

Final Project 1 Solution

1. Finding and Analyzing Data:

- Locate an open-source dataset.
- Write technical and non-technical observations about the dataset (schema, data quality, potential improvements).

The dataset I have picked for this project is from UCI Machine Learning Repository

[Home - UCI Machine Learning Repository](#) with the subject area “Climate and Environment”. Here is the link to find the data set [Forest Fires - UCI Machine Learning Repository](#),

[Chase the Fire \(kaggle.com\)](#)

This predicts the burned area of forest fires in Portugal, it does not have any missing values

About the dataset and Observations:

Technical Observations:

- **Schema:** The dataset consists of 13 columns. The columns are "X", "Y", "month", "day", "FFMC", "DMC", "DC", "ISI", "temp", "RH", "wind", "rain", and "area".
- **Data Types:** The dataset contains a mix of numerical and categorical data. "X", "Y", "FFMC", "DMC", "DC", "ISI", "temp", "RH", "wind", "rain", and "area" are numerical, while "month" and "day" are categorical.
- **Missing Values:** From the given data, it appears there are no missing values. However, a more thorough check would be needed on the full dataset.
- **Data Quality:** There are no obvious inconsistencies in the data. Outliers would need to be checked for in the numerical data.

Non-Technical Observations:

- **Purpose of Data/insights:** The data seems to be related to forest fires and includes various weather and fire danger index conditions. It could be used to understand the relationship between these conditions and the occurrence or severity of forest fires.
- **Limitations/Biases:** The data seems to be limited to specific locations (as indicated by the X and Y coordinates). If these locations are not representative of the area as a whole, this could introduce bias. Also, many fires seem to have an area of 0, which could indicate an imbalance in the data.
- **Temperature:** The ‘temp’ column seems to vary quite a bit, indicating a range of weather conditions.
- **Rain:** The ‘rain’ column has a lot of zero values, suggesting that rain is not a common occurrence in these observations.

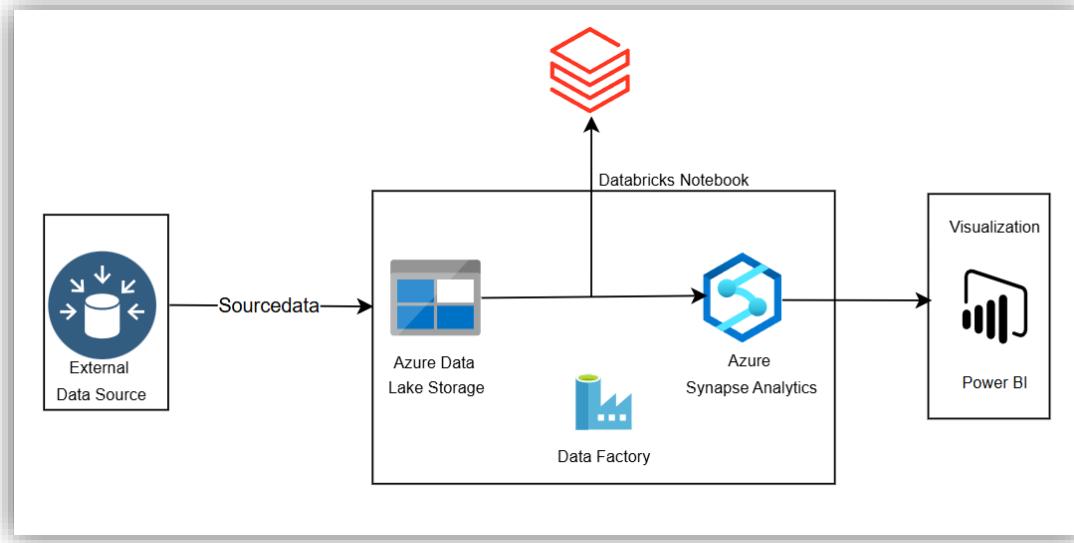
- **Area:** The ‘area’ column also has many zero values. This could indicate areas where no forest fire was observed.
- **Month and Day:** The dataset includes observations from different months and days, indicating it spans across different seasons and days of the week.

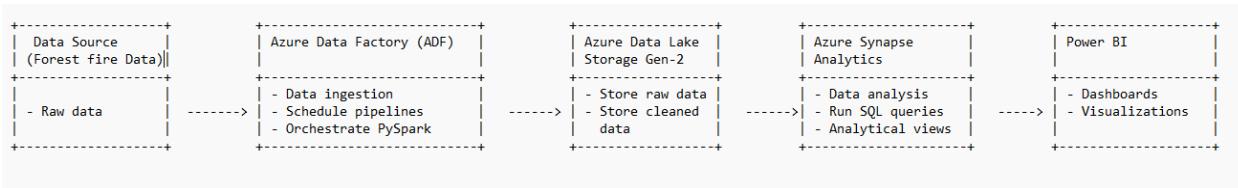
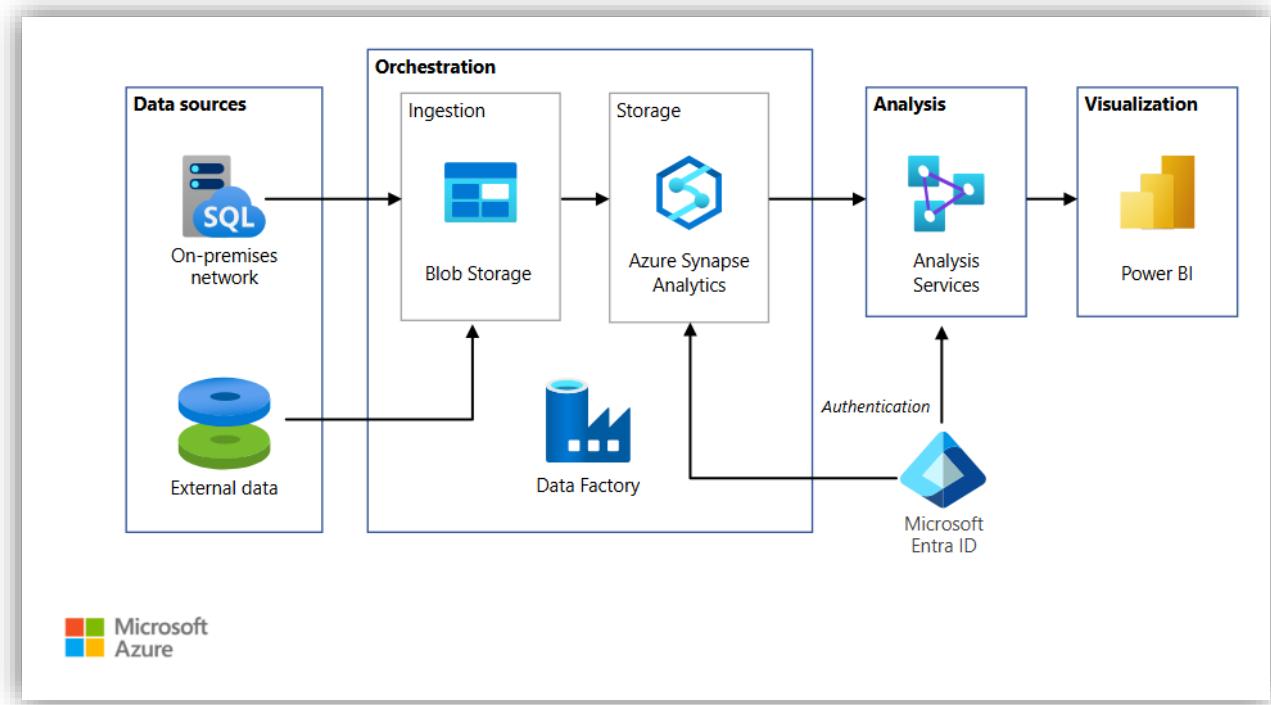
Potential Improvements:

- **Encoding:** The categorical variables ('month' and 'day') could be encoded to numerical values for use in certain types of analysis or machine learning models.
- **Outliers:** It would be beneficial to check for outliers in the numerical data which may skew analysis.
- **Feature Engineering:** New features could be created from the existing ones to provide more insights. For example, a 'season' feature could be engineered from the 'month' column.
- **Data Scaling:** The numerical data spans different ranges. Scaling the data to a standard range might be beneficial for certain types of analysis or models.
- **Adding Additional Data Sources:** Additional data such as more detailed geographical data, historical weather patterns, or human activity could enrich the dataset and provide more predictive power.

2. Architectural Diagram:

- Create an architectural diagram showing how services will be integrated.
- Specify each service's role (e.g., ADF for ingestion, ADLS for storage).





3. Data Pipeline Creation:

- Ingest data into the chosen source.
- Explain the choice of ingestion method.
- Store data in ADLS.
- Apply data transformations using PySpark.

The screenshot shows a Microsoft Excel spreadsheet titled "forestfires (1)". The data consists of 39 rows and 14 columns. The columns are labeled A through W, and the rows are numbered 1 through 39. The data includes various meteorological and geographical variables such as FFMC, DMC, ISI, temperature, relative humidity, wind, rain, and area. The first few rows of data are as follows:

	X	Y	month	day	FFMC	DMC	ISI	temp	RH	wind	rain	M	N	O	P	Q	R	S	T	U	V	W
2	7	5	mar	fri	86.2	26.2	94.3	5.1	8.2	51	6.7	0	0	0	0	0	0	0	0	0	0	
3	7	4	oct	tue	90.6	35.4	669.1	6.7	18	33	0.9	0	0	0	0	0	0	0	0	0	0	0
4	7	4	oct	sat	90.6	43.7	686.9	6.7	14.6	33	1.3	0	0	0	0	0	0	0	0	0	0	0
5	8	6	mar	fri	91.7	33.3	77.5	9	8.3	97	4	0.2	0	0	0	0	0	0	0	0	0	0
6	8	6	mar	sun	89.3	51.3	102.2	9.6	11.4	99	1.8	0	0	0	0	0	0	0	0	0	0	0
7	8	6	aug	sun	92.3	85.3	488	14.7	22.2	29	5.4	0	0	0	0	0	0	0	0	0	0	0
8	8	6	aug	mon	92.3	88.9	495.6	8.5	24.1	27	3.1	0	0	0	0	0	0	0	0	0	0	0
9	8	6	aug	mon	91.5	145.4	608.2	10.7	8	86	2.2	0	0	0	0	0	0	0	0	0	0	0
10	8	6	sep	tue	91	129.5	692.6	7	13.1	63	5.4	0	0	0	0	0	0	0	0	0	0	0
11	7	5	sep	sat	92.5	88	698.6	7.1	22.8	40	4	0	0	0	0	0	0	0	0	0	0	0
12	7	5	sep	sat	92.5	88	698.6	7.1	17.8	51	7.2	0	0	0	0	0	0	0	0	0	0	0
13	7	5	sep	sat	92.8	73.2	713	22.6	19.3	38	4	0	0	0	0	0	0	0	0	0	0	0
14	6	5	aug	fri	63.5	70.8	665.3	0.8	17	72	6.7	0	0	0	0	0	0	0	0	0	0	0
15	6	5	sep	mon	90.9	126.5	686.5	7	21.3	42	2.2	0	0	0	0	0	0	0	0	0	0	0
16	6	5	sep	wed	92.9	133.3	699.6	9.2	26.4	21	4.5	0	0	0	0	0	0	0	0	0	0	0
17	6	5	sep	fri	93.3	141.2	713.9	13.9	22.9	44	5.4	0	0	0	0	0	0	0	0	0	0	0
18	5	5	mar	sat	91.7	35.8	80.8	7.8	15.1	27	5.4	0	0	0	0	0	0	0	0	0	0	0
19	8	5	oct	mon	84.9	32.8	664.2	3	16.7	47	4.9	0	0	0	0	0	0	0	0	0	0	0
20	6	4	mar	wed	89.2	27.9	70.8	6.3	15.9	35	4	0	0	0	0	0	0	0	0	0	0	0
21	6	4	apr	sat	86.3	27.4	97.1	5.1	9.3	44	4.5	0	0	0	0	0	0	0	0	0	0	0
22	6	4	sep	tue	91	129.5	692.6	7	18.3	40	2.7	0	0	0	0	0	0	0	0	0	0	0
23	5	4	sep	mon	91.8	76.5	724.3	9.2	19.1	38	2.7	0	0	0	0	0	0	0	0	0	0	0
24	7	4	jun	sun	94.3	96.3	200	56.1	21	44	4.5	0	0	0	0	0	0	0	0	0	0	0
25	7	4	aug	sat	90.2	110.9	537.4	6.2	19.5	43	5.8	0	0	0	0	0	0	0	0	0	0	0
26	7	4	aug	sat	93.5	139.4	594.2	20.3	23.7	32	5.8	0	0	0	0	0	0	0	0	0	0	0
27	7	4	aug	sun	91.4	142.4	601.4	10.6	16.3	60	5.4	0	0	0	0	0	0	0	0	0	0	0
28	7	4	sep	fri	92.4	117.9	668	12.2	19	34	5.8	0	0	0	0	0	0	0	0	0	0	0
29	7	4	sep	mon	90.9	126.5	686.5	7	19.4	48	1.3	0	0	0	0	0	0	0	0	0	0	0
30	6	3	sep	sat	93.4	145.4	721.4	8.1	30.2	24	2.7	0	0	0	0	0	0	0	0	0	0	0
31	6	3	sep	sun	93.5	149.3	728.6	8.1	22.8	39	3.6	0	0	0	0	0	0	0	0	0	0	0
32	6	3	sep	fri	94.3	85.1	692.3	15.9	25.4	24	3.6	0	0	0	0	0	0	0	0	0	0	0
33	6	3	sep	mon	88.6	91.8	709.9	7.1	11.2	78	7.6	0	0	0	0	0	0	0	0	0	0	0
34	6	3	sep	fri	88.6	69.7	706.8	5.8	20.6	37	1.8	0	0	0	0	0	0	0	0	0	0	0
35	6	3	sep	sun	91.7	75.6	718.3	7.8	17.7	39	3.6	0	0	0	0	0	0	0	0	0	0	0
36	6	3	sep	mon	91.8	78.5	724.3	9.2	21.2	32	2.7	0	0	0	0	0	0	0	0	0	0	0
37	6	3	sep	tue	90.3	80.7	730.2	6.3	18.2	62	4.5	0	0	0	0	0	0	0	0	0	0	0
38	6	3	oct	tue	90.6	35.4	669.1	6.7	21.7	24	4.5	0	0	0	0	0	0	0	0	0	0	0
39	7	4	oct	fri	90	41.5	682.6	8.7	11.3	60	5.4	0	0	0	0	0	0	0	0	0	0	0

