



# 中华人民共和国国家标准

GB/T XXXXX—XXXX

## 数据匿名化流通实施及评估指南

Guidelines for the implementation and evaluation of data anonymization circulation

（征求意见稿）

（本稿完成时间：2026 年 3 月 20 日）

XXXX - XX - XX 发布

XXXX - XX - XX 实施

国家市场监督管理总局  
国家标准化管理委员会 发布



# 目 次

前言 .....	III
引言 .....	IV
1 范围 .....	1
2 规范性引用文件 .....	1
3 术语和定义 .....	1
4 数据匿名化流通总则 .....	3
4.1 数据匿名化流通原则 .....	3
4.2 数据匿名化流通目标 .....	3
4.3 数据匿名化流通框架 .....	4
4.4 匿名化流通相关主体 .....	5
5 数据匿名化流通过程 .....	5
5.1 概述 .....	5
5.2 数据及场景分析 .....	6
5.3 数据匿名化流通准备 .....	7
5.4 数据匿名化流通实施 .....	8
5.5 结果审计/效果评估 .....	8
5.6 数据流通风险管理 .....	8
6 去标识化处理技术措施 .....	8
7 流通环境安全技术措施 .....	9
7.1 概述 .....	9
7.2 基本安全技术能力 .....	9
7.3 安全计算技术能力 .....	10
8 数据匿名化流通管理措施 .....	10
8.1 数据提供方管理要求 .....	10
8.2 数据需求方管理要求 .....	11
8.3 数据开发方管理要求 .....	12
8.4 数据流通环境提供方管理要求 .....	12
9 匿名化处理过程规范性评估概述 .....	12
9.1 评估构成要素 .....	12
9.2 数据流通场景风险 .....	13
9.3 数据去标识化处理程度 .....	13
9.4 匿名化管理措施落实 .....	13
9.5 流通环境安全技术措施保障 .....	13
9.6 匿名化处理过程规范性评估流程 .....	13

9.7 重新评估情形..... 14

10 匿名化处理过程规范性评估过程..... 14

10.1 评估准备..... 14

10.2 信息调研..... 14

10.3 初步评估..... 16

10.4 细化评估..... 16

10.5 评估总结..... 18

附录 A（资料性）基于场景的数据匿名化流通示例..... 19

附录 B（资料性）常见去标识化技术和匿名模型..... 23

附录 C（资料性）数据处理协议中匿名化条款参考示例..... 27

附录 D（资料性）不同风险场景下的匿名模型参数建议..... 29

附录 E（资料性）重识别风险度量指标计算方法..... 30

## 前 言

本文件按照GB/T 1.1—2020《标准化工作导则 第1部分：标准化文件的结构和起草规则》的规定起草。

本文件由全国数据标准化技术委员会提出并归口。

本文件起草单位：北京赛西科技发展有限责任公司、中国电子技术标准化研究院、北京航空航天大学北京市大数据科学与脑机智能高精尖创新中心、浙江大学、复旦大学、蚂蚁科技集团股份有限公司、北京腾云天下科技有限公司、云基华海信息技术股份有限公司、中电数据产业集团有限公司、阿里巴巴（中国）有限公司、北京百度网讯科技有限公司、北京快手科技有限公司、北京火山引擎科技有限公司、北京小桔科技集团有限公司、华控清交信息科技（北京）有限公司、联通数据智能有限公司、中国科学院信息工程研究所、瓴羊智能科技有限公司、中兴通讯股份有限公司等。

本文件主要起草人：胡影、周晨炜、徐羽佳、李建欣、邵振赢、任奎、刘健、杨哲慙、黄丽华、白晓媛、潘无穷、黄洋成、张亚东、鲁胜强、胡成盛、顾伟、鲁艳、徐艺澈、落红卫、王昕、蔡权伟、刘笑岑、孙晓鹏、靳晨、钟康维、朱雪峰、苏丹、李海东、高晨涛、林海、刘洋、侯雨桥、贾紫薇、姚栋、梅傲婷、徐敏等。

## 引 言

数据作为生产要素，其市场化和价值化潜力备受瞩目。然而，如何平衡流通利用与安全保护，成为困扰个人数据流通的重要因素之一。2025年1月6日，《关于完善数据流通安全治理 更好促进数据要素市场化价值化的实施方案》正式印发，为个人数据合规高效流通利用指出了可行的方案，即个人数据流通应当依法依规取得个人同意，或者经过匿名化处理。

数据是否为匿名信息的判定属于法律定性判断，不在标准化文件的范畴内，匿名化一方面是一个综合的处理过程，另一方面也包含处理后的数据无法识别特定自然人且不能复原的结果，因此技术层面的匿名化判定，需要同时对过程的规范性和结果的有效性两方面进行判断，本标准认为，匿名化过程的规范性评估是匿名化结果有效性判断的前置条件，同时考虑到与TC260相关匿名化标准的有效衔接，本文件聚焦于匿名化处理过程，针对数据流通场景，提出匿名化处理过程的规范性基准，明确匿名化的操作规范、技术指标和流通环境要求，TC260相关匿名化标准则作为后置的结果判断标准依据。

当前，我国尚未形成广泛共识的匿名化规则，业界对于去标识化的应用效果，匿名化处理后数据的流通和利用价值，匿名化处理后数据的流通责任等成为争议焦点。与此同时，数据分析技术的不断发展和数据流通交易的逐步推进，多渠道数据的汇聚融合与关联分析增加了重新识别特定自然人的风险。

为此，本文件提出，在数据流通环节实施个人信息匿名化处理，要结合数据流通利用的特点和需求，既要突出流通环境的安全保障作用，从数据处理、环境安全、管理措施等方面综合实施匿名化处理，也要在数据处理中同步考虑数据接收方对数据有用性、质量、使用价值等方面的需求。

为提升数据流通场景下匿名化处理过程的规范性，本文件一方面提出了面向流通的匿名化处理实施框架，该框架基于数据流通利用典型场景，通过去标识化处理、流通环境安全保障等环节实施数据匿名化处理，关注基于流通环境的数据匿名化处理过程。另一方面，本文件构建了涵盖数据流通场景风险、数据去标识化处理程度、流通环境安全保障能力及匿名化管理措施落实情况的综合评估体系，评估匿名化处理过程是否达到规范性基准要求，从而为进一步判定其匿名化效果提供必要的前置基线，促进个人信息在安全合规的基础上高效利用。本文件适用于指导组织在数据供给、流通和利用过程中对个人信息进行匿名化流通，以及匿名化处理过程规范性的自评估与第三方评估，也可为主管监管部门进行数据流通监督管理提供参考。

# 数据匿名化流通实施及评估指南

## 1 范围

本文件提出了面向数据流通的匿名化处理实施框架和规范性评估方法，规定了数据匿名化流通的原则、目标和流程，给出了典型场景数据匿名化流通实施示例，数据匿名化流通相关技术、环境和管理措施建议，以及匿名化处理过程规范性评估的概述、流程、评估过程具体步骤。

本文件适用于指导组织在数据供给、流通和利用过程中对个人信息进行匿名化流通，以及匿名化处理规范性的自评与第三方评估，也可为主管监管部门进行数据流通监督管理提供参考。

本文件不适用于非数据流通场景下的匿名化实施与匿名化处理规范性评价。

## 2 规范性引用文件

下列文件中的内容通过文中的规范性引用而构成本文件必不可少的条款。其中，注日期的引用文件，仅该日期对应的版本适用于本文件；不注日期的引用文件，其最新版本（包括所有的修改单）适用于本文件。

GB/T 25069—2022	信息安全技术	术语
GB/T 35273—2020	信息安全技术	个人信息安全规范
GB/T 37964—2019	信息安全技术	个人信息去标识化指南
GB/T 39335—2020	信息安全技术	个人信息安全影响评估指南
GB/T 42460—2023	信息安全技术	个人信息去标识化效果评估指南

## 3 术语和定义

GB/T 25069—2022、GB/T 35273—2020、GB/T 37964—2019、GB/T 39335—2020、GB/T 42460—2023界定的以及下列术语和定义适用于本文件。

### 3.1

**去标识化** de-identification

个人信息经过处理，使其在不借助额外信息的情况下无法识别特定自然人的过程。

### 3.2

**匿名化** anonymization

个人信息经过处理无法识别特定自然人且不能复原的过程。

### 3.3

**数据提供方** data provider

产生、持有或控制数据，并在流通中出售、提供数据的组织或个人。

**注：**数据加工处理流通情形可能涉及多个提供方对数据进行联合处理或融合计算，针对非本方提供的数据而言，提供方也可能属于数据接收方。

### 3.4

#### 数据需求方 data requester

在数据流通中收集、接收、购买或使用数据的组织或个人。

### 3.5

#### 数据接收方 data recipient

在数据流通利用过程中可能接触流通数据的相关方。

注：数据直接流通情形时，数据接收方就是数据需求方；数据加工处理流通情形时，数据接收方可能涉及非本方数据的提供方、数据需求方、数据开发方、流通环境提供方等。

### 3.6

#### 数据流通环境 data circulation environment

支持数据在不同主体之间进行传输、存储、交换、共享、交易、计算等流通利用活动的硬件、软件设施和服务的集合。

注1：数据流通，是指数据在不同主体之间流动的过程，包括数据开放、共享、交易、交换等。

注2：本文件所称“不同主体”指具有不同法人资格的组织。对于同一集团公司内部不同子公司、或具有独立运营和数据控制权的业务部门之间的数据提供，参考本文件执行。

### 3.7

#### 标识符 identifier

单独或结合其他属性可以实现对个人唯一标识的数据属性。

注1：标识符可用于识别特定自然人，可分为直接标识符和准标识符。

注2：直接标识符（direct identifier），是指特定环境下可单独识别个人的标识符，如姓名、公民身份号码、护照号、驾照号、详细住址、电子邮件地址、电话号码等。

注3：准标识符（quasi-identifier），是指特定环境下结合其他属性可唯一识别个人的标识符，如性别、出生日期或年龄、民族、职业、婚姻状况、国籍等。

### 3.8

#### 附加信息 additional information

不包含在去标识化数据集内的，结合去标识化数据集后能够帮助重标识个人信息主体的信息。

### 3.9

#### 数据还原 data reversion

将经过去标识化处理后的流通数据，重新关联到个人信息主体或者重新恢复原始数据的过程。

注：数据还原主要包括数据重标识、数据复原两方面。其中数据重标识，是指将处理后的数据重新关联到个人或一组个人的过程；数据复原是指将处理后的数据重新恢复为原始数据的过程。

