, 2023.2.

② 사용자 교육¹⁰⁾

- ▶ 사용자에 대하여 웹 3.0 서비스 이용의 이점과 함께 개인정보보호와 관련된 위험 및 책임 등에 관한 교육이 필요
- **(서비스의 복잡성)** 웹 3.0의 작동방식, 안전하고 효과적인 사용법, 개인정보 및 신원(ID) 보호 방법 등은 이해하기 어려울 수 있음
 - 개인정보보호 표준이나 요구사항이 상이한 다양한 플랫폼, 애플리케이션, 프로토콜과 상호작용하는 과정에서 일관성 결여는 사용자에게 혼란을 유발
- **(낮은 신뢰성)** 서비스의 복잡성으로 인해 사용자는 웹 3.0 기술을 불신하게 되고, 서비스 이용에 저항감을 가질 수 있음
- **(책임과 의무)** 개인 키나 디지털 지갑에 대한 접근 권한을 상실하거나, 해킹·피싱 공격에 노출될 경우의 책임과 위험부담 등에 대해 올바른 이해가 필요
 - 사용자에게 개인정보 유출에 따른 손해 발생 시의 구제 수단이나 보호조치 등에 대하여 사전에 주지시켜야 함

③ 사용자 통제

- ▶ 웹 3.0 환경에서 사용자는 개인정보에 대해 강화된 통제권을 갖게 되지만, 사용자가 주도적으로 통제권을 활용하기 어려워 개인정보 오남용이 발생할 수 있음
- ▶ **(사용자 보안 및 책임)** 웹 3.0 플랫폼에서 사용자는 중앙집중형 중개자나 기관에 의존하지 않고 스스로 자신의 개인정보, 지갑, 키 등을 안전하게 관리해야 하는 과제에 직면
 - 즉, 사용자는 직접 강력한 비밀번호 생성, 키 백업, 피싱 공격 방지, 거래 확인 등 보안 및 개인정보보호와 관련된 다양한 조치를 이해하고 실행해야 함
 - 탈중앙화 시스템은 중개자 없이 사용자가 자신의 데이터를 소유·제어·공유할 수 있는 P2P 네트워크와 분산원장에 의존하고 있어 사용자 개개인에 자신의 데이터를 관리하고 보안 및 개인정보보호를 보장해야하는 책임이 귀속
 - 예를 들어, 블록체인 기반 시스템은 공개키 암호화로 거래와 신원을 확인하여 보안 위협에 대응하지만, 해당 시스템의 사용자는 해커의 공격에 노출되거나 공개키를 분실하지 않도록 주의를 기울여야 함
- ▶ **(사용자 익명성 및 책임)** 웹 3.0 플랫폼에서 사용자는 익명성(anonymity)과 가명성(pseudonymity)을 유지할 수 있으나, 이를 통해 얻을 수 있는 장점과 제반 위험 사이에서 균형이 필요

10) Forbes, Web 3.0: How To Prepare For A Privacy-Driven Future, 2021.10.28.

- 익명성과 가명성은 사용자가 자신의 실제 신분을 드러내지 않고 웹 3.0 공간 내에서 타인들과 상호작용할 수 있게 함으로써 사용자의 프라이버시와 표현의 자유를 향상시킴
- 반면 이러한 익명성과 가명성은 사이버 범죄, 테러, 혐오 발언, 잘못된 정보와 같은 악의적인 활동을 가능하게 하고, 가해자를 추적해 책임을 묻기 어렵게 만들 수 있음
- ▶ **(개인정보보호 규정)** 각국의 개인정보보호 규정은 개인정보 컨트롤러와 프로세서에게 개인정보보호 의무를 부과하여 허가받지 않은 접근이나 사용 또는 공개로부터 사용자 개인정보를 보호하고 있음
- 그러나 웹 3.0 플랫폼에서는 분산화된 특성과 글로벌 규모의 서비스로 인해 공공 원장 상에서 개인정보 주체의 권리 행사에 대한 규정이 모호하거나 국가 간 관련 규정이 상이해 침해가 발생할 가능성이 존재
- ▶ **(데이터 공유 촉진)** 웹 3.0 플랫폼 상에서 개인의 자발적 데이터 공유는 연구, 혁신, 개인화 등 사회적 공익이나 경제적 이익 차원에서 다양한 형태의 응용이 가능하나, 재식별, 프로파일링, 차별 등과 같은 위험을 수반¹¹⁾
- 웹 3.0에서 사용자 개인정보보호와 동시에 데이터 공유를 촉진하기 위하여 암호화, 영지식 증명, 차등 개인정보보호와 같은 PET 기술 활용에 대한 검토가 필요
- 한편 웹 3.0 플랫폼은 토큰, 스마트 계약, 평판 시스템 등과 같은 메커니즘을 활용하여 사용자의 선호도와 동의 여부에 따라 데이터 공유에 대한 보상을 제공할 수 있음

3. 요약 및 결론

- ▶ 웹 3.0은 시맨틱, 탈중앙화, 인텔리전스 기반 개인화, 사용자 중심성 등을 특징으로 하며 현재 관련 논의가 활발
- 이러한 일련의 특징들은 상호 직접적인 연계 속에서 통일성 있는 발전을 이뤄왔다기 보다는 기술 진화 양상이 웹과의 실험적 결합을 통해 다양한 양상으로 변화하고 있으며, 이 점에서 웹 3.0은 개념 자체가 추상적이면서 동태적인 특성을 보유
- ▶ 웹 3.0의 주요 특징 중 지능화와 관련해서는 AI 시스템에서 일반적으로 언급될 수 있는 수준의 개인정보보호 이슈들에 주목해야 함
- 향후 웹 3.0 관련 생태계 내 다양한 주체들의 참여가 예상됨에 따라 이와 관련된 데이터 수집 활동에 대한 효율적인 통제체계 수립 논의가 선행될 필요가 있음

11) Jessica Groopman, Web 3.0 security risks: What you need to know, 2023.2.

