

# IBM Journey to Cloud and AI

## Analytics Modernization Workshop

*Workbook*



Author:

Vikram S Khatri  
Executive I/T Specialist  
[vikram.khatri@us.ibm.com](mailto:vikram.khatri@us.ibm.com)

Dated:10-20-2018: Last Updated: 11-17-2018

## Acknowledgements

- **Duane Almeter** and **Craig Maddux** for their exemplary leadership to bring best of the technical experts together to help create this multidiscipline asset and bring multiple technical sellers for them to deliver this workshop in their respective markets. Without your leadership, this would not have been possible – Thank you.
- **Anjali Shah, Rui Fan, Daniel Kikuchi** and **Ben Chard** for bringing their data science expertise to this asset and creating the AI-based predictive modeling for this asset.
- **John Van Buren** for his expertise in Catalog, Governance and Integration to enrich this asset.
- **Kent Rubin** for his expertise in business dashboards.
- **Jeff Tuck** for writing Predictive Modeling microservice to enrich the stock trader application written by **John Alcorn** and **Greg Hintermeister** under the guidance of **Eric Herness**, IBM Fellow and **Bratt Coffmann** for his UI assistance.
- **Sriram Srinivasan**, Chief Architect of the IBM Cloud Private for Data and his team – **Ashwin Dev, Prashant Patel, Sahil Shah, Parth Komperla** from the IBM development labs to provide technical details and help.
- **Dean Compher, Dominic Farrar, Gary Brunell, Karen Groski, Kyle Talish, Neal Finkelstein, Patrick Pitre, Paul Betts, Ryan Kather, S Suh, Dale Mumper** and **David Solomon** for active participation in testing the asset and for taking the role of proctor to facilitate delivery of this workshop.
- **Claus Huempel** from IBM Germany to learn and deliver this asset in Europe.
- **Eric Watson** for helping to refine the content.
- **Kristen Mucci** for planning and organizing the events throughout North America.

Thank you - **Beth Friday**, VP, North America Technical sales for being our Technical Executive Sponsor and **Hemanth Manda**, Director, Platform Offerings and **Sampada Basarkar**, Director of Development for providing help from the development teams.

---

# Contents

<b>LAB 01</b>	<b>SETUP AND INTRODUCTION .....</b>	<b>6</b>
1.1	IBM JOURNEY TO CLOUD AND AI: ANALYTICS MODERNIZATION WORKSHOP .....	6
1.2	WHAT IS IBM CLOUD PRIVATE? .....	6
1.3	WHAT IS IBM CLOUD PRIVATE FOR DATA? .....	6
1.4	AUDIENCE FOR THIS IBM WORKSHOP .....	6
1.5	LAB ENVIRONMENT .....	7
1.6	LET'S GET STARTED! .....	7
1.7	EXPLORE IBM CLOUD PRIVATE .....	10
1.8	EXPLORE HADOOP PLATFORM .....	11
1.9	EXPLORE HOME PAGE .....	13
1.10	USER MANAGEMENT: PERSONA-BASED ROLES AND TEAMS .....	13
1.11	INTEGRATION WITH SOURCE CONTROL .....	14
1.12	WHY MAKE IBM CLOUD PRIVATE FOR DATA YOUR PLATFORM FOR DATA .....	16
<b>LAB 02</b>	<b>EXECUTIVE DEMO .....</b>	<b>17</b>
2.1	IMPORT PROJECT .....	17
2.2	CONNECT TO A DATA SOURCE .....	18
2.3	CHECK RUNNING ENVIRONMENT .....	22
2.4	STOCK TRADER OPENING BELL ANALYSIS DASHBOARD .....	22
2.5	STOCK TRADER ANALYSIS CLOSING BELL ANALYSIS .....	24
2.6	RUN STOCK-TRADER BEFORE APPLICATION .....	25
2.7	RUN STOCK-TRADER AFTER APPLICATION .....	29
<b>LAB 03</b>	<b>COLLECT .....</b>	<b>36</b>
3.1	MULTIPLE SOURCES .....	36
3.2	INTEGRATION WITH THE HADOOP PLATFORM .....	36
3.3	DATA VIRTUALIZATION FROM HADOOP PLATFORM .....	37
3.4	VIRTUALIZATION OF Db2 IN ICP-D .....	39
3.5	VIRTUALIZATION OF ORACLE IN ICP-D .....	40
3.6	TEST DATA VIRTUALIZATION .....	41
<b>LAB 04</b>	<b>ORGANIZE .....</b>	<b>43</b>
4.1	CREATE CONNECTIONS .....	43
4.2	CREATE BUSINESS GLOSSARY .....	45
4.3	CREATE GOVERNANCE POLICIES AND RULES .....	48
4.4	DISCOVER ASSETS .....	50
4.5	SHOPPING FOR DATA .....	56
4.6	TRANSFORM DATA .....	57
<b>LAB 05</b>	<b>ANALYZE .....</b>	<b>70</b>
5.1	BUSINESS ANALYST .....	70
5.2	DATA SCIENTISTS – APPROACH TO A SOLUTION .....	87
<b>LAB 06</b>	<b>DEPLOY .....</b>	<b>102</b>
6.1	INTEGRATION WITH SOURCE CONTROL .....	102
6.2	DEPLOY THE MODEL .....	103
6.3	TEST DEPLOYMENT .....	107
6.4	CONSUME DEPLOYMENT IN STOCK TRADER APPLICATION .....	109
6.5	CONCLUSION .....	113
<b>APPENDIX A.</b>	<b>STOCK – OPENING BELL DASHBOARD .....</b>	<b>114</b>
8.1	BUILD DASHBOARD .....	118
<b>APPENDIX B.</b>	<b>INTEGRATED DATA PLATFORM .....</b>	<b>124</b>
8.2	APPLICATION WITH ANALYTICS MODERNIZATION .....	125
8.3	COLLECT .....	126
8.4	ORGANIZE .....	127
8.5	ANALYZE .....	129
8.6	DEPLOY .....	130
8.7	BUSINESS IMPACT .....	131
<b>APPENDIX C.</b>	<b>INTEGRATION OF HADOOP CLUSTER WITH ICP-D .....</b>	<b>132</b>
8.8	ACCESS TO HADOOP BEHIND FIREWALL .....	132
8.9	DOWNLOAD SOFTWARE .....	132

8.10	PREPARE HADOOP CLUSTER FOR ICP-D INTEGRATION .....	132
8.11	CREATE <b>DSXHI</b> USER .....	134
8.12	INSTALL SOFTWARE IN THE DOCKER CONTAINER .....	134
8.13	CREATE <b>DSXHI</b> CONFIGURATION FILE.....	135
8.14	INSTALL <b>DSXHI</b> .....	135
8.15	CREATE <b>DSXHI</b> SYSTEMD SERVICE .....	136
8.16	COMMIT DOCKER CONTAINER TO SAVE THE INSTALLATION.....	137
8.17	TROUBLESHOOTING .....	137
	<b>APPENDIX AA. NOTICES .....</b>	<b>139</b>
	<b>APPENDIX BB. TRADEMARKS AND COPYRIGHTS .....</b>	<b>141</b>

THIS PAGE INTENTIONALLY LEFT BLANK

---

## Lab 01 Setup and Introduction

### 1.1 IBM Journey to Cloud and AI: Analytics Modernization Workshop

This session provides hands-on experience with modernizing your microservices applications by enriching them with machine learning and artificial intelligence. The Journey to Artificial Intelligence requires a strong information architecture that supports self-service capabilities and balances the needs of both the agility required by lines of business and the ‘Enterprise class’ delivery required by IT. This journey can move significantly faster if you use an integrated platform to: **Collect, Organize** and **Analyze** your data. Let’s start this journey by defining some key concepts.

### 1.2 What is IBM Cloud Private?

IBM Cloud Private is a cloud-computing model run solely for one organization. It can be managed internally or by a third party. It can be hosted on premises behind your company’s firewall or externally on a third-party provider’s cloud. Private cloud offers the benefits of a public cloud including: rapid deployment and scalability, ease of use, elasticity. It also provides: greater control, increased performance, predictable costs, tighter security and flexible management options. It can be customized to meet your unique needs and security requirements.

### 1.3 What is IBM Cloud Private for Data?

IBM Cloud Private for Data (ICP-D), built on the foundation of IBM Cloud Private, is an integrated end-to-end platform designed to help make data more accessible and trusted. It provides access to many analytical tools to help you gain insights from your data.

ICP-D provides the data platform which enables you to climb the ladder to AI faster. With it, you can quickly build, train and deploy machine learning (ML) and artificial intelligence (AI) models. You can also: ‘Inventory and catalog’ your data sources, provide ‘self-service’ shopping for data, and provide data integration and refinement capabilities. High quality, trusted data can be more easily prepared and assembled through this integrated platform.

### 1.4 Audience for this IBM Workshop

This IBM workshop is mainly designed for line-of-business professionals who are tasked to gain new insights from all available data – regardless of its type and origin. The following personas will greatly benefit from this workshop.

Role	Capabilities
Business analyst	Business analyst delivers value by taking data, using it to answer questions, and communicating the results to help make better business decisions.
Data scientist	Data scientist brings expertise in statistics and the process of building ML/AI models to make predictions and answer key business questions.

Data steward	Data steward brings integration and governance to the data.
Data engineer	Data engineer builds and optimizes the systems to allow data scientists and analysts to perform their work. The data engineer ensures that any data is properly received, transformed, stored, and made accessible to other users.

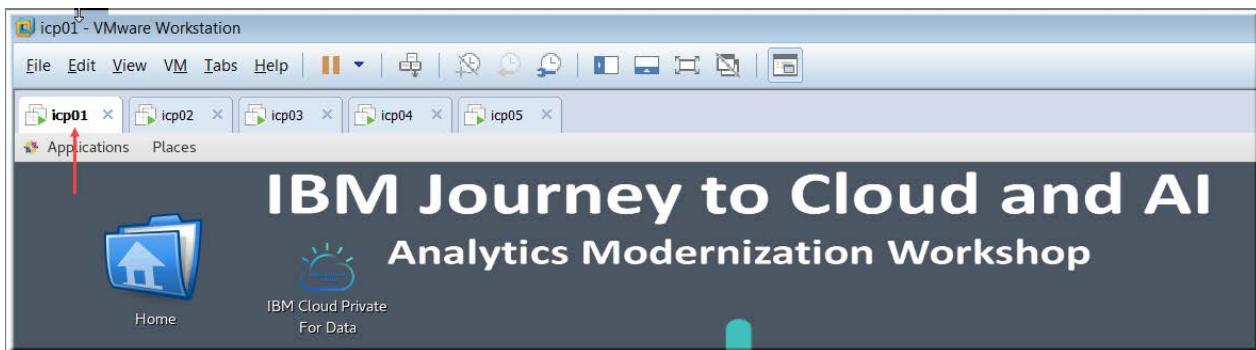
## 1.5 Lab environment

We are using eleven Virtual Machines (VMs) for this workshop. This software environment was built to the following specifications:

- IBM Cloud Private 3.1 as foundational technology
- IBM Cloud Private for Data 1.1.0.2 as an Integrated Data Platform

## 1.6 Let's get started!

- 1. On your laptop/workstation/Cloud VM, multiple virtual machines are running. If you are not already in [icp01](#) VM, click the [icp01](#) tab to switch to the first virtual machine.



- 2. Double-click [IBM Cloud Private for Data](#) from the desktop to open the web browser to start the IBM Cloud Private for Data web console. Maximize your browser window.



- 3. The web UI console displays as shown. Please type **admin** and **password** for the *Username* and *Password* and click **Submit**.

The screenshot shows the sign-in interface for IBM Cloud Private for Data. At the top, there are 'Sign in' and 'Sign up' buttons. Below them is a 'Username' field containing 'admin'. Underneath is a 'Password' field with several red asterisks. A red arrow points to the password field. To the right of the password field, the text 'Password is password' is displayed. At the bottom is a 'Submit' button.

- 4. Review the **Let's get started!**

The screenshot shows the 'Let's get started!' page. It includes a pie chart, a 'WELCOME, ADMIN!' message, a large 'Let's get started!' heading, and a video player for an 'IBM Cloud Private for Data: Overview' video. Below the video are links for 'Collect and organize' and 'Analyze'.

- 5. Even if you only spend a few minutes on the **Let's get started!** page, be sure look at the following:

- a. One menu that provides capabilities to **Collect**, **Organize** and **Analyze** data using a single end-to-end platform.



- b. **Collect** – You can collect data from a myriad of sources, no matter where it resides. For example, you can collect data from IBM Db2®, Oracle, Teradata, Hadoop, flat files and so on. It does not matter if data is stored in a cloud or on-premises.

- \_\_\_c. **Organize** – Use embedded machine learning to automate the process of assigning business terms to newly discovered data sources. Enforce governance policies as you search, find and catalog your data. Create an enterprise-level data catalog through which you can search for the data you need. Collect and prepare the data for analysis that you need to perform. Use embedded security to provide visibility to all data and still control who accesses it.
  - \_\_\_d. **Analyze** – Build predictive and prescriptive models using open source programming languages and/or simpler graphical user interfaces. With the click of a button, deploy these models into production and publish an Open REST API that can be consumed by any application which can be running either on IBM Cloud Private or any other platform. Build dashboards and interact with your data to gain insight and business understanding.
  - \_\_\_e. The key personas' roles are divided into different categories. If you are a data engineer / data steward / data analyst, you will focus areas on the **Collect** and **Organize** category.
  - \_\_\_f. The data scientists and business analysts will use **Analyze** category. All personas will collaborate using an integrated single platform.
  - \_\_\_g. The **Administer** category provides tools to organize teams, define roles and provide access and control.
- \_\_\_6. IBM Cloud Private for Data is most useful in these four scenarios.



- \_\_\_a. **Manage All Your Data**: Use Discovery to automate the process of assigning business terms to technical assets. Use enterprise catalog to secure, govern and control access to your data regardless where it resides. Use a combination of virtualization and transformation to prepare your data for analysis.
- \_\_\_b. **Build Your Ladder to AI**: Build a strong information architecture to help you realize the value of leveraging ML/AI.
- \_\_\_c. **Modernize Your Data and Analytics Workloads**: Modernize your data platform and provide data scientists and application developers the ability to quickly add AI to your applications. ICP is the foundation of ICP-D. Use ICP to build cloud-native, microservice-based applications and use ICP-D to enrich them with machine learning and artificial intelligence.
- \_\_\_d. **Compliance Readiness** – Our strong governance capabilities allow you to provide regulatory compliance.

We will be demonstrating all of the above capabilities in this workshop from Collect, Organize and Analyze to demonstrate AI-driven model development and consumption of those APIs into a microservice application.

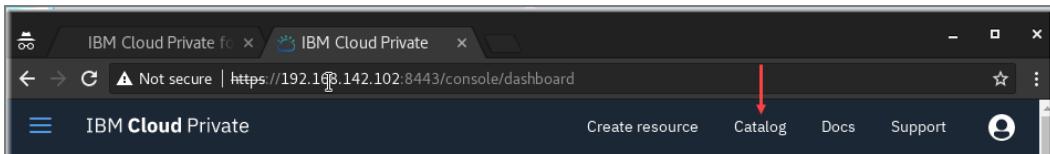
## 1.7 Explore IBM Cloud Private

- \_\_7. ICP-D foundation is IBM Cloud Private. On your [icp01](#) VM's desktop, double-click **IBM Cloud Private Web UI**.



- \_\_8. Log in using [admin](#) and [admin](#).

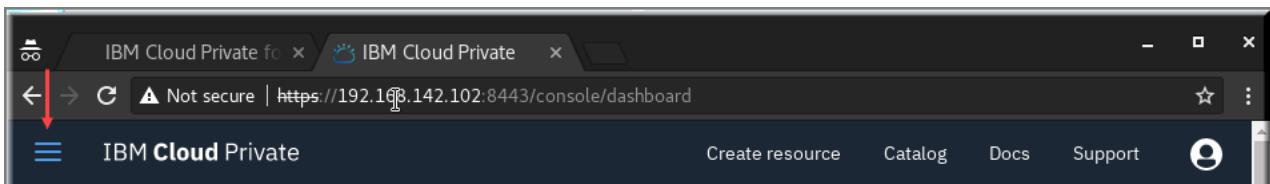
- \_\_9. Click [Catalog](#).



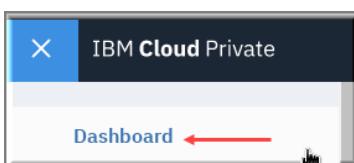
- \_\_10. This is where you can find IBM and Open Source software that you can install in the IBM Cloud Private platform. These are IBM-vetted software products that can be installed in the IBM Cloud Private platform.

**IBM is bringing power of cloud platform within the bounds of your firewall.**

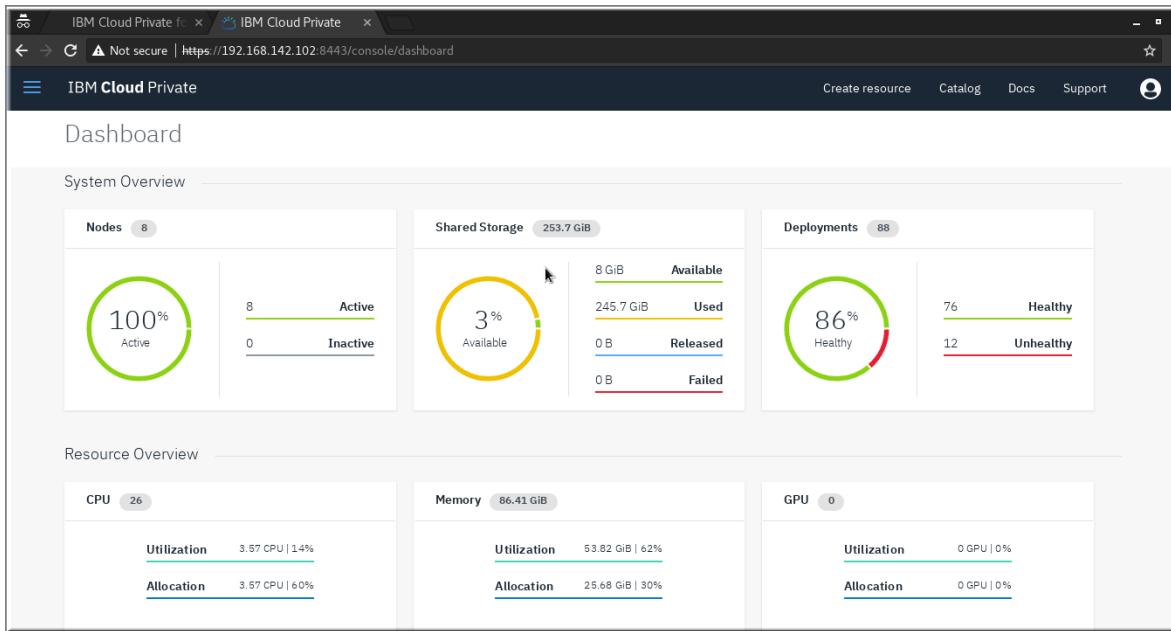
- \_\_11. Click the three horizontal bar symbol (or hamburger).



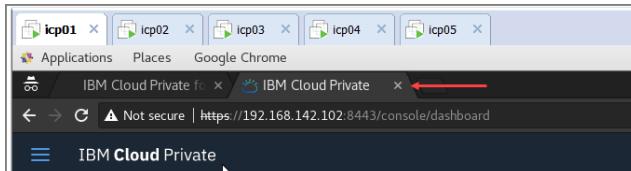
- \_\_12. Click [Dashboard](#).



- 13. The [System Overview](#) dashboard gives you a picture of the overall health of the IBM Cloud Private cluster through single pane of glass.



- 14. Close the [IBM Cloud Private](#) tab.

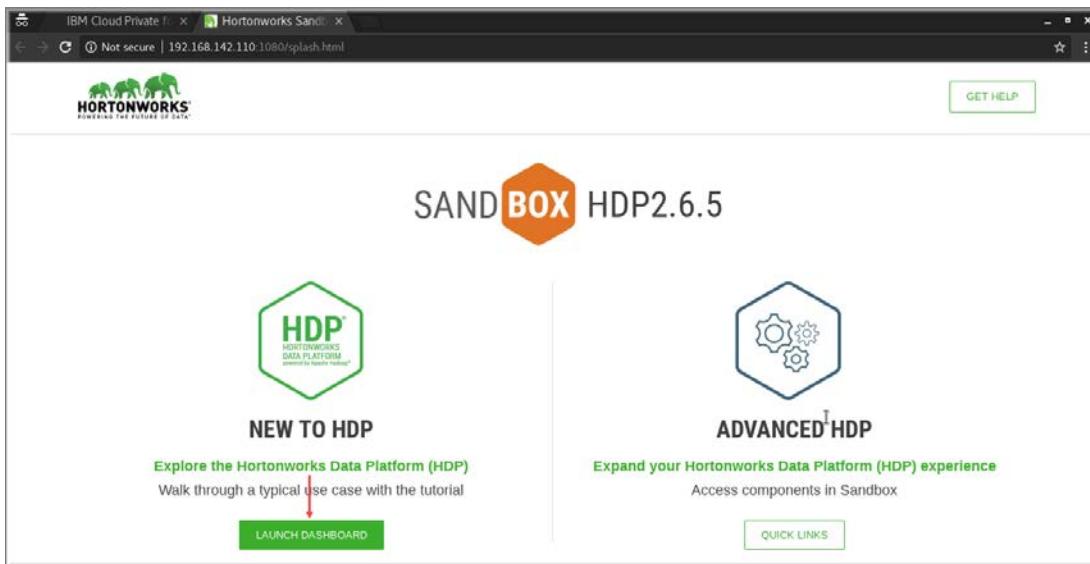


## 1.8 Explore Hadoop Platform

- 15. In our journey of the integrated data platform, we will show [Collect](#) – which is about collection of data regardless where it resides.
- 16. On your [icp01](#) VM's desktop, double-click [Hadoop Web Console](#) to launch Hortonworks Hadoop Data Platform.



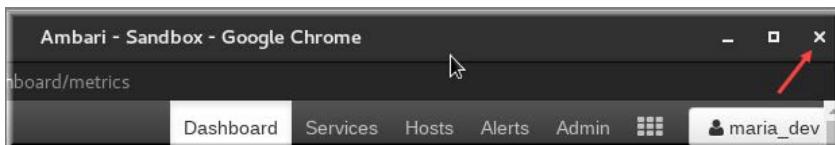
\_17. Click [Launch Dashboard](#)



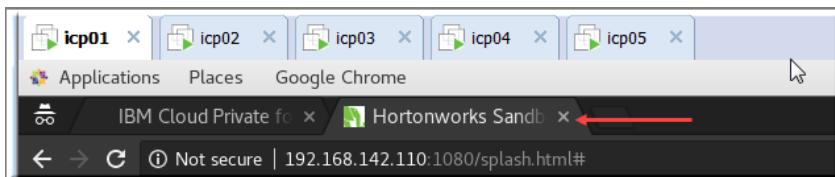
\_18. Type in Username: [admin](#) and Password: [password](#). Click [Sign in](#).

\_19. The Apache Ambari dashboard opens in a new window. We will return to this in later lab exercises.

\_20. Click the [Ambari](#) browser window.



\_21. Close the [Hortonworks Sandbox](#) tab.



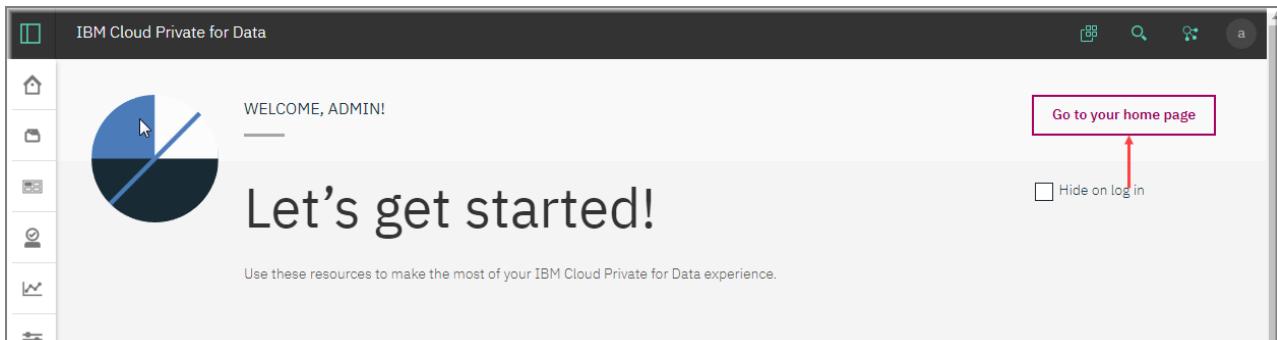
\_22. ICP-D allows integration of Hadoop edge node (or gateway) to transfer data to and from the Hadoop cluster. Livy integration also provides the ability to execute jobs remotely on the Hadoop cluster. This can be done using [Administer](#) ⇒ [Hadoop Integration](#). Note: We will use the Hadoop integration later lab exercises.

	<b>Note:</b> By design, Hadoop is not running (nor is it intended to run) on the ICP-D cluster. Many companies use Hadoop as an enterprise data repository. Our approach with ICP-D is to leave the data in Hadoop, and access it as needed for specific projects.
--	--

\_23. Switch back to the ICP-D web console.

## 1.9 Explore Home Page

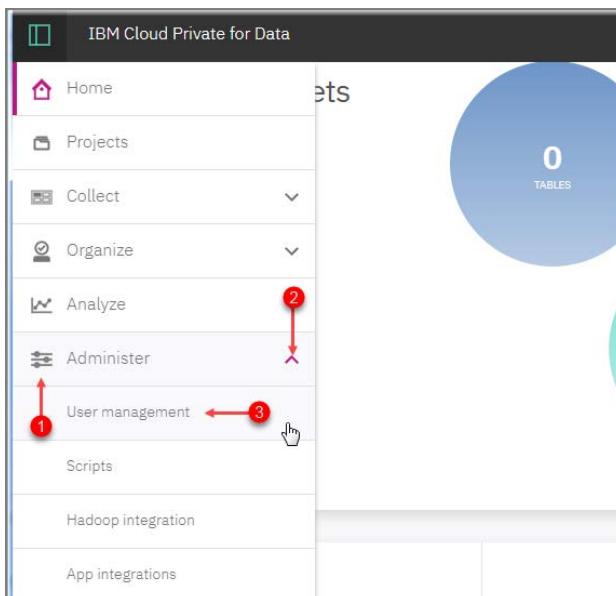
- \_24. Click [Go to your home page](#).



- \_25. The home page shows the inventory of the integrated data platform for Governance, Terms, Policies, Assets, Governed Assets, Asset Distribution and Transformation jobs. The intent of this page is to provide individual customizations on daily tasks for different personas.
- \_26. Scroll through the page to see inventory and asset types. Since we are starting fresh, we will visit the home page later to see the assets and various lists.

## 1.10 User Management: Persona-based Roles and Teams

- \_27. This section explores data governance and user authorizations for the various stages of the data analytics pipeline.
- \_28. Hover your mouse on the navigation bar on the left side of the page to pull down the menu and click the down arrow on [Administer](#) to bring the drop-down menu and click [User management](#).



- \_\_29. ICP-D can interact with your LDAP for the user management.

NAME	STATUS	USERNAME	DATE ADDED	USER ID
admin	Approved	admin	--	999

- \_\_30. Click **Roles** and notice Business Analyst, Data Engineer, Data Scientist, and Data Steward roles. These roles are customizable for different types of permissions such as Catalog, Policy, Provision Database, Manage Catalog, Manage Governance Policies and Virtualize and Transform data.

ROLE	DESCRIPTION	LAST UPDATED	ENABLED PERMISSIONS
Administrator	Administrator role		Administrator
Business Analyst	Business analyst role		Catalog Access
Data Engineer	Data engineer role		Catalog Access + 2 more
Data Scientist	Data scientist role		Catalog Access
Data Steward	Data steward role		Manage Catalog + 3 more

## 1.11 Integration with Source Control

- \_\_31. Click . This is located towards the right on the menu bar and click **Settings**.

- Signed in as: admin
- Getting Started
- Settings
- Sign Out

- \_\_32. Click [Git Integrations](#) and click [Add Token](#).

The screenshot shows a user interface for managing access tokens. At the top, there are tabs for 'Profile', 'Permissions', and 'Git Integrations'. The 'Git Integrations' tab is selected and highlighted with a pink underline. Below the tabs, the section title 'Access Tokens' is displayed. A large button at the bottom right of the section is labeled 'Add token' with a small circular icon containing a plus sign. A red arrow points directly at this button.

- \_\_33. Notice that you can manage repository with [GitHub](#), [GitHub Enterprise](#), [BitBuket](#) and [BitBucket Server](#).

The screenshot shows the 'Add token' page within the 'Git Integrations' tab. The title 'Add token' is at the top. Below it, a note says 'Visit [GitHub personal access tokens](#), select repository scope and generate a token.' A section titled 'Platform\*' contains four options: 'GitHub' (selected with a red dot), 'GitHub Enterprise', 'BitBucket', and 'BitBucket Server'. The 'GitHub' option is highlighted with a red circle.

- \_\_34. This will allow you to integrate ICP-D projects with your current CICD (Continuous Integration and Continuous Delivery) pipe line to automate delivery of the artifacts. You can use capabilities from IBM Cloud Platform to build cloud native microservices applications which are tied to the ML/AI model development and delivery pipeline.

- \_\_35. Click [Add-on](#) icon.



- \_\_36. Notice that [Db2 Event Store](#) and [Db2 Warehouse](#) applications available and ready to be installed in the IBM Cloud Private for Data platform with click of a button. You will see more IBM and IBM Business Partners applications in the Add-on as time progresses. Click [View Deployed Databases](#).

The screenshot shows the 'Add-ons' page. The title 'Add-ons' is at the top. Two add-on cards are visible: 'Db2 Event Store' and 'Db2 Warehouse'. The 'Db2 Event Store' card has a 'Get started' button with a red arrow pointing to it. The 'Db2 Warehouse' card has a 'View deployed databases' button with a red arrow pointing to it. Both cards include a brief description of their functions and an 'IBM' logo.

- \_\_37. In this local ICP-D cluster – we have already installed Db2 Warehouse. Click the three vertical dots (AKA **More Vert**) and click **Details**.

Your database is available.

db2wh

Database name	BLUDB
Database type	db2whsmp
Database software version	2.11.0
Created on	2018-10-02
Status	Available

Node(s)

HOSTNAME	CPU	MEMORY
10.1.238.159	1 cores	3 GB

Storage

Claim name zen-db2wh

Access Information

Username	user999	edit
Password	#OV?xq17K81ur?L*	edit
JDBC Connection URL	jdbc:db2://192.168.142.10...	edit



**Note:** By design DB2 Warehouse is running on the ICP-D cluster. The intent of this is to provide 'sandbox' capability for projects on the platform. Clients can also choose to use this capability to run production data marts and data warehouses.

- \_\_38. Note that applications and databases available through **Add-on** are native to the ICP-D. IBM collaborates and partners with [MongoDB](#), [portworx](#), [Datameer](#), [Lightbend](#), [Senzing](#), [Prolifics](#) and others to build an ecosystem around ICP-D. The details can be found at <https://www.ibm.com/products/cloud-private-for-data/partners>

## 1.12 Why make IBM Cloud Private for Data your platform for data

- \_\_39. IBM understands data and provides an integrated, end-to-end data platform that enables enterprises to:
- Collect relevant data and make it simpler and more accessible
  - Use federation, virtualization and/or transformation to combine and refine data sets
  - Organize data so it can be trusted
  - Analyze insights on demand
  - Embed machine learning everywhere
- \_\_40. Let's begin this journey with the following lab exercises to demonstrate all of the above.

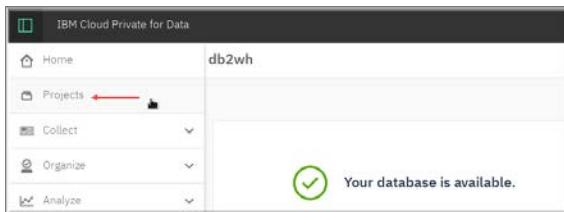
**\*\* End of Lab 01 - Setup and Introduction.**

## Lab 02 Executive Demo

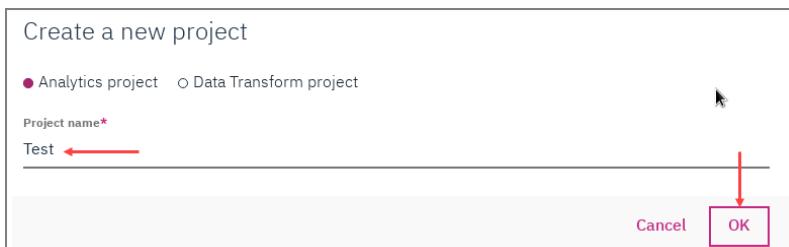
Before we proceed with the lab exercises, let's demonstrate some of the key functions of ICP-D and the microservices-based application before and after implementation of the machine learning models.

### 2.1 Import Project

- 1. ICP-D is a platform that allows a team to work collaboratively using tools and a collection of related assets in a project with access and view control.
- 2. We will work with a project that we have created for the purpose of this workshop. You will create your own projects and explore features in the following lab exercises.
- 3. Switch to the ICP-D web console.
- 4. Hover your mouse on the left navigation pane and it will pull the menu towards right. Click **Projects**.



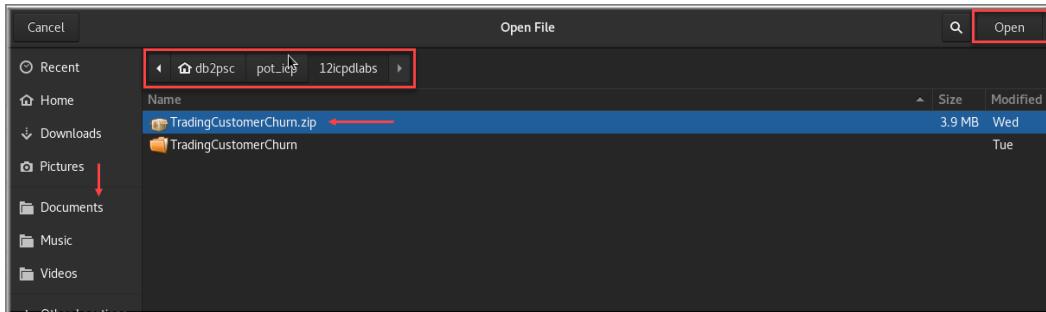
- 5. Click **New Project**.
- 6. Give the project the name **Test** and click **OK**.



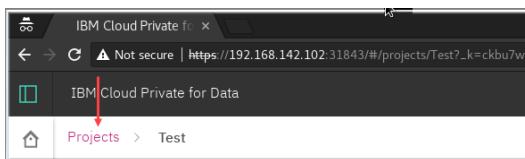
- 7. Click the middle tab **From File** and then click **browse**.



- \_\_8. On the left pane, double-click **Documents**. On the right pane, double-click **pot\_icp**, double-click **12icplabs** and click **TradingCustomerChurn.zip**.



- \_\_9. Click **Create** (towards the bottom right).
- \_\_10. The project creation for the first time is an involved process that launches many resources in the ICP-D cluster. It may take a minute or two.
- \_\_11. Click **Projects** (top left-hand corner).



- \_\_12. Notice that we imported Project **TradingCustomerChurn** from a zip file.
- \_\_13. Click **TradingCustomerChurn**.

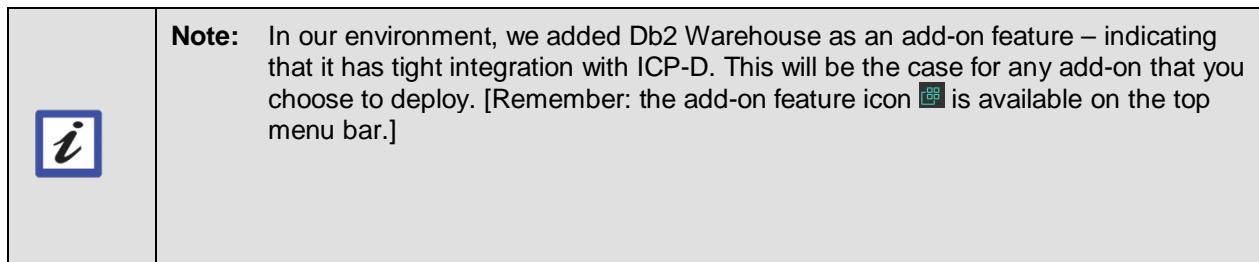
NAME	PROJECT TYPE
TradingCustomerChurn	Standard
dsx-samples	-

- \_\_14. It displays the number of assets, environments, data sources, jobs and collaborator.

## 2.2 Connect to a Data Source

- \_\_15. Click **Data Sources** and then click **Db2Warehouse**.

NAME
Db2Warehouse



- \_\_16. Scroll down and specify *Username*: db2psc and *Password*: password and click **Save** [Located in bottom right corner.]

Username \*  
db2psc

Password \*  
password

Shared

Password is password

- \_\_17. Note: There was no test connection facility. Let's see how this can be done easily.

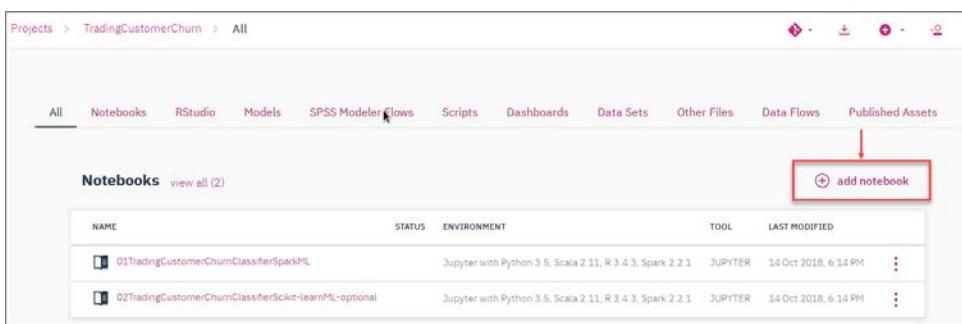
- \_\_18. Click **TradingCustomerChurn** [middle breadcrumb in the top menu bar.]



- \_\_19. Click **Assets**. [Note: Your number may be different.]



- \_\_20. Click **+ add notebook**.



- \_\_21. Enter Name: **Test** and select the *Environment*: **Jupyter with Python 3.5, Scala 2.11, R 3.4.3, Spark 2.2.1**.

Blank From File From URL

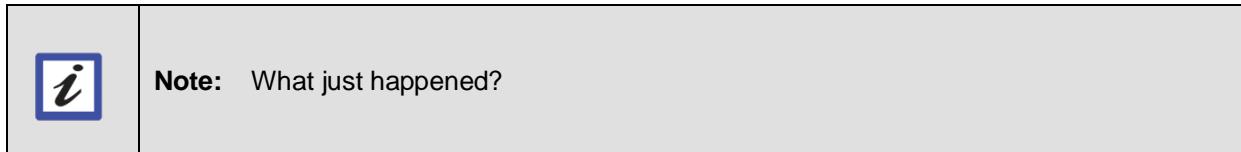
Name\*  
Test 46  
This name is valid

Description  
Type your description here

Environment\*  
Jupyter with Python 3.5, Scala 2.11, R 3.4.3, Spark 2.2.1 500

Language\*  
Python 3.5

\_\_22. Click **Create** [bottom right corner.]



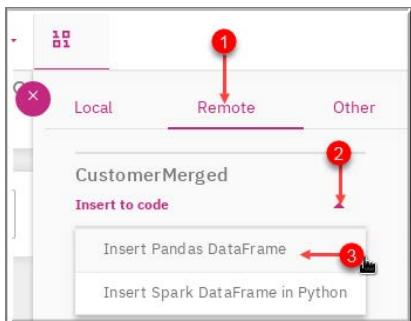
\_\_23. Wait for the Jupyter notebook to launch and notice what happens.

- A pod was instantiated – which means loading a complete compute Jupyter notebook environment (7+ GB) with all the artifacts from the private ICP-D registry.
- IBM Cloud Private schedules this pod on any VM – wherever CPU and memory resources are available.
- IP addresses and connections are all configured automatically.
- The same working environment can be used by multiple users. If a single pod's resources are not sufficient, another environment is created automatically.
- When the number of users grow, you can add more machines to the ICP-D cluster and scheduling of resources is handled automatically.
- ICP-D's scale-out model is pretty effective.
- You no longer have to wait days or even weeks to get the compute resources.
- IBM Cloud Private makes 40 – 60 percent more efficient use of compute resources. This means more users can be accommodated with same compute capacity. As one task completes, its resources are freed up to work on next one.

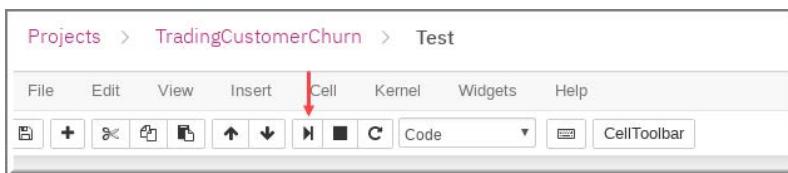
\_\_24. Click the **Binary Number** icon on the menu bar (last icon).



