# Enable Personal Data Governance for Sustainable Compliance

Vincent Rejany; Bogdan Teleuca, SAS Institute Inc. Cary, NC

## ABSTRACT

In the context of European Union's General Data Protection Regulation (EU GDPR), one of the main challenges for data controllers and data processors is to demonstrate compliance by documenting all their data processing activities and where appropriate, to assess the risk of these processes for the individuals. Such requirements cannot be achieved without being able to build an efficient data governance program combining Legal driven top-down activities through personal data compliance and IT driven bottom-up operations through personal data mapping including personal data categories definition and discovery.

We use several processes developed in SAS® Data Management Studio to identify the personal data and update the governance view within SAS® Business Data Network and SAS® Lineage. We demonstrate several features in other products such as the Personal Data Discovery Dashboard in SAS® Visual Analytics, and SAS® Personal Data Compliance Manager as it applies to Records of Processing Activities and the Data Protection Impact Assessment.

## INTRODUCTION

Data governance is not an old concept; at SAS we have been pitching data governance benefits for years. However, it is often seen as something that is nice to have, even though it is a recognized method for mitigating risk, increasing operational efficiency, and enabling innovation.

As of this writing, we are close to the deadline to implement the requirements brought by the European Union's new General Data Protection Regulation (GDPR). On 25 May 2018, not only will GDPR be enforced across the European Union, but it will also have a global impact on all organizations that deal with the information of EU citizens.

The objective of the regulation is to give citizens more control over their data and to create a uniform set of rules to enforce across the continent. The main priority for organizations will be to show accountability by regaining control of their data processes, especially the processes and reasons for collecting, processing, updating, archiving, and deleting personal data records. To achieve such a task, being able to size the effort and discover the type and location of personal data is essential. Such a perspective is a key element for addressing data protection impact assessment when one process represents a high risk to the rights and freedoms of individuals.

GDPR breathes data governance and calls for discipline, integrity, and trust. "In the middle of difficulty lies opportunities" said Albert Einstein, so GDPR should be embraced as an opportunity to create value for your business, to gain a competitive edge, to innovate, to reinvent the way you manage your customer relationship, and to start doing more with personal data. Knowing the information that you hold, its quality, the reason that you hold it, and the length of time you can retain it is key in terms of operational efficiency. Organizations that support this idea of transparency will gain one competitive advantage by differentiating themselves in the market.

### WHAT IS PERSONAL DATA?

Personal data is any information that enables one person to be identified, directly or indirectly, by reference to an identifier such as a name, an identification number, location data, an online identifier, but also to one or more factors specific to the physical, physiological, genetic, mental, economic, cultural, or social identity of that natural person. Different pieces of information, which collected together can lead to the identification of a person, also constitute personal data. These are examples of personal data:

- name and surname
- home address

- email address such as name.surname@company.com
- identification card number such as VISA, American Express, or a loyalty card
- location data
- network identifiers such as IP addresses, even if they are dynamic
- cookie ID
- advertising identifier of your phone

Personal data that has been rendered anonymous in such a way that the individual is not or is no longer identifiable is no longer considered personal data. for data to be truly anonymized, the anonymization must be irreversible.

The regulation also defines the concept of special categories of data, for which specific safeguards and requirements are specified, such as a higher level of consent. These special categories relate to personal data that are "particularly sensitive in relation to fundamental rights and freedoms" and, therefore, "merit specific protection." These categories include data "revealing racial or ethnic origin, political opinions, religious or philosophical beliefs, or trade-union membership, and the processing of genetic data, biometric data to uniquely identify a natural person, data concerning health or data concerning a natural person's sex life or sexual orientation."

## UNDERSTANDING THE NEED FOR PERSONAL DATA GOVERNANCE

How can you comply with GDPR? By applying **G**ood **D**ata **PR**actices! The question is not how to comply, but how to be compliant and remain compliant. The aim of data protection regulations such as GDPR is to change behaviors and mindsets. Taking that perspective, the accountability principle (in Article 5 of the GDPR) makes the data controller to be the one responsible for demonstrating compliance with these GDPR principles:

- Lawfulness, fairness, and transparency must exist in processes that manage personal data.
- Limitation of purpose. Personal data must be collected for specified, explicit, and legitimate purposes.
- Data minimization. There should be no reason to use more data than necessary for the defined purpose.
- Accuracy. Data quality must be ensured and personal data be kept up-to-date.
- Storage limitation. Personal data must be processed for no longer than is necessary.
- Integrity and confidentiality. Appropriate security measures must be taken.

So how do you prove and show accountability? Under GDPR, accountability can be proven by nominating a data protection officer, drafting your privacy notice, and responding to requests (such as access requests) from individuals. However, the main action that you must take is to document internally all your processing activities, and to make this documentation available to supervisory authorities upon request. This "record of processing activities" is required by GDPR Article 30, and will facilitate the compliance with the other principles.

Data controllers must also carry out data protection impact assessments (DPIAs) when data processes could represent a high risk to individuals' rights and freedoms, particularly when new technologies are involved. The DPIA is required by Article 35 of the GDPR, and contains information about how a new or modified application might affect the privacy of personal information processed by or stored within the application.

Remember that any large organization has hundreds of systems, data assets, and processing activities, and thousands of personal data types to review daily, weekly, or monthly. Describing these items is a significant effort, but maintaining an up-to-date view of them is even more time-consuming and is prone to errors.