- ▶ 분산 네트워크를 기반으로 한 탈중앙화의 특징은 웹 3.0 고유의 개인정보보호 관련 이슈를 제기
 - 웹 3.0에서는 중앙의 플랫폼 관리 주체나 별도의 데이터 컨트롤러 및 프로세서가 존재하지 않기 때문에 GDPR 등 개인정보보호 규정에서 정한 정보주체의 권리 보장이 취약
 - 또한 IPFS와 같이 개별 기기가 노드가 될 수 있는 웹 환경에서는 분산 네트워크 내 모든 참여자들에게 콘텐츠에 대한 접근성이 개방되므로 개인정보보호 측면에서 인프라 설계 기술 개발을 위한 산업계의 관심이 요구됨
 - 이에 더하여 전 세계를 아우르는 네트워크는 국경 간 데이터 이동을 전제로 하는 만큼 관련 규정 준수와 관련하여 복잡성을 확대
- ▶ 웹 3.0은 중앙의 통제를 받지 않는 특성으로 인해 개인정보보호를 위해 사용자 개인의 통제 수단과 조치에 크게 의존하고 있음
 - 따라서 사용자의 정보보호 인식을 제고하고 개인정보 유출에 따른 위험 및 책임 등에 관한 적극적인 교육과 홍보를 실시하는 등 사전예방 노력을 기울이는 것이 중요

표 _ 웹 3.0의 주요 특징별 개인정보보호 이슈 요약

구분	주요 특징	개인정보보호 이슈
탈중앙화	<ul style="list-style-type: none"> • 데이터가 암호화 및 개인 키를 통해 피어(peer) 간에 분산됨에 따라 중앙 집중형 조직에 의한 데이터 유출 또는 오용의 위험을 회피 • 사용자의 데이터 접근·공유 권한 강화 	<ul style="list-style-type: none"> • 분산 네트워크에서 개인정보 유출 위험 • 데이터 보안·무결성, 개인정보보호·기밀성 등을 위한 개인정보보호 기술(PET) 적용 • 규제 준수
지능화	<ul style="list-style-type: none"> • 데이터가 구조화되고 의미론적 주석이 추가되어 보다 정확하고 관련성 높은 쿼리 및 응답이 가능 • 사용자는 불필요한 정보를 노출하지 않고 어떤 데이터를 공개하거나 확인할 지 선택할 수 있음 	<ul style="list-style-type: none"> • 웹 3.0 생태계 내 다양한 주체들의 데이터 수집 활동에 대한 통제 어려움 • AI 시스템이 취합한 개인정보의 편향성 극복 • AI 시스템용 데이터의 기밀성 강화
사용자 중심성	<ul style="list-style-type: none"> • 중앙 조직의 개입 없이 사용자가 직접 디지털 ID를 생성, 소유, 통제 및 검증할 수 있는 자기주권적 신원 관리(SSI, Self Sovereign Identity) 실현 • 사용자가 온라인에서 자신의 신원을 보호하기 위해 가명 또는 익명성 사용 	<ul style="list-style-type: none"> • ID 관리 • 사용자 통제권 이슈 • 사용자 교육 필요성

출처: 넥스텔리전스(주)

Reference

1. Forbes, Web 3.0: How To Prepare For A Privacy-Driven Future, 2021.10.28.
2. Heleen Janssen, et al, Decentralized data processing: personal data stores and the GDPR, 2020.12.28.
3. Jessica Groopman, Web 3.0 security risks: What you need to know, 2023.2.
4. NIST, Differential Privacy for Privacy-Preserving Data Analysis: An Introduction to our Blog Series, 2020.7.27.
5. Primavera De Filippi, The Interplay between Decentralization and Privacy: The Case of Blockchain Technologies, 2016.9.14.

사우디아라비아의 데이터·프라이버시 규제 샌드박스 추진현황

[요약]

1. 추진배경과 목적

사우디아라비아 데이터·인공지능청(SDAIA)이 새로운 개인정보보호법의 시행을 맞아 개인 정보를 다루는 중소기업과 스타트업의 준법 이행을 지원하기 위함

2. 데이터·프라이버시 규제 샌드박스의 추진현황

- (1) **(도입 목표)** ▲사우디아라비아의 데이터와 프라이버시 규제 강화 ▲혁신적 비즈니스 모델과 서비스의 지원 ▲설계 단계부터 프라이버시를 고려한 제품과 서비스의 개발
- (2) **(기대 효과)** 참여 기업은 제품과 서비스의 준법 여부 확인 및 SDAIA의 지원을 확보하고 규제당국은 혁신 증진과 시장 경쟁 활성화 효과를 기대할 수 있음
- (3) **(지원 대상)** 데이터와 프라이버시 관련 제품이나 서비스를 보유한 사우디 소재 기업
- (4) **(평가 기준)** ▲사우디의 혁신 지원 ▲사우디의 사회경제 환경 및 사우디 비전 2030에 대한 가치 증대 ▲고객 중심의 솔루션 ▲비즈니스 모델의 준비도, 테스트 계획 등
- (5) **(진행 절차)** 신청서 제출 → 평가 → 온보딩 → 테스트 → 종료의 다섯 단계로 진행
- (6) **(기타 사항)** 샌드박스 참여 시 관련 규제에 대한 지침과 지원 제공 및 규제 유예 가능

3. 시사점

규제당국과 기업의 협력을 통해 데이터와 프라이버시를 활용한 신규 서비스 도입 과정에서 개인정보 침해의 위험을 최소화할 수 있을 것을 기대

1. 추진배경과 목적

- ▶ 사우디아라비아 데이터·인공지능청(SDAIA)*이 데이터의 가치를 활용하는 동시에 소비자의 개인정보 권리를 보호하기 위한 데이터·프라이버시 규제 샌드박스를 출범('23.2.6.)

* Saudi Data and Artificial Intelligence Authority: 무하마드 빈 살만 왕세자의 직속 기구로 2019년 설립되어 사우디의 디지털 전환을 주도

- 이 샌드박스는 데이터와 프라이버시 보호 규정을 준수하면서 혁신과 실험 허브로서 사우디아라비아를 자리매김하기 위한 역내 최초 시도로서, 현지 기업들이 개인정보보호 규제가 자사의 제품과 서비스에 미치는 영향을 시험할 수 있도록 지원
- ▶ 사우디아라비아는 '21.9월 개인정보보호법(PDPL)*을 처음 공포했으며, 원래 '22.3월 시행 예정이었으나 1년 연기하여 올해 3월 17일부터 시행
 - * Personal Data Protection Law
- 사우디아라비아 내 개인정보를 처리하는 모든 기업에 적용되는 PDPL의 목적은 개인정보와 프라이버시를 보장하고 개인정보 공유를 규제하며 개인정보의 남용을 방지하기 위함
- PDPL은 목적 제한 및 데이터 최소화, 데이터 처리기록의 등록과 유지관리 등의 컨트롤러의 의무와 정보주체의 권리 및 법률 위반 시의 처벌 등을 포괄하며, 특정 경우를 제외하고는 개인정보의 국외이전을 금지
- PDPL 발효로부터 처음 2년간 법률 집행을 담당하는 관할기관은 SDAIA이며, 이후 데이터 보호 환경의 발전에 따라 SDAIA 산하의 국가데이터관리기관(NDMO)*으로 이관될 예정
 - * National Data Management Office
- ▶ SDAIA는 PDPL 시행을 앞두고 애드테크(AdTech)나 레그테크(RegTech) 또는 프라이버시 강화 기술(PET) 솔루션을 사업모델로 하는 영세기업과 중소기업 및 스타트업의 준법을 지원하기 위해 규제 샌드박스를 도입
- SDAIA는 규제 샌드박스를 통해 개인정보를 다루는 기업들이 서비스 출시에 앞서 데이터와 프라이버시 보호를 우선할 수 있도록 제품과 서비스의 최적화에 필요한 지침과 자문, 전문 지식을 제공할 계획

