

Hands-on lab guide

Citing your sources: Explaining generative AI output

Eric Martens

emartens@us.ibm.com

Learning Content Development, Data and AI



Loucas Loumakos

loucas.loumakos@ibm.com

Learning Content Development, Data and AI

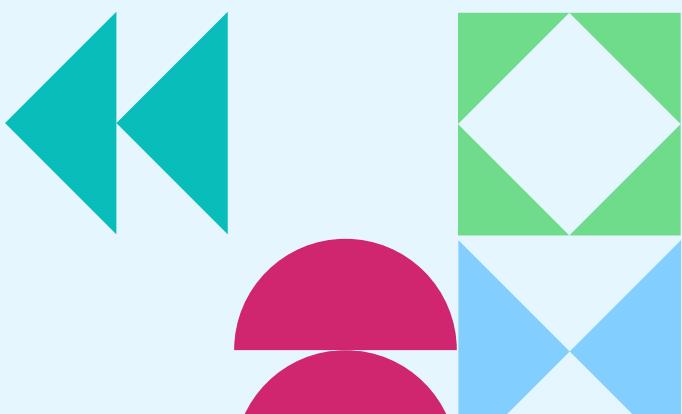
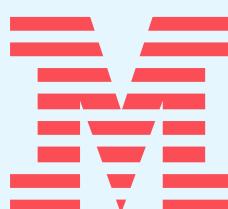


Table of Contents

1	Introduction	5
1.1	About this hands-on lab	5
2	Setting up the environment	5
2.1	Sign into IBM Cloud	5
2.2	Generate an API key	6
2.3	Sign into watsonx	8
2.4	Create a watsonx project	9
2.5	Associate a machine learning service	10
2.6	Get the project ID	11
3	Executing the Jupyter notebook	12
3.1	Create the notebook	12
3.2	Install the necessary libraries in the notebook environment	14
3.3	Update the notebook with your credentials	15
3.4	Adding and querying your own text	15
4	Congratulations	19

Notices and disclaimers

© 2024 International Business Machines Corporation. No part of this document may be reproduced or transmitted in any form without written permission from IBM.

U.S. Government Users Restricted Rights – use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM.

This document is current as of the initial date of publication and may be changed by IBM at any time. Not all offerings are available in every country in which IBM operates.

Information in these presentations (including information relating to products that have not yet been announced by IBM) has been reviewed for accuracy as of the date of initial publication and could include unintentional technical or typographical errors. IBM shall have no responsibility to update this information.

This document is distributed “as is” without any warranty, either express or implied. In no event, shall IBM be liable for any damage arising from the use of this information, including but not limited to, loss of data, business interruption, loss of profit or loss of opportunity. IBM products and services are warranted per the terms and conditions of the agreements under which they are provided. The performance data and client examples cited are presented for illustrative purposes only. Actual performance results may vary depending on specific configurations and operating conditions.

IBM products are manufactured from new parts or new and used parts.

In some cases, a product may not be new and may have been previously installed. Regardless, our warranty terms apply.”

Any statements regarding IBM's future direction, intent or product plans are subject to change or withdrawal without notice.

Performance data contained herein was generally obtained in a controlled, isolated environments. Customer examples are presented as illustrations of how those customers have used IBM products and the results they may have achieved. Actual performance, cost, savings or other results in other operating environments may vary.

References in this document to IBM products, programs, or services does not imply that IBM intends to make such products, programs or services available in all countries in which IBM operates or does business.

Workshops, sessions and associated materials may have been prepared by independent session speakers, and do not necessarily reflect the views of IBM. All materials and discussions are provided for informational purposes only, and are neither intended to, nor shall constitute legal or other guidance or advice to any individual participant or their specific situation.

Notices and disclaimers (Continued)

It is the customer's responsibility to ensure its own compliance with legal requirements and to obtain advice of competent legal counsel as to the identification and interpretation of any relevant laws and regulatory requirements that may affect the customer's business and any actions the customer may need to take to comply with such laws. IBM does not provide legal advice or represent or warrant that its services or products will ensure that the customer follows any law.

Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products. IBM does not warrant the quality of any third-party products, or the ability of any such third-party products to interoperate with IBM's products. **IBM expressly disclaims all warranties, expressed or implied, including but not limited to, the implied warranties of merchantability and fitness for a purpose.**

The provision of the information contained herein is not intended to, and does not, grant any right or license under any IBM patents, copyrights, trademarks or other intellectual property right.

IBM, the IBM logo, and ibm.com are trademarks of International Business Machines Corporation, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at "Copyright and trademark information" at

[Learn more →](#)

1 Introduction

The release of ChatGPT in November of 2022 sparked an explosion of interest in generative AI and large language models. The power and flexibility of the new models sparked a surge of creativity as organizations sought to apply new techniques to existing business problems. One of the most popular use cases for generative AI is retrieval-augmented generation (RAG), in which the model accesses additional authoritative information outside its training data set to answer questions and provide information. For example, an organization could provide its health care policy as the external dataset, and then allow the model to field employee questions about policy specifics. RAG offers the ability to customize and refine a general-purpose model for a specific use case without the time and expense of training a new model.

Additionally, the watsonx.governance platform allows organizations to generate explanations of the information provided by the RAG model, which can be critical for verifying the accuracy of the answers the model provides. In the health care policy example, the provided explanations can point to specific portions of the policy used to generate the response, allowing the user to trust that the model output is grounded in factual information, and is not a hallucination.

1.1 About this hands-on lab

In this lab, you will explore this capability by creating a project, associating it with the proper services, and running a Jupyter notebook that will allow you to see this service in action. Finally, you will explore how the Python SDK used in the lab can be easily integrated into an application that allows for the surfacing of the explanation to users.

