

# SECURITY



技术版 ▶▶ 与安全人士分享技术心得 Share technique experience with security professionals



**动态可控**  
业务数据

**动态可控**  
运维数据

**静态可知**  
数据

**动态可控**  
办公数据

★ 本期焦点

监管合规下的个人信息安全保障

数据安全产品分析

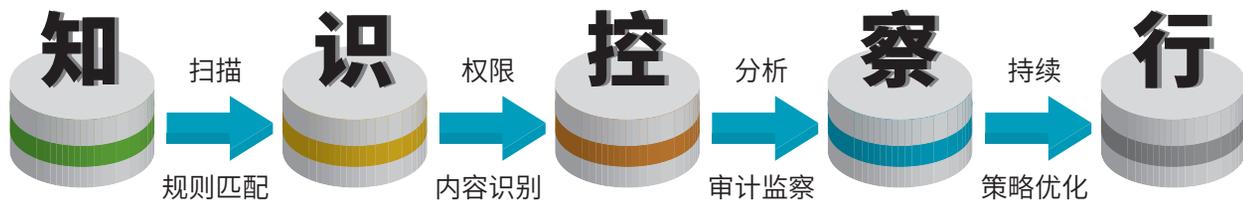
绿盟数据安全整体解决方案

基于大数据分析的敏感数据检测及响应方案

绿盟科技官方微信



数据安全治理方法论



■ 定义敏感数据

■ 发现敏感数据

■ 控制敏感数据

■ 监督敏感数据

■ 持续运营服务



主办: 绿盟科技  
策划: 解决方案中心  
地址: 北京市海淀区北洼路4号益泰大厦三层  
邮编: 100089  
电话: (010) 6843 8880-5438  
传真: (010) 6843 7328  
网址: www.nsfocus.com

数据安全专刊

欢迎您扫描封面左下角的二维码, 关注绿盟科技官方微信, 分享您的建议和评论, 或者来信nsmagazine@nsfocus.com 与我们交流。



© 2019 绿盟科技

本刊图片与文字未经相关版权所有人书面批准, 一概不得以任何形式、方法转载或使用。本刊保留所有版权。

数据安全专刊是绿盟科技的专用商标。

需要获取更多信息, 请访问WWW.NSFOCUS.COM

▪ 卷首语	李晨	4
▪ 市场解读		5-20
监管合规下的个人信息安全保障	孙昌卫	5
浅谈数据保护中的隐私问题	刘潇奕	9
数据安全产品研究	施岭	13
隐私影响评估 (PIA) 介绍	刘宇	16
▪ 解决方案		21-44
绿盟数据安全整体解决方案	施岭	21
数据安全咨询服务介绍	贾晓萍 刘宇	28
大数据安全的解决思路	孙叶	32
金融行业数据治理方案	李迪	38
▪ 技术 & 应用		45-79
基于大数据分析的敏感数据检测及响应方案	梁莎 李景 皮靖 吴天昊	45
特权访问管理下的数据安全如何保护	许德昭	51
数据库审计产品的技术运用	梁步庭	56
浅析 AWS 数据安全保护措施	刘弘利	61
数据内容识别技术深度剖析	施岭	65
大数据环境安全管控浅析	王豪	70
大数据安全之敏感数据发现	肖春亮	75

在数字化时代，Gartner 将网络安全重新定义为“数字安全”，是指数字化转型中的网络安全。

数据安全 (Data Security) 有两方面的含义，一是数据本身的安全，主要是指采用现代密码算法对数据进行主动保护，如数据保密性、数据完整性、双向强身份认证等。二是数据防护的安全，主要是采用现代信息识别手段对数据进行主动防护，如通过协议识别、类型识别、内容识别、机器学习等智能化手段保证数据的安全。

尽管《网络安全法》、《GDPR》等法规政策相继出台，但过去的 2018 年遭媒体曝光的数据泄露事件数量仍然远超过往的年份，个人敏感数据遭泄露的人数需要以亿计算。“数据安全”已连续三年蝉联 RSA 热词榜冠军，而国内的数据安全市场仍处于成长期，很多投资商依然在观望，但按照数据泄露事件数量激增、性质不断恶化的形势发展下去，国内数据安全市场井喷之日很快就会到来，据中商产业研究院整理统计，预计 2020 年中国数据安全产业规模预计达 70.2 亿元。

随着数据的激增与共享，智慧城市的不断发展，云计算、大数据、物联网得到了更为广泛的应用，大数据的安全也成为了企业关注的焦点，IDC 估测，在大数据时代，数据将以每年 50% 的速度增长。

绿盟科技一直对数据安全有着持续的研究，本刊中将我们现阶段的研究成果做了汇总与阐述，希望可以与大家共同探讨数据安全的最佳实践。

绿盟科技 副总裁 李晨

# 监管合规下的个人信息安全保障

解决方案中心 孙昌卫

摘要：当今社会，数据的重要性和价值被越来越重视，个人信息因为和公众密切相关而被社会各界聚焦关注，本文介绍了国内外典型的个人数据法规、标准，并提出安全建议，供大家参考。

据《中国互联网络发展状况统计报告》统计显示，截至 2018 年 6 月，中国网民规模为 8.02 亿，较 2017 年末增加 3.8%，互联网普及率达 57.7%，其中手机网民规模达 7.88 亿，在上网人群中的占比达 98.3%。

高比例的上网人群催生出来大量的应用程序，这些应用程序为广大网民提供各式各样的便捷服务，但其中也出现了对个人信息违法收集、滥用、泄露等问题，这不仅会导致个人隐私数据泄露，同时还严重影响个人生命和财产安全。另外随着国家政务系统的集约化，提供公共服务的应用系统不断增加，数据共享使用、融合存储，数据资源集中后高价值明显，公民个人信息在此过程中面临的安全风险也在日益剧增。

网络安全是一个全球性的问题，其中公民个人信息保护问题面临的挑战更加严峻，也更加复杂。各国政府一直致力于公民个人信息保护，但依然不断出现个人隐私数据泄露事件，例如 2018 年 3 月 Facebook 出现的数据泄露事件，导致 5000 万用户信息被第三方公司 Cambridge Analytica 用于大数据分析，并引发连锁反应。2018 年 8 月华住集团“5 亿条个人敏感信息泄露”，全部数据泄露资料更是高达 141.5GB，9 月份警方将犯罪嫌疑人抓获归案，同时也将依法查处涉案的主体单位。可见在落实网络安全主体责任和安全防护措施方面我们丝毫不能懈怠，需要切实加强安全防护保障措施，在防护效果方面达到“进不来、看不懂、拿不走、改不了、赖不掉”。

从法规保护层面而言，目前各个国家国情不同、重视程度不一样，立法的进度和保障措施也不尽相同，其中以欧盟为代表的组织和国内在个人信息保护方面存在着明显差异。

2018年5月25日，欧盟通用数据保护条例（General Data Protection Regulation, GDPR）正式实施，在全球范围内产生广泛影响。GDPR提出了个人数据保护六大原则：合法合理透明性原则、目的限制原则、最小化数据处理原则、数据准确性原则、限制存储期限原则、数据的完整性和保密性原则，在六大原则的指导下，通过一系列严格的问责机制，从系统设计和默认设置着手的隐私保护（Data protection by design and by default）、保留处理活动记录、实施安全保障措施、数据泄露报告与通知、数据保护影响评估、事先协商、设置数据保护官等措施，对数据主体的知情权、访问权、纠正权、删除权（被遗忘权）、限制处理权、可移植权（可携带权）、拒绝权和与自动化个人决策相关权利进行保护。

GDPR要求数据控制者和处理者实施适当的技术和管理措施，并在必要时进行审查和更新，尽管全文并未给出详细的控制措施实施要求，但仍指出需对以

下因素进行着重考虑：

- 个人数据的匿名化和加密；
- 数据系统持续保持保密性、完整性、可用性以及恢复的能力；
- 在发生自然事故或者技术事故的情况下，个人信息的及时获取以及再存储能力；
- 对技术性及管理性措施的有效性定期进行测试、访问、评估，以确保处理过程的安全性。

需特别注意的是，GDPR关于“同意”的认定标准较以往更为严格，且对儿童个人信息的保护更为注重。

### **国内数据安全法律法规、政策标准**

我国在多项法律中关注和强化对个人信息的保护。如2012年全国人大常委会通过了《关于加强网络信息保护的決定》；2015年《中华人民共和国刑法修正案（九）》中明确了对个人信息保护的规定；2016年《中华人民共和国网络安全法》确定了个人信息保护的基本规则。2017年《中华人民共和国民法总则》中也明确规定了自然人的个人信息受法律保护。2019年《中华人民共和国电子商务法》中也纳入了保护消费者个人信息等规定。除此以外《网络安全等级保护条例（征求意见稿）》和《关键信息基础设施安全保护条例（征

求意见稿)》中也对个人信息保护进行了相关规定,要求网络运营者落实重要数据和个人信息安全保护制度,采取保护措施,保障数据和信息在收集、存储、传输、使用、提供、销毁过程中的安全。

其中《中华人民共和国网络安全法》在第四章网络信息安全部分,较大篇幅的对个人信息安全进行了规定,并且明确规定了“任何个人和组织不得窃取或者以其他非法方式获取个人信息,不得非法出售或者非法向他人提供个人信息”。为了落实《中华人民共和国网络安全法》《消费者权益保障法》,2019年1月25日,中央网信办、工业和信息化部、公安部、市场监管总局正式对外发布《关于开展App违法违规收集使用个人信息专项治理的公告》。从中我们也看出政府治理违规收集使用个人信息方面的决心,专项治理公告中明确要求:

- App运营者收集使用个人信息时要严格履行《中华人民共和国网络安全法》规定的责任义务,对获取的个人信息安全负责,采取有效措施加强个人信息保护。

- 相关部门依法编制大众化应用基本业务功能及必要信息规范,并对用户基数大、民众密切相关的

App隐私政策和个人信息收集使用情况进行评估。

- 主管部门将加大对违法违规收集使用个人信息行为的监管和处罚力度,包括依法暂停相关业务、停业整顿、吊销相关业务许可证或者吊销营业执照等处罚手段。

- 开展App个人信息安全认证,同时也鼓励运营者通过App个人信息安全认证。

在个人信息安全标准方面,2018年5月1日,信安标委组织制定的《信息安全技术 个人信息安全规范》正式实施,并于2019年1月30号公布二次修订草案。这是国内在个人信息安全保障方面的提升,在这部重磅的个人信息安全标准中明确了个人信息安全的基本原则,个人信息的收集、保存、使用、委托处理、共享、转让、公开披露的要求,个人信息安全事件处置和对组织的管理要求。同时相关的配套法规、标准也在陆续制定当中,相信推出时间指日可待。

虽然国内针对个人信息安全保护有一系列的法律法规、政策标准。但是总体而言法规、标准依然不完善,缺少总体规划,碎片化明显,同时存在落地执行不到位等情况。需要我们在立法层面加强顶层设计,统一

规划，紧密衔接，加强监管和处罚措施，不断完善现有管理机制，从国家层面为个人信息安全提供法律保障。

### **对个人信息安全的几点建议**

个人信息安全不单纯是技术问题或者法律问题，需要统筹结合，上下联动，综合防御。各组织单位、行业客户需要梳理清楚自身数据资产，识别个人敏感信息，加强内部安全管理措施，数据在哪里，安全保障就要覆盖到哪里。

在个人信息安全防护方面，行业单位需要落实国家网络安全等级保护制度。在安全合规建设中及时整改、落实管理，构建动态防御体系。加强数据安全动态监测预警机制，加强信息收集、分析研判，个人信息安

全事件发生时及时预警、迅速处置。对接触到个人敏感信息的人员，进行鉴权控制、行为审计，并定期组织安全培训，强化素质教育、安全意识教育、安全技能教育，减少敏感数据从内部泄露的风险。敏感数据定期备份，备份信息定期离线保存，避免数据勒索事件发生。加强普法教育，个人信息安全从身边的小事做起，不随意丢弃快递单、不随便参加扫描抽奖活动、使用 App 谨慎放权。

面对个人信息安全，没有谁可以独善其身，加强个人信息安全保护不仅需要国家立法保障，更加需要行业单位加强落地执行，不断完善监督、检查、处罚机制，同时也需要公民积极参与。只有社会各界共同努力，才能构建风清气朗的网络空间。

# 浅谈数据保护中的隐私问题

解决方案中心 刘潇奕

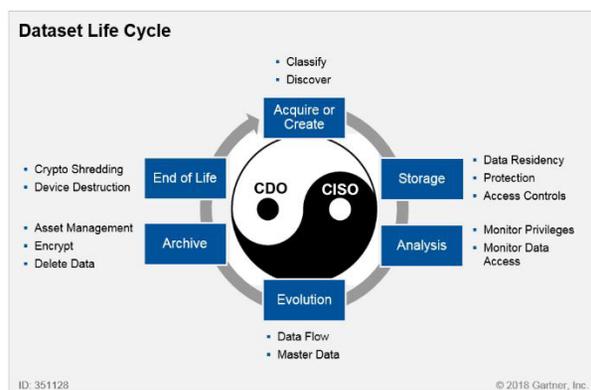
关键词：数据安全生命周期、GDPR、隐私数据、数据安全治理

摘要：隐私及敏感数据是数据保护的关键，因此首席信息安全官(CISO)发挥着越来越重要的作用。

CISO 的职责包括在数据生命周期的各个阶段保护数据，并验证第三方策略和实践的合规性。CISO 需要合作以确保满足客户需求和期望，并且应与执行职能部门协同工作。同时，GDPR 对隐私的要求需要 CISO 和数据保护官 (DPO) 之间的重要合作，并要求进行数据保护影响评估 (DPIA) 来解决这些隐私问题。

## 一、数据集生命周期的管理以及 GDPR 对隐私问题的影响

每个数据集的生命周期将从几分钟到几十年不等，并对生命周期管理产生不同的需求。在 Gartner 题为《How to Use the Data Security Governance Framework》的报告中对数据集生命周期的介绍：



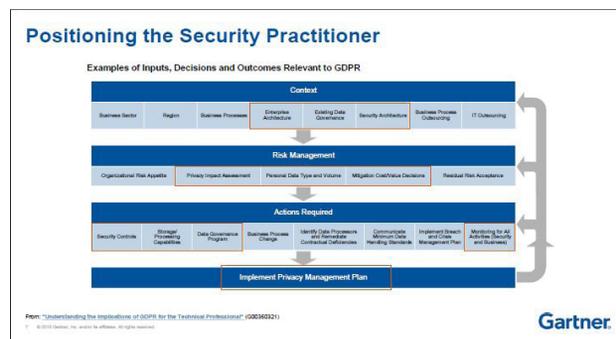
由上图可见，数据集的生命周期依次为获取或创

建、存储、分析、演进、归档、销毁这六个阶段。其中，数据集的获取或创建包括分类、发现；存储包括数据留存、保护、访问控制；分析包括监控权限、监控数据访问；演进包括数据流、主要的数据库；归档包括资产管理、加密、删除数据；一旦加密或设备损坏生命周期即销毁。每个数据集的生命周期将从几分钟到几十年不等，并对生命周期管理有不同的需求。然而，这也为首席数据官(CDO)和首席信息安全官(CISO)团队之间的合作创造了紧迫性，因为许多生命周期问题都是重叠的。国际上通过数据驻留以及数据保护和隐私法（例如欧盟的通用数据保护法规 (GDPR) 的影响来创造风险。根据 Gartner 的总结，GDPR 数据处理风险由低至高分别为处理个人数据、处理大量个人数据、处理敏感个人数据、处理大量敏感数据、（试

图去) 预测主体的行为或行动、对主体进行预测或做决定、对所做的(预测性的)决定使用人工智能。值得一提的是, 隐私及敏感数据是数据保护的关键, 因此首席信息安全官(CISO)发挥着越来越重要的作用。CISO 的职责包括在数据生命周期的各个阶段保护数据, 并验证第三方策略和实践的合规性。CISO 需要合作以确保满足客户需求和期望, 并且应与执行功能部门协同工作。同时, GDPR 对隐私的要求需要 CISO 和数据保护官(DPO) 之间的重要合作, 需要采取行动和政策来保护与个人权利有关的个人数据, 这可能是与业务目标相矛盾的。

GDPR 中的第 35 条提到要求进行数据保护影响评估(DPIA) 来解决这些隐私问题。DPIA 作为设计保护不可或缺的一部分。使用的技术包括数据剖析、信息权限管理(IRM) 产品隐私扩展和专业软件。DPIA 提供需要的先决条件去记录哪些个人数据被处理。在 GDPR 授权的 DPIA 中, CISO 必须与业务利益相关者和 DPO 合作, 记录对客户和员工的潜在影响, 以及相应减轻措施。这意味着必须将 DPIA 纳入业务影响评估(BIA)。任何已识别的隐私风险都将影响业务选择和后续业务风险。在 Gartner 题为《Approaches to Data Security in a Regulated World》的报告中, 关

于安全从业者的配置, 以 GDPR 的投入、决策和产出为例, 其中包括环境、风险管理、所需行动、实施隐私管理计划这四个层面, 这些层面之间直接或间接相互作用。其中, 环境主要包括企业架构、存在的数据监管、安全架构, 风险管理主要包括隐私影响评估、个人数据类型和容量、减缓成本/价值决定。



对于数据安全领导者面临的挑战, Gartner 建议 CISO 与 CDO 建立企业赞助关系并与主要利益相关者合作; 建立利润中心来管理数据, 例如使用信息经济学来评估风险并明确其优先级; 针对人、数据、分析, 映射存在的安全策略; 制定计划以通过管理控制台编排策略; 寻找像 DCAP 一样可提供统一方法的供应商, 这可能需要在预算、员工技能和可用性方面对于实际部署进行权衡。

## 二、敏感数据是数据保护的关键要素

缺乏关于敏感数据的情景感知,可能会让组织暴露到重大风险中。关键是要确定敏感数据是否存在于 Hadoop 中,它位于何处并且继而触发适当的数据保护措施,例如数据屏蔽、数据校订、标记化或加密。

对于进入 Hadoop 的结构化数据,例如来自数据库的有关数据,或者逗号分隔值(CSV)、JavaScript 对象符号(JSON)格式化文件,敏感数据的位置和分类可能已经知道。在这种情况下,保护那些列或者字段可以以编程方式发生。

对于非结构化数据,对敏感数据的位置、数值和分类的掌握变得难得多。敏感数据的发现和定位,成为数据保护重要的第一步。

## 三、选择供应商来处理用于安全检测的数据时的注意事项

有效地管理、理解和有效且规律地使用检测提供者和技术会付出比较大的代价。处理组织生成的大量数据,以及对其特定安全用例要求缺乏了解,这是一项艰巨的任务。

安全与风险管理(SRM)领导者应通过有效识别、确定优先级并使风险与这些数据保持一致,以此来认

识活动数据如何有助于实现业务目标的安全追求。

在选择技术供应商或服务提供商来获取、处理和存储此类日志或网络传输数据以用于安全检测功能时,领导者应考虑:在用例范围内调整时,需要哪些数据来满足服务或技术的需求;安全性概述的要求可能需要哪些其它数据,以及需要它们的人如何访问它们;哪些合规性要求会影响要存储的数据类型,数据必须存储多长时间以及存储位置。

