

IBM Guardium
Discover and Classify

SUPERVISED AI OVERVIEW

IBM GUARDIUM DISCOVER AND CLASSIFY

VERSION 4.2.2

TABLE OF CONTENTS

Table of Contents	2
Glossary of Terms and Acronyms	3
Purpose	4
Use cases	5
Possible actions	5
Main scenarios	5
Structure	8
Landing page	8
RDAs & Virtual Views page	8
Create a Virtual View/RDA	8
Data Training - Virtual View Modification page	9
Virtual view visualization (right part).	9
Mapping (left part)	9
Checking existing mapping	10
Preview Sample Data	10
Manage mappings	10
Advanced mapping configuration	12
Candidates Page	12

GLOSSARY OF TERMS AND ACRONYMS

Table 1: Glossary of Terms

TERM	DEFINITION
Appliance name	Name of the IGDC analytic appliance or CM appliance.
Bubble	Visualization of a group of candidates combined by common properties - data source and virtual view.
Candidate	Personal data instances that were not confirmed against any of the RDAs and therefore cannot be considered "trusted" (verified).
Candidate Virtual View Tab	Tab for the identified virtual view structure modification and adjustment of the column mapping to the Supervised AI data elements (Supervised AI > Data Training).
Constraint	Single data element or a combination of data elements that specify a unique person in the IBM Guardium inventory of personal data.
Data source	Network element that stores data in a structured or unstructured format. The IGDC uses them as a data source to create an inventory of the personal information.
Data subject	Instance of retrieved personal data confirmed as a unique data subject due to match with RDA.
Database	Structured data storage serving as a source of personal data.
Database name	Database vendor name. For example, Oracle, MySQL, MariaDB.
Group (complex)constraint	Combination of data elements that specify a unique person in IBM Guardium inventory of personal data.
Hostname	Data source URL, which format depends on the data source type and vendor.
Last time analysis	Date and time of the last data source analysis.
Mandatory field	Attribute of a virtual view column specifying it as mandatory for personal information instance retrieval. The rows with a mandatory blank field (potential data subject record) will be ignored during virtual view analysis.
RDA/VV list	List of root data assets and approved virtual views created or modified in the Data Training page.
Relationship Tab	Tab for virtual view creation from scratch, including table links and column-data element mappings (Supervised AI > Data Training).
Root Data Asset (RDA)	Set of structured data used by IBM Guardium as a source for comparison with detected personal information.
Sample Data	Table with data elements of 30 data subjects retrieved by the modified virtual that facilitate verification of the column-data element mappings.
Schema	Logical collection of objects (tables, views, indexes) in a database.
Supervised AI	<p>Wizard-based machine learning tool enabling a "non-data scientist" to build, modify, and train AI models on how to identify personal and sensitive information based on discovered results.</p> <p>It assists users in the analysis of the identified candidates based on database virtual views and training the system to accept or reject the candidates as trusted data, creation of RDAs, and applying the system-detected virtual views as RDAs.</p>
Supervised AI data element	Name used by IGDC software for specific personal data types. For example, CC_NUMBER for credit card numbers. The Supervised AI data elements are identical to the data elements across the IGDC modules like Personal Information Search, Advanced Search and others. You can change the list of data elements and modify data element recognition in the Data Element Configurator (CM UI > Settings > Data Recognition > Data Element Configuration).
Training mode	Mode triggered by changes in a virtual view when the system trains IGDC analyzers to identify Information according to the new structure. The training process is sequential, meaning that the system begins processing a changed virtual view only after finishing training the previous group of candidates (bubble).
Virtual View (VV)	Set of database tables or schema used as a source of a specific group of candidates forming a collection of identical data elements designated to the subject candidate group.

PURPOSE

Supervised AI is a wizard-based machine learning tool enabling "non-data scientists" to build, modify, and train AI models to identify personal and sensitive information based on automatically discovered results.

The Supervised AI module visualizes the possible sources of personal data (candidate virtual views) and assists in training the system to accept or reject the candidates as trusted data. You can also create root data assets and apply the system-detected virtual views as root data assets.

Data training using the Supervised AI tool improves the quality of data recognition during IGDC discovery and allows the customization of the system according to your network and business specifics.