2 Setting up the environment

In this section, you will sign in to the watsonx environment, where you can build, test, deploy, and collaborate on data science and AI projects. Watsonx offers full access management, integration with dozens of databases and data connections, and a full suite of cutting-edge data science and AI tools from prompt labs to no-code rapid prototyping to industry-standard Jupyter notebook runtimes. For this lab, you will create a watsonx project, which organizes data, code, and other assets, and integrates with watsonx.ai and watsonx.governance services.

2.1 Sign into IBM Cloud

1. To begin, you will need to sign in to [watsonx](#) using the provided credentials. From your lab reservation screen, scroll down and locate the **Username** and **Password**. Copy and paste these values into a text file.
2. Click on the **IBM Cloud Login** link. The login window opens.

The screenshot shows the 'Reservation Details' page for generating an API key. It includes fields for 'Username' (user_7jcx), 'Password' (l07iwowskxz23mc), and 'IBM Cloud API key'. A red box highlights the 'Username' field, and a red circle labeled '1' is on the 'Password' field. A red box highlights the 'IBM Cloud API key' field, and a red circle labeled '2' is on the 'IBM Cloud Login' link.

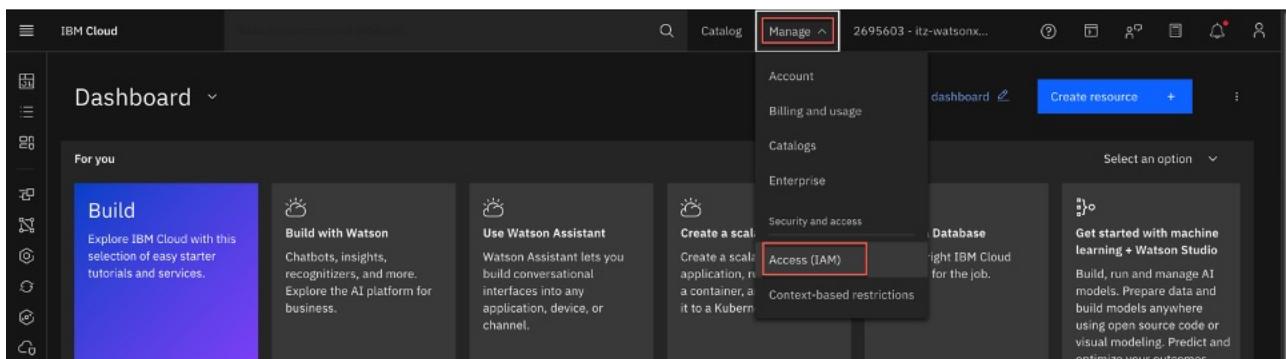
3. Enter the username and password into the appropriate fields and click the **Sign in** button. After you are signed in, you will be taken to the **IBM Cloud** dashboard.



2.2 Generate an API key

To connect to services in your reserved account, you will need an API key, which can be generated via the user interface.

1. Click on the **Manage** button to open the menu, then click on the **Access (IAM)** menu item.



2. If necessary, click on the **Manage identities** menu item to expand it, then click on the **API keys** menu item.

The screenshot shows the IBM Cloud Identity and Access Management interface. On the left, there's a sidebar with 'IAM' selected. Under 'API keys', the 'API keys' option is highlighted with a red box. The main content area has a title 'IBM Cloud Identity and Access Management' and a sub-section 'Securely authenticate users for platform services and control access to resources.' It includes buttons for 'Invite users' and 'Create access group'. To the right is a graphic illustrating security concepts like a shield, a key, and a fingerprint. At the bottom, there's a 'Whats new?' section and a note about 'Limit access with resource attribute-based conditions'.

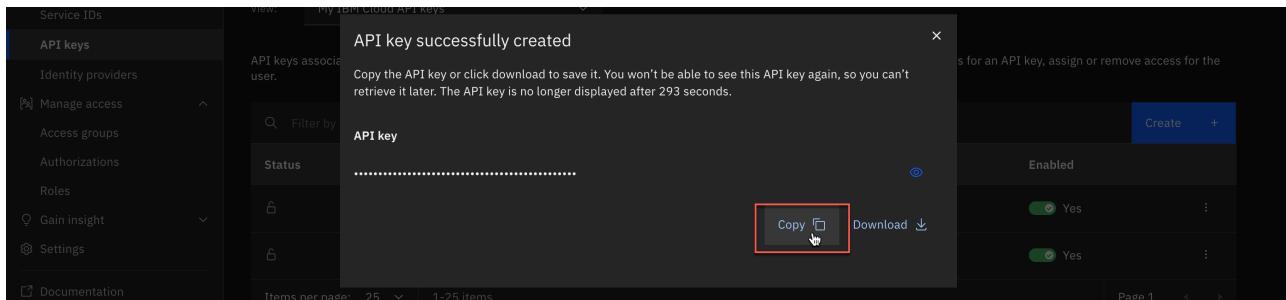
- Click on the **Create** button. The **Create IBM Cloud API key** window opens.

This screenshot shows the 'Create IBM Cloud API key' dialog box. In the 'Name' field, 'RAG API key' is entered. The 'Create' button at the bottom right is highlighted with a red box. The background shows the 'API keys' list from the previous screenshot.

