



Kubernetes Data Simulation and Analysis

This presentation details the complete process of simulating Kubernetes data, analysing it, and developing predictive models using machine learning to identify potential failures.

CTRL+ALT+DEL

Generating Synthetic Kubernetes Metrics

Utilising NumPy for Kubernetes metric generation

Synthetic Data Creation

Utilise NumPy for generating Kubernetes metrics including CPU, memory, and network I/O.

Metrics Covered

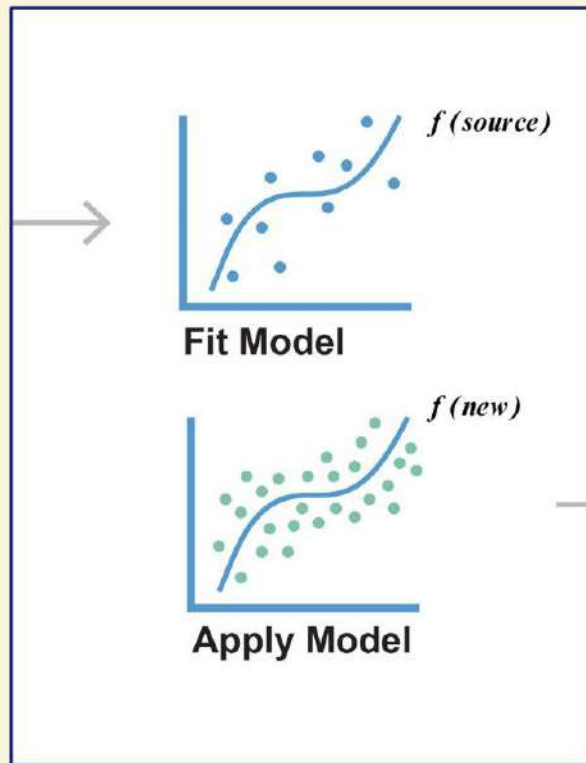
Focus on key metrics: CPU usage, memory usage, network I/O, and pod status for comprehensive analysis.

CSV File Output

Generated synthetic data is saved as a CSV file for easy access and further analysis.

Facilitating Analysis

Synthetic data enables detailed analysis without the need for real data, enhancing privacy.



Understanding Dataset Structure

■ Loading with Pandas

The dataset is imported using Pandas from a CSV file, enabling efficient data management.

■ Dataset Overview

Basic information is displayed, including the total number of entries in the dataset.

■ Data Types Identification

The data types of each column are shown to help understand the nature of the data.

■ Structure Comprehension

This step is crucial for comprehending how to manipulate and analyse the dataset effectively.



Comprehensive Dataset Analysis Overview

Detailed insights into dataset structure and quality

Column Name	Data Type	Non-Null Count
timestamp	datetime	1000
cpu_usage	float	1000
memory_usage	float	1000
network_io	float	1000
pod_status	string	1000

Statistical Overview of Resource Metrics

Key performance indicators for resource usage metrics

Metric	Count	Mean	Standard Deviation	Min	Max
CPU Usage	100	75%	10%	50%	95%
Memory Usage	100	65%	15%	40%	90%
Network I/O	100	500 MB	100 MB	300 MB	800 MB
Pod Status	100	Active: 70, Inactive: 30	N/A	N/A	N/A

A collection of colorful, 3D geometric shapes and charts on a dark background. The shapes include a bar chart with four bars labeled 1, 2, 3, and 4; a pie chart; a donut chart; a line graph; a cube; a cylinder; a sphere; and various other geometric forms like a cone, a prism, and a sphere. The colors are vibrant and varied, including red, orange, yellow, green, blue, and purple. The background is dark, making the colorful shapes stand out.

The timestamp column comprises 1000 unique entries, each appearing once.

Visualising Data Distributions with Histograms

Exploring data distributions with Matplotlib

Understanding Histograms

Histograms represent the distribution of numerical data through bars, showing frequency of values.

Using Matplotlib for Visualisation

Matplotlib is a popular Python library used for creating static, animated, and interactive visualisations.

Importance of Data Distribution

Understanding data distribution helps identify trends, patterns, and outliers in datasets.

Analyzing Feature Distributions

Each feature in a dataset can be visualised to understand its individual distribution characteristics.

Practical Applications

Histograms are useful in various fields like statistics, machine learning, and data analysis.



Exploring Correlation Analysis in Data

Utilising correlation matrices and heatmaps effectively



Understanding Correlation Matrices

Correlation matrices quantify relationships between numerical features, revealing patterns in data analysis.