USE CASES

POSSIBLE ACTIONS

1. Visualizing the identified candidates in groups by data source and virtual view. See [landing page](#).
2. Modifying the field mappings to the supervised AI data elements based on system mappings within the candidate virtual view. See [Candidate Virtual View tab](#).
3. Modifying the virtual view by adding relationships with out-of-virtual-view tables from the same data source. See [Add non-detected table links to virtual view](#).
4. Reviewing a list of candidates. See [Candidates list](#).

MAIN SCENARIOS

1. **Ignore a group of candidates to skip them in future analysis.**

Click **Ignore the VV** under virtual view visualization to mark this virtual view as invalid; the system will flag this virtual view as "ignored," and it will not be analyzed by IGDC in the future analysis.

2. **Accept a group of candidates as "trusted" data.**

Click **Valid VV** under virtual view visualization to accept this virtual view as "trusted" data; the system will flag this virtual view as "valid". In the next data source analysis, this personal information will be either mapped to the existing RDAs (can happen because of mapping modification) or mapped to the default root data asset (list of all "trusted" data subjects generated by the system).

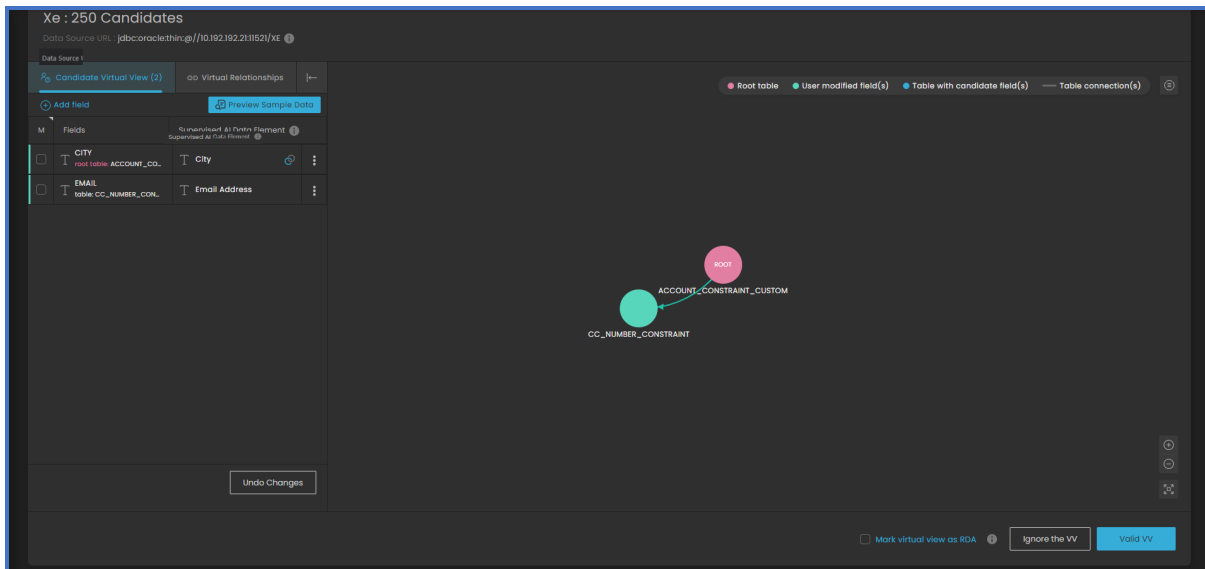


Figure 1: Data training

3. **Create an RDA based on the virtual view.**

Check **Mark that virtual view as RDA** and then click **Valid VV** under virtual view visualization. In the popup that open, enter the RDA name and description. The system will create an RDA entry with all the data subjects from this virtual view.

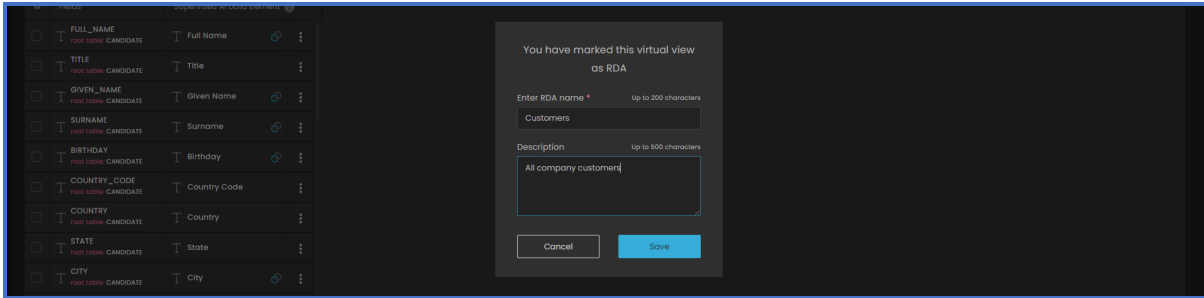


Figure 2: Virtual view-based RDA popup

4. Training mode

After scenarios 1, 2 or 3, the system switches to training mode. The system then trains the IGDC analyzers to identify information according to the new structure.