2. 데이터·프라이버시 규제 샌드박스의 추진현황

(1) 도입 목표

- ▶ 사우디아라비아는 석유 중심의 산업구조를 탈피하기 위한 국가 혁신 전략 '사우디 비전 2030'을 통해 디지털 경제발전을 추진하고 있으며, SDAIA는 데이터의 가치를 포착하고 극대화하기 위한 데이터 AI 국가 전략(NSDAI)을 수립('20.10월)
 - * National Strategy for Data & AI
- ▶ 국가데이터관리기관(NDMO)은 데이터·프라이버시 규제 샌드박스 도입을 통해 데이터의 가치를 극대화한다는 SDAIA의 비전을 더욱 발전시키고, 혁신 기업과 규제당국이 협력해 한정된 기간에 제품과 서비스를 테스트하기 위한 안전한 환경을 조성하는 것을 목표로 함

- 규제 샌드박스의 전략적 목표는 ▲사우디아라비아의 데이터와 프라이버시 규제 강화 ▲혁신적 비즈니스 모델과 서비스의 지원 ▲설계 단계부터 프라이버시를 고려한 제품과 서비스의 개발임

(2) 기대 효과

표 _ 사우디아라비아 데이터·프라이버시 규제 샌드박스의 기대 효과

구분	주요 내용
참여 기업	<ul style="list-style-type: none"> • SDAIA가 보유한 전문가 및 생태계 협력사들에 직접 접근해 자문과 지원 확보 • 개인정보보호법(PDPL)과 관련 규제를 포함한 사우디아라비아의 규제수단에 대한 가이드 확보 • 규제 샌드박스의 범위 안에 있는 규제 수단에 대응해 제품이나 서비스의 준법 여부를 테스트 • PDPL과 NDMO 규제와 관련해 규제 유예의 가능성
규제당국	<ul style="list-style-type: none"> • 사우디의 디지털 역량을 극대화하기 위한 국가적 자산으로서 데이터를 활용 • 규제 수단에 개선이 필요한 영역을 확인하기 위한 중요한 통찰을 확보 • 규제당국의 감독 아래 설계 단계부터 프라이버시를 고려한 제품 개발을 목표로 중소기업의 신규 솔루션 탐색을 통해 국가적 혁신을 증진 • 중소기업의 경제 기여 확대로 시장 경쟁 활성화
소비자	<ul style="list-style-type: none"> • 혁신적 제품과 서비스의 데이터 프라이버시, 안전 및 신뢰 수준 제고 • 테스트 건본집단이 되어 피드백을 제출함으로써 혁신 솔루션의 개발에 참여 • 사우디아라비아의 데이터 프라이버시 요구사항에 부합하는 기업들의 혁신적 제품과 서비스 향유

출처: SDAIA 자료를 토대로 넥스텔리전스 정리

(3) 지원 대상

▶ 규제 샌드박스에 참여하기 위해서는 아래 조건을 모두 충족해야 함

- 사우디아라비아 소재의 기업이나 기업가
- 프라이버시 강화 기술(PET)을 포함한 데이터와 프라이버시 관련 솔루션/서비스/비즈니스 모델을 보유
- 지원자가 제안한 활용사례가 PDPL 및 관련 현행 법률의 범위 안에 속하여 해당 법률에 대응한 테스트가 필요해야 함
- 참가자는 개발 단계에 있는 혁신적인 기술 솔루션을 시연해야 하며, 해당 솔루션은 명확한 활용사례를 갖추고 규제 샌드박스에서 테스트를 진행할 수 있어야 함 (구상이나 설계 단계의 솔루션은 참여 불가능)

(4) 평가 기준

- ▶ 규제 샌드박스의 지원자는 코호트 단위로 허가되며, 지원자의 평가 기준은 ▲사우디의 혁신 지원 ▲사우디의 사회경제 환경 및 사우디 비전 2030에 대한 가치 증대 ▲고객 중심의 솔루션 ▲비즈니스 모델의 준비도, 테스트 계획, 출구 전략 등임
- 규제 샌드박스에 참여하는 솔루션은 사우디 내 상용화된 제품이나 서비스와 차별화되는 제품이나 솔루션, 비즈니스 모델이어야 함
- 사우디의 경제와 비전 2030 달성에 긍정적인 영향을 미쳐야 하며, 사용자의 고객 경험이나 삶의 질을 개선해야 함
- 지원자는 구체적인 비즈니스 모델과 테스트 사례 및 방법론, 주요 이정표 등 실행방안을 포함한 테스트 계획과 명확한 출구 전략을 제시해야 함

표 _ 규제 샌드박스 신청 기업의 평가 기준(예시)

구분	참여 가능 사례	참여 불가 사례
혁신 지원	<ul style="list-style-type: none"> 제안한 비즈니스 모델이 시장에 없다는 시장조사 결과를 제시 샌드박스 팀이 시장 내 유사한 제품을 확인하지 못함 	<ul style="list-style-type: none"> 시장 내 대체 가능한 솔루션이 다수 존재 지원자가 비즈니스 모델의 부가 가치를 입증 불가
사회경제 환경 및 사우디 비전 2030에 대한 가치 증대	<ul style="list-style-type: none"> 제안 솔루션이 경제에 미치는 영향을 예상한 연구 결과를 제시 지원자가 제품 출시와 운영의 결과로 창출되는 예상 일자리 수를 제시 비즈니스 모델에 해당 산업 부문의 예상 성장치를 제시 	<ul style="list-style-type: none"> 지원자가 보유 솔루션의 사회경제적 영향에 대한 뚜렷한 비전을 제시하지 못함
고객 중심의 솔루션	<ul style="list-style-type: none"> 제안 솔루션이 고객 경험을 개선한다는 증거를 제시 샌드박스 팀이 데스크 리서치와 분석을 통해 제안 솔루션이 고객의 프라이버시를 개선한다고 판단 	<ul style="list-style-type: none"> 지원자가 제안한 솔루션이 고객에게 가치를 더한다는 충분한 증거를 제시하지 못함 샌드박스 팀이 제안 솔루션이 소비자 경험에 부정적 영향을 미칠 가능성을 확실히 배제할 수 없음
비즈니스 모델의 준비도, 테스트 계획, 출구 전략	<ul style="list-style-type: none"> 지원자가 샌드박스 내 테스트를 위한 적절한 예산이 할당되었음을 입증 가능 지원자가 명확한 테스트 기준과 예상 결과를 포함한 구체적인 테스트 시나리오를 제시 지원자가 제품 출시와 연속성에 대한 명확한 로드맵을 갖고 출구 전략을 제시 	<ul style="list-style-type: none"> 지원자가 불충분한 테스트 계획이나 비즈니스 모델을 제시 지원자가 테스트를 위한 충분한 예산을 할당했음을 입증하지 못함 지원자가 샌드박스 이후 제품 출시할 준비가 되었다는 충분한 증거를 제시하지 못함

출처: SDAIA 자료를 토대로 넥스텔리전스 정리

(5) 진행 절차

- ▶ 규제 샌드박스는 신청서 제출 → 평가 → 온보딩 → 테스트 → 종료의 다섯 단계로 진행되며, 2월 6일부터 3월 31일까지 신청서를 제출받음

그림 _ 규제 샌드박스의 진행 절차



출처: SDAIA, Data and Privacy Regulatory Sandbox Guidelines, 2023.2.6.