\_\_25. Click **Remote** [Middle]. Click **down arrow** for the first item and click **Insert Pandas Data Frame**.



- \_\_\_26. Let's run this Jupyter Notebook cell. [Remember **Shift-Enter** is the keyboard shortcut to run a selected cell.]



- \_\_\_27. During the execution of a cell in a Jupyter notebook, an asterisk (\*) displays in the square bracket which changes to a number when execution for that cell completes.
- \_\_\_28. If the connection is successful, you get the top few rows of the DB2 warehouse table. Skip the next step – which is only required if your cell did not execute.
- \_\_\_29. If your cell did not execute successfully, don't worry. Minimize the browser window and double-click [Create Workshop User in Db2 Warehouse](#) to create a temporary Db2 Warehouse user.

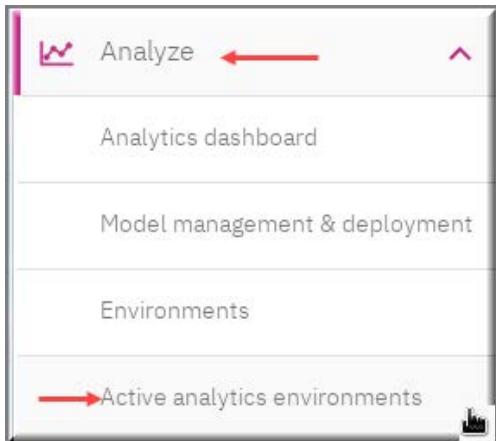


- \_\_\_30. Switch back to the browser window and press **Shift-Enter** to execute the cell again. You should now see the first few rows from the Db2 Warehouse table as shown below.

ID	CHURNRISK	GENDER	STATUS	CHILDREN	ESTINCOME	HOMEOWNER	AGE	TOTALDOLLARVALUETRADED	TOTALUNITSTRADED	LARGESTSINGLETRANS/
0	Low	F	S	1	38000.00	N	24	59755.98	206	29877.990
1	Low	M	M	2	29616.00	N	49	29782.98	45	14891.490
2	Low	M	M	0	19732.80	N	51	24812.48	22	12406.240
3	High	M	S	2	96.33	N	56	26132.61	32	13066.305
4	High	F	M	2	52004.80	N	25	5030.50	23	1257.625

## 2.3 Check Running Environment

- \_\_31. Hover the mouse on the left menu bar and click [Analyze](#)  $\Rightarrow$  [Active analytics environments](#).



- \_\_32. Notice the Jupyter runtime environment and check the status. It should be green.

Runtimes 1									All Projects	All Runtime Types
Name	Runtime Type	User	Project	Job Name (Run ID)	Date Started	CPU (cores)	Gpus	Memory (GB)	Status	
Jupyter with Python 3.5, Scala 2.11, R 3.4.3, Spark 2.2.1	Environment	admin	TradingCustomerChurn		14 Oct 2018, 6:53 PM	-	-	-	<span style="color: green;">●</span>	



**Note:** Make sure that you have only one environment active at any given time in this resource-constrained environment. If more than one environment is running, the system may become sluggish.

## 2.4 Stock Trader Opening Bell Analysis dashboard

- \_\_33. Boatswain Trading is aware that revenues are declining. Let's discover the actual trend and what the data presents. In our scenario, the business analyst has requested aggregated stock transaction data. It contains data such as the historical total number of individual customer visits to our website and number of trades per customer for the past year.
- \_\_34. Let's take a look at the dashboard first. We will build same dashboard in a later exercise to learn to use the ICP-D integrated data platform to gain an insight into our business.
- \_\_35. Hover your mouse on the left menu bar and click [Projects](#). Click [TradingCustomerChurn](#) project.

NAME	PROJECT TYPE
TradingCustomerChurn	Analytics
dsx-samples	-

- \_\_36. Click **Assets** and then click **Dashboards**.

Projects > TradingCustomerChurn > All

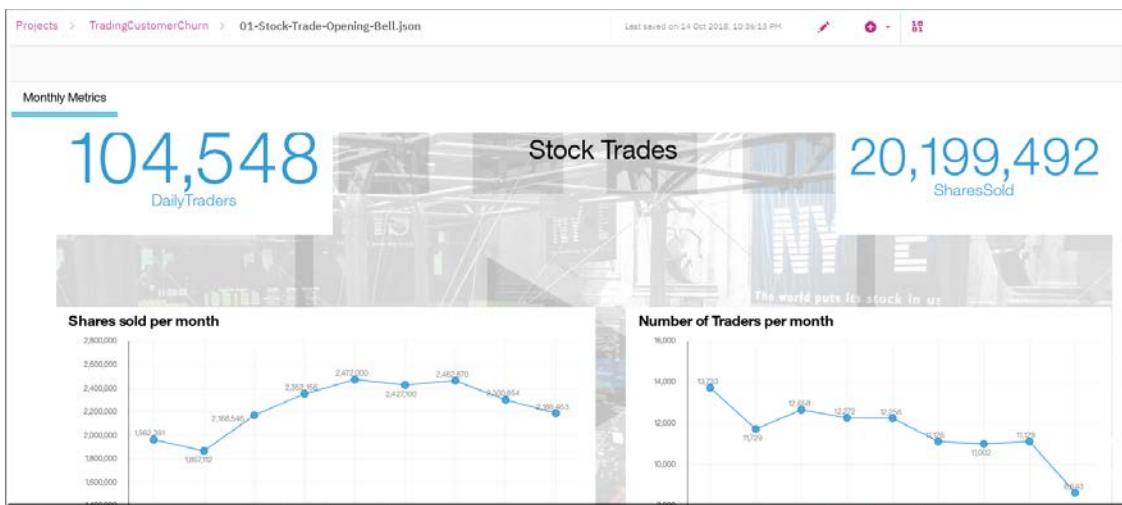
All Notebooks RStudio Models SPSS Modeler Flows Scripts Dashboards Data Sets

- \_\_37. We have prebuilt three dashboards. [Note: You will get a chance to build these dashboards in later exercises.]
- \_\_38. The first dashboard shows a visualization of the 'Stock's Before' analysis.
- \_\_39. The second dashboard shows the results of the analysis of the customer demographics data that was collected from various sources, such as Data Warehouse and Hadoop environment, after curation of data. This information is used to evaluate which measures (AKA factors or variables) have the greatest impact on the customer attrition.
- \_\_40. The third dashboard shows the analysis after implementation of the AI models and the impact that it had on results of the business.
- \_\_41. Click **01-Stock-Trader-Opening-Bell.json**.

Dashboards (3)

NAME
03-Stock-Trader-Closing-Bell.json
02-Stock-Trader-Demographic-Discovery.json
01-Stock-Trader-Opening-Bell.json

- \_\_42. The dashboard shows the consolidated report of our business. You can see that the number of shares sold per month is relatively flat and daily trades are declining.



- \_\_43. The data for this report is stored at a central place and is governed and managed through the team collaboration between business analysts, data engineer, data scientist and data steward.

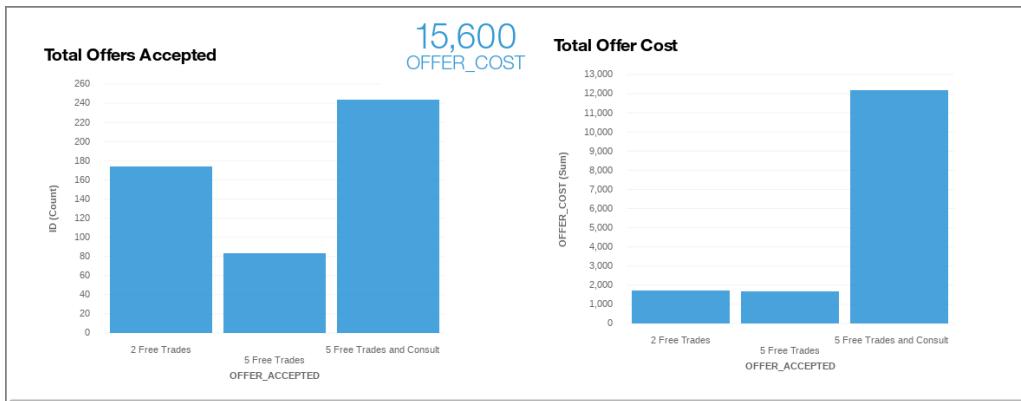
## 2.5 Stock Trader Analysis Closing Bell Analysis

The purpose of our lab exercises is to show the journey of analytics modernization from our opening dashboard to closing bell dashboard.

- \_\_44. Let's look at the finished product, and then dive into the details.
- \_\_45. Go back to [Project](#) ⇒ [TradingCustomerChurn](#) ⇒ [Assets](#) ⇒ [Dashboard](#) and click [03-Stock-Trader-Closing-Bell.json](#)
- \_\_46. The dashboard is built from the data sets present in the project.
- \_\_47. Notice the trends of the Shares Sold per Month and Number of Traders per Month and compare with the Opening Bell dashboard.
- \_\_48. The graph should display as the following figure:



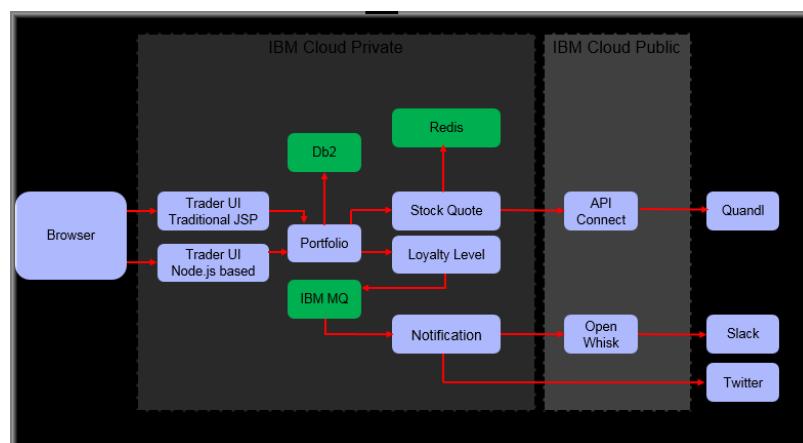
- 49. After implementation of the predictive models built and scored using different machine learning algorithms, changes to the modern microservice application were made to offer personalized promotions based on customer's history and risk of separation. The results after a few months show an increase in the number of shares sold per month and number of stock traders per month.
- 50. Scroll down the dashboard and to see the total number of retention offers accepted and the total cost of those offers. Given the increased number of 'shares sold' and 'traders per month' from the prior graph, the USD 15,600 cost (that is, lost revenue from free trades and wealth consulting) seems like a very good investment.



- 51. This dashboard uses data for a full year from a .csv file and the second data set uses the customer offers and costs to business from a Db2 Warehouse table.

## 2.6 Run Stock-Trader Before Application

- 52. The software development and distribution methods and deployment have changed rapidly since 2015. The application that we show is built on a microservices-based architecture that allows loosely coupled services to work together. This allows components to upgraded incrementally, as opposed to traditional, monolithic application development.
- 53. The sample stock trader application architecture is displayed below showing interaction between various microservices.



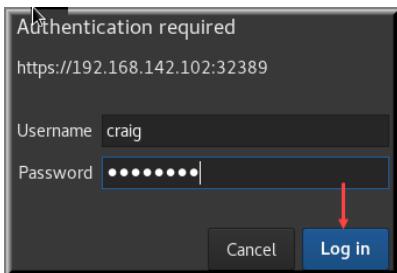
- \_\_54. The infrastructure services Db2, MQ and Redis are shown in green boxes whereas microservices are shown in purple.
- \_\_55. The main advantage of microservice-based application is the ability to modify, deploy, update and scale services dynamically and individually.
- \_\_56. In a modern microservice mesh architecture, it is possible to run multiple versions of the same service in production and control the traffic dynamically based on rules that can be defined outside of the service itself, eliminating the need to need to modify any of the microservice code.

	<b>Note:</b> Start talking with your developers about the concept of a Service Mesh and implement it on IBM Cloud Private platform. Consult us if you would like more information.
---	--

- \_\_57. The application shown here is the demo application that does not address all features and functions of a commercial application but gives you insight into the concept of application modernization. The IBM Cloud Private platform provides tools to help you migrate your monolithic applications to a microservices-based architecture, free of cost [ML1].
- \_\_58. If your journey on this path has not yet begun, share this application with your development teams. It will provide you with the framework you need to start building your own microservice-based applications. You can find it at: <https://github.com/IBMStockTrader>
- \_\_59. Let's launch this application.
- \_\_60. Minimize your browser window.
- \_\_61. Double-click **Stock Before**.



- \_\_62. A login popup text box displays – Enter Username: **craig** and Password: **password**.



- \_\_63. NOTE: The first time the application is launched it may take a little longer to start due to limited system resources of the workshop environment.
- \_\_64. Select **Create a new portfolio** and click **Submit**.

Welcome to IBM Trader powered by ICP for Data

<input checked="" type="radio"/> Create a new portfolio	<input type="radio"/> Retrieve selected portfolio
<input type="radio"/> Update selected portfolio (add stock)	<input type="radio"/> Delete selected portfolio

Owner	Total	Loyalty Level

**Submit** **Change User**

**Note:** One advantage of a microservices-based application is that it allows polyglot services to be written in the language of choice by the developers. This allows services and applications to work together easily.

Developers are no longer forced to change their preferred language or learn a new language due to monolithic legacy applications.

\_\_65. Type **TechStocks** as name of the portfolio and click **Submit**.

Welcome to IBM Trader powered by ICP for Data

Owner	TechStocks
-------	------------

**Submit**

\_\_66. Select **Update selected portfolio** and click **Submit**.

<input checked="" type="radio"/> Create a new portfolio	<input type="radio"/> Retrieve selected portfolio
<input type="radio"/> Update selected portfolio (add stock)	<input type="radio"/> Delete selected portfolio

Owner	Total	Loyalty Level
TechStocks	\$0	Basic

**Submit** **Change User**

\_\_67. Type **IBM** and **1000** and click **Submit** to buy 1000 stocks of IBM. [Note: You can choose any other, if you know the stock symbol.]

## Add Stock

Welcome to IBM Trader powered by ICP for Data

Owner	
Stock Symbol	IBM ←
Number of Shares	1000 ←
<b>Submit</b>	



**Note:** If Internet access is available – the stock quotes are taken from [quandl.com](https://quandl.com) using IBM API Connect through IBM Cloud.

The cached values from Quandl are stored in redis.

- 68. Notice that the loyalty level changed to Gold.

## Summary

Welcome to IBM Trader powered by ICP for Data

<input type="radio"/> Create a new portfolio		
<input checked="" type="radio"/> Retrieve selected portfolio		
<input type="radio"/> Update selected portfolio (add stock)		
<input type="radio"/> Delete selected portfolio		
Owner	Total	Loyalty Level
TechStocks	\$115,670	Gold ←
<b>Submit</b>	<b>Change User</b>	

- 69. The Ingress rules are defined such that the [notification-service](#) can trigger a notification to either [Slack](#) or [Twitter](#).

- 70. If [notification-service](#) sent it to Twitter, you can check the message at <https://twitter.com/ibmstocktrader>.



**IBM Stock Trader** @IBMSockTrader  
The IBM Stock Trader sample for IBM Cloud: [github.com/IBMSockTrader](https://github.com/IBMSockTrader). Notifications posted here when stock

**Tweets** 445    **Following** 4    **Followers** 7

**Tweets** **Tweets & replies**

IBM Stock Trader @IBMSockTrader · 4m  
On Monday, October 15, 2018 at 2:30 PM UTC, TechStocks changed status from Basic to Gold. #IBMSockTrader

- \_\_71. The goal is to enhance this application with analytics modernization to increase revenue and profits. We will name this application as "[Stock - After](#)".
- \_\_72. The modern Stock Trader application built on modern IBM Cloud Private platform uses agile technologies with a modern CI/CD pipeline to implement changes rapidly.

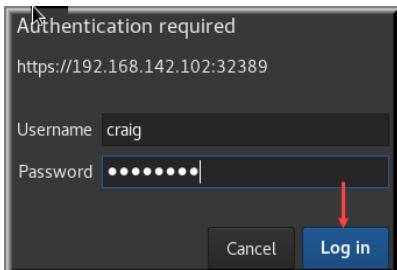
	<b>Note:</b> Gone are the days of "Code Freeze" – Welcome to the modern CI/CD world.
---	--

## 2.7 Run Stock-Trader After Application

- \_\_73. Companies that are able to keep pace with the technological breakthroughs will have a competitive advantage over companies that do not.
- \_\_74. Modernization of your application is the essence, but ICP-D goes one step further – it allows you to use your data to modernize your analytics capabilities.
- \_\_75. The key here is to speed the organization of the data and make it ready for machine learning and artificial intelligence to provide analytics modernization.
- \_\_76. Let's "begin with the end in mind" and take a look at the finished product. Then we will dive into the steps that lead us to [Analytics Modernization](#).
- \_\_77. Double-click [Stocks After](#).



- \_\_78. Log in using Username: [craig](#) and Password: [password](#).



- \_\_79. Notice the offer "[no processing fee for next 5 trades](#)".

Welcome to IBM Trader powered by ICP for Data

- Create a new portfolio
- Retrieve selected portfolio
- Update selected portfolio (add stock)
- Delete selected portfolio

Owner	Total	Loyalty Level
TechStocks	\$115,670	Gold

Though looking simple - a lot has gone through to provide machine learning predictive model scoring service.

no processing fee for next 5 trades

Advertisement

**IBM Cloud Private for Data**

- Cloud agile
- Lightning fast
- AI-ready

No assembly required!

- \_80. A predictive model that was built (a journey through which we will go in subsequent lab exercises) looks for the score of the user "[craig](#)" and returns a number of parameters. Our new predictive analysis microservice uses a single parameter and transforms it into an offer to the customer based on the "separation risk" that the model predicts.
- \_81. This will become clear as you work through details in the [Analyze](#) lab exercise.
- \_82. Of course, models are only as good as the data they use. We will use [Organize](#) lab to classify, catalog, curate and establish the lineage. This will provide the best relevant and necessary data to the [Analyze](#) phase – where models are built, tested, scored and continuously improved.
- \_83. Before we [Organize](#) the data, we first must collect it – regardless where it resides. The [Collect](#) exercise shows you the process of getting data.

i

**Note:** Remember the term "Borderless Data".

- \_84. Click [Change User](#).

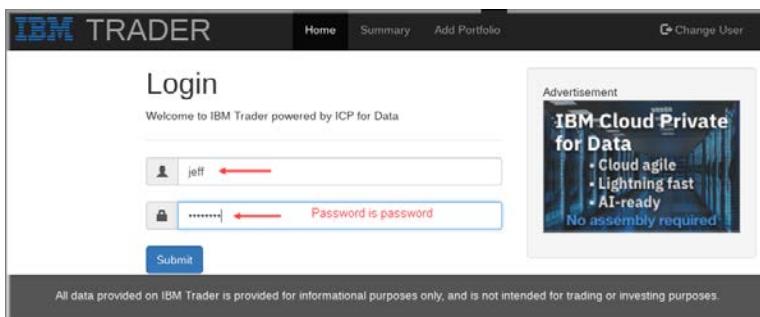
Change User

no processing fee for next 5 trades

Advertisement

**IBM Cloud Private**

- \_85. Specify credentials: [jeff](#) and [password](#).



\_\_86. Notice the personalized award offer.

The screenshot shows the IBM TRADER summary page. The top navigation bar includes Home, Summary (which is selected), Add Portfolio, Predictive Analysis, and Change User. The main content area has a heading "Summary" and a note "Welcome to IBM TRADER". Below this are four radio button options: "Create a new portfolio", "Retrieve selected portfolio" (which is selected), "Update selected portfolio (add stock)", and "Delete selected portfolio". A table below shows a single portfolio entry: "Owner" (TechStocks), "Total" (\$115,670), and "Loyalty Level" (Gold). At the bottom are "Submit" and "Change User" buttons. A red callout bubble points to the text "Personalized offer based upon customer demographics" which is displayed above a box containing the offer: "no processing fee for next 5 trades + 1 free wealth management session with certified planner". To the right of the summary is an advertisement for "IBM-Cloud Private for Data" with the tagline "No assembly required".

- \_\_87. Let's dive deep in behind the scene ML model prediction. Click [Predictive Analysis](#). This screen is built to directly interact with the ML model. Change variables to see how those specific sets of variables impact predictions.

The screenshot shows the 'Customer Churn Predictor' page. At the top, there are navigation links: Home, Summary, and Add Portfolio. Below the title, a welcome message reads: 'Welcome to IBM Trader powered by ICP for Data'. The form contains the following fields:

Age:	26
Gender:	<input type="radio"/> Male <input checked="" type="radio"/> Female
Marital Status:	<input type="radio"/> Married <input checked="" type="radio"/> Single
Number of Children:	0
Home Owner:	<input type="radio"/> Yes <input checked="" type="radio"/> No
Estimated Income:	55000
Net Realized Gains YTD:	1500
Net Realized Losses YTD:	0
Smallest Single Transaction:	150
Largest Single Transaction:	1500
Total Dollar Value Traded:	8500
Total Units Traded:	350

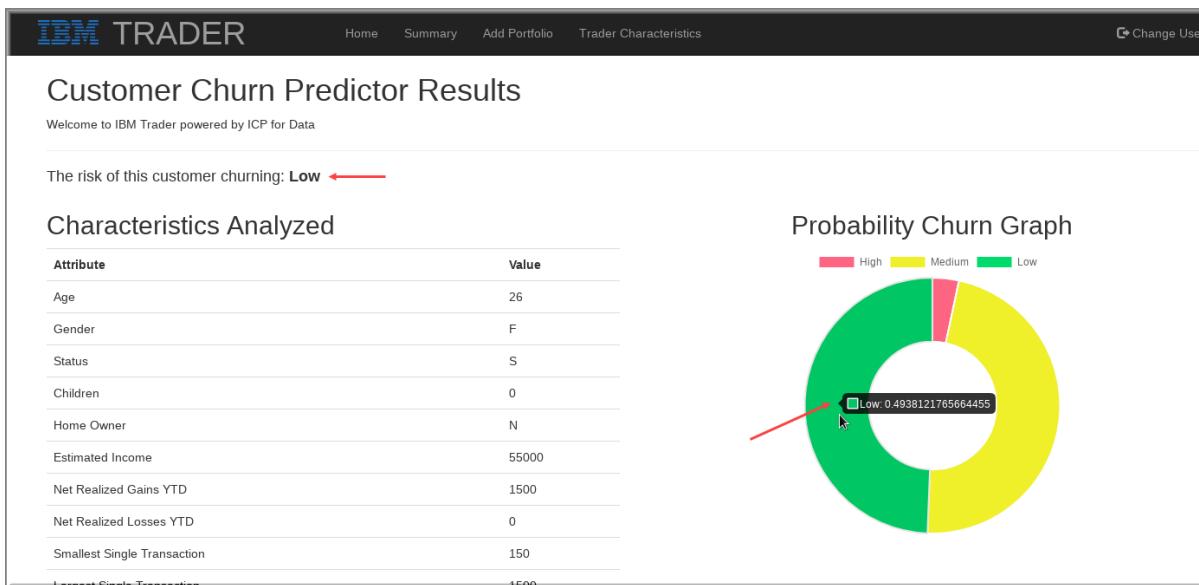
- \_\_88. Scroll down and click [Submit](#) to assess the separation risk based on default parameters.

The screenshot shows a form with the following fields:

Total Units Traded:	350
Days Since Last Login:	2
Days Since Last Trade:	4
Percentage Change Calculation:	50

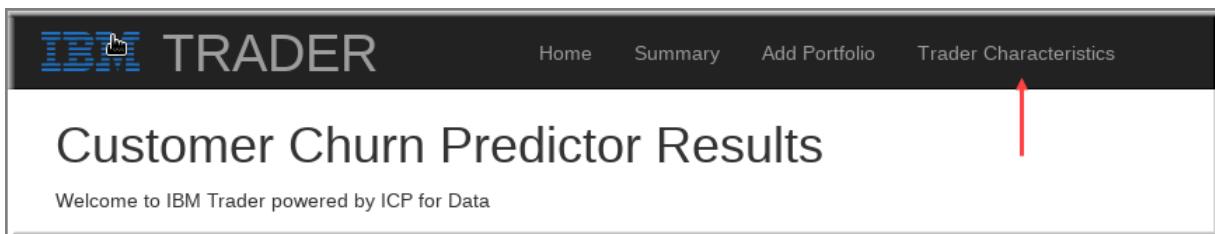
A red arrow points to the 'Submit' button at the bottom left of the form.

- \_\_89. This is the result using the default values of separation risk predictor. In this case, there is a *Low* to moderate risk of this customer leaving. This customer has a profitable portfolio.



- \_\_90. Let's change some parameters.

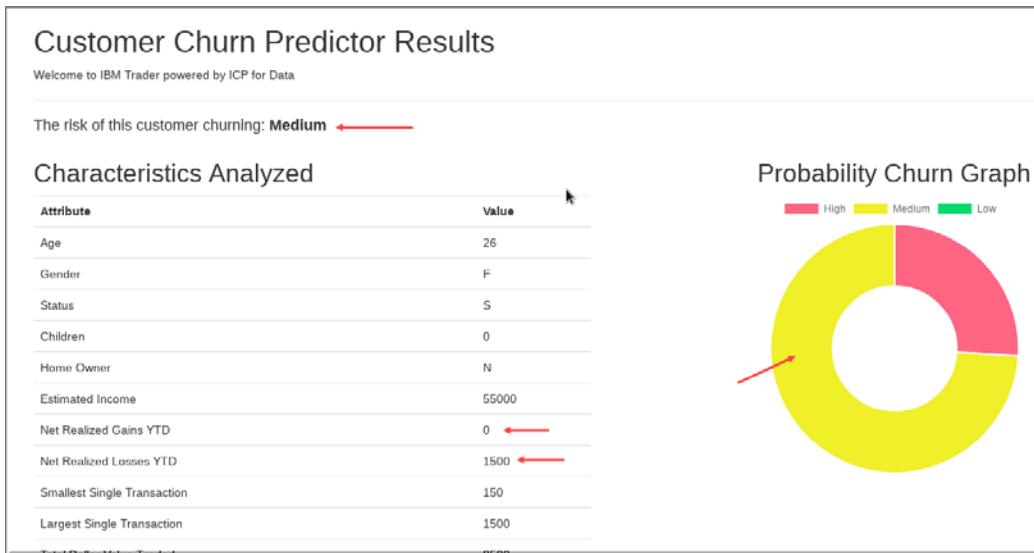
- \_\_91. Click [Trader Characteristics](#).



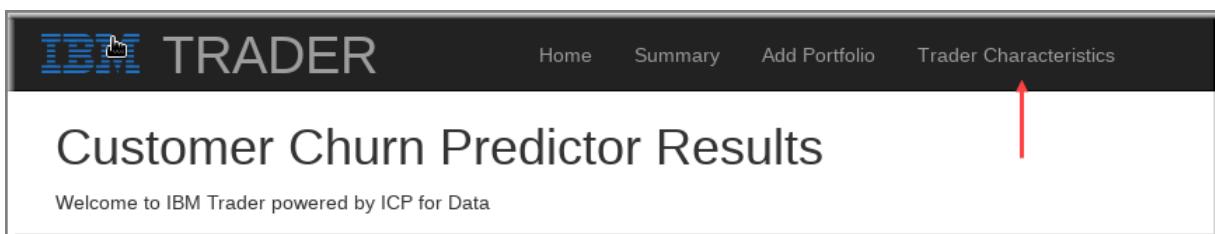
- \_\_92. Change [Net Realized Gains YTD](#) from **1500** to **0** and change [Net Realized Losses YTD](#) to **1500** from **0**. Basically, making this portfolio as not profitable. Scroll to the bottom and click [Submit](#).

The screenshot shows the 'Trader Characteristics' form. It contains fields for Number of Children (0), Home Owner (Yes), Estimated Income (55000), Net Realized Gains YTD (0), Net Realized Losses YTD (1500), Smallest Single Transaction (150), and Largest Single Transaction (1500). Red arrows point to the 'Net Realized Gains YTD' and 'Net Realized Losses YTD' fields, with text indicating changes from 1500 to 0 and 0 to 1500 respectively.

- \_\_93. Notice the probability of separation risk of this customer. The risk has changes from **Low** to **Medium** given the fact that portfolio of this customer is no longer profitable.



- \_\_94. Let's make some more modifications. Click **Trader Characteristics**.



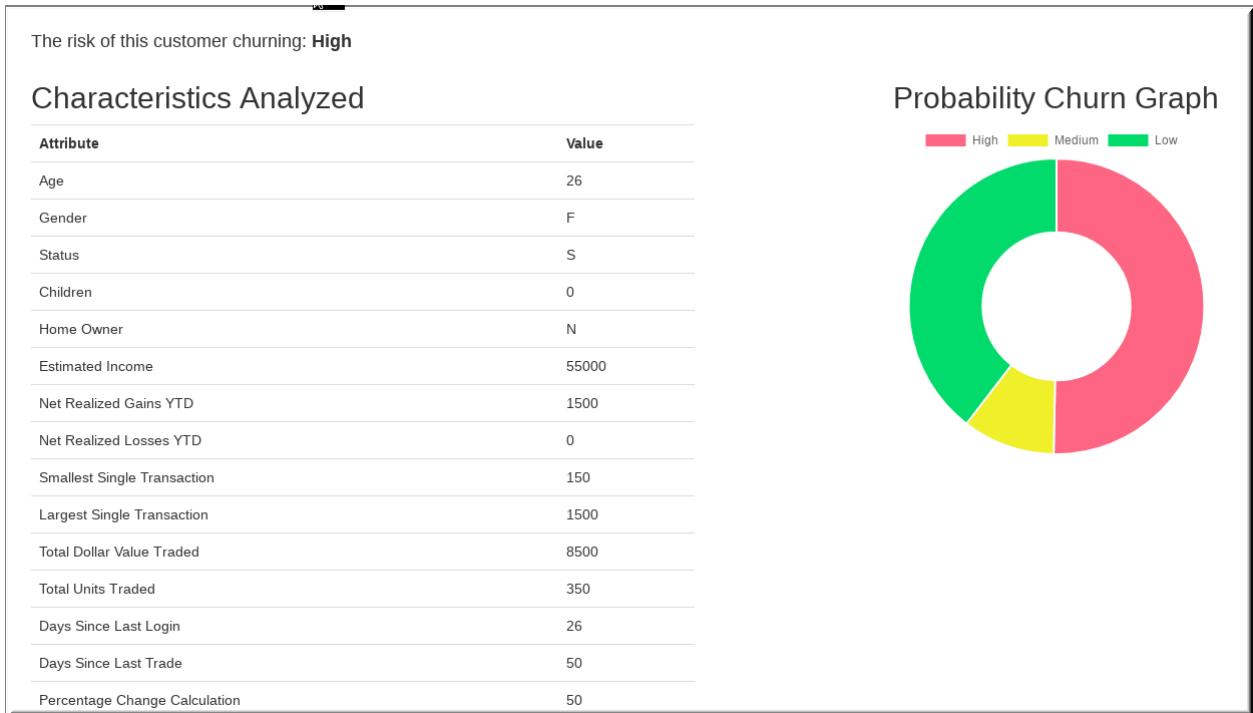
- \_\_95. The default parameters are restored and the portfolio is profitable again. Scroll down and let's change the **Days Since Last Login** from **2** to **26** and **Days Since Last Trade** from **4** to **50**. Click **Submit**.

The screenshot shows the "Trader Characteristics" form. It contains several input fields and a "Submit" button at the bottom.

- "Total Units Traded": 350
- "Days Since Last Login": 26 (with a red arrow and note: "Change to 26 from 2"))
- "Days Since Last Trade": 50 (with a red arrow and note: "Change to 50 from 4"))
- "Percentage Change Calculation": A slider set to 50.

A red arrow points to the "Submit" button at the bottom left.

- \_\_96. Notice that the separation risk changes to **High**. It's been a while since we've seen this customer and there is little account activity. In this case, there is a high probability this customer might leave.



- \_\_97. Regardless of your problem classification domain, ICP-D provides the tools needed to access, clean, shape and govern your data in preparation for ML model development, release and continuous improvement framework.
- \_\_98. From automatic ML model generation (no coding) to comprehensive tool sets to develop complex models (with coding), ICP-D provides the tools necessary to produce your specific ML model classification.
- \_\_99. Once ML models are deployed, they are easily consumed by applications as RESTful-compliant services. The deployed models are easily scalable to meet usage demands as they are packaged as independent scalable docker containers.
- \_\_100. As probabilities are consumed from one or more models, application logic can be applied to allow applications to make intelligent business decisions based on a set of predefined business rules and workflows.
- \_\_101. Let's deep dive in **Collect, Organize, Analyze** and **Deploy** features of the ICP-D as an integrated data platform.

	<b>Note:</b> Do you talk about getting value from the "Dark Data"?
---	--

## \*\* End of Lab 02: Executive Demo

## Lab 03 Collect

Collect data regardless where it resides. This is the key aspect of data virtualization to access, govern and analyze it without duplication and massive data movement.

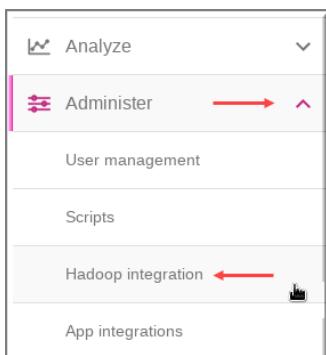
One step further is to use streaming data from IoT devices, banking transactions and more. The Db2 Event Store's add-on feature in ICP-D provides this capability.

### 3.1 Multiple Sources

- 1. ICP-D has add-in capabilities for Db2 OLTP, Db2 Warehouse and Db2 Event Store to install and upgrade at any time to in the platform itself. The add-on capability will support additional data sources from MongoDB, DataMeer, and others in the future.
- 2. The support for data sources available outside ICP-D platform are available for [BigSQL](#), [Oracle](#), [Db2 for z/OS](#), [Hive](#) for Hortonworks and Cloudera, [HDFS](#) and [Informix](#). Additional support for data sources is available through the [Organize](#) category. The list of connected data sources is designed to grow.

### 3.2 Integration with the Hadoop Platform

- 3. The Hadoop cluster has several services which are insecure (using HTTP) and hence the Hadoop cluster is behind the firewall of an organization.
- 4. ICP-D provides integration to Hadoop cluster for Hortonworks and Cloudera through a secured gateway built and customized using Apache KNOX.
- 5. The integration with the Hadoop cluster requires an additional software package that needs to be installed on one of the Hadoop edge nodes (AKA Gateway).
- 6. This has already been done in this lab environment. You can refer to Appendix-C for the necessary details for integration software to be installed on Hadoop edge node for the integration with ICP-D.
- 7. Switch the ICP-D web console tab in your browser. If you do not have web console running, double-click [IBM Cloud Private for Data](#) on the desktop.
- 8. Hover your mouse over the left menu bar and click the **Administer** down arrow to expand it and then click [Hadoop Integration](#).



\_\_9. Determine whether our existing Hadoop cluster is registered.

NAME	SERVICE USER ID	URL
hadoop	dsxhi	<a href="https://sandbox-hdp.hortonworks.com:8443/gateway/192.168.142.102">https://sandbox-hdp.hortonworks.com:8443/gateway/192.168.142.102</a>

\_\_10. If not, click **Add Registration**. [top, right-hand corner.]

\_\_11. Type:

- \_\_a. Display Name: [Hadoop](#)
- \_\_b. Service URL: <https://sandbox-hdp.hortonworks.com:8443/gateway/192.168.142.101>
- \_\_c. Service User ID: [dsxhi](#)
- \_\_d. Click **Add**.

Add Registration

Display Name \*  
Hadoop

Service URL \*  
<https://sandbox-hdp.hortonworks.com:8443/gateway/192.168.142.101>

Service User ID \*  
dsxhi

Cancel Add

\_\_12. You should see the registration success message.

Hadoop Integration		
Successfully registered the system.		
NAME	SERVICE USER ID	URL
hadoop	dsxhi	<a href="https://sandbox-hdp.hortonworks.com:8443/gateway/192.168.142.102">https://sandbox-hdp.hortonworks.com:8443/gateway/192.168.142.102</a>

### 3.3 Data Virtualization from Hadoop Platform

\_\_13. Let's verify a file in Hadoop platform that we will use for data virtualization.

#### 3.3.1 Copy a File to Hadoop

\_\_14. Double-click **Hadoop Web Console** from the VM desktop.

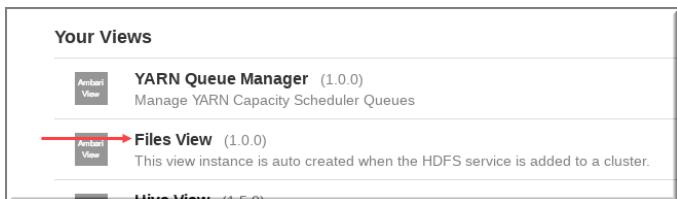
\_\_15. Click **Launch Dashboard** when the HDP Web UI appears.

\_\_16. Type Username: [admin](#) and Password: [password](#).

\_17. From main menu bar, click the 3x3 checkerboard .



\_18. Click **Files View**.



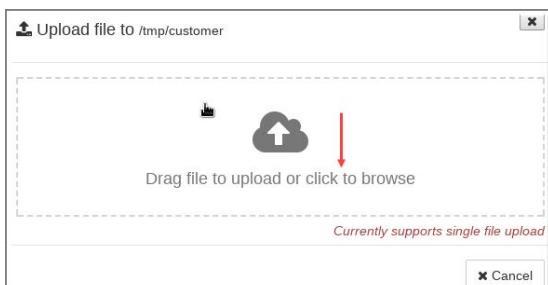
\_19. Click **tmp** folder.

\_20. Click **customer**.

\_21. Maximize your browser screen, if not done already. Click **Upload**.



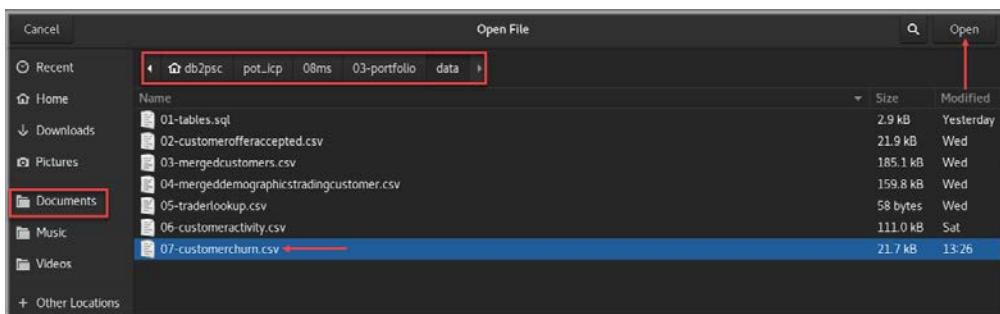
\_22. Click to browse.



\_23. Click **Documents** in the left pane and click **pot\_icp** in right pane.

\_24. Click each **08ms**  $\Rightarrow$  **03-portfolio**  $\Rightarrow$  **data**  $\Rightarrow$  **07-customerchurn.csv**

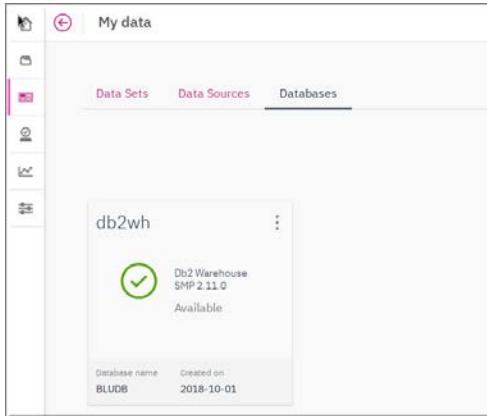
\_25. Click **Open**



\_26. Close the Hadoop browser window and the Hortonworks tab from the main browser window.

### 3.4 Virtualization of Db2 in ICP-D

- \_27. Note that ICP-D has an add-on for Db2 Warehouse. This means that the Db2 Warehouse data from within ICP-D platform is already available.
- \_28. Click **Collect** ⇒ **My Data** ⇒ **Databases**



- \_29. Let's virtualize Db2 OLTP database which we did not add as an add-on in this case.
- \_30. Click **Projects** ⇒ **TradingCustomerChurn**
- \_31. Click **Data Sources**.
- \_32. Click **add data source**. [Towards right side of Data Sources].
- \_33. Type the data source name **db2oltp**, Select **Db2** from the drop-down menu, and type the following JDBC URL: [No typos please.]  
**jdbc:db2://ibmdb2-ibm-db2oltp-dev.stocktrader.svc.cluster.local:50000/PSDB**
- \_34. Type Username: **db2psc** and Password: **password**. Click **Add remote data set**

The screenshot shows the 'Add data source' form with the following field values and state:

- Data source name \***: db2oltp (1)
- Description**: Type your description here
- Data source type \***: Db2 (2)
- JDBC URL \***: jdbc:db2://ibmdb2-ibm-db2oltp-dev.stocktrader.svc.cluster.local:50000/PSDB (3)
- Username \***: db2psc (4)
- Password \***: ..... (5)
- Shared**:  Shared (6)
- Add remote data set** button (7)

- \_\_35. Type *Remote data set name* STOCK, *Schema* DB2PSC and *Table* STOCK. [Far right bottom of screen.] (The names should be in uppercase.)

