# Lab Guide

# Hands-on-Lab
# Virtualizing Db2 Warehouse data with data virtualization

Shivam R Solanki
Data Scientist
Shivam.raj.solanki@ibm.com

IBM

For decades, companies have tried to break down silos by copying data from different operational systems into central data stores for analysis, such as data marts, data warehouses, and data lakes. This is often costly and prone to error. Most struggle to manage an average of 33 unique data sources, which are diverse in structure and type, and are often trapped in data silos that are hard to find and access.

With *data virtualization,* you can query data across many systems without having to copy and replicate data, which reduces costs. It also can simplify your analytics and make them more up to date and accurate because you're querying the latest data at its source.

In this tutorial, we're going to learn how to virtualize Db2® Warehouse data with data virtualization on IBM Cloud Pak for Data to make queries across multiple data sources.

## Learning objectives

In this tutorial, you will learn how to:

- [Setup the project on IBM Cloud Pak for Data](#)
- [Add datasets to IBM Cloud Pak® for Data.](#)
- [Add a data source for data virtualization.](#)
- [Virtualize the data and create a joined view.](#)
- [Assign virtualized data to a project.](#)
- [Add roles to users and perform admin tasks.](#)

# Steps

## Step 1. About the dataset

The dataset used for this tutorial contains information about fraud auto insurance claims for an insurance company. The data was collected in three CSV files.

1. ***claims.csv***

Some of the important attributes in this file are

- Claim Id

- Capital Gains *($)*
- Capital Loss *($)*
- Incident Type *(Single Vehicle collision, Vehicle Theft etc.)*
- Collison Type *(Read Collision, Side Collision etc.)*
- Incident Severity *(Minor Damage, Total Loss etc.)*
- Authorities Contacted *(Ambulance, Police etc.)*
- Incident Hour of the day
- Number of vehicles involved
- Witnesses
- Total claim amount *($)*
- Fraud reported *(Yes, No)*

2. *customer.csv*

Some of the important attributes in this file are

- Customer Id
- Insured Sex *(Male, Female)*
- Insured Occupation *(Craft repair, sales etc.)*
- Insured Hobbies *(Chess, Cross fit etc.)*

3. *policies.csv*

Some of the important attributes in this file are

- Claim Id
- Policy Id
- Coverage *(Basic, Premium etc.)*
- Policy Annual Premium *($)*
- Auto Make *(Dodge, Chevrolet etc.)*

## Step 2. Set up the project on IBM Cloud Pak for Data

**Log in to IBM Cloud Pak for Data**

1. Launch a browser and navigate to the IBM Cloud Pak for Data url.

**Create a new IBM Cloud Pak for Data project**
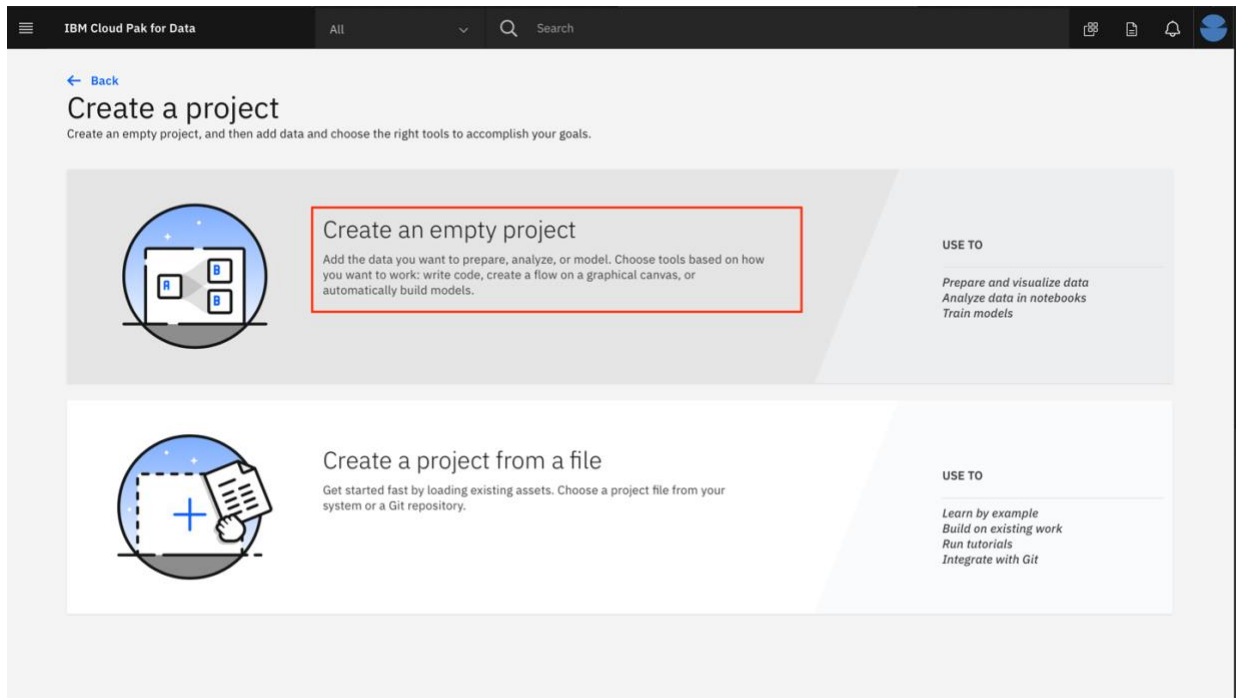
1. Go the hamburger (☰) menu and click **Projects**.