### 3.10

#### 目标属性 target attribute

流通数据集中提供主要分析利用价值的数据属性或特征。

注：由于目标属性提供主要分析利用价值，通常需要保留原始属性值或做尽可能少的修改；如果目标属性涉及敏感



个人信息或私密个人信息，目标属性也称为敏感属性；切断个人与敏感属性的联系，防止敏感个人信息泄露、被推断，也是匿名化处理的目标之一。

### 3.11

#### 假名化技术 pseudonymization technique

一种使用假名替换标识符的去标识化技术。

注：假名化有可逆和不可逆两种。可逆假名化要求对于假名化过程中产生的附加信息（如密钥，身份映射表等）需要安全独立保管，不和假名化后的数据一起发布给数据需求方。

### 3.12

#### $k$ -匿名 $k$ -anonymity

一种形式化隐私度量模型，确保数据集的每条记录与其他至少 $k-1$ 条记录在所有准标识符上具有相同的值。

注1： $k$ -匿名确保数据集中的每条记录都在一个大小至少为 $k$ 的等价类中，从而将个体隐藏在一定大小的群体中，使得该个体与其他至少 $k-1$ 个个体不可区分。

注2：等价类是指数据中所有准标识符属性值相同的记录构成的集合。对于满足 $k$ -匿名模型的数据，其每个等价类所包含的记录个数均不少于 $k$ 个。

注3： $l$ -多样性是在 $k$ -匿名基础上的扩展隐私保护概念，是指等价类在所选属性上的取值，至少具有 $l$ 个良好表示的值。

注4： $t$ -接近性是对 $l$ -多样性的改进技术，指等价类中某个选定属性的值分布与整个数据集中该属性的值分布之间的距离不超过阈值 $t$ 。

### 3.13

#### 差分隐私 differential privacy

一种形式化隐私度量模型，确保无论数据集中是否包含特定个体的数据，统计分析结果的概率分布差异都不会超过一个预先设定的值。

## 4 数据匿名化流通总则

### 4.1 数据匿名化流通原则

数据匿名化流通依据以下原则开展：

- a) 合法合规原则：遵守法律、行政法规等有关规定要求，尊重社会公德和伦理道德，不危害国家安全、公共利益、组织或个人合法权益，保障个人信息主体享有的法定权利；
- b) 平衡效用原则：统筹考虑数据流通利用和安全保护需求，在安全合规的基础上强调数据利用价值和质量管理要求，探索可兼顾数据利用和安全合规需求的数据匿名化流通方式；
- c) 分类分级原则：加强数据匿名化流通分类分级管理，对数据流通利用场景、去标识化处理程度、流通环境保障能力进行分类分级，丰富典型场景的数据合规流通利用方式；
- d) 风险防控原则：考虑当前合理可能的匿名化技术水平和数据还原风险，采用数据匿名化处理和流通环境安全保障相结合的方法，实现数据流通利用风险可控并持续管理风险。

### 4.2 数据匿名化流通目标

数据匿名化流通的目标包括：

- a) 无法识别：接收方在具体的流通场景和环境中，不能通过其已知的数据识别出特定自然人，不能将不同数据集中关于相同个人信息主体的信息关联起来，难以使用其他属性以较高概率推导出原本未知的敏感个人信息；
- b) 不能复原：接收方无法复原出未经去标识化处理的原始数据，例如通过泛化、随机化等不可逆去标识化技术进行处理，或者采用可逆假名化技术及流通环境保障使得接收方无法复原数据；
- c) 最少够用：在满足合法合规和风险可控的前提下，结合业务目标和数据特性，选择合适的去标识化技术和模型，确保处理后的数据质量可用，且仅流通满足使用目的最少够用的个人信息；
- d) 范围限制：数据需求方遵守与提供方的数据流通协议和环境安全责任约束，在法律规定、协议约定或个人信息主体同意授权的必要范围内使用，不超过约定的数据流通、使用的目的、方式和范围；
- e) 环境可控：综合采取技术和管理措施保障流通环境可信可控，能够防范数据重识别、复原、滥用、泄露、篡改、破坏、非法流通利用等风险，确保数据流通利用全过程可记录、可审计、可追溯，同时建立相应管理制度和保障措施，并定期跟踪评估和持续改进。

4.3 数据匿名化流通框架

本文件提出面向数据流通的匿名化处理实施框架，如图1所示。该框架采用数据处理、环境保障相结合的方式，通过数据及场景分析、数据去标识化处理、数据流通环境安全保障、结果审计/效果评估等环节实现个人信息流通利用，其核心思路是将去标识处理后的数据使用，约束在一个可信的数据流通环境，确保个人信息流通后的数据还原风险足够低，使得数据接收方无法识别特定自然人且不能复原。

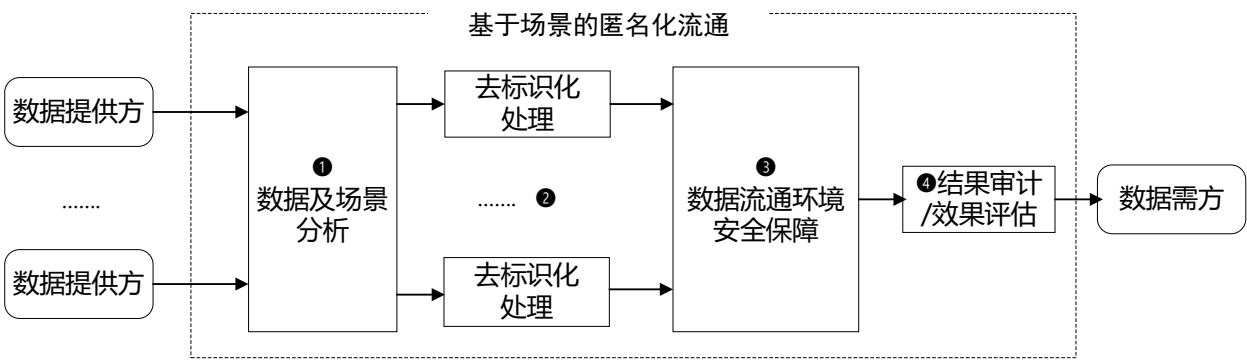


图 1 面向数据流通的匿名化处理实施框架

按照数据是否基于环境进行流通，将数据流通分为数据直接流通、加工处理流通两种情形。不同情形在上述框架的基础上，可采取以下方式进行个人信息流通利用。

- a) 数据直接流通：也称单方处理流通，数据提供方将数据处理后直接传输提供给数据需求方，如数据的内部共享、交互交换、接口访问、交易售卖、数据开放、公开发布等。数据直接流通情形可对个人信息进行匿名化处理，处理后的结果数据不属于个人信息，可自由流通给数据需求方。

注1：数据直接流通情形的个人信息匿名化处理，见《数据安全技术 个人信息匿名化处理指南及评价方法》。

- b) 加工处理流通：也称合作处理流通，两方或多方依托数据流通环境，对数据进行联合处理、融合计算等，协同完成数据处理目的后输出结果数据给需求方，如联合建模、联合统计等。加工处理流通情形可采取数据匿名化流通方式实现个人信息流通利用：

注2：融合计算是指将来自多个参与方的原始数据传输至统一的中央平台或可信环境进行物理汇聚，以对汇集后的

数据集进行综合性处理与分析的计算模式：联合处理是指各参与方的原始数据保留在本地不发生物理转移，仅在参与方之间交换模型参数或中间计算结果，以协同完成特定分析或计算任务的模式。

- 1) 对个人信息进行去标识化处理，并采用具有相应保障能力的环境进行流通利用，使得个人信息在该流通场景和环境中的数据接收方无法识别特定自然人且不能复原；
- 2) 经过匿名化处理后的结果数据，如果可指向特定自然人或者使用方可还原个人信息，需按照相关要求获得数据提供方或数据主体授权同意。

#### 4.4 匿名化流通相关主体

数据匿名化流通参与主体主要涉及组织机构（例如企业和企业之间进行联合建模），个人信息主体不在此标准的数据匿名化流通各主体范围内。数据匿名化流通参与主体的角色包括：

- a) 数据提供方：原始数据的拥有者，数据流通环境的用户之一。数据提供方根据不同的数据使用场景，选择合适的去标识化技术，对所需的原始数据进行去标识化处理后直接流通或输入数据流通环境中。
- b) 数据需求方：需求方期望从提供方获取数据，或者对提供方数据进行加工处理后获取数据结果。数据需求方向数据提供方或数据流通环境发起数据请求，并接收来自数据提供方或数据流通环境提供的数据。数据需求方也可以是数据提供方之一，或者也可以提供数据流通环境。
- c) 数据流通环境提供方：数据流通环境具有相应的安全保障能力，为处理后的流通数据提供撮合、数据交付等服务。在本文件提出的数据匿名化流通场景中，流通数据的具体使用场景在数据流通利用环境中完成开发，仅将结果输出给数据需求方。数据提供方原始数据及去标识化处理后的数据均不直接交付给数据需求方。
- d) 数据开发方：数据开发方按场景完成数据开发利用，开发利用方在数据流通环境中无法识别特定自然人或复原原始数据。数据开发方也可以是数据需求方。

### 5 数据匿名化流通过程

#### 5.1 概述

组织开展数据匿名化流通，主要包括以下实施阶段和步骤，如图2所示：

- a) 数据及场景分析：数据流通前，数据供需方分析参与流通的数据特点，涉及个人信息的直接标识符、准标识符以及目标属性，同时分析数据流通场景，明确涉及的流通目的、方式和范围、相关方、数据使用需求等，识别可能面临的流通合规要求和安全风险。本文件提出了统计分析、联合建模、大模型训练等典型场景，典型场景见附录 A；
- b) 数据匿名化流通准备：在数据及场景分析基础上，制定数据匿名化流通策略，选择合适的去标识化处理技术和模型，明确流通环境的安全保障措施，并在实施前对策略进行审核评估；
- c) 数据匿名化流通实施：按照数据匿名化流通策略，不同数据提供方需根据数据特征及使用场景，选择相应去标识化技术和模型，对流通的个人信息进行去标识化处理，方可开展流通。处理后的数据可采用安全可信的环境进行流通利用，相关方按约定在数据流通环境中进行加工处理，接收方均无法从环境中读取或获取输入数据，数据加工完后，数据流通环境自动销毁输入的数据和处理过程产生的临时数据，避免数据还原；
- d) 结果审计/效果评估：经数据流通环境加工完后，环境对输出的结果数据进行审计，确保输出数据与场景描述一致，并作必要的安全处理，留存审计结果，输出数据给需求方；对于输出数据可能包含个人信息的情形，可开展数据匿名化效果评估，判断是否达到匿名化效果；

注：匿名化效果评估方法见《数据安全技术 个人信息匿名化处理指南及评价方法》。

- e) 数据流通风险管理：采取日志记录、合规审计、风险监测、合同协议等措施，在数据流通前、流通中、流通后对数据流通利用安全风险进行持续管理，将数据流通安全风险控制在可接受风险水平。

