

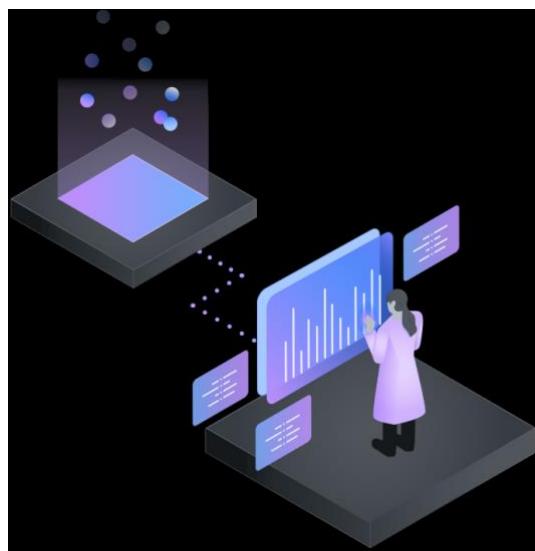
IBM Journey to Cloud and AI

Analytics Modernization Workshop

Featuring: IBM Cloud Pak for Data



Core Lab Workbook



Lab workbook and IBM Cloud Pak for Data workshop design and environment by:
Burt Vialpando, IBM Executive Analytics Architect
Kent Rubin, IBM Solution Architect

October 26, 2020

Acknowledgements

- **Duane Almeter** and **Eric Watson** for their leadership on the project
- **Daniel Kikuchi** for the CPD 3.0 cluster install environments, and other technical input and support
- **Ed Duhe, Rich Russo** and **Mitchell Odum** for providing the ESX server development platform
- **John Lucas** for product testing, workbook publishing, and project management support
- **John Van Buren** for the Organize lab designs
- **Rajesh Kartha** for Data Virtualization caching lab and DV z/OS lab
- **Benjamin Herta** for the AutoAI non-AVX work around
- **Sidney Phoon** for the last Notebook updates
- **Eric Martens** for the OpenScale lab design
- **Owais Hashmi** for OpenScale storyboarding and lab assistance
- **Rohit Gargate** for OpenScale Auto setup assistance
- **Ben Chard** for providing Analyze lab resources and assistance
- **Tom Konchan** for SPSS/IPS and Decision Optimization lab
- **Daniel Hancock** for the NPS Getting Started lab
- **David Trotter** for the z/OS work on the DV z/OS lab
- **Linda Snow** for providing assistance with complete workshop review



Table of Contents

LAB 01	GETTING STARTED	8
1.1	IBM Journey to Cloud and AI: Analytics Modernization Workshop	8
1.2	Cloud Pak for Data: defined	8
1.3	OpenShift: defined	8
1.4	Lab workshop environment	8
1.5	Audience for this IBM workshop	9
1.6	Getting started	10
1.7	User management: Persona-based roles and teams	13
1.8	Managing platform options	15
1.9	Reviewing the profile settings	17
1.10	Reviewing CPD services	18
1.11	Reviewing instances	20
1.12	Customizing branding	21
1.13	Lab conclusion	23
LAB 02	BUSINESS USE CASE: CUSTOMER CHURN	24
2.1	Lab overview	24
2.2	Persona represented in this lab	24
2.3	Logging into the CPD web client (if you have not already done so)	25
2.4	Reviewing the dashboard: Monthly Metrics	25
2.5	Reviewing the dashboard: Demographics Discovery	27
2.6	Devising a plan	31
2.7	Reviewing the dashboard: Monthly Metrics After AI	34
2.8	Lab conclusion	35
LAB 03	COLLECT: CONNECTIONS	36
3.1	Lab overview	36
3.2	Persona represented in this lab	36
3.3	Db2 data overview – Transforming for analytics	37
3.4	MongoDB data overview – Virtualizing for analytics	42
3.5	Lab conclusion	47

LAB 04	ORGANIZE	48
4.1	Lab overview.....	48
4.2	Persona represented in this lab	48
4.3	Logging into the CPD web client (if you have not already done so)	49
4.4	Reviewing a data asset in the project.....	49
4.5	Reviewing a business glossary	59
4.6	Reviewing a Governance Policy and a Rule.....	62
4.7	Reviewing Classifications, Data Classes, and Reference Data	64
4.8	Searching for Data.....	67
4.9	Reviewing the catalog	71
4.10	Transforming Data	73
4.11	Lab conclusion.....	83
LAB 05	COLLECT: VIRTUALIZE	84
5.1	Lab overview.....	84
5.2	Persona represented in this lab	84
5.3	Logging into the CPD web client (if you have not already done so)	85
5.4	Remove older Data Virtualization sources.....	85
5.5	Adding Data Virtualization data sources.....	86
5.6	Virtualizing the Db2 data	88
5.7	Virtualizing the MongoDB data	91
5.8	Joining the virtualized tables	93
5.9	Lab conclusion.....	98

LAB 06 ANALYZE: AutoAI.....	99
6.1 Lab overview.....	99
6.2 Persona represented in this lab	99
6.3 Logging into the CPD web client (if you have not already done so).....	100
6.4 Setting up the AutoAI experiment.....	100
6.5 Running the AutoAI experiment	106
6.6 Running a Notebook.....	109
6.7 Reviewing the AutoAI results.....	115
6.8 Saving the model.....	119
6.9 Lab conclusion.....	123
LAB 07 DEPLOY	124
7.1 Lab overview.....	124
7.2 Persona represented in this lab	124
7.3 Logging into the CPD web client (if you have not already done so)	124
7.4 Reviewing the notebook deployment space.....	125
7.5 Deploying and testing the AutoAI model	138
7.6 Lab conclusion.....	144
Lab 08 Infuse: Watson OpenScale	145
8.1 Lab overview.....	145
8.2 Persona represented in this lab	145
8.3 Logging into the CPD web client (if you have not already done so).....	145
8.4 Credit Risk built-in demo	146
8.5 Monitoring the AutoAI model in OpenScale	167
8.6 Lab conclusion.....	178

LAB 09 INFUSE: COGNOS ANALYTICS - INTRODUCTION	179
9.1 Lab overview.....	179
9.2 Persona represented in this lab	179
9.3 Logging into the CPD web client (if you have not already done so).....	179
9.4 Logging into the IBM Cognos provisioned instance	180
9.5 The IBM Cognos Analytics User Interface	181
9.6 Importing the CSV file as an Exploration.....	183
9.7 Cleaning up the data	185
9.8 Exploring Data Relationships	191
9.9 Creating Exploration Cards	194
9.10 Embedding Cognos Content.....	207
9.12 Lab conclusion.....	208
LAB 10 WRAP-UP	209
10.1 Lab overview.....	209
10.2 Data Scientist wrap-up	209
10.3 Data Engineer wrap-up.....	210
10.4 Administrator wrap-up.....	212
10.5 Workshop conclusion	214
<i>Back Page: Notices</i>	<i>215</i>
<i>Back Page: Trademarks and Copyrights.....</i>	<i>217</i>

[This page left intentionally blank]

Lab 01 GETTING STARTED

1.1 IBM Journey to Cloud and AI: Analytics Modernization Workshop

This workshop provides hands-on experience with Cloud Pak for Data that will show you how to modernize your microservices applications by enriching them with Machine Learning (ML) and Artificial Intelligence (AI).

The Journey to AI requires a strong information architecture (IA) that supports self-service capabilities and balances the needs of both the agility required by lines of business as well as the “Enterprise class” delivery required by IT. This journey can move significantly faster and with more efficiency when you use a single integrated platform like Cloud Pak for Data. It is the world’s leading platform that allows you to **Collect**, **Organize** and **Analyze** data, and then **Deploy** the results to **Infuse** your applications with AI.

1.2 Cloud Pak for Data: defined

Cloud Pak for Data (CPD) is an integrated end-to-end data and analytics platform designed to help make data more accessible and trusted, as well as to provide access to many analytical tools to help your organization gain insights from your data.

CPD provides the data platform that accelerates the journey up the “AI Ladder.” With it, you can quickly build, train, deploy, and manage machine learning (ML) models to create applications with Artificial intelligence (AI). CPD provides inventory and cataloging of your data sources, self-service shopping for data, and data integration and refinement capabilities. Thus, high quality and trusted data can be more easily prepared, assembled and used in one modern, integrated, collaborative and scalable platform.

Cloud Pak for Data is installed on the foundation of OpenShift for the cluster this workshop uses.

1.3 OpenShift: defined

Red Hat OpenShift is an open, hybrid cloud Kubernetes platform used to build, run, and scale container-based applications. OpenShift itself is built upon a foundation of Red Hat Enterprise Linux. OpenShift includes everything you need to manage your development lifecycle, including standardized workflows, support for multiple environments, continuous integration, and release management.

Cloud Pak for Data can be installed and managed on public cloud platforms as well, including IBM Cloud, AWS, Microsoft and more.

1.4 Lab workshop environment

We are using a CPD cloud cluster for this workshop. This software environment was built with the following key software components:

- [Red Hat OpenShift Container Platform \(RHOCUP\) 4.3.21](#) as the foundational cloud-native technology platform of Kubernetes and Docker, as well as other open-source tools.
- [Cloud Pak for Data v3.0.1](#) as the microservice-built, integrated data and analytics platform with various add-ons installed and enabled.

1.5 Audience for this IBM workshop

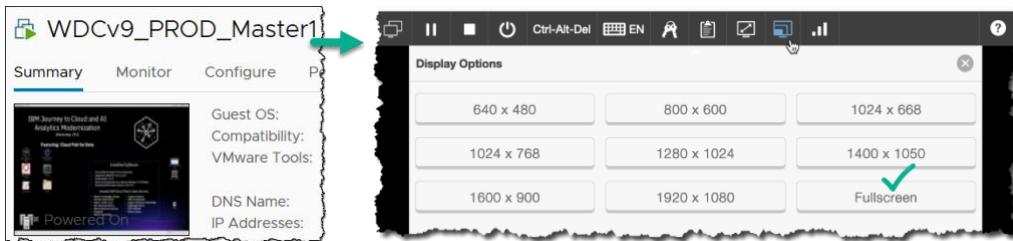
This IBM workshop is aimed at the line-of-business professionals who are tasked to gain new insights from all available data – regardless of its type and origin. The following personas who will be represented in the various labs will greatly benefit from this workshop:

Persona (Role)	Capabilities
 Administrator	<p>Administrators set up and maintain the CPD environment itself.</p> <p>Note: while some of the Admin work can be done in the CPD web client, most of the Admin work on the cluster would be done in OpenShift which is outside the scope of this workshop.</p> <p>The exercises in this first lab represent some typical CPD Administrator activities.</p>
 Data Engineer	<p>Data Engineers build and optimize the systems to allow data scientists and business analysts to perform their work. The Data Engineer ensures that any data is properly received, transformed, stored, and made accessible to other users.</p>
 Data Steward	<p>Data Stewards bring integration and transformation of the data as well as providing governance, lineage and classification of the data.</p>
 Data Quality Analyst	<p>Data Quality Analysts perform advanced curation of the data and analyze the quality of the data.</p>
 Business Analyst	<p>Business Analysts deliver value by analyzing data to answer questions and communicating the results to help make better business decisions.</p>
 Data Scientist	<p>Data Scientists bring expertise in statistics and the process of building ML/AI models to make predictions and answer key business questions.</p>
 Developer	<p>Developers create and maintain the end-user applications that utilize the output from all the other personas on the CPD platform.</p>

1.6 Getting started

1. To launch your CPD cluster, you will select the virtual machine with the label **Master1**.

Click on this node (VM) to expand the image and then click on the **display** option in the top bar and select **Fullscreen**.



You are now in the first (headed) virtual machine image of a cluster of virtual machines that comprise your Cloud Pak for Data cluster. Everything you will be doing in the labs will be driven from this first image.

If a screensaver function has locked the screen, hit **[Enter]** to get to the desktop.

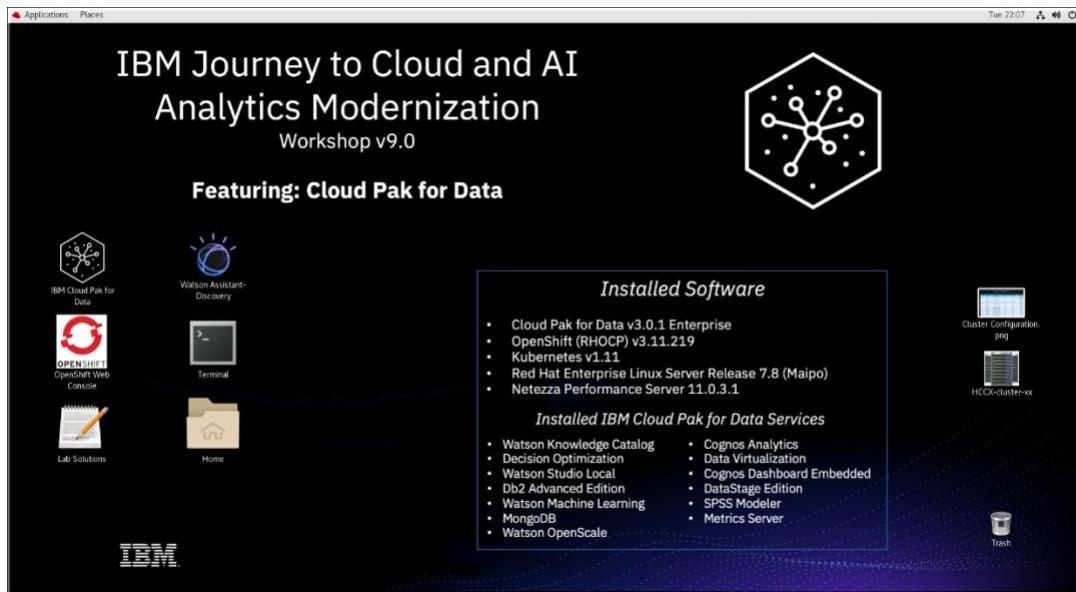




Admin

If necessary, you can log back into the Linux OS with: User **ibmdemo**, Password **passw0rd**.

2. In this VM, notice the lab desktop looks like this:



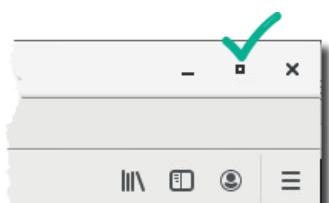
- __3. Make sure to run in “Fullscreen” mode to make the most of your computer screen’s real estate.



- __4. Double-click the icon IBM [Cloud Pak for Data](#).



- __5. After launching, maximize the browser window.



- __6. The CPD web client GUI displays as shown below.

Use [cpduser](#) and [cpdaccess](#) for the *Username* and *Password* and click [Sign in](#).

SIGN IN

IBM Cloud Pak for Data

Username
cpduser

Password
.....cpdaccess

Sign in →

- __7. You should now be at the [Home Page](#).

- __8. Scroll down to review the quick navigation and resources links on this page.

You will be exploring many of these in more detail in this workshop, so don't follow these links right now.

- __9. If you do happen leave this page by clicking on a link and you want to return to it, you can do so by clicking:

[Navigation Menu](#) ("hamburger" icon) ⇔ [Home](#)

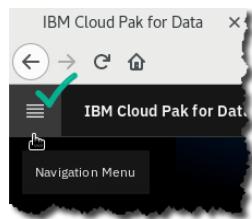
- __10. You can, of course, also use the browser back and forward arrow keys to navigate through main screens in the CPD web client.

 Admin	In this workshop, we will demonstrate the Collect , Organize and Analyze capabilities to create a machine learning model that you can Deploy and then Infuse into a microservices application.
-------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

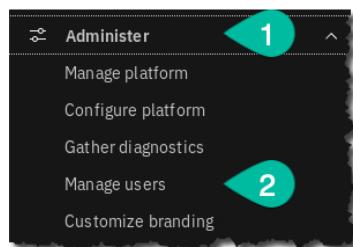
1.7 User management: Persona-based roles and teams

This section explores user authorizations for the various stages of the data analytics pipeline.

- __11. Click on the [Navigation Menu](#) (“hamburger” icon) at the top left of the screen.



- __12. Click on the [Administer](#) Group (to display the drop-down menu) then click [Manage users](#).



- __13. Review the [Users](#) first, then click on the [Roles](#) section to review the various personas that can be represented by any given user. A user can be granted more than one role if needed.

[Hover](#) over the permissions of any role to see what individual permissions exist for that role.

Role	Description	Modified on	Enabled permissions
Administrator	Administrator role	Jun 26, 2020 10:09 AM	Access advanced governance capabilities + 16 more
Business Analyst	Business Analyst role	Jun 24, 2020 12:57 PM	Access catalogs + 2 more
Data Engineer	Data engineer role	Jun 24, 2020 12:57 PM	Access catalogs +
Data Quality Analyst	Data quality analyst role	Jun 24, 2020 12:57 PM	Access catalogs + 5 more



These existing roles can be edited (customized) or new ones created, to suit your organization's needs. CPD is very much “persona driven” in that each user can play their particular part in your organization’s journey to AI. Each user (acting as one or more personas) can hand off and/or share their work with other users/personas, for a totally collaborative environment.

- 14. Click back to the [Users](#) section and then click on [Configure LDAP](#). You can review the fields required to do this here.



The screenshot shows the 'Manage users' interface in the IBM Cloud Pak for Data web client. The 'Users' tab is selected. In the top right corner, there is a green checkmark icon followed by the text 'Configure LDAP'. Below this, the 'Configure LDAP' dialog box is open, titled 'Configure LDAP'. It contains fields for 'LDAP server information': 'LDAP protocol' (set to 'ldaps://'), 'LDAP hostname' ('sample.com'), 'LDAP port' ('636'), 'Domain search user (optional)' ('A user that can perform lookups in the LDAP server'), 'Domain search base' ('dc=sample,dc=com'), 'Domain search password (optional)' ('The password for domain search user'), and 'User search field' ('For example, cn, uid, or sAMAccountName'). A 'Cancel' button is visible in the top right of the dialog box.



Admin

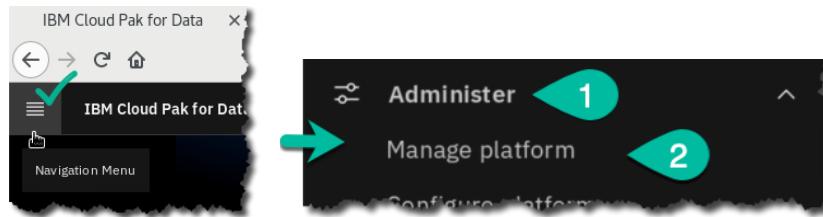
For the sake of simplicity, you will remain logged into the CPD web client throughout this workshop as the user [cpduser](#), which has been granted all persona roles. This was done so that you will not be required to log off and log on again as different users to represent the varying personas as you make your way through each lab.

In your organization, however, it is likely that once you have a mature CPD environment set up, separation of duties will be defined by persona where different users will be assigned one or more personas to do their particular tasks.

1.8 Managing platform options

You can view the underlying OpenShift services and pods by doing the following:

- _15. Click [Navigation Menu](#) \Rightarrow [Administer](#) \Rightarrow [Manage platform](#).



- _16. In the Search area, type **db2**.

- _17. Click the deployment link [Db2 Advanced Edition](#).

A screenshot of the "Manage platform" interface. At the top, it says "Manage platform: cp4d (primary)" and "Last updated: 6/23/2020 6:21 PM". Below that is a "Deployments" section with a sub-instruction: "See the deployments running in each Red Hat OpenShift project on the platform. You can see how many resources each deployment is currently using compared to the number of resources that the deployment has reserved." There are filter options "Filter by: All types" and "Clear all". A search bar contains "db2" with a checkmark. Below is a table with columns "Name", "Type", and "Installed". It lists three entries: "Total" (Base, 6/19/2020), "Db2 Advanced Edition" (Base, 6/19/2020 with a checkmark next to it), and "db2wh catalog" (Base, 6/18/2020).

- _18. Under the tab [Fixed resources](#), notice the Deployment CPU and Memory usage.

Click on any Deployment to review the pods for it.

A screenshot of the "Db2 Advanced Edition deployment" fixed resources page. At the top, it says "Db2 Advanced Edition deployment" and "Last updated: 6/23/2020 7:00 PM". Below is a navigation bar with tabs "Fixed resources" (which is selected and highlighted in blue), "Dynamic runtimes", and "Service instances". A note states: "The Db2 Advanced Edition deployment reserves a certain amount of resources to run required services. These amounts represent the expected baseline usage for the deployment." A search bar "Find fixed resources" is present. A table lists resources with columns "Name" and "Pods". It shows three rows: "Total" with 8 pods, "db2oltp-1592545856832-ibm-unified-console-api" with 5 pods (which has a checkmark next to it), and "db2oltp-1592545856832-ibm-unified-console-influxdb" with 1 pod.

- __19. The next screen shows the underlying OpenShift/Kubernetes pods for this deployment. (Your pod names will differ.)

Pods	SCC	Service account	Created on	Status ⓘ	vCPU
Total					0.19 of 1.00
db2oltp-1592545856832-ibm-unified-console-api-599c57d76...	db2u-scc	db2u	6/19/2020 1:51 AM	●	0.19 of 1.00
db2oltp-1592545856832-ibm-unified-console-collector-1592...	db2u-scc	db2u	6/23/2020 7:05 PM	●	0.00 of 0.00
db2oltp-1592545856832-ibm-unified-console-collector-1592...	db2u-scc	db2u	6/23/2020 7:03 PM	●	0.00 of 0.00

- __20. Click on **Fixed resources** in the bread crumb line (or use the back arrow).



- __21. Click on tab **Service Instances**.

This shows how many instances of the Db2 Advanced Edition have been created using the Db2 Advanced Edition deployment. In our case, it is only one.

We will review this instance in a different way later in this lab.

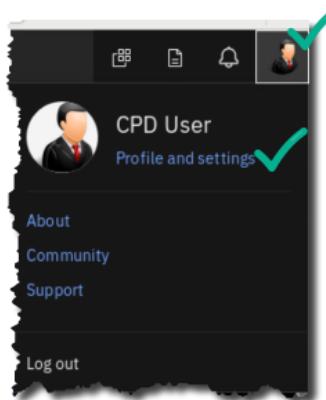
Name	Type	Created by	Created on	Users	Pods	Status	vCPU	Memory (GB)
Total				1	6	0.00 of 5.50	0.02 of 28.30	
Db2 Advanced Edition	db2oltp	user1001	6/19/2020 1:51 AM	1	6	●	0.00 of 5.50	0.02 of 28.30

1.9 Reviewing the profile settings

- _22. Click the top right circle of your web client screen that has your user icon on it. This provides a drop down.

Choose [Profile and settings](#).

(Note: this is also the location where you can [Log out](#) of the web client)



- _23. Review the things you can change in your [Profile](#), then review [Permissions](#).

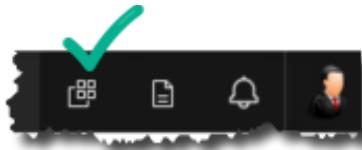
Note: in Permissions, your user has all permissions to allow you to do anything required in the workshop lab exercises. In the real world, your permissions would be more limited and controlled by an administrator.

- _24. Now review [Git Integrations](#) ↴ [Generate API key](#) and [New token](#).

These allow you to configure Git with CPD, which allows you to integrate CPD projects with your current CI/CD (Continuous Integration and Continuous Delivery) pipeline to automate delivery of the artifacts you create in the CPD platform. You can use the capabilities from the underlying OpenShift platform to build cloud native microservice applications which are tied to the ML / AI model development with a delivery pipeline.

1.10 Reviewing CPD services

- _25. Click the **Services** icon (four little squares over one bigger square) on the top right corner of your screen.



- _26. This will bring up all available services for Cloud Pak for Data.

Click through the various categories to see what services you can install on CPD.

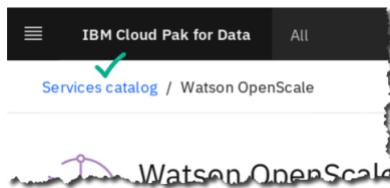
For example, click on category **AI** and notice what is available there. Those with “Enabled” are accessible by this CPD cluster right now.

The screenshot shows the Services catalog interface. On the left, a sidebar lists categories such as AI, Analytics, Dashboards, Data governance, Data sources, Developer tools, Industry solutions, Storage, Pricing, Source, Status, and Industry accelerators. The AI category is selected, indicated by a green checkmark. The main area displays various AI services with their icons, names, descriptions, and status (IBM or Premium). Two services are highlighted with red boxes: Watson Machine Learning (Enabled) and Watson OpenScale (Enabled). Other visible services include Watson ATOps, Watson Assistant, Watson Discovery, Watson Knowledge Studio, Watson Language Translator, Watson Speech to Text, Watson Studio, and Watson Text to Speech.

- _27. Now, click on the **Watson Machine Learning** service tile to get more details on what this service can do.

The screenshot shows the detailed view of the Watson OpenScale service. At the top, it says "Watson OpenScale" and "Enabled". Below that is the "Insights Dashboard" which displays metrics for Model Monitors, Deployments Monitored, Quality Alerts, Fairness Alerts, and Drift Alerts. The dashboard includes three cards: "GeneralCreditRiskModelICP", "GeneralCreditRiskModelChurn", and "GeneralCreditRiskModelProd". To the right, there is a "Version" section (3.0.1), a "Description" section (explaining that Watson OpenScale is an enterprise-grade environment for AI infused applications), and a "About the developer" section (IBM).

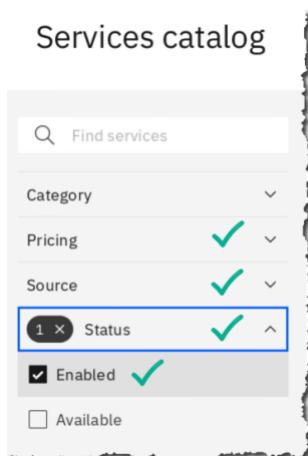
- _28. Click the browser back arrow to return to the [Services catalog](#) screen.
 (Or you can click on the [Services catalog](#) link in the breadcrumb trail itself)



- _29. Another convenient way to review these Services is to filter them in the Services Catalog options.

Click through [Pricing](#), [Source](#) and [Status](#).

In [Status](#), check [Enabled](#)



- _30. This will show all the enabled services in your cluster (as long as you did not check a filter for Categories, Pricing or Source.)

Category	Service	Status	Description
AI	Watson Machine Learning	IBM Enabled Premium	Build and train machine learning models with tools for all skill levels. Deploy and manage models at scale.
	Watson OpenScale	IBM Enabled Premium	Infuse your AI with trust and transparency. Understand how your AI models make decisions to detect and mitigate bias.
	Watson Studio	IBM Enabled Premium	Unleash the power of your data. Build custom models and infuse your business with AI and machine learning.
	Industry accelerators	Jump-start your analysis of common business problems with sample data science assets.	
Analytics	Data Refinery	IBM Enabled	Simplify the process of preparing large amounts of raw data for analysis.
	Decision Optimization	IBM Enabled Premium	Evaluate millions of possibilities to find the best solution to any given problem.
	SPSS Modeler	IBM Enabled Premium	Create flows to prepare and blend data, build and manage models, and visualize the results.

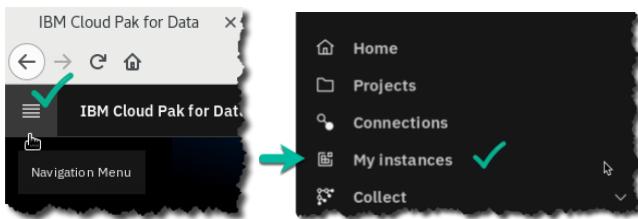
- __31. Explore through this Services page to find a few capabilities that you might find useful for your organization.



IBM continues to provide more services with each release of Cloud Pak for Data. Some are part of the base offering, others are purchasable as IBM “cartridges,” and still other are purchasable through a 3rd party vendor.

1.11 Reviewing instances

- __32. Click [Navigation Menu](#) \Rightarrow [My instances](#).



- __33. Click tab [Provisioned Instances](#) and then click the twistie to sort the instances that were provisioned for this workshop.

The screenshot shows the 'My instances' page in the IBM Cloud Pak for Data interface. The 'Provisioned instances' tab is selected. The table lists four provisioned instances:

Name	Type	Created by	vCPU (Cores)	Memory (GB)	Users	Status	Created on
MongoDB-1	mongodb	user1001	1	4	1	green	Jun 23, 2020
cognos-analytics-app	cognos-analytics-app	cpduser	9	22.59	1	green	Jun 19, 2020
data-virtualization	dv	user1001	8	32	2	green	Jun 19, 2020
Db2 Advanced Edition	db2oltp	user1001	5	27	1	green	Jun 19, 2020



The term “instance” in this context means a copy of a persistent data store within the CPD platform. These instances are stateful Kubernetes services like Db2, MongoDB, Streams, Data Virtualization and even Cognos Analytics.

- __34. Check the tabs for [Environments](#) and [Jobs](#). Note: the cluster may not have any of these running at this time, so these pages could be empty.

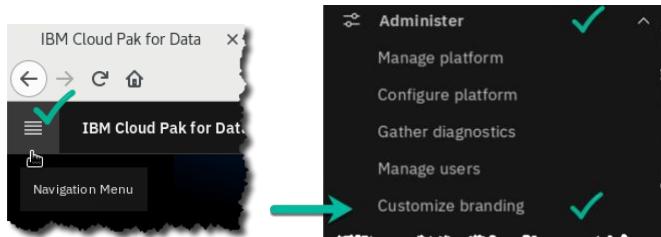


The term “environment” in this context means a copy of an analytics runtime that is running (taking up resources) on the cluster. These can be Jupyter/Python, Zeppelin/Anaconda, R Studio, Decision Optimization, etc.

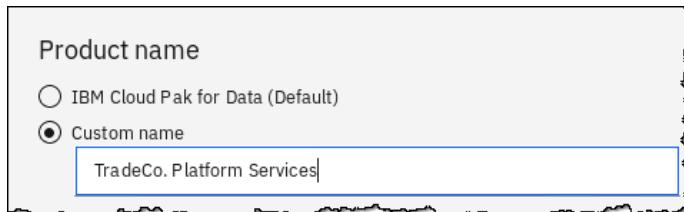
The term “job” in this context means a task scheduled within the platform. These can be analytics related (like a scheduled batch scoring job) or they can be a scheduled ETL or Streams job, etc.

1.12 Customizing branding

- __35. Click **Navigation Menu** \Rightarrow **Administer** \Rightarrow **Customize branding**.

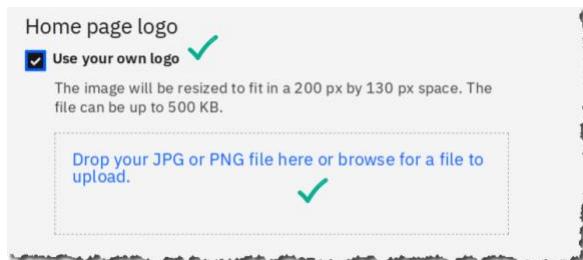


- __36. In section **Product name**, click **Custom name**, then fill in **TradeCo. Platform Services**.



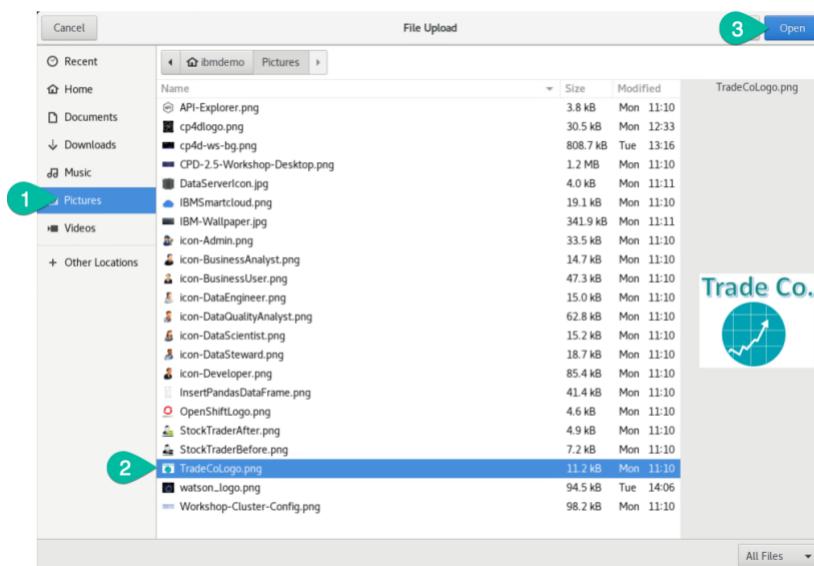
- __37. In section **Home page logo**, check **Use your own logo**.

Click on the box: **Drop your JPG or PNG file here or browse for a file to upload**.

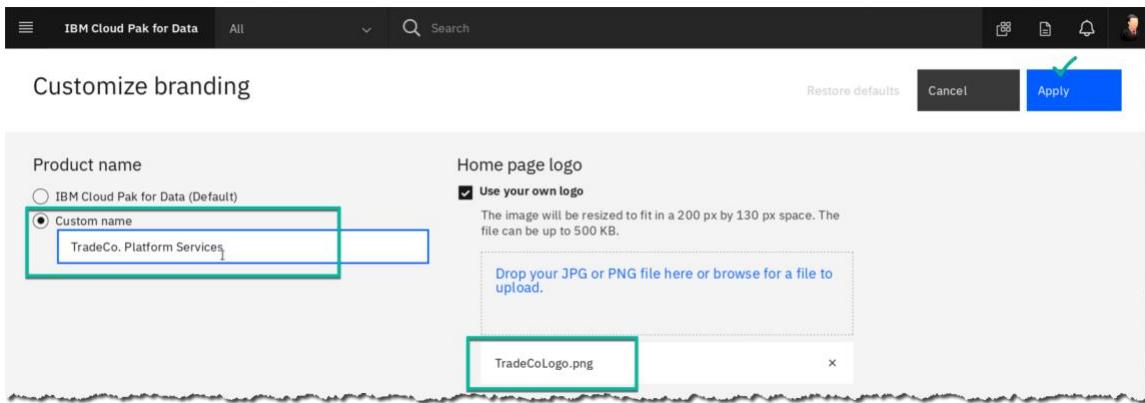


- __38. For Browser Only: Download TradeCoLogo.png from <https://ibm.biz/BdqhHa>.

For Unified Desktop: Under directory **Pictures**, select file **TradeCoLogo.png** \Rightarrow **Open**.



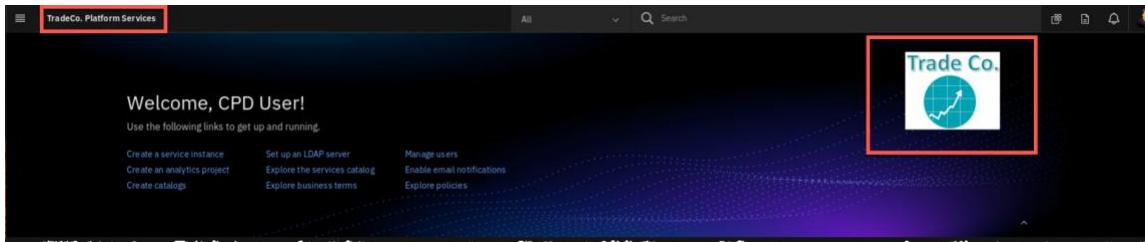
__39. Click button **Apply**.



__40. Click **Navigation Menu** \Rightarrow **Home**.

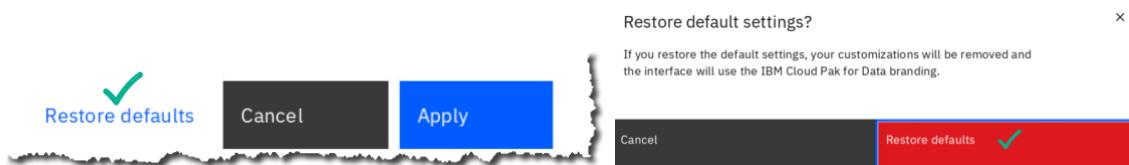


__41. Notice the Navigation bar is customized, and the Home Page can have a company logo on it.
Note: You may have to refresh your browser if it does not immediately show.



__42. Let's reset this to the default setting so the rest of the lab workbooks will be consistent with your environment:

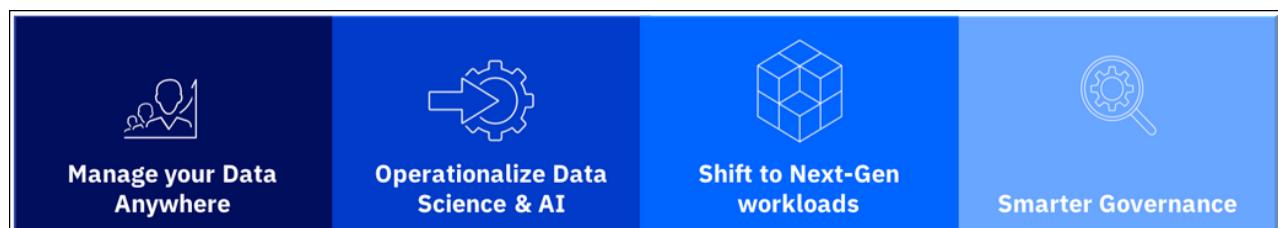
Click **Navigation Menu** \Rightarrow **Administer** \Rightarrow **Customize branding**, then **Restore defaults** \Rightarrow **Restore defaults**.



1.13 Lab conclusion

Cloud Pak for Data is useful in the following business macro use-case scenarios:

1. **Manage Your Data Anywhere:** Use data virtualization, streaming, cataloging, governance and more to prepare your data for analysis.
2. **Operationalize Data Science & AI:** Build, deploy, manage & govern models and data at scale to improve business outcomes like controlling customer churn, cross selling and up selling, predictive maintenance and more.
3. **Shift to Next-Gen workloads:** Shift to Cloud Native to be able to provision and scale in minutes, build once and deploy anywhere with multi-cloud support, and use built-in automation and collaboration to increase productivity.
4. **Smarter Governance:** Enable self-service analytics with auto-discovery of meta data, implementing governance rules and policies, enforcement of privacy to mitigate risk and to ensure compliance for regulatory requirements like GDPR.



Make Cloud Pak for Data your platform for data and analytics. Why? Because IBM understands data and provides an integrated, end-to-end data platform that enables enterprises to:

- Collect relevant data and make it simple and accessible
- Use federation, virtualization and/or transformation to combine and refine data sets
- Organize data so it can be trusted
- Analyze insights on demand
- Infuse machine learning into your applications

All of the above and much more will be demonstrated in the following workshop labs.

** End of Lab 01 – Getting Started

Lab by Burt Vialpando and Kent Rubin, IBM

Lab 02 BUSINESS USE CASE: CUSTOMER CHURN

2.1 Lab overview

Trade Co. is experiencing a decline in review, which appears to be due to losing customers. The current process of predicting customer churn seems to be ineffective as well. The company's executives have asked their senior Business Analyst to help them understand why and to help find a solution.



Trade Co. Challenges

-  Customer retention problem leading to declining revenue
-  Underperforming rules-based system to identify separation (churn) risk
-  Lack of centralized, vetted, and reliable data to ensure accuracy of analytics
-  Disparate analytical tools for reporting and model development
-  No simple way to infuse machine learning models into the customer facing Stock Trader Application

2.2 Persona represented in this lab

When embarking on machine learning projects, many organizations engage **Business Analysts** to help gain insight into their data. This persona can use tools like Analytics Dashboards to build visualizations to help the organization more clearly understand their business challenges.

The **Business Analyst** persona is the likely role to perform the exercises in this lab.

Persona (Role)	Capabilities
 Business Analyst	Business Analysts deliver value by taking data, using it to answer questions, and communicating the results to help make better business decisions.

2.3 Logging into the CPD web client (if you have not already done so)

- __1. If you are starting this lab stand-alone (without going through previous labs) do the following:
- __2. Double-click the desktop icon: [Cloud Pak for Data Web Client](#).



- __3. The CPD web client GUI displays as shown. Use `cpduser` and `cpdaccess` for the *Username* and *Password* and click [Sign in](#).

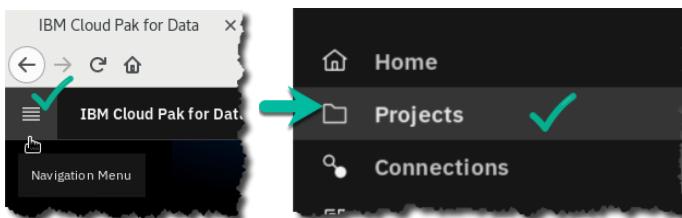
The screenshot shows the "SIGN IN" page for "IBM Cloud Pak for Data". It features a decorative header image of people interacting with data visualization and server icons. Below the image are two input fields: "Username" containing "cpduser" and "Password" containing "cpdaccess". A blue "Sign in" button is at the bottom, with a green checkmark icon to its left.

2.4 Reviewing the dashboard: Monthly Metrics

The Business Analyst has used the Cloud Pak for Data Analytics Dashboarding service to analyze Trade Co. issues. He has come up with a Monthly Metrics dashboard that demonstrates the company's concerns.

This first dashboard was built with the recent trading information delivered to him for the year. It was placed into a CSV file in the project the Business Analyst shares with his team.

- __4. In the CPD web client, click the [Navigation Menu](#) ("hamburger" icon) \Rightarrow [Projects](#).



- __5. Select the project: [CPD Workshop Analytics Project](#).

The screenshot shows a table listing projects. The columns are "Name" and "Project type". One row is selected, indicated by a green checkmark icon next to the project name. The project details are partially visible below the table.

Name	Project type
CPD Workshop Analytics Project	Analytics
CPD Workshop Data Transformation Project	Data transformation

- __6. Under tab **Assets**, scroll down until you find **Dashboards**.

Click **Monthly Metrics - Trade Co.**

My projects / CPD Workshop Analytics Project

Overview **Assets ✓** Environments

Dashboards ✓

Name	Shared	Last editor
Demographics Discovery - Trade Co.	CPD User	
Monthly Metrics - Trade Co. ✓	CPD User	
Monthly Metrics - After AI Trade Co.	CPD User	



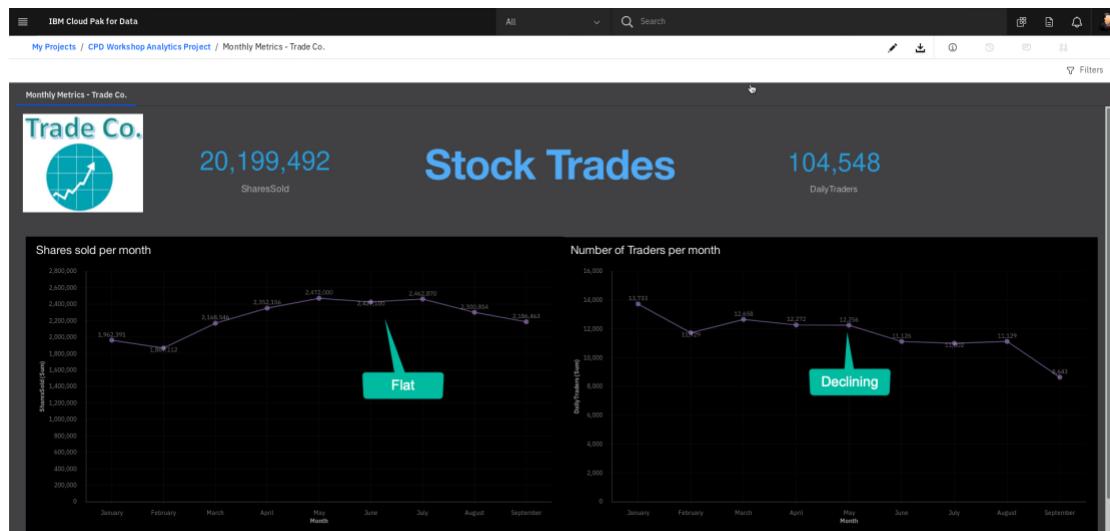
Business Analyst

Use [Ctrl] – and [Ctrl] + keys to adjust the zoom on the dashboard to best suit your screen.
You can also use [Ctrl] [Mouse-Scroll-Button].

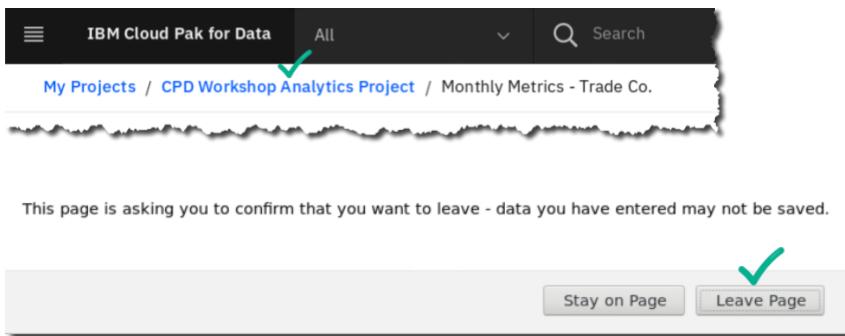
- __7. Notice that from January to September, **Shares sold per month** remains flat while **Number of Traders per month** are declining.

This verifies that customer churn is in fact occurring, even though the remaining traders are trading more each month, keeping **shares sold** nearly even.

If Trade Co. could somehow find a way to retain their customers, shares sold per month would go up, thus driving revenue up.



- __8. Leave this dashboard by clicking on the breadcrumb trail to navigate to the project again.



 Business Analyst	<p>If you would like to see how this dashboard was built, you can build it yourself by doing the following lab: (Deeper Dive) – Cognos Dashboard Embedded</p>
-------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

2.5 Reviewing the dashboard: Demographics Discovery

The Business Analyst next asked for the company's customer demographics data, joined together with the customer churn data, and joined again with the recent activity data.

This data (representing three different data sets together) was placed in the team project and the Business Analyst created the following dashboard to better understand the situation.

- __9. Under tab [Assets](#) find [Dashboards](#).

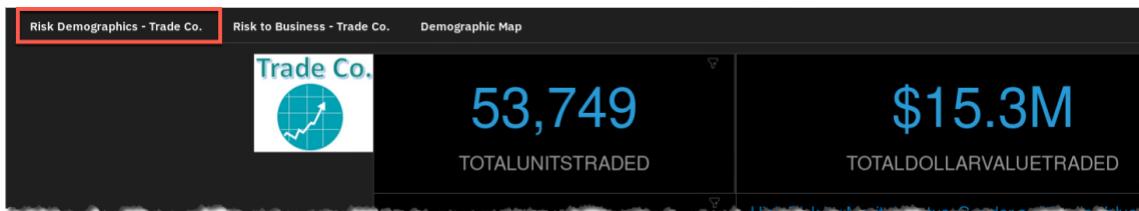
Click [Demographics Discovery – Trade Co.](#)

Name	Shared	Last editor
Demographics Discovery - Trade Co.	✓	CPD User
Monthly Metrics - Trade Co.		CPD User
Monthly Metrics - After AI Trade Co.		CPD User

- __10. Notice there are two tabs in this dashboard; you are currently positioned at tab:

[Risk Demographics - Trade Co.](#)

The top visualizations give [Total Units Traded](#) and [Dollar Value Traded](#). You will see later that these are interactive with the rest of the dashboard.



- __11. Review the top left filter (funnel) icon and notice that a filter has been set to only see information for those considered [High](#) in [Churn Risk](#).



- __12. Click on the aqua blue portion of the pie chart (which represents [Female](#) customers).

Notice the top visualizations change in value. For example, the new number displayed in the [Total Units Traded](#) visualization would now indicate [High Risk Females](#).

Click again on the aqua portion of the pie chart to deselect the values for [Female](#).



- __13. The next visualization is also filtered by High Risk, and then sub-divides the data into two charts between **Females** (on the left) and **Males** (on the right).

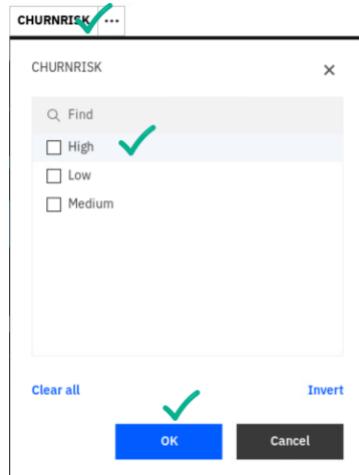
Further, each chart shows bars (**Dollars**) and lines (**Units**) for the Marital Status of Divorced(**D**), Married(**M**) and Single(**S**) customers.

This is a complex visualization, but from it you can see that Married Females are the largest group.



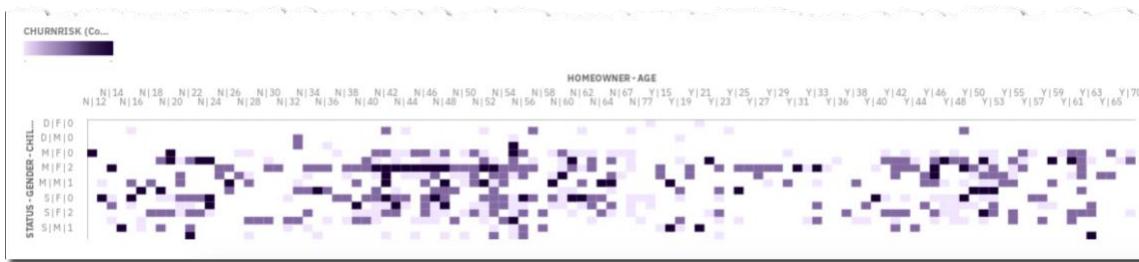
- __14. In the top filters section, click on the filter for **CHURNRISK**.

Deselect the filter for **High** and click **[OK]**.



- __15. With no filters in place for this tab on the dashboard, scroll down to the bottom visualization.

- 16. The bottom visualization on this dashboard is a Heat Map of many different demographics data points along with ChurnRisk. On this map, the darker the square, the more the risk.



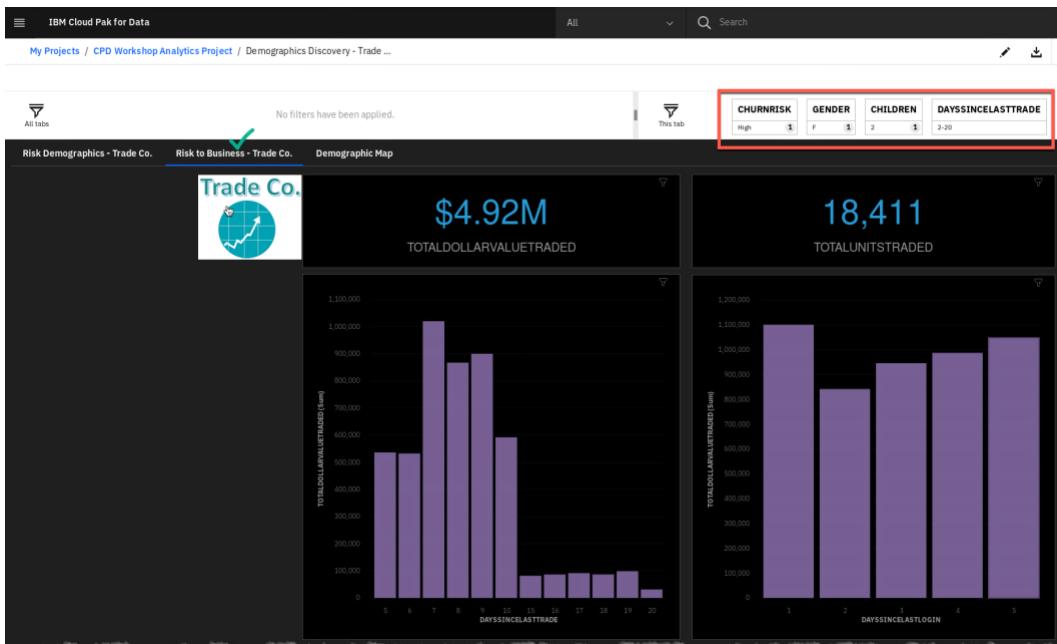
- 17. Hover over darker purple square on the top left of the Heat Map. Notice it represents Married Females with zero children, who are not homeowners and are 20 years of age.



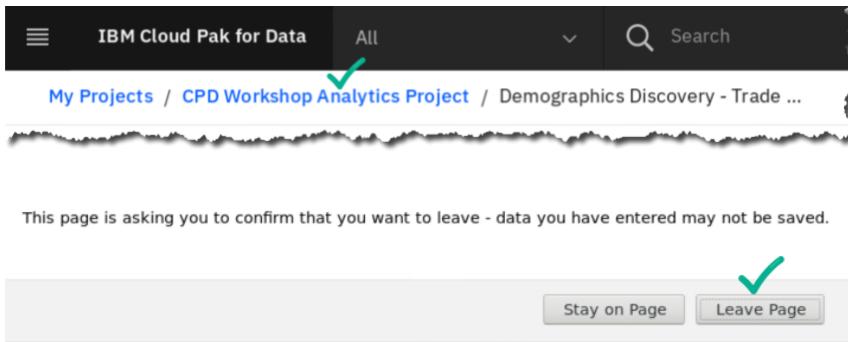
- 18. Click on the second tab in this dashboard called [Risk to Business – Trade Co.](#).

This analyzes Dollars and Units traded on the top, filtered by the demographics you found in the first tab visualization that appear to be the higher risk. (Note: these filters can be changed at any time in Edit mode to give flexibility to the user.)

This tab also has bottom visualizations that demonstrate that even though this demographic appears to be high risk, they are still fairly active.



- __19. Leave this dashboard by clicking on the breadcrumb trail to navigate to the project again.



2.6 Devising a plan

Armed with his findings, the Business Analyst brings this to the Data Scientists on the team. Together they come up with a plan to fix their current Rules Based system of predicting churn.

Separation(Churn) Risk: Current Rules Based System

Built Using Limited Data
Rules are developed using a single source of data that contains customer demographic information.

Manual Process to Develop Rules
Rules are manually developed based on the past experience of the marketing team. Rules are only updated once a year.

Low Overall Predictive Accuracy
Low overall predictive accuracy. We are both missing identifying customers who ultimately separate and incorrectly assigning high risk to customers who ultimately stay.

The intent is to use a better approach: Leverage Cloud Pak for Data to build a data driven Machine Learning model to infuse into their Stock Trader application.

Separation (Churn) Risk: New Data Driven Approach



Incorporate Multiple Data Sources

Use vetted centralized transactional data along with customer demographics to understand separation behavior. Also, include the outcomes of the rules based system for each customer where an accurate prediction was rendered.



Data Driven Process to Develop Machine Learning Models

Develop predictive models for separation risk that automatically discover and incorporate all the patterns in the data including interactions and contingencies.

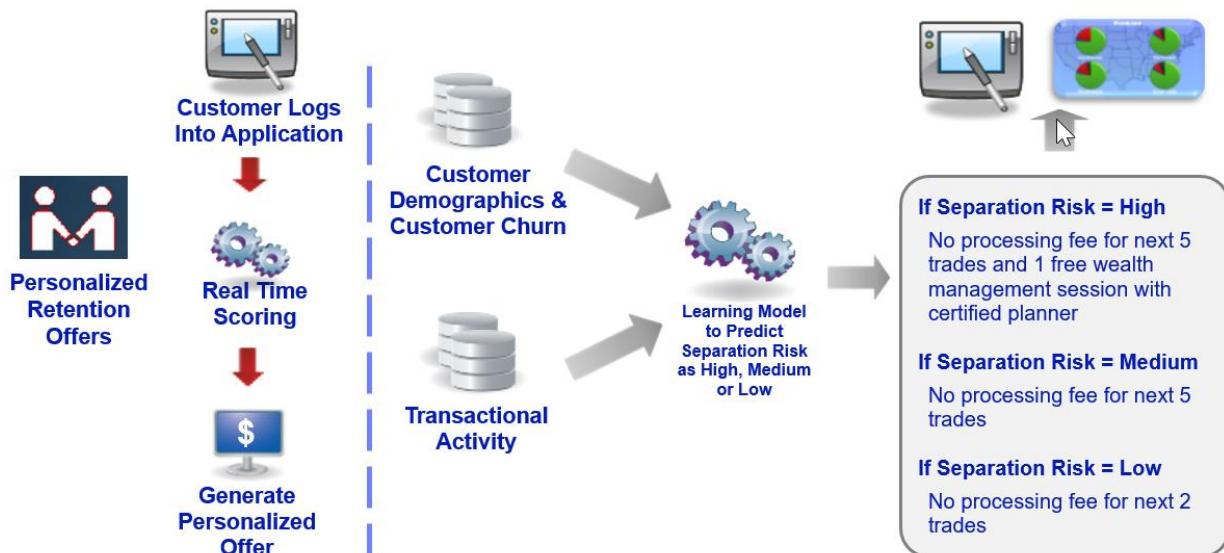


High Accuracy from Adaptive Machine Learning

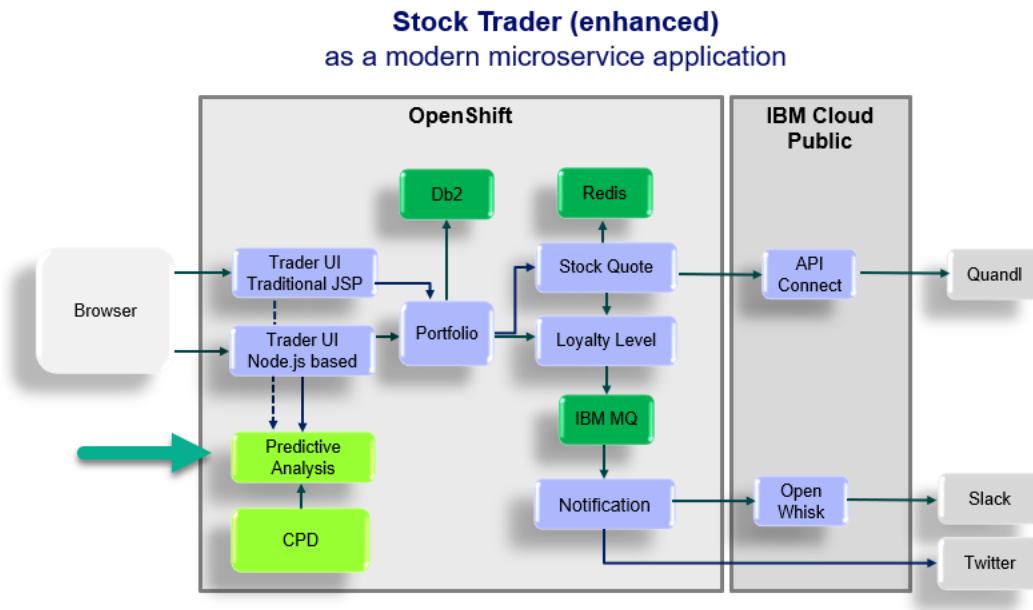
Models will classify separation risk with a higher overall accuracy and will adapt to changing patterns in risk to maintain that accuracy. Machine Learning models will incorporate all the understanding from the rules-based system and build on that to develop highly complex set of predictive conditions.