2.  Click on **New project**.



3.  Select **Create an empty project**.

IBM

4. Provide a name and optional description for the project and click **Create**.
   *(Note – For the purpose of this lab, please add your initials at the end of the project name to avoid conflict since all the participants are working in the same environment. For ex: Fraud-Claim-srs)*

## Step 3. Set up the Db2 Warehouse on IBM Cloud

We'll need a place to store our data. It is suggested to use [Db2 Warehouse on IBM Cloud](#) in order to conserve resources on the CPD cluster.

Note that IBM Cloud Pak for Data can work with any database with a JDBC connector, so Db2 warehouse is only one of many choices.

Steps below are shown for fetching connection details for Db2 Warehouse on IBM Cloud. *(Note - This step has already been done for you so you can use the credentials provided by the instructor. The screenshots shown here are for reference. You can skip to the Step 4 now)*

1. Go to Service Credentials and click **New credential +**. Click the **Copy to clipboard** icon and save the credentials for later.

2. Now go to Manage and click **Open Console**.



**Get SSL certificate for Db2 Warehouse on IBM Cloud**

You will need an SSL cert for IBM Cloud Pak for Data to use the IBM Cloud Db2 Warehouse instance.

IBM

1. In the Db2 Warehouse console, from the upper-left hamburger (☰) menu,
   click **Administration → Connections**, then **Download SSL Certificate**.



2. You will need to convert the SSL certificate from .crt to a .pem file using OpenSSL. Run
   the following command:

```bash
openssl x509 -in DigiCertGlobalRootCA.crt -out DigiCertGlobalRootCA.pem -outform PEM -inform DER
```

**Seed the Db2 Warehouse on IBM Cloud**

Steps below are shown for uploading data to Db2 Warehouse on IBM Cloud. *(Note - The 3 data files claims.csv, customer.csv and policies.csv have been uploaded to Db2 Warehouse on IBM Cloud for you so you can skip to the Step 4 now)*

IBM

1. From the upper-left hamburger (☰) menu, click **LOAD** → **Load data**.



2. Click **browse files** and select the claims.csv file after downloading it from the link, then click **Next**.

3. Choose Schema **INSURANCE** and click **+ New table**. Under Create a new Table, provide CLAIMS as the name of the table and click **Create**, then **Next**.



4. Accept the defaults and click **Next**. On the next screen, click **Begin Load**.



3. Repeat for the policies.csv file, naming the table `POLICIES` and the customers.csv table `CUSTOMERS`.

Now the Db2 warehouse has been set up on IBM Cloud.

IBM

## Step 4. Add a new data source connection

**Add the new data source**

1.  To add a new data source, go to the hamburger (☰) menu and click on
    the **Connections** option.



2.  At the overview, click **New connection +**.

**Enter connection details for Db2 Warehouse on IBM Cloud**

1.  Start by giving your new connection a name, and select **Db2 Warehouse on Cloud** as your connection type. More fields should appear. Fill in the new fields with the credentials for your Db2 Warehouse connection. Click the checkbox for **Use SSL**.

    *(Note – Please add your initials at the end of the connection name with a dash. For example, db2-warehouse-srs in order to avoid conflict)*

2.  Enter the following service credentials to create a connection to the Db2 warehouse on IBM Cloud.

    "**hostname**": "db2w-jtveemh.us-south.db2w.cloud.ibm.com",
    "**password**": "GxduaIWj_f4aZq82Ah@4@KqmHAkEF",
    "**port**": 50001,
    "**db**": "BLUDB",
    "**username**": "bluadmin",