Remote data set name \*  
STOCK

Description  
Description of dataset

Schema  
DB2PSC

Table \*  
STOCK

- \_\_36. Click [Create](#)
- \_\_37. Click [db2oltp](#) data source again, scroll down and click [Add data set](#).
- \_\_38. Type *Remote data set name* PORTFOLIO, *Schema* DB2PSC and *Table* PORTFOLIO and click [Create](#).
- \_\_39. Click [Save](#).

### 3.5 Virtualization of Oracle in ICP-D

- \_\_40. Click [Add data source](#)
- \_\_41. Type the data source name [oracle](#), Select [Oracle](#) from the dropdown, and type the following JDBC URL: [No typos please.]
- <jdbc:oracle:thin:@192.168.142.106:1521:XE>
- Notice: Next, you need to add : and @.
- \_\_42. Type Username: [scott](#) and Password: [tiger](#). Click [Add remote data set](#)

Data source name \*  
oracle

Description  
Type your description here

Data source type \*  
Oracle

JDBC URL \*  
<jdbc:oracle:thin:@192.168.142.106:1521:XE>

Username \*  
scott

Password \*  
tiger  
Password is tiger

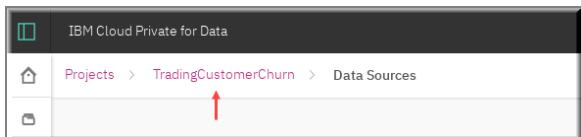
Shared

+ Add remote data set

- \_\_43. Type: *Remote data set name* ALL\_TABLES, *Schema* SYS and *Table* ALL\_TABLES and click [Create](#).

### 3.6 Test Data Virtualization

\_\_44. Click **TradingCustomerChurn** project.



\_\_45. Click **Assets**  $\Rightarrow$  **Notebooks**

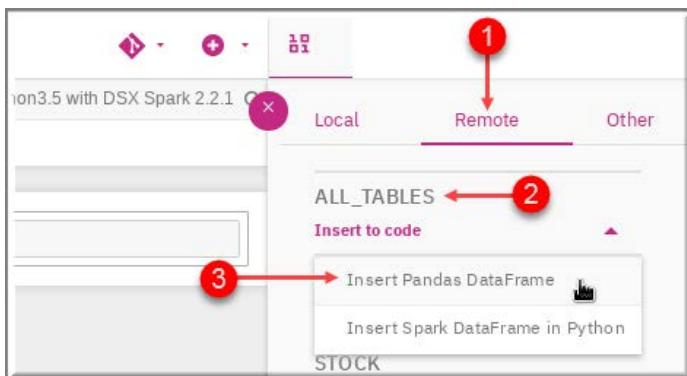
\_\_46. Click **Test** to open.

\_\_47. Click **scissor** icon on the menu bar to delete all cells.

\_\_48. Click **binary number** icon.



\_\_49. Select middle tab **Remote**. Insert **ALL\_TABLES** (From Oracle 11g database) Panda DataFrame to a cell.



\_\_50. Press **Shift-Enter** to execute the cell.

\_\_51. You should see the **ALL\_TABLES** data from the Oracle database.

```
In [12]:
```

```
import dsx.core.utils.requests, jaydebeapi, os, io, sys
from pyspark.sql import SparkSession
import pandas as pd
df1 = None
dataSet = dsx.core.utils.get_remote_data_set_info('ALL_TABLES')
dataSource = dsx.core.utils.get_data_source_info(dataSet['datasource'])
if (sys.version_info >= (3, 0)):
    conn = jaydebeapi.connect(dataSource['driver_class'], dataSource['URL'], [dataSource['user'], dataSource['password']])
else:
    conn = jaydebeapi.connect(dataSource['driver_class'], [dataSource['URL'], dataSource['user'], dataSource['password']])
query = 'select * from ' + (dataSet['schema'] + '.' if (len(dataSet['schema'].strip()) != 0) else '') + dataSet['table']
df1 = pd.read_sql(query, con=conn)
df1.head()
```

```
Out[12]:
```

	OWNER	TABLE_NAME	TABLESPACE_NAME	CLUSTER_NAME	IOT_NAME	STATUS	PCT_FREE	PCT_USED	INI_TRANS	MAX_TRANS	... SKIP_CORRUPT	MONITORING	CLUSTER_OWNER	DEPENDENCIES	COMPRE:
0	SYS	ICOL\$	SYSTEM	C_OBJ#	None	VALID	0.0	0.0	0.0	0.0	... DISABLED	YES	SYS	DISABLED	DISABLED
1	SYS	CONS\$	SYSTEM	None	None	VALID	10.0	40.0	1.0	255.0	... DISABLED	YES	None	DISABLED	DISABLED
2	SYS	UNDOS\$	SYSTEM	None	None	VALID	10.0	40.0	1.0	255.0	... DISABLED	YES	None	DISABLED	DISABLED
3	SYS	PROXY_ROLE_DATA\$	SYSTEM	None	None	VALID	10.0	40.0	1.0	255.0	... DISABLED	YES	None	DISABLED	DISABLED
4	SYS	FILE\$	SYSTEM	None	None	VALID	10.0	40.0	1.0	255.0	... DISABLED	YES	None	DISABLED	DISABLED

\_\_52. Add Db2 **PORTFOLIO** in a cell, press **Shift-Enter** to get the data from Db2 data source.

```
In [13]:
```

```
df2 = None
dataSet = dsx.core.utils.get_remote_data_set_info('PORTFOLIO')
dataSource = dsx.core.utils.get_data_source_info(dataSet['datasource'])
if (sys.version_info >= (3, 0)):
    conn = jaydebeapi.connect(dataSource['driver_class'], dataSource['URL'], [dataSource['user'], dataSource['password']])
else:
    conn = jaydebeapi.connect(dataSource['driver_class'], [dataSource['URL'], dataSource['user'], dataSource['password']])
query = 'select * from ' + (dataSet['schema'] + '.' if (len(dataSet['schema'].strip()) != 0) else '') + dataSet['table']
df2 = pd.read_sql(query, con=conn)
df2.head()
```

```
Out[13]:
```

	OWNER	TOTAL	LOYALTY
0	TechStock	141500.0	Gold

\_\_53. Similarly Add **STOCK** data set in a cell and press **Shift-Enter** to execute the cell.

```
In [14]:
```

```
df3 = None
dataSet = dsx.core.utils.get_remote_data_set_info('STOCK')
dataSource = dsx.core.utils.get_data_source_info(dataSet['datasource'])
if (sys.version_info >= (3, 0)):
    conn = jaydebeapi.connect(dataSource['driver_class'], dataSource['URL'], [dataSource['user'], dataSource['password']])
else:
    conn = jaydebeapi.connect(dataSource['driver_class'], [dataSource['URL'], dataSource['user'], dataSource['password']])
query = 'select * from ' + (dataSet['schema'] + '.' if (len(dataSet['schema'].strip()) != 0) else '') + dataSet['table']
df3 = pd.read_sql(query, con=conn)
df3.head()
```

```
Out[14]:
```

	OWNER	SYMBOL	SHARES	PRICE	TOTAL	DATEQUOTED
0	TechStock	IBM	1000	141.5	141500.0	2018-10-15



**Note:** The data virtualization gives you the ability to **Collect** your hard to reach data in a single ICP-D platform.

**\*\* End of Lab 03: Collect**

## Lab 04 Organize

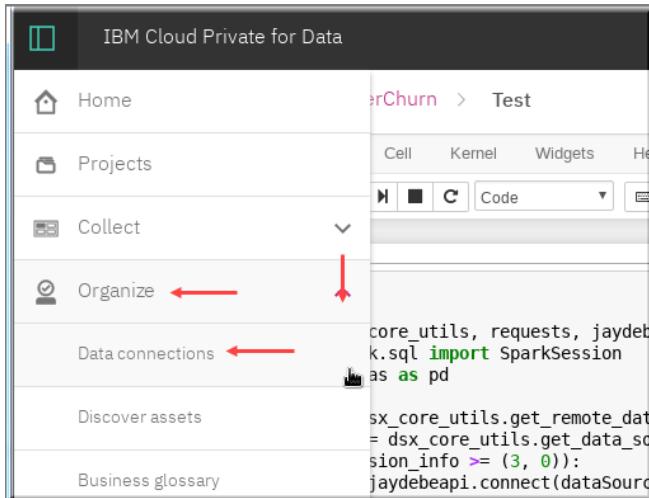
Most organizations find it difficult to understand their own data because it is dispersed across many silos and controlled by different teams.

The lab shows you how to uncover hidden data and to build a lineage that is otherwise difficult to establish. ICP-D helps you move from manual processes to establish relationships between data to automation, aided by machine learning.

Let's begin to understand how this task is simplified by ICP-D.

### 4.1 Create Connections

- \_\_1. Switch to the [ICP-D](#) Web UI.
- \_\_2. Hover your mouse on the left-menu bar and click [Organize](#), then click the down arrow and select [Data connections](#).



- \_\_3. Click [Create Connections](#).



- \_\_4. Type: Name **BLUDB**, select **Db2** from the drop-down menu and type JDBC URL exactly as shown below:

**jdbc:db2://db2whsmp-1538416690.zen.svc.cluster.local:50000/BLUDB**

- \_\_5. [No typos please.]

\_\_6. Enter Username: **db2psc** and Password: **password**.

Name \*  
BLUDB

Description  
Description

Choose connection  
Db2

JDBC URL \*  
jdbc:db2://db2whsmp-1538416690.zen.svc.cluster.local:50000/BLUDB

Username \*  
db2psc

Password \*  
.....

\_\_7. Click **Test Connection**.

\_\_8. After successful connection, click **Save Connection**.



\_\_9. Under **Organize** ⇒ **Data connections**, click **Create Connection**.

\_\_10. Enter **Name hadoop**, keep **connection HDFS**, and keep **File System WebHDFS**. Enter **Username: admin**, **Password: password**, **Host: 192.168.142.110** and **Port: 50070**.

Name \*  
hadoop

Description  
Description

Choose connection  
HDFS

File System  
WebHDFS

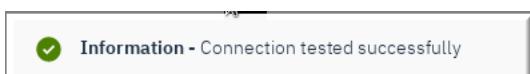
Host \*  
192.168.142.110

Port  
50070

Username \*  
admin

Password \*  
.....

\_\_11. Click **Test Connection**.



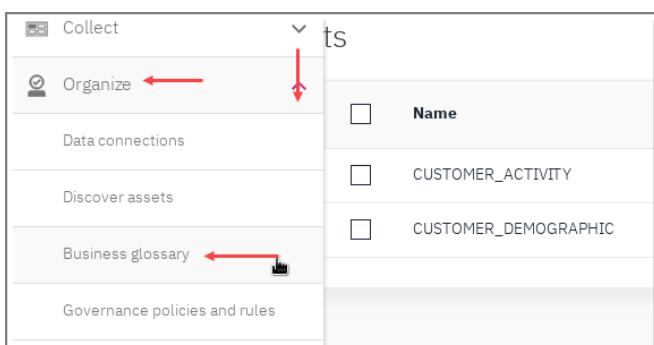
\_\_12. After successful connection, click [Save Connection](#).

\_\_13. Two connections are created.

NAME	CONNECTION TYPE
BLUDB	JDBCConnector
hadoop	HDFSFileConnector

## 4.2 Create Business Glossary

\_\_14. Navigate on the left menu to [Organize](#) ⇒ [Business Glossary](#)



i

**Note:** The business terms that we enter help users understand the data. These will automatically be associated with our data sources in a later step.

Generally speaking, these terms would either already exist within a client's environment, or there would be a governance board tasked with creating a glossary, prior to onboarding new data sources.

\_\_15. Click [Create Category](#).

Categories
Terms

0 Categories available

NAME	CREATED	MODIFIED
No Categories found		

Import Categories
Create Category

- \_\_16. Type Name: **Customer Churn** and description: **Contains terms related to customer churn.**

Create Category

**Name\***

Customer Churn ←

241

**Parent Category**

Type to find and add → Keep blank as this is the parent

**Short Description**

Contains terms related to customer churn ←

215

- \_\_17. Click **Save** [Towards right bottom of the screen.]
- \_\_18. Go to the second tab **Terms** and click **Create Term**.

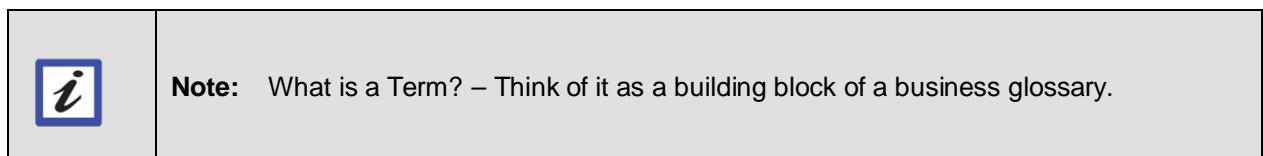
Categories Terms ←

Search for Terms in the catalog

0 Terms available

Import Terms Create Term

NAME	CREATED	MODIFIED
------	---------	----------



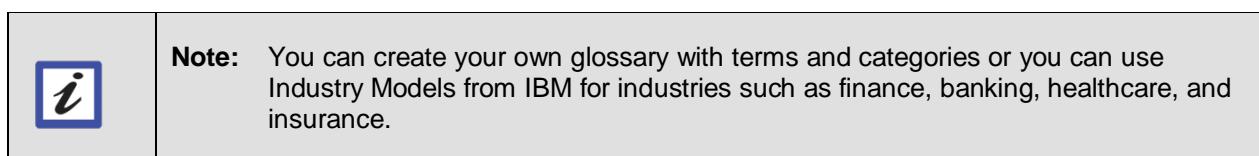
- \_\_19. Create a term **Income** with a parent category **Customer Churn** (Type **C** in the box), provide a description as: **Yearly income of the customer** and choose **Status** as **Standard** from the drop-down menu.

Name\*  
Income 1

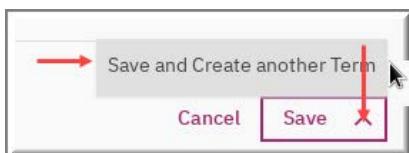
Parent Category\*  
Customer Churn 2

Short Description  
Yearly income of the customer

Status\*  
Standard 3



- \_\_20. Click the down arrow next to **Save** and click **Save and Create another Term**.



- \_\_21. Now, repeat above exercise to create the following **Terms**. [Remember to use **Save and Create another Term**.] (The description for each term is found in the next table.)

Name	Category	Status
Home Owner	Customer Churn	Standard
Net Realized Losses	Customer Churn	Standard
Net Realized Gains	Customer Churn	Standard
Gender	Customer Churn	Standard
Days Since Last Trade	Customer Churn	Standard
Total Dollar Value Traded	Customer Churn	Standard

- \_\_22. Take the **Description** from the second column for the **Term** being entered.

Name	Description
Home Owner	Flag indicating whether the customer owns a home.
Net Realized Losses	The net dollar value of losses realized for a customer, usually a year to date metric.
Net Realized Gains	The net realized gains for a customer, usually a year to date metric.
Gender	The gender of the customer, either 'M' or 'F'
Days Since Last Trade	The number of days since the customer last executed a trade.
Total Dollar Value Traded	The total amount the customer has traded with us since onboarding.

- \_\_23. For the last term, just click **Save** instead of Save and Create another Term
- \_\_24. Notice the categories **Customer Churn** that you just created and the **Terms** associated with that category. Review the **Categories** and **Terms** tab.



#### 4.3 Create Governance Policies and Rules

- \_\_25. Navigate on the left menu to **Organize** ⇒ **Governance policies and rules**
- \_\_26. Click **Policies** and **Create Policy**.



- \_\_27. Type Name: **Net Gains and Net Losses Are Mutually Exclusive** and Description: **A customer can only have a value in Net Gains or Net Losses, but not in both.**

The screenshot shows a 'Create Policy' form. The 'Name\*' field contains 'Net Gains and Net Losses Are Mutually Exclusive' with a character count of 208. The 'Short Description' field contains 'A customer can only have a value in Net Gains or Net Losses, but not in both.' with a character count of 178. A red arrow points to the end of the name entry, and another red arrow points to the end of the short description entry.

- \_\_28. Click **Save**.

- \_\_29. We just created a policy and now create a rule. Click **Create Rule** under **Rules**.

The screenshot shows a 'Rules' catalog page. At the top, there are tabs for 'Rules' (which is selected) and 'Policies'. Below the tabs is a search bar with the placeholder 'Search for Rules in the catalog'. In the center, it says '0 Rules available'. On the right side, there are buttons for 'Import Rules' and 'Create Rule'. A red box highlights the 'Create Rule' button, and a red arrow points to it from below.

- \_\_30. Type Name: **Net Realized Gains and Losses Validity Check**, type 'n' to select policy and type Description as **If Net Realized gains > 0 then Net Realized Losses = 0 else if Net Realized Losses > 0 then Net Realized Gains = 0.**

The screenshot shows a 'Create Rule' form. The 'Name\*' field contains 'Net Realized Gains and Losses Validity Check' with a character count of 211. The 'Referencing Policies' section shows a single policy 'Net Gains and Net Losses Are Mutually Exclusive' selected. The 'Short Description' field contains 'If Net Realized gains > 0 then Net Realized Losses = 0 else if Net Realized Losses > 0 then Ne' with a character count of 141. A red arrow points to the end of the name entry, a red arrow points to the selected policy in the referencing policies list, and a red arrow points to the end of the short description entry.

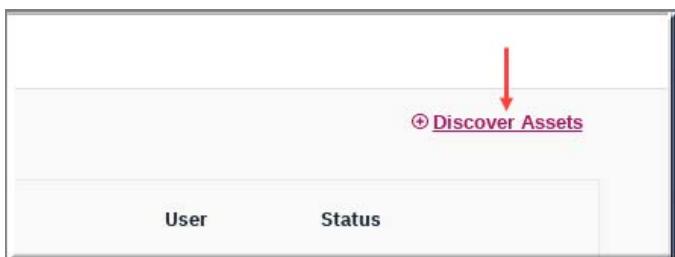
- \_\_31. Click **Save**.

## 4.4 Discover Assets

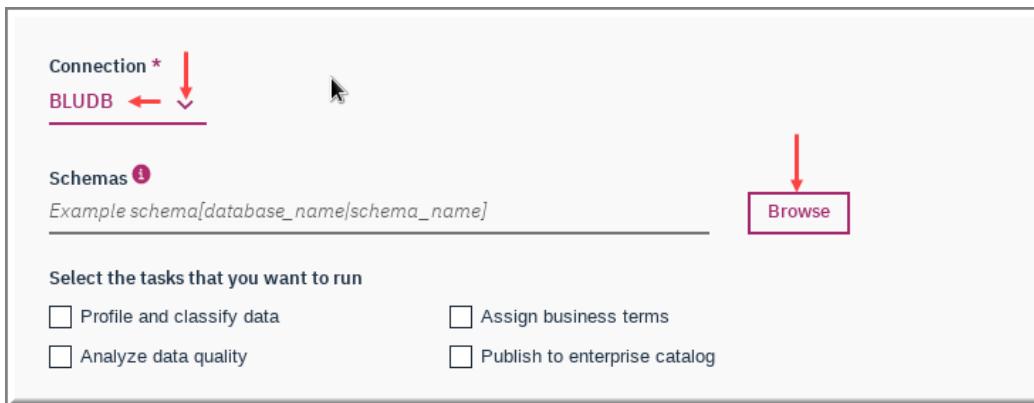
- \_\_32. Click **Organize** and expand it and click **Discover assets**.



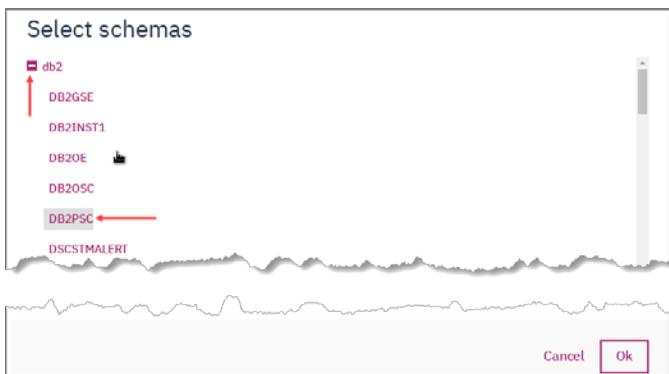
- \_\_33. Click **Discover Assets**. [Towards top right]



- \_\_34. Select **BLUDB** from the drop-down menu and click **Browse**.



- \_\_35. Expand **db2** and click **DB2PSC** to select it and click **Ok**.



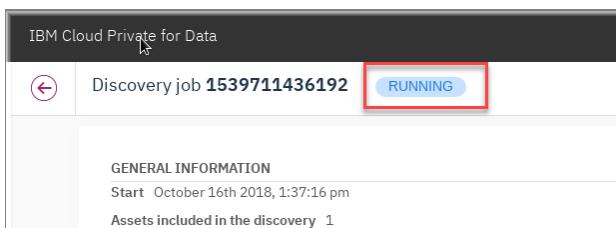
- \_\_36. Check the boxes for **Profile and classify data**, **Assign business terms**, **Analyze data quality**, and **Publish to enterprise catalog**.



- \_\_37. Click **Discover**. [Far bottom right of the screen.]



- \_\_38. A job is created and begins to execute. It may take few minutes to complete.



- \_\_39. Click **Refresh** to see if the job is still running and wait for this to complete.



- \_\_40. The Discovery phase has two stages – 1. **Import** and 2. **Analyze**.

- \_\_41. When it is finished, you will see that the 'Analyze' phase has a flag marked **FINISHED**.

- \_\_42. As you refresh the page, you can see the percentage complete in each phase.

DISCOVERED ASSETS INFORMATION			
Asset name		Asset type	Status
			Phase IMPORT
DB2PSC	Schema	Start October 17th 2018, 8:49:13 am End October 17th 2018, 8:49:23 am <span style="background-color: #2e7131; color: white; border-radius: 15px; padding: 2px 5px;">FINISHED</span>	Phase ANALYZE Done 100.00% Successful 100.00% Cancelled 0.00% Failed 0.00% Start October 17th 2018, 8:49:39 am End October 17th 2018, 8:54:44 am <span style="background-color: #2e7131; color: white; border-radius: 15px; padding: 2px 5px;">FINISHED</span>

- \_\_43. Wait for this to complete. It may take 4-6 minutes.

	<p><b>Note:</b> What just happened? – An automated discovery, assignment a quality score and assignment of terms if there was a match with a TERM that we created.</p> <p>This process is incredibly powerful. What just happened is that, simply by pointing to a data source[ML1]:</p> <p>ICP - D has automatically profiled the data to highlight data quality, classified the data so that the business purpose is easily understood, that is, do columns contain PII data? Are some columns more significant to our analytics project than others?</p> <p>And finally, it has automatically assigned these data columns to business terms that we previously added to our glossary. This allows business users to easily find or 'shop' for data that they want to research or exploit.</p>
---	--

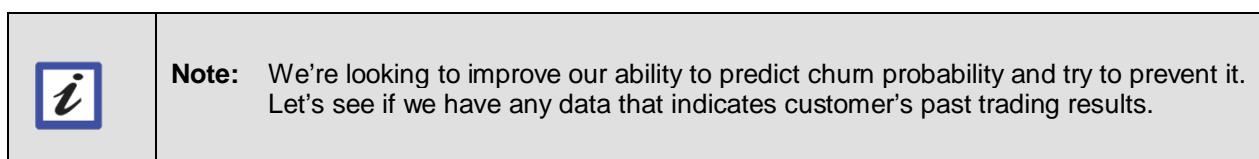
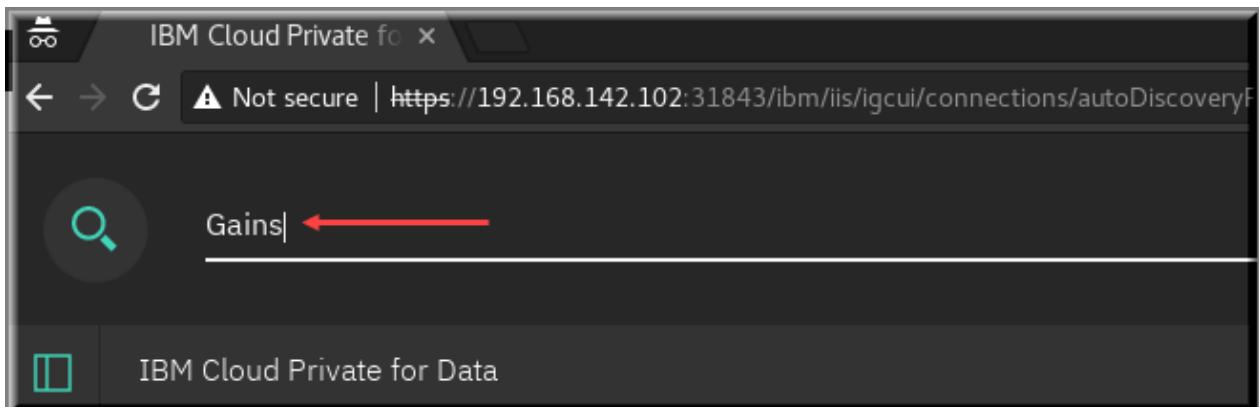
- \_\_44. When it completes, you will see that the **Analyze** phase has a flag marked **FINISHED**.
- \_\_45. Click **Review** under *Review results*.
- \_\_46. Expand **CUSTOMER\_ACTIVITY** and notice the quality score assigned to a table, individual data items and match with a known data class and possible match with a Term.

Assets					
	Name	Quality	Data class	Assigned terms	Actions
>	<input type="checkbox"/> CUSTOMER_ACTIVITY	88.76 %	-	Income 40.41% ← +1 Suggested	  
>	<input type="checkbox"/> CUSTOMER_DEMOGRAPHIC	99.86 %	-	Income 40.41% ← +1 Suggested	  

- \_\_47. The discovered assets are published to the catalog so that they are searchable.
- \_\_48. Click **Search** [Middle icon on the top menu bar towards the right.]



\_\_49. Type **Gains** and click **Enter**.

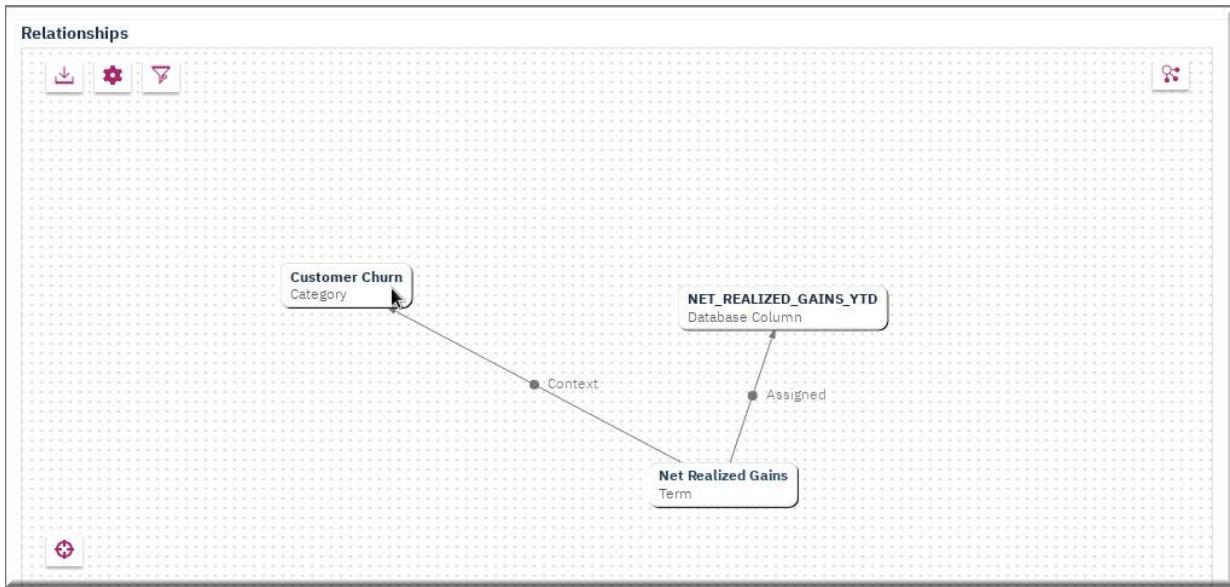


\_\_50. You see number of results and one of them is **Net Realized Gains**.

TERM	GROUP	CLASS	LAST MODIFIED
Net Realized Gains	Glossary and Governance	GlossaryExtensions.BusinessTerm	10-17-2018

\_\_51. Click **Relationship Graph**.

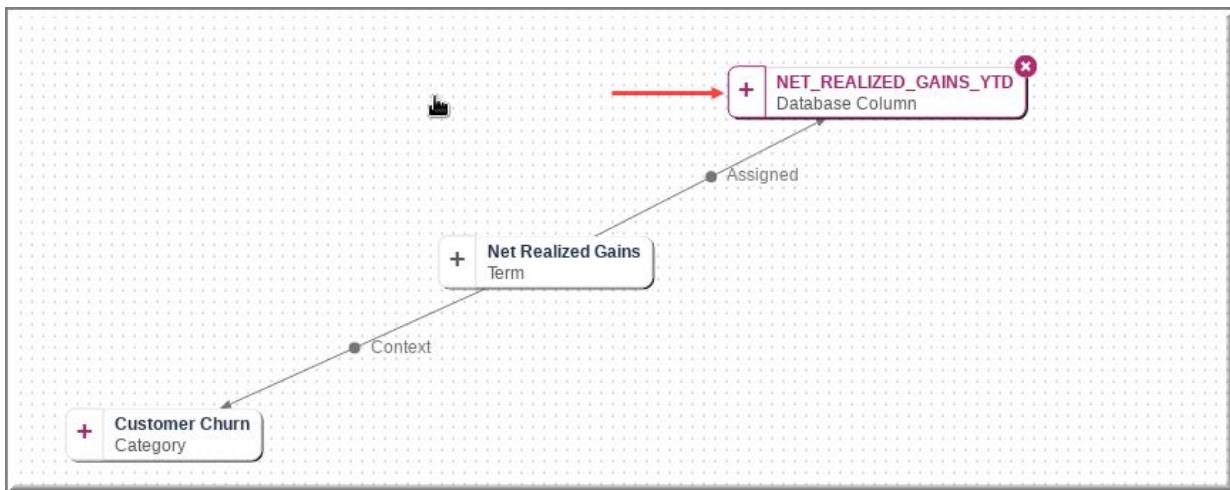
\_\_52. The relationship graph between **Category**, **Term** and the associated **Database Column** displays.



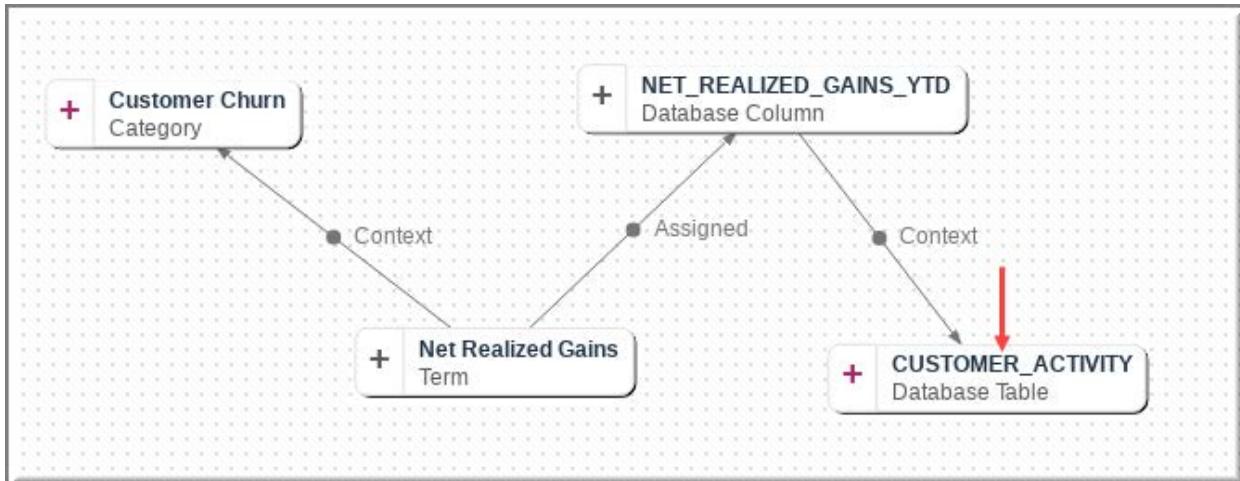
\_\_53. Click **Explore relationships** icon at the top-right corner of the screen.



\_\_54. Hover your mouse over **NET\_REALIZED\_GAINS\_YTD** and click + sign.



- \_\_\_55. What did you just learn? A visual relationship matrix between a business TERM and the database table and column it finds its association.



	<b>Note:</b> Deriving this relationship manually is a very cumbersome, labor-intensive job, should you have thousands of tables. Let ICP-D platform do this for you with a degree of confidence that you can see, validate and collaborate with comments and star ratings.
--	--

- \_\_\_56. Click the **CUSTOMER\_ACTIVITY** box, to display the details on the right-side with the connection information.
- \_\_\_57. A solid discovery and classification process of your data is a recipe to find your **dark data** and build machine learning models that we will cover in next lab exercise.

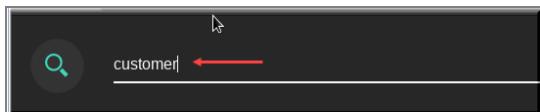
	<b>Note:</b> Gartner defines <b>dark data</b> as the information assets organizations collect, process and store during regular business activities, but generally fail to use for other purposes (for example, analytics, business relationships and direct monetizing).
--	---

## 4.5 Shopping for Data

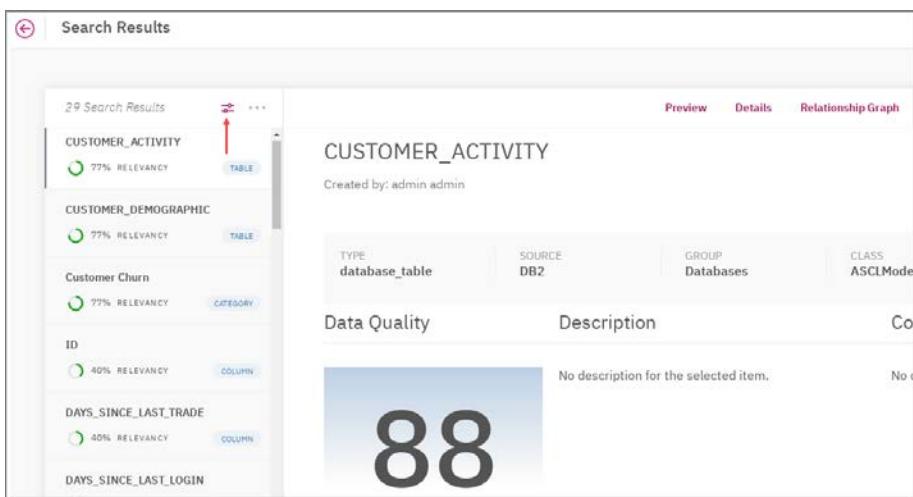
\_58. Click **Search** icon.



\_59. Type **customer** and click **Enter**.

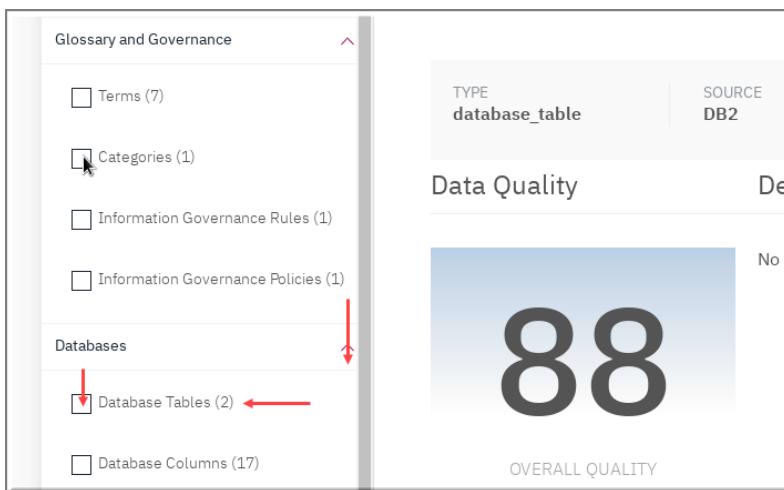


\_60. Assets related to **customer** display. Click



\_61. Expand **Databases**.

\_62. Check **Database Tables**.



- \_\_63. Click **Properties** icon, **CUSTOMER\_ACTIVITY** and then click **Add to cart**.

**CUSTOMER\_ACTIVITY**

Created by: admin admin

TYPE database\_table SOURCE DB2 GROUP Databases CLASS ASCLModel.DatabaseTable LAST MODIFIED 10-17-2018

Data Quality Description Comments

- \_\_64. In the same way, click **CUSTOMER\_DEMOGRAPHIC** and click **Add to cart**.



**Note:** We've now chosen two tables to extract. From here, we build an ETL job to extract, join and transform this data. We will also add previous customer churn data from Hadoop to create our final data set for analysis.

## 4.6 Transform Data

- \_\_65. Click **Transform**.

**CUSTOMER\_DEMOGRAPHIC**

Created by: admin admin

TYPE database\_table SOURCE DB2 GROUP Databases CLASS ASCLModel.DatabaseTable LAST MODIFIED 10-17-2018

Data Quality Description Comments

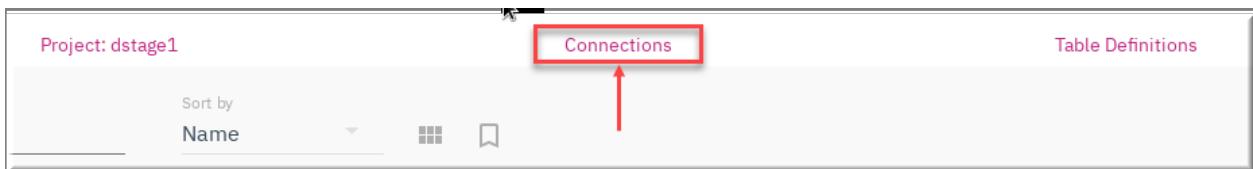
- \_\_66. A list of ETL projects displays, click **dstage1**.

Name	Collaborators
ANALYZERPROJECT	(AA, AI)
DataClick	(AA, AI)
<b>dstage1</b>	(AA, AI)



**Note:** After a few seconds, the screen repaints and displays the ETL canvas, with our two data sources that we just added to our cart placed onto the canvas, so that we may further develop our ETL job.

\_67. Click **Connections**



\_68. Click **Down arrow** [toward the right.]

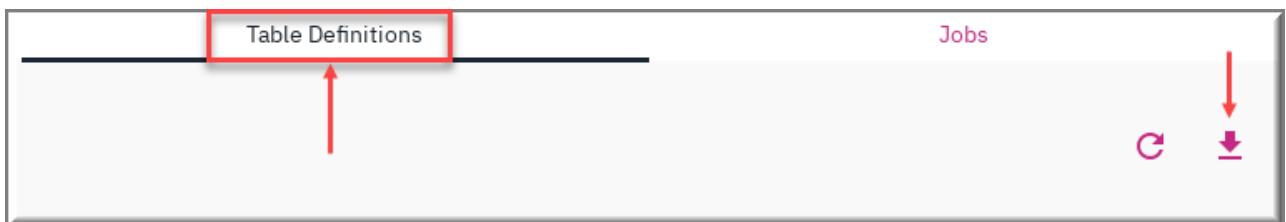


\_69. Check **hadoop** data source and click **Import**.

\_70. You should see two connections imported.

Name	Description	Modified on	Category
BLUDB	Imported from BLUDB on 2018-10-1	2018-10-17 16:34:16	\Stage Types\Parallel\Database
hadoop	Imported from hadoop on 2018-10-1	2018-10-17 16:34:19	\Stage Types\Parallel\Database

\_71. Next click **Table Definitions** and click **Down arrow**.



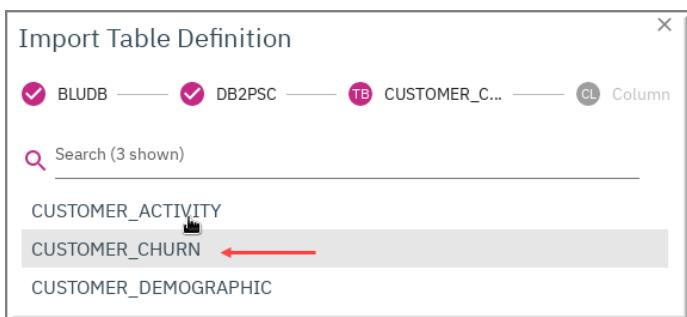
\_72. The connection **BLUDB** is selected, click **Next**.