For data professionals (such as data owners, data stewards, and data controllers), the typical manual or semi-automatic steps no longer stand a chance when facing GDPR requirements.

## SAS APPROACH FOR PERSONAL DATA GOVERNANCE

The challenges that companies face in complying with GDPR are both external (toward supervisory authorities) and internal. When supervisory authorities perform a review of a company's data protection status, they will require a company-wide overview of all data sources and processes. And data sources aren't just systems or databases – they include network drives, files on PCs, and everywhere else personal data can be stored!

To address this challenge, there is one fair and straightforward method: "Say what you do and do what you say," which matches the classic data governance top-down and bottom-up analysis. Personal data governance (Figure 1) embeds both personal data compliance and mapping efforts.



**Figure 1. Personal Data Governance**

Personal data compliance aims at addressing the legal requirements, mentioned previously, to document data sources, to record processing activities, and to list the potential risks of these activities. Understanding data processes outside of IT is critical to capturing risks and potential control gaps. Such activities require advanced capabilities in the areas of structured methodologies shared across the organization, versioning, workflow for process management of approvals and reviews, and notifications. Because organizations can have thousands of data sources and hundreds of processes, relying on spreadsheets is not an option.

Personal data mapping is an approach that is proposed to facilitate the governance efforts and to significantly reduce the amount of time and effort needed to have the latest view of personal data. The location of personal data is essential information, to be as exhaustive as possible in your documentation and to show to the supervisory authority that you have established the processes needed to handle personal data. Moreover, recording the locations where personal data is stored will help your organization easily locate the information when an individual exercises his or her rights.

## PERSONAL DATA COMPLIANCE

### SAS® PERSONAL DATA COMPLIANCE MANAGER

SAS® Personal Data Compliance Manager is a recently released product based on the new SAS® Risk Governance Framework. SAS® Personal Data Compliance Manager is a workflow-driven and regulator-facing solution that automates the management of governance, risk, and compliance data. The product facilitates the entry, collection, transfer, storage, tracking, and reporting of operational losses, gains, and recoveries that are drawn from multiple locations across an organization.

The goal of SAS® Personal Data Compliance Manager is to provide organizations with customizable templates and workflows to document personal data processes and assess the risk of these processes.

The software has been developed based on GDPR requirements, as well as supervisory authorities' guidelines such as Working Party 29 (WP29), the CNIL (France), ICO (UK) and CPP (Belgium). SAS® Personal Data Compliance Manager is not specific to GDPR, and intends to address all data protection regulations, which are often strongly like the European law.

The product is intended to work in conjunction with the SAS® Data Management components and SAS® Visual Analytics. Together they provide a technology platform for participating firms to deal with the EU regulation on personal data protection

In this first release, SAS® Personal Data Compliance Manager can also be used to perform these tasks:

- document and maintain data processing activities
- define data controllers, processors, and data subject categories
- conduct data protection impact assessment
- describe and maintain data assets and systems
- define controls and security measures
- manage incidents, data breaches, and data subject correspondence

## RECORDS OF PROCESSING ACTIVITIES

### Legal Background

Data processing activity is defined in Article 4 (definitions 2 and 6) of the GDPR. Processing covers a wide range of operations performed on personal data, including both manual and automated processes. It includes the collection, recording, organization, structuring, storage, adaptation, alteration, retrieval, consultation, use, disclosure by transmission, dissemination or otherwise making available, alignment or combination, restriction, and erasure or destruction of personal data. The records shall be in writing, including in electronic form.

According to Article 30, the documentation of processing activities must include the following information:

- the name and contact details of the controller and, where applicable, the joint controller, the controller's representative, and the data protection officer;
- the purposes of the processing. In appendix, Table 1 lists some examples of activities, as recently provided by the Belgian supervisory authority.
- the basis of the processing that could be
    - necessary for the performance of a contract
    - A legal obligation
    - protection of the vital interests of the data subject
    - A task carried out in the public interest or in the exercise of official authority
    - legitimate interests pursued by the controller or by a third party
    - a result of data subject consent
- a description of the categories of data subjects and of the categories of personal data
- the categories of recipients to whom the personal data have been or will be disclosed including recipients in third countries or international organizations
- where applicable, transfers of personal data to a third country or an international organization, including the identification of that third country or international organization and, in the case of certain transfers the documentation of suitable safeguards
- where possible, the envisaged time limits for erasure of the different categories of data
- where possible, a general description of the technical and organizational security measures referred to in GDPR Article 32(1)

**Recording data processing activities with SAS® Personal Data Compliance Manager**

*Systems and Data Assets Definition*

Within SAS® Personal Data Compliance Manager, the recording of processing activities starts with the definition of systems and data assets. One system is a general identifier for one application, such as ERP, Finance, CRM. One data asset is more precise and allows to differentiate software from file system or databases. One processing activity can cover more than one data asset, and one data asset can of course be linked to more than one processing. Display 1 illustrates the creation of one data asset "Finance".



**Display 1. Personal Data Compliance – Data Asset Definition**

## Processing activities definition

SAS® Personal Data Compliance Manager supports the recording of processing activities through one dedicated web form. By default, the definition of processing activities is workflow driven so each step can be validated by approved users. One processing activity can be linked to one or more data assets.

The different sections have been defined based on the Article 30 of EU GDPR and the recommendations of various European supervisory authorities such as the ICO (UK), CNIL (France) and CPP (Belgium CBPL/CPVP). Display 2 shows the first tab of the form with the details about the processing activity.