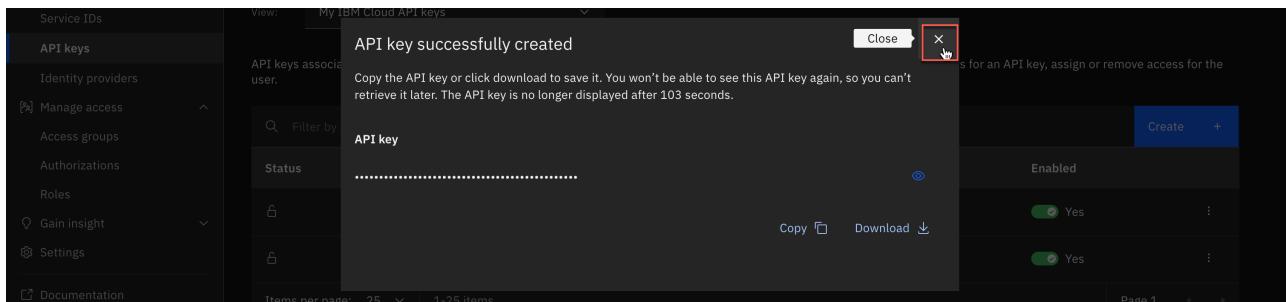
- Give your API key a name such as **RAG API key**, then click the **Create** button.

This screenshot shows the 'Create IBM Cloud API key' dialog box. The 'Name' field contains 'RAG API key'. The 'Leaked action' section has 'Disable the leaked key' selected. The 'Session creation' section has 'No' selected. The 'Create' button at the bottom right is highlighted with a red box. The background shows the 'API keys' list from the previous screenshots.

- Once the key has been created, click the **Copy** button to copy it to your clipboard, then paste the value into the text file that contains your login credentials. When you edit the Jupyter notebook later in the lab, this value will be used as the **API_KEY**.



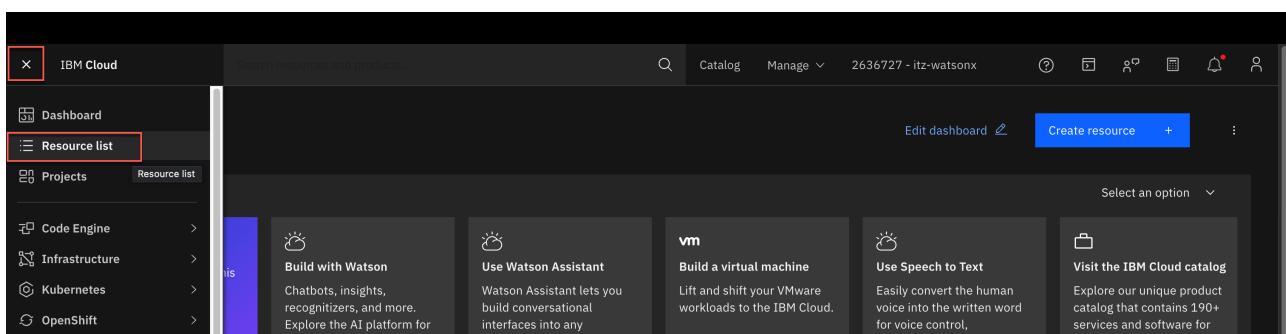
- Click on the X button to close the API key window.



You now have the credentials necessary to proceed with the lab. In the next step, you will sign into the watsonx platform.

2.3 Sign into watsonx

- Click on the hamburger menu in the upper left to open the menu, then click on the **Resource list** menu item. The **Resource list** page opens.

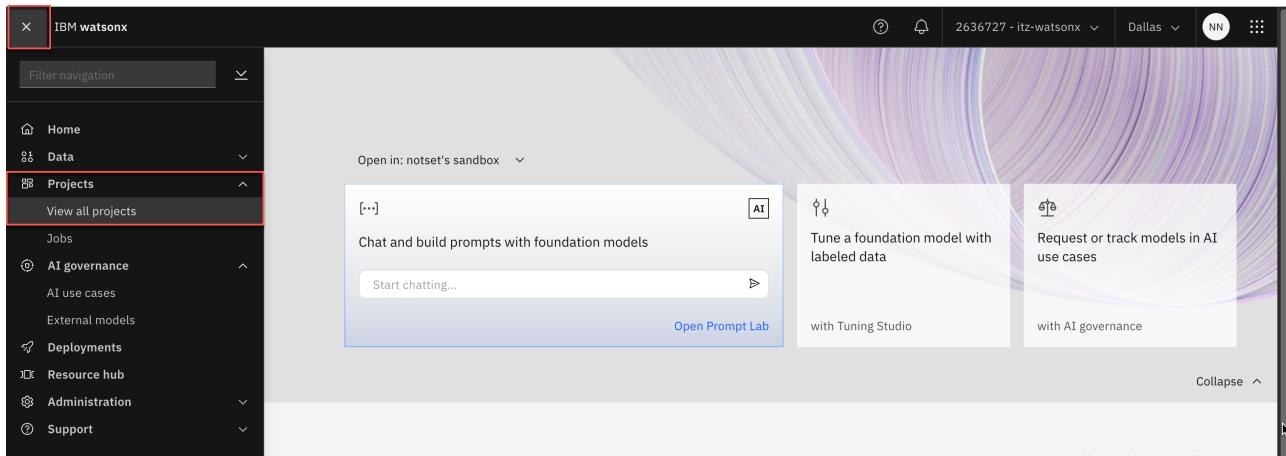


- Locate the **AI / Machine Learning** item from the list and click on it to expand it. Locate the **Watson OpenScale** item in the expanded section and click on it. Note that the name may differ slightly from the one in the screenshot below.