阶段	具体工作
数据及场景分析	<ul style="list-style-type: none"><li>1. 流通数据分析</li><li>2. 流通场景分析</li><li>3. 使用需求分析</li><li>4. 流通风险分析</li></ul>
数据匿名化流通准备	<ul style="list-style-type: none"><li>1. 制定数据处理方案</li><li>2. 制定环境保障方案</li><li>3. 流通方案审核评估</li></ul>
数据匿名化流通实施	<ul style="list-style-type: none"><li>1. 数据最小化处理</li><li>2. 数据去标识化处理</li><li>3. 流通环境安全保障</li></ul>
结果审计/效果评估	<ul style="list-style-type: none"><li>1. 策略实施验证</li><li>2. 输出结果审计</li><li>3. 流通匿名化效果评估</li><li>4. 评估结果记录和留存</li></ul>
数据流通风险管理	<ul style="list-style-type: none"><li>1. 日志记录</li><li>2. 合规审计</li><li>3. 风险监测</li><li>4. 合同协议</li><li>.....</li></ul>

图 2 面向数据流通的匿名化处理实施流程

## 5.2 数据及场景分析

数据及场景分析，主要分析流通数据、数据流通场景、数据使用需求、输出数据特征和流通安全风险，主要包括以下步骤：

- a) 流通数据分析：梳理流通数据集包含的数据属性（如数据字段、数据项等）和相关来源合法性，识别涉及的个人信息，明确涉及的个人信息标识符、敏感个人信息、目标属性等。
- 1) 梳理数据属性：识别流通数据集包含哪些数据属性，分析数据属性的格式、关系、取值、规模等，明确主要提供分析利用价值的目标属性（如病症描述、薪资总额等）。
  - 2) 识别直接标识符：识别流通数据包含的个人信息直接标识符，直接标识符可单独识别特定自然人（如姓名、身份证号、电子邮件地址、电话号码等）；
  - 3) 识别准标识符：识别流通数据包含的个人信息准标识符，准标识符可结合其他属性识别特定自然人（如性别、出生日期、民族、职业等）；
  - 4) 识别敏感个人信息：识别流通数据包含的敏感个人信息，如个人生物识别、宗教信仰、特

定身份、医疗健康、金融账户、行踪轨迹、不满 14 周岁未成年人个人信息等；

- b) 流通场景分析：识别数据流通场景类型，明确该场景属于数据直接流通还是加工处理流通，分析数据流通利用的实现方式、业务模式、相关方、流转过程等。
  - 1) 识别流通场景类型：梳理数据流通场景，明确该场景属于数据直接流通还是加工处理流通情形；
  - 2) 识别实现方式：识别数据流通的实现方式，包括传输交换、接口访问、融合计算、联合处理等；
  - 3) 识别业务模式：识别数据流通的业务模式，包括数据开放、数据共享、数据交易、公共数据授权运营等；
  - 4) 识别流通相关方：识别数据流通场景中涉及的相关方和责任划分情况，分析数据接收方的动机能力和已经掌握的数据情况；
  - 5) 识别数据流转过程：识别数据流转过程的步骤环节，以及各环节在不同相关方之间流转的数据类型。
- c) 使用需求分析：分析流通数据中各数据属性的可用性和质量需求，明确需方对输出数据的特征需求，如是否与使用需求相关、是否要保留真实性、是否接受精度损失、是否需要关联、输出数据是否含标识符、输出数据涉及的属性特征等。
- d) 流通风险分析：识别该流通场景需满足的安全合规要求，分析在满足数据使用需求的前提下，可能存在的重识别、复原、安全风险和对个人信息权益的影响，综合确定场景风险水平。

### 5.3 数据匿名化流通准备

数据匿名化流通准备，主要在数据及场景分析基础上，制定数据匿名化流通策略，包括数据去标识化处理策略、流通环境保障策略。

- a) 制定数据处理策略：根据数据流通场景、数据使用需求等，分析可能的重识别攻击类型，选择合适的去标识化处理技术和模型，制定数据去标识化处理策略，主要包括以下步骤：
  - 1) 必要性可行性分析：综合考虑数据流通场景、数据使用需求等因素，分析对数据进行匿名化处理的必要性和可行性（例如是否与初步评估中的基线要求相悖）；
  - 2) 确定数据处理对象：确定需要进行处理的数据属性，包括需要消除的直接标识符、需要处理的准标识符、目标属性等；
  - 3) 数据处理技术选择：根据数据属性特点和数据使用需求，确定采用的去标识化处理技术。常见技术如泛化、抑制、假名化等技术，技术选择见第 6 章；
  - 4) 匿名模型选择：根据需要处理的数据属性，选择相应的匿名模型，例如  $k$ -匿名、差分隐私等，并根据场景风险水平设定模型参数（例如  $k$ -匿名的  $k$  值）。
- b) 制定环境保障策略：针对加工处理流通情形，按照数据匿名化流通的原则和目标，在数据处理策略的基础上，制定环境保障策略，明确流通环境安全技术要求，具体见第 7 章。
  - 1) 环境安全需求分析：结合数据流通风险分析和数据处理策略，分析为达到无法识别、不能复原、最少够用、环境可控等目标，流通环境需要满足的安全保障需求，例如基本安全保障、原始数据不出域、可用不可见、可用不可复制、可算不可识、可控可计量、可溯可审计等。
  - 2) 选择安全技术措施：根据数据流通场景和流通环境安全需求，选择相应流通环境安全技术，包括访问控制、权限管理、安全审计等基本安全技术措施，以及数据沙箱、安全多方计算、可信执行环境、联邦学习等安全计算技术。

- c) 流通策略审核评估：在匿名化流通实施前，按照第 9~10 章对匿名化处理策略进行评估。评估内容主要包括：所采取数据去标识化处理程度、环境保障措施等是否与场景风险相适应。

#### 5.4 数据匿名化流通实施

数据匿名化流通实施，主要按照数据匿名化流通策略，开展数据最小化处理、数据去标识化处理，采取流通环境安全保障，包括以下步骤：

- a) 数据最小化处理：按照最少够用目标，删除流通数据中与数据使用需求无关且涉及个人信息的数据属性。
- b) 数据去标识化处理：按照数据处理策略，采用去标识化技术和模型对数据属性进行处理。
  - 1) 处理直接标识符：根据数据处理策略中对直接标识符的处理方式，对直接标识符进行删除或假名化处理。
  - 2) 处理准标识符：根据数据处理策略中对准标识符需满足的匿名模型和参数，采取相应的去标识化技术，直至处理后的数据满足匿名模型要求。
  - 3) 处理敏感个人信息：根据数据处理策略对敏感个人信息进行去标识化处理，尤其对不必原始保留的敏感个人信息进行泛化、抑制等处理，使得敏感个人信息不易被推断。
- c) 流通环境安全保障：根据数据流通环境安全保障方案，部署相应安全技术措施，如采用数据沙箱、可信执行环境、安全多方计算等安全计算环境，配套访问控制、身份鉴别、权限管理、安全审计等基本技术措施，实现数据匿名化流通的无法识别、不能复原、最少够用、范围限制、环境可控等目标，防范数据重识别、复原、泄露、篡改、破坏、滥用、非法流通利用等风险。

#### 5.5 结果审计/效果评估

结果审计/效果评估，旨在验证数据流通处理过程是否按照既定策略实施，并对输出结果是否包含个人信息进行审计。视情通过匿名化效果评估确认是否达到匿名化效果，并对审计/评估过程和结果进行记录，具体步骤包括：

- a) 策略实施验证：验证实际开展的数据去标识化处理、流通环境安全保障措施，是否与数据匿名化流通策略保持一致；
- b) 输出结果审计：数据经流通利用平台加工完后，数据流通环境提供方或独立第三方对计算的结果进行审计，确保输出数据与场景描述一致，并作必要的安全处理，留存审计结果，输出数据给需求方。
- c) 匿名化效果评估：对于输出数据可能包含个人信息的情形，宜开展匿名化效果评估；  
注：匿名化效果评估方法见《数据安全技术 个人信息匿名化处理指南及评价方法》。
- d) 评估结果记录和留存：记录策略实施验证、输出结果审计、匿名化效果评估的过程和结果，保存相应记录和报告至少 3 年。

#### 5.6 数据流通风险管理

数据流通风险管理，主要采取日志记录、合规审计、风险监测、合同协议等措施，在数据流通前、流通中、流通后对数据流通利用安全风险进行持续管理，将数据流通安全风险控制在可接受风险水平，见第8章。

### 6 去标识化处理技术措施

根据流通数据的属性特点、使用需求等，选择适合的去标识化处理技术和模型，不同技术方法也可组合使用，选择策略包括但不限于：

- a) 根据使用需求是否要关联数据，可对直接标识符进行删除或假名化处理，如对数据进行假名化处理，应对附加信息进行安全保存。
  - 1) 如假名化采用加密的方式，应确保密钥的安全；
  - 2) 如假名化采用映射表的方式，应对映射表进行隔离和加密存储，并设置严格的访问权限；
- b) 根据使用需求是否要保留数据真实性，对准标识符选择泛化或随机化类处理技术，如需保留数据真实性且可接受精度损失的，可选择泛化类处理技术；
- c) 针对群体进行统计、挖掘分析的场景，可采用随机化、泛化等技术进行去标识化处理，在满足差分隐私模型要求下输出。
- d) 针对数据精度要求不高的开发测试场景，可采用数据合成技术代替原始数据，合成数据保留与原始数据相符的特征；
- e) 结构化数据集可采用  $k$ -匿名等匿名模型，可针对典型数据流通场景，按照场景风险水平，设定  $k$ -匿名模型参数（如  $k$  值、 $l$  值、 $t$  值等），例如准标识符满足 5-anonymity，敏感属性 A 需满足的 2-diversity、敏感属性 B 需满足 3-diversity 等。

常见去标识化技术和匿名模型见附录B。

## 7 流通环境安全技术措施

### 7.1 概述

数据流通环境，可能涉及数据流通交易环境、数据加工利用环境。其中，数据流通交易环境撮合数据流通交易，提供流通数据登记、数据流通合同协议等服务，并将数据使用场景做准确描述；数据加工利用环境提供数据分析处理的环境，数据开发方根据合同协议描述场景和数据输入，自定义开发数据加工处理程序，数据开发利用环境为数据加工处理程序提供执行环境。数据加工利用环境可支持数据库操作、实时计算引擎、离线计算引擎、隐私保护计算、可信数据空间等环境。

### 7.2 基本安全技术能力

数据流通环境可采取硬件或软件环境，应具备以下基本安全技术能力：

- a) 身份认证：采用多因素身份鉴别措施，对相关方身份进行认证；
- b) 访问控制：从系统功能权限和数据权限方面对相关方进行管理，限制未经授权的数据访问行为。
- c) 安全隔离：确保不同接收方的数据在逻辑或物理上进行安全隔离；
- d) 加密保护：对敏感个人信息进行加密存储，采用的密码技术应符合国家密码管理相关要求；
- e) 安全传输：数据在互联网传输时进行加密，并支持安全的数据传输措施（如 VPN 或专线接入）；
- f) 数据销毁：计算任务完成后，删除环境中的原始数据和中间结果；数据安全加工环境数据流通结束后，所有输入数据，运行过程产生的临时数据都会被立刻删除。
- g) 数据防泄漏：采取数据防泄漏技术措施，防范数据泄露或未授权访问；
- h) 附加信息保护：应采取措施对可还原数据的附加信息进行加密和隔离保护，除数据提供方外其他方应禁止访问附加信息；
- i) 境内存储：从事境内数据流通利用服务的环境和基础设施，应部署在我国境内；
- j) 接口安全：应对数据接口访问进行身份鉴别，对不安全输入参数进行限制和过滤，为接口提供异常处理能力；