3.  Download the SSL certificate and upload it to create connection.


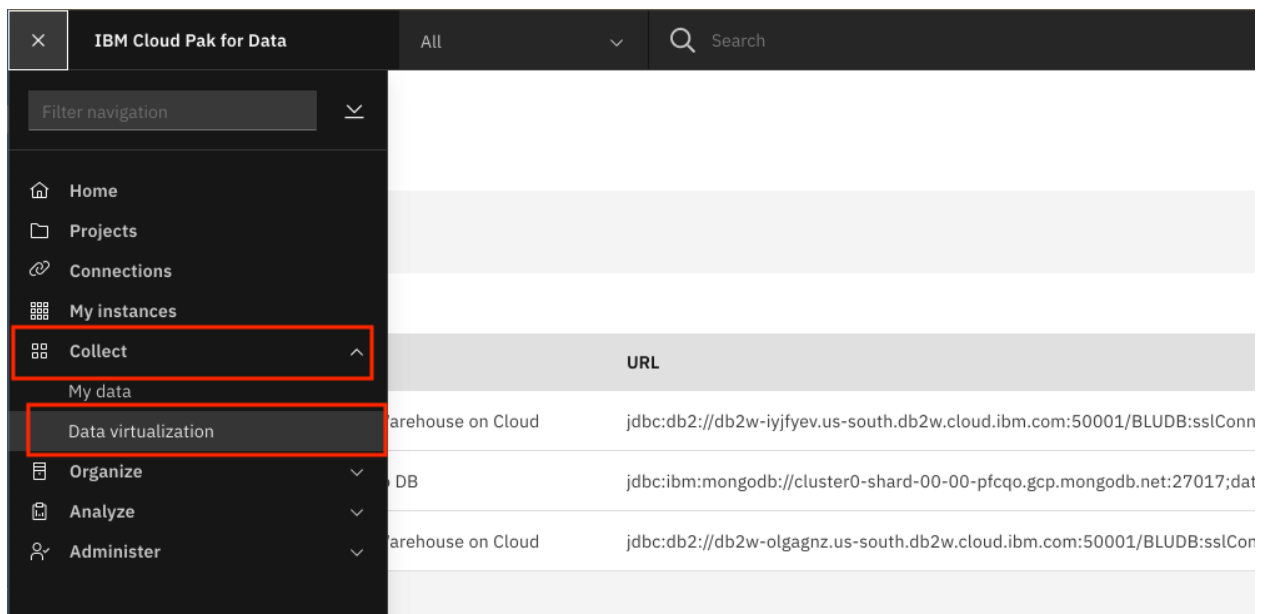
The new connection will be listed in the overview.

## Step 5. Virtualize Db2 data with data virtualization

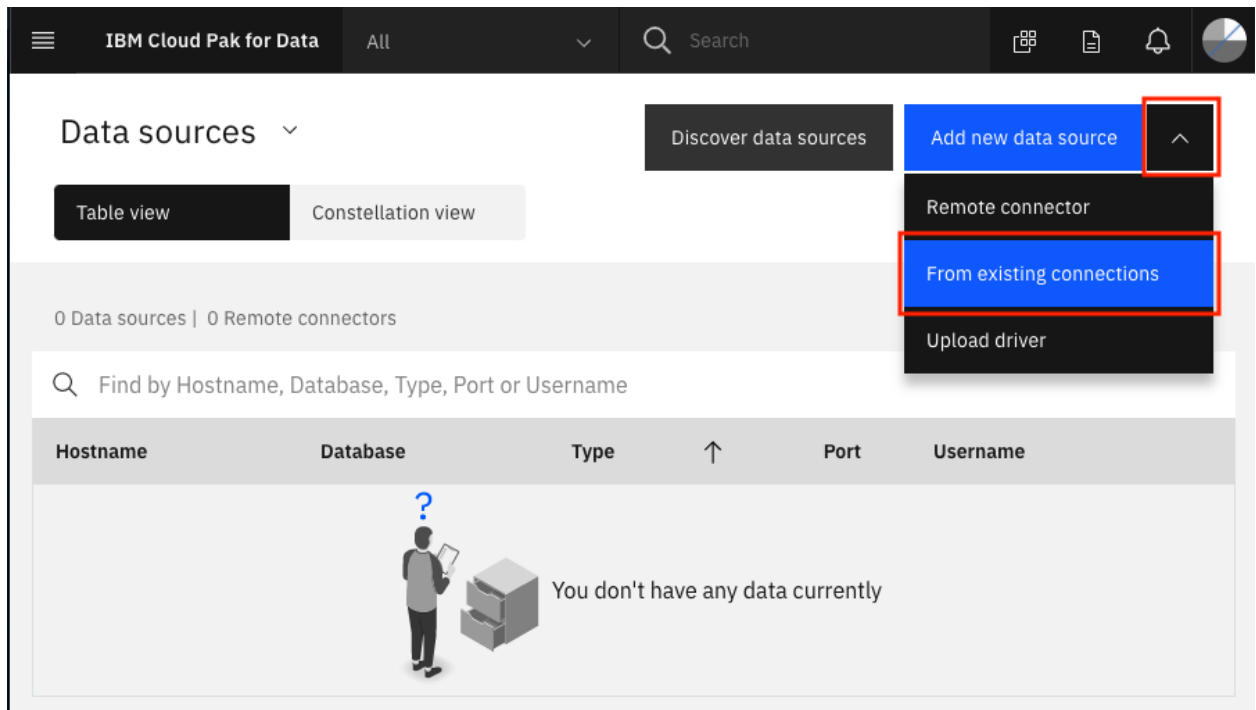**NOTE**: This section requires Admin user access to the IBM Cloud Pak for Data cluster.