3. Click on the **Launch watsonx.governance** button. If you are asked to provide any additional information, wait a moment, then click the **Continue** button when it becomes enabled. Note that you may need to wait for up to two minutes for the watsonx.governance services to be prepared.
4. Read the terms and policies if you wish, then check the box to agree, and click the **X** button to close the popup window without taking the tour.

2.4 Create a watsonx project

1. Click on the **hamburger menu** in the upper left to expand it, then click on the **Projects** item to expand it. Finally, click on the **View all projects** menu item.



2. Click the **New project** button to create a new project. The **Create a project** screen opens.
3. Fill out the project details, including a **Name** and optional **Description**.
4. Use the **Select storage service** dropdown to select the cloud object storage instance for your environment. If you have multiple options, select the one that most resembles the one in the screenshot below.

The screenshot shows the 'Create a project' form. On the left, there's a '+ New' button, 'Local file' option, and 'Sample' option. The main area is titled 'Define details' and contains fields for 'Name' (with 'RAG explanations' entered) and 'Description (optional)' (with 'Explain the output of RAG queries'). Below these are 'Tags (optional)' and 'Add tags' fields. Under 'Define storage', there's a 'Select storage service' dropdown with 'Target Cloud Object Storage Instance' checked and 'itzcos-110000b3qc-kwzdy' selected. A note at the bottom says 'Project includes integration with Cloud Object Storage for storing project assets.' Red numbers 3 and 4 are overlaid on the screenshot to indicate specific steps: 3 points to the 'Name' field, and 4 points to the 'Select storage service' dropdown.

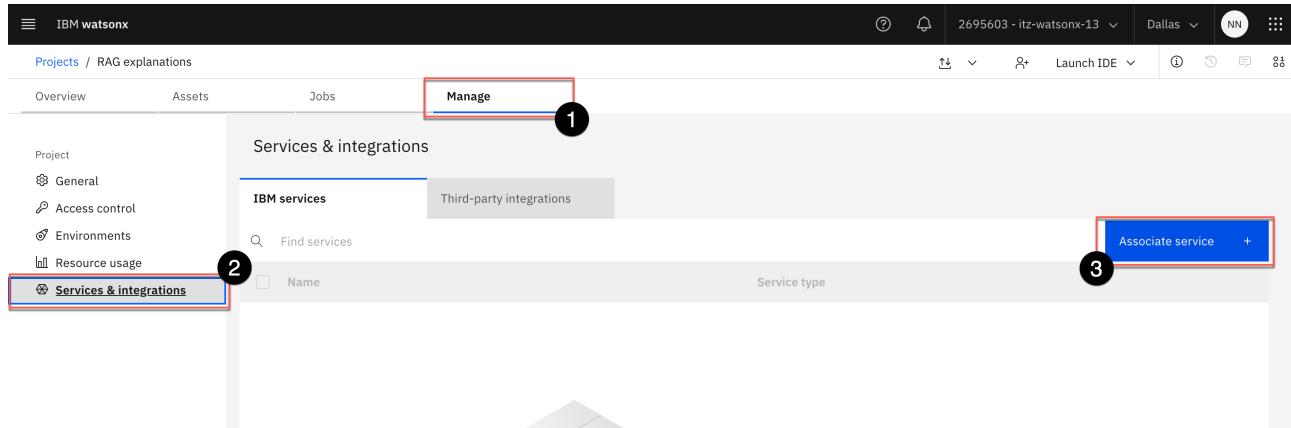
5. Click the **Create** button to create your project.

2.5 Associate a machine learning service

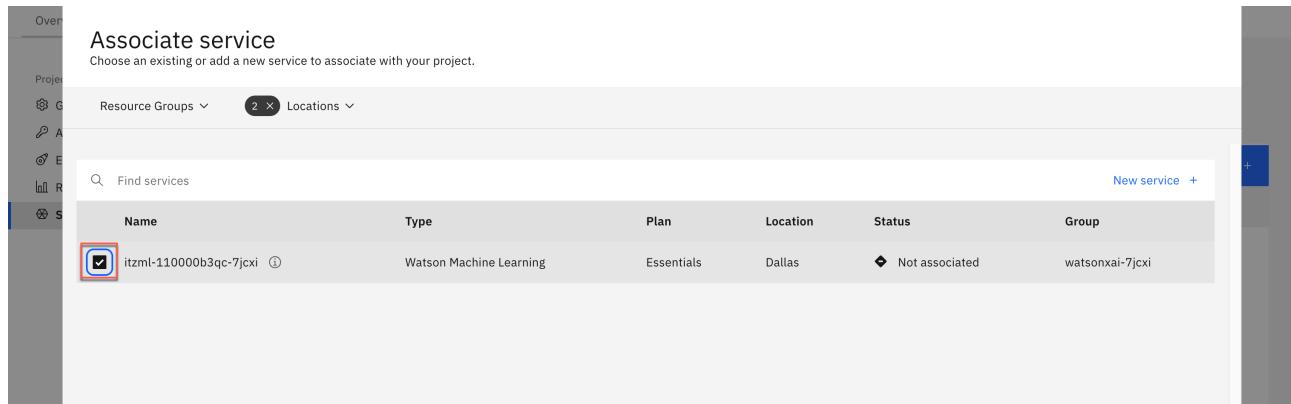
For project assets such as Jupyter notebooks to access the IBM Foundation models to perform RAG queries, you must associate a machine learning service with the project. Your lab environment comes with a provisioned machine learning service.

1. From your project screen, click on the **Manage** tab.

2. From the list on the left, click on the **Services & integrations** menu item.
3. Click the blue **Associate service** button. The **Associate service** window opens.