- k) 风险监测：开展数据流向监控和异常行为检测，并在发现安全缺陷、漏洞时，立即采取补救措施；
- l) 安全审计：对数据流通、访问或操作进行记录，并定期开展安全审计，及时发现未授权或风险操作行为。

### 7.3 安全计算技术能力

#### 7.3.1 基础级安全计算技术能力

数据流通环境的基础级安全计算技术能力包括：

- a) 提供容器化、虚拟化或物理硬件的隔离环境，确保不同接收方的处理任务间的隔离。例如采用数据沙箱技术实现数据隔离计算；
- b) 数据流通环境应采取以下管控措施，确保环境可控可信：
  - 1) 环境应具备基本的抗攻击能力，能够阻断接收方和外部人员的攻击行为。包括：窥探、干扰隔离环境内的执行情况；将去标识化处理后的数据重识别到特定自然人；将去标识化处理后的数据复原为原始数据；
  - 2) 环境应只执行场景规定的计算逻辑，且数据提供方能够对此进行验证；
  - 3) 环境应能够防止非预期的数据进入到隔离环境；
  - 4) 环境应能够防止非预期的数据从隔离环境中输出。
  - 5) 任何人无法通过环境获取输入数据或计算的中间结果数据。
- c) 确保计算环境有独立边界和严格管控，同时抵御外部黑客和内部人员攻击，如采用独立系统并经过加固；
- d) 建立完整的操作日志记录机制，对数据接入、计算任务启动、参与方行为、结果输出等关键操作进行记录，支持事后追溯与审计。日志应防篡改并保留一定期限。

#### 7.3.2 增强级安全计算技术能力

数据流通环境的增强级安全计算技术能力如下：

- a) 确保计算过程中，原始数据始终处于密态不可见的状态，抵御特权管理员攻击，例如采用可信执行环境、多方安全计算、同态加密等技术；
- b) 应通过技术手段防止数据被非法复制、下载或持久化存储，如采用全流程密态计算等技术；
- c) 如采用可信执行环境，可信执行环境应具备可信度量、远程认证、隔离等能力，保证可信执行环境中的数据和代码的安全性；
- d) 在需要对密文数据进行持久化存储的场景，宜支持密钥分片存储，防止单点泄露；
- e) 应建立完整、不可篡改的数据使用链路追踪机制，支持从数据源、计算任务、参与方到结果输出的数据流通利用全过程全链路追溯与审计，如采用数据水印、数据血缘、区块链等技术。

## 8 数据匿名化流通管理措施

### 8.1 数据提供方管理要求

数据提供方管理要求包括：

- a) 对数据需求方的使用场景、目的和处理过程进行充分了解和审核；
- b) 提供满足数据使用场景所需的最小范围数据；
- c) 对数据使用场景和数据匿名化处理的合规性和安全性进行评估；



- d) 采用合同协议等具有约束效力的形式，明确数据流通相关方的责任和义务，包括但不限于：
  - 1) 数据使用场景、使用目的和范围限制；
  - 2) 数据安全保护义务和责任；
  - 3) 禁止使用方对匿名化流通的数据进行重标识或将其与其他数据链接；
  - 4) 是否允许使用方将数据进行二次提供及其条件；
  - 5) 数据泄露事件的通知和应急响应机制；
  - 6) 违反数据流通利用约定的处罚措施。

注：数据处理协议中匿名化条款参考示例见附录 C。

- e) 确保提供的数据符合法律相关要求；
- f) 明确数据匿名化处理的人员和职责，定期对负责人员进行相关技术培训；
- g) 建立数据匿名化流通管理制度，明确数据流通场景、采取的去标识化处理技术、剩余风险分析、安全保障措施等；
- h) 制定并公开数据匿名化流通处理规则，告知数据流通的目的、处理方式、数据范围、相关方，以及采取哪些措施减少数据匿名化流通的安全风险；
- i) 在个人信息流通前，制定数据匿名化流通策略，并对策略开展个人信息保护影响评估；
- j) 建立数据匿名化流通统一台账，并纳入个人信息保护影响评估范围；
- k) 制定并定期演练数据匿名流通应急预案，发生数据还原、泄露、滥用等安全事件时，立即启动预案，采取措施防止危害扩大，消除安全隐患，并按照规定向有关部门报告；
- l) 跟踪相关技术、法律法规和业界实践发展，定期审核更新数据匿名化流通的制度、处理规则和策略。
- m) 持续监控内外部风险变化，并在条件满足时，组织并完成重新评估工作，确保匿名化效果的持续有效性。

## 8.2 数据需求方管理要求

数据需求方管理要求包括：

- a) 按照最少够用和目的限制原则，申请满足实际使用目的的最小必要数据，限制匿名化流通的数据的访问权限；
- b) 与数据提供方等签订数据使用合同协议，对数据使用场景、使用目的和处理流程进行如实说明，合同协议应满足8.1中d)要求；
- c) 按照事先约定的数据使用场景和使用目的，对数据进行开发利用；
- d) 对匿名化流通的数据的使用目的进行评估，核查自身的使用目的是否包含识别特定个人的目的；
- e) 对匿名流通处理的数据使用进行监控，持续监督内部人员数据使用范围是否超出约定场景；
- f) 对会接触到匿名化流通的数据的人员进行培训，并签订保密协议；
- g) 采取措施禁止任何试图重新识别特定自然人、复原原始数据、超范围使用或滥用数据的行为，并确保采取适当措施销毁任何意外重新识别的个人信息；
- h) 对将其他数据与匿名化流通的数据进行关联分析的能力进行控制，以管理因关联而产生的数据还原风险；
- i) 如结果数据未达到匿名化效果，按照相关规定在数据使用前告知个人信息主体使用目的、方式和范围，并获得个人授权同意；
- j) 对匿名化流通的数据的开发利用操作进行日志记录，并定期对操作行为进行安全审计；

- k) 达成使用目的后，将匿名化流通的数据进行删除或匿名化处理。
- l) 就数据使用方式、环境变化或任何可能影响匿名化效果的情况，及时向数据提供方通知。

### 8.3 数据开发方管理要求

数据开发方管理要求包括：

- a) 严格按照需求方和供应方签署的合同场景和处理流程开发数据加工利用程序；
- b) 不采取技术手段在加工利用程序中做违反合同约定操作，例如未经供需方同意查看、输出、重定向、存储输入数据、输出数据或中间结果数据；
- c) 不通过技术手段将信息隐藏在输出结果中。

### 8.4 数据流通环境提供方管理要求

数据流通环境提供方管理要求包括：

- a) 提供满足7.2要求的基本安全技术能力，按照场景需求具备满足7.3要求的安全计算技术能力；
- b) 向数据流通利用相关方，告知数据流通环境提供的安全技术能力和安全管理措施；
- c) 定期对所提供的流通环境安全技术措施进行安全评估，确保其提供的安全保障能力满足数据匿名化流通安全需求，不会对数据进行泄露、未授权访问、重识别、复原等；
- d) 按照最小授权原则对数据流通环境设置严格的访问控制策略，并确保访问控制策略的实施；
- e) 对数据提供方、数据需求方进行身份审核和鉴别；
- f) 对供需双方签署合同及开发方的直接相关代码进行审核；
- g) 对数据提供方、需方、开发利用方在环境内的操作行为进行日志记录，并定期对敏感操作行为进行安全审计；
- h) 制定并定期演练数据匿名流通应急预案，发生安全事件时及时响应，采取措施消除安全隐患防止危害扩大，并向相关方及有关部门报告。
- i) 当其平台环境的技术架构或安全能力发生可能影响重识别风险的实质性变更时，通知数据流通场景相关方。

## 9 匿名化处理过程规范性评估概述

### 9.1 评估构成要素

数据流通场景下的匿名化处理过程规范性评估构成要素包括：

- a) 数据流通场景风险：刻画数据流通场景的重识别风险高低，一旦遭到重识别对个人权益的影响；
- b) 数据去标识化处理程度：刻画对原始数据的直接标识符、准标识符进行转换处理后，结果数据的可识别性强弱；
- c) 匿名化管理措施落实：处理者实施匿名化应配套落实的管理措施或管理体系；
- d) 流通环境安全能力保障：降低数据流通安全风险各类措施，包括基本安全技术措施，以及能够进一步提升数据流通安全的设施、软硬件环境、增强技术等安全计算技术措施。

其中，数据流通场景的风险决定了数据处理、管理措施、安全环境等措施的综合强度要求，场景风险越高，对应的措施要求越强；对原始数据的转换处理是数据匿名化的核心，直接反映数据本身的可识别性高低；匿名化管理措施是确保匿名化处理的规范性、持续有效的保障；通过提高流通环境的安全性，降低数据流通安全风险，能够辅助性地、间接地降低结果数据遭到重识别的总体风险、综合可能性。

## 9.2 数据流通场景风险

数据流通场景风险分为：

- a) 高风险场景：场景涉及人群的规模大、数据敏感度高、流通范围较广、流通链条长、接收方数量大等，导致此类场景下结果数据面临的初始重识别风险较高。例如公开发布、数据交易场景。
- b) 中风险场景：场景涉及人群的规模、数据敏感度、流通范围、流通链条、接收方数量适中，导致此类场景下结果数据面临的初始重识别风险适中。例如跨机构共享场景。
- c) 低风险场景：场景涉及人群的规模小、数据敏感度低、流通范围有限、流通链条短、接收方数量小等，因此此类场景下结果数据面临的初始重识别风险较低。例如内部共享场景。

## 9.3 数据去标识化处理程度

数据去标识化处理程度分为：

- a) 完全处理：结果数据不包含任何标识符。
- b) 充分处理：消除了直接标识符，对准标识符进行了处理，在相应的数据流通场景风险下，能够提供较高程度重识别防范作用。
  - 1) 充分消解单独挑出、关联、推断风险；
  - 2) 充分抵御身份披露攻击、属性披露攻击、统计推断攻击。
- c) 部分处理：消除了直接标识符，对准标识符进行了处理，在相应的数据流通场景风险下，能够提供一定程度重识别防范作用。
  - 1) 部分消解单独挑出、关联、推断风险；
  - 2) 部分抵御身份披露攻击、属性披露攻击、统计推断攻击。
- d) 不足处理：结果数据仍包含直接标识符。

## 9.4 匿名化管理措施落实

匿名化管理措施包括组织建设、匿名化策略管理、规则公开、记录留存、接收方评估、相关方约束、应急机制、监督审计等方面，见第8章。

## 9.5 流通环境安全技术措施保障

流通环境安全技术措施包括：

- a) 基本安全技术措施：包括身份认证、访问控制、存储隔离、加密保护、操作审计、数据销毁等。
  - b) 安全计算环境技术措施：包括数据沙箱技术、隐私保护计算技术等。
- 具体措施见第7章。