- ▶ **(신청)** 지원자는 코호트 접수 발표 후 30일 이내 SDAIA/NDMO 포털에서 제공되는 지원 양식을 작성해 비즈니스 모델을 뒷받침하는 보충 자료와 함께 제출
- 샌드박스 팀은 규정된 자격 기준에 따라 신청서와 제출된 자료를 검토하여 지원자에게 결과를 통보하고 합격자는 평가 단계로 이동
- ▶ **(평가)** 샌드박스 팀은 상기 평가 기준에 따라 구체적인 평가를 진행하며, 이 단계는 최대 60일까지 소요될 수 있음
- 이 단계에는 규제 유예 영역의 특징을 포함해 필요한 샌드박스 도구의 확인 작업도 포함되며, 합격자에게는 임시 허가서와 온보딩 초대장이 제공됨
- ▶ **(온보딩)** 샌드박스 팀은 합격한 지원자에게 온보딩 교육을 제공하며 제출한 서류에 대한 자세한 검토를 통해 테스트 단계를 준비하며, 이 단계는 최대 60일까지 소요될 수 있음
- 온보딩 과정에는 행정 서류 서명, 담당자 배정 등의 온보딩 관련 운영 활동, 샌드박스과 PDPL 및 관련 규정에 대한 설명을 제공하는 Q&A 세션, 테스트 계획과 출구 전략의 마무리 지원 등이 포함

- 샌드박스 팀은 온보딩에 필요한 행정 서류를 처리하고 구체적인 참여 약관이 담긴 최종 허가서를 제공하며, 지원자가 제출한 문서가 기준에 충족되지 않거나 테스트 준비가 부족하다는 결론이 나오면 최종 허가서 발급을 거부할 수 있음
- ▶ **(테스트)** 이 단계에서 참가자는 승인된 테스트 계획에 따라 최대 6개월 동안 솔루션을 테스트할 수 있음
 - 참가자는 샌드박스 팀에게 정해진 보고 양식에 따라 정기적으로 상태 업데이트를 제공해야 하며, 테스트 결과 기록을 유지하고 최종 허가서에 나온 약관에 의거해 데이터를 수집해야 함
 - 샌드박스 팀은 제출받은 정기 보고서를 검토해 테스트 결과를 감독 및 평가하고 사전에 합의된 테스트 운영 KPI의 준수 여부를 확인
 - 샌드박스 팀은 참가자가 최종 허가서의 약관을 준수하지 않았다고 판단되는 경우, 추가 분석과 조사 시까지 테스트 중단을 요구할 수 있으며, 약관의 심각한 위반이 확인된 경우 테스트를 무기한 중단할 수 있음
 - 참가자가 예상치 않은 기술적·비기술적 상황으로 인해 테스트를 중단해야 하는 경우 테스트 기한 연장을 요청할 수 있으며, 샌드박스 팀은 사안에 따라 요청을 검토
 - 테스트 기간은 최종 허가서 수령일로부터 6개월 후 종료되며, 참가자는 종료 시 테스트 계획에 적시된 KPI에 부합하는 전체 결과를 담은 보고서를 제출
- ▶ **(종료)** 샌드박스 팀은 종료 보고서와 이전에 설정된 KPI의 준수 여부를 평가하여, 참가자의 샌드박스의 성공적인 이수를 승인할 권리가 있음
 - 샌드박스 과정을 성공적으로 마친 경우, 참가자는 SDAIA 웹사이트의 성공한 참가자 명단에 등재됨
 - 테스트 관계가 종료된 이후 참여자는 온보딩 단계에서 수립한 출구 전략을 이행해야 하며, 출구 전략은 참가자의 결정에 따라 솔루션 개발 중단이나 제품의 생산 개시 등 다양한 결과를 도출할 수 있음

(6) 기타 사항

- ▶ 온보딩 단계에서 참가자들은 테스트에 필요한 샌드박스 도구와 양식을 요구할 수 있으며, 샌드박스 팀이 제공하는 도구의 종류는 아래와 같음
 - (보고 양식) 온보딩 단계 완료 시 참가자에게 제공

- (규제 지침 및 생태계 지원) 샌드박스 팀은 참가자에게 PDPL과 관련 규정에 대한 안내를 제공하며, 참가자 요청에 따라 SDAIA 생태계/네트워크에 속한 기타 기관과 참가자 간 연결을 지원
- (규제 유예) SDAIA와 NDMO는 규제 샌드박스 내의 테스트 기간 동안 참가자가 사우디 PDPL이나 NDMO의 규정을 일부 위반해도 책임을 묻지 않으며, 규제 유예 여부는 사안별로 샌드박스 팀의 재량에 따라 결정됨
- ▶ 참가자들은 샌드박스 절차의 일환으로 정보주체의 권리와 이익을 보호하기 위한 일련의 보호조치를 취해야 하며, 조치는 개별 테스트 시나리오에 따라 달라질 수 있으나 아래의 내용을 포함해야 함
- 규제 샌드박스에서 솔루션의 테스트에 참여하는 정보주체와의 합의(개인정보의 이용에 대한 명시적 동의 포함)
- 규제 샌드박스에서 참가자의 활동으로 인해 발생할 수 있는 정보주체의 손실에 대한 보상 약속
- 규제 샌드박스에서 참가자의 테스트 목적에 따라 개인정보를 제공해야 하는 정보주체의 수에 대한 구체적인 제한의 설정

3. 시사점

- ▶ 사우디아라비아에는 최근까지 개인정보 보호법이 부재하고 사이버범죄방지법, 전자상거래법 등 부문별 규정에 개인정보보호 조항이 산재해 있었으나 '19.8월 SDAIA의 설립과 함께 개인정보보호 규제에서도 많은 진전을 보임
- SDAIA는 '20.10월 데이터 AI 국가전략 수립과 함께 EU GDPR과 상당 부분 유사한 국가 데이터 거버넌스 임시규정*을 발표했으며, '21.9월에 데이터 컨트롤러·프로세서, 정보주체 등 GDPR의 핵심 개념을 포함한 PDPL이 공포되어 올해 3월 17일에 발효
- * National Data Governance Interim Regulations
- ▶ SDAIA는 PDPL의 시행과 함께 데이터와 프라이버시 관련 제품이나 서비스를 보유한 중소기업과 스타트업 등의 준법 지원과 해당 분야의 혁신 증진을 위한 규제 샌드박스를 출범했으며, 이는 역대 최초 시도이며 글로벌 차원에서는 유사한 사례가 있음
- 일례로 싱가포르에서는 '22.7월 정보통신미디어개발청(IMDA)와 개인정보보호위원회(PDPC)가 프라이버시 강화 기술(PET) 시범프로젝트를 지원하기 위한 규제 샌드박스를 출범