**Display 2. Personal Data Compliance – Processing Activity Definition**

One of the most important part of the processing activity creation is the specification of the personal data categories involved in the processing as well as the data subject categories concerned. Display 3 presents the related section within SAS® Personal Data Compliance Manager Data Process.



**Display 3. Personal Data Compliance – Processing Activity - Data Information details**

## DATA PROTECTION IMPACT ASSESSMENT

### Legal background

Data Protection Impact Assessment (DPIA), also known as Privacy Impact Assessment (PIA) is one the specific process required by the GDPR. Impact assessments are not new, as similar risk assessments are also required by ISO/IEC 27001. DPIAs aim at identifying the risks related to the use of personal within one or several processing activities by evaluating them versus the GDPR principles mentioned formerly. For each risk identified safeguards and security measures must be defined. Article 35 of EU GDPR defines three conditions for which one DPIA must be conducted:

- When evaluating a natural person using automated processing (including profiling) to make decisions or have legal impacts on the subject. The use of new technologies, i.e. Big Data or Artificial Intelligence over personal data, is typically that situation
- When processing large quantities of special categories of data, or personal data relating to criminal convictions and offences.
- When systematically monitoring a publicly accessible area on a large scale, i.e. CCTV

According to the UK Information Commissioner (ICO) one DPIA should contain:

- A description of the processing operations and the purposes, including, where applicable, the legitimate interests pursued by the controller.
- An assessment of the necessity and proportionality of the processing in relation to the purpose.
- An assessment of the risks to individuals.
- The measures in place to address risk, including security and to demonstrate that you comply.

### Assessing data protection risk with SAS® Personal Data Compliance Manager

Within SAS® Personal Data Compliance Manager, one DPIA pre- assessment is required over each processing activity. This pre- assessment is key as the absence of DPIA would have to be justified to the supervisory authority in case of an audit.

Information collection has been structured through three tabs and is workflow driven. The first tab "Details" in Display 4 provides a view of the processing activities involved in the DPIA as well as an overview of the Privacy Risk Assessment. Depending on the risk severity of one assessment, regular reviews will be required. SAS® Personal Data Compliance Manager also supports the ability to define and schedule notifications for ensuring such review.



**Display 4. Personal Data Compliance – DPIA Definition**

The Data Protection Assessment Tab, Display 5, proposes to conduct the DPIA through a series of questions covering topics such as: Collection, Use, Retention, Sharing and Transfer, Access and Security, Data Privacy Assessment. Relying on the answers provided, the risk (inherent and residual) will be calculated by SAS® Personal Data Compliance Manager. Documentation can be attached to the DPIA.



**Display 5. SAS® Personal Data Compliance Manager – Data Protection Impact Assessment**

## PERSONAL DATA MAPPING

Personal data mapping is an essential task for complying efficiently with personal data protection regulations. Under GDPR, organizations must be fully accountable for all data flows, but visualizing and mapping data flows without the appropriate tools is very challenging. It is one key activity to demonstrate accountability by proving the reliability of your documentation effort to authorities, especially toward the recording of processing activities. Efficiency is important as organizations can have thousands of data assets, avoiding manual process through automated processes for discovering and mapping personal data attributes is required.

Data mapping is at the cornerstone of data governance, and data protection regulation require higher visibility regarding data processes. Data mapping enables you to get back control of your data management activities by supporting reverse engineering of your different data assets. You can then gain valuable insights for improvements in design and privacy compliance.

### PERSONAL DATA MAPPING METHODOLOGY

Four steps have been defined to support Personal Data Mapping through the SAS® for Personal Data Protection end-to-end approach. SAS® for Personal Data Protection is an accelerator containing pre-built assets to identifying, governing and protecting personal data, which relies on SAS® Data Management, an industry-leading solution built on a data quality platform that helps you improve, integrate and govern your data.

**Figure 2. Personal Data Mapping Methodology**

**Document Data Assets & Personal Data Categories:** From a risk and compliance point view this step has been fulfilled within SAS® Personal Data Compliance Manager. Therefore, these definitions are sourced from SAS® Personal Data Compliance Manager and presented into SAS® Business Data Glossary to be accessible for Business and IT people. Metadata related to these assets must be linked to applications or data owners.

**Identify Personal Data**: Personal data discovery is the key part of personal data mapping. It aims at finding, cataloging, and analyzing personal data attributes, across data assets. Results of the personal data discovery process are surfaced into a SAS® Visual Analytics Dashboard.

**Associate & Review Personal Data**: Once personal data have been identified within one data asset, related columns metadata will be automatically associated with their corresponding personal data category business term within SAS® Business Data Network. It will allow to document that such personal data attribute is present in this data asset and to get a clear mapping between metadata, personal data business terms, data assets and processes. Associations with a medium or low confidence are submitted for manual review

**Visualize Data Flows & Relationships**: Having centralized the definitions of the data assets, processes, data owners, personal data categories linked to databases and data jobs metadata, it now possible to surface a complete picture of the relationships between the objects. Such data flow is exposed into the processing activities definition with SAS® Personal Data Compliance Manager and the related personal data categories are provided by the personal data discovery step.