## 9.6 匿名化处理过程规范性评估流程

数据匿名化处理过程规范性评估按照如下流程开展：

- a) 评估准备：确定评估场景、评估范围、组建团队、制定方案等。
- b) 信息调研：涉及数据流通场景、流通数据情况、匿名化技术策略、匿名化管理措施、数据流通安全措施。
- c) 初步评估：数据匿名化基线要求符合性评价。
- d) 细化评估：在通过初步评估的基础上，围绕数据流通场景风险、数据去标识化程度、匿名化管理措施落实、流通环境安全技术措施保障开展进一步评估。
- e) 评估总结：编制评估报告、评估报告复核、评估过程记录存档等。

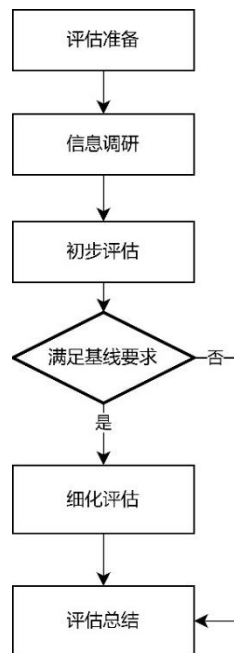


图 3 匿名化处理过程规范性评估流程

9.7 重新评估情形

匿名化评估后，出现以下情形时需要重新开展匿名化评估：

- a) 原始数据发生重大更新（如新增关键字段），或结果数据被用于与评估未覆盖的其他数据集进行关联、融合；
- b) 数据流通的目的、接收方范围或地域范围超出了初始评估所界定的边界；
- c) 数据流通环境的安全控制措施发生变更；
- d) 发生与该数据集相关的重识别事件、数据泄露事件；
- e) 相关的法律法规、国家标准或行业最佳实践发生变化。

10 匿名化处理过程规范性评估过程

10.1 评估准备

评估准备阶段包括：

- a) 确定评估目标：包括数据匿名化处理过程**基线达标**、数据匿名化处理过程**规范性达标**；
- b) 评估对象范围：包括数据流通场景、待评估的数据集、数据流通环境等；
- c) 组建评估团队：包括个人信息保护专家、匿名化技术专家、数据安全专家、法律合规专家；
- d) 制定评估方案：包括评估范围、内容、方法、人员、计划等；
- e) 准备评估工具，包括但不限于评估报告模板、用于重识别风险分析的技术工具等。

10.2 信息调研

10.2.1 数据流通场景

调研数据流通的目的、方式、范围，涉及的相关方、数据流转等情况，包括：

- a) 流通利用目的：
    - 1) 场景目的类型：包括但不限于数据交易、商业运营、公共服务、科学研究、技术开发等；
    - 2) 为实现数据流通利用目的是否需要识别、关联、推断；
    - 3) 结果数据的使用是否仍针对个人；
    - 4) 数据匿名化处理目的：包括但不限于尽责性个人信息权益保障、免于取得同意等合法性基础、应当开展匿名化处理时履行自身义务等。
  - b) 实现方式：包括但不限于直接传输、接口访问、协同计算等。
  - c) 业务模式：包括但不限于数据开放、数据共享、数据交易、公共数据授权运营等。
  - d) 流通范围：包括但不限于公开发布、对外提供、内部共享等。
- 注：本文件不涉及跨境场景。
- e) 场景相关方：数据提供方、数据接收方、环境提供方、数据开发方。
  - f) 流通环境：使用的数据流通环境基本情况，例如数据基础设施。

### 10.2.2 流通数据情况

调研提供方及接收方的数据情况，包括：

- a) 提供方数据情况：原始数据和结果数据的规模、包含的属性项、是否最小化处理、是否包含敏感个人信息。
- b) 接收方数据情况：拥有的与所提供数据共同的属性项、可访问的其他数据情况。
- c) 提供方、接收方数据的属性类别梳理：
  - 1) 目标属性梳理：哪些属性属于为实现流通利用目的所必需，需要原始保留的属性；
  - 2) 直接标识符梳理：哪些属性属于与实现流通利用目的无关，仅具备标识作用的属性；

注 1：常见直接标识符示例 GB/T 42460—2023 附录 A。

- 3) 准标识符梳理：哪些属性属于既与实现流通利用目的相关但不必要，又能通过组合其他属性具备标识作用的属性。

注 2：准标识符可接受一定程度精度损失，通过泛化等处理，可以在降低数据可识别性的同时保留一定的使用价值。常见准标识符示例 GB/T 42460—2023 附录 B。

### 10.2.3 匿名化技术策略

调研实施匿名化所采取的数据匿名化技术处理方案，包括：

- a) 对各属性的类别判断结果，是否遗漏对直接标识符、准标识符的识别。
- b) 采用的数据转换处理技术、匿名模型：
  - 1) 采取的数据转换处理技术及其可逆性：包括对直接标识符的转换处理技术（例如删除、假名化等）、准标识符的转换处理技术（例如泛化等）；
  - 2) 是否采用匿名模型、以及匿名模型的参数设置：包括但不限于  $k$ -anonymity、 $l$ -diversity、 $t$ -closeness、差分隐私等模型及其参数设置。
- c) 选取的匿名模型与场景、数据的适配性，包括但不限于：
  - 1) 需要抵御身份披露风险时，是否使用  $k$ -anonymity 等模型；
  - 2) 需要抵御属性披露风险时，是否使用  $l$ -diversity、 $t$ -closeness 等模型；
  - 3) 需要抵御统计推断攻击时，是否使用差分隐私等模型。
- d) 匿名化技术实施的规范性，包括但不限于：
  - 1) 是否设置了各个准标识符的泛化层次：包括数值型属性泛化层次、分类型属性泛化层次。

- 2) 是否对于结果数据能否满足设定的匿名模型进行验证;
- 3) 是否设置了提高结果数据有用性的数据效用度量, 包括专用指标、通用指标等;
- 4) 是否采用了常见的匿名化实现算法;
- 5) 是否使用了匿名化技术工具。

#### 10.2.4 匿名化管理措施

调研实施的匿名化相关管理措施情况, 包括但不限于以下方面:

- a) 组织建设: 包括是否涉及专门岗位职责、相关人员是否具备匿名化专业技术能力、是否开展相关教育培训等。
- b) 匿名化策略管理: 包括匿名化策略制定、审核、实施、验证、持续改进等机制是否完善。
- c) 规则透明: 是否制定匿名化处理规则并及时更新, 是否在个人信息保护政策中进行必要的公开声明等。
- d) 记录留存: 是否留存匿名化策略、规则的制定、审核、更新等相关记录。
- e) 接收方评估: 是否对数据接收方的重识别动机、能力, 已经掌握的数据情况等进行分析评估。
- f) 相关方约束: 是否通过合同、协议等, 约束结果数据的使用目的范围、禁止开展重识别等。
- g) 应急机制: 是否建立匿名化相关应急预案、处置机制等。
- h) 监督审计: 是否要求接收方提供结果数据的后续处理记录, 并定期进行审查。

#### 10.2.5 流通环境安全技术措施

调研采取的流通环境安全技术措施, 包括但不限于:

- a) 基本安全技术措施: 是否采取身份认证、访问控制、存储隔离、加密保护、操作审计、数据销毁等技术措施。
- b) 安全计算环境技术措施: 是否采取数据沙箱、安全多方计算、可信执行环境等技术措施。

### 10.3 初步评估

初步评估用于快速识别不符合以下基线要求的情形, 包含以下评估项:

- a) 流通目的的实现未建立在识别、关联、推断之上, 与匿名化的保护目的不相悖;
- b) 结果数据的使用不针对个人, 即结果数据不与个人相关联;
- c) 最小化处理: 原始数据、结果数据, 不包含超出必要范围的属性(与流通利用目的无关的属性);
- d) 不可直接识别: 结果数据不存在直接标识符;
- e) 不可稳定链接: 不存在能将不同数据集中关于同一个人关联起来的标识符, 例如链接用 ID;
- f) 采用的处理技术不可逆: 采取泛化、抑制等不可逆的处理技术; 采用可逆假名化技术处理的, 对密钥、映射表进行安全存储;
- g) 相关方均满足第 8 章中的匿名化相关管理措施要求;
- h) 数据流通环境满足 7.2 中的基本安全技术措施要求。

以上任一要求不满足的, 则匿名化处理过程**不达标**; 全部满足的, 则匿名化处理过程**基线达标**。

### 10.4 细化评估

#### 10.4.1 场景风险评估

场景风险评估综合以下因素给出结论, 具体包括:

- a) 数据涉及个人的规模;
- b) 包含属性项的数量;

- c) 是否包含敏感个人信息；
- d) 相关方（特别是数据接收方）的数量；
- e) 接收方重识别动机强弱；
- f) 流通目的类型属于公益性或商业性；
- g) 流通范围的大小、数据流通链条的长短；
- h) 如果发生重识别对个人权益的影响。

评估结论分为：高风险、中风险、低风险。

#### 10.4.2 数据去标识化处理程度评估

判断去标识化处理程度需结合接收方可访问的数据情况，在识别结果数据与接收方可访问数据的共同属性的基础上，对数据去标识化处理程度的评估可采取以下方案，包括但不限于：

- a) 标准标识符识别法。基于标准标识符集，识别结果数据中是否包含标准标识符集中的属性。标准标识符集包括：

- 1) 标准直接标识符集；
- 2) 标准准标识符集。

注 1：通过标准标识符识别评估的，数据去标识化处理程度相当于部分处理。

注 2：标准标识符集，可能根据数据流通场景所涉及行业领域的不同有所区别。

注 3：标准标识符识别法，不适用于高风险数据流通利用场景。

- b) 匿名模型及参数法。针对单独挑出、可链接性、属性推断等，验证匿名化处理后的数据，是否满足所声称的匿名模型及其参数，包括但不限于：

- 1)  $k$ -anonymity 模型及参数  $k$ ；
- 2)  $l$ -diversity 模型及参数  $l$ ；
- 3)  $t$ -closeness 模型及参数  $t$ ；
- 4) 差分隐私模型及参数  $\epsilon$ 。

注 4：基于数据流通利用场景和流通数据情况，可能有不同的匿名模型要求及参数阈值。例如低风险场景下，以  $k$ -anonymity 模型为基线要求，充分处理、部分处理对应的阈值可能分别取 3、2；高风险场景下，对应阈值提高到 20、5。不同风险场景下的匿名模型参数建议见附录 D。

- c) 重识别风险指标法。可采用的风险指标包括但不限于：

- 1) 检察官、记者、营销者重识别风险指标；
- 2) 记录唯一性指标（例如唯一记录占比）。

注 5：基于风险阈值设定，重识别风险法可用于判断是否数据去标识化处理程度相当于充分处理、部分处理。例如设定最大检察官风险不超过 20%、平均检察官风险不超过 10%为部分处理的风险阈值；最大检察官风险不超过 10%，平均检察官风险不超过 5%为充分处理的风险阈值等。重识别风险度量指标计算方法见附录 E。