- 국내에서는 '20.8월 출범한 개인정보보호위원회가 개인정보보호 규제 유예 제도를 통해 규제 존재 여부의 신속 확인, 신기술 제품이나 서비스의 시험·검증을 허용하는 실증 특례, 일정 기간 임시로 서비스 운영을 허가하는 임시 허가 등으로 133건에 달하는 사안을 처리
- ▶ SDAIA의 규제 샌드박스는 현재 신청서를 접수 중인 단계로 성공 여부를 속단하기 어려우나, 규제당국과 기업의 협력을 통해 데이터와 프라이버시를 활용한 신규 서비스 도입 과정에서 개인정보침해의 위험을 최소화할 수 있을 것을 기대
- ▶ 우리나라에서도 9월 15일부터 이용자의 개인정보 전송 요구권 도입을 골자로 하는 개인정보 보호법 개정안의 시행을 앞둔 가운데, 규제 변화로 인한 신규 서비스의 준법을 지원하기 위한 규제 샌드박스 등의 제도 도입을 모색 가능

Reference

1. DataGuidance, Saudi Arabia: SDAIA launches data and privacy regulatory sandbox, 2023.2.10.
2. DataGuidance, EU - Saudi Arabia: GDPR v. Saudi Arabia's PDPL, 2022.3.
3. SDAIA, Data and Privacy Regulatory Sandbox, 2023.2.6.
4. SDAIA, Data and Privacy Regulatory Sandbox Guidelines, 2023.2.6.

개인정보보호와 활용성 강화 기술로 주목받는 재현데이터 기술의 특성과 활용 과제

[목 차]

1. 배경

2. 재현데이터의 부상

3. 개인정보보호 측면에서 재현데이터의 특성과 장단점

- (1) 재현데이터의 특성
- (2) 재현데이터 생성 방법
- (3) 재현데이터의 장점
- (4) 재현데이터의 단점

4. 결론 및 시사점

1. 배경

- ▶ 4차 산업혁명의 원동력으로서 '데이터'의 중요성이 커지면서 데이터 활용을 촉진하는 한편, 개인정보 활용으로 침해될 수 있는 개인정보 주체의 정보주권을 보장하고 균형을 확보하는 것이 주요 과제로 부상
- ▶ 재현데이터(synthetic data)는 실제 세계를 직접 관찰해 수집한 데이터가 아니라 인위적으로 생성한 데이터로, 개인식별정보(PII, Personally Identifiable Information)를 포함하지 않아 각종 인공지능(AI) 서비스 이용을 촉진하면서도 개인정보보호가 가능
- 재현데이터는 실제 데이터에서 통계적 특성을 파악해 모델을 만드는 통계적 방법이나 생성적 적대 네트워크(GAN) 등 여러 방법으로 생성

- ▶ 특히 ChatGPT 등 인공지능 기반 서비스가 빠르게 확산하면서 인공지능 서비스에 적용할 수 있는 개인정보보호 기술(PET, Privacy Enhancing Technologies) 중 하나로 재현데이터를 안전하게 이용할 수 있는 방안에 대한 관심과 연구가 활발
- ▶ 현재 보건의료, 사회 복지, 교통 및 물류, 교육, 금융 등 여러 산업 분야에서 재현데이터 활용이 확대
 - 특히 민감한 정보를 취급하는 건강 및 의료, 금융 부문을 중심으로 재현데이터 활용에 대한 기대감이 높고 활용이 활발

2. 재현데이터의 부상

- ▶ 시장조사기관 가트너(Gartner)는 2022년 발표한 '인공지능 하이프사이클(2022 Gartner Hype Cycle for Artificial Intelligence)'¹²⁾과 '개인정보보호 하이프사이클(Hype Cycle for Privacy 2022)'¹³⁾ 분석에서 모두 재현데이터를 최상위 관심 기술로 선정
- ▶ 가트너는 인공지능의 여러 관련 기술 중에서 재현데이터에 대한 관심이 최고조에 달하고 있으며, 향후 다양한 산업 분야에서 재현데이터 도입이 급증할 것으로 전망¹³⁾

그림 _ 가트너의 2022 인공지능 하이프사이클 분석



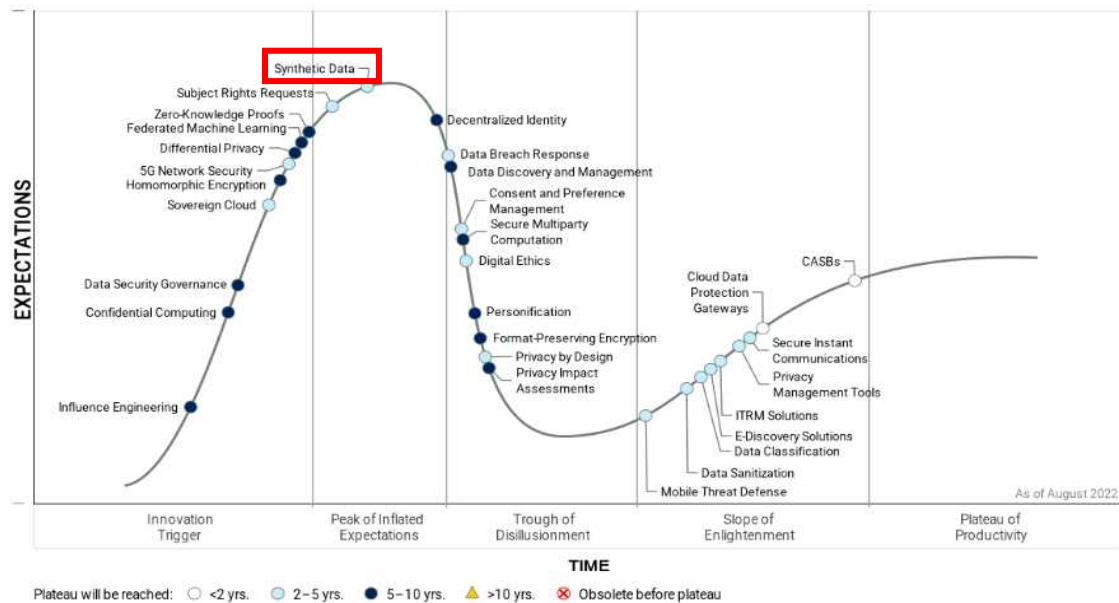
출처: Gartner, What is New in Artificial Intelligence from the 2022 Gartner Hype Cycle, 2022.9.