The screenshot shows the Microsoft Azure Storage Container named "sourcecontainer". The container has one blob named "forestfires.csv". The blob details are as follows:

- Name:** forestfires.csv
- Modified:** 14/06/2024, 12:14:27
- Access tier:** Hot (Inferred)
- Archive status:** Not yet archived
- Blob type:** Block blob

Microsoft Azure Search resources, services, and docs (G+)

keerthanaakkula@gmail.. DEFAULT DIRECTORY

Home > Default Directory | App registrations >

RakshithaServicePrinciple

Search Delete Endpoints Preview features

Overview Quickstart Integration assistant

Manage

- Branding & properties
- Authentication
- Certificates & secrets
- Token configuration
- API permissions
- Expose an API
- App roles
- Owners
- Roles and administrators
- Manifest

Support + Troubleshooting

- Troubleshooting
- New support request

Copied

^ Essentials

Display name: RakshithaServicePrinciple

Application (client) ID: 22b45241-70c5-4ea2-a23a-fb7efc40d683

Object ID: c4ab1157-42f1-4502-8c83-57728ee7b258

Directory (tenant) ID: cc876379-1061-49ab-964b-7e5d840d62a1

Supported account types: My organization only

Client credentials: 0_certificate_1_secret

Redirect URLs: Add a Redirect URI

Application ID URI: Add an Application ID URI

Managed application in local directory: RakshithaServicePrinciple

Welcome to the new and improved App registrations. Looking to learn how it's changed from App registrations (Legacy)? [Learn more](#)

Starting June 30th, 2020 we will no longer add any new features to Azure Active Directory Authentication Library (ADAL) and Azure Active Directory Graph. We will continue to provide technical support and security updates but we will no longer provide feature updates. Applications will need to be upgraded to Microsoft Authentication Library (MSAL) and Microsoft Graph. [Learn more](#)

Get Started Documentation

Build your application with the Microsoft identity platform

The Microsoft identity platform is an authentication service, open-source libraries, and application management tools. You can create modern, standards-based authentication solutions, access and

Microsoft Azure Search resources, services, and docs (G+)

keerthanaakkula@gmail.. DEFAULT DIRECTORY

Home > Default Directory | App registrations > RakshithaServicePrinciple

RakshithaServicePrinciple | Certificates & secrets

Search Got feedback?

Overview Quickstart Integration assistant

Manage

- Branding & properties
- Authentication
- Certificates & secrets
- Token configuration
- API permissions
- Expose an API
- App roles
- Owners
- Roles and administrators
- Manifest

Got a second to give us some feedback? →

Credentials enable confidential applications to identify themselves to the authentication service when receiving tokens at a web addressable location (using an HTTPS scheme). For a higher level of assurance, we recommend using a certificate (instead of a client secret) as a credential.

Application registration certificates, secrets and federated credentials can be found in the tabs below.

Certificates (0) Client secrets (1) Federated credentials (0)

A secret string that the application uses to prove its identity when requesting a token. Also can be referred to as application password.

+ New client secret

Description	Expires	Value ⓘ	Secret ID
This is use to read the d...	12/09/2024	P5o8Q~EadmMkLGn9...	921d91e2-3574-4db6-...

Support + Troubleshooting

- Troubleshooting
- New support request

Service Principle: Client ID (username) :22b45241-70c5-4ea2-a23a-fb7efc40d683

Secret value (password): P5o8Q~EadmMkLGn99uGv7St9sMw9luUUdE_kEcgg

The screenshot shows the Microsoft Azure Access Control (IAM) interface for the storage account 'adlsfinalstorage'. The 'Role assignments' tab is selected. It displays 2 role assignments out of a maximum of 4000. A search bar at the top allows filtering by name or email. The results table includes columns for Name, Type, Role, Scope, and Condition. One assignment is listed for the user 'HEM MARYADA'.

Name	Type	Role	Scope	Condition
HEM MARYADA	User	Owner	Subscription (Inherited)	None

The screenshot shows the Microsoft Azure Access Control (IAM) interface for the storage account 'adlsfinalstorage'. The 'Current role assignments' tab is selected. A search bar at the top allows filtering by name. The results table shows 0 role assignments for the service principal 'RakshithaServicePrinciple'. Other tabs include 'Eligible assignments', 'Role assignments (0)', 'Deny assignments (0)', and 'Classic administrators (0)'.

Role assignments	Eligible assignments
(0)	(0)

Microsoft Azure Search resources, services, and docs (G+)

Home > adlsfinalstorage | Access Control (IAM) > Add role assignment

Role Members Conditions Review + assign

Selected role Storage Blob Data Contributor

Assign access to User, group, or service principal Managed identity

Members + Select members

Name	Object ID	Type
No members selected		

Description Optional

Review + assign Previous Next Select Close

Select members

Select ⓘ RakshithaServicePrinciple

No users, groups, or service principals found.

Selected members:

RakshithaServicePrinciple Remove

mPrivate Microsoft Add role assignment - Microsoft

https://portal.azure.com/#view/Microsoft_Azure_AD/AddRoleAssignmentsLandingBlade/scope/%2Fsubscriptions%2Fbf9ac98f-4827-4552-bd2c-930ad6cd590b%2FresourceGroups%2Frakresourcegroup/providers/Microsoft.Storage/storageAccounts/adlsfinalstorage

Microsoft Azure Search resources, services, and docs (G+)

keerthanaakkula@gmail... DEFAULT DIRECTORY

Home > adlsfinalstorage | Access Control (IAM) > Add role assignment

Role Members Conditions Review + assign

Role Storage Blob Data Contributor

Scope /subscriptions/bf9ac98f-4827-4552-bd2c-930ad6cd590b/resourceGroups/rakresourcegroup/providers/Microsoft.Storage/storageAccounts/adlsfinalstorage

Members

Name	Object ID	Type
RakshithaServicePrinciple	46b612a3-55ec-40c4-8145-3577dac41d28	App

Description No description

Condition None

Review + assign Previous Next Feedback

Microsoft Azure Search resources, services, and docs (G+/-) keerthanaakkula@gmail... DEFAULT DIRECTORY

Home > adlsfinalstorage

adlsfinalstorage | Access Control (IAM)

Storage account

Search Add Download role assignments Edit columns Refresh Remove Feedback

Overview Activity log Tags Diagnose and solve problems Access Control (IAM) Data migration Events Storage browser

Number of role assignments for this subscription 3 4000

All Job function (2) Privileged (1)

Search by name or email Type : All Role : All Scope : All scopes Group by : Role

3 items (1 Users, 2 Service Principals)

Name	Type	Role	Scope	Condition
HEM MARYADA keerthanaakkula@gmail.com#EXT#@keerthan...	User	Owner	Subscription (Inherited)	None
RakshithaServicePrinciple	App	Storage Blob Data Contributor	This resource	Add
RakshithaServicePrinciple	App	Storage Blob Data Contributor	This resource	Add

InPrivate (2) Create a key vault - Microsoft Az... + https://portal.azure.com/#create/Microsoft.KeyVault

Microsoft Azure Search resources, services, and docs (G+/-) keerthanaakkula@gmail... DEFAULT DIRECTORY

Home > Create a resource > Marketplace >

Create a key vault

Azure Key Vault is a cloud service used to manage keys, secrets, and certificates. Key Vault eliminates the need for developers to store security information in their code. It allows you to centralize the storage of your application secrets which greatly reduces the chances that secrets may be leaked. Key Vault also allows you to securely store secrets and keys backed by Hardware Security Modules or HSMs. The HSMs used are Federal Information Processing Standards (FIPS) 140-2 Level 2 validated. In addition, key vault provides logs of all access and usage attempts of your secrets so you have a complete audit trail for compliance.

Project details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription * Azure subscription 1

Resource group * rakresourcegroup Create new

Instance details

Key vault name * rkvaulttest

Region * East US

Pricing tier * Standard

Recovery options

Soft delete protection will automatically be enabled on this key vault. This feature allows you to recover or permanently delete

Previous Next Review + create Give feedback

The screenshot shows the Microsoft Azure Deployment Overview page for a deployment named 'rakvaulttest'. The deployment status is marked as 'complete' with a green checkmark icon. Key details listed include:

- Deployment name: rakvaulttest
- Subscription: Azure subscription 1
- Start time: 14/06/2024, 13:19:02
- Correlation ID: bae8be31-a247-435f-b...
- Resource group: rakresourcegroup

Below the main summary, there are sections for 'Deployment details' and 'Next steps'. A prominent blue button labeled 'Go to resource' is located at the bottom left. To the right of the main content, there are two side cards:

- Cost management**: A card with a dollar sign icon, encouraging users to stay within their budget and prevent unexpected charges. It includes a link to 'Set up cost alerts >'.
- Microsoft Defender for Cloud**: A card with a shield icon, securing apps and infrastructure. It includes a link to 'Go to Microsoft Defender for Cloud >'.

The screenshot shows the Microsoft Azure Key Vault Secrets page for the 'rakvaulttest' key vault. The left sidebar lists various management options like Overview, Activity log, Access control (IAM), Tags, and more. The 'Secrets' option is currently selected and highlighted in grey. The main content area displays a table with columns for Name, Type, Status, and Expiration date. A message at the top states: 'The operation is not allowed by RBAC. If role assignments were recently changed, please wait several minutes for role assignments to become effective.' Below this, a note says: 'You are unauthorized to view these contents.'

In the above screenshot as you can see that the you have to assign the admin role to this key vault account.

Microsoft Azure Search resources, services, and docs (G+/)

Home > rakvaulttest

rakvaulttest | Access control (IAM)

Key vault

Search Add Download role assignments Edit columns Refresh Remove Feedback

Overview Activity log Access control (IAM) Tags Diagnose and solve problems Access policies Events

Objects Keys Secrets Certificates

Settings Access configuration Networking Microsoft Defender for Cloud

Role assignments Roles Deny assignments Classic administrators

Number of role assignments for this subscription 4 4000

All Job function (1) Privileged (1)

Search by name or email Type : All Role : All Scope : All scopes Group by : Role

2 items (2 Users)

Name	Type	Role	Scope	Condition
HEM MARYADA keerthanaakkula_g...	User	Owner	Subscription (Inherited)	None
HEM MARYADA keerthanaakkula_g...	User	Key Vault Administrator	This resource	None

Open data bricks and open this link

<https://<databricksinstance>#secrets/createScope>

Microsoft Azure Upgrade Search resources, services, and docs (G+/)

Home > rakvaulttest

rakvaulttest | Properties

Key vault

Save Discard changes Refresh

Access policies Events

Objects Keys Secrets Certificates

Settings Access configuration Networking Microsoft Defender for Cloud Properties Locks

Monitoring Alerts Metrics Diagnostic settings Logs Insights

Name rakvaulttest

Sku (Pricing tier) Standard

Location eastus

Vault URI https://rakvaulttest.vault.azure.net/ Copied

Resource ID /subscriptions/bf9ac98f-4827-4552-bd2c-930ad6cd590b/resourceGroups/rakresourcegroup/providers/...

Subscription ID bf9ac98f-4827-4552-bd2c-930ad6cd590b

Subscription Name Azure subscription 1

Directory ID cc876379-1061-49ab-964b-7e5d840d62a1

Directory Name Default Directory

Soft-delete Soft delete has been enabled on this key vault

Days to retain deleted vaults 90

Purge protection

Disable purge protection (allow key vault and objects to be purged during retention period)
 Enable purge protection (enforce a mandatory retention period for deleted vaults and vault objects)

The screenshot shows the 'Create Secret Scope' dialog box on the Microsoft Azure Databricks platform. The dialog has a title bar 'Create Secret Scope' with 'Cancel' and 'Create' buttons. Below the title, a descriptive text states: 'A store for secrets that is identified by a name and backed by a specific store type. [Learn more](#)'. The 'Scope Name' field contains 'rakshithaSecretScope'. The 'Manage Principal' dropdown is set to 'All workspace users'. Under the 'Azure Key Vault' section, the 'DNS Name' field is 'https://rakvaulttest.vault.azure.net/' and the 'Resource ID' field starts with '/subscriptions/bf9ac98f-4827-4552-bd2c-930ad6cd590b/resourceGroups/rakresour'. On the left, there is a sidebar with various icons.

The screenshot shows the same 'Create Secret Scope' dialog box, but now it displays a modal confirmation message: 'The secret scope named `rakshithaSecretScope` has been added.' Below the message, it says 'Manage secrets in this scope in Azure KeyVault with manage principal = users'. At the bottom right of the modal is an 'OK' button. The background of the dialog is dimmed.

Under secret value put the Service principal name

You have make sure that you the admin/permissions of that key vault to create the secret or any certification, if you not admin/have relevant permissions you can't cannt create an secret.

The screenshot shows the 'Create a secret' dialog box in the Microsoft Azure portal. The 'Name' field is set to 'secretforServicePrinciple'. The 'Secret value' field contains a masked password. The 'Enabled' button is set to 'Yes'. At the bottom are 'Create' and 'Cancel' buttons.

