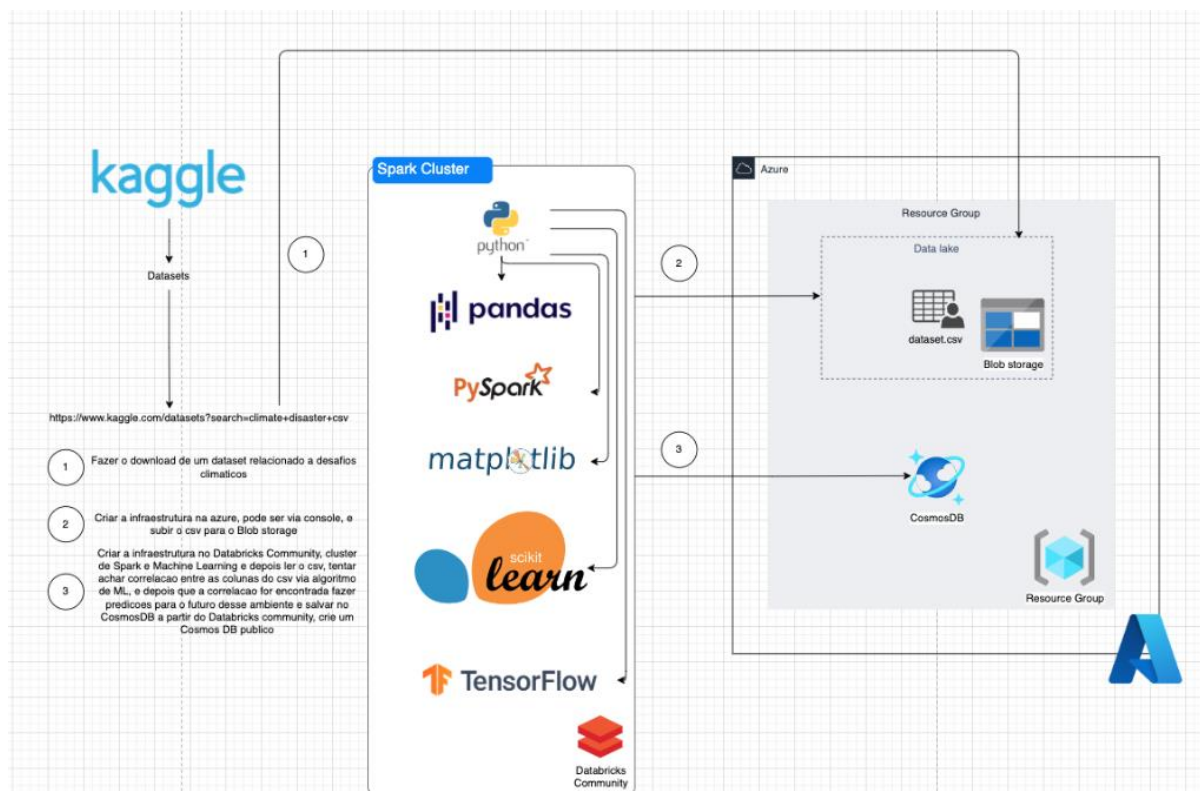


Global Solutions

| RM | Nome |
|------------------------------|--------|
| Caiky de Ávila Pereira Dutra | 99274 |
| Lucas Pinheiro de Souza | 99061 |
| Yann Dantas | 550783 |

Arquitetura solicitada para tarefa



Passo a passo da criação da nossa infraestrutura

O dataset escolhido para essa tarefa refere-se ao consumo de energia disponibilizado através desse link do Kaggle:

<https://www.kaggle.com/datasets/slawarechaniv/powerconsumption/data>

Para acesso do código fonte através do Github basta clicar neste link:

https://github.com/lucaspinheiro27/Globalsolutions_cloud

1. Criação de um Resource Group

O Resource Group é o contêiner lógico que agrupa todos os recursos da Azure para facilitar o gerenciamento. Ele permite organizar e gerenciar os recursos necessários para o projeto em um único lugar

```
Your Cloud Shell session will be ephemeral so no files or system changes will persist beyond your current session.
mr99061 [ ~ ]$ az group create --name GSResource --location brazilsouth
{
  "id": "/subscriptions/4c6283da-b703-409f-9293-b062565d0ccf/resourceGroups/GSResource",
  "location": "brazilsouth",
  "managedBy": null,
  "name": "GSResource",
  "properties": {
    "provisioningState": "Succeeded"
  },
  "tags": null,
  "type": "Microsoft.Resources/resourceGroups"
}
mr99061 [ ~ ]$
```