The training process is **sequential**, i.e. the system begins processing a changed virtual view only after completing the training of the previous group of candidates (bubble). However, you can return to the landing page and train other candidate virtual views at any point in the process.

Upon accepting or rejecting the virtual view, the changes will be applied after the next scheduled analysis of the relevant data source.

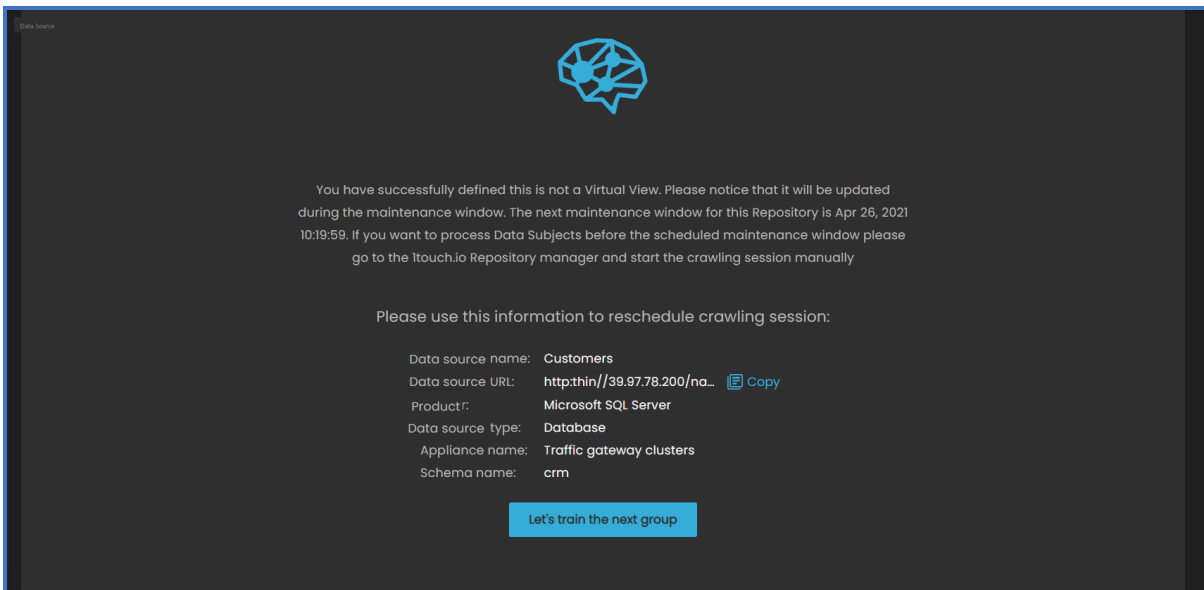


Figure 3: Training mode screen

5. Review the flow.

On the **RDAs & Virtual View** page, you can review your work history, and edit/delete any entry: valid or ignored virtual views, as well as created RDAs.

Supervised AI | RDAs / Virtual Views

Filter Collapse All Reset

0 Filters active

RDA / Virtual View List

Search by data source name, data source URL, RDA name Create RDA/Virtual View

RDA/VV Status	Date	Data Source Name	Data Source Type	Data Source URL	Schema	Created	Processing Status
> VV Ignored	31.08.2022 02:35:51AM	N/A	N/A	jdbc:oraclethin://10.192.192.21:1521/XE	CANDIDATES_250	Manually	Valid
> VV Approved	31.08.2022 02:18:15AM	N/A	N/A	jdbc:db2://10.192.192.4:50000/testdb	FIAT_IT_40_FULL	Manually	Changed
> RDA	31.08.2022 02:07:17AM	N/A	N/A	jdbc:oraclethin://10.192.192.21:1521/XE	CANDIDATES_250	Auto	Changed

Figure 4: Flow

6. Create an RDA from scratch. See [Create a virtual view/RDA](#).
7. Create a virtual view from scratch. See [Create a virtual view/RDA](#).