The screenshot shows the 'rakvaulttest | Secrets' page in the Microsoft Azure portal. A single secret named 'secretforServicePrinciple' is listed, showing it is enabled. The table has columns for Name, Type, Status, and Expiration date.

Name	Type	Status	Expiration date
secretforServicePrinciple		✓ Enabled	

dbutils.secrets.get(scope = "", key = ""). I got an error while I was running the databricks notebook. You cannt access the key vault using a role-based access you need to enable the Vault access policy from access configuration

Home > rakvaulttest

rakvaulttest | Access configuration

Key vault

Access policies Refresh

Access control (IAM) Access policies

Settings Access configuration

Permission model

Grant data plane access by using a [Azure RBAC](#) or [Key Vault access policy](#)

Azure role-based access control (recommended) ⓘ

Vault access policy ⓘ

[Go to access policies](#)

Resource access

WARNING: You are changing the permission model. This may immediately change users and services that are allowed to access this key vault. You may proceed if this key vault is new, not used in production workloads, or if you are undoing a previous change. Otherwise it's strongly recommended that you perform this action in the beginning of your own planned maintenance event, during which you can test the new configuration and undo if necessary.

Choose among the following options to grant access to specific resource types

Azure Virtual Machines for deployment ⓘ

Azure Resource Manager for template deployment ⓘ

Azure Disk Encryption for volume encryption ⓘ

[Apply](#) [Discard changes](#)

Microsoft Azure |  databricks Search data, notebooks, recents, and ... CTRL + P mywkrakdatabricks ⓘ H

LogicTransforma... Python ⌂

File Edit View Run Help L Provide f Run all HEM MARYADA's Cluster Schedule Share

New Assistant: OFF ⓘ

Assistant

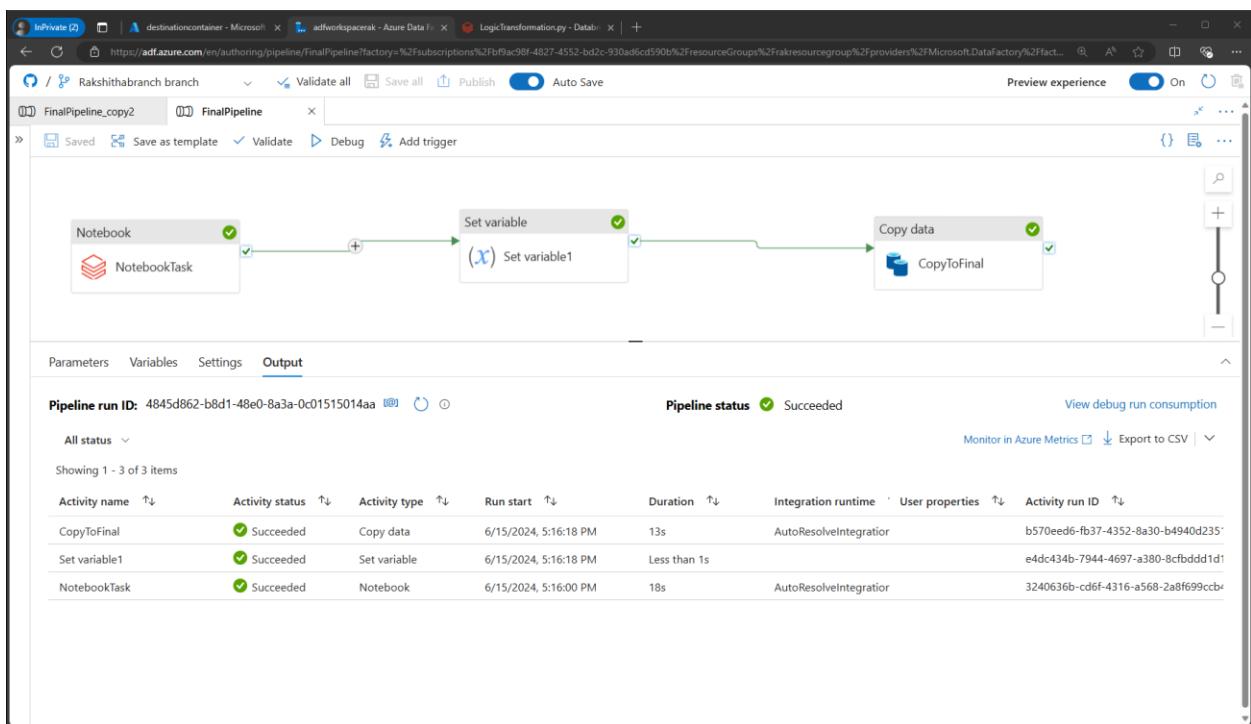
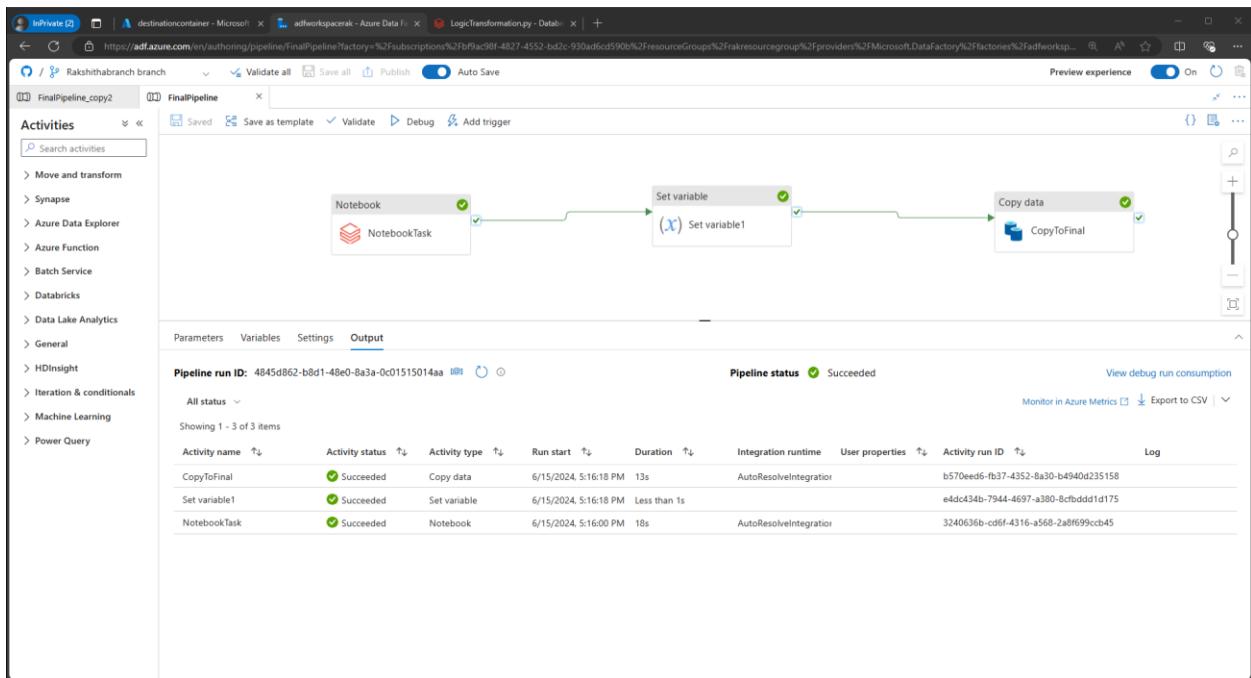
The error you're encountering indicates a permissions issue with accessing the Azure Key Vault from Databricks. Specifically, Databricks is not authorized to perform the action on the resource, which in this case is retrieving a secret from your Azure Key Vault. To resolve this issue, you need to ensure that the Databricks service principal has the necessary permissions on the Azure Key Vault.

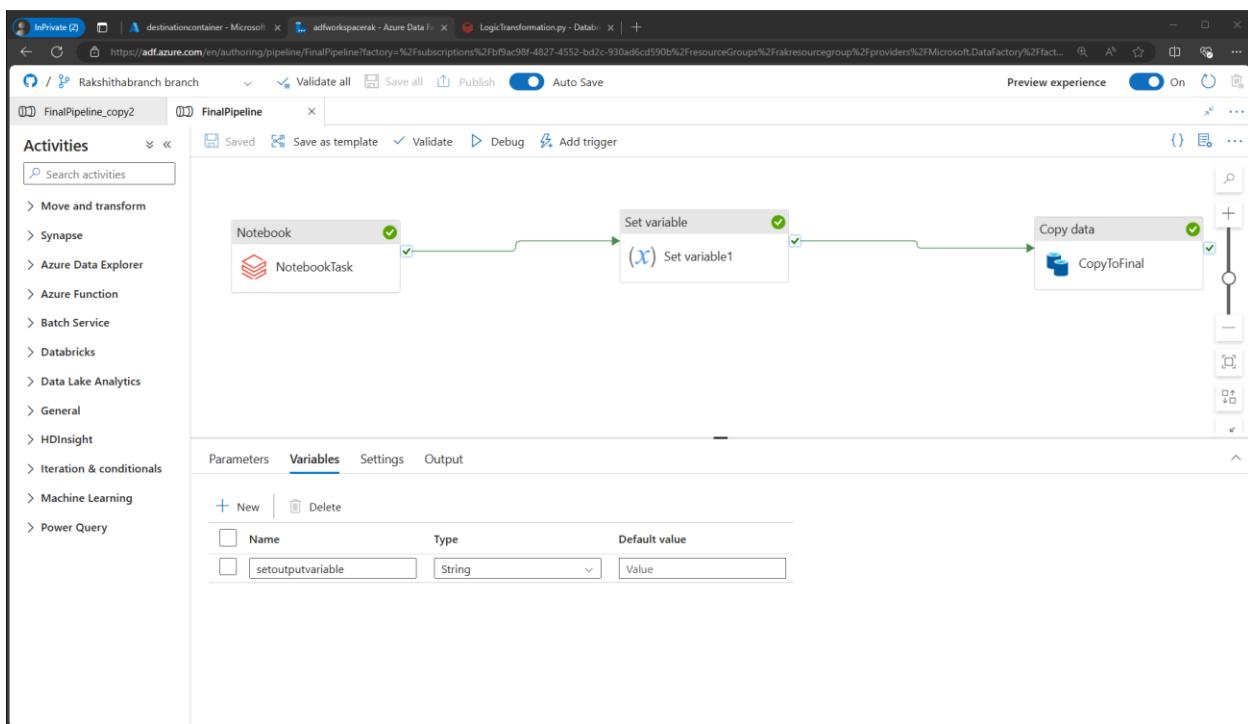
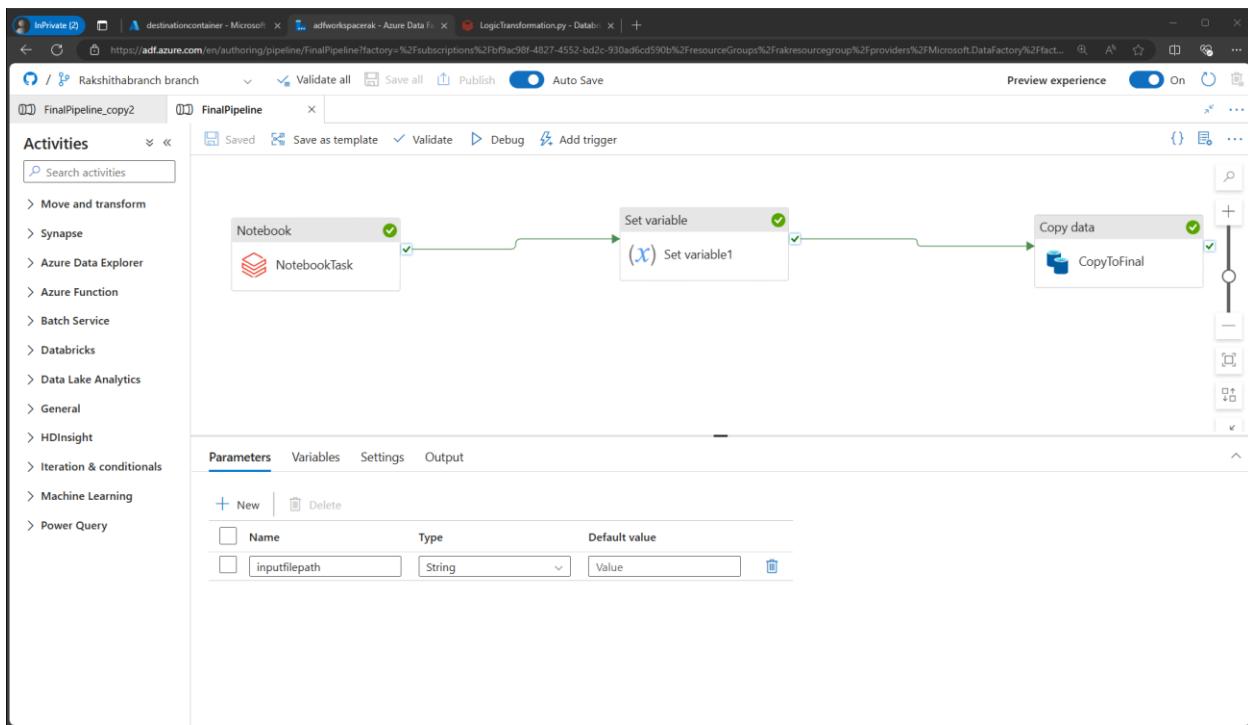
Here are the steps to fix the issue:

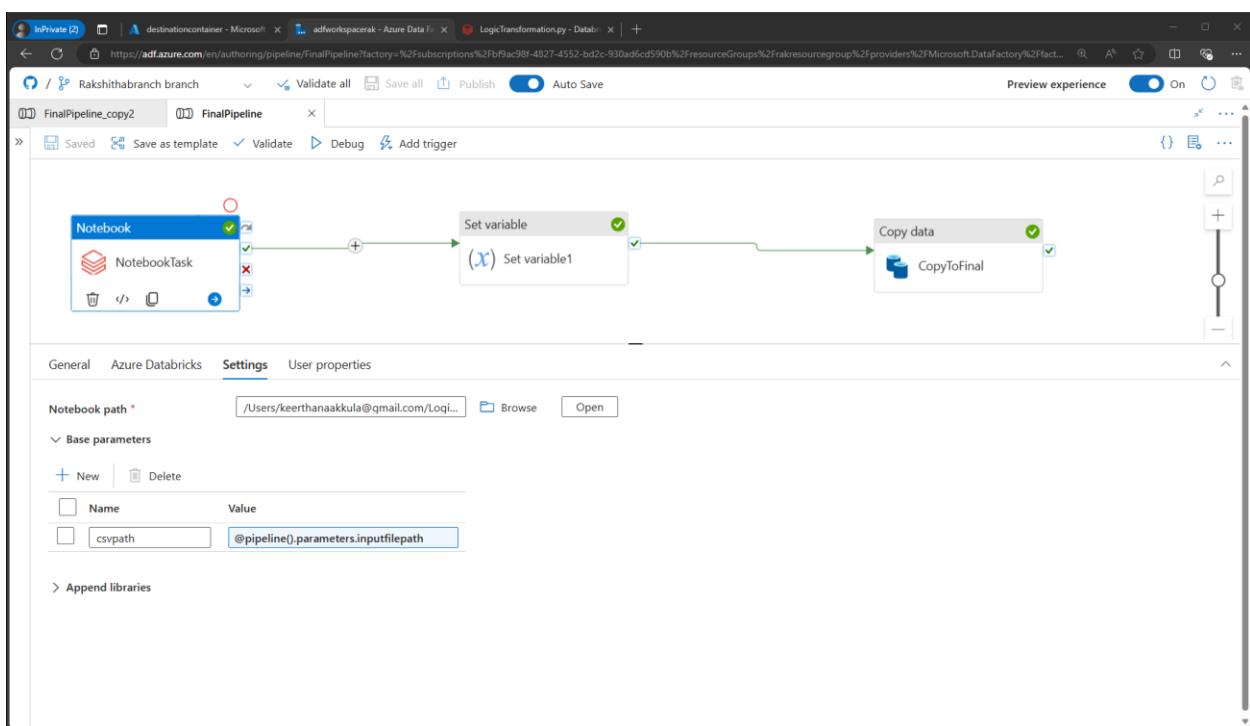
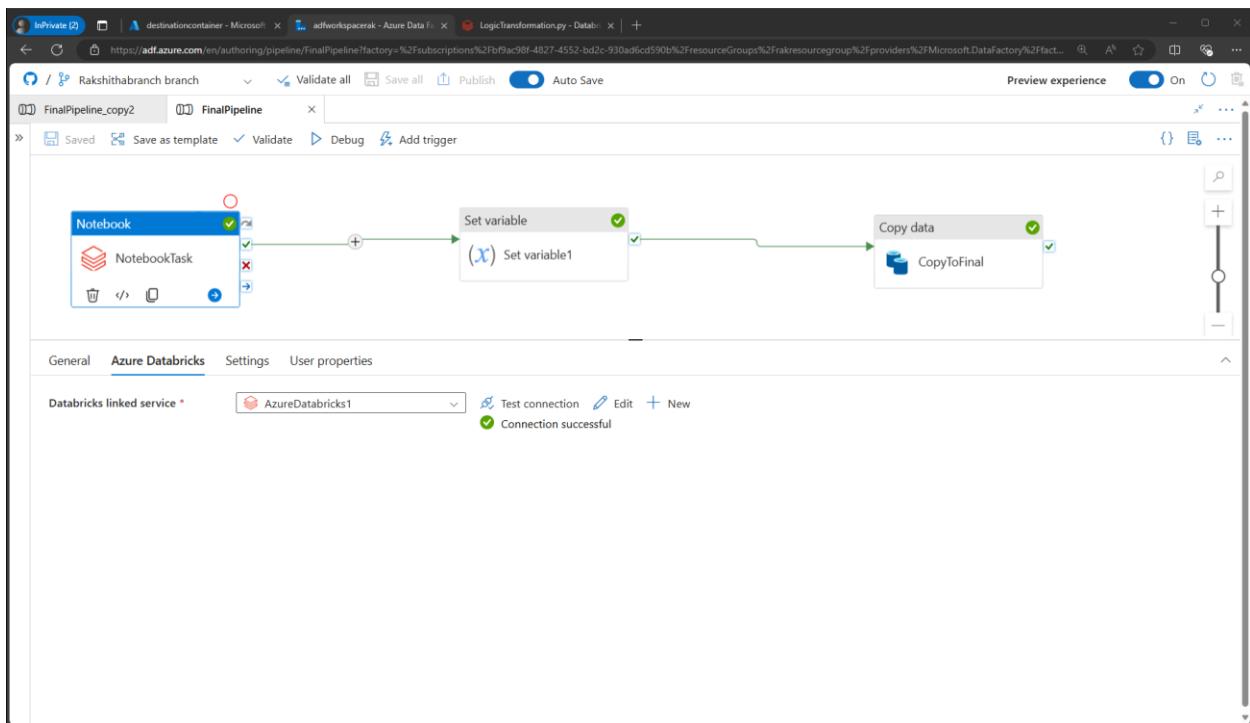
- Set Permission Model to Vault Access Policy:** Ensure your Azure Key Vault is using the Vault access policy permission model, not the Azure role-based access control (RBAC) model. Databricks currently supports the Vault access policy model for secret scopes.
- Add Access Policy for Databricks Service Principal:**
 - Go to your Azure Key Vault in the Azure Portal.
 - Navigate to the "Access policies" section.
 - Click on "Add Access Policy".
 - Select the "Secret Management" template or manually select the "Get" and "List" permissions under "Secret permissions".
 - Click on "Select principal" and add the Databricks service principal. You can find the Application (client) ID for the Databricks service principal in the error message you provided (`appid=2ff814a6-3304-4ab8-85cb-cd0e6f879c1d`).
 - Click "Add" to add the policy, then click "Save" to apply the changes.
- Allow Trusted Microsoft Services:**
 - Still in the Azure Key Vault settings, navigate to the "Networking" section.
 - Ensure that "Allow trusted Microsoft services to bypass this firewall" is checked. This setting is crucial for services like Databricks to access your Key Vault.

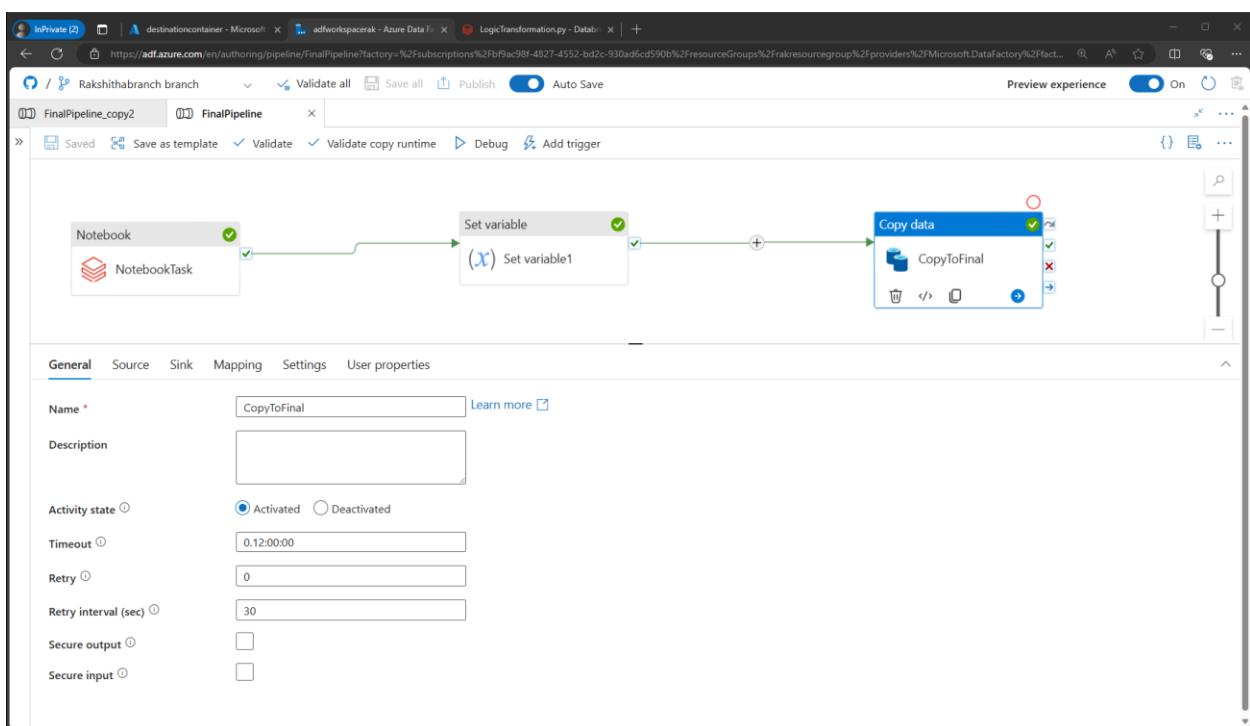
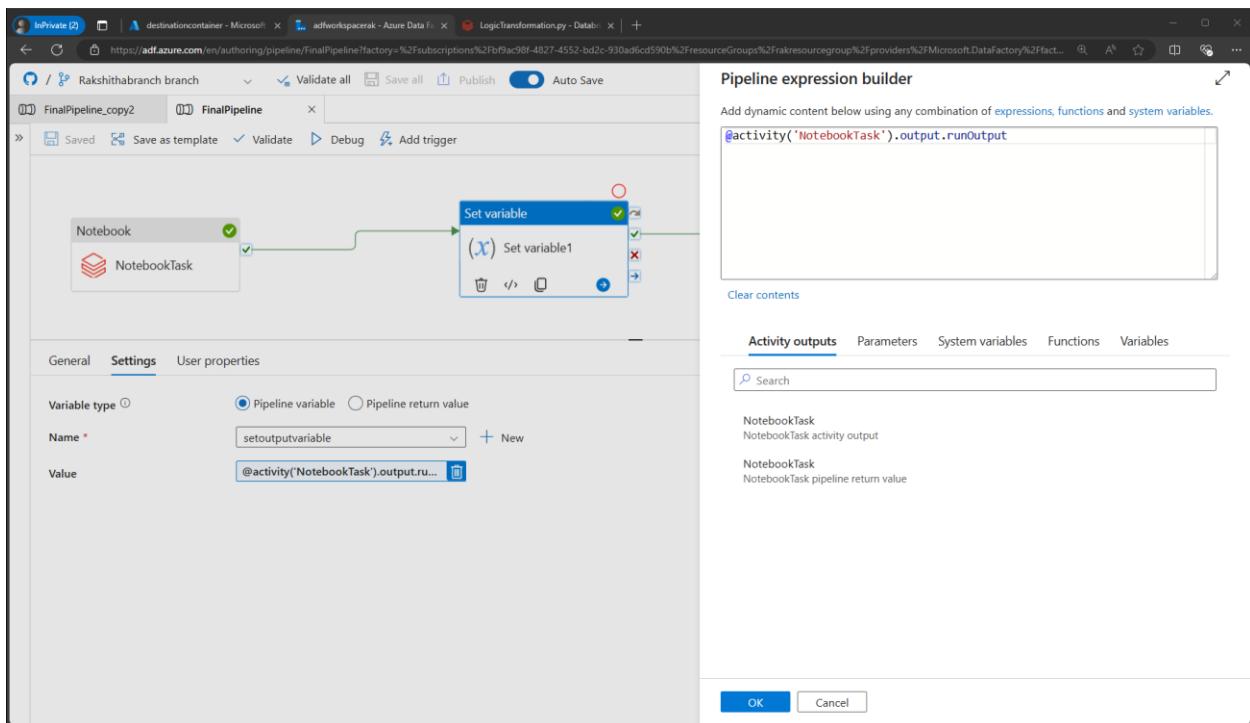
After applying these changes, wait a few minutes for the permissions to propagate, then try accessing the secret from

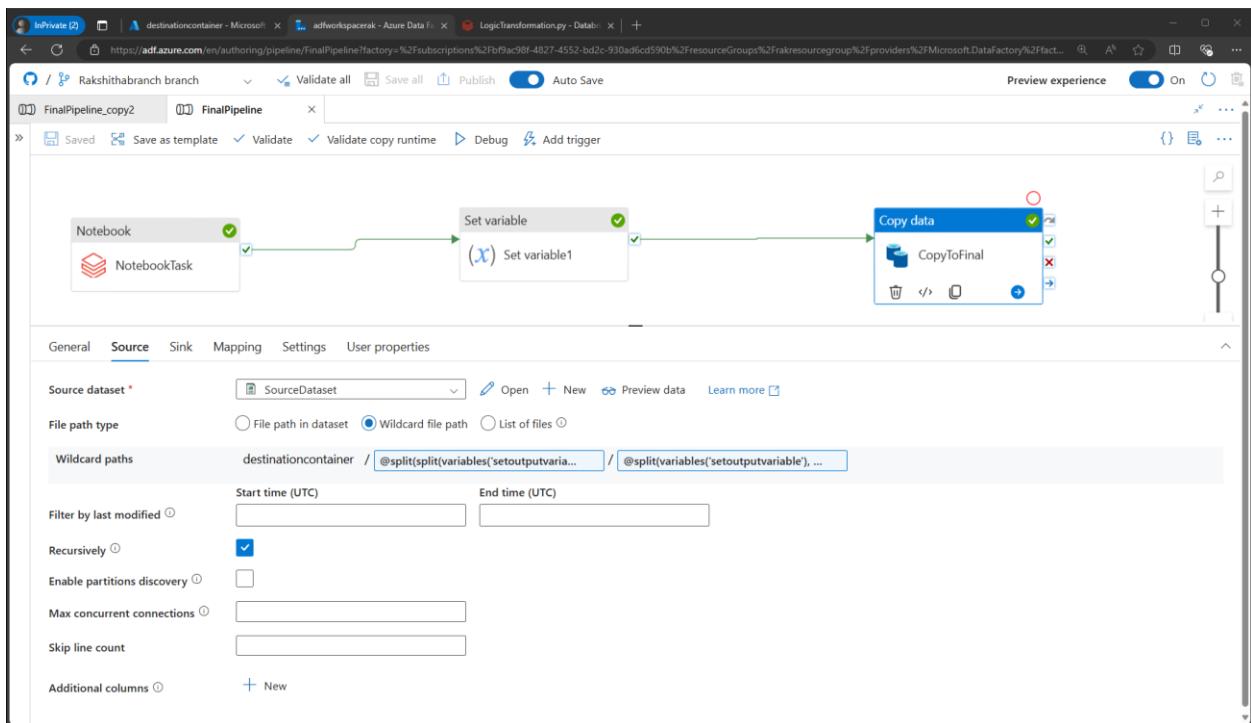
[Fix my code](#)











The screenshot shows the Azure Data Factory dataset editor for "SourceDataset". The "Connection" tab is selected, showing the following configuration:

- Linked service:** AzureDataLakeStorage1
- File path:** destinationcontainer / Directory / File name
- Compression type:** Select...
- Column delimiter:** Comma (,)
- Row delimiter:** Default (\r\n, or \n)
- Encoding:** Default(UTF-8)
- Quote character:** Double quote ("")
- Escape character:** Backslash (\)
- First row as header:** Checked.
- Null value:** Unchecked.