2. Configuração do Blob Storage

O Azure Blob Storage será usado como um data lake, onde o powerconsumption.csv será carregado. Este serviço oferece armazenamento econômico para grandes volumes de dados.

```
mr99061 [ ~ ]$ az storage account create --name globalsolutionstorage --resource-group GSResource --location brazilsouth --sku Standard_LRS
{
  "accessTier": "Hot",
  "accountMigrationInProgress": null,
  "allowBlobPublicAccess": false,
  "allowCrossTenantReplication": false,
  "allowSharedKeyAccess": null,
  "allowedCopyScope": null,
  "azureFilesIdentityBasedAuthentication": null,
  "blobRestoreStatus": null,
  "creationTime": "2024-11-21T00:56:12.633374+00:00",
  "customDomain": null,
  "defaultToOAuthAuthentication": null,
  "dnsEndpointType": null,
  "enableExtendedGroups": null,
  "enableHttpsTrafficOnly": true,
  "enableNfsV3": null,
  "encryption": {
    "encryptionIdentity": null,
    "keySource": "Microsoft.Storage",
    "keyVaultProperties": null,
    "requireInfrastructureEncryption": null,
    "services": {
      "blob": {
        "enabled": true,
        "keyType": "Account",
        "lastEnabledTime": "2024-11-21T00:56:12.727195+00:00"
      },
      "file": {
        "enabled": true,
        "keyType": "Account",
        "lastEnabledTime": "2024-11-21T00:56:12.727195+00:00"
      },
      "queue": null,
      "table": null
    }
  },
  "queue": null,
  "table": null
},
```

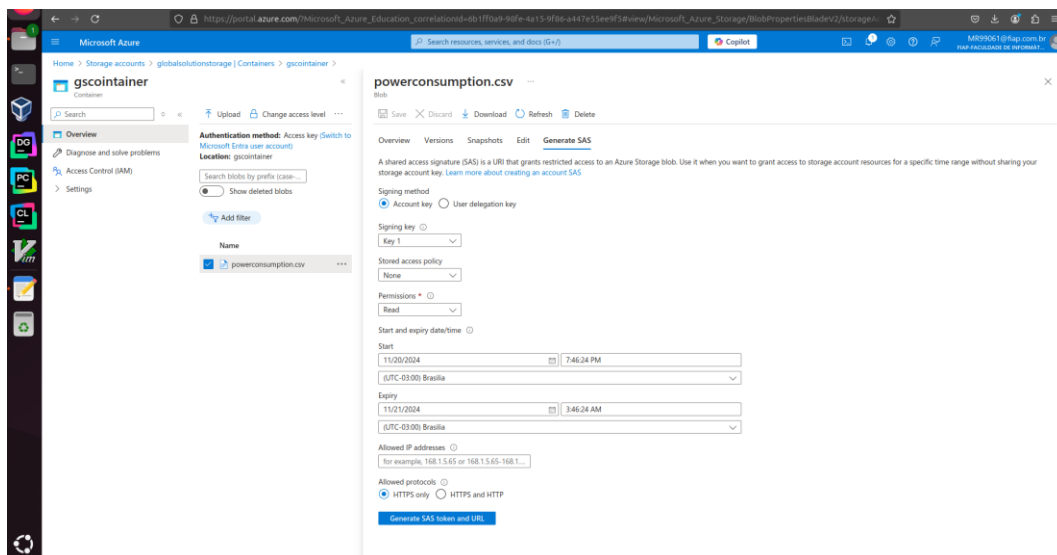
```
mr99061 [ ~ ]$ az storage container create --account-name globalsolutionstorage --name gscointainer

There are no credentials provided in your command and environment, we will query for account key for your storage account.
It is recommended to provide --connection-string, --account-key or --sas-token in your command as credentials.

You also can add '--auth-mode login' in your command to use Azure Active Directory (Azure AD) for authorization if your login account is assigned required RBAC roles.
For more information about RBAC roles in storage, visit https://docs.microsoft.com/azure/storage/common/storage-auth-aad-rbac-cli.

In addition, setting the corresponding environment variables can avoid inputting credentials in your command. Please use --help to get more information about environment variable usage.
{
  "created": true
}
mr99061 [ ~ ]$
```

Após o carregamento do arquivo, não esqueça de gerar o token SAS para o arquivo ser carregado no Databricks posteriormente.