For this section, we'll use the data virtualization tool to import the data from Db2 Warehouse, which is now exposed as a connection in IBM Cloud Pak for Data.
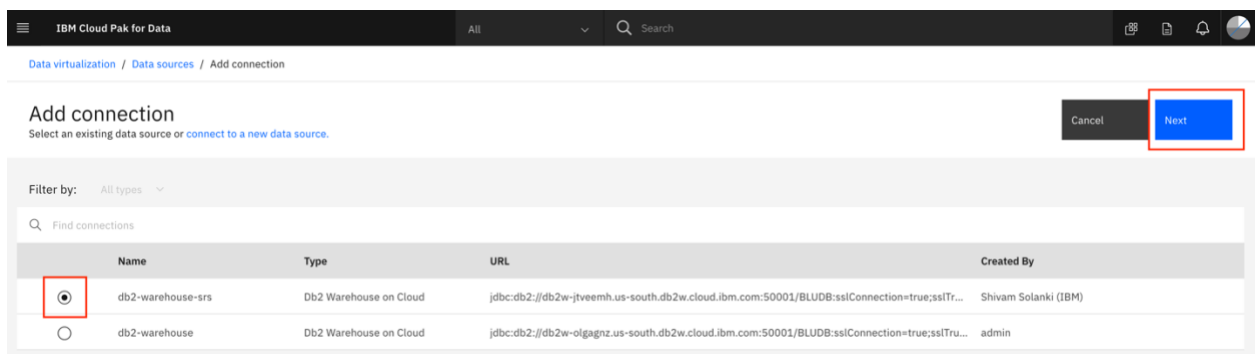
**Add a data source to Data Virtualization**

1.  To launch the data virtualization tool, go the hamburger (≡) menu and click **Collect** and then **Data Virtualization**.

2. At the empty overview, click the drop-down next to Add new data source and select **From existing connections**



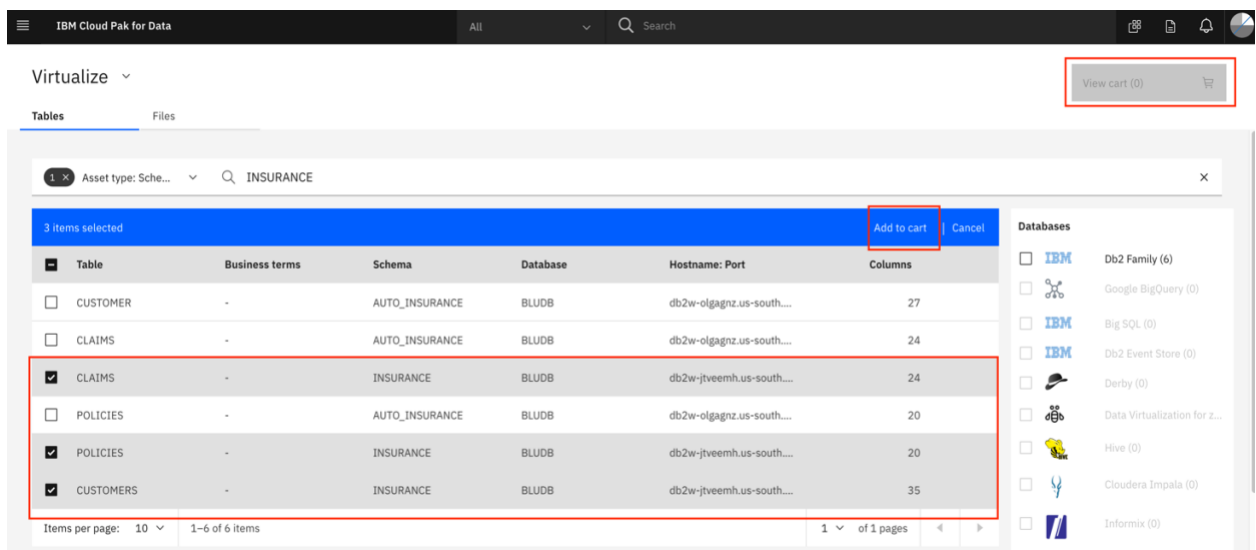3. Select the data source we made in the previous step and click **Next**.



**Start virtualizing data**

Because we now have access to the Db2 Warehouse data, we can virtualize the data to our IBM Cloud Pak for Data project.

