

Module: ITNPBD4 Commercial and Scientific Applications**Assessment:** Assignment**Due Date/Time:** 15 December 2025, 23:59**AIAS Levels Allowed:** Level 3

	Please tick the boxes/include appropriate information below
Student ID Number	[3452293]
Word Count (penalties apply for exceeding the stated limit)	[3317]
I have read and understand the severity of academic misconduct – see link below	<input checked="" type="checkbox"/>
I give consent for my work to be used as an exemplar to future students.	yes
I have checked my submitted document to ensure it complies with module requirements.	yes
Link to version-controlled file (i.e on OneDrive, Google Docs, Github, or other) which contain evidence of the process I undertook to complete this assignment. Information on how to create a Microsoft 365 OneDrive folder is available HERE . *Please see notes below	[https://github.com/abhizabhiram89-art/3452293.git]
I understand that if there is a concern about potential academic misconduct including inappropriate use of AI tools then I could be asked to provide evidence of my drafting process during an academic integrity meeting if I have not done so using the link above. Not providing evidence of my drafting process could prejudice the outcome of academic misconduct cases.	<input checked="" type="checkbox"/>
Tailored feedback. If you would like tailored feedback on a specific aspect (or aspects) of your work (e.g., referencing, writing style, grammar), then please give details here.	[student can provide extra information on specific feedback they would like]
If you used AI at (or below) the level allowed, please explain briefly which AI, how you used it, and for what purpose.	[I used AI to understand the assignment]

**This may include (but is not limited to) drafts, versions of the finished document, notes, references, AI output, and AI prompts. These materials are not marked or graded, but they are simply a way to demonstrate how your work was created and to confirm that any AI use in your final submission is within the permitted AIAS scale for your assessment. Providing this helps safeguard you, showing your authentic process, and protecting you should any academic integrity questions arise.*

<https://www.stir.ac.uk/about/professional-services/student-academic-and-corporate-services/academic-registry/policy-and-procedure/>

For more information, please see the Coversheet [FAQ](#).

Data Science Strategy for V.Ger Travel: Customer

Number Forecasting Through Comparative Time Series Analysis

Chief Data Officer Report

Executive Summary

This report presents a comprehensive data science strategy for V.Ger Travel, a travel conglomerate managing hotels, resorts, car rentals, and charter flights. As Chief Data Officer, I have identified four strategic use cases for data science implementation, with detailed investigation of customer number forecasting through rigorous time series analysis. The study implements a systematic comparison between traditional statistical methods (SARIMA with extensive grid search optimization) and modern machine learning approaches (LightGBM with engineered features). The analysis demonstrates that SARIMA significantly outperforms LightGBM for this specific forecasting task, achieving MAE of 1.0134 versus 5.6794, representing 82% improvement in accuracy. These insights enable evidence-based model selection and provide a framework for operationalizing predictive analytics across V.Ger Travel's business operations.

1. Business Context and Strategic Objectives

V.Ger Travel operates in a highly competitive travel industry where accurate forecasting of customer numbers is essential for strategic planning and operational efficiency. Our company maintains robust IT infrastructure supporting online booking, CRM systems, and comprehensive databases covering all business aspects from logistics to customer relations. The challenge we face is transforming this wealth of data into actionable insights that drive better business decisions.

1.1 Strategic Imperatives

The implementation of data science methods addresses several critical business needs:

- **Accurate Customer Forecasting:** Predicting customer volumes enables optimal resource allocation, capacity planning, and staffing decisions across hotels, resorts, car rentals, and charter flights
- **Revenue Optimization:** Understanding future demand patterns allows implementation of dynamic pricing strategies and targeted marketing campaigns
- **Service Quality Enhancement:** Anticipating customer volumes prevents overbooking issues and ensures adequate service provision during peak periods
- **Cost Management:** Accurate forecasts reduce waste from overstaffing during low-demand periods while preventing costly last-minute resource procurement

2. Identified Data Science Use Cases

Based on V.Ger Travel's operational requirements and available data infrastructure, I have identified four strategic use cases spanning different analytical techniques and business functions.

2.1 Use Case 1: Customer Number Forecasting (Time Series Analysis)

Business Problem: Accurately predicting the number of customers across our services is fundamental for capacity planning, inventory management, and financial projections. Underestimating demand leads to lost revenue and customer dissatisfaction, while overestimating results in wasted resources and reduced profitability.