3. Cosmos DB para Armazenamento de Resultados

O Azure Cosmos DB será utilizado para armazenar os dados processados e os resultados das análises. Ele permite que os dados estejam disponíveis para consultas em tempo real e oferece uma integração com MongoDB API, facilitando o acesso.

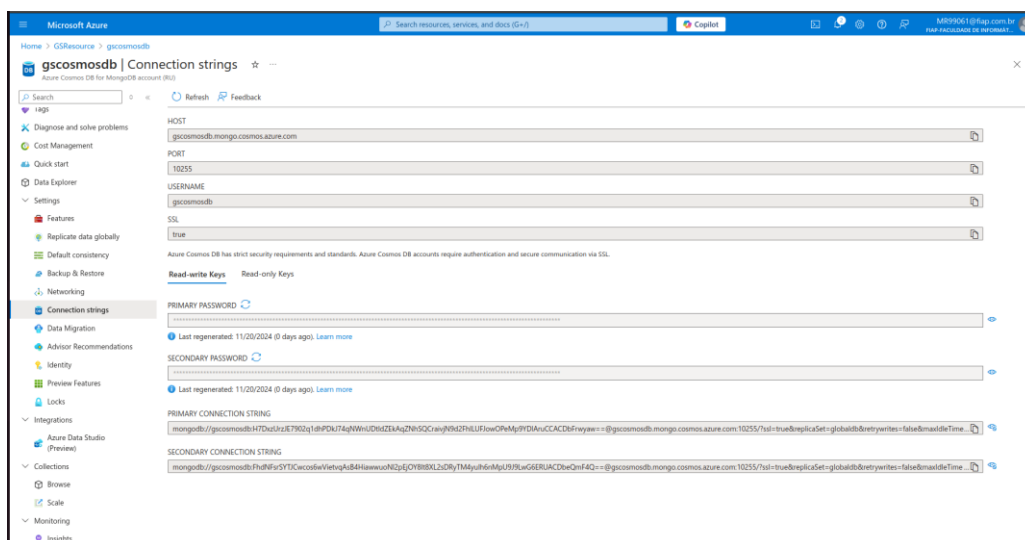
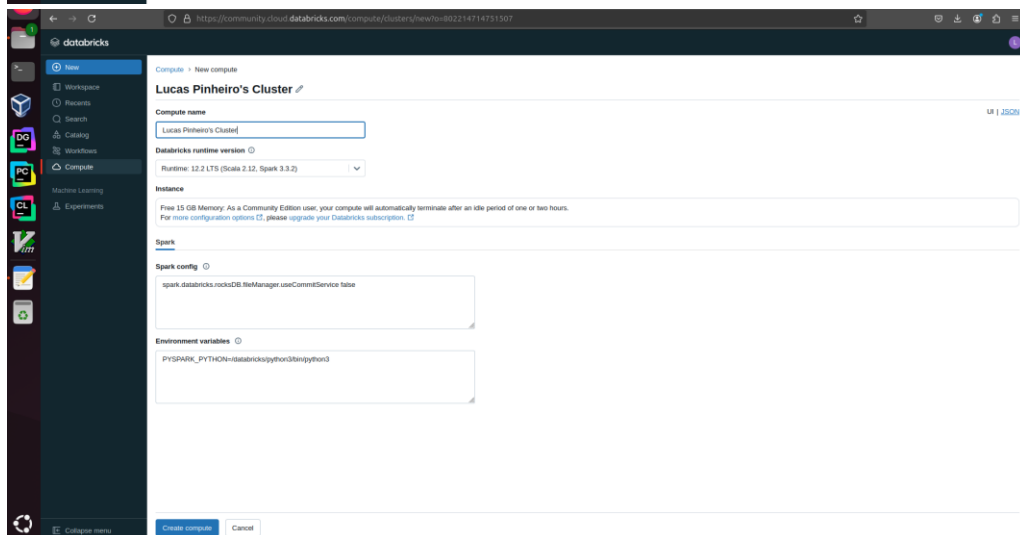
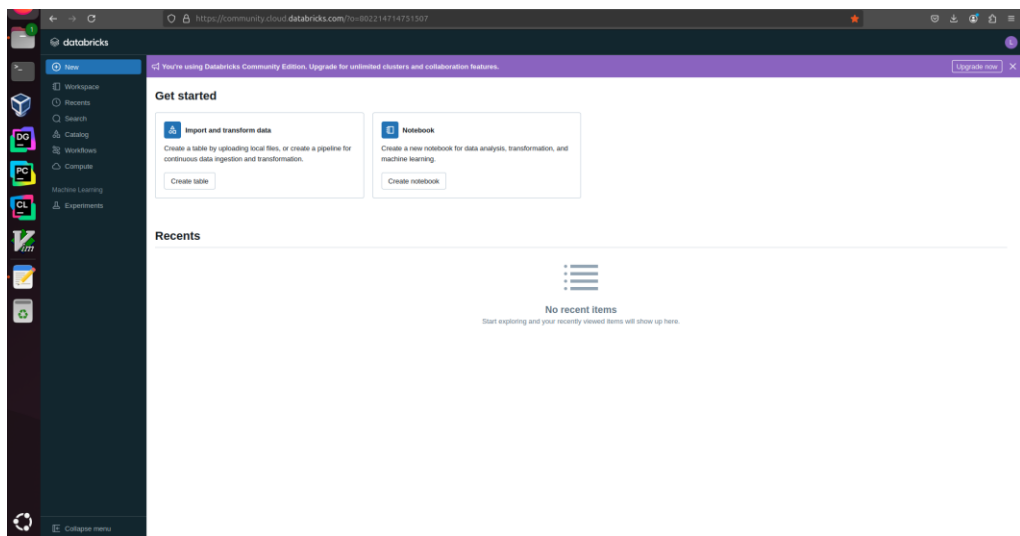
```
az cosmosdb create --name gscosmosdb --resource-group GSResource --kind MongoDB --locations regionName=brazilsouth failoverPriority=0 isZoneRedundant=False
{
  "analyticalStorageConfiguration": {
    "schemaType": "FullFidelity"
  },
  "apiProperties": {
    "serverVersion": "3.6"
  },
  "backupPolicy": {
    "migrationState": null,
    "periodicModeProperties": {
      "backupIntervalInMinutes": 240,
      "backupRetentionIntervalInHours": 8,
      "backupStorageRedundancy": "Geo"
    },
    "type": "Periodic"
  },
  "capabilities": [
    {
      "name": "EnableMongo"
    }
  ],
  "capacity": null,
  "connectorOffer": null,
  "consistencyPolicy": {
    "defaultConsistencyLevel": "Session",
    "maxIntervalInSeconds": 5,
    "maxStalenessPrefix": 100
  },
  "cors": [],
  "createMode": null,
  "customerManagedKeyStatus": null,
  "databaseAccountOfferType": "Standard",
  "defaultIdentity": "FirstPartyIdentity",
  "disableKeyBasedMetadataWriteAccess": false,
  "disableLocalAuth": false,
  "tags": {}
}
```

```
az cosmosdb mongodb database create \