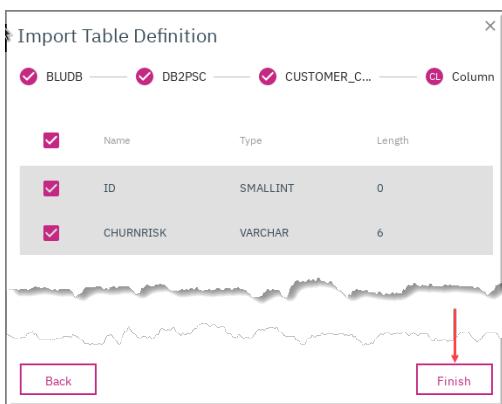
\_\_73. Select **DB2PSC** schema and click **Next**.



\_\_74. Select **CUSTOMER\_CHURN** and click **Next**.



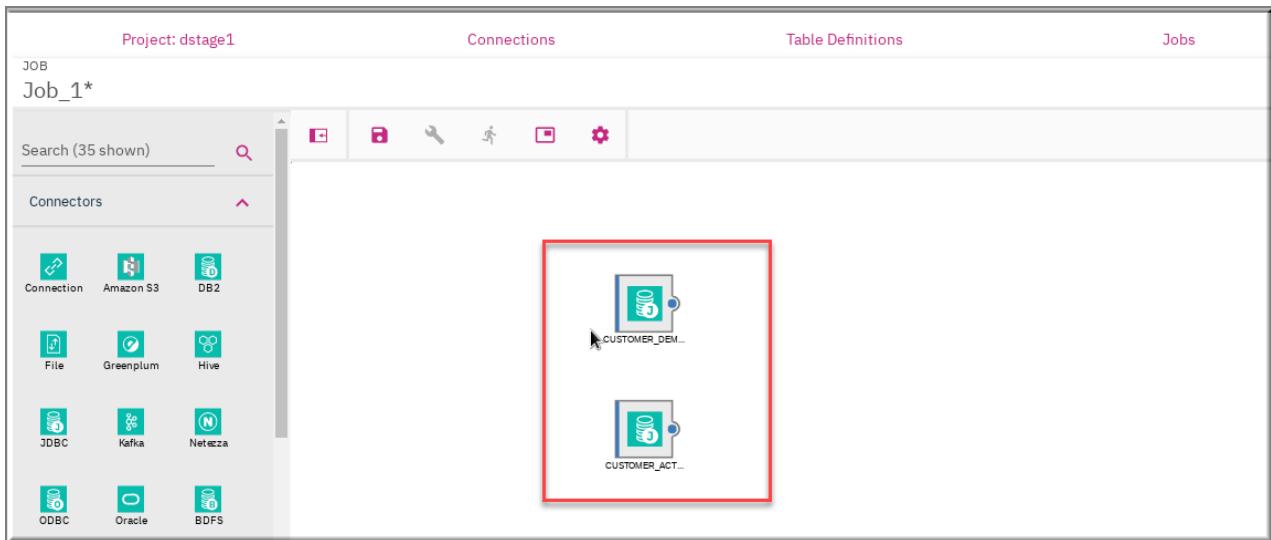
\_\_75. Column names from the table display. Click **Finish**.



\_\_76. Click **Save**.

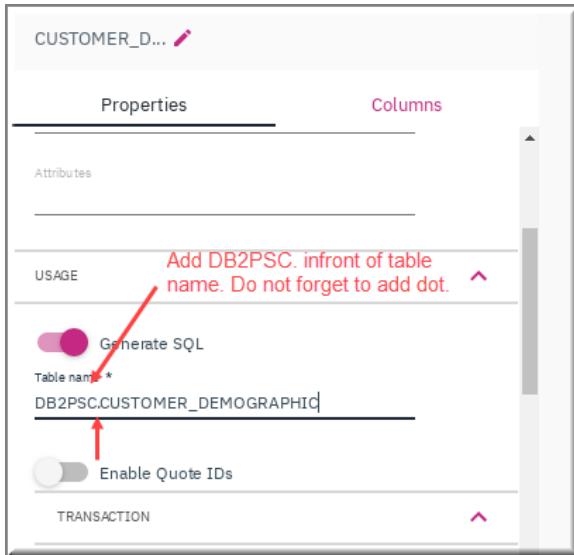
\_\_77. Click **Jobs\_1\***.

- \_\_78. Since we had two tables **CUSTOMER\_ACTIVITY** and **CUSTOMER\_DEMOGRAPHICS** already in the shopping cart, those tables are added automatically to the **Job\_1\***.

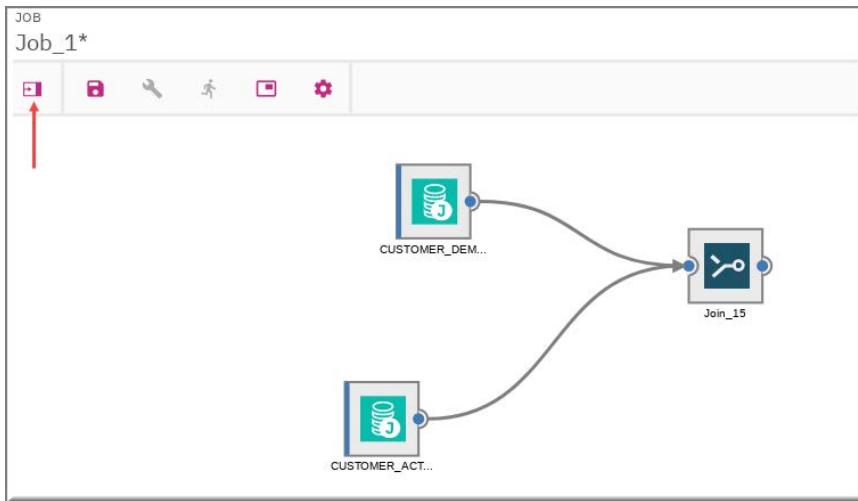


- \_\_79. Click the **CUSTOMER\_DEMOGRAPHICS** icon on the canvas.
- \_\_80. Navigate to the palette on the left and do the following:
- Find the **Join** icon (in the Stages section.)
  - Click it once
  - Hold your mouse button down and drag the stage to the canvas
  - A link (shown as a grey line) should be drawn from **CUSTOMER\_DEMOGRAPHIC** to the join stage. [If you do not see the link being drawn, remove the Join and ensure that you click **CUSTOMER\_DEMOGRAPHIC** first.] When you drag a new stage onto the canvas (Join, in this case), the interface automatically draws a link from whatever stage on the canvas is highlighted (in this case **CUSTOMER\_DEMOGRAPHIC**) to the stage you've just placed it on.
  - Next, click the **CUSTOMER\_ACTIVITY** stage
  - Click and hold on the stage and drag a link to the join stage
- \_\_81. You should see the join between **CUSTOMER\_DEMOGRAPHIC** and **CUSTOMER\_ACTIVITY**.
- \_\_82. If you do not get it right it the first time, delete the **Join** from canvas and try again.
- \_\_83. Double-click **CUSTOMER\_DEMOGRAPHIC** to edit some attributes. The right pane will slide-in.

- \_\_84. Go to the **Table name** field and add the schema name **DB2PSC.** (include the **dot** after) in front of the table name **CUSTOMER\_DEMOGRAPHIC.**

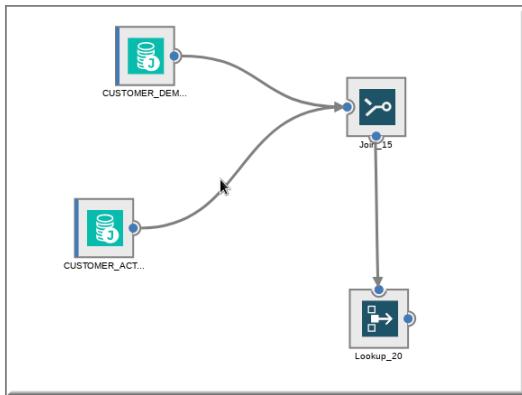


- \_\_85. Scroll all the way down and click **OK**.
- \_\_86. Double-click **CUSTOMER\_ACTIVITY** and add **DB2PSC.** in front of the table name and click **OK**.
- \_\_87. Remember – you can hide/unhide the palette to use different stage artifacts, using the icon shown by the red arrow below.



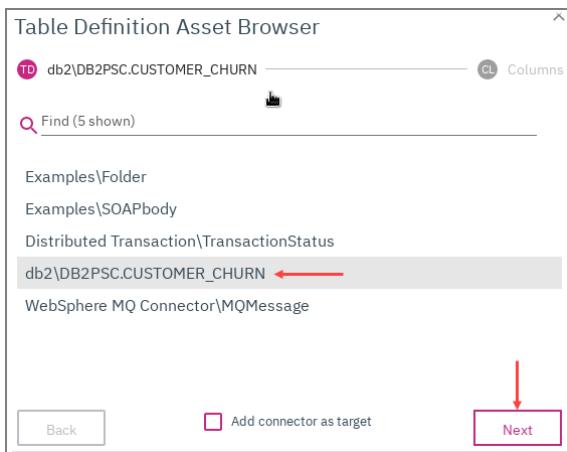
- \_\_88. Click **Join** on the canvas to select it. [Make sure you select it before proceeding to the next step.]
- \_\_89. On the palette, click **Lookup** to select and then double-click it. This draws the link between **Join** stage and the **Lookup** stage.

- \_\_90. Re-arrange the different stages so that **Lookup** stage is underneath the **Join** stage.

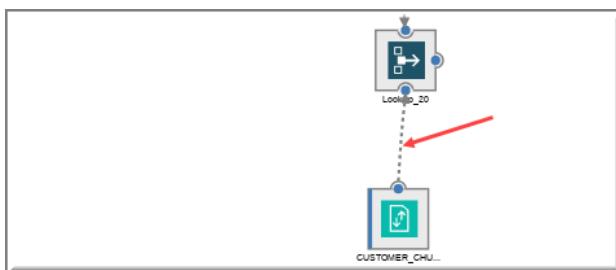


#### 4.6.1 Join with lookup Data from HDFS

- \_\_91. **Make sure you have not selected any node.** In the Connectors palette, click **File** (not sequential file) and drag it below the **Lookup** stage. This opens a window.  
 \_\_92. Select **CUSTOMER\_CHURN** and click **Next**.



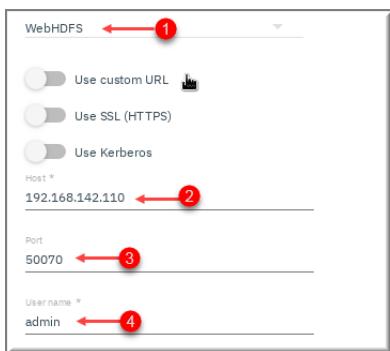
- \_\_93. Click **Add to Job**.  
 \_\_94. The **File** stage has a little circle link on the right side. Click that little circle and the cursor will change to a line and then click **Lookup** stage to draw the connection between two.  
 \_\_95. This should be a dotted line as shown – which means that it is a *Reference*.



\_\_96. Double-click this File stage to open edit box on the right.

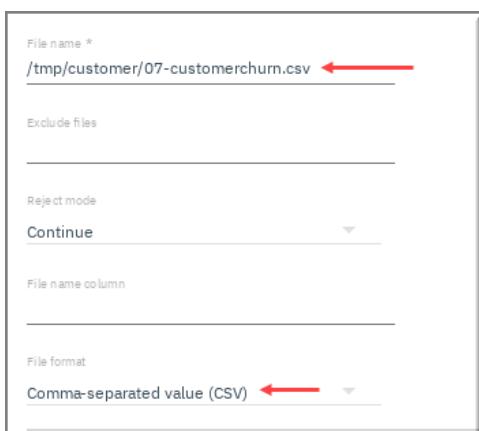
\_\_97. Under the Connection property:

- \_\_a. Select **WebHDFS** from the drop-down menu.
- \_\_b. Type Host: **192.168.142.110**
- \_\_c. Port: **50070**
- \_\_d. Username: **admin**



\_\_98. Under Usage property:

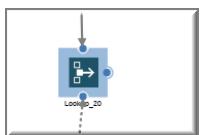
- \_\_a. Type file name: **/tmp/customer/07-customerchurn.csv** [all lower case.]
- \_\_b. Select file format **Comma-Separated value (CSV)** from the drop-down menu.
- \_\_c. Slide the slider towards right to indicate that the **First row is header**.



\_\_99. Scroll down and click **OK**.

	<b>Note:</b> Now we need to do a few things to finish the properties in our job. First, we need to let ICP-D know what columns to use from Hadoop as a key to match on.
---	---

\_\_100. On the canvas, click the **Lookup** stage to select it.

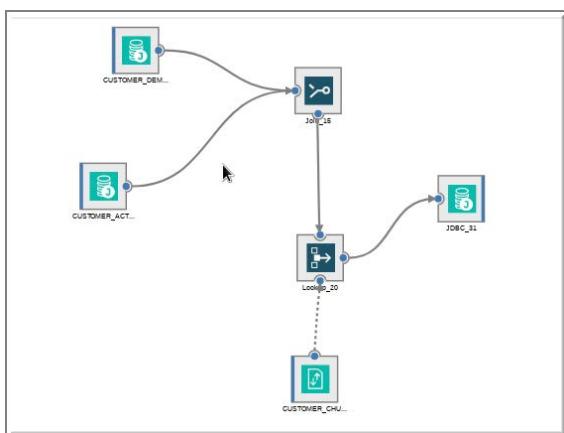


\_\_101. From the palette, drag a **Connection** stage and drop it to the right side of the lookup stage.

\_\_102. Check **Add selected connection as target**. Click **Add to Job**.



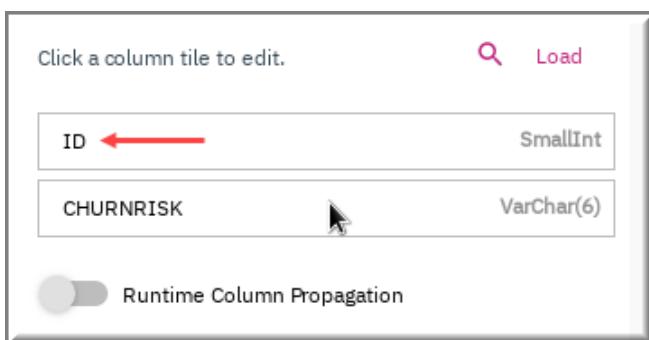
\_\_103. This should link to the **Lookup** stage. Your flow should display the screenshot below.



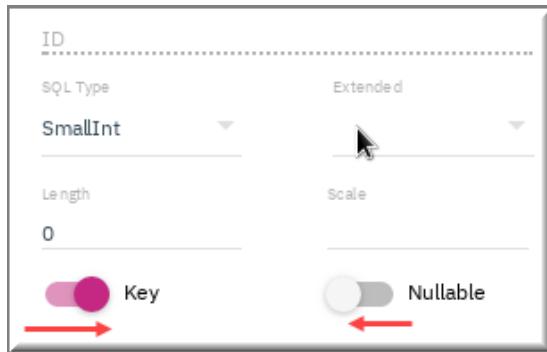
\_\_104. Double-click **CUSTOMER\_CHURN** stage.

\_\_105. Select **Columns** tab from the right-hand pane.

\_\_106. Click the **ID** field.



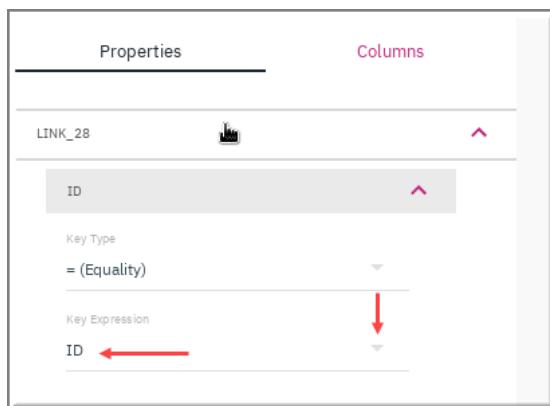
\_107. Slide the **Key** slide to the right to make it a key and slide the **Nullable** slide to the left to make it non-nullable.



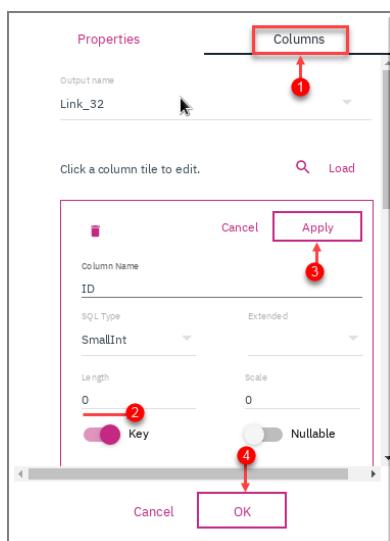
\_108. Click **Apply** in the *column tile* and then click **OK**.

\_109. Double-click the **Lookup** stage.

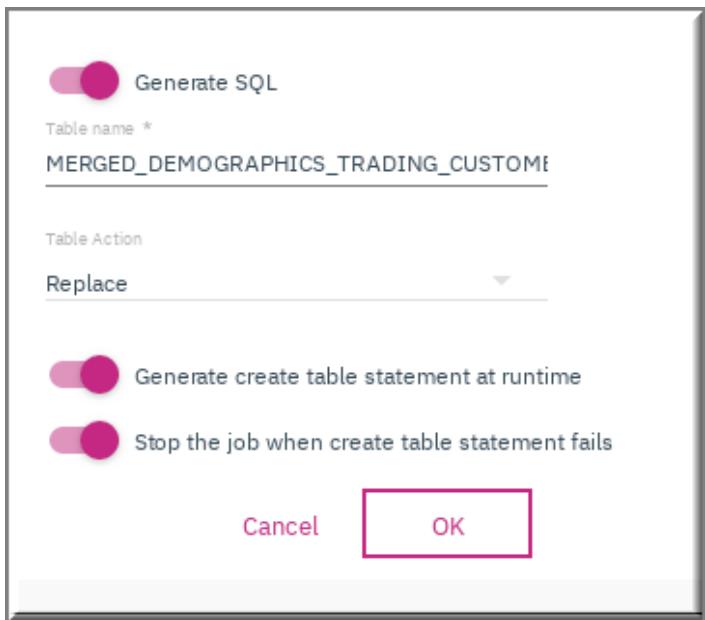
\_110. On **Properties**, click key expression drop-down menu and select **ID** to join.



\_111. Click **Columns** and click **ID field**, slide key slider to the right, click **Apply** and the **OK**.

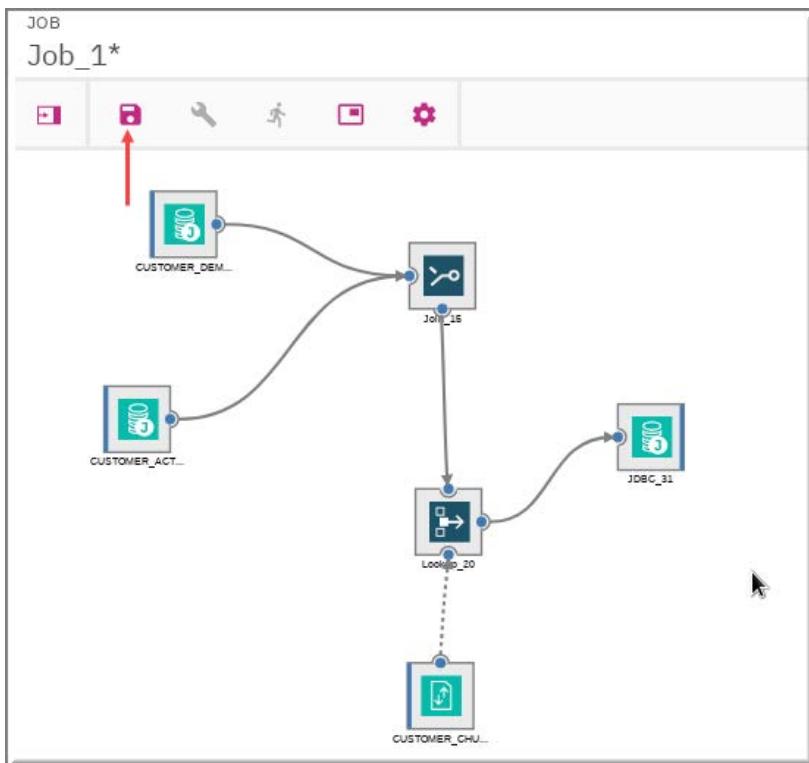


\_112. Double-click the Target stage [Look for the name as JDBC\_XX]. Navigate down to the table name and type MERGED\_DEMOGRAPHICS\_TRADING\_CUSTOMER and select table action Replace from the drop-down menu. Click OK to save.

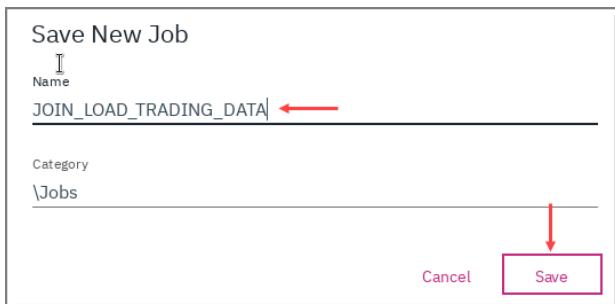


#### 4.6.2 Save, Compile and Run Transformation Job

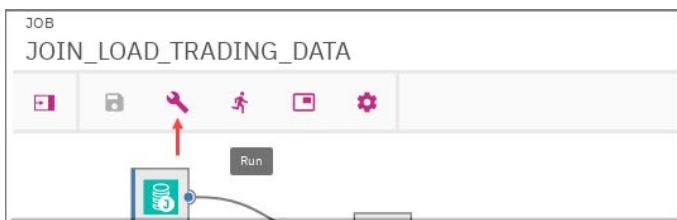
\_113. Click **Save**.



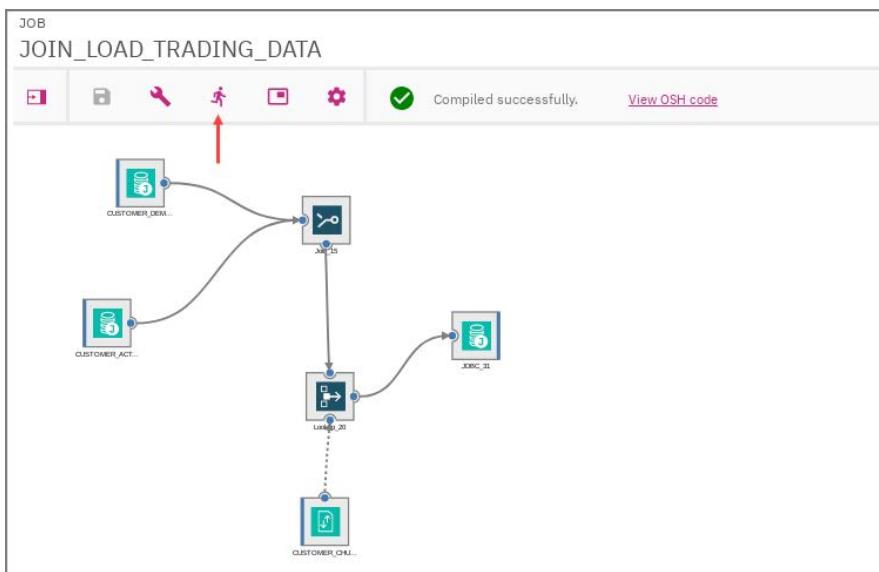
\_114. Type name: JOIN\_LOAD\_TRADING\_DATA and click Save.



\_115. Click the wrench to compile the job.



\_116. Click the runner icon to Run the job.



\_117. Click Run to confirm the run.

\_118. It may take a minute or two to run the job. You can click the refresh icon to see the status.

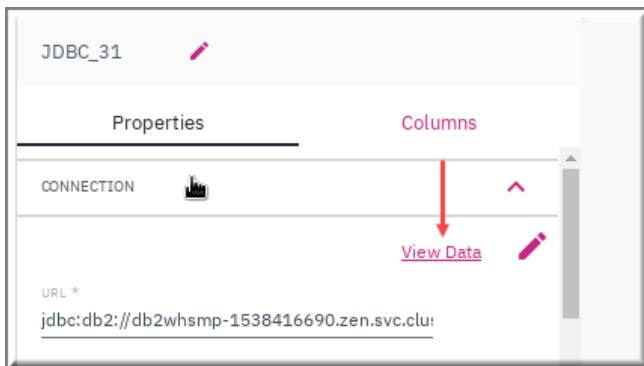


\_119. Click Refresh.



\_120. Double-click Target stage [JDBC\_XX] to open detailed pane on right side.

\_121. Click View Data.



\_122. The data transformation job between two Db2 Warehouse tables and data from Hadoop/HDFS was joined and created new table in Db2 Warehouse.

**View Data**

ID	GENDER	STATUS	CHILDREN	ESTIMATED_INCOME	HOME_OWNER	AGE	TOTAL_DOLLAR_VALUE_TRADED	TOTAL_UNIT
1	M	M	2	29416.00	N	49	29782.98	45
2	M	M	0	19732.80	N	51	24812.48	22
3	M	S	2	96.33	N	56	26132.61	32
4	F	M	2	52004.80	N	25	5030.50	23
5	M	M	2	53010.80	N	19	12451.25	46
10	F	M	2	47902.00	N	26	17441.74	49
11	M	M	1	7545.96	Y	17	22392.24	178
12	F	S	0	78851.30	N	48	370.04	28
13	F	S	1	17540.70	Y	63	22172.22	13
14	F	M	0	83891.90	Y	61	28922.89	45

[Copy to Clipboard](#) OK

—123. We will use Analyze the data from this table in next lab exercise. We will build a business dashboard from ML2 and then build machine learning models that we will use to modernize our application from analytics perspective.

	<p><b>Note:</b> We have seen the value in applying a glossary, terms and classification to make data searchable so that data scientist, data engineers or business analyst can shop for data.</p> <p>The ICP-D platform reduces the time-consuming data organization task, making it easier for data scientist. We will see that value in next lab exercise.</p> <p>This lab has shown you how to automatically discover data sources, automatically classify those sources with business classifications using ML methods, automatically assign business terms to those data sources using ML and fuzzy matching methods.</p> <p>It shows how to shop for specific data, and join and move that data to our analytics warehouse for further use by data scientists.</p> <p>The steps covered here used to take months, and sometimes years, to complete using traditional manual methods.</p> <p>ICP-D automates them so that you can accelerate the time to value of your analytics projects.</p>
---	---

Credit for Lab: Jan Van Buren, Business Unit Exec, NA Info Integration and Governance Technical Sales

**\*\* End of Lab 04: Organize**

## Lab 05 Analyze

When embarking on machine learning projects, most organizations immediately engage their data scientists hoping to gain insight into their data. They quickly realize that the data that they need is not available in the format they require.

[Analyze](#) is our third pillar in our Integrated Data Platform. This is where data engineers (who prepare data), business analysts (who make sense out of data), data scientists (who build machine learning and artificial intelligence models) join forces to gain insights from their data.

The data stewards help [Organize](#) the data (as shown in the previous lab) through discovery, business glossary, policy, rules and governance. This is the foundation on which business analysts and data scientists build.

Now, let's look at the role – business analysts and data scientists would play in using IBM Cloud Private for Data to accomplish Analytics Modernization.

### 5.1 Business Analystist

We will assume the role of two personas who are key to driving business decisions. Let's start by examining the role that a business analyst plays in the discovery of business problems. We will seek solutions in later sections.

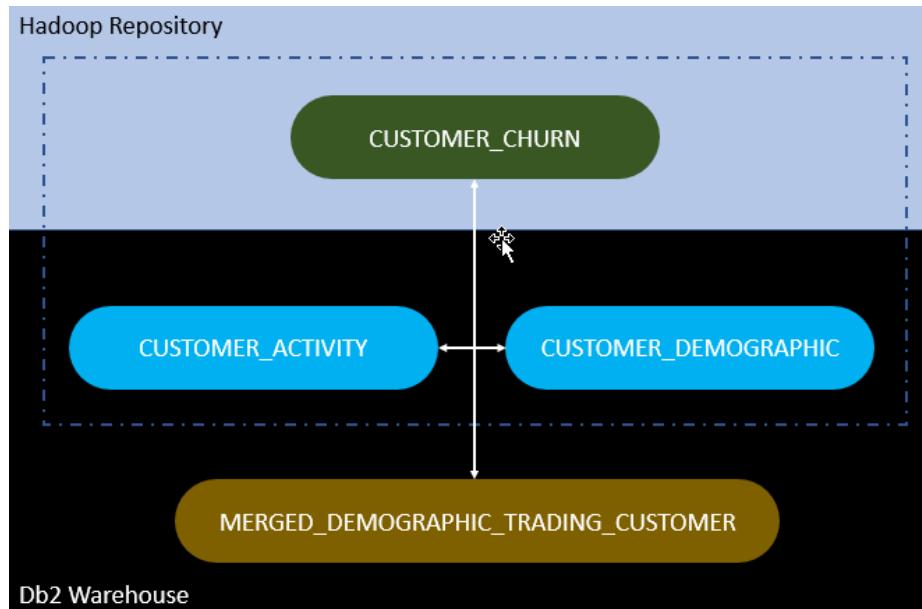
#### 5.1.1 Dashboard – Identify Focus Area

- 1. In our Executive Demo lab, you were presented with Stock Opening Bell business analysis that painted the picture of how the business is doing. You noticed a gradual decline in revenue and flat share volume.
- 2. Refer to Appendix-'A' if you would like to learn how to build this dashboard. You will find step-by-step directions to analyze current trends of customer visits and daily trade volume generated through the Stock Trader application.
- 3. This is our Stock Trade Opening Bell dashboard derived from the data.



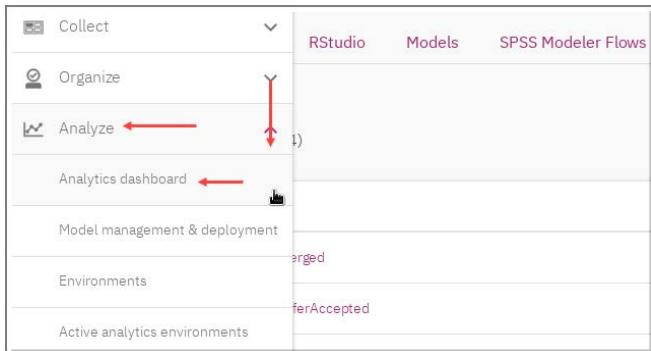
### 5.1.2 Data Preparation

- 4. In previous lab, we collected Customer Risk data. In our fictitious scenario, the Customer Risk legacy application assigns a **Low**, **Medium** and **High** risks to customers based on the company's past complaint history. This data is stored in a Hadoop repository.
- 5. In that lab, you'll recall that we joined **CUSTOMER\_ACTIVITY** and **CUSTOMER\_DEMOGRAPHIC** data from Db2 Warehouse with **CUSTOMER\_CHURN** data from Hadoop and produced the **MERGED\_DEMOGRAPHICS\_TRADING\_CUSTOMER** table.



- 6. The table **MERGED\_DEMOGRAPHIC\_TRADING\_CUSTOMER** has data that business analysts will visualize with the business dashboard. Then, the data scientists will apply ML/AI capabilities over the data.
  - 7. Let's begin with business analysts role to analyze customer demographics.
- ### 5.1.3 Business Dashboard – Analyze Customer Demographics
- 8. Click **Projects** from the left menu bar.
  - 9. Click **TradingCustomerChurn** project.

\_\_10. Click **Analyze**  $\Rightarrow$  **Analytics Dashboard** from the left menu bar.



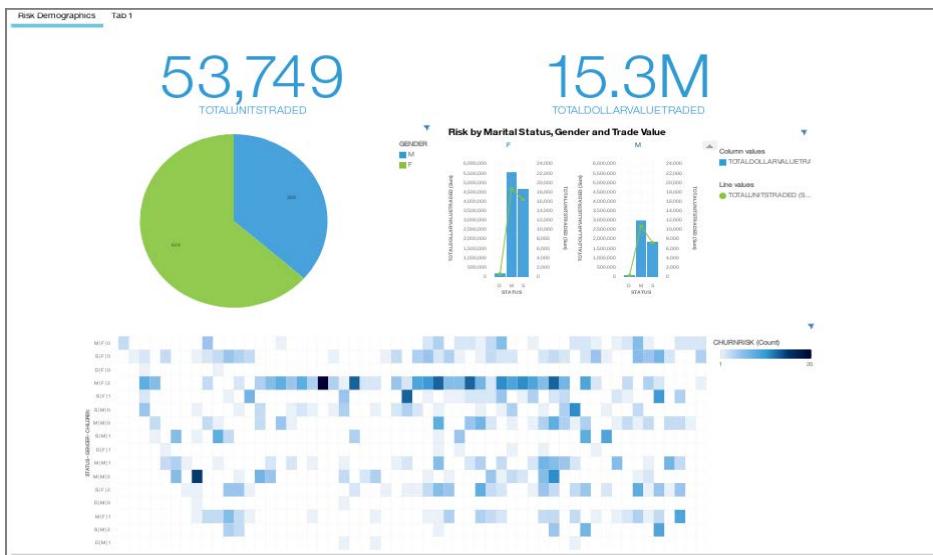
\_\_11. We pre-built **02-Stock-Trader-Demographic-Discovery** dashboard.

A screenshot of the dashboard list interface. At the top, there is a search bar labeled 'Search by dashboard'. Below it, a message says '3 dashboards available'. On the right, there is a 'Create dashboard' button. The list shows three entries:

NAME	PROJECT	MODIFIED
01-Stock-Trade-Opening-Bell.json	TradingCustomerChurn	17 Oct 2018, 8:40 AM
03-Stock-Trader-Closing-Bell.json	TradingCustomerChurn	17 Oct 2018, 8:40 AM
<b>02-Stock-Trader-Demographic-Discovery.json</b>	TradingCustomerChurn	17 Oct 2018, 8:40 AM

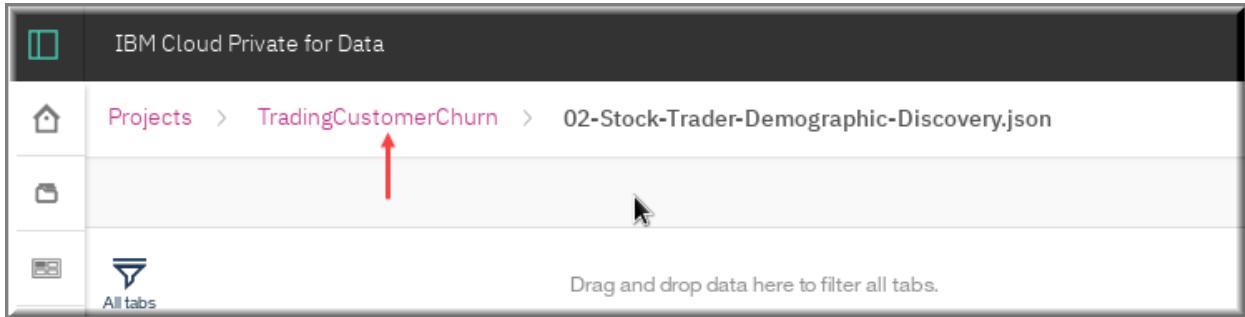
A red arrow points to the '02-Stock-Trader-Demographic-Discovery.json' entry in the list.

\_\_12. Click **02-Stock-Trader-Demographic-Discovery** to open the dashboard.



\_\_13. Notice the metrics that are present in this dashboard. **1.** Know your customers (male or female), **2.** Separation Risk by marital status, gender and trade value are used to generate a heat map to help us understand which type of customers are more likely to separate from us.

- \_\_14. This dashboard will paint a high-level picture of the customer demographic we need to target to reverse our declining revenue trend.
- \_\_15. The dashboards allow you to apply a filter **Low**, **Medium** and **High** to the Risk score to help you understand the data.
- \_\_16. Let's build this dashboard from the beginning.
- \_\_17. Click Project **TradingCustomerChurn**.



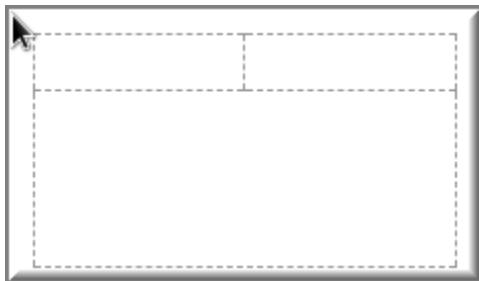
- \_\_18. Click **Assets**, then click **Dashboard** and then click **add dashboard**.



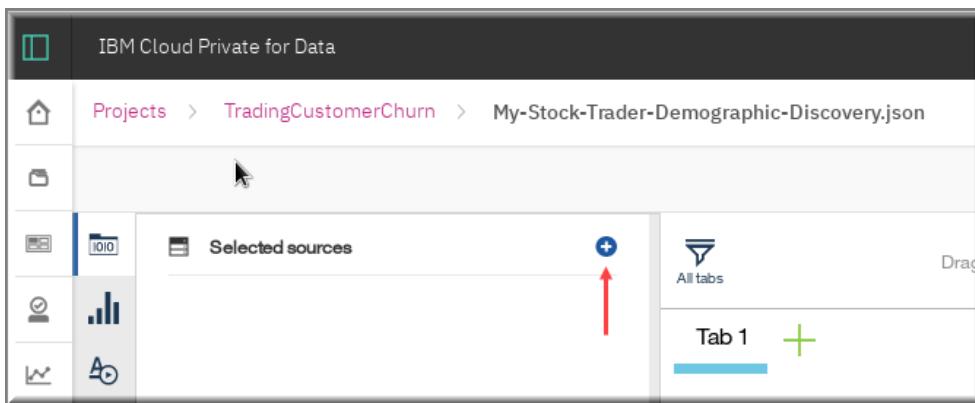
- \_\_19. Enter **My-Stock-Trader-Demographic-Discovery** in the Name field and click **Create**. [bottom right.]

The screenshot shows the "Create Dashboard" dialog box. At the top, it says "Create Dashboard". Below that, there are two options: "Blank" and "From File", with "From File" being selected. A red arrow points down to the "Name\*" field. The name "My-Stock-Trader-Demographic-Discovery" is entered into the field. To the right of the field, there is a green checkmark icon and the number "13". Below the name field, a message says "This name is valid". Underneath the name field, there is a "Description" section with a text input field containing the placeholder "Type your description here".

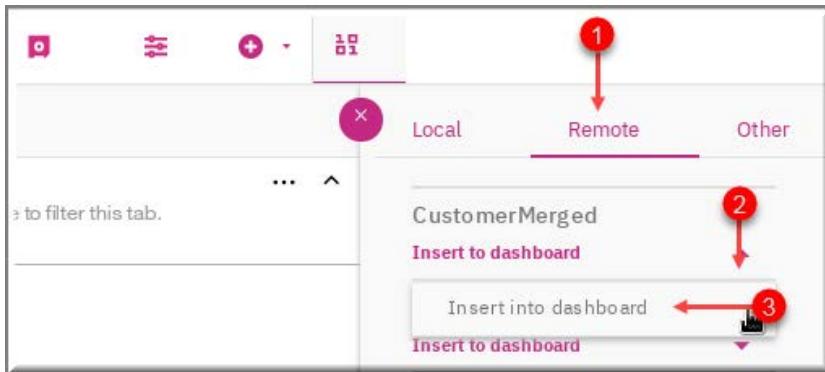
\_20. Select this template and click **Ok**.



\_21. Click the + next to **Selected Sources**



\_22. Click **Remote** on the flyout menu on the right side. Expand **CustomerMerged** and click **Insert to dashboard**.



\_23. This **CustomerMerged** Db2 warehouse table is inserted in the **Selected sources**. Click to select it.



\_\_24. Expand the table details by clicking the **down arrow** / caret.

The screenshot shows the 'CustomerMerged' table structure. On the left, a tree view lists columns under 'MERGED\_DEMOGRAPHIC\_CUSTOMER'. A red arrow points down to the 'ID' column, indicating where to click to expand. The right side shows a tabbed interface with 'Tab 1' selected, which is currently empty.

\_\_25. Click **ID** and then click ellipse (right side).

The screenshot shows the 'CustomerMerged' table structure with the 'ID' column selected. A red arrow points to the ellipsis button (three dots) located to the right of the 'ID' column header.

\_\_26. Click **Properties**.

The screenshot shows the 'Properties' dialog box for the 'ID' column. A red arrow points to the 'Properties' button at the bottom right of the dialog.

\_\_27. Select **Identifier** from the *Usage* drop-down menu and **Count** from the *Aggregate* drop-down menu. Click **Close**.

The screenshot shows the 'Properties' dialog box for the 'ID' column. The 'Usage' dropdown is set to 'Identifier' and the 'Aggregate' dropdown is set to 'Count'. Both dropdowns have a red arrow pointing to them, indicating they need to be changed.