The screenshot shows the Azure Data Factory Pipeline expression builder interface. A pipeline step named "NotebookTask" is selected, and its "Source" tab is active. In the "Source dataset" dropdown, "SourceDataset" is chosen. Under "File path type", the "Wildcard file path" option is selected, and the "Wildcard paths" field contains the expression `destinationcontainer / @split(split(variables('setoutputvariable'), '/')[0][1], '/')[1]`. This expression is used to dynamically generate a file path based on the value of the variable `setoutputvariable`, which is split into segments and the second segment is taken as the directory name. The "Activity outputs" pane on the right lists the output of the NotebookTask activity.

This screenshot shows a similar configuration in the Azure Data Factory Pipeline expression builder. The "Source dataset" is "SourceDataset", and the "File path type" is "Wildcard file path". The "Wildcard paths" field contains the expression `destinationcontainer / @split(split(variables('setoutputvariable'), '/')[sub(length(split(variables('setoutputvariable'), '/')), 1)], '/')`. This expression uses the `sub` function to get the last segment of the variable `setoutputvariable` and combines it with the rest of the path. The "Activity outputs" pane shows the output of the NotebookTask activity.

InPrivate (2) Microsoft.Azure.SynapseAnalytics - Microsoft Azure adfworkspacera - Azure Data ... LogicTransformation.py - Data ... Copilot

https://portal.azure.com/#view/HubsExtension/DeploymentDetailsBlade/~/overview/id/%2Fsubscriptions%2Fb9ac98f-4827-4552-bd2c-9...

Microsoft Azure Search resources, services, and docs (G+/-)

Home > Microsoft.Azure.SynapseAnalytics-20240616174929 | Overview

Deployment

Delete Cancel Redeploy Download Refresh

✓ Your deployment is complete

Deployment name : Microsoft.Azure.SynapseAnalytics-20240616174929
Subscription : Azure subscription 1
Resource group : rakresourcegroup
Start time : 16/06/2024, 17:50:17
Correlation ID : 5b067cd6-9227-4464-b0f0-c84a324b1b23

> Deployment details
Next steps

Go to resource group

Give feedback
Tell us about your experience with deployment

InPrivate (2) Microsoft.Azure.SynapseAnalytics - Microsoft Azure synapsecworkspace - Microsoft Azure adfworkspacera - Azure Data ... LogicTransformation.py - Data ... Copilot

https://portal.azure.com/#view/HubsExtension/DeploymentDetailsBlade/~/overview/id/%2Fsubscriptions%2Fb9ac98f-4827-4552-bd2c-9...

Microsoft Azure Search resources, services, and docs (G+/-)

Home > Microsoft.Azure.Synapse.SqlPoolOnExistingWorkspace | Overview

Deployment

✓ Deployment succeeded

Deployment 'Microsoft.Azure.Synapse.SqlPoolOnExistingWorkspace...' to resource group 'rakresourcegroup' was successful.

Go to resource Pin to dashboard

✓ Your deployment is complete

Deployment name : Microsoft.Azure.Synapse.SqlPoolOnExistingWorkspace_a51fef71083d4
Subscription : Azure subscription 1
Resource group : rakresourcegroup
Start time : 16/06/2024, 17:56:11
Correlation ID : a43939e5-7c89-4a97-bdc9-a0fee40e027f

> Deployment details
Next steps

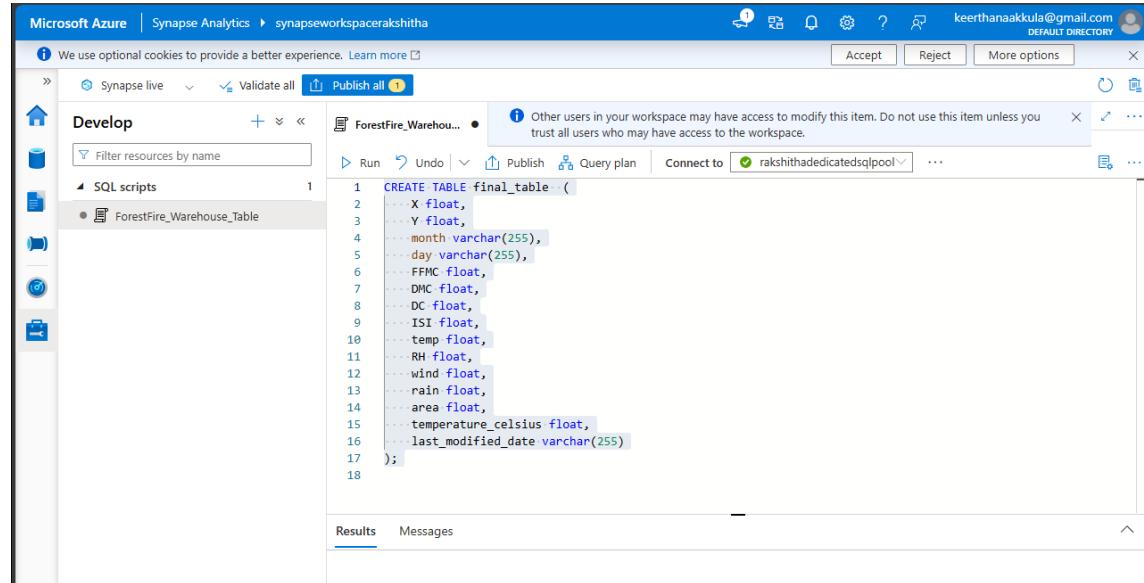
Go to resource

Give feedback
Tell us about your experience with deployment

```

CREATE TABLE your_table_name (
    X float,
    Y float,
    month varchar(255),
    day varchar(255),
    FFMC float,
    DMC float,
    DC float,
    ISI float,
    temp float,
    RH float,
    wind float,
    rain float,
    area float,
    temperature_celsius float,
    last_modified_date varchar (255)
);

```

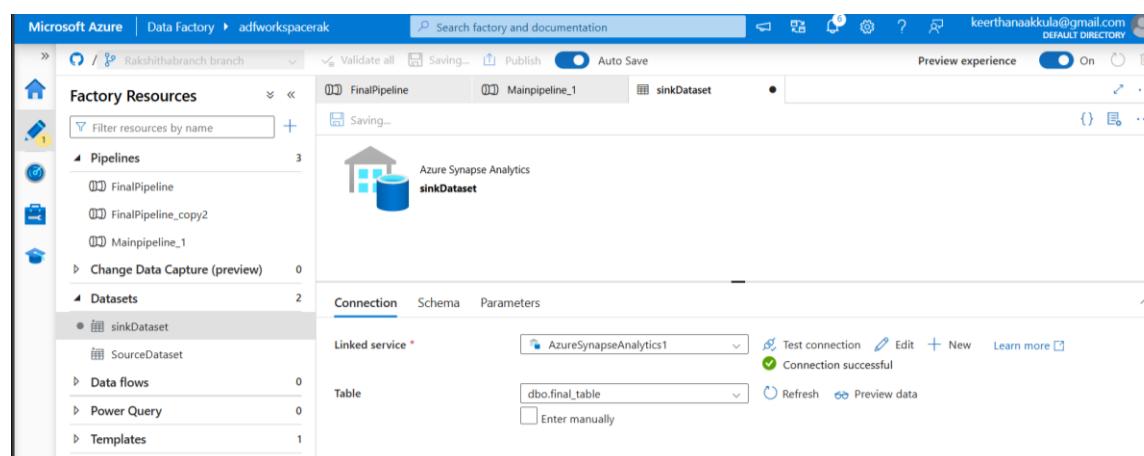


The screenshot shows the Microsoft Azure Synapse Analytics workspace interface. The top navigation bar includes 'Microsoft Azure', 'Synapse Analytics', 'synapseworkspace:rakshitha', and a user email 'keerthanaakkula@gmail.com'. The main area is titled 'Develop' and contains a 'SQL scripts' section. A specific script named 'ForestFire_Warehouse_Table' is selected. The code editor displays the following SQL statement:

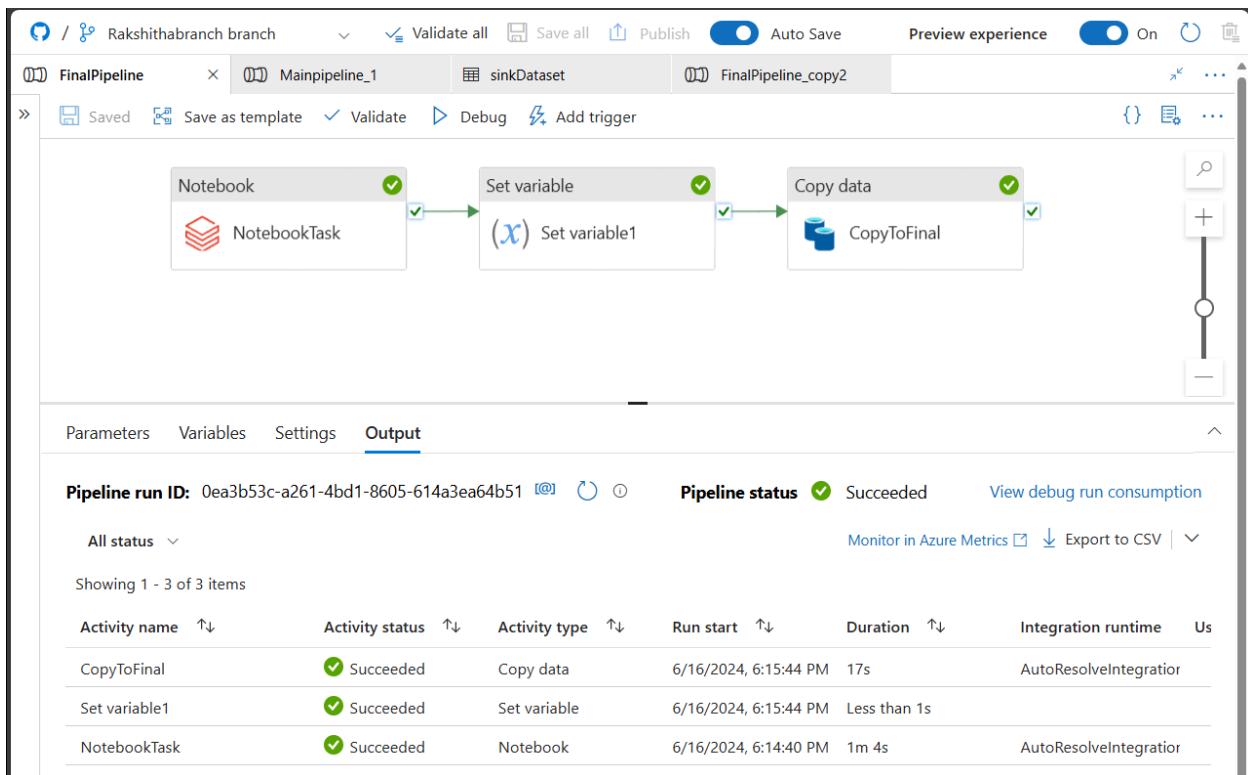
```

CREATE TABLE final_table (
    X float,
    Y float,
    month varchar(255),
    day varchar(255),
    FFMC float,
    DMC float,
    DC float,
    ISI float,
    temp float,
    RH float,
    wind float,
    rain float,
    area float,
    temperature_celsius float,
    last_modified_date varchar (255)
);