--account-name gscosmosdb \
--name mydatabase \
--resource-group GSResource
{
  "id": "/subscriptions/4c6283da-b703-409f-9293-b062565dbccf/resourceGroups/GSResource/providers/Microsoft.DocumentDB/databaseAccounts/gscosmosdb/mongodatabases/mydatabase",
  "location": null,
  "name": "mydatabase",
  "options": null,
  "resource": {
    "createMode": null,
    "etag": null,
    "id": "mydatabase",
    "restoreParameters": null,
    "rid": null,
    "ts": null
  },
  "resourceGroup": "GSResource",
  "tags": null,
  "type": "Microsoft.DocumentDB/databaseAccounts/mongodatabases"
}
```

```
az cosmosdb mongodb collection create --account-name gscosmosdb --database-name mydatabase --name mycollection1 --resource-group GSResource --throughput 400
{
  "id": "/subscriptions/4c6283da-b703-409f-9293-b062565dbccf/resourceGroups/GSResource/providers/Microsoft.DocumentDB/databaseAccounts/gscosmosdb/mongodatabases/mydatabase/collections/mycollection1",
  "location": null,
  "name": "mycollection1",
  "options": null,
  "resource": {
    "analyticalStoragePct": null,
    "createMode": null,
    "etag": null,
    "id": "mycollection1",
    "indexes": [
      {
        "key": {
          "keys": [
            "id"
          ]
        },
        "options": null
      }
    ],
    "restoreParameters": null,
    "rid": null,
    "shardKey": null,
    "ts": null
  },
  "resourceGroup": "GSResource",
  "tags": null,
  "type": "Microsoft.DocumentDB/databaseAccounts/mongodatabases/collections"
}
```