1. Click on the **Data Sources** drop-down and choose **Virtualize**.



2. Several tables will appear (many are created as sample data when a Db2 Warehouse instance is provisioned) in the table. Find the tables that was created earlier, the previous instructions suggested naming them CUSTOMERS, POLICIES, and CLAIMS. Or you can search them by filtering on Schema = INSURANCE. Once selected, click **Add to cart** and then **View Cart**.



3. The next panel prompts the user to choose which project to assign the data to. Choose **My virtualized data** and *uncheck* the box that says Submit to catalog. Click **Virtualize** to start the process.

4. You'll be notified that the virtual tables have been created. Let's see the new virtualized data from the data virtualization tool by clicking **View my virtualized data**.



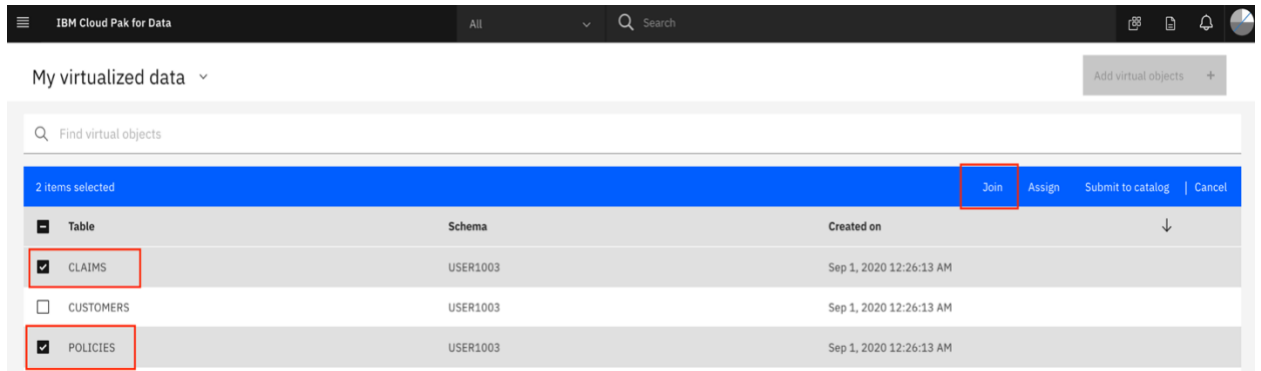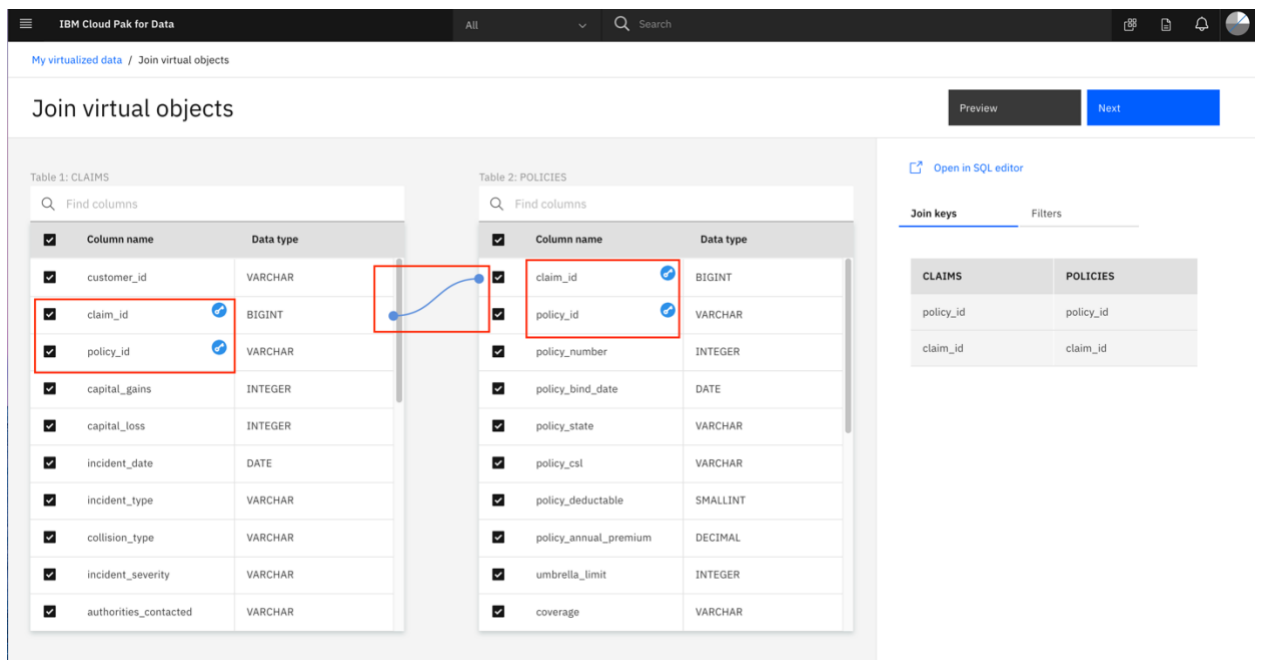**Join the virtualized data**

Now we are going to *join* the tables we created so we have a merged set of data. It will be easier to do it here rather than in a notebook where we would have to write code to handle three different datasets.