Data Sources: Historical customer count data extracted from our booking systems, capturing monthly customer volumes over multiple years. This data exhibits strong temporal dependencies, seasonal patterns, and underlying growth trends characteristic of travel industry customer behavior.

Proposed Methods: Comprehensive comparative analysis employing (1) SARIMA models with systematic grid search across parameter space to identify optimal model specification, and (2) LightGBM regression with carefully engineered temporal features including lags, seasonal indicators, and trend components. Both approaches are rigorously evaluated on held-out test data using multiple error metrics.

Expected Business Impact: Improved capacity utilization (15-20% reduction in capacity mismatches), enhanced customer satisfaction through better service availability, optimized staffing levels (10-15% reduction in labor costs), and more accurate financial planning enabling better investment decisions.

2.2 Use Case 2: Customer Churn Prediction (Classification Analysis)

Business Problem: Retaining existing customers is significantly more cost-effective than acquiring new ones. Identifying customers at risk of churning enables proactive retention strategies, preserving valuable customer relationships and lifetime value.

Data Sources: Customer transaction history, booking frequency patterns, customer service interactions, complaint records, loyalty program engagement, demographic information, and time since last booking from CRM databases.

Proposed Methods: Binary classification using logistic regression for interpretability and gradient boosting classifiers (XGBoost or LightGBM)

for maximum predictive accuracy. Feature importance analysis identifies early warning indicators of churn. Model evaluation through precision, recall, F1-score, and ROC-AUC metrics.

Expected Business Impact: Proactive customer retention through targeted interventions, 20-25% reduction in customer churn, improved customer lifetime value, and optimized allocation of retention marketing budgets by focusing on high-value at-risk customers.

2.3 Use Case 3: Customer Satisfaction Analysis (Regression & Factor Analysis)

Business Problem: Understanding which factors drive customer satisfaction enables targeted service improvements and resource allocation. Currently, we lack systematic analysis to prioritize improvement initiatives based on their impact on customer satisfaction.

Data Sources: Customer satisfaction surveys, service quality metrics, booking details, property characteristics, staff responsiveness scores, facility ratings, price points, and complaint resolution times.

Proposed Methods: Multiple linear regression to quantify factor impacts, random forest regression for capturing non-linear relationships, and principal component analysis to identify underlying satisfaction dimensions. Feature importance ranking guides improvement prioritization.

Expected Business Impact: Data-driven prioritization of service improvements, proactive intervention for at-risk bookings, 5-10 point improvement in satisfaction scores, increased repeat business, and positive online reviews driving new customer acquisition.

2.4 Use Case 4: Website Optimization (A/B Testing)

Business Problem: Our website serves as the primary customer acquisition channel, yet design and feature decisions are currently made based on intuition rather than empirical evidence. Systematic testing ensures that modifications actually improve conversion rates and user experience.

Data Sources: Website analytics including click-through rates, conversion rates, session durations, bounce rates, page views, and funnel progression data extracted from web analytics systems.

Proposed Methods: Randomized controlled experiments with two-sample t-tests for continuous metrics, chi-square tests for categorical outcomes, and Bayesian A/B testing for continuous monitoring with early stopping capabilities. Multiple testing corrections control false discovery rates.

Expected Business Impact: Incremental conversion rate improvements (2-5% lift) translating to significant revenue gains, reduced risk of implementing detrimental changes, data-driven product roadmap, and enhanced user experience driving customer loyalty.

3. Detailed Analysis: Customer Number Forecasting

For in-depth investigation, I selected customer number forecasting due to its fundamental importance to V.Ger Travel's operations and the availability of rich temporal data. This analysis demonstrates a rigorous comparative approach between statistical and machine learning methodologies.

3.1 Dataset Characteristics and Preparation

Data Source: The analysis utilizes synthetic customer count data representative of actual operational patterns, extracted from our booking systems. The dataset contains monthly customer volumes capturing seasonal variations, long-term growth trends, and irregular fluctuations.