STRUCTURE

The Supervised AI module is a part of the IGDC platform. It consists of six screens.

LANDING PAGE

The landing page consists of two parts: **main window** and **sidebar**. The main window visualizes groups of candidates with identical attributes in "bubbles".

Table 1: Attributes for grouping candidates in a bubble

ATTRIBUTE	EXPLANATION
Data source	All candidates retrieved from the same data source (databases).
Virtual view	All candidates are associated with the same virtual view

The main window shows 27 random bubbles of three sizes - large, medium, and small (nine of each size) with brief information regarding the candidates group in each bubble.

The sidebar filters allow you to modify the "bubbles" shown in the main window as follows:

- Candidates' group size - large, medium, small;
- Database vendor - MySQL, Oracle, etc.;
- Appliance - name of the IGDC analytic appliance that identified and retrieved the group of candidates.

You can also search for the desired bubbles by database or hostname.

Hover over any bubble to see auxiliary information - number of personal information instances, database vendor, type of processing, appliance name and time of the data source last analysis.

In hover mode, click Train to start working with the group of candidates.

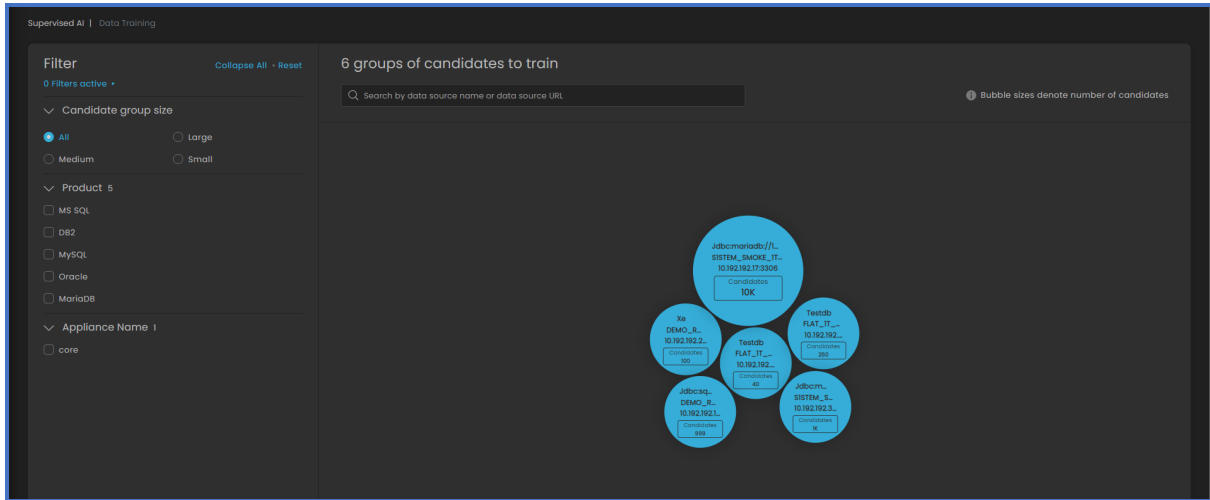


Figure 2: Supervised AI landing page

RDAS & VIRTUAL VIEWS PAGE

On the **RDAs & Virtual View** page, you can review, edit, and delete the history of your work - all valid and ignored virtual views, as well as created RDAs.

CREATE A VIRTUAL VIEW/RDA

In supervised AI, you can create an RDA or virtual view from scratch in 3 steps:

1. Click **Create RDA/Virtual View** on the **RDAs & Virtual View** page and select the desired data source.

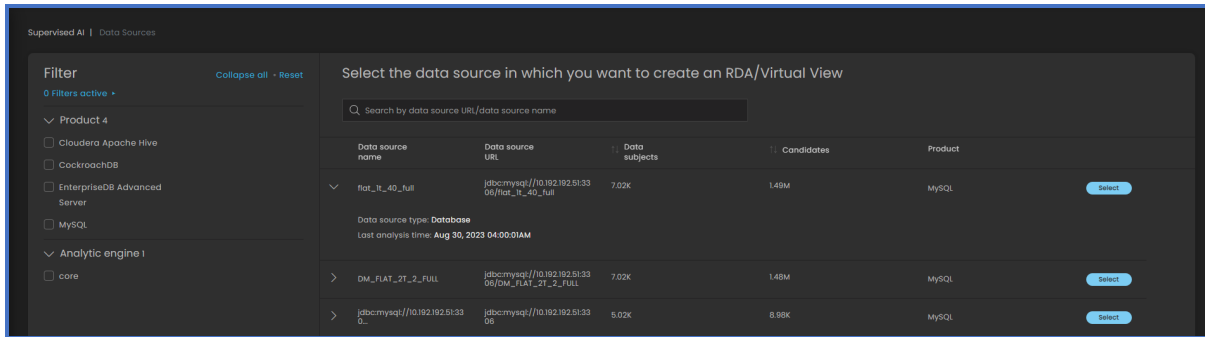


Figure 1: Selecting data source for RDA

2. Define relationships between tables and fields in the relationships tab in the relationships tab. The same logic as for [virtual view modification](#).
3. Define mappings for the desired fields in the Candidate Virtual View tab. The same logic as for [virtual view modification](#).
4. Click **Valid VV** to save a virtual view OR check **Mark that virtual view as RDA** and then click **Define as data subject** to create an RDA.

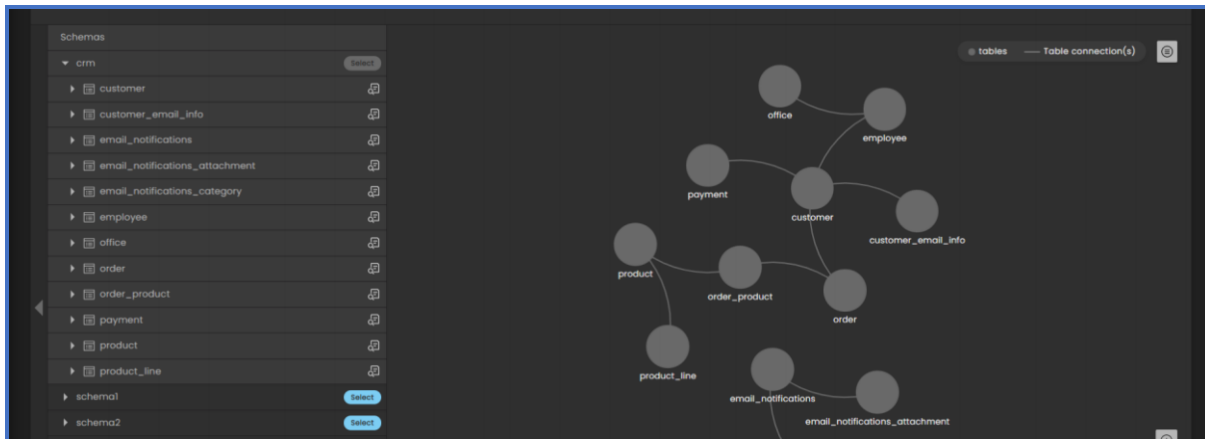


Figure 2: Table relations definition

DATA TRAINING - VIRTUAL VIEW MODIFICATION PAGE

VIRTUAL VIEW VISUALIZATION (RIGHT PART).

Circles represent tables that were detected as part of a virtual view (**blue** - non-modified; **green** - modified, **pink** - root table). The table name is shown under the circle. If the circles are linked by a grey line, it means there is a relationship (connection) between them.

MAPPING (LEFT PART)

Mapping means the association of a field from a virtual view table with an entity type supported by Inventa. For example, a field with *ssn* column name can be mapped with the *ssn_number* entity type.

Checking Existing Mapping

The **Candidate Virtual View** tab is the "light" mapping mode based on table relationships built by the system. A user can ensure that the mapping is correct and modify some mappings if necessary, as well as configure fields as mandatory (i.e. the system will skip a potential PII candidate if this field is empty).

A virtual view mapping must contain at least one constraint (🔗 - field is a constraint; 🔗 - the field is a part of a complex constraint like given name + phone number).