## DOCUMENT DATA ASSETS & PERSONAL DATA CATEGORIES

After the initial work has been performed within SAS® Personal Data Compliance Manager, there is the need to share this perspective with a larger audience, beyond the Legal or Compliance departments. SAS® Business Data Network is a business data glossary, part of SAS® Data Management that enables collaboration between business, technical, and data steward users. SAS® Business Data Network can be used as a single-entry point for all data consumers to better understand and govern their data asset through the definition and the maintenance of business terms.

Business terms can be organized through hierarchies and relationships, and can be linked to different roles such as data or business owner or data steward. Different types of terms can be defined according to the information that needs to be documented.

**Display 6: SAS® Business Data Network Main View**

In the context of the personal data governance approach, seven term types have been defined and more than 700 terms related to data privacy have been already integrated into SAS® Business Data Network. Table 1, below, presents the term types implemented

| Business Term Type | Description |
|---|---|
| Subject Area | Master entity identifying the subject or the project being described, such as personal data protection |
| Domain | Defines the second-level object types, such as data assets, processing activities, and personal data categories |
| Data Asset | Describes the databases or application potentially containing personal data, such as CRM, ERP, and data warehouse |
| Business Process | Defines business processes such as customer onboarding or marketing campaign, which is often the purpose of the processing activity |
| Processing Activity | Defines the processing activity including the information requested by the regulator |
| Personal Data Category | Contains the definition of the personal data to be identified, such as E-mail address or Network Address |
| Wiki | A generic term type used to enter glossary-like definitions, such as PII, consent, or data transfer |

**Table 1. Term Types Defined in SAS® Business Data Network**

## Personal Data Protection Bridge

The personal data protection bridge (PDP Bridge) is an asset developed by SAS® Data Management experts, tha enables information exchange between SAS® Personal Data Compliance Manager and SAS® Business Data Network. It consists of a database and processes for extracting information to be synchronized between the two products. Objects defined within SAS® Personal Data Compliance Manager such as data assets and processing activities are automatically populated into SAS® Business Data Network and are available for use. Once the personal data categories are identified and associated with one data asset in SAS® Business Data Network, the categories are synchronized from SAS® Business Data Network to SAS® Personal Data Compliance Manager through the PDP Bridge.

The PDP Bridge uses the data export and import features available in SAS® Personal Data Compliance Manager and the SAS® Business Data Network REST API.



**Figure 3: Personal Data Protection Bridge Overview**

## Data Asset Definition

Information about data assets can either be entered manually or imported from SAS® Personal Data Compliance Manager using the PDP Bridge. If you use either SAS® Federation Server or SAS® Data Quality as a personal data discovery engine, related connection information must be provided. This information is critical, because by default these connections are unknown in the SAS metadata.

**Display 7: Data Asset Example in SAS® Business Data Network**

Next, you need to add the connection as a library type to the database, in the list of associated items of the term. Defined SAS libraries (Display 8) are available from SAS® Lineage, which is integrated with SAS® Business Data Network.



**Display 8: SAS Library Association**

## Personal Data Categories Definition

A personal data category, such as **Phone**, **Payment Card Number**, **Delivery Address**, and so on, is the term that regroups the personal data that is identified in different sources. This also ensures that users understand the definitions behind the personal data categories.



**Display 9: Personal Data Category Example within SAS® Business Data Network**

Firstly, you need to point all Personal Data Categories, such as an e-mail address, to a SAS® Quality Knowledge Base personal data token.

Secondly, you need to set a confidence level for the identification.  The Personal Data identification uses fuzzy matching and the Quality Knowledge Base definitions. The usage of a match percentage is to link automatically all data with a match percentage over the high limit; send in a remediation queue (a manual confirmation process) all data with a match percentage between the low limit and the high limit and reject everything under the low limit. The SAS® Quality Knowledge Base is explained in the next paragraph.

You set, for example the QKB Personal Data Token to E-mail for the "E-mail address", the High limit to: 85, the Low limit to: 65. In the section "Associate & Review Personal Data" we will see how these thresholds are used.



**Display 10: Personal Data Category Definition**

## IDENTIFY PERSONAL DATA

Personal data discovery is the core part of the data mapping exercise. This process is about the categorization of personal data in structured data sources. The discovery report provides an overview of locations (databases, directories, columns etc.), types of personal data, indication of amount and other information. Unstructured data sources can also be analyzed but require another approach using text analytics. The approach is not covered in this paper.

### Enable Personal Data Discovery with the SAS® Quality Knowledge Base

You can use the SAS® Quality Knowledge Base which consists of a repository that contains literally thousands of pre-built data quality rules applicable across the full range of data subject areas, including Customer and Product. The SAS® Quality Knowledge Base contains all the SAS® Data Management algorithms to standardize, identify, extract, fuzzy match data. Each locale offered by SAS has built-in data quality algorithms for typical data types (such as name and address conventions, phone standards and so forth). The algorithms vary from region to region. So far, 35 countries are currently covered.

SAS® Data Management Studio is used to identify the personal data and update the governance view in the business glossary.

