

EV Charging Station Load Prediction Using PySpark

Leveraging Big Data for Smarter Grid Management and Enhanced User Experience



Forecasting EV Charging Demand

The rapid adoption of electric vehicles (EVs) introduces significant challenges for existing electric grids and charging infrastructure operators. Accurate load forecasting is critical for smarter grid management, optimising operational costs, and improving the overall user experience at charging stations.

This project documents an end-to-end Big Data solution that uses Apache PySpark and machine learning techniques to predict EV charging station load in kilowatts (kW), based on a combination of temporal, environmental, and operational factors.



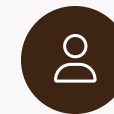
Grid Management

Enabling reliable and stable power distribution.



Cost Optimization

Minimising peak demand charges for operators.



User Experience

Reducing wait times and ensuring charger availability.

Literature Review: Existing Load Forecasting Systems

Current methodologies for EV charging load forecasting face limitations in scalability, linearity assumptions, or proprietary integration. Our PySpark pipeline offers a scalable, open-source alternative.

Statistical Time-Series Models

Models like ARIMA/SARIMA are well-suited for short-term forecasts but are constrained by linear assumptions and struggle with non-linear EV demand patterns.

Deep Learning Architectures

LSTM networks effectively capture complex temporal dependencies but necessitate vast amounts of training data and high computational resources.

Hybrid Methods

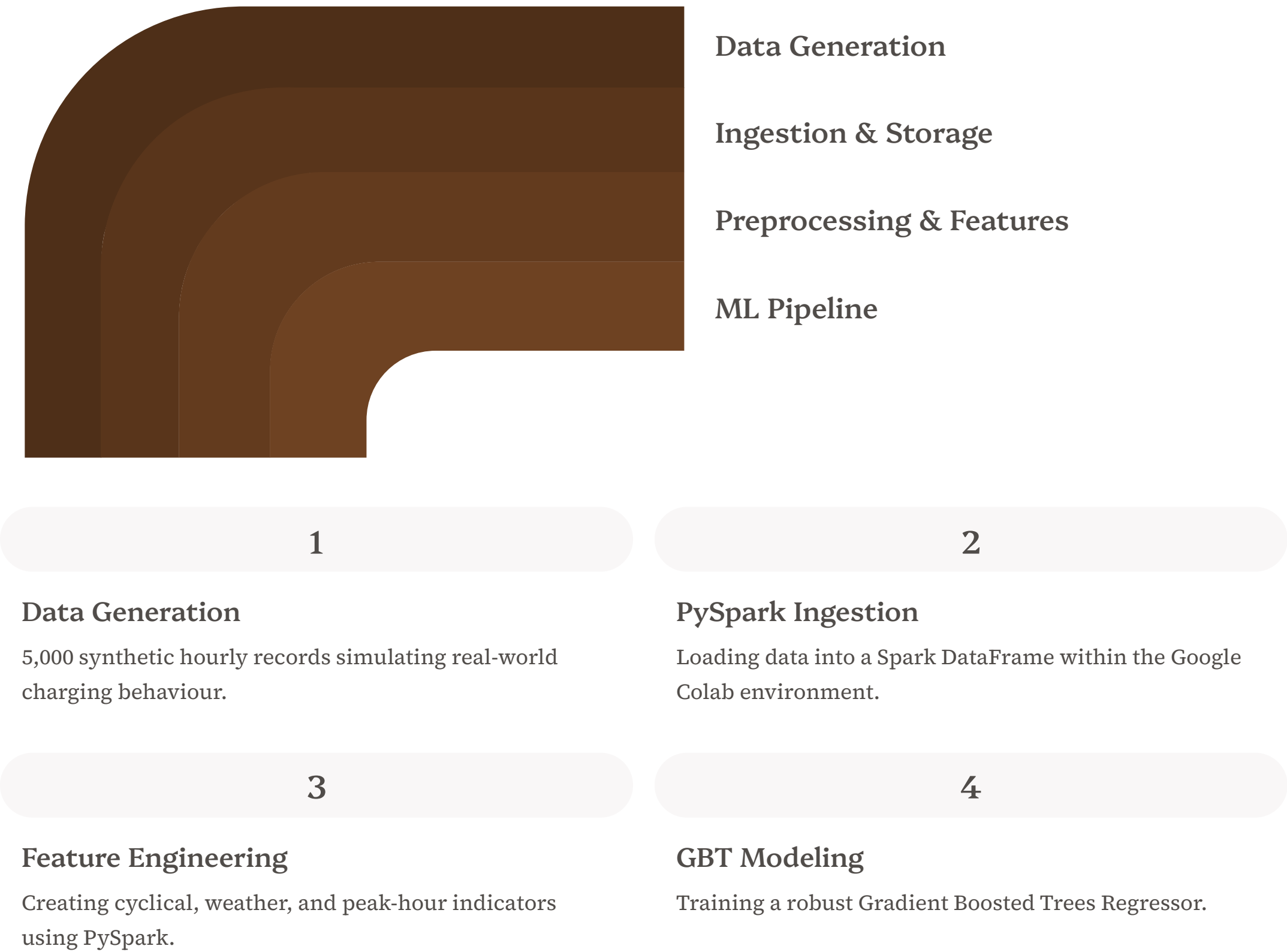
Combining weather data and calendar effects improves accuracy, but these solutions often lack the inherent scalability required for a rapidly expanding EV network.