## 四、我们需要关注物联网中的数据隐私和安全问题

据 Gartner 预测,到 2021 年,关键基础设施的监管合规性将使全球物联网安全支出达到 10 亿美元,高于今天的不到 1 亿美元。

对于想尽办法入侵企业网络或家庭网络的犯罪分子来说,物联网设备是一条全新的攻击途径。如果一个保护不当的物联网设备或传感器连接到企业网络,可能为攻击者提供了进入网络的新方法,进而他们有可能找到想非法获取的宝贵数据。物联网中的威胁包括窃取或窥探、恶意软件注入、物理干涉、供应链问题等,而涉及的数据安全隐私监管多与 GDPR、E- 隐私、HIPAA 法案等有关。为保护隐私,从分类、文件数据流、修改数据存取政策,到加密、令牌化、掩码和其它隐

私保护技术，有关人员在数据安全政策和执行方面已做出了不少努力。

整体来看，有关人员在实践时需要遵守三个基本原则：考虑运营的必要性，在适当的地方直接识别数据；输入分析时，在可操作的情况下尽量使用假名；输出也应该是匿名的，此时只与数据有关而与个人无关，同时也需要小心重新识别的风险。

因网络安全和业务风险的融合不断提升，隐私和敏感数据泄露的风险持续加大，对新一代安全领导者建立具有风险意识的文化与管理体制，建立智能的安全运维和快速威胁响应机制，以及确保一个安全易恢复的网络等提出了与时俱进的需求。安全领导者需要建立风险意识管理系统来进行识别并针对业务威胁和风险进行优先级排序，建立一项风险意识战略和企业架构，并为需要的技能与能力平衡人力资源。

网络安全是一个全球性的问题，其中公民个人信息保护问题面临的挑战更加严峻，也更加复杂。各国政府一直致力于公民个人信息保护，但依然不断出现个人隐私数据泄露事件，例如 2018 年 3 月 Facebook 出现的数据泄露事件，导致 5000 万用户信息被第三方公司 Cambridge Analytica 用于大数据分析，并引发

连锁反应。2018 年 8 月华住集团“5 亿条个人敏感信息泄露”，全部数据泄露资料更是高达 141.5GB，9 月份警方将犯罪嫌疑人抓获归案，同时也将依法查处涉案的主体单位。可见在落实网络安全主体责任和安全防护措施方面我们丝毫不能懈怠，需要切实加强安全防护保障措施，在防护效果方面达到“进不来、看不懂、拿不走、改不了、赖不掉”。

# 数据安全产品研究

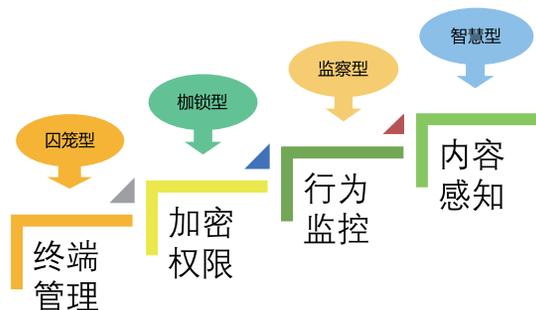
解决方案中心 施岭

关键词：数据安全、囚笼型、枷锁型、监察型、智慧型

摘要：“数字安全”的战场中涌现出了不同类型的产品和技术，每一类都有其独有的特点，专注点也各有不同。随着科技的发展，各种技术开始互相之间融合和搭配，如环境隔离、加密授权、审计监控、内容识别、机器学习等技术，融合后产品将更好地应对已知和未知的风险。

自 1986 年第一只蠕虫进入互联网开始，网络安全愈演愈烈。近十年来，由于利益的驱使，网络安全的战场中涌现出了以数据为核心的安全问题。为此，Gartner 将网络安全重新定义为“数字安全”及数字化转型中的网络安全。

在此“数字安全”的战场中涌现出了不同类型的产品和技术，主要包括囚笼型、枷锁型、监察型、智慧型四类。每一类都有其独有的特点，专注点也各有不同。随着科技的发展，机器学习技术已经逐步融入其中。



下面就详细介绍一下四种类型的数据安全产品的特点

## ▪ 囚笼型 - 数据安全产品

囚笼型产品主要特点为 **设备强管控，采用逻辑隔离手段，构建安全隔离容器。**

自 2000 年后，国外的安全管理产品相继涌入中国，刚开始是概念式引导，慢慢地转化为产品。有名的产品厂商包括 Symantec、LANDesk，2005 年至 2008 年他们在中国的市场占有率已经到了 80%。2008 年以后，随着发展，国内产品开始大量进入市场，至今国外终端管理类产品已经被国内产品大量替换。虽然市场已经呈现出饱和状态，但由于技术的不断革新、新功能的融入，再加上终端准入的重要性，每年还会有很多的用户将资金用于终端管理产品。

终端管控产品为管理带来了方便，严格的控制保障了数据只能在内部使用，但随着万物互联的共享时代的到来，数据的交互成了问题。如何在严控的同时让交互也方便，后面的类型中会找到答案。

- **枷锁型 - 数据安全产品**

枷锁型产品主要表现为 **数据强管控，提供内容源头级纵深防御能力；数据的分类、分级、加密、授权与管理。**

与终端管理不同，数据加密与权限控制产品已经将关注点从设备变化成了具体的数据文件，控制方式更加细粒度化，保密方式更优秀。从 2007 年开始至今，市场中涌现出很多有实力的优秀厂商，因为国家的监管要求，加密类产品只能在获得相关保密资质、密码认证后才可以国内使用，所以使得国外产品无法在国内大面积的销售，加密和权限类产品至今为止每年的需求量还是很大，各个行业都有数据防护的需求。虽然市场竞争激烈，但使用者还是担心数据会被加密绑架，而且是全局范围内的。不过还好目前所有产品都很成熟，很稳定。而且由于这两年的数据勒索问题层出不穷，加密厂商还有针对性的提出了防御勒索的方案，得到了市场的一致好评。

- **监察型 - 数据安全产品**

监察型的产品是 **行为强审计，利用准确关键字对数据操作行为的审计，文档的新建、修改、传输、存储、删除的行为监察。**

行为审计，分为网络行为审计和终端行为审计，网络行为审计可以有效的监控员工工作时间内的网络访问行为，而终端行为审计可以更有针对性的完成对关键数据文件的操作行为。审计产品与其他网络和终端产品共存，可以互相补充，至今市场占有率依然很高，不过随着发展，很多网络和终端产品的不断完善和提升，单独行为审计产品已经无法顺利的存活，多元化开始受到客户青睐。

当今时代数据风险的态势感知受到更多的关注，通过对各种行为的监控，将攻击行为、操作行为、流转路径、存储位置、外发途径等所有日志进行分析关联，让未知的风险提前暴露出来，在事前就可以得到合理的预防。

- **智慧型 - 数据安全产品**

智慧型产品追求 **数据智能管控，可识别、可发现、可管理，提供共性管控能力。**

为了更加全面的对数据进行管控，终端管理产品与加密权限类产品做了很多组合的方案，但都是属于全局强管控，有一定的局限性，无法应用到更加复杂

的数据环境中。在这种情况下，世界各地又不断发生着各种各样的数据泄密事件，人们对数据的重视程度就落在了内容上。这时，内容感知型DLP产品应运而生，通过内容来识别数据的重要性，通过内容来为数据进行分类，通过内容来对数据进行级别划分，智能化的管控方式也带来了便利性和灵活性。

自2013年以来，国内大力推动国产DLP产品的生产和应用，在金融行业和运营商行业更是掀起了一个潮流，但国内产品还处于一种萌芽阶段，产品的不成熟和不稳定为DLP国产化的道路带来了阻力，很多终端、加密和审计厂商开始转型，DLP产品所具有的内容识别算法得到了广泛的应用，对保护数据提供了有力的技术手段。

▪ **总结：**

数据已然成为企业、社会、个人的核心资产，在遵循法律法规的前提下对数据的共享和开放，是提高数据使用效率的基础，也是趋势使然。从安全的角度看，要保障数据共享开放，基础的要求包括数据内容真实可靠、完整唯一。

数据产生的行业、环境、位置以及角色存在巨大差异，数据安全问题极难一劳永逸，没有哪一两款产

品就能在数据全生命周期中做到面面俱到、游刃有余。随着《网络安全法》逐步推行和开展，数据安全以及数据安全产品提供商也会迎来新的机遇和挑战，技术融合也势必成为加快数据安全产品落地的新方案。

数据安全提供商应为企业建设数据内容价值评估体系，对数据的价值、敏感度进行标定，同时通过技术和法规的手段来提高监管约束，保障数据不被滥用。

# 隐私影响评估 (PIA) 介绍

安全服务部 刘宇

关键词：ISO/IEC 29134、隐私保护、PIA、风险评估

摘要：本文详细介绍了 ISO/IEC 29134 标准中实施隐私影响评估 (PIA) 的具体步骤。简单分析了当前社会环境中对隐私保护的需求并介绍了隐私保护评估 (PIA) 所涉及的相关概念。从企业的角度阐述了进行隐私影响评估 (PIA) 所能够带来的收益。通过将隐私影响评估 (PIA) 实施中的具体内容归纳为三个阶段十四个步骤，逐一介绍每个步骤中的输入内容、输出内容和实施方法。有效的为读者提供直观了解隐私影响评估 (PIA) 实施方法的途径。

## 一、介绍

我国网络环境下的隐私侵犯事件层出不穷，随着国家的重视和公民的关注，任何业务与个人信息相关的企业都应当重视其内部可能存在的隐私影响风险并采取有效手段进行处理。而用于界定隐私影响程度，判定隐私风险和提供隐私风险处理方式的有效参考办法即为隐私影响评估 (PIA)。本文所阐述的隐私影响评估 (PIA) 方法均参考由国际标准化组织所发布的 ISO/IEC 29134 标准，该标准作为国际通用标准为 PIA 的实施提供了详细且全面的指南。

当前市场环境下，诸多企业亟需通过隐私影响评估来提升自身的隐私保护能力，改变自身的隐私保护政策以及隐私保护文化。2019 年 2 月 17 日下午，京东金融就其 App 涉嫌侵犯用户隐私事件向全体用户道

歉。在该事件中，用户发现京东金融 App 客户端的本地文件夹会储存用户的手机截图，由此用户质疑京东金融 App 通过上传用户截图的行为来非法收集用户的隐私信息。京东金融对此解释为其储存功能的仅为方便客服与客户沟通。从用户的角度来看，不论其截图的使用目的为何，用户都会感受到其隐私正在遭受侵犯。该事件暴露出的问题是京东对用户隐私保护的漠视，企业在设计产品时一味注重功能并未考虑到对用户隐私安全的保护。该事件无疑对京东金融造成了极大的负面影响。京东金融 App 对用户隐私侵犯的问题并不是个例，根据南方都市报发布的 2018 个人信息保护年度报告，诸多企业都存在着用户隐私保护缺失或不规范的情况。此类企业的用户隐私问题亟待解决，而隐私影响评估 (PIA) 正是帮助企业有效避免

类似京东金融事件发生的有效手段。以京东金融为例，如果在 App 发布之前针对其 App 实施隐私影响评估 (PIA)，便可有效地发现该问题，并且基于隐私风险评估 (PIA) 中关于风险处置的方法，结合通过设计保护隐私 (privacy by design) 的思路，就可以在产品发布之前对其进行整改，最终避免因隐私保护失责导致的用户信任体系崩塌。因此，对于尚未暴露或已经暴露出隐私保护问题的企业，隐私影响评估 (PIA) 服务拥有着广阔的市场前景。

## 二、实施隐私风险评估 (PIA) 的收益

了解实施隐私风险评估 (PIA) 的各种不同收益可以使我们在进行服务实施时更好的明确服务的交付目标。

对企业实施隐私影响评估能够为企业发现、评估并处理其内部业务系统涉及个人识别信息 (PII) 的所有数据处理行为中潜藏的隐私风险。发现并处理这些隐私风险能够为企业带来多方面的收益，其主要收益体现在以下五个方面：

- 针对企业中存在的隐私风险和问题提供早期预警。

- 在企业发生隐私安全事件之后为企业提供尽责证明。
- 帮助企业获得公民对其的信任。
- 为企业管理层提供可用于决策的可靠信息。
- 在监管部门之前发现自身存在的隐私风险问题，及时处理规避处罚。

隐私风险评估 (PIA) 不仅能够为企业识别现存的隐私风险并加以处理，它还可在项目开始的早期阶段进行实施并对项目本身加以影响，以保证项目符合通过设计保证隐私 (privacy by design) 的要求。

### 隐私影响评估 (PIA) 的实施流程

隐私影响评估 (PIA) 的实施可分为三个阶段：准备阶段、实施阶段和发布评审阶段，每个阶段又分为数个具体的实施步骤如下图所示：



图 1. 隐私影响评估实施流程图

接下来我们针对隐私影响评估 (PIA) 的具体实施方法进行介绍：

### 1) 识别 PII 信息流

目标：梳理所评估目标系统的 PII 信息流；

实施方法：PIA 团队应当与企业内部以及外部人员进行访谈来收集梳理 PII 信息流的必要信息，这些信息包括但不限于：

- PII 的来源以及收集方式；
- 企业 PII 处理的负责人；
- PII 处理的目的；
- PII 的处理方式；
- PII 的保留和销毁策略；
- PII 的管理和修改策略；
- 企业对 PII 的保护方式；
- 被传输到低安全等级区域的 PII；

PIA 团队应当考虑 PII 信息流即 PII 生命周期中所存在的风险点，并将这些风险点用于后续的隐私风险评估。

### 2) 分析用例含义

目标：识别潜在的用户行为，辨别出存在隐私风险的用户行为；

实施方法：PIA 团队应分析典型用户类别的使用案

例，分析提取其中可能存在的隐私安全风险点。典型的高危用户行为有：

- 用户对正在运行设备的操作系统安全设置进行错误更改；
- 用户倾向于丢失移动设备和智能芯片；
- 对设备的误操作和对应用设置的错误理解；
- 用户对非法行为缺少防范意识；

### 3) 明确隐私保护需求

目标：明确被评估系统的隐私保护需求；

实施方法：PIA 团队应当确保被评估系统满足相应的法律条例以及合同要求的隐私保护需求；

PIA 团队应当：

- 被 PII 处理相关地法律、法规、以及使用的合同条款；
- 制定相关的信息管理系统管理标准 (如 ISO/IEC 27001 系列标准)；
- 辨认出目标系统的相关隐私要求；
- 描述已经实施或在计划中的隐私保护措施；
- 收集之前与隐私保护相关的项目中的可用信息；
- 检查 PII 主体是否被正确告知并同意企业使用其 PII 的目的；

• 检查 PII 主体是否被赋予获取、修改以及收回 PII 的能力；

#### 4) 评估隐私风险

目标: 识别、分析、以及评测目标系统中的隐私风险；

实施方法: 实施隐私风险评估需要识别、分析、和评测三个步骤。

识别隐私风险需要 PIA 团队选取符合其目标和可行性的隐私风险识别工具和技术。

隐私风险分析需要 PIA 团队针对已识别的隐私风险分析其可能造成的影响，分析其可能性与严重性。隐私风险分析可以根据 PIA 团队的信息和资源拥有情

况，以定性、半定量或定量的方式来展开

隐私风险评测要求 PIA 团队从隐私风险发生的可能性和其发生后的严重性两个维度对隐私风险进行综合评测，并根据评测结果将已识别的隐私风险定位在隐私风险图中。PIA 团队应当综合考虑企业的实际情况（包括风险承受能力）来对隐私风险进行判定，并向企业提供处置隐私风险的优先级。下图所示为隐私风险图的一个例子。

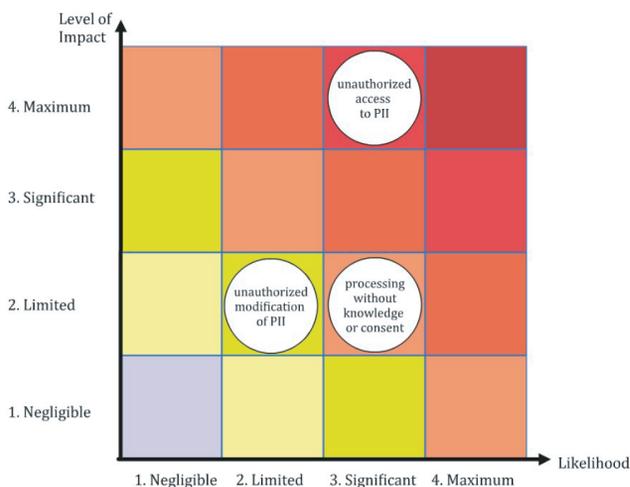


图 2. 隐私风险图

### 5) 隐私风险处理准备

目标：为已评估的隐私风险选择合适的处置方式；

实施方法：典型的风险处置方式有风险削弱、风险保留、风险规避、和风险转移。PIA 团队需要根据不同隐私风险的危害性和其特征选择合适的风险处置选项。

#### **发布评审阶段**

作为向企业管理层和利益相关人进行汇报的工具，隐私影响评估 (PIA) 报告中的内容为实施隐私影响评估各步骤中所产生结果的详细记录。由于 PIA 实施的过程信息可能包含导致企业信息系统受损的机密信息，所以产生的完整版 PIA 报告需要做保密处理。可供公开发布的 PIA 报告中不应出现与企业机密相关的信息。

### 三、总结

当前社会环境下，隐私保护是社会和公民所关注的重点问题。国家相关机构陆续颁布新的法规与标准逐步完善我国的隐私保护体系。隐私影响评估 (PIA) 是隐私保护落实的关键，它不仅能够为信息系统提供隐私风险早期预警等收益，还可以帮助企业从设计层面构建隐私保护体系 (privacy by design)。通过三个阶段数十个步骤，完整实施的隐私影响评估 (PIA) 能够有效地发现评估对象存在的隐私风险，并对隐私风险进行评估和处理。同时为企业提供的持续服务让隐私风险评估 (PIA) 成为长期保障隐私的有效机制。

# 绿盟数据安全解决方案

解决方案中心 施岭

关键词：数据安全、数据安全治理、数据共享、数据泄露、大数据

摘要：绿盟科技为数据安全设计了全面可信的防御体系，提出“知”、“识”、“控”、“察”、“行”的数据安全治理方法论，以及包括数据梳理、运维数据监管、业务数据监管、办公数据监管，以及数据的可视化的完整解决方案，有效保护数据在全生命周期过程中的安全，达到合法采集、合理利用、静态可知、动态可控的防护目标。

## 序

人类经历了三次数据量的跃升，Web 1.0 时代以门户网站为主，Web 2.0 时代用户原创内容带来“数据爆炸”，物联网时代数据上 TB、PB 级，进入“大数据时代”，IDC 估测，数据以每年 50% 的速度增长。

这是一个互联，且不断保持“在线”的社会，数据只有被共享才更有价值。

当今数据的来源复杂而多样，云计算、大数据、物联网、移动互联网、车联网、手机、平板电脑、个人电脑 (PC) 以及遍布地球各个角落的各种各样的传感器，无一不是数据来源或者承载的方式。