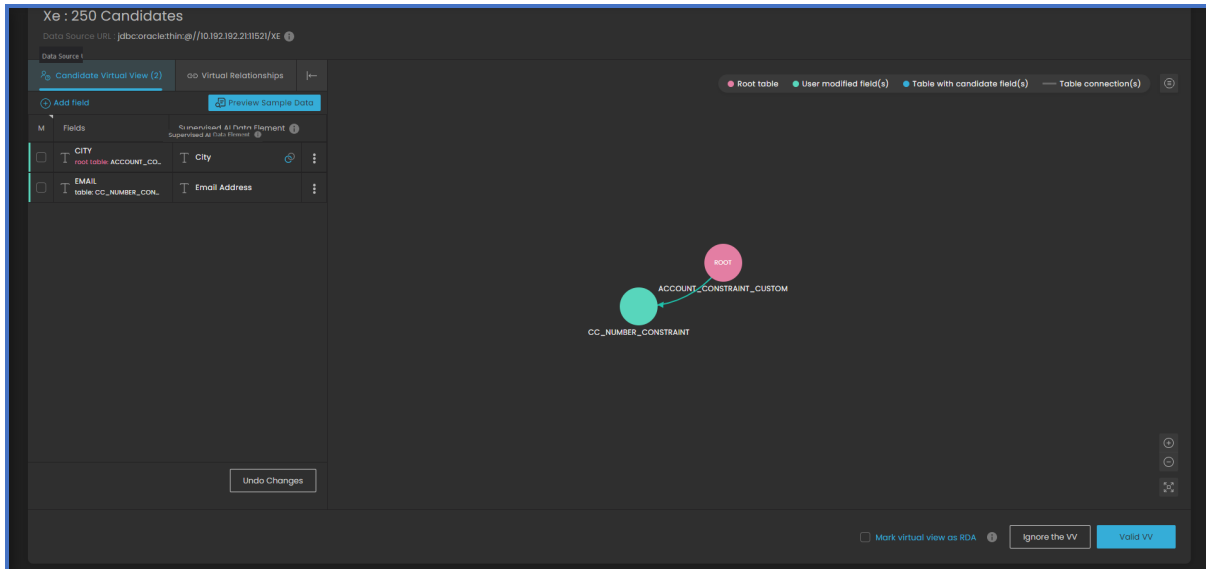


Figure 1: Candidate Virtual View tab

Preview Sample Data

If the information in the tab is insufficient to decide if the mapping is correct, click **Preview Sample Data**. You will see 30 rows with samples of PII candidates in the right part.

If the mappings have not been modified, the system will provide true data in all fields per PII candidate. After mapping is modified, sample data is not related to a specific candidate.

Manage Mappings

In the **Candidate Virtual View** tab, you can add a mapping by clicking **Add Field**. You need to select a table and field from the current virtual view (left column) and then map it to the Inventa entity type (right column).

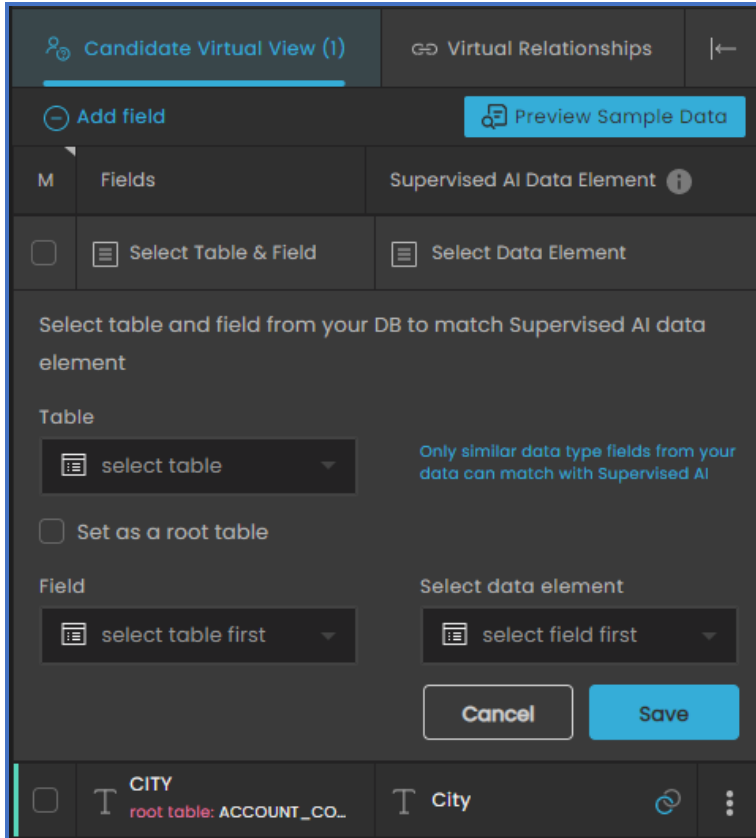


Figure 2: Adding new mappings in Candidate Virtual View tab

You can also edit and delete the mappings.