4. Locate and check the box to the left of the Watson Machine Learning service in the list.



5. Click the **Associate** button to associate the service.

2.6 Get the project ID

Getting the ID of the project will allow you to query watsonx.ai Foundation models using the service you just associated.

1. Beneath the **Project** header on the left, click on the **General** menu item.

The screenshot shows the IBM WatsonX interface. At the top, there's a navigation bar with 'IBM watsonx', a search bar, and various dropdowns and icons. Below it, a breadcrumb trail says 'Projects / RAG explanations'. The main area has tabs for 'Overview', 'Assets', 'Deployments', 'Jobs', and 'Manage', with 'Manage' being the active tab. On the left, a sidebar titled 'Project' has sections for 'General' (which is highlighted with a red box), 'Access control', 'Environments', 'Resource usage', and 'Services & integrations'. The main content area is titled 'Services & integrations' and shows a sub-section 'IBM services (1)'. It lists a single service entry: 'Name: itzml-110000b3qc-7jcx1' and 'Service type: Watson Machine Learning'. There's also a 'Find services' search bar and an 'Associate service' button.

2. In the **Details** tile, locate the **Project ID** item and click the **Copy to clipboard** icon to copy the value. Paste it into the text file you have been using to store credentials. In the Jupyter notebook, this value will be the **PROJECT_ID**.

The screenshot shows the 'Details' tile for the project 'RAG explanations'. It includes sections for 'Name' (RAG explanations), 'Description' (Explain the output of RAG queries), 'Tags' (Add tags to make projects easier to find), and 'Project ID' (da573099-8b14-4743-95bc-783f450f70f). A 'Copy to clipboard' button is located below the Project ID field, and both the Project ID field and the 'Copy to clipboard' button are highlighted with red boxes. To the right of the tile, there's a 'Storage' section showing 0 Bytes used and a bucket named 'ragexplanations-donotdelete-pr-4ovfk8enoocpqr'.

Your environment is now fully configured, and you have all of the credentials necessary to access Foundation models. In the next step, you will create a Jupyter notebook to query the models.

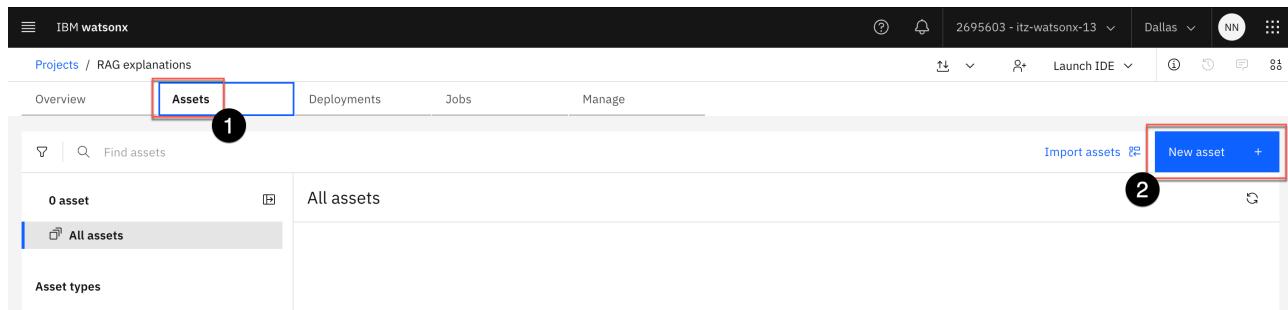
3 Executing the Jupyter notebook

In this section of the lab, you will load a Jupyter notebook from a GitHub repository, input your credentials, and run it.