数据是未来最大的资产，用好数据不仅可以提高企业自己的产品和服务，也可以攫取大量利润。一旦没有守好数据，那么很有可能成为下一个负面信息的主角。

数据安全已经成为全世界瞩目的焦点问题，连续三年蝉联 RSA 热词榜冠军。

## 1. 数据安全风险

移动互联网的普及使得人们越来越多的在网络上留下信息，这些信息如果被分析和利用，将对个人隐私和安全形成极大的威胁，同时海量数据也增加了信息保护的难度。

历年来，数据泄露事件愈演愈烈，透过事件我们看到了问题的本质：

- 对数据的违规使用

非法收集：

外部：漏洞攻击、木马注入、弱配置、APT

内部：越权盗窃、离职

数据滥用：

诱导、贩卖、敲诈

- 数据泄露带来的风险：

访问：认证、权限

共享：业务（门户、调用测试）、人员交互

外发：跨区、第三方（网络、邮件）

外带：出差、回家

## 2. 政策上的措施

针对不断涌现的数据泄露问题，数据和隐私保护政策陆续出台：

- 我国于 2017 年 6 月 1 日正式施行《中华人民共和国网络安全法》，规定了公民使用网络服务需要实名认证，任何网络侵入、干扰和窃取网络数据都是违法的，个人信息安全得到真正的法律保护，从此确立了公民个人信息保护的基本法律制度，促进经济社会信息化健康发展。

- 我国《网络安全等级保护条例》提出对信息进行收集、存储、传输、交换、处理的系统进行不同等级的保护要求。对定级不准确不合理的网络运营者，应准确履行自己的网络安全义务，工作不到位的网络运营者主要负责人以及网络安全相关负责人将受到相应的处罚。

- 我国《关键信息基础设施安全保护条例》提出对保护范围内的单位运行、管理的网络设施和信息系统，一旦遭到破坏、丧失功能或者数据泄露，可能严重危害国家安全、国计民生、公共利益的，都应受到网络安全法的处罚，为企业带来了合规挑战。

- 我国于 2018 年 5 月 1 日正式实施《个人信息安全规范》，规范个人信息控制者在收集、保存、使用、

共享、转让、公开披露等信息处理环节中的相关行为，进一步强调了个人敏感信息被泄露、非法提供或滥用可能危害人身、财产安全，致使个人名誉、身心健康受到损害或歧视性待遇等严重后果，遏制个人信息非法收集、滥用、泄漏等乱象，最大程度地保障个人的合法权益和社会公共利益。

- 欧盟于 2018 年 5 月 25 日正式施行《通用数据保护条例》，简称 GDPR，被称为史上最严数据保护法，其最高惩罚代价为暂停使用个人数据。

### 3. 数据安全防护思路

#### 3.1 数据安全防护目标

绿科技为数据安全设计了全面可信的防御体系，有效保护数据在全生命周期过程中的安全，达到合法采集、合理利用、静态可知、动态可控的防护目标。

- 合法采集：利用大数据分拣技术，使企业在法律约束范围内合法采集敏感数据；
- 合理利用：通过建立数据模型，以及对数据的敏感级别进行划分，设立不同的访问层级，在数据被开发利用前做好防护措施，杜绝非法滥用；
- 静态可知：对存储中的静态数据进行扫描发

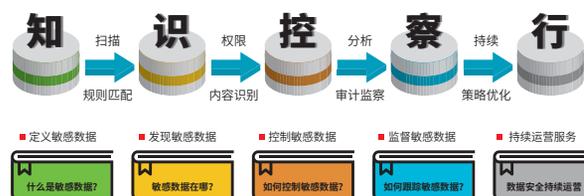
现，并展示数据的分布；

- 动态可控：对流动的数据进行监控，防止数据在交互、共享中有意无意的泄露。

#### 3.2 数据安全治理设计思路

数据安全就是对数据的安全治理，是从政策到数据的全生命周期的监察与保护。

绿盟科技结合客户的需求，以及对实际环境的调研了解，总结出了一套完整又科学的数据安全治理方法，及“知”、“识”、“控”、“察”、“行”。

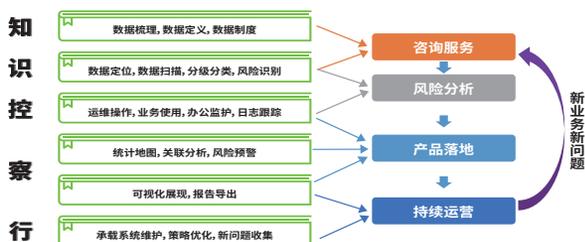


- 知：分析政策法规、梳理业务及人员对数据的使用规范，定义敏感数据；
- 识：根据定义好的敏感数据，利用工具对全网进行敏感数据扫描发现，对发现的数据进行数据定位、数据分类、数据分级。
- 控：根据敏感数据的级别，设定数据在全生

命周期中的可用范围，利用规范和工具对数据进行细粒度的权限管控。

- 察：对数据进行监督监察，保障数据在可控范围内正常使用的同时，也对非法的数据行为进行了记录，为事后取证留下了清晰准确的日志信息。
- 行：对不断变化的数据做持续性的跟踪，提供策略优化与持续运营的服务。

将数据安全治理方法“知”、“识”、“控”、“察”、“行”应用于实际项目中，利用咨询服务发现数据风险，通过产品落地实现对数据的可视化监控、风险点排除，及时预警、及时阻止对数据的非法使用行为，最后对数据进行持续运营服务，让数据始终处于被监控的安全状态，当有新的业务上线时，可根据此数据治理方法快速的实现新数据的安全监控。



#### 4. 数据安全解决方案

绿盟科技针对数据安全提出了完整的解决方案，

包括数据梳理、运维数据监管、业务数据监管、办公数据监管，以及数据的可视化，全面对数据在各种场景中的全生命周期安全进行了阐述。



#### 4.1 数据梳理与风险评估

- 形成全局数据分布情况，奠定分级分类管理基础
- 全方位了解数据安全现状与所面临的安全威胁
- 提供专业整改建议，促进系统整改
- 建立标准化、规范化、专业化数据安全管理体系
- 提升数据安全技术防护水平
- 针对性地、有序地进行各项日常数据安全工作和开展数据安全建设

#### 4.2 运维数据安全防护

企事业单位 IT 系统不断发展，网络规模迅速扩大，设备数量激增，建设重点逐步从网络平台建设，转向以深化应用、提升效益为特征的运行维护阶段，IT 系统运维与安全管理正逐渐走向融合。信息系统的运行直接关系企业效益，构建一个强健的 IT 运维安全管理体系对企业信息化的发展至关重要，对运维的安全性也提出了更高要求。



#### 日常工作监管

通过对运维数据安全防护可以对运维人员和合作伙伴的日常操作情况做保护和记录，方便监管。



#### 法律法规遵循

许多企业和单位需要满足国家或者行业监管部门的法律法规要求(如等级保护、企业内控制度等)。



#### 事后追溯

通过事后审计信息可以对发生的安全事件进行追溯，保留一份不能修改不可删除的“证据”。

堡垒机做为专业的运维监管系统，提供了先进的运维安全管控与审计能力，是运维数据防护的第一道防线，其目标是帮助企业转变传统 IT 安全运维被动响应的模式，建立面向用户的集中、主动的运维安全管控模式，降低人为安全风险，对运维数据的访问行为实现全面的审计，满足合规要求，保障企业效益。

数据库访问疏于监管带来了巨大危害，数据库审计

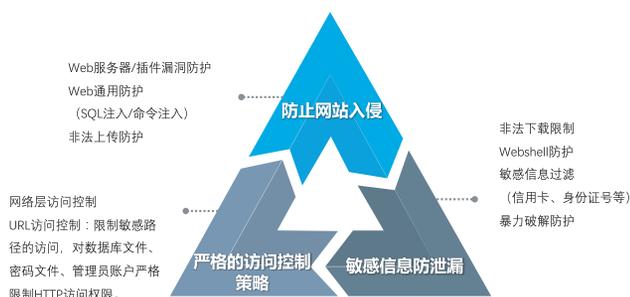
与防护产品受到空前关注。将数据库审计与数据库防火墙结合使用，就实现了具有集数据库 IPS、IDS 和审计功能为一体的综合安全防护能力。

大数据的思想及其初步应用已经惠及人们的日常生活，与大数据相互依存的云计算技术、物联网、智慧城市等新的应用模式同时印证了其在信息化时代的重要地位。随着数据的价值越来越重要，大数据的安全稳定也逐渐被重视，在大数据时代，无论对于数据本身的保护，还是对于由数据而演变的一些信息的安全，都对大数据环境提出了更高的要求。虽然大数据安全与大数据业务是相对应的，但对于业务和环境的维护将更为重要，大数据安全运维将包含风险检查、控制防护、审计监控三个方面。

### 4.3 业务数据安全防护

随着信息技术日新月异的发展，近些年来，企业利用计算机网络技术与各重要业务系统相结合，可以实现无纸办公。有效地提高了工作效率，如外部门户网站系统、内部网站系统、办公自动化系统等。然而信息化技术给我们带来便利的同时，各种网络与信息系统安全问题也逐渐暴露出来，业务数据泄露事件频发。因此要对业务系统进行安全防护，既要保障业务系统

中的静态数据、动态数据的一致性和业务的关联性，又能保证被输出的敏感数据得到脱敏处理。



的措施来消除这些威胁，降低整体安全风险，确保内部办公环境下的数据安全，具体防护方案可以归纳为以下几个方面：

- 数据外发途径全面防控
- 关键数据加密使用防护
- 数据离网交换安全防护
- 用户上网行为监管
- 网络数据外发监控
- 邮件数据安全防护

#### 4.4 办公数据安全防护

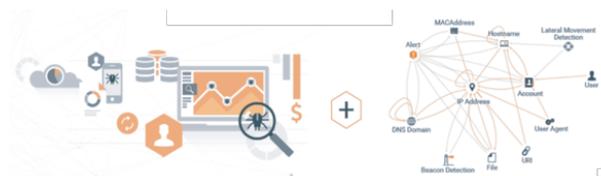
无纸化办公的普遍应用，办公环境中每位员工都拥有独立的办公终端，企业的研发源代码，财务、政券信息，运营资料，人资等敏感数据以不同的文件形式存在于每个人的终端电脑及存储服务器中，员工通过办公终端接入网络，连接到各种业务系统、各类服务器进行数据的使用与交换，很难对敏感数据进行准确的定位和分析发现。

数据在办公环境中的类型和存在的形态多种多样，数据操作和传输途径更是纷繁复杂。

当前无论是终端、网络还是管理方面，都对敏感数据的泄漏存在有很大的风险。因此，需要采取相应

#### 4.5 数据可视化展现

通过数据溯源和用户实体行为分析来完成可视化的展现，使用户可以更全面更具体的了解数据在全生命周期过程中的状态及风险，为后续防护措施提供有力支撑。



### 5. 方案价值

### 1) 满足合规要求

现如今，国家对数据安全已经出台了多项法规，通过本方案的实施，可以对法规中提到的鉴别信息数据、重要个人信息、重要业务数据做到针对性的监控与保护，使企业在发现数据风险前及时做出响应，避免因数据丢失造成的危害。

### 2) 权限划定清晰

责权不清一直都是最根本的问题，通过本方案的实施，将数据合理的进行级别划分，再结合管理与业务的需要对数据的访问、使用，进行清晰的权限管控，做到权责分离，事后还可以通过审计结果明确事故责任方，避免了责任不清出现的推诿扯皮。

### 3) 数据生命周期全面掌控

掌握数据的全生命周期是对数据风险的提前预知，利用本方案对数据的生命周期中各个环节做监控，掌握数据的动态，了解数据的流向，提前对可能发生的数据泄露风险进行预警，保障数据在安全的可控范围内流转、使用与存储。

### 4) 降低数据泄露风险

通过对数据的扫描与跟踪，利用内容识别、UEBA、机器学习等技术，及时发现数据所承载的系统、业务、网络、终端中的安全威胁，提前做好防范措施，

让泄密风险看得见、使数据泄漏防得住。

### 5) 提高数据使用者的安全意识

绿盟数据安全解决方案的应用，让数据使用者了解数据的重要程度，规范数据使用者的操作行为，从潜意识里指导与帮助人们正确使用资源，合理利用资源，保护数据的安全。

## 6. 总结

绿盟数据安全解决方案为客户提供了全面可信的数据风险识别与防护体系，将个人隐私数据、企业敏感数据、鉴别类信息进行有效的分拣区分，从数据治理到合规监管，从及时预警到风险态势，对不同场景提供有效的数据安全保障服务。

# 绿盟科技数据安全咨询服务介绍

安全服务部 贾晓萍 刘宇

关键词：数据安全、个人信息安全、咨询服务

摘要：绿盟科技根据国内外数据安全政策法规及相关标准要求，并借助已有的成熟技术评估服务流程体系及评估工具，为客户提供数据安全专项评估、数据安全治理、数据安全认证及数据安全整体防护方案等适用于多种数据安全场景的咨询服务。

随着社会的进步和科技的发展，信息已经成为我国实现经济转型升级的基础性资源，数据安全则是信息化持续推进的基本前提，但当前的数据保护情况不容乐观，数据泄露、数据滥用、个人信息交易等现象时有发生，数据安全问题日渐凸显，数据安全保护已经成为影响国家安全、社会秩序以及公民利益的焦点问题。

根据国内外数据安全政策法规及相关标准要求，并借助绿盟科技成熟的技术评估服务流程体系及评估工具，通过数据安全专项评估、数据安全治理、数据安全认证咨询及数据安全整体防护方案等多方面咨询服务，为客户提供适用于数据安全多种场景的咨询服务。



## 一、数据安全咨询服务介绍

### 1.1 数据安全专项评估服务

#### (1) 数据梳理咨询

绿盟科技数据梳理咨询包括数据分类分级和数据映射两部分服务，一方面，通过对数据的识别与收集，对数据进行分类与分级，为数据的分类分级的管控做铺垫，另一方面，对收集到的数据进行数据映射，协助

▶▶ 解决方案

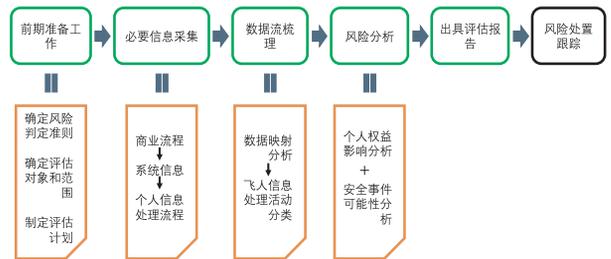
客户清晰了解数据的分布状态，为后期数据管控合规性评估做铺垫。



(2) 个人数据影响评估服务

个人信息安全影响评估服务参考《信息安全技术 个人信息安全影响评估指南（征求意见稿）》，通过数据映射分析对客户个人信息处理行为进行梳理与分类，并对客户个人信息处理行为中的个人信息安全风险从危害性和可能性两个维度进行评估，来实现对客户个

人信息安全影响的评估工作。根据评估结果，绿盟科技还将提供风险处置的建议，并在需要的情况下对风险处置的结果进行跟踪。



(3) 数据泄露安全评估

面对严峻的数据泄露形势，为积极防范客户由于外部及内部原因导致的数据泄露行为发生，结合其监管机构各项数据安全风险防范工作要求，绿盟科技通过提供基础安全评估工作、业务安全测试、移动 APP 测试及专项 API 测试工作，为客户提供数据泄露安全评估服务，进行数据防护手段建设。



#### (4) 数据生命周期安全评估

绿盟科技数据生命周期安全评估服务从数据生命周期通用安全和各阶段安全两个方面对数据生命周期进行安全评估，以发现客户数据生命周期安全管理及管控措施方面存在的安全问题，提出整改意见。

数据生命周期各阶段安全要求					
数据采集	数据传输	数据存储	数据处理	数据交换	数据销毁
<ul style="list-style-type: none"> <li>· 数据收集和获取情况</li> <li>· 数据分类情况</li> <li>· 数据分级情况</li> </ul>	<ul style="list-style-type: none"> <li>· 传输保密性控制措施</li> <li>· 传输完整性控制措施</li> <li>· 网络边界安全</li> <li>· 网络可用性管理</li> </ul>	<ul style="list-style-type: none"> <li>· 数据存储加密</li> <li>· 数据备份和恢复</li> <li>· 数据库安全防护</li> <li>· 数据访问控制</li> <li>· 数据冗余与时效性</li> </ul>	<ul style="list-style-type: none"> <li>· 数据正当使用</li> <li>· 数据脱敏</li> <li>· 数据权限管理</li> </ul>	<ul style="list-style-type: none"> <li>· 数据交换的必要性</li> <li>· 数据交换的合法性</li> </ul>	<ul style="list-style-type: none"> <li>· 介质使用管理</li> <li>· 数据销毁处置</li> <li>· 介质销毁处置</li> </ul>
数据生命周期通用安全要求					
<ul style="list-style-type: none"> <li>· 数据安全策略与制度</li> <li>· 组织及人员管理</li> <li>· 合规性管理</li> </ul>	<ul style="list-style-type: none"> <li>· 系统建设</li> <li>· 系统运维</li> <li>· 网络和数据安全</li> </ul>	<ul style="list-style-type: none"> <li>· 第三方服务</li> <li>· 基础环境</li> <li>· 设备与计算安全</li> </ul>	<ul style="list-style-type: none"> <li>· 数据安全事件应急处置</li> <li>· 安全审计</li> </ul>		

## 1.2 数据安全治理咨询服务

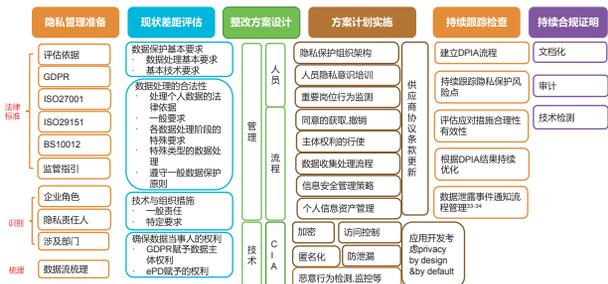
### (1) 通用数据安全治理咨询

绿盟科技通用数据安全治理咨询服务是在对客户数据进行有效理解和分析下，对数据进行不同类别和密级的分类分级工作及数据映射梳理；在对数据分级分类的基础上，了解这些数据在生命周期中的安全管控情况，并辅助以渗透测试、漏洞扫描、配置核查等技术检测工作；最后针对不同数据的安全需求，在满足数据正常使用的目标下，完成相应安全方案的设计、实施及优化服务。



### (2) 隐私数据治理 (GDPR) 安全咨询

隐私数据治理 (GDPR) 咨询服务，通过详细解读 GDPR 的相关法规要求，结合国际成熟的信息安全管理体系框架 (ISO27001, ISO27018, BS 10012:2017) 等，以及先进的技术评估工具，从人员、流程、技术和隐私治理多个方面，帮助客户快速发现不合规的领域并提供应对措施，以督促企业尽快符合 GDPR 的要求，并通过实施 GDPR 合规项目，协助客户逐渐形成成熟的隐私数据治理与保护体系。