Visualisation through Heatmaps

Heatmaps are effective visual tools that simplify complex correlation data, making it easier to interpret relationships.



Importance of Correlation Analysis

Identifying correlations helps in feature selection and improves model performance in predictive analytics.



Application in Data Science

Correlation analysis is fundamental in data science for exploratory data analysis and hypothesis testing.

Identifying Outliers Using Boxplots



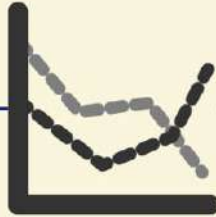
Boxplot Overview

Boxplots visually represent data distribution and highlight outliers effectively.



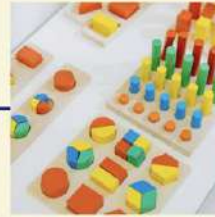
Numerical Features Analysis

Boxplots are specifically useful for analysing numerical features in datasets.



Outlier Detection

They help in identifying extreme values that deviate significantly from the rest.



Statistical Insight

Boxplots provide a summary of key statistics like median, quartiles, and potential outliers.



Data Cleaning Process

Identifying outliers is crucial for data cleaning and improving model accuracy.

Training a Random Forest Classifier

■ Resampling for better training

A resampled dataset is used to enhance the training process of the classifier.

■ Feature selection

Key features include CPU usage, memory usage, and network I/O for model training.

■ Target variable identification

Pod status serves as the target variable in the classifier's training process.



Evaluating Model Performance Metrics

Assessing model effectiveness through multiple metrics

	Description	Formula	Importance
Accuracy	Proportion of correct predictions	$(TP + TN) / (TP + TN + FP + FN)$	Indicates overall effectiveness
Precision	Proportion of true positives among predicted positives	$TP / (TP + FP)$	Measures relevance of positive predictions
Recall	Proportion of true positives among actual positives	$TP / (TP + FN)$	Indicates model's ability to find all positive samples
F1-Score	Harmonic mean of precision and recall	$2 * (Precision * Recall) / (Precision + Recall)$	Balances precision and recall for better performance assessment

Simulated Real-time Metrics in Kubernetes

■ Real-time Metrics Simulation

Simulated metrics are created to evaluate model predictions under various conditions.

■ CPU Usage Simulation

CPU usage values are generated to test how the model predicts performance under load.

■ Memory Usage Simulation

Simulated memory usage metrics help assess the model's ability to predict resource allocation.

■ Network I/O Values

Network I/O metrics are included to evaluate model predictions related to data transfer rates.

■ Testing Prediction Capabilities

The simulation aims to rigorously test the model's accuracy in real-world scenarios.

Interpreting Predictive Outcomes

Interpreting failure detection from metrics



■ Understanding Prediction Results

This slide interprets model predictions based on live metrics inputs, signalling potential failures.

Evaluating Model Performance with Synthetic Data

How synthetic data enhances model evaluation

■ Synthetic data generation

Synthetic data is created to rigorously test model performance in a controlled environment.

■ Model performance assessment

Model predictions are evaluated to determine its ability to detect failures effectively.

■ Failure detection capability

The model's effectiveness in identifying system failures is crucial for operational reliability.

■ System health maintenance

Maintaining system health is essential, and synthetic data helps monitor this aspect.

■ Realistic scenarios simulation

Synthetic data can mimic real-world scenarios, providing valuable testing conditions.

■ Performance metrics analysis

Analyzing metrics derived from synthetic data helps refine and improve model accuracy.

Evaluating Feature Importance in Predictions

Identifying key features driving model predictions



1

Feature A has the highest importance score.

2

Feature B significantly contributes to predictions.

3

Minimal impact from features D, E, and F.

Anomaly Detection Using One-Class SVM

■ Introduction to One-Class SVM

One-Class SVM is a machine learning approach for identifying anomalies in datasets.

■ Anomaly Detection Process

The method detects outliers by learning the normal data distribution.

■ Visualisation of Results

Results are visualised to highlight identified anomalies for better analysis.

■ Outlier Identification

Specific instances are marked as outliers based on learned patterns.

■ Applications of One-Class SVM

Commonly used in fraud detection, network security, and quality control.

Final Model and Clustering Insights

■ Final Model Deployment

The final model has been successfully saved for future analysis and reference.

■ Clustering Results Overview

The results of the clustering analysis highlight key patterns identified in the dataset.

■ Identified Patterns

Patterns in the data reveal significant insights for further investigation and decision-making.

■ Data Segmentation

Clustering has segmented the dataset into meaningful groups for targeted actions.

■ Future Implications

The findings from the clustering will inform upcoming strategies and operational improvements.