1. Click on any two tables (CLAIMS and POLICIES, for instance), then click the **Join** button.

IBM

2. To join the tables, we need to pick a key that is common to both datasets. Here we choose to map *claim_id* and *policy_id* from the first table to *claim_id* and *policy_id* on the second table. Do this by clicking one and dragging it to the other. When the line is drawn, click **Next**.



3. Next, you have a chance to edit column names, but we will keep them as-is. Click **Next**.

4. In the next panel, we'll give our joined data a unique name such as POLICIESCLAIMS (to be consistent with SQL standards, pick an uppercase name). Under **Assign to**, choose **My virtualized data** and uncheck the box that says Submit to catalog. Click **Create view** to start the process.
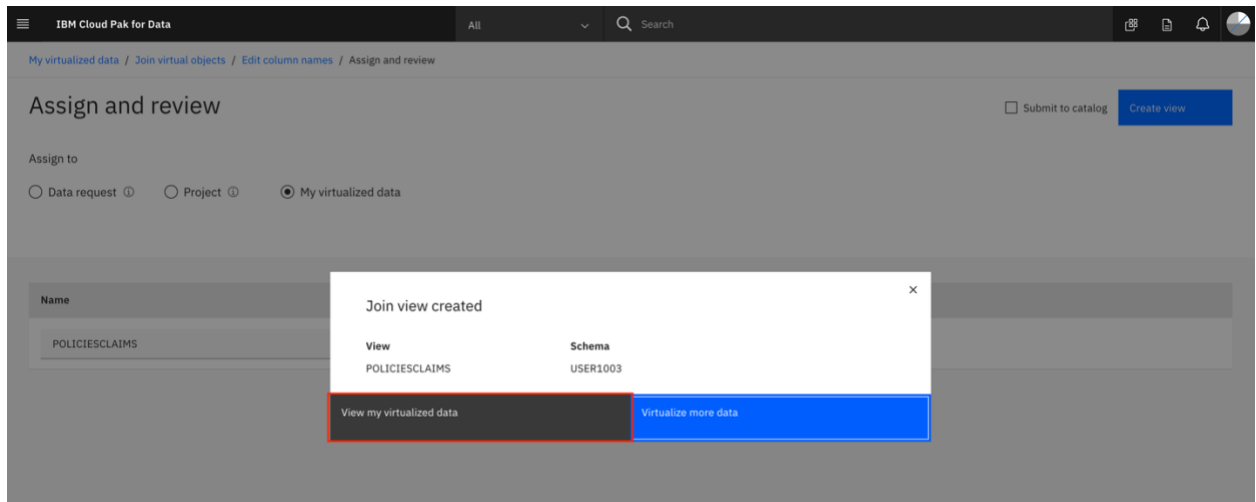
5.  You'll be notified that the join has succeeded. Click on **View my virtualized data** to go back and see all your virtualized data.



6.  **IMPORTANT** Now join the new joined view (POLICIESCLAIMS) and the last virtualized table (CUSTOMERS) to create a new joined view that has all three tables; let's call it POLICIESCLAIMSCUSTOMERS. Switching back to the **My virtualized data** screen should show all three virtualized tables and two joined tables. Do not go to the next section until this step is performed.



**Grant access to the virtualized data**

For other users to have access to the data you just virtualized, you need to grant it. Follow these steps to make your virtualized data visible to them:

1.  Go to **Data Virtualization** from the hamburger (☰) menu. Click on **Menu → My virtualized data**.
2.  Click on the virtualized data you've created, then click the three vertical dots to the right, and choose **Manage access.**

3. Click the **Specific users** radio button, then **Add user +**.

4. Select the users you wish to grant access to and click **Add users**.



Repeat the above steps for the remaining tables and views.

**Assign the Engineer role to the users**

IBM Cloud Pak for Data users that need to use data virtualization functions must be assigned specific roles based on their job descriptions. These roles are Admin, Engineer, User, and Steward. You can learn more about these roles on the IBM Cloud Pak for Data product hub.

Let's assign the Engineer role to some users:

IBM

1. From the hamburger (☰) menu, choose the **Data Virtualization** option, then click **My virtualized data → User management**.