## ▶▶ 解决方案

### 1.3 数据安全领域认证咨询

近年来，随着对数据安全保护的重视度的提升，数据安全领域的认证也引起了很多企业的关注，绿盟科技目前可提供 ISO 29151 认证咨询服务，主要参考《ISO/IEC 29151:2017 信息技术 - 安全技术 - 个人信息保护实践规则》、《ISO/IEC 29100:2011 信息技术安全技术 隐私框架》要求，对客户相关认证咨询服务，配合客户规范个人信息收集、存储、处理、使用和披露等各个环节中数据操作的相关行为，提高业务流程的安全性和可靠性，降低 IT 运营过程中的个人可识别身份信息风险，并协助客户获得 ISO 29151 认证证书。



### 1.4 数据安全整体防护方案

绿盟科技数据安全整体防护方案从组织建设、人员能力、制度流程、技术工具四方面进行阐述，为客户提供数据安全整体防护方案，协助客户全方位提升

数据安全管理与控制水平。



### 二、数据安全咨询服务客户收益

绿盟科技通过提供数据安全咨询服务，协助客户达到保护数据资产、管理敏感信息、风险规避以及满足政策合规的目的。



# 大数据安全的解决思路

BSG产品管理部 孙叶

关键词：大数据平台、数据安全

摘要：本文从大数据平台下的安全问题，引出大数据安全的法规标准、防护思路及解决方案。

## 引言

随着互联网、物联网、云计算等技术的快速发展，全球数据量出现爆炸式增长；根据 IDC 研究的“大数据摩尔定律”表明，人类社会产生的数据一直在以每年 50% 的速度增长，也就是说，每两年就增加一倍。在大数据不断向各个行业渗透、深刻影响国家的政治、经济、民生和国防的同时，其安全问题也将对个人隐私、社会稳定和国家安全带来巨大的潜在威胁与挑战。

政务信息化的推进，电信、金融、互联网等行业的平台升级，加速推进大数据安全和隐私保护需求。为了应对这些需求，我国正在开展的全国网络安全执行大检查行动中，首次开展针对大数据安全的整治工作，具体包括大数据的采集、传输、存储、处理、交换、

销毁等全生命周期的监控与保护。

## 一、大数据平台下的安全问题

大数据平台涉及到的内容比较广泛，安全问题可以从这 5 个维度去考虑，安全管理、平台安全、数据安全、运维安全、业务安全。



### 1、安全管理

安全管理是指大数据平台安全管理方面的要求，包括管理制度、机构和人员管理、系统建设管理、运维管理等内容及配套管理流程。安全防护离不开管理与技术协同，国家、政府、行业自上而下应该有安全管理制度和管理流程，指导具体安全工作的开展和实施。

### 2、平台安全

平台安全指平台主机、系统、组件自身的安全和身份鉴别、访问控制、接口安全、多租户管理等安全问题，是对大数据平台传输、存储、运算等资源的安全防护要求。企业大多数都使用基于社区化、开源化组件的Hadoop平台，缺乏安全方面的考虑。

### 3、数据安全

数据属于一种资产，有6个生命周期阶段：采集、传输、存储、处理、交换、销毁；数据安全要保障数据在任何阶段下都是安全的。围绕数据全生命周期考

虑数据安全问题，例如：数据采集阶段的分类分级、清洗比对、质量监控；数据传输阶段的安全管理；数据存储阶段的安全存储、访问控制、数据副本、数据归档、数据时效性；数据处理和交换阶段的分布式处理安全、数据加密、数据脱敏、数据溯源；数据交换阶段的数据导入导出、共享、发布、交换监控；数据销毁阶段的介质使用管理、数据销毁、介质销毁等安全问题。

### 4、运维安全

运维人员的权限相对较大，运维人员直接对数据库进行操作，涉及的数据量非常大，数据的安全难以保障。例如：内部人员的误操作导致数据丢失或不可用，蓄谋恶意行为导致数据泄露。

### 5、业务安全

业务安全跟业务强相关，跟应用场景和业务流量特征有关，一般的防护手段很难发现，涉及到业务学习和行为分析。例如：缓慢少量攻击、共谋、在噪音

中隐身、持续渗漏尝试、长期潜伏者等。

## 二、大数据安全法规标准

大数据时代是万物互联的时代，数据在共享中体现价值，因此，国内外法律法规也终将完善大数据安全领域的防护和技术要求，助力大数据安全建设。

国家、政府、各行业相继出台大数据平台安全和数据安全相关的国标、行标、企标、地标，推动大数据产业的良性发展。《中华人民共和国网络安全法》、《中华人民共和国计算机信息系统安全保护条例》、《等保2.0》、《个人信息安全规范》、《GDPR》、《电信网与互联网大数据平台安全防护技术要求》等。

## 三、大数据安全防护思路与解决方案

数据共享是必然需求，大数据安全的防护目标要在保障业务正常的前提下，以合理成本，保护大数据平台下数据的安全。业务需求与风险并存，防护要在业务需求与风险之间寻求平衡，对不同价值和属性的数据，在不同业务需求下，实施不同级别的防护措施，控制防护成本。

### 1、防护思路

大数据安全防护方案可按层次考虑，平台安全、数据安全、运维安全、业务安全，层层深入，逐步提升安全性。

- **平台安全**

数据的存储和流转依托大数据平台和各业务系统，平台自身安全是第一步，通过平台各组件与系统的漏洞扫描管理、规范化的基线核查管理、平台态势感知，确保大数据平台的安全运行。

- **数据安全**

关注数据的安全存储，数据梳理，掌握数据全景图，让数据风险可量化；关注数据在处理、交换、使用时安全，身份认证、访问控制、数据加密、数据脱敏，防止非法或越权访问数据，对数据访问进行管控、数据审计。

- **运维安全**

收敛大数据平台的数据访问途径，对运维人员访

## ► 解决方案

问大数据平台的操作行为进行操作管控、操作审计。

### ▪ 业务安全

机器学习建模，对敏感数据的访问行为和敏感业务进行机器学习，对用户行为进行分析，感知和预测业务安全风险。

## 2、国外厂商方案

Gartner 在《Market Trends: Database Security, Worldwide, 2017》报告中列出几个 Big Data 厂商，我们挑选 Informatica 和 Dataguise 这 2 个厂商的方案进行简单介绍。

### ▪ Informatica

该厂商的大数据安全解决方案支持对结构化、非结构化数据做发现、分类、风险评分，数据访问和操作监控，发现可疑或未授权操作，数据保护，敏感数据扩散跟踪，让风险跟踪处置形成闭环。

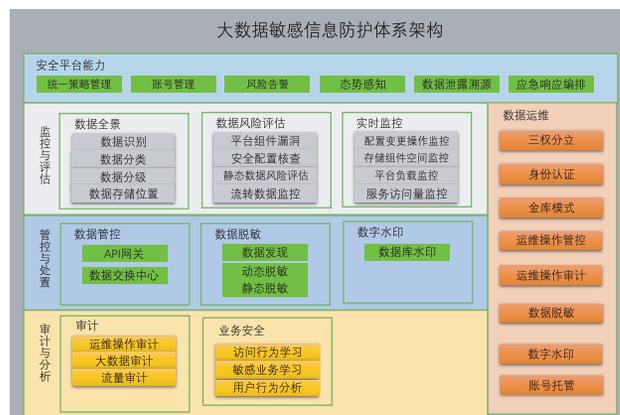
### ▪ Dataguise

该厂商提供全局敏感数据管理解决方案，产品旁

路部署在大数据集群边界，可实时检测、审计、保护和监视敏感数据资产。

### ▪ 绿盟大数据安全防护方案

绿盟结合大数据安全防护思路和实操性，从安全平台能力、运维安全管控、监控与评估、管控与处理、审计与分析这几方面去考虑，设计出一套大数据安全防护方案，覆盖平台安全、数据安全、运维安全、业务安全，提供数据发现、分类、分级、评估、监控、保护、审计、溯源、态势感知一整套大数据安全防护方案。



安全平台对接设备，集中管理、日志收集、智能分析，

拥有统一策略管理、账号管理、风险告警、态势感知、数据泄露溯源、应急响应编排能力。

监控与评估，从监控的角度入手，首先对大数据平台上的数据做梳理，数据识别、分类、分级、存储位置定位，生成数据全景图，为数据细粒度的访问授权提供依据，同时对动态数据做跟踪，监控数据流转、使用是否符合预期；从评估的角度出发，结合平台自身的组件漏洞、配置安全性及数据在平台上的存储、流转情况做综合评估，全面剖析数据风险，量化数据风险，为大数据平台的态势监测与防护提供有力支撑。

管控与处置，结合客户行业数据特征，提供行业数据分类分级模板；结合客户实际业务需求，灵活提供细粒度访问控制、数据加解密和数据脱敏方案，根据实时应用和业务流量监控，及时处置异常行为，避免进一步风险。

审计与分析，审计指对行为操作和数据访问做审计，为事件问题定位、溯源和大数据分析提供依据；业务安全指基于用户画像和异常行为分析做业务安全风险的感知和预测，及时给出处置策略。

运维安全管控，对平台上所有运维操作进行统一

管控和审计，利用运维操作审核机制，防止内部人员私自、独立对平台配置和数据进行操作。（等保 2.0 要求：大数据平台的管理流量和系统业务流量分离，因此，运维流量和业务流量也要分开做管控）

#### 四、大数据安全面临的挑战

大数据安全与传统数据安全相比，存在一些差异，大数据环境的特点是分布式、组件多、接口多、类型多、数据量大，这些特性给大数据安全引入了技术难点。



主流开源大数据组件二十多款，还有大量第三方封

装的组件，不同组件使用的交互接口不同，安全产品面对这么多组件接口，在监控、防护、溯源的方案设计和技术实现上都有难度。

大数据平台要存储和处理的数据量庞大，IDC 预计，到 2020 年全球数据总量将超过 40ZB，面对持续膨胀的数据量，安全产品不仅要提高单机产品的处理性能，还要考虑产品扩容和延展性。

大数据平台要存储和处理的数据类型众多，结构化数据、半结构化数据、非结构化数据。要对非结构化数据做识别、分类分级和脱敏处理，有一定技术难度。

## 五、大数据安全未来发展方向

由于政务大数据覆盖了自然人、法人、企业、政府机构等，同时和医疗、教育、民生服务等各个部门相关；因此，解决了政务大数据安全问题，就能有效解决其他行业大数据安全问题，有力支撑国家治理体系和治理能力现代化目标的实现。从企业层面来看，国家将统一标准规范，避免行业交流繁杂、数据所有权混乱、开发成本高等一系列问题。统一的数据管理平台，统

一的数据存储，统一的数据标准，进行统一的数据资产管理，统一进行授权管理，这是未来探索的一个方向。

# 金融行业数据治理方案

亿赛通 安全服务总监 李迪

关键字：数据安全、数据治理、数据资产、分级分类

摘要：亿赛通根据十六年数据安全实践经验以及对金融行业的深刻理解，提出金融行业数据治理方案，以数据安全防护为核心，围绕数据生命周期，从组织建设、制度流程、技术工具和人员能力等四个方面进行建设，在实施层面按照需求梳理、分级分类、策略制定、技术落地、优化改进五个方面进行数据治理工作的落地开展。

## 一. 金融行业数据安全背景分析

金融行业作为国家的经济重要领域，数据资产庞大，使用角色繁杂，数据共享和分析的需求强烈，但目前金融机构在数据管理方面仍存在较多问题，数据处理过程中大量的用户信息及用户业务使用信息等个人隐私数据管控机制不足，面临违规越权使用或被用于非法用途等数据泄漏安全风险，对员工有意或无意的敏感数据泄漏缺乏检测与防护手段。

近几年主管部门相继出台了多项规章制度，《网络安全法》于2017年6月1日起施行，银保监会于18年5月发布《银行业金融机构数据治理指引》，证监会于2018年9月发布《证券期货业数据分类分级指引》，从监管层面不断完善数据安全工作，指导金融机构加

强数据治理，提高数据质量，以数据驱动金融机构发展。

## 二. 金融行业数据治理方法论

亿赛通数据治理专业服务以数据安全防护为核心，以合规与业务需求为导向，围绕数据生命周期，从组织建设、制度流程、技术工具和人员能力等四个方面进行能力建设。

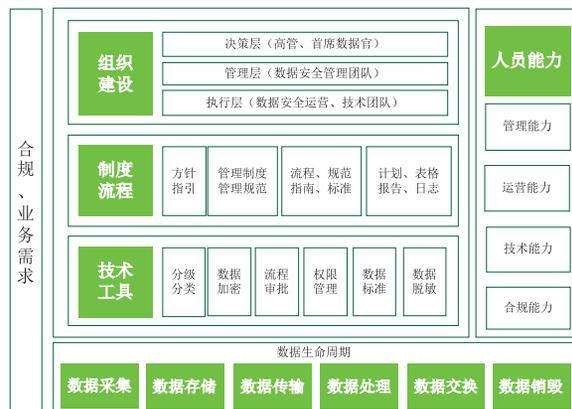


图 2.1 数据治理体系架构

## 2.1 合规和业务需求

在合规监管制度层面，银证监已经下发了《银行业金融机构数据治理指引》，从业务需求方面，数据向第三方平台共享的安全管控需求，以及数据大集中资产梳理的业务需求驱使推动数据治理建设开展。

## 2.2 组织架构

传统金融行业安全均由科技部门负责，随着数据治理工作的深入开展，业务部门要深入参与数据资产梳理以及分级分类工作，因此原有的组织架构和项目模式无法支撑数据治理的深入开展，需要自上而下形

成高层牵头、跨业务部门、数据全覆盖的组织架构。

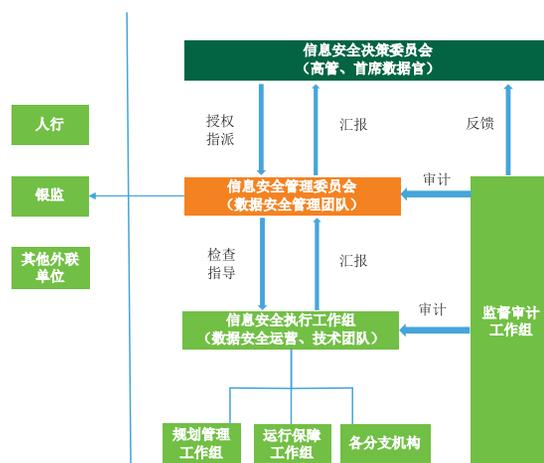


图 2.2 数据治理组织架构

## 2.3 制度流程

目前金融行业大多有较完整的安全规范，如分级分类规定，保密规定等，但一方面没有独立的数据安全规范，可执行性不强，另一方面缺乏技术监管手段，落地执行较难。因此需要制定独立的数据安全管理文件，按不同级别分期建设，逐步落地。



图 2.3 数据安全管理制度体系

## 2.4 技术工具

数据安全项目真正地执行不仅需要管理制度的规范，更需要技术工具的管控。根据数据分级分类进行安全环境保障、边界管控及合规监管，保障数据的保密性、完整性和可用性。

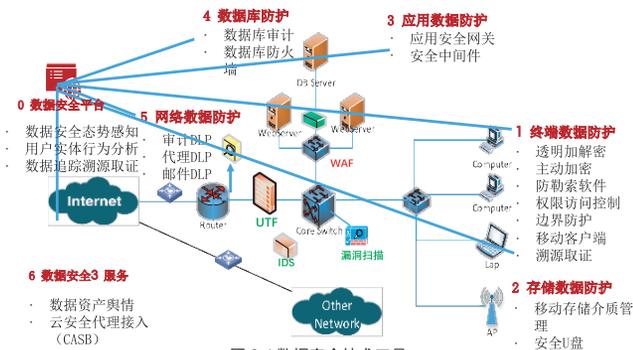


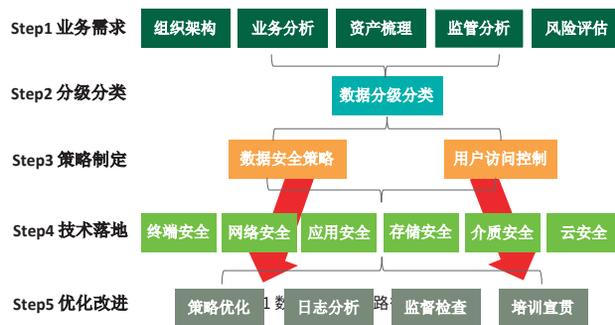
图 2.4 数据安全技术工具

## 2.5 人员能力

传统安全人员的技术能力大多以网络安全和信息安全为基础，而在数据安全层面需要既懂金融业务，又懂数据安全体系的复合型人才，对数据治理人员的培养和管理制度的宣贯需形成常态化机制，提高数据安全人员能力。

## 三 . 金融行业数据安全治理实施路径

亿赛通数据安全治理服务按照业务需求、分级分类、策略制定、技术落地、优化改进五个步骤进行数据安全实施工作。



### 3.1 确定业务需求

开展数据治理首先要在确定组织架构的基础上，通过资产扫描工具结合人工业务调研，对企业数据资产进行全覆盖梳理。同时进行全生命周期风险评估和监管政策对标分析，以此确定业务需求和目标。

### 3.2 数据分级分类

分级分类是数据治理的前提，也是工作量最繁重的环节，在数据资产全覆盖梳理的基础上，首先对条线进行业务细分，确定管理主体和数据治理覆盖范围。其次根据业务调研结果，进行数据资产的数据归类，确定数据类别，完成数据分类工作，然后按照数据损坏丢失可能造成的影响程度进行数据定级，基于数据分级分类对不同级别的数据实行差异化安全控制手段。

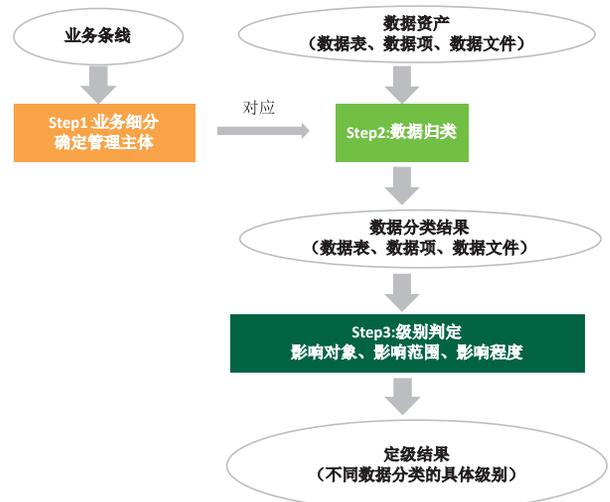


图 3.2 分级分类流程

### 3.3 策略制定

在数据治理的策略制定上，从用户权限控制和数据安全策略两个方向考虑如何实施数据安全治理，在针对“人”的权限控制上，要明确数据的访问者、访问对象、访问行为，尤其在对第三方开放数据共享时，要严格控制系统开放的权限；在针对“数据”的安全防护上，要基于不同级别的数据制定有针对性的数据安全策略，对核心商密和普通商密数据进行加密管控，对内部公开数据进行安全审计管控，形成整体化全生

命周期的数据策略体系。

### 3.4 技术落地

亿赛通针对数据安全治理提供全生命周期安全管控平台，可以对客户的主机和服务器数据进行全面扫描和梳理，智能识别用户数据安全资产。根据分级分类的落地应用进行标密定密，对于核心数据资产进行加密防护。通过在终端、网络、应用、介质等数据传输通道的全面监控，建立数据安全边界，将技术工具与管理有效结合，实现对用户核心信息资产的全方位保护。

