

**CLOUD APPLICATION DEVELOPMENT**  
**BIG DATA ANALYSIS WITH IBM CLOUD DATABASE**  
**PROJECT SUBMISSION : PHASE 4**

**PROBLEM STATEMENT:**

Continue building the big data analysis solution by applying advanced analysis techniques and visualizing the results.

Apply more complex analysis techniques, such as machine learning algorithms, time series analysis, or sentiment analysis, depending on the dataset and objectives.

Create visualizations to showcase the analysis results. Use tools like Matplotlib, Plotly, or IBM Watson Studio for creating graphs and charts.

**SOLUTION:**

**(i) Machine Learning Algorithms:**

When working with IBM Cloud for big data analysis, we can leverage various machine learning algorithms depending on our specific data and objectives. IBM Cloud provides services and tools that make it easier to implement machine learning algorithms at scale. One popular service is IBM Watson Machine Learning, which integrates with IBM Cloud and supports various machine learning libraries and frameworks. Here's an example of a machine learning algorithm we can use in this context:

**Algorithm:** Random Forest

**Why Random Forest:**

Random Forest is a versatile and powerful machine learning algorithm often used for big data analysis because it is robust, scalable, and can handle a wide range of data types, including structured and semi-structured data.

**How to Implement Random Forest on IBM Cloud:**

➤ **Data Preparation:**

Prepare the dataset in a format that can be ingested by IBM Watson Machine Learning. This might involve loading data from the IBM Cloud database instance and transforming it into a suitable format (e.g., a CSV file).

➤ **Data Splitting:**

Divide the dataset into training and testing sets. This is essential to assess the performance of the Random Forest model.

➤ **Model Training:**

Use IBM Watson Machine Learning to train a Random Forest model. We can use popular Python libraries like scikit-learn, which are supported by IBM Cloud services.

Configure hyperparameters such as the number of trees, maximum depth, and minimum samples per leaf to optimize the model.

➤ **Model Evaluation:**

Evaluate the Random Forest model's performance using appropriate metrics, such as accuracy, precision, recall, and F1-score. IBM Watson Machine Learning provides tools for model evaluation and performance monitoring.

➤ **Hyperparameter Tuning:**

Fine-tune the model by adjusting hyperparameters to improve its accuracy or efficiency.

➤ **Deployment:**

Deploy the trained Random Forest model on IBM Cloud. We can use Watson Machine Learning for this purpose, making your model accessible via API endpoints.

➤ **Scalability:**

IBM Cloud provides scalability options, allowing you to handle large volumes of data efficiently. You can scale up your computational resources as needed.

### ➤ **Monitoring and Management:**

Continuously monitor the deployed model's performance and retrain it as needed to keep it up-to-date with evolving data.

Random Forest is just one example, and many other algorithms like Gradient Boosting, Support Vector Machines, or Neural Networks can also be used. IBM Cloud's services and tools provide a robust environment for implementing and scaling machine learning solutions for big data analysis.

### **RANDOM FOREST ALGORITHM:**

```
# Import necessary libraries
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Load our data and split it into features and labels
# Replace 'X' and 'y' with our feature matrix and target variable.
X = ...
y = ...

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Create a Random Forest Classifier
rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)

# Train the model on the training data
rf_classifier.fit(X_train, y_train)

# Make predictions on the test data
y_pred = rf_classifier.predict(X_test)

# Evaluate the model's accuracy
accuracy = accuracy_score(y_test, y_pred)

print(f'Accuracy: {accuracy * 100:.2f}%')

# Once we're satisfied with our model's performance, we can deploy it on IBM
Cloud using Watson Machine Learning.
```

## **(ii) TIME SERIES ANALYSIS:**

### **➤ Data Preparation:**

Start by accessing our data stored in the IBM Cloud database, making sure it contains a timestamp or date column.

Extract the necessary data for our time series analysis.

### **➤ Load Data:**

Retrieve the relevant data from our database using SQL queries or database connectors in our preferred programming language (e.g., Python or Java).

### **➤ Time Series Visualization:**

Use a data visualization library like Matplotlib or Plotly to create time series plots. Visualize our time series data to observe trends, seasonality, and anomalies.

### **EXAMPLE CODE:**

```
# Import necessary libraries

import matplotlib.pyplot as plt


# Assuming we have a time series dataset with two columns: 'Date' and 'Value'

# Replace 'date_data' and 'value_data' with our own data.

date_data = ['2023-01-01', '2023-01-02', '2023-01-03', '2023-01-04', '2023-01-05']

value_data = [100, 120, 110, 130, 95]
```

```
# Convert 'Date' column to datetime format

date_data = pd.to_datetime(date_data)


# Create a time series plot

plt.figure(figsize=(12, 6))

plt.plot(date_data, value_data, marker='o', linestyle='-')

plt.title('Time Series Plot')

plt.xlabel('Date')

plt.ylabel('Value')

plt.grid(True)


# Show the plot

plt.show()
```

### ➤ **Time Series Decomposition:**

Decompose the time series data into its individual components, typically trend, seasonality, and residuals, using techniques like additive or multiplicative decomposition.

### ➤ **Statistical Analysis:**

Perform statistical tests to identify stationarity, autocorrelation, and seasonality in the time series. Tools like the Augmented Dickey-Fuller test or the Autocorrelation Function (ACF) plot can be helpful.

➤ **Model Selection:**

Select an appropriate time series forecasting model. Common models include ARIMA (AutoRegressive Integrated Moving Average), Exponential Smoothing methods, or Seasonal Decomposition of Time Series (STL).

➤ **Model Fitting:**

Fit our chosen model to the time series data using statistical packages or libraries like statsmodels in Python. Tune model hyperparameters as necessary.

➤ **Forecasting:**

Use the fitted model to make future predictions. We can specify the number of future data points to forecast.

➤ **Model Evaluation:**

Assess the model's performance by comparing the predicted values with the actual values. Common evaluation metrics include Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE).

➤ **Visualization of Forecasts:**

Plot the forecasted values along with the actual data to visualize how well our model's predictions align with reality.

➤ **Anomaly Detection:**

Implement anomaly detection algorithms, such as the Z-score method, to identify unusual data points that deviate significantly from the expected time series pattern.

## Z-SCORE ALGORITHM:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

# Load and preprocess our time series data

# Replace 'data' with our own time series data.

mean = data.mean() # Calculate the mean and standard deviation of the time series
std = data.std()

threshold = 2.0 # Define a threshold (usually in terms of standard deviations) to
identify anomalies

z_scores = (data - mean) / std # Calculate the Z-scores for each data point

anomalies = np.abs(z_scores) > threshold # Identify anomalies by comparing Z-scores
to the threshold

# Plot the time series data with identified anomalies

plt.figure(figsize=(12, 6))

plt.plot(data, label="Time Series Data")

plt.scatter(data[anomalies].index, data[anomalies], c='red', label="Anomalies")

plt.xlabel("Time")

plt.ylabel("Value")

plt.legend()

plt.title("Time Series Anomaly Detection with Z-Score Method")

plt.show()

print("Anomalies:")

print(data[anomalies])
```

➤ **Reporting and Interpretation:**

Summarize our findings and insights from the time series analysis. Share our results with stakeholders and decision-makers.

When working with an IBM Cloud database, we may use the database's native query and data retrieval capabilities to extract the relevant time series data. Additionally, we can leverage IBM Cloud's scalability and processing power for handling large volumes of time series data efficiently.