12) 가트너는 인공지능 등 신기술이나 주목해야 할 이슈를 중심으로 성장 단계 및 트렌드를 분석·예측하는 하이프 사이클 보고서를 매년 발간

13) Gartner, Hype Cycle for Privacy, 2022, 2022.8.

- ▶ 또한 2025년까지 대규모 조직의 60%가 분석, 인텔리전스 또는 클라우드 영역에서 개인정보보호 강화 컴퓨팅 기술(Privacy-enhancing computation techniques)을 사용할 것으로 예측¹⁴⁾
- 가트너는 개인정보보호 강화 컴퓨팅 기술 분야 중에서 재현데이터를 최고의 관심 분야로 선정하였고, 그 외 정보주체 권한 요청, 영지식 증명, 동형 암호화 등을 관심도가 높은 기술로 지목

그림 _ 가트너의 2022 개인정보보호 하이프사이클 분석



출처: Gartner, Hype Cycle for Privacy, 2022, 2022.8.

- 가트너는 재현데이터의 중요성을 강조하면서, 비즈니스적 영향, 촉진/장애 요인, 관련 권장사항, 주요 공급업체 등을 다음과 같이 분석해 제시

표 _ 재현데이터의 중요성과 특성

구분	주요 내용
재현데이터의 중요성	<ul style="list-style-type: none"> • AI 모델이 효과적으로 훈련될 수 있도록 실제 데이터를 획득하고 레이블을 지정하는 것은 AI 개발의 주요 과제 중 하나 • 그러나 실제 데이터를 획득해 활용하는 데에는 여러 부담이 수반되는데, 재현 데이터를 이용하면 시간과 비용을 크게 줄일 수 있음 • 또한 재현데이터는 개인식별정보(PII)를 효과적으로 제거
비즈니스 영향	<ul style="list-style-type: none"> • 컴퓨터 비전 및 자연어 응용 프로그램에서의 사용과 함께 다양한 산업 분야에서 재현데이터 채택과 활용이 급증할 전망 <ul style="list-style-type: none"> - 원본 데이터의 합성·변형 또는 데이터 일부의 합성 대체를 통해 머신러닝 모델 학습에서 PII 사용을 피할 수 있음 - 저렴한 비용으로 신속하게 필요한 데이터를 확보할 수 있으므로 머신러닝 개발에 필요한 비용 및 시간 절약이 가능 - 학습 데이터가 많을수록 학습 결과가 향상되므로 머신러닝 성능도 향상

14) Gartner, What is New in Artificial Intelligence from the 2022 Gartner Hype Cycle, 2022.9.

구분	주요 내용
촉진 요인	<ul style="list-style-type: none"> AI 학습 데이터에서 재현데이터로 개인정보를 보호할 수 있으므로 의료와 금융 분야를 중심으로 재현데이터에 대한 관심이 증가 재현데이터 수요에 부응하기 위해 신규 및 기존 공급업체들이 관련 제품을 출시하면서 공급업체 환경이 확장되고 재현데이터 채택이 촉진 재현데이터는 데이터 수익화, 외부 분석 지원, 플랫폼 평가 및 테스트데이터 개발 등 활용 영역이 점차 확장 시뮬레이션 기술 채택이 증가함에 따라 재현데이터의 활용이 가속화
장애 요인	<ul style="list-style-type: none"> 재현데이터에는 편향 위험이 존재하며, 개발이 어렵거나, 기존의 실제 데이터에 새로운 정보를 제공하지 못할 수 있음 데이터 품질이 데이터를 개발하는 인공지능 모델에 따라 달라질 수 있음 다른 데이터 파이프라인 도구와 함께 재현데이터 관련 기술을 사용하는 시기와 방법에 대한 명확한 결정이 어려움 재현데이터라 하더라도 조직의 민감한 정보를 노출할 위험은 여전히 존재 실제 데이터에 비해 열등하거나 가짜인 데이터로 인식될 수 있음
사용자를 위한 권장사항	<ul style="list-style-type: none"> 의료, 금융 등 여러 산업 분야에서 재현데이터를 사용할 때에는 규제에 주의를 기울이고 규칙을 준수 개인정보가 데이터셋에 포함되는 경우에는 원본 데이터의 재현 변형(synthetic variations) 또는 데이터 부분의 재현 교체(synthetic replacement) 등을 활용 샘플링 접근방식, 데이터 과학자 등을 활용하여 통계적 유효성과 재현데이터의 배포를 보장 기술 발전에 맞춰 전문 공급업체를 활용
주요 공급업체	<ul style="list-style-type: none"> Bitext, Datagen, Rendered.ai, Diveplan, Hazy, LeapYear, MOSTLY AI, Neuromation, Statics, Tonic

출처: Gartner, Hype Cycle for Privacy, 2022, 2022.8.

3. 개인정보보호 측면에서 재현데이터의 특성과 장단점¹⁵⁾

(1) 재현데이터의 특성

- ▶ 재현데이터는 인공적으로 생성된 데이터로, 실제 세계를 직접 관찰한 결과가 아니며 실제 데이터의 통계적 특성을 보존하는 것을 목적으로 하는 알고리즘에 의해 생성
 - 재현데이터는 운영 또는 생산 데이터를 모방하고, 수학적 모델을 검증하고, 머신러닝 모델을 교육하는 등 다양하게 활용
- ▶ 재현데이터의 생성은 다음과 같은 네 가지 특징적 요소들과 관련
 - ① **소스 데이터**) 재현데이터가 재현하는 통계적 특성을 가진 원본 데이터
 - 소스 데이터가 개인과 관련되면 '준식별자'(성별, 인종, 아동, 흡연 여부, 성적 취향 등) 또는 '직접 식별자'(얼굴 이미지, 유전적 프로필 등)와 같은 개인정보를 포함할 가능성이 높음

15) Data Guidance, International: Is synthetic data the future of privacy?, 2023.2.,
Data Guidance, International: What is synthetic data and how is it generated?, 2023.2.