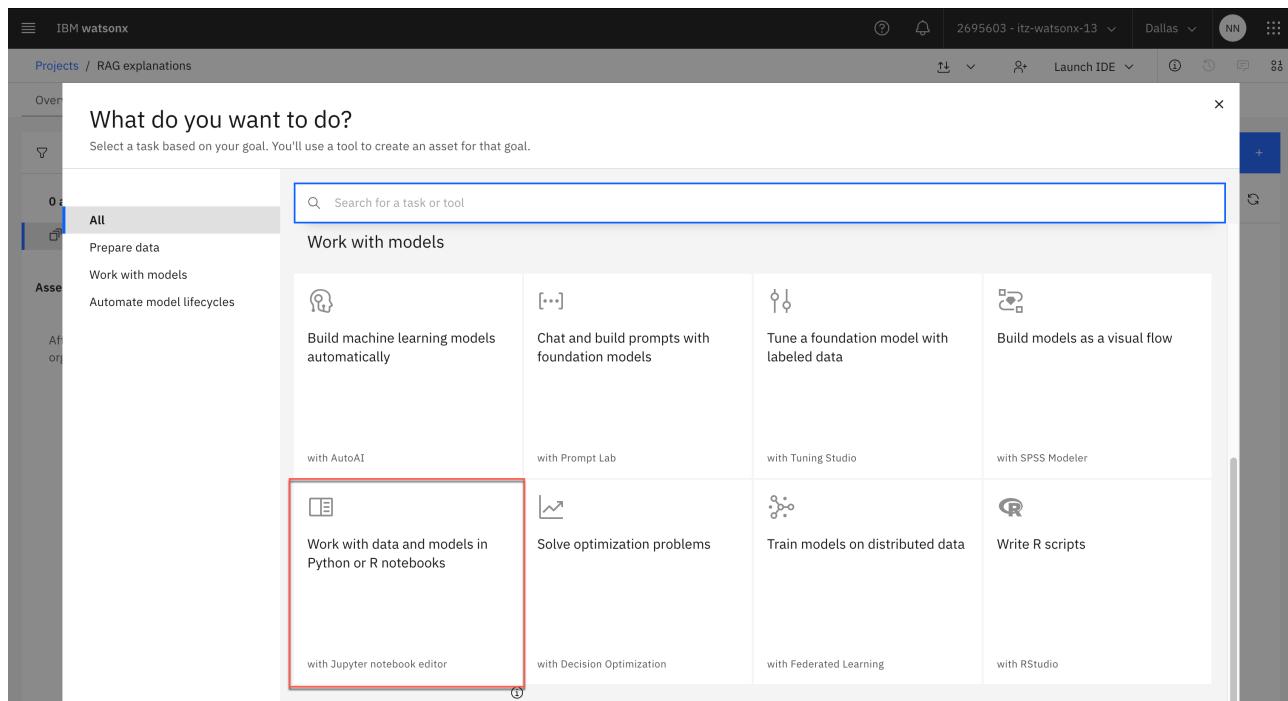
3.1 Create the notebook

1. Click on the **Assets** tab of your project.
2. Click on the **New asset** button. The **What do you want to do?** window opens.

IBM TechXchange

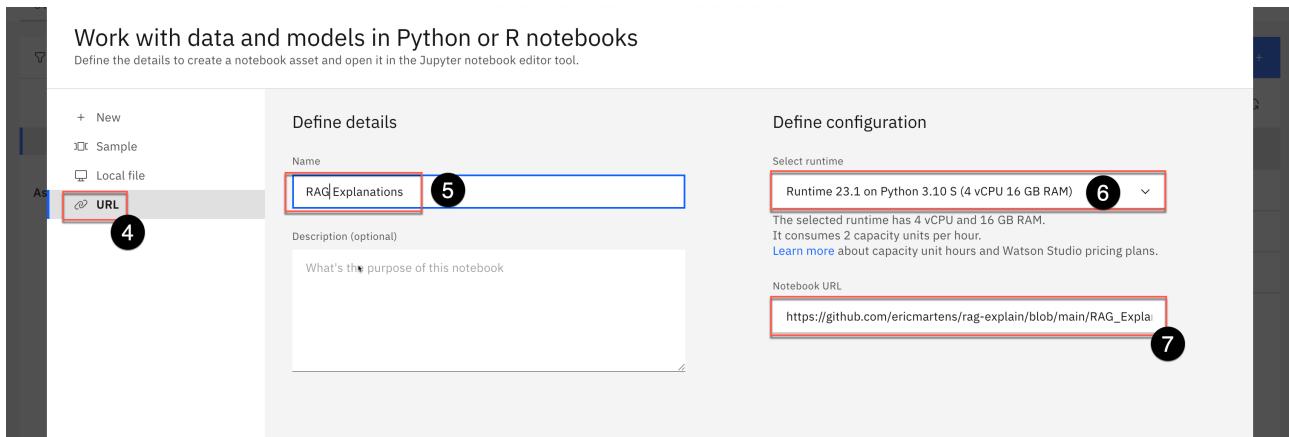


3. Scroll down to the **Work with models** section and click on the **Work with data and models in Python or R notebooks** tile.



4. Click on the **URL** menu option on the right.
5. Give your notebook a name such as **RAG Explanations**.
6. Click on the **Select runtime** dropdown and select **Runtime 23.1 on Python 3.10 S...** from the list.
7. Copy and paste the following value into the **Notebook URL** field:

```
https://github.com/ericmartens/rag-explain/blob/main/RAG_Explanations_TXC_2024.ipynb
```



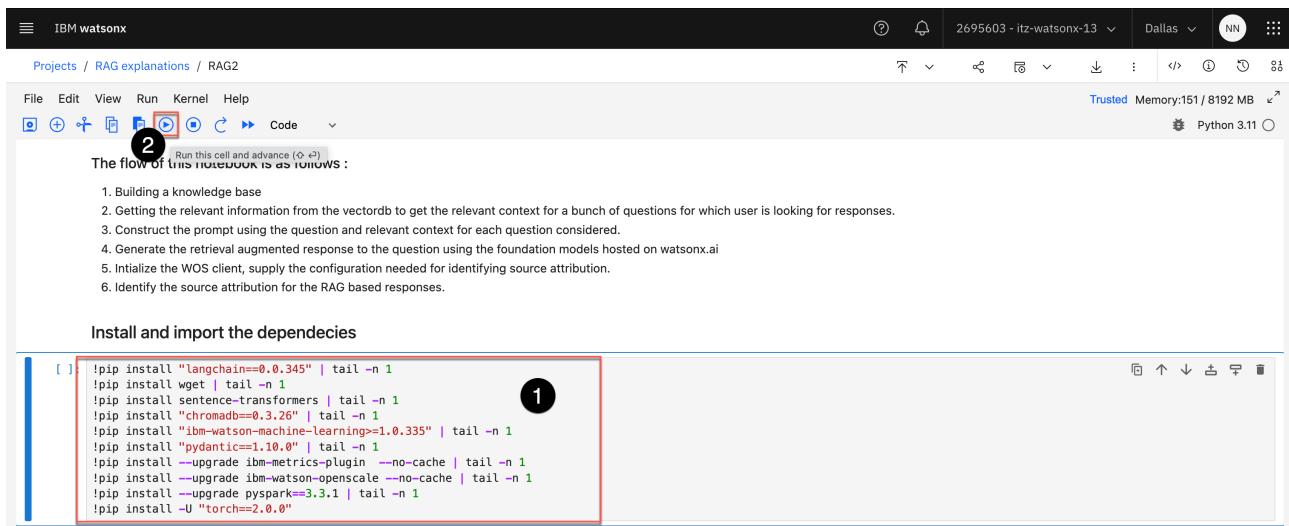
- Click the **Create** button. The notebook will be created in your project, which can take up to two minutes to complete. Note that if the notebook creation step is taking too long, refreshing the page can often fix the problem.