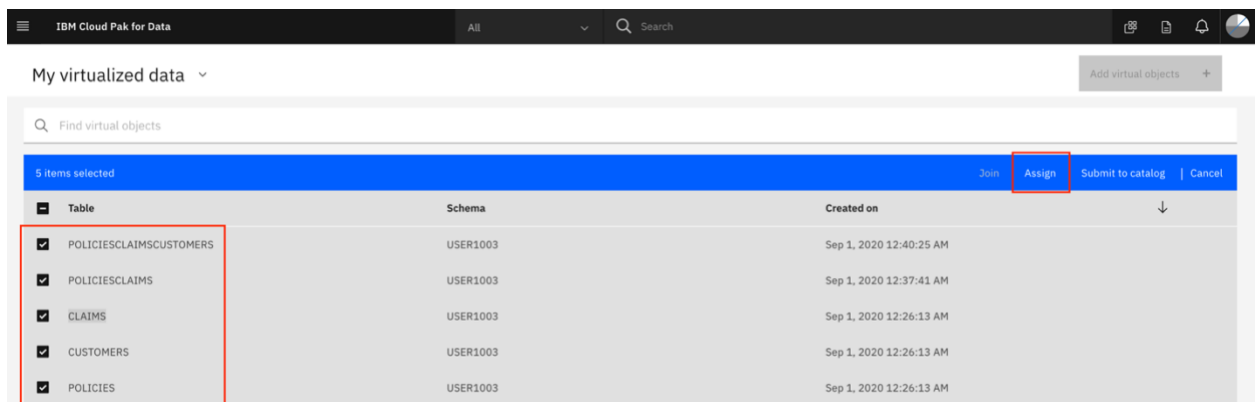2. Click on **Add users +** and update the role of the user to Engineer.
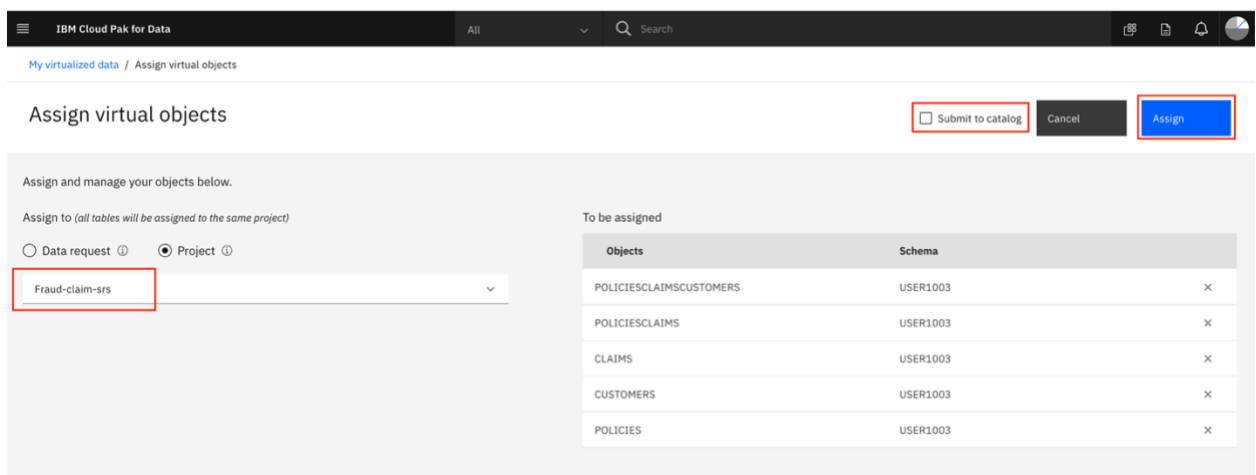
## Step 6. Users assign virtualized data

Now let's look at how a user who has access to virtualized data can assign the data to their project — how to add the virtualized data as asset to a project.
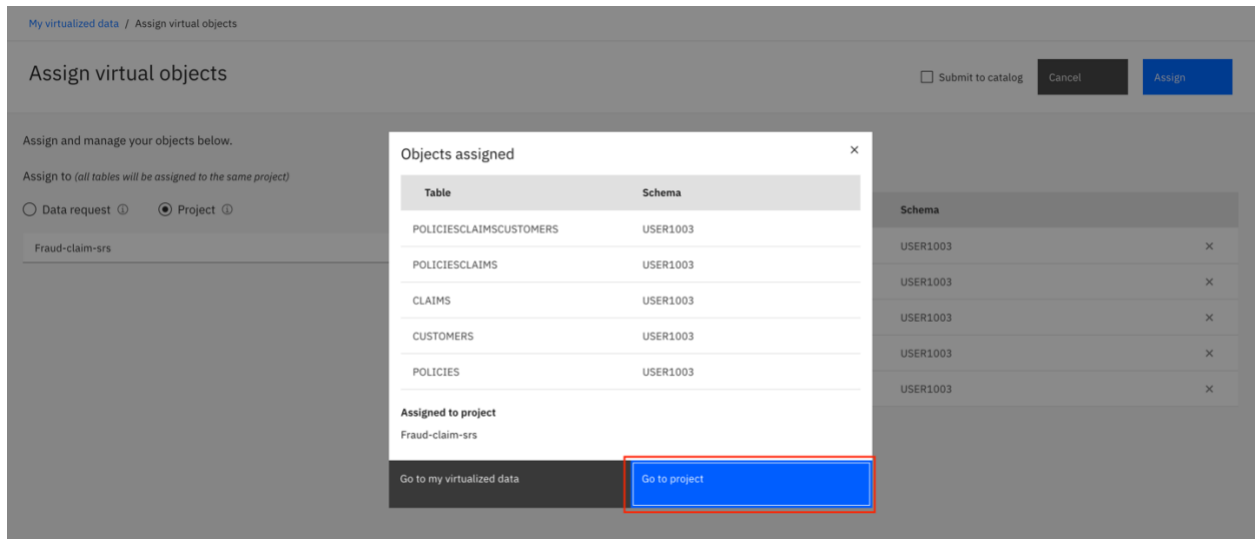
**Assign the data to your project**