With this better model that more accurately predicts churn, the developers will infuse it into their Stock Trader application that will act upon this information in real time. Once a customer (trader) logs into the application, the new ML model will predict the risk of churn for this customer and will then make one of three offers designed to retain them.

Deployment: Stock Trader App with Integrated AI



In this case, infusing the new ML model into the Stock Trader application will be as simple as writing one extra microservice that is invoked upon sign in. It will be called the “Predictive Analysis” service (indicated by the green arrow below) and will make the offer after sign in.



2.7 Reviewing the dashboard: Monthly Metrics After AI

The model was created and infused into the Stock Trader application, which ran for three months. New activity data was captured for this period. The Business Analyst created another dashboard similar to the first that displays the results.

- _20. Still in the same project, under [Assets](#), scroll down until you find [Dashboards](#).

Click [Monthly Metrics - After AI Trade Co.](#)

The screenshot shows the 'Assets' tab selected in the navigation bar. Below it, a dropdown menu titled 'Dashboards' is open, listing three items: 'Demographics Discovery - Trade Co.', 'Monthly Metrics - Trade Co.', and 'Monthly Metrics - After AI Trade Co.'. The third item is highlighted with a green checkmark icon.

- _21. This dashboard is similar to the first one in the two visualizations:

[Shares sold per month](#) and [Number of Traders per month](#)

Notice however, new data is present for October, November, and December.



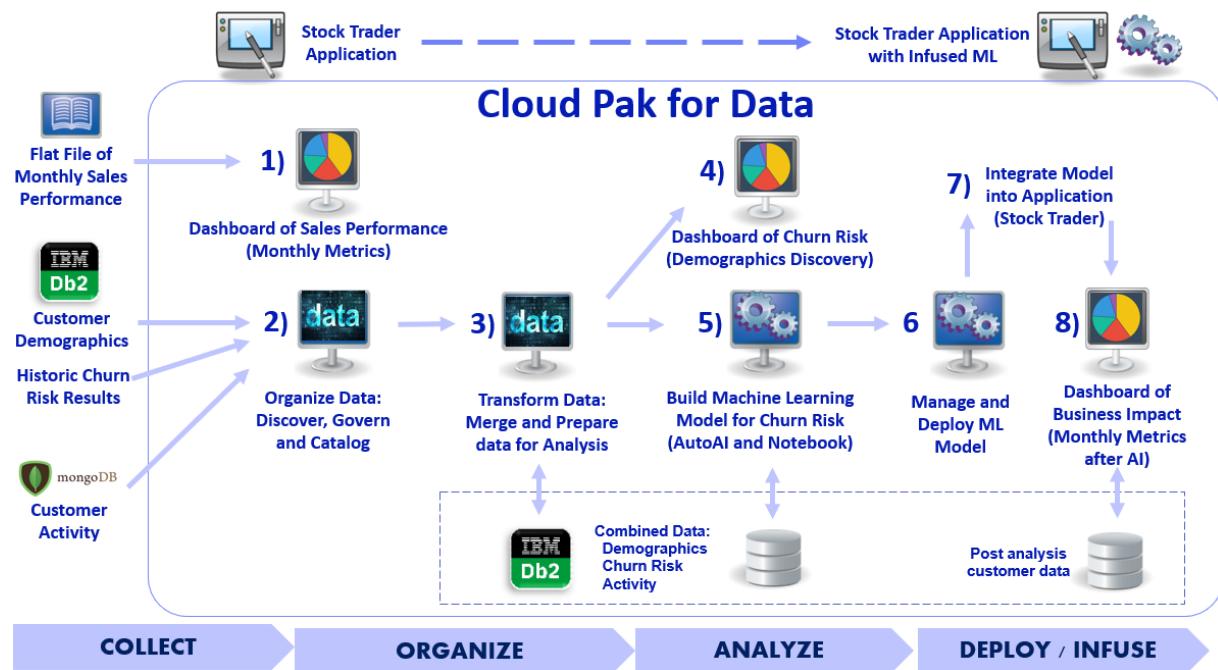
The data is very encouraging: both data points have gone up!

Trade Co. is back in business with their successful customer retention offer driven by a smarter application infused with a machine learning model that was built and maintained on Cloud Pak for Data.

2.8 Lab conclusion

This workshop will walk you through how Trade Co. was able to pull off their business success using the Cloud Pak for Data platform.

Each subsequent lab will walk you through one of the steps (Collect, Organize, Analyze, Deploy, Infuse) taken by Trade Co. in their Journey to Cloud and AI using this amazing, industry leading Analytics Modernization platform called Cloud Pak for Data. You will also take on the roles of the various personas involved along the way. The steps in the journey are depicted in the illustration below.



** End of Lab 02 – Business Use Case: Customer Churn

Lab by Burt Vialpando and Kent Rubin, IBM

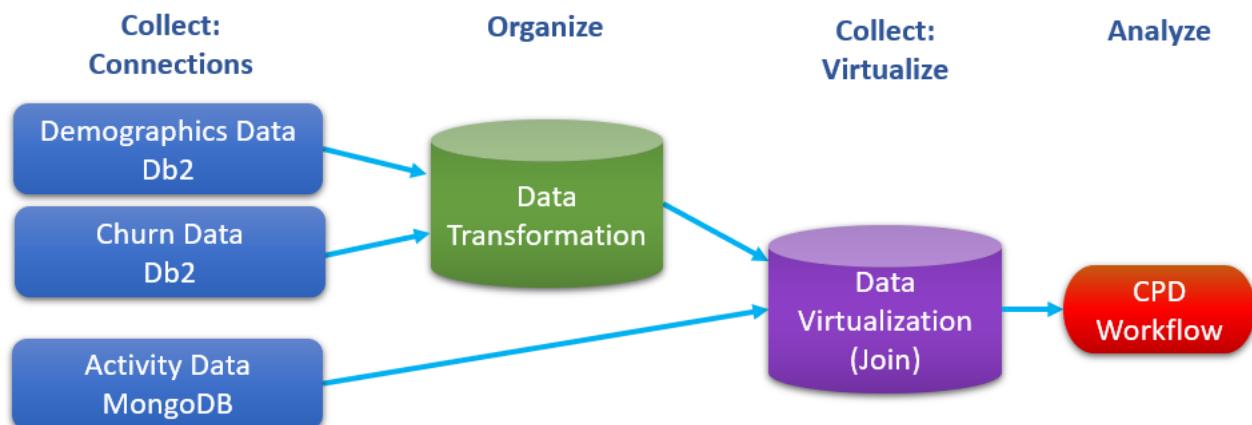
Lab 03 COLLECT: CONNECTIONS

3.1 Lab overview

The [Collect](#) capability of CPD means accessing your organization's data regardless of where it resides, whether that be an in-cluster data source (e.g. the Db2 Advanced Edition) or a native connection to a remote data source. You can even use the Db2 Event Store or Streams services to provide streaming access that is best suited for Internet of Things (IoT) processing.

Additionally, Data Virtualization and Data Transformation are available to streamline the access, performance, and formatting of the data for use in later steps of the CPD analytics workflow.

In this lab you will explore [Connections](#) for the [Collect](#) process. In later labs you will explore [Data Transformation](#) and [Data Virtualization](#).



3.2 Persona represented in this lab

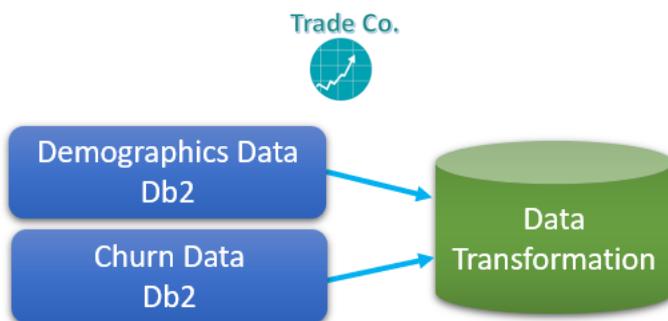
The [Data Engineer](#) persona is the likely role to perform the various [Collect](#) tasks shown in this lab.

Persona (Role)	Capabilities
 Data Engineer	Data Engineers build and optimize the systems to allow data scientists and business analysts to perform their work. The Data Engineer ensures that any data is properly received, transformed, stored, and made accessible to other users.

3.3 Db2 data overview – Transforming for analytics

The Db2 data in our Trade Co. scenario (Demographics and Churn) have two key factors that make an appropriate scenario for [Data Transformation](#):

- The data only changes once a month (it is relatively static). Thus, copying it and/or changing it for downstream processing in our CPD analytics workflow is OK because we are not required to have the latest data to get the results we need.
- The data must be changed before it can be processed in our CPD analytics workflow.



Although you will be doing the actual [Data Transformation](#) in the next lab, you will prepare for it in this one by reviewing some steps in this [Connections](#) lab.

3.3.1 Logging into the CPD web client (if you have not already done so)

- __1. If you are starting this lab stand-alone (without going through previous labs) do the following:
- __2. Double-click the desktop icon: [Cloud Pak for Data Web Client](#).



- __3. The CPD web client GUI displays as shown.

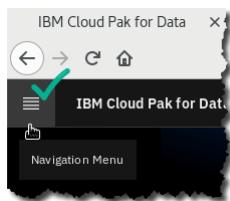
Use [cpduser](#) and [cpdaccess](#) for the *Username* and *Password* and click [Sign in](#).



3.3.2 Reviewing the Db2 Advanced Edition data

- 4. To review this data, you will be launching the *Db2 Advanced Edition Data Management Console*.

In the CPD web client, click the [Navigation Menu](#) (“hamburger” icon).



- 5. Click [Collect](#) ⇒ [My Data](#) ⇒ [Databases](#) ⇒ [Db2 Advanced Edition](#) ⇒ [ellipsis](#) ⋮ ⇒ [Open database](#) (Note: Please wait for the green check to appear next to Db2 Advanced Edition first.)

- 6. At the top left of the Db2 Advanced Edition console, click on: [Summary](#) ⇒ [Explore](#) ⇒ [Tables](#).

__7. Select schema [CPDUSER](#).

Schemas	
Name	Type
USER1001	User
<input checked="" type="checkbox"/> CPDUSER ✓	User
SOLUTIONS	User

__8. Click [CUSTOMER_CHURN](#) and click on it to bring up the table definitions view.

Click on [View Data](#).

The screenshot shows a database interface with a sidebar titled "Schemas" containing three entries: "USER1001", "CPDUSER" (selected), and "SOLUTIONS". The main area displays the "Tables" section with two tables listed: "CUSTOMER_CHURN" (selected) and "CUSTOMER_DEMOGRAPHIC". The "CUSTOMER_CHURN" table definition is shown, indicating approximately 2066 rows (64 KB) last updated on 2020-06-22 17:57:33. The table has two columns: "ID" (SMALLINT) and "CHURNRISK" (VARCHAR). The data view shows five rows of data:

ID	CHURNRISK
0	Low
1	Low
2	Low
3	High
4	High

This data represents churn risk that was determined through a manual methodology by Trade Co. It will be enhanced by a machine learning methodology created in this workshop.

__9. Click on the [Back](#) icon, then table [CUSTOMER_DEMOGRAPHICS](#) ↴ [view Data](#).

The screenshot shows a database interface with a sidebar titled "Schemas" containing three entries: "USER1001", "CPDUSER" (selected), and "SOLUTIONS". The main area displays the "Tables" section with two tables listed: "CUSTOMER_CHURN" and "CUSTOMER_DEMOGRAPHICS" (selected). The "CUSTOMER_DEMOGRAPHICS" table definition is shown, indicating approximately 2066 rows (320 KB) last updated on 2020-06-22 17:58:22. The table has seven columns: ID, GENDER, STATUS, CHILDREN, ESTINCOME, HOMEOWNER, and AGE. The data view shows four rows of data:

ID	GENDER	STATUS	CHILDREN	ESTINCOME	HOMEOWNER	AGE
0	F	S	1	38000.00	N	24
1	M	M	2	29616.00	N	49
2	M	M	0	19732.80	N	51
3	M	S	2	96.33	N	56

This data represents the known characteristics of the Trade Co. customers themselves. The machine learning methodology of this workshop will take advantage of this data to determine churn risk more accurately.

3.3.3 Reviewing the Db2 connection

- _10. In order to create a connection to a CPD data source, you would first need the connection information. To review what that looks like for the in-cluster Db2 Advanced Edition, click:

[Navigation Menu](#) \Rightarrow [Collect](#) \Rightarrow [My Data](#) \Rightarrow [Databases](#) \Rightarrow [Db2 Advanced Edition](#) \Rightarrow [ellipsis](#) \vdots \Rightarrow [Details](#)

- _11. This takes you to the details page for this data source.

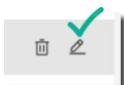
On the right you will find the [Access Information](#) section which gives you the information you need to create a connection.

Note: [Username](#) and [Password](#) are explicitly shown (and easily copied using the icons) while the JDBC Connection URL string tells us the [Host:Port/Database](#).

- _12. Now you can review the connection to the Db2 Advanced Edition that was already created for this workshop.

Go to: [Navigation Menu](#) \Rightarrow [Connections](#) \Rightarrow [Db2 Advanced Edition](#)

- __13. Hover over the left part of the connection and click the [Edit](#) (pencil) icon.



- __14. The Db2 Advanced Edition [Connection](#) edit screen looks like this:

Db2 Advanced Edition

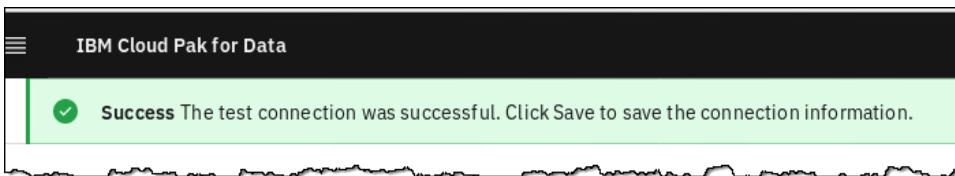
Connection name	Db2 Advanced Edition
Description (optional)	Enter a description for the connection
Connection type	Db2
Host	worker5.clusterw9
Port	32030
Database	BLUDB
Username ⓘ	user1001
Password	*****

Each data source type (Db2, Mongo, Hive, S3, Drop Box, and so on) has their own connection page format. Many of the relational database sources are similar to this one, while others require different credentials altogether.

- __15. At the bottom of the screen, click [Test connection](#).

Username ⓘ	user1001	Password	*****
Cancel		Test connection	Save

- __16. At the top left of the screen, you should see this Success message.



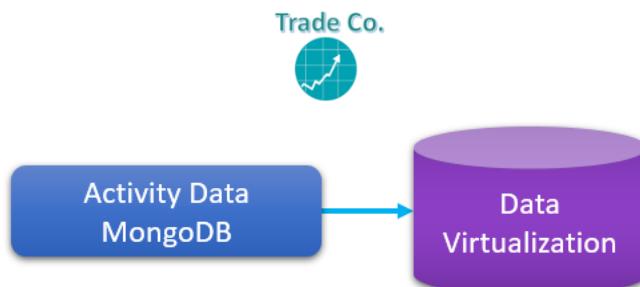
- __17. To exit the edit screen, click [Cancel](#). DO NOT SAVE this connection.

 Data Engineer	<p>Note: if you do NOT have a successful connection test as shown above, you should fix this by launching another CPD web client from the desktop and use the second web client to return to the Details screen of the Db2 Advanced Edition. Use both screens to make sure all the credentials match between the Details screen and the Connection screen, especially Username/Password.</p>
------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

3.4 MongoDB data overview – Virtualizing for analytics

The Mongo data in our scenario (Activity) has two key factors that makes for an appropriate scenario for Data Virtualization:

- The data is constantly changing (it is dynamic)
- Our analytics processing requires the latest data



We have chosen semi-structured data in this scenario to demonstrate the power of CPD Data Virtualization – it even works with JSON documents. But we could just have easily chosen a structured database source for it, like Db2, Informix, SQL Server, and so on.

Semi-structured (and even unstructured) data is commonly used in many systems of engagement applications (like Trade Co.'s Stock Trader) so this was another reason we chose this as an example data source type.

Our scenario presumes that this data comes from a mobile application that will be constantly changing the data, and that we would require the latest data for our analytics workflow.

3.4.1 Reviewing the MongoDB data

The MongoDB database was provisioned after the MongoDB Enterprise service was installed, which is also located in [My data](#).

__18. Start at the [Navigation Menu](#).

__19. Click [Collect](#) ⇒ [My data](#) ⇒ [Databases](#) ⇒ [MongoDB-1](#) ⇒ [ellipsis :](#) ⇒ [Details](#).

__20. Examine the section [MongoDB Ops Manager](#) on the right.

Here you can review the username for the Ops Manager (which is cpduser).

Copy the [Password](#) using the icon.

First Ops Manager user	cpduser	<input type="checkbox"/>
Password	*****	
Replica set	jdbc:ibm:mongodb://mongodb-1592925755860-replica:27017/	Copy password

- _21. Now return to the [My data](#) screen and choose [Databases](#) ⇒ [MongoDB-1](#) ⇒ [ellipsis](#) : ⇒ [Open Database](#).

The very first time you log into this *MongoDB Ops Manager* console it may not ask you for user and password credentials. We have saved these for you in a browser setting.

However, if the credentials were not saved and should you encounter this Login screen (as above) you can use the password copied from the [Details](#) screen.

Note: if the copy/paste of the password does not work for you using [Details](#) screen icon, then try this technique:

1. In the [Details](#) screen, perform the “copy” by highlighting the password and [right-click](#) ⇒ [copy](#).
2. In the [Ops Manager](#) login screen use [right-click](#) ⇒ [paste](#).



Data
Engineer

- __22. In the *MongoDB Ops Manager* console tab, click **Deployment** on the top left of the screen.
In the **Processes** tab, click on the **Replica Set** link,

The screenshot shows the MongoDB Ops Manager interface. The top navigation bar includes 'Access Manager' and 'Support'. Below it, a breadcrumb path shows 'MONGODB-1592925755860-REPLICA-SET > MONGODB-1592925755860-REPLICA-SET'. The main area is titled 'Deployment' with a green checkmark above it. The 'Processes' tab is selected, indicated by a green underline. Other tabs include 'Servers', 'Agents', and 'Security'. A search bar and a 'VIEW' dropdown with 'CLUSTERS' and 'LIST' options are also present. At the bottom, a list shows a single replica set: 'mongodb-1592925755860-replica-set' with version 4.2.6, backup status 'Disabled', and MongoDB count '1 Mongod'. A green checkmark is placed over the 'mongodb-1592925755860-replica-set' link.

- __23. Click on the **Data** tab to review the JSON documents in database **mongodb.activity01**.

The screenshot shows the 'Data' tab for the 'mongodb-1592925755860-replica-set' deployment. The top navigation bar includes 'Access Manager' and 'Support'. The main area is titled 'mongodb-1592925755860-replica-set' with a green checkmark above it. The 'Data' tab is selected, indicated by a green underline. Other tabs include 'Overview', 'Real Time', 'Metrics', 'Profiler', and 'Performance Advisor'. A 'FILTER' button with the value '("filter": "example")' is shown. The 'mongodb.activity01' database is selected, showing collection size 283.43KB, total documents 1000, and index size 44KB. The 'Find' tab is selected. A red box highlights a specific JSON document in the 'QUERY RESULTS 1-20 OF MANY' section:

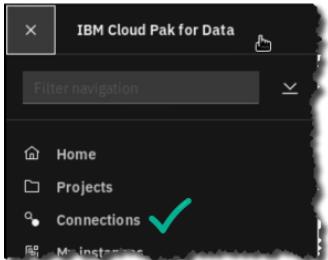
```
_id: ObjectId("5ef22f6d16a80b97d03a2a00")
ID: 0
TotalDollarValueTraded: 59755.98
TotalUnitsTraded: 206
LargestSingleTransaction: 29877.99
SmallestSingleTransaction: 2987.799
PercentChangeCalculation: 51.5
DaysSinceLastLogin: 3
DaysSinceLastTrade: 10
NetRealizedGains_YTD: 2987.799
NetRealizedLosses_YTD: 0
```

- __24. When finished, close the tab **Data | MongoDB Ops Manager**.

The screenshot shows a browser window with two tabs: 'IBM Cloud Pak for Data' and 'Data | MongoDB Ops Manager'. The 'Data | MongoDB Ops Manager' tab has a red X icon indicating it is closed. The address bar shows the URL 'cp4d-cpd-cp4d.apps.clusterw9:3274'.

3.4.2 Reviewing the MongoDB connection

_25. Go to [Navigation Menu](#) ⇒ [Connections](#)



_26. Click on connection [MongoDB](#) line and choose the [Edit](#) icon.

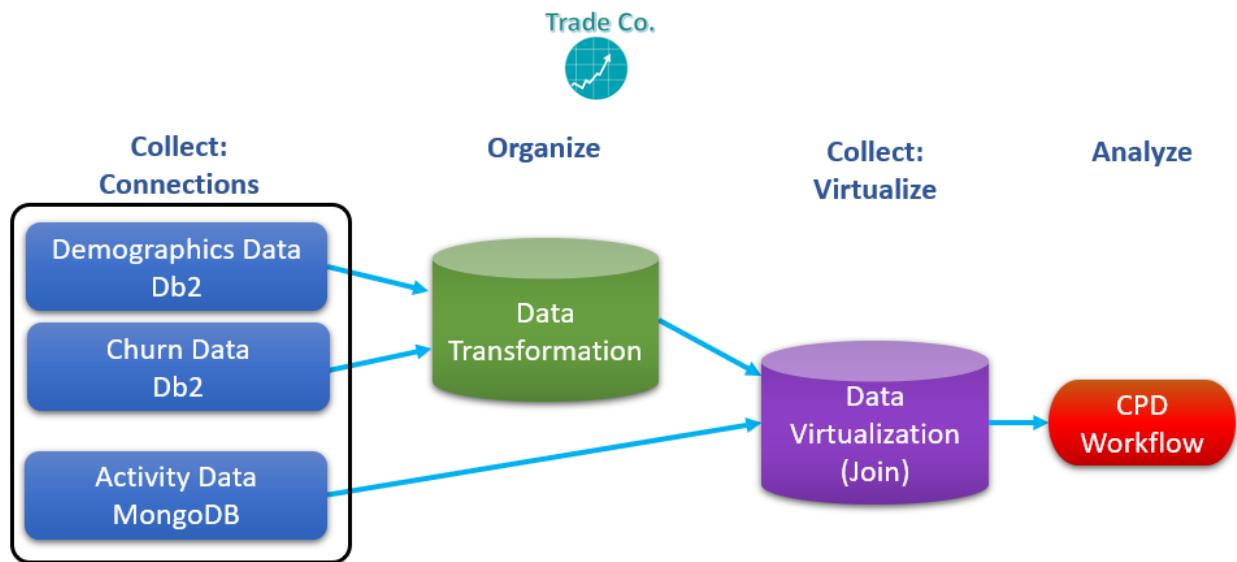


_27. Note that the credentials information is taken from the Details screen you reviewed earlier.

This screenshot shows the "Edit" screen for the "MongoDB" connection. It includes fields for Connection name (MongoDB), Description (optional), Connection type (Mongo DB), Host (mongodb-1592925755860-replica-set-0.mongodb-1592925755860), Port (27017), Database (MONGODB), Username (cpduser), and Password (redacted). A callout bubble points to the connection type field with the text "This information is taken from the MongoDB Details screen". At the bottom are "Cancel" and "Save" buttons.

3.5 Lab conclusion

In this Collect: Connections lab, you reviewed connections for the Db2 Advanced Edition and MongoDB data. This prepares you for the Organize and Collect: Virtualize labs in this workshop.



In the [Organize](#) lab coming up, we will be transforming the Db2 data sets into one.

After that, we will finish up the [Collect](#) processing by virtualizing the results from both the Db2 Transformation output and the MongoDB data together.

**** End of Lab 03 - Collect: Connections**

Lab by Burt Vialpando and Kent Rubin, IBM

Lab 04 ORGANIZE

4.1 Lab overview

Many organizations find it difficult to understand their own data because it originates from many sources, is dispersed across many silos, and is controlled by different teams.

This [Organize](#) lab will show you how to uncover the hidden value in your organization's data and how to build a lineage that is otherwise difficult to establish. Cloud Pak for Data helps your organization move from the manual processes required to establish relationships between data to an automated one aided by the platform's built-in machine learning capabilities.

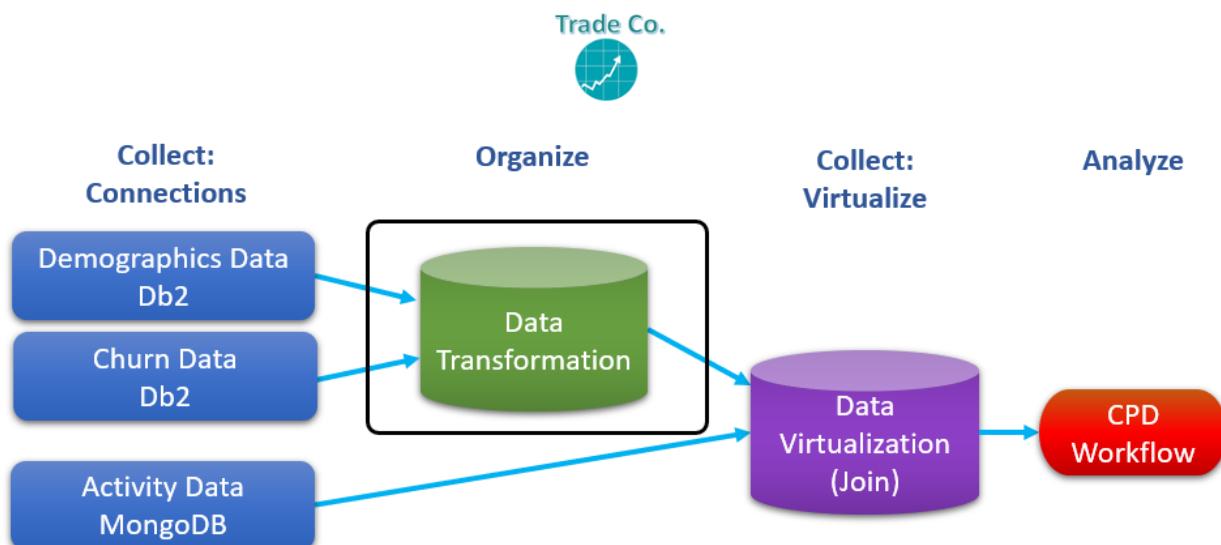
4.2 Persona represented in this lab

The [Data Steward](#) persona is the likely one to perform most of the [Organize](#) tasks shown in this lab.

Persona (Role)	Capabilities
 Data Steward	Data Stewards integrate and transform data as well as provide governance, lineage and classification of the data.

Note: this persona often works closely with the [Data Engineer](#) because they both work with the data to prepare it for analytics processing by other personas. For example, in this lab one of activities the [Data Engineer](#) will do is to create a Transformation job with the Db2 data. In the next lab, the [Data Engineer](#) will go on to use that output table to create a final virtualized view of all the data sources joined together.

Note: The Data Steward persona also works closely with the [Data Quality Analyst](#) persona.



Before we start transforming data, let's first explore the other crucial aspects of the CPD [Organize](#) capabilities.

4.3 Logging into the CPD web client (if you have not already done so)

- 1. If you are starting this lab stand-alone (without going through previous labs) do the following:
- 2. Double-click the desktop icon: [Cloud Pak for Data Web Client](#).



- 3. The CPD web client GUI displays as shown. Use `cpduser` and `cpdaccess` for the *Username* and *Password* and click **Sign in**.

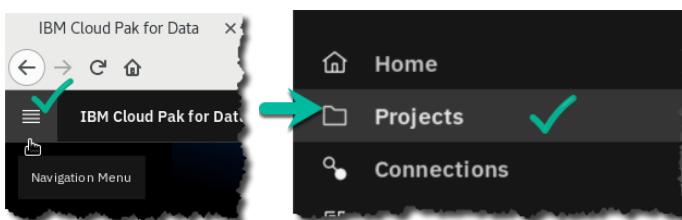


4.4 Reviewing a data asset in the project

Some of the more interesting “Organize” capabilities can be done on a data asset that has been added to a project. In this Trade Co. scenario, the team has created a Project from which they will all collaborate and work from together.

4.4.1 Data asset overview

- 4. In the CPD web client, click the [Navigation Menu](#) (“hamburger” icon) \Rightarrow [Projects](#).



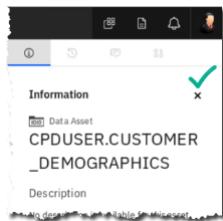
- 5. Select the project: [CPD Workshop Analytics Project](#).

Name	Project type
CPD Workshop Analytics Project ✓	Analytics
CPD Workshop Transform Project	Data transform

- __6. Under tab **Assets** ⇒ **Data assets**, click **CPDUSER.CUSTOMER.DEMOGRAPHICS**.

Name	Type	Created by
CPDUSER.CUSTOMER.CHURN	Data Asset	CPD User
CPDUSER.CUSTOMER_DEMOGRAPHICS	Data Asset	CPD User
customer_demochurn_activity_analyze.csv	Data Asset	CPD User

- __7. Close (x) the Information window on the top left.



- __8. In the **Preview** section, general information about the data asset is displayed, as well as column specific information and sample data.

ID	GENDER	STATUS	CHILDREN	ESTINCO...	HOM
1	Female	Active	0	Not cl...	T...

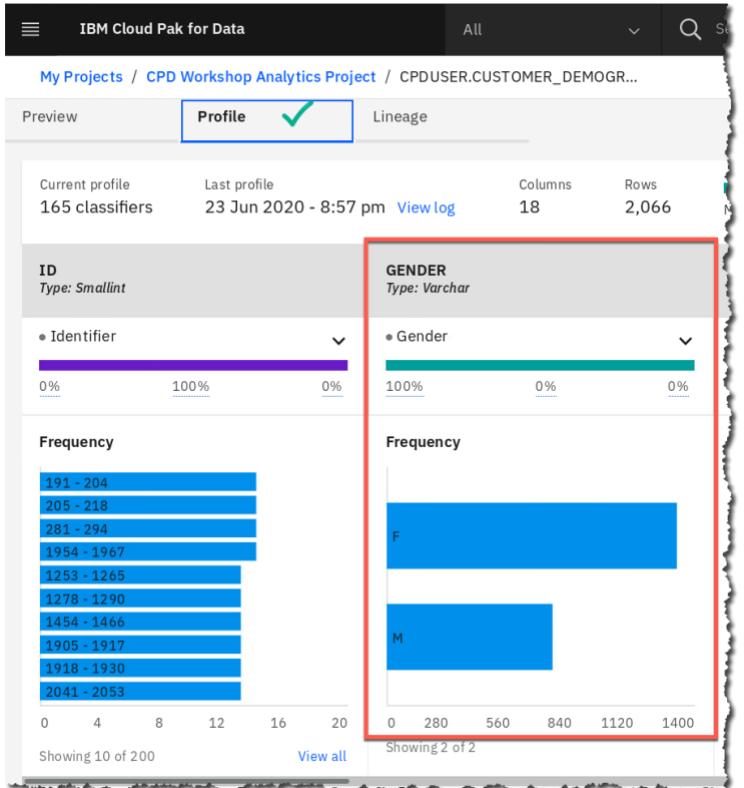
- ___9. Click on the down arrow for column **TAXID** to see how Profiling categorized the data.

The identifier is: **US Social Security Number**.

ID	GENDER	STATUS	CHILDREN	ESTINCO...	HOMEOW...	AGE	TAXID	CREDITO
	Smallint	String	String	Smallint	Decimal	String	Smallint	Char
Identifier	Gender	Code	Code	Not cl...	Indica...	Code	US So...	Identifier
0	F	S	1	38000	N	24	US Social Se...	6549061
1	M	M	2	29616	N	49	Canadian So...	6436360
2	M	M	0	19732.8	N	51	Not classifie...	4849378
3	M	S	2	96.33	N	56		2926742
4	F	M	2	52004.8	N	25	Profiling Tra...	141013706
								4132500

- ___10. Click **Profile** to see more details on the profiling information. (It may take a minute or two to render – be patient.)

Notice column **GENDER** was assigned Identifier Gender, which has 2 of 2 values (M and F) and a displayed frequency for each.



- __11. Scroll over to column **TAXID** and click the twisty to see how this data is profiled.

TAXID ✓
Type: Varchar

• US Social Security Number ✓ ^

- ✓US Social Security Num... 93% matches
- Canadian Social Insura... 9% matches
- Not classified 5% matches
- Routing Transit Number 3% matches
- Individual Taxpayer Ide... 2% matches

[View all](#)

189017262 - 194603113
176441136 - 181073923
170039718 - 176163659
139489175 - 144122501
106198864 - 110729933
111063740 - 115332481

4.4.2 Refine visualizations

- __12. Click back to **Preview** then click **Refine**

IBM Cloud Pak for Data All Search

My Projects / CPD Workshop Analytics Project / CPDUSER.CUSTOMER_DEMOGR...

Preview ✓ Profile Lineage

Schema: 18 Columns | 2066 rows Preview: 1000 rows Last refresh: just now Refine ✓

ID Smallint	GENDER String	STATUS String	CHILDREN Smallint	ESTINCO... Decimal	HOMEOW... String	AGE Smallint	TAXID String	CREDITC... Char	DOB Date	ADDRESS... String	ADDRESS... String	CITY String
Identifi... v	Gender v	Code v	Code v	Not cl... v	Indica... v	Code v	US So... v	Identifi... v	Date o... v	Text v	Not cl... v	City
0	F	S	1	38000	N	24	147889187	6549061697939	1947-11-11	159 HUTTON ST	ABS	
1	M	M	2	29616	N	49	113772166	6436360484417	1992-03-17	31 WOODLAND F	SAI	
2	M	M	0	19732.8	N	51	132420919	4849378808118	1907-09-08	1910 COCHRAN	KEA	
3	M	S	2	96.33	N	56	700548452	2926742654852	1980-04-29	187 HAYES MILL	RUS	

- __13. Click **Visualizations ↴ Map**

(Hint: you may have to select the double down arrow to see **Map**.)

IBM Cloud Pak for Data All Search

My Projects / CPD Workshop Analytics Project / CPDUSER.CUSTOMER_DEMOGR... / Refine data

Operation + Code an operation to cleanse and shape your data

Data Profile Visualizations

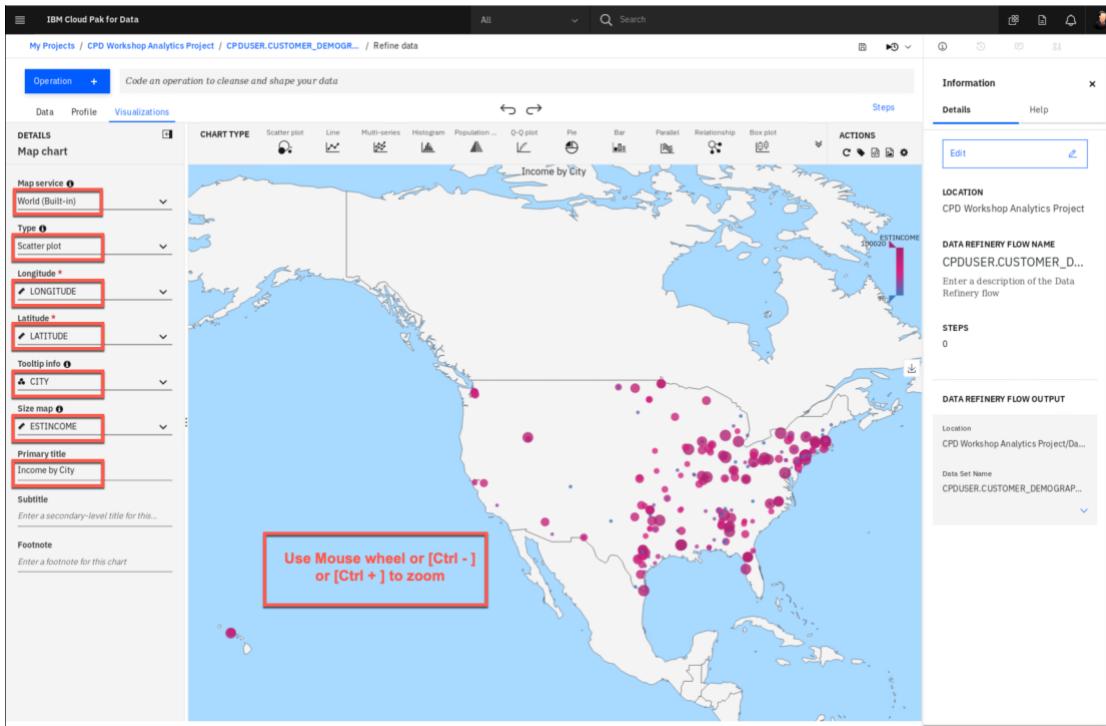
CHART TYPE Scatter plot Line Multi-series Histogram Population ... Q-Q plot Pie Bar Parallel Relationship Box plot Heat map t-SNE Word cloud

ALL CHARTS (28)

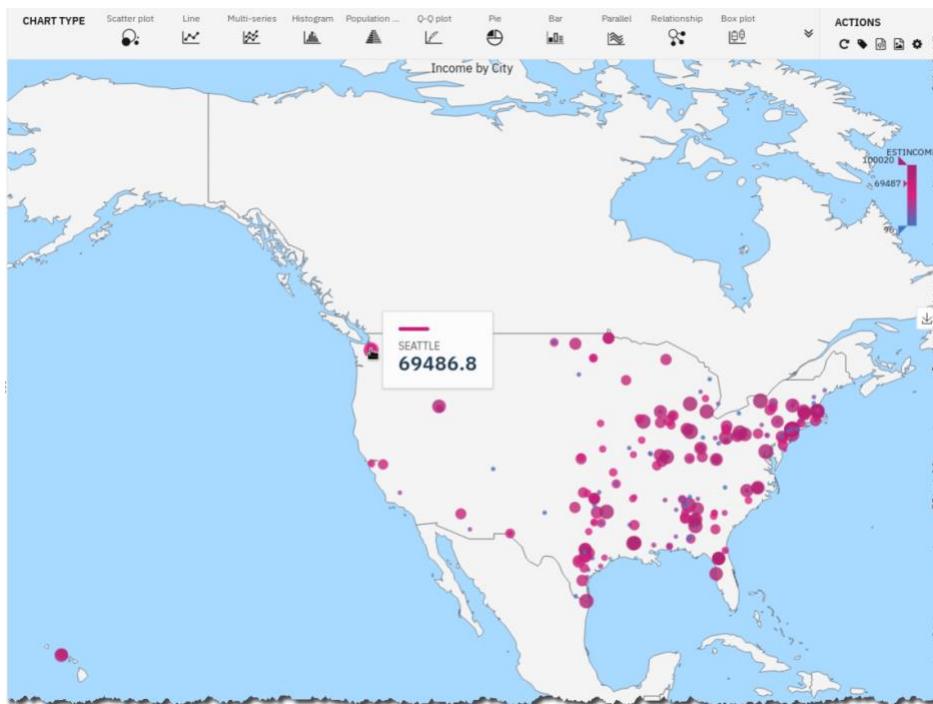
Scatter plot Line Multi-series Histogram Population ... Q-Q plot Pie Bar Parallel Relationship Box plot Heat map t-SNE Word cloud

Map Treemap Radar Theme River Circle packi... Bubble 3D Error bar Scatterplot ... Candlestick Multi-chart Dual Y-axes Time plot Customized

- __14. Fill in the Details as World, Scatter plot, LONGITUDE, LATITUDE, CITY, ESTINCOME and Income by City

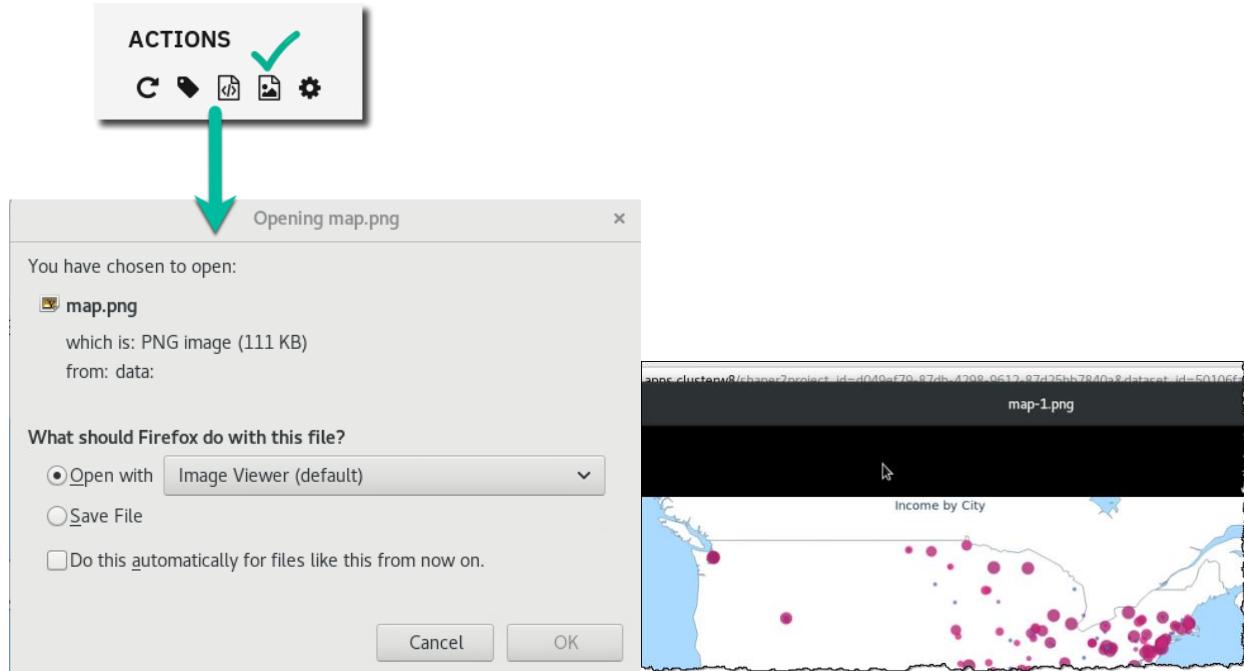


- __15. Use your mouse's wheel to Scroll up and center the United States on the visualization.
Use the [Ctrl] – and [Ctrl] + to zoom in and out to properly size the visualization.
- __16. Hover over any circle on the map to see the Estimated Income for that City.

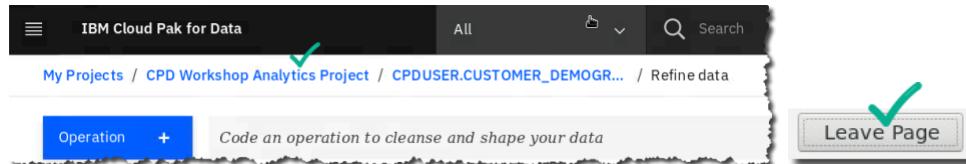


- 17. If the Data Scientist or Business Analyst wants to keep a visualization, they can download it by clicking on the **image icon** shown.

They can save or review it by downloading it from the **ACTIONS** section and then using **Open with Image Viewer**.



- 18. Close the Image viewer (if you have it open) and click on the breadcrumb back to the link for the **CPD Workshop Analytics project** then **Leave Page**.



4.4.3 Data refinery flows

You should be back at the project [CPD Workshop Analytics Project](#). If you are not there, navigate there as you did at the beginning of this lab.

- __19. Under **Assets**, scroll down to find **Data Refinery flows**

Click [CPDUSER.CUSTOMER.DEMOGRAPHICS_flow](#)

Name	Type
CPDUSER.CUSTOMER_DEMOGRAPHICS_flow	Data Refinery flow

- __20. This Refinery flow was pre-created for you to shorten this lab. It was created in the Refine screen you were in earlier. Click the **Steps** Box (1).

There is only one step in this job, which is a “Text” operation on the ZIP4 column.

Hover over the step to see the icons for it.

Click the **edit** (pencil) icon (2) to review the details of the step.

1 Steps

2

1 Steps

Data Source

CPDUSER.CUSTOMER_DEMOGRAPHICS

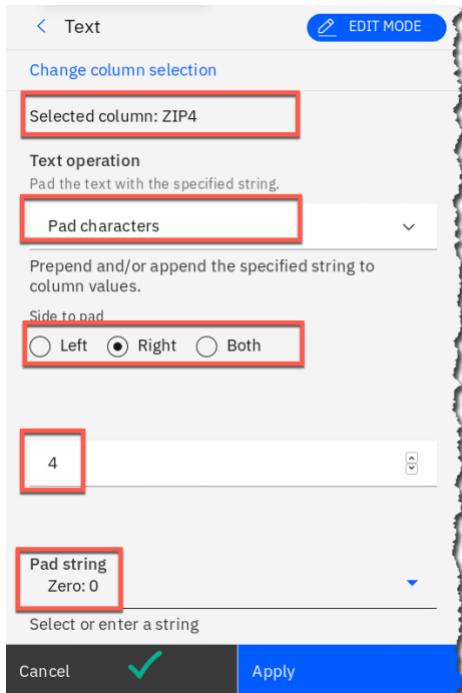
Text

Padded text in ZIP4 with 0 on the right for total length of 4 into ZIP4

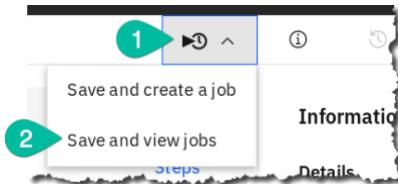
__21. This step refines the ZIP4 column by padding 4 characters of “0” to the Right.

This means if the column is already filled in, it will not do anything. If it is empty or filled with 0, it will fill it in with 0000.

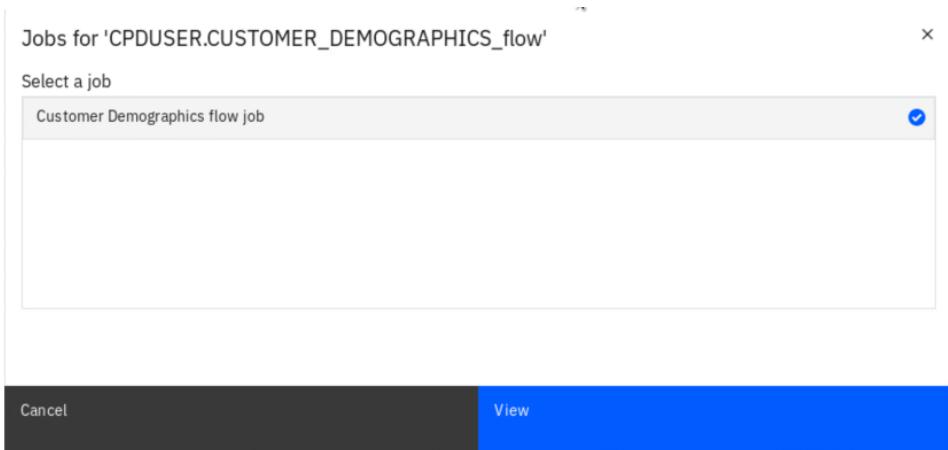
Click [Cancel](#) (the job is OK as it is, but you could have changed it here if you needed to).



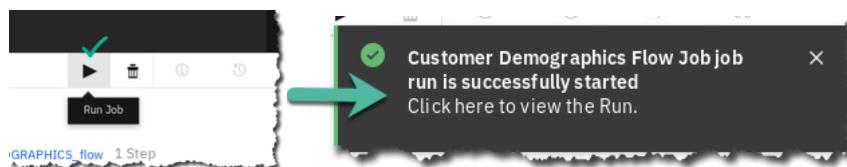
__22. Click on the [Jobs](#) (arrow by a clock) icon, then [Save and view jobs](#).



__23. Click the [Customer Demographics Flow Job](#) shown, then [View](#).



- __24. Click the [Run Job](#) icon.



Note: If the job fails for any reason, reselect the source csv file to prime it and rerun the job.

- __25. The job takes about a minute to complete. You can see it Running on the screen.
Hit the [refresh](#) icon to see status more quickly.



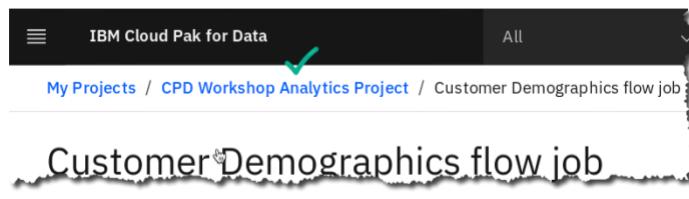
- __26. When [Completed](#), you can check the logs of the job if you want to. This can be especially useful if the job is complex.

Click on it to review the details of the job.



 Data Steward	If the job run fails, simply click the ellipses on that job run and delete it, then run it again.
--------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------

- __27. Click on the breadcrumb trail to return back to the project.



- __28. A new [Data Asset](#) was created as output from the Flow (refine job).
Scroll down to find and then click on [Customer Demographics_shaped.csv](#)



__29. Scroll over to find the ZIP4 column.

Notice the data is padded with four zeroes if it was less than four zeroes. If the row already had a four-digit value, it was left alone.

Schema: 18 Columns										
Preview: 1000 rows										
HOMEOW...	AGE	TAXID	CREDITC...	DOB	ADDRESS...	ADDRESS...	CITY	STATE	ZIP	ZIP4
N	49	113772166	6436360484417	1992-03-17	31 WOODLAND R		SAINT LOUIS	MO	63121	0000
N	56	700548452	2926742654852	1980-04-29	187 HAYES MILL		RUSTON	LA	71270	0000
N	19	163371244	2231773884473	1992-12-06	7850 45TH AVE N		CHESTER	MA	01011	0000
N	60	206227068	5553618912566	1912-11-23	515 KENSINGTO		ISSAQAH	WA	98075	0000
N	33	119762649	6813572896826	1977-02-09	188 W OLYMPIC I		EL PASO	TX	79925	0000
N	26	817366094	2046608099384	1905-02-01	21579 LARAMIE		PHILADELPHIA	PA	19104	0000
N	48	124158559	8979358234254	1933-06-14	4716 SW VIOLA (TENAFLY	NJ	07670	1057
Y	61	165912006	6325263828540	1947-01-25	4220 BARDSTOW		ARENA	ND	58494	0000
N	16	730825728	8111200911782	1916-10-05	D18 CALLE 5		HOUSTON	TX	77008	4708

The power of **Refine** can be used by more than just the Data Steward because it is launched from a Project, which enables many more personas the ability to shape and refine data assets to which the Data Steward may have given them access. This gives anyone the ability to explore and shape data on their own with self-service capabilities.

If the need is to make a permanent change to the data, for example, to write it back to a database from where it came, that is where **Transform** comes in. We will explore Transform later in this lab.

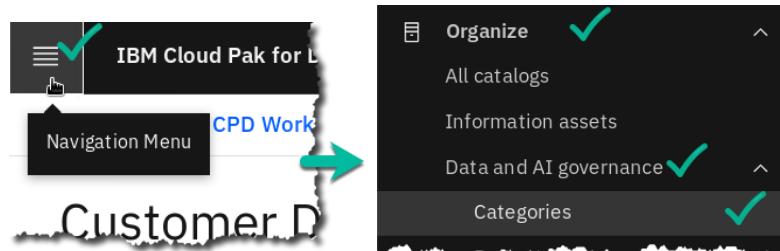
4.5 Reviewing a business glossary

A **business glossary** consists of **categories** and **terms**.

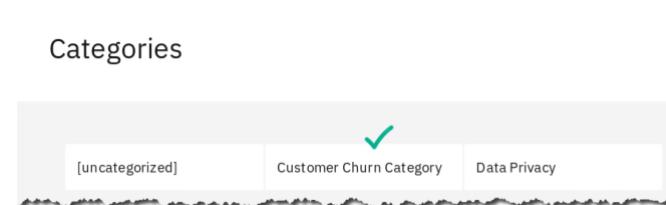
Categories provide the logical structure for the glossary so that you can browse and understand the relationships among terms and categories in the glossary. Categories can be organized in a hierarchy based on their meaning and relationships to one another.

A **Business term** is a word or phrase that describes a characteristic of the enterprise. Business terms are the fundamental building blocks of the glossary. Each Business term has a parent Category, but it can also be referenced by other Categories. When you create a Business term, you need to provide a meaningful name. Business terms can be assigned to other Business terms, and to other asset types as well.

- __30. Click **Navigation Menu** \Rightarrow **Organize** \Rightarrow **Data and AI governance** \Rightarrow **Categories**.



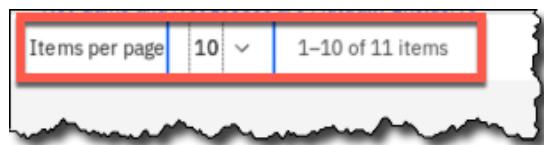
- __31. Review a category already created for you: **Customer Churn Category**.



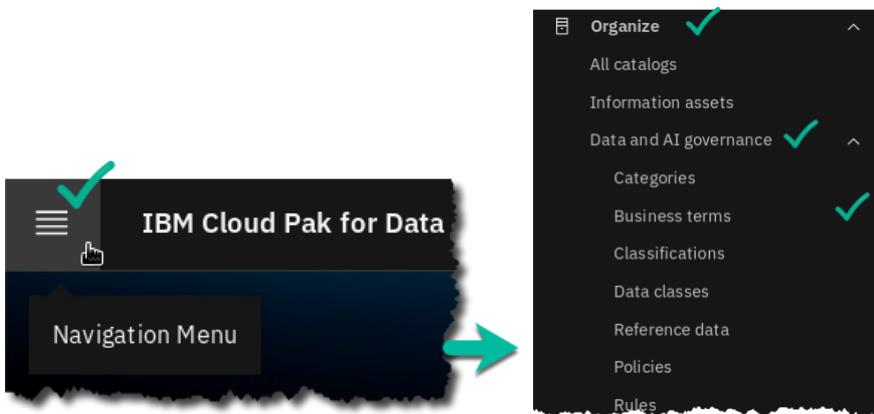
- __32. Review the Business terms and Policy for the Category.

Name	↑	Description	Type
Days Since Last Trade		Number of days since the customer executed a trade on our platform	Business term
Gender		Customer Churn Category	Business term
Home Owner		Flag indicating whether a customer owns a home	Business term
Income		Annual household income of the customer	Business term
Net Gains and Net Losses are Mutually Exclusive		A customer can only have a value in Net Gains or Net Losses	Policy

Note: You can view more than 5 items per page.



- __33. Click **Navigation Menu** \Rightarrow **Organize** \Rightarrow **Data and AI governance** \Rightarrow **Business terms**.



- __34. Here you can edit or add new Business terms that are either in Published or Draft mode.
Click on one to review it in more detail.

The screenshot shows the "Business terms" page. At the top, there are tabs for "Published" (which is selected) and "Draft". Below the tabs is a search bar with placeholder text "Find business terms" and a "Sort by: Name" dropdown. The main area displays a list of business terms:

- Days Since Last Trade** ✓
Number of days since the customer executed a trade on our platform
 Customer Churn Category
Last modified: Jun 22, 2020
- Gender
 Customer Churn Category

Below this, a detailed view of the "Days Since Last Trade" term is shown in a modal window:
Days Since Last Trade Published
Overview Related content
^ General
Description Number of days since the customer executed a trade on our platform
Primary category Customer Churn Category

Data
Steward

You can create your own Glossary with Categories and Business terms manually, or import them from a file. This workshop was prepared using .csv files in in the following directory on your OS:

```
[root@master1 organize]# pwd
/workshop/labs/organize
[root@master1 organize]# ls *.csv
-rw-rw-r--. 1 ibmdemo ibmdemo 980 Aug  4 16:39 organize-businessterms.csv
-rw-rw-r--. 1 ibmdemo ibmdemo 845 Aug  4 16:39 organize-businessterms-fixed.csv
-rw-rw-r--. 1 ibmdemo ibmdemo 264 Aug  4 16:39 organize-categories.csv
-rw-rw-r--. 1 ibmdemo ibmdemo 329 Aug  4 16:39 organize-policies.csv
-rw-rw-r--. 1 ibmdemo ibmdemo 1245 Aug  4 16:39 organize-referencedatasets.csv
-rw-rw-r--. 1 ibmdemo ibmdemo 2438 Aug  4 16:39 organize-rules.csv
[root@master1 organize]#
```

In addition, you can import Industry Models from IBM for industries such as finance, banking, healthcare, and insurance and import them into CPD.

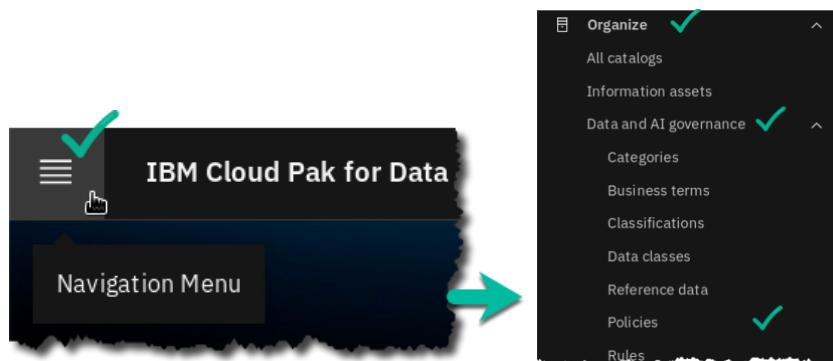
See the services screen then [Industry Accelerators](#).

4.6 Reviewing a Governance Policy and a Rule

An information governance **Policy** is a natural-language description of a governance subject area. It can contain multiple information governance sub-policies or reference one or more information Governance Rules. It must fulfill a business objective and be relevant and understandable to all users of the policy.

An information Governance **Rule** is a natural-language description of the criteria that are used to determine whether information assets are compliant with business objectives. Generally, information governance rules are derived from information governance policies and are more specific. The Rules define the actions to take in specific situations to implement the Policy. An information Governance Rule can be referenced by one or more information Governance Policies.

- 35. Click **Navigation Menu** \Rightarrow **Organize** \Rightarrow **Data and AI governance** \Rightarrow **Policies**.



Note: the first time using the page may take a minute or two to render – be please patient.