Commercial Smart-Grid Solutions

These integrate real-time telemetry but rely on proprietary platforms, limiting customization and open-environment development.

This project differentiates itself by implementing a **scalable PySpark ML pipeline** using interpretable tree-based models and realistic synthetic data for open-source development.

End-to-End Big Data Architecture



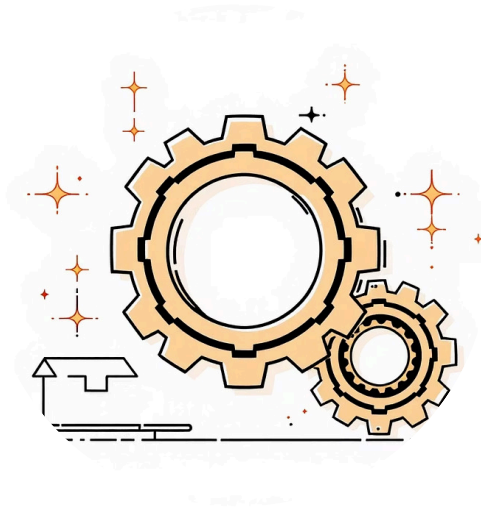
Key Pipeline Modules and Functionality

Each module is implemented using PySpark to ensure distributed processing capability, making the pipeline scalable for real-world Big Data volumes.



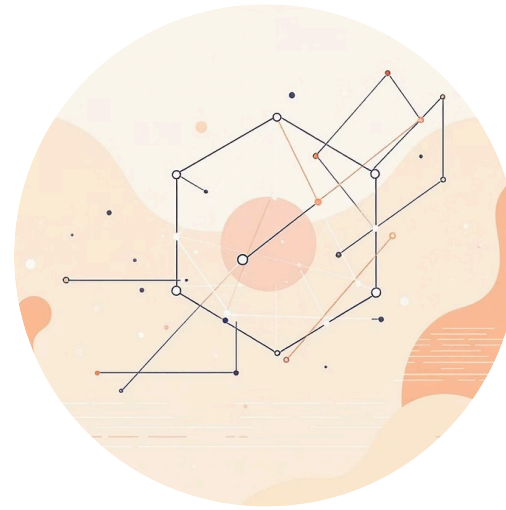
Data Generation

A Python script simulates 5,000 hourly timestamps, randomising inputs like temperature, price, and charger availability, then computes **charging_load_kW** based on defined seasonal patterns plus noise.



Feature Engineering

PySpark transformations create crucial predictors, including **hour_sin** and **hour_cos** for cyclical time encoding, a squared temperature feature, and a **peak_hour** indicator for rush periods.



Modeling & Training

A Spark MLlib pipeline stages the features with **VectorAssembler** and **StandardScaler**, culminating in a **GBRegressor** (100 trees, depth 5, step size 0.1).



Evaluation & Reporting

Data is split 80/20. We compute key metrics (RMSE, MAE, R^2) and generate a summary report, including a persistence step to save the trained model for future deployment.

Implementation Environment and Workflow

The entire solution is developed and executed within a Google Colab environment, leveraging the power of PySpark for distributed computation on the synthetic dataset.

Setup and Initialization

Installation of `pyspark` and `findspark` libraries; initialisation of the `SparkSession` to enable distributed processing.

Pipeline Construction

Building the ML pipeline stages: defining the `VectorAssembler`, setting up the `StandardScaler`, and instantiating the `GBT Regressor` for the final prediction step.

Model Persistence

The final trained `GBT` model is saved to disk, ensuring it can be loaded for future inference without retraining.

Data Loading and Schema Inspection

The synthetic dataset CSV is loaded using `spark.read.csv`, followed by a check of the inferred schema and total record count (5,000).

Training and Prediction

The pipeline is fitted on the training split. Predictions are then generated on the testing split, and results are converted to `Pandas` for plotting.

Project Results and Performance Metrics

The Gradient Boosted Trees model achieved high predictive accuracy on the simulated data, demonstrating the efficacy of distributed feature engineering.

5,00,000

Dataset Records

Total hourly simulation records used for training and testing.

10

Input Features

Number of features engineered from temporal, operational, and weather data.

0.92

R² Score

Coefficient of determination, indicating high variance explained by the model.

10 kW

RMSE

Root Mean Squared Error, the typical difference between actual and predicted load.

Top 5 Feature Importances

Temporal and operational factors were the primary drivers of charging load variability.

- **hour_of_day** (Cyclical time encoding)
- **peak_hour** (Weekday rush hour indicator)
- **temperature** (Environmental factor)
- **num_available_chargers** (Operational constraint)
- **electricity_price** (Market influence)

Conclusion & Future Roadmap

The PySpark-based GBT model successfully established a scalable and interpretable pipeline for EV charging load prediction, delivering high accuracy and valuable feature insights. The project provides a strong foundation for operational deployment.

→ Scalability Demonstrated

The distributed nature of PySpark confirms viability for processing petabytes of real-world data from large charging networks.

→ Interpretable Insights

The tree-based model provides clear feature importances, guiding infrastructure investment and pricing strategies.

→ Operational Benefits

Accurate forecasts support grid reliability, inform power purchasing, and improve overall service quality for EV users.

Future Enhancements

R\$

Integrate real-world charging telemetry data.

📶

Extend to real-time predictions using Spark Structured Streaming.

⚡

Add live weather API integration for external inputs.

🏠

Expand to multi-station geospatial analysis for network-scale optimization.