- **(② 알고리즘)** 재현데이터를 생성하는 데 사용되는 통계 모델
 - 재현데이터를 생성하는 방법은 여러 가지가 있으나 현재는 딥러닝 기술을 포함한 인공지능/머신러닝 도구가 가장 일반적으로 활용
- **(③ 재현데이터)** 재현데이터 생성 모델을 바탕으로 생성된 데이터
- **(④ 개인정보보호 및 유용성 척도)** 소스 데이터에 포함된 개인의 정보가 재현데이터에도 동시에 존재하는 경우 개인정보 위험이 있을 수 있으므로 소스 데이터와 재현데이터의 분포 또는 통계 속성 간의 유사성이나 차이 정도를 평가

(2) 재현데이터 생성 방법

- ▶ 재현데이터는 딥페이크처럼 머신러닝 기술을 통해 생성되며, 2단계의 접근 방식을 취함
 - **(1단계)** 실제 데이터를 기반으로 머신러닝 모델을 훈련하는 단계로, 이 단계에서 머신러닝 모델은 실제 데이터에 존재하는 패턴을 학습하게 되는데 이러한 머신러닝 모델을 생성 모델이라고 함
 - **(2단계)** 머신러닝 모델에서 새로운 데이터(재현데이터)를 생성하는 단계로, 생성된 재현 데이터는 원본 데이터의 많은 특성을 유지하고 있어 해당 데이터를 활용할 경우 실제 데이터를 활용할 때와 동일한 결론 도출이 가능
- ▶ 재현데이터를 생성하는 데에는 다양한 머신러닝 기법이 활용될 수 있는데 구체적인 방법은 데이터셋의 특성에 따라 달라지므로 각 상황에 적합한 모델링 도구를 사용
 - 예를 들어 하나의 스프레드시트로 표출되는 설문조사 데이터는 통계적 방법을 사용하여 적절하게 모델링할 수 있는 반면, 병원의 전자 의료기록 시스템의 복잡한 데이터를 모델링할 때에는 보다 정교한 딥러닝 기술이 필요
- ▶ 생성된 가상의 데이터와 실제 데이터는 일대일(1:1) 매칭되지 않으므로 재현데이터는 강력한 개인정보보호 특성을 보유
 - 생성된 재현데이터는 실제 데이터의 주체인 개인과 매칭하기 어려움
- ▶ 재현데이터는 비식별 정보의 형태로서, 개인정보 이용 시 필요한 정보주체의 동의 획득과 같은 절차를 거치지 않고 이용 가능
 - 일반적으로 개인정보 비식별화 조치를 수행하면 개인정보 활용에 대한 사전 동의 획득과 같은 의무가 줄어들고, 해당 데이터의 이용 및 공유가 용이해짐

- 재현데이터의 개인정보 노출 위험을 객관적으로 평가하기 위한 지표들이 개발되어 있는데, 이러한 지표들은 재현데이터에서 발생할 수 있는 다양한 유형의 유출 위험을 파악하고 성공 가능성을 정량화하여 제시
- ▶ 일반적으로 재현데이터의 이용자는 '개인정보보호 커뮤니티'와 '데이터 실사용자'라는 2가지 유형의 이해관계자로 구분할 수 있으며, 재현데이터를 활용하는 조직은 이들 이해관계자의 요구사항을 모두 고려해야 함
- **(개인정보보호 커뮤니티)** 재현데이터 생성을 개인정보보호 강화를 위한 도구로 간주하며, 개인정보 공개 위험을 축소하는 재현데이터의 특성에 관심
- **(데이터 실사용자)** 데이터 분석가를 비롯하여 재현데이터를 활용하는 실제 사용자로, 재현데이터의 높은 유용성에 관심
- ▶ 재현데이터에서는 개인정보보호와 유용성 사이의 균형이 필요하므로, 이러한 기준을 동시에 최적화하면서 수용 가능한 수준에서 생성 모델의 균형을 맞추는 것이 중요
- 재현데이터는 개인정보보호 기능이 높을수록 유용성이 적어지고, 반대로 유용성이 클수록 개인정보보호 기능이 낮을 수 있음
- 개인정보보호와 유용성 사이에서 적절한 균형을 맞추기 위하여 대체로 개인정보보호에 대한 임계치를 설정한 후, 개인정보보호 위험이 임계치 미만인 범위 내에서 유용성을 최대화 확대
- ▶ 일반적으로 재현데이터가 생성되면 이와 함께 개인정보보호 보고서와 유용성 보고서를 각각 생성
- 개인정보보호 보고서는 생성된 재현데이터에 대한 개인정보보호 지표를 문서화하고, 유용성 보고서는 재현데이터의 품질이 얼마나 양호한지를 문서화함
- 이러한 보고서들은 재현데이터 생성 방법 프로세스의 일환으로 자동적으로 만들어져야 하며, 생성된 데이터를 적용하는 데 필요한 정보를 사용자에게 제공해야 함
- ▶ 재식별 공격(re-identification attack)¹⁶⁾이 정기적으로 보고되면서 전통적인 비식별화 및 익명화 방법에 대한 우려가 존재하는 가운데, 재현데이터는 데이터를 보호하기 위한 보다 안정적이고 현대적인 접근방식으로 간주되고 있음
- 더욱이 재현데이터 생성은 대부분 자동화될 수 있는 반면 기존의 비식별화 및 익명화 기술을 제대로 적용하기 위해서는 상당한 전문 지식이 필요

16) 재식별을 목적으로 비식별화된 데이터 세트의 데이터를 원래의 데이터 주체와 연관시키기 위하여 수행하는 공격 행위를 의미

- ▶ 한편, 재현데이터 생성과 관련해서는 다음과 같은 몇몇 한계들에 대한 고려가 필요
 - 빈번히 발생하지 않는 극히 드문 사례는 생성 모델에서 포착하기 어려움
 - 이는 일반적인 모델링 문제로 생성 모델에만 국한되지 않음
 - 재현데이터로 개인정보 공개 위험을 완벽하게 차단하기는 어려움
 - 이는 모든 개인정보보호 강화 기술에 해당하는 제한사항이나, 재현데이터 생성이 제대로 수행함으로써 개인정보 공개 위험성을 크게 축소할 수는 있음

(3) 재현데이터의 장점

- ▶ 재현데이터는 여러 산업 부문에서 활발하게 활용되고 있으며, ▲개인정보보호 ▲품질 ▲비용 절감 ▲확장성 측면에서 상당한 이점을 제공
- (① 개인정보보호) 이론적으로 재현데이터는 건강 또는 생체 정보와 같은 민감한 데이터 사용과 관련된 개인정보보호 문제를 해결
 - 고유 식별자를 제거하는 대신 통계적으로 유사한 정보를 데이터에 추가하므로 개인의 신원을 감출 수 있으며 재식별 위험을 완화
 - 단, 데이터를 재현하여 가상으로 생성하더라도 개인정보 침해 위험을 완전하게 해소하는 것은 불가능
- (② 품질) 재현데이터는 품질과 데이터의 균형성 및 다양성 측면에서 실제 데이터보다 더욱 향상된 데이터 제공이 가능
 - 실제 데이터는 관련 법규로 인해 활용할 수 없는 경우가 많고, 모든 조건이나 이벤트를 고려할 수 없으므로 오류, 불균형, 부정확성, 편향성 등의 문제를 내포
- (③ 비용 절감) 실제 데이터는 비싸고 수집하기 복잡하지만, 재현데이터를 생성하면 실제 데이터를 수집하는 것보다 훈련 데이터를 훨씬 빠르고 저렴하게 확보 가능
- (④ 확장성) 머신러닝 처리에는 엄청난 양의 데이터가 필요한데, 예측 모델링에 필요한 충분한 규모의 실제 데이터를 확보하기는 어려운 데 반해 재현데이터는 누락되었거나 불충분한 정보를 채우며 원하는 조건에 맞도록 조건을 조정하여 정확하고 확장 가능한 인공지능 모델을 생성
 - 재현데이터를 통해 편향을 줄이거나, 실제로는 아직 확인되지 않는 조건으로 데이터를 표출할 수 있으며, 재현데이터와 실제 데이터를 결합시켜 더욱 향상된 데이터셋을 만들어낼 수 있음