- 36. Review the Published Policy that has been associated with our previously reviewed Category.
Policies

The screenshot shows the Policies page with the following details:

- Published** (highlighted with a checkmark) and **Draft** buttons.
- Find policies** search bar.
- Sort by: Name** button.
- Data Privacy**: Company-wide data privacy policy for securing private data. Last modified: Jun 22, 2020.
- Net Gains and Net Losses are Mutually Exclusive** (highlighted with a checkmark): A customer can only have a value in Net Gains or Net Losses. Last modified: Jun 23, 2020.
- Showing 2 of 2 policies.

- __37. Scroll down until you find the Rule that is associated with this Policy.
Notice you could add another rule to this policy here. (Don't do this now.)
Review the Rule by clicking on it.

The screenshot shows a software interface for managing rules. At the top, there is a header with a 'Rules' section and a green checkmark icon. Below this is a table with columns for 'Name' and 'Description'. A single row is selected, showing the name 'Net Realized Gains and Losses Validity Check' and a green checkmark icon. To the right of the table is a large button labeled 'Add rules +'. Below the table, a detailed view of the selected rule is displayed. The title is 'Net Realized Gains and Losses Validity Check' with a 'Published' status indicator. The 'Overview' tab is selected, showing a brief description: 'If Net Realized Gains > 0 then Net Realized Losses = 0 if Net Realized Losses > 0 then Net Realized Gains = 0'. The 'Related content' tab is also visible. Under the 'General' section, there is a 'Description' field containing the same logic as the overview. The 'Primary category' is listed as 'Customer Churn Category'.

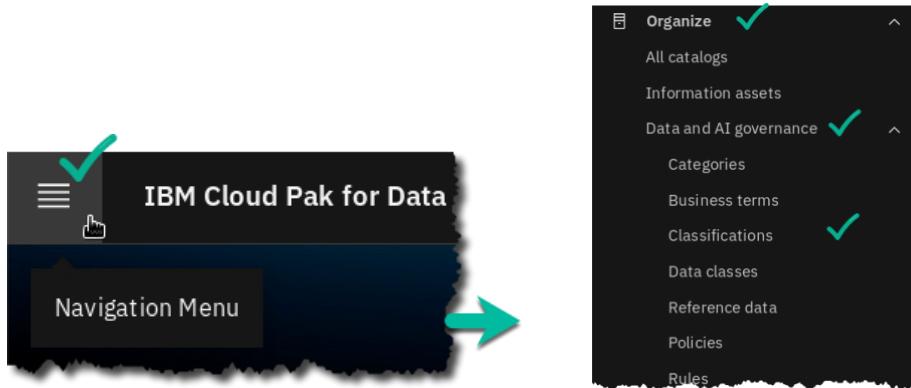
 Data Steward	A Data Dictionary contains a Business Glossary (Categories and Business terms) as well as information Governance Policies and Rules to ensure data compliance with business objectives. In this lab we have a beginning sample of these items, but in reality, a Data Dictionary for any organization is quite large and can and should be updated as frequently as new data sources, regulations, and other criteria require it.
--------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

4.7 Reviewing Classifications, Data Classes, and Reference Data

4.7.1 Classifications

A [Classification](#) describes the sensitivity level of data. In catalogs, a classification describes the sensitivity of a whole data asset to help catalog members understand the asset. You can use classifications to describe Business Terms, Data Classes, Reference Data Sets, and Governance Rules. In data protection Rules, you can include Classifications in conditions to identify the type of data to restrict.

- _38. Click [Navigation Menu](#) ⇒ [Organize](#) ⇒ [Data and AI governance](#) ⇒ [Classifications](#).



- _39. Scroll to find Classification [Confidential](#) and click on it

Classifications

This screenshot shows the 'Classifications' page. At the top, there are tabs for 'Published' (which is selected) and 'Draft'. Below that is a search bar with 'Find classifications' placeholder text, a 'Sort by: Name' dropdown, and a 'Show: All' button. The main area displays a single classification entry: 'Confidential' with a checkmark. The description below it reads: 'Confidential data is data that if compromised in some form, is likely to result in significant and/or long-term harm to the institution and/or individuals whose data it is. Access to confidential information is restricted to those who have a...'. Underneath the description is a category link '[uncategorized]' and a timestamp 'Last modified: Jun 18, 2020'.

- _40. The Classification is described here. You could also add the primary Category here, but there is no need to do so now.

This screenshot shows the details page for the 'Confidential' classification. At the top, it says 'Confidential' and 'Published'. Below that is a navigation bar with 'Overview' (selected) and 'Related content'. The 'Overview' section contains a 'General' summary and a 'Description' section. The 'Description' section includes a detailed text about confidential data and a 'Show more' link. At the bottom, there's a 'Primary category' section with a note '(uncategorized)'.

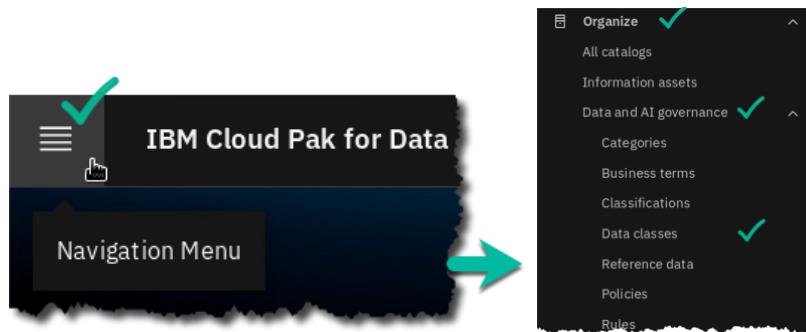
4.7.2 Data classes

Data classes describe the contents of data in a column in a relational or structured data set. Data classes are assigned to columns when profiling a structured data asset and shown on the Profile page in a Catalog or Project.

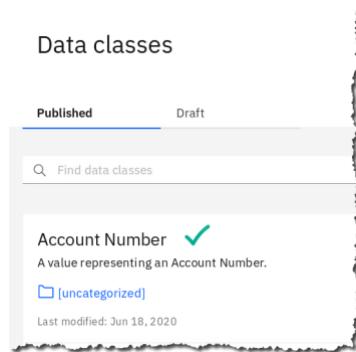
Watson Knowledge Catalog provides a predefined set of Data Classes. Some Data Classes are categorized into groups, for example:

- If you select **Date**, it also includes **Date of Birth**.
- If you select **Driver's License**, it also includes all driver licenses listed below.

__41. Click [Navigation Menu](#) ⇨ [Organize](#) ⇨ [Data and AI governance](#) ⇨ [Data classes](#).



__42. Scroll to find [Account Number](#) and click on it.



__43. The Data Class is described here.

From here you could edit the description, examples or add the primary Category.

The screenshot shows the 'Account Number' data class details page. At the top, it says 'Account Number' (Published). Below that, there are tabs for 'Overview' (selected) and 'Related content'. Under 'General', there is a 'Description' field containing 'A value representing an Account Number.' and an 'Examples' field containing '123456'. There are three green circles with edit icons overlaid on the 'Description', 'Examples', and 'Primary category' fields. The 'Primary category' field contains '[uncategorized]'.

4.7.3 Reference Data Set

[Reference Data Sets](#) define list of permissible values that are allowed for use within a data field and may be referenced by Business Terms, Policies, Rules and Data Classes in Watson Knowledge Catalog.

You can capture, manage, and socialize reference data — setting it up once and re-using the reference data in other places.

- __44. Click [Navigation Menu](#) ⇒ [Organize](#) ⇒ [Data and AI governance](#) ⇒ [Reference data](#).



- __45. Scroll to find Reference Data Set [State and Province Codes](#) and click on it.



- __46. Scroll down to review the Reference Data Set rows (the data).

The screenshot shows the data table for the 'State and Province Codes' reference data set. The table has columns for 'Code' and 'Value'. The data rows are:

	Code	Value
▼	AA	Armed Forces (the) Americas
▼	AB	Alberta
▼	AE	Armed Forces Europe
▼	AK	Alaska
▼	AL	Alabama
▼	AP	Armed Forces Pacific
▼	AR	Arkansas

- __47. Click [Related content](#) and review how you could relate Data classes and Rules to it.



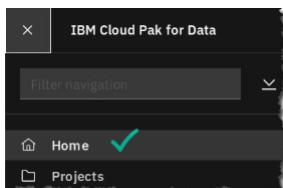
4.8 Searching for Data

The Data Scientist and Business Analyst personas may not always know what has been made available to them by the Data Engineers and Data Stewards in CPD. In fact, individual Data Engineers and Data Stewards may not always know what other users of the same persona have made available through their Collect and Organize activities.

This is remedied by the search functionality provided in CPD.

- 48. Start the search here: [Navigation Menu](#) \Rightarrow [Home](#).

(Doing this gives your search a neutral context.)



- 49. Click in the [Search](#) bar.



- 50. Type [churn](#) and hit [Enter](#).



- 51. Notice in the absence of a filter (or context) that multiple kinds of assets related to the term [churn](#) are found.

- 52. Select from [Any type – Data Asset](#).

 A screenshot of the search results page titled 'Search results for churn'. The search filters are set to 'Type' (selected), 'Any source' (selected), and 'Steward/Owner' (selected). Below the filters, a 'Data asset' button is shown with a green X. The results table shows three items:

Name	Type
CPDUSER.CUSTOMER_CHURN All projects > CPD Workshop Analytics Project	Data asset
CPDUSER.CUSTOMER_CHURN All catalogs > CPD Workshop Catalog	Data asset
customer_demo churn activity analyze.csv	

 On the left side of the search interface, there is a sidebar with a dropdown menu for 'Any type' (selected) and several other options like 'Business term', 'Category', 'Data asset' (selected), and 'machine-lear...' (partially visible).

Notice that one Data asset is in a catalog and another is in a project.

- __53. Click on the Data asset called **CPDUSER.CUSTOMER.CHURN** which is located in **All projects > CPD Workshop Analytics Project**.

The screenshot shows a search interface for 'churn'. The search filters are set to 'Type: Data asset', 'Any source', and 'Steward/Owner'. The results table has columns 'Name' and 'Type'. The first row, 'CPDUSER.CUSTOMER.CHURN' (located in 'All projects > CPD Workshop Analytics Project'), is highlighted with a red box.

Name	Type
CPDUSER.CUSTOMER.CHURN All projects > CPD Workshop Analytics Project	Data asset
CPDUSER.CUSTOMER.CHURN All catalogs > CPD Workshop Catalog	Data asset
customer_demo churn activity analyze.csv	

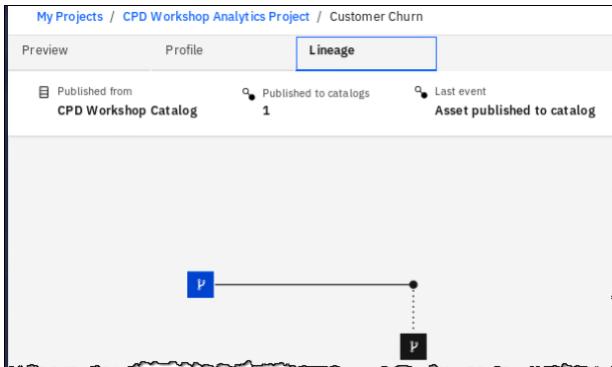
- __54. Here you can review a sample of the data.

Next, click **Lineage**.

The screenshot shows the 'Lineage' tab selected for the data asset 'CPDUSER.CUSTOMER.CHURN'. The schema is listed as '2 Columns'. The table displays two columns: 'ID' (Smallint) and 'CHURNRISK' (String). Two rows are shown: one with ID 0 and CHURNRISK 'Low', and another with ID 1 and CHURNRISK 'Low'.

ID	CHURNRISK
0	Low
1	Low

- ___55. The Lineage shown here is relatively simple for this Data asset; however, this could show a more complex lineage should this Data asset have been the result of a Data Flow Design or refine job flow output from a virtualized view of a join of two tables in two different servers.





Data Steward

The results of this search and exploration of the Data asset tells us that we want to utilize this asset for a Data Flow Design job later in this lab.

At this point, we could go to the [Overview](#) page of the [Catalog View](#) and click [Add to Project](#) for easy access later. However, this has already been done, so there is no need to do so now.



- ___56. Try the [Search](#) again, this time using the word [demographics](#).



- ___57. Choose the second Data asset named [CPDUSER.Customer_Demographics](#) which is under [All projects > CPD Workshop Analytics Project](#).

The screenshot shows a search results page titled 'Search results for demographics'. It lists six items, each with a name, type, and a small description. The second item, 'CPDUSER.CUSTOMER_DEMOGRAPHICS', is highlighted with a red box.

Name	Type
Customer Demographics Flow Job	Job
Demographics Discovery - Trade Co.	Dashboard
CPDUSER.CUSTOMER_DEMOGRAPHICS All catalogs > CPD Workshop Catalog	Data asset
CPDUSER.CUSTOMER_DEMOGRAPHICS All projects > CPD Workshop Analytics Project	Data asset
CPDUSER.CUSTOMER_DEMOGRAPHICS_flow	Data flow
CPDUSER.CUSTOMER_DEMOGRAPHICS_shaped.csv All projects > CPD Workshop Analytics Project	Data asset

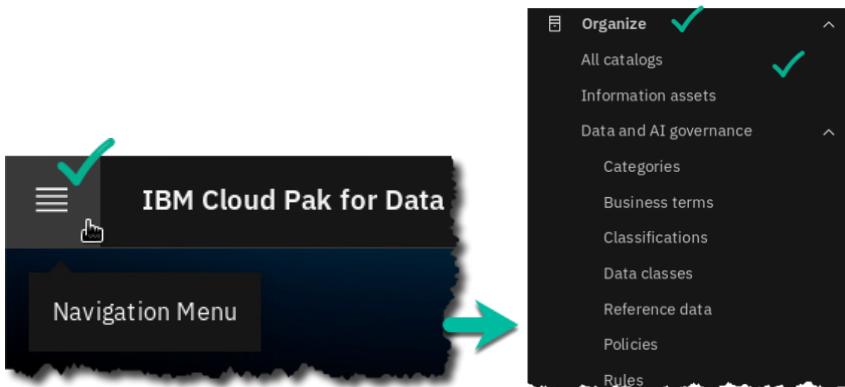
- __58. This takes us to the [CPDUSER.Customer_Demographics](#) Data asset we reviewed in an earlier exercise at the beginning of this lab.

The screenshot shows the IBM Cloud Pak for Data interface. At the top, it says "IBM Cloud Pak for Data" and "All". Below that, it shows "My Projects / CPD Workshop Analytics Project / CPDUSER.CUSTOMER_DEMOGRA...". There are three tabs: "Preview" (which is selected), "Profile", and "Lineage". Under "Preview", it says "Schema: 18 Columns | 2066 rows" and "Preview: 1000 rows". The data table has columns: ID (Smallint), GENDER (String), STATUS (String), CHILDREN (Smallint), ESTINCO... (Decimal), HOMEOW... (String), AGE (Smallint), TAXID (String). The first row of data is: 0, F, S, 1, 38000, N, 24, 14788.

 Data Steward	You will be using this data in the Data Flow Design exercise to follow. To recap: you have searched for, found and selected two tables, from which you will build a job to join and transform this data to create another permanent table.
------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

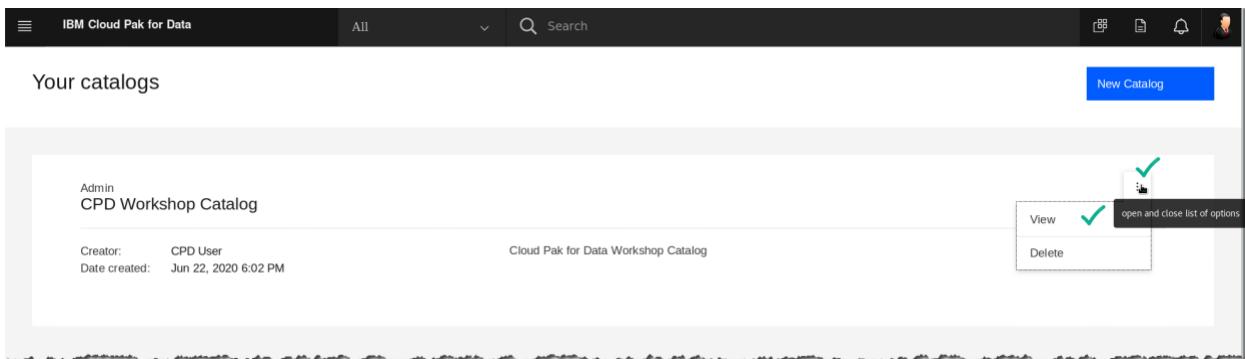
4.9 Reviewing the catalog

- _59. Click **Navigation Menu** ⇒ **Organize** ⇒ **All catalogs**.

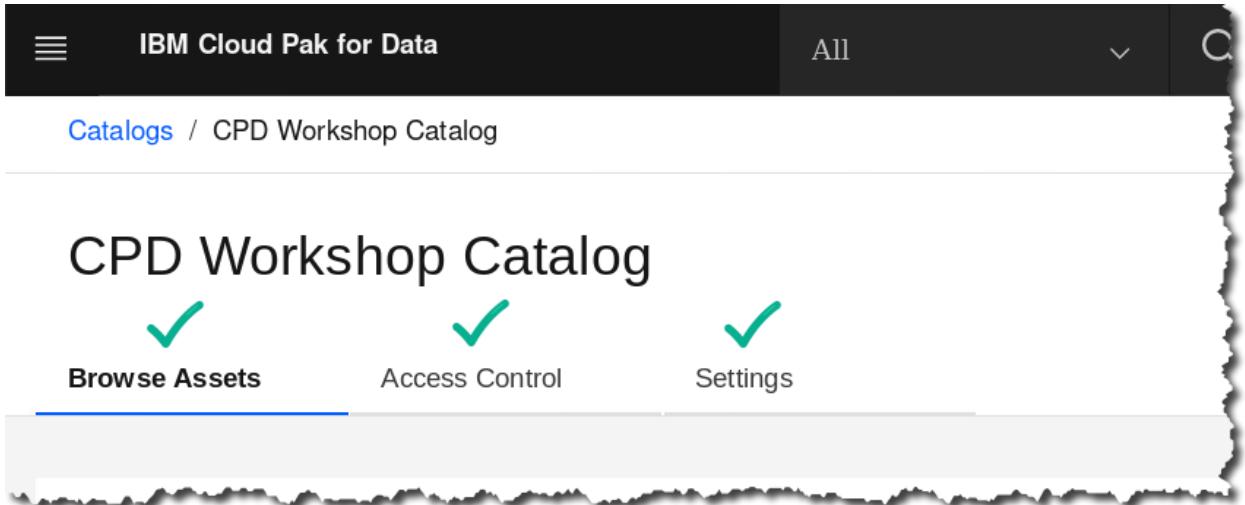


- _60. You will see a catalog named **CPD Workshop Catalog**.

Review it by clicking **ellipses** ⋮ ⇒ **View**. (Tip: You can also click on the catalog name.)



- _61. Review each of the Catalog sections: **Browse Assets**, **Access Control**, and **Settings**.



- _62. Return to **Browse Assets**.

- __63. Click the ellipses to the right of the [CPDUSER.Customer_Demographics](#) tile.
(Tip: You can either hover over the name for a flyout that shows the complete name or zoom out.)

The screenshot shows the IBM Cloud Pak for Data interface. In the top navigation bar, 'IBM Cloud Pak for Data' is selected. Below it, 'Catalogs / CPD Workshop Catalog' is shown. The main title 'CPD Workshop Catalog' is at the top. Under 'Watson Recommends', there are two data asset cards. The first card for 'CPDUS' has a tooltip with options: 'Open' (with a checkmark), 'Add to project', and 'Remove'. The second card for 'SOLUTION' has a 5-star rating. A green checkmark icon is overlaid on the screenshot near the tooltip.

- __64. Click the tab [Review](#).

The screenshot shows the 'Customer Demographics' page under the 'Review' tab. The top navigation bar says 'DATA ASSET'. The main title is 'Customer Demographics'. Below it, there are tabs: 'Overview', 'Access', 'Review' (which is highlighted with a blue border and a green checkmark icon), 'Profile', and 'Lineage'. The 'Review' tab is active.

- __65. This is where you can find review ratings and comments for the data, which is another aid in helping you find the best data for your projects.

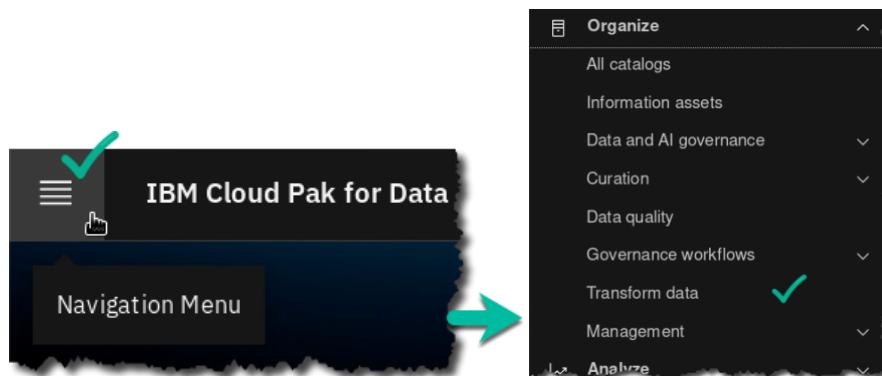
The screenshot shows the 'Customer Demographics' page under the 'Review' tab. The top navigation bar says 'DATA ASSET'. The main title is 'Customer Demographics'. Below it, there are tabs: 'Overview', 'Access', 'Review' (highlighted with a blue border), 'Profile', and 'Lineage'. The 'Review' tab is active. On the left, there's a 'Review summary' section with a 5.0 rating and 1 review. On the right, there's a 'Other Reviews' section showing one review from a 'Data Scientist' on May 30, 2020, with a 5-star rating and the comment: 'This is the best demographics data for our TradeCo clients.'

4.10 Transforming Data

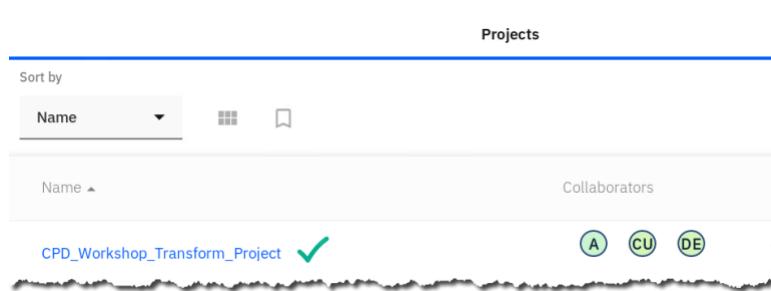
In this set of exercises, you will use the CPD built-in Data Flow Designer (DFD) to build a job that can transform your data. Note: it is the Data Engineer who would do this.

Persona (Role)	Capabilities
 Data Engineer	Data Engineers build and optimize the systems to allow data scientists and business analysts to perform their work. The Data Engineer ensures that any data is properly received, transformed, stored, and made accessible to other users.

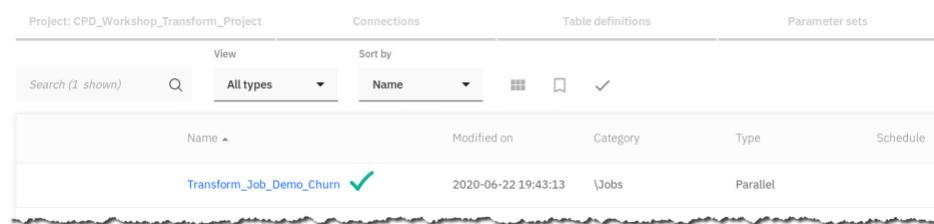
- __66. Click **Navigation Menu** ⇒ **Organize** ⇒ **Transform data**.



- __67. Select the project: [CPD_Workshop_Transform_Project](#).

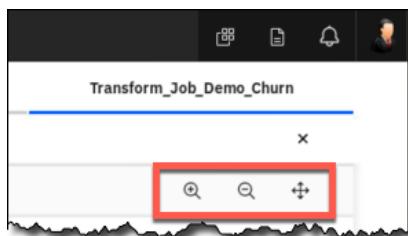


- __68. Select the job [Transform_Job_Demo_Churn](#).



 Data Engineer	<p>Note: if there is a lock icon next to this job  then click on the lock icon to unlock the job, then you can select it.</p> <p>Just make sure you are not in two different web client sessions at the same time doing this because that might be causing your locked job issue.</p>
---------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

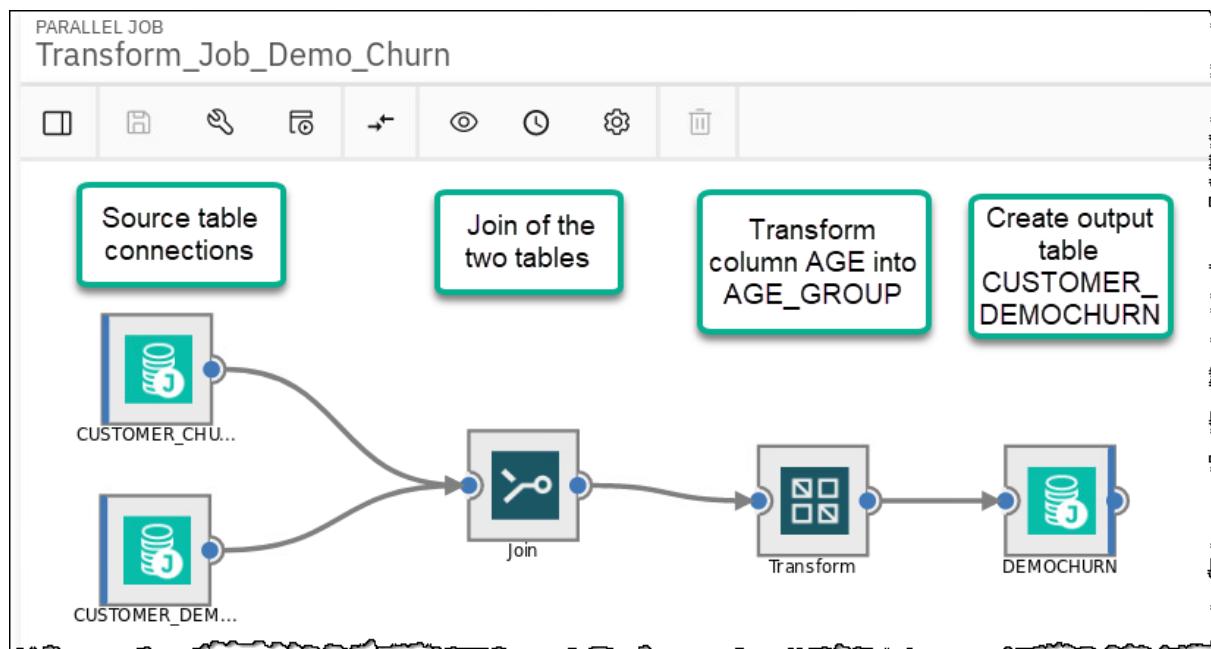
- __69. Use the zoom icons to get the best view of the entire transform job on your screen.



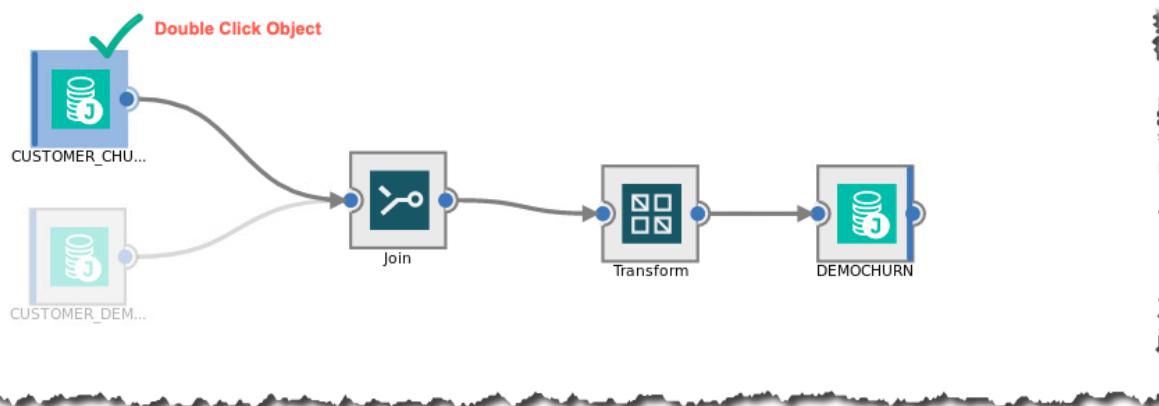
- __70. This job consists of four stages as shown below. Its purpose is to join two tables (represented by the Data assets you previously searched for) [CUSTOMER_CHURN](#) and [CUSTOMER_DEMOGRAPHICS](#).

It then adds one new column called [AGE_GROUP](#) which is derived (transformed) from the column [AGE](#).

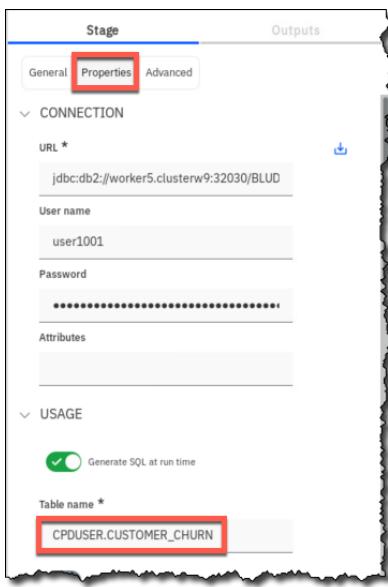
Finally, it writes out the results to a single table [CUSTOMER_DEMOCHURN](#).



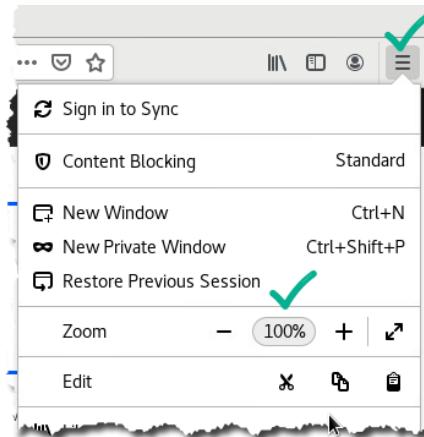
- __71. Review the connection of CUSTOMER_CHURN (the top one) by double clicking on it.



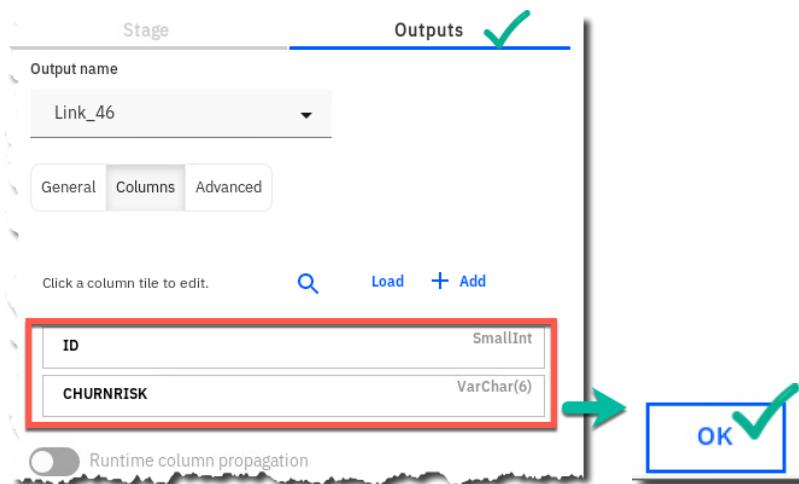
- __72. The **Properties** for this connection will appear. Review the connection and other properties.



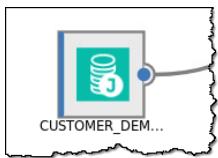
(Note: if this does not appear, resize your browser zoom as shown below)



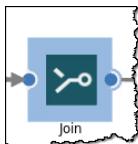
- __73. Review the tab **Outputs** to see the table columns being passed on to the next stage of the job. Then click **OK**.



__74. Do the same with the connection for **CUSTOMER_DEMOGRAPHICS** (the bottom one).



__75. Next, double click on the stage **Join**.



__76. Notice in this stage's **Properties**, the two tables are joined by the column: **ID**.

Stage Properties

JOIN KEYS

Key: ID

OPERATIONS

Join Type: Inner

OK

__77. In this Join stage, select section: **Outputs**.

Notice that the columns now include both tables.

Join

Outputs

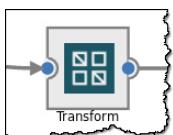
Output name: Link_52

Click a column tile to edit.

ID
CHURNRISK
GENDER
STATUS
CHILDREN
ESTINCOME
HOMEOWNER
AGE
TAXID
CREDITCARD
DOB
ADDRESS_1

OK

- __78. Double click stage [Transform](#).



- __79. Select the section [Inputs](#). These are simply the inputs from the Join stage.

Transform - Transformer stage

Properties		Inputs		Outputs		
Input name	Link_52	Jump to	Columns			
<input style="width: 100%;" type="text" value="Search input columns (19...)"/> Q				<input checked="" type="checkbox"/> Description		
Column name	Key	SQL type	Extended	Length	Scale	Nullable
ID	false	SMALLINT	true	0	0	true
CHURNRISK	false	VARCHAR	false	6	0	true
GENDER	false	VARCHAR	false	1	0	true

- __80. Review the section [Outputs](#). Notice each column from Inputs is given a coded “link” prefix before the column name.

Transform - Transformer stage

Properties		Inputs		Outputs		
Output name	Link_57	Jump to	Columns			
Constraint		Abort after rows	0			
<input checked="" type="checkbox"/> Otherwise/Log				<input type="checkbox"/> Description Load + Add 		
<input type="checkbox"/> Derivation		Column name	Key	SQL type	Extended	Length
<input type="checkbox"/> Link_52.ID		ID	false	SMALLINT	true	0
<input type="checkbox"/> Link_52.CHURNRISK		CHURNRISK	false	VARCHAR	false	6

- __81. Scroll down the list of output columns (using both scroll bars) to find the last output column, which is a derived (transformed) column called **AGE_GROUP**.

Notice it is assigned a SQL type (data type) or VARCHAR, Length 11.

Click on this derived column to review its formula.

	Derivation	Column name	Key	SQL type	Extended	Length	Scale	Nullable
<input type="checkbox"/>	Link_52.STATE	STATE	false	CHAR	false	2	0	true
<input type="checkbox"/>	Link_52.ZIP	ZIP	false	VARCHAR	false	5	0	true
<input type="checkbox"/>	Link_52.ZIP4	ZIP4	false	VARCHAR	false	4	0	true
<input type="checkbox"/>	Link_52.LONGITUDE	LONGITUDE	false	DECIMAL		9	6	true
<input type="checkbox"/>	Link_52.LATITUDE	LATITUDE	false	DECIMAL		9	6	true
<input checked="" type="checkbox"/>	If Link_52.AGE < 18 THEN "Child" ELSE IF Link_52.AGE < 30 THEN "Young adult" ELSE IF Link_52.AGE < 65 THEN "Adult" ELSE "Senior"	AGE_GROUP	false	VARCHAR	false	11	0	false

Runtime column propagation

Cancel OK

- __82. The Derivation Builder screen shows how it transforms AGE into AGE_GROUP.

Derivation Builder - AGE_GROUP VARCHAR(11)

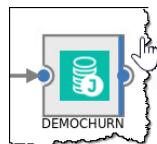
Derivation

```
If Link_52.AGE < 18 THEN "Child" ELSE IF Link_52.AGE < 30 THEN "Young adult" ELSE IF Link_52.AGE < 65 THEN "Adult" ELSE "Senior"
```

- __83. Click **Cancel** and **Cancel** again to ensure you have not inadvertently changed something.



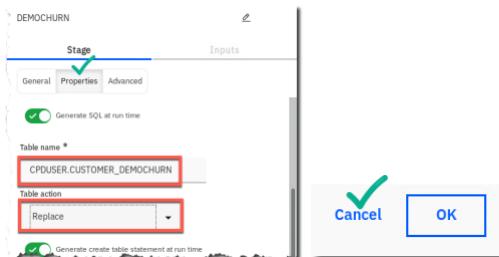
- __84. Double click the last stage, DEMOCHURN



__85. Scroll down in the Properties to find **Table name** and **Table action**.

This indicates that `CPDUSER.CUSTOMER_DEMOCHURN` will be replaced with each run of this job.

Close without changing anything in this stage by clicking **Cancel**.

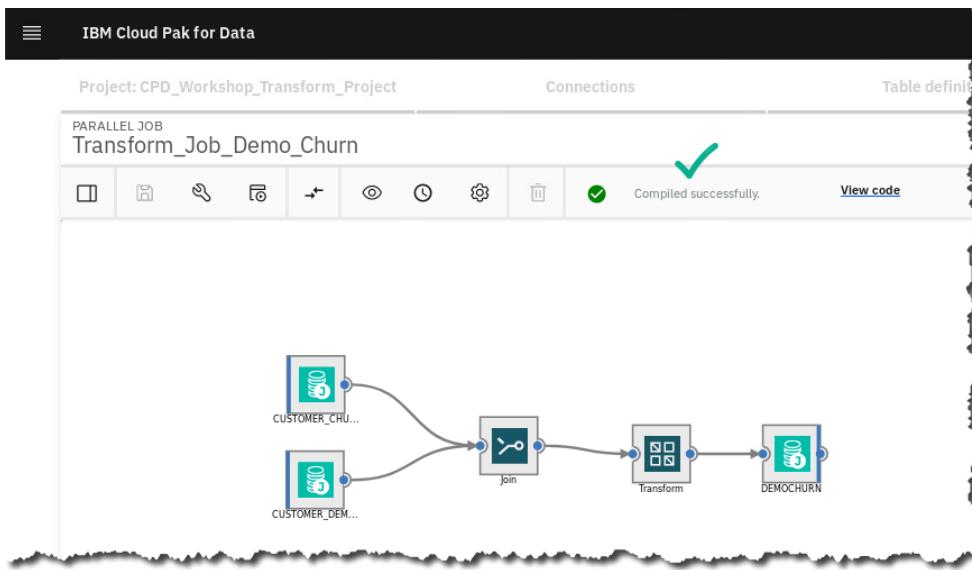


__86. Click on the **Save** then **Compile (wrench)** icon.



 Data Engineer	<p>The compile action does the following:</p> <ul style="list-style-type: none"> <i>Primary Input.</i> If you have more than one input link to a Transformer stage, the compiler checks that one is defined as the primary input link. <i>Reference Input.</i> If you have reference inputs defined in a Transformer stage, the compiler checks that these are not from sequential files. <i>Key Expressions.</i> If you have key fields specified in your column definitions, the compiler checks that there are key expressions joining the data tables. <i>Transforms.</i> If you have specified a transform, the compiler checks that this is a suitable transform for the data type.
---------------------------------------------------------------------------------------------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

__87. The job should complete successfully.

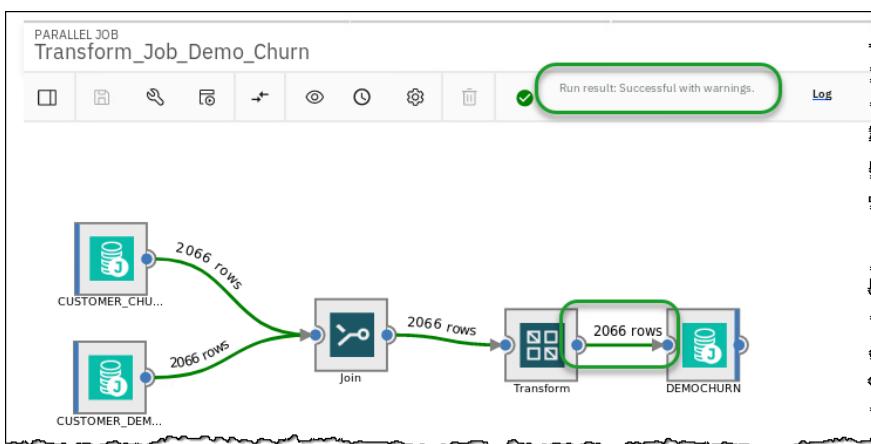


__88. Click the **Run** icon. When the Job run options box appears, select **Run**.



__89. Watch the job run...

When it completes, it should look like this:



Data
Engineer

- You can select the Log to view any warnings or errors during this process. This is helpful when troubleshooting. In this case, you notice a warning that the script attempted to drop a table that didn't exist. No need to worry.

__90. Review the data by opening the Data Server Manager (DSM) for the Db2 Warehouse.

Navigation Menu \Rightarrow Collect \Rightarrow My Data \Rightarrow Databases \Rightarrow Db2 Warehouse \Rightarrow ellipsis \vdots
 \Rightarrow Open database

The screenshot shows the 'My data' section of the IBM Cloud Pak for Data interface. The 'Databases' tab is selected. A search bar and filter dropdown are at the top. Below, two databases are listed: 'Db2 Advanced Edition' and 'MongoDB-1'. A context menu is open over the 'Db2 Advanced Edition' entry, with the 'Open database' option highlighted. Other menu items include Details, Configure, Submit connection for approval, Manage access, and Delete.

__91. At the top left of the Db2 DSM console, click on: Summary \Rightarrow Explore \Rightarrow Tables.

The screenshot shows the Db2 DSM interface. The top navigation bar shows 'My data: Databases / Db2 Advanced Edition'. The main menu on the left has three options: 'Summary' (marked with a green circle 1), 'Explore' (marked with a green circle 2, which is highlighted with a blue border), and 'Tables' (marked with a green circle 3). The 'Explore' option is currently active.

_92. Select schema CPDUSER \Rightarrow table CUSTOMER_DEMOCHURN \Rightarrow View Data.

The screenshot shows the IBM Data Studio interface. On the left, under 'Schemas', the 'CPDUSER' schema is selected (indicated by a green checkmark). In the center, under 'Tables', the 'CUSTOMER_DEMOCHURN' table is selected (also indicated by a green checkmark). On the right, the 'Table definition' pane is open, showing the structure of the CUSTOMER_DEMOCHURN table with columns like ID, CHURNRISK, GENDER, etc. A blue button labeled 'View data' with a green checkmark is visible at the bottom right of the table definition pane.

Name	Type	Tables	Name	Schema	Properties	Table definition
USER1001	User	218	CUSTOMER_CHURN	CPDUSER		CUSTOMER_DEMOCHURN
<input checked="" type="checkbox"/> CPDUSER	User	3	<input checked="" type="checkbox"/> CUSTOMER_DEMOCHURN	CPDUSER		
SOLUTIONS	User	3	<input type="checkbox"/> CUSTOMER_DEMOCHURN	CUSTOMER_DEMOCHURN		

Total: 3, selected: 1 Total: 3, selected: 1 View data

_93. The new column AGE_GROUP contains the derived data, and the output is also the join of the two tables in this schema.

CPDUSER.CUSTOMER_DEMOCHURN			
AGE_GROUP	ID	CHURNRISK	GENDER
Adult	7	High	M
Adult	141	Low	F
Adult	138	High	M

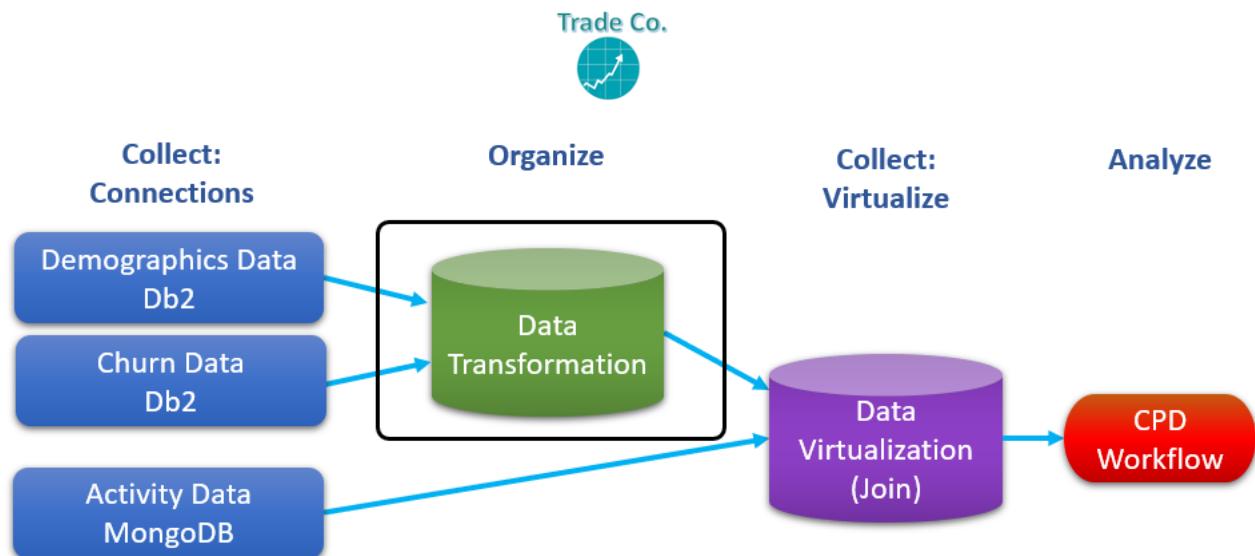
Data Engineer

If you want to try your hand at building a transformation job from scratch, please see the Organize Deeper Dive lab, section “Transformation Data – Creating a job.”

4.11 Lab conclusion

We have seen the value in creating a [Data Dictionary](#) by creating a [glossary of categories](#) and [terms, policies](#) and [rules](#), to make data searchable so that data scientist, data engineers and business analysts can “shop for data.”

This lab showed you how data can be [profiled, visualized, refined, searched](#) and [transformed](#).



The steps covered here could normally take many weeks, months, and sometimes even years, to complete using traditional manual methods. Cloud Pak for Data automates these things so that you can accelerate the time to value of your organization’s analytics projects.

** End of Lab 04 - Organize

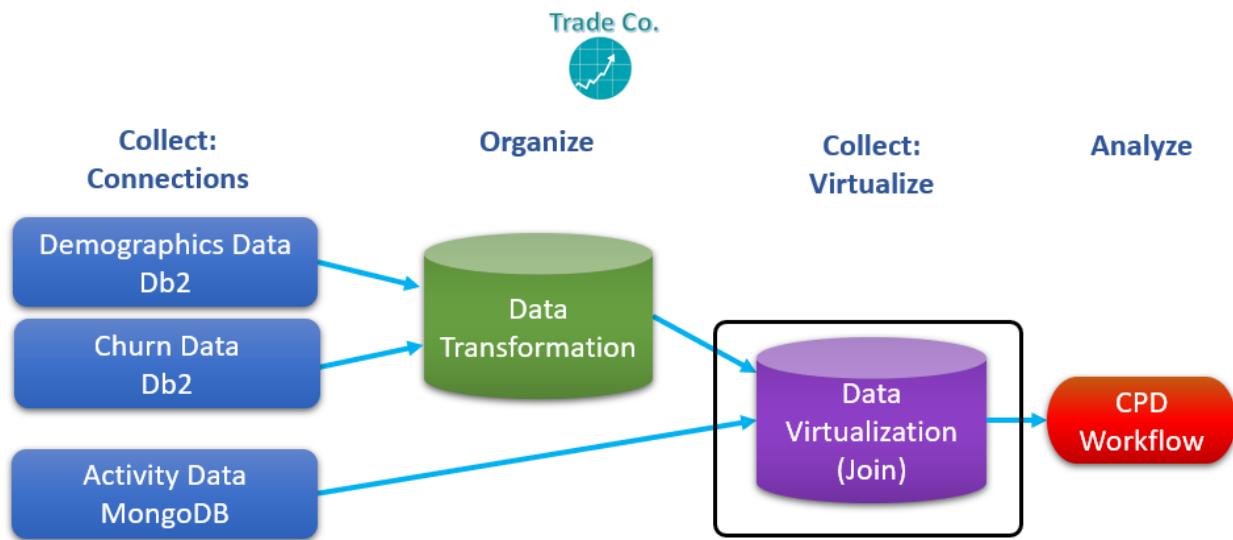
Lab by Burt Vialpando, Kent Rubin and John Van Buren, IBM

Lab 05 COLLECT: VIRTUALIZE

5.1 Lab overview

In this lab, you will learn about [Data Virtualization](#) to complete the [Collect](#) tasks by creating a virtualized view of the transformed Db2 Demographics and Churn data, joined with the MongoDB Activity data.

The team from Trade Co. wants to use data virtualization to easily join the disparate data sources into one view for easier consumption in their analytics work.



5.2 Persona represented in this lab

The [Data Engineer](#) persona will be performing the various [Collect](#) tasks shown in this lab.

Persona (Role)	Capabilities
 Data Engineer	Data Engineers build and optimize the systems to allow data scientists and business analysts to perform their work. The Data Engineer ensures that any data is properly received, transformed, stored, and made accessible to other users.

5.3 Logging into the CPD web client (if you have not already done so)

- 1. If you are starting this lab stand-alone (without going through previous labs) do the following:
- 2. Double-click the desktop icon: [Cloud Pak for Data Web Client](#).



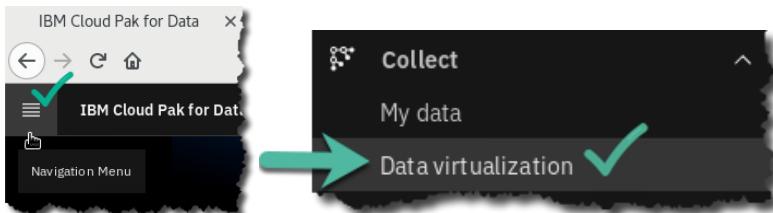
- 3. The CPD web client GUI displays as shown. Use `cpduser` and `cpdaccess` for the *Username* and *Password* and click [Sign in](#).



5.4 Remove older Data Virtualization sources

This step is needed to remove older data sources that are no longer required for this lab.

- 4. Start at the [Navigation Menu](#) ("hamburger" icon) ⇒ [Collect](#) ⇒ [Data virtualization](#).



- 5. In the Data sources screen, click the [ellipse](#) on the MySQL Community Edition and [Remove](#)

The screenshot shows the "Data sources" screen in the IBM Cloud Pak for Data Workshop. It lists three data sources in a table view:

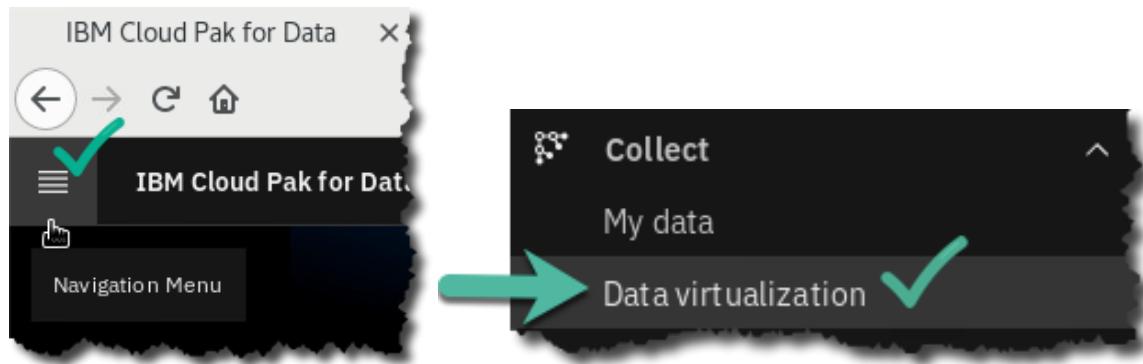
Hostname: Port	Database	Type	Username
sl-us-south-1-portal.55.dblayer.com: 21568	STOCKS	MySQL Community Edition	admin
db2w-naxmdbf.us-south.db2w.cloud.ibm.com: 50000	BLUDB	Db2 Family	dvuser
db2w-dplyzso.us-south.db2w.cloud.ibm.com: 50000	BLUDB	Db2 Family	dvuser

A green arrow points to the "Edit connection" and "Remove" buttons for the first row. A message at the bottom says "Deleted the data source successfully."

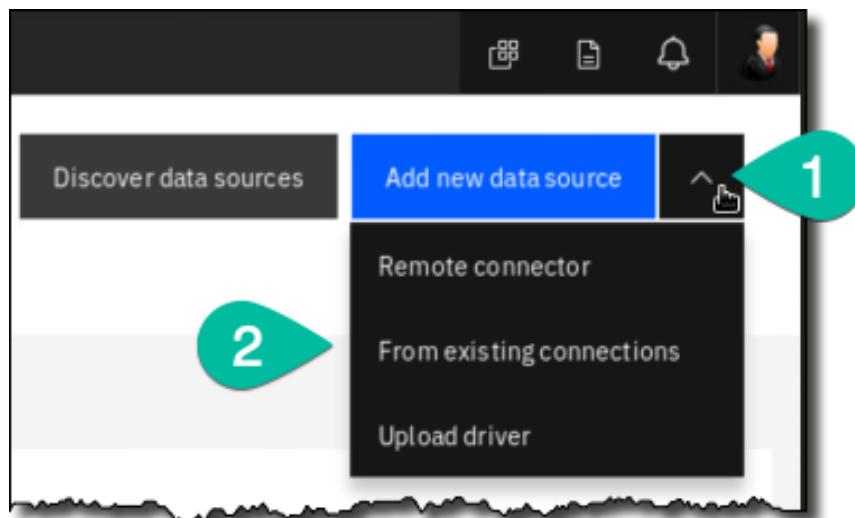
5.5 Adding Data Virtualization data sources

Let's explore the Data Virtualization process by adding new data sources to virtualize.

- _6. Start at the [Navigation Menu](#) ("hamburger" icon) \Rightarrow [Collect](#) \Rightarrow [Data virtualization](#).



- _7. In the Data sources screen, click the [Arrow](#) next to 'Add new data source' \Rightarrow [From existing connections](#).



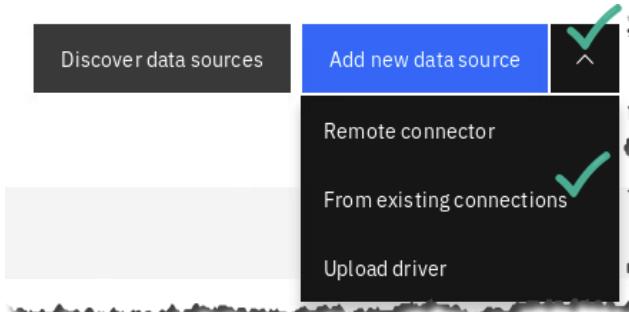
__8. Select (click on) Db2 Advanced Edition \Rightarrow Next.



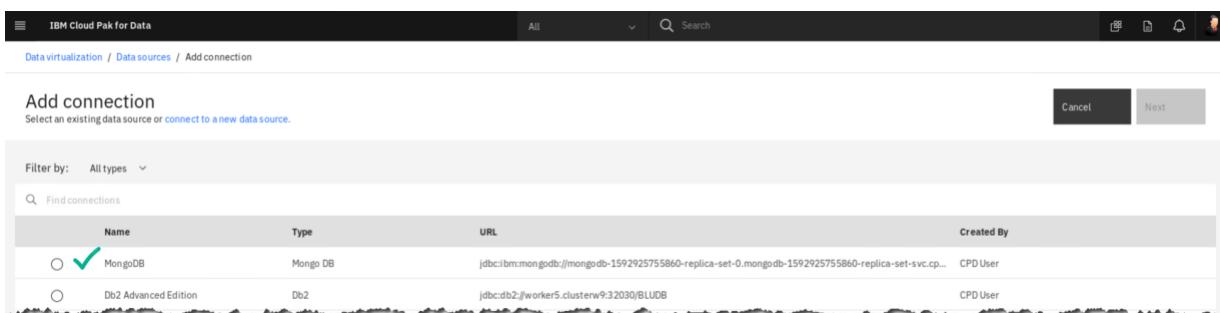
__9. You will notice a new Data source hostname created by Username user1001 (port 32030)



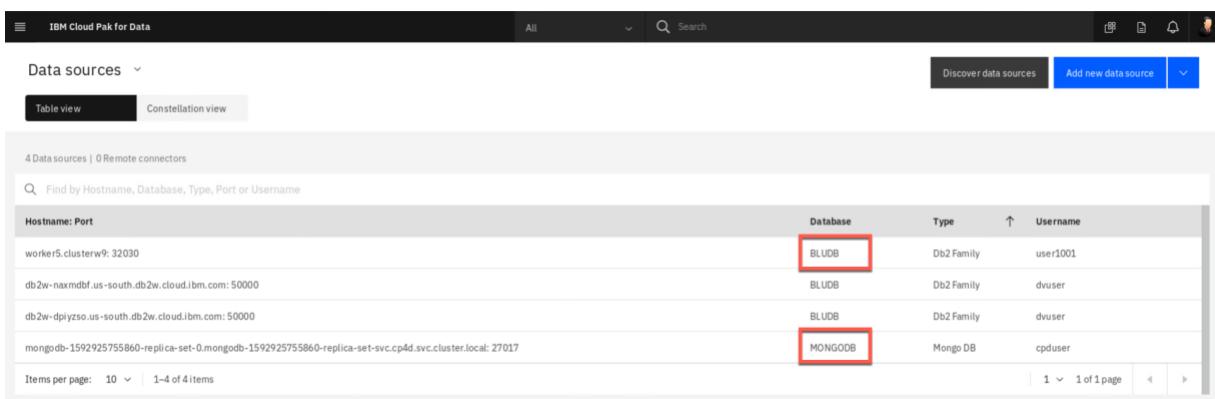
__10. Repeat this again: Add new data source \Rightarrow From existing connections.



__11. Select (click on) MongoDB \Rightarrow Next.



__12. You should now see the following Data virtualization Data sources.



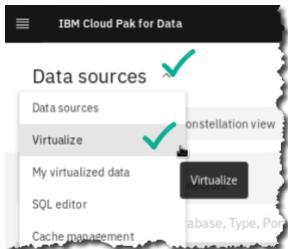
Data
Engineer

In this exercise, you used previously created connections to create Data sources in the Data virtualization instance. You could have created the connection dynamically on this screen as well.

5.6 Virtualizing the Db2 data

In the next exercise, you will virtualize the Db2 table [CUSTOMER_DEMOCHURN](#).

- __13. At the top left of the screen, click [Data sources](#) \Rightarrow [Virtualize](#).