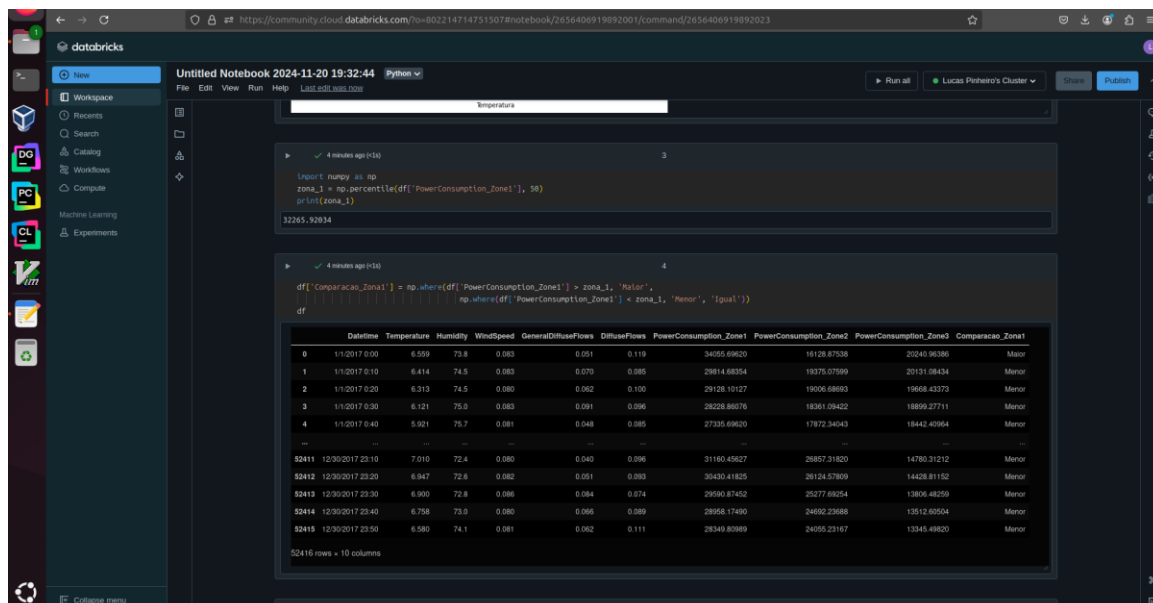
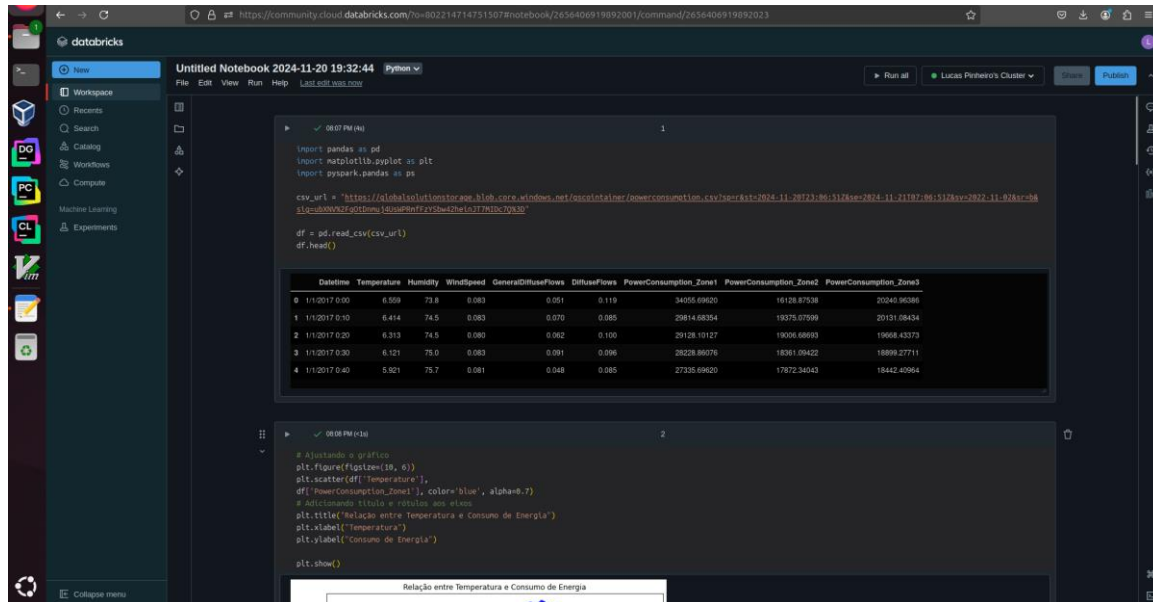
Fase 4 - Configuração do Databricks Community para o Cluster Spark:

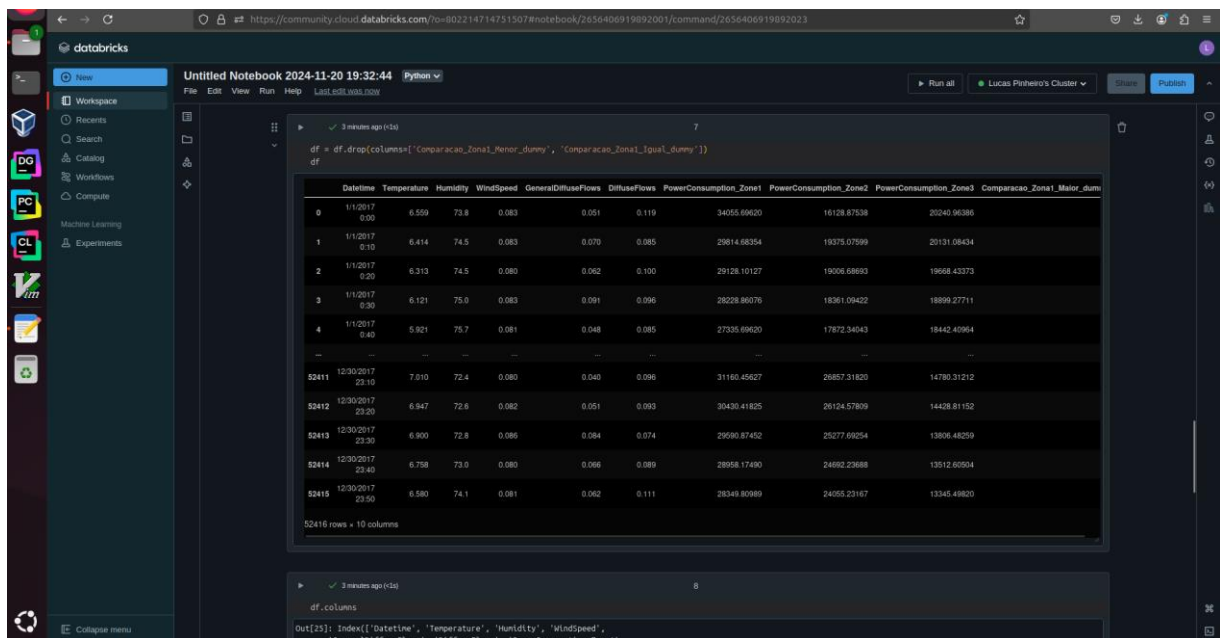
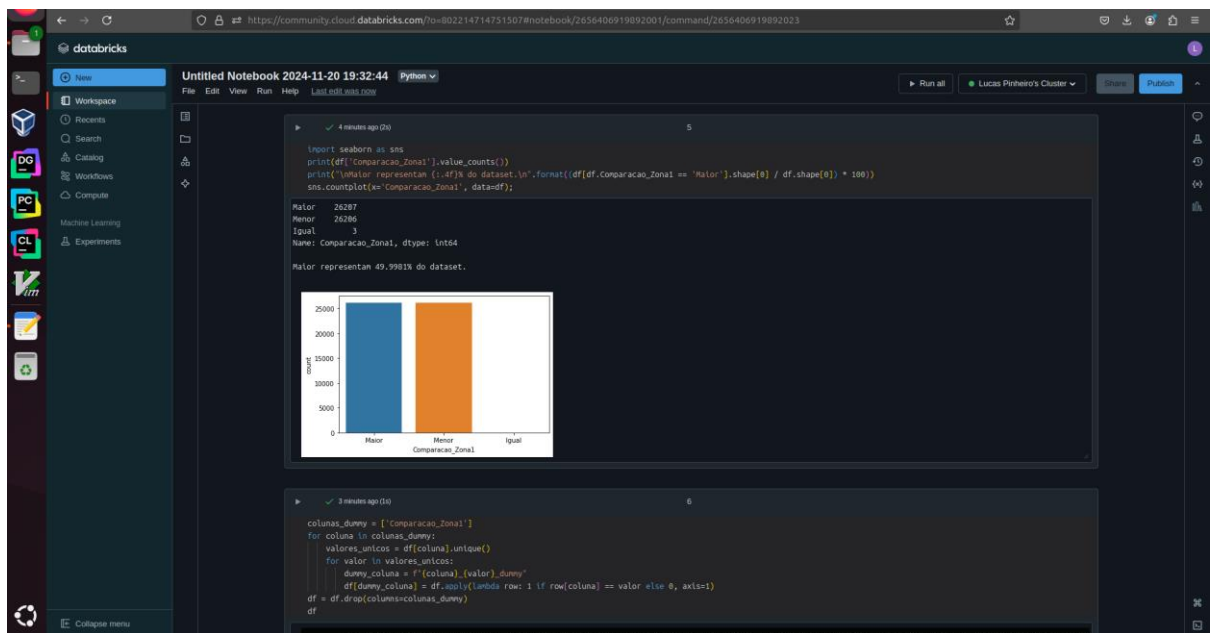
Não esqueça de copiar a connection string do cosmosdb presente no container criado anteriormente