```



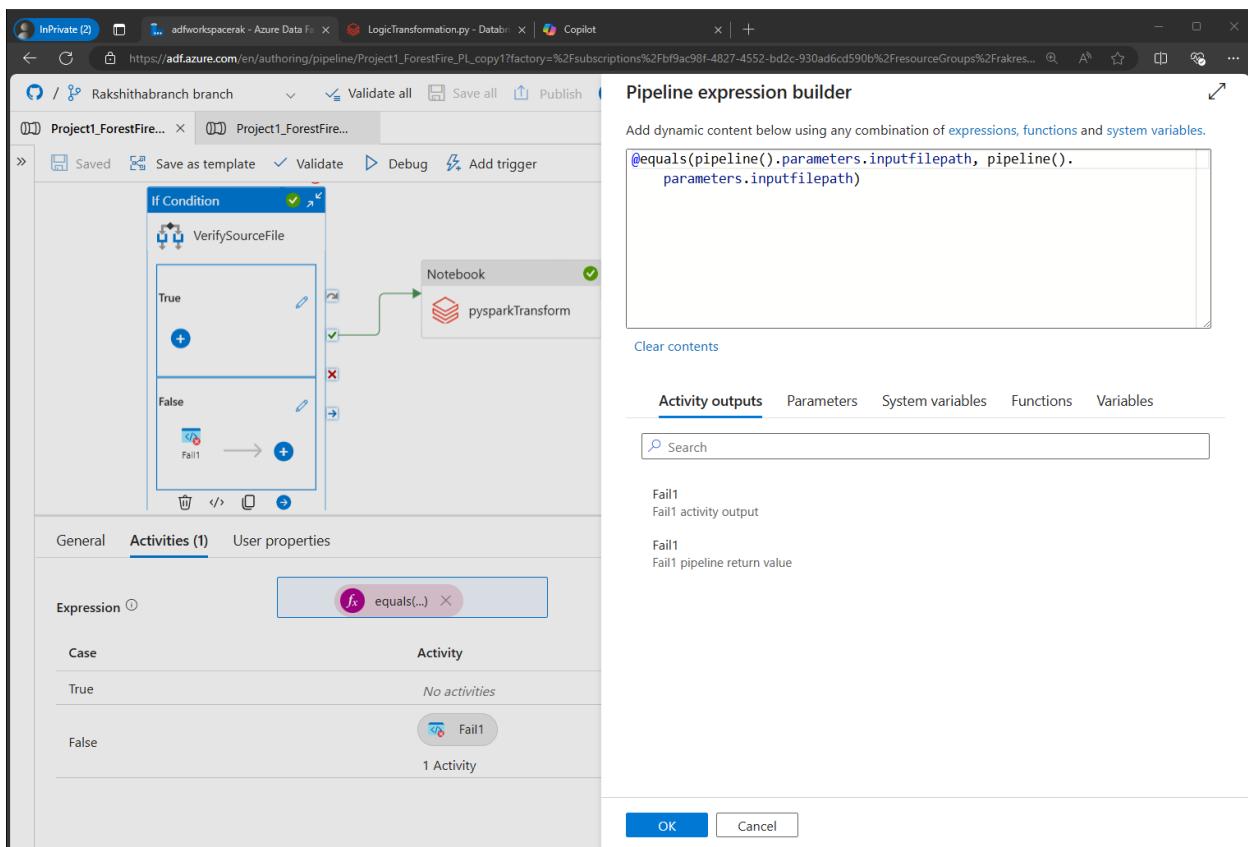
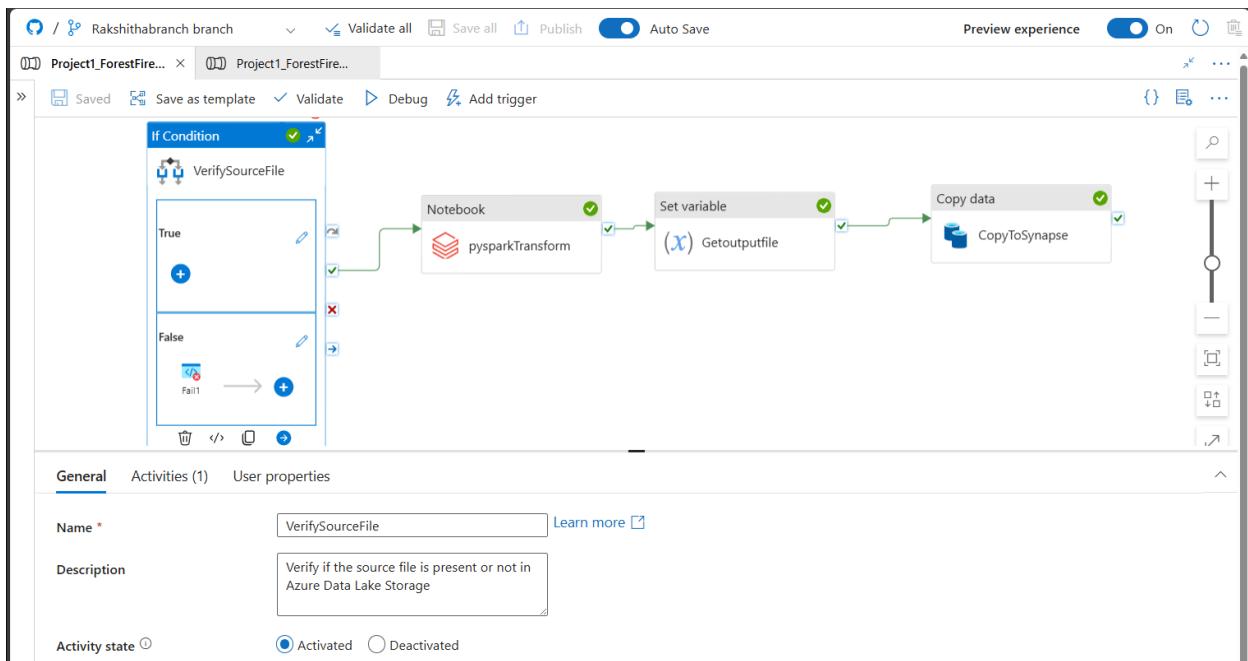
The screenshot shows the Microsoft Azure Data Factory workspace interface. The top navigation bar includes 'Microsoft Azure', 'Data Factory', 'adfworkspacera', and a user email 'keerthanaakkula@gmail.com'. The left sidebar lists 'Factory Resources' such as Pipelines, Datasets, Data flows, Power Query, and Templates. In the main pane, a pipeline named 'FinalPipeline' is selected. Under the 'sinkDataset' section, it shows an 'Azure Synapse Analytics' icon. Below this, the 'Connection' tab is active, showing a linked service 'AzureSynapseAnalytics1' with a 'Test connection' status of 'Connection successful'. The 'Table' dropdown is set to 'dbo.final_table'.

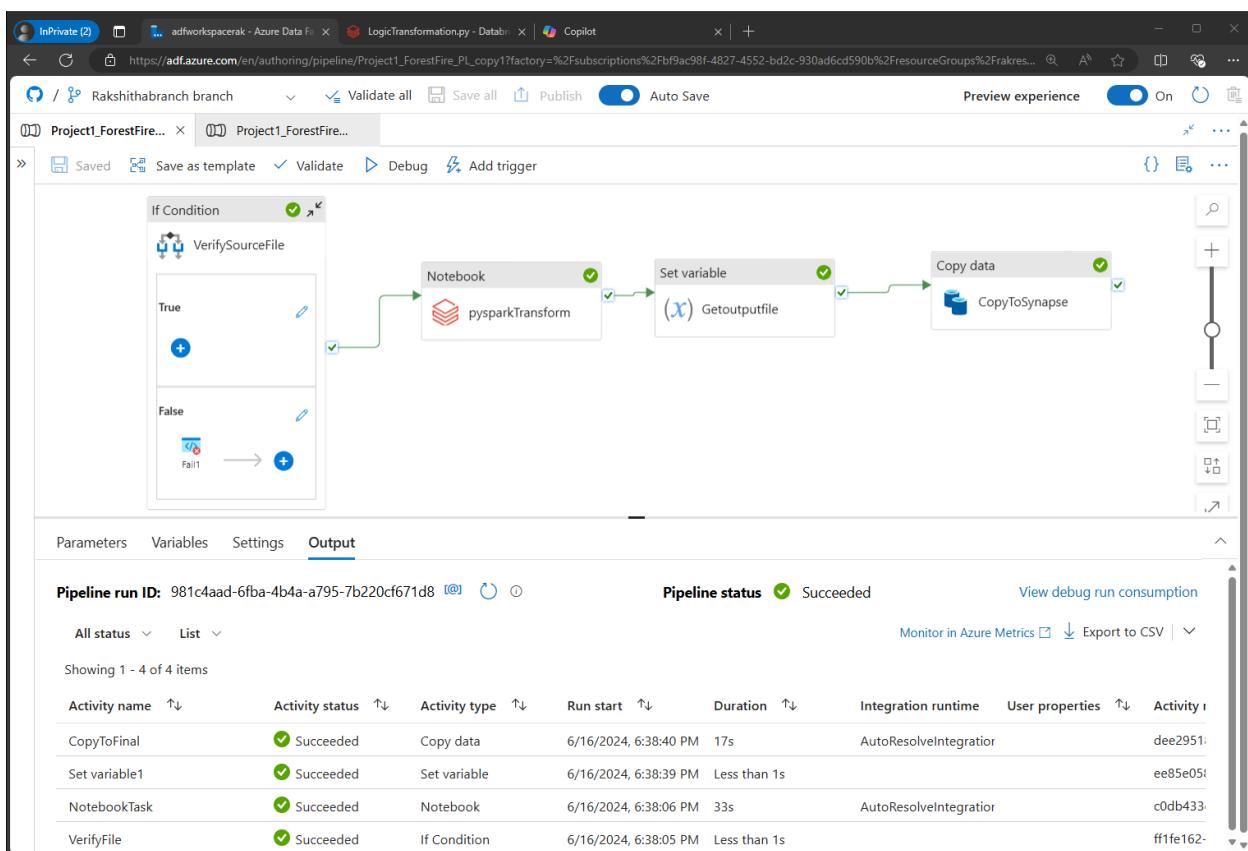
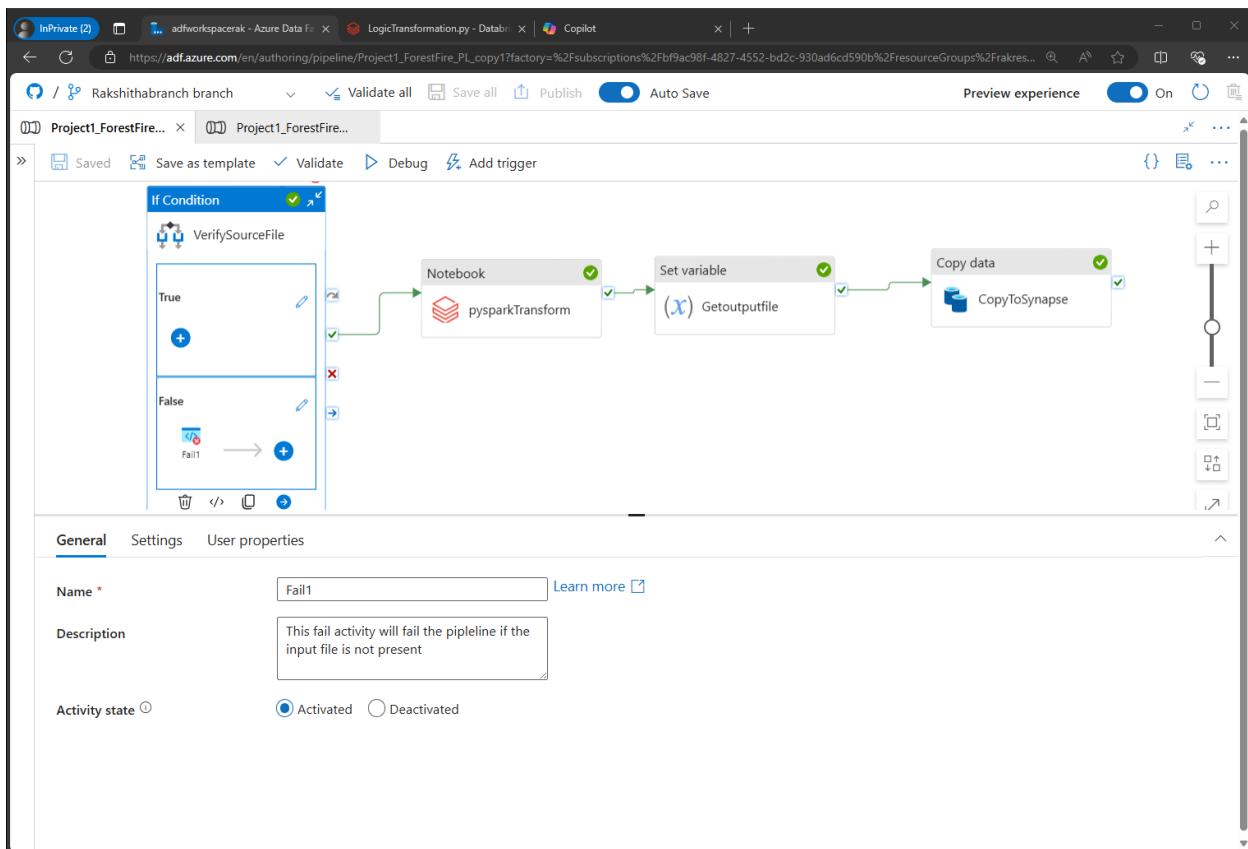


```

CREATE TABLE final_table (
    X float,
    Y float,
    month varchar(255),
    day varchar(255),
    FFMC float,
    DMC float,
    DC float,
    ISI float,
    temp float,
    RH float,
    wind float,
    rain float,
    area float,
    temperature_celsius float,
    last_modified_date varchar(255)
);
select count(*) from final_table
  
```

Added if-else condition to check if the source file is present or not





Adding schedule and alerts

The screenshot shows the Azure Data Factory studio interface. On the left, there is a pipeline canvas with a flow from 'If Condition' to 'Notebook'. The 'If Condition' step has a 'VerifySourceFile' activity. The 'Notebook' step has a 'pysparkTransform' activity. Below the canvas, there is a 'Parameters' tab with a table for defining parameters: Name (inputfilepath), Type (String), and Default value (Value). On the right, a modal dialog titled 'Add triggers' is open, showing a list of triggers: 'TriggerDaily' (Type: Schedule, Parameters: 1). A 'Choose trigger...' dropdown is also present.