3.2 Install the necessary libraries in the notebook environment

Once the notebook has been created and appears on your screen, you must install additional libraries into the environment.

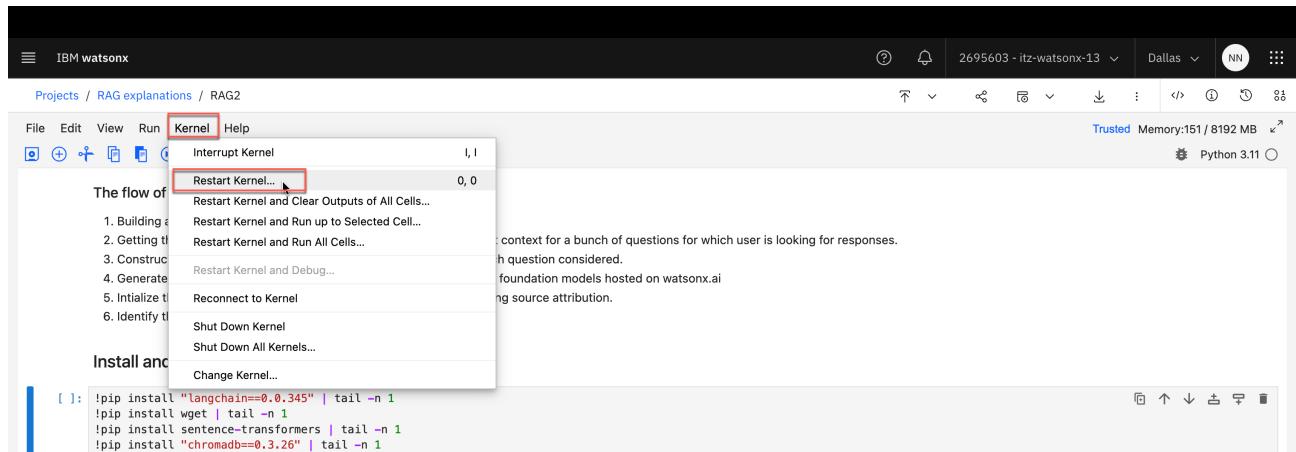
- Scroll down, locate and click on the first code cell beneath the **Install and import dependencies** header. A blue highlight box will appear around the active cell.
- Run the cell by clicking the **Run** button from the notebook menu, or by using the **shift + return** keyboard shortcut. The notebook will use the pip utility to install a variety of libraries into the environment for the notebook to use. Note that this cell may generate some warning and error messages; these can be ignored.

As the cell is running, an asterisk (*) will appear in the box to the left of the code in the cell. When the install process is finished, the asterisk will be replaced by a number, indicating that the cell has completed.



- When the cell has completed running, you will need to restart the kernel to ensure that the environment is updated. Click on the **Kernel** menu item to open it, and then click on the **Restart kernel** menu option. A

warning notification window will ask you to confirm your choice. Click the **Restart** button to confirm. The environment's kernel will restart.



- As the kernel restarts, a **Kernel restarting** status message will display in the upper right of the notebook screen. When the restart finishes, this status will disappear, and you may proceed with the rest of the lab.

3.3 Update the notebook with your credentials

When the kernel finishes restarting, you will need to add the API_KEY and PROJECT_ID credentials you gathered in previous steps.

- Locate and click on the first code cell beneath the **Edit the two values below...** cell to make it active.
- Copy and paste the **API_KEY** and **PROJECT_ID** values that you obtained in previous steps into the cell between the quotation marks on each line. When you are finished, your cell should resemble the one in the screenshot below:



Run the notebook

At this point, you may run through the notebook one cell at a time, beginning with the credentials cell you just edited. Notebook cells can be run by clicking on the **Run** button from the menu, or by using the **shift + return** keyboard shortcut. Note that some warning messages may appear; these are normal and can be ignored.

When you reach the **Custom information** section of the notebook, you may return to this guide for information on adding your own data to the project for use in the notebook.

3.4 Adding and querying your own text

It is easy to add your own custom text for queries and explanations. Prior to completing this step, you will need to save the text to a file on your machine so that it can be imported into the notebook.

IBM TechXchange

1. Click on the empty code cell beneath the **Custom content** title cell. Your cursor should appear in the empty code cell, and a blue highlight box will appear around the cell.
2. Click on the **Code snippets** button on the top right of the screen to open the **Code snippets** panel.
3. Click on the **Read data** tile.

The screenshot shows the IBM WatsonX interface with the 'RAG explanations / RAG2' project selected. In the top right, there's a 'Code Snippets' button with a red box around it and the number '2'. Below it is a 'Code Snippets' panel with a heading 'Data Ingestion'. Inside the panel, there's a 'Read data' section with a red box around it and the number '3'. The main workspace shows some Python code in a code cell:

```
import warnings
warnings.filterwarnings("ignore")
results = client.ai_metrics.compute_metrics(configuration=config_json,data_frame=data)

metrics = results.get("metrics_result")
results = metrics.get("explainability").get("protodash")

import json
for idx, entry in enumerate(results):
    print(f"====idx:{idx}: Question:{questions[idx]} Response:{data['generated_text'][idx]}====")
    print(json.dumps(entry,indent=4))
```

Below the code cell is an 'Explanation' section with a note about source attribution. To the right of the code cell is a 'Data Ingestion' panel with sections for 'Read data' and 'Prompt engineering'.