评估结论分为：完全处理、充分处理、部分处理、不足处理。其中，经以上任一方法评估结论为部分处理的，最终结论可为部分处理；经以上任一方法评估结论为充分处理的，最终结论可为充分处理；经以上所有方法评估结论为充分处理时，最终结论可为完全处理。

#### 10.4.3 安全计算技术能力评估

评估安全计算环境技术措施达到的能力水平：

- a) 基础级安全计算技术能力：实现了基本的可用不可见、可用不可复制、可溯可审计等能力。
- b) 增强级安全计算技术能力：实现了较强的可用不可见、可用不可复制、可溯可审计等能力。

评估结论分为：不具备安全计算技术能力、具备基础级安全计算技术能力、具备增强级安全计算技术能力。

10.4.4 符合度结论

根据表1判断，在该场景风险下，数据去标识化处理程度、匿名化管理措施、环境保障措施是否达到相应水平要求。均达到要求则该场景、目的、环境、管理措施等约束下，认为匿名化处理过程规范性达标，否则认为匿名化处理过程基线达标。

表 1 不同场景风险下数据处理、管理与环境保障措施组合的规范性要求

数据流通场景风险	数据去标识化处理程度	匿名化管理措施	流通环境保障能力	
			基本安全技术能力	安全计算技术能力
高/中/低风险	完全处理	√	√	—
	充分处理	√	√	基础级
	部分处理	√	√	增强级
注：不同场景风险下，充分处理、部分处理对应的参数建议值不同，匿名模型及参数法的不同场景风险参数建议见附录 D。				

10.5 评估总结

评估总结阶段包括：

- a) 编制评估报告。评估报告内容包括但不限于：
  - 1) 评估目的；
  - 2) 评估依据；
  - 3) 评估对象和范围；
  - 4) 评估团队；
  - 5) 评估过程；
  - 6) 评估方法和工具；
  - 7) 风险分析；
  - 8) 评估结论；
  - 9) 改进建议。
- b) 评估报告复核；
- c) 评估过程记录存档。



## 附录 A

### (资料性)

#### 基于场景的数据匿名化流通示例

#### A.1 统计分析场景

##### A.1.1 场景描述

统计分析场景主要对数据进行汇总、整理和分析，以获取统计特征、揭示整体分布规律或趋势，如群体分析、业务看板、社会调查、市场预测等。数据需求方获得统计分析结果。

- a) 场景相关方包括：数据提供方（如电商平台、医疗机构等）、统计分析方（数据需求方，如咨询公司、科研机构等）、数据流通环境提供方。
- b) 涉及的数据类型包括：
  - 1) 直接标识符：如身份证号、手机号码、用户账号等；
  - 2) 准标识符：如年龄、性别、地理位置、职业、学历等；
  - 3) 目标属性：如消费记录、医疗记录、行为特征等。
- c) 按照数据流通方式，统计分析场景分为以下两种情形：
  - 1) 直接流通情形：数据提供方将去标识化处理后的数据集直接提供给需求方，由需求方自行完成统计分析；需求方接触的是记录级数据，输入侧去标识化处理是主要保护机制；
  - 2) 加工处理流通情形：数据提供方将数据输入流通环境，统计分析在流通环境内完成，需求方仅获得统计结果；原始数据不出域，输出侧差分隐私保护是主要保护机制。
- d) 统计分析场景涉及多个数据提供方时，各方应以约定的类别型分析维度（如年龄段、性别、地区等）对本方数据进行分组，分别计算各组的统计值（如计数、均值、占比等），仅将组级统计结果用于跨方比较或合并计算，不涉及个体级记录的跨方传输或对齐；需要识别并对齐不同数据集中同一自然人记录后再进行联合统计的，不满足匿名化基线要求，不适用本场景。

##### A.1.2 数据去标识化处理

- a) 直接流通情形，数据提供方在流通前对数据集实施完整的去标识化处理，措施包括但不限于：
  - 1) 与统计目的无关的属性，在流通前删除，不纳入流通数据集；
  - 2) 对直接标识符进行删除或遮蔽处理；
  - 3) 对准标识符，根据统计目的是否需要保留属性真实值进行选择：需要保留真实值或可接受精度损失的，采用泛化处理（如数值型属性区间化、地理位置层级泛化等）；无需保留真实值的，采用随机化或抑制处理；
  - 4) 对目标属性中涉及敏感个人信息的，进行泛化或抑制处理，确保敏感个人信息不易被推断；
  - 5) 对处理后的结构化数据集，采用  $k$ -匿名等匿名模型验证准标识符的处理结果是否满足设定参数要求，参数建议见附录 D。
- b) 加工处理流通情形，数据提供方在流通前应进行最小化处理，删除与统计目的无关的属性及非必要直接标识符；输入侧去标识化处理作为内部人员防护的纵深措施，可参照 a)1)~4)实施。

##### A.1.3 流通环境安全保障

- a) 直接流通情形，流通环境满足 7.2 基本安全技术能力要求；
- b) 加工处理流通情形：
  - 1) 单一数据提供方的，流通环境满足 7.2 基本安全技术能力要求；

- 2) 涉及多个数据提供方的，流通环境在满足 1)的基础上，具备 7.3.1 基础级安全计算技术能力；
- c) 加工处理流通情形下，对统计查询输出结果采用差分隐私技术进行保护，参数  $\epsilon$  的取值建议见附录 D。

#### A. 1.4 结果审计

- a) 直接流通情形，输出数据集不应包含直接标识符；对输出数据集开展重识别风险评估，风险度量方法见附录 E；
- b) 加工处理流通情形，输出统计结果不应包含个人信息标识符；对输出统计结果开展差异攻击测试（Differencing Attack），评估攻击者能否通过比对多次查询结果反推个人信息；测试发现风险的，应补充差分隐私保护或限制查询次数；
- c) 记录并留存上述审计过程和结果，保存期限不少于 3 年。

### A. 2 联合建模场景

#### A. 2.1 场景描述

联合建模场景指多方利用各自数据资源，共同参与模型训练过程，联合构建预测或分类模型，如多机构联合建立风控模型、多家医院共同开发疾病预测模型等。数据需求方获得训练完成的机器学习模型或模型的预测结果。模型预测结果针对个人的情形不满足匿名化基线要求，不适用本场景。

- a) 场景相关方包括：数据提供方、模型开发方（数据开发方）、数据流通环境提供方、模型使用方（数据需求方）。
- b) 涉及的数据类型包括：
  - 1) 训练数据：各方持有的特征数据、标签数据、样本数据等；
  - 2) 模型数据：模型参数、权重、结构等；
  - 3) 中间结果：训练过程中产生的梯度、损失值、评估指标等。

#### A. 2.2 数据去标识化处理

联合建模场景的数据去标识化处理，针对各方持有的训练数据实施，措施包括但不限于：

- a) 对训练数据中的直接标识符进行删除、杂凑或随机化等处理；
- b) 对准标识符进行泛化处理；
- c) 对涉及敏感个人信息的属性，进行泛化或抑制处理，使敏感个人信息不易被推断；
- d) 涉及敏感个人信息的模型训练，可采用差分隐私技术（DP-SGD）对训练梯度添加随机噪声，参数  $\epsilon$  的取值建议见附录 D。

#### A. 2.3 流通环境安全保障

流通环境安全保障措施包括但不限于：

- a) 各参与方的训练数据在进入流通环境前，应按照 A.2.2 完成去标识化处理；
- b) 联合处理情形（各方数据不发生物理转移，仅交换模型参数或中间计算结果），流通环境应满足 7.2 基本安全技术能力要求，并具备 7.3.1 基础级安全计算技术能力；
- c) 融合计算情形（各方数据汇聚至统一环境进行处理），流通环境应在满足 b)的基础上，具备 7.3.2 增强级安全计算技术能力，确保原始数据在计算全程处于密态；
- d) 建模完成后，所有输入数据及训练过程产生的中间结果数据应及时销毁。

#### A. 2.4 结果审计

- a) 联合建模场景输出的机器学习模型，不应包含可识别特定自然人的个人信息；
- b) 对输出模型开展以下测试：
  - 1) 成员推断攻击测试（Membership Inference Attack）：评估攻击者能否通过查询模型推断某条记录是否出现在训练数据中；
  - 2) 模型逆向攻击测试（Model Inversion Attack）：评估攻击者能否通过模型输出反推训练数据中的个人信息；
- c) 测试发现显著个人信息泄露风险的，应采取模型遗忘、差分隐私再训练等措施处置，或限制模型对外提供范围；
- d) 记录并留存上述审计过程和结果，保存期限不少于 3 年。

### A.3 大模型训练场景

#### A.3.1 场景描述

大模型训练和调优阶段，使用大量数据集开展模型预训练或优化训练，以提升模型性能或开展安全对齐。预训练阶段使用海量无标注语料；优化训练阶段通过有监督微调（SFT）、基于人类反馈的强化学习（RLHF）等方式进行模型优化或价值观对齐。

- a) 按训练阶段划分，大模型训练场景分为以下三种子类型：
  - 1) 预训练阶段：使用网页文本、书籍、百科、专业数据集等无标注语料开展基础模型训练；
  - 2) 监督微调阶段（SFT）：使用指令-回答对数据开展模型能力优化训练，数据可能来源于真实用户对话或业务记录；
  - 3) 对齐训练阶段（RLHF）：使用标注人员对模型输出的偏好评价数据开展价值观对齐训练，以及基于用户反馈进行优化训练的数据。
- b) 场景相关方包括：数据提供方、数据标注方、模型训练方（数据需求方）、环境提供方等。

#### A.3.2 数据去标识化处理

大模型训练场景的数据去标识化处理措施包括但不限于：

- a) 对训练数据中直接标识符的识别，可采用基于规则的正则匹配与基于命名实体识别（NER）模型的自动识别相结合的方式，辅以人工抽样校验；
- b) 对具有固定格式的直接标识符（如手机号、身份证号、电子邮件地址等），采用删除或占位符替换处理；
- c) 对非结构化文本中的直接标识符（如人名、详细地址等），采用删除或同类虚构值替换处理；
- d) 对个人信息密度过高、无法有效处理的文档，进行记录级过滤删除；
- e) 对 SFT 数据及用户反馈数据中涉及的个人信息，按照 b)~d) 进行处理；
- f) 对标注人员身份信息进行假名化处理；
- g) 数据去标识化处理完成后，应进行抽样人工复核，记录处理质量指标。

#### A.3.3 流通环境安全保障

流通环境安全保障措施包括但不限于：

- a) 满足 7.2 规定的基本安全技术能力要求；训练完成后，及时销毁输入数据及训练过程产生的临时数据；
- b) 对不同类型的训练数据采用相应的加密保护：
  - 1) 标注类文件数据，采用文件级透明加密；

