

```
1 # =====
2 # AQUA-ANALYTICS: PREDICTIVE MODEL SCRIPT
3 # This script will connect to Google Sheet, train a model, find the key
4 # drivers, and save the results with predictions.
5 # =====
6
7 # Step A: Import necessary libraries
8 import pandas as pd
9 from sklearn.model_selection import train_test_split
10 from sklearn.ensemble import RandomForestRegressor
11 from google.colab import auth
12 import gspread
13 from google.auth import default
14
15 # Step B: Authenticate and connect to Google Sheet
16 print("--> Step 1: Connecting to Google Sheets...")
17 try:
18     auth.authenticate_user()
19     creds, _ = default()
20     gc = gspread.authorize(creds)
21     print("Authentication successful!")
22 except Exception as e:
23     print(f"Authentication failed. Please try running this cell again. Error:
24         {e}")
25
26 # --- !!! IMPORTANT CONFIGURATION !!! ---
27 # THE FULL URL of Google Sheet.
28 SPREADSHEET_URL = "https://docs.google.com/spreadsheets/d/
29 1rOp2pRQBmPPhenSPoOgE7nJHZbm8XUr6EZ7BHngfTYE/edit?gid=460422318#gid=460422318"
30 # The name of the worksheet that contains my final, clean data.
31 WORKSHEET_NAME = "Master Sheet"
32 # The name of the column I want to predict.
33 TARGET_COLUMN = "Certification_Deficit_Pct"
34 # --- END OF CONFIGURATION ---
35
36 try:
37     # To open the spreadsheet and the specific worksheet
38     worksheet = gc.open_by_url(SPREADSHEET_URL).worksheet(WORKSHEET_NAME)
39
40     # To get all the data from the sheet and convert it into a pandas DataFrame
41     data = worksheet.get_all_records()
42     df = pd.DataFrame(data)
43     print(f"Successfully loaded {len(df)} rows from '{WORKSHEET_NAME}'")
44     print("\nData Preview:")
45     print(df.head())
46 except Exception as e:
47     print(f"Failed to load data. Please check your URL and Worksheet Name.
48         Error: {e}")
49
50 # Step C: To prepare the data for the Machine Learning model
51 print("\n--> Step 2: Preparing data for the model...")
52 # To define the features (inputs) and the target (what I want to predict)
53 features = df.drop(columns=['State/UT', TARGET_COLUMN]) # Using all columns
54 except the state name and the target itself
55 target = df[TARGET_COLUMN]
56
57 # Handle any non-numeric data by converting it to numbers (one-hot encoding)
58 features = pd.get_dummies(features)
59 print("Data preparation complete.")
60
61 # Step D: Training the Predictive Model
62 print("\n--> Step 3: Training the predictive model...")
63 # I am using RandomForestRegressor, which is a powerful and reliable model for
64 this kind of task.
65 model = RandomForestRegressor(n_estimators=100, random_state=42)
66 model.fit(features, target)
67 print("Model training complete!")
68
69 # Step E: Analyzing the Key Drivers (Feature Importance)
70 print("\n--> Step 4: Analyzing the Key Drivers of Implementation Risk...")
71 # Getting the importance of each feature from the trained model
72 importances = model.feature_importances_
73 feature_names = features.columns
74
75 # Creating a DataFrame to display the results clearly
76 feature_importance_df = pd.DataFrame({'Feature': feature_names, 'Importance':
77 importances})
78 feature_importance_df = feature_importance_df.sort_values(by='Importance',
79 ascending=False)
```

```
76
77 print("\n=====")
78 print("
79 print("=====")
80 print("This table shows which factors have the biggest impact on predicting the
81 Certification Deficit.")
82 print("A higher importance score means the factor is more influential.")
83 print(feature_importance_df)
84 print("=====")
85
86 # Step F: Making Predictions and Save the Results
87 print("\n--> Step 5: Making predictions and preparing the final output...")
88 # Using the trained model to make predictions on the entire dataset
89 predictions = model.predict(features)
90
91 # Add the predictions as a new column to the original DataFrame
92 df['Predicted_Certification_Deficit'] = predictions
93
94 # Preparing the file for download
95 output_filename = "predictions_output.csv"
96 df.to_csv(output_filename, index=False)
97
98 print(f"\nSUCCESS! A file named '{output_filename}' is ready for download.")
99 print("This file contains all of the original data plus the new AI predictions.
100 ")
101 feature_importance_df.to_csv("feature_importance_output.csv", index=False)
102 print(f"\nSUCCESS! A file named '{feature_importance_output.csv}' is ready for
103 download.")
```



--> Step 1: Connecting to Google Sheets...
Authentication successful!
Successfully loaded 34 rows from 'Master Sheet'.

Data Preview:

	State/UT	Overall_FHTC_Pct	Village_Inequality_Index	\
0	Andaman and Nicobar Islands	1.00		0.00
1	Andhra Pradesh	0.74		0.11
2	Arunachal Pradesh	1.00		0.00
3	Assam	0.81		0.15
4	Bihar	0.96		0.03
	PWS_Infrastructure_Gap	Literacy_Rate_2011	Certification_Deficit_Pct	
0	0.00	85.0		0.00
1	0.00	60.0		0.23
2	0.00	60.0		0.00
3	0.02	69.0		0.44
4	0.02	60.0		1.00

--> Step 2: Preparing data for the model...
Data preparation complete.

--> Step 3: Training the predictive model...
Model training complete!

--> Step 4: Analyzing the Key Drivers of Implementation Risk...

```
=====
*** KEY DRIVERS ANALYSIS ***
=====
This table shows which factors have the biggest impact on predicting the Certification Deficit.
A higher importance score means the factor is more influential.
Feature Importance
3 Literacy_Rate_2011 0.419329
0 Overall_FHTC_Pct 0.216872
1 Village_Inequality_Index 0.188521
2 PWS_Infrastructure_Gap 0.175278
=====
```

--> Step 5: Making predictions and preparing the final output...

SUCCESS! A file named 'predictions_output.csv' is ready for download.
This file contains all of the original data plus the new AI predictions.

SUCCESS! A file named feature_importance_output.csv is ready for download.