NOTE: If you experience an error message after a minute or two and see this message:

The screenshot shows the 'Tables' tab selected in the navigation bar. A prominent red error message box at the top states: 'No table found or discovered, please try it later.' Below this, the search bar shows 'Asset type: All' and a placeholder 'Find tables by name, schema, column, or business term'. The main area displays 'Available tables: 0' and three filter buttons: 'Table', 'Business terms', and 'Schema'.

Simply wait 5 minutes then refresh the browser and the tables should appear.

- __14. Under [Databases](#), check box [Db2 Family](#).

Fill in the [Search](#) bar with [CUSTOMER_DEMOCHURN](#).

The screenshot shows the 'Tables' tab selected. In the search bar, the text 'CUSTOMER_DEMOCHURN' is entered and highlighted with a red box.

- __15. Check the able [CPDUSER.CUSTOMER_DEMOCHURN](#) from the [CPDUSER](#) schema.

If. You didn't do the Organize/Transform lab, choose [SOLUTIONS.CUSTOMER_DEMOCHURN](#).

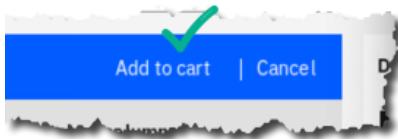
The screenshot shows the 'Virtualize' interface with the 'Tables' tab selected. A search bar at the top contains the text 'customer_demochurn'. Below the search bar, a blue header bar indicates '1 item selected'. The main table lists three items:

Table	Business terms	Schema	Database
<input checked="" type="checkbox"/> CUSTOMER_DEMOCHURN	-	CPDUSER	BLUDB
<input type="checkbox"/> CUSTOMER_DEMOCHURN	-	SOLUTIONS	BLUDB
<input type="checkbox"/> CUSTOMER_DEMOCHURN	-	SOLUTIONS	BLUDB

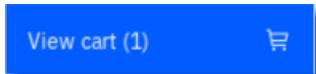
At the bottom left, there is a 'Items per page:' dropdown set to 10, and a status message '1-3 of 3 items'.

Note: the above tables are identical, but the first one was only created if you did the Organize / Transform lab preceding this lab. The pre-created table is useful for jumping right into this lab. Also, use your browser functionality to zoom out if you cannot see the full table names clearly.

- __16. Click [Add to Cart](#).



- __17. Click [View cart](#).



- __18. Select *Project CPD Workshop Analytics Project*.

Click [Virtualize](#)

The screenshot shows the 'Virtualize' interface with the 'Review cart and virtualize tables' step selected. The 'Assign to' dropdown is set to 'CPD Workshop Analytics Project', which is highlighted with a red box. The 'Virtualize' button at the top right is also highlighted with a green checkmark.

Below the dropdown, there is a table titled 'Objects to be virtualized' showing the details of the selected table:

Table	Schema	Source schema	Databases/File Path	Hostname:Port	Grouped tables
CUSTOMER_DEMOCHURN	USER1001	SOLUTIONS	BLUDB	worker5.clusterw9: 32030	1

- __19. You have now just created a virtualized table and placed it in a project.

Click [Virtualize more data](#).

Virtual tables created X

1 of 1 tables successfully virtualized.

Table	Schema	Status
CUSTOMER_DEMOCHURN	USER1001	✓ success

Assigned to project

CPD Workshop Analytics Project

✓

[Virtualize more data](#) [View my virtualized data](#) [Go to project](#)

5.7 Virtualizing the MongoDB data

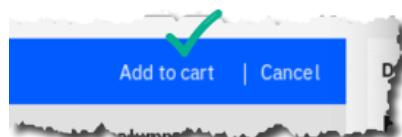
You should be at the **Virtualize** screen again, after clicking **Virtualize more data** from the previous step.

If you are not there, you can always get there by doing the following: **Navigation menu** \Rightarrow **Collect** \Rightarrow **Data virtualization** \Rightarrow **Menu** \Rightarrow **Virtualize**.

20. In the *Virtualize* screen Click the *Databases* box **Mongo DB**, select **ACTIVITY01**.

The screenshot shows the 'Virtualize' screen in the IBM Cloud Pak for Data interface. On the left, there's a search bar and a dropdown for 'Asset type: All'. Below it, there are tabs for 'Tables' and 'Files', with 'Tables' selected. A search bar below the tabs says 'Find tables by name, schema, column, or business term'. In the main area, a table lists selected items. One row for 'ACTIVITY01' is highlighted with a green circle containing the number '2'. To the right of the table is a sidebar titled 'Databases' with a list of available databases, one of which is 'Mongo DB (1)' with a green circle containing the number '1' next to it. At the bottom right of the main area is a 'View cart (0)' button.

21. Click **Add to cart**.



22. Click **View cart**.



23. Select *Project CPD Workshop Analytics Project*.

Click **Virtualize**.

The screenshot shows the 'Review cart and virtualize tables' screen. At the top, there's a message 'Assign to (all tables will be assigned to the same project)'. Below it, there are three radio buttons: 'Data request' (unchecked), 'Project' (checked), and 'My virtualized data' (unchecked). A dropdown menu shows 'CPD Workshop Analytics Project'. At the bottom, there's a table titled 'Objects to be virtualized' with one row for 'ACTIVITY01'. The table columns are: Table, Schema, Source schema, Databases/File Path, Hostname: Port, and Grouped tables. The 'Grouped tables' column shows the value '1'. There are also buttons for 'Submit to catalog', 'Empty cart', and a green checkmark icon next to the 'Virtualize' button.

__24. You have now created a virtualized table of Mongo data in a project.

Click [View my virtualized data](#).

Virtual tables created ×

1 of 1 tables successfully virtualized.

Table	Schema	Status
ACTIVITY01	USER1001	success

Assigned to project

CPD Workshop Analytics Project

[Virtualize more data](#) [View my virtualized data](#) [Go to project](#)

5.8 Joining the virtualized tables

You should be at the [My Virtualize data](#) screen, after clicking [View my virtualized data](#) in the previous exercise.

If you are not, you can always navigate there by doing the following: [Navigation menu](#) \Rightarrow [Collect](#) \Rightarrow [Data virtualization](#) \Rightarrow [Menu](#) \Rightarrow [My Virtualized data](#).

- 25. Before joining any of the virtualized tables you have created, preview each of them to make sure they are working properly.

For each table, click on [ellipses](#) \vdots \Rightarrow [Preview](#).

Table	Schema	Created on
ACTIVITY01	USER1001	Jun 23, 2020 11:16:18 PM
CUSTOMER_DEMOCHURN	USER1001	Jun 23, 2020 11:10:44 PM

PERCENTCHANGECALCULATION	LARGESTSINGLETRANSACTION	TOTALUNITSTRADED	NETREALIZEDGAINS_YTD
51.5	29877.99	206	2987.799
3	13246.325	12	1324.6325
44.5	11196.12	178	0
4.2	92.51	28	0
5.5	5770.575	22	577.0575
14	6660.665	56	0
21.75	6595.66	87	659.566

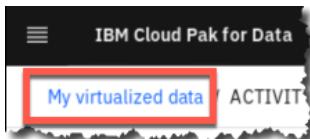


Data
Engineer

Check for data in ACTIVITY01 and CUSTOMER_DEMOCHURN.

If any of these virtualized tables are not displaying data, you can fix it by performing a Remove action on the virtualized table (use the ellipses \vdots) and then re-create it.

_26. Return to My virtualized data by selecting breadcrumb back.



_27. Choose tables **ACTIVITY01** and **CUSTOMER_DEMOCHURN**, then click **Join**.

Table	Schema	Created on
ACTIVITY01	USER1001	Jun 23, 2020 11:16:18 PM
CUSTOMER_DEMOCHURN	USER1001	Jun 23, 2020 11:10:44 PM

_28. Scroll down on the table **CUSTOMER_DEMOCHURN** until you can see the Column **ID**.

Join the virtualized tables by clicking on Column **ID** on the table **ACTIVITY01** and dragging the cursor to the Column **ID** on the table **CUSTOMER_DEMOCHURN**.

If done correctly, you will see blue key icons next to each column in the virtualized tables.

Column name	Data type
PERCENTCHANGECALCULATION	DOUBLE
LARGESTSINGLETRANSACTION	DOUBLE
TOTALUNITSTRADED	INTEGER
NETREALIZEDGAINS_YTD	DOUBLE
TOTALDOLLARVALUETRADED	DOUBLE
DAYSSINCIELASTLOGIN	INTEGER
NETREALIZEDLOSSES_YTD	DOUBLE
_ID	VARCHAR
DAYSSINCIELASTTRADE	INTEGER
ID	INTEGER
SMALLESTSINGLETRANSACTION	DOUBLE

Column name	Data type
AGE_GROUP	VARCHAR
ID	SMALLINT
CHURNRISK	VARCHAR
GENDER	VARCHAR
STATUS	VARCHAR
CHILDREN	SMALLINT
ESTINCOME	DECIMAL
HOMEOWNER	VARCHAR
AGE	SMALLINT
TAXID	VARCHAR
CREDITCARD	CHAR
DOB	DATE
ADDRESS_1	VARCHAR
ADDRESS_2	VARCHAR

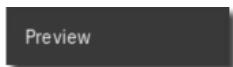
__29. On table **ACTIVITY01** find the Column **_ID**.

This column is different from the join column you just used called **ID**. The **_ID** field is a Mongo generated column that is not useful for our Analytics processing.

Click on the box next to the Column Name to uncheck Column **_ID**. This will indicate that you do not want this column included in the joined virtualized view of the two virtualized tables. You may have noticed that the arrow between ACTIVITY01 and CUSTOMER_DEMOCHURN disappeared. That's ok. The key still indicates a join relationship.

Table 1: ACTIVITY01		
<input type="checkbox"/>	Column name	Data type
<input checked="" type="checkbox"/>	PERCENTCHANGECALCULATION	DOUBLE
<input checked="" type="checkbox"/>	LARGESTSINGLETRANSACTION	DOUBLE
<input checked="" type="checkbox"/>	TOTALUNITSTRADED	INTEGER
<input checked="" type="checkbox"/>	NETREALIZEDGAINS_YTD	DOUBLE
<input checked="" type="checkbox"/>	TOTALDOLLARVALUETRADED	DOUBLE
<input checked="" type="checkbox"/>	DAYSSINCELASTLOGIN	INTEGER
<input checked="" type="checkbox"/>	NETREALIZEDLOSSES_YTD	DOUBLE
<input checked="" type="checkbox"/>	_ID	VARCHAR
<input checked="" type="checkbox"/>	DAYSSINCELASTTRADE	INTEGER
<input checked="" type="checkbox"/>	ID	INTEGER
<input checked="" type="checkbox"/>	SMALLESTSINGLETRANSACTION	DOUBLE

__30. Click **Preview**.

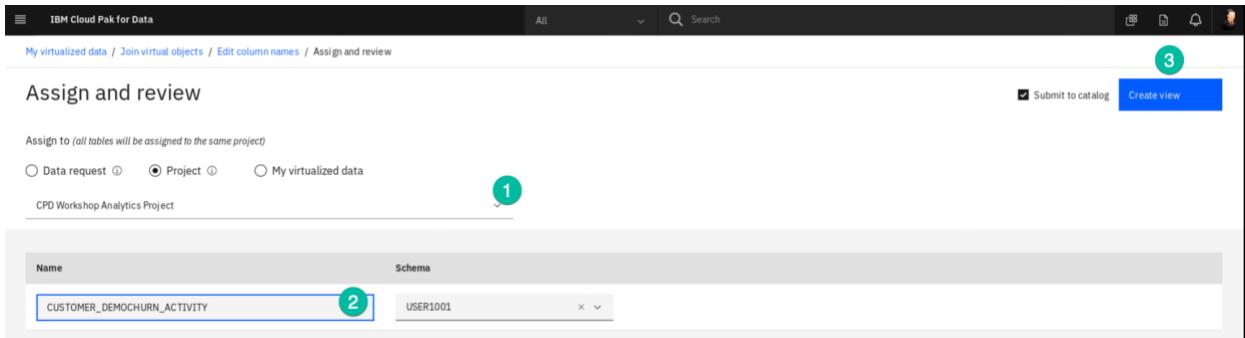


__31. Scroll to the right to make sure data columns from both tables are represented.

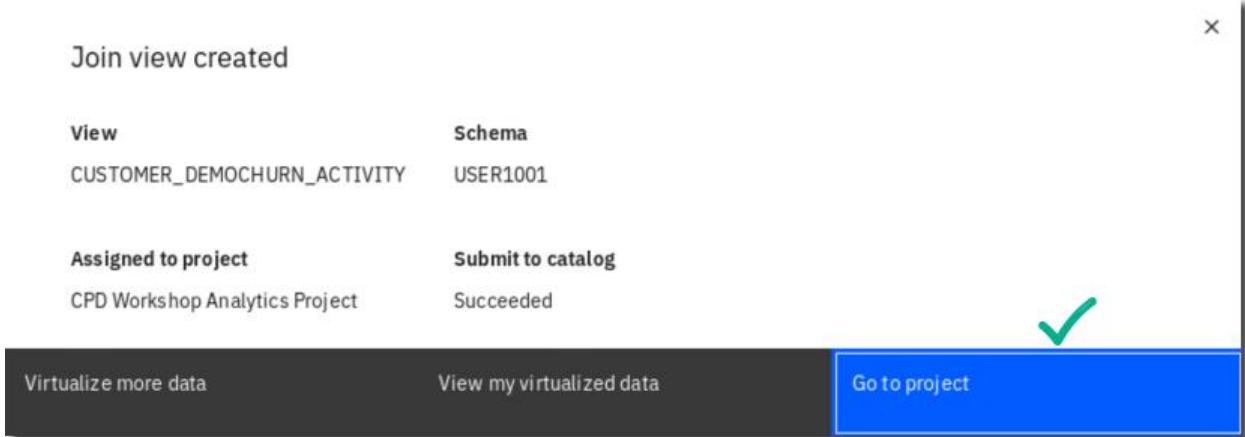
Click **x** to close.

New join preview				
	ID	SMALLESTSINGLETRANSACTION	AGE_GROUP	CHURNRISK
	0	2987.799	Young adult	Low
	1	1489.149	Adult	Low
	2	1240.624	Adult	Low
	3	1306.6305	Adult	High
	4	125.7625	Young adult	High
	5	622.5625	Young adult	High

- __32. Click **Next**. Now you can Edit column names if you like.
- __33. Click **Next again**.
- __34. In the Assign and review screen, select Project **CPD Workshop Analytics Project** (1), then type view name: **CUSTOMER_DEMOCHURN_ACTIVITY** (2) and Click **Create View** (3)



- __35. Click **Go to project**.



__36. You will be navigated to My Projects \Rightarrow CPD Workshop Analytics Project \Rightarrow Assets.

Find and click on your new virtualized view (of the two virtual tables)
USER1001.CUSTOMER_DEMOCHURN_ACTIVITY.

Name	Type	Created by	Last modified
USER1001.CUSTOMER_DEMOCHURN_ACTIVITY	Data Asset	CPD User	Jul 22, 2020, 11:59 AM

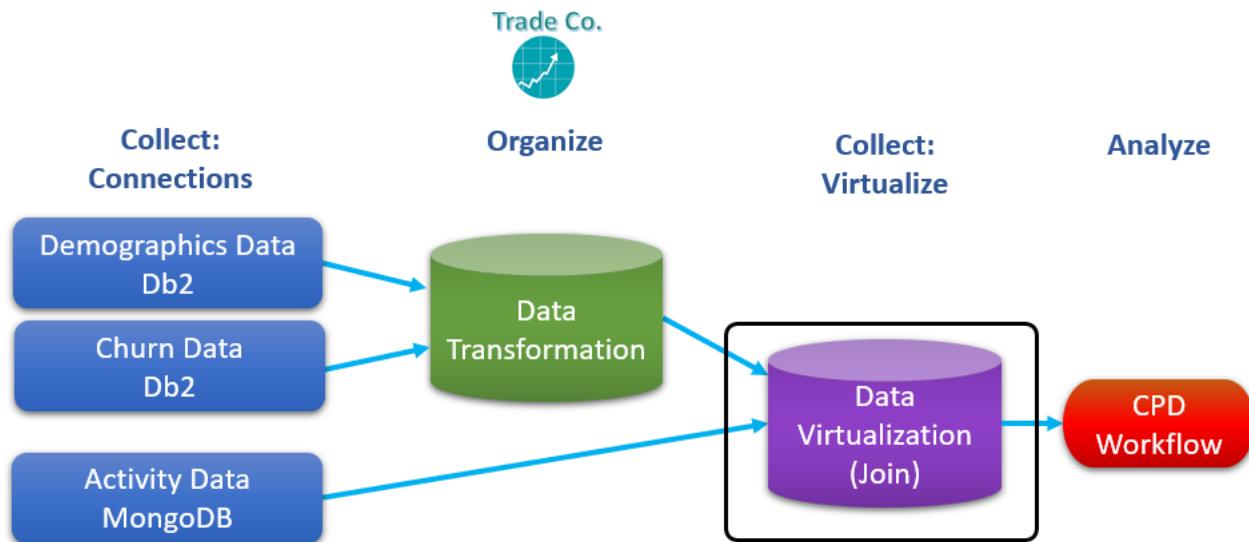
__37. Scroll through the data in the preview mode.

Notice at this point you could refine the data even further. (We will do that in another exercise.)

PERCENTCHANGECALCULATED	LARGESTSINGLETRANSACTION	TOTALUNITSSTRAIGHT	NETREALIZEDGAINS	TOTALDOLLARVALUEUTR	DAYSSINCELASTL	NETREALIZEDLOSSES	DAYSSINCELASTY	ID	SMALLEST
51.5	29877.99	206	2987.799	59755.98	3	0	10	0	2987.799
11.25	14891.49	45	1489.149	29782.98	3	0	9	1	14891.49
5.5	12406.24	22	1240.624	24812.48	1	0	9	2	12406.24
8	13066.309	32	0	26132.61	3	1306.6309	5	3	13066.309
3.45	1257.625	23	0	5030.5	2	251.525	19	4	1257.625
11.5	6225.625	46	0	12451.25	2	622.5625	8	5	6225.625
8	13261.325	32	0	26522.65	4	1326.1325	10	6	13261.325
3.25	10111.01	13	0	20222.02	3	1011.101	10	7	10111.01
5.7	2345.235	38	0	9380.94	2	234.5235	15	8	2345.235
0.9	2412.7425	6	0	9650.97	3	241.27425	11	9	2412.7425
12.25	8720.87	49	872.087	17441.74	2	0	8	10	8720.87
44.5	11196.12	178	0	22392.24	1	1119.612	5	11	11196.12
4.2	92.51	28	0	370.04	5	9.251	13	12	9.251
3.25	11086.11	13	0	22172.22	4	11086.11	9	13	11086.11
11.25	14461.445	45	0	28922.89	5	14461.445	6	14	14461.445
3	13246.325	12	13246.325	26492.65	1	0	9	15	13246.325
21.75	6595.66	87	659.566	13191.32	1	0	7	16	6595.66
14	6660.665	56	0	1321.33	3	666.0665	7	17	6660.665
5.5	5770.575	22	577.0575	11541.15	2	0	9	18	5770.575
6.25	7415.74	25	0	14831.48	5	741.574	7	19	7415.74
13.65	1197.62	91	0	4790.48	3	119.762	13	20	1197.62
20.25	8395.84	81	0	16791.68	1	839.584	6	21	8395.84
1.65	820.0825	11	0	3280.33	4	164.0165	15	22	820.0825
12.5	5525.555	50	552.5555	11051.11	1	0	7	23	5525.555
73.25	10686.07	293	10686.07	21372.14	3	0	9	24	10686.07
0	0	0	0	0	0	0	0	25	0
40.75	6185.62	163	6185.62	12371.24	2	0	5	26	6185.62
0	0	0	0	0	0	0	0	27	0

5.9 Lab conclusion

In this lab, you learned about [Data Virtualization](#) to complete the [Collect](#) tasks by creating a virtualized view of the transformed Db2 Demographics and Db2 Churn data, joined with the MongoDB (JSON) Activity data.



 Data Engineer	<p>The SQL query runs on the remote data sources (Db2 and MongoDB) when you call the data set from the project. This capability provides the absolute latest “current state” information from the activity (i.e. transactional data store). If you want to learn even more about data virtualization, see:</p> <p>(Deeper Dive) Data Virtualization (Caching)</p> <p>Final note: Cloud Pak for Data refers to joins of individual virtualized tables as “virtualized views.”</p>
------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

** End of Lab 05 - Collect: Virtualize

Lab by Burt Vialpando and Kent Rubin, IBM

Lab 06 ANALYZE: AUTOAI

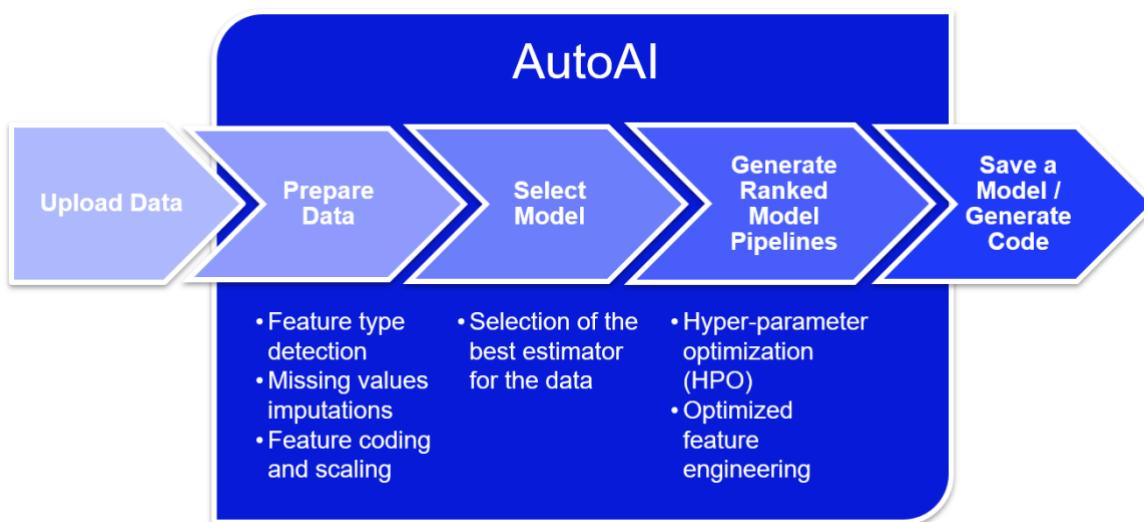
6.1 Lab overview

Analyze is the third phase in the Cloud Pak for Data platform and workflow. This is where data scientists and business analysts can join forces to gain insights from their organization's data. Analyze capabilities can be delivered with a number of different CPD services. Examples:

- AutoAI (included with the Watson Machine Learning service)
- Notebook model creation (included with the Watson Studio service)
- SPSS Modeler
- Decision Optimization
- Cognos Analytics (or Cognos Analytics Embedded)
- IBM Streams (which could be “Collect” or “Analyze” depending on how it is used)
- RStudio model creation...and others

In this lab you will explore AutoAI, which helps simplify the Machine Learning model AI lifecycle by automating the following:

- Data preparation
- Model development
- Feature engineering
- Hyper-parameter optimization



6.2 Persona represented in this lab

The **Data Scientist** persona is the most likely role to perform the **Analyze** tasks in this lab, that is, to create a machine learning model with AutoAI that can be deployed and infused into an AI application.

Persona (Role)	Capabilities
 Data Scientist	Data Scientists bring expertise in statistics and the process of building ML/AI models to make predictions and answer key business questions.

6.3 Logging into the CPD web client (if you have not already done so)

- 1. If you are starting this lab stand-alone (without going through previous labs) do the following:
- 2. Click the desktop icon: [Cloud Pak for Data Web Client](#).



- 3. The CPD web client GUI displays as shown. Use `cpduser` and `cpdaccess` for the *Username* and *Password* and click [Sign in](#).



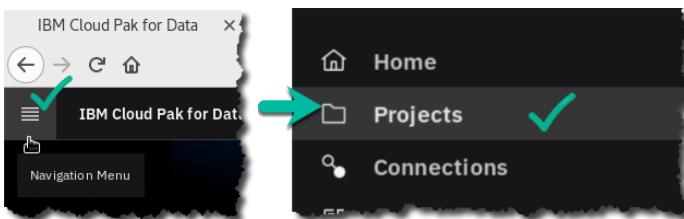
6.4 Setting up the AutoAI experiment

In the lab, you will create an AutoAI experiment that will be used to automatically create the machine learning model that best fits the data to provide the desired outcome. One only needs to provide general guidance, and AutoAI will do the rest of the work.

In our scenario, Trade Co. data scientists accelerate their time to value using this powerful tool.



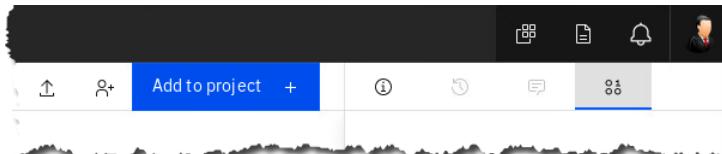
- 4. In the CPD web client, click the [Navigation Menu](#) ("hamburger" icon) ⇨ [Projects](#).



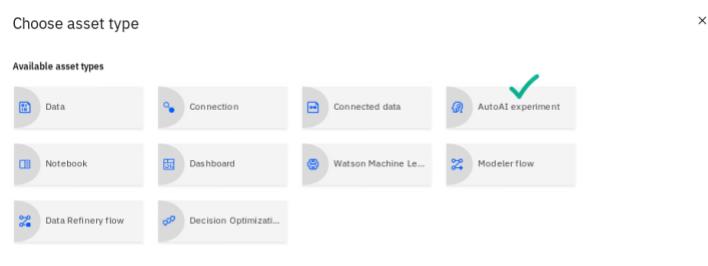
- __5. Select the project: [CPD Workshop Analytics Project](#).

Name	Project type
CPD Workshop Analytics Project ✓	Analytics
CPD Workshop Data Transformation Project	Data transformation

- __6. Once the project is opened, at the top right corner of the screen click on: [Add to project +](#).



- __7. In the next screen hover over and then click on the tile: [AutoAI experiment](#).



- __8. In **Name**, enter [ChurnRisk AutoAI experiment](#).

Fill in anything you like in **Description**.

Leave the defaults for **Compute configuration**.

Click [Create](#).

New AutoAI experiment

Define details

Name * ChurnRisk AutoAI Experiment

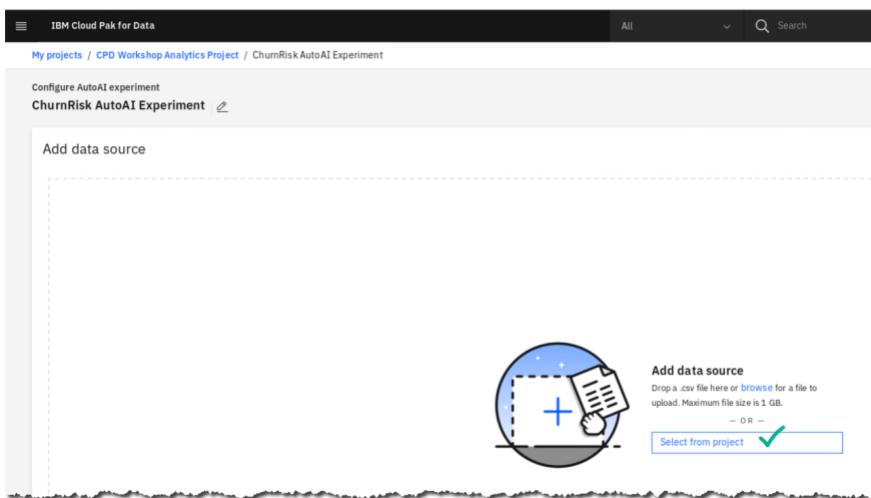
Compute configuration * 4 vCPU and 16 GB RAM

Description CPD Workshop Auto AI

Create

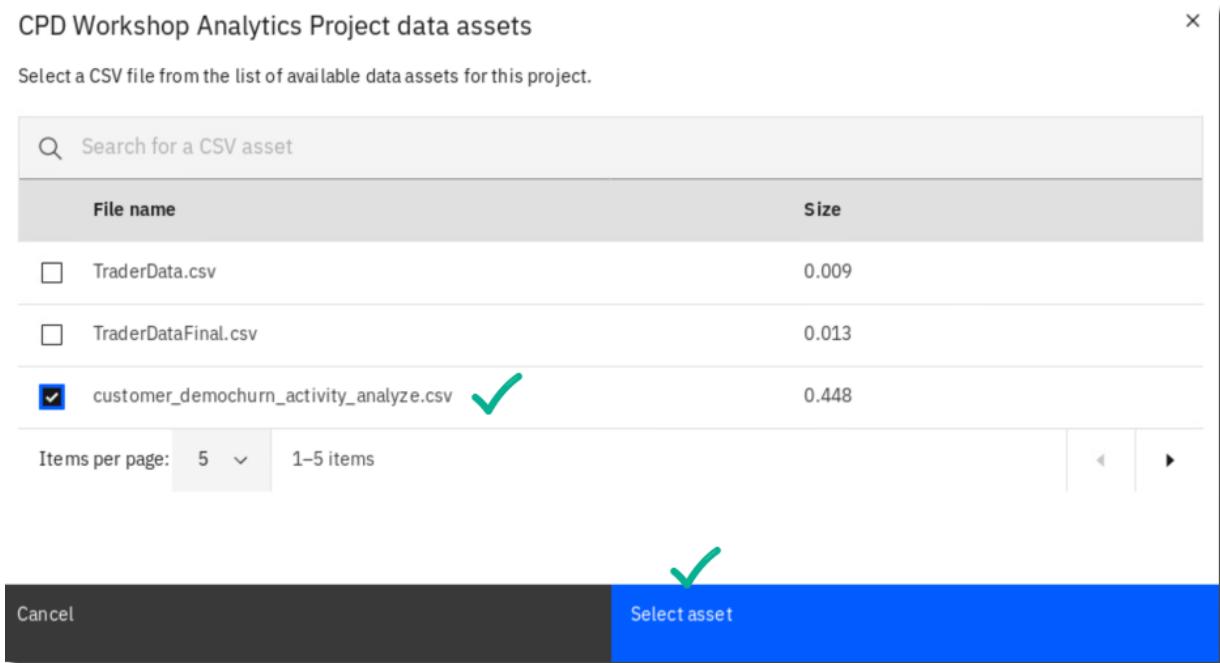
Note: If you have to wait a long time for the spinning circle to complete, simply return to the CPD Workshop Analytics Project and select ChurnRisk AutoAI Experiment from the list of Assets.

_9. In the screen **Add data source**, click **Select from project**.



_10. Select File name: **customer_demochurn_activity_analyze.csv**.

Click **Select asset**.



Cancel

Select asset



The CSV file used in this AutoAI experiment is a join of Db2 CUSTOMER_CHURN, Db2 CUSTOMER_DEMOGRAPHICS and MongoDB CUSTOMER_ACTIVITY data.

If you have been doing all the labs so far in this workshop, you would have completed the previous Data Flow Designer and Data Virtualization labs that transformed and joined these tables together as one virtualized view.

Since AutoAI requires a file as input, this virtualized view was exported to a CSV file (**customer_demochurn_activity_analyze.csv**) to be used as input for this lab.

__11. In the screen [Select prediction column](#), select Column name: CHURNRISK.

The screenshot shows the 'Select prediction column' interface. On the left, there's a sidebar for 'Add data source' with a file upload area and a 'Select from project' button. The main panel lists columns from a data source named 'customer_demo churn_activity_analyze.csv'. The 'CHURNRISK' column is highlighted with a green checkmark and a red box around the selection checkbox. Other columns listed include ID, AGE_GROUP, GENDER, STATUS, CHILDREN, ESTINCOME, HOMEOWNER, AGE, and TAVIS. At the bottom, it says 'Prediction column: CHURNRISK', 'PREDICTION TYPE: Multiclass Classification', and 'OPTIMIZED METRIC: Accuracy'.

__12. Notice that the bottom of this screen now fills in when CHURNRISK is select as the Prediction column.

AutoAI has determined for you a PREDICTION TYPE = Multiclass Classification.

Click [Experiment settings](#).

This is a close-up view of the 'Experiment settings' section. It shows 'Prediction column: CHURNRISK' at the top. Below it, 'PREDICTION TYPE' is set to 'Multiclass Classification' (with a red box around it), and 'OPTIMIZED METRIC' is set to 'Accuracy'. At the bottom, there are two buttons: 'Experiment settings' (with a checkmark icon) and 'Run experiment'.

- __13. The first Experiment setting you can change from the values chosen for you is [Data source](#)
 This allows you to change the [Holdout data split](#) (for testing vs. training)

The screenshot shows the 'Data source settings' section of the experiment configuration. It includes fields for the prediction column ('CHURNRISK'), column data type ('String'), and data source ('customer_demo'). A 'Subsample' toggle is off. The 'Training data split' slider is set to 85% (90% - 3 folds) and the 'Holdout data split' is set to 10%. A callout box highlights the 'Holdout data split' value of 10%.

If you scroll down you can see that you can optionally choose to select which columns to include in the experiment as well.

Leave these settings as-is and click on the section: [Prediction](#).

- __14. The second setting you can change is [Prediction](#).

The first section in the Prediction settings is the [Prediction type](#).

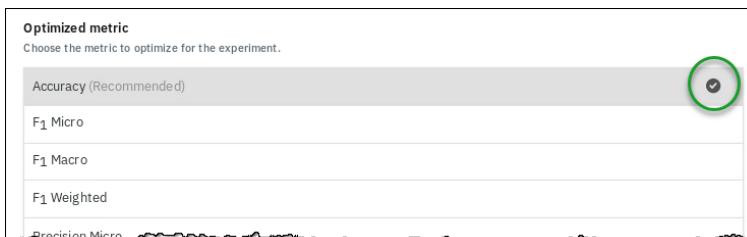
Read what all three Prediction types are best suited for. The one chosen for you was derived from the column [CHURNRISK](#), which, by the way, has three values - High, Medium and Low. Thus, AutoAI determined that Multiclass classification is best suited for this data.

Leave these setting as-is .

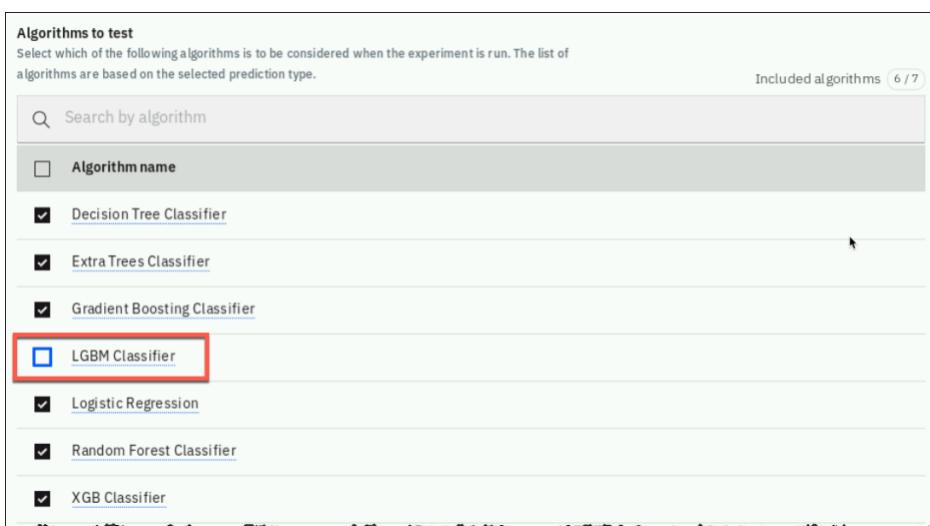
The screenshot shows the 'Prediction type' section of the experiment configuration. It includes descriptions for 'Binary classification' (classifying data into two categories) and 'Multiclass classification' (classifying data into multiple categories). A green arrow points to the 'Multiclass classification' section, which is highlighted with a checked checkbox.

- __15. Scroll down to review the next section: [Optimized metric](#). Notice that AutoAI chose **Accuracy** for you. Hover over a few of the other metrics you could choose if you so desired.

A Data Scientist would best determine if and when to deviate from this recommended metric, but you will leave this choice as-is.

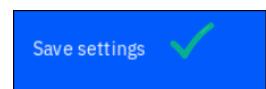


- __16. Scroll to review third section which allows you to choose which algorithms to test. [Remove the LGBM Classifier](#) from the list by deselecting it.



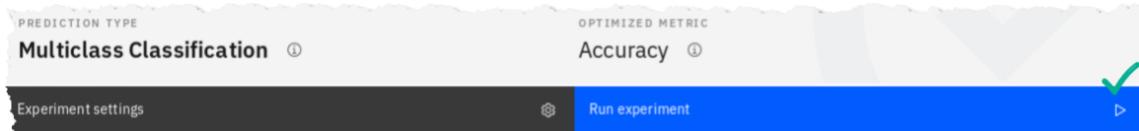
- __17. Click the [Runtime](#) setting to review the last settings for your experiment.

- __18. Click [Save settings](#) to save your changes.

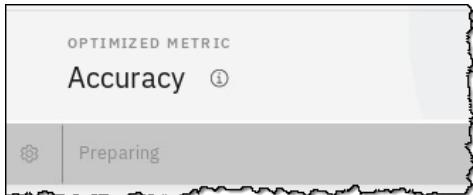


6.5 Running the AutoAI experiment

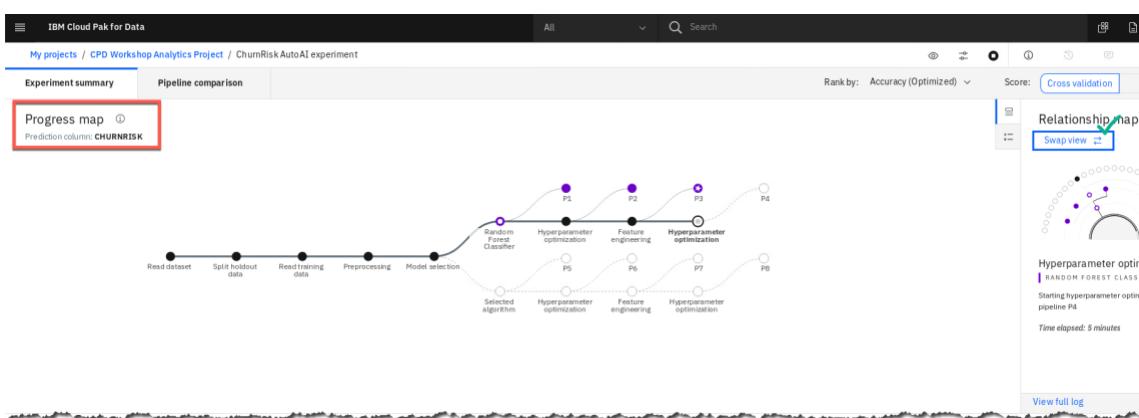
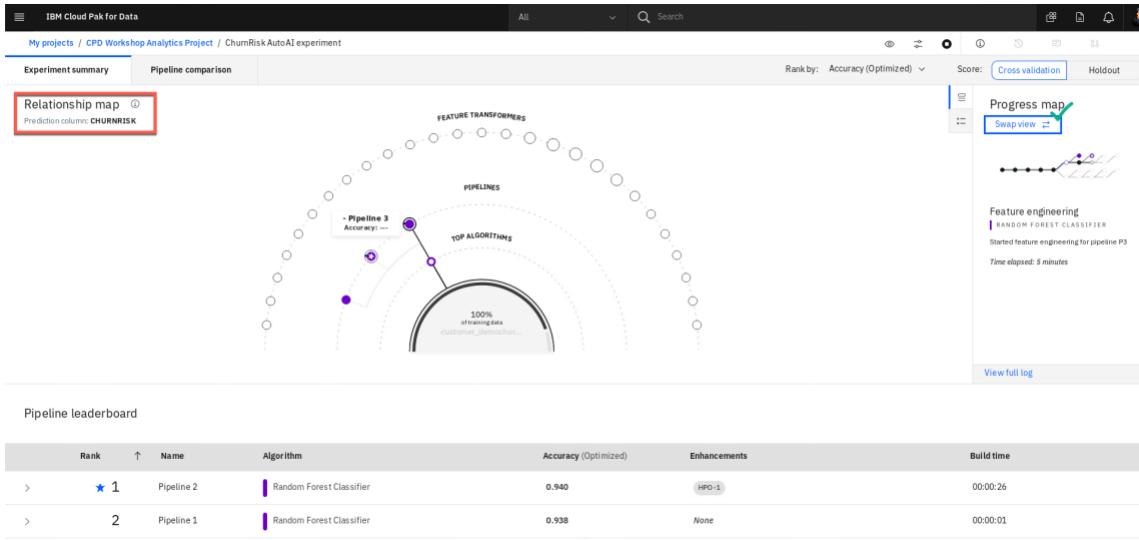
- __19. At the bottom of the screen, click Run experiment.



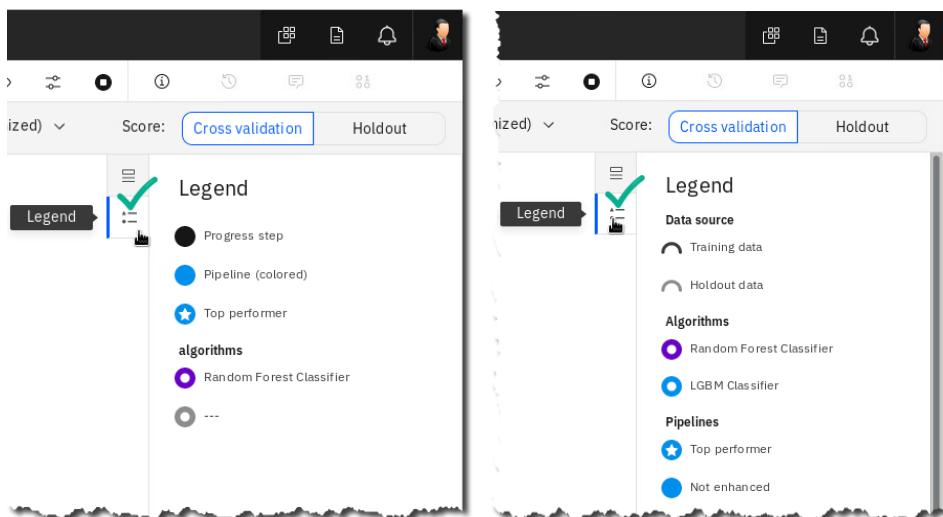
- __20. You will first see AutoAI go into a Preparing mode... then it will be Pending.



- __21. Once it is running, you can click on [Swap view] to see each of the two different infographics: Relationship map and Progress map.



__22. Review the [Legend](#) for each infographic map.



__23. On either infographics map, you can scroll down to see the [Pipeline](#) leaderboard.

AutoAI begins building several different pipelines (possible solutions for the best model) ranking them as it goes along. You will see the pipelines complete and rank as the process continues.

Pipeline leaderboard						
Rank	↑	Name	Algorithm	Accuracy (Optimized)	Enhancements	
>	★ 1	Pipeline 4	Random Forest Classifier	0.950	HPO-1	FE HPO-2
>	2	Pipeline 8	Decision Tree Classifier	0.944	HPO-1	FE HPO-2
>	3	Pipeline 3	Random Forest Classifier	0.943	HPO-1	FE
>	4 *	Pipeline 2	Random Forest Classifier	0.940	HPO-1	
>	5	Pipeline 7	Decision Tree Classifier	0.939	HPO-1	FE
>	6	Pipeline 6	Decision Tree Classifier	0.938	HPO-1	
>	7	Pipeline 1	Random Forest Classifier	0.938		None
>	8	Pipeline 5	Decision Tree Classifier	0.891		None

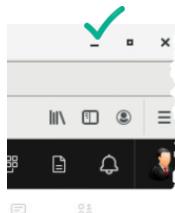
24. As AutoAI does its work, watch it flow through its various steps:

- **Read dataset:** Reads the data set you provided for the experiment.
- **Split holdout data:** Splits the data between testing and training.
- **Read training data:** Reads training data to prepare for preprocessing.
- **Preprocessing:** Most data sets contain different data formats and missing values, but standard ML algorithms work with numbers and no missing values. AutoAI applies various algorithms (estimators) to analyze, clean, and prepare your raw data for machine learning. It automatically detects and categorizes features based on data type, such as categorical or numerical. It determines the best combination of strategies for missing value imputation, feature encoding, and feature scaling for your data.
- **Model selection:** AutoAI uses a novel approach that enables testing and ranking candidate algorithms against small subsets of the data, gradually increasing the size of the subset for the most promising algorithms to arrive at the best match by ranking large numbers of candidate algorithms. This approach saves time without sacrificing performance.
- **Selected estimator:** Shows the estimator chosen from the model selection step.
- **Hyperparameter optimization:** Refines the best performing model pipelines by using a novel hyperparameter optimization algorithm optimized for costly function evaluations such as model training and scoring that are typical in machine learning. This approach enables fast convergence to a good solution despite long evaluation times of each iteration.
- **Feature engineering:** Attempts to transform the raw data into the combination of features that best represents the problem to achieve the most accurate prediction. AutoAI uses a novel approach that explores various feature construction choices in a structured, non-exhaustive manner, while progressively maximizing model accuracy using reinforcement learning. This results in an optimized sequence of transformations for the data that best match the algorithms of the model selection step.

6.6 Running a Notebook

While the AutoAI experiment is running, open another CPD web client to perform a parallel exercise.

- _25. Minimize the CPD web client browser to be able to get to the desktop.

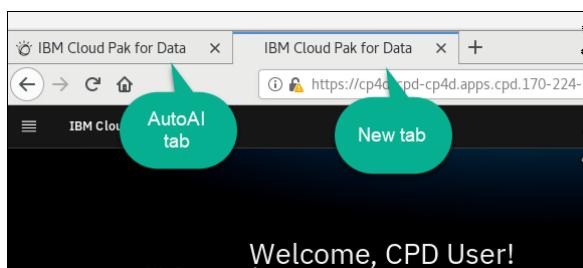


- _26. Double click the desktop icon: [Cloud Pak for Data Web Client](#).

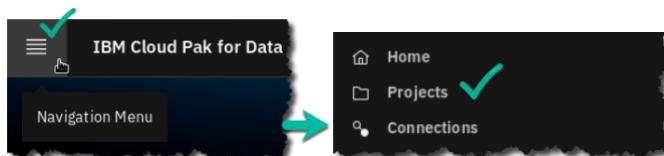


- _27. This will open a second CPD Web client browser tab.

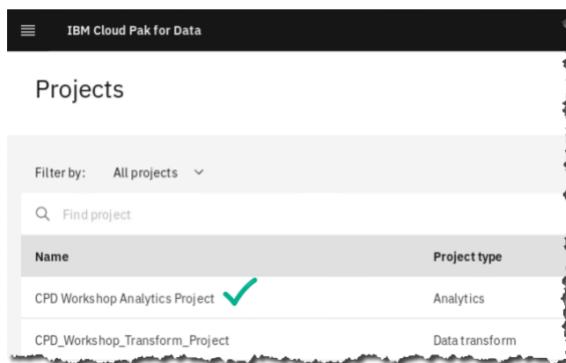
The AutoAI session is running in the first tab **so don't close it!**



- _28. In the CPD web client, click the [Navigation Menu](#) \Rightarrow [Projects](#).



- _29. Select the project: [CPD Workshop Analytics Project](#).



__30. In the section [Assets](#), scroll down to find [Notebooks](#).

Click [TradingCustomerChurnClassifier-Py36](#).

__31. You will be presented with the opened notebook.

__32. Click the [Edit](#) (pencil) icon to put the Notebook in edit mode.

(Note: If you are returning to the Notebook, it may already be in edit mode.)

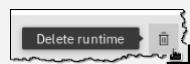


Note: if you received an error like this “403: forbidden.”

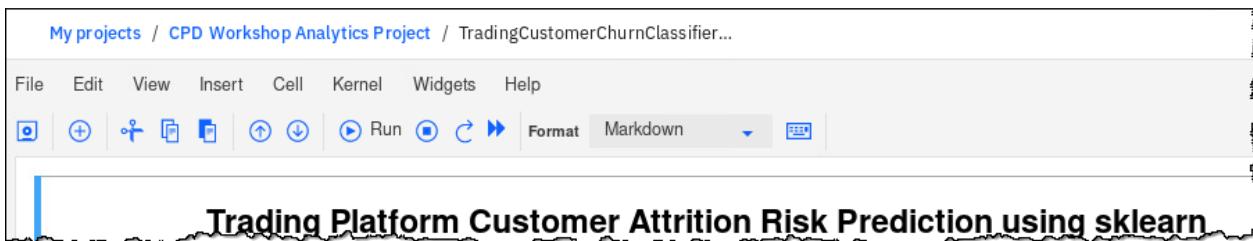


Data Scientist

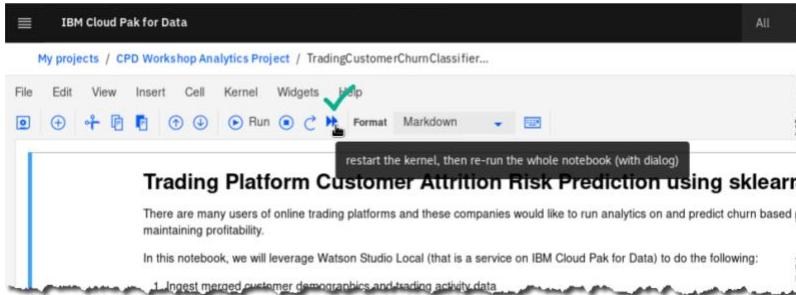
Fix the problem by leaving this screen and going to [My Instances](#) \Rightarrow [Environments](#) and then delete the Runtime environment for Python 3.6 that is currently running. Try opening the notebook again after returning to the project.



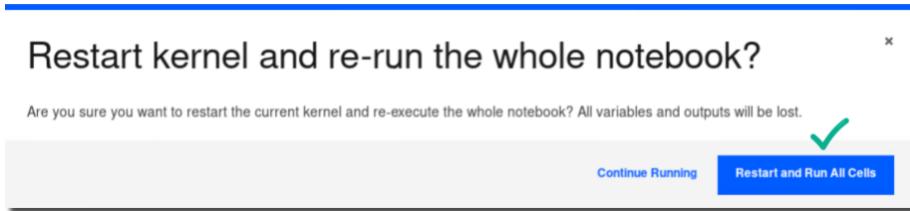
- __33. The Notebook will start a runtime and present a screen with the top left looking like this:



- __34. Click on the **Fast Forward** (double arrow) icon to re-run the whole notebook.



- __35. Click **Restart and Run All Cells**.



- __36. The notebook will now start running each cell in sequential order. A number will be placed next to each cell after it has executed.

Scroll up to the first cells under **1. Load libraries**.

 A screenshot of the Jupyter Notebook showing the first cell of the '1. Load libraries' section. The code cell contains:


```
In [1]: #Uncomment and run once to install the package in your runtime environment
!pip install scikit-learn==0.22
!pip install sklearn-pandas
```

 Below the cell, a progress bar indicates the download of 'scikit-learn-0.22-cp36-cp36m-manylinux1_x86_64.whl' (7.0 MB) from a local source.



Note: You can run the notebook cell by cell or all at once. Either way will give you the same result. Any cell that has not yet run is indicated like this: [*]

- __37. Scroll to the 4th cell (under [2. Load data example](#)) to see the input file for this notebook. Notice it is the same input file from the same project that AutoAI is using.

In [4]:

```
df_churn_pd = pd.read_csv('/project_data/data_asset/customer_demochurn_activity_analyze.csv')
df_churn_pd.head()
```

Out[4]:

	ID	AGE_GROUP	CHURNRISK	GENDER	STATUS	CHILDREN	ESTINCOME	HOMEOWNER	AGE	TAXID	...	LATITUDE
0	0	Young adult	Low	F	S	1	38000.00	N	24	147889187	...	NaN
1	1	Adult	Low	M	M	2	29616.00	N	49	113772166	...	38.687261
2	2	Adult	Low	M	M	0	19732.80	N	51	132420919	...	NaN
3	3	Adult	High	M	S	2	96.33	N	56	700548452	...	32.531971
4	4	Young adult	High	F	M	2	52004.80	N	25	141013706	...	33.593192

Scroll through this section of the notebook to see various visualizations of the data.

- __38. Review the cells in [3. Data preparation](#).

3. Data preparation

[Top](#)

Data preparation is a very important step in machine learning model building. This is because

During this process, we identify categorical columns in the dataset. Categories needed to

- __39. The next section is self-explanatory: [4. Build Random Forest classification model](#).

4. Build Random Forest classification model

[Top](#)

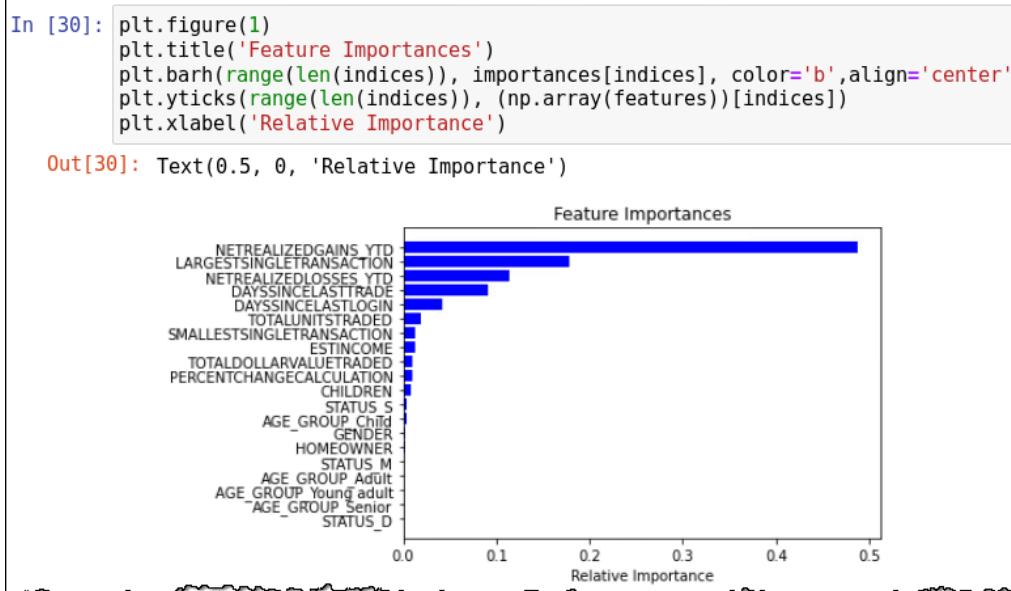
We instantiate a decision-tree based classification algorithm, namely, RandomForestClassifier. Next we define how to combine multiple algorithms into a single pipeline, or workflow.

We split original dataset into train and test datasets. We fit the pipeline to training data and apply the trained

- __40. Scroll to find [Model Results](#), then find the Accuracy.

The screenshot shows a Jupyter Notebook cell with the title "Model results". Below it, a text cell contains the Python code: `In [25]: print('Accuracy: ', sklearn.metrics.accuracy_score(y_test, y_prediction))`. The output of the code, "Accuracy: 0.9387096774193548", is highlighted with a red box and a green arrow pointing to it.

- __41. Cell [30] shows the “Feature Importances,” that is, the data columns that affected the model the most



- __42. Finally, review the name listed in [5. Save the model into WML Deployment Space](#).

We will be referring to this Deployment Space in a later lab.

5. Save the model into WML Deployment Space

[Top](#)

Before we save the model we must create a deployment space. Watson Machine Learning pro
space

The screenshot shows a Jupyter Notebook cell with the code: `In [32]: # Specify a names for the space being created, the saved model and the model deployment
space_name = 'deployment-space-analytics-project-workshop'
model_name = 'churn_risk_model'
deployment_name = 'churn_risk_model-deployment'`. The variable `space_name` is highlighted with a red box and a green arrow pointing to it.

- 43. Also note at the very end of the notebook in the last two cells that two files are created for batch scoring and evaluation.

```

Write test data into .csv files for batch scoring and model evaluations

In [48]: # Write the test data a .csv so that we can later use it for batch scoring
write_score_CSV=X_test
write_score_CSV.to_csv('/project_data/data_asset/model_batch_score.csv', sep=',', index=False)

In [49]: # Write the test data to a .csv so that we can later use it for Evaluation
write_eval_CSV=X_test
write_eval_CSV.to_csv('/project_data/data_asset/model_eval.csv', sep=',', index=False)

```

Your notebook has NOT finished until you see that the last two code cells (above) have a number from the run.