- 2) 数据仓库类语料数据，采用数仓级透明加密；
- 3) 数据库类日志数据，采用数据库级透明加密。
- c) 涉及敏感个人信息的模型训练，应采用差分隐私技术（DP-SGD）对训练梯度添加随机噪声，参数  $\epsilon$  的取值建议见附录 D；
- d) 涉及对第三方大模型精调的，应在安全可控环境中训练，建立训练数据清单，采用数据加密、可信计算等技术保护训练过程数据；
- e) 接入检索增强（RAG）能力的，应对检索库数据实施加密保护与访问控制，检索库中的个人信息数据应按 A.3.2 处理后方可存入；
- f) 应对数据提供方、数据标注方实施分类分级、权限管控、操作审计等管理措施；应采用数据血缘追踪、数字水印等技术，支持全流程追溯与审计。

#### A.3.4 结果审计

结果审计应同时覆盖数据维度和模型维度，包括：

- a) 对去标识化处理后的训练数据进行抽样检查，记录个人信息残留情况；残留率超出可接受水平的，应重新开展数据处理；
- b) 对训练完成的模型开展以下测试：
  - 1) 成员推断攻击测试（Membership Inference Attack）：评估模型对训练数据的整体记忆程度；
  - 2) 训练数据提取攻击测试（Training Data Extraction Attack）：评估模型是否会复现训练数据中含有个人信息的文本；
  - 3) 模型输出个人信息泄露抽样审计：对模型在典型输入下的输出内容进行抽样，检测是否包含可识别特定自然人的信息；
- c) 模型维度审计发现显著个人信息泄露风险的，应采取模型遗忘（Machine Unlearning）、差分隐私再训练等措施处置，或限制模型对外部署范围；
- d) 记录并留存上述审计过程和结果，保存期限不少于 3 年。

附 录 B  
(资料性)  
常见去标识化技术和匿名模型

## B.1 常见去标识化技术

常见数据去标识化技术如表B.1所示，具体技术描述、适用场景及注意事项如下：

- a) 记录抑制 (Record Suppression) :
  - 1) 描述：删除或置空整个记录，防止敏感数据暴露。
  - 2) 适用场景：去除异常值或无法满足匿名化标准的记录。
  - 3) 注意事项：可能对数据集的统计属性（如平均值、中位数、方差等）造成影响。
- b) 属性抑制 (Attribute Suppression) :
  - 1) 描述：删除或置空整个属性（列），通常是直接标识符或不必要的字段。
  - 2) 适用场景：删除对数据分析无用但可能泄露隐私的直接标识符（如姓名、身份证号）。
  - 3) 注意事项：应确保数据源在导出时就不包含不必要的属性。
- c) 字符遮掩 (Character Masking) :
  - 1) 描述：用特定符号（如“\*”或“#”）替换部分字符，遮掩敏感信息，如“123456”变为“123\*\*\*”。
  - 2) 适用场景：当属性的部分信息足以满足用途，且遮掩部分能提供所需的匿名化程度。常用于展示电话号码、身份证号等敏感字段。
  - 3) 注意事项：需要考虑原始数据的长度是否会泄露信息。
- d) 假名化 (Pseudonymization) :
  - 1) 描述：用生成的假名替换标识符，可分为可逆和不可逆两种。
  - 2) 适用场景：需要保持数据记录间的关联性（如用户行为数据分析），但无需识别真实身份。例如将用户姓名替换为唯一的编码。
  - 3) 注意事项：需安全存储映射关系，防止泄露导致隐私风险。
- e) 泛化 (Generalization) :
  - 1) 描述：降低数据值的精度，将具体值替换为范围或类别，例如“25 岁”变为“20-30 岁”。
  - 2) 适用场景：当数据精确值对分析不关键，而泛化后仍能满足用途的场景。常用于地理位置（如具体地址替换为城市）或数值数据（如收入区间）。
  - 3) 注意事项：泛化范围需合理选择，避免过度泛化导致数据失去实用性。泛化层级过低可能无法有效保护隐私，过高则可能削弱数据分析效果。
- f) 置换 (Swapping) :
  - 1) 描述：在不同记录之间交换属性值，从而打破属性与个体的直接关联。
  - 2) 适用场景：数据分析只需属性层面的统计，而不需保留个体级别的完整关联性。例如交换患者病历中的诊断结果和地理位置。
  - 3) 注意事项：确保置换过程的随机性，防止被推断还原。注意保留原始数据的统计规律，避免置换引入偏差。
- g) 扰动 (Perturbation) :
  - 1) 描述：通过向数据值添加微小随机变化，增加隐私保护，如在收入数据上加减随机数。
  - 2) 适用场景：当数据的精确值不重要，而允许一定程度的误差。适用于分析中对数据分布更关注的场景，如聚类或回归模型。

3) 注意事项：扰动幅度需与数据属性的范围相适应，过大可能失真，过小则保护不足。可能影响某些分析模型的精度，需权衡隐私与分析质量。

h) 聚合（Aggregation）：

- 1) 描述：将多个记录的数据汇总为统计值，如总数、平均值或中位数。
- 2) 适用场景：当数据分析仅需汇总信息，而不需要个体记录的具体内容。常用于公开数据的发布，如人口普查数据的统计结果。
- 3) 注意事项：无法聚合的属性可能需要单独处理或删除。聚合过程中可能需要创造新的属性以保存统计结果，确保数据完整性。

表 B.1 常见去标识化处理技术

技术	描述	适用场景	注意事项/限制	输出数据类型	数据记录级保真性	适用数据类型	适用属性类型
记录抑制	删除或置空整个记录	去除异常值或无法满足匿名化标准的记录	可能对数据集的统计属性（如平均值、中位数、方差等）造成影响	微数据	是	所有	所有
属性抑制	删除或置空整个属性（列）	删除对数据分析无用但可能泄露隐私的直接标识符（如姓名、身份证号）	应确保数据源在导出时就不包含不必要的属性	微数据	是	所有	直接标识符
字符遮掩	用特定符号（如“*”或“#”）替换部分字符	当属性的部分信息足以满足用途，且遮掩部分能提供所需的匿名化程度	需要考虑原始数据的长度是否会泄露信息	微数据	是	所有	直接标识符/准标识符
假名化	用生成的假名替换标识符	需要保持数据记录间的关联性，但无需识别真实身份	需安全存储映射关系，防止泄露导致隐私风险	微数据	是	所有	直接标识符
泛化	降低数据值的精度	当数据精确值对分析不关键，而泛化后仍能满足用途的场景	泛化范围需合理选择，避免过度泛化导致数据失去实用性	微数据	是	连续数据/分类数据	准标识符
置换	在不同记录之间交换属性值	数据分析只需属性层面的统计，而不需保留个体级别的完整关联性	确保置换过程的随机性，防止被推断还原	微数据	否	所有	所有
扰动	通过向数据值添加微小随机变化	当数据的精确值不重要，而允许一定程度的误差	扰动幅度需与数据属性的范围相适应	微数据	否	连续数据	准标识符
聚合	将多个记录的数据汇总为统计值	当数据分析仅需汇总信息，而不需要个体记录的具体内容	无法聚合的属性可能需要单独处理或删除	统计数据	否	所有	所有

## B.2 常见匿名模型

匿名模型旨在通过对数据进行合理的变换或约束，减少隐私泄露的风险，常见匿名模型如表B.2所示，具体模型描述、应用及其限制如下：

### a) $k$ -匿名 ( $k$ -anonymity)：

- 1) 描述：数据集中的每个记录在其准标识符（如性别、年龄、邮编等）的取值上与至少  $k-1$  个其他记录相同，即任何一个记录无法通过准标识符被唯一定位，从而防止单独挑出和关联攻击。
- 2) 应用：作为隐私保护的基本风险阈值，用于衡量数据集抵御再识别攻击的能力。广泛应用于需要分享或公开的敏感数据中，例如医疗数据或用户行为数据。
- 3) 限制：无法防止同质性、背景知识攻击。

注1：同质性攻击是指当一个等价类中的敏感属性值过于单一时，攻击者仍然可能通过推断确定目标的敏感信息。

注2：背景知识攻击是攻击者利用额外的背景知识（如已知某人属于某等价类）推断敏感信息。

### b) $l$ -多样性 ( $l$ -diversity)：

- 1) 描述：扩展  $k$ -匿名，要求每个等价类中目标敏感属性具有至少  $l$  种不同的值，从而增加保护强度，防止属性披露。
- 2) 应用：防止属性披露，即使攻击者知道等价类，也无法明确推测目标属性的具体值。适用于需要发布包含多种敏感信息的场景，例如金融数据和医疗数据的匿名化。
- 3) 限制：无法防止语义相似性攻击、近邻攻击、偏斜性攻击。

注3：语义相似性攻击是指即使有  $l$  种不同的属性值，但这些值在语义上高度相似（如“低收入”、“微薄收入”），攻击者仍可能推断敏感信息。

注4：近邻攻击是指当等价类中的连续型属性的值彼此间的距离非常接近时，攻击者可以通过细微差别推测目标属性。

### c) $t$ -接近性 ( $t$ -closeness)：

- 1) 描述：进一步扩展  $l$ -多样性，通过确保每个等价类中目标敏感属性的分布与整体数据集的分布接近，其差异度（通常用 EMD 或 KL 散度衡量）不超过预定阈值  $t$ ，从而提高保护精度。
- 2) 应用：防止推断攻击（Inference Attack），确保即使攻击者拥有较多知识，也难以通过属性分布差异推测出目标敏感信息。常用于隐私要求较高且数据分布要求严格一致的场景，例如人口统计学和精准医疗领域的数据共享。
- 3) 限制：当目标属性的分布本身极为不均匀时，确保  $t$ -接近性可能导致等价类分组困难，或者需要对数据进行过度修改；数据实用性下降，当  $t$  值设置得过小以确保严格保护时，可能导致数据失真过大，从而影响数据分析的实用性；计算复杂性，在多维目标属性的场景下，计算等价类分布与总体分布的接近性会显著增加算法的复杂性。

### d) 差分隐私 (Differential Privacy)：

- 1) 描述：向查询结果中添加随机噪声，确保单个记录的存在与否不会显著影响输出结果，从而保障个体隐私。噪声的大小由隐私预算参数  $(\epsilon, \delta)$  控制， $\epsilon$  值越小，隐私保护程度越高，但数据效用也越低； $\delta$  通常设置为一个非常小的值，表示可接受一定概率的隐私泄露。
- 2) 适用场景：适用于统计查询，例如计数、求和、平均值等。在发布统计数据时，可以采用差分隐私技术，在保证数据可用性的同时，保护个体隐私。

- 3) 注意事项：需要仔细选择隐私预算参数( $\epsilon, \delta$ )，平衡隐私保护和数据效用。 $\epsilon$  值过小会导致数据效用过低，而  $\epsilon$  值过大会削弱隐私保护效果。