For facing such a challenge in an efficient way, you rely on the capabilities present within the SAS Quality Knowledge Base for identifying personal data. One identification definition aims at categorizing one string based on vocabulary and structure analysis. Each possible token defined is score against the input string and the best score is returned. Figure 4 illustrates this principle.
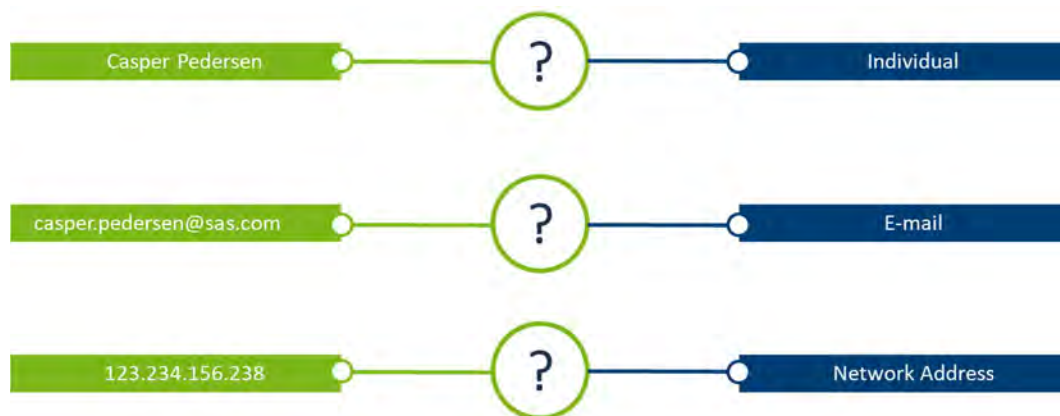


**Figure 4: Personal Data Discovery**

Thanks to the work already done by SAS, we are already able to identify more than 15 types of personal data categories in multiple languages. Depending on the country and language this list can contain additional personal data categories.
Table 2 presents the list of attributes identifiable in the context of the Personal Data Protection initiative.

| Personal Data Tokens | |
|---|---|
| Country, Country (Iso2), Country (Iso3) | Individual |
| Date, Date/Time | Organization |
| E-Mail | Phone |
| Geographical Point | Postal Code |
| IBAN | National Id |
| Network Address | Vehicle Registration |
| Payment Card Number | Identity Card Number |
| URL | Passport Number |
| City | Tax ID |
| Delivery Address | |

**Table 2. PDP – Personal Data (Core) Identification Definition Tokens**

The SAS Quality Knowledge Base features are available in multiple SAS products such Base SAS®, SAS® Data Quality, SAS® Data Integration Studio, and SAS® Federation Server. It is also available in Hadoop, Teradata, and SAS® Event Stream Processing.

Each locale can be extended to meet the needs of any organization by using the customize function of SAS® Data Management Studio to enhance or add algorithms as necessary. SAS works closely with its international offices to constantly refine existing locales and develop new versions

**Automate Personal Data Discovery using SAS® DataFlux Data Management Studio**

SAS® DataFlux Data Management Studio is a data management platform that combines data quality and data integration. It provides a process and a technology framework to deliver a single, accurate and consistent view of your enterprise data. With SAS® DataFlux Data Management Studio, you can establish an effective data governance through a wide range of activities, including personal data discovery.

SAS® DataFlux Data Management Studio offers multiple way for using the SAS® Quality Knowledge Base discovery capabilities presented formerly, including:

- Data Exploration
- Data Quality jobs

In the example below, we use Data Exploration to discover personal data in a table BP000.



**Display 11: Data Exploration – Data Source configuration**

Display 12. Shows the configuration of the Data Exploration job to discover personal data based the field on two definitions, "PDP – Personal Data (Core)" and "PDP – Personal Data (Patterns)". This last definition contains industry specific regular expressions to identify codes like VIN numbers, IMEI, or IMSI.



**Display 12: Data Exploration – Identification Definition configuration**

The identification analysis produces an overview of the results, by personal data category. You also get a sample view of the data that have been identified in the process. The field in the source, FIELD9, was identified a Payment Card Number with a 96% match percentage.
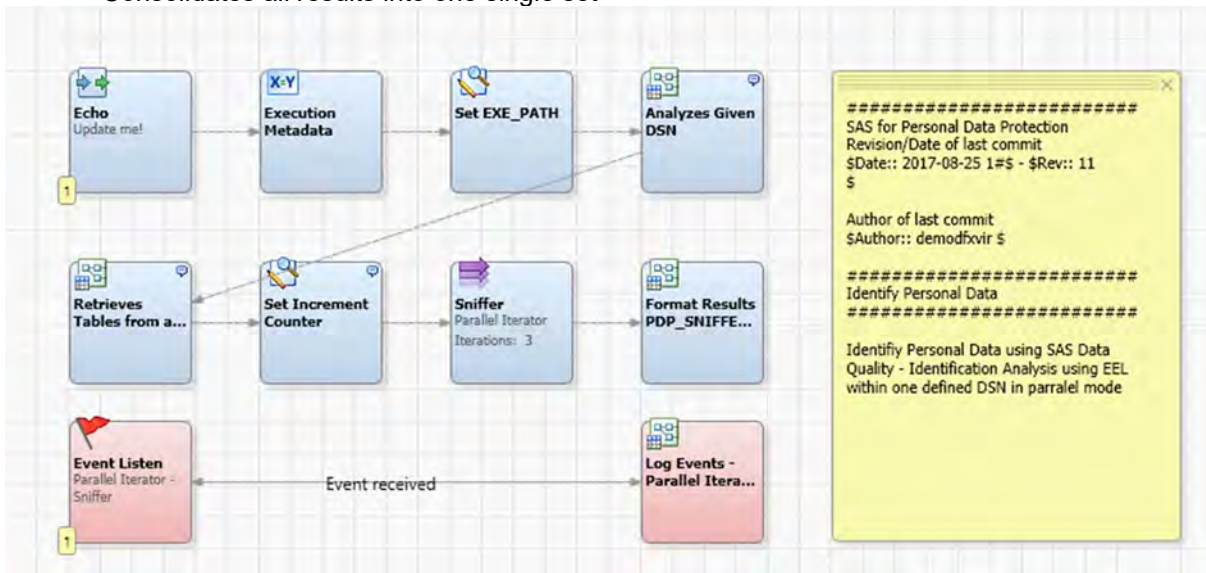
**Display 13: Data Exploration – Identification Report**