工具进行事后审计，对企业内核心数据资产进行分布统计、合规检查，并对违规行为进行趋势分析和合规预警。打造“事前防御、事中控制、事后审计”的全方位防护体系。

### 3.5 优化改进

数据安全建设需要长期不断进行管理和技术的调整优化。数据安全平台管控策略要随管理制度细化不断完善优化，也要进行各部门的自查和监督部门的结果性审计检查，同时要做好公司内的培训宣贯，规避常见办公安全风险，贯彻数据安全管理规定。

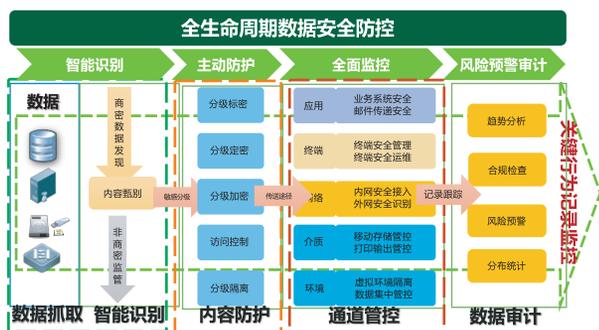


图 3.3 全生命周期数据安全防控技术手段

在数据安全防控的基础上，可通过数据安全稽核

## 四 . 安全治理案例

### 4.1 Q 银行

亿赛通协助 Q 银行开展业务梳理、资产梳理和分级分类工作，在管理方面制定了数据资产安全管理规章制度，在技术方面建设数据安全治理平台，对办公网内的非结构化电子文档数据共享和移动存储进行加密防护、权限管理和外发管控。

## ▶▶ 解决方案



图 4.1 Q 银行数据安全治理平台

在行内开展安全运营，对内网终端进行全方位审计，并运用大数据分析技术形成信息安全态势分析，提供信息安全事件追踪溯源，统计分析，管理控制等辅助决策技术手段。

Q 银行科技处与亿赛通共同完成的《城商行数据安全治理平台的研究及实践》课题获得 2018 年度银行业信息科技风险管理课题研究成果二类奖项。

### 4.2 Z 银行

亿赛通协助 Z 银行对总行及全国各分支行进行数据治理工作，进行关键数据梳理、风险评估、制定分级分类标准及相关管理规定，从管理和技术上实现对重要数据资产的分级防护，打造全行数据安全防护体系。



图 4.2 Z 银行数据治理方法

在项目中利用态势感知技术实现可视化管理、风险预警和溯源取证。实现“授权用、带不走、无法读、留痕迹”的总体目标。

## 五 . 结语

数据治理是随着金融行业资产数据化和数据资产化，从大数据时代到云时代发展转变的产物，企业数据治理的最终落脚点在于数据应用和价值实现，从数据管控向数据价值转变，实现数据驱动业务发展。目前国内的数据治理方兴未艾，仅在法律法规层面有了方向性的指引，尚缺乏可执行的监管标准和业界最佳实践。亿赛通将不断从理论和实践层面完善数据治理水平，打造成成熟的数据治理业界最佳实践，为金融客户数据安全保驾护航。

# 基于大数据分析的敏感数据检测及响应方案

ESM技术部 梁莎 李景 皮靖  
TRG产品部 吴天昊

关键词：敏感数据泄露、数据安全、数据防护方案、敏感数据检测及响应

摘要：数据化新时代悄然而至，新型的数据风险也随之而来。笔者所在的 ESM 技术团队通过分析敏感数据泄露的成因，推出基于大数据分析的敏感数据检测及响应方案。这是业界第一个主攻“过失泄露”方向的敏感数据防护方案。该方案整合了机器学习等技术，通过采集多维度的数据，对服务、主机、账号持续画像，监控与分析，在数据泄露发生之前，“早观察、早发现、早知晓”，化被动为主动，有效保证了企业敏感数据的安全。

随着互联网技术的应用越来越成熟，众多企业组织已经步入了数据化新时代。信息化建设有利于企业的健康成长，但同时也存在较高的数据风险，特别是敏感数据（即企业运营数据、客户信息、个人行为隐私等重要信息）泄露的风险。近年频频发生的用户敏感数据泄露事件给企业与组织造成了巨大的损失。2018年3月17日有媒体报道，一家名为“剑桥分析”的英国公司，在未经用户许可的情况下获取 Facebook 上 5000 万名用户个人信息数据。受此影响，Facebook 股价大跌，两日市值蒸发 500 亿美元。这个事件也影响了 Twitter、Snapchat 等社交媒体公司的股价，导致其大跌。加强敏感数据的防护已成为企业数据安全管理工作中的重中之重。

## 一、敏感数据泄露成因

《2017 年数据泄露 QuickView 报告》显示，黑客攻击仍然是数据泄露的主因之一，但其对数据泄露的影响下降到第二位。自 2008 年以来，无意的数据泄露和其他数据处理错误比恶意入侵网络导致更多的数据丢失。根据 Verizon《2018 数据泄漏报告调查》统计，2017 年关于企业信息安全的大小事件共发生了大约 53,000 起，其中确定为数据泄漏的事件有 2,216 起。据统计，68% 的数据泄露事件是“过失泄露”，即企业内部人员因为过失导致公司数据泄露；22% 的数据泄露是“恶意违规”，即攻击者或内鬼的恶意行为导致数据泄露；10% 的数据泄露是“账号攻陷”，即存在漏洞被攻击者攻陷。

## 二、方案介绍

针对这一情况，笔者所在的ESM技术团队推出面向“过失泄露”的解决方案，这是业界第一个主攻“过失泄露”方向的敏感数据防护方案。本解决方案通过大数据和机器学习的手段，以基线、画像与持续监控的方式，学习账号与系统的正常行为，及时发现系统内的异常行为，对偏离正常行为的动作进行及时告警，因此可以有效分析出安全问题的源头，助力企业与组织解决业务数据安全问题。本项目已帮助数家大企业客户做到了智能化的分析，保障了敏感数据的安全，带来了良好的市场收益。

本解决方案通过多维度数据采集，梳理敏感数据异常场景，如图1所示，初步实现了敏感数据防泄漏的模式。

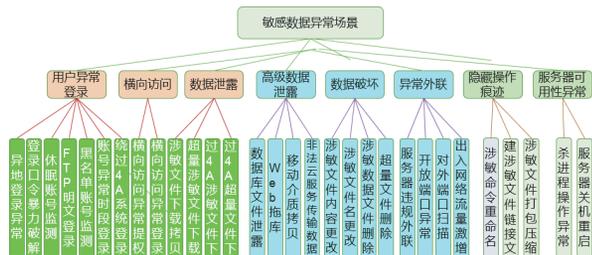


图1 敏感数据异常场景

解决方案的实现细节，如图2：



图2 敏感数据防护实现

### 2.1 多维度采集数据

本解决方案实现了数据的多维度采集，如网络流量数据、业务系统数据、设备告警数据、其他终端等数据的采集，提取了三个层面的敏感信息：

1. 账号：为习得账号的正常行为基线，发现其异常行为，采集如4A系统账号，OA账号以及其他业务系统账号的活动数据。
2. 主机：为了学习主机的正常模式，发现其异常行为，采集如主机终端的进程活动日志、网络活动日志、文件操作日志等。
3. 服务：为了获取服务的正常行为基线，发现异

常,采集如服务系统(如FTP服务、数据库服务、客户服务等)日志、服务访问日志、服务所在网络的网络全流量日志等。

通过多维度数据采集,对用户异常行为实现有效监控,基于用户异常行为分析基线,针对涉及敏感数据的业务,能有效的检测到用户、主机、服务以及相关数据的异常行为。

## 2.2 对服务、主机、账号等持续画像与监控

同时,实现了对服务、主机、账号等的多维度持续画像,监控与分析,如图3所示:



图3 敏感数据多维度监控

有了合适的数据源,接下来使用基线进一步对账号、主机、服务等进行刻画与描述。

1. 数值型基线:对于历史流入流出流量,历史访问端口行为,历史访问主机行为等可以量化的指标,使用数值型基线。

2. 标称型基线:对于账号常用登录区域,账号惯常登录时间,账号常用IP,主机历史开放端口等不能量化的指标,使用标称型基线。

3. 图基线:特别的,对于账号访问路径这种可以用有向图表示的情况,使用图基线。

对于黑名单账号活动、休眠账号活动、敏感文件访问、内网主机外联动作等进行持续监控;对时序性发生的反常行为构成故事线,进行时序性分析;对个体风险评估在特定时间间隔内陡增,进行风险激增分析;对主机与主机包含服务的风险关系,进行关联分析。

通过以上手段,能够有效针对移动重点业务服务器(DPI)进行流量画像,对访问信息、活跃时间段、存活情况、活跃服务数、应用/协议类型、流量特征等进行分析,发现多起服务器非法外联,异常时间异常操作,服务异常等。

### 三、方案的优势

#### 3.1 业界第一个主攻“过失泄露”方向的敏感数据防护方案

传统的敏感数据防护方案大部分是防范恶意攻击者造成的数据泄露，对数据泄露事件的认识与理解还不够充分。针对这一情况，本项目推出面向“过失泄露”的解决方案，研究成果可以迅速的在绿盟全流量系统中落地实现，有助于提升绿盟在数据安全与敏感数据防护领域的竞争力。

#### 3.2 在数据泄露发生之前，“早观察、早发现、早知晓”，化被动为主动

传统的数据防护方案重在加强资产管理，外设管理等管理方面，一般要等到数据泄露发生了才能发出有效的告警。本项目推出的解决方案研究发现，数据泄露通常伴随着大量的异常动作，异常用户行为（如异常时间段登录、异地登录、流出流量激增等）。因此，本方案从服务、主机和账号等多个维度，学习历史行为模式，感知正常的模式，得到正常行为基线，利用算法生成行为基线模型，当实时的动作行为发生改变或偏离基线时进行跟踪分析，及早告警。

#### 3.3 整合机器学习等多种技术

为了发现可疑的异常行为，除了依靠传统的统计及特征的方法外，本方案还利用大数据，机器学习等先进技术和算法，依靠建立的行为模型，维护正常的访问行为模式，形成用户的行为基线，从而发现偏离历史惯常行为的动作，及时告警。此外，为验证用户身份及检测异常行为，本方案发明了一种基于用户在某站点系统历史访问轨迹建立用户行为图谱，并利用行为图谱持续验证用户身份以及判断是否存在异常访问行为的方法。该方法已经在国内申请专利。

### 四、案例说明

#### 4.1 时间线故事场景示例 -- 数据窃取

下面通过一个时间线故事来描述基于大数据分析的敏感数据检测及响应方案如何帮助客户发现并制止高危事故，如数据窃取等。

某公司员工小王：

1. 时刻一，收到了来自陌生人的邮件，里面包含了疑似恶意附件，这些异常行为触发了相应的告警，风险积分也随之增加；

2. 时刻二，小王的办公 PC 主动对境外一个陌生 IP 发起了 TCP 连接，这个异常行为触发了主动外联境外陌生主机的告警，风险评分继续增加；

3. 时刻三，系统发现小王的办公 PC 利用小王的账号频繁尝试访问用户系统的数据库，并登录成功，系统对这个异常行为给出了告警，并增加了风险积分；

4. 时刻四，系统发觉小王的办公 PC 利用小王的账号在用户系统数据库执行 dump table 操作，马上触发了相关的告警，风险评分增加；

一直到时刻五，系统发现小王的办公 PC 有对外网的大量 DNS 隐蔽信道通信，触发了相关告警，风险评分继续增加，此时，小王累计的风险评分已经达到了 190 分，风险排名位居第一位，引起了安全人员的重视，在事态恶化之前，得到了及时处理。

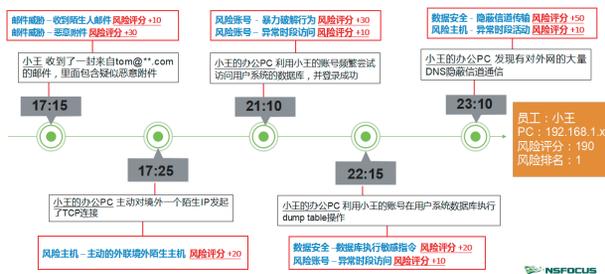


图 4 时间线故事场景示例 -- 数据窃取

基于大数据分析的敏感数据检测及响应方案通过对账号、主机、服务等进行多维度持续的画像与监控，可以在第一时刻发现异常行为，并给出相关告警，和相关的风险评分。当累计的风险积分不断增加时，一方面，这说明有风险的行为越来越多，系统出现安全事故的可能性在增加；另一方面，这个用户引起安全人员重视的可能性也在增加。这样，安全人员能及时发现可疑用户，并进行处理，从而保障了系统的安全。

#### 4.2 真实案例 -- 内部人员账号失窃导致数据外泄

将基于大数据分析的敏感数据检测及响应方案部署到客户的系统中，帮助客户发现并阻止了很多高危事故。在某个重点客户局点：

首先，系统发现了账号异常行为，在某一时刻，账号在不常用的 IP 登入，触发了账号异常的告警；

紧接着，系统发现该账号访问了不常访问的路径，触发了账号异常的告警；

最后，系统发现实时传输的数据量超过历史基线的 5 倍以上，触发了数据传输异常的告警。

至此，敏感数据防护系统的账号失窃导致数据外泄的 playbook 被触发，系统判定为高危事故，主动

向绿盟态势感知平台发出封堵消息，一键封堵子系统作出封堵动作，及时止损，避免了事态的进一步恶化。



图 5 真实案例 -- 发现数据外泄行为

基于大数据分析的敏感数据检测及响应方案利用机器学习的方法进行基线检测，在训练过程中先描摹出账号的正常行为，从而实时鉴别出账号的异常行为，做出相应的告警。将这些告警串接起来，可以为鉴别高危事故提供相应的线索与证据，形成一条完整的故事线。

## 五、结束语

数据安全领域正在快速发展，技术供应商会越来越多。数据的完整性与安全性得到越来越高的重视，致力于用户的数据安全诉求，更智能化的提供敏感数据防护的方案，会更有助于产品的定位与市场策略的制定。ESM 技术团队的基于大数据分析的敏感数据检测及响应方案将数据安全转化成了业务安全，发挥了“看清已知，检测未知”的重要作用，能为企业业务数据的安全防护作出更大的贡献。

# 特权访问下数据安全防护方案

BSG产品管理部 许德昭

关键词：特权访问管理、数据安全、最小权限原则、三权分

摘要：据预测未来 30 年，数据将成为生产资料，将会是企业核心重要资产，如何对具有系统特权用户访问业务数据进行权限管控和安全审计成为企业管理中的关键问题，本文通过对特权访问特点和潜在安全风险的分析，提出一种特权访问下数据安全防护方案。

## 一、特权访问与安全风险

### 1.1、特权访问

特权访问是指具有业务系统特权的人对业务系统进行配置更改或者对业务数据访问操作等行为，此类访问行为可控制组织资源、修改安全策略以及访问大量敏感数据，常常和系统运行维护操作相关，系统运行维护访问操作是最为典型的特权访问行为。例如在 Linux 系统上以 root 权限登录系统修改系统参数、停止核心业务服务或执行系统关机操作；亦或在交换机上以管理员权限登录修改交换机路由配置添加访问白名单放行外部用户访问内部敏感系统和数据等。特权访问具有隐秘性强、可执行权限高和影响范围广的特点。

- 隐秘性强

业务系统为了权限分配管理和系统运行维护需要，系统中必定存在特权帐号。业务系统上线后由于企业组织主要精力在业务运行上，往往忽略了对特权帐号的管控，特别是对第三方人员对业务系统的运维操作，甚至于客户还不知晓一些特权帐号的存在。

同时通过网络远程访问的特权操作大量存在，由于网络通信的不安全性，大量远程特权访问行为通信数据都被加密后进行传输，因此特权访问行为更难以详细展示出来。

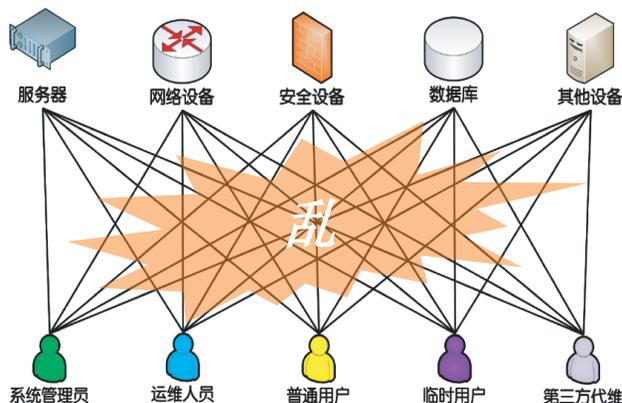
- 可执行权限高

特权帐号相比普通用户拥有更多的数据读写权限，以及系统操作执行权限。特权访问时可访问业务系统中的核心关键数据信息，可对其进行编辑修改，更新

关键业务流程步骤等操作。

▪ 影响范围广

当特权用户使用特权帐号访问业务系统后，由于具备高执行权限，可修改关键业务数据和配置，一旦修改成功后会业务和系统都具有广泛影响。例如修改企业边界网络防火墙访问白名单时，可设置所有人均可访问网络资源；自动更新企业核心数据库系统组件等，此类操作都可能导致企业业务出现中断或者崩溃。



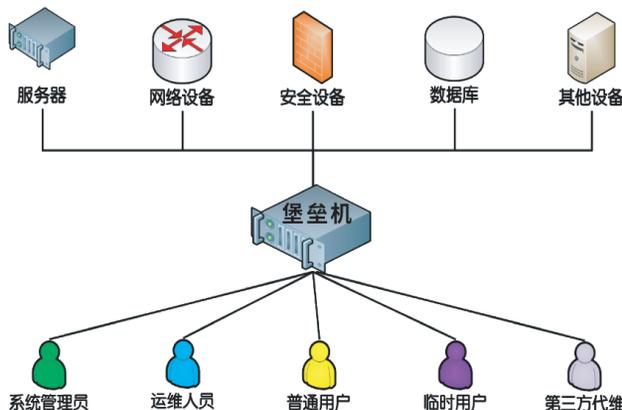
二、特权访问下数据安全解决方案

将分散、混乱特权访问现状进行集中统一管控是有效的解决办法，实现特权身份和访问权限进行集中管理，并对访问行为进行全程实时记录，为事后安全审计提供有力证据。

1.2、安全风险

“未来 30 年，数据将成为生产资料，计算会是生产力，互联网将成为一种生产关系。”（引用自第四届世界互联网大会上马云在开幕式致辞内容），数据将会是企业核心重要资产，这些重要资产在访问控制上往往缺乏有效的控制管理，常常存在如下的安全风险。