表 B.2 常见匿名模型

模型	描述	适用场景	注意事项/限制	输出数据类型	适用属性类型
$k$ -匿名	数据集中的每个记录在其准标识符上与至少 $k-1$ 个其他记录相同，防止身份披露	防止链接攻击	无法防止同质性攻击、背景知识攻击	微数据	准标识符
$l$ -多样性	每个等价类中目标敏感属性具有至少 $l$ 种不同的值，防止属性披露	防止同质性攻击	无法防止语义相似性攻击、近邻攻击、偏斜性攻击	微数据	敏感属性
$t$ -接近性	每个等价类中目标敏感属性的分布与整体数据集的分布接近，防止属性披露	防止推断攻击	当目标属性的分布本身极为不均匀时，确保 $t$ -接近性可能导致等价类分组困难	微数据	敏感属性
差分隐私	向查询结果中添加随机噪声，确保单个记录的存在与否不会显著影响输出结果	适用于统计查询，例如计数、求和、平均值等	需要仔细选择隐私预算参数，平衡隐私保护和数据效用	统计数据	所有



## 附 录 C

(资料性)

## 数据处理协议中匿名化条款参考示例

适用场景：个人信息流通过程中，如合作双方决定对数据进行匿名化处理，建议参考如下模版，在相关的数据处理协议中增加相应条款。

【条款一】双方为实现本协议目的对相关个人信息进行处理时，应当严格遵守数据保护法的规定。为确保个人信息处理的合法性与安全性，双方一致同意采取相应技术手段，对本协议所涉及的个人信息进行匿名化处理。

【条款二】本协议项下双方拟处理的个人信息及匿名化相关情况：

个人信息处理基本情况

- 数据类型/字段：如手机号、设备标识符
- 数据量级：
- 数据处理目的：
- 数据处理方式：
- 数据传输方式：
- 数据展示方式：如个体粒度还是群体粒度
- 是否会引入第三方进行数据处理：如是，需要继续填写
  - 第三方的名称；
  - 第三方在数据处理活动中的角色/地位或作用；
    - 第三方处理数据的类型/字段；
    - 第三方处理数据的目的；
    - 第三方处理数据的方式。

匿名化处理情况及说明

- 拟进行匿名化处理的数据类型/字段：
- 拟采取的匿名化技术：
- 拟配备的安全计算环境：
- 拟实现的匿名化效果：

【条款三】：双方承诺

• 双方承诺对匿名化处理后的数据，将严格限制在本协议约定的数据类型、处理目的范围之内进行处理。

• 双方承诺对匿名化处理后的数据，在合理考虑现有技术、实施成本的前提下，不会以任何方式破坏其匿名化效果，不会通过任何方式利用匿名化处理后的数据识别特定自然人，亦不会以任何方式对匿名化处理后的数据进行复原。

• 双方承诺已建立严格的数据访问及处理的安全/合规评估机制、权限管控机制以及数据隔离机制，对匿名化数据与各自持有的可能与匿名化数据建立关联并可能导致破坏其匿名化效果的数据进行严格的隔离，并通过上述评估与权限管控机制防止匿名化效果被破坏。

• 双方承诺针对匿名化处理后的数据，已经建立了相应的匿名化效果监督机制以及风险应急响应与处置机制。双方会定期对数据的匿名化效果开展评估，一旦发现数据的匿名化效果受到减损，将及时启动响应机制并寻求有效的处置措施，以便恢复/加固数据的匿名化效果。

- 双方承诺对各自员工定期进行个人信息保护、数据安全以及匿名化相关的培训，不断提升双方员工的数据合规意识。

**【条款四】：**违约责任：可与主协议的违约责任进行关联

附 录 D  
(资料性)  
不同风险场景下的匿名模型参数建议

不同风险场景下的匿名模型参数建议见表D.1。

表 D. 1 不同风险场景下的匿名模型参数建议

场景	高风险场景			中风险场景			低风险场景		
处理程度	完全处理	充分处理	部分处理	完全处理	充分处理	部分处理	完全处理	充分处理	部分处理
$k$ -匿名 ( $k$ )	—	$\geq 20$	$\geq 5$	—	$\geq 5$	$\geq 3$	—	$\geq 3$	$\geq 2$
$l$ -多样性 ( $l$ ) (如适用)	—	$\geq 5$	$\geq 3$	—	$\geq 3$	$\geq 2$	—	$\geq 2$	—
$t$ -接近性 ( $t$ ) (如适用)	—	$\leq 0.05$	$\leq 0.10$	—	$\leq 0.10$	$\leq 0.20$	—	$\leq 0.20$	—
差分隐私 $\varepsilon$ (如适用)	$\leq 0.5$	$\leq 1$	$\leq 2$	$\leq 1$	$\leq 2$	$\leq 8$	$\leq 2$	$\leq 8$	$\leq 10$

## 附 录 E

### （资料性）

### 重识别风险度量指标计算方法

#### E.1 概述

本附录旨在为评估匿名化处理后数据的重识别风险提供三种量化度量指标。这三种指标分别对应不同的攻击者动机和攻击场景，评估数据主体身份被重新识别的可能性。指标的选择应基于对数据流通场景、数据接收方能力以及潜在攻击动机的综合判断。

$D_S$ : 经匿名化处理后，拟对外披露的数据集（样本集）。

$n$ : 披露数据集  $D_S$  中的记录总数。

$D_P$ : 攻击者可能掌握的、包含身份信息的外部识别数据库（可视为总体或背景知识库，如选民名单、公开登记册等）。

$N$ : 外部识别数据库  $D_P$  中的记录总数。

$j$ : 特定等价类的索引。

$f_j$ : 在披露数据集  $D_S$  中，第  $j$  个等价类包含的记录数量。

$F_j$ : 在外部识别数据库  $D_P$  中，第  $j$  个等价类包含的记录数量。

#### E.2 检察官风险

##### E.2.1 场景定义

检察官风险（Prosecutor Risk）场景假设攻击者（如检察官）具有明确的攻击目标。攻击者事先知道某个特定个体（如被告）的信息一定存在于披露数据集  $D_S$  中，其唯一目的是在该数据集中准确地找到属于该特定个体的记录。

##### E.2.2 计算方法

检察官风险衡量的是，在已知目标存在的情况下，成功匹配该目标记录的最高概率。该风险取决于披露数据集中规模最小的等价类。

##### E.2.3 计算公式

检察官风险（ $R_{prosecutor}$ ）的计算公式如下：

$$R_{prosecutor} = \frac{1}{\min(f_j)}$$

其中， $\min(f_j)$  是披露数据集  $D_S$  中所有等价类规模的最小值。

##### E.2.4 说明

该风险完全基于披露数据集  $D_S$  本身进行计算，无需外部数据库信息。它代表了数据集中最脆弱记录的“最坏情况”下的重识别风险。对于满足  $k$ -匿名（ $k$ -anonymity）要求的数据集，由于所有等价类的规模  $f_j \geq k$ ，因此其检察官风险  $R_{prosecutor} \leq 1/k$ 。该指标适用于攻击目标明确、且能确认目标存在于数据集内的场景。

#### E.3 记者风险

##### E.3.1 场景定义

记者风险（Journalist Risk）场景假设攻击者（如记者）的目标不是寻找特定个体，而是证明该匿名

化数据集存在隐私泄露风险。攻击者试图通过将披露数据集 $D_S$ 与其掌握的外部识别数据库 $D_P$ 进行链接，成功重识别出任意一个个体，并以此作为新闻报道来质疑数据发布方的可信度。

### E.3.2 计算方法

记者风险衡量的是，攻击者通过链接外部数据库能成功重识别任意个体的最高概率。该风险取决于外部识别数据库（总体）中规模最小的等价类。

### E.3.3 计算公式

记者风险（ $R_{journalist}$ ）的计算公式如下：

$$R_{journalist} = \frac{1}{\min(F_j)}$$

其中， $\min(F_j)$ 是在外部识别数据库 $D_P$ 中，与 $D_S$ 存在交集的等价类里，规模最小的等价类的记录数。

### E.3.4 说明

计算记者风险需要了解或估计总体数据 $D_P$ 的分布情况。

在实践中，数据发布方通常无法获取 $D_P$ ，因此 $F_j$ 的值需要通过统计方法（如泊松分布、对数线性模型等）基于样本数据集 $D_S$ 进行估计。

该指标反映了数据集中最独特个体（相对于总体而言）的重识别风险，是比检察官风险更严格的度量。

适用于攻击者可访问外部背景知识库，并试图通过重识别任意个体来制造影响的场景。

## E.4 营销者风险

### E.4.1 场景定义

营销者风险（Marketer Risk）场景假设攻击者（如营销人员）的目标是最大化重识别的记录总数，而不是专注于某个高风险的个体。攻击者希望将披露数据集 $D_S$ 与其掌握的外部识别数据库 $D_P$ 进行匹配，以获取尽可能多数据主体的身份信息用于后续营销活动。攻击者关心的是整体匹配的成功率。

### E.4.2 计算方法

营销者风险衡量的是披露数据集中能够被正确重识别的记录的期望比例。它是个体记录被正确匹配概率的加权平均值，而不是最坏情况下的风险。

### E.4.3 计算公式

营销者风险（ $R_{marketer}$ ）的计算公式如下：

$$R_{marketer} = \frac{1}{n} \sum_j \frac{f_j}{F_j}$$

对于 $D_S$ 中等价类 $j$ 里的任意一条记录，它在 $D_P$ 中对应的等价类 $j$ 里有 $F_j$ 条记录。随机匹配时，匹配正确的概率是 $1/F_j$ 。等价类 $j$ 里共有 $f_j$ 条记录，因此在等价类 $j$ 中期望能正确匹配的记录数是 $f_j \times (1/F_j)$ 。

### E.4.4 说明

与记者风险类似，营销者风险的计算也需要关于总体数据 $D_P$ 的信息（ $F_j$ ），通常需要进行统计估计。该指标评估的是整个数据集的平均风险水平，而非单个记录的风险。即使数据集中不存在唯一的记录（即记者风险较低），如果存在大量小型等价类，营销者风险也可能很高。适用于攻击者目标是批量获取身份信息，对整体成功率敏感，且能容忍一定错误匹配成本的场景。

E. 5 重识别风险指标对比

重识别风险指标对比见表 E. 1。

表 E. 1 重识别风险指标对比

指标类型	攻击者目标	计算所需信息	公式	风险说明
检察官风险	识别已知存在的特定个体	仅披露数据集 ( $D_S$ )	$1/\min(f_j)$	数据集内最脆弱个体的最坏情况风险
记者风险	识别任意一个个体以证明风险	披露数据集 ( $D_S$ ) 和外部数据库 ( $D_P$ )	$1/\min(F_j)$	相对总体最独特个体的最坏情况风险
营销者风险	最大化重识别的记录总数	披露数据集 ( $D_S$ ) 和外部数据库( $D_P$ )	$\frac{1}{n} \sum_j \frac{f_j}{F_j}$	整个数据集可被成功重识别的期望比例（平均风险）