## *Data Quality Jobs for Personal Data Identification*

Data quality jobs are the main way to process data in SAS® DataFlux Data Management Studio. Each data job specifies a set of data-processing operations that flow from source to target. These data processing operations can address multiple objectives such classic ETL operations, data quality operations like data standardization, matching and clustering, as well as personal data identification.

To industrialize and accelerate the personal data discovery process over one database, we have designed one dedicated data quality process job, named the "Personal Data Sniffer". For one connection and schema, the "Sniffer":

- Creates one folder for storing all the process output results and logs
- Analyzes the connection and parameters provided (sampling, randomization, definitions to apply)
- Lists all the tables contained in the connection and schema configured
- Initializes one parallel iterator to process multiple tables at the same time
- Executes the personal data identification over each table
- Consolidates all results into one single set



**Display 14: Data Quality Job – Personal Data Sniffer**

## Report on personal data discovery in SAS® Visual Analytics

Results delivered by Data Exploration or Data Quality Job executions are challenging to analyze as the volume of information can be quite significant. In that perspective we have designed in SAS® Visual Analytics the Personal Data Discovery Dashboard to review the results of the different personal data discovery processes.

The dashboard provides a clear inventory of the personal data identified in any of the sources analyzed. The personal data categories are displayed by source, data asset, in which tables and columns they are present. Data Stewards and other users have a clear overview of the personal data contained in their applications, reducing their crypticity and making them transparent.



**Display 15: Personal Data Discovery Dashboard – Full Report Tab**

The "Explorer" tab allows you to investigate why a certain field has been identified as personal data category by the SAS® Quality Knowledge Base, in which table and with what match percentage.



**Display 16: Personal Data Discovery Dashboard – Full Report Tab**

## ASSOCIATE & REVIEW PERSONAL DATA

### Link Personal Data Discovery results with Personal Data Category Terms

In the first step of the methodology we have defined the personal data categories and data assets terms into SAS® Business Data Network. For each data asset we know the ODBC connection as well as the corresponding SAS® Library. The personal data discovery process provides one personal data attribute for each table and column analyzed. In the second step, the field metadata is associated to their corresponding personal data category term in SAS® Business Data Network.



**Figure 5: Metadata Linker Process – Overview**

This process named the "Metadata Linker" connects data assets with the personal data categories automatically if personal data is identified in a table or file within the connection defined in the data asset. The level of confidence and the mapping with the SAS® Quality Knowledge Base token helps you validating automatically this association.

Display 17 illustrates the different steps of the process.

- Retrieve all defined personal data categories and data assets from SAS® Business Data Network
- Map Personal Data Discovery results with Personal Data Category Terms considering only the scores with the highest confidence (Figure 5 – Step 1)
- Search for the tables and columns into SAS® Lineage for each data asset according to the SAS Library mapped
- Associate column metadata with Personal Data Category Terms (Figure 5 – Step 2)
- Link identified Personal Data Category Terms to the related Data Asset Term



**Display 17: Metadata Linker Process – Detailed view**

Associations are automatically posted to SAS® Business Data Network using its REST API. For the Data Asset Terms, the results of the Metadata Linker process are visible in the Related Terms section. The "Finance" data asset presented in Display 7 now contains personal data categories (Display 18)

**Display 18: Data Asset Term associated with Personal Data Category Terms**

The full list of columns metadata identified as "E-mail" has been added automatically by the Metadata Linker to the Personal Data Category "E-mail". You now have a complete view (identification, definition, governance, lineage) with just one process.



**Display 19: Column Metadata associated automatically to a Personal Data Category Term**

## Remediate Low Confidence Association in SAS® Data Remediation

In SAS® Business Data Network, for each personal data category you had to specify a QKB personal data token and set a high limit and a low limit for the match percentage. After you specified these items, the personal metadata linker performs these tasks:

- links data with a match percentage over the high limit
- sends in a remediation queue (a manual confirmation process) all data with a match percentage between the low limit and the high limit
- rejects everything under the low limit

You can use SAS® Data Remediation to review propositions made by the personal metadata linker.



**Display 20: SAS® Data Remediation – Association Remediation main screen**

When you open a proposition, you are being asked: "Associate Column with Term", then, "Do you want to associate the column FIELD4 from the table TABLE_US in library ORACLE with term Phone?)". The term Phone is the term from the Personal Data Category in SAS® Business Data Network. You might have no clue what the FIELD4 contains, but you can review the content behind. A stored process behind opens a window with 20 sample records from the table TABLE_US, for the FIELD4 column.