- 特权身份冒用、滥用
- 访问权限管理混乱
- 缺乏有效的安全审计，无法满足安全监管要求
- 数据传输泄露和威胁分析能力不足



## 2.1、特权身份集中管理

### ▪ 主帐号集中管理

把具有特权身份自然人抽象定义为主帐号，所有可访问业务系统帐号密码信息抽象定义为从帐号，将所有主帐号和从帐号统一管理起来是特权访问管理的前提。通常采用三权分立原则对主帐号进行管理，可以划分为特权身份管理员、特权审计员和系统维护员三类角色权限，其中特权身份管理员负责对主帐号新建、编辑、权限分配、注销等一系列全生命周期管理；特权审计员负责对主帐号操作行为、从帐号使用情况进行审计分析，并对审计结果进行统计报表等；系统维护员负责对特权身份管理系统的配置、更新和维护等。三类权限相互牵制，防范特权权限监管真空区。同时结合双因素或多因素认证方式对主帐号进行身份鉴别，解决特权身份混用、冒用问题，也为安全事件指证和定则提供可靠依据。并引入身份鉴别防护机制，例如对暴力尝试破解密码行为进行锁定登录，静默会话自动注销，不能使用重复密码，帐号密码信息加密存储等等安全机制保护主帐号信息。

### ▪ 从帐号集中管理

把所有业务系统抽象定义为目标设备。将目标设备中的所有从帐号进行集中管理形成从帐号分布全景图，等同于管理好了访问企业信息资产保险库的“金钥匙”。基于全景图的基础上管理好“金钥匙”的分发和使用情况，同时也要做好周期性巡查工作，及时发现企业中未纳管的目标设备和从帐号信息。例如通过 SSO(单点登录) 技术使得主帐号用户在不知道从帐号密码的条件下也可访问业务系统和数据。周期性扫描 IDC 机房中存活的业务系统以及发现从帐号信息。

定期检查从帐号密码状态，及时发现异常情况，保管好“金钥匙”。例如周期性对从帐号密码有效性进行验证，可及时发现从帐号密码泄露或失窃，发现特权访问时越权改密操作行为。周期性对从帐号密码进行改密，使得密码满足强密码规范要求，解决从帐号密码泄露和失窃问题。周期性检测从帐号状态，可及时发现非法植入的幽灵(后门)帐号，因员工离职后未及时注销的孤儿(长期不用的)帐号等异常从帐号情况。对于核心业务系统“金钥匙”最好是能够改造升级鉴别机制，升级到双因素认证方式(即支持可知因素和不可知因素的双因素认证)，例如从帐号鉴别通过固定密码和动态密码组合方式进行认证，可彻底

解决密码丢失、窃取和周期更新问题。

## 2.2、访问权限集中管控

### ▪ 最小访问权限原则

将访问权限尽可能划分为最小粒度，仅赋予特权访问所需的最小权限集合，统一集中分配特权访问时的权限，形成特权访问权限全景图，清晰描述哪些自然人能够访问哪些业务系统，具备哪些访问权限，尽可能减少特权访问中权限滥用或越权行为发生。例如数据库从帐号按查询和编辑权限划分为两类帐号 user1 和 user2，当仅需要查询操作时分配 user1 即可，防范误删除数据。也可以按服务器应用特点，将权限划分为上传和下载权限来进行管控，例如文件服务器等。

### ▪ 金库模式

对于访问高价值业务系统和高危级别操作时，应采用实时金库模式进行管控，即配置“操作 - 监管”的双岗位模式对特权访问进行管理，实行高价值业务系统“一访问一审批”，高危级别操作“一操作一审批”，并对访问操作过程专人专岗实时管理。例如访问网络边界出入口交换机和防火墙时，修改访问控制配置或重启设备操作时，都应进行操作审批和确认。

## 2.3、全程集中安全审计

事后事件分析的主要内容是谁在什么时间，什么地点对哪个业务系统进行了什么操作，具备什么权限，进一步可以提升到操作者是谁管理的，谁导入到运维环境中的，事件中的业务系统主管单位或者主管人员是谁，访问权限分配是否合理，访问权限都是有谁分配和审核，经过了哪些调整。这些问题都可以通过安全审计的方式完整记录下来。事后分析中更重要的是能够完整还原事件的过程，准确评估事件的风险和损失。

## 2.4、数据加密和威胁分析

### ▪ 通信协议加密保护

加密数据是解决网络嗅探和监听的最好方式。对特权访问通信的数据流进行数据加密，可有效防范监听和流量还原导致的数据泄露情况。例如将文件传输 FTP 协议更新为 SFTP，TELNET 更新为 SSH，VNC 更新为 RDP 等等。

### ▪ 威胁分析和检测

业务系统被特权访问后留下的数据是否对业务系统稳定性、业务核心组建的安全影响有多大，是否存

在安全威胁? 这些问题时刻困扰着管理员和 CISO 们。由于特权访问的强隐秘性, 传统安全检测手段 (例如 IDS, 网络审计或安全沙箱等) 难以发现安全威胁。若是在传统安全检测技术基础上增加协议代理或数据摆渡技术可以有效解决特权访问过程中数据威胁分析, 提高数据安全能力。

### 三、总结

随着信息化建设的不断深入, 云计算、大数据等新技术不断发展成熟, 数据将成为生产资料, 将会是企业核心重要资产, 构建并提高数据访问管控能力至关重要。在特权访问管理的“识别”、“控制”、“审计”和“防护”四个层面上进行安全能力建设, 可有效降低特权访问下数据安全风险, 提高企业数字资产运行效率, 杜绝特权身份滥用、冒用和混用等问题, 防范威胁数据对数字资产的攻击, 满足法规监管安全审计要求。

# 数据库审计产品的技术运用

合作产品部 梁步庭

关键字：数据库审计、存在问题、UEBA、EDM 机器学习

摘要：数据库审计过程存在许多审计和运维问题。如何更好地解决? 本文以运用 UEBA、EDM 技术为例，着重介绍运用这两种技术解决审计结果达不到预期和预警信息运维的难题。

## 一、引言

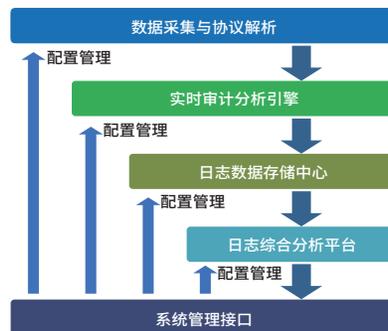
近几年随着数字化社会进程加快，数据库本身作为数据核心载体之一，且数据量庞大。数据的多样性、敏感性等特点凸显，也是数据安全泄露的高发地。一时间数据库安全成为数据安全中不可或缺的一环。早年数据安全概念已经不能和现在同日而语，数据安全已经被提高到了整个国家高度，被投以青眼越来越受到重视。这也就决定了数据库审计产品除了完成本职工作外，还需要肩负起数据安全这杆大旗，承担更多额外的安全使命。

所以如何运用成熟技术来修炼内功，夯实审计能力是当务之急。而本文主要探讨数据库审计除现有技术运用外，还可以运用哪些技术来解决数据库审计现

存的问题或用户需求。

## 二、数据库审计工作原理简述

简单回顾目前数据库审计已经运用的技术，以及审计过程，了解数据库审计产品的工作原理。



图一：审计原理图

通过主动(Agent 模式) 或被动方式(流量镜像模式) 采集网络数据, 通过协议识别、解析从被采集数据中筛选出 SQL 会话, 剥离数据中的通讯交互讯息、SQL 语句、返回结果等信息。

实时审计分析引擎将经过剥离的待审数据按照系统预设配置进行策略匹配, 完成审计结果、风险等级、告警等信息的输出展现。然后对数据库操作日志记录、会话、事件、统计信息等审计结果预计历史告警记录, 建立日志索引表, 完成全部日志的本地压缩、归档。日志本地保存同时开方对外日志接口, 实现日志备份。

日志归档后通过日志综合分析平台对历史审计日志进行检索, 统计、综合分析及价值挖掘等离线查询工作以及多格式的报表导出。系统管理端展现审计分析引擎和日志存储中心提供的实时审计日志与历史归档查询日志。

### 三、存在问题与解决方向

数据库审计产品的最终输出信息实质是一个个 SQL 语句和返回结果, 审计严重等级划分后通过或独立、或关联形式展现。但当数据量成指数级的增长,

原来传统的审计内容与结果的展现思路, 越来越不能满足当下的安全要求, 问题也随之暴露出来:

#### 第一个问题: 审计结果不达预期的问题。

数据采集开始到最终呈现审计结果和触发告警的前提, 必须与预设策略完全匹配, 而策略预设需要审计人员具备起码的 SQL 知识并了解实际业务, 还需要根据业务变化进行审计策略的调整, 对审计人员的学习和应变能力有一定技术要求。例如: 某企业对业务数据库只允许进行查询操作, 且返回结果不能超过 N 行。按此要求分析翻译后得到: 除 SELECT 以外的全部操作命令为风险语句, 且 SELECT 命令下的返回行数多余 N 行时即时告警的审计策略。回看这条审计要求, 在传统审计能力上限阈值不超过一定行数的要求该怎么设定, 显然只能对某个固定数值进行约束, 比如 100 行。但是, 固定值少了告警数量激增, 多了告警无法触发, 严重影响审计质量。同时阈值需求本质是希望对数据泄露进行一定预警和防范, 但固定阈值的方法只能解决一时。只要进行试探性尝试, 阈值上限完全可以通过语句控制进行规避, 比如只返回 99

行。显然传统策略能力得出的审计结果不达预期，以及存在漏报误报问题。

UEBA 通过机器学习来发现高级威胁，实现了自动化的建模，在发现异常用户行为、用户异常行为等方面有了非常高的“命中率”。引入用户实体行为分析 UEBA 技术，也就是将人对数据库合法操作行为和习惯纳入到策略范围，原本死板的策略能力就有了较大缓冲空间，对提升审计结果准确度、改善审计结果问题有较大帮助。

将自动化建模能力，纳入到 DAS 策略机制中，优先通过对合法行为进行分析。学习所得策略集，配合人工去重、查漏、补缺，形成最终行为策略。再将行为策略与普通策略结合，互相补充完成审计工作。大大节省了策略维护成本，同时还提高了审计精度，解决审计结果不达预期的问题。

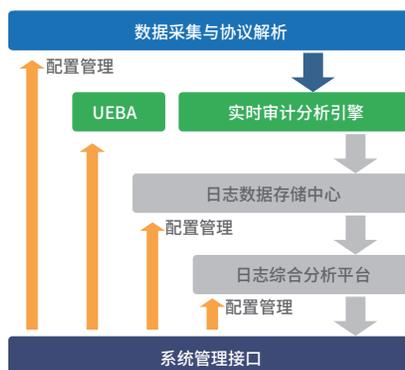


图 2：UEBA 原型

建立 UEBA 原型，整个数据库审计将变为两部分：

第一部分：学习与建模

当数据被采集完成协议解析后，待审计数据优先经过 UEBA 引擎进行学习，引擎通过定向定期学习（学习时长视流量大小与业务轻重而定），完成数据库操作行为的建模。建模完成形成策略集向引擎自动输出，交由审计人员进行人为干预（查阅、校准）。完成最终管理策略集制定并下发至实时审计分析引擎。

第二部分：审计与展现

回归原始审计处理流程。

## 第二个问题，预警信息难运维的问题

数据预警结果输出主要依靠图表统计、信息告警、详情列举等形式进行展现，但信息建立的基础全部来自于一个个抽象独立的 SQL 语句。例如图三、四，从产品特性和审计对象角度，这种审计结果符合预期。但就预警结果而言，预警结果不显性，阅读门槛高的问题暴露无遗。尤其当审计人员缺乏专业知识，面对大量 SQL 语句无法快速定位风险，还都需要再做进一步排查，投入大量时间成本。而被审计数据来源单一是主要原因之一，传统的数据提供方式主要依靠网络镜像或 Agent 提取，数据获取前提必须由客户端发起，数据库审计系统只能单方面接收和分析。哪些库、表、字段需要重点监控，什么数据，一概不知。



图三：告警信息

序号	SQL语句	执行时间	数据库名称	风险类型	风险名称	风险等级	数据库用户名	客户端IP	客户端名称
1	use information_schema	2019-01-28 18:05	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
2	SHOW COLUMNS FROM information_sche	2019-01-28 18:05	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
3	SHOW TABLE STATUS LIKE 'COLUMN'	2019-01-28 18:05	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
4	SHOW CREATE TABLE 'information_sche	2019-01-28 18:04	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
5	use information_schema	2019-01-28 18:04	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
6	SELECT * FROM 'information_schema'.	2019-01-28 18:04	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
7	use information_schema	2019-01-28 18:04	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
8	SHOW COLUMNS FROM 'information_sche	2019-01-28 18:04	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
9	SHOW TABLE STATUS LIKE 'COLLATION'	2019-01-28 18:04	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
10	SHOW CREATE TABLE 'information_sche	2019-01-28 18:04	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
11	use information_schema	2019-01-28 18:04	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
12	SELECT * FROM 'information_schema'.	2019-01-28 18:04	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	
13	use information_schema	2019-01-28 18:04	01_8_63_mysql	风险操作	高风险	中风险	root	192.168.8.140	

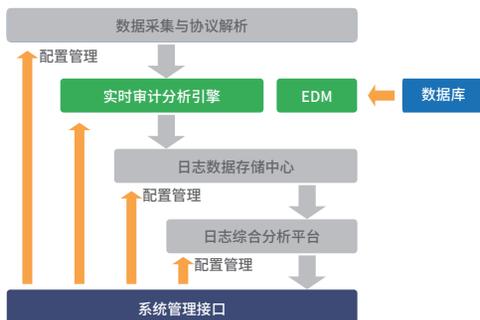
图四：详情页

精确数据比对技术 (EDM)，可用于分析、保护结构化格式的数据。允许根据特定数据列中的任何数据栏组合进行检测，例如在特定记录中检测 M 个字段中的 N 个字段。能够在“值组”或指定的数据类型集上触发，例如可接受名字与身份证号这两个字段的组合，但不接受名字与手机号这两个字段的组合，还支持相近逻辑以减少可能的误报情形。

EDM 可以主动对数据库、表进行扫描，分析数据类型了解被保护数据，形成库、表的指纹信息和表类型标识。配合正则表达式、业务化语言等能力实现对 SQL 会话的归类、翻译和去重，为风险告警提供关键信息依据，精炼预警信息改变告警信息的呈现方式。

EDM 技术的运用，对提炼、整合告警信息起到了决定作用。提高信息可读性，实现风险定位的最后一步成为可能。经过调整告警信息将会变为：

【警告! \*\* 用户 \*\*IP 正在 (或在 \*\* 时间段) 对 \*\* 实例下的 \*\* 类型表进行的风险操作，受影响行数为 \*\*，包含 \*\* 类数据，风险等级为 \*\*，累计操作时长 \*\*。详情中展现受影响数据节选和相关 SQL 会话。】



图五：EDM 原型

建立 EDM 原型，整个数据库审计流程将变为三部分：

第一部分：学习与建模

向 EDM 学习引擎告知一个数据库地址与合法身份，引擎主动发起一次对数据库全库表的扫描和分析，完成数据库各库表的初级指纹的建立，并对结其进行关键性描述（如：对库表是否携带银行卡、身份证、手机等信息列进行告知）。

第二部分：指纹筛选与建立

完成初级指纹建立后自动输出，交由审计人员进行人为干预对其进行矫正与过滤。形成关键性指纹，与

常规审计策略进行组合合并送至实时审计分析引擎中。

第三部分：审计输出

回归原始审计处理流程。

结语

探讨“新”技术在数据库审计中的运用，并不是为了让数据库审计更加“专家化”，而是让数据库审计更加“傻瓜化”。通过运用新技术和新方法降低数据库审计产品使用门槛，让越来越多的非专业人士也能轻松掌握和使用。这不仅仅是对产品的一次提升，同时也是对数据库安全、乃至数据安全的一次重要提升。

# 浅析亚马逊AWS数据安全的措施

解决方案中心 刘弘利

**摘要：**公有云的使用越来越普遍，用户利用第三方云计算平台，数据安全和隐私保护格外重要。本文以亚马逊 AWS 为例，分析云计算平台数据安全保护的措施。企业在自建的数据中心，私有云等平台，可以借鉴 AWS 的安全实践。

云计算有很多优势，越来越多的企业迁移到公有云开展业务。按照服务类型，云计算分为三种：基础设置即服务 (IaaS)，平台即服务 (PaaS) 以及软件即服务 (SaaS) 三种。应用和业务系统部署在云端，不管是哪一种云计算服务，应用系统的数据都会驻留在云端，数据安全性是云端租户需要考虑的重要问题。

云计算的数据安全，需要云服务提供商和租户一起分担。不同类型的云计算服务，二者的责任范围和比例有所区分。事实上，云计算平台的数据安全，更多的是租户的责任。

表列举不同类型云计算的网络安全责任划分，对于用户访问和数据安全，云计算租户的角色都是主责人。

私有云 用户自建数据中心	IaaS (基础设施及服务)	PaaS (平台及服务)	SaaS (软件及服务)
用户访问	用户访问	用户访问	用户访问
数据安全	数据安全	数据安全	数据安全
应用	应用	应用	应用
操作系统	操作系统	操作系统	操作系统
网络流量	网络流量	网络流量	网络流量
虚拟化环境	虚拟化环境	虚拟化环境	虚拟化环境
基础设施	基础设施	基础设施	基础设施
物理设备	物理设备	物理设备	物理设备

租户责任

云计算提供商责任

即将发布的等级保护 2.0 中，对云计算安全扩展的要求中，也有类似的规定。核心有两点，其一是对云计算平台自身安全的防护要求，其二是对云计算平台上租户的安全防护。

亚马逊 AWS，微软 Azure 和阿里云是公有云业界三强。据报道，亚马逊 AWS 云计算发展早，技术积

累深厚，一家公司就占了 51% 的全球市场份额<sup>1</sup>。以 AWS 为例，了解公有云计算平台数据安全保护的措  
施，企业在构建私有云或者使用其他公有云，可以借  
鉴 AWS 的数据安全保护措施。

## 一、AWS 数据安全

AWS 建议用户采用“良好架构设计”<sup>2</sup>，在开始上  
云时就注重整体架构。这是在 AWS 构建和应用的最  
佳实践，帮助用户构建稳定、安全、高效、低廉的系统。

良好架构设计的 5 大支柱如下表所示。

5大支柱	描述
Operational Excellence	运行和监控系统以提供业务价值，能够持续优化流程。
Security	通过风险评估和缓解策略提供商业价值，同时提供保护信息、系统和资产的能力。
Reliability	从基础架构或服务中断中恢复的能力，弹性获取计算资源以满足需求，并缓解诸如错误配置或瞬时断网问题。
Performance Efficiency	能够有效的使用计算资源，满足系统要求，并随着需求变化和技术的发展保持效率。
Cost Optimization	能够以最低价格运行系统，交付业务价值。