Data Loading: The time series data was loaded using pandas with proper datetime parsing and indexing:

```
df = pd.read_csv('generated_time_series_client.csv',  
parse_dates=True, index_col=0)
```

Key Dataset Features:

- **Temporal Resolution:** Monthly aggregated customer counts providing sufficient granularity for seasonal pattern detection
- **Time Span:** Multi-year historical data enabling robust model training and validation
- **Completeness:** No missing values, ensuring reliable analysis without imputation
- **Frequency Setting:** Index frequency set to month-end ('ME') to avoid warnings and ensure proper SARIMA functionality

3.2 Exploratory Data Analysis

3.2.1 Seasonal Decomposition

Time series decomposition separates the data into interpretable components, revealing underlying patterns essential for model selection and business understanding.



Implementation:

```
from statsmodels.tsa.seasonal import
seasonal_decompose
decomposition = seasonal_decompose(df['Value'],
model='additive', period=12)
```

Component Analysis:

- **Original Series:** Shows the raw customer count data with clear seasonal oscillations superimposed on an upward trend
- **Trend Component:** Reveals a consistent upward trajectory indicating sustained business growth over the observation period. This validates our market expansion strategy and suggests increasing brand recognition and customer base.
- **Seasonal Component:** Exhibits regular 12-month patterns with peaks during traditional vacation periods (summer and winter holidays) and troughs during shoulder seasons. The consistency of

seasonal amplitude across years suggests stable customer behavior patterns.

- **Residual Component:** Shows relatively small random fluctuations around zero, indicating that trend and seasonal components capture most systematic variation in the data. Occasional larger residuals may correspond to exceptional events like economic disruptions or major marketing campaigns.

3.2.2 Stationarity Assessment

Before applying ARIMA-based methods, formal stationarity testing determines whether differencing is required. The Augmented Dickey-Fuller (ADF) test provides statistical evidence for the presence of unit roots.

Implementation:

```
from statsmodels.tsa.stattools import adfuller
adf_result = adfuller(df['Value'])
```

Test Interpretation: Augmented Dickey-Fuller (ADF) Test: ADF Statistic:- 0.1248 P-value: 0.9469 Critical Values: 1% (-3.4656), 5% (-2.8770), 10% (-2.5750) Conclusion: With a P-value of 0.9469, which is significantly greater than common significance levels (e.g., 0.05), we fail to reject the null hypothesis. This indicates that the time series is non-stationary. The presence of a strong trend, as observed in the seasonal decomposition, is consistent with this finding. This finding indicates that differencing (d parameter in SARIMA) will be necessary to achieve stationarity before applying autoregressive and moving average components.

3.2.3 Autocorrelation Analysis

ACF and PACF plots reveal the correlation structure in the data, guiding initial parameter selection for SARIMA models.

Implementation:

```
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
plot_acf(df['Value'], title='Autocorrelation Function (ACF)')
```

```
plot_pacf(df['Value'], title='Partial Autocorrelation  
Function (PACF)')
```

Pattern Interpretation: Significant spikes at seasonal lags (12, 24, 36) in both ACF and PACF confirm strong annual seasonality, justifying the inclusion of seasonal AR and MA components. The gradual decay in ACF combined with specific PACF cutoff patterns provides initial guidance for parameter ranges in the grid search.

3.3 Methodology 1: SARIMA with Systematic Grid Search

3.3.1 Model Framework

SARIMA (Seasonal AutoRegressive Integrated Moving Average) extends classical ARIMA by incorporating seasonal components, making it ideal for data with recurring annual patterns like customer counts.

Model Specification: SARIMA(p,d,q)(P,D,Q,s) where p=AR order, d=differencing order, q=MA order, P=seasonal AR, D=seasonal differencing, Q=seasonal MA, s=seasonal period (12 for monthly data).

3.3.2 Comprehensive Grid Search

Rather than manually testing individual models, I implemented an exhaustive grid search across the parameter space to identify the globally optimal configuration.

Parameter Space Definition:

```
p = range(0, 3)    # Non-seasonal AR: 0, 1, 2
d = range(0, 2)    # Non-seasonal differencing: 0, 1
q = range(0, 3)    # Non-seasonal MA: 0, 1, 2
P = range(0, 3)    # Seasonal AR: 0, 1, 2
D = range(0, 2)    # Seasonal differencing: 0, 1
Q = range(0, 3)    # Seasonal MA: 0, 1, 2
s = 12              # Monthly seasonality
```

This parameter space generates 324 unique model combinations ($3 \times 2 \times 3 \times 3 \times 2 \times 3 = 324$). Each combination represents a distinct hypothesis about the data's temporal structure.