- \_\_28. Repeat same for the following columns, setting their Usage and Aggregate values as follows:

Column Name	Usage	Aggregate
CHURNRISK	Attribute	None
GENDER	Attribute	None
STATUS	Attribute	None
CHILDREN	Attribute	None
ESTINCOME	Measure	None
HOMEOWNER	Attribute	None
AGE	Attribute	None
TOTALDOLLARVALUETRADED	Measure	Total
TOTALUNITSTRADED	Measure	Total
LARGESTSINGLETRANSACTION	Measure	None
SMALLESTSINGLETRANSACTION	Measure	None
PERCENTCHANGECALCULATION	Measure	None
DAYSSINCELASTLOGIN	Measure	None
DAYSSINCELASTTRADE	Measure	None
NETRALIZEDGAINS_YTD	Measure	None

- \_\_29. Drag **TOTALUNITSTRADED** towards the text **Drop here to manage** at top left of the screen. The icon color will change to light blue. Drop the field so that it occupies the full space.



- \_\_30. The result should display the total units traded as follows:



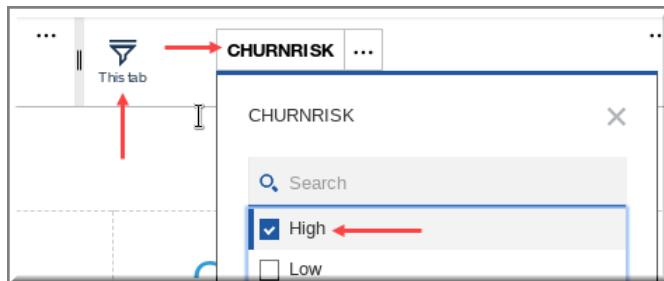
- \_\_31. Similarly drag **TOTALDOLLARVALUETRADED** to the top right and drop it when the icon turns light blue as you hover the text *Drag here*.

- \_\_32. The result should display the following.

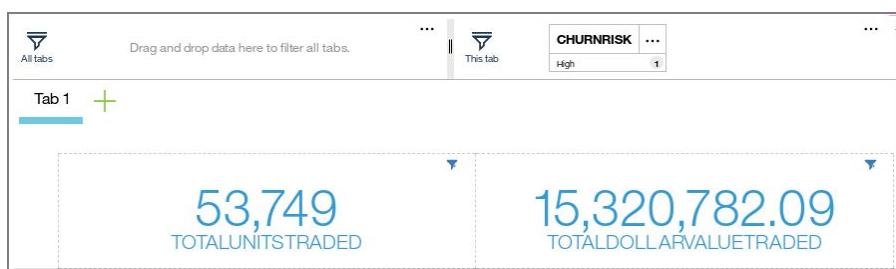


- \_\_33. Now that we have scored a historical churn risk of existing customers, let's discover traits of who those customers are.

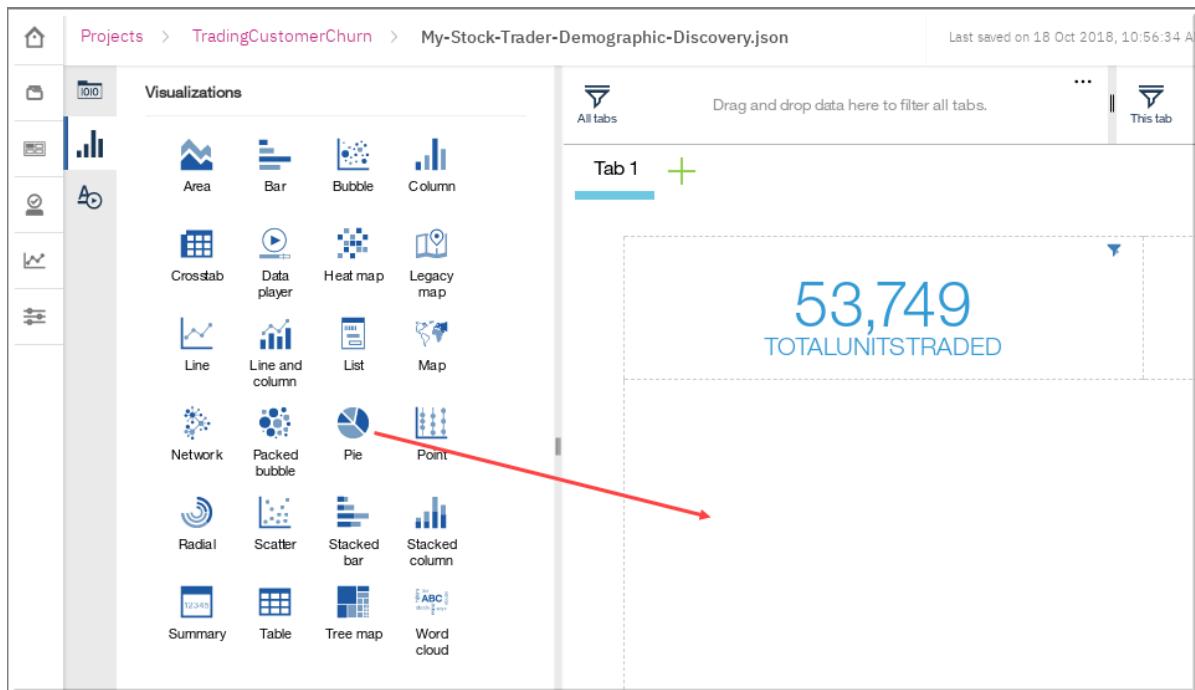
- \_\_34. Drag CHURNRISK and drop it on the **This Tab** text at the top right. Click CHURNRISK and check **High**. Click **OK**.



- \_\_35. Notice the amount changes after applying a filter of **High** risk customers.



- \_\_36. Click Middle Chart (left menu bar) icon and drag the **Pie chart** from the flyout menu to the bottom left.

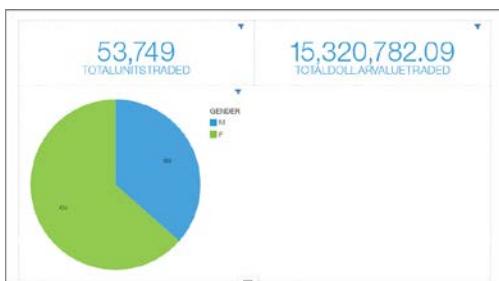


\_\_37. Drag **GENDER** to the top **Segment** section (next to the pie icon) and drag **ID** to the **#** section.

\_\_38. Click **Collapse** (top right corner of the pie chart.)



\_\_39. Your dashboard should display the following. Notice the relative proportion of female and male customers.

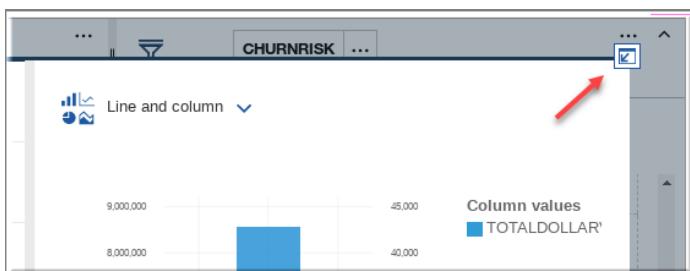


\_\_40. Let's find out who makes the most trades for the most amount of trade value.

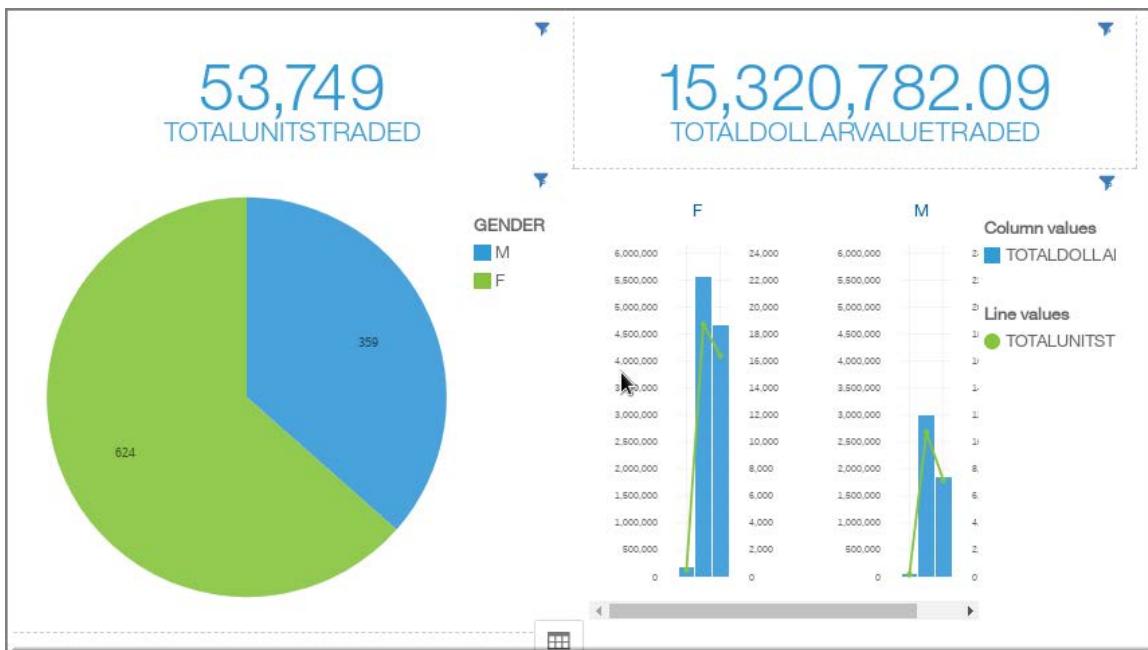
\_\_41. From the **Chart** palette, drag **Line and Column** to bottom right of the canvas.

- \_\_42. Drag each of the numbered items to the right to the position indicated (that is, Drag status to the top, then Length and so on.)

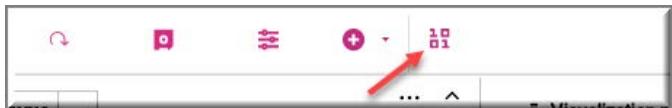
- \_\_43. Collapse the chart.



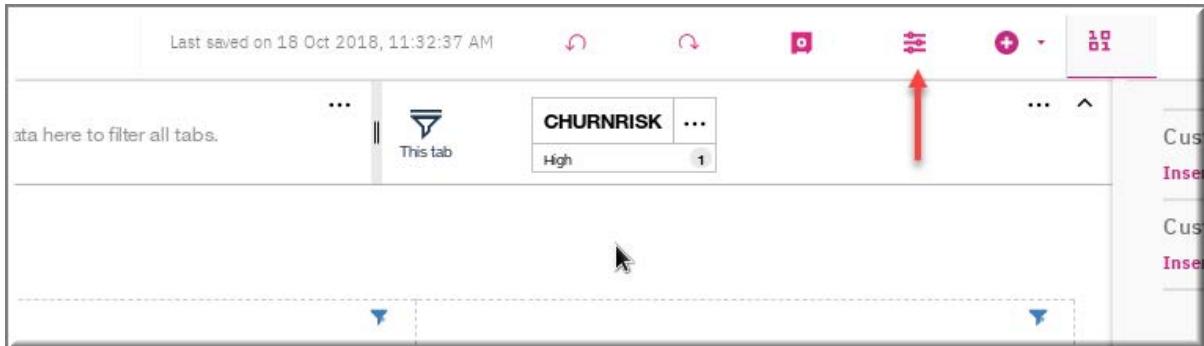
- \_\_44. The dashboard should display the following figure:



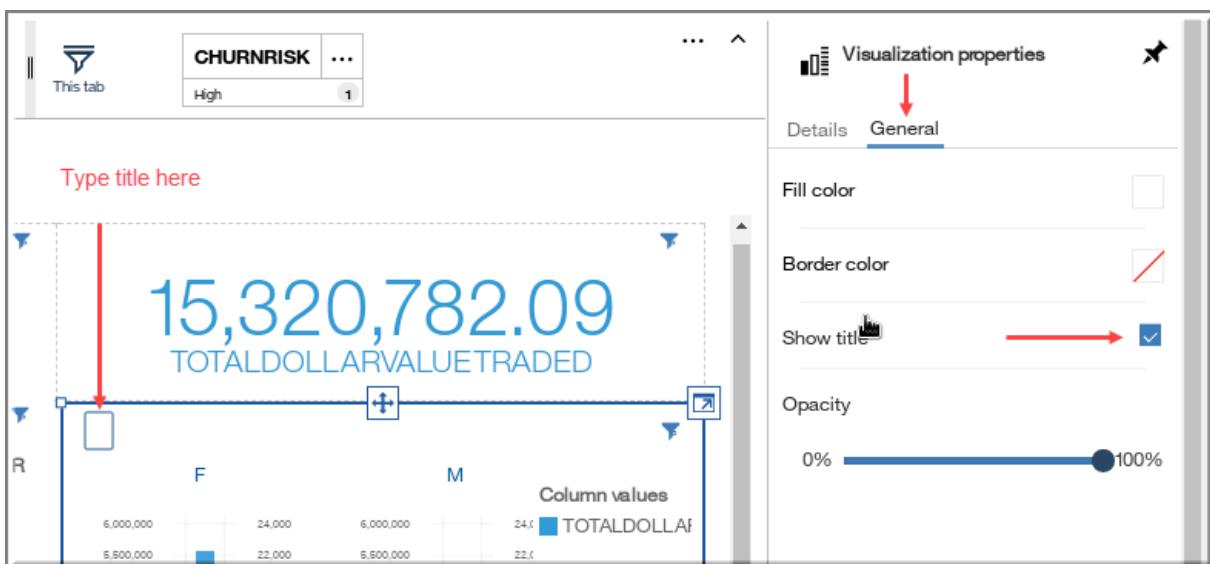
- \_\_45. Make some room available by hiding the right palette by clicking **Find Data** icon.



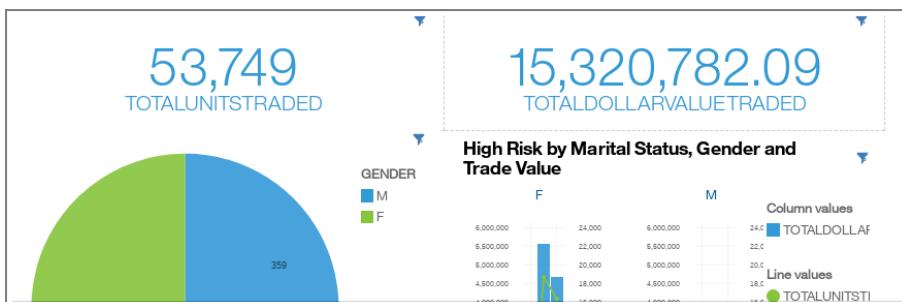
- \_\_46. Select the line column chart that you just added and click **Properties** from top menu bar.



- \_\_47. Click the **General** tab. Check **Show Title** and type: **High Risk by Marital Status, Gender and Trade Value**

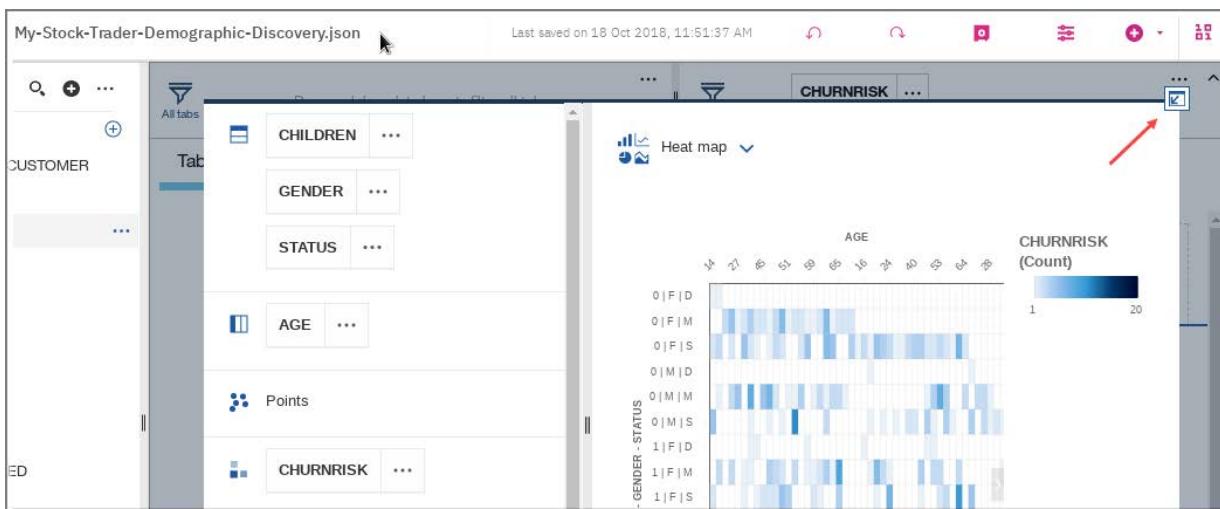


- \_\_48. Your dashboard is beginning to take shape. It should display as follows:



- \_\_49. It appears from the chart that the greatest number of high-risk customers who trade and spend the most with us are Married and Single Females.
- \_\_50. Is there any other insight we can gain? For example, does the number of children or home ownership plays a part?
- \_\_51. From the **Chart** palette, **drag Heat Map** to the bottom of the pie chart.
- \_\_52. Drag the **STATUS**, **GENDER**, and **CHILDREN** individually to the **Rows** section. Drag **AGE** to the **Column** section and drag **CHURNRISK** to the **Heat**. [As indicated in the number tags below]

- \_\_53. Collapse the **Heat Map** Chart.

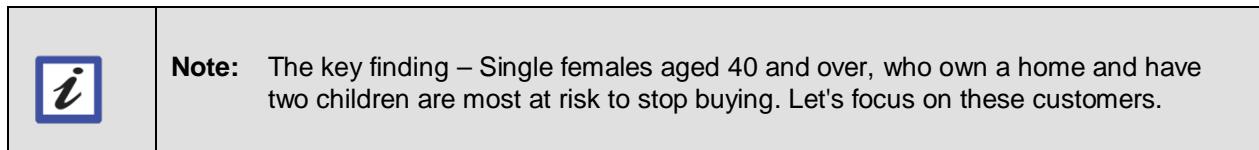
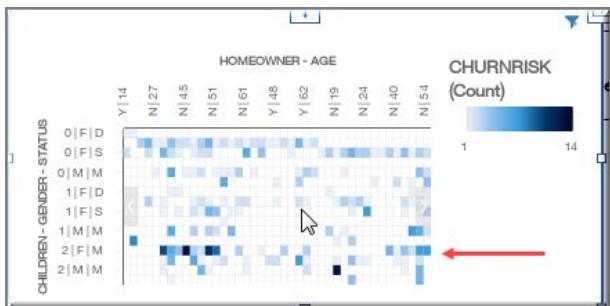


- \_\_54. The resultant Heat Map Chart shows that married couple with two children account for most of our customers.
- \_\_55. Does Home Ownership play a part?

- \_\_\_56. Expand the [Heat Map](#) chart and drag **HOMEOWNER** above **AGE** (Note: Do not replace Age).



- \_\_\_57. Collapse the chart.
- \_\_\_58. The [Heat Map](#) visual shows that single female homeowners over 40 years old with two children are most likely to be high risk.



- \_\_\_59. Change the title of the first tab to [Risk Demographics](#).

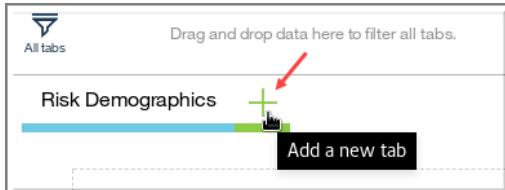


- \_\_\_60. Adjust the different visuals in the first tab. For example: Place the heat map below the pie and column charts.

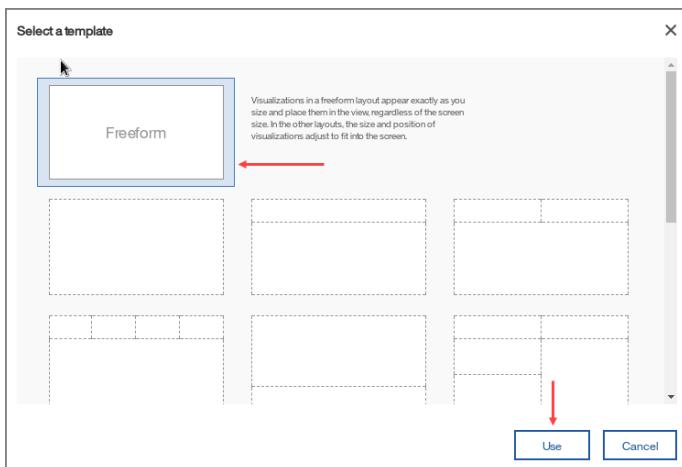


#### 5.1.4 Analyze Risk to the Business

- \_\_61. Let's create another tab to see the risk to the business.
- \_\_62. Click **+** to add a tab.



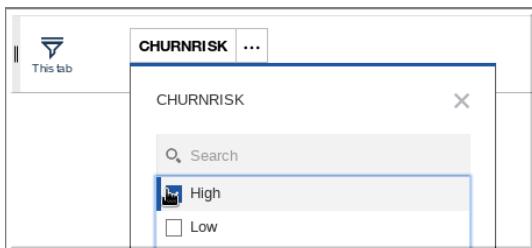
- \_\_63. Select the **Free Form** template and click **Use**.



- \_\_64. Change the title to **Risk to the Business**. [Click the **pencil** in the tab.]

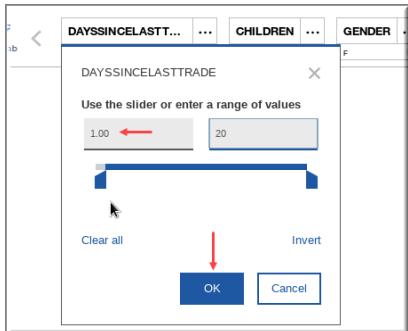


- \_\_65. Let's filter our results so that we focus on this group.
- \_\_66. Drag column **CHURNRISK** to the *Drag and Drop Data to Filter THIS tab*. [Top right of the canvas.] Click it and check **High** and click **OK**.



- \_\_67. Drag **GENDER** and click after you drop in the filter area and check **Female**. Click **OK**.
- \_\_68. Drag **CHILDREN** and click after you drop in the filter area and select **2** and click **OK**.

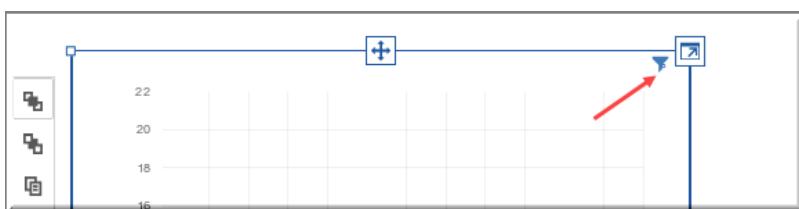
- \_\_69. Drag **DAYSSINCELASTTRADE** and click after you drop in the filter area and select **2** for the lower bound and keep the upper bound to **20**. Click **OK**.



- \_\_70. Now we have applied filters as per our previous analysis, let's build the dashboard.
- \_\_71. Drag **TOTALDOLLARVALUETRADED** to the top left area of the canvas. Resize the chart so it appears in the top left corner. This creates a single value field that is filtered based upon our tab filter.



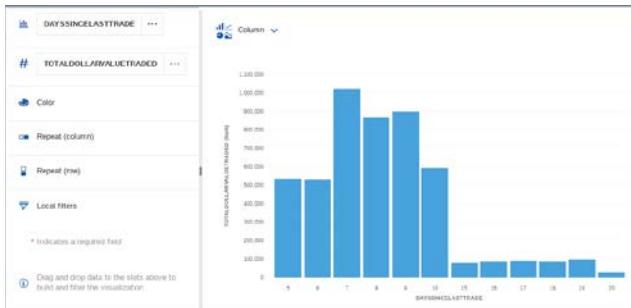
- \_\_72. It looks like this group trades roughly \$5M.
- \_\_73. Pull over **TOTALUNITSTRADED** over to the top right area of the canvas. Resize the chart so that it appears in the top right corner.
- \_\_74. Let's see the **TOTALDOLLARVALUE** at risk by this group based upon the days since they last traded with us and the date they last logged in to our system.
- \_\_75. Hold control key down and click **DAYSSINCELASTTRADE** and **TOTALDOLLARVALUETRADED** and drag to the right.
- \_\_76. Select **Expand**.



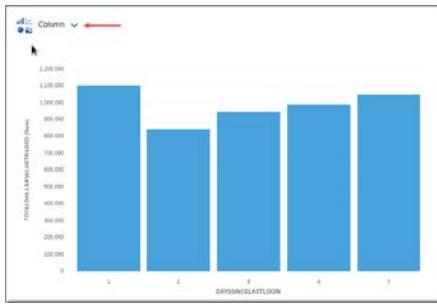
- \_\_77. Change the chart type to **Column**. [Click and select from the chart palette.]



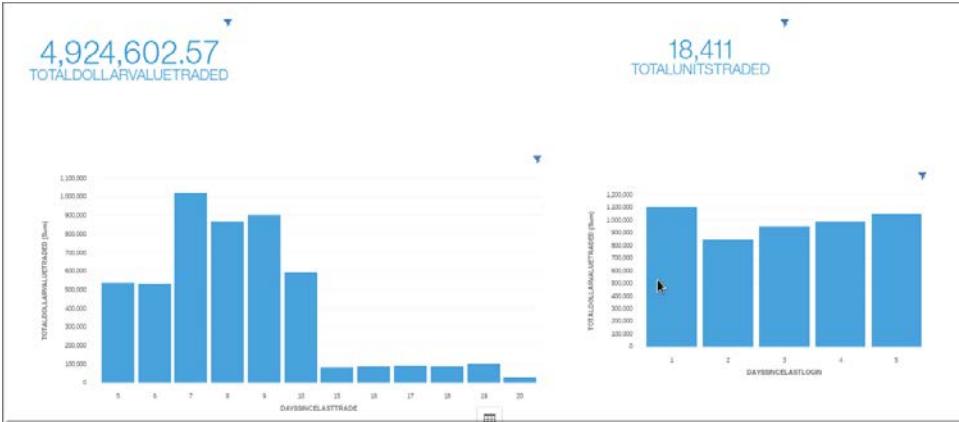
- \_\_78. The chart gets drawn as shown below:



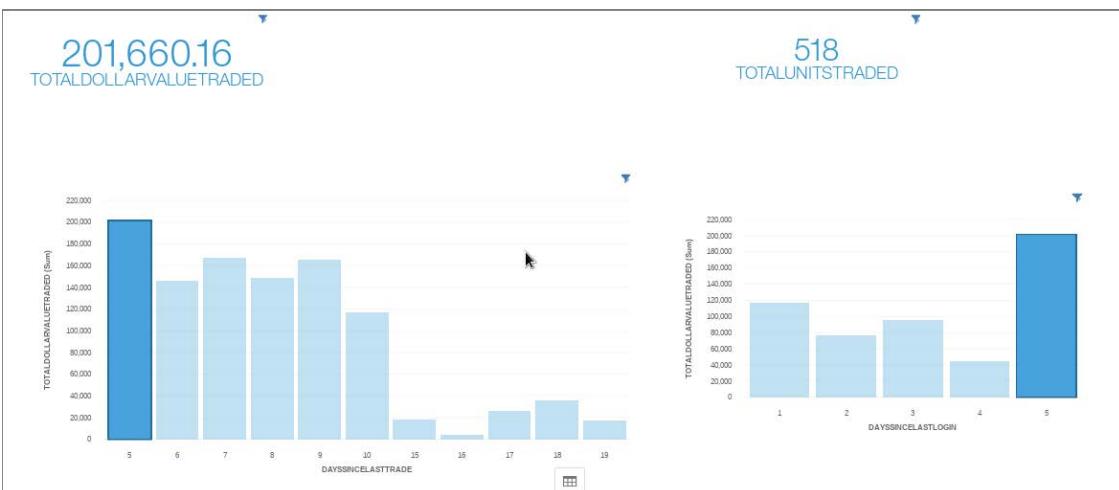
- \_\_79. Collapse the screen by clicking the collapse icon to the top right of the frame.
- \_\_80. Resize this chart and fit into the bottom left of the canvas using the chart movement arrows shown on the border of the chart.
- \_\_81. Hold control key down and click **DAYSSINCEDLASTLOGIN** and **TOTALDOLLARVALUETRADED** and drag to the right.
- \_\_82. Expand the chart.
- \_\_83. Change the Chart type to **Column**. It should display the following:



- \_\_84. Collapse the chart by clicking the top right icon on the chart and use resize arrows to fit in the bottom right of the screen.
- \_\_85. The second tab will display as shown below.



86. You can click any column from each chart to show the aggregate value for total dollar traded and total units traded. For example: Click first column on first chart and last column on second chart to see the values.



	<b>Note:</b> As a business analyst, you identified the group likely to be at risk to separate from the business.  An intelligent scientific solution needs to be crafted to mitigate the risks.  Data scientists are now engaged to move our analysis forward.
--	--

**Credit for the Lab:** Kent Rubin, Analytics and Cognitive Solution Architect

## 5.2 Data Scientists – Approach to a Solution

As seen in the previous labs, the business analysts have visualized the data so that it can be understood. The search of a solution begins now by engaging data scientists from the organization.

We have uncovered few key points along the way.

- Data resides everywhere and collecting the data in one lake (or platform) is not realistic.
- Data virtualization within ICP-D expands the capability to collect the data.
- Data governance is key to comply with rules and regulations but is difficult to implement, which hinders discovery and use of data.
- Dark Data is what leads many organizations to not take an action when it was necessary.
- Data preparation consumes too much of the data scientists' time today. They often spend more time preparing data than developing models.

We have gone through the journey of automation, machine learning while discovering and cataloging of data, it is now time to allow data scientists to provide solutions.

\_87. Let's begin through the next phase of Analytics Modernization.

### 5.2.1 Analyzing Data and Toolset

\_88. The recent contribution of universities, research institutes, corporations and individuals to Open Source has led to the development of innovative tools in data science. Most notable of them are:

- Jupyter and Spark Notebooks
- Zeppelin and Spark Notebooks
- R-Studio

IBM Cloud Private for Data supports all three.

- \_89. Python has become a language a choice for many data scientists mainly due to a very extensive ecosystem of libraries built for the Python language platform. It is no surprise that today Python commands 35 percent of the market share amongst all languages.
- \_90. We will use Jupyter notebooks and Spark for this lab exercise.
- \_91. ICP-D platform provides sample notebooks that show real-world scenarios and contain many useful examples of computations and visualizations that will jumpstart your analytics projects.

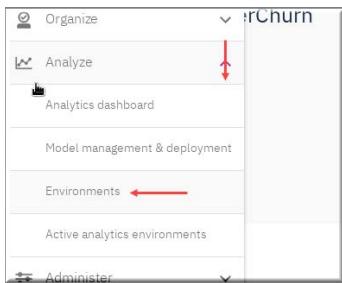


**Note:** The non-obvious value of the IBM Cloud Platform is to provide an agile "cloud-native" infrastructure within your firewall.

## 5.2.2 Walk Through the Model

- \_\_92. Developing, training, and deploying a machine learning lab from scratch is beyond the scope of this introduction to ICP for Data, so we will walk through a Jupyter notebook that represents an example of the model development process. If you are familiar with the process you will recognize the tools used and will learn how ICP for Data speeds and facilitates collaboration throughout the process. If you're new to machine learning it will serve as a gentle introduction to the process. We will develop, test and score a model that addresses our business problem.

- \_\_93. Switch back to the ICP-D Web UI and click [Analyze](#)  $\Rightarrow$  [Environments](#).



- \_\_94. We need to stop all [Running](#) environments to free-up resources. Click **three vertical dots** to the right of the [Running](#) environment.

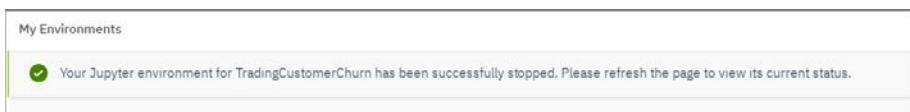
My environments							
NAME	TYPE	PROJECT	STATUS	CPU (CORES)	GPUS	MEMORY (GB)	DATE STARTED
Decision Optimization		TradingCustomerChurn	Stopped	2.0	0	1.0	
Jupyter with Python 2.7, Scala 2.11, R 3.4.3, Spark 2.0.2	Jupyter	TradingCustomerChurn	Stopped	—	—	—	
Jupyter with Python 3.5 for GPU	Jupyter	TradingCustomerChurn	Stopped	—	—	—	
Jupyter with Python 3.5, Scala 2.11, R 3.4.3, Spark 2.2.1	Jupyter	TradingCustomerChurn	Running	—	—	—	24 Oct 2018, 4:04 PM
RStudio with R 3.4.3	RStudio	TradingCustomerChurn	Stopped	—	—	—	

- \_\_95. Click [Stop](#).

Running	—	—	—	24 Oct 2018, 4:04 PM	⋮
Stopped	—	—	—		<a href="#">Stop</a> $\leftarrow$
Stopped	—	—	—		<a href="#">Edit settings</a>
Stopped	—	—	—		<a href="#">Open</a>

- \_\_96. Repeat same for any other environment, which is [Running](#) in your case.

- \_\_97. You will see the message of the environment getting stopped.



- \_\_98. Click [Projects](#) from the left menu bar.

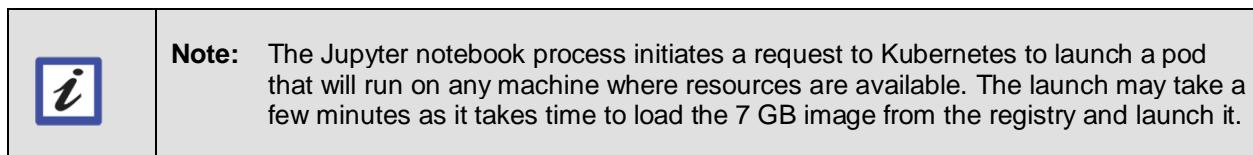
\_\_\_99. Click [TradingCustomerChurn](#) to open it.

NAME	PROJECT TYPE	USER ROLE	LAST UPDATED	ACTIONS
TradingCustomerChurn	Analytics	Admin	17 Oct 2018, 8:40 AM	[More]
dsx-samples	-	Viewer	1 Oct 2018, 10:39 AM	[More]
ANALYZERPROJECT	Data Transform	Admin	5 Sep 2018, 5:22 PM	[More]
DataClick	Data Transform	Admin	5 Sep 2018, 5:23 PM	[More]
dstage1	Data Transform	Admin	5 Sep 2018, 5:13 PM	[More]

\_\_\_100. Click [Assets](#). [Your assets number might be different than shown below].



\_\_\_101. Click Jupyter Notebook – [01TradingCustomerChurnClassifierSparkML](#) to open it.



\_\_\_102. If you see the free resources message, please ignore it. It is normal due to the restricted resources in the lab environment.



\_\_\_103. If you are an experienced user of Jupyter notebooks, you may skip steps that are review for you.

\_\_\_104. A notebook a collection of compute cells that can be run individually, or as a group. You can pause execution between steps to review results, make adjustments to your code and re-run. One of the features of notebooks is that your notes and markdown cells can be interspersed with your code. This allows notebooks to be self-documenting.

\_\_\_105. A notebook can use environment of your choice. For example, we chose to use Python 3.5 with Spark 2.2.1. Other choices will be available, depending on your cluster configuration and hardware. For example, other levels of Python, Spark, R, GPUs, might be available. These options are chosen at the time of notebook creation.

The screenshot shows the IBM Cloud Private for Data interface. At the top, there's a navigation bar with 'IBM Cloud Private for Data'. Below it, a breadcrumb path shows 'Projects > TradingCustomerChurn > 01TradingCustomerChurnClassifierSpa'. A toolbar with various icons is visible above the main content area. The main content area contains a title 'Trading Platform Customer Attrition Risk Prediction using SparkML' and a brief introduction about predicting customer churn based on user activity.

\_106. Read through the introduction to learn what this notebook provides.

In this notebook, we will leverage Data Science Experience Local to do the following:

1. Ingest merged customer demographics and trading activity data
2. Visualize merged dataset and get better understanding of data to build hypotheses for prediction
3. Leverage SparkML library to build classification model that predicts whether customer has propensity to churn
4. Expose SparkML classification model as RESTful API endpoint for the end-to-end customer churn risk prediction and risk remediation application

The diagram illustrates a machine learning workflow. It starts with 'Data Preparation' and 'Data Visualization' leading to a 'Training set'. This set feeds into two parallel models: 'Model 1 All Features' and 'Model 2 Top 10 Important Features', both of which have arrows pointing to a central 'learn' node. The 'learn' node then points to a 'Test set', which leads to 'Evaluation' and finally 'Deployment'. The Apache Spark logo is at the bottom, and the IBM Data Science Experience logo is at the bottom right.

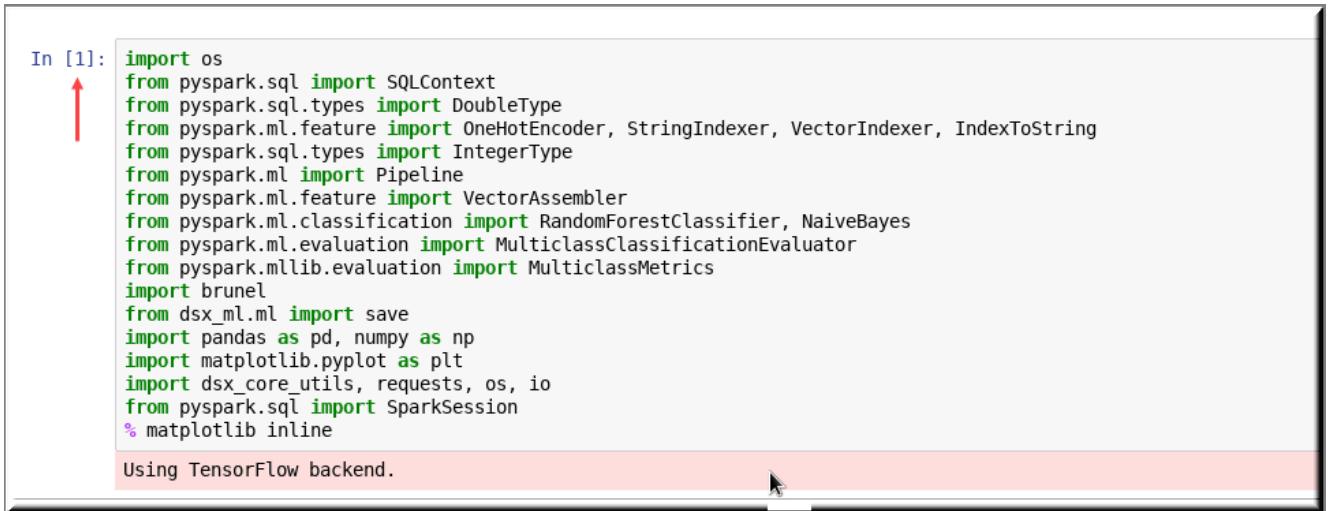
\_107. The notebook explains two type of cells – markdown and code. The documentation is done using markdown (easy to learn) and code cells are actual code that will execute when we want them to.

\_108. To execute a cell (even including markdown cell), you can either press menu icon [Run cell select below](#) or use keyboard [Shift-Enter](#) to execute a cell.

The screenshot shows the Jupyter notebook toolbar. The 'Cell' menu icon is highlighted with a red arrow. Other icons include File, Edit, View, Insert, Kernel, Widgets, Help, and CellToolbar.

\_109. You are on a first cell, press [Shift-Enter](#) to execute the cell.