安全在 AWS 是 5 大支柱之一，可见安全对于云端  
业务系统重要性。良好架构设计，在安全方面建议了  
指导原则和安全最佳实践。数据保护是 5 个安全最佳

1 <https://tech.sina.com.cn/it/2018-11-29/doc-ihpevhcm3225635.shtml>

2 <https://aws.amazon.com/architecture/well-architected/>

实践之一。

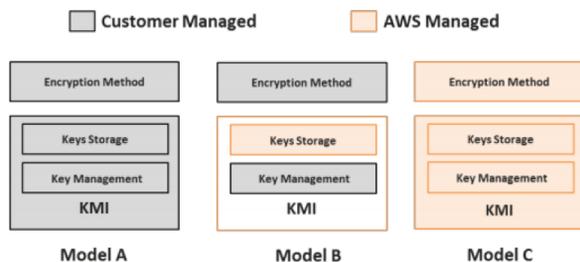
## 二、AWS 数据安全防护措施

### 1、数据分类

租户根据敏感程度以及重要性对数据进行分  
类，IAM、AWS KMS 和 AWS CloudHSM 进 而 依  
靠标记执行不同的管理策略。例如敏感数据存储  
在某个 Amazon S3 桶里，可以为该桶设置类似  
“DataClassification=CRITICAL” 的 标 记 ， 关 联 到  
AWS KMS 服务，就能透明的对存储和计算中的数据进  
行加密和解密。除了人工对数据进行分类，AWS 还提  
供了 Amazon Macie 来帮助租户发现、分类和保护存  
储在 Amazon S3 中的敏感数据。

### 2、加密和脱敏

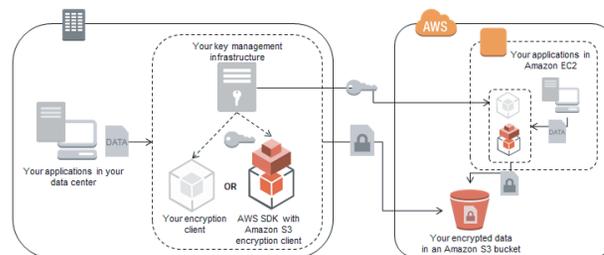
AWS 提供了 AWS KMS 和 AWS CloudHSM 两种密  
钥管理服务，创建和管理密钥，并且能够和 AWS 上绝  
大部分服务无缝集成。AWS 提供了三种密钥管理体系  
KMI(key management infrastructure)，租户可以根据  
业务情况进行选择加密方法、密钥存储和密钥管理的  
责任方。



数据脱敏 (Tokenization), 租户自定义无意义的 Token, 用来代替敏感数据 (如客户的信用卡号码), Token 和敏感信息的对应关系加密存储在 Amazon RDS 等数据库服务之中, 供应用程序调用。

### 3、持久化数据保护

AWS 建议对租户的持久化数据 (如对象存储、数据库等) 进行加密和访问控制, 从而降低未授权访问的风险。租户可以自行加密, 然后将数据放到 Amazon S3 存储上, 也可以利用 AWS 提供的 AWS KMS 服务, 在服务端进行加密。如下图为一种客户端加密的过程



### 4、传输过程中数据保护

对于传输中的数据, AWS 建议实用加密方式传输。其中, 对于管理员日常维护管理, AWS 建议利用 VPN 方式远程接入; 对于业务数据, AWS 建议采用 HTTPS 方式传输。

下表汇总 AWS 服务中, 数据从“ A ”点到“ B ”点流动的保护方式。

Point "A"	Point "B"	Data flow protection
Enterprise data sources	Amazon S3	Encrypted with SSL/TLS; S3 requests signed with AWS Sigv4
Amazon S3	Amazon EMR	Encrypted with SSL/TLS
Amazon S3	Amazon Redshift	Encrypted with SSL/TLS
Amazon EMR	Clients	Encrypted with SSL/TLS; varies with Hadoop application client
Amazon Redshift	Clients	Supports SSL/TLS; Requires configuration
Apache Hadoop on Amazon EMR		<ul style="list-style-type: none"> <li>Hadoop RPC encryption</li> <li>HDFS Block data transfer encryption</li> <li>KMS over HTTPS is not enabled by default with Hadoop KMS</li> <li>May vary with EMR release (such as Tez and Spark in release 5.0.0+)</li> </ul>

### 5、数据备份

最后一道防线是数据备份，不管是什么原因的数据丢失、损坏、误删除，都可以通过恢复备份的数据，将损失减少到最低。

### 三、AWS 数据安全保护的启示

AWS“良好架构设计”中，有两条原则与数据安全直接相关。一是保护传输中的数据与静态数据，二是建立让人远离数据的机制和工具，避免直接访问数据，减少数据泄露或者人为错误的风险。

对 AWS 数据安全措施的分析，我们得到以下启示：

- 数据安全不是孤立的，需要多层次安全措施
- 数据安全保护是一个整体，从数据分类分级打标记，实现自动的安全防护措施
- 既要考虑静态数据，也要考虑网络传输中的数据
- 数据备份是业务持续性的保障

数字化转型大潮中，云计算在其中扮演了重要角色，越来越多的企业采用公有云构建应用。数据和隐

私在云计算环境下的保护，是企业采用公有云的一大障碍。亚马逊 AWS，作为公有云行业的领导者，在数据安全的保护为云上租户提供服务。这既是为客户提供了数据保护的安全服务，也是在市场竞争中建立壁垒。本文浅述 AWS 数据安全的措施，感兴趣的读者，可以以此为契机，深入探索。



绿盟科技作为国内安全厂商中为数不多的 AWS 进阶技术类合作伙伴，目前已经为 AWS 客户提供 VPC 边界防护、运维安全、等保合规等场景的解决方案，和 AWS 共同保障客户云上业务和数据的安全。

---

参考资料：

1. AWS Well-Architected Framework
2. Security Pillar — AWS Well-Architected Framework
3. Data Lake Security Best Practice

# 数据内容识别技术深度剖析

解决方案中心 施岭

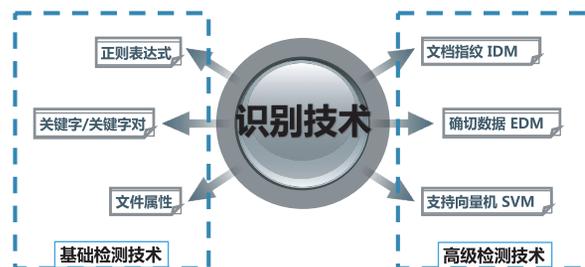
关键词：数据安全、内容识别、数据丢失

摘要：内容识别技术具备的识别能力包括关键字、正则表达式、文档指纹、确切数据源（数据库指纹）、支持向量机等。

## 一、引言

以前要管控数据，大多是强管控，直接全部隔离，或者全部加密，我们称之为囚笼、枷锁式的管控，在实际的数据生产、使用、流转中带来了许多不必要的麻烦，人们需要更加灵活的方式来处理数据，此时，智能化的数据安全管控应运而生，通过内容识别技术实现对数据的检索和分级，企业管理员可以按照数据的重要程度有针对性的对数据进行控制。

下面我们就来了解一下内容识别技术的能力。内容识别应该具备的识别能力具体来说有正则表达式、关键字/关键字对、文档属性、文档指纹、确切数据（数据库指纹）、支持向量机。



为了预防数据丢失，无论数据的存储、复制或传输位置在哪里，都必须准确地检测所有类型的机密数据。如果没有准确的检测，数据安全系统就会生成许多误报（将并未违规的消息或文件标识为违规）以及漏报（未将违反策略的消息或文件标识为违规）。误报会大量耗费进行进一步调查和解决明显事故所需的时间和资源。漏报会掩盖安全漏洞，导致数据丢失、潜在财务损失、法律风险并有损组织声誉。因此需要准

确的检测技术来做保障。为了确保最高的准确性，数据内容识别采用了三种基础检测技术和三种高级检测技术。

## 二、基础检测技术

基础检测技术中通常有三种方式，正则表达式检测(标示符)、关键字和关键字对检测、文档属性检测。

### 1. 正则表达式

正则表达式又称规则表达式，属于最常用的检测技术的一种，正则表达式通常被用来检索符合某个模式或规则的文本，此方法可以对明确的敏感信息内容进行检测。

### 2. 关键字 / 关键字对

关键字匹配同样属于最常用的检索方法，此方法可以准确地对敏感信息内容进行检索；

关键字对，则可以将被穿插了其他字符而故意混淆的一对关键字进行准确地检索，从而提高检测的准确性。

### 3. 文档属性

文档属性检测主要是针对文档的类型、文档的大小、文档的名称进行检测，其中文档的类型的检测是基于文件格式进行检测，不是简单的基于后缀名检测，对于修改后缀名的场景，文件类型检测可以准确的检测出被检测文件的类型，理论上讲可以识别任何类型的文件。

## 三、高级检测技术

高级检测技术中也有三种方式，精确数据比对(EDM)、指纹文档比对(IDM)、向量分类比对(SVM)。EDM 用于保护通常为结构化格式的数据，例如客户或员工数据库记录。IDM 和 SVM 用于保护非结构化的数据，例如 Microsoft Word 或 PowerPoint 文档。对于 EDM、IDM、SVM 而言，敏感数据会先由企业标识出来，然后再利用具有内容识别的系统来判别其特征，以进行精准的持续检测。判别特征的流程包括访问和检索文本及数据、予以正规化，并使用不可逆的打乱方式进行保护。

内容识别技术是以实际的机密内容为基础，而非根据文件本身。因此，内容识别技术不只能检测敏感数据的检索项或衍生项，而且能够标识文件格式与特征信息格式不同的敏感数据。例如，如果已经判别出机密 Microsoft Word 文档的特征，DLP 就能够在相同的内容以 PDF 附件的方式通过电子邮件进行提交时，将其准确检测出来。

## 1. 文档指纹 (IDM)

“指纹文档” (IDM) 可确保准确检测以文档形式存储的非结构化数据，例如 Microsoft Word 与 PowerPoint 文件、PDF 文档、财务、并购文档，以及其他敏感或专有信息。IDM 会创建文档指纹特征，以检测原始文档的已检索部分、草稿或不同版本的受保护文档。

IDM 首先要进行敏感文件的学习和训练，拿到敏感内容的文档时，IDM 采用语义分析的技术进行分词，然后进行语义分析，提出来需要学习和训练的敏感信息文档的指纹模型，然后利用同样的方法对被测的文档或内容进行指纹抓取，将得到的指纹与训练的指纹进行比对，根据预设的相似度去确认被检测文档是否

为敏感信息文档。这种方法可让 IDM 具备极高的准确率与较大的扩展性。

## 2. 确切数据 (EDM)

确切数据 (EDM) 可保护客户与员工的数据，以及其他通常存储在数据库中的结构化数据。例如，客户可能会撰写有关使用 EDM 检测的策略，以在消息中查找“名字”、“身份证号”、“银行帐号”或“电话号码”其中任意三项同时出现的情况，并将其映射至客户数据库中的记录。

EDM 允许根据特定数据列中的任何数据栏组合进行检测；也就是在特定记录中检测 M 个字段中的 N 个字段。它能够在“值组”或指定的数据类型集上触发；例如，可接受名字与身份证号这两个字段的组合，但不接受名字与手机号这两个字段的组合。

由于会针对每个数据存储格存储一个单独的打乱号码，因此只有来自单个列的映射数据才能触发正在查找不同数据组合的检测策略。例如，有个 EDM 策略请求“名字 + 身份证号 + 手机号”的组合，则“张三” + “1333333333”“110001198107011533”可触发此策略，

但是即使“李四”也位于同一数据库中，“李四”+“1333333333”“110001198107011533”也不能触发此策略。EDM 也支持相近逻辑以减少可能的误报情形。对于检测期间所处理的自由格式文本而言，单个特征列中所有数据各自的字数均必须在可配置的范围，方可视为匹配项。例如，依默认，在检测到的电子邮件正文的文本中，“张三”+“1333333333”“110001198107011533”各自的字数必须在选定的范围内，才会出现匹配项。对于含有表式数据（例如 Excel 电子表格）的文本而言，单个特征列中所有数据都必须位于表式文本的同一行上，方可视为匹配项，以减少整体误报情形。

### 3. 支持向量机 (SVM)

支持向量机 (Support Vector Machines) 是由 Vapnik 等人于 1995 年提出来的。之后随着统计理论的发展，支持向量机也逐渐受到了各领域研究者的关注，在很短的时间就得到很广泛的应用。支持向量机是建立在统计学习理论的 VC 维理论和结构风险最小化原理基础上的，利用有限的样本所提供的信息对模型的复杂性和学习能力两者进行了寻求最佳的折中，以获得最好的泛化能力。SVM 的基本思想是把训练数据非线性的映射到一个更高维的特征空间 (Hilbert 空

间) 中，在这个高维的特征空间中找到一个超平面使得正例和反例两者间的隔离边缘被最大化。SVM 的出现有效的解决了传统的神经网络结果选择问题、局部极小值、过拟合等问题。并且在小样本、非线性、数据高维等机器学习问题中表现出很多令人瞩目的性质，被广泛地应用在模式识别，数据挖掘等领域。

SVM 比对算法适合那些具有微妙的特征或很难描述的数据，如财务报告和源代码等。使用过程中，先将文档按照内容细化分类，每一类文档集合有属于本类的意义，经过 SVM 比对，确定被检测的文档属于哪一类，并取得此类文档的权限和策略。同时，针对 SVM 的特点，可以进行终端或服务器上的文档按照分类含义进行分类数据发现。

IDM 和 SVM 的比对区别是，IDM 将待检测文件的指纹和训练模型中的每一个文件进行指纹比对；而 SVM 是将待检测文件向量化，并归属到某一类训练集所建立的向量空间。

#### 四、内容识别技术的应用

现如今，内容识别技术已经得到了广泛的应用，比如数据的扫描挖掘、分类分级、交换共享、访问控制、泄露防护等。通过对不同检测技术的测试，基础检测技术速度更快，效率更高，而高级检测技术则效率偏低，而且对计算性能要求较高，因此，一般情况下会将内容识别技术以模块方式嵌入产品中，可以根据客户的需要，自由的选择在不同的场景应用合适的内容检测技术。

对于静态数据的扫描，基础检测技术与高级检测技术都应用较多，对于动态数据的检测，为了不影响业务连续性，及时发现问题并快速响应和处置，则基础检测技术应用性更高。不管是哪种检测技术，在应用时都是需要先预设检测规则，然后再进行检测，内容识别技术只对数据内容进行检测，如果想达到响应与防护的效果，必须与协议识别、应用识别、驱动识别等技术共同组合才能达到理想的目标结果。

总之，利用内容识别技术，在当今纷繁复杂的大数据时代下，从各类数据中快速精准有效的分拣出个人隐私数据或企业敏感数据，在提高效率的同时，让企业了解到数据的价值，并制定针对性的数据防护策略。

# 大数据环境安全管控浅析

SPG研发部 王豪

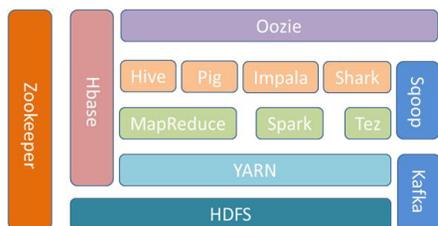
关键字：数据库审计、存在问题、UEBA、EDM 机器学习

摘要：数据库审计过程存在许多审计和运维问题。如何更好地解决? 本文以运用 UEBA、EDM 技术为例，着重介绍运用这两种技术解决审计结果达不到预期和预警信息运维的难题。

## 一、大数据环境和安全特性

大数据系统作为现在流行的数据存储、计算和分析平台，被越来越多的应用在各个行业。在数据为王的时代，由于大数据平台存储了大量的数据，对大数据平台及其存储数据的安全需求也越来越紧迫，在本文中，我们主要针对 Hadoop 的大数据环境的安全接管控进行分析。了解现有的大数据管控技术和思路。

典型的 Hadoop 大数据架构如下图，包含基础存储、资源调度、分布式计算框架、分析工具、数据采集等，图中包含了现有的基础组件和部分主流组件。



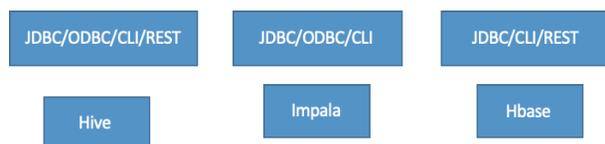
大数据环境相比传统数据环境在安全管控方面三大特性需要关注：

1、分布式，传统数据库虽然也有分布的概念，但应用更多的还是 hadoop 等大数据平台，由于 Hadoop 的特性，对单机的硬件性能要求不像传统数据库那么高，但带来的数据存储和访问的分布式，大量的集群节点使得防护的边界变得大更分散，增加对其进行管控的难度和复杂度；

2、多样性，由于 Hadoop 在近些年得快速使用，基于 Hadoop 平台衍生了大量的开源存储、计算等组件，从而使得安全人员需要了解大量组件的原理和部署来进行相关防护，而各种组件得更新和增加速度大大加快，也带来了更多的学习成本；

3、多接口，各种组件提供了多个访问接口，如 Hive 有 CLI、JDBC、ODBC 等接口，从而需要考虑多接口下的安全接入，避免遗漏。

在进行大数据系统的安全管控前，需要分析哪些组件主要提供对外接口，哪些组件会涉及到资源及数据的操作，这样才能了解边界，做到有的放矢，下图是 Hadoop 环境下的一些查询和存储组件的接口。



从上图可以看到这几个组件都提供了大量的操作接口，因此在了解各个组件提供了哪些接口后，在大数据业务上线前需要优先分析业务场景，确定各个组件需要开放的接口，对于不需要的接口，则要关闭相应服务和端口，避免接口对外暴露，减少安全风险。

## 二、大数据安全管控

在大数据安全中，安全管控处于非常重要的一环，哪些用户有权限访问，能做哪些操作，能访问哪些数据决定了平台和数据的安全性。一般来说，安全管控主要划分为认证管理、账号管理、权限管理和审计。

认证管理是对用户的接入进行安全认证，保证其有权限接入，账号管理是对所有的账号进行统一管理，保证账号安全性和易用性，权限管理是对资源进行管理，保证用户能且只能访问有权限的资源，审计是对用户操作、资源操作进行审计，满足合规要求，做到事后分析追踪。

## 三、认证管理

Hadoop 的 1.0 版本及之前，并没有安全的认证机制。默认集群中所有节点都是可信的，这导致恶意用户能够伪装成真正的用户或者服务端入侵集群，危害平台和数据安全。在之后的版本中，Hadoop 加入了 Kerberos 认证机制，集群中的节点使用密钥认证，只有认证通过的接口才能正常访问。

Kerberos 支持 Hadoop 环境下多个组件认证，支持 Unix、LDAP 等相关用户对接，使用 Kerberos 后将有效提升集群个节点的安全性。

## 四、账号管理

主要根据业务场景而定，而用户所在系统通常已有账号管理系统，如 LDAP，而 Kerberos 等也提供了和

其的对接，功能易于实现。

## 五、权限管理

由于大数据环境在安全管控方面的几大特性，对系统资源和数据的相关访问需要重点关注，从现有技术分析，主要有代理和插件两种方法来实现权限管理。

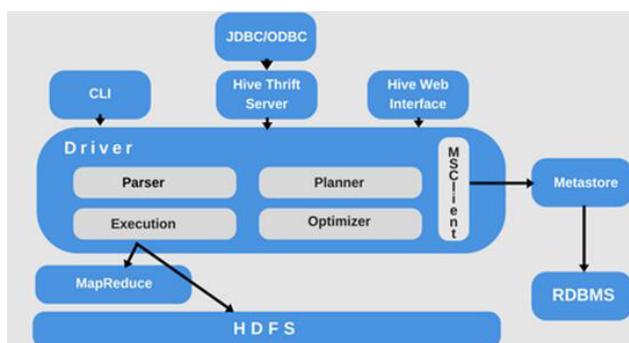
**代理方式：**代理方式主要针对通过网络方式访问的行为，在服务节点前设置代理，对上下行流量分析来进行权限管理。