- 44. One key data point to take particular note of from this notebook run is found in cell 25, which is the accuracy of the Random Forest model created by this notebook.

```

In [25]: print('Accuracy: ', sklearn.metrics.accuracy_score( y_test, y_prediction ))
Accuracy:  0.9387096774193548

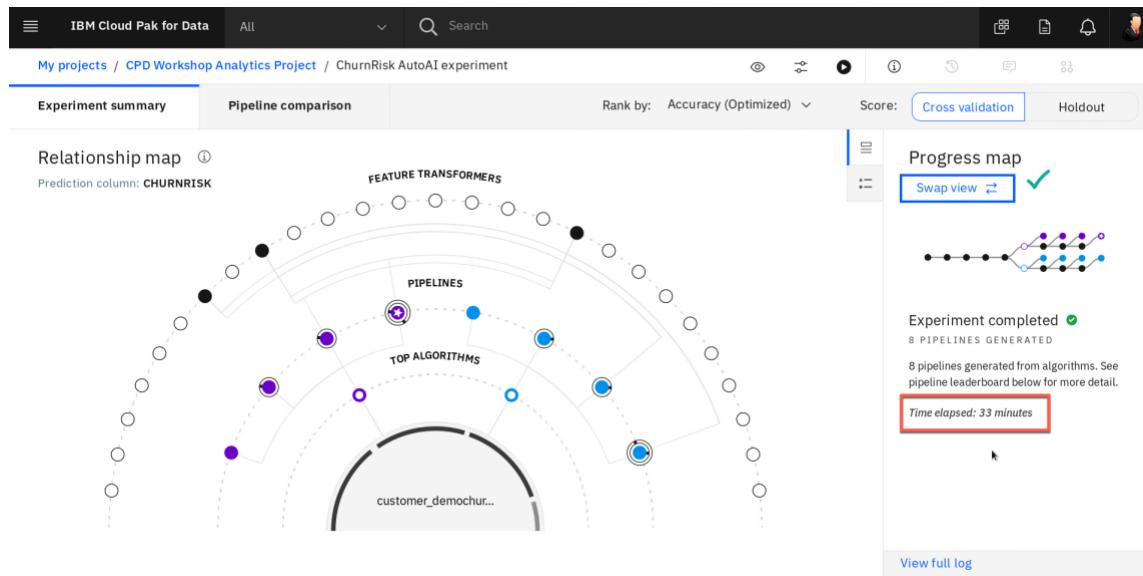
```

 Data Scientist	<h3 style="text-align: center;">Ten Reasons Why I Like my Jupyter Notebook</h3> <ol style="list-style-type: none"> All in one place: The Jupyter Notebook is a web-based interactive environment that combines code, rich text, images, videos, animations, mathematical equations, plots, maps, interactive figures and widgets, and graphical user interfaces, into a single document. Easy to share: Notebooks are saved as structured text files (JSON format), which makes them easily shareable. Easy to convert: Jupyter comes with a special tool, nbconvert, which converts notebooks to other formats such as HTML and PDF. Language independent: The architecture of Jupyter is language independent. The decoupling between the client and kernel makes it possible to write kernels in any language. Easy to create kernel wrappers: Jupyter brings a lightweight interface for kernel languages that can be wrapped in Python. Wrapper kernels can implement optional methods, notably for code completion and code inspection. Easy to customize: Jupyter's interface can be used to create an entirely customized experience in the Jupyter Notebook (or another client application such as the console). Extensions with custom magic commands: Create IPython extensions with custom magic commands to make interactive computing even easier. Many third-party extensions and magic commands exist, for example, the %%cython magic that allows one to write Cython code directly in a notebook. Stress-free Reproducible experiments: Jupyter notebooks can help you conduct efficient and reproducible interactive computing experiments with ease. It lets you keep a detailed record of your work. Also, the ease of use of the Jupyter Notebook means that you don't have to worry about reproducibility; just do all of your interactive work in notebooks, put them under version control, and commit regularly. Don't forget to refactor your code into independent reusable components. Effective teaching-cum-learning tool: The Jupyter Notebook is not only a tool for scientific research and data analysis but also a great tool for teaching. Interactive code and data exploration: The ipywidgets package provides many common user interface controls for exploring code and data interactively.
-------------------------------------------------------------------------------------------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

6.7 Reviewing the AutoAI results

- 45. Return to your first CPD Web client browser tab.

If the AutoAI experiment has not completed yet, let it finish. You can tell if it is finished by looking at the [Relationship map view](#).



Data
Scientist

Note: this experiment may take 25 minutes or more to complete, but this is not normal. We have turned off AVX/AVX2 processor support for this workshop so that the Cloud Pak for Data workshop image can run on servers that do not have AVX nor AVX2 processors. If we did not do that and the workshop was run on servers without AVX/AVX2 processors, the experiment would fail. So, we took this route to make sure the AutoAI experiment would run under any circumstances on any server, with or without AVX/AVX2 processors.

That said, with a CPD cluster on AVX/AVX2 supported processors this experiment completes on average, in around 3 minutes.

Servers with AVX/AVX2 processors were released in 2011 and are mostly ubiquitous. Your organization probably has them.



Data
Scientist

Note: the results from your AutoAI experiment may vary from the illustrations in this workbook. This is especially true for Feature Transformations.

- 46. AutoAI chooses the best model from the various pipeline leaderboard options as Rank #1. You can scroll down in either map infographic screen to see the leaderboard.

Notice the Algorithm chose (Random Forest classifier) and the most accurate result.

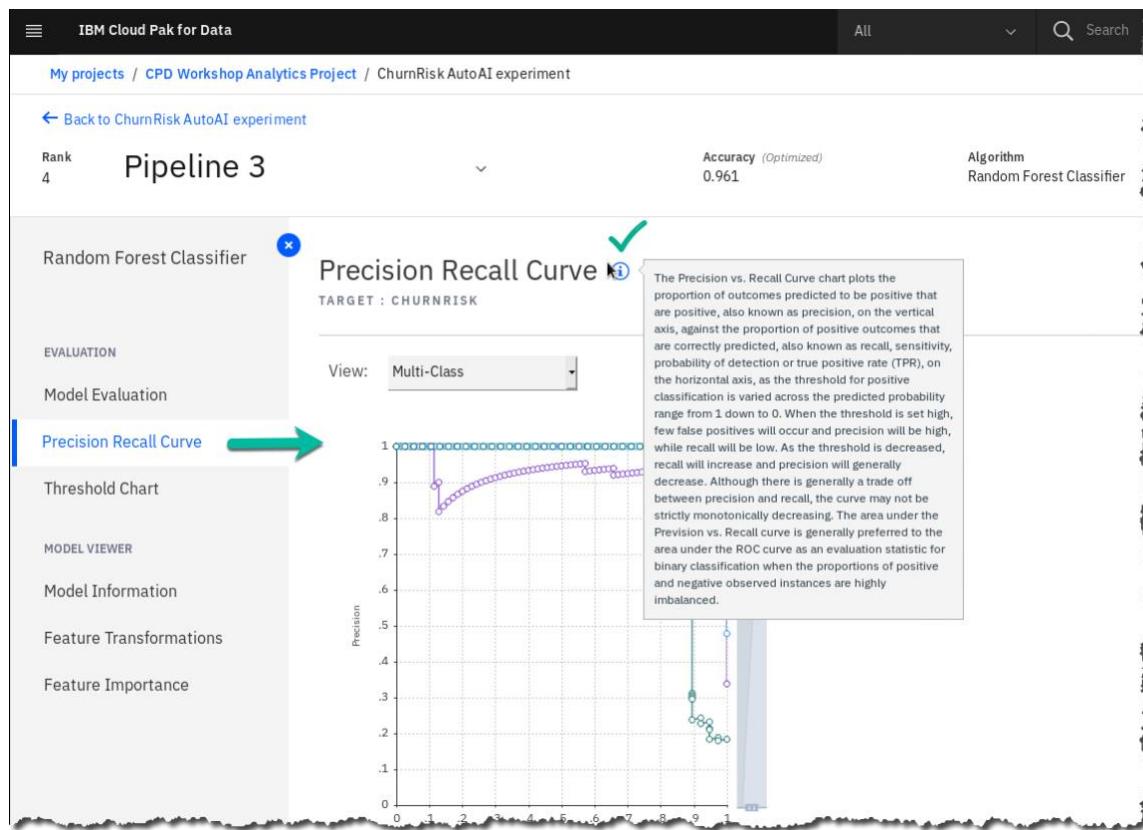
Click on the #1 ranked Pipeline 4.

Rank	Name	Algorithm	Accuracy (Optimized)	Enhancements
> 1	Pipeline 4	Random Forest Classifier	0.950	HPO-1 FE HPO-2
> 2	Pipeline 8	Decision Tree Classifier	0.944	HPO-1 FE HPO-2
> 3	Pipeline 3	Random Forest Classifier	0.943	HPO-1 FE
> 4	Pipeline 2	Random Forest Classifier	0.940	HPO-1
> 5	Pipeline 7	Decision Tree Classifier	0.939	HPO-1 FE

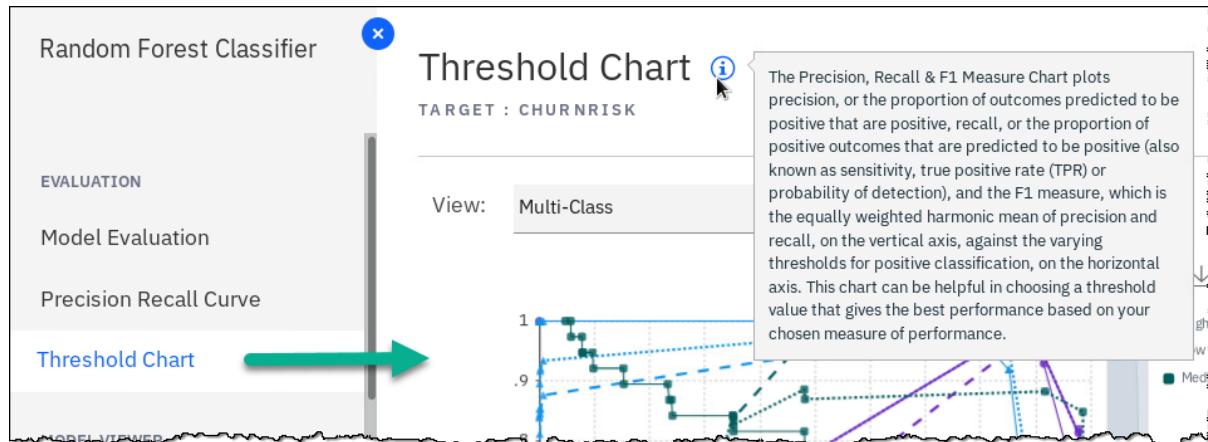
- 47. Model Evaluation shows the various evaluation accuracy figures. Click through them but make sure to return to the original view - Multi-Class.

The screenshot shows the IBM Cloud Pak for Data interface. At the top, it displays "My projects / CPD Workshop Analytics Project / ChurnRisk AutoAI experiment". Below this, a navigation bar includes "Back to ChurnRisk AutoAI experiment". The main content area is titled "Pipeline 3" with a rank of 4 and an accuracy of 0.961. A sidebar on the left lists "Model Evaluation", "Precision Recall Curve", "Threshold Chart", "Model Information", "Feature Transformations", and "Feature Importance". The "Model Evaluation" option is selected. A dropdown menu titled "View:" is open, showing "Multi-Class" (which is checked), "High (One v. Rest)", "Low (One v. Rest)", and "Medium (One v. Rest)". To the right of the dropdown, there is a "ROC Curve" plot with a dashed diagonal line and several colored data points representing different thresholds. The overall interface has a modern design with dark mode elements.

—48. **Precision Recall Curve** Shows the tradeoff between precision and recall.



—49. **Threshold Chart** helps choose a threshold for best performance based on the chosen measure of performance.



_50. [Model Information](#) gives info on type of model fitted.

The screenshot shows the 'Model Information' section of a dashboard. On the left, there's a sidebar with options: Random Forest Classifier, Precision Recall Curve, Threshold Chart, MODEL VIEWER, Model Information (which is selected and highlighted with a green arrow), and Feature Transformations. The main area has a title 'Model Information' with a help icon. Below it, it says 'TARGET : CHURNRISK'. A callout box states: 'This table contains information on the type of model fitted, identifies the target field and the number of input features.' The table itself has three rows:

Label (Target)	CHURNRISK
Model Type	Random Forest Classifier
Number of Features	45

_51. [Feature Transformations](#) displays new features created by AutoAI. This is a powerful capability of AutoAI because this is not something a Data Scientist may intuitively do themselves in a notebook. (Note: your results may vary from this screen shot.)

The screenshot shows the 'Feature Transformations' section of a dashboard. The sidebar on the left includes: Random Forest Classifier, Precision Recall Curve, Threshold Chart, MODEL VIEWER, Model Information, Feature Transformations (selected and highlighted with a green arrow), and Feature Importance. The main area has a title 'Feature Transformations' with a help icon. Below it, it says 'TARGET : CHURNRISK'. A callout box explains: 'Displays any new features created during pipeline building, along with the transformation function(s) and the original feature(s) transformed. Hover over the transformation name to get more details.' The table lists two new features:

New Feature	Original Feature	Transformation
NewFeature_11	All	pca(All)
NewFeature_9	All	pca(All)

_52. [Feature Importance](#) displays the relative importance of the feature in predicting the target.

The screenshot shows the 'Feature Importance' section of a dashboard. The sidebar on the left includes: Random Forest Classifier, Precision Recall Curve, Threshold Chart, MODEL VIEWER, Model Information, Feature Transformations, and Feature Importance (selected and highlighted with a green arrow). The main area has a title 'Feature Importance' with a help icon. Below it, it says 'TARGET : CHURNRISK'. A callout box states: 'Shows the relative importance of each feature in predicting the target, based on an averaging of nine different importance measures.' The chart is a horizontal bar chart with three bars:

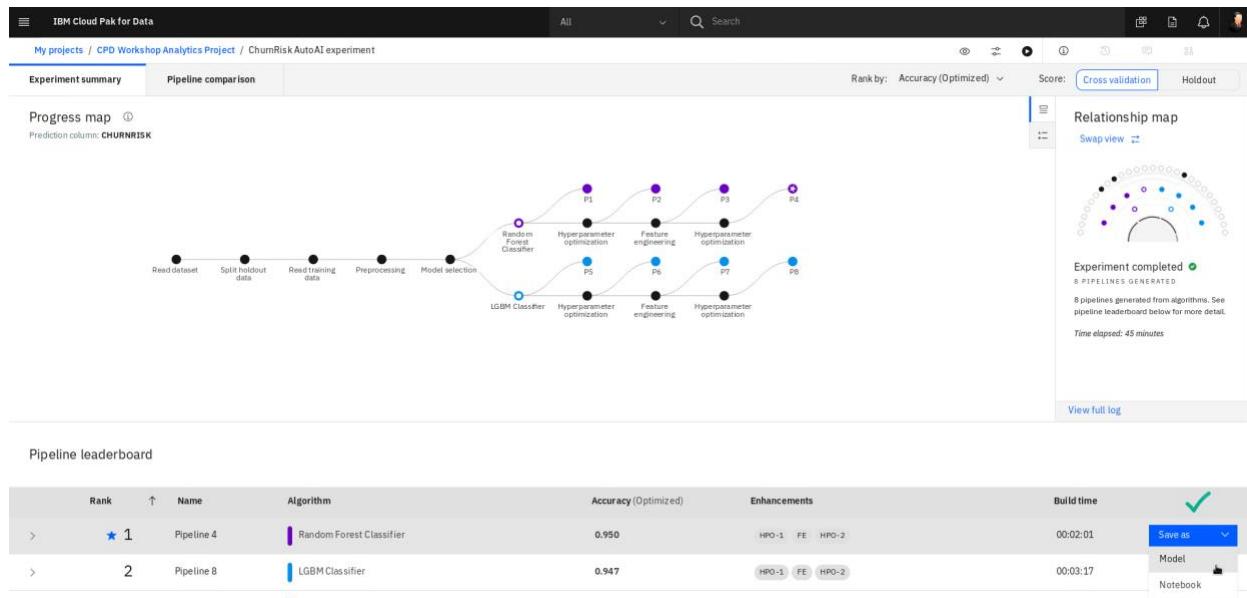
Feature	Importance
NETREALIZEDG...	0.40
NewFeature_11	0.24
NETREALIZEDL...	0.20

6.8 Saving the model

AutoAI gives you two options to save the model you decide is the best one for your application. You can either save the model as a model or save the model as a Notebook.

6.8.1 Saving the model as a model

_53. Hover your mouse over the right side of your Top Rank model, click **Save as** .



Rank	Name	Algorithm	Accuracy (Optimized)	Enhancements	Build time	
> 1	Pipeline 4	Random Forest Classifier	0.950	HPO-1 FE HPO-2	00:02:01	 Save as
> 2	Pipeline 8	LGBM Classifier	0.947	HPO-1 FE HPO-2	00:03:17	Model Notebook

_54. In the screen **Save as model**:

Model name: keep it as is

Description: **CPD Workshop AutoAI experiment**

Click **Save**.

Save as model

Save this model as a project asset so you can deploy, train, and test it.

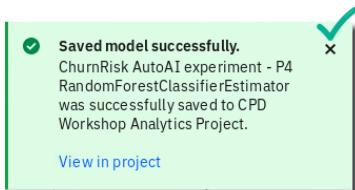
Model name

Description (optional)

Associated project

Cancel
Save 

__55. Exit the confirmation pop up by clicking **x**.



6.8.2 Saving the model as a notebook

__56. On the top right corner of the screen, click **Save as ↗ Notebook**.

Rank	Name	Algorithm	Accuracy (Optimized)	Enhancements	Build time
1	Pipeline 4	Random Forest Classifier	0.950	HPO-1 FE HPO-2	00:01:45
2	Pipeline 8	LGBM Classifier	0.947	HPO-3 FE HPO-2	00:02:15

__57. In the screen *New notebook*:

Name: keep it as is

Description: [CPD Workshop AutoAI – Notebook creation](#)

Click [Create notebook](#).

__58. You will be taken to the new notebook in edit mode in the project.

Take a moment to examine the new auto-generated notebook you just created.

Exit the notebook by clicking **CPD Workshop Analytics Project** on the breadcrumb trail.

```
#####
#Licensed Materials - Property of IBM
#(C) Copyright IBM Corp. 2020
#US Government Users Restricted Rights - Use, duplication disclosure restricted
#by GSA ADP Schedule Contract with IBM Corp.
#####
The auto-generated notebooks are subject to the International License Agreement for Non-Warranted Programs (or equivalent) and License Terms. Specifically, the Source Components and Sample Materials clause included in the License Information document for Watson Studio AutoAI generated notebooks. By using this notebook, you agree to the License Terms. http://www14.software.ibm.com/cgi-bin/weblap/lap.pl?li\_formnum=L-AMCU-BHU2B7&title=IBM+Cloud+Pak+for+Data
```

IBM AutoAI Auto-Generated Notebook v1.11.10

Note: Notebook code generated using AutoAI will execute successfully. If code is modified or reordered, there is no guarantee it will succeed. Different data will result in different output. For different data, please consider returning to AutoAI Experiments to generate a new pipeline. Please read our documentation for more information. (Cloud Pak For Data) https://www.ibm.com/support/knowledgecenter/SSONUZ_3.0.0/wsj/analyze-data/autoai-notebook.html.

Before modifying the pipeline or trying to re-fit the pipeline, consider: The notebook converts dataframes to numpy arrays before fitting them. This may result in different output for categorical values during fit of the preprocessing pipeline. Delete its members before re-fitting.

Representing Pipeline from run: Pipeline_4 from run e3247cb03b924c728fb25a4f1c45b58

6.8.3 Reviewing the project

__59. You should be back in the project and under the Assets tab. (If you are not, navigate there now.) Scroll down to find the Notebooks section. There you will find the AutoAI generated notebook you were just in.

Name	Shared	Scheduled	Status
TradingCustomerChurnClassifier-Py36			Green circle
ChurnRisk AutoAI experiment - P4 notebook			Green circle

- __60. Scroll down farther to find the model you created from the AutoAI save feature. Select the **ChurnRisk AutoAI experiment - P4 RandomForestClassifierEstimator**.

Name	Type	Software specification
ChurnRisk AutoAI experiment - P4 RandomForestClassifierEstimator	wml-hybrid_0.1	hybrid_0.1

- __61. Review your new AutoAI generated model.

Model
ChurnRisk AutoAI Experiment - P4 RandomForestClassifierEstimator

Overview

Summary

Model Type: wml-hybrid_0.1

Software specification: hybrid_0.1

- __62. Close any extra browser tabs you may have open, leaving only one.

IBM Cloud Pak for Data

IBM Cloud Pak for Data

https://cp4d-cpd-cp4d.apps.c

IBM Cloud Pak for Data

6.9 Lab conclusion

The AutoAI generated model had greater accuracy than the one built by the notebook. AutoAI can test many variations of a model in minutes that a Data Scientist would take days, weeks or even months to try out manually using a notebook approach. And AutoAI does not require any coding at all, making it an available solution to more people.

Data Scientists will always love their notebooks for the reasons listed in this lab, but AutoAI is an industry unique IBM differentiator that will enable the Data Scientist to be more productive. A data scientist can accelerate model creation using AutoAI and then tweak it using a Jupyter notebook if desired. In minutes AutoAI generates this code that would normally take hours or days to create manually.

To fulfill the promise of AI, organizations are tackling skill-set gaps, deployment, and governance processes today. In particular, businesses are seeking an alternative where citizen data scientists can quickly get started, while expert data scientists can speed experimentation time from weeks and months to minutes and hours.

Both groups need a multimodal data science and AI environment where data and analytics specialists collaborate with other experts and optimize model performance end-to-end.

AutoAI in the Cloud Pak for Data Watson Machine Learning service solves this requirement.

Read here how AutoAI won the 2019 Alconics award for “*Best Innovation in Intelligent Automation*”: <http://ibm.biz/The-Alconics-Awards>

Award for IBM's AutoAI



** End of Lab 06 - Analyze: AutoAI

Lab by Burt Vialpando, Kent Rubin, Anjali Shah and Sidney Phoon - IBM

Lab 07 DEPLOY

7.1 Lab overview

Note: This lab requires that you have completed Lab 06 Analyze: AutoAI.

In the previous lab, you created a model from both a notebook and from AutoAI. You will learn how to Deploy these models in this lab.

In our scenario, Trade Co. uses these steps to deploy their machine learning models into production.



7.2 Persona represented in this lab

The [Developer](#) persona is the likely role to perform the various [Deploy](#) tasks in this lab. However, the Data Scientist persona could perform these tasks as well.

Persona (Role)	Capabilities
 Developer	Developers create and maintain the end-user applications that utilize the output from all the other personas on the CPD platform.

7.3 Logging into the CPD web client (if you have not already done so)

- __1. If you are starting this lab stand-alone (without going through previous labs) do the following:
- __2. Double-click the desktop icon: [Cloud Pak for Data Web Client](#).



- __3. The CPD web client GUI displays as shown. Use `cpduser` and `cpdaccess` for the *Username* and *Password* and click [Sign in](#).

7.4 Reviewing the notebook deployment space

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 access and consume the model through the same interface. While this is a manual process in most organizations, Cloud Pak for Data can automate deploying and maintaining models without writing a single line of code.

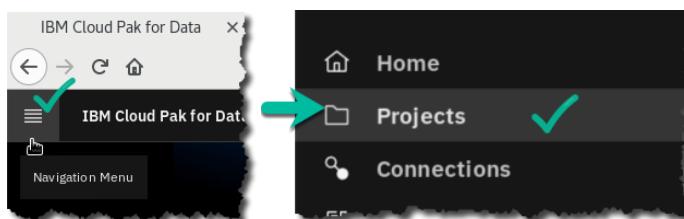
CPD eliminates the need for to do the following:

- Writing code to perform the above capability and using a runtime to deploy it.
- Creating a runtime on bare metal machines that require OS installation, network, storage, etc.
- Creating a runtime on a virtual machine in VMware on Intel, or IBM POWER VM® on a POWER platform.
- Creating a runtime in Docker or CRI-O requires someone to build the image and deploy it on one of the above platforms.

Each of the above requires manpower and machine resources. Using CPD, you can bypass this and quickly harvest insight from your data in a repetitive manner by integrating it with your CI/CD pipeline.

7.4.1 Associating a deployment space

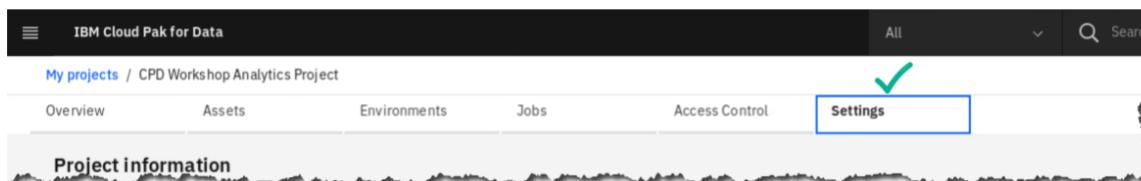
- 4. In the CPD web client, click the [Navigation Menu](#) a Projects.



- 5. Select the project: [CPD Workshop Analytics Project](#).

Name	Project type
CPD Workshop Analytics Project ✓	Analytics

- 6. Click section [Settings](#).



7. Scroll down to find [Associated deployment space](#),
 Notice it says this project is currently not associated with a deployment space.
 Click [Associate a deployment space](#),

- Click section [Existing](#),
 Click Deployment Space [deployment-space-analytics-project-workshop](#),
 Click [Associate](#),

Note: this deployment space was generated by running the notebook in the previous lab.

```
# Specify a names for the space being created, the saved model and the model deployment
space_name = 'deployment-space-analytics-project-workshop'
model_name = 'churn_risk_model'
deployment_name = 'churn_risk_model-deployment'
```

5.1 Create Deployment Space

```
# create the space and set it as default
space_meta_data = {
    client.spaces.ConfigurationMetaNames.NAME : space_name
}
```

However, you could have created a deployment space through the CPD web client as well.

Developer

- __8. The project should now show that it is associated with this deployment space. (You may need to scroll down to see this.)

The screenshot shows the 'Settings' tab of a project in IBM Cloud Pak for Data. Under the 'Project information' section, there is a 'Project name' field containing 'CPD Workshop Analytics Project'. Below it is a 'Description' field with placeholder text 'Project description'. In the 'Associated deployment space' section, there is a 'Name' field containing 'deployment-space-analytics-project-workshop', which is highlighted with a red rectangular box. The 'Integrations' section is partially visible below.

- __9. Scroll up and click **Assets**,

The screenshot shows the 'Assets' tab selected in the navigation bar of the IBM Cloud Pak for Data interface. The 'My projects' section shows 'CPD Workshop Analytics Project'. The 'Assets' tab has a green checkmark icon next to it. The 'Project information' section is visible below, showing a 'Project name' field.

- __10. If the data asset add screen is open, close it.

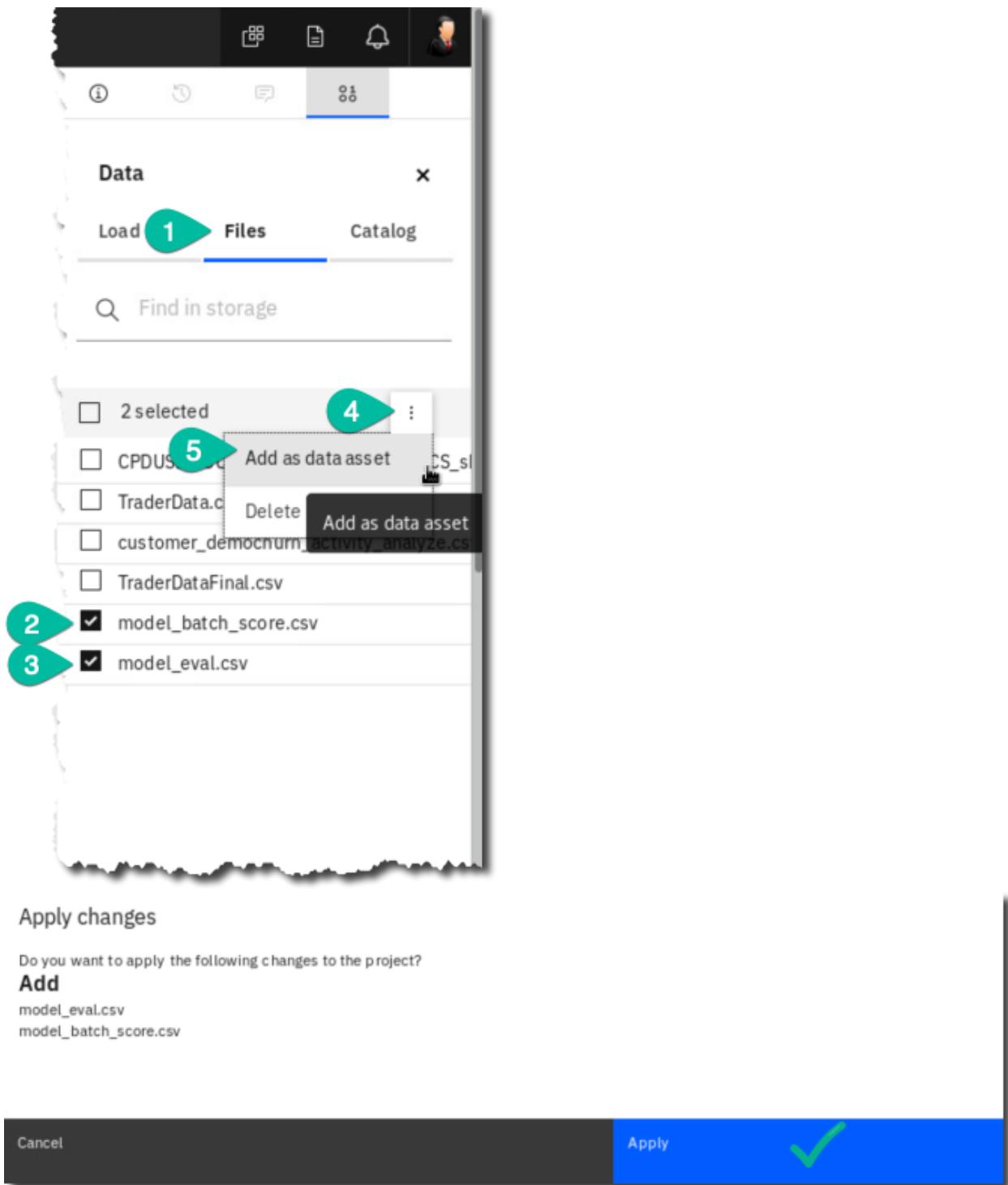
The screenshot shows the 'Data' tab selected in the 'New data asset' add screen. The 'Load' tab is active, and there is a dashed box for dropping files or browsing for files to upload. A green checkmark icon is present in the top right corner of the screen.

- __11. Click **+ New data asset** (This is how you open the data asset screen.)

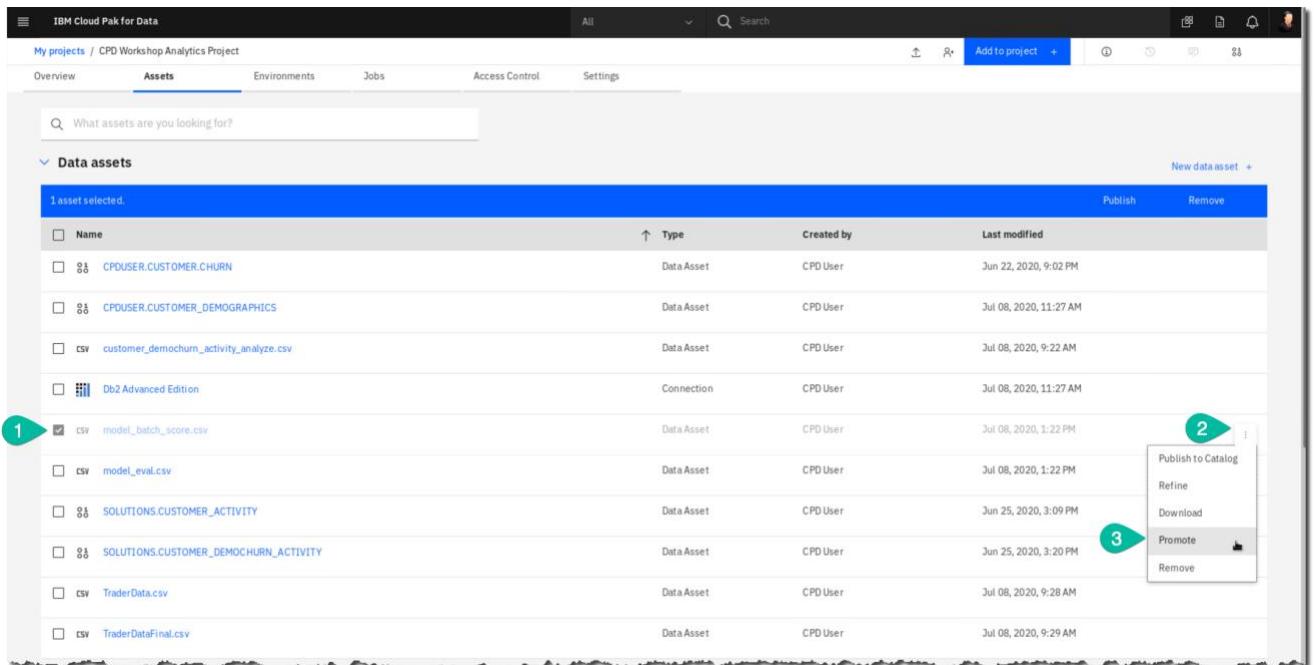
The screenshot shows a blue button labeled 'New data asset' with a white plus sign icon and a green checkmark icon to its right, used to open the data asset screen.

__12. Click **Files**, then select the two files: [model_batch_score.csv](#) and [model_eval.csv](#).

(Note: these files were created by the notebook), then click ellipses a [Add as data asset](#) and [Apply](#).



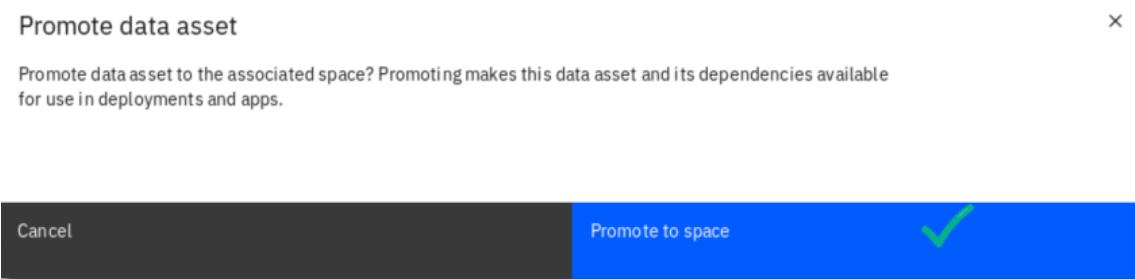
__13. Find the newly added data asset [model_batch_score.csv](#), click on ellipses  and Promote.



The screenshot shows the 'Assets' tab in the IBM Cloud Pak for Data interface. A search bar at the top has 'What assets are you looking for?' entered. Below it, a section titled 'Data assets' shows a list of 1 asset selected. The selected item is 'model_batch_score.csv'. To the right of the list is a context menu with three numbered steps: 1. 'Publish' (disabled), 2. 'Refine' (disabled), 3. 'Download' (disabled), 4. 'Promote' (highlighted with a green circle and checkmark), and 5. 'Remove' (disabled).

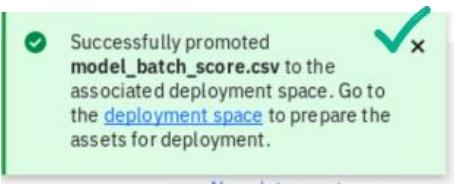
Name	Type	Created by	Last modified
CPUSER.CUSTOMER.CHURN	Data Asset	CPD User	Jun 22, 2020, 9:02 PM
CPUSER.CUSTOMER_DEMOGRAPHICS	Data Asset	CPD User	Jul 08, 2020, 11:27 AM
customer_demochurn_activity_analyze.csv	Data Asset	CPD User	Jul 08, 2020, 9:22 AM
Db2 Advanced Edition	Connection	CPD User	Jul 08, 2020, 11:27 AM
model_batch_score.csv	Data Asset	CPD User	Jul 08, 2020, 1:22 PM
model_eval.csv	Data Asset	CPD User	Jul 08, 2020, 1:22 PM
SOLUTIONS.CUSTOMER_ACTIVITY	Data Asset	CPD User	Jun 25, 2020, 3:09 PM
SOLUTIONS.CUSTOMER_DEMOCHURN_ACTIVITY	Data Asset	CPD User	Jun 25, 2020, 3:20 PM
TraderData.csv	Data Asset	CPD User	Jul 08, 2020, 9:28 AM
TraderDataFinal.csv	Data Asset	CPD User	Jul 08, 2020, 9:29 AM

__14. Click [Promote to space](#).

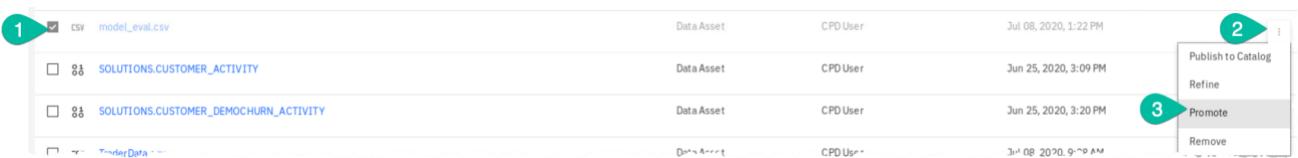


The screenshot shows a modal dialog titled 'Promote data asset'. It contains a message: 'Promote data asset to the associated space? Promoting makes this data asset and its dependencies available for use in deployments and apps.' At the bottom are two buttons: 'Cancel' and 'Promote to space', which is highlighted with a green checkmark.

__15. The data asset is now promoted to the associated Deployment Space.



__16. Repeat for the data asset [model_eval.csv](#).



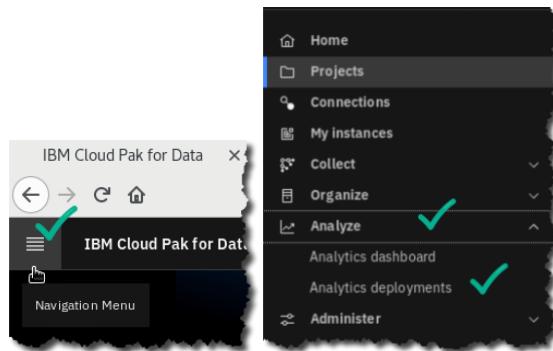
The screenshot shows the 'Assets' tab in the IBM Cloud Pak for Data interface. A search bar at the top has 'What assets are you looking for?' entered. Below it, a list of assets shows 'model_eval.csv' selected. To the right of the list is a context menu with three numbered steps: 1. 'Publish' (disabled), 2. 'Refine' (disabled), 3. 'Promote' (highlighted with a green circle and checkmark), and 4. 'Remove' (disabled).

Name	Type	Created by	Last modified
model_eval.csv	Data Asset	CPD User	Jul 08, 2020, 1:22 PM
SOLUTIONS.CUSTOMER_ACTIVITY	Data Asset	CPD User	Jun 25, 2020, 3:09 PM
SOLUTIONS.CUSTOMER_DEMOCHURN_ACTIVITY	Data Asset	CPD User	Jun 25, 2020, 3:20 PM
TraderData	Connection	CPD User	Jul 08, 2020, 9:29 AM

7.4.2 Working in the deployment space

Now you can go to the deployment space to work in it:

- 17. Click [Navigation menu](#) ⇒ [Analyze](#) ⇒ [Analytics deployments](#)



- 18. In this screen, note that you could create a new deployment space if you wanted to. You will not need to, however, because the one we need was created by the notebook.

Instead, click [deployment-space-analytics-workshop](#).

Name	Last updated	Associated project
deployment-space-analytics-project-workshop	Jul 8, 2020 1:17 PM	CPD Workshop Analytics Project

- 19. In the Assets section, there is one Model: [churn_risk_model](#) (created by the notebook) and there are two Data Assets (the ones you just promoted).

Name	Type	Software specification	Last modified
churn_risk_model	scikit-learn_0.22	scikit-learn_0.22-py3.6	Jul 8, 2020 1:02 PM

Name	Type	Last modified
model_eval.csv	Data Asset	Jul 8, 2020 1:27 PM
model_batch_score.csv	Data Asset	Jul 8, 2020 1:25 PM

- _20. In the Deployments section, there is one deployment: [churn_risk_model-deployment](#) (which uses model [churn_risk_model](#)).

The screenshot shows the 'Deployments' tab selected in the navigation bar. Below it, a table lists a single deployment entry:

Name	Type	Status	Asset	Last modified
churn_risk_model-deployment	Online	Deployed	churn_risk_model	Jul 8, 2020 1:02 PM



Developer

The notebook created all of these CPD assets – part of the code that did it is shown here.

```
# Specify a names for the space being created, the saved model and the model deployment
space_name = 'deployment-space-analytics-project-workshop'
model_name = 'churn_risk_model'
deployment_name = 'churn_risk_model-deployment'
```

- _21. Return to section [Assets](#) (click on it).

Click on the model: [churn_risk_model](#).

The screenshot shows the 'Assets' tab selected in the navigation bar. Below it, a table lists a single model entry:

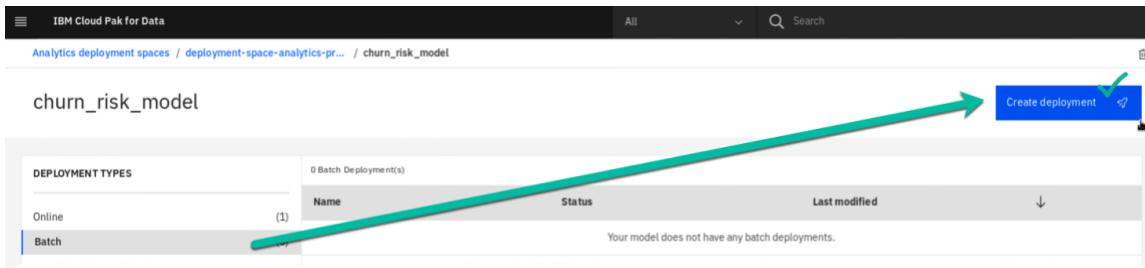
Name	Type	Software specification	Last modified
churn_risk_model	scikit-learn_0.22	scikit-learn_0.22-py3.6	Jul 8, 2020 1:02 PM

- _22. It shows one Online deployment (as seen previously in the Deployments section).

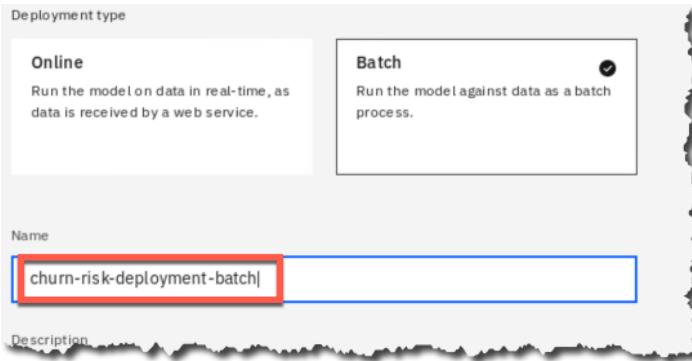
The screenshot shows the 'DEPLOYMENT TYPES' section. It indicates 1 Online Deployment(s). The table below shows the deployment details:

Name	Status
churn_risk_mod...	Deployed

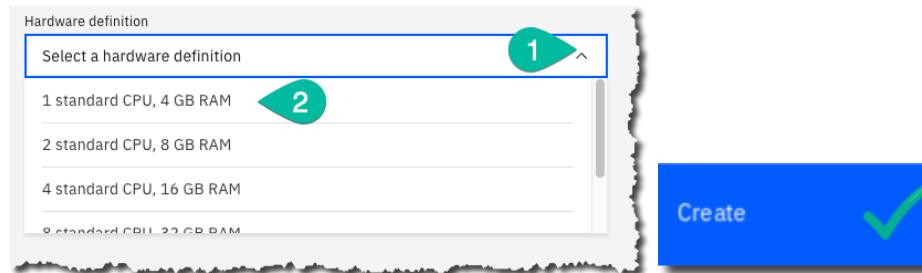
_23. Click **Batch** ⇒ **Create Deployment**



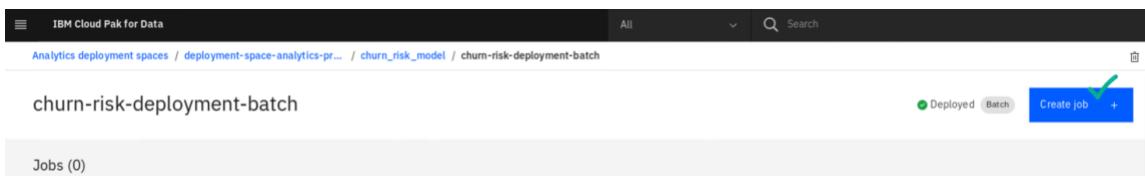
_24. Fill in name: **churn-risk-model-deployment-batch**.



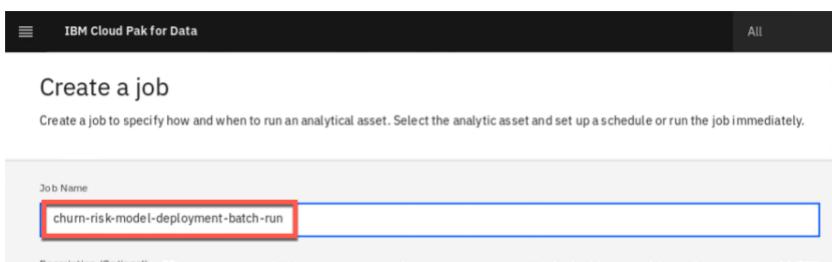
_25. Scroll down and find Select a Hardware definition. Select **1 standard CPU,4 GB RAM**, and click **Create**.



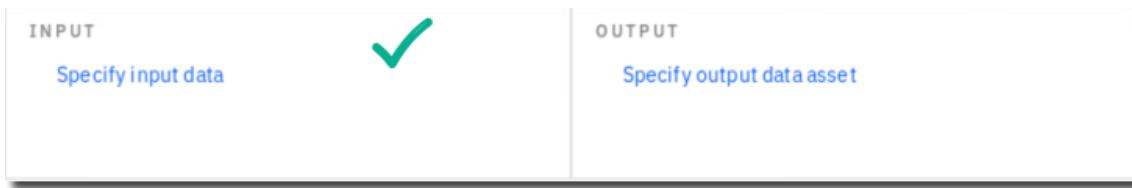
_26. You will be taken to the Jobs screen. Click **Create job +**.



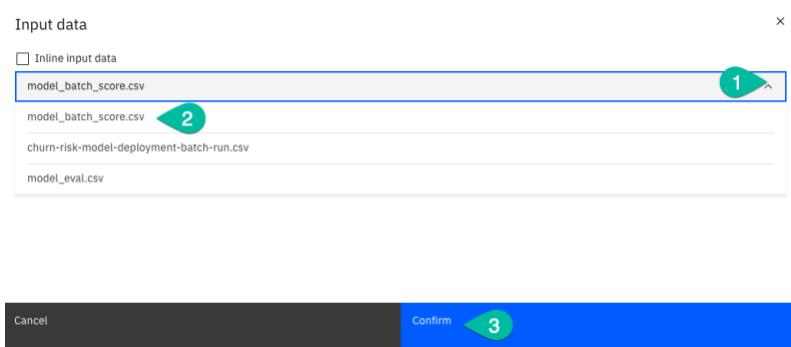
_27. Fill in the Job Name: **churn-risk-model-deployment-batch-run**.



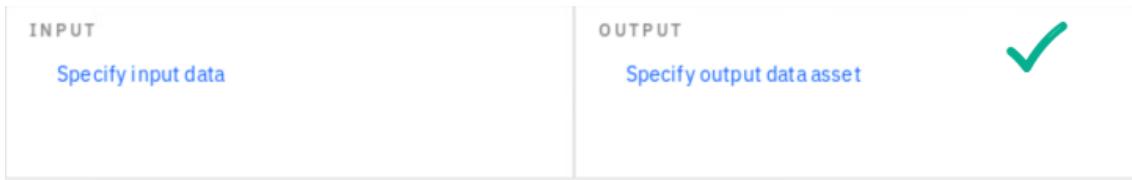
- __28. Click on sections for **INPUT – Specify input data.**



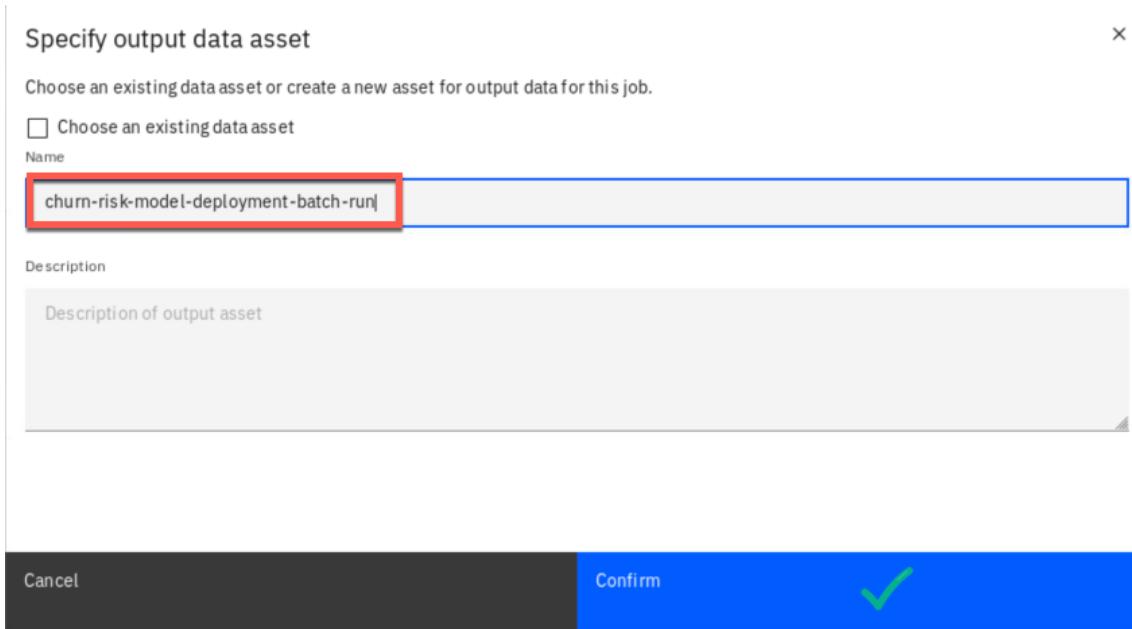
- __29. Set the input value to **model_batch_score.csv** and Confirm.



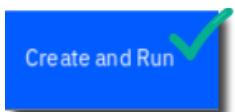
- __30. Select the **OUTPUT – Specify output data asset.**



- __31. In this case, we will output to a new dataset. Enter **churn-risk-model-deployment-batch-run** for name and **Confirm**.



__32. Select **Create and Run**.



__33. The job will queue, run, and complete.

Note: You may need to refresh your browser to see the “Completed” result. You may have to wait around 5-10 minutes for the job to complete.



__34. Click on the **Start Time value** when Completed.

Runs (1)

Start Time	Status
Jul 08, 2020 1:47:05 PM	Completed

__35. Review the job – select **Show More** at bottom right.

Status	Duration (s)	Started by	Environment
Completed	421	CPD User	1 standard CPU, 4 GB RAM

Log tail | Total 35 lines [Download log](#)

```
{
  "deployment": {
    "href": "/v4/deployments/c4f2a00-3960-4ac5-abf1-3c47b35cfcd9",
    "id": "c4f2a00-3960-4ac5-abf1-3c47b35cfcd9"
  },
  "hardware_spec": {
    "id": "f3beb7d-0a75-41bc-bd48-e931a28cc4c5"
  },
  "scoring": {
    "input_data_references": [
      {
        "location": {
          "href": "/v2/assets/920nf7e5-a3dd-4d41-adce-d841de2fedec7space_id=7e0905f1-a608-4d59-a2a7-fc169c9728ac"
        },
        "type": "data_asset",
        "connection": {}
      }
    ],
    "output_data_reference": {
      "location": {
        "href": "/v2/assets/55456e18-0f17-4ee7-804f-722277f2ee77space_id=7e0905f1-a608-4d59-a2a7-fc169c9728ac",
        "name": "churn-risk-model-deployment-batch-run"
      },
      "type": "data_asset",
      "connection": {}
    }
  },
  "status": {
    "completed_at": "2020-07-08T17:54:07.473825Z",
    "running_at": "2020-07-08T17:54:09.182138Z",
    "state": "Completed"
  },
  "space_id": "7e0905f1-a608-4d59-a2a7-fc169c9728ac"
}
```

Show less ^

__36. After reviewing the job, click on the breadcrumb to return to the deployment space.

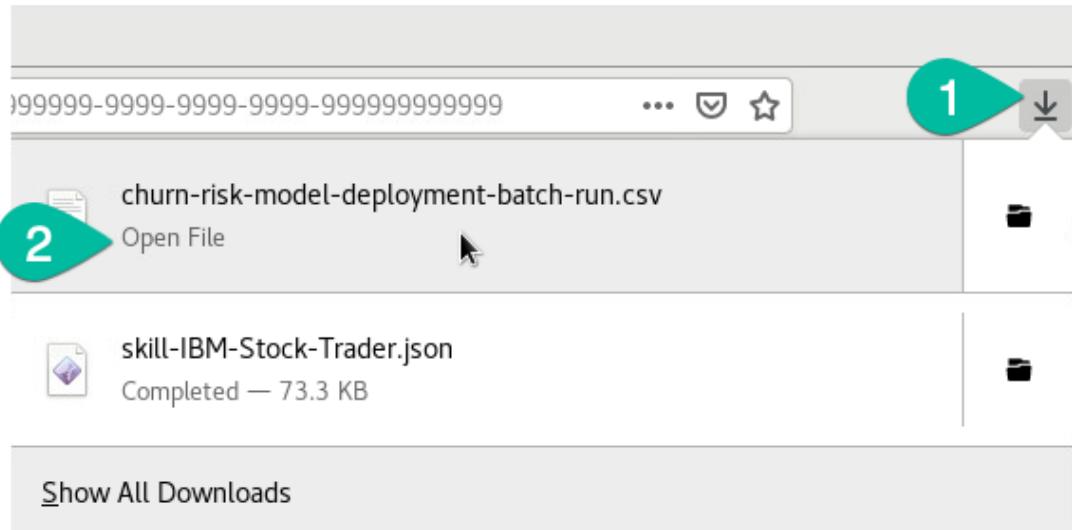
- __37. In the section **Assets**, find the new output file under **Data Assets**.
 The file name is **churn-risk-model-deployment-batch-run.csv**.
- __38. On that file line, click **Hover** to uncover the download (down arrow) icon : and **Click it**.

Name	Type	Software specification	Last modified
churn_risk_model	scikit-learn_0.22	scikit-learn_0.22-py3.6	Jul 8, 2020 1:02 PM
model_eval.csv			Jul 8, 2020 1:27 PM
model_batch_score.csv			Jul 8, 2020 1:25 PM

- __39. Save your File to the Documents folder and Save.

Name	Size	Modified
churn-risk-model-deployment-batch-run.csv	13.9 kB	14:53
CustomerChurnDemographics-CA_Lab.csv	277.5 kB	24 Jun 11:15

- __40. Review the downloaded file by clicking on the download (down arrow) icon in your Firefox browser, then select the file.



Note: this file is located on the operating system under /home/ibmuser/Downloads.

- __41. The batch run output is available for review.

```
prediction,probability
1, "[0.0, 0.99, 0.01]"
0, "[0.79, 0.02, 0.19]"
0, "[0.72, 0.01, 0.27]"
0, "[0.97, 0.01, 0.02]"
0, "[0.97, 0.0, 0.03]"
2, "[0.04, 0.0, 0.96]"
2, "[0.03, 0.0, 0.97]"
1, "[0.01, 0.9408333333333333, 0.04916666666666664]"
1, "[0.0, 1.0, 0.0]"
```

- __42. Close the file edit windows after reviewing.

- __43. Click [Deployments](#).

Click Online deployment: [churn_risk_model_deployment](#).

Name	Type	Status
churn-risk-deployment-batch	Batch	Deployed
churn_risk_model-deployment	Online	Deployed

- 44. Under section [API reference](#), review the Endpoint for this model. This is what your applications can use to call this model.

The screenshot shows a user interface for an API reference. At the top, there are two tabs: "API reference" (which is selected, indicated by a green checkmark) and "Test". Below the tabs, there is a section titled "Direct link" under "Endpoint". A red box highlights the URL: `https://cp4d-cpd-cp4d.apps.clusterw9/v4/deployments/93e932b4-0d6d-40e3-981d-b1a18b9813f6/predictions`.

- 45. [Review the code snippets](#) for each language that can aid a Developer to easily embed the model into an application.

The screenshot shows a "Code snippets" section for Java. It includes a "curl" example and Java code. The Java code is as follows:

```
# TODO: manually define and pass values to be scored below
curl -X POST --header 'Content-Type: application/json' --header 'Accept: application/json' --header "Authorization: Bearer $WML_AUTH_TOKEN" -d '{ "array_of_input_fields": [ "field1", "field2", "field3" ] }' https://cp4d-cpd-cp4d.apps.clusterw9/v4/deployments/93e932b4-0d6d-40e3-981d-b1a18b9813f6/predictions
```

The screenshot shows a "Code snippets" section for Java. The Java tab is selected, indicated by a green checkmark. The Java code is identical to the one shown in the previous screenshot.

The screenshot shows a "Code snippets" section for Python. The Python tab is selected, indicated by a green checkmark. The Python code is as follows:

```
import urllib3, requests, json

# NOTE: you must construct mltoken based on provided documentation
header = {'Content-Type': 'application/json', 'Authorization': 'Bearer ' + mltoken}

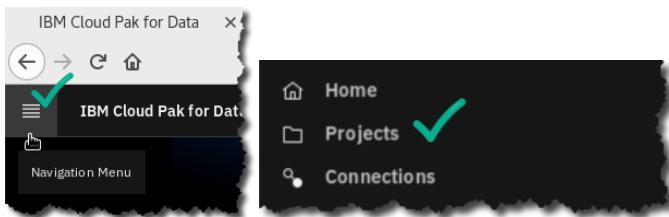
# NOTE: manually define and pass the array(s) of values to be scored in the next line
payload_scoring = {"input_data": [{"fields": [array_of_input_fields], "values": [array_of_values_to_be_scored, another_array_of_values]}]}

response_scoring = requests.post('https://cp4d-cpd-cp4d.apps.clusterw9/v4/deployments/93e932b4-0d6d-40e3-981d-b1a18b9813f6/predictions', headers=header, json=payload_scoring)
print("Scoring response")
print(json.loads(response_scoring.text))
```

7.5 Deploying and testing the AutoAI model

7.5.1 Deploying the AutoAI model

- 46. In the CPD web client, click the [Navigation Menu](#) ⇒ [Projects](#).



- 47. Select the project: [CPD Workshop Analytics Project](#).

Name	Project type
CPD Workshop Analytics Project ✓	Analytics
CPD Workshop Transform Project	Data transform

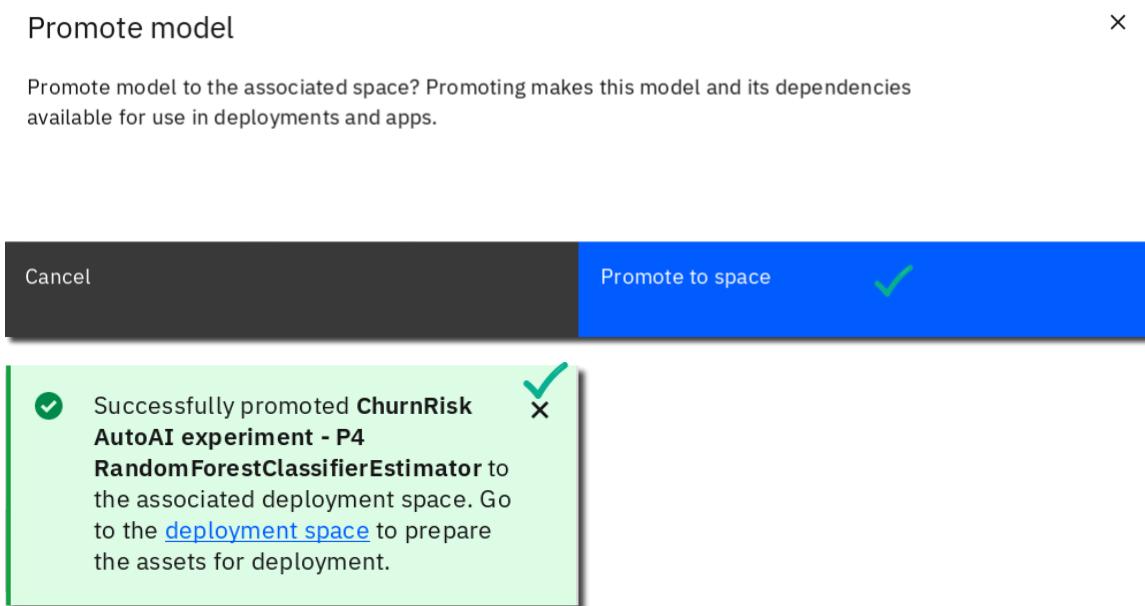
- 48. Under [Assets](#), scroll down to [Models](#).