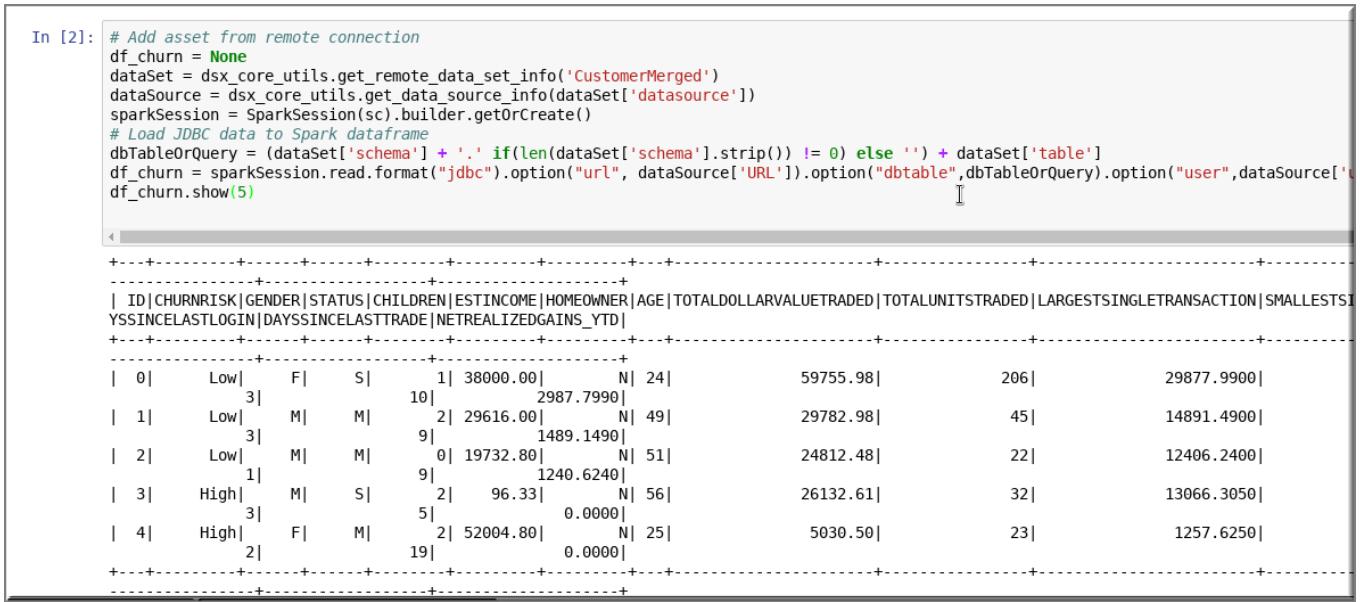
- 110. When you execute a cell, you will see an \* in the cell square bracket. This is an indication that the cell is currently running. You can execute cell one at a time, in which case execution pauses at the next cell. You can also "Run All" cells which causes a serial execution from cell to cell without pause.
- 111. When a cell completes its execution, you will see a number displayed in its cell boundary as shown.



```
In [1]: import os
from pyspark.sql import SQLContext
from pyspark.sql.types import DoubleType
from pyspark.ml.feature import OneHotEncoder, StringIndexer, VectorIndexer, IndexToString
from pyspark.sql.types import IntegerType
from pyspark.ml import Pipeline
from pyspark.ml.feature import VectorAssembler
from pyspark.ml.classification import RandomForestClassifier, NaiveBayes
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
from pyspark.mllib.evaluation import MulticlassMetrics
import brunel
from dsx_ml.ml import save
import pandas as pd, numpy as np
import matplotlib.pyplot as plt
import dsx_core_utils, requests, os, io
from pyspark.sql import SparkSession
%matplotlib inline

Using TensorFlow backend.
```

- 112. The next cell is a document cell that explains the procedure to insert a data frame which can be either a Pandas or Spark data frame.
- 113. Press **Shift-Enter** to execute a cell and watch its output and then go to the next cell by pressing **Shift-Enter** again and keep repeating it until we reach to the end of this notebook.
- 114. The next data cell is the inserted Panda data frame for connecting to Db2 Warehouse.



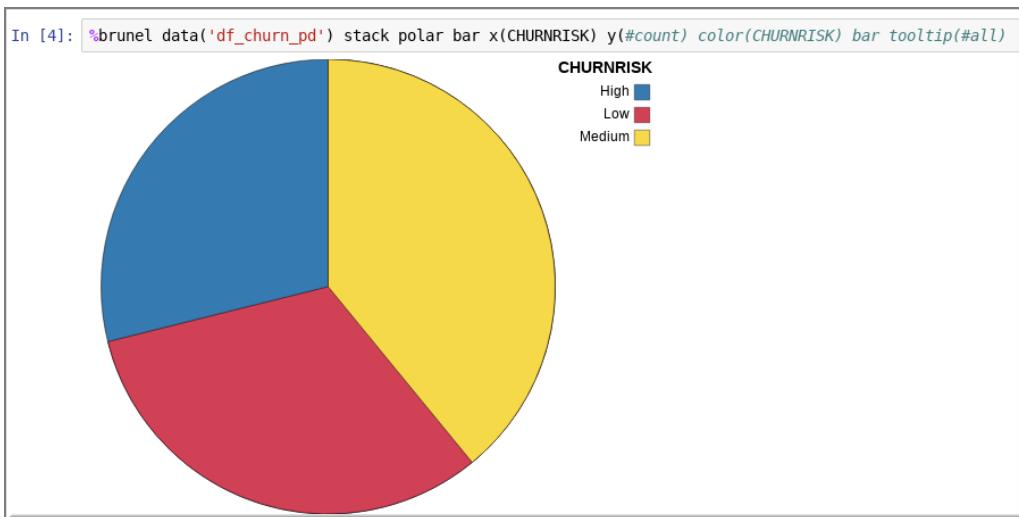
```
In [2]: # Add asset from remote connection
df_churn = None
dataSet = dsx_core_utils.get_remote_data_set_info('CustomerMerged')
dataSource = dsx_core_utils.get_data_source_info(dataSet['datasource'])
sparkSession = SparkSession(sc).builder.getOrCreate()
# Load JDBC data to Spark dataframe
dbTableOrQuery = (dataSet['schema'] + '.' if len(dataSet['schema'].strip()) != 0 else '') + dataSet['table']
df_churn = sparkSession.read.format("jdbc").option("url", dataSource['URL']).option("dbtable", dbTableOrQuery).option("user", dataSource['u'])
df_churn.show(5)
```

ID	CHURNRISK	GENDER	STATUS	CHILDREN	ESTINCOME	HOMEOWNER	AGE	TOTALDOLLARVALUETRADED	TOTALUNITSTRADED	LARGESTSINGLETRANSACTION	SMALLESTSINGLETRANSACTION	SINCEREGISTRATION	SINCERLASTLOGIN	DAYSSINCERLASTTRADE	NETREALIZEDGAINS_YTD	
0	Low	F	S	1	38000.00		24	59755.98	206	29877.9900						
1	Low	M	M	2	29616.00		49	29782.98	45	14891.4900						
2	Low	M	M	0	19732.80	N	51	24812.48	22	12406.2400						
3	High	M	S	2	96.33	N	56	26132.61	32	13066.3050						
4	High	F	M	2	52004.80	N	25	5030.50	23	1257.6250						

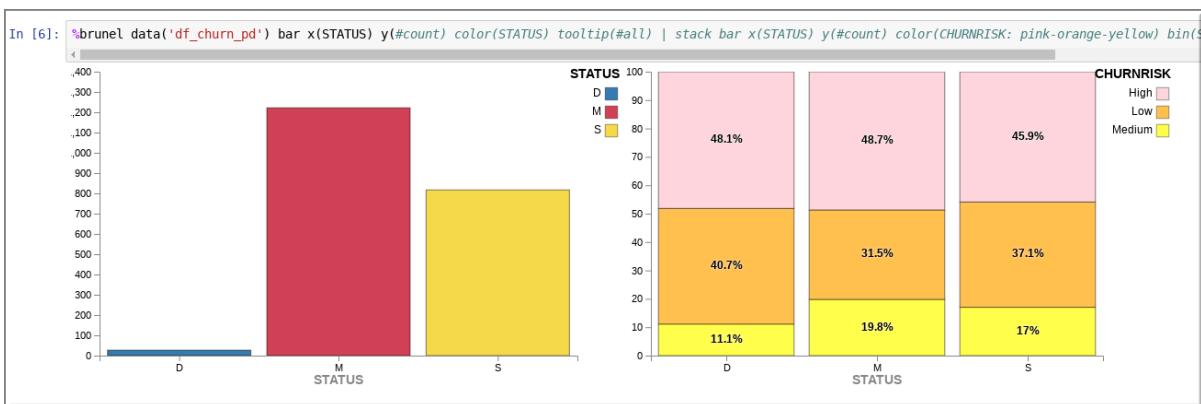
- \_\_115. The data above shows the output from `MERGED_DEMOGRAPHICS_TRADING_CUSTOMER` table which was the result of the previous lab exercise and the focus of the business analysis that we did in [Business Analyst](#) section.

	<p><b>Note:</b> You might wonder – do I have to write this much code to connect to my data source?</p> <p>The answer is No – you just define a data source and drag it here from your remote data set and the code gets generated automatically.</p>
---	--

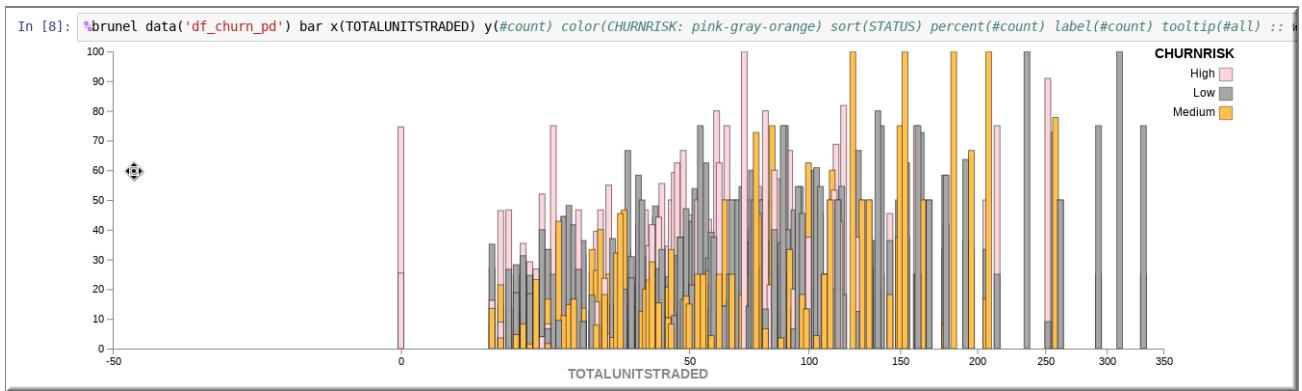
- \_\_116. The next cell contains code to visualize the data using the Brunel library. Brunel is just one of many open source visualization libraries available for Python and Jupyter. You may also have heard of matplotlib, plotly, Pixiedust and others. New libraries are constantly being released.



- \_\_117. The next visualization shows the separation risk of different category and status.



\_\_118. Next visualization is total units traded by each separation risk category.



\_\_119. The next stage in model development is data preparation, which involves the indexing of categories and assembling them into a feature vector to train models.

### 5.2.3 Build Spark ML Model using Random Forest Classification

\_\_120. These steps show the definition of a pipeline chains together transformations and estimators for machine learning algorithms.

```
# Split original dataset into train and test datasets. We fit the pipeline to training data and apply the trained model to transform test data.
In [17]: # instantiate a random forest classifier, take the default settings
rf=RandomForestClassifier(labelCol="label", featuresCol="features")

# Convert indexed labels back to original labels.
labelConverter = IndexToString(inputCol="prediction", outputCol="predictedLabel", labels=labelIndexer.labels)

stages += [labelIndexer, assembler, rf, labelConverter]

pipeline = Pipeline(stages = stages)

In [18]: # Split data into train and test datasets
train, test = df_churn.randomSplit([0.7,0.3], seed=100)
train.cache()
test.cache()

Out[18]: DataFrame[ID: int, CHURNRISK: string, GENDER: string, STATUS: string, CHILDREN: int, ESTINCOME: int, HOMEOWNER: int, LARGESTSINGLETRANSACTION: int, SMALLEST_SINGLETRANSACTION: int, PERCENTCHANGECALCULATION: int, DAYSSINCELASTTRANSACTION: int, predictedLabel: string, prediction: float, probability: array]
```

\_\_121. The next process builds the models using the pipeline, transforms and shows the results.

```
In [22]: results = model.transform(test)
results.select(results["ID"],results["CHURNRISK"],results["label"],results["predictedLabel"],results["prediction"],results["probability"])
results.toPandas().head(6)

Out[22]:
```

	ID	CHURNRISK	label	predictedLabel	prediction	probability
0	4	High	0.0	High	0.0	[0.825155165927, 0.000936872651527, 0.17390796...]
1	7	High	0.0	High	0.0	[0.776170343948, 0.0975125901263, 0.126317065926]
2	8	Medium	2.0	High	0.0	[0.508695939111, 0.000646174977108, 0.49065788...]
3	9	Medium	2.0	Medium	2.0	[0.239374989112, 0.0114018087855, 0.749223202103]
4	15	Low	1.0	Low	1.0	[0.0377812370278, 0.922980122156, 0.0392386408...]
5	18	Low	1.0	Low	1.0	[0.134921727231, 0.837544808518, 0.0275334642508]

\_\_122. The model results and interpretation are described in the documentation cell.

### Model results

In a supervised classification problem such as churn risk classification, we have a true output and a model-generated predicted output for each data point. For this reason, the results for each data point can be assigned to one of four categories:

1. True Positive (TP) - label is positive and prediction is also positive
2. True Negative (TN) - label is negative and prediction is also negative
3. False Positive (FP) - label is negative but prediction is positive
4. False Negative (FN) - label is positive but prediction is negative

These four numbers are the building blocks for most classifier evaluation metrics. A fundamental point when considering classifier evaluation is that pure accuracy (i.e. was the prediction correct or incorrect) is not generally a good metric. The reason for this is because a dataset may be highly unbalanced. For example, if a model is designed to predict fraud from a dataset where 95% of the data points are not fraud and 5% of the data points are fraud, then a naive classifier that predicts not fraud, regardless of input, will be 95% accurate. For this reason, metrics like precision and recall are typically used because they take into account the type of error. In most applications there is some desired balance between precision and recall, which can be captured by combining the two into a single metric, called the F-measure.

\_\_123. The model precision for this is 91 percent.

```
In [24]: print('Model Precision = {:.2f}'.format(results.filter(results.label == results.p
Model Precision = 0.91.
```

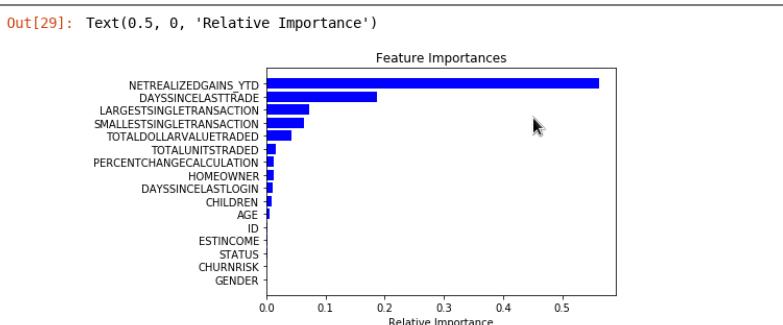
### 5.2.4 Evaluate Different Models

\_\_124. The process of building model is incomplete unless it can be evaluated with a test data to see its accuracy.

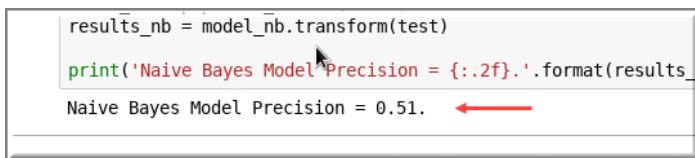
```
Overall Statistics
Model F-measure = 0.9096774193548387

Statistics by Class
Class 0.0 F-Measure = 0.9177631578947368
Class 1.0 F-Measure = 0.9864253393665159
Class 2.0 F-Measure = 0.7052631578947368
```

\_\_125. The next cell shows the “Feature Importance” which shows how much of an impact each feature has on prediction.



\_\_126. Next is Naïve Bayes classifier trained on the training data set.



```

results_nb = model_nb.transform(test)
print('Naive Bayes Model Precision = {:.2f}'.format(results_nb))
Naive Bayes Model Precision = 0.51. ←

```

\_\_127. The above indicates that Random Forest Classifier shows a high F-measure upon evaluation, which shows strong performance.

### 5.2.5 Save the Model in Repository

\_\_128. After comparison of two models, the Random Forest model is saved into the repository.



```

In [31]: save(name='TradingChurnRiskClassificationSparkML',
          model=model,
          test_data = test,
          algorithm_type='Classification',
          description='This is a SparkML Model to Classify Trading Customer Churn Risk')
Out[31]: {'path': '/user-home/999/DSX_Projects/TradingCustomerChurn/models/TradingChurnRiskClassificationSparkML/1',
          'scoring_endpoint': 'https://dsxl-api/v3/project/score/Python35/spark-2.2/TradingCustomerChurn/TradingChurnRiskClassificationSparkML/1'}

In [32]: # Write the test data without label to a .csv so that we can later use it for batch scoring
write_score_CSV=test.toPandas().drop(['CHURNRISK'], axis=1)
write_score_CSV.to_csv('../datasets/TradingCustomerSparkMLBatchScore.csv', sep=',', index=False)

In [33]: # Write the test data to a .csv so that we can later use it for Evaluation
write_eval_CSV=test.toPandas()
write_eval_CSV.to_csv('../datasets/TradingCustomerSparkMLEval.csv', sep=',', index=False)

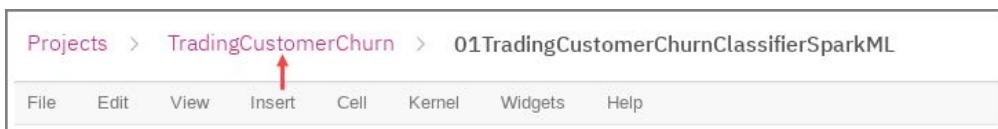
```



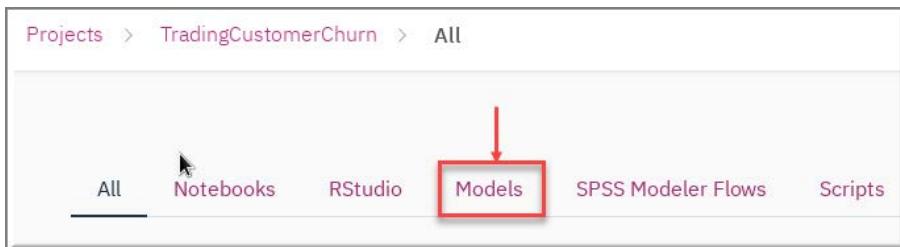
**Note:** The model is saved in the platform repository, providing governance to the data and model.

### 5.2.6 Test, Batch Score and Evaluate Saved Model

\_\_129. Click [TradingCustomerChurn](#) Project.



\_\_130. Click [Models](#).



\_\_131. Click to open **Trading Custom Churn** model that we saved in the previous section.

The screenshot shows a table titled 'Models (1)'. It has three columns: 'NAME', 'TYPE', and 'STATUS'. There is one row containing the model 'TradingChurnRiskClassificationSparkML' with version 'v1'. A red arrow points to the 'v1' button. The 'TYPE' column shows 'spark-2.2' and the 'STATUS' column shows 'trained'.

NAME	TYPE	STATUS
TradingChurnRiskClassificationSparkML v1	spark-2.2	trained

\_\_132. We noticed that the Random Forest Classification gave us 91 percent and Bayes classifier was at 51 percent accuracy. Click **Test**.

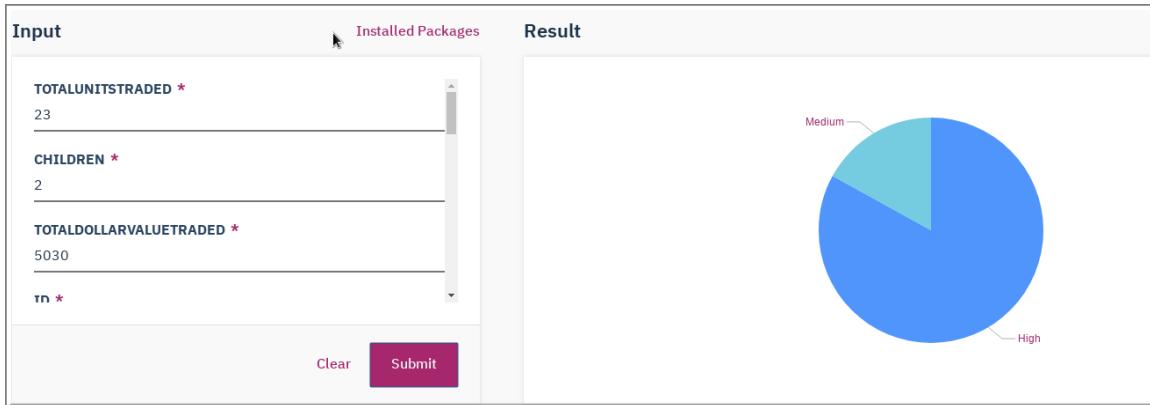
The screenshot shows the 'Test' tab selected in the navigation bar. Below it, the title is 'TradingChurnRiskClassificationSparkML v1'. A red arrow points to the 'Test' tab. On the left, there's a 'Accuracy' section with a green donut chart labeled '91%' and a 'Accuracy history' section with a line graph. A red arrow points to the '91%' label on the chart.

\_\_133. Scroll down to the input section where sample input data can be entered to test the model. Click **Submit**.

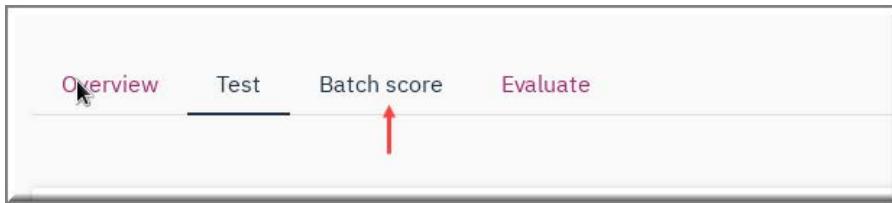
The screenshot shows an input form with three sections: 'Input', 'Installed Packages', and 'Result'. The 'Input' section contains several fields with sample data: 'TOTALUNITSTRADED \*' (23), 'CHILDREN \*' (2), 'TOTALDOLLARVALUETRADED \*' (5030), 'ID \*' (4), and 'NETDEALTZEROCOUNT VTD \*' (0). A red arrow points to the 'Submit' button at the bottom of the input section. The 'Installed Packages' and 'Result' sections are currently empty.

Input	Installed Packages	Result
TOTALUNITSTRADED *		
23		
CHILDREN *		
2		
TOTALDOLLARVALUETRADED *		
5030		
ID *		
4		
NETDEALTZEROCOUNT VTD *		
0		
<input type="button" value="Clear"/>	<input type="button" value="Submit"/>	

- \_\_134. The model online score returns the results indicating percentage of high or medium separation risk.



- \_\_135. Scroll up and click **Batch Score**.



- \_\_136. Select Input data set [TradingCustomerSparkMLBatchScore.csv](#) and type the output file name [batchscoreresults.csv](#). Click **Generate batch script**.

The figure shows the 'Batch scoring script inputs' dialog. It has two main sections: 'Batch scoring script inputs' on the left and 'Result' on the right. In the 'Batch scoring script inputs' section, there are three fields: 'Execution Type \*' set to 'DSX', 'Input data set \*' set to 'TradingCustomerSparkMLBatchScore.csv', and 'Output data set \*' with a radio button selected for 'Local file' and the value 'batchscoreresults.csv'. A red arrow points from the 'Input data set' field down to the 'Generate batch script' button at the bottom. Another red arrow points from the 'Output data set' field down to the same button. The 'Result' section on the right shows a dark, mostly black preview area with a small number '1' in the top-left corner.

\_\_137. Click **Run now**.

```

1 #!/usr/bin/python
2
3 import sys, os
4 from pyspark.sql import SparkSession
5 from pyspark.ml import Pipeline, Model, PipelineModel
6 from pyspark.sql import SQLContext
7 import pandas
8 import dsx_core_utils, re, jaydebeapi
9 from sqlalchemy import *
10 from sqlalchemy.types import String, Boolean
11
12 # define variables
13 args = {'source': '/datasets/TradingCustomerSparkMLBatchScore.csv', 'output_type': 'Localfile', 'target': '/datasets/batchscore'}
14 input_data = os.getenv("DSX_PROJECT_DIR") + args.get("source")
15 output_data = os.getenv("DSX_PROJECT_DIR") + args.get("target")
16 model_path = os.getenv("DSX_PROJECT_DIR") + "/models/TradingChurnRiskClassificationSparkML/1/model"
17
18 # create spark context
19 spark = SparkSession.builder.getOrCreate()
20 sc = spark.sparkContext
21
22

```

\_\_138. Scroll down and you will see that a batch job is running.

ID	NAME	TARGET HOST	TRIGGERED BY	STARTED AT	DURATION (S)	RESULT
1599899137-999	Run 1	Local instance	admin	18 Oct 2018, 5:45 PM	In progress	Running

\_\_139. Wait for this to complete and you should see the **Success** message.

DURATION (S)	RESULT
54	Success

\_\_140. Click **TradingCustomerChurn** project [Top left corner].

\_\_141. Click **Assets** **Data Sets**.

- \_\_142. Click [batchscoreresults](#) to open and preview it.
- \_\_143. Scroll to the right and check the prediction for the sample data for different IDs.

probability	prediction	predictedLabel
[0.825155165927,0.000936872651527,0.173907961422]	0.0	High
[0.776170343948,0.0975125901263,0.126317065926]	0.0	High
[0.508695939111,0.000646174977108,0.490657885912]	0.0	High
[0.239374989112,0.0114018087855,0.749223202103]	2.0	Medium
[0.0377812370278,0.922980122156,0.0392386408163]	1.0	Low
[0.134921727231,0.837544808518,0.0275334642508]	1.0	Low

i

**Note:** Based upon demographics information of customers and some information from legacy risk applications, the model predicts the probability of separation.

Models must be periodically tested, trained and redeployed.

This is where you will see the value of ICP-D as an integrated data platform.

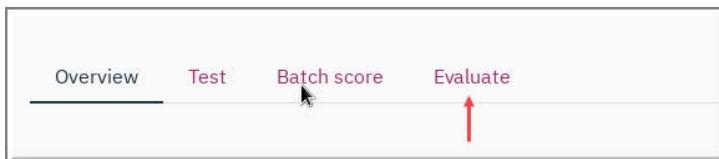
- \_\_144. Click [Close](#) to preview the batch results.

## 5.2.7 Create an Evaluation

- \_\_145. Click [Models](#)

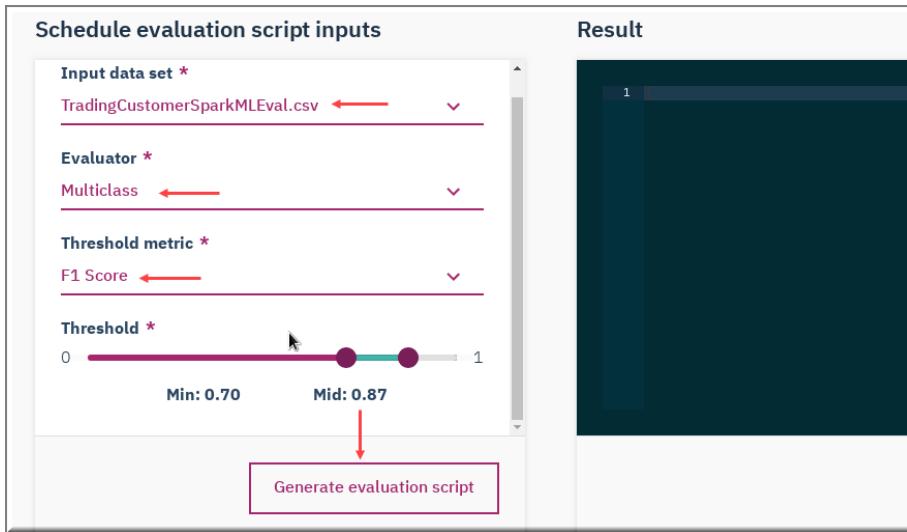
- \_\_146. Click to open the model that we just previewed in previous steps.

- \_\_147. Click [Evaluate](#)



- \_\_148. Scroll down and select input data set [TradingCustomerSparkMLEval.csv](#) from the drop-down menu.
- \_\_149. Select [Multiclass](#) as the evaluator, [F1 score](#) as the threshold metric.

\_150. Change the **Threshold** criteria between **0.71** to **0.89**. Click **Generate evaluation script**.



\_151. The script is generated and click **Run now**.

```

1 #!/usr/bin/python
2
3 import pandas as pd
4 import json
5 import time, sys, os, shutil, glob, io, requests
6 from pyspark.sql import SparkSession
7 from pyspark.ml import Pipeline, Model, PipelineModel
8 from pyspark.sql import SQLContext
9 import dsx_core_utils
10
11 # define variables
12 args = {"evaluator_type": "multiclass", "threshold": {"mid_value": 0.87, "metric": "f1Score", "min_value": 0.7}, "dataset": "TradingCustomerSparkMLEval.csv", "model_path": os.getenv("DSX_PROJECT_DIR") + "/models/ChurnRiskClassificationSparkML/1/model"}
13 # create spark context
14 spark = SparkSession.builder.getOrCreate()
15 sc = spark.sparkContext
16
17 # load the input data
18 input_data = os.getenv("DSX_PROJECT_DIR") + args.get("dataset")
19 dataframe = SQLContext(sc).read.csv(input_data, header="true", inferSchema = "true")
20
21
22

```

\_152. The evaluation job is scheduled. Scroll down and see the job **Running**.

ID	NAME	TARGET HOST	TRIGGERED BY	STARTED AT	DURATION (S)	RESULT
1539900387-999	Run 1	Local instance	admin	18 Oct 2018, 6:06 PM	In progress	<span style="color: green;">●</span> Running

\_153. Wait for this job to complete and the Result should show as **Success**.



- \_\_154. Click [TradingCustomerChurn](#) project [top left corner].
- \_\_155. Click [Assets ⇒ Data Sets](#).
- \_\_156. Click [TradingCustomerSparkMLEval](#) to open and preview it.

ID	CHURNRISK	GENDER	STATUS	CHILDREN	ESTINCOME	HOMEOWNER	AGE
4	High	F	M	2	52004	N	25
7	High	M	M	0	19749	N	60
8	Medium	M	S	1	57626	Y	44
9	Medium	M	M	2	20078	N	33
15	Low	F	M	2	28220	N	39
18	Low	M	S	2	89459	N	53

- \_\_157. Notice the high-risk category prediction – which matches with the business analysis through the heat map.

	<p><b>Note:</b> The journey does not end here but is the beginning. The success is going to be through a repeatable process of prepare, test, evaluate and deploy the model.</p> <p>The platform provides agility that you need from ICP-D.</p> <p>The next lab exercise will show the deployment and the repeatable process that can fit into your CICD pipe line of DevOps.</p>
---	---

**Lab Credits: Rui Fan and Anjali Saha, IBM Data Scientists.**

**\*\* End of Lab 05: Analyze**

## Lab 06 Deploy

In our previous lab exercise, we created a model, tested it and scored it using the ICP-D Web UI. Business value is achieved when a successful model is deployed and easily updated.

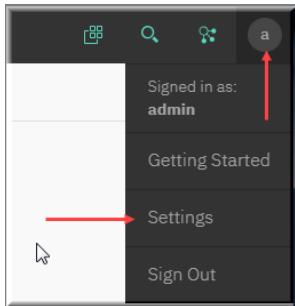
In this lab exercise, we describe the procedure of deploying our model so that it is available for use by application developers.

We will have a much better appreciation of the capabilities that ICP-D platform provides once we finish this lab exercise. We will show automation in harvesting the wealth of information from our data through **Analytics Modernization**.

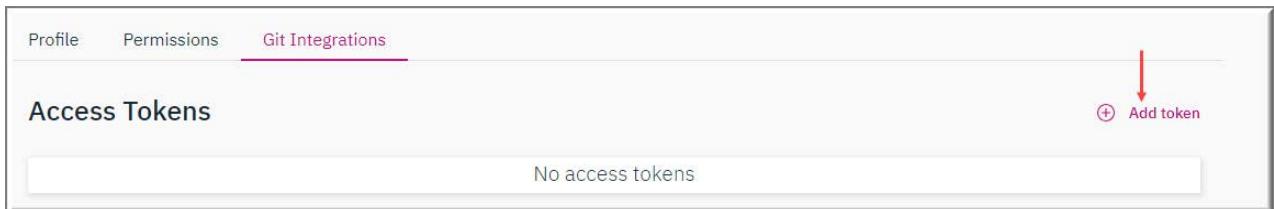
Let's complete this journey with the last pillar of our ICP-D platform.

### 6.1 Integration with Source Control

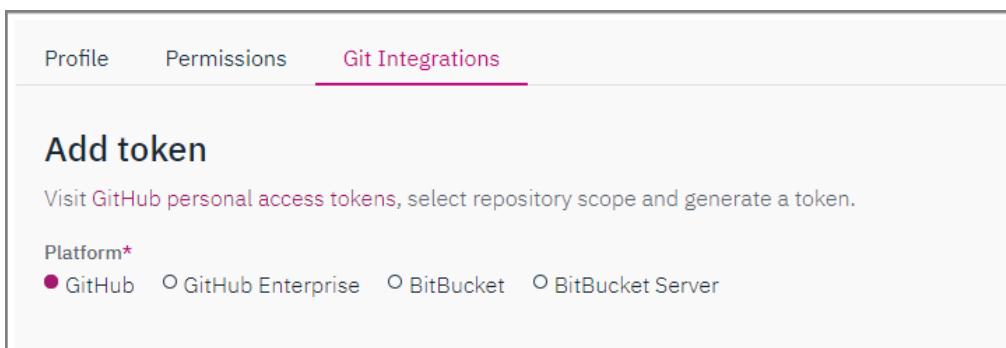
- 1. Click . This is located towards right on the menu bar, then click **Settings**.



- 2. Click **Git Integrations** and click **Add Token**.

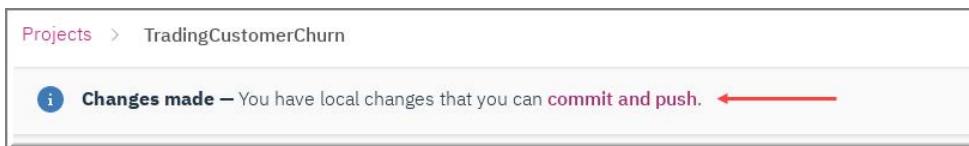


- 3. Notice that you can manage repository with **GitHub**, **GitHub Enterprise**, **BitBucket** and **BitBucket Server**.

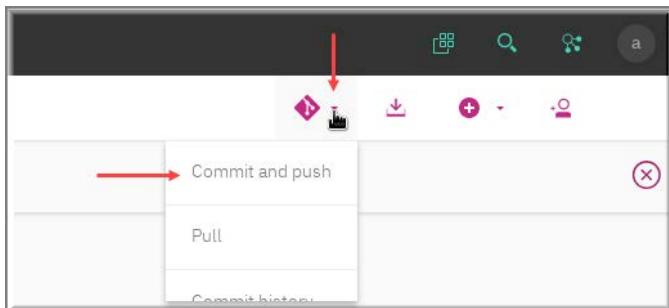


## 6.2 Deploy the Model

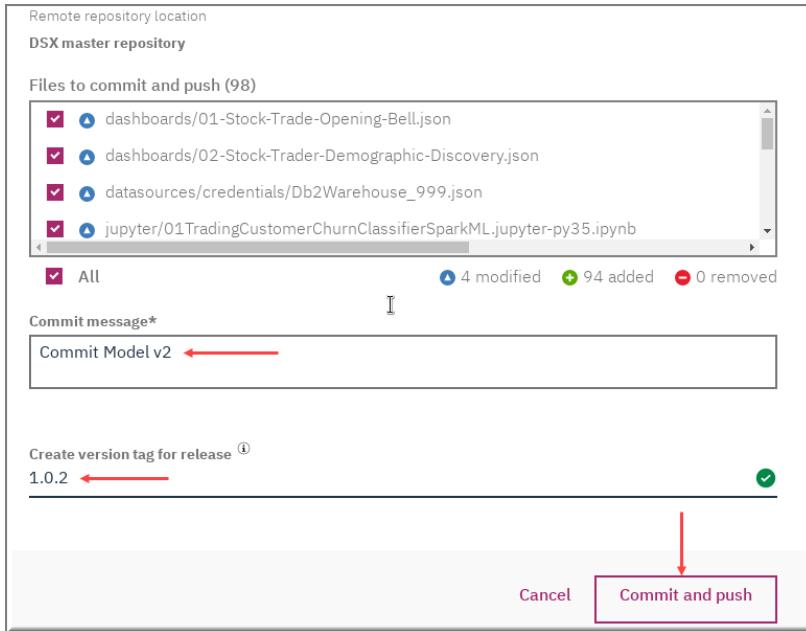
- 4. You can think of model deployment as the equivalent of writing a self-service application that takes the model and makes it available through a REST API interface. Application developers will access and consume the model through the same interface.
- 5. This is a manual process in most organizations. But ICP-D automates deploying and maintaining models without writing a single line of code.
- 6. The IBM Cloud Private platform provides an easy method to deploy the model by eliminating the need for the following:
  - Write the code to perform this capability and use a runtime to deploy it.
  - The runtime could be bare metal machine that requires OS installation, provisioning, network, storage and so on.
  - The runtime could be spinning up a virtual machine in a VMware on Intel or an IBM POWER VM® on a POWER platform.
  - The runtime could be deploying a Docker image – which requires someone to build the image and deploy it on some platform.
- 7. Each of the above requires manpower and machine resources. Using the IBM ICP-D platform, you can bypass this and quickly harvest insight from your data in a repetitive fashion by integrating it with your CICD.
- 8. Let's examine the automated delivery and node deployment which you can integrate with your DevOps implementation.
- 9. Click [Projects](#) from the left menu bar.
- 10. Click [TradingCustomerChurn](#) to open it.
- 11. You might see this message.



- 12. If you do, click [Git](#) icon and click [Commit and push](#).



- \_\_13. Type: *Commit message Commit Model v2*, type version tag as: **1.0.2**, and click **Commit and push**.



- \_\_14. It may take a minute or two. Wait for it to complete. When completed. The following message displays:

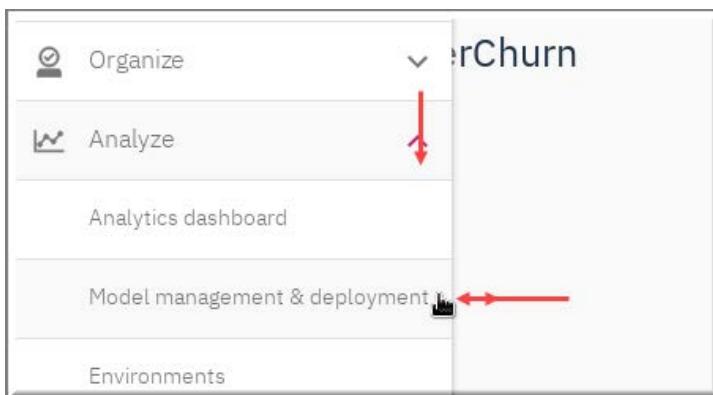


- \_\_15. Click **View**.

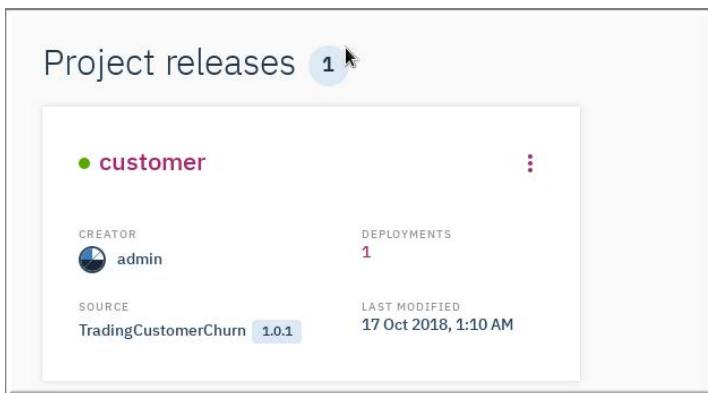
- \_\_16. Scroll down to see the collaboration team who has worked on this including: admin (you), Anjali, Rui, Kikuchi, Rubin and so on.

## 6.2.1 Model Management and Deployment

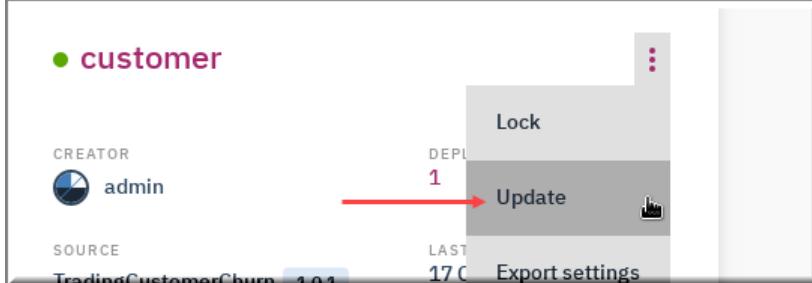
- \_\_17. Click **Analyze** ⇒ **Model management and deployment**



- \_\_18. You would notice that a customer release is already up and running.



- \_\_19. For the purpose of the Executive Demo, we deployed the model already. We now use the same [Project Release](#) to deploy the second version.
- \_\_20. Click the three vertical dots and click [Update](#).



- \_\_21. Select [TradingCustomerChurn](#) from the drop-down menu and Tag [1.0.2](#).