**Display 21: SAS® Data Remediation – Stored process showing column values**

At this point you might recognize a telephone number pattern and decide that the proposed association is legitimate. You decide to associate the Field4 with the term Phone in SAS® Business Data Network. From SAS® Data Remediation, one REST API call is executed over SAS® Business Data Network to post the changes. The new associated field is visible in the "Phone" Personal Data Category Term (Display 22).
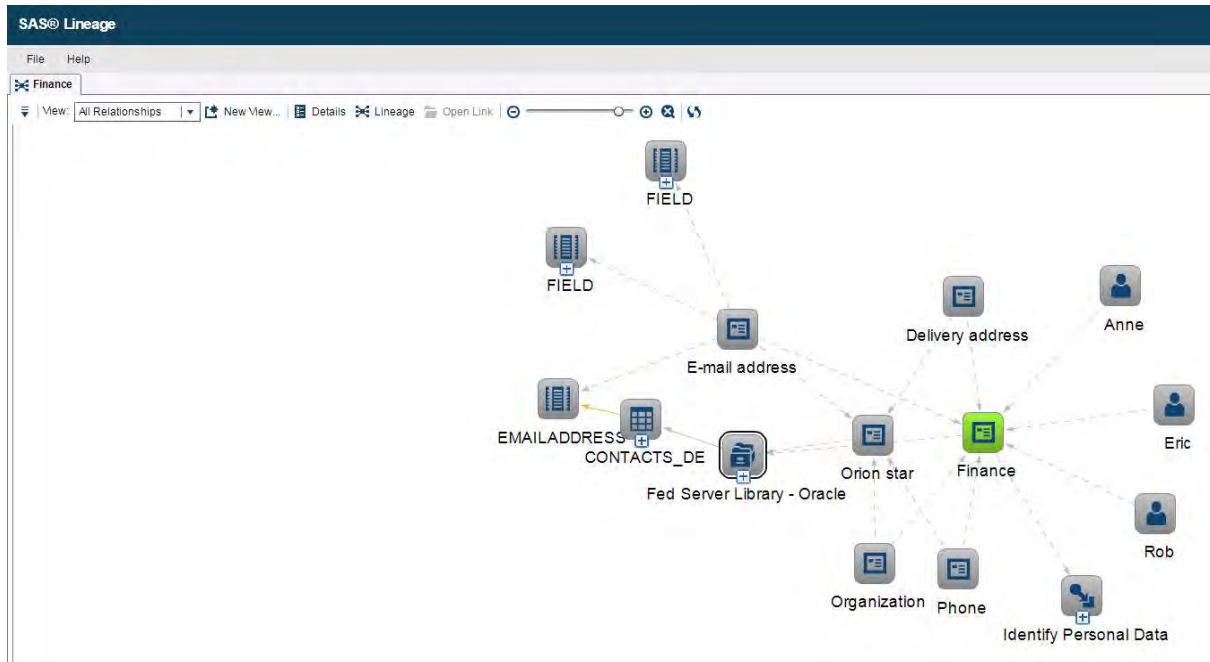


**Display 22: SAS® Data Remediation – Manual Review**

## VISUALIZE DATA FLOWS & RELATIONSHIPS

### SAS® Lineage

Once the former steps completed you can have a complete view of all the data assets and the personal data categories through SAS® Lineage. Lineage capabilities tie business and technical information together in a single and cohesive information store. It is a strategic component of a data governance plan, in the quest to understand the data assets in your organization.
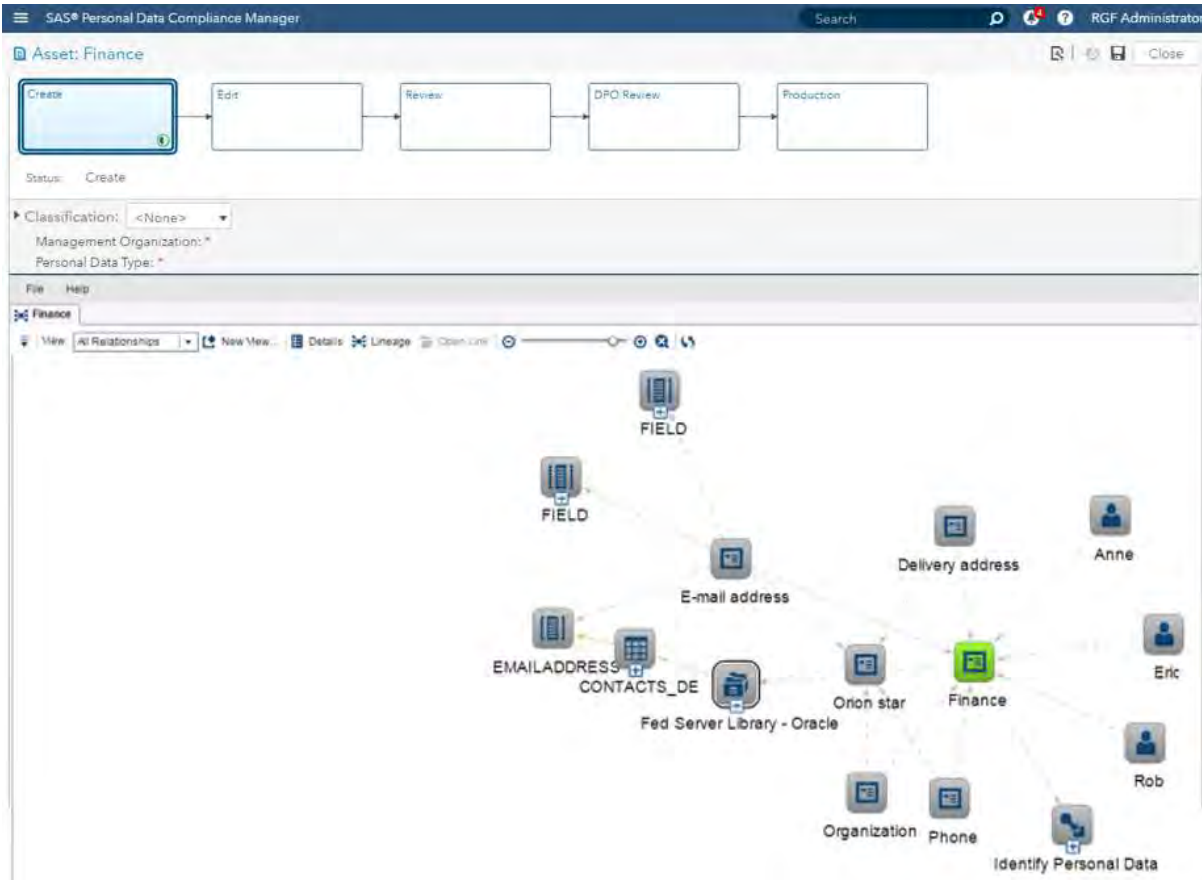
SAS® Lineage supports the management and analysis of object and metadata relationships, including dependencies and lifecycle. This management and analysis process reveals where data comes from, how it is transformed, and where it is going. It surfaces all the metadata present in the SAS® Platform including libraries, tables columns, but also SAS® Data Integration jobs, SAS® Data Quality jobs and SAS® Visual Analytics reports. Impact analysis can be activated to reveal which objects are affected when another object is changed or deleted.