Iremos utilizar a URL do arquivo importado para um notebook do Databricks, depois iremos encontrar algumas análises, incluindo preparação e limpeza dos dados. Neste exemplo podemos notar que há uma correlação entre a temperatura e a Zona 1, a partir do percentil de 50, podemos

fazer a classificação dessa zona se é maior, menor ou igual. Verificamos se o há um balanceamento dos dados e faremos a dummização dessa coluna para futura predição do dataset. Nesta predição utilizamos as colunas 'Temperatura', 'Humidade' e 'Velocidade do vento'. Utilizamos os algoritmos de árvore de decisão floresta aleatória





databricks

Untitled Notebook 2024-11-20 19:32:44 Python v

File Edit View Run Help Last edit was 1 minute ago

Run all Lucas Pinheiro's Cluster Share Publish

df.columns

Out[23]: Index(['Datetime', 'Temperature', 'Humidity', 'WindSpeed', 'GeneralDiffuseFlows', 'DiffuseFlows', 'PowerConsumption_Zone1', 'PowerConsumption_Zone2', 'PowerConsumption_Zone3', 'Comparacao_Zonal_Maior_dummy'], dtype='object')

columns = ['Temperature', 'Humidity', 'WindSpeed']

from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(df[columns],
df['Comparacao_Zonal_Maior_dummy'], test_size=0.5, random_state=42)

from sklearn.tree import DecisionTreeRegressor
from sklearn.tree import DecisionTreeClassifier
from sklearn import metrics
classificador = DecisionTreeClassifier()
classificador.fit(x_train, y_train)
y_pred = classificador.predict(x_test)
print("Acurácia:", metrics.accuracy_score(y_test, y_pred))

Acurácia: 0.6957417582417582

from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

databricks

Untitled Notebook 2024-11-20 19:32:44 Python v

File Edit View Run Help Last edit was 1 minute ago

Run all Lucas Pinheiro's Cluster Share Publish

from sklearn.tree import DecisionTreeClassifier
from sklearn import metrics
classificador = DecisionTreeClassifier()
classificador.fit(x_train, y_train)
y_pred = classificador.predict(x_test)
print("Acurácia:", metrics.accuracy_score(y_test, y_pred))

Acurácia: 0.6957417582417582

from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
rf_classifier = RandomForestClassifier(random_state=42)
rf_classifier.fit(x_train, y_train.values.ravel())
rf_predictions = rf_classifier.predict(x_test)
rf_accuracy = accuracy_score(y_test, rf_predictions)
print(f"Acurácia do Random Forest: {rf_accuracy}")

Acurácia do Random Forest: 0.73468827838278

Python

[Shift+Enter] to run and move to next cell
[Ctrl+Shift+F] to open the command palette
[Esc+F] to see all keyboard shortcuts

New

Workspace

Recents

Search

Catalog

Workflows

Compute

Machine Learning

Experiments

Untitled Notebook 2024-11-20 19:32:44 Python

File Edit View Run Help Last edit was 5 minutes ago

Interrupt 08:40 Lucas Pinheiro's Cluster Share Publish

```
import pymongo
import json
import pandas as pd

CONNECTION_STRING = "mongodb://gscosmosdb:FhdWfSr5Y73CwcosdwVletvqk8B4HLawuON12pEJOY81tEXL2iDHyTMeyuLhnmPU929LwG6RIUACDbeYnf4Q==@gscosmosdb-mongo.cosmos.azure.com:10255/?ssl=true&replicaSet=globaldb&retrywrites=false&maxIdleTimeMS=120000&appName=gscosmosdbg"

DB_NAME = "mydatabase1"
UNSHARDED_COLLECTION_NAME = "nl-collection"

def save_documents(collection, documents):
    """Salva uma lista de documentos na coleção."""
    if not documents:
        print("No documents to insert.")
        return

    for document in documents:
        try:
            collection.insert_one(document)
            print("Inserted document:", document)
        except Exception as e:
            print(f"Error inserting document {document}: {e}")

    print("Saved {} documents to the collection".format(len(documents)))

def main():
    client = pymongo.MongoClient(CONNECTION_STRING)

    try:
        client.server_info()
    except pymongo.errors.ServerSelectionTimeoutError:
        raise TimeoutError("Invalid API for MongoDB connection string or timed out when attempting to connect")
```