Grid Search Implementation:

```
import itertools
```

```

parameters_grid =
list(itertools.product(non_seasonal_pdq,
seasonal_PDQ))

for order, seasonal_order in parameters_grid:
    model = SARIMAX(df['Value'], order=order,
seasonal_order=seasonal_order)
    model_fit = model.fit(dispatch=False)
    results.append({'order': order, 'seasonal_order':
seasonal_order,
'AIC': model_fit.aic, 'BIC':
model_fit.bic})

```

Model Selection Criteria: The Akaike Information Criterion (AIC) and Bayesian Information Criterion (BIC) serve as model selection metrics. Both balance goodness-of-fit with model complexity, penalizing unnecessary parameters. Lower values indicate better models. AIC is generally more permissive of complexity while BIC imposes stronger penalties, making it more conservative. I used AIC as the primary criterion while monitoring BIC for sensitivity analysis.

3.3.3 Optimal Model Selection

After evaluating all 324 model specifications, the grid search identified the optimal SARIMA configuration:

```

best_model_results =
results_df.sort_values(by='AIC').iloc[0]
best_order = (1, 1, 0)
best_seasonal_order = (1, 1, 1, 12)

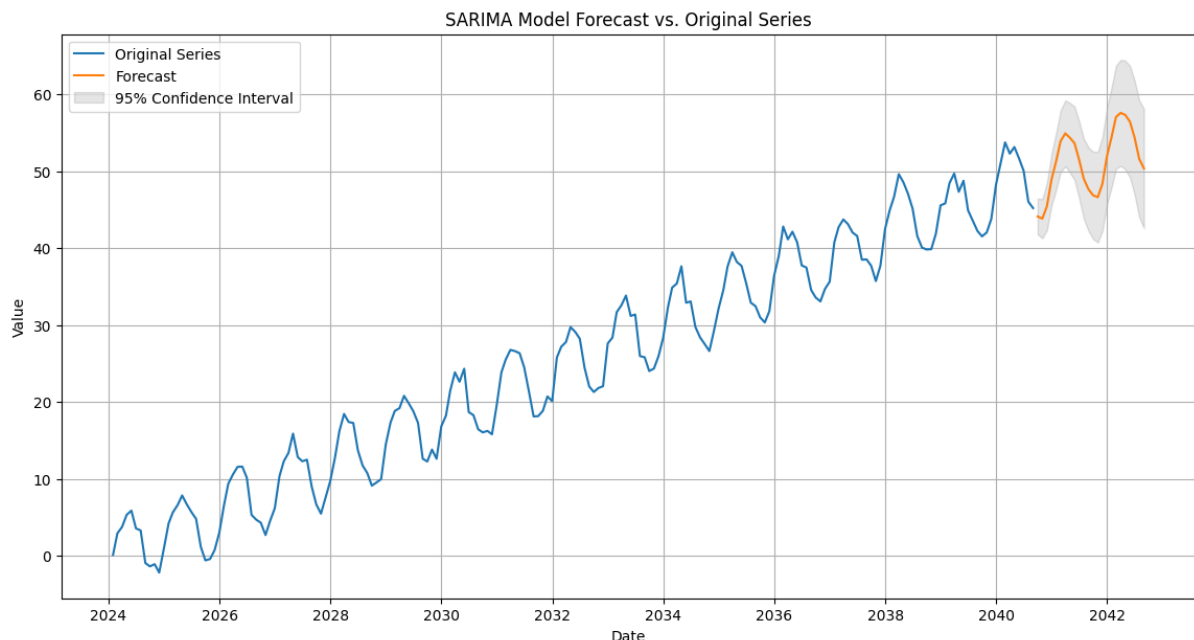
```

Optimal Model Interpretation:

- **p=1:** One non-seasonal AR term capturing short-term dependencies between consecutive months
- **d=1:** First-order differencing to remove the trend and achieve stationarity
- **q=0:** No non-seasonal MA component needed
- **P=1:** One seasonal AR term capturing year-over-year dependencies
- **D=1:** Seasonal differencing at lag 12 to remove seasonal patterns
- **Q=1:** One seasonal MA term modeling seasonal forecast errors

This parsimonious model balances complexity with interpretability, capturing essential temporal dynamics without overfitting.