**插件方式：**由于 Hadoop 环境中很多组件提供了 hook，因此在 hook 点上添加插件，能够实现相应的权限管理。

如对 Hive 的权限管理，Hive 是 Hadoop 的数据仓库，对存储在 HDFS 中的数据进行分析和管理的。它主要功能是将结构化数据映射为数据库表，提供类 SQL 的查询功能，将 SQL 语句转化为 MapReduce 任务进行运行。Hive 主要是针对传统的数据工程师、分析师等提供一种相似 SQL 的查询方式，减少学习成本，提供工作效率。

Hive 的基本架构如下，Hive 提供了 CLI、JDBC、

ODBC、HWI 的访问接口。其中 CLI 是 shell 访问方式，常用在运维等场景，JDBC、ODBC 主要用在客户端访问等场景。



各种接口的权限管理：

### JDBC/ODBC

在 JDBC/ODBC 的访问中，首先会连接到 HiveServer2，再通过 HiveServer2 进行后续的操作，因此考虑在 HiveServer2 中作为一个安全接入点。由于该接口用户网络访问接入，支持代理和插件的方式。

**代理方式：**在 HiveServer2 前使用代理，客户端首先连接代理，代理再将数据转发给 HiveServer2，代理可和策略中心联动，实时获取策略或将数据上传至策略中心，进行权限管理。

插件方式：插件功能是在 Hive 的 hook 点上对数据进行监听，从而实现权限管理，管理策略方式和代理类似。

### CLI

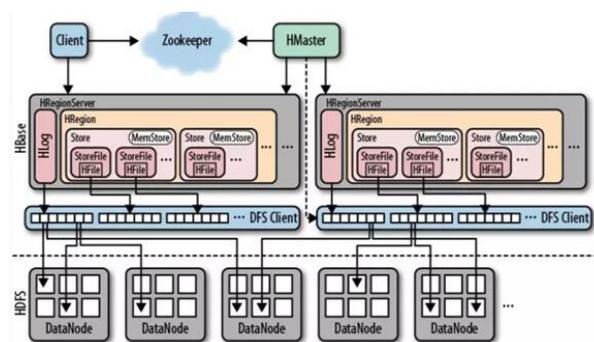
CLI 由于不通过 HiveServer2，因此通过插件的方式进行权限控制。分析 CLI 的访问原理，在通过 CLI 进行数据查询时，会从 Metastore 中获取到相关的元数据，从而去 HDFS 上读取相关数据，因此，在 Metastore 上实现对操作行为的监听，实现权限控制。

从上面 Hive 的分析中可以看出，实现组件的权限控制主要从以下几步入手：

- 1、 组件的功能分析，是否涉及资源或数据操作
- 2、 组件的接口和原理分析，是否通过代理或插件方式找到相应接入点
- 3、 组件的操作流程分析，结合集群部署状态分析接入点的部署形式

如 Hive 的 JDBC/ODBC 的权限控制只需要在 HiveServer2 节点上部署，这样就能对所有对 Hive 的操作进行控制。而在 Hbase 中，参考下图的 Hbase 架构图，由于客户端可直接通过 HRegionServer 访问

数据，因此在所有 HRegionServer 节点上都需要进行权限控制。



权限管理是对被访问资源的管控，因此需要做最大化的细粒度，现有技术中，能够做到对 Hive、Hbase 等的库、表、行、列、字段等细粒度控制，最大化满足用户需求。

## 六、审计

审计包括操作审计和数据审计。操作审计包括用户登录、退出等。数据审计是对资源、数据的访问行为记录，同时对行为进行分析，达到对主客体、时间、空间等的详细审计。而审计主要是通过代理或插件将审计日志上传至日志平台进行存储和分析。

综上所述，Hadoop 大数据环境中，大数据的防护管控不仅针对 Hadoop，更是对其生态圈中的大量组件要做到安全管控，特别是流行的查询、存储、计算和资源调度等组件，这样才能实现对整个大数据环境和数据的全方位控制。

## 七、大数据安全管控市场和技术分析

Hadoop 环境的安全管控现已有多个开源项目或国内外厂商提供相应产品或解决方案，如 Ranger、Sentry 等 Apache 顶级开源项目，很多大数据环境已使用这些项目进行管控，满足其安全需求。

但开源项目仍然存在相应的问题。比如 Ranger 所提供的功能只能简单满足客户的部分需求，在支持组件粒度、使用便捷性上都存在一定的问题，很多客户在此基础上自行开发或者第三方厂商基于此进行二次开发来满足更多的安全和使用需求。

在数据安全和大数据安全中，敏感数据发现、数据分级分类、加密、脱敏等都是必不可少的。大数据安全管控中仍然需要这些功能，用户在访问数据资源时，需要根据数据的安全级别和配置策略进行相应的

脱敏和加密，保证敏感数据内容不被泄露。而敏感数据发现和数据分级分类是数据安全最前期的需求，只有了解数据才能做到有策略的防护，同时数据的梳理对后续的管控和脱敏加密更能做到易用性和实时跟踪。

## 八、总结

在大数据安全平台中，加入认证管理、帐号管理，能够很好地守护大数据平台的大门，能够严格控制访问者的访问资格，将没有权限的访问者拒之门外。权限管理，能对用户的权限进行细分，什么帐号能访问什么资源，不能访问什么资源，都能进行控制。审计能够跟踪访问者进行的操作及访问的数据，当发生用户恶意操作（如删除数据表，篡改文件）时，能够进行追溯。从上面三个方面进行管控以后，能够极大地保护用户大数据平台及数据，实现大数据平台及数据安全的目标。

# 大数据安全之敏感数据发现

SPG研发部 肖春亮

关键字：结构化数据、非结构化数据、敏感数据、数据分级分类

摘要：本文会涉及到在大数据环境下，如何去区别结构化，半结构化，非结构化的数据；如何对结构化数据进行敏感数据发现；同时对非结构化数据的分级思考，实现非结构化数据中的敏感数据定位及发现。

## 一、引言

自国务院发布《关于促进大数据发展的行动纲要》以来，大数据在我国的发展与应用上升到国家战略层面，随着国家大数据战略的推进和大数据应用的逐步深化，大数据分析和应用得到了各个行业的关注，人们试图从大量数据中发现蕴含的模式和规律，进而产生更多的价值，“数据”作为分析对象在这个过程中所起到的作用是决定性的。不确定数据在哪里，不确定数据的价值以及敏感程度，也不确定数据访问权限以及数据的共享情况，是很多企业面临的问题。

## 二、传统数据与大数据的区别

传统数据与大数据主要存在 3 个方面的不同：

### ▪ 数据采集面不同

传统数据一般采集数据都是存储必要的的数据，如张三去吃饭，老板一般会记录如下信息：

客户	类型	食物	消费金额(元)
张三	早餐	牛肉面	15

而大数据采集数据会更加的全面，可能采集数据的数据如下：

客户	入场时间	入场方式	坐的位置	点餐时间	吃饭时间	食物	离开时间	消费金额(元)
张三	2017-01-01 9:30:27	自行车	A2	2017-01-01 9:32:27	2017-01-01 9:35:20	牛肉面	2017-01-01 9:50:20	15

### ▪ 存储方式不同

传统数据库一般都是用关系型数据库存储，如：oracle, mysql, postgres 等；

大数据存储一般都是使用 hadoop 生态进行存储，如 hive, impala, hbase 等。

▪ 数据类型不同

传统数据主要指结构化数据，半结构化数据。

大数据包括结构化数据、半结构化数据、非结构化数据。

大数据与传统数据相比的主要特点可以概括为：数据量“大而全”、数据类型“复杂”、数据价值“无限”。

三、大数据的类型

由于大数据存储的数据类型多而复杂，各种形式的数据可以归结为结构化数据、半结构化数据和非结构化数据三大种类。



▪ 结构化数据

(黄一，21，1989/06/09，201706092382)

▪ 半结构化数据

(姓名：黄一，年龄：21，出生日期：1989/06/09，学号：201706092382)

▪ 非结构化数据

黄一今年21岁，出生于1989/06/09，学号为201706092382

可见，结构化数据、半结构化数据、非结构化数据的最主要区别在于是否存在预先定义好数据模型，更确切的说是概念数据模型。

结构化数据能够用统一的某种结构加以表示，离开了这种结构，数据就没有意义；

半结构化数据具有某种结构，即数据的结构信息同数据混在一起，常见的有xml、html、json文件等等。

非结构化数据没有概念数据模型形式的限制，可以自由表达；非结构化数据不仅包含了文本，而且还包含图象、声音、影视、超媒体等典型信息，在互联网上的信息内容形式中占据了很大比例。由于非结构化数据中没有限定结构形式，表示灵活，蕴含了丰富的信息。因此，对于非结构化数据的敏感数据发现将会是更具挑战性的问题。

传统数据存储主要用于存储结构化数据，大数据存储可以存储结构化、半结构化、非结构化数据。

#### 四、数据的分级分类

在大数据中，不同的行业会有不同行业的敏感数据分级分类规则。一般情况下，可以将数据分级定义为4个级别：低敏感级、较敏感级、敏感级、极敏感；数据分类由于各个行业可能各不相同，分类级数可能也不一致。

在产品中，为了方便客户使用分级分类规则，一般会内置相关行业的规则，同时用户也可以复制内置的行业规则，形成自己行业的规则，并可以编辑自定义行业的规则。

用户可以使用对应的行业规则对其数据进行整体扫描、分级、分类，并根据分级分类结果做进一步的安全防护如细粒度访问控制等。

#### 五、大数据中敏感数据发现

大数据其数据量大而杂，并且无规律可寻。如何从大数据中识别出敏感数据是一个非常有挑战性的一个问题。



下面分别从结构化、半结构化、非结构化数据进行讨论，从这些数据类型中寻找敏感数据的方法。

##### 1. 结构化数据

结构化数据一般是存储于传统的关系型数据库或大数据平台的hive, impala等组件中，由于数据格式是固定的，所以对于敏感数据的发现是比较轻松的，同时能够确认数据的分类分级类型。比如，某运营商敏感规则中，有如下规则：

数据类别	一级子类	二级子类	对应数据	数据分级 4: 极敏感 3: 敏感级 2: 较敏感级 1: 低敏感级
A类: 用户身份相关数据	A1: 用户身份和标识信息	A1-1: 自然人身份标识	性别	3

当数据命中规则性别时，就认为该数据为敏感数据，其分级为：敏感级；分类为：【A类用户身份相关数据】【A1: 用户身份和标识信息】【A1-3: 用户基本资料】

在进行自动化敏感数据发现时，可以使用脚本程序去连接传统数据库或大数据平台的组件；首先通过API函数去获取有那些库，表，字段；然后对库，表，字段的值进行随机抽样，比如取第1-10条，第990-1000条，对取出来的值进行敏感规则匹配，若没有匹配上，则说明不是敏感数据；若匹配上，则说明是对

应分类及分级的敏感数据。

## 2. 半结构化数据

半结构化数据由于没有固定的数据格式，传统关系型数据库和大数据平台组件的 hive, impala 等组件都不方便存储，在进行敏感数据发现时，将其划分为非结构化数据中的文本的敏感数据发现。

## 3. 非结构化数据

非结构化数据一般存储在大数据组件的 HDFS 或 HBase 等中。

由于其没有固定的数据格式，比如文本数据，可能同时匹配多个敏感规则，就不能确定具体属于哪种分类，所以不方便对非结构化数据进行分类。所以对于非结构化数据，对其进行分级是一种比较合理的方式，若同时命中多个级别的规则，则为最高级别的敏感数据。

非结构化数据的分类分级是实现数据安全与合规的基础，对于不同类型的数据，需要采取不同的分类方法。不恰当的工具可能会产生不必要的业务问题，带来安全风险，增加额外的实施与维护成本。

### 3.1 非结构化数据分级方法

非结构化数据分级方法可以分为用户驱动和自动化完成两种，这两种方式是互为补充的。例如，自动分级能够简化用户驱动型分级的过程，用户驱动型分级也可以用于纠正自动分级中可能存在的错误。

#### ▪ 用户驱动型分级

对用户进行有关数据分级的培训，由用户主导对电子表格、报告、邮件、影视、超媒体等非结构化数据进行分级。该操作需要与相关文档应用协同开展，在操作系统、文件系统等层面完成。

用户驱动型分级工具通常通过自动化提高用户分级的效率，向用户建议分级并尽可能地减小工作量。这些工具一般还具备监控并执行数据分级策略的功能，例如：

- 要求用户在文件保存之前先对文件进行分级，或在发送电子邮件之前对邮件进行分级；
- 发现或阻止未经授权的分级更改行为。

#### ▪ 自动化分级

我们可以使用以下两种方法实现非结构化数据的

自动分级：

- 内容感知分级方法

该方法依赖于对非结构化数据内容的自动分析来确定分级，其中涉及了很多技术（正则表达式、数据字典技术，部分或完整指纹识别、机器学习等），应用的数据类型或应用程序不同，各类技术的适用性也有所不同。

对于信用卡号、社保号数据，正则表达式技术能够简便地检测到信用卡号、社保号等信息。

添加数据字典技术用以检测名称、地址或医疗条件可以提高正则表达式的准确性，但也会增加此类分级规则的复杂性。

指纹识别技术对于检测某些特定文档可能有效，但部分指纹识别需要持续性的维护，因为新的敏感信息还在不断地产生。

机器学习技术对于难以用用户定义模式描述的文档非常有用，但所产生的模型可能是不透明的，而且还会产生在定义阈值时难以解释的相似指数。

- 情境感知分级方法

该方法依赖于现有的分级知识库；利用广泛的情景上下文属性，同现有分级知识库对比，若相似度达到 70% 及以上，则将需要分级的数据划分到对应的

分级知识库中，从而确定数据的分级类型。

这种分级方法适用于静态数据。

理论上来说，用户驱动和自动化分级这两种分级方式可以处理所有数据类型，但结果的准确性则取决于用户的知识背景、勤奋度以及分级关系的情境确定性。当内容感知技术不够准确时，其它方法则加以补充。

## 六、总结

大数据的敏感数据发现是一个非常庞大而又有挑战性的工作，有的可以使用程序自动化发现，有的需要有用户的参与；要能够精确的发现大数据平台中的敏感数据，需要用户及安全厂商紧密合作，共同协作完成对企业大数据的敏感数据整理及定位。

# 数据安全小常识

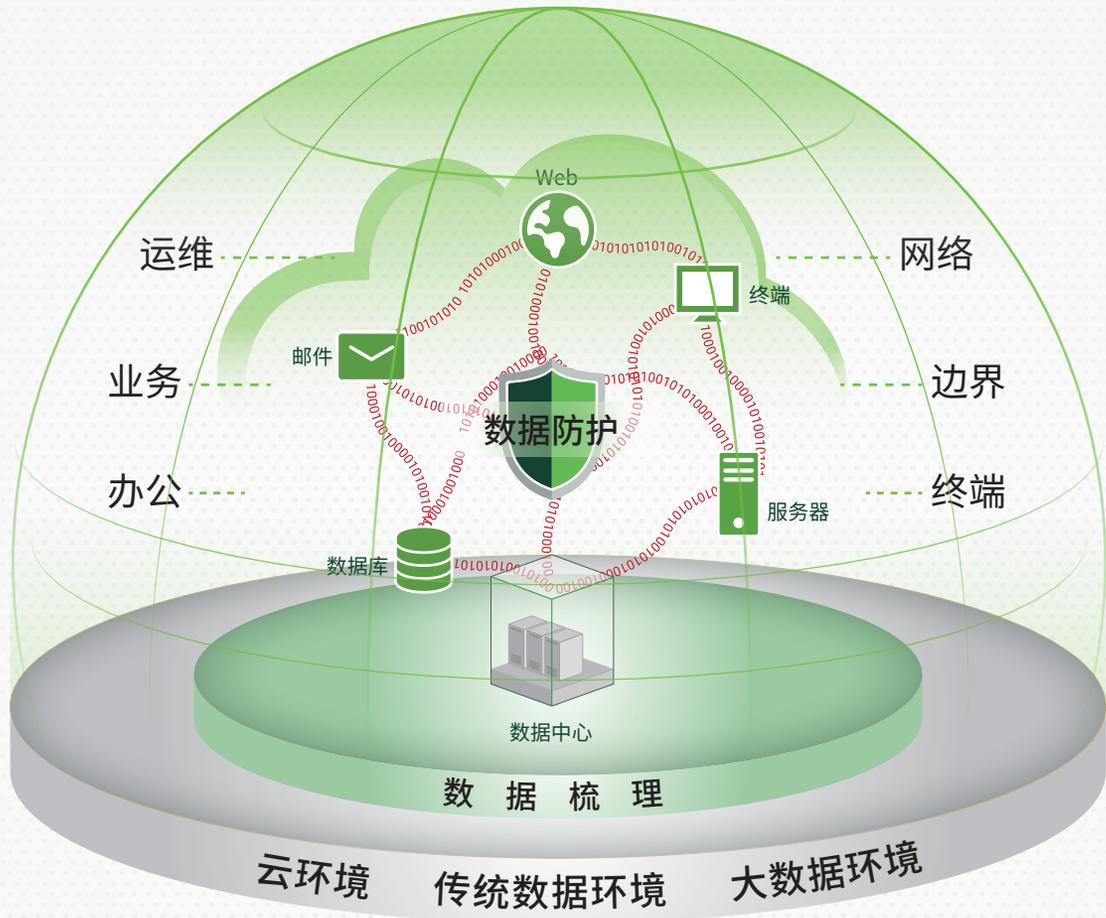


个人电脑安全口诀

- 敏感文件要加密，邮件文件莫忘记
- 软件请到官网下，MD5 要对比
- 文件删除要彻底，硬盘移交须脱密
- 密码设置要复杂，分级安全又好记
- 系统补丁及时打，不怕黑客常惦记
- 数据备份要定期，备份文件须加密

# 绿盟数据安全解决方案

## NSFOCUS DATA SECURITY SOLUTIONS



**THE EXPERT  
BEHIND GIANTS**  
巨人背后的专家

多年以来，绿盟科技致力于安全攻防的研究，  
为政府、运营商、金融、能源、互联网以及教育、医疗等行业用户，提供具  
有核心竞争力的安全产品及解决方案，帮助客户实现业务的安全顺畅运行。  
在这些巨人的背后，他们是备受信赖的专家。

# NSFOCUS IDR

绿盟敏感数据发现与风险评估系统

NSFOCUS Insight For Discovery And Risk

## 大数据安全治理必备武器



集多种扫描于一身



绿盟科技官方微信

**THE EXPERT  
BEHIND GIANTS  
巨人背后的专家**

多年以来，绿盟科技致力于安全攻防的研究，为政府、运营商、金融、能源、互联网以及教育、医疗等行业用户，提供具有核心竞争力的安全产品及解决方案，帮助客户实现业务的安全顺畅运行。在这些巨人的背后，他们是备受信赖的专家。

 **NSFOCUS** 绿盟科技