(4) 재현데이터의 단점

- ▶ 재현데이터의 품질은 사용된 모델과 개발된 데이터셋의 품질에 따라 달라지며 위험과 한계가 존재
- ▶ 따라서 재현데이터를 활용하려는 조직은 ▲실제 데이터의 처리 ▲법적 모호성 ▲재식별 위험 ▲사실성 ▲편향 등 여러 도전과제들에 대한 고려가 필요
- **(① 실제 데이터 처리)** 재현데이터에 대한 현실적인 매개변수를 결정하기 위해서는 부분적으로 실제 데이터 처리가 요구됨
 - 실제 데이터는 개인을 식별하거나 식별 가능한 있는 정보를 포함할 수 있으며, 이 때 재현데이터 생성은 관련 개인정보보호법의 규제 대상이 될 수 있음
- **(② 법적 모호성)** 각국 정부 또는 개인정보보호 관련 규제 기관들은 개인정보, 비식별화 또는 익명화를 각각 다르게 정의하고 있으며, 재현데이터의 활용은 이러한 다양한 법규의 영향을 받을 수 있음
- **(③ 재식별 위험)** 원본 데이터 파일이 재현데이터에 포함된 경우 재식별이 가능
 - 원본 데이터로부터 개인에 대한 민감한 정보를 수집하지 않으면서 재현데이터가 실제 데이터의 통계적 속성을 정확하게 반영하도록 만드는 것은 쉽지 않음
 - 그 결과 일부 재현데이터는 원본 데이터와 매우 유사하게 생성될 수 있으며, 이를 통해 재식별이 가능할 수 있음
 - 예를 들어 알고리즘 모델이 소스 데이터의 통계적 속성을 지나치게 유사하거나 정확하게 학습하는 상황 즉, 훈련 데이터에 '과적합(overfitting)'되는 경우에는 소스 데이터를 복사하여 재식별이 더 쉬워지며, 과적합이 없더라도 우연한 파일 복제가 발생하기도 함
 - 그러나 재현에 사용되는 방법이나 알고리즘에 따라 어떤 데이터 특성이 유지되고 어떤 패턴이 억제될지 예측하는 것은 불가능
- **(④ 사실성)** 재현데이터는 기밀성 의무가 적용되는 소스 데이터를 올바르게 반영해야 하는데, 개인정보를 노출하지 않으면서 실제 현실을 반영한 데이터를 생성하기 위해서는 개인정보보호와 데이터의 유용성 사이에 균형이 필요
 - 재현데이터의 정밀도가 부족하면 모델의 학습 또는 테스트 프로세스가 적절히 이뤄지지 않아 재현데이터의 유용성이 제한될 수 있음

- (⑤ 편향) 재현데이터는 훈련 데이터 내 편향을 제한할 수 있는 기술로 각광받고 있으나, 때로는 편향을 줄이기보다 오히려 강화
 - 소스 데이터 내 오류가 있는 경우 이를 잘못 해석하여 확대하거나 인공지능 및 머신러닝 시스템의 내재적 또는 과거 편향을 모방하면서 공정성과 정확성이 저해될 수 있음

4. 결론 및 시사점

- ▶ 개인정보 침해 위험을 고려할 때 재현데이터 기술을 단독으로 사용하는 것보다는 여러 개인정보보호 강화 기술을 결합하여 사용하는 것이 효과적일 수 있음
- 차등 개인정보보호 기술과 같이 프라이버시 위험을 줄이기 위해 여러 단계에서 익명화나 가명화 기술의 적용을 고려하는 것이 권장됨
- ▶ 재현데이터는 재식별 공격 등 다양한 위험에서 완전히 자유롭지 않고, 재현데이터로 인해 정보주체의 기본권 침해가 발생하거나, 인공지능 및 머신러닝에서 편견이 제거되지 않고 영속화될 수 있으므로 재현데이터의 한계에 대한 이해와 대응이 필요
- 개인정보보호법에 명시된 신중한 접근방식을 견지하면서, 개인정보 영향평가 등을 통해 개인정보 침해 위험을 예측하고 완화하는 것이 중요

Reference

1. Data Guidance, International: Is synthetic data the future of privacy?, 2023.2.
2. Data Guidance, International: What is synthetic data and how is it generated?, 2023.2.
3. Gartner, Hype Cycle for Privacy, 2022, 2022.8.
4. Gartner, What is New in Artificial Intelligence from the 2022 Gartner Hype Cycle, 2022.9.

〈2023년 개인정보보호 월간 동향 보고서 발간 목록〉

번호	호수	제 목
1	1월 01	주요국 개인정보보호 강화기술 정책동향 분석 및 시사점
2	1월 02	EU 인공지능법(안)과 GDPR의 상호작용 분석
3	1월 03	해외 아동 개인정보 보호 침해 관련 행정처분 사례 분석
4	2월 01	해외 경쟁법 관련 개인정보보호 이슈 분석
5	2월 02	미국의 개인정보보호 법제 입법 동향
6	2월 03	디지털 자산과 개인정보보호의 관련성 및 고려사항
7	3월 01	웹 3.0 시대 도래로 부상하는 개인정보보호 이슈 분석
8	3월 02	사우디아라비아의 데이터·프라이버시 규제 샌드박스 추진현황
9	3월 03	개인정보보호와 활용성 강화 기술로 주목받는 재현데이터 기술의 특성과 활용 과제

2023

개인정보보호 월간동향분석 제3호

발 행 2023년 3월 29일

발행처 한국인터넷진흥원

전라남도 나주시 진흥길 9

Tel: 061-820-1865

1. 본 보고서는 개인정보보호위원회 출연금으로 수행한 사업의 결과입니다.
2. 본 보고서의 내용을 발표할 때에는 반드시 한국인터넷진흥원 사업의 결과임을 밝혀야 합니다.
3. 본 보고서의 판권은 한국인터넷진흥원이 소유하고 있으며, 허가 없이 무단전재 및 복사를 금합니다.

※ 본 보고서의 내용은 한국인터넷진흥원의 공식 입장과는 다를 수 있습니다.