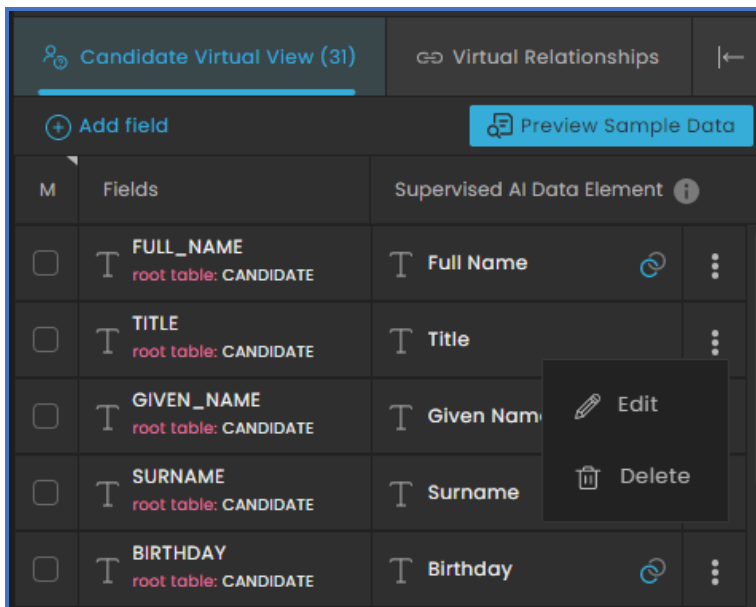


Figure 3: Mapping edition and deletion

To return to default settings, click **Undo Changes**.

Advanced Mapping Configuration

Relationships tab - advanced ERD-based mode for adding non-detected tables to the virtual view or building a virtual view from scratch.

Add non-detected tables to virtual view

If the system has not automatically detected all virtual view tables, you can add a link to any table from the repository in the Relationships tab. You can also preview sample data from both the virtual view fields and the fields not related to the virtual view.

Enter table name in Search and review all its relationships with other tables within the repository, including the links not related to the current virtual view.

You can add relationships on both the level of tables and the level of fields in tables. After adding a link, go to the **Candidate Virtual View** tab and map it to the entity type.

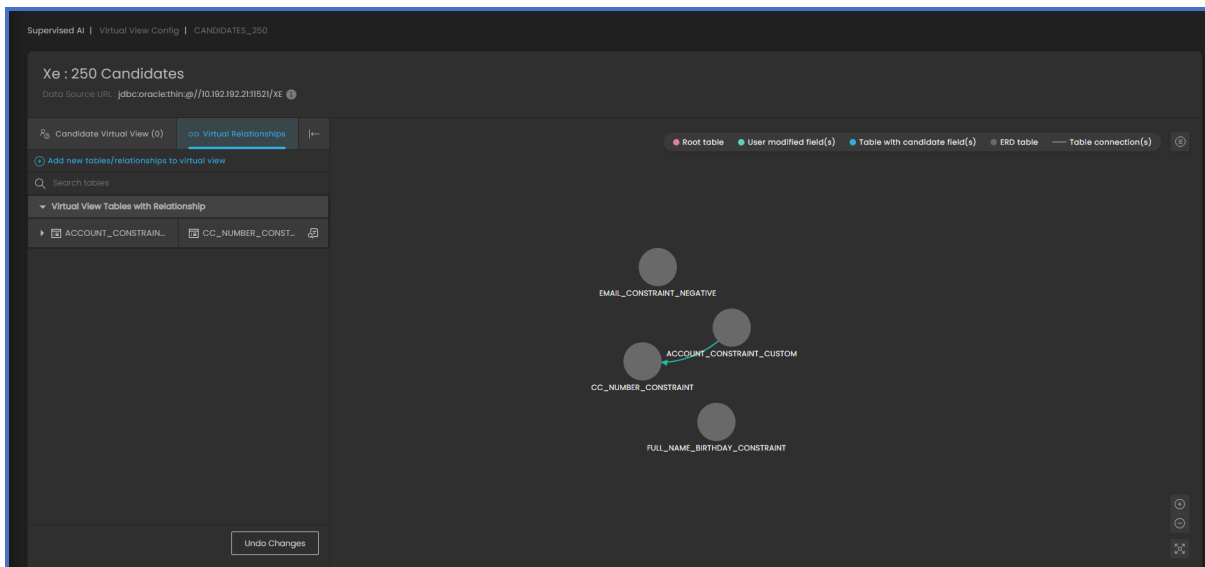


Figure 4: Relationships tab

CANDIDATES PAGE

The **Candidates** page allows you to view, filter, and export candidates data - data pertaining to instances of personal data retrieved from the data source but not confirmed against the root data asset (RDA). The page shows candidates discovered in all data sources supported by IGDC, including file shares, databases, data lakes, cloud storages, etc.