4. Click on the **Upload a data file** button.

The screenshot shows the same WatsonX interface as above, but the 'Upload a data file' button in the 'Data Ingestion' panel is highlighted with a red box and the number '4'. The 'Data Ingestion' panel also has a magnifying glass icon and a note about no data assets found. The main workspace shows the same Python code as before.

5. Click on the **Drop data files here...** link to browse to a file on your machine, and select the custom text you wish to work with.

The screenshot shows the WatsonX interface with the 'Upload data files' panel open. The 'Drop data files here or browse for files to upload' area is highlighted with a red box and the number '5'. The main workspace shows the Python code from earlier steps.

IBM TechXchange

6. Once the file has finished uploading, it has been added as an asset in your project, and can be used by the notebook. Click on the **Code snippets** button once again.
7. Click on the **Read data** tile.

The screenshot shows the IBM Watsonx interface with the 'Code Snippets' panel open. The 'Data Ingestion' section is highlighted with a red box and contains a 'Read data' tile. The tile has a sub-tile 'Generate a code snippet to load data from a data asset or connection into your notebook.' A red box highlights the 'Select data from project' button in this sub-tile. The main notebook area shows some Python code related to RAG explanations and protodash explainers.

8. Click on the **Select data from project** button. The **Select data from project** window opens.

The screenshot shows the 'Select data from project' window open in the foreground. The 'Data assets' section is highlighted with a red box and shows a single item: 'custom_data.txt'. The main notebook area shows some Python code related to RAG explanations and protodash explainers.

9. Click on the **Data asset** category, then click on the file you imported in a previous step and click on the **Select** button.

The screenshot shows the 'Select data from project' window with the 'Selected assets' panel open. It displays the details for the selected asset: 'custom_data.txt'. The asset is a 'Data asset' type, 5 KB in size, with a mime type of 'text/plain'. It was created at 2024/10/22 10:14:45 and last updated at the same time. A red box highlights the 'Selected assets' panel.

10. Click on the **Insert code to cell** button.

`print(json.dumps(entry, indent=4))`

Explanation

Source attribution can be understood using the weights (the attribution/contribution factor) and the prototypes (the relevant context/source) which has attributed to the response by the foundation model behind the scenes. For example a weight: 1.0 indicate that that a single paragraph of the context has attributed for response by foundation model. Likewise weights : 0.6,0.3,0.1 indicate that 3 paragraphs have attributed for response by foundation model behind the scenes. The prototype values are the paragraphs supplied as part of the relevant context.

Custom content

Please refer to your lab guide for instructions on using custom text in the lab.

[]: Insert code to cell

11. The notebook will create a code snippet to import the contents of the text file as a StreamingBody object into to the notebook. Run the cell to execute it, making note of the variable name used to hold the contents of the file.

Custom content

Please refer to your lab guide for instructions on using custom text in the lab.

[]: Generate a code snippet to load data from a data asset or connection into your notebook.

Selected data: `custom_data.txt`

Load as: `StreamingBody object`

```
import os, types
import pandas as pd
from botocore.client import Config
import ibm_boto3

def __iter__(self): return 0

# @hidden_cell
# The following code accesses a file in your IBM Cloud Object Storage. It includes your credentials.
# You might want to remove those credentials before you share the notebook.

cos_client = ibm_boto3.client(service_name='s3',
    ibm_api_key_id='QnRadVnE9hmwgWL1tFrnddBAByroG0dBE7IMlws7Kt',
    ibm_auth_endpoint='https://iam.cloud.ibm.com/identity/token',
    config=Config(signature_version='oauth'),
    endpoint_url='https://s3.direct.us-south.cloud-object-storage.appdomain.cloud')

bucket = 'ragexplanations-donotdelete-pr-4ovfk8enoocpqr'
object_key = 'custom_data.txt'

# load data of type "text/plain" into a botocore.response.StreamingBody object.
# Please read the documentation of ibm_boto3 and pandas to learn more about the possibilities to load the data.
# ibm_boto3 documentation: https://ibm.github.io/ibm-cos-sdk-python/
# pandas documentation: http://pandas.pydata.org/

streaming_body_1 = cos_client.get_object(Bucket=bucket, Key=object_key)['Body']
```

12. Note that each time you use the code snippet to insert a new set of text, the variable name will increment (`streaming_body_1`, `streaming_body_2`...). You will need to ensure that the variable used in the next code cell is updated to match the variable used to store the imported data.

ibm_auth_endpoint='https://iam.cloud.ibm.com/identity/token',
config=Config(signature_version='oauth'),
endpoint_url='https://s3.direct.us-south.cloud-object-storage.appdomain.cloud'

bucket = 'ragexplanations-donotdelete-pr-4ovfk8enoocpqr'
object_key = 'custom_data.txt'

load data of type "text/plain" into a botocore.response.StreamingBody object.
Please read the documentation of ibm_boto3 and pandas to learn more about the possibilities to load the data.
ibm_boto3 documentation: https://ibm.github.io/ibm-cos-sdk-python/
pandas documentation: http://pandas.pydata.org/

`streaming_body_1 = cos_client.get_object(Bucket=bucket, Key=object_key)['Body']`

Store the custom text in a variable

[9]: `custom_text = streaming_body_1.read()`

StreamingBody object

13. You may now proceed with the rest of the notebook, returning to the custom text loading cell to load new files if you wish.

4 Congratulations

You have explored the RAG explanation capabilities of watsonx.ai and watsonx.governance. Thank you for attending this session, and enjoy the remainder of your conference.