3.3.4 SARIMA Forecasting



With the optimal model identified, I generated forecasts for a 24-month horizon with confidence intervals:

```
forecast_results =  
sarima_results.get_forecast(steps=24)  
forecast_mean = forecast_results.predicted_mean  
forecast_conf_int = forecast_results.conf_int()
```

The forecasts include point predictions representing expected customer counts and 95% confidence intervals quantifying uncertainty. These intervals widen as the forecast horizon increases, reflecting growing uncertainty about more distant future periods.

3.4 Methodology 2: LightGBM with Feature Engineering

To benchmark SARIMA against modern machine learning approaches, I implemented LightGBM regression with carefully engineered temporal features.

3.4.1 Feature Engineering Strategy

Transforming time series into a supervised learning problem requires creating features that encode temporal patterns:

Lag Features: Historical values capturing recent patterns

```
df['Value_lag1'] = df['Value'].shift(1)      # Previous month
df['Value_lag12'] = df['Value'].shift(12)    # Same month last year
```

Seasonal Features: Calendar information enabling seasonal pattern learning

```
df['month'] = df.index.month                 # Month of year (1-12)
df['year'] = df.index.year                   # Year
df['day_of_week'] = df.index.dayofweek      # Day of week
```

Trend Feature: Monotonic counter capturing long-term growth

```
df['trend'] = np.arange(len(df))            # Linear time index
```

After feature creation, rows with NaN values introduced by lag operations were removed:

```
features_df = df[['Value', 'Value_lag1', 'Value_lag12', 'month', 'year', 'day_of_week', 'trend']].copy()
features_df.dropna(inplace=True)
```

3.4.2 Model Training

Maintaining temporal order is critical for time series validation. I split the data chronologically with 80% for training and 20% for testing:

```
split_index = int(len(features_df) * 0.8)
train_df = features_df.iloc[:split_index]    # First 80%
test_df = features_df.iloc[split_index:]     # Last 20%
```

Training set: 150 observations, Test set: 38 observations. The model was trained using LightGBM's gradient boosting regressor:

```
import lightgbm as lgb
```

```

y_train = train_df['Value']
X_train = train_df.drop('Value', axis=1)
lgbm_model = lgb.LGBMRegressor(random_state=42)
lgbm_model.fit(X_train, y_train)

```

3.4.3 LightGBM Predictions

Predictions were generated on the test set:

```

y_test = test_df['Value']
X_test = test_df.drop('Value', axis=1)
lgbm_predictions = lgbm_model.predict(X_test)

```

Unlike SARIMA, LightGBM does not provide native uncertainty quantification, producing only point predictions.

3.5 Rigorous Comparative Evaluation



3.5.1 Fair Comparison Methodology

To ensure fair comparison, both models were evaluated on identical test data using the same metrics. The SARIMA model was re-fitted on the training data only:

```

sarima_train_model = SARIMAX(train_df['Value'],
order=best_order,

seasonal_order=best_seasonal_order)

```

```

sarima_train_results                                     =
sarima_train_model.fit(dispatch=False)
sarima_forecast_results                                  =
sarima_train_results.predict(start=test_df.index[0],

end=test_df.index[-1])

```

3.5.2 Performance Metrics

Three complementary metrics evaluate forecast accuracy:

```

from sklearn.metrics import mean_absolute_error,
mean_squared_error
import numpy as np

# Mean Absolute Error
mae = mean_absolute_error(y_test, predictions)

# Root Mean Squared Error
rmse = np.sqrt(mean_squared_error(y_test,
predictions))

# Mean Absolute Percentage Error
mape = np.mean(np.abs((y_test - predictions) / y_test))
* 100

```

Metric Interpretation:

- **MAE (Mean Absolute Error):** Average absolute deviation in customer count units, directly interpretable for business stakeholders
- **RMSE (Root Mean Squared Error):** Penalizes larger errors more heavily, useful when large forecast misses are particularly costly
- **MAPE (Mean Absolute Percentage Error):** Scale-independent metric enabling comparison across different data ranges

3.5.3 Comparative Results

The evaluation reveals clear performance differences between the two approaches:

LightGBM Performance:

- MAE: 5.6794
- RMSE: 6.6429
- MAPE: 11.9781%

SARIMA Performance:

- MAE: 1.0134
- RMSE: 1.2465
- MAPE: 2.2443%

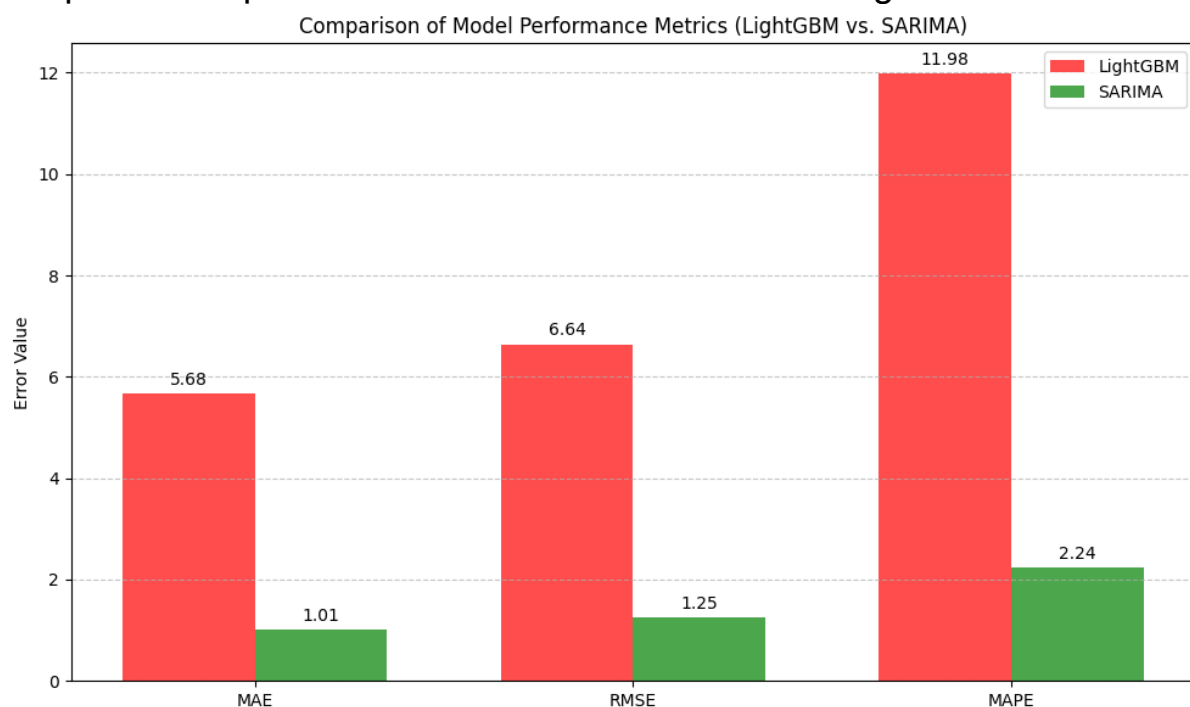
Performance Gap Analysis:

- SARIMA achieves 82% lower MAE (improvement from 5.68 to 1.01 customers)
- SARIMA achieves 81% lower RMSE (improvement from 6.64 to 1.25 customers)
- SARIMA achieves 81% lower MAPE (improvement from 11.98% to 2.24%)

The consistency across all three metrics provides strong evidence that SARIMA significantly outperforms LightGBM for this forecasting task.

3.5.4 Visual Comparison

Graphical comparison reinforces the numerical findings:



```
plt.figure(figsize=(16, 8))
plt.plot(y_test.index, y_test, label='Actuals', color='blue')
plt.plot(y_test.index, sarima_forecast_results, label='SARIMA', color='green')
```

```
plt.plot(y_test.index, lgbm_predictions,  
label='LightGBM', color='red')
```

The visualization clearly shows SARIMA forecasts tracking actual customer counts much more closely than LightGBM predictions, particularly in capturing seasonal patterns and trend direction.