1. From the hamburger (☰) menu, click on **Collect → Data Virtualization**. You will be brought to the My virtualized data section. Here you should see the data you can access (or that the administrator has assigned to you). Select the checkbox next to our original tables (CLAIMS, POLICIES, CUSTOMERS) and the joined tables (POLICIESCLAIMS, POLICIESCLAIMS CUSTOMERS), and click the **Assign** button to import them into your project.
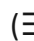
2. On the Assign virtual objects screen, choose the project to assign the data. If there is a **Submit to catalog** checkbox on the top right, uncheck it and click the **Assign** button to add the data to your project.
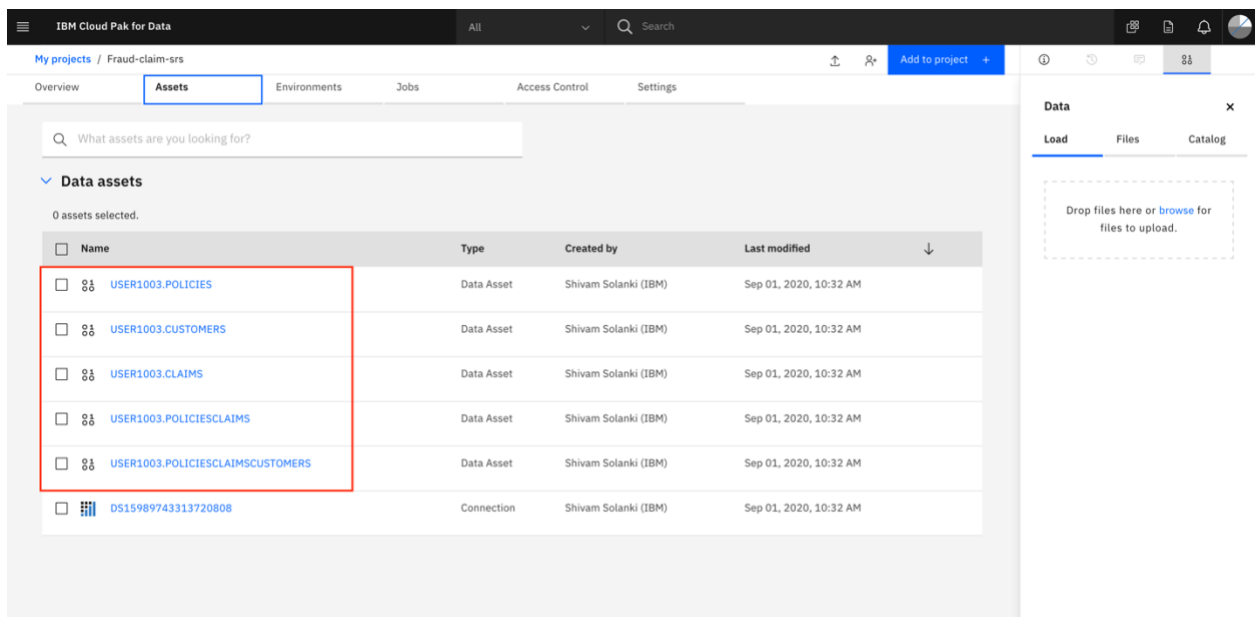
3. In the pop-up panel, you will receive a confirmation that the objects have been assigned to your project. Click the **Go to project** button.



Alternatively, close the modal and go to your projects by clicking on the hamburger (☰) menu, then choosing **Projects**.

4. On the project page, click on the **Assets** tab to display the virtualized tables and joined tables that are now in your project.

**Summary**

This lab tutorial explained how to virtualize Db2 Warehouse data with data virtualization on IBM Cloud Pak for Data to make queries across multiple data sources. To continue the series and learn more about IBM Cloud Pak for Data, take a look at the next tutorial, Data visualization with data refinery.