New

Workspace

Recents

Search

Catalog

Workflows

Compute

Machine Learning

Experiments

Untitled Notebook 2024-11-20 19:32:44 Python

File Edit View Run Help Last edit was 6 minutes ago

Interrupt 08:51 Lucas Pinheiro's Cluster Share Publish

```
( Interrupt 08:51 15 Python
```

```
print("Saved {} documents to the collection".format(len(documents)))

def main():
    client = pymongo.MongoClient(CONNECTION_STRING)

    try:
        client.server_info()
    except pymongo.errors.ServerSelectionTimeoutError:
        raise TimeoutError("Invalid API for MongoDB connection string or timed out when attempting to connect")

    db = client[DB_NAME]
    collection = db[UNSHARDED_COLLECTION_NAME]

    documents = json.loads(df.to_json(orient='records'))

    save_documents(collection, documents)

if __name__ == '__main__':
    main()

Inserted document: {'Datetime': '1/12/2017 22:28', 'Temperature': 15.18, 'Humidity': 83.3, 'WindSpeed': 0.88, 'GeneralDiffuseFlows': 0.877, 'DiffuseFlows': 0.137, 'PowerConsumption_Zone1': 38649.11392, 'PowerConsumption_Zone2': 23919.75884, 'PowerConsumption_Zone3': 23334.93978, 'Comparacao_Zonal_Maior_dummy': 1, '_id': ObjectId('673e7a16231c81e8a3a6baa6')}
Inserted document: {'Datetime': '1/12/2017 22:38', 'Temperature': 15.19, 'Humidity': 83.8, 'WindSpeed': 0.88, 'GeneralDiffuseFlows': 0.851, 'DiffuseFlows': 0.126, 'PowerConsumption_Zone1': 37676.96283, 'PowerConsumption_Zone2': 23354.48729, 'PowerConsumption_Zone3': 22854.93978, 'Comparacao_Zonal_Maior_dummy': 1, '_id': ObjectId('673e7a16231c81e8a3a6baa7')}
Inserted document: {'Datetime': '1/12/2017 22:48', 'Temperature': 15.2, 'Humidity': 83.1, 'WindSpeed': 0.88, 'GeneralDiffuseFlows': 0.87, 'DiffuseFlows': 0.184, 'PowerConsumption_Zone1': 36771.64557, 'PowerConsumption_Zone2': 22697.87234, 'PowerConsumption_Zone3': 22282.48964, 'Comparacao_Zonal_Maior_dummy': 1, '_id': ObjectId('673e7a16231c81e8a3a6baa8')}
Inserted document: {'Datetime': '1/12/2017 22:58', 'Temperature': 15.18, 'Humidity': 83.3, 'WindSpeed': 0.879, 'GeneralDiffuseFlows': 0.862, 'DiffuseFlows': 0.141, 'PowerConsumption_Zone1': 35671.89873, 'PowerConsumption_Zone2': 22147.11246, 'PowerConsumption_Zone3': 21744.57831, 'Comparacao_Zonal_Maior_dummy': 1, '_id': ObjectId('673e7a16231c81e8a3a6baa9')}
Inserted document: {'Datetime': '1/12/2017 23:08', 'Temperature': 15.18, 'Humidity': 83.6, 'WindSpeed': 0.875, 'GeneralDiffuseFlows': 0.84, 'DiffuseFlows': 0.093, 'PowerConsumption_Zone1': 34553.92485, 'PowerConsumption_Zone2': 21384.88243, 'PowerConsumption_Zone3': 21235.66265, 'Comparacao_Zonal_Maior_dummy': 1, '_id': ObjectId('673e7a16231c81e8a3a6baaa')}
Inserted document: {'Datetime': '1/12/2017 23:18', 'Temperature': 15.17, 'Humidity': 83.8, 'WindSpeed': 0.874, 'GeneralDiffuseFlows': 0.844, 'DiffuseFlows': 0.152, 'PowerConsumption_Zone1': 33435.04891, 'PowerConsumption_Zone2': 20625.99841, 'PowerConsumption_Zone3': 20535.66265, 'Comparacao_Zonal_Maior_dummy': 1, '_id': ObjectId('673e7a16231c81e8a3a6baab')}
```