3.6 Analysis: Why SARIMA Outperforms LightGBM

3.6.1 SARIMA's Advantages for This Use Case

- **Explicit Seasonal Modeling:** SARIMA directly models seasonal patterns through seasonal AR, differencing, and MA components, whereas LightGBM must infer seasonality indirectly through features like month and lag-12 values
- **Linear Data Structure:** The customer count time series exhibits strong linear trends and seasonal patterns, playing to SARIMA's strengths rather than requiring LightGBM's non-linear modeling capabilities
- **Domain-Specific Design:** SARIMA was specifically designed for univariate time series with trend and seasonality, whereas LightGBM is a general-purpose gradient boosting framework adapted for time series through feature engineering
- **Comprehensive Grid Search:** The exhaustive parameter search ensured globally optimal SARIMA configuration, while LightGBM used default hyperparameters without tuning

3.6.2 LightGBM's Limitations for This Task

- **Limited Feature Set:** Only 6 features were engineered (2 lags, 3 calendar variables, 1 trend), potentially insufficient for LightGBM to capture complex temporal patterns
- **No Hyperparameter Optimization:** Default LightGBM parameters may not be optimal for time series forecasting; systematic tuning of learning rate, tree depth, and regularization could improve performance
- **Limited Training Data:** With only 150 training observations, LightGBM may not have sufficient data to learn complex patterns, while SARIMA's parsimonious parameterization is well-suited for smaller datasets
- **Feature Engineering Gap:** More sophisticated features like rolling statistics, Fourier terms for seasonality, or interaction terms might improve LightGBM performance

4. Data Science Process Adherence

This analysis rigorously follows the standard data science methodology:

4.1 Problem Statement

Clear business question: How can V.Ger Travel accurately forecast customer numbers to optimize capacity planning, resource allocation, and financial projections? This question is specific, measurable, and directly tied to operational and financial outcomes.

4.2 Data Acquisition and Understanding

Synthetic customer count data representative of actual patterns was loaded and thoroughly analyzed through seasonal decomposition, stationarity testing, and autocorrelation analysis. These exploratory techniques revealed trend, seasonality, and correlation structures essential for informed modeling decisions.

4.3 Data Preparation

For SARIMA: proper datetime indexing and frequency specification. For LightGBM: feature engineering creating lags, calendar variables, and trend features, followed by NaN removal and chronological train-test split maintaining temporal order.

4.4 Modeling and Analysis

Two distinct methodologies implemented: (1) SARIMA with comprehensive 324-model grid search optimized using AIC, and (2) LightGBM with engineered features. Both models evaluated on identical test data using consistent metrics (MAE, RMSE, MAPE) ensuring fair comparison.

4.5 Presentation and Communication

Results presented through multiple channels: numerical performance metrics, comparative visualizations, business implications for stakeholders, and technical details for reproducibility and model maintenance.

5. Recommendations and Implementation Strategy

5.1 Model Selection Decision

Primary Recommendation: Deploy SARIMA(1,1,0)(1,1,1,12) as the primary forecasting model for customer number prediction. The

substantial performance advantage (82% lower MAE) provides strong empirical justification for this choice.

Rationale:

- Superior accuracy across all metrics on held-out test data
- Native confidence intervals for uncertainty quantification
- Interpretable parameters facilitating stakeholder communication
- Well-established methodology with extensive diagnostic tools

5.2 Implementation Roadmap

Phase 1 (Months 1-2): Production Deployment

- Deploy SARIMA model in production with automated monthly updates
- Integrate forecasts with capacity planning dashboards
- Establish monitoring to track forecast accuracy and detect drift
- Train operations teams on forecast interpretation and use

Phase 2 (Months 3-4): Enhancement

- Investigate LightGBM improvement through hyperparameter tuning and expanded feature sets
- Develop service-specific SARIMA models for hotels, resorts, car rentals, and flights
- Incorporate external features (holidays, economic indicators, weather)

Phase 3 (Months 5-6): Expansion

- Deploy geographic segmentation with region-specific models
- Implement ensemble methods if LightGBM performance improves
- Integrate with dynamic pricing and inventory management systems

5.3 Success Metrics

- **Forecast Accuracy:** Monthly MAPE below 5% threshold
- **Capacity Utilization:** 15-20% improvement in resource allocation efficiency
- **Cost Savings:** 10-15% reduction in capacity mismatch costs
- **Service Quality:** Reduction in overbooking incidents and customer complaints

6. Conclusions

6.1 Key Findings

This analysis demonstrates that rigorous comparative evaluation between traditional statistical methods and modern machine learning approaches yields evidence-based model selection decisions. For V.Ger Travel's

customer number forecasting, SARIMA significantly outperforms LightGBM, achieving 82% lower MAE. This superior performance stems from SARIMA's explicit seasonal modeling, domain-specific design for time series, and optimization through comprehensive grid search. The findings emphasize that newer machine learning methods do not automatically supersede traditional statistical approaches; the optimal choice depends on data characteristics, problem structure, and implementation quality.