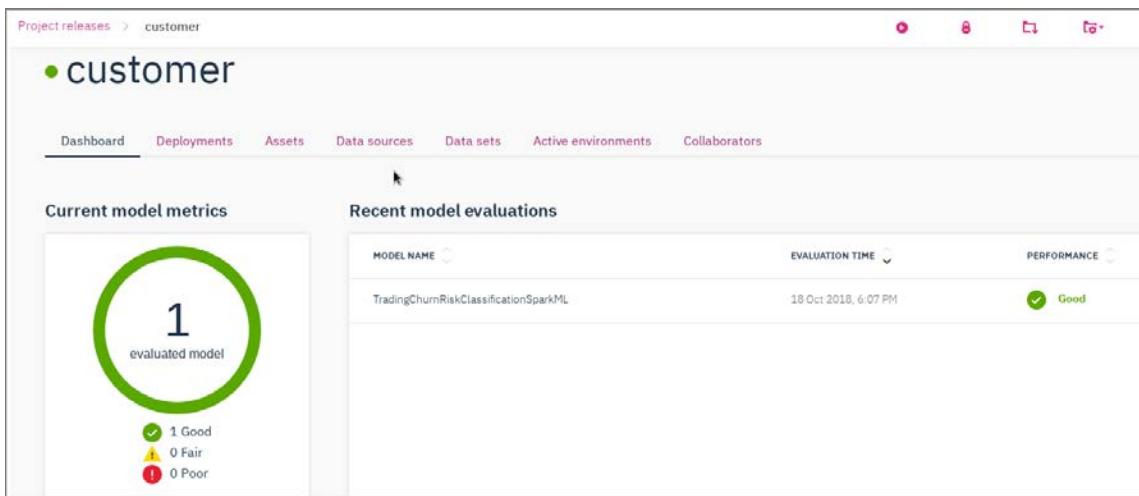
Name  
customer

Route  
award

Source project \*  
TradingCustomerChurn

Tag \*  
1.0.2

- \_\_22. Notice the *Name* **customer** and *Route* **award**. These are the name and route we used to deploy the model. The name identifies the model and route is used by the proxy within ICP-D to route the request.
- \_\_23. Click **Next** to see the list of assets to be updated.
- \_\_24. Click **Update**.
- \_\_25. You will see current model metrics and the evaluations that we ran in our previous lab exercises.



- \_\_26. Click **Deployments**

The screenshot shows a dashboard titled 'customer'. At the top, there are tabs for Dashboard, Deployments (which is selected), Assets, Data sources, Data sets, Active environments, and Collaborators. A search bar is located at the top right. Below the tabs is a table with columns: NAME, ASSET, TYPE, VISIBILITY, DATE STARTED, and AVAILABILITY. One row in the table is highlighted, showing 'customer' as the name, 'TradingChurnRiskClassificationSparkML' as the asset, 'Web service' as the type, and 'Enabled' as the availability status. An arrow points to the 'customer' entry in the table.

- \_27. The model is deployed and enabled.



**Note:** ICP-D platform builds a Docker container and using the underlying IBM Cloud Private, the container and its Kubernetes pod are deployed onto the cluster. Once deployed they are ready to service the API requests and provide results back from the model.

This is an example of the automation that our ICP-D platform provides in support of integrating data and application development pipelines.

- \_28. Click [customer](#).

- \_29. The [Overview](#) section shows the [POST](#) REST API that an application can consume.

The screenshot shows the 'Overview' section of a deployment page. It features two tabs: 'Overview' (which is selected) and 'API'. Below the tabs, there is a large button labeled 'POST' in white text, followed by a URL: 'https://192.168.142.102:31843/dmodel/v1/award/pyscript/customer/score'. To the right of the URL is a small icon of a browser window.

### 6.3 Test Deployment

- \_30. The model deployment through ICP-D is secured with an authentication token. Any applications using the model need to authenticate themselves by providing this token.
- \_31. Scroll down to the **Function** metrics. These give an overview of the REST API consumption metrics.
- \_32. Click [API](#) and scroll down to the [Request](#) section. Look for [INSERT\\_VALUE](#) and change it to [100](#). Click [Submit](#).

## Request

**Function name \***



---

**Body \***

```
{"input_json_str": [{"TOTALUNITSTRADED":23,"CHILDREN":2,"TOTALDOLLARVALUETRADED":5030,"ID":4,"NETREALIZEDGAINS_YTD":"INSERT VALUE","DAYSSINCELASTTRADE":19,"SMALLESTSINGLETRANSACTION":125,"GENDER":"F","DAYSSINCELASTLOGIN":2,"ESTINCOME":52004,"LARGESTSINGLETRANSACTION":1257,"PERCENTCHANGECALCULATION":3,"AGE":25,"STATUS":"M","HOMEOWNER":"N"}]}
```

Change to 100

Clear
Submit

- \_\_33. You can see the **Response** from the model (which was deployed). The **JSON** output is returned by the **REST API**.

```

{
  "result": [
    {"predictions": [High], probabilities: [[0.572096381241254, 0.24983552631578942, 0.17806809244295652]], classes: [High, Low, Medium]}",
    "",
    0
  ],
  "stdout": [
    "+---+-----+-----+-----+-----+-----+-----+-----+
    "|AGE|CHILDREN|DAYSSINCELASTLOGIN|DAYSSINCELASTTRADE|ESTINCOME|GENDER|HOMEOWNER| ID|LARGESTSINGLETRANSACTION
    "+---+-----+-----+-----+-----+-----+-----+-----+
    "| 25|          2|              2|             19|      52004|     F|       N|  4|           1257
    "+---+-----+-----+-----+-----+-----+-----+-----+
    +-----+-----+-----+-----+-----+-----+-----+-----+
    |NETREALIZEDGAINS_YTD|PERCENTCHANGECALCULATION|SMALLESTSINGLETRANSACTION|STATUS|TOTALDOLLARVALUETRADED|TOTALUNITTRADED|
    +-----+-----+-----+-----+-----+-----+-----+-----+
    |          100|                  3|             125|      M|      5030|          23|
    +-----+-----+-----+-----+-----+-----+-----+-----+
    +-----+-----+-----+-----+-----+-----+-----+-----+
    GENDERIndex|GENDERclassVec|STATUSIndex|STATUSclassVec|HOMEOWNERIndex|HOMEOWNERclassVec|         features| rawPrediction|
    +-----+-----+-----+-----+-----+-----+-----+-----+
    | 0.0| (1,[0],[1.0])| 0.0| (2,[0],[1.0])| 0.0| (1,[0],[1.0])|[1.0,1.0,0.0,1.0,...|[11.4419276248250...|
    +-----+-----+-----+-----+-----+-----+-----+-----+
    +-----+-----+-----+-----+
    probability|prediction|predictedLabel|,
    +-----+-----+-----+
    |[0.572096381241254...|      0.0|      High|,"
    +-----+-----+-----+
  ],
  "stderr": []
}

```

- \_\_34. Notice within the above response, the prediction probability is returned, based on the input JSON that is given to the model.

\_\_35. Click [generate code](#)

The screenshot shows a deployment interface with the following sections:

- ALLOCATED CPU:** Unallocated
- ALLOCATED MEMORY:** Unallocated
- REQUESTS:** 7
- Response:** A code editor containing a JSON response object.
- generate code:** A button with a double-headed arrow icon located above the response code editor.

```

1 {
2   "result": [
3     {"\predictions": ["High"], "\probabilities": [[0.572096381241254, 0.2498355
4     "",
5     0
6     ],
7     "stdout": [

```

- \_\_36. The attached curl command is generated. This can be used to test the REST API and can be given to the application developers for them to implement the API in their application.

The screenshot shows a modal dialog box titled "Generate code" containing a code editor with a curl command:

```

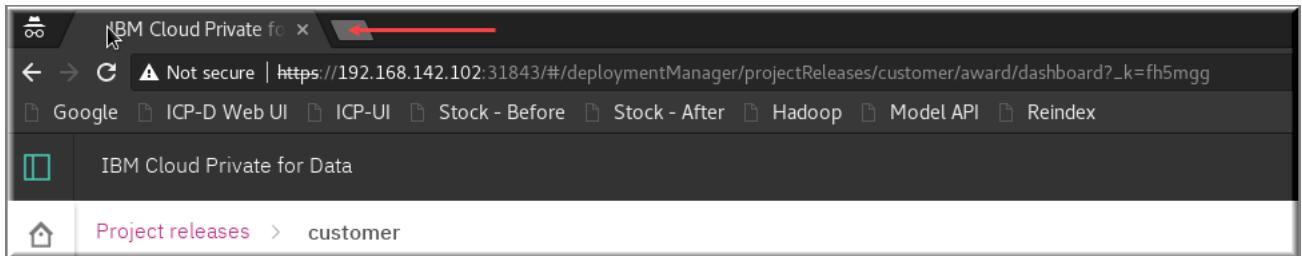
1 curl -k -X POST \
2   https://192.168.142.102:31843/dmodel/v1/award/pyscript/customer/score
3   -H 'Authorization: Bearer eyJhbGciOiJSUzI1NiIsInR5cCI6IkpXVCJ9eyJ1c2
4   -H 'Cache-Control: no-cache' \
5   -H 'Content-Type: application/json' \
6   -d '{"args":{"input_json_str":"[{\\"TOTALUNITSTRADED\\":23,\\"CHILDREN\\"

```

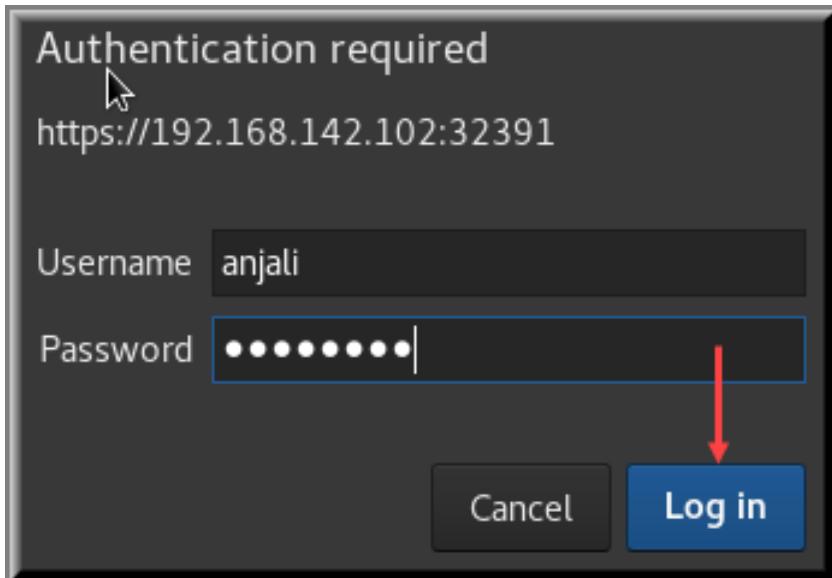
A red arrow points from the "generate code" button in the previous screenshot to this modal. Another red arrow points from the "Close" button in the modal to the "Close" button in the bottom right corner of the slide.

## 6.4 Consume Deployment in Stock Trader Application

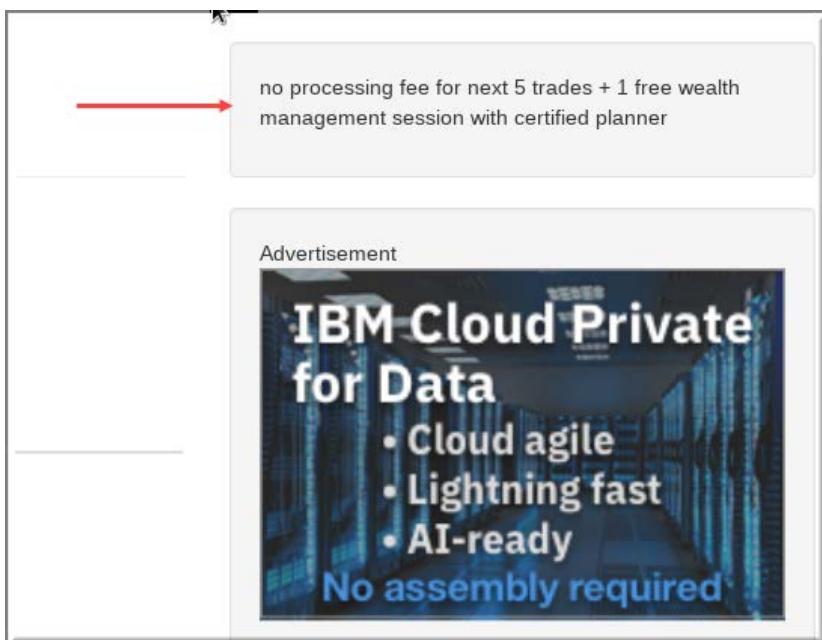
- \_\_37. The machine learning model was created by data scientists and deployed using ICP-D platform. Now, let's see how a modern microservices-based application can consume it.
- \_\_38. We have shown the [Stock – Before](#) and [Stock – After](#) application in the Executive Demo lab.
- \_\_39. Let's revisit it more detail.
- \_\_40. Open a new browser window tab.



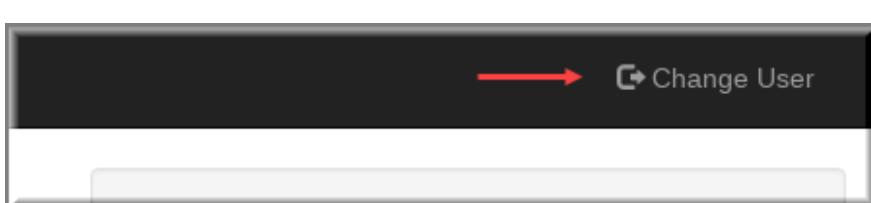
- \_\_41. Click **Stock – After** desktop icon.
- \_\_42. Log in as *Username*: **anjali** and *Password*: **password**. Click **Log in**.



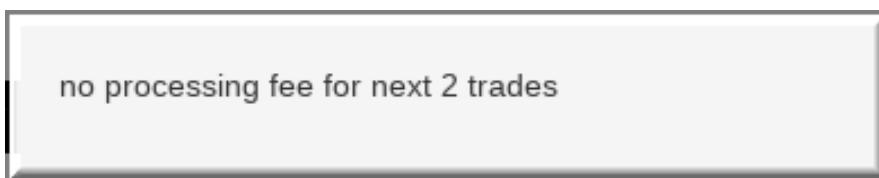
- \_\_\_43. Notice the offer box that shows after logging in to the application.



- \_\_\_44. Based on user ID of the login, we get our demographic information about that user from our database, provide that information to a REST API call from the backend of this microservice. The demographic information is given to the model running on the ICP-D platform. The model returns a response (churn prediction value), based on the response business logic, is invoked to present a retention offer to the user. All of this is performed in real time, the user then logs into the application.
- \_\_\_45. Click [Change User](#).



- \_\_\_46. Type User: [dan](#) and Password: [password](#) and see the offer.



- \_\_\_47. The offer is based on the risk assessment of the ID of the user.
- \_\_\_48. Let's go further into the REST API call.
- \_\_\_49. Open one more browser tab and click [Model API](#) desktop icon.

\_50. Click **GET** to test the REST API.



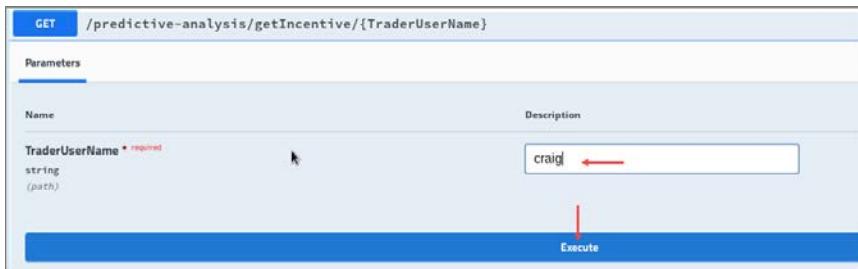
\_51. Select the HTTPS version of the **REST API** as **http** will be rejected by the proxy server of ICP-D.



\_52. Click **Try it out**.



\_53. Type **craig** and click **Execute**.



- \_\_54. Notice the responses: [curl](#) command with GET verb, 200 response code, and the response body which is consumed by the application.

The screenshot shows a curl command interface. The 'Request URL' field contains: `curl -X GET "https://predictive-analysis-service.stocktrader.svc.cluster.local:9443/predictive-analysis/getIncentive/craig" -H "Accept: application/json"`. The 'Server response' section shows a 'Code' of 200 and a 'Response body' containing: `{ "offer": "no processing fee for next 5 trades" }`. The 'Response headers' section includes: `content-language: en-US  
content-length: 47  
content-type: application/json  
date: Fri, 19 Oct 2018 19:07:55 GMT  
x-powered-by: Servlet/3.1`.

- \_\_55. Try user [foo](#). You will notice no response as we have no demographic information for this user.

The screenshot shows a curl command interface. The 'TraderUserName \* required' field has 'foo' entered. The 'Request URL' field contains: `curl -X GET "https://predictive-analysis-service.stocktrader.svc.cluster.local:9443/predictive-analysis/getIncentive/foo" -H "Accept: application/json"`. The 'Server response' section shows a 'Code' of 200 and a 'Response body' containing: `{ "offer": "" }`.

## 6.5 Conclusion

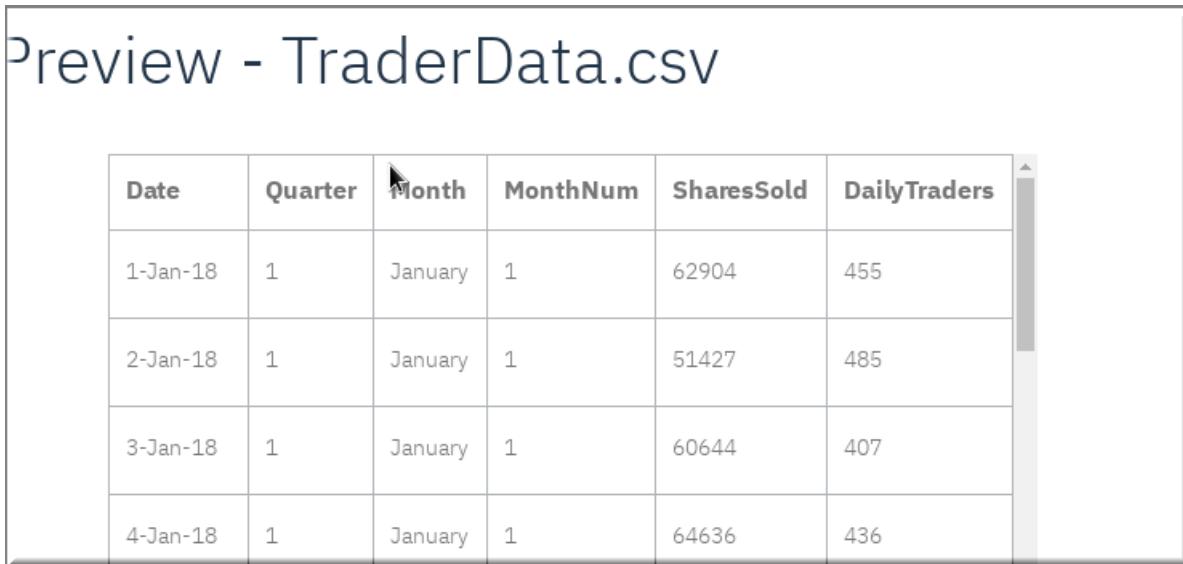
- \_\_56. As an organization, we collect, and have access to, more data than we can possibly analyze.
- \_\_57. Data is spread all over, in many silos, and is hard to find. Locating that one person in an organization who has the breadth and depth of knowledge of our data is like looking for a unicorn.
- \_\_58. IBM Cloud Private for Data can be that unicorn. It provides an integrated data platform where data scientists can shop for data, Data engineers help to prepare the data. Business analysts provide visual insights of data to the stakeholders. Data stewards can discover, classify, provide governance, rules and policies to meet regulations.
- \_\_59. Turning your data into a corporate asset with end-to-end data integration on modern, cloud-native platform. We have completed the IBM Journey to Cloud and AI: Analytics Modernization.

**\*\* End of Lab 06: Deploy**

---

## Appendix A. Stock – Opening Bell Dashboard

- \_\_1. Let's begin by analyzing current trends of customer visits and daily trades in our Stock Trader application.
- \_\_2. We requested data engineers to provide a file with historical totals of visits and trades for the past year. This file was provided and deposited in our project where we all collaborate with each other.
- \_\_3. The data steward guidelines of the company do not allow data to be migrated to individual's laptop for analysis.
- \_\_4. Let's build a dashboard with this data to see if we see any trends.
- \_\_5. Click [Projects](#) from the left menu bar.
- \_\_6. Click [TradingCustomerChurn](#) project.
- \_\_7. Click [Assets](#) and then click [Data Sets](#).
- \_\_8. Notice [TraderData.csv](#) available in our project. We will use this data file to build the [Stock – Opening Bell](#) Dashboard.
- \_\_9. Click the link [TraderData.csv](#) to preview the data.

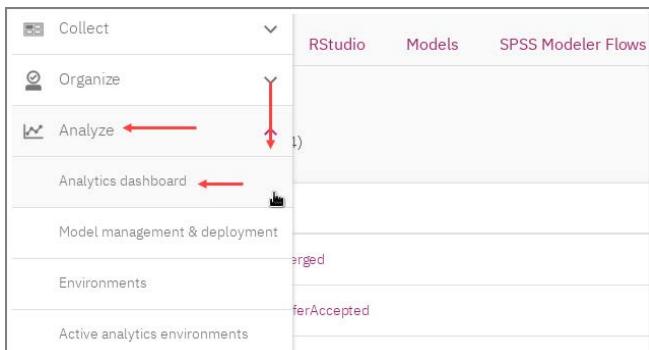


The screenshot shows a 'Preview - TraderData.csv' window. The table has the following data:

Date	Quarter	Month	MonthNum	SharesSold	DailyTraders
1-Jan-18	1	January	1	62904	455
2-Jan-18	1	January	1	51427	485
3-Jan-18	1	January	1	60644	407
4-Jan-18	1	January	1	64636	436

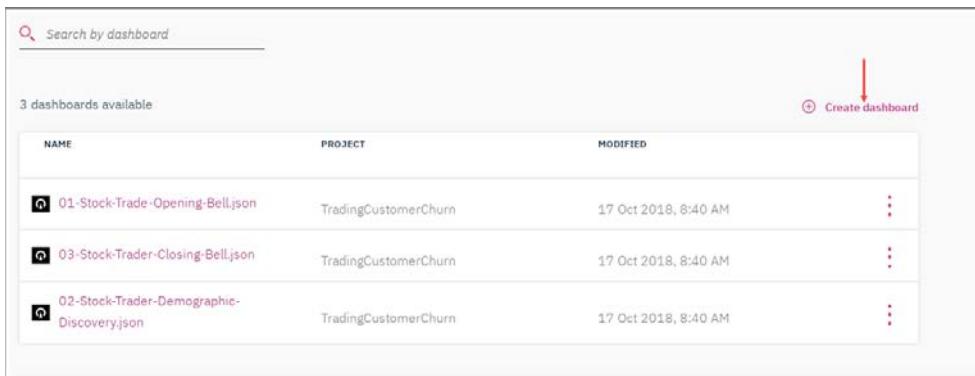
- \_\_10. Click [Close](#) after preview.

\_\_11. Click **Analyze**  $\Rightarrow$  **Analytics dashboard** from the left menu bar.



\_\_12. We used the **01-Stock-Trade-Opening-Bell** dashboard in the Executive Demo lab and we will build this exact same dashboard here.

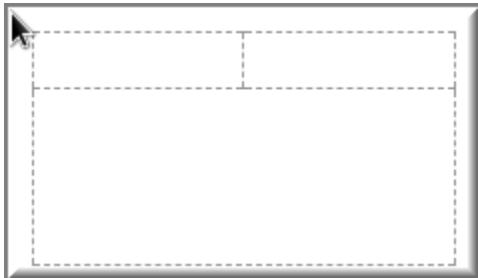
\_\_13. Click **Create dashboard**.



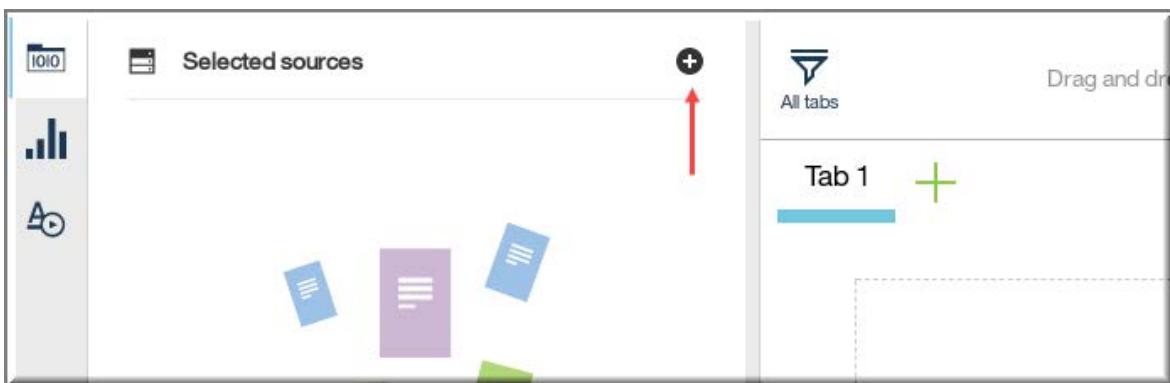
\_\_14. Type name: **My-Own-Stock-Opening-Bell** and choose **TradingCustomerChurn** project from the drop-down menu.

Name*	My-Own-Stock-Opening-Bell	19
Description	Type your description here	
Project*	TradingCustomerChurn	

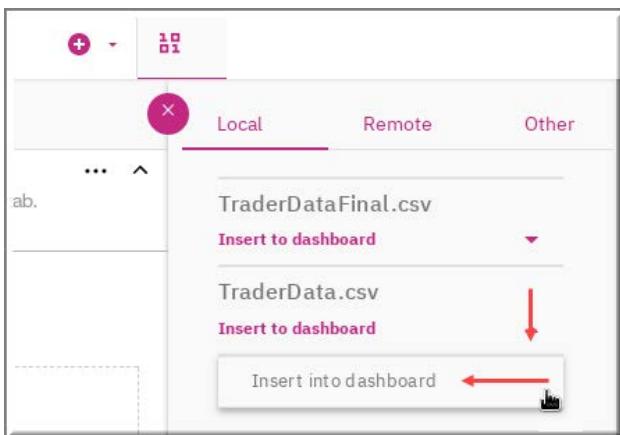
- \_\_15. Click **Create**. [bottom-right corner.]
- \_\_16. You are now presented with a choice of canvas templates. Select the one that looks as shown:



- \_\_17. Click **OK**.
- \_\_18. From the **Selected sources** area near the top left, click



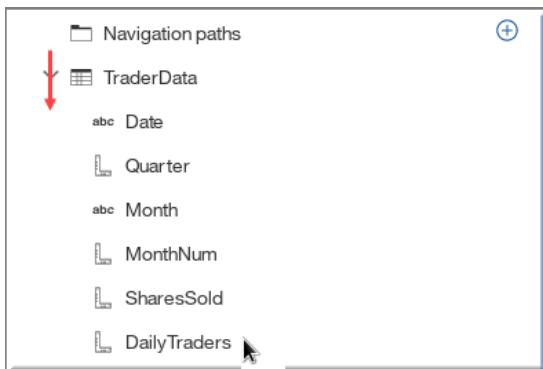
- \_\_19. The right pane will slide into the canvas. Click the down arrow of **TraderData.csv** and click **Insert Into dashboard**.



\_20. You will now see **TraderData** in the **Selected sources** area.

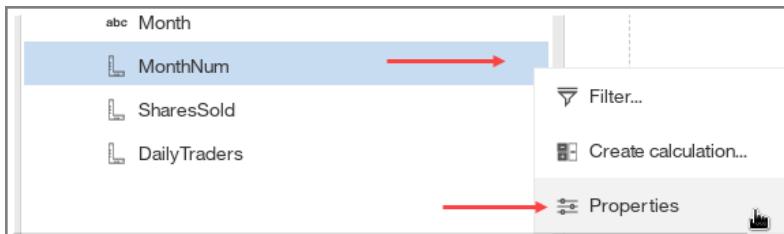


\_21. Click **TraderData** and expand it to show all the data items in that data source file.

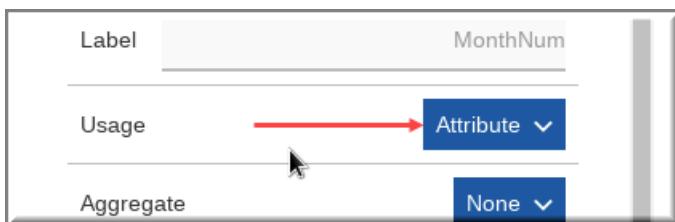


\_22. Let's do some housekeeping so that our data is represented correctly within our dashboard.

\_23. Click the **MonthNum** ellipse, then select **Properties** from the flyout menu.

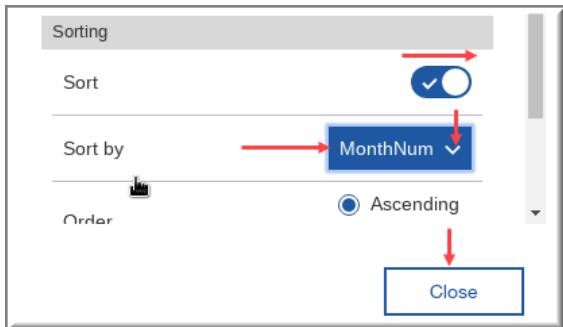


\_24. **MonthNum** in this case is not a measure. It's simply an attribute on which we sort our months. Change **MonthNum** usage to be an **Attribute** and click **Close**.

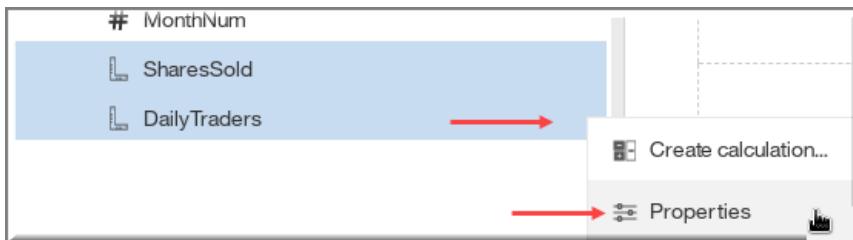


\_25. Now, let's select the Properties of **Month** by selecting the ellipse next to **Month**.

- \_\_26. Slide the **Sort** icon to turn it on and sort by **MonthNum** in **Ascending** order, then click **Close**.



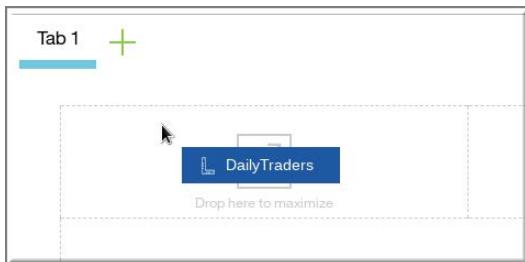
- \_\_27. Next, we set the aggregation of our measures for totaling.
- \_\_28. Hold the shift key down and click **SharesSold** and **DailyTraders** to select both. Click the ellipse and then click **Properties**.



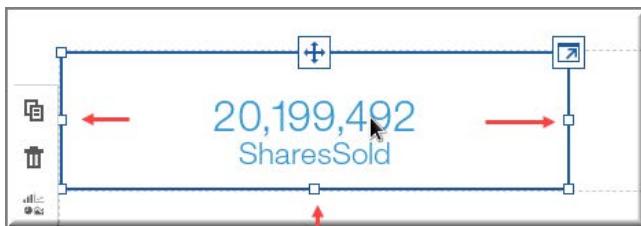
- \_\_29. Change Aggregate to **Total** from the drop-down menu and click **Close**.

## 8.1 Build Dashboard

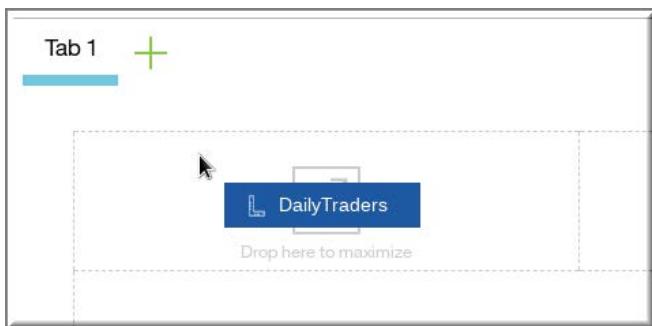
- \_\_30. Drag **SharesSold** to the top left box, hovering over the *Drop here to maximize* area when it turns blue. This gives us a total of **Shares sold** over all time.



- \_\_31. If you did not drop at the right place, adjust the guide to match with the outline of the left box.



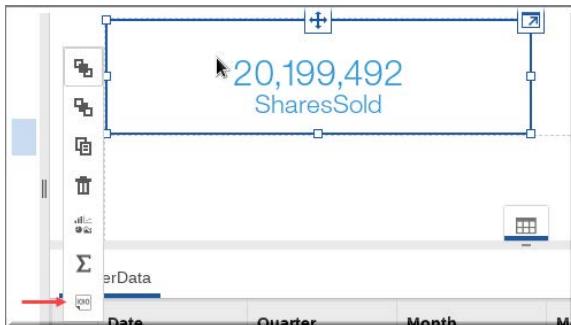
- \_\_32. Drag **DailyTraders** to the top right box, hovering over the *Drop here to maximize* area when it turns blue as well. This maximizes this metric in this box. This gives us a total of trades over all time.



- \_\_33. After you complete both top boxes, the dashboard should display as shown:



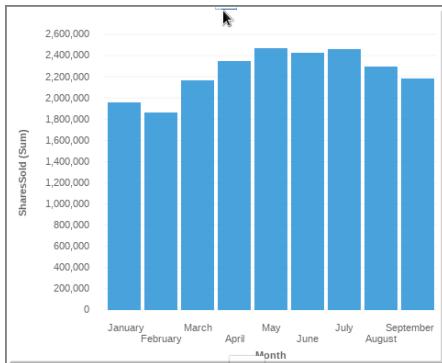
- \_\_34. Click **SharesSold** metric and select the **format** from the flyout menu.



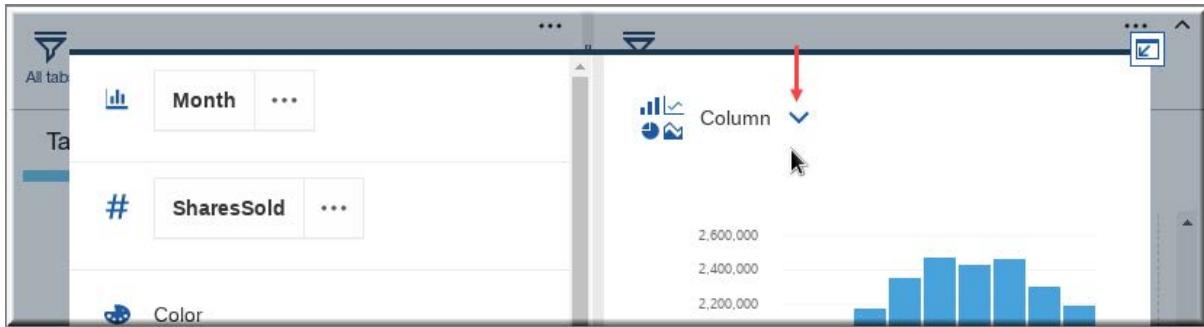
- \_\_35. Click **Abbreviate**.



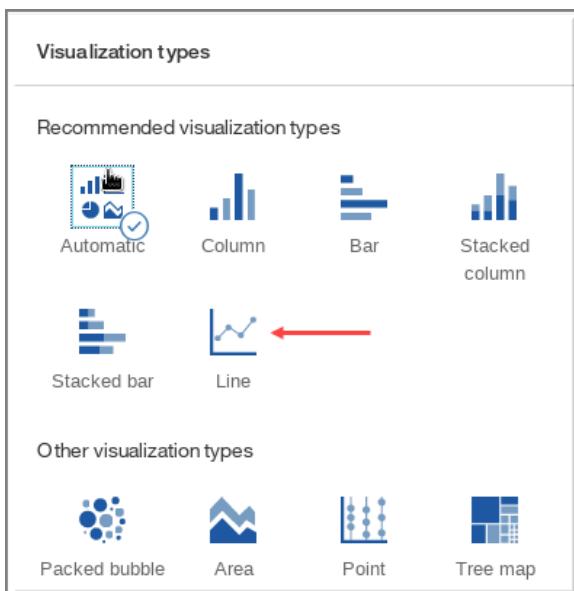
- \_\_36. Hold the shift key down, click **Month** and **SharesSold** and drag the two onto the lower canvas area. [Do not drop in the *Drop here to maximize* area this time.]



- \_\_37. Click (on top right of the chart). This will expand the chart and allow you to make adjustments.
- \_\_38. Click the drop-down menu to change the type of chart from Column to **Line Chart**.



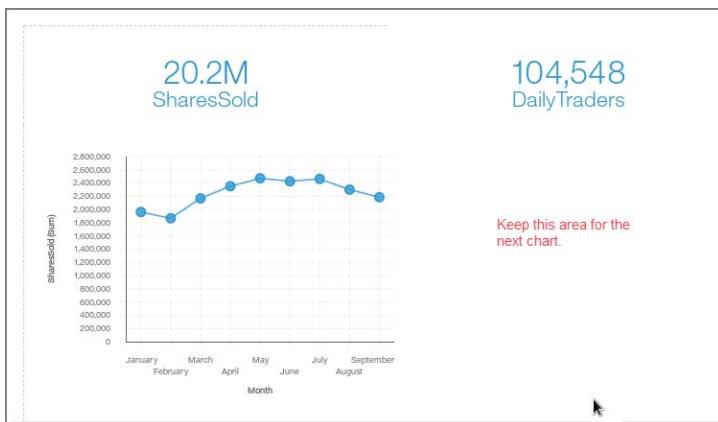
- \_\_39. Select **Line** visualization type.



- 40. Now select the **Expand** button again to return it to its original size.

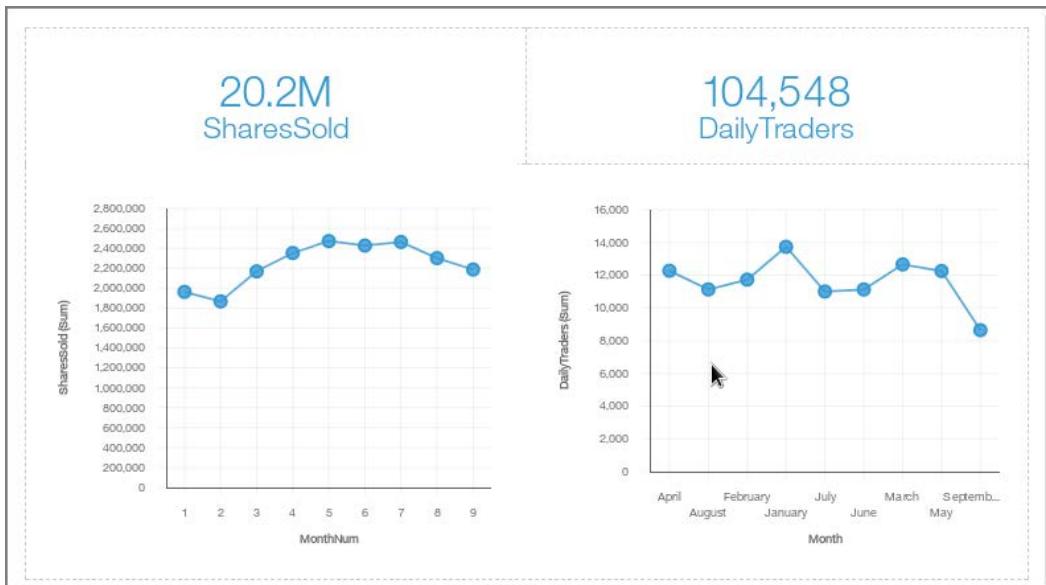


- 41. Drag the chart so it aligns with the left side of the template area and squeeze in the right side of the template so it aligns with the middle of the screen.

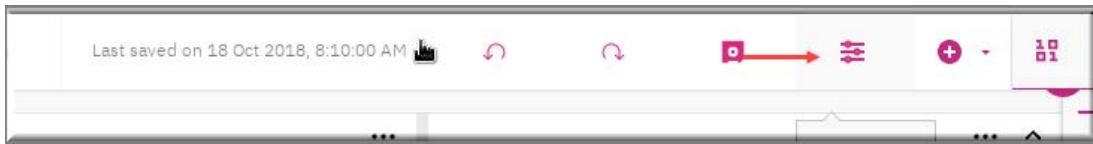


- 42. Hold the shift key down, click **Month** and **DailyTraders** and drag the two onto the lower canvas area. [Do not drop in the *Drop here to maximize* area this time.]
- 43. Change the visualization from **Column** to **Line** as we did with the first graph. [Click **Expand**, change the type and click the same button to bring its in original size.]
- 44. Adjust the top, bottom, left and right of the chart boundaries so that all the boxes are aligned.

- 45. Your final dashboard should display as below.



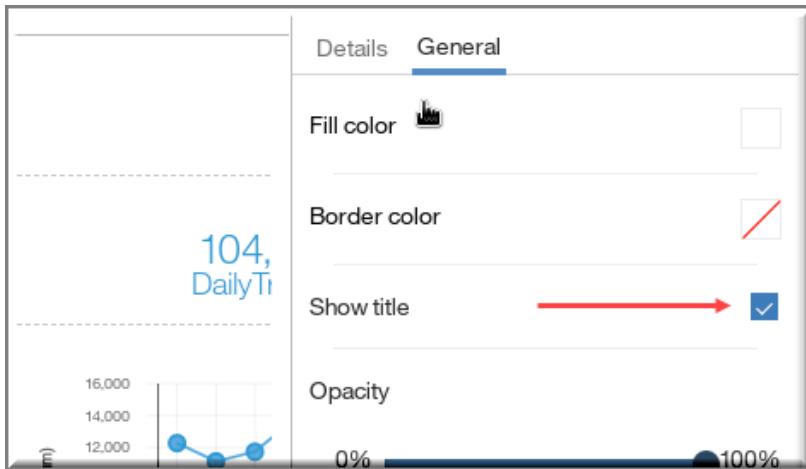
- \_\_46. Select the bottom left chart and then select the **Properties** button at the top.