On the model created from the AutoAI experiment (ChurnRisk AutoAI experiment – P4 RandomForestClassifierEstimator), click on [ellipses :](#) then [Promote](#) and [Promote to space](#).

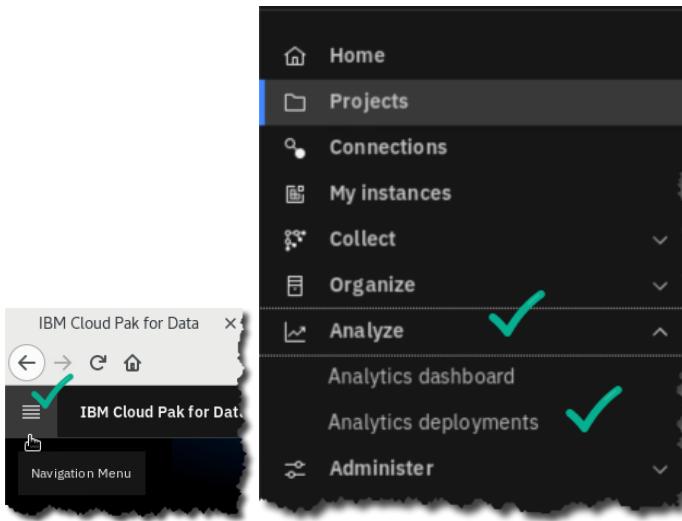
 A screenshot of the CPD web client showing the 'Models' section. At the top, there's a navigation bar with 'My projects / CPD Workshop Analytics Project' and tabs for 'Overview', 'Assets' (with a green checkmark), and 'Environments'. Below is a 'Project information' section with a 'Project name' field. The main area shows a table of models. One row is selected, showing details: 'Name' is 'ChurnRisk AutoAI experiment - P4 RandomForestClassifierEstimator', 'Type' is 'wml-hybrid_0.1', 'Software specification' is 'hybrid_0.1', and 'Last modified' is 'Jul 08, 2020'. To the right of the table is a 'Promote' button with a green arrow icon. A context menu is open over this button, with item 1 labeled 'Promote' and item 2 labeled 'Delete'. Below the table is a 'Functions' section. At the bottom, a modal dialog titled 'Promote model' contains the message: 'This model was already promoted to the associated space. Promoting it again will create a copy of the model in the space. Are you sure you want to promote this model and its dependencies?'. It has 'Cancel' and 'Promote to space' buttons, with the latter being highlighted in blue.

Name	Type	Software specification	Last modified
ChurnRisk AutoAI experiment - P4 RandomForestClassifierEstimator	wml-hybrid_0.1	hybrid_0.1	Jul 08, 2020

__49. This promotes the model to the associated Deployment Space



__50. Click [Navigation menu](#) \Rightarrow [Analyze](#) \Rightarrow [Analytics deployments](#).



__51. Click [deployment-space-analytics-project-workshop](#).

IBM Cloud Pak for Data All Search

Analytics deployment spaces

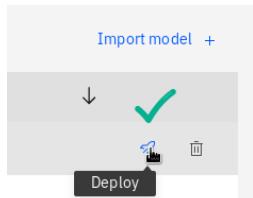
Which deployment space are you looking for?

Name	Last updated	Associated project
deployment-space-analytics-project... ✓	Jul 8, 2020 1:17 PM	CPD Workshop Analytics Pro

__52. Under section **Assets**, find the model you just promoted from your project. (*Don't click on it.*)

Name	Type	Software specification
ChurnRisk AutoAI experiment - P4 RandomForestClassifierEstimator	wml-hybrid_0.1	hybrid_0.1
No description provided.	scikit-learn_0.22	scikit-learn_0.22-py3.6

__53. At the end of the row for this model, click **Hover** and click on the **Deploy** (rocket) icon.



__54. Choose **Online**.

_55. Name: ChurnRisk AutoAI Experiment deployment a Create.



_56. You will see the online model being deployed.

1 Online Deployment(s)	
Name	Status
ChurnRisk Auto...	C In progress

_57. While you are waiting for it to deploy, review the model [Schema](#).

Analytics deployment spaces / deployment-space-analytics-pr... / ChurnRisk AutoAI experiment - ...		All
ChurnRisk AutoAI experiment - P4 RandomForestClassifierEstimator		
Deployments Schema ✓		
Input (1)	Column	Type
	ADDRESS_1	object
	ADDRESS_2	float64
	AGE	int64
	AGE_GROUP	object
	CHILDREN	int64
	CITY	object

__58. Select **Deployments**. The Online model is now deployed in short time with zero coding.

DEPLOYMENT TYPES				1 Online Deployment(s)		
Online	Batch	(1)	(0)	Name	Status	Last modified
				ChurnRisk Auto...	Deployed	Jul 8, 2020 4:12 PM

__59. Click on the deployment.

1 Online Deployment(s)	
Name	Status
ChurnRisk Auto...	Deployed

7.5.2 Testing the AutoAI model

__60. Click on tab **Test**.

Fill in the first three input data fields as shown below.

Enter input data	
ID	1
AGE_GROUP	Adult
GENDER	Female
STATUS	

__61. Scroll down and click on Predict.

(Note: filling in more fields would make the prediction more accurate, but we only did three fields to show you the process.)



__62. Your machine learning prediction and probability is returned.

Result

```
0 {
1   "predictions": [
2     {
3       "fields": [
4         "prediction",
5         "probability"
6       ],
7       "values": [
8         [
9           "High",
10          [
11            0.7005350837370795,
12            0.03857883761495348,
13            0.2608860786479671
14          ]
15        ]
16      ]
17    }
18  ]
19 }
```

7.6 Lab conclusion

Today, organizations need better, faster and more integrated methods for model deployment that can include online, batch and virtual testing. The Cloud Pak for Data deploy capabilities help serve this need so they can more easily infuse AI into their applications.

In our scenario, the Trade Co. developer used the platform to easily deploy and test models.



** End of Lab 07 - Deploy

Lab by Burt Vialpando and Kent Rubin, IBM

Lab 08 Infuse: Watson OpenScale

8.1 Lab overview

Business value is achieved when a model is infused into an application, regularly monitored and updated. But most organizations need help with that. They want to know when the accuracy of a model begins to decline so it can automatically be corrected. They also want to know how their models make decisions, and make sure they eliminate bias.

The first part of this lab will leave the Trade Co. use case to focus on a Credit Risk use case demo that is built into Watson OpenScale. It shows how OpenScale can detect and mitigate bias while providing explainable output for the model's decision.



In this [Infuse](#) lab you will learn how Watson OpenScale assists you in this.

8.2 Persona represented in this lab

The [Developer](#) persona is the likely role to perform the various [Infuse](#) tasks in this lab. However, the Data Scientist persona could perform these tasks as well.

Persona (Role)	Capabilities
 Developer	Developers create and maintain the end-user applications that utilize the output from all the other personas on the CPD platform.

8.3 Logging into the CPD web client (if you have not already done so)

- __1. If you are starting this lab stand-alone (without going through previous labs) do the following:
- __2. Double-click the desktop icon: [Cloud Pak for Data Web Client](#).



- __3. The CPD web client GUI displays as shown. Use [cpduser](#) and [cpdaccess](#) for the *Username* and *Password* and click [Sign in](#).



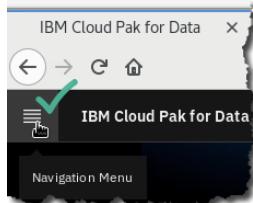
8.4 Credit Risk built-in demo

8.4.1 Run OpenScale Auto setup

Watson OpenScale provides a quick setup utility that will automatically set up a data mart, create a model, and record seven days' worth of measurements into the OpenScale monitors. This creates the Credit Risk demo that we will go through in this lab.

- 4. Start by first retrieving the Db2 Advanced Edition connection information.

Click on the [Navigation Menu](#) (“hamburger” icon) at the top left of the screen.



- 5. Click [Collect](#) ⇒ [My Data](#) ⇒ [Databases](#) ⇒ [Db2 Advanced Edition](#) ⇒ [ellipsis :](#) ⇒ [Details](#).

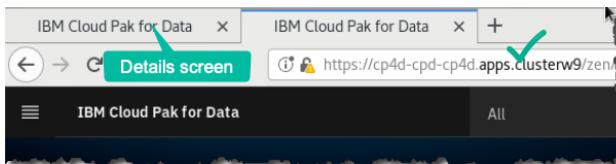
- 6. Scroll to find the Access Information. Notice Username, Password, Host and Database.

Access information	
Username	<input type="text" value="user1001"/>
Password	<input type="password" value="REDACTED"/>
JDBC Connection URL	<input type="text" value="idbc:db2://worker5.clusterw9:32030//BLUBB"/>

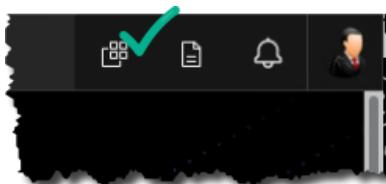
- 7. If you are using the remote desktop, click on the icon for [Cloud Pak for Data Web Client](#); otherwise, duplicate your existing browser tab.



- 8. There should now be two web client browser tabs open. The first tab has the [Details Screen](#) for the Db2 Advanced Edition just reviewed. The second is the new web client browser page just launched.

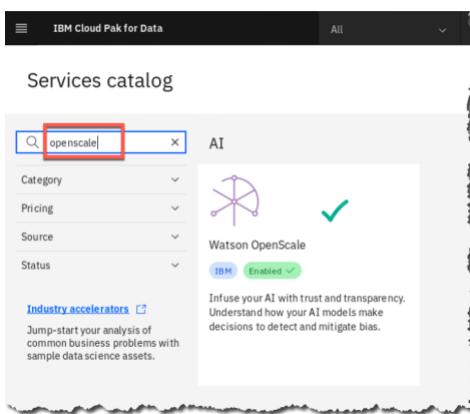


- 9. From the second (new) web client browser tab, click the icon for [Services](#) (at the top right corner of the screen)

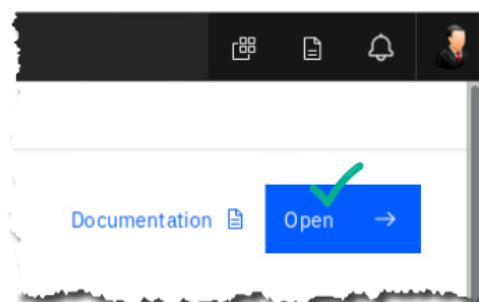


- 10. In the search window, type [openscale](#).

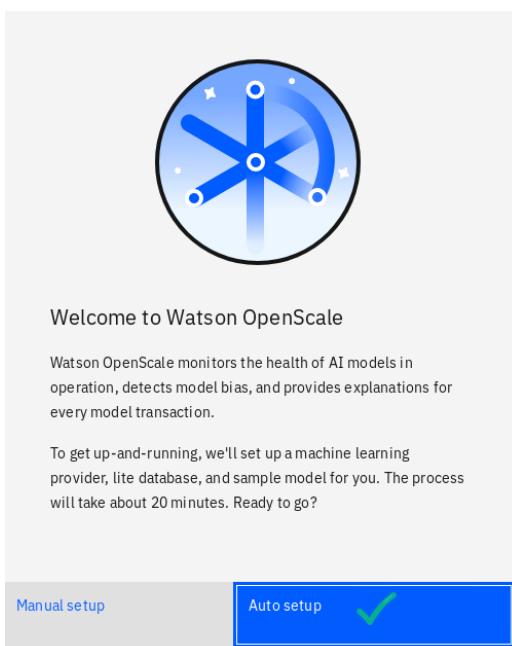
Click on tile [Watson OpenScale](#).



- 11. Click [Open](#).



__12. Click **Auto setup** ⇨ **Next.**



Connect to Watson Machine Learning

Watson OpenScale will deploy the sample model in your Watson Machine Learning service. Follow the link to learn how to obtain your service credentials. [Learn more.](#)

Use Watson Machine Learning instance on the local environment

Service Credentials

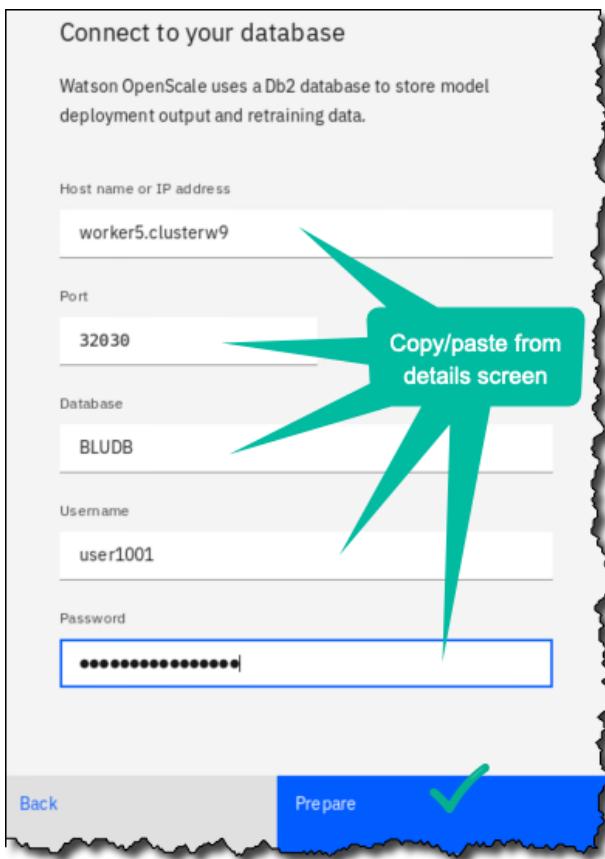
```
{
  "apikey": "<<APIKEY>>",
  "iam_apikey_description": "<<APIKEY_DESCRIPTION>>",
  "iam_apikey_name": "<<APIKEY_NAME>>",
  "iam_role_crn": "<<ROLE_NAME>>",
  "iam_serviceid_crn": "<<SERVICE_ID_CRN>>",
  "instance_id": "<<INSTANCE_ID>>",
  "password": "<<PASSWORD>>",
  "url": "https://us-south.ml.cloud.ibm.com",
  "username": "<<USERNAME>>"
}
```

Back

Next

__13. Fill in the connection information here from the **Db2 Advanced Edition Details** screen you have opened in the other tab. (Note: Your hostname may vary.)

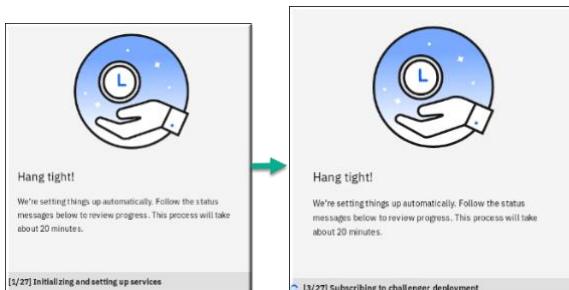
When all filed in, click **Prepare**.



The best technique to use to copy and paste the password is:

1. Highlight the password
2. Right click - Copy
3. Right click - Paste

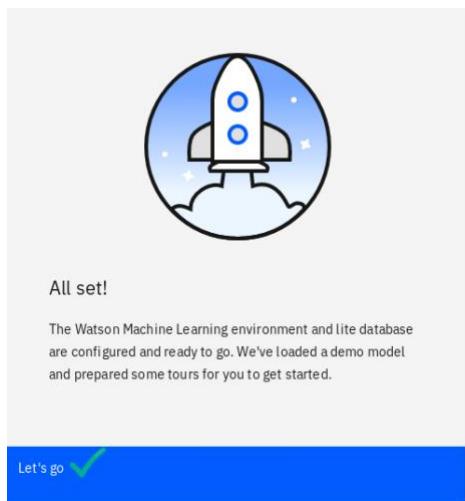
14. The Auto setup will take 20 minutes or so to run. During this time, read about the OpenScale monitors and the Credit Risk scenario it is creating as it goes through the 27-step setup.



 Developer	<p>For this demo, you will be monitoring a model that attempts to predict credit risk based on demographic data as well as credit history, residence information, age, employment status, and more. The scenario and model use synthetic data based on the UCI German Credit dataset.</p> <pre> graph LR CL[Credit Lender] --> DS[Data Science Team] DS --> RM[Risk Model] RM --> WOS[Watson OpenScale] WOS --> LA[Loan Applicants] WOS <--> DS WOS <--> CL </pre>
---------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

 Developer	<p>Credit lenders can monitor risk models for performance, bias and explainability to limit risk exposure from regulations and create more fair and explainable outcomes for customers.</p> <p>Traditional lenders are under pressure to expand their digital portfolio of financial services to a larger and more diverse audience, which requires a new approach to credit risk modeling. Their data science teams currently rely on standard modeling techniques - like decision trees and logistic regression - which work well for moderate datasets and make recommendations that can be easily explained. This satisfies regulatory requirements that credit lending decisions must be transparent and explainable.</p> <p>To provide credit access to a wider and riskier population, applicant credit histories must expand beyond traditional credit, like mortgages and car loans, to alternate credit sources like utility and mobile phone plan payment histories, plus education and job titles. These new data sources offer promise, but also introduce risk by increasing the likelihood of unexpected correlations which introduce bias based on an applicant's age, gender, or other personal traits.</p> <p>The data science techniques most suited to these diverse datasets, such as gradient boosted trees and neural networks, can generate highly accurate risk models, but at a cost. Such "black box" models generate opaque predictions that must somehow become transparent, to ensure regulatory approval such as Article 22 of the General Data Protection Regulation (GDPR), or the federal Fair Credit Reporting Act (FCRA) managed by the Consumer Financial Protection Bureau.</p>
---------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

- __15. When the Auto setup has completed, click [Let's go](#). You need not go through the Guided Tour, but you can if you have the time.



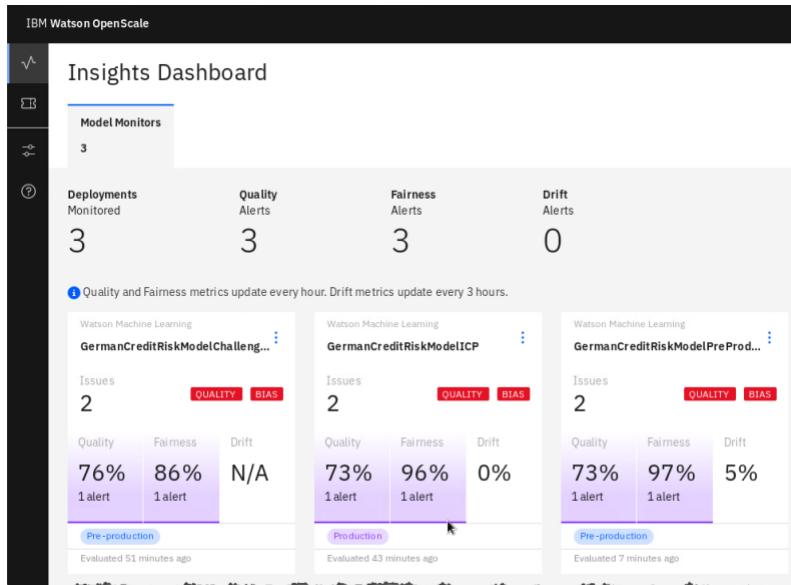
8.4.2 Explore the fairness and quality monitors

Watson OpenScale provides two types of monitors: application monitors, and model monitors.

We will begin with the model monitors. The model monitors section of the Insights Dashboard provides an overview of all the models being monitored by OpenScale.

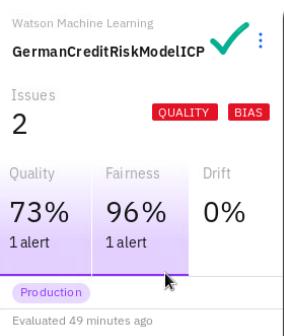
You can monitor models hosted with Watson Machine Learning on public or private clouds, Microsoft Azure, Amazon SageMaker, or custom models that provide JSON prediction output.

- __16. The dashboard shows how many deployments are currently being monitored, as well as an overview of the alerts from those models. Below, each model is represented by a tile showing the machine learning provider and alerts for that model. The left side of the screen shows all of the active monitors for the models, divided into sections for [Fairness](#), [Quality](#), and [Drift](#).

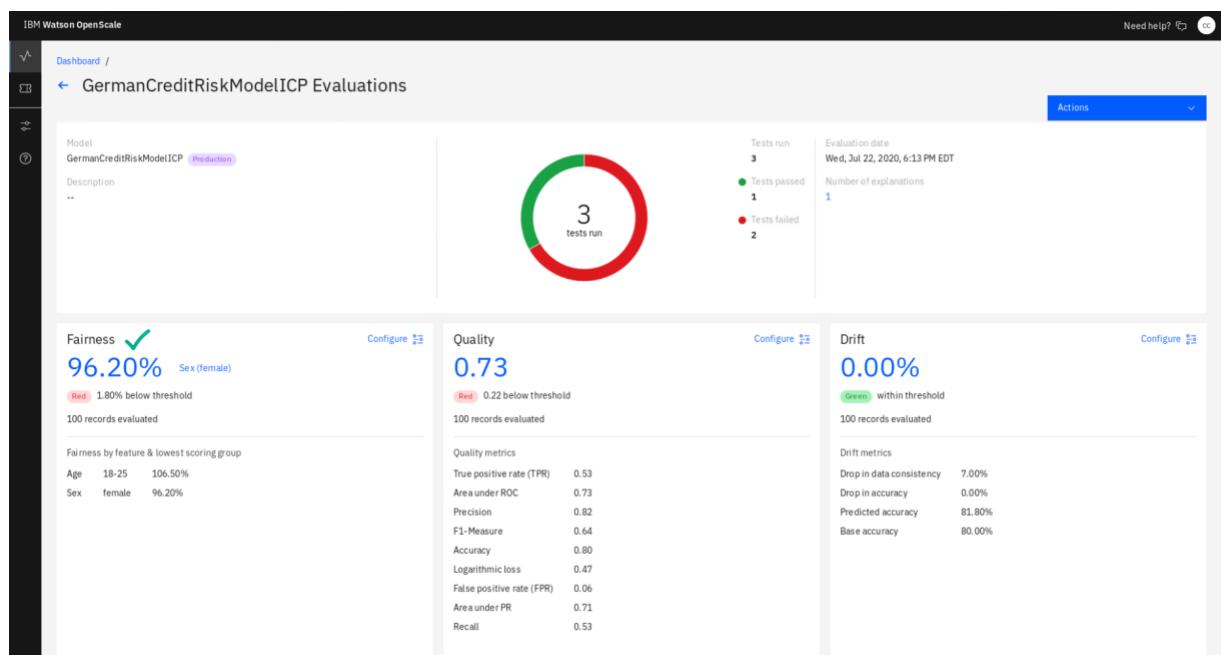


NOTE: Your actual values may vary slightly from what is shown in these lab screen shots as the data source has some randomness built into it.

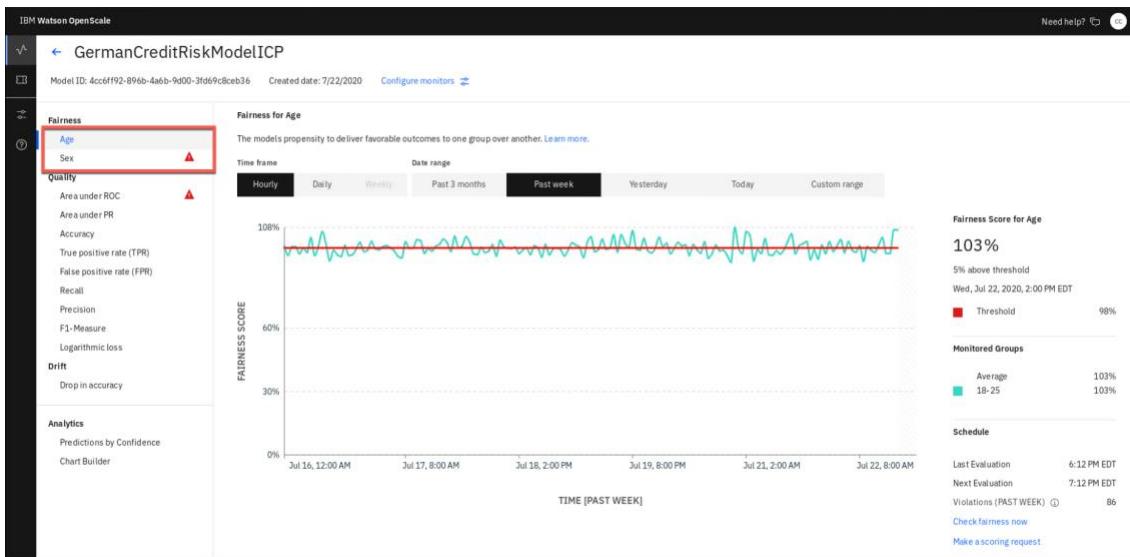
__17. Click on the tile [GermanCreditRiskModelICP](#).



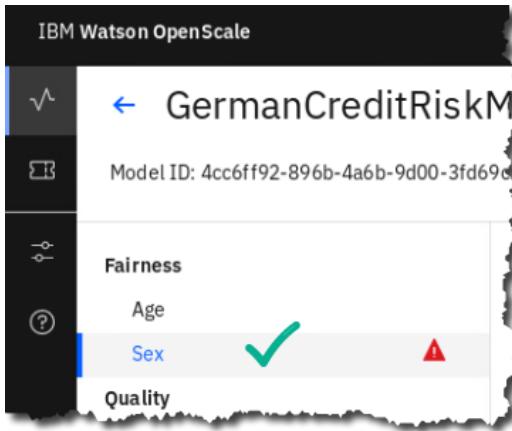
__18. You are presented with a dashboard of your model evaluations. Select the [Fairness percentage value link](#) in the Fairness tile.



- __19. In the Fairness section, you can see that we have chosen to model two features, Age and Sex, for fairness. Additionally, you can see that we have an alert for the Sex feature.

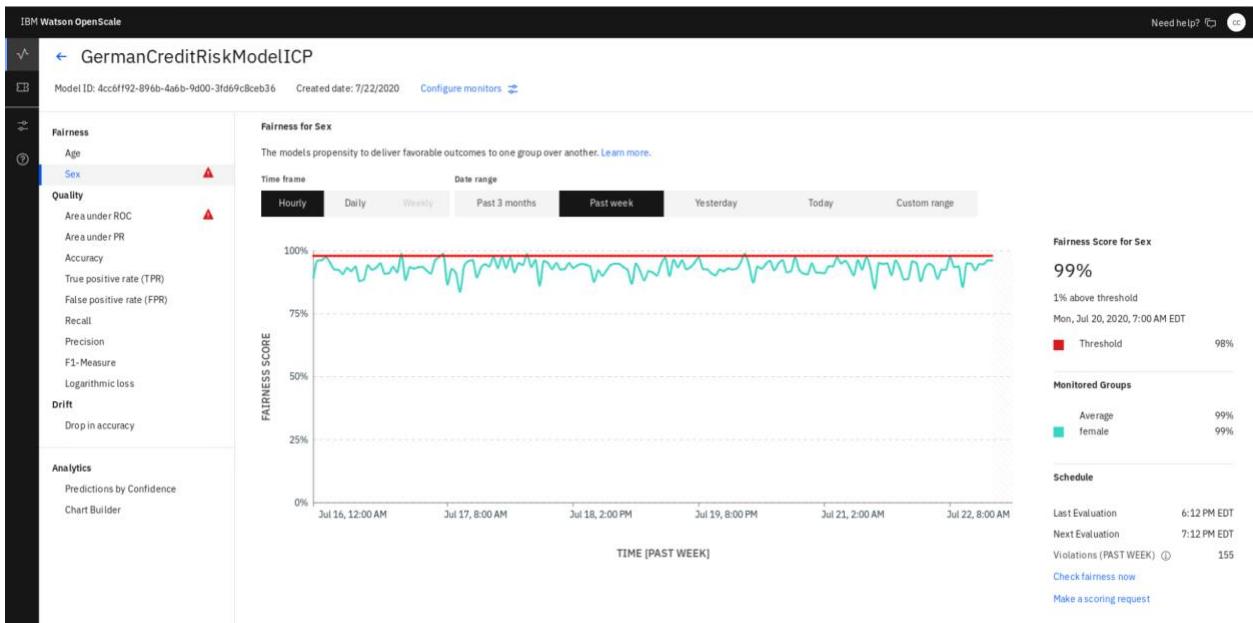


- __20. Click on the fairness monitor [Sex](#).

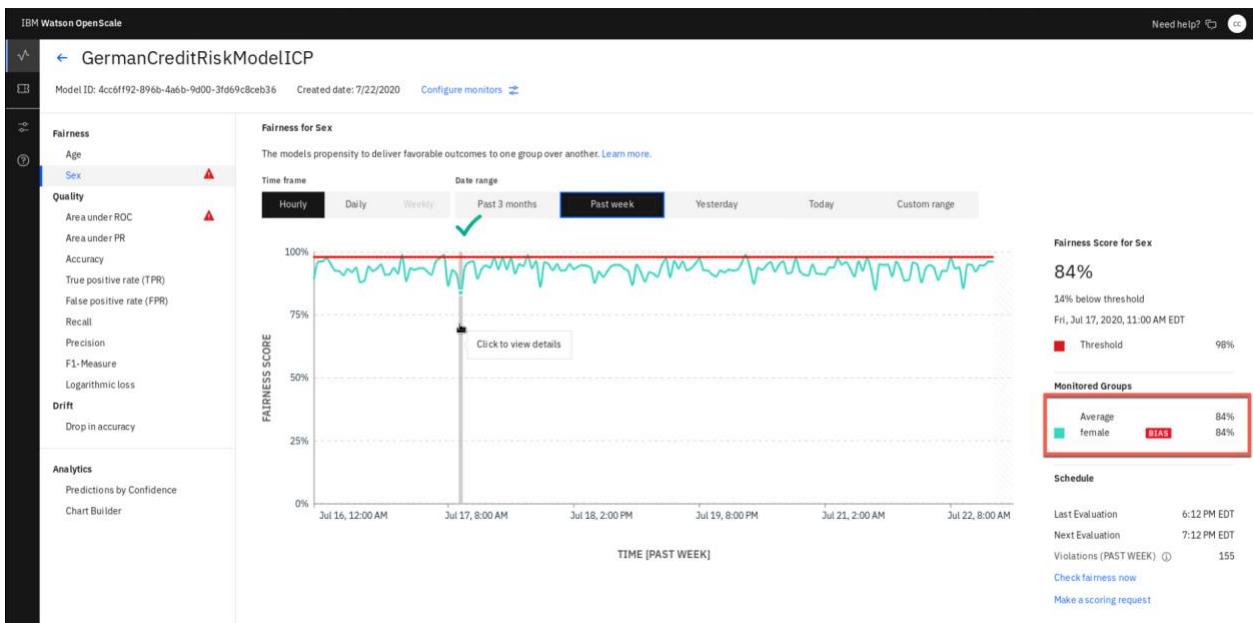


- __21. Note the time frame above the graph. We are looking at hourly data from the past week, but the time frame can be adjusted as necessary. The graph shows the fairness score for females as a green line. The threshold we have set for an alert is shown by the red line.

As you can see, the fairness score has dropped beneath the threshold consistently over the past week, alerting us to a potential unfair bias issue with the model. As you **move your cursor inside that date range**, values are updated to show the Fairness Score for that date and time.



- __22. Move your mouse on the chart to the point where fairness score was at its lowest (**84%**)



This screen shows us fairness details for this particular time period. OpenScale calculates its fairness score using a combination of actual predictions (payload data) and perturbed data,

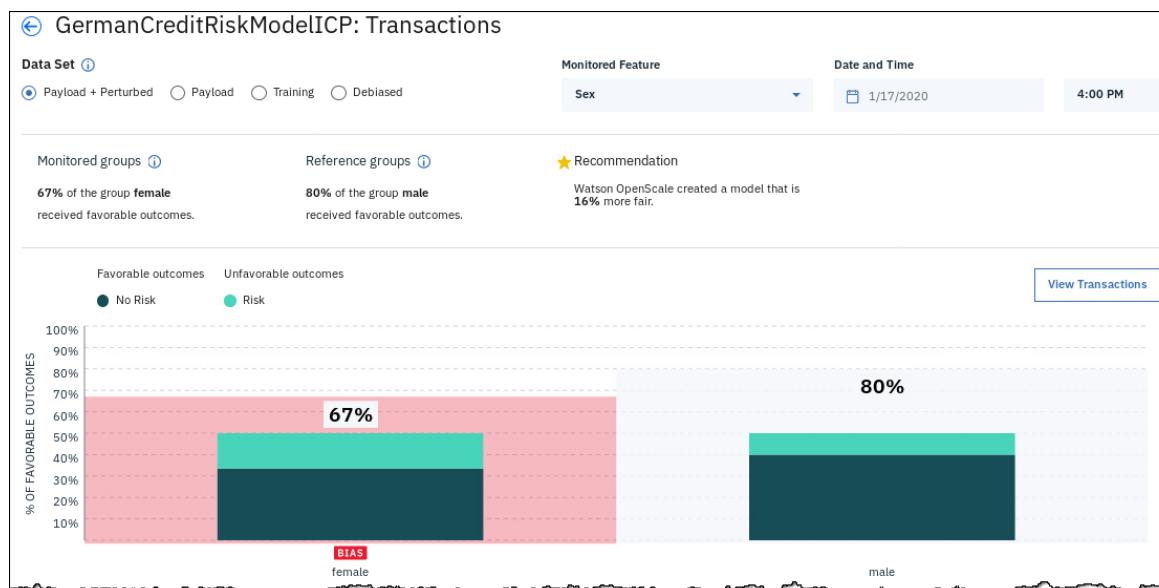
generated when the prediction probability is close to 50%. OpenScale will flip the monitored feature to see how it affects the prediction outcome.

- __23. Click on that lowest score line when it says: [Click to view details](#).



- __24. The fairness score is reached by dividing the percentage of positive outcomes for the monitored group (females, 67%) by the percentage of positive outcomes for the reference group (males, 80%).

The graph shows the breakdown of positive and negative predictions for our two groups. You can use the radio button at the top to view payload and perturbed data, actual payload (prediction) data, and training data.



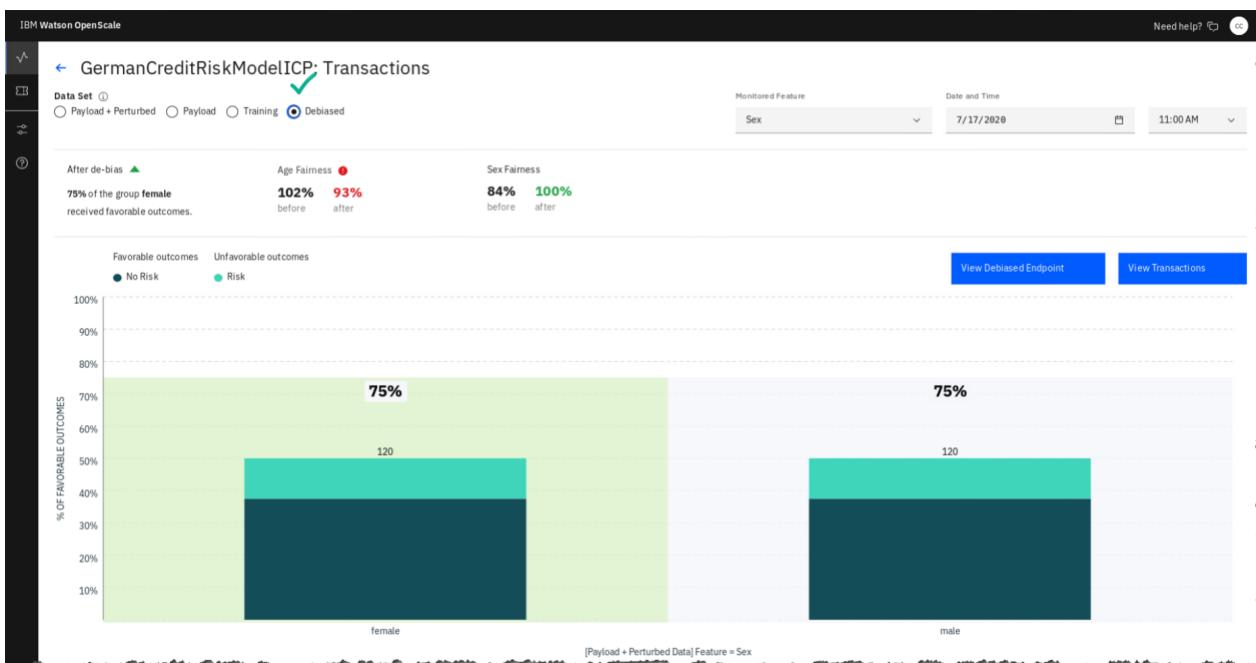
- __25. Click data set (radio button) [Training](#) to view the training data breakdown.

As you can see, our training data had significantly more records for males than for females, which may be a potential source of our unfair bias against females.

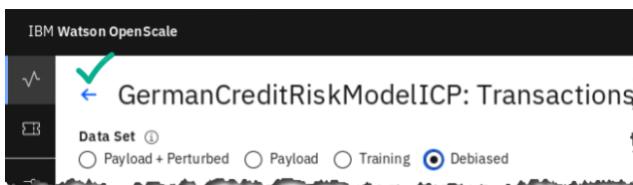


_26. Click the data set (button) **Debiased**.

OpenScale can create a sort of “corrective lens” to reduce or remove unfair model bias. It does this by training another model to predict when an outcome of the production model is likely to be unfairly biased, flipping the feature value from the monitored group (female) to the reference group (male) and returning this prediction. On this screen, you can see how using this model will affect the fairness scores for other features.



_27. Click the **back arrow** to return to the model dashboard.



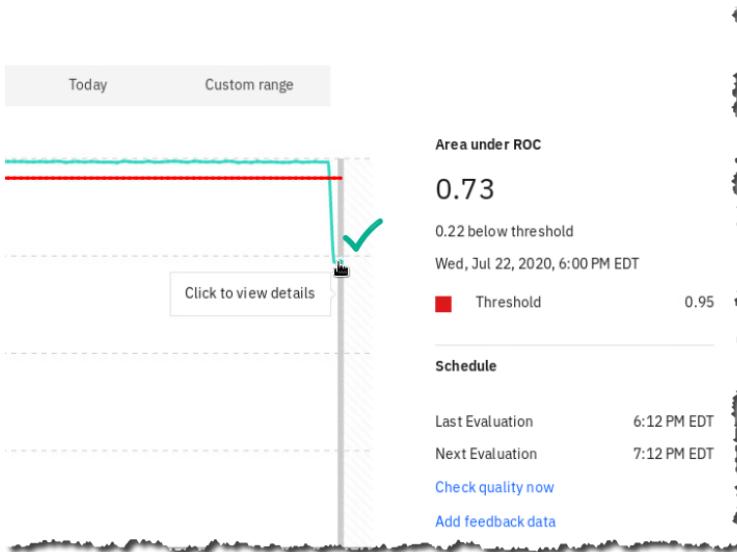
_28. Click the monitor **Area** under **ROC** in the section **Quality**.

- 29. OpenScale provides several different quality measurements. For our binary classification model, *Area under ROC* provides the best standard for model quality. These scores are generated by providing ground truth feedback data to the model, either via CSV upload or using a RESTful endpoint provided by OpenScale.

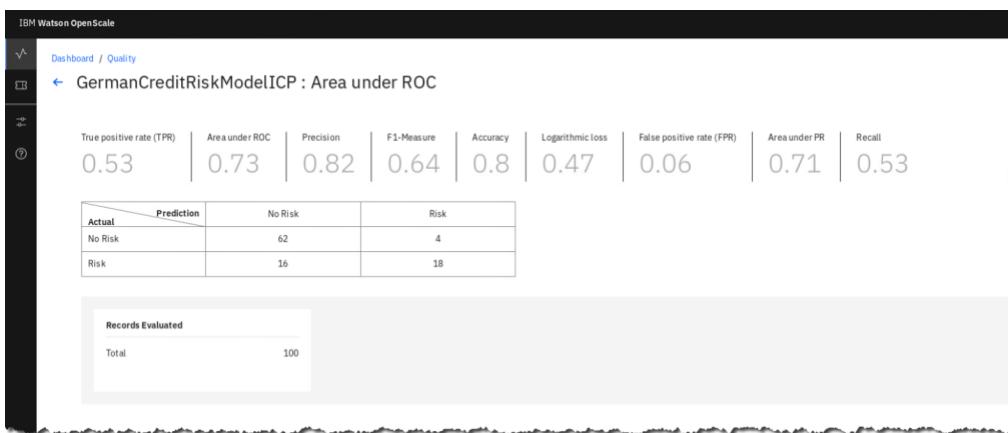


As with the fairness monitor, the chart in the middle of the screen shows model performance over an adjustable time window, with the relevant measurement shown as the green line and the alert threshold represented by the red line. As you can see, our model quality has consistently been above the threshold until the most recent measurement, represented by the right-most portion of the chart.

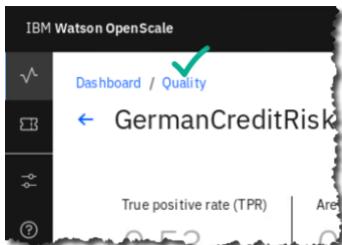
- 30. Click on the right-most portion of the chart, where Area under ROC drops to **0.73**.



- __31. Here we can see a further breakdown of the feedback data and the various accuracy scores, and the number of feedback records evaluated.



- __32. Click on the [Quality](#) link to return to the GermanCreditRiskModelICP [Dashboard](#).

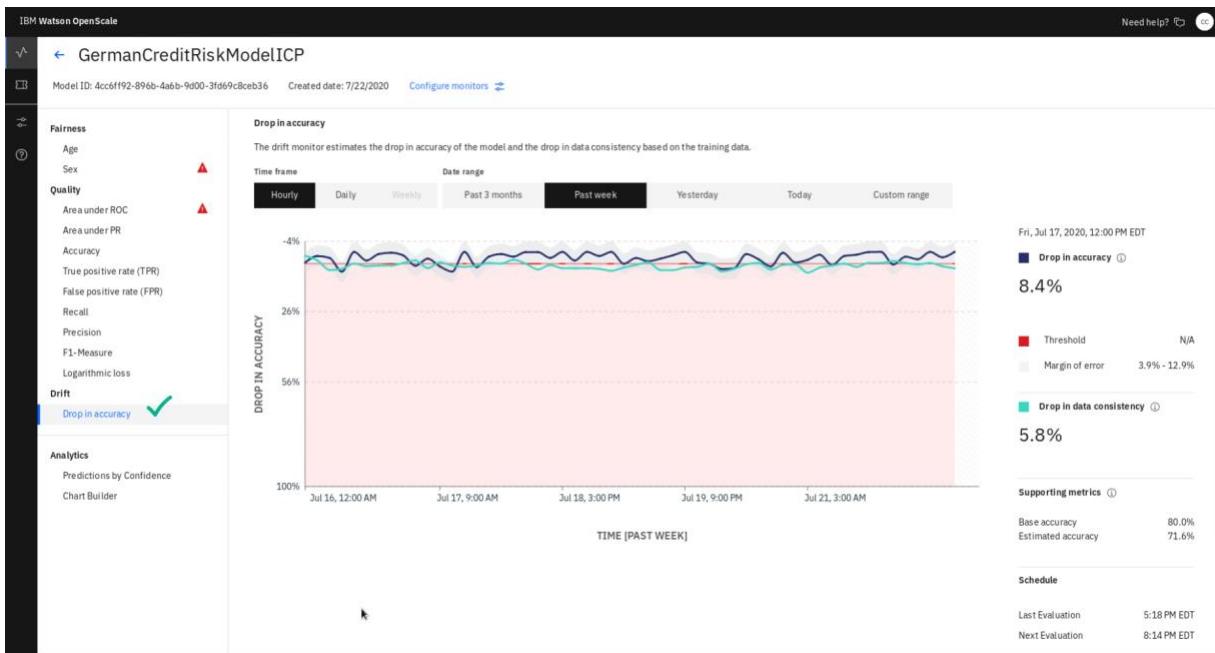


- __33. The OpenScale drift monitor is a separate linear regression drift model, trained to determine which types of data the production model struggles to correctly predict. This drift model allows OpenScale to forecast potentially costly drops in model accuracy without requiring additional feedback data.

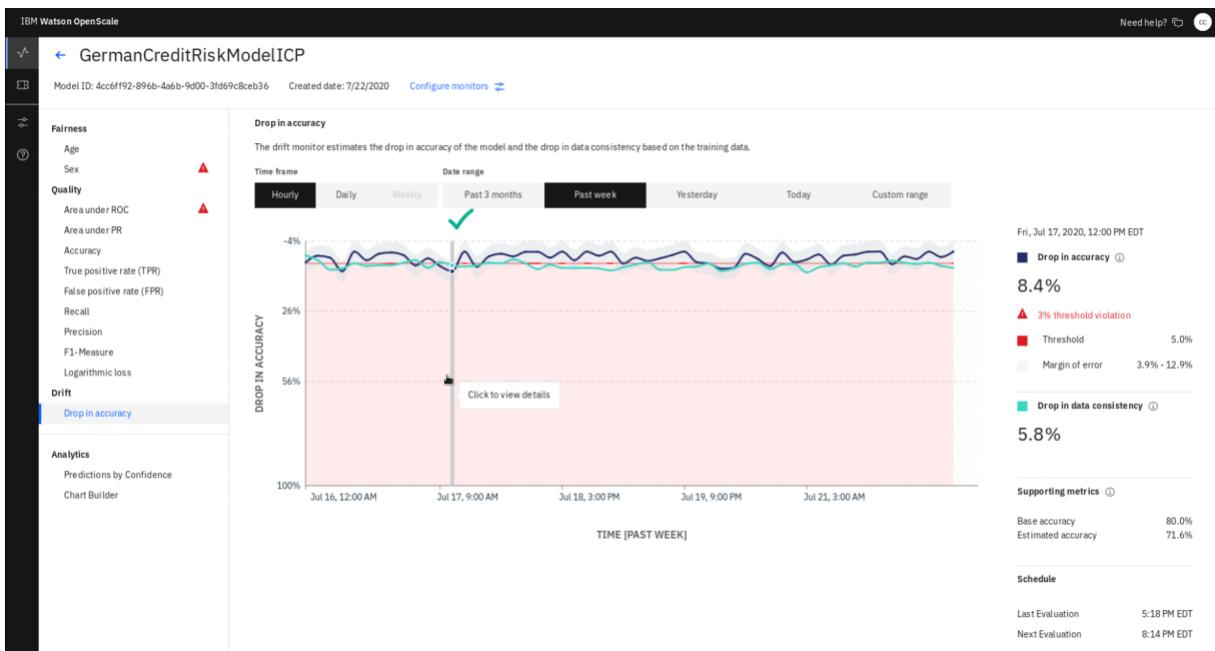
Additionally, the drift monitor compares incoming prediction requests with the training data to identify changes in data consistency that also may affect model output.

These two measurements are shown on the drift monitor screen. Estimated drop in accuracy is represented by the black line, drop in data consistency by the green line, and alert threshold by the red line.

__34. Select the [Drop in accuracy](#) link.

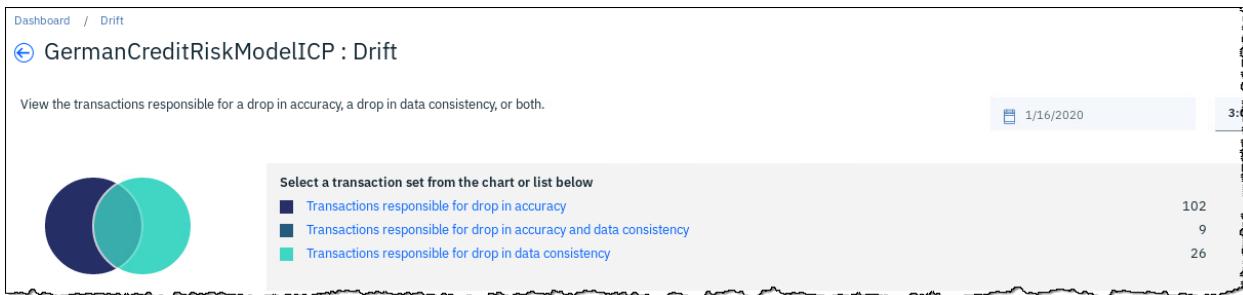


__35. Click the chart where the drop in accuracy is at its greatest (8.4%).

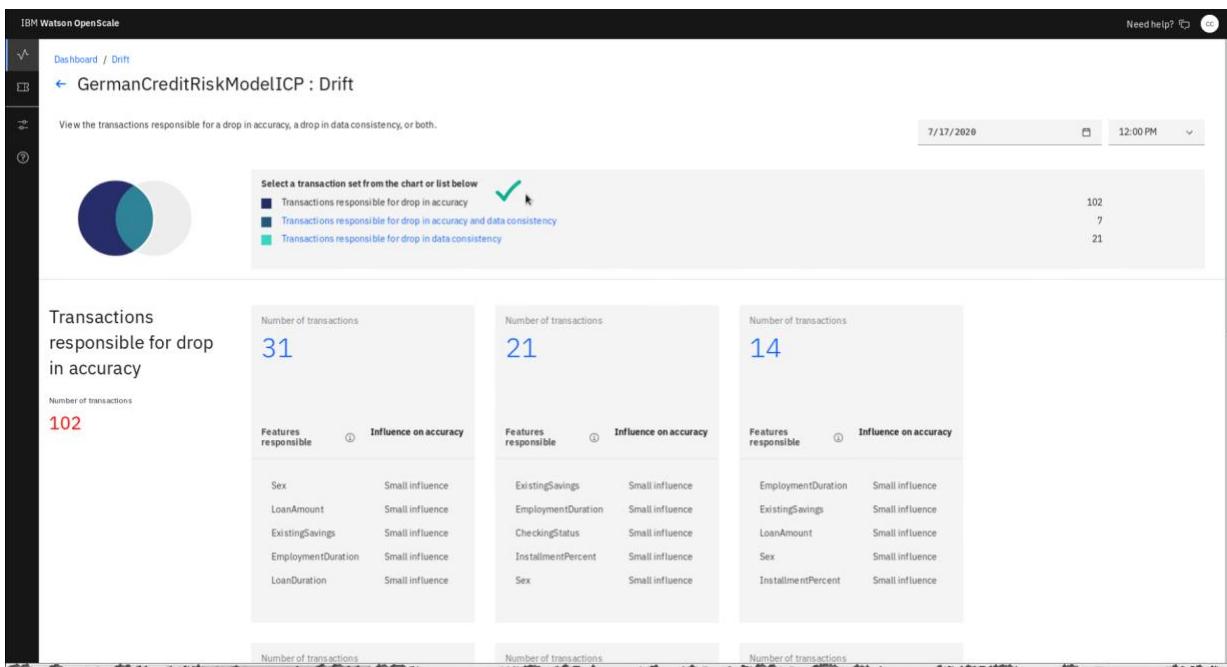


NOTE: Your values may be slightly different than shown due to randomness in data used for training.

- __36. Here, we can get a detailed view of the transactions responsible for estimated drops in accuracy, data consistency, or both.

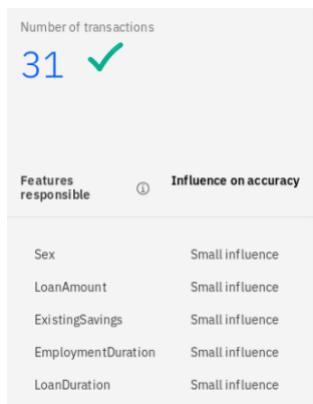


- __37. Click [Transactions responsible for drop in accuracy](#).



OpenScale divides transactions (predictions) that are affecting model accuracy into groups that share feature characteristics, providing a view of which feature values are causing our drift and how much influence each is having.

- __38. Click on one of the tiles for a [transaction grouping](#).



- __39. Here, OpenScale provides a detailed summary of how the values are affecting the model, as well as recommendations for how to address the issues with corrected training data.

OpenScale lists the predictions and feature values for this grouping.

The screenshot shows the IBM Watson OpenScale interface for a 'GermanCreditRiskModelICP' drift analysis. The top navigation bar includes 'Dashboard / Drift / Drifted Transactions'. Below this, a section titled 'Reason for drop in accuracy' explains that Sex, LoanAmount, ExistingSavings, EmploymentDuration, and LoanDuration feature values are causing a drop in accuracy. A 'Recommendation' section suggests reviewing these features. The main content is a table titled 'Transactions responsible for drop in accuracy' with columns: Transactions, Timestamp, Model Output, Sex, LoanAmount, ExistingSavings, and EmploymentDuration. The table lists five transactions with their corresponding details.

Transactions	Timestamp	Model Output	Sex	LoanAmount	ExistingSavings	EmploymentDuration
715f4f5fd88145c898d0b9d886ea88ca-1	Jul 17, 2020, 1:14:01 PM	No Risk	male	3580	less_100	1_to_4
e7109a7591364a4db72fc5dec18b27bc-1	Jul 17, 2020, 1:14:01 PM	Risk	male	5123	greater_1000	greater_7
0191390dd21d465ea8ce2e18db3adc56-1	Jul 17, 2020, 1:14:01 PM	No Risk	male	4626	greater_1000	4_to_7
a060bc91fd1747869670be0078765a65-1	Jul 17, 2020, 1:14:01 PM	Risk	male	3577	500_to_1000	greater_7
98f49157ce404c3d8f907c995f3e350a-1	Jul 17, 2020, 1:14:01 PM	Risk	male	933	500_to_1000	greater_7

8.4.3 Explain an individual prediction

Using a variety of open source algorithms, OpenScale can provide highly detailed explanations of the predictions your model has made.

- __40. Click the [Explain prediction](#) link in the [Actions](#) column of the transactions table.

The screenshot shows a table with three columns: 'LoanDuration', 'Confidence', and 'Actions'. The first row has a 'LoanDuration' of 22, a 'Confidence' of 96.1%, and an 'Actions' column containing a blue link labeled 'Explain prediction' with a green checkmark icon. The second row has a 'LoanDuration' of 29, a 'Confidence' of 98.5%, and an 'Actions' column containing a blue link labeled 'Explain prediction'.

LoanDuration	Confidence	Actions
22	96.1%	Explain prediction ✓
29	98.5%	Explain prediction

- 41. The OpenScale explanation feature works by slightly perturbing the feature values from the original prediction, sending these values to the production model, and measuring the impact the changes have on the outcome. By sending thousands of perturbed requests, OpenScale can gain a detailed picture of feature importance for not only relatively simple models like linear regression or decision tree classifiers, but also complex neural networks and image recognition models. (It can take a minute or two for the details to be calculated, so please be patient.)

The screenshot shows the OpenScale interface for a transaction ID of 240208c96baa4b46b41eb827fad0fda7-1. The 'Details' section includes:

- Model Name:** GermanCreditRiskModelICP
- Type:** Original
- Deployment:** GermanCreditRiskModelICP
- Transaction:** 240208c96baa4b46b41eb827fad0fda7-1

The 'Minimum changes for Risk outcome' section shows:

Feature	Value
Age	74.0
LoanDuration	59.51652
OwnsProperty	unknown

The 'Maximum changes allowed for the same outcome' section shows:

Feature	Value
CheckingStatus	no_checking
LoanDuration	23.00000
CreditHistory	prior_payments_delayed

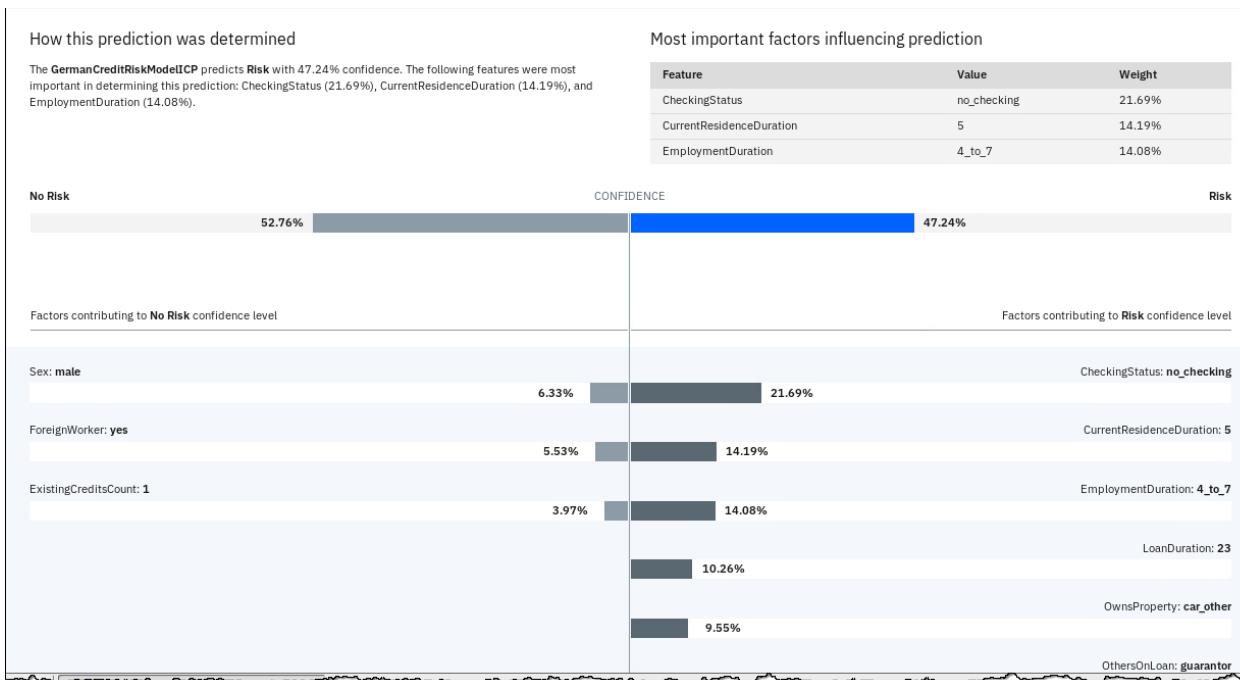
The upper portion of the screen shows information on the model and original prediction, as well as IBM's contrastive explanation technology. The **Minimum changes for No Risk outcome** show the Pertinent Negative values, or the least amount the feature values for the transaction can be changed to get a different outcome.