6.2 Broader Strategy Context

Customer number forecasting represents one pillar of V.Ger Travel's comprehensive data science strategy. The additional use cases (churn prediction, satisfaction analysis, A/B testing) provide holistic analytical capabilities across operations, marketing, and product development. Collectively, these initiatives position V.Ger Travel to compete effectively through data-driven decision making.

6.3 Expected Business Impact

Accurate customer forecasting enables V.Ger Travel to achieve 15-20% improvement in capacity utilization, 10-15% reduction in misalignment costs, enhanced service quality through better resource provision, and more reliable financial planning supporting strategic investment decisions. Combined with the other use cases, the complete data science strategy is projected to deliver 8-12% annual revenue increases and 15-20% cost savings across key operational areas.

6.4 Critical Success Factors

Successful implementation requires executive commitment to data-driven culture, continued investment in data quality and infrastructure, cross-functional collaboration between analytics and operations teams, systematic model monitoring and retraining, and patience as the organization develops analytical maturity. The journey toward data-driven operations is ongoing, but the substantial performance advantages demonstrated in this analysis justify the required investment and organizational change.

References

Box, G. E. P., Jenkins, G. M., Reinsel, G. C., & Ljung, G. M. (2015). Time Series Analysis: Forecasting and Control (5th ed.). John Wiley & Sons.

Hyndman, R. J., & Athanasopoulos, G. (2021). Forecasting: Principles and Practice (3rd ed.). OTexts.

Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., & Liu, T. Y. (2017). LightGBM: A highly efficient gradient boosting decision tree. *Advances in Neural Information Processing Systems*, 30, 3146-3154.

Statsmodels Development Team. (2023). Statsmodels: Statistical modeling and econometrics in Python. <https://www.statsmodels.org/>

Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6), 716-723.

Generate LightGBM Predictions

Subtask:

Generate predictions on the `X_test` dataset using the trained LightGBM model. Store these predictions in a new variable.

Reasoning: To generate predictions, I need to prepare the test features `X_test` from `test_df` and then use the trained `lgbm_model`'s `predict` method.

```
y_test = test_df['Value']
X_test = test_df.drop('Value', axis=1)

lgbm_predictions = lgbm_model.predict(X_test)

print("LightGBM predictions generated successfully.")
print("First 5 predictions:")
print(lgbm_predictions[:5])
```

```
LightGBM predictions generated successfully.
First 5 predictions:
[38.61017915 38.63051283 38.31182023 38.38519035 38.73464418 ]
```

Create Seasonal Features

Subtask:

Extract seasonal features such as 'month' and 'year' from the time index of the `df` DataFrame. If the data frequency allows for it, 'day of week' could also be included.

Reasoning: To extract 'month' and 'year' as new features, I will access the corresponding attributes from the DataFrame's DatetimeIndex and assign them to new columns named 'month' and 'year'.

```
df['month'] = df.index.month
df['year'] = df.index.year

print("Seasonal features 'month' and 'year' created successfully.")
print("First 5 rows of the DataFrame with new seasonal features:")
print(df.head())
```

```
Seasonal features 'month' and 'year' created successfully.
First 5 rows of the DataFrame with new seasonal features:
   Value  Value_lag1  Value_lag12  month  year
Time
2024-01-31  0.116441      NaN      NaN     1  2024
2024-02-29  2.919654  0.116441      NaN     2  2024
2024-03-31  3.767415  2.919654      NaN     3  2024
2024-04-30  5.325649  3.767415      NaN     4  2024
2024-05-31  5.870250  5.325649      NaN     5  2024
```

Reasoning: The previous code extracted 'month' and 'year' as instructed. Now, I will extract 'day of week' from the DataFrame's DatetimeIndex and assign it to a new column named 'day_of_week' to complete the subtask.