- \_\_47. Check **Show Value Labels** from the Details tab.



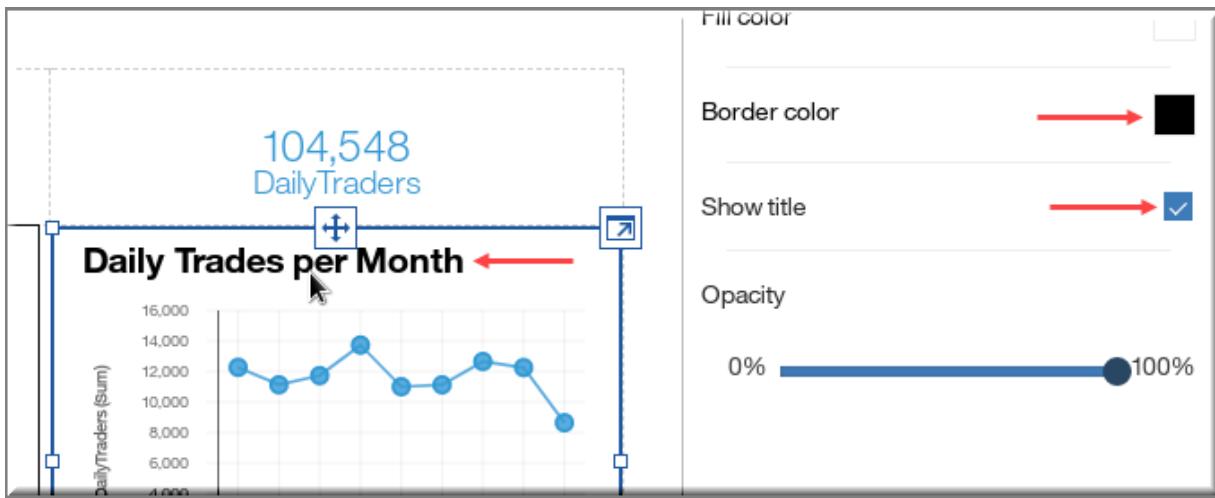
- \_\_48. Now select the **General** tab and check to **Show title**.



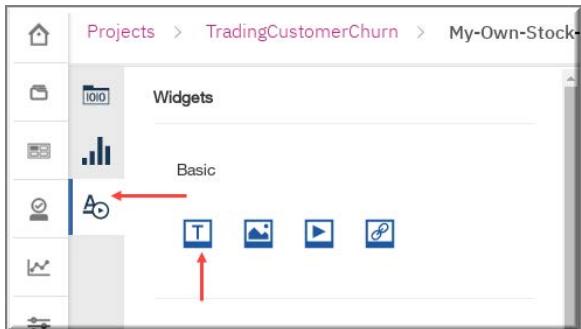
- \_\_49. Enter **Shares sold per month** for the title of the left chart and change the border color to **Black**.



- \_\_50. Click the bottom right-hand chart and enter **Daily Trades per month** for the title of the right chart and change the border color to **Black**.

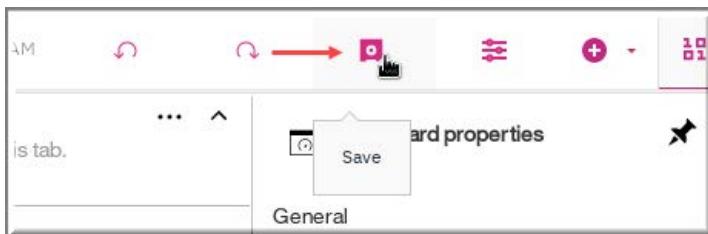


- \_\_51. On the left side of the canvas, click the **Widgets** menu and drag a **Text box** between the top two charts.



- \_\_52. Title it **Stock Trades**. Adjust the box so the title stands out.

- \_\_53. Click **Save**.



- \_\_54. You notice from examining the charts that **Shares sold** is relatively flat and daily trades are falling off. We need to use IBM Cloud Private for Data to discover the WHY behind this trend.

## Appendix B. Integrated Data Platform

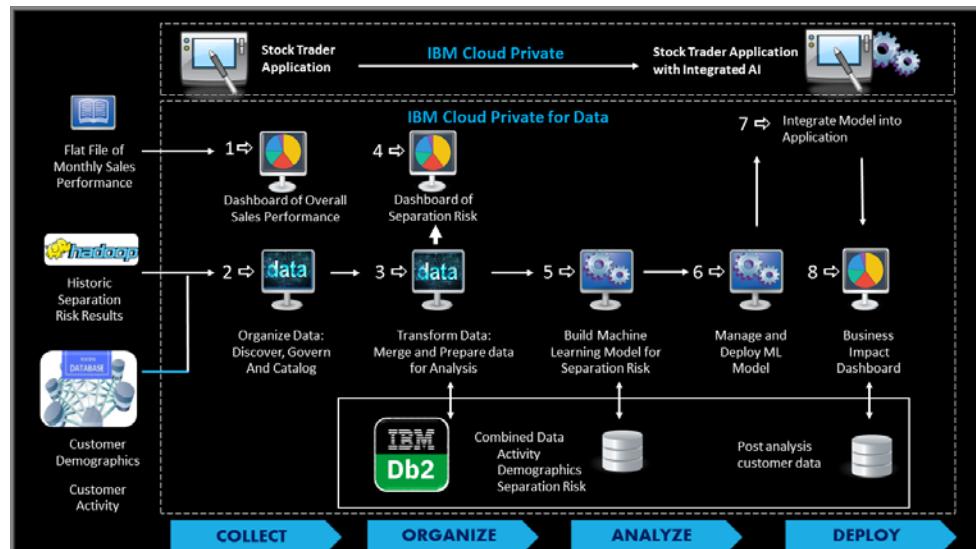
This workshop covers end-to-end analytics modernization with examples to help you on this journey.

Regardless of the role you play in your organization, you are probably aware of difficulties faced by people throughout the Lines of Business and IT. Some of these may sound familiar to you:

- Our data is very siloed. It is a huge challenge to search for it, understand the 'business context' of it, determine its quality and if it can be trusted, and knowing the right people across our business to resolve these challenges.
- The data owner can't give me access to production data, so we have to copy full data sets into our own environment. This is time consuming and a duplication of effort. Also, the data is stale the moment I get it. Note: A case study found that one bank had 82 data warehouses and 63 percent of its data is a duplicate of other data.<sup>1</sup>
- The data we get from different silos is unorganized and it is a complex effort to clean and prepare it. It takes 80 percent of our time to just prepare it. For some projects, this takes many months. In one case study, it was found that 73 percent of the time is spent in Extract Transform and Load (ETL), working with DBAs and infrastructure support teams to get the data.
- We do not have borderless data and we do not know the lineage of the data.
- Building good models takes good data. I'm not always confident in the data I'm using because it's curated by tribal knowledge. Once I do find a good source, or spend effort curating data, it's not discoverable and readily available to others.
- Once I create a well-performing, predictive model, I struggle with the operational aspects of making the model easily consumable by application developers and keeping the production model tuned over time. Multiply this by dozens (or hundreds) of models from multiple lines of business across our enterprise and you see our struggle.
- New programming languages and their resulting tools for data scientists are becoming available annually. We need a way to efficiently support all of this tooling, and the multiple versions of each individual tool.

Let's start working with an example and see how ICP-D helps address these challenges.

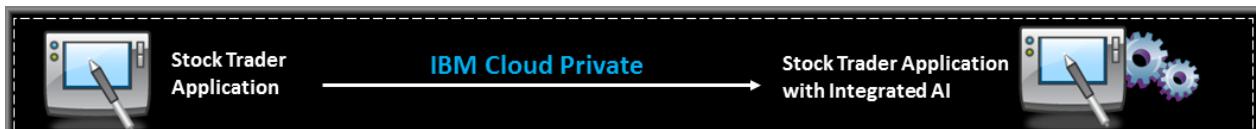
For the purpose of this workshop, let's examine the following workflow involving **Business Analyst**, **Data Engineer**, **Data Scientist** and **Data Steward**. Each one of these personas has a role in each of the four lab exercises of **Collect**, **Organize**, **Analyze** and **Deploy**.



<sup>1</sup> Eric Watson, IBM Corporation

## 8.2 Application with Analytics Modernization

- 55. Boatswain is the hypothetical company for which we work. It is an on-line stock trading company. We have approximately 100,000 customers who maintain their portfolios and buy/sell stocks through our website. One of our bigger challenges is retaining customers. Their cost of switching to a new provider is very low. As a business, we need to understand our 'Customer Separation Risk'. Specifically, how many customers are we losing to competition, how much is this costing us, what are the main drivers and/or customer profiles who leave, what actions will we take to retain these individuals, and what are the costs and benefits of our corrective actions?
- 56. Our corporate strategy is to offer 'Individually customized' retention offers and other services so that we differentiate our website and customer experience from our competition. With this approach we will 'Market to a Segment of One Customer' versus marketing to groups or subsegments of our overall customer base.
- 57. As a first step in our journey, we successfully ported our core customer-facing application from a 'legacy monolithic architecture' to a 'cloud-native, microservices architecture' running in Docker containers on a Kubernetes cluster. The next step in our journey is to add artificial intelligence so that we can effectively 'Market to a customer segment of one customer'.
- 58. The name of our application is [Stock Trader](#). At the start of this workshop, the application does not yet include AI and we are calling the application "Stock – Before". This application can be deployed in any Kubernetes-based cloud environment. For this workshop, we are running the application in [IBM Cloud Private](#).
- 59. The outcome of our workshop is to enhance our application with a machine learning model (artificial intelligence) to drive increased revenue and profit by significantly improving our customer retention problem. We will name this application "Stock – After".
- 60. This is the top box in our workflow.



	<p><b>Note:</b> <a href="#">Stock – Before</a> – there are no individual retention offers.</p> <p><a href="#">Stock – After</a> – we get <b>high</b>, <b>medium</b> or <b>low</b> risk separation probability through our machine learning prediction model for the logged-in user. The offers are made 'in real time' based upon risk.</p> <table style="margin-left: 20px;"> <tr> <td>Low</td><td>– 3 free trades</td></tr> <tr> <td>Medium</td><td>– 5 free trades</td></tr> <tr> <td>High</td><td>– 5 free trades plus a free consult from a wealth management advisor</td></tr> </table>	Low	– 3 free trades	Medium	– 5 free trades	High	– 5 free trades plus a free consult from a wealth management advisor
Low	– 3 free trades						
Medium	– 5 free trades						
High	– 5 free trades plus a free consult from a wealth management advisor						

## 8.3 Collect

The purpose of the Collect lab is to show ways to get access to data that is locked in different silos. To get started, we simply upload a local .csv file to build our first dashboard.

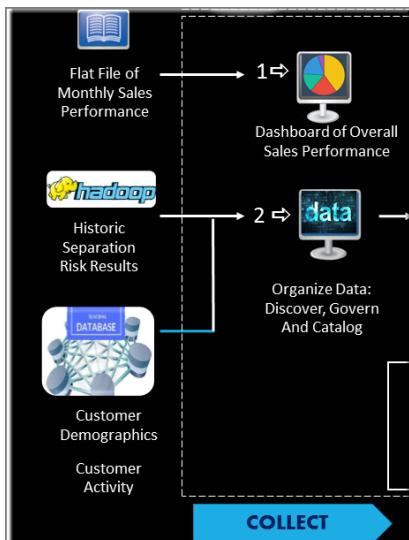
'Passing of data between users' is a common practice in many enterprises that can lead to unmanageability. When data is 'passed around' on thumb drives, it is generally not secure nor governed. It may include Personally Identifiable Information (PII) such as: customer name, age, tax id number, address, and more. This practice can lead to security breaches, tarnishing a company's reputation in the marketplace, and potentially legal action.

ICP-D provides a better way to protect and govern your data, while at the same time making it easily accessible to business analysts and data scientists.

Through this workshop, we will connect to different data sources including: Hadoop and DB2. Note that ICP-D can currently connect to additional sources in a "Secure and Governed way" including: Oracle, SQL Server, Teradata, and Event Store.

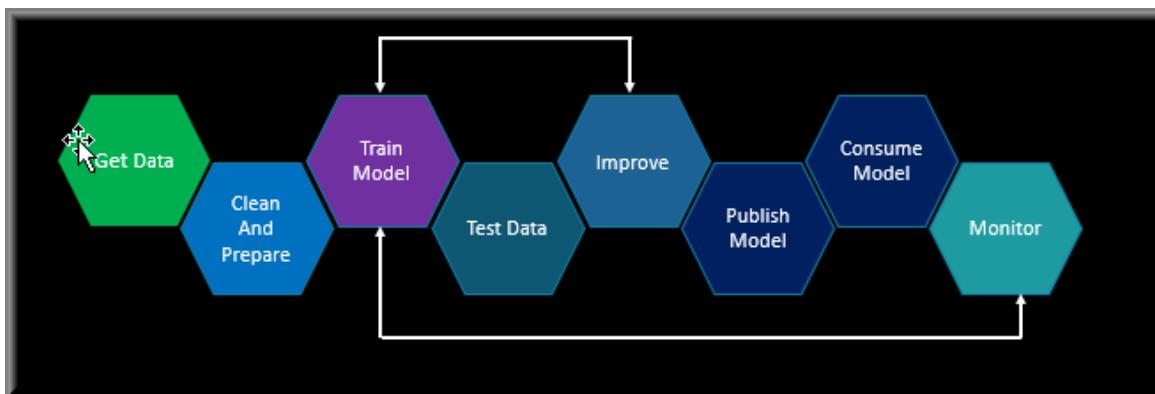
Stay tuned for the next version of this workshop. It will contain new capability on 'virtualizing SQL Queries' across these data sources.

- 61. Boatswain has realized that revenue and profits are declining. This will be highlighted in a Sales Performance Dashboard that you will build. Throughout our workshop, we justify additional investment in our project based upon hidden insight that can be uncovered from the voluminous data dispersed in different silos.
- 62. The business analyst takes the aggregated historical transactional data from our website and displays a dashboard of declining revenue and number of daily traders. We will see this in the lab exercise entitled [Collect Stock – Opening Bell](#) business dashboard.
- 63. As part of this exercise, you log in to the [Stock Trader – Before](#) application and work with the application to understand the workflow and service that we provide to our customers.
- 64. In the Work Map Index – the [Collect](#) process is described as below:



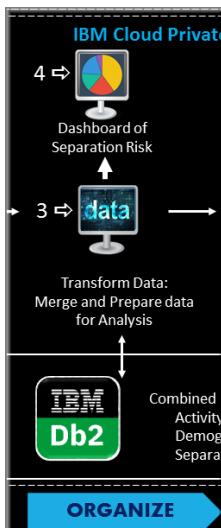
## 8.4 Organize

- 65. The purpose of the **Organize** lab is to show how ICP-D helps to govern data and make it accessible to business analysts and data scientists. One of the biggest challenges companies face when trying to govern their data is to define the inventory of their data. To accomplish this, companies need to execute the painstaking process of looking at all the columns in all of the database tables, understanding the quality of what is in each column, correctly mapping business terms to each column, and making these business terms (and their respective columns) available to business analysts and data scientists in a secure way.
- 66. Once the right data is located, the next step is to combine multiple sources into one file. It has been stated that “Simple models based on a lot of data is better than more complex models based on less data.” Said differently leveraging all available, relevant data sources will produce the best performing model
- 67. ICP-D helps to accelerate these challenges. The Organize Category has a very important role. The repetitive workflows of **Find** ⇒ **Classify** ⇒ **Quality Score** ⇒ **Load to Catalog**, and **Get Data** ⇒ **Combine Sources** ⇒ **Transform** ⇒ **Load**



- 68. A model can be made successful one time with herculean efforts, but this process is not scalable without a strong **Organize** capability. A model can only be as good as the quality of data that it uses.
- 69. Building a good model requires a great deal of time and effort. The **Organize** capability of ICP-D helps you to scale and govern the process and ensures your modeling pipeline is being fed by high quality data.
- 70. In this lab exercise, we will go through the process:
- Gathering of data across disparate sources
  - Classification of data
  - Assignment of terms to data
  - Application of quality rules to data
  - Transformation of data

- \_\_71. The above process is critical for a successful machine learning model development. Without automation, most organizations struggle to successfully implement a clean, end-to-end machine learning pipeline.
- \_\_72. The success of the machine learning model depends on connecting to data, regardless of where it is stored. Classifying data manually is cumbersome. ICP-D provides an automated way of classifying the data based on more than 200 prebuilt classifications.
- \_\_73. Confidence levels are assigned to the data based on automated classifications. This automation process leads to a much faster way of getting to know the data and where it resides.
- \_\_74. Quality rules are applied to the data. For example: Whether or not a credit card number is actually chargeable or valid. Run those rules against the columns, if it matches with the classification of data and if it matches associated business terms.
- \_\_75. Most organizations have been unsuccessful in properly classifying the data simply due to the extensive manual efforts required. These efforts are reduced significantly in this platform.
- \_\_76. The ICP-D platform assigns business [Terms](#) that exist already in the business glossary to the data that it has classified automatically to the data elements. This process is about simply assigning a column name of Tables to the descriptive names in the business glossary automatically.
- \_\_77. Finally - moving that data to the target using transformation capabilities so that it is available to data scientists.
- \_\_78. This capability of [Organize](#) is the strongest link for a successful and repetitive model deployment for the data scientist to [Search](#) and [Shop for Data](#) – a capability most organizations have been lacking due to the enormous amount of manual effort required.
- \_\_79. IBM studies with clients have shown a 300 percent average, 5-year ROI with adoption of automation in [Organize](#) process.
- \_\_80. In the Work Map Index – the [Organize](#) process is described as below:



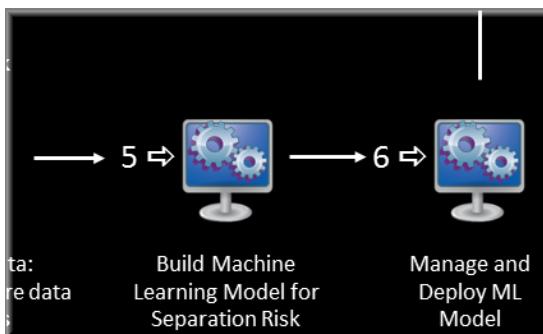
- \_\_81. We will examine the process in detail as described above in [Organize](#) Lab.

## 8.5 Analyze

The purpose of the analyze lab is to show how ICP-D helps data scientists can: accomplish their job using language and their tooling of choice, improve their access to trusted, governed data sources, and be positioned with additional tooling to make their models available in a production environment for consumption by application developers.

	<b>Note:</b> We are using a Jupyter Notebook with a Python runtime to build the machine learning model for this workshop. ICP-D also supports other user-friendly GUI interfaces for both coders and non-coders R-Studio, Scala, Spark, Zeppelin with Anaconda and more.
---	--

- 82. Once the data has been collected and organized properly, with search and classification capabilities, it is much easier for data scientists to shop for the data.
- 83. The ICP-D platform reduces that 80 percent manual labor that is spent on data preparation to 10-20 percent.
- 84. In the Analyze lab exercise, we will develop, test and validate a machine learning model to better predict which customers are at the highest risk to separate from us.
  - The model is built using Jupyter notebook and imports the merged customer table which was a result of hybrid data coming from relational and Hadoop data platforms.
  - The data for the model will be prepared for the machine learning algorithm.
  - 'Separation Risk' will be identified as the critical target.
  - The proper variables will be identified as input predictors and a supervised learning algorithm using data as inputs will be executed.
  - The model will be tested and validated for performance ensuring that we select the best model.
  - The model will be saved to the model repository.
- 85. The ICP-D platform provides an integrated platform for data scientists to work together and collaborate instead of each data scientist working on their own set of analytical tools without strong governance of the data.
- 86. Collaboration using GitHub or BitBucket is integrated into the platform, which brings a cohesiveness to the work culture and helps to automate the CICD pipeline.
- 87. In the Work Map Index – the [Analyze](#) process is described as below:



- 88. We will examine the process in detail as described above in [Analyze](#) Lab.

## 8.6 Deploy

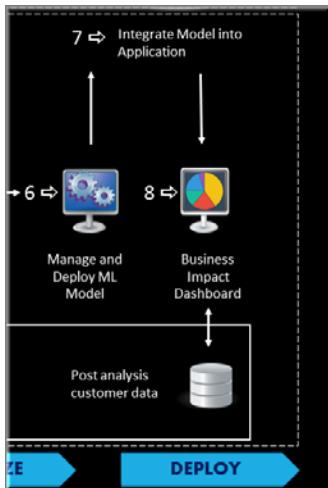
- 89. The purpose of the Deploy lab is to make the model available in a production environment for application developers to consume. We will enhance our Stock Trader application to consume the model. Here you will work with the '[Stock – After](#)' application and witness our new AI capability.



**Note:** The final result of the end-to-end, integrated data platform with machine learning capabilities is demonstrated by consuming the REST API into our existing Stock Trader application.

- 90. When a user logs into the application, in real time, we look up the current demographics of the user from our master file, run the demographics through our machine learning model to get a 'churn prediction' score, then based on that score, we provide the user with one of the following four retentions offers
- No offer
  - Next 3 trades free
  - Next 5 trades free
  - Next 5 trades free and a free wealth advisor consultation
- 91. As a company, we are differentiating ourselves by offering individualized retention offers in real time on our website. Based on our cloud-native platform, we can easily scale this application based on user/business need and can support tens or hundreds of thousands of users.
- 92. The final result of the end-to-end integrated data platform with ML/AI capabilities are demonstrated by consuming the REST API into our existing Stock Trader application.
- 93. In this [Deploy](#) lab exercise, we will show:
- a. Integrate model into Stock Trader Application
  - b. Predict separation risk for high value customers through the REST API's endpoint
  - c. Customized promotions offered to customers based on their loyalty and separation risk
- 94. The impact of the model deployment to the business is then highlighted. The business analyst takes the data after a three-month period and analyzes the impact of embedding AI into the Stock Trader application model. They see that it has led to a significant change in our business' overall performance.

- \_\_\_95. In the Work Map Index – the **Deploy** process is described as below:



## 8.7 Business Impact

- \_\_\_96. The purpose of our last lab is to provide a dashboard quantifying the business impact of our project over time. We will turn the clock forward 90 days and measure business impact. Specifically, how many retention offers were made, how many customers took advantage of them, what was the cost of the offers (that is lost revenue for free trades and free wealth consulting sessions) AND what was the incremental revenue we brought to our company by significantly improving our 'customer separation' problem.
- \_\_\_97. The goal of this workshop has been to demonstrate analytics modernization and enhance applications with machine learning models, which drive our organization's journey to cloud and AI.



---

## Appendix C. Integration of Hadoop Cluster with ICP-D

- 98. Note: These instructions are for IT professionals who will integrate their ICP-D with a Hadoop cluster.
- 99. Data virtualization is the key strength of IBM Cloud Private for Data.
- 100. The Collect of ICP-D is to get data regardless where it resides. If data is in a Hadoop cluster of either Hortonworks or Cloudera, it can be accessed without actually moving it.
- 101. The integration of ICP-D with Hadoop cluster will be useful for the following tasks:
  - a. Accessing the HDFS file system to read and write files
  - b. Connecting to the Hive database to access table data
  - c. Using Big SQL on the Hortonworks data platform to access table data
  - d. Push down Spark processing to the Spark2 server on the HDP

### 8.8 Access to Hadoop Behind Firewall

- 102. When the Hadoop services behind a firewall, the authentication via Knox and/or Kerberos is implemented. In such situation, it is recommended to implement the ICP-D Hadoop Integration package (DSX-HI) on one of the edge nodes of the cluster.
- 103. In a secure Hadoop cluster, only authentication services are exposed to the outside world. The insecure services such as native HDFS, WebHDFS, WebHCat, Livy, and so on. are behind a firewall and this can only be accessed through a DSX-HI gateway.
- 104. The DSX-HI gateway is built on Apache Knox.

### 8.9 Download Software

- 105. The Hadoop Integration software can be downloaded from Passport Advantage. At the time of writing of this lab document, we used ICP 3.1 with ICPD 1.1.0.2. The part number of the *Hadoop integration software for IBM Cloud Private for Data 1.1.0.2* is [CNW0BEN](#).
- 106. In our lab environment, Hortonworks Sandbox (Hadoop VM) is available in the tenth VM.
- 107. Log in to the Hadoop VM and use [root](#) password: [hadoop](#).

```
# ssh 192.168.142.110
```

### 8.10 Prepare Hadoop Cluster for ICP-D Integration

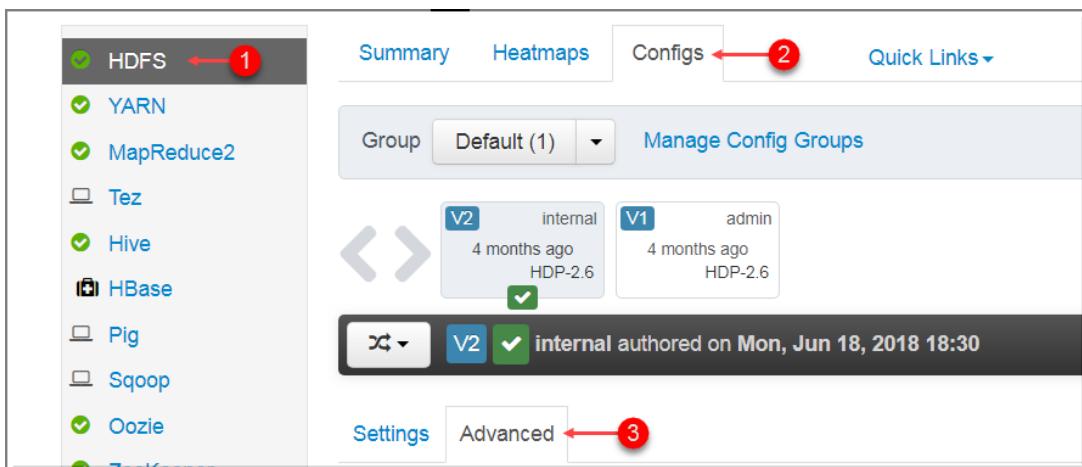
- 108. The DSX-HI services connect to the Hadoop cluster with special privileges to impersonate other users to ensure that YARN jobs are submitted on behalf of the user and that HDFS authorization settings are respected. We will use user [dsxhi](#) as a proxy user to do jobs on behalf of requests coming from ICP-D.

\_\_109. Log in to the Hortonworks Sandbox. Open URL <http://192.168.142.110:1080> using the Chrome browser from VM **icp01**.

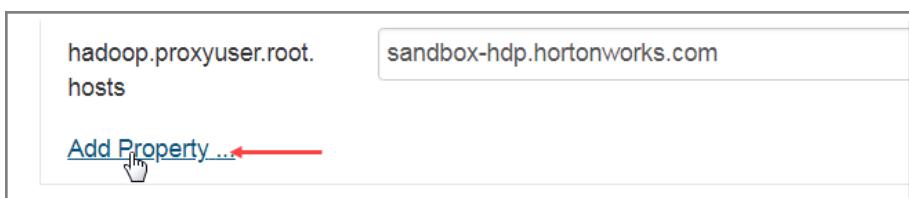
\_\_110. Click **Launch Dashboard**.

\_\_111. Login using Username: **admin** and Password: **password**.

\_\_112. Click **HDFS** ⇒ **Configs** ⇒ **Advanced**.



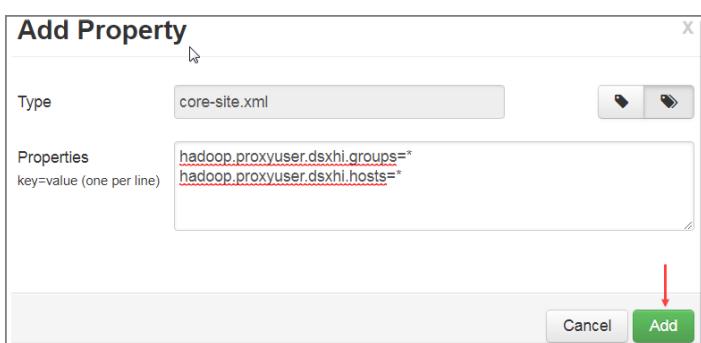
\_\_113. Scroll down to **Custom core-sites** and expand it. Go to the end of property list and click **Add Property**.



\_\_114. Type in the box, the following properties.

```
hadoop.proxyuser.root.hosts
hadoop.proxyuser.dsxhi.groups=*
hadoop.proxyuser.dsxhi.hosts=*
```

\_\_115. Click **Add**.



- \_\_116. Scroll to the top and click **Save** twice and click **Proceed Anyway** and click **OK**.
- \_\_117. Click **Hive**  $\Rightarrow$  **Configs**  $\Rightarrow$  **Advanced**.
- \_\_118. Scroll down to **Custom webhcat-site**.
- \_\_119. Expand it and go to the end of list and click **Add Property**.
- \_\_120. Add the following two lines in the Properties box and click **Add**.

```
webhcat.proxyuser.dsxhi.hosts=*
webhcat.proxyuser.dsxhi.groups=*
```

- \_\_121. Scroll to top and click **Save** twice and click **Proceed Anyway** and click **OK**.
- \_\_122. Click **Spark2**  $\Rightarrow$  **Configs** and scroll down to **Custom livy2-conf**. Expand it and click **Add Property**. Add the following line and click **Add**.

```
livy.superusers=dsxhi
```

- \_\_123. Scroll to top and click **Save** twice and click **Proceed Anyway** and click **OK**.

## 8.11 Create **dsxhi** User

- \_\_124. Switch back to the command line window in the Hadoop VM to run the following commands as root.
- \_\_125. Login to the Docker container **sandbox-hdp-security**

```
# docker exec -it sandbox-hdp-security bash
```

- \_\_126. Run the following to create dsxhi user

```
# useradd -u 3001 -g hdfs dsxhi
```

- \_\_127. Create **dsxhi** HDFS home directory and change permissions.

```
# hdfs dfs -mkdir /user/dsxhi
# hdfs dfs -chown dsxhi:hdfs /user/dsxhi
```

- \_\_128. Create a user for every user registered in ICP-D. In our case, we are using only one user **admin**.

```
# useradd -u 3002 admin
```

- \_\_129. Change permissions.

```
# hdfs dfs -chown admin:admin /user/admin
```

## 8.12 Install Software in the Docker Container

- \_\_130. You can **scp** RPM to the Docker container or put it on any web server and use **rpm -ivh http://...** to install it.

\_\_131. For example: I copied the RPM to my local web server at <https://www.zinox.com/wp-content/uploads/dsxhi-icpdata-ee-1.0.1.0-49.noarch.rpm>

\_\_132. Install `dsxhi` on Hadoop edge node.

```
# rpm -ivh dsxhi-icpdata-ee-1.0.1.0-49.noarch.rpm
Preparing...                                              [100%]
Updating / installing...
 1:dsxhi-1.0.1.0-49                                     [100%]
```

## 8.13 Create `dsxhi` Configuration File.

\_\_133. Switch to `conf` directory.

```
# cd /opt/ibm/dsxhi/conf
```

\_\_134. Copy the Hortonworks template

```
# cp dsxhi_install.conf.template.HDP dsxhi_install.conf
```

\_\_135. Modify the config file to specify service user, webchat URL, ICP-D server name and so on, in `dsxhi_install.conf`.

\_\_136. The `dsxhi_install.conf` as per our configuration is as shown.

```
# sed -e '/^$/d' -e '/^#.*/d' dsxhi_install.conf
dsxhi_license_acceptance=a
dsxhi_serviceuser=dsxhi
dsxhi_serviceuser_group=hdfs
dsxhi_gateway_port=8443
dsxhi_rest_port=8082
cluster_manager_url=http://192.168.142.110:8080
cluster_admin=admin
exposed_hadoop_services=webhdfs,webhcat
existing_webhcat_url=http://192.168.142.110:50070
dsxhi_livyspark_port=8998
dsxhi_livyspark2_port=8999
known_dsx_list=https://192.168.142.101:31843
package_installer_tool=
packages=lapack
cluster_nodes=
cluster_ssh_user=
cluster_ssh_key_path=
```

## 8.14 Install DSXHI

\_\_137. Save the existing keystore

```
# cd /etc/pki/ca-trust/extracted/java
# mv cacerts cacerts.orig
```

138. Install DSXHI

```
# cd /opt/ibm/dsxhi/bin
# ./install.py --dsxhi_gateway_master_password=password --
dsxhi_java_cacerts_password=password --password=password
```

139. The installation does prerequisite check and if they pass, it will configure [knox](#) gateway and start all services.

```
IBM Cloud Private for Data Enterprise Edition Hadoop Integration Service
-----
Terms and Conditions: http://www14.software.ibm.com/cgi-bin/weblap/lap.pl?la_formnum=&li_formnum=L-KMRY-B2632F&title=IBM+Cloud+Private+for+Data+-+Enterprise+Edition+V1.1.0.1+(bundles+ICP+Foundation)&l=en
--Determining properties
--Running the prechecks

--Configure gateway
--Create template for gateway
--Install dsxhi_rest
--Start all services
--Setting up known DSX Local clusters
--Install finished. Check status in /var/log/dsxhi/dsxhi.log
```

140. Check if Gateway, HDFSWeb and file system services ports are in use.

```
# netstat -nlp | grep 8443
tcp        0      0 0.0.0.0:8443        0.0.0.0:*      LISTEN      75285/java
# netstat -nlp | grep 50070
tcp        0      0 172.18.0.2:50070    0.0.0.0:*      LISTEN      6772/java
# netstat -nlp | grep 8020
tcp        0      0 172.18.0.2:8020    0.0.0.0:*      LISTEN      6772/java
```

141. The modified [knox](#) gateway should be used and the default knox gateway should be kept disabled.**8.15 Create [DSXHI](#) Systemd Service**142. The [dsxhi](#) gateway should start automatically whenever the Hadoop edge (gateway) node starts.143. Create a service as per the following.

```
cat << 'EOT' > /usr/lib/systemd/system/icpdhi.service
[Unit]
Description=ICPD DSXHIO Service
After=ambari-server.service ambari-agent.service mysql.service dbus.service

[Service]
Type=oneshot
ExecStart=/opt/ibm/dsxhi/bin/start.py
TimeoutStartSec=1200
```

```
[Install]
WantedBy=multi-user.target
EOT
```

- \_\_144. Enable the service so that the gateway is started automatically whenever the Docker container is restarted.

```
# systemctl enable icpdhi
```

## 8.16 Commit Docker Container to Save the Installation

- \_\_145. Exit from container

```
# exit
```

- \_\_146. Commit the image with same tag.

```
ID=$(docker ps --format '{{.ID}}' --filter Name=sandbox-hdp-security)
IMAGE=$(docker ps --format '{{.Image}}' --filter Name=sandbox-hdp-security)
docker commit $ID $IMAGE
```

- \_\_147. The image is saved and the [dsxhi](#) install will not disappear if the container is stopped or removed.

## 8.17 Troubleshooting

- \_\_148. When the registration of the Hadoop edge node fails, [Administer](#) ⇒ [Hadoop Integration in ICP-D Web UI](#).

- \_\_149. The reason for this failure is the ICP-D [utils-api](#) pod is unable to resolve the name of the Hadoop edge node.

- \_\_150. Make sure that the name of the Hadoop edge server is defined and resolvable through the name server that ICP-D is accessing.

- \_\_151. The troubleshooting can be done by looking at the log files from the [utils-api](#) pod in ICP-D and the log files from the Hadoop edge node.

```
/var/log/dsxhi/gateway
/var/log/dsxhi
```

- \_\_152. Commands to log onto Hadoop sandbox in our cluster.

```
# ssh Hadoop
# docker exec -it sandbox-hdp-security bash
# cd /opt/ibm/dsxhi/bin
```

- \_\_153. The commands [start.py](#), [status.py](#) and [stop.py](#) can be used to start, status and stop the [dsxhi](#) gateway server.

\_\_154. Run command `manage_known_dsx.py --list` to find the service URL that we need to use in ICP-D.

```
# ./manage_known_dsx.py --list
DSX Local Cluster URL          DSXHI Service URL
https://192.168.142.102:31843  https://sandbox-hdp.hortonworks.com:8443/gateway/192.168.142.102
https://192.168.142.101:31843  https://sandbox-hdp.hortonworks.com:8443/gateway/192.168.142.101
```

\_\_155. Check the logs of the `utils-api` pod

```
# kubectl -n zen get pods | grep utils
utils-api-5966956c44-t5qcn  1/1    Running   0  1h
# kubectl -n zen logs utils-api-5966956c44-t5qcn
```

Credit: Based on document prepared by Frank Ketelaars, Technical Lead for IBM Cloud Private for Data.

**\*\* End of Appendix:**

---

## Appendix AA. Notices

This information was developed for products and services offered in the USA.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing  
IBM Corporation  
North Castle Drive  
Armonk, NY 10504-1785  
USA

For license inquiries regarding double-byte (DBCS) information, contact the IBM Intellectual Property Department in your country or send inquiries, in writing, to:

IBM World Trade Asia Corporation  
Licensing  
2-31 Roppongi 3-chome, Minato-ku  
Tokyo 106-0032, Japan

**The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law: INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.** Some states do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this IBM product and use of those Web sites is at your own risk.

IBM may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

Any performance data contained herein was determined in a controlled environment. Therefore, the results obtained in other operating environments may vary significantly. Some measurements may have

been made on development-level systems and there is no guarantee that these measurements will be the same on generally available systems. Furthermore, some measurements may have been estimated through extrapolation. Actual results may vary. Users of this document should verify the applicable data for their specific environment.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

All statements regarding IBM's future direction and intent are subject to change or withdrawal without notice and represent goals and objectives only.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to the names and addresses used by an actual business enterprise is entirely coincidental. All references to fictitious companies or individuals are used for illustration purposes only.

#### COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrate programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs.

---

## Appendix BB. Trademarks and copyrights

The following terms are trademarks of International Business Machines Corporation in the United States, other countries, or both:

IBM	AIX	CICS	ClearCase	ClearQuest	Cloudscape
Cube Views	Db2	developerWorks	DRDA	IMS	IMS/ESA
Informix	Lotus	Lotus Workflow	MQSeries	OmniFind	
Rational	Redbooks	Red Brick	RequisitePro	System i	
System z	Tivoli	WebSphere	Workplace	System p	

Adobe, Acrobat, Portable Document Format (PDF), and PostScript are either registered trademarks or trademarks of Adobe Systems Incorporated in the United States, other countries, or both.

Cell Broadband Engine is a trademark of Sony Computer Entertainment, Inc. in the United States, other countries, or both and is used under license therefrom.

Java and all Java-based trademarks and logos are trademarks of Sun Microsystems, Inc. in the United States, other countries, or both. See Java Guidelines

Microsoft, Windows, Windows NT, and the Windows logo are registered trademarks of Microsoft Corporation in the United States, other countries, or both.

Intel, Intel logo, Intel Inside, Intel Inside logo, Intel Centrino, Intel Centrino logo, Celeron, Intel Xeon, Intel SpeedStep, Itanium, and Pentium are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

UNIX is a registered trademark of The Open Group in the United States and other countries.

Linux is a registered trademark of Linus Torvalds in the United States, other countries, or both.

ITIL is a registered trademark and a registered community trademark of the Office of Government Commerce and is registered in the U.S. Patent and Trademark Office.

IT Infrastructure Library is a registered trademark of the Central Computer and Telecommunications Agency which is now part of the Office of Government Commerce.

Other company, product and service names may be trademarks or service marks of others.

## NOTES

## NOTES



---

© Copyright IBM Corporation 2018.

The information contained in these materials is provided for informational purposes only and is provided AS IS without warranty of any kind, express or implied. IBM shall not be responsible for any damages arising out of the use of, or otherwise related to, these materials. Nothing contained in these materials is intended to, nor shall have the effect of, creating any warranties or representations from IBM or its suppliers or licensors, or altering the terms and conditions of the applicable license agreement governing the use of IBM software. References in these materials to IBM products, programs, or services do not imply that they will be available in all countries in which IBM operates. This information is based on current IBM product plans and strategy, which are subject to change by IBM without notice. Product release dates and/or capabilities referenced in these materials may change at any time at IBM's sole discretion based on market opportunities or other factors and are not intended to be a commitment to future product or feature availability in any way.

IBM, the IBM logo and ibm.com are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both. If these and other IBM trademarked terms are marked on their first occurrence in this information with a trademark symbol (® or ™), these symbols indicate U.S. registered or common law trademarks owned by IBM at the time this information was published. Such trademarks may also be registered or common law trademarks in other countries. A current list of IBM trademarks is available on the web at "Copyright and trademark information" at [ibm.com/legal/copytrade.shtml](http://ibm.com/legal/copytrade.shtml)

Other company, product and service names may be trademarks or service marks of others.



Please Recycle

---