The **Maximum changes allowed for the same outcome** show the Pertinent Positive values, or how much the feature values can change and still have the model make the same prediction.

- 42. Scroll down to see the factors that influenced the prediction.

Here, you can see the confidence the model has in its prediction as well as a quick summary of the most important factors that led the model to make a prediction that this particular loan application represents a risky one.

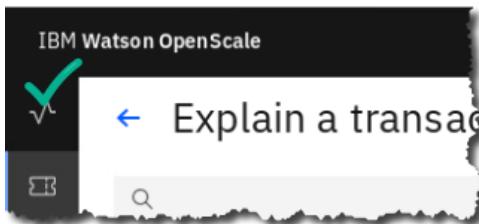
The chart at the bottom of the screen shows the feature values for this prediction, whether they contributed to a No Risk or Risk prediction, and how much they influenced the model.



This detailed information allows you to ensure that your models are making predictions grounded in reality, as well as providing full explainability for predictions in case of an external audit or internal review of the model.

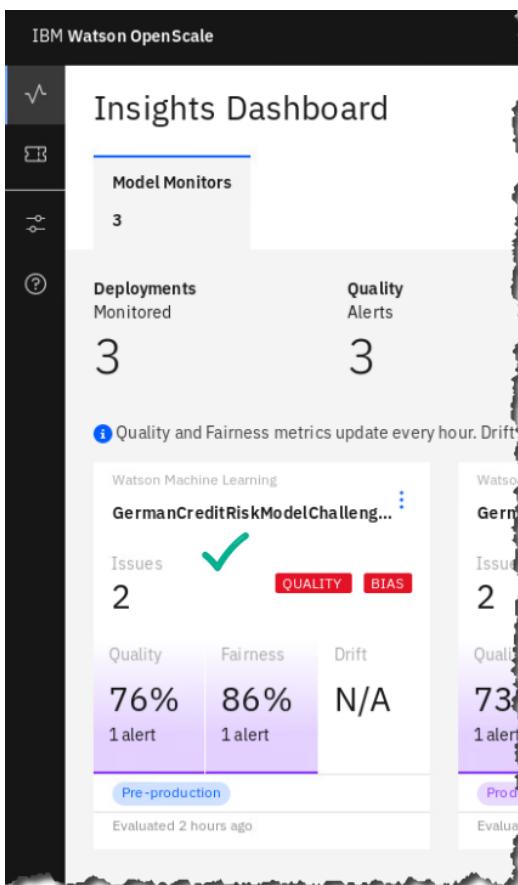
Finally, the data from this explanation is stored in the OpenScale data mart, where it, along with all other metrics, can be retrieved via API and surfaced to business users or even customers via dashboards or other applications.

- 43. Return to the [Insights Dashboard](#) ⇨ [Model Monitors](#).

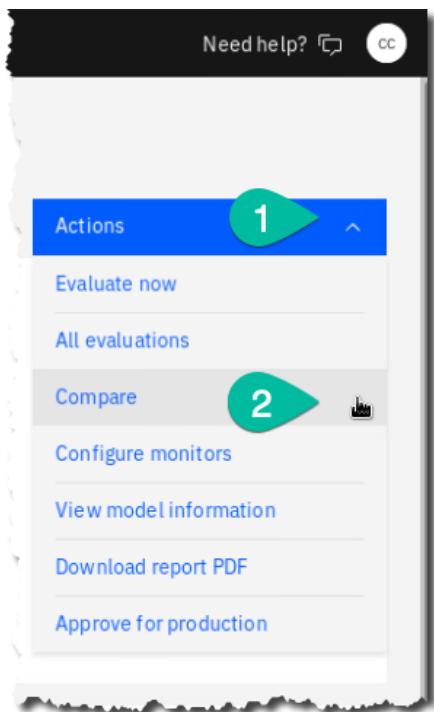


8.4.4 Comparing PreProduction Models

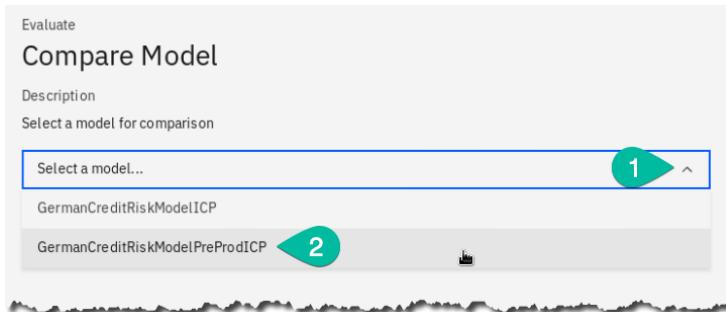
- 44. You can also compare how pre-production models are performing. To do this, select the [GermanCreditRiskModelChallengerICP](#) box.



_45. Once opened, select **Actions – Compare** options from the list.



_46. Select the other pre-production model [GermanCreditRiskModelPreProdICP](#) from the list.



- __47. Examine the comparative results of both models.

The screenshot shows the 'Evaluate' section of the Watson OpenScale interface with the title 'Compare Model'. A dropdown menu is open, showing 'GermanCreditRiskModelPreProdICP' as the selected model. The main area displays a table comparing two models across various metrics. The table has columns for 'Model' (which is currently empty), 'Fairness', 'Quality', and two unnamed columns representing different models. The 'Fairness' row contains entries for 'Age' and 'Sex'. The 'Quality' row contains entries for 'True positive rate (TPR)', 'Area under ROC', 'Precision', 'F1-Measure', 'Accuracy', 'Logarithmic loss', and 'False positive rate (FPR)'. The last two columns show numerical values for each metric. The entire table is highlighted with a red border.

Model			
Fairness	Age	100.00	107.80
	Sex	86.10	97.40
Quality	True positive rate (TPR)	0.57	0.53
	Area under ROC	0.76	0.73
	Precision	0.87	0.82
	F1-Measure	0.69	0.64
	Accuracy	0.82	0.80
	Logarithmic loss	0.43	0.47
	False positive rate (FPR)	0.05	0.06

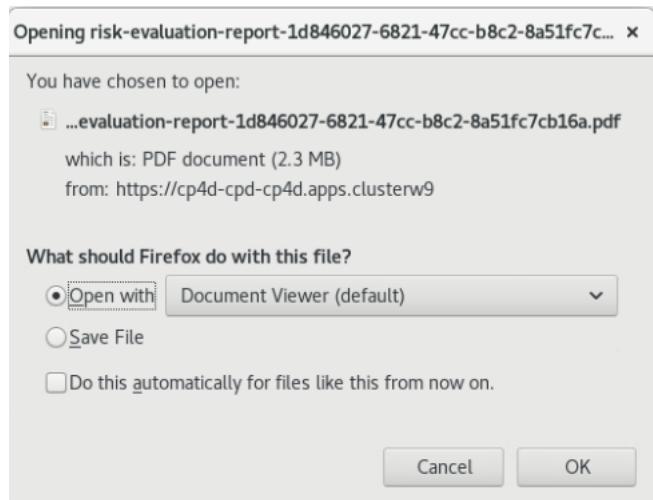
- __48. Select the X to close the comparison.



- __49. One final step you can do is to deliver everything you have generated in a single PDF report. From the [Actions](#) menu, select [Download report PDF](#). It may take a minute to prepare.

The screenshot shows the 'Actions' menu of the Watson OpenScale interface. The menu items are: 'Evaluate now' (highlighted with a green circle labeled '1'), 'All evaluations', 'Compare', 'Configure monitors', 'View model information', 'Download report PDF' (highlighted with a green circle labeled '2'), and 'Approve for production'. The 'Download report PDF' option is the one intended for selection.

- __50. If you are using the remote desktop experience, Select to [Open the document](#) with the default Document Viewer. If you are using a browser, simply download and open on your own machine.



- __51. Browse through the report you have created!

The screenshot shows a browser window displaying the 'IBM Openscale Metric Summary' report. The left sidebar contains a table of contents with the following structure and page numbers:

	Page
IBM Watson OpenScale	1
Metric Summary	1
IBM Watson OpenScale	1
GermanCreditRiskMode...	1
July 22, 2020	1
Report prepared by ...	1
Overview	2
GermanCreditRiskMo...	2
Report Details	2
Model Details	2
Training data details	2
Metrics	3
Metric details	3
Summary	3
RED BREACH	3
Age	3
Sex	3
RED BREACH	3
Statistics	4
Appendix	5
Appendix	6
Area under ROC	6
Area under PR	6
Accuracy	7
True positive rate (T...	7
False positive rate (F...	7
Recall	8
Precision	8
F1-Measure	9
Logarithmic loss	9
Fairness	10
Drop in accuracy	11
Drop in data consiste...	11
Estimated accuracy	12
Base Accuracy	12
Throughput	13

The main content area of the report includes:

- IBM Watson OpenScale Metric Summary
- IBM Watson OpenScale
- GermanCreditRiskModelChallengerICP Evaluation Report
- July 22, 2020

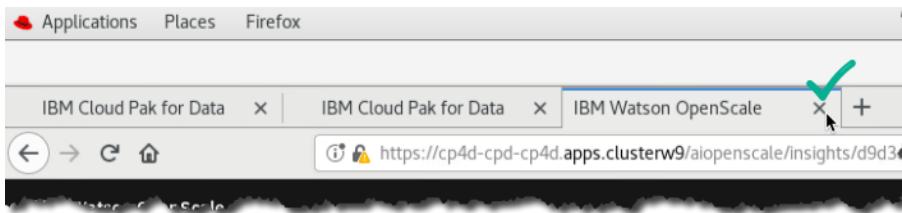
8.5 Monitoring the AutoAI model in OpenScale

Note: this section can only be completed as-is if you have previously completed the Analyze: AutoAI and Deploy labs.

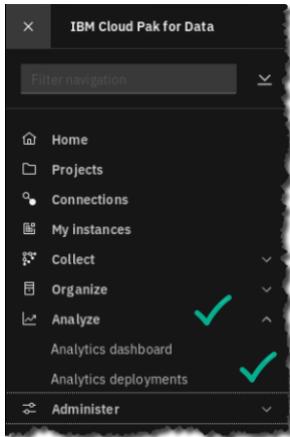
8.5.1 Associate the model with the new deployment space

Before you can monitor the AutoAI model in OpenScale that you created in a previous lab, you need to do a bit of housekeeping first to associate that model to the OpenScale deployment.

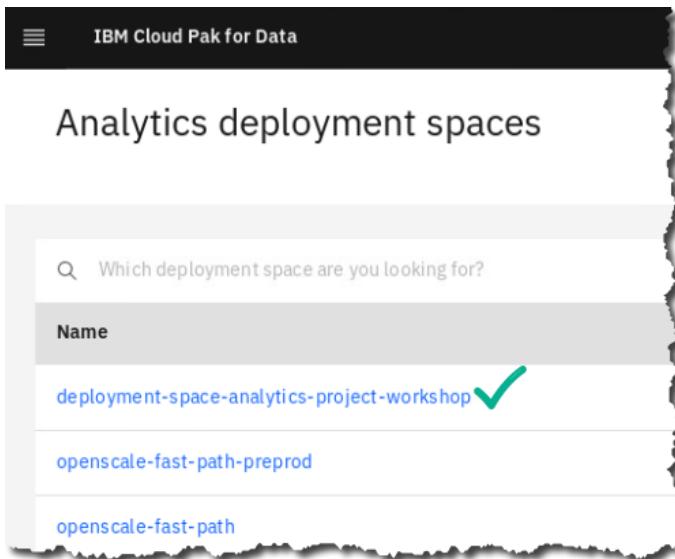
- __52. Close the current browser tab and select your open IBM Cloud Pak for Data tab.



- __53. Click Navigation menu \Rightarrow Analyze \Rightarrow Analytics deployments,



- __54. Select deployment space deployment-space-analytics-project-workshop.



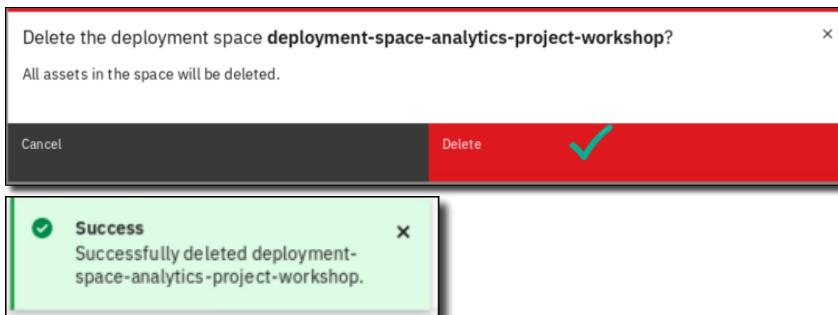
- __55. Select **Deployments** tab, then select the **ellipses** next to every deployment and delete them all.

- __56. Return to **Analytics deployment spaces** (using the breadcrumb on the page).

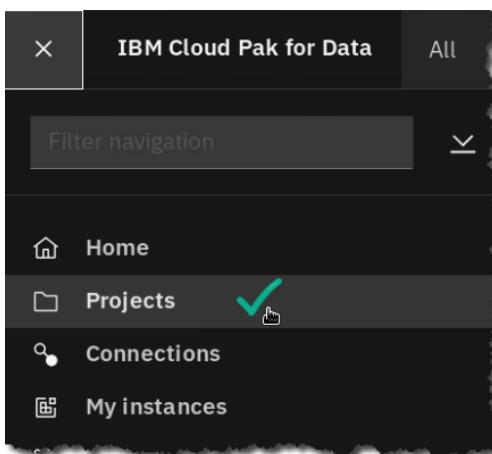
- __57. With no deployment in it, you can now delete the deployment space **deployment-space-analytics-project-workshop**.

Click on the deployment, then **ellipses** \Rightarrow **Delete**.

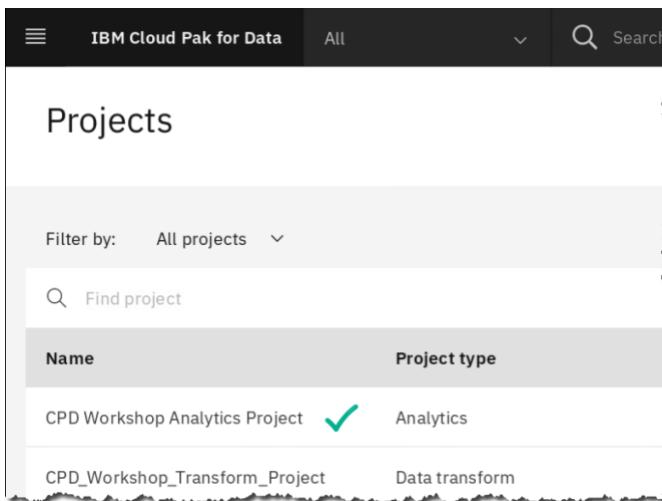
__58. Click **Delete** (to confirm). You should see the message as below.



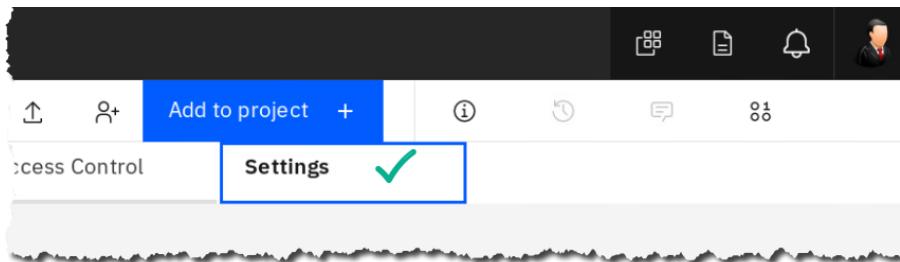
__59. Go to the **Navigation Menu** ⇒ **Projects**.



__60. Select the project: **CPD Workshop Analytics Project**.



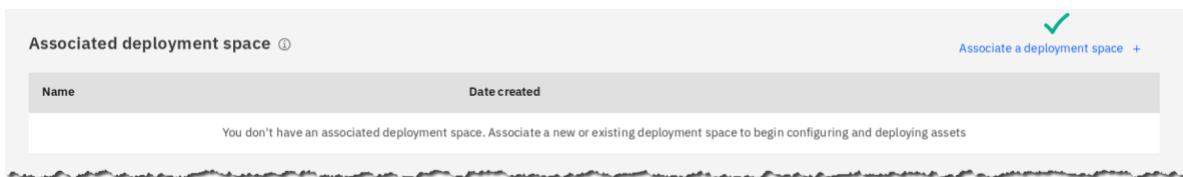
_61. Click tab **Settings**.



_62. Scroll down to find **Associated deployment space**.

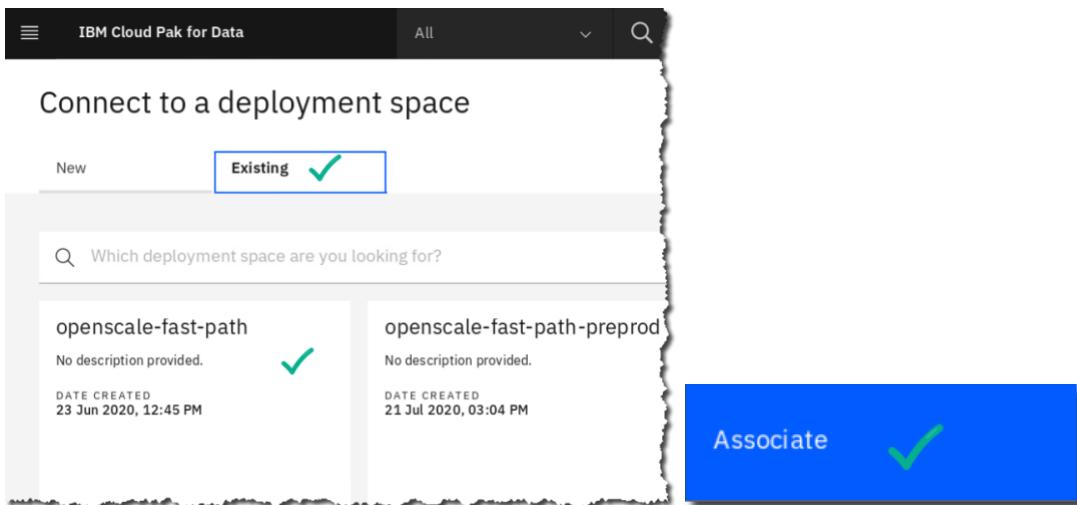
Notice it says this project no longer is associated with a deployment space.

Click [Associate a deployment space](#).

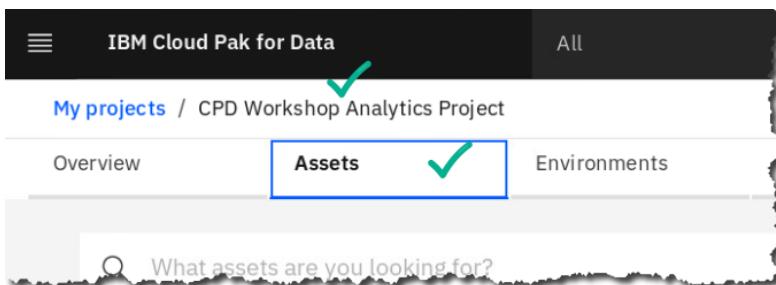


_63. Click section **Existing**.

Click on Deployment Space [openscale-fast-path](#) \Rightarrow [Associate](#).



_64. While remaining in this project, go the tab **Assets**,



- __65. Scroll to find section **Models**.
- __66. Select the ellipses next to the model **ChurnRisk AutoAI Experiment - P4 RandomForestClassifierEstimator** – then click **Promote**.

The screenshot shows the 'Watson Machine Learning models' interface. A table lists a single model entry:

Name	Type	Software specification	Last modified
ChurnRisk AutoAI experiment - P4 RandomForestClassifierEstimator	wml-hybrid_0.1	hybrid_0.1	Jul 22, 2020

A context menu is open over the model row, with the 'Promote' option highlighted. Other options visible in the menu are 'Delete' and three dots (...).

- __67. Choose **Promote to space**.

The screenshot shows a 'Promote model' dialog box. It contains the following text:

Promote model
Promote model to the associated space? Promoting makes this model and its dependencies available for use in deployments and apps.

At the bottom, there are two buttons: 'Cancel' and 'Promote to space'. The 'Promote to space' button is highlighted with a blue background and a green checkmark icon.

- __68. The model is now promoted to the OpenScale deployment space the project is associated with.

The screenshot shows a green success message dialog box with the following text:

Successfully promoted ChurnRisk AutoAI experiment - P4 RandomForestClassifierEstimator to the associated deployment space. Go to the [deployment space](#) to prepare the assets for deployment.

8.5.2 Deploy the model in the new deployment space

Now you can go to that deployment space to deploy the model.

- __69. Click **Navigation menu** ⇒ **Analyze** ⇒ **Analytics deployments** ⇒ **openscale-fast-path**.

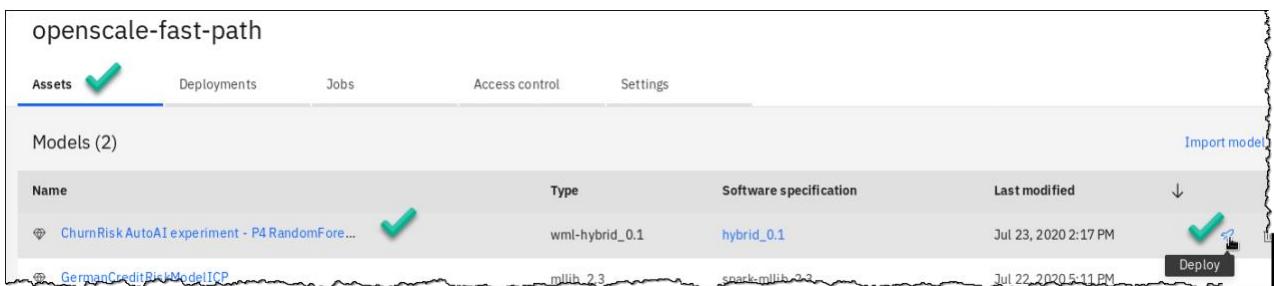
The screenshot shows the 'IBM Cloud Pak for Data' navigation menu. The path taken is:

- Navigation menu
- Projects
- Connections
- My instances
- Collect
- Organize
- Analyze
- Analytics dashboard
- Analytics deployments
- Administrator

Each step in the path is highlighted with a green checkmark icon. A separate window titled 'Name' shows the value 'openscale-fast-path' with a green checkmark icon.

- __70. Under the tab **Assets**, hover on the **launch** (rocket) icon for the model **ChurnRisk AutoAI experiment - P4 RandomForestClassifierEstimator**.

Click **Deploy**.



- __71. Select **Online**.

Fill in Name: **ChurnRisk AutoAI Experiment Deployment**.

Click **Create**.

The screenshot shows the 'Create a deployment' form. At the top, it says 'Associated asset: ChurnRisk AutoAI experiment - P4 RandomForestClassifierEstimator'. Below that, 'Deployment type' has 'Online' selected with a green checkmark. The 'Name' field contains 'ChurnRisk AutoAI Experiment Deployment' with a green checkmark. The 'Description' field is empty. At the bottom right, the 'Create' button is highlighted with a blue box and a green checkmark.

- __72. You will see it as an [Online](#) deployment that is first: [In-progress](#).

DEPLOYMENT TYPES				1 Online Deployment(s)		
Online	(1)	Name	Status	Last modified		
Batch	(0)	Churn Risk Auto...	In progress	Jul 23, 2020 2:28 PM		

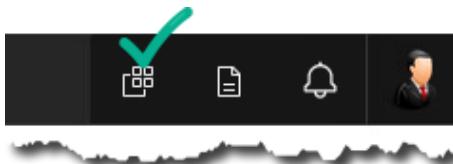
...in a minute or two, it will show as: [Deployed](#).

DEPLOYMENT TYPES				1 Online Deployment(s)		
Online	(1)	Name	Status	Last modified		
Batch	(0)	Churn Risk Auto...	Deployed	Jul 23, 2020 2:28		

8.5.3 Monitor the model in OpenScale

Now you will be able configure the monitoring of that model in Watson OpenScale.

- __73. From the web client, click the icon for [Services](#) (at the top right corner of the screen).



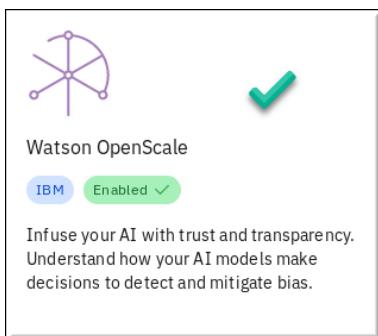
- __74. Click [Services Category](#) \Rightarrow [AI](#).

Find services

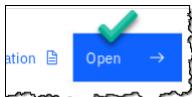
1 × Category

 AI

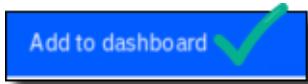
__75. Find and click on tile Watson OpenScale.



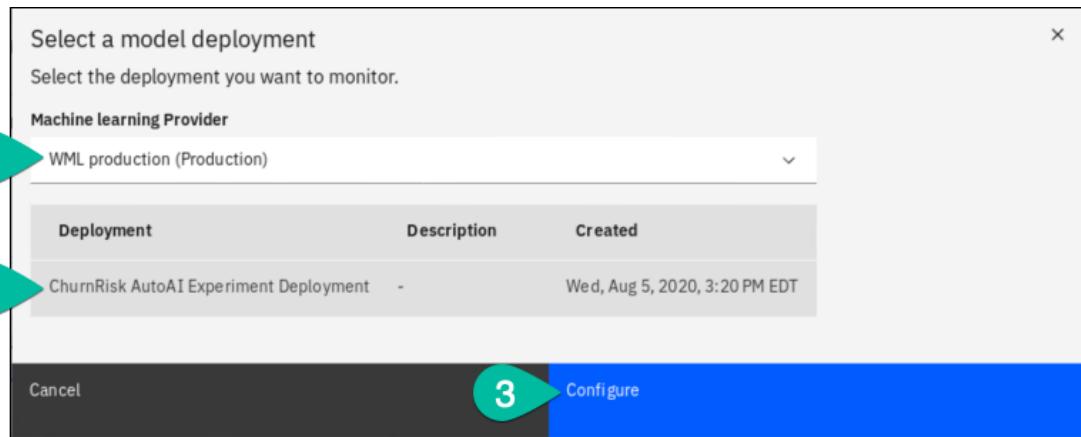
__76. Click (top right corner) [Open](#).



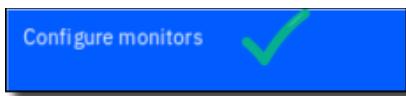
__77. Click (top right corner) [Add to dashboard](#).



__78. Select the Deployment: ChurnRisk AutoAI Experiment Deployment \Rightarrow Configure.



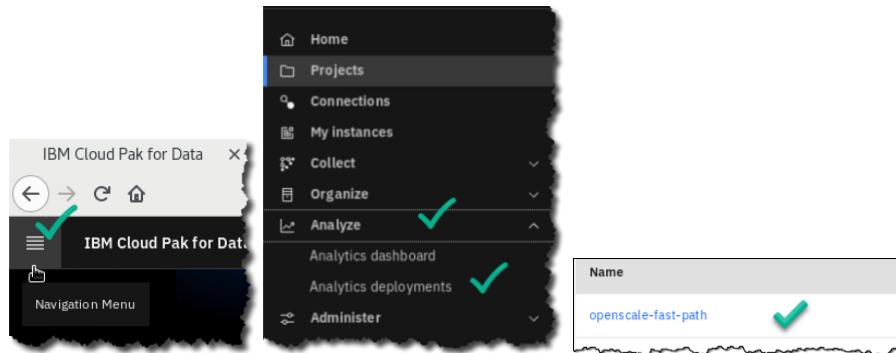
__79. Click [Configure monitors](#).



_80. Next you will need to do at least one test (or payload log action) to activate the monitor.

If you are using the remote desktop, click the [Cloud Pak for Data web client](#) to open a new browser tab or duplicate your browser tab. If you are using your browser, duplicate the tab.

From that tab return to [deployment spaces](#) ⇒ [openscale-fast-path](#).



_81. Choose deployment [ChurnRisk AutoAI Experiment Deployment](#).

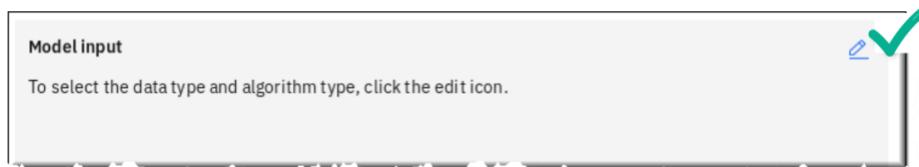


- __82. Under tab **Test** enter the information shown below Click **Predict**.

(Note: you don't have to fill in all the features (input data) for this test, just the four values shown. This provides at least one payload log action... the prediction itself is not relevant here.)

The screenshot shows the 'Churn Risk AutoAI Experiment deployment' interface. The 'Test' tab is selected. The 'Enter input data' section contains four fields with values: ID (81), AGE_GROUP (Adult), GENDER (F), and STATUS (S). A blue box highlights the 'Predict' button at the bottom right.

- __83. Return to the OpenScale browser tab and click the **edit** (pencil) icon in the tile **Model Input**.



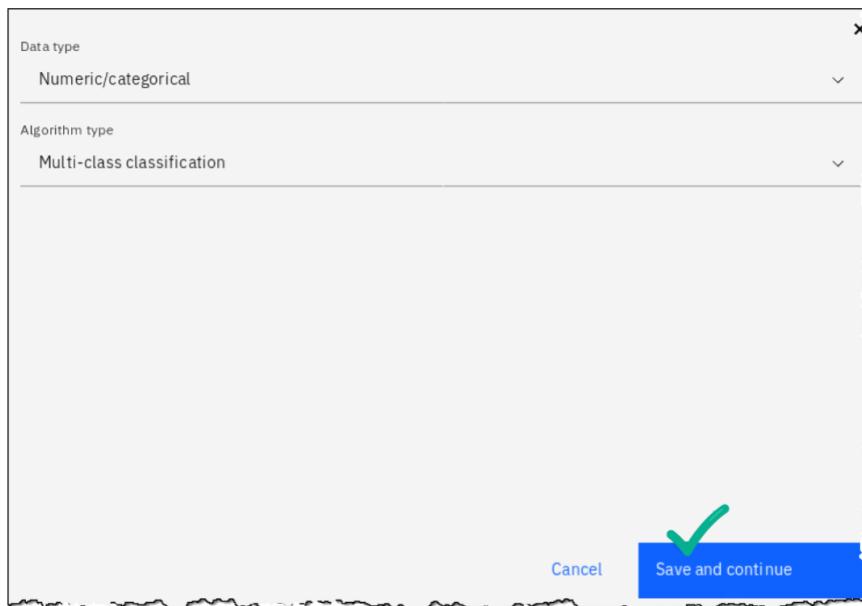
- __84. Under **Select data type**, choose **Numerical/categorical**.



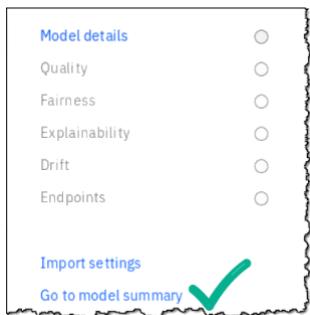
- __85. Under **Select algorithm type**, choose **Multi-class classification**.



__86. Click [Save and continue](#).



__87. Click [Go to model summary](#).



For the remainder of the lab, try to complete the rest of the configuration tabs yourself.

The purpose of this exercise was to show you how to get started with model monitoring in Watson OpenScale with your own model. How to configure all of the options is probably best suited for a future Deep Dive lab.

8.6 Lab conclusion

Watson OpenScale is an enterprise-grade environment to help your organization infuse your AI with trust and transparency at scale, delivering a quick return on investment.

This Cloud Pak for Data service understands your models, detects and mitigates drift, and delivers explainable outcomes that are free from bias.

With this feature, Trade Co. is able to stay ahead of the curve.



**

Read here about how Watson OpenScale won a 451 Firestarter award in Q2 2019.

http://ibm.biz/openscale_award

IBM Watson OpenScale wins 451 Firestarter Award for making AI outcomes fair and explainable

JUNE 27, 2019

** End of Lab 08 – Infuse: Watson OpenScale

Lab by Eric Martens, Burt Vialpando and Kent Rubin - IBM

Lab 09 INFUSE: COGNOS ANALYTICS - INTRODUCTION

9.1 Lab overview

In this lab, we will upload a CSV file given to us by our IT department to see if there is any interesting information and if we can use this information to better target market our potential high-risk churn customers.

For this lab, we will use the power of the IBM Cloud Pak for Data platform and the IBM Cognos Analytics cartridge. Installation of the IBM Cognos Analytics cartridge has already been done for you.

9.2 Persona represented in this lab

For this lab, we will assume the roles of both a Data Engineer and a Business Analyst.

Persona (Role)	Capabilities
 Business Analyst	Business Analysts deliver value by taking data, using it to answer questions, and communicating the results to help make better business decisions.

9.3 Logging into the CPD web client (if you have not already done so)

- __1. If you are starting this lab stand-alone (without going through previous labs) do the following:
- __2. Double-click the desktop icon: [Cloud Pak for Data Web Client](#).



- __3. The CPD web client GUI displays as shown. Use `cpduser` and `cpdaccess` for the *Username* and *Password* and click [Sign In](#).

 A screenshot of the IBM Cloud Pak for Data sign-in page. The page title is "SIGN IN" above "IBM Cloud Pak for Data". Below the title is a graphic showing a person interacting with a server and databases. The sign-in form contains fields for "Username" (with "cpduser" entered) and "Password" (with "cpdaccess" entered). A blue "Sign in" button with a checkmark icon is at the bottom.

9.4 Logging into the IBM Cognos provisioned instance

- __4. Select the [Navigation Menu](#) icon at the top of the screen.



- __5. Select [My Instances](#) from the menu.

- __6. Select [Provisioned instances](#) from the tab at the top.

Here you will see a list of instances that have been provisioned for you already.

- __7. Select the [ellipsis](#) at the end of the [cognos-analytics-app](#) and select [Open](#).

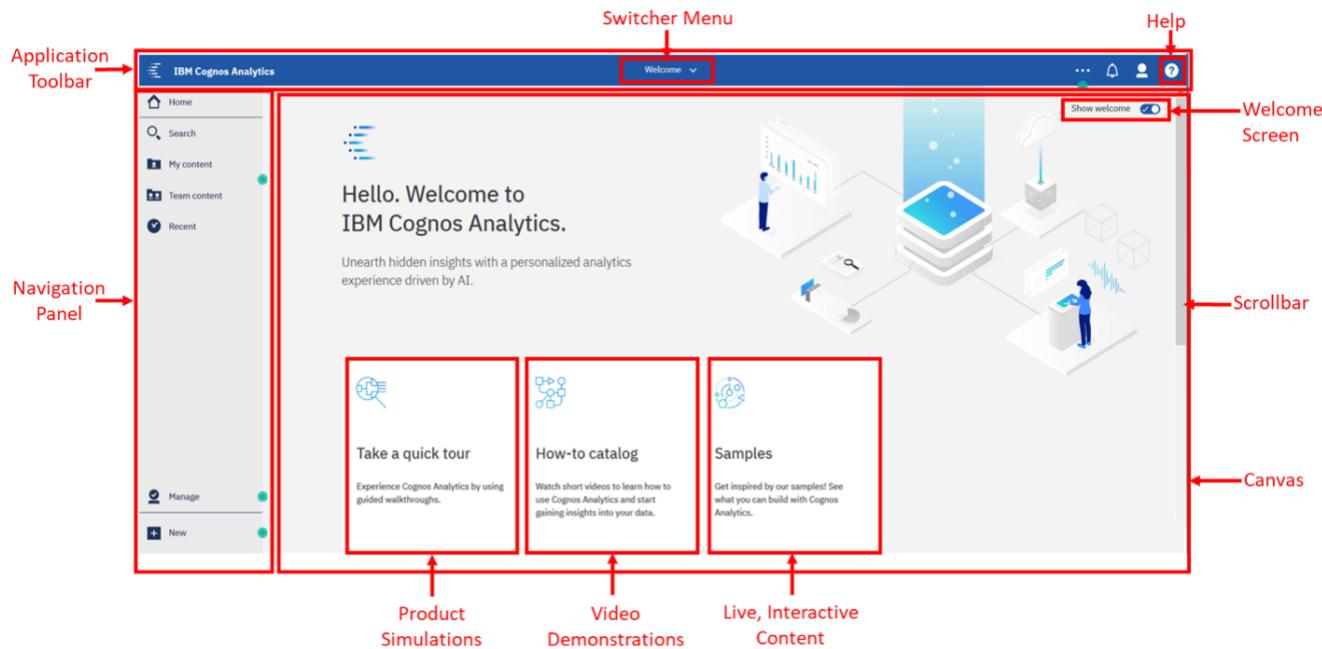
You will now be automatically directed and logged into the IBM Cognos Analytics instance on your IBM Cloud Pak for Data platform.

9.5 The IBM Cognos Analytics User Interface

The purpose of the User Interface (UI) is to provide Users with a streamlined way to get started using Cognos Analytics and view content and activities pertinent to them.

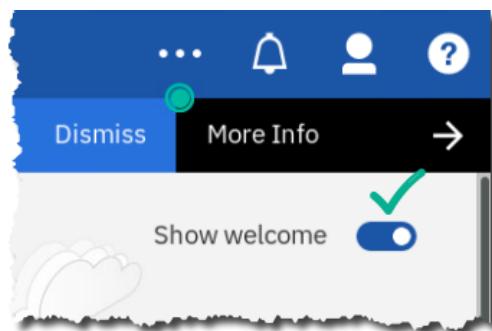
The User Experience brings you directly into the completely redesigned IBM Cognos Analytics User Interface (UI).

All IBM Cognos Analytics Users begin their navigation here at the [Welcome Screen](#).

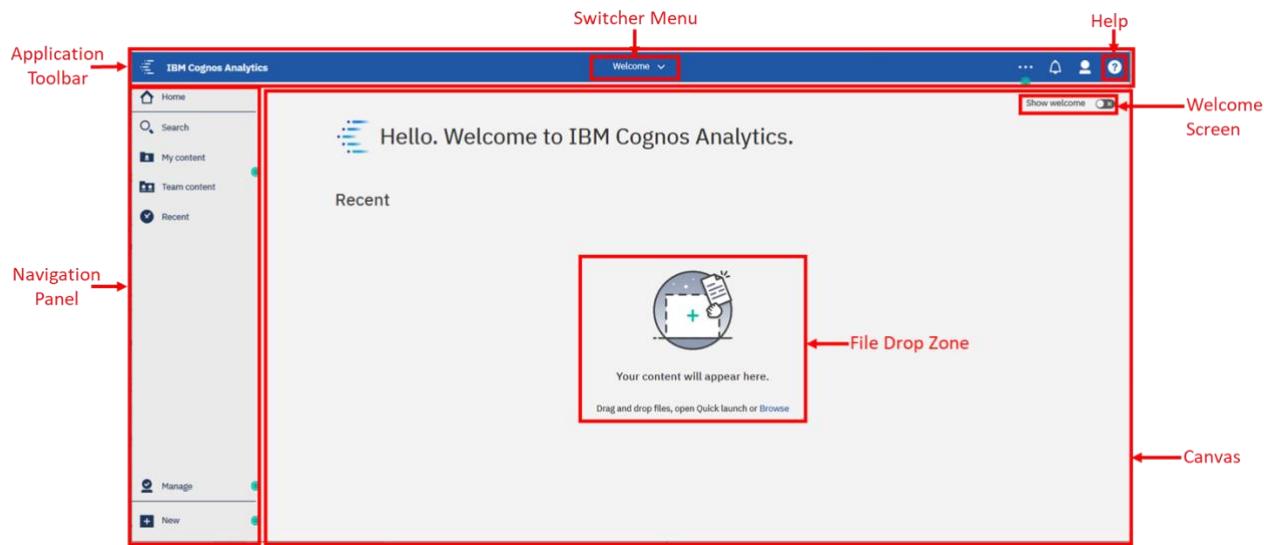


TECH TIP: Not all [Welcome Screen Getting Started](#) tiles may be available in your workshop instance. The Welcome screen getting started tiles presented are based on which have been configured to run in the environment.

- __8. Click on the [Welcome Screen toggle button](#) on the upper left of the canvas to collapse the Welcome Screen's Getting Started content.



- 9. The canvas now shows the Recently used files, if any, in the **Recent** section, along with the **File drop zone** where users can easily upload their data files.

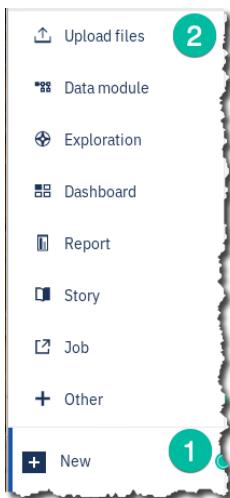


- 10. Once you begin working with content, the canvas will update with your recently used items. In your Cognos Analytics instance, you may see recent content on the canvas.

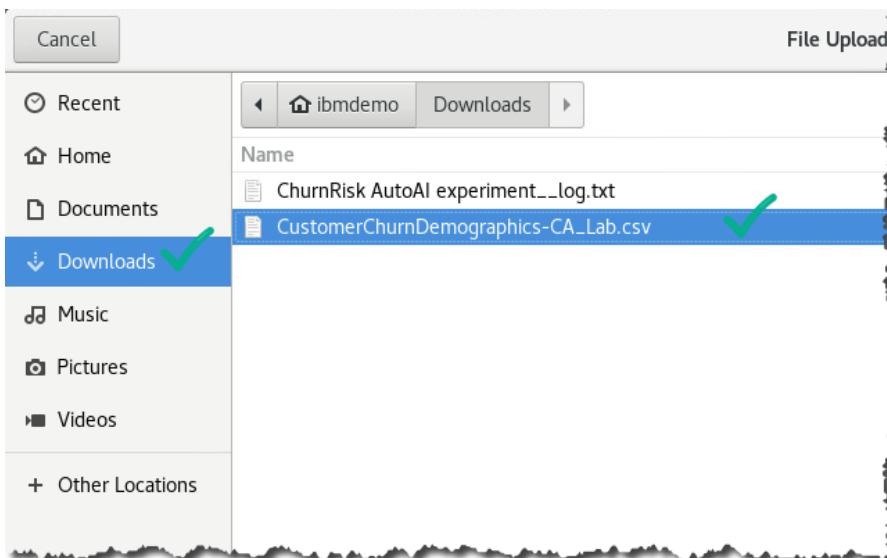
Now that we are in the IBM Cognos Analytics user interface, let's make a connection to the IBM Cloud Pak for Data – Virtualized Data. Remember, our business use case here is to analyze and create a dashboard quickly from this new set of data we virtualized. We did not have to wait for a Data Warehouse or Data Lake to be completed. Let's do some analytics at the speed of thought!

9.6 Importing the CSV file as an Exploration

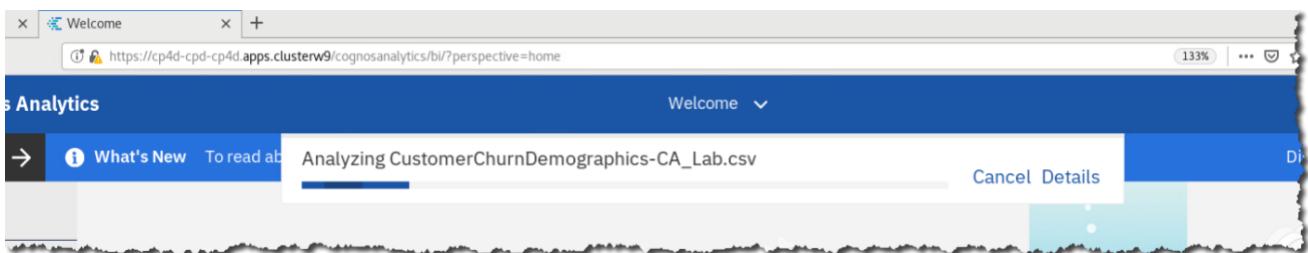
- _11. From your Cognos Analytics Home screen, select [New – Upload](#).



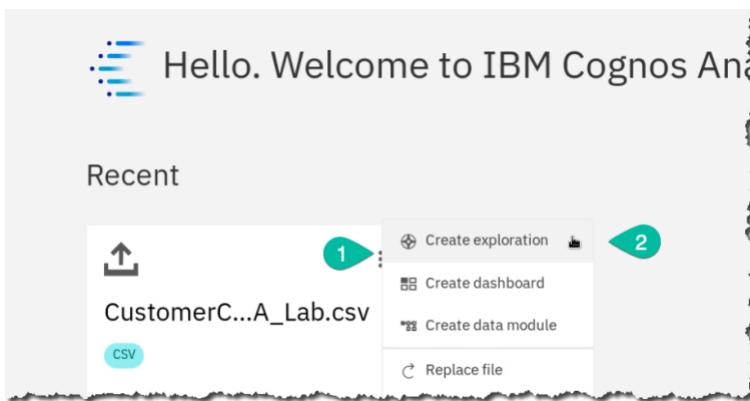
- _12. Select [CustomerChurnDemographics-CA_Lab.csv](#) from [/home/ibmdemo/Downloads](#) if you are accessing via remote desktop. If you are accessing via browser, you may download the file from <http://ibm.biz/Churn-Demographics>.



- __13. Cognos Analytics will then begin to import and analyze the data in the CSV file.



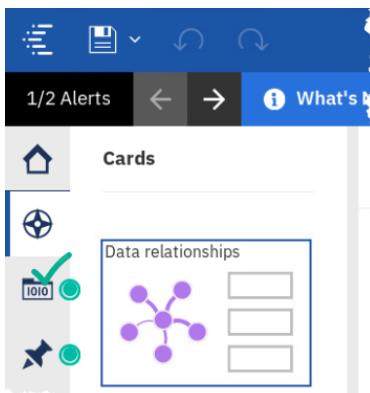
- __14. In the Recent area, in the upper left corner of the CustomerChurn CSV card, select [ellipses](#) \Rightarrow [Create exploration](#).



9.7 Cleaning up the data

When IBM Cognos Analytics imports CSV files, it does its best to understand the data type, usage and aggregation from reading your data. Depending upon the data you import, you may need to help IBM Cognos Analytics out by defining some data properties.

- _15. From the left-hand menu, select the icon for the [Data tab](#).



Note: You will notice icons next to each field showing a representation of the data items usage as recognized by IBM Cognos Analytics. Measures are represented by a ruler, numeric values as a hash, text values as ABC, time values as a clock and geographic points as a pin.

- _16. Because this data is a snapshot total of values for a given date, you will need to change the aggregation properties of your measures, so they roll up correctly and do not give you misleading or false information.

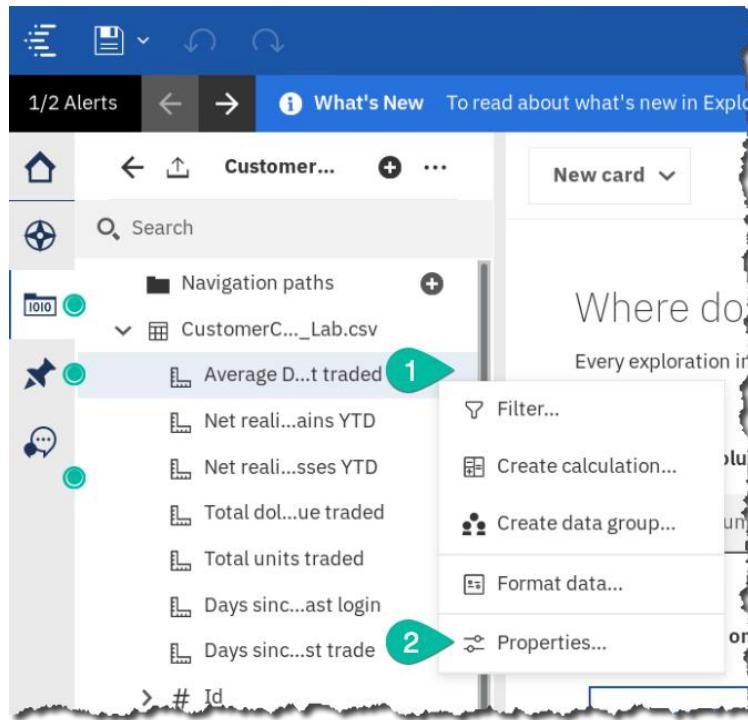


Business Analyst

TIP: If you cannot read the full name of the data item shown on your screen, simply hover over the name and a fly out will display the full name. Another option is to use the zoom out feature of your browser.

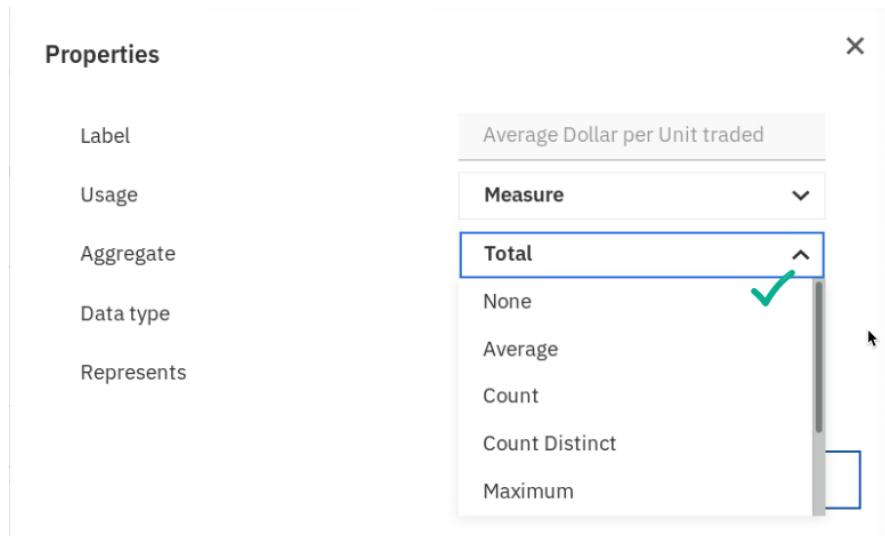
- __17. For the following measures, change the Properties of the Data Item by selecting the ellipses at the end of each item and choosing [Properties](#) from the menu.

For example, the following displays the properties for [Average Dollar per Unit traded](#):



- __18. You will be presented with the [Properties](#) of that measure.

Change the [Aggregate](#) of [Average Dollar per Unit traded](#) measure to [None](#), then [Close](#).

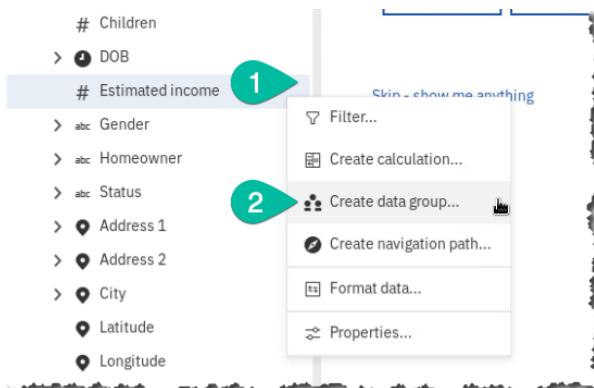


- __19. Change the other **Data Items** with the following property values by following the previous steps.

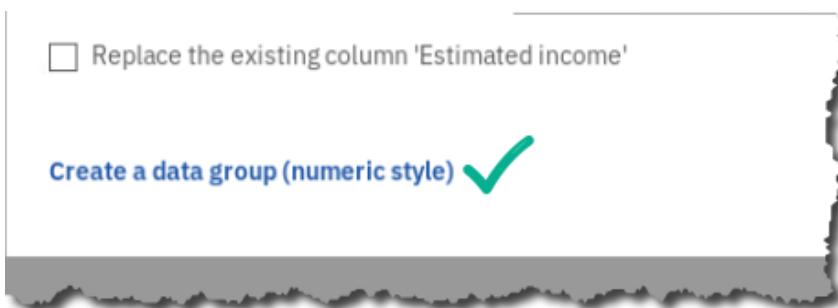
Data Item	Usage	Aggregate
Average Dollar per Unit traded	Measure	None
Net realized gains YTD	Measure	None
Net realized losses YTD	Measure	None
Days since last login	Attribute	Count Distinct
Days since last trade	Attribute	Count Distinct
Age	Attribute	Count Distinct
Children	Attribute	Count Distinct
Estimated Income	Attribute	None

- __20. Because there are a lot of **Estimated Income** values, it would be better to group this data in buckets for easier understanding.

By **Estimated income**, select **ellipses** \Rightarrow **Create data group**.



- __21. Choose **Create a data group (numeric style)**.



_22. You will be presented with the following auto grouping:

The screenshot shows the 'Create a data group (numeric style)' dialog box. The 'Name' field is set to 'Estimated income (Group)'. The 'Groups' dropdown is set to 5. The 'Range border values' section shows five ranges: '96019.2660 and above', '72038.5320 to < 96019.2660', '48057.7980 to < 72038.5320', '24077.0640 to < 48057.7980', and 'less than 24077.0640'. Each range has a corresponding 'Higher' and 'Lower' value listed below it. A checkbox for 'Group NULL values as' is present. At the bottom are 'Create' and 'Cancel' buttons.

_23. Change the number of **Groups** to 4.

The screenshot shows the 'Create a data group (numeric style)' dialog box. The 'Name' field is set to 'Estimated income (Group)'. The 'Groups' dropdown is set to 4, indicated by a green checkmark. The 'Range border values' section shows four ranges: '90024.0825 and above', '60048.1650 to < 90024.0825', '30072.2475 to < 60048.1650', and 'less than 30072.2475'. Each range has a corresponding 'Higher' and 'Lower' value listed below it. A checkbox for 'Group NULL values as' is present. At the bottom are 'Create' and 'Cancel' buttons.

Note: You could manually enter group values or use Cognos to automatically do it for you!

_24. Select [Create](#).

You will notice you have created a new grouping at the top of your Data Items list.

The screenshot shows the IBM Cognos Analytics interface. On the left, there is a blue sidebar with a 'Create' button containing a green checkmark. In the center, a dropdown menu is open under 'abc Estimate... (Group)'. The menu items are: 'less th...72.2475', '30072.2...8.1650', '60048.1...4.0825', and '90024.0... above'. A red box highlights this dropdown menu. A green circle with the number '1' is positioned to the left of the 'Create' button.

_25. When you are finished, return to the Exploration canvas by selecting the [Exploration](#) tab.

The screenshot shows the IBM Cognos Analytics interface. The 'Exploration' tab is highlighted with a green checkmark. The sidebar on the left also has a green checkmark next to the 'Exploration' icon.

IBM Cognos Analytics Exploration asks you: [Where do you want to start?](#)

You can see that IBM Cognos Analytics suggests starting with some of the measures it found in your Data Module. This is helpful in discovering what are the possible drivers behind each of your metrics. IBM Cognos Analytics uses a powerful AI engine to assist you in discovering possible drivers between measures and dimensions, strength of relationships between dimensions, and also possible ways of examining this data for you by suggesting highly visual graphs of this data by dimension.

When IBM Cognos Analytics suggests: [Try starting with one of these...](#) it is providing you an easy way to begin discovering and exploring what it finds as potentially important (what is driving our metrics?) based on a measure or dimension.

We are not limited to these options, though.

- __26. Select the area where it asks you to [Enter data column](#) and you will see that all of your Data Module Data Items are presented.

The screenshot shows a user interface for exploring data. At the top left is a button labeled "New card ▾". Below it is a question "Where do you want to start?". A note says "Every exploration includes a starting card." A text input field contains the placeholder "Start with any column in 'CustomerChurnDemographics-CA_Lab.csv'. You can always change it later". Below this is a dropdown menu with the heading "Enter data column. Not sure? Try Average Dollar per Unit traded, Days since last trade". The menu lists several items from a dataset named "CustomerChurnDemographics-CA_Lab.csv":

- CustomerChurnDemographics-CA_Lab.csv
 - Estimated income (Group)
 - Average Dollar per Unit traded
 - Net realized gains YTD
 - Net realized losses YTD
 - Total dollar value traded
 - Total units traded
 - # Days since last login
 - # Days since last trade



Business Analyst

TIP: If you're really adventurous, or do not know the data you are exploring very well, you could select the [Skip – show me anything](#) link and IBM Cognos Analytics will select any measure as a starting point.

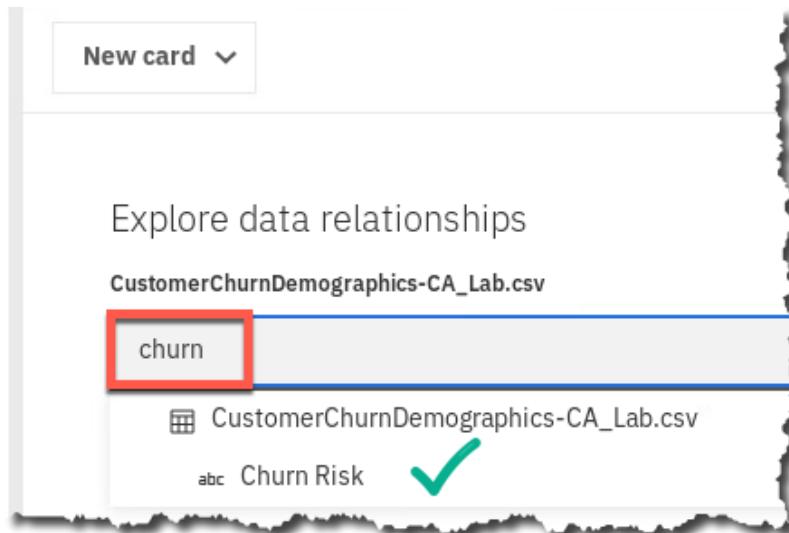
[Do not do this for this Workshop!](#) It's good to know you can use the AI help from IBM Cognos Analytics when you need.