The screenshot shows the Azure Data Factory studio interface. The pipeline canvas is identical to the previous one. On the right, a modal dialog titled 'Edit trigger' is open. The 'TriggerDaily' configuration includes:

- Name:** TriggerDaily
- Description:** (empty)
- Type:** ScheduleTrigger
- Start date:** 6/17/2024, 7:15:00 PM
- Time zone:** Eastern Time (US & Canada) (UTC-4)
- Recurrence:** Every 1 Day(s)
- Advanced recurrence options:**
 - Execute at these times:** (empty)
 - Hours:** (empty)
 - Minutes:** (empty)
 - Schedule execution times:** 19:15
 - Specify an end date:** (checkbox)

At the bottom of the dialog are 'OK' and 'Cancel' buttons.

Edit trigger

Trigger Run Parameters

Parameters that are not provided a value will not be included in the trigger.

Name	Type	Value
inputfilepath	string	abfs://sourcecontainer@adlsf...

Make sure to "Publish" for trigger to be activated after clicking "Save"

Save **Cancel**

Alerts

ALERT	ENABLED	RESOURCE TYPE	RESOURCES	ACTIONS
NewAlert	On	Pipeline	1	

Microsoft Azure | Data Factory > adfworkspacerak

Search factory and documentation

Alerts & metrics

ALERT ENABLED

NewAlert On

Edit alert rule

Alert rule name *: NewAlert

Description:

Severity *: Sev0

Target criteria

Whenever Pipeline Failed Runs metric is Greater

Add criteria

Notifications

Action group type

AlertGroup 1 Email

Configure notification

Enable rule upon creation

Update alert rule Cancel

This screenshot shows the 'Alerts & metrics' section of the Azure Data Factory portal. A new alert rule named 'NewAlert' is being created. The alert is set to trigger whenever the 'Pipeline Failed Runs' metric is greater than a certain threshold. It is configured to send one email notification to the specified action group.

Microsoft Azure | Data Factory > adfworkspacerak

Search factory and documentation

Alerts & metrics

ALERT ENABLED

NewAlert On

Configure alert logic

Selecting the dimension values will help you filter to the right time series.

Dimension

Name: Project1_ForestFire_PL

FailureType: 3 selected

Values

Alert logic

Condition *: Greater than

Time aggregation *: Total

Threshold count *: 1

Evaluate based on

Period *: Over the last 1 minutes

Frequency *: Every 1 minute

Update criteria Cancel

This screenshot shows the configuration of the 'NewAlert' alert rule's logic. It specifies that the alert should trigger whenever the 'Project1_ForestFire_PL' dimension value is selected and the failure type is greater than 1. The alert is set to evaluate over the last 1 minute at a frequency of every 1 minute.

Notifications	Action group type	Actions
AlertGroup	1 Email	

[+ Configure notification](#)

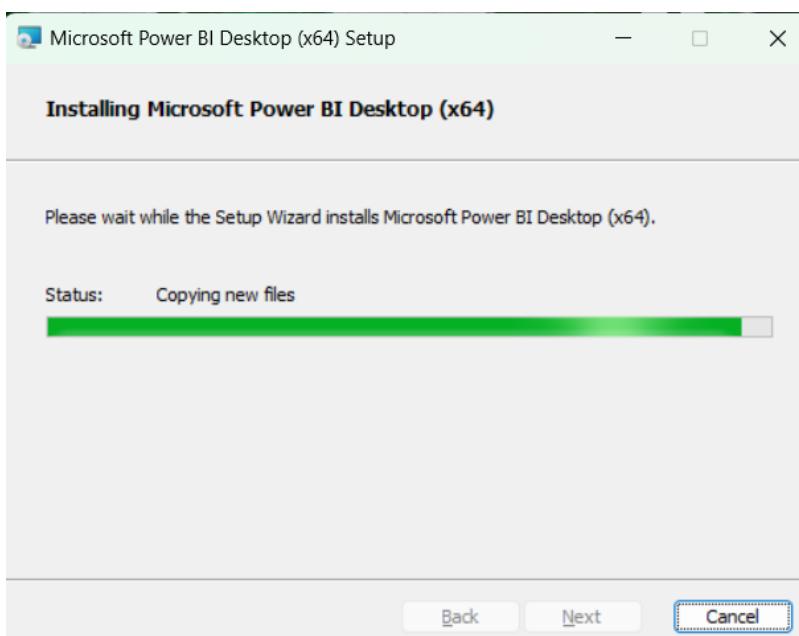
GitHub repo URL

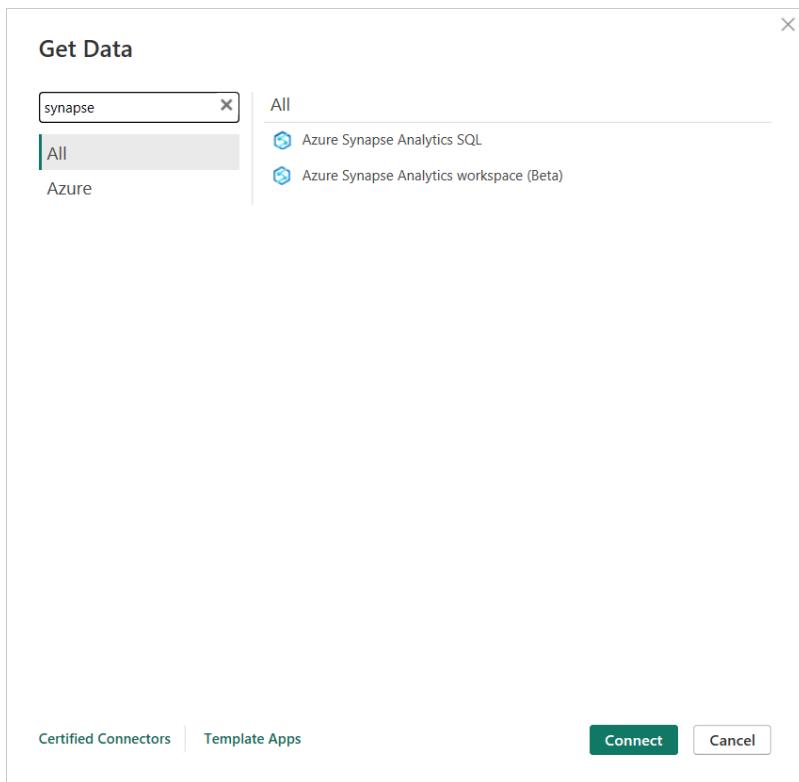
https://github.com/Rakshitha-apple/ETL_Pipeline_repo.git

The screenshot shows the GitHub repository page for 'forestfire_repo'. At the top, there's a yellow banner indicating an activity from 'adf_publish'. Below the banner, the repository structure is shown with several commits listed. On the right side, there are sections for 'About', 'Releases', and 'Packages'. The 'About' section contains a brief description of the project: 'This project is building a data pipeline using Azure services and unlocking the power of open-source data analysis'. The 'Releases' section shows 'No releases published' and a link to 'Create a new release'. The 'Packages' section shows 'No packages published' and a link to 'Publish your first package'.

The screenshot shows the GitHub repository page for 'forestfire_repo'. The left sidebar shows the file structure with a 'main' branch selected. The main content area shows a commit from 'Rakshitha-apple' updating a pipeline file. The commit message is 'Updating pipeline: Project1_ForestFire_PL'. The commit was made 37 minutes ago. The right sidebar shows the GitHub sidebar with various icons for notifications, issues, pull requests, etc.

Visualize results and prepare a dashboard in Power BI.





Microsoft Azure

keerthanaakkula@gmail.com

Home > rakswk

Search resources, services, and docs (G+ /)

Overview

Activity log

Access control (IAM)

Tags

Diagnose and solve problems

Settings

Microsoft Entra ID

Properties

Locks

Analytics pools

SQL pools

Apache Spark pools

Data Explorer pools (preview)

Resource group ([move](#))
[rakresourcegroup](#)

Status
Succeeded

Location
East US

Subscription ([move](#))
[Azure subscription 1](#)

Subscription ID
bf9ac98f-4827-4552-bd2c-930ad6cd590b

Managed virtual network
No

Managed Identity object ID
b01bd055-4430-466e-8fce-16cc7af6db05

Workspace web URL
<https://web.azuresynthesize.net?workspace=%2bsubscriptions%2fbf9ac...>

Networking
[Show firewall settings](#)

Primary ADLS Gen2 account URL
<https://rakshithastorageadls.dfs.core.windows.net>

Primary ADLS Gen2 file system
projectcontainer

SQL admin username
sqladminuser

SQL Microsoft Entra admin
<live.com#keerthanaakk> Copied

Dedicated SQL endpoint
rakswk.sql.azuresynthesize.net

Serverless SQL endpoint
rakswk-on-demand.sql.azuresynthesize.net

Development endpoint
<https://rakswk.dev.azuresynthesize.net>

Tags ([edit](#))
[Add tags](#)

Microsoft Azure | Microsoft 365 | Power BI | Copilot | Copilot log | New InPrivate

https://portal.azure.com/#@keerthanaakkula@gmail.com/resource/subscriptions/bf9ac98f-4827-4552-bd2c-930ad6cd590b/resourceGroups/rakresourcegroup/providers...

Microsoft Azure | Microsoft 365 | Power BI | Copilot | Copilot log | New InPrivate

keerthanaakkula@gmail.com | DEFAULT DIRECTORY

Home > raskwk

raskwk | SQL pools

Synapse workspace

Search | Refresh | Assign tags | Delete

Access control (IAM) Tags Diagnose and solve problems

Settings Microsoft Entra ID Properties Locks

Analytics pools

- SQL pools (selected)
- Apache Spark pools
- Data Explorer pools (preview)

Security

Encryption

Name	Type	Status	Size
Built-in	Serverless	N/A	Auto
dedicatedsqlpool	Dedicated	Online	DW300c

SQL Server database

Server ⓘ
raskwk.sql.azuresynapse.net

Database (optional)
dedicatedsqlpool

Data Connectivity mode ⓘ
 Import
 DirectQuery

Advanced options

OK Cancel

Navigator

Display Options ▾

- ▶ rakswk.sql.azuresynapse.net: dedicatedsqlpool...
 - fire_data
 - fire_data_new
 - ForestFireData_Final

No items selected for preview

Navigator

Display Options ▾

- rakswk.sql.azuresynapse.net: dedicatedsqlpool...
 - fire_data
 - fire_data_new
 - ForestFireData_Final** (selected)

ForestFireData_Final

x	y	month	day	FFMC	DMC	DC	ISI
7	5	mar	fri	86.2	26.2	94.3	
8	6	mar	fri	91.7	33.3	77.5	
8	6	mar	sun	89.3	51.3	102.2	
5	5	mar	sat	91.7	35.8	80.8	
6	4	mar	wed	89.2	27.9	70.8	
4	4	mar	tue	88.1	25.7	67.6	
4	4	mar	mon	87.2	23.9	64.7	
4	4	mar	mon	87.6	52.2	103.8	
2	2	mar	sun	89.3	51.3	102.2	
2	2	mar	sun	89.3	51.3	102.2	
4	5	mar	fri	91.7	33.3	77.5	
4	5	mar	fri	91.2	48.3	97.8	
5	4	mar	fri	91.7	33.3	77.5	
1	3	mar	mon	87.6	52.2	103.8	
6	5	mar	sat	91.7	35.8	80.8	
8	6	mar	fri	91.7	35.8	80.8	
3	4	mar	sat	69	2.4	15.5	
4	5	mar	fri	85.9	19.5	57.3	
4	5	mar	thu	91.4	30.7	74.3	
4	4	mar	fri	85.9	19.5	57.3	
3	4	mar	fri	91.7	33.3	77.5	
3	4	mar	tue	88.1	25.7	67.6	
3	5	mar	tue	88.1	25.7	67.6	

Select Related Tables Load Transform Data Cancel

Untitled - Power BI Desktop

File Home Insert Modeling View Optimize Help

Cut Copy Paste Format painter Clipboard Get data from workbook data hub OneLake Data Enter Dataverse Recent sources Data Transform Refresh data New visual Text box More visual Calculations Sensitivity Share Copilot

Build visuals with your data
Select or drag fields from the Data pane onto the report canvas.