**Display 23: SAS® Lineage – Finance Asset Overview**

## SAS® Personal Data Compliance Manager integrated with SAS® Lineage

The Personal Data Protection Bridge (PDP Bridge) makes the Lineage View accessible in the Asset screen from SAS® Personal Data Compliance Manager. This opens the door to the synchronization of personal data categories discovered with the view available to a Data Controller.



**Display 24: SAS® Personal Data Compliance Manager – SAS® Lineage Integration**

## CONCLUSION

We explained the legal background set by the General Data Protection Regulation, the provisions of the Article 30 for the records of processing activities and personal data definition. We then looked how data processing activities, data assets and data protection impact assessments are defined in the SAS® Personal Data Compliance Manager.

We introduced the need for a Personal Data Mapping Methodology. We demonstrated the need for SAS® Business Data Network and SAS® Lineage as a central point to maintain the personal data categories definitions, the sources of data where these can be found and the link with the data assets via related terms.

We proposed a personal data identification process with SAS® Data Management Studio to identify personal data using the SAS® Quality Knowledge Base definitions for personal data protection and highlighted the need for a Personal Data Discovery Dashboard in SAS® Visual Analytics to explore the results. Lastly, we introduced the usage of SAS® Data Remediation to solve low confidence proposals and stressed the importance of SAS® Lineage to visualize the links and associations.

The benefit for having the SAS platform at the core of the personal data protection initiatives is to reduce considerably the effort and the time to answer the regulator or the customers: what personal data is available in the sources, what are the processing activities and the risk posed by those to the liberties and rights of data subjects.

## REFERENCES

- European Parliament and the Council of the European Union. 2016. *REGULATION (EU) 2016/679 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data*. Brussels, Belgium: European Union. Available: http://eur-lex.europa.eu/legal-content/EN/TXT/?uri=uriserv:OJ.L_.2016.119.01.0001.01.ENG

- *Information Commissioner's Office. 2017. Preparing for the General Data Protection Regulation (GDPR): 12 Steps to Take Now. Cheshire, England: Information Commissioner's Office. Available: https://ico.org.uk/media/1624219/preparing-for-the-gdpr-12-steps.pdf*

- EU General Data Protection Regulation (GDPR): An Implementation and Compliance Guide - *ITGP Privacy Team*

## ACKNOWLEDGMENTS

## RECOMMENDED READING

- Casper Pedersen, 2016. "SAS Solution for Personal Data Protection, Enabling compliance with new regulation". Available https://www.sas.com/content/dam/SAS/en_us/doc/other1/solution-for-personal-data-protection-108517.pdf

- Rineer, Brian, 2015. "Garbage In, Gourmet Out: How to Leverage the Power of the SAS® Quality Knowledge Base." Proceedings of the SAS Global Forum 2015 Conference. Cary, NC: SAS Institute Inc. Available http://support.sas.com/resources/papers/proceedings15/SAS1852-2015.pdf

- Hoffritz, Cecily, 2017. "I Spy PII: Detect, Protect, and Monitor Personal Data with SAS® Data Management." Proceedings of the SAS Global Forum 2017 Conference. Cary, NC: SAS Institute Inc. Available: http://support.sas.com/resources/papers/proceedings17/SAS0639-2017.pdf

## CONTACT INFORMATION

Your comments and questions are valued and encouraged. Contact the authors at:

Vincent Rejany
Domaine de Grégy
Grégy-sur-Yerres
77257 Brie Comte Robert Cedex
SAS Institute, Inc.
+33 (0)6 40 54 17 99
vincent.rejany@sas.com
http://www.sas.com

Bogdan Teleuca
Hertenbergstraat 6,
3080 Tervuren, Belgium
SAS Institute, Inc.
+32 475 36 33 81
bogdan.teleuca@sas.com
http://www.sas.com