9.8 Exploring Data Relationships

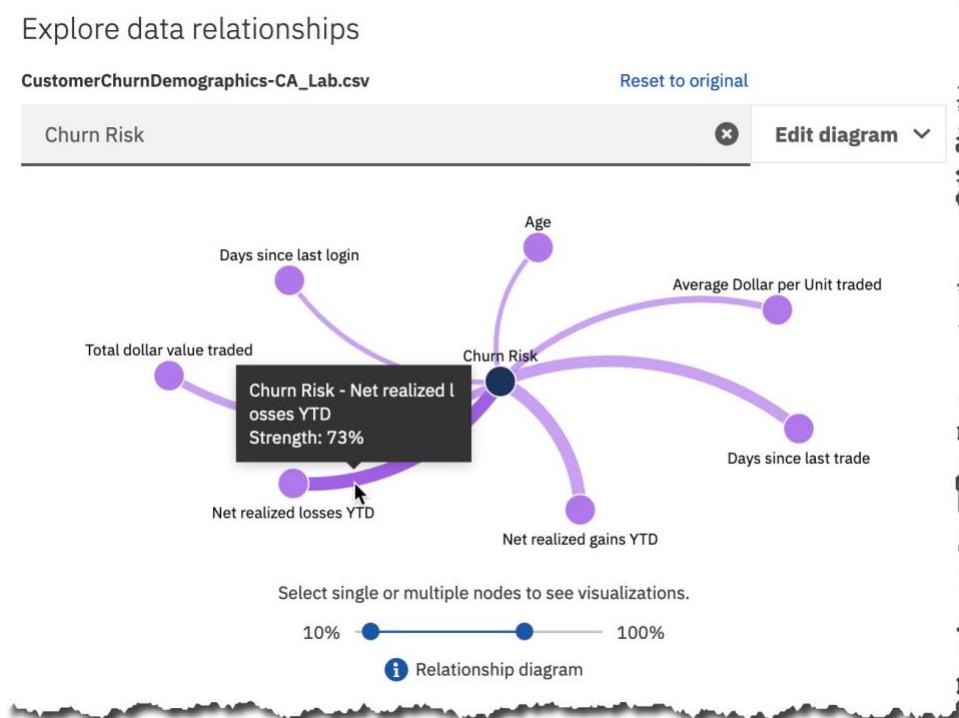
For this workshop, let's have IBM Cognos Analytics help us understand our customers better from this data set.

Let's first find out which customers are spending the most on trades with us.

- _27. From the dropdown, type [Churn](#), then select [Churn Risk](#).



- _28. By hovering over each of my measures associated with [Churn Risk](#), I can get a quick understanding of relationship strength that is driving [Churn Risk](#).

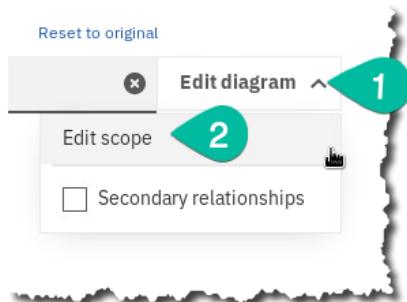


For example, there seems to be a 73% relationship strength between [Net realized losses YTD](#) and [Churn Risk](#).

__29. Hover over other measures to see how they relate as well.

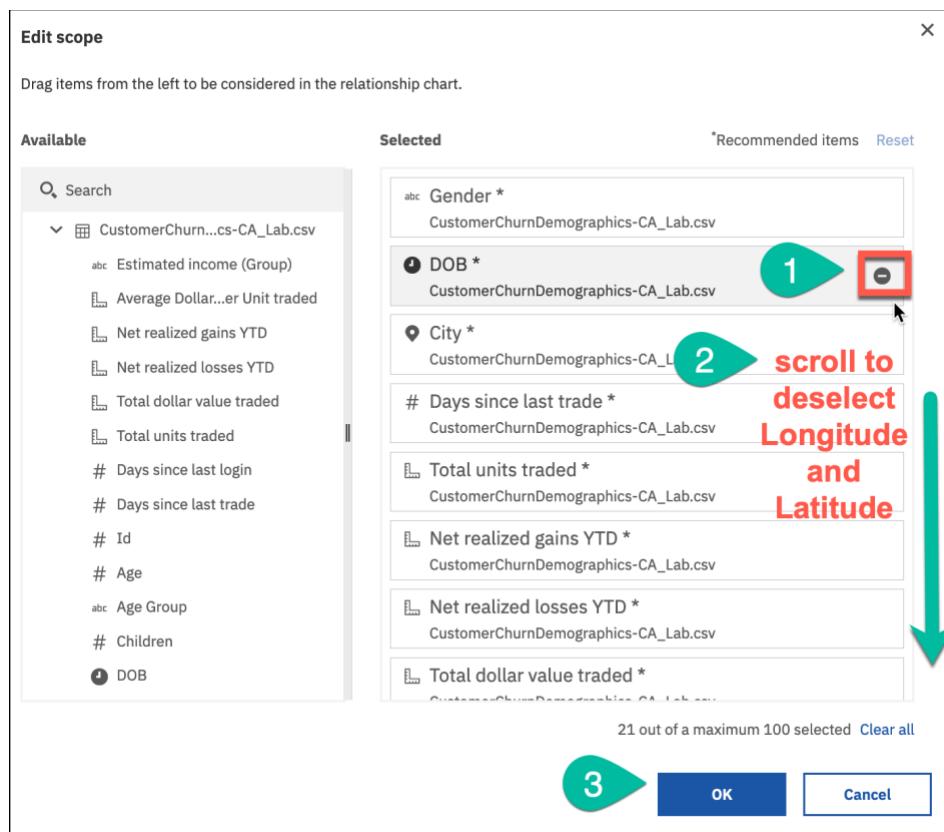
You can see that IBM Cognos Analytics has selected the scope of fields to show relationship strength between. You are not limited to this, however. You can change the scope of this relationship examination to remove fields that you know may interfere with getting a proper result.

__30. Select the drop-down menu to [Edit diagram](#) and select [Edit scope](#).



Fields such as [DOB](#), [Longitude](#) and [Latitude](#) are probably Data Items that do not bring any value when examining your data relationships and may cause static when exploring your data.

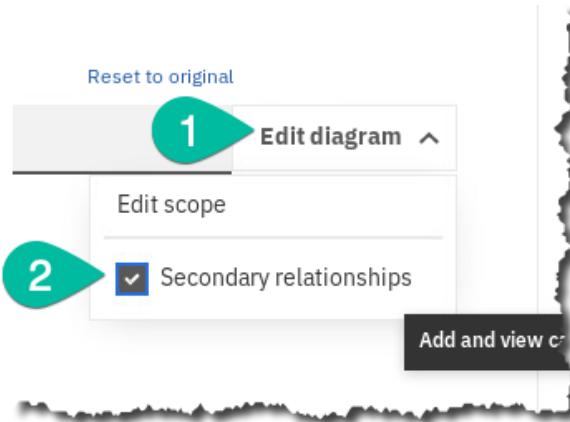
__31. You can remove those by hovering over, then selecting the next to [DOB](#), [Longitude](#) and [Latitude](#), then [OK](#).



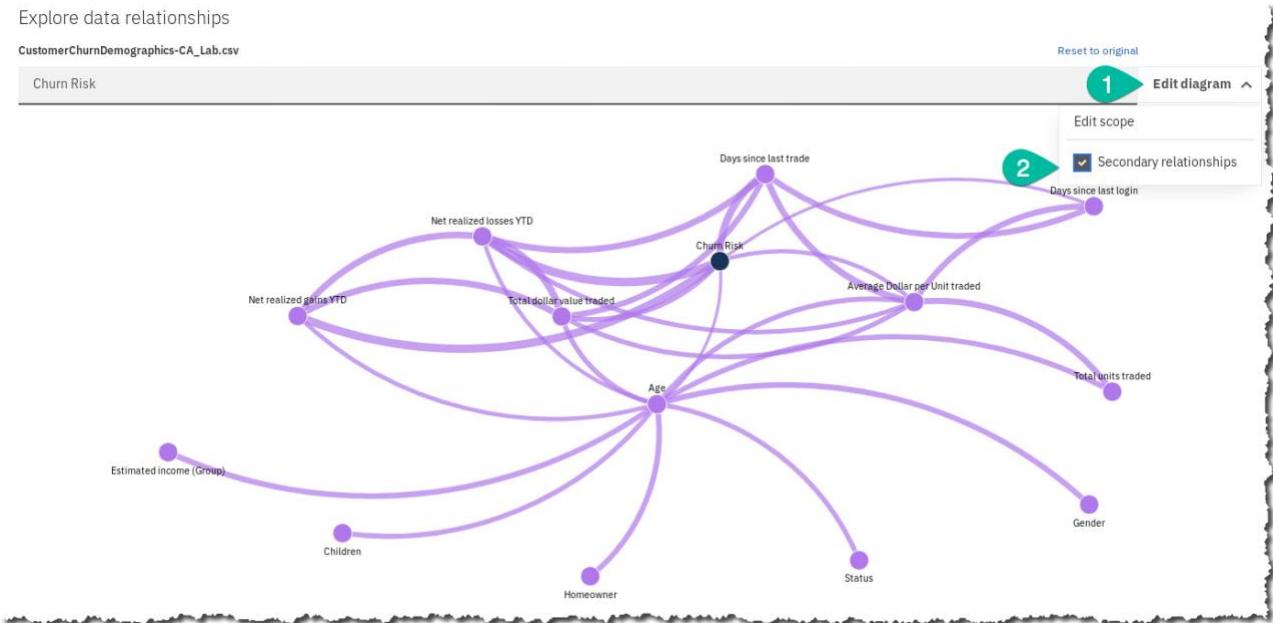
IBM Cognos Analytics gives you the flexibility to remove fields that you may not require when exploring your data and may interfere (static) with your analysis.

- __32. You can also investigate secondary relationships to see the strength of the relationship between all fields.

Select **Edit diagram** ⇒ **Secondary relationships**.



- __33. Hover around each of the fields to see view the relationship strength IBM Cognos Analytics discovered between them.

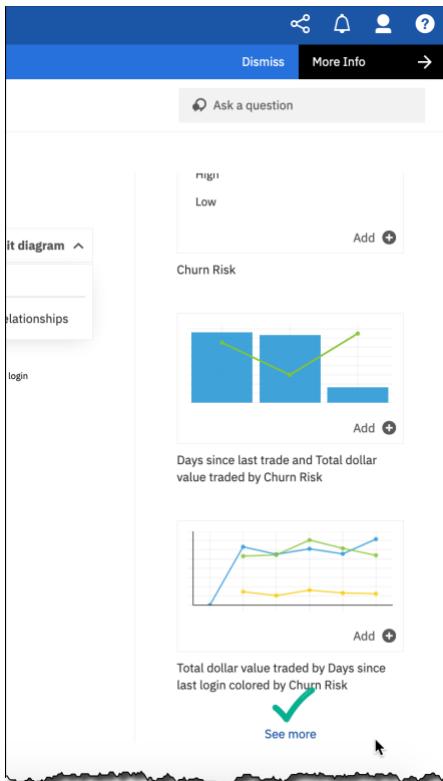


9.9 Creating Exploration Cards

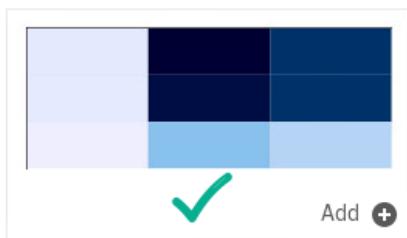
When IBM Cognos Analytics examines relationships, the powerful AI engine will begin to surface charts and graphs that may be of interest to you. This is especially helpful if you are unfamiliar with the data you are working with.

Scroll down on the far-right portion of your screen, then notice that IBM Cognos Analytics has presented you with three visualizations. You can begin with one of these info graphs or 'Cards' or see more Cards.

- __34. Select link [See more](#).



- __35. You are now able to see more auto-generated cards via the power of IBM Cognos Analytics AI. IBM Cognos Analytics is using its powerful AI engine to discover new possible drivers or interesting information from your data source.
- __36. Select inside the card: [Total dollar value traded by Churn Risk and Status](#).



Total dollar value traded by Churn Risk and Status

- __37. Notice that it has added the Card to the left side in the Exploration Cards area.

The screenshot shows the IBM Watson Analytics interface. On the left, there's a sidebar with icons for Home, Cards, Data relationships, and Alerts. The main area is titled 'Cards' and shows a card titled 'Total dollar value traded by Churn Risk'. This card has a red border around its title and a small preview image. To the right of the card, there's a detailed view of the same card with the title 'Total dollar value traded by Churn Risk' and some numerical data below it. The overall interface has a clean, modern design with a blue header and a white background.

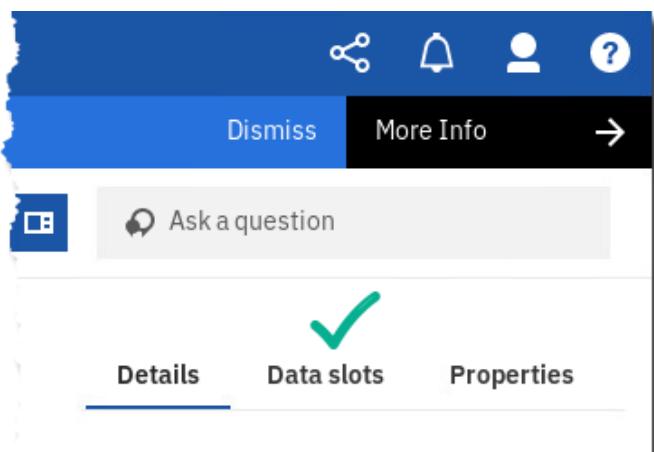
- __38. We can see that our High and Medium Churn Risk customers are mostly married. You can view more interesting discoveries by viewing the [Details](#) tab.

The screenshot shows the 'Details' tab for the 'Total dollar value traded' card. The tab has three tabs: 'Details' (which is selected), 'Data slots', and 'Properties'. The content area contains several discovery points:

- For **Total dollar value traded**, High and Low are the most important categories of **Churn Risk** with a total value of 29,995,149.55 (89.9 % of the total).
- The value of **Total dollar value traded** is unusually low when **Churn Risk** is Medium.
- For **Total dollar value traded**, M and S are the most important categories of **Status** with a total value of 32,836,463.79 (98.5 % of the total).
- The value of **Total dollar value traded** is unusually low when **Status** is D.
- The sum of **Total dollar value traded** for all values of **Churn Risk** and **Status** is 33,347,734.92.

Let's now see if Gender plays a role in how much money they trade with us.

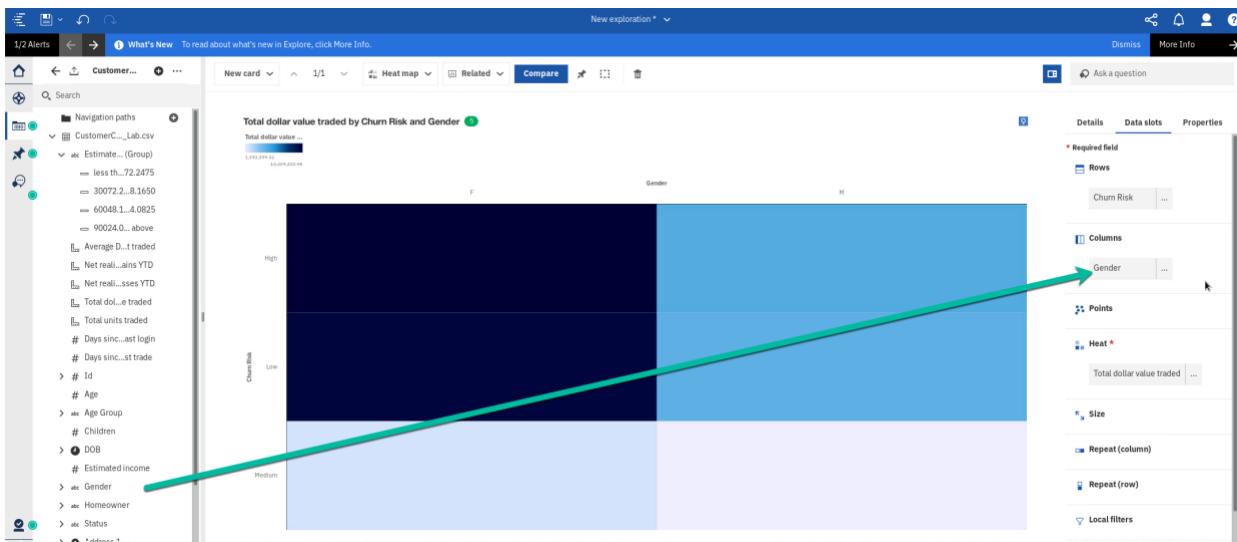
- __39. Select the Data slots tab at the top next to Details.



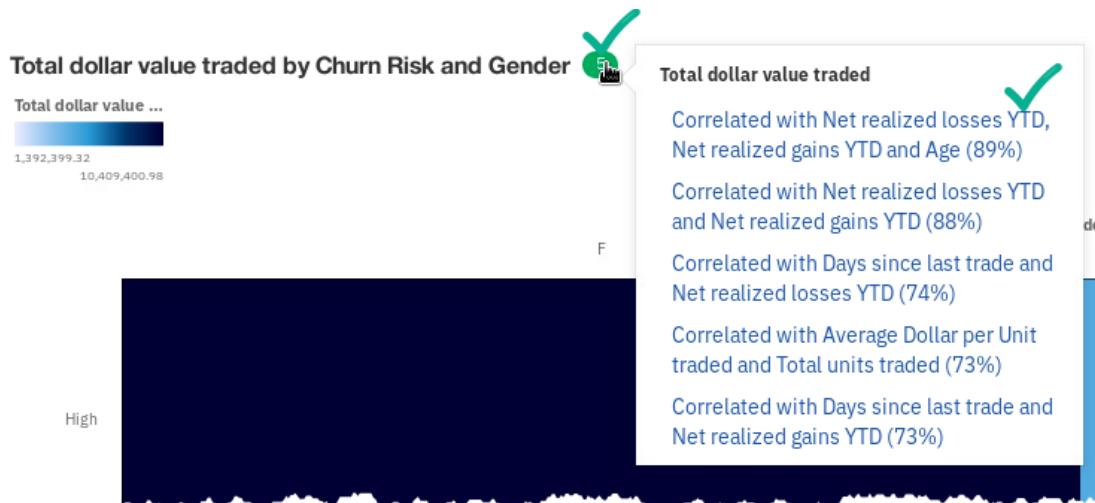
- __40. Cognos Analytics provides you interesting insight to your data automatically, but you can take control, too. Select the data tab on the far-right side.



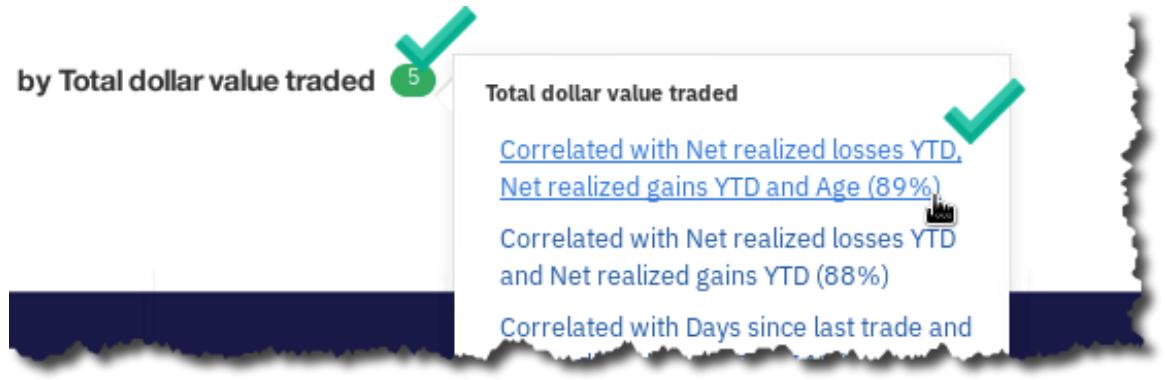
- __41. From here, drag Gender over to on top of Status.



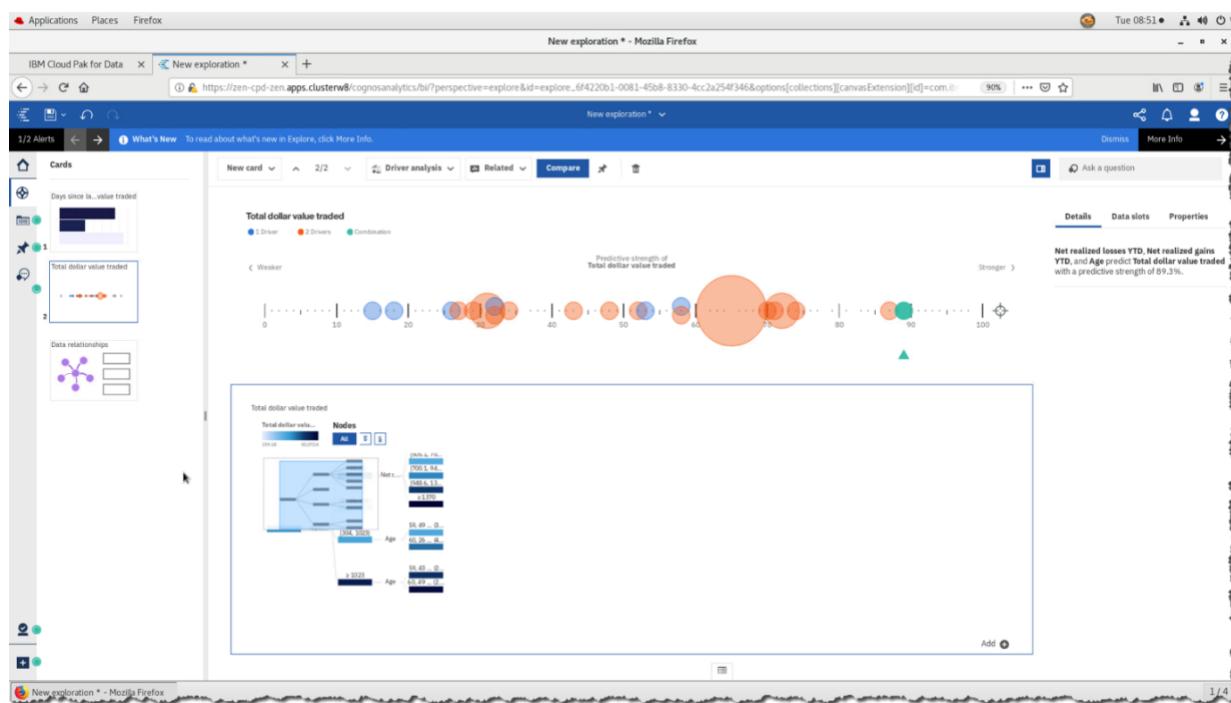
- __42. This chart is showing that our high-risk customers who trade the most dollars with us are Female and Married. Cognos Analytics also provides hints at other data points which seem to be related. Cognos Analytics found 5 other points of interest for you. Select the 5 at the top and choose Correlated with Net realized losses YTD, Net realized gains YTD and Age (89%).



- __43. Next to the title for the fly out, select the 5 \Rightarrow Correlated with Net realized losses YTD, Net realized gains YTD and Age (89%).



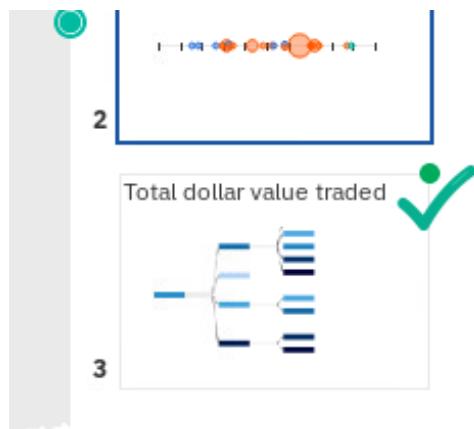
___44. Choose the **Explorations** tab on the left  and you will now see the following:



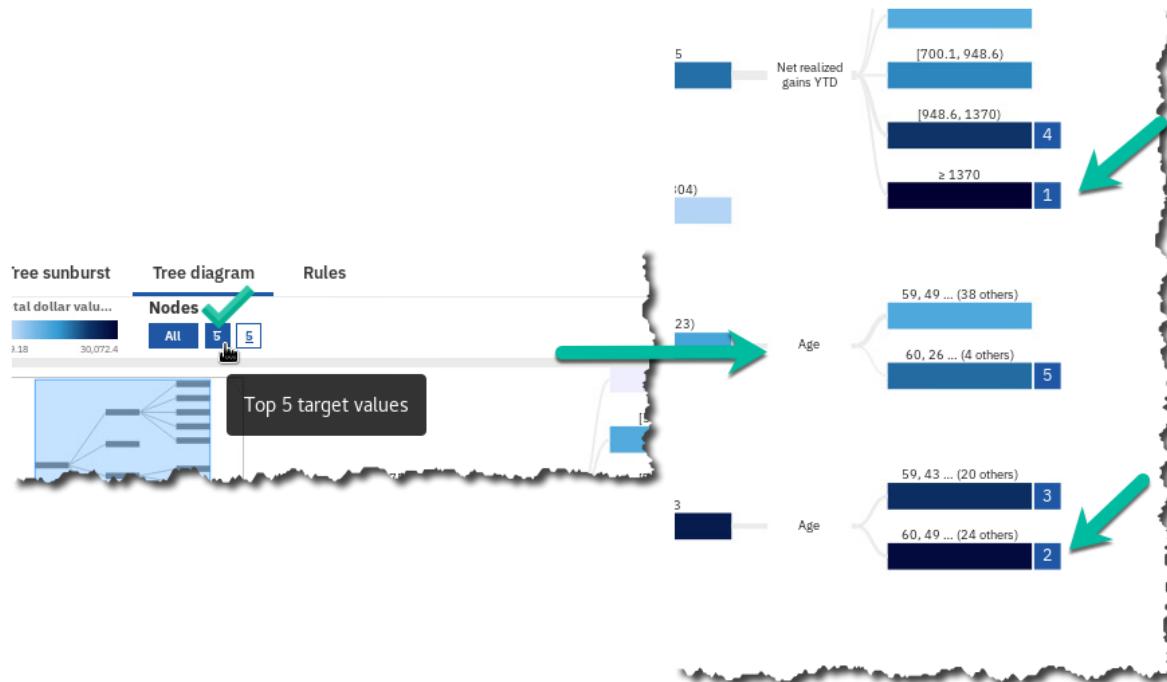
- ___45. Notice the strength of the predictor is higher when a combination of drivers is applied.
- ___46. Examine the Decision Tree at the bottom provided by IBM Cognos Analytics.
- ___47. At the bottom right, select **Add button** (to add this Decision Tree as a new Exploration card).



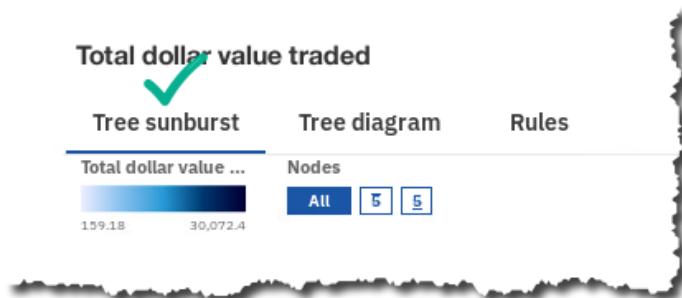
___48. Select the **Decision Tree** card at the left.



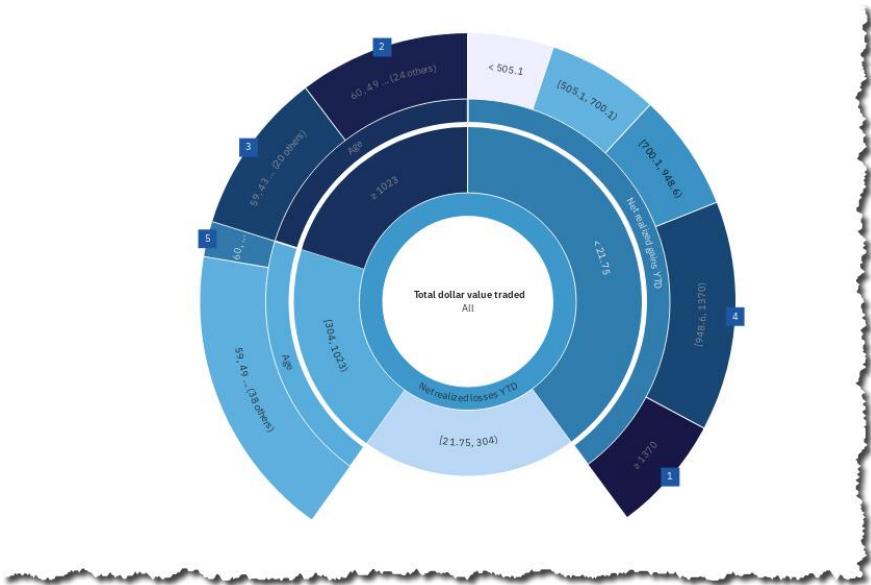
- __49. Now you are easily seeing combinations of drivers and their relationship strength in terms of Total dollar value traded. Select the Top 5 Nodes to show your top 5 nodes in the tree.



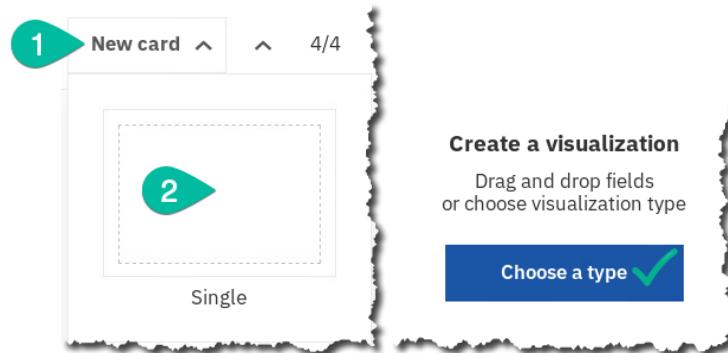
- __50. But you are not limited there. You can view other ways of looking at this data.
Choose [Tree sunburst](#) from the top tab.



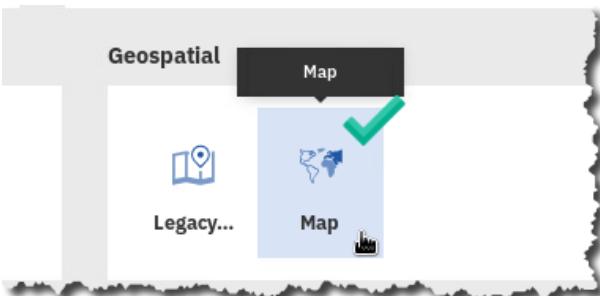
- __51. You will now be presented with the Sunburst diagram.



- __52. So far, IBM Cognos Analytics has been guiding us and creating charts based upon a metric we were looking for. We also have the ability to create exploration cards on our own.
- __53. From the top menu, select **New card** \Rightarrow **Single – Choose a Type**.



- __54. Choose **Map** under the group **Geospatial**.

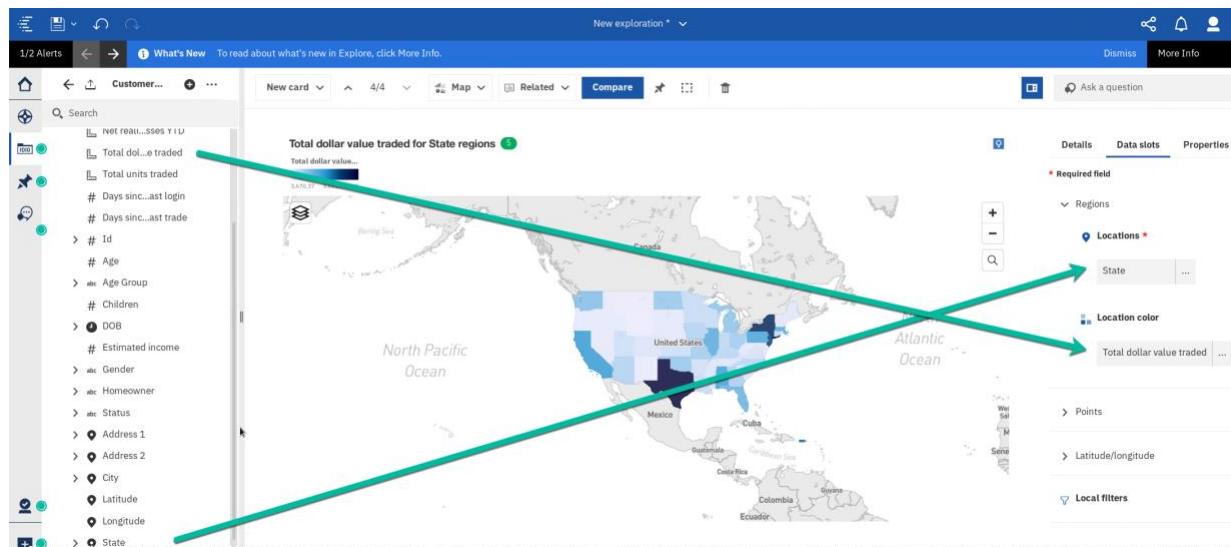


__55. Select the **Data Tab**, then under **Data slots**, expand **Regions**.



Drag **Total dollar value traded** to **Location color**.

Drag **State** to **Locations**.



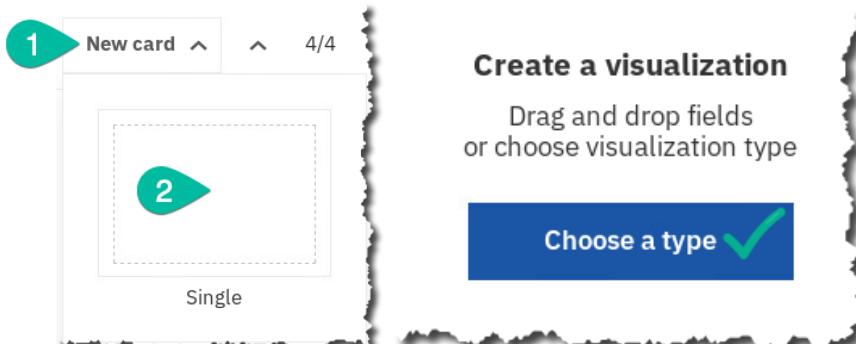
__56. Select the **Details** tab to get more information on this map.

It looks like most of our big spenders are in Texas and New York.

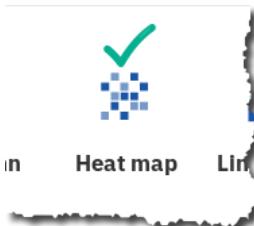
Zoom in on the U.S. and center it to focus.

__57. Let's create one more Exploration card.

From the top tab, select **New card** \Rightarrow **Single** \Rightarrow **Choose a type**.



__58. This time let's choose a **heat map** from the **Relationships** group.



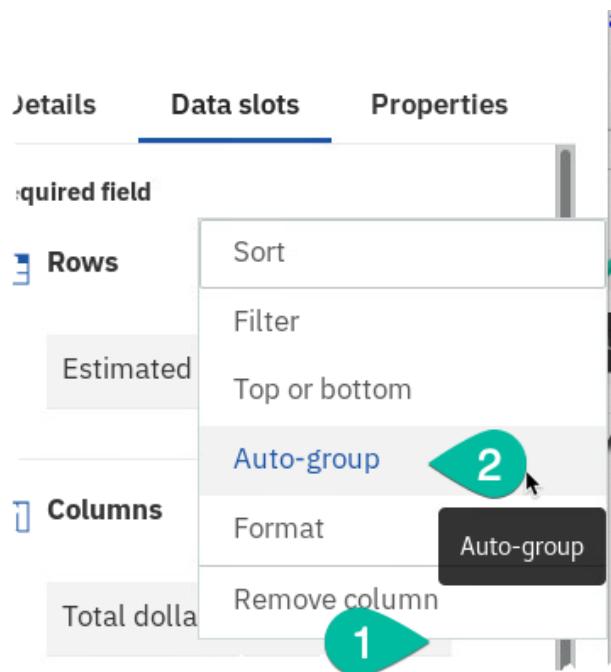
__59. Drag **Estimated income (Group)** to **Rows**.

Drag **Total dollar value traded** to **Columns**.

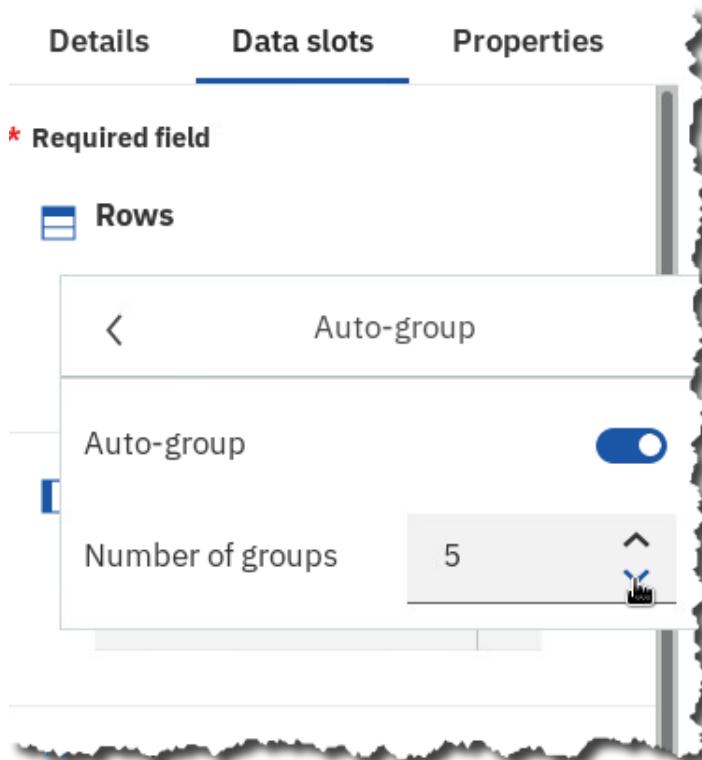
Drag **Total units traded** to **Heat**.

- _60. That looks really noisy! Let's now group that **Total dollar value traded**. Lucky for us, IBM Cognos Analytics allows you to Auto-group on the fly.

Select the **ellipses** next to **Total dollar traded** \Rightarrow **Auto-group**.

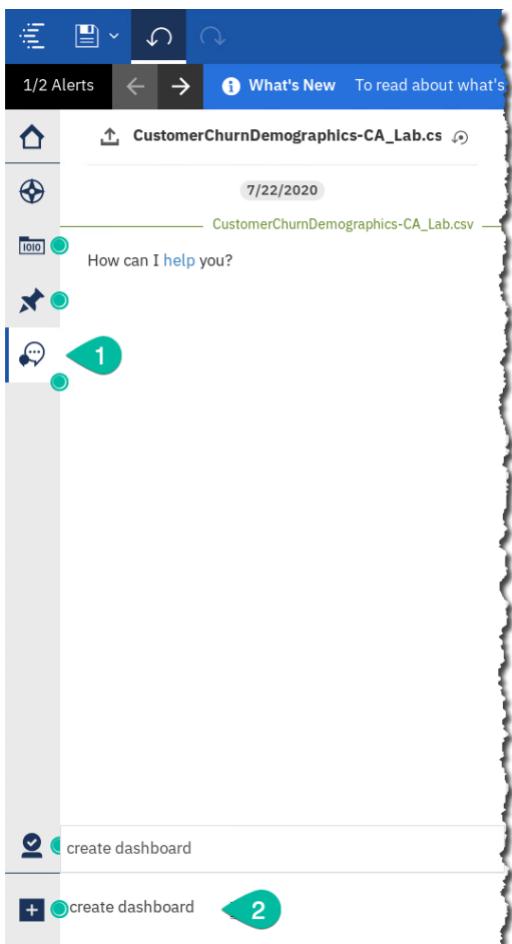


- _61. Turn on **Auto-group** and set **Number of groups** to 5.

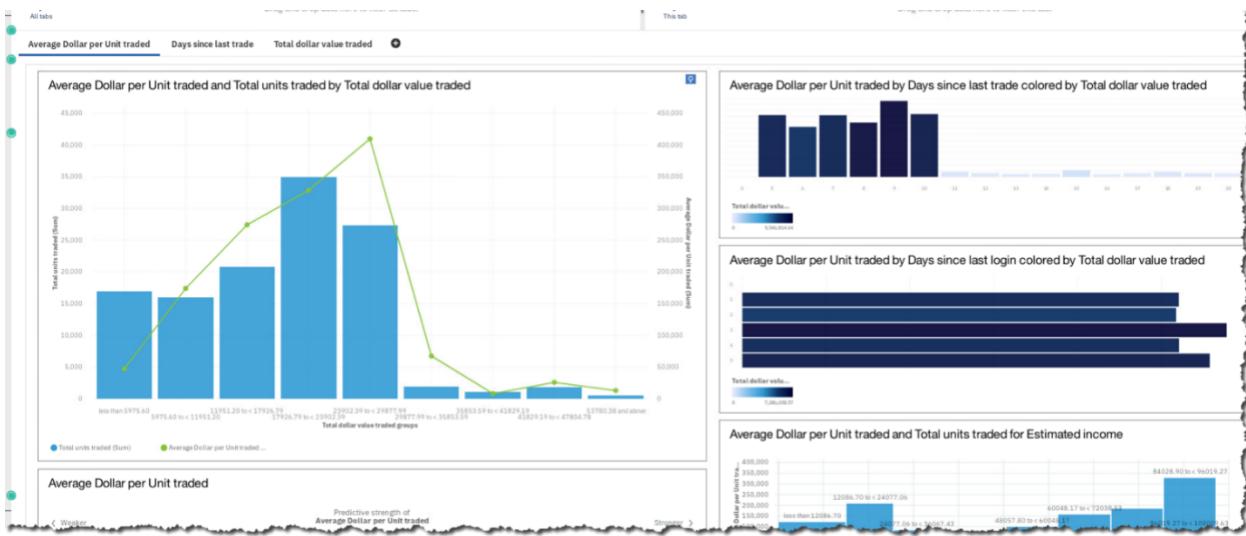


- __62. That's easier to read! Visualizing that sweet spot is now so much easier.
- __63. Now we will have IBM Cognos Analytics Assistant create a Dashboard for us.

Select the [Assistant](#) link on the left, then choose the suggestion of [create dashboard](#).



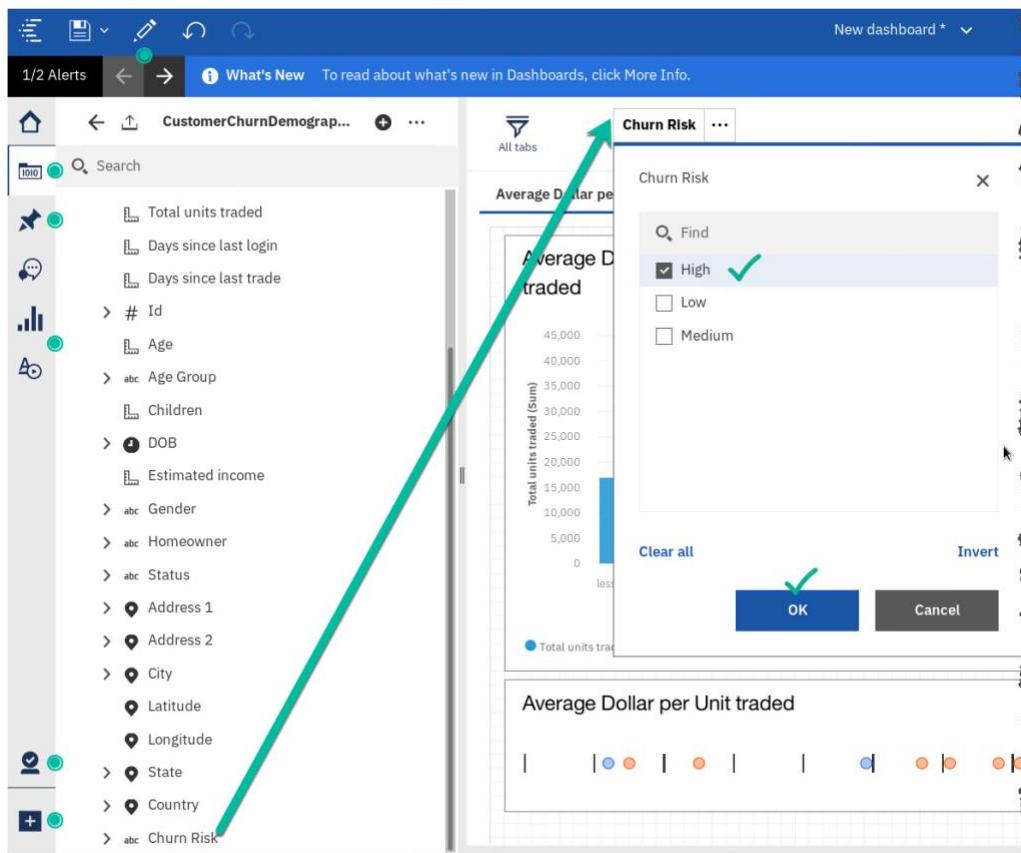
- __64. You should now see the following:



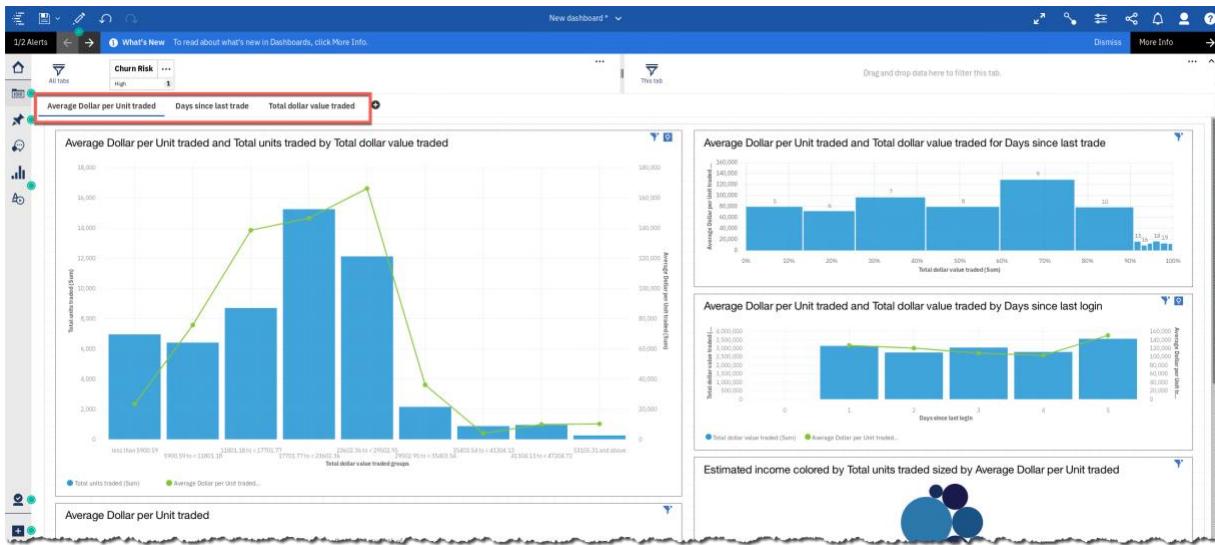
IBM Cognos Analytics Assistant created a dashboard for you!

- __65. You can investigate this dashboard by selecting the tabs at the top or hover over any of the charts and graphs. We wish to focus in on our high churn customers. So, let's filter this dashboard.

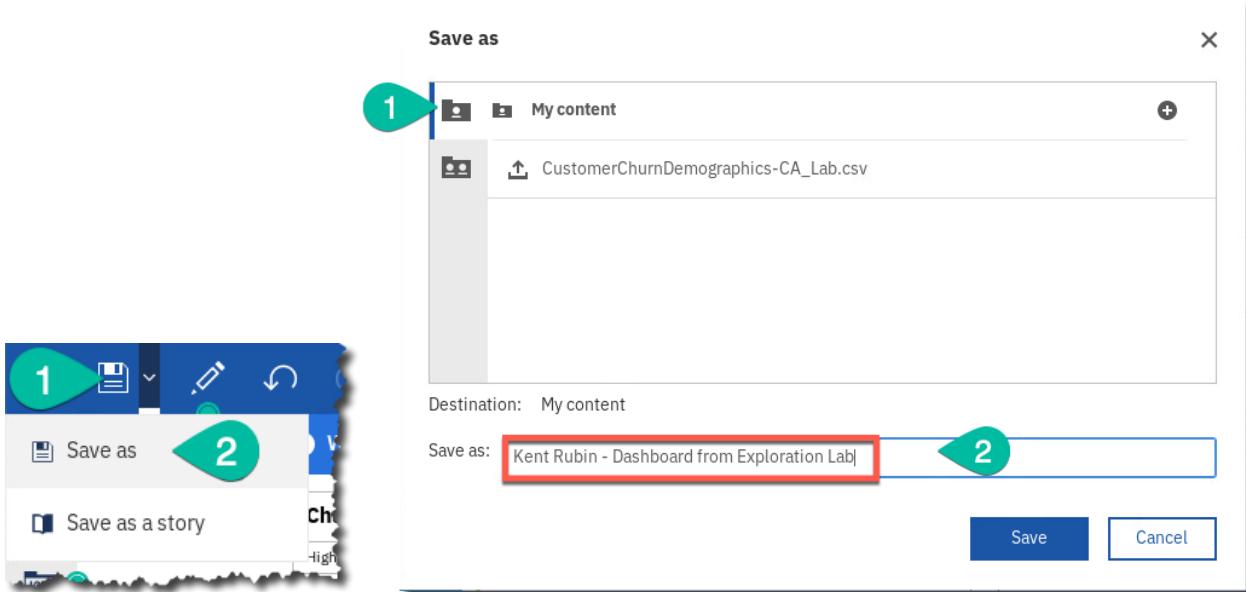
From the **Data** tab , drag **Churn Risk** to the Filter **All tabs** location .



- __66. Close the **Data items** tab by selecting it again and you should see all of your charts and tabs filtered on just your high churn risk customers. Building dashboards is easy in Cognos Analytics!



- __67. Save your Auto-created Dashboard in your Personal Folder as Your name – Dashboard from Exploration Lab.



9.10 Embedding Cognos Content

- _68. Now that you have created and saved your dashboard, you can Share your information to others. There are two options when sharing a link. You can either provide an HTML link to your dashboard or you can view the embedded code. You can also export your dashboard to a PDF file. To view all these options, choose the **Share** button at the top left.



- _69. You will then see the Link options in the flyout



- _70. Use this information to embed and link your Cognos Analytics dashboard.

9.12 Lab conclusion

The results of this Exploration and Dashboard give us a good understanding of our high churn risk customers, how much they spend and where they are located. I can make marketing decisions in these areas to ensure my business is not at risk of losing these customers.

You have now witnessed the ease at which IBM Cognos Analytics can create highly visual dashboards for greater business insight...all without need to code any expressions in a secure, governed manner.

You're off to a good start. Keep in mind, IBM Cognos Analytics offers insight and infused AI using several governed, easy to use components. Through this lab, you have experienced working with Data Modules, Explorations and Dashboards. But IBM Cognos Analytics offers other areas of insight as well, such as Reports and other areas of analysis that could not fit in this lab.

We encourage you to delve deeper into all these other avenues of insight by taking on other online IBM Cognos Analytics workshops available to you.

TIP: For a more detailed workshop on IBM Cognos Analytics Dashboards, please see:

<http://ibm.biz/Cognos-Dashboards>

TIP: For a more detailed workshop on IBM Cognos Analytics Reporting, please see:

<http://ibm.biz/Cognos-Reporting>

**** End of Lab 09 – Infuse: Cognos Analytics - Introduction**

Lab by Kent Rubin, IBM

Lab 10 WRAP-UP

10.1 Lab overview

Let's do some wrap up tasks the various personas might do after completing the entire set of "core" labs for this workshop.

Note: if you have not completed all the "core" labs in this workshop, doing some of these exercises may not make sense.

10.2 Data Scientist wrap-up



Data Scientist

10.2.1 Saving an analytics project

The Data Scientist saves her work by [Exporting](#) the project she just worked on to a file on the server in a location that the OS administrators consistently backup.

- _1. Click [Navigation Menu](#) \Rightarrow [My Projects](#) \Rightarrow [CPD Workshop Analytics Project](#).
Click on the [Export to Desktop](#) (up arrow) icon (at the top of the screen).

- _2. Click the first box to select all the assets and then [Export](#) \Rightarrow [Continue export](#).

__3. Click [Save file](#) ⇒ [OK](#).

(Note: the file will be in

10.3 Data Engineer wrap-up

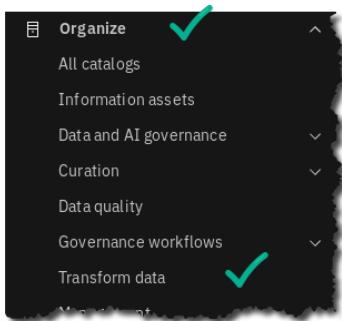


Data Engineer

10.3.1 Saving a transform project

The Data Engineer saves her work by [Exporting](#) the project she just worked on to a file on the server in a location that the OS administrators consistently backup.

__4. Click [Navigation Menu](#) ⇒ [Organize](#) ⇒ [Transform data](#).



__5. Click on [CPD_Workshop_Transform_Project](#).



- __6. Click ellipses ⇒ Export.

- __7. Here is where you choose which jobs in the project to Export. You can click them all with the top checkbox. Since there is only one job in this project, checking either box works the same.

Click Export.

Export jobs - CPD_Workshop_Transform_Project

- __8. The export is kicked off in the background. When it finishes, the file will be in [/home/ibmdemo/Downloads](#) if you are accessing via remote desktop. If using browser, you can save wherever you wish.

- __9. Review it in a terminal window.

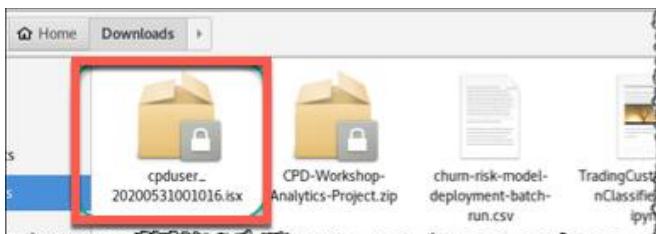
```

/home/ibmdemo$ pwd
/home/ibmdemo/Downloads
/home/ibmdemo$ ls -l
total 528
-rw-rw-r--. 1 ibmdemo ibmdemo 13929 Jun  2 12:28 churn-risk-model-deployment-batch-run.csv
-rw-r--r--. 1 root    root   51505 Jun  2 15:18 cpduser_20200531001016.lsx
-rw-r--r--. 1 root    root   444913 Jun  2 15:17 CPD-Workshop-Analytics-Project.zip
-rw-rw-r--. 1 ibmdemo ibmdemo 24172 Jun  1 17:49 TradingCustomerChurnClassifier-Py36.ipynb
[root@bastion Downloads]#

```

transform project export

_10. You can also review it in your desktop browser.



10.4 Administrator wrap-up

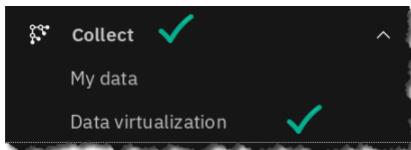


Administrator

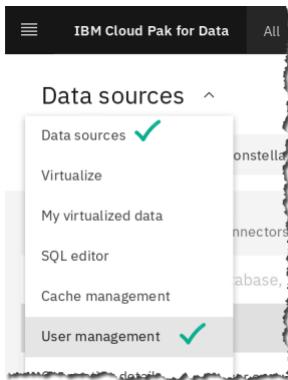
10.4.1 Granting Data Virtualization access

The Administrator has been asked to open up Data Virtualization capabilities to a trusted user called **Data Steward** which allows the user to create them himself.

_11. The **Navigation Menu** \Rightarrow **Collect** \Rightarrow **Data virtualization**.



_12. Click the **Menu dropdown** (which now says: **Data sources**) \Rightarrow **User management**.



_13. Click **Add users** + .



- __14. Select Name: **Data Steward** ⇒ Role: **Steward** ⇒ Add.

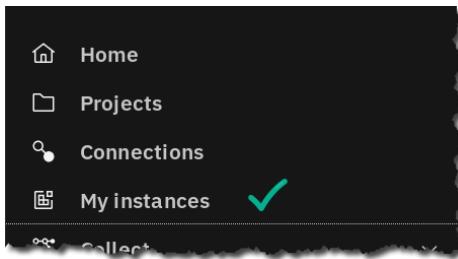
<input type="checkbox"/> Name	Username	Role
<input type="checkbox"/> Data Scientist	datascientist	User
<input checked="" type="checkbox"/> 1 Data Steward	datasteward	2 User ^
<input type="checkbox"/> Developer	developer	Admin
<input type="checkbox"/> User	user	Engineer
		3 Steward
		User

Cancel Add 4

10.4.2 Environment management

The Administrator may need to reclaim resources for environments that are long running and not in use.

- __15. Click **Navigation Menu** ⇒ **My Instances**.



Click the Environments tab.

- __16. Here you can delete any runtime environment that has been left running by a user to free up resources on the cluster.

Name	Type	ID	Project	Created by	Created on	Pods	Status	vCPU	Memory (GB)
Total						0		0.00 of 1.00	0.17 of 2.00
Default Python 3.6 Enviro...	jupc...	CPD Work...	cpduser	Jun 1, 2020	1	1	Green	0.00 of 1.00	Delete runtime

10.5 Workshop conclusion

With this workshop, you can now see how Cloud Pak for Data turns your organization's data into a critical corporate asset with end-to-end data integration and collaboration on a modern, cloud-native platform.

You have now completed your IBM Journey to Cloud and AI: Analytics Modernization using Cloud Pak for Data.

**** End of Lab 10 - Wrap-up**

Lab by Burt Vialpando and Kent Rubin, IBM

Back Page: 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.

Back Page: 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 2019.

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:

<https://www.ibm.com/legal/us/en/copytrade.shtml>

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



Please Recycle