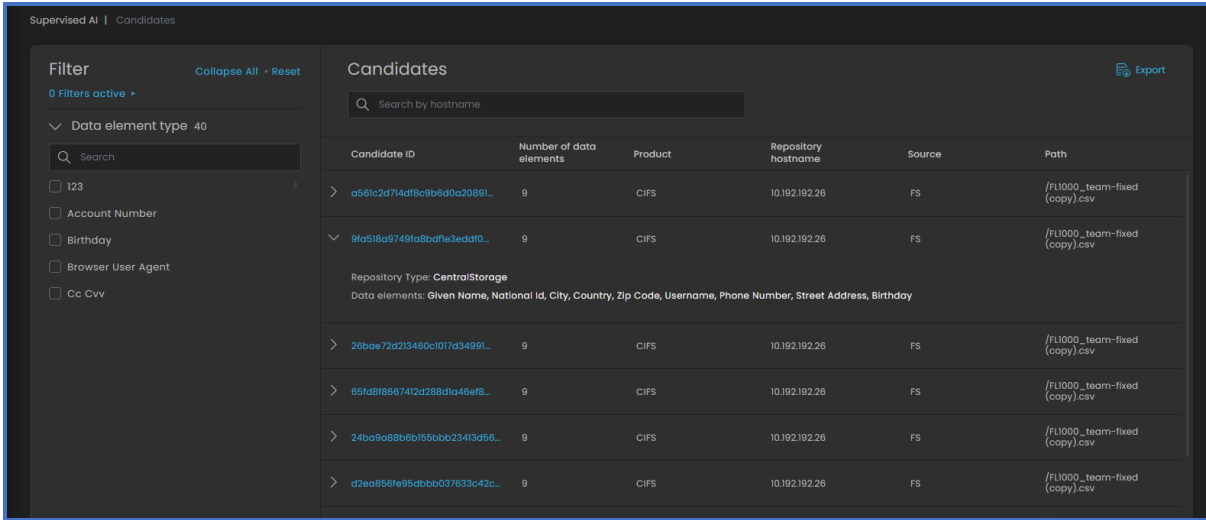


Figure 1: Candidates page

You can also view all the data discovered for a specific candidate.

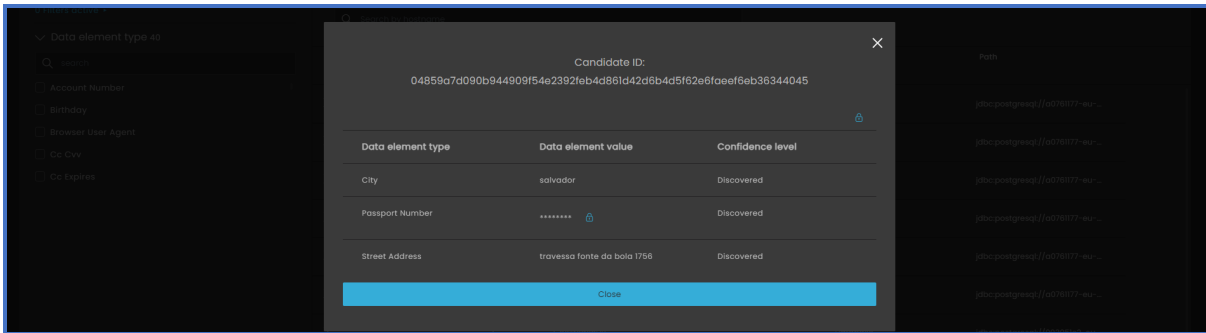


Figure 2: Candidate profile

IBM, the IBM logo, and IBM Guardium Discover and Classify are trademarks or registered trademarks of International Business Machines Corporation, in the United States and/or other countries. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on ibm.com/trademark.