Filters

Visualizations

Values

Add data fields here

Drill through

Cross-report

Keep all filters

Add drill-through fields here

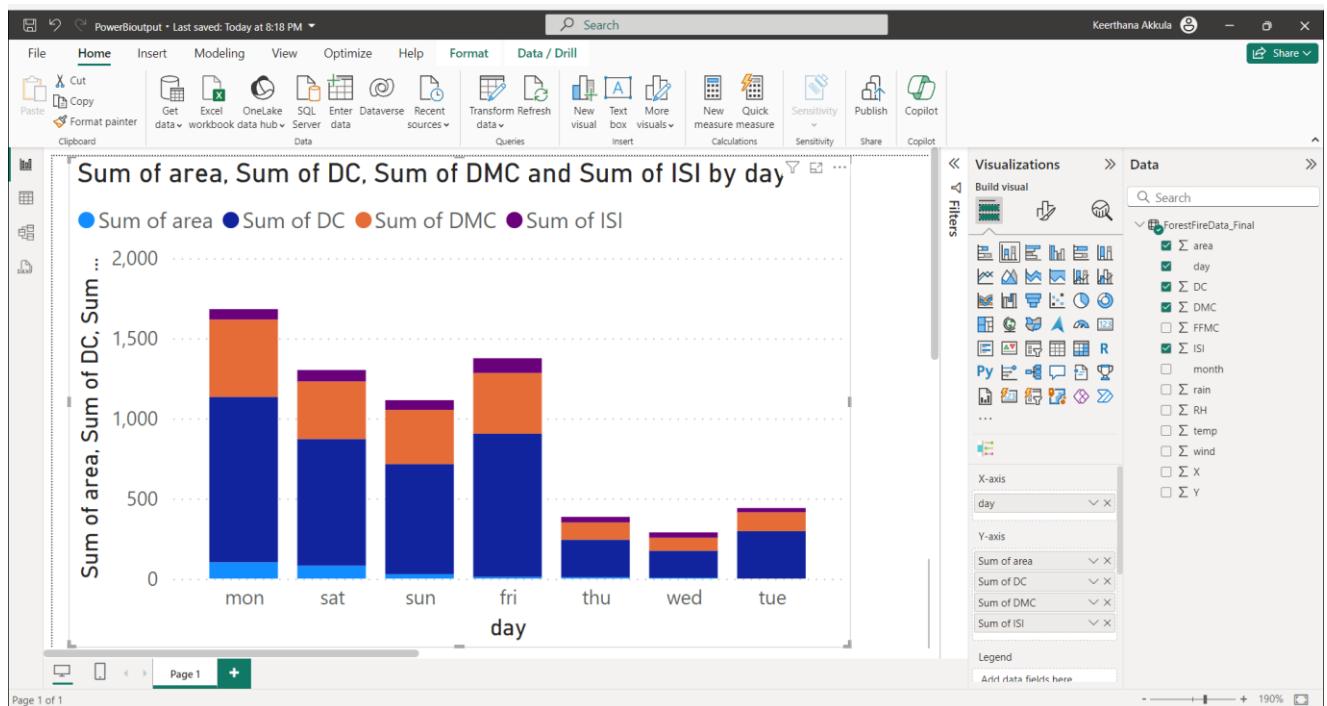
Data

Search

ForestFireData_Final

- area
- day
- \sum DC
- \sum DMC
- \sum FFMC
- \sum ISI
- \sum month
- \sum rain
- \sum RH
- \sum temp
- \sum wind
- \sum X
- \sum Y

Page 1 of 1



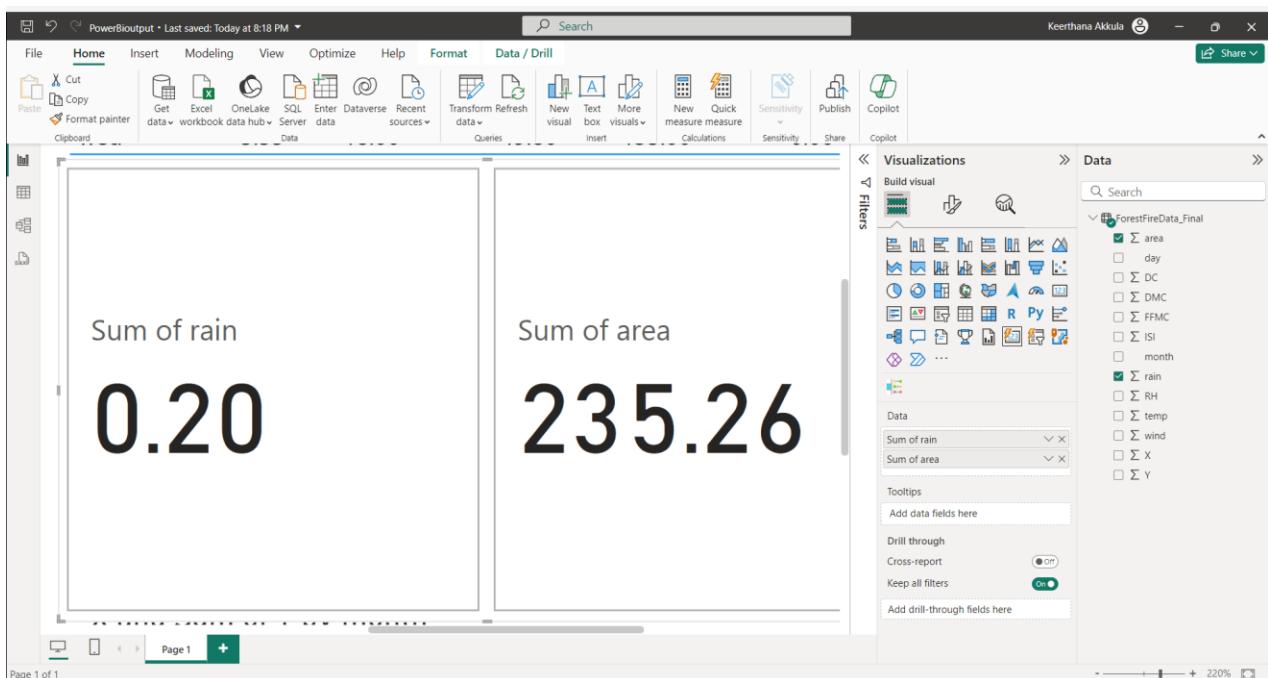
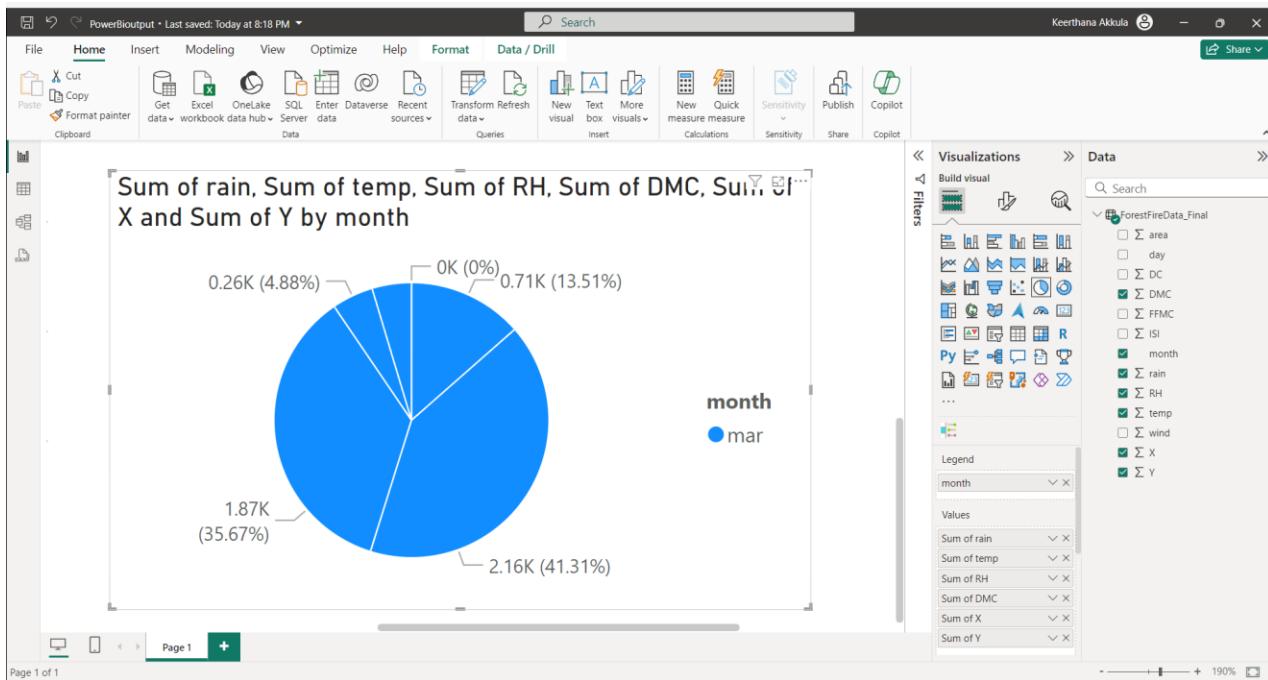
PowerBI Output - Last saved: Today at 8:18 PM

day	Sum of area	Sum of X	Sum of temp	Sum of RH	Sum of rain
fri	10.84	59.00	154.00	422.00	0.20
mon	101.94	54.00	141.40	504.00	0.00
sat	80.65	45.00	154.30	324.00	0.00
sun	27.53	37.00	74.40	396.00	0.00
thu	8.75	27.00	57.10	225.00	0.00
tue	0.00	15.00	75.50	154.00	0.00
wed	5.55	18.00	49.80	135.00	0.00
Total	235.26	255.00	706.50	2,160.00	0.20

Visualizations pane:

- Build visual
- Filters:
 - Σ area
 - day
 - Σ DC
 - Σ DMC
 - Σ FFMC
 - Σ ISI
 - month
 - Σ rain
 - Σ RH
 - Σ temp
 - Σ wind
 - Σ X
 - Σ Y
- Columns:
 - day
 - Sum of area
 - Sum of X
 - Sum of temp
 - Sum of RH
 - Sum of rain
- Drill through:
 - Cross-report

Page 1 of 1



Create the visualization using the data



PowerBIOutput • Last saved: Today at 8:18 PM

File Home Insert Modeling View Optimize Help Format Data / Drill

Paste X Cut Copy Get data workbook data hub OneLake SQL Server Data Enter Dataverse Recent sources Data Transform Refresh Queries New visual Insert More visuals Text box Quick measure measure Calculations Sensitivity Share Publish Copilot

Sum of area, Sum of DC, Sum of DMC and Sum of ISI by day

Sum of area, Sum of DC, Sum of DMC, Sum of ISI

day	Sum of area	Sum of DC	Sum of DMC	Sum of ISI
mon	~1000	~500	~200	~100
sat	~1000	~500	~200	~100
sun	~1000	~500	~200	~100
fri	~1000	~500	~200	~100
thu	~500	~200	~100	~100
wed	~500	~200	~100	~100
tue	~500	~200	~100	~100

Sum of rain 0.20 **Sum of area** 235.26

Sum of area, Sum of DC, Sum of DMC, Sum of ISI, Sum of FFMC, Sum of rain, Sum of RH, Sum of temp, Sum of wind, Sum of X and Sum of Y by day and month

month ●mar

Sum of rain, Sum of temp, Sum of RH, Sum of DMC, Sum of X and Sum of Y by month

month	Sum of rain
mar	~0.26K (4.88%)
okt	~0.71K (13.51%)
apr	~0.24K (4.47%)
mai	~1.87K (35.67%)
jun	~2.18K (41.31%)

Visualizations Data

Build visual

Filters

Search

ForestFireData_Final

- ✓ Σ area
- ✓ Σ day
- Σ DC
- Σ DMC
- Σ FFMC
- Σ ISI
- month
- ✓ Σ rain
- ✓ Σ RH
- ✓ Σ temp
- Σ wind
- ✓ Σ X
- ✓ Σ Y

Columns

day, Sum of area, Sum of X, Sum of temp, Sum of RH, Sum of rain

Drill through

Cross-report

Page 1 of 1

PowerBIOutput • Last saved: Today at 8:18 PM

File Home Insert Modeling View Optimize Help Format Data / Drill

Paste X Cut Copy Get data workbook data hub OneLake SQL Server Data Enter Dataverse Recent sources Data Transform Refresh Queries New visual Insert More visuals Text box Quick measure measure Calculations Sensitivity Share Publish Copilot

Sum of area, Sum of DC, Sum of DMC and Sum of ISI by day

Sum of area, Sum of DC, Sum of DMC, Sum of ISI

day	Sum of area	Sum of DC	Sum of DMC	Sum of ISI
mon	~1000	~500	~200	~100
sat	~1000	~500	~200	~100
sun	~1000	~500	~200	~100
fri	~1000	~500	~200	~100
thu	~500	~200	~100	~100
wed	~500	~200	~100	~100
tue	~500	~200	~100	~100

Sum of rain 0.20 **Sum of area** 235....

Sum of area, Sum of DC, Sum of DMC, Sum of ISI, Sum of FFMC, Sum of rain, Sum of RH, Sum of temp, Sum of wind, Sum of X and Sum of Y by day and month

month ●mar

Sum of rain, Sum of temp, Sum of RH, Sum of DMC, Sum of X and Sum of Y by month

month	Sum of rain
mar	~0.26K (4.88%)
okt	~0.71K (13.51%)
apr	~0.24K (4.47%)
mai	~1.87K (35.67%)
jun	~2.18K (41.31%)

Visualizations Data

Filters

Search

ForestFireData_Final

- ✓ Σ area
- ✓ Σ day
- Σ DC
- Σ DMC
- Σ FFMC
- Σ ISI
- month
- ✓ Σ rain
- ✓ Σ RH
- ✓ Σ temp
- Σ wind
- ✓ Σ X
- ✓ Σ Y

Columns

day, Sum of area, Sum of X, Sum of temp, Sum of RH, Sum of rain

Drill through

Cross-report

Page 1 of 1

Conclusion:

- Document the advantages and outcomes of the data engineering process.
- Mention improvements in data quality and the benefits for the organization.

The Importance and Benefits of Data Engineering

1. Centralizing Data:

- Data engineering helps gather and integrate data from different sources, creating a single view of the organization's operations.
- This unified data allows key stakeholders to access the information they need easily and securely.

2. Ensuring Data Security:

- Data engineers implement strong security measures to protect data throughout its lifecycle.
- Proper handling and protection of data prevent unauthorized access and cyber attacks.

3. Better Decision-Making:

- When data is accurate and well-organized, it helps analysts and executives make informed decisions.
- Organizations can use these reliable insights to optimize processes, spot trends, and adapt to market changes.

4. Improving Data Quality:

- Data engineering involves cleaning and transforming data to ensure its quality.
- High-quality data leads to better analysis, more accurate predictions, and improved business results.

5. Organizational Benefits:

1. Efficient data engineering has several advantages for organizations, including:
2. Efficiency: Automated data pipelines save time and reduce manual work.
3. Scalability: Well-designed systems can handle large amounts of data smoothly.
4. Innovation: Reliable data supports advanced analytics, like AI and machine learning.
5. Cost Savings: Effective data management reduces operational costs.
6. Competitive Edge: Organizations that use data well gain a competitive advantage.