

Bike Sharing Analysis Report

1. Executive Summary

This analysis focuses on predicting hourly bike utilization using the Capital Bikeshare dataset. After exploring multiple models—including Linear Regression, Gradient Boosting, and Support Vector Regression—we selected the RandomForestRegressor as the preferred model. The RandomForest model achieved a mean absolute deviation (MAD) of **32.86**, which is significantly lower than the MADs of the other tested models (LinearRegression: 74.11, GradientBoosting: 57.03, SVR: 94.48). A lower MAD indicates more accurate predictions on average, making the RandomForest approach more reliable for decision-making in a business context.

2. Model Selection Rationale

Several factors influenced our decision to choose the RandomForest model:

- **Performance:** With a MAD of 32.86, RandomForest provided the lowest error among all models, ensuring more dependable forecasting.
- **Robustness:** RandomForest inherently handles nonlinear relationships and interactions between features, which are common in real-world datasets like bike sharing.
- **Interpretability:** Despite being an ensemble method, RandomForest offers insights into feature importance, aiding in business decisions.
- **Scalability:** The model can be efficiently trained on large datasets and is less prone to overfitting, particularly when dealing with high-dimensional categorical data that was transformed using one-hot encoding.

3. Exploratory Data Analysis (EDA) and Visualizations

The report includes several plots to understand the data better:

- **Distribution of Total Bike Rental Count:** A histogram with a KDE overlay illustrates the distribution of bike rentals.
- **Hourly Average Rentals:** A time-series plot showing the average bike rentals per hour.
- **Bike Rentals by Weather Situation:** A box plot that reveals how different weather conditions affect rental counts.
- **Correlation Heatmap:** Displays correlations among key features (temperature, humidity, windspeed) and the target variable.

These visualizations not only support our choice of the model by confirming the presence of complex interactions in the data but also guide further feature engineering steps if needed.

4. Conclusion

In summary, the RandomForestRegressor is the most suitable model for this business case, achieving a MAD of 32.86. Its strong performance, ability to model complex nonlinearities, and interpretability make it a robust choice for operationalizing bike utilization forecasts. The following section contains the complete source code that underpins this analysis.

Relevant Source Code:

The relevant code snippet that the report is based on, includes these key functions. such as:

- **load_data:** Reads the CSV file.
- **perform_eda:** Performs exploratory data analysis and saves plots.
- **preprocess_data:** Prepares and transforms the data for modeling.
- **train_model:** Trains the RandomForestRegressor and computes the Mean Absolute Deviation.
- **main:** Runs the workflow (including unit tests).

load_data():

```
def load_data(filepath: str) -> pd.DataFrame:
    try:
        data = pd.read_csv(filepath)
        logging.info("Data loaded successfully from %s", filepath)
        return data
    except Exception as e:
        logging.error("Failed to load data: %s", e)
        raise
```

perform_Eda():

```
def perform_eda(data: pd.DataFrame, output_dir: str = "plots") -> None:
    os.makedirs(output_dir, exist_ok=True)
    # Plot creation code for distribution, hourly average, box plot, and heatmap
    ...
    logging.info("EDA plots saved in the '%s' directory.", output_dir)
```

preprocess_data():

```
def preprocess_data(data: pd.DataFrame) -> pd.DataFrame:
    features = ['season', 'yr', 'mnth', 'hr', 'holiday', 'weekday',
                'workingday', 'weathersit', 'temp', 'atemp', 'hum', 'windspeed']
    target = 'cnt'
    df = data[features + [target]].copy()
    for col in ['season', 'yr', 'mnth', 'hr', 'holiday', 'weekday', 'workingday', 'weathersit']:
        df[col] = df[col].astype('category')
    logging.info("Data preprocessing completed.")
    return df
```

train_model()

```
def train_model(data: pd.DataFrame, random_state: int = 42) -> (RandomForestRegressor, float):
    X = data.drop(columns='cnt')
    y = data['cnt']
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=random_state)
    X_train_enc = pd.get_dummies(X_train, drop_first=True)
    X_test_enc = pd.get_dummies(X_test, drop_first=True)
    X_test_enc = X_test_enc.reindex(columns=X_train_enc.columns, fill_value=0)
    model = RandomForestRegressor(n_estimators=100, random_state=random_state)
    model.fit(X_train_enc, y_train)
    predictions = model.predict(X_test_enc)
    mad = mean_absolute_error(y_test, predictions)
    logging.info("Model trained successfully. MAD: %.2f", mad)
    return model, mad
```

Appendices

Appendix A – Contains the plots generated during the EDA (saved in the "src/plots" directory):

- **cnt_distribution.png**
- **hourly_avg.png**
- **cnt_by_weathersit.png**
- **correlation_heatmap.png**

These visualizations are attached as below:

1. **cnt_distribution.png:**

This image shows the distribution of the total bike rental counts (the "cnt" variable). It's a histogram (with a kernel density estimate overlay) that illustrates how frequently different rental count values occur. This helps in understanding the spread, central tendency, and skewness of the rental count data.

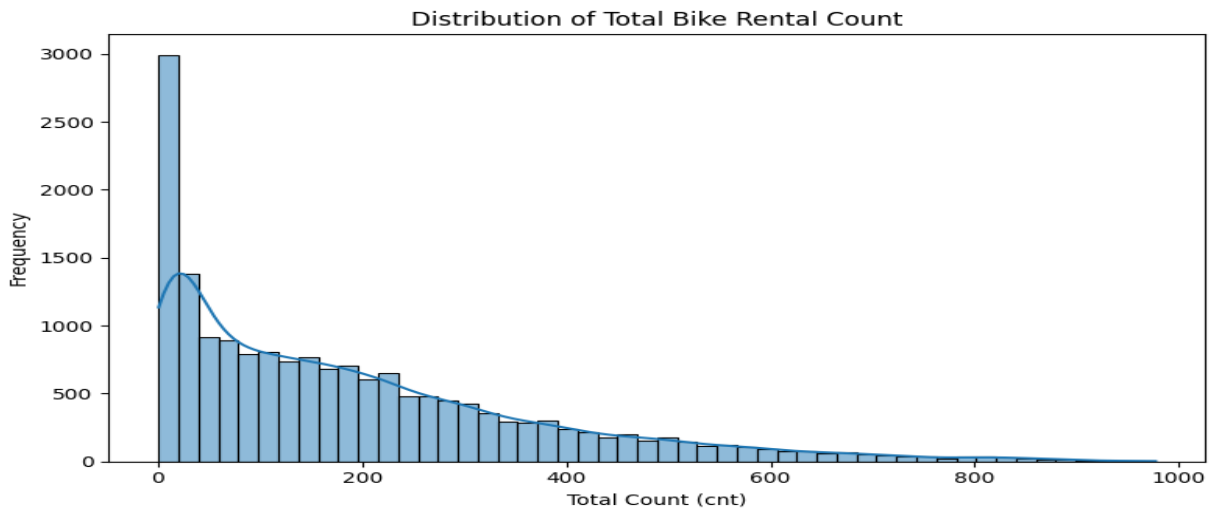


Figure 1: cnt_distribution.png

2. **hourly_avg.png:**

This plot visualizes the average bike rental count for each hour of the day. By grouping the data by the "hr" variable and calculating the mean "cnt" for each hour, the chart highlights trends in bike usage over the course of a day. It is useful for identifying peak hours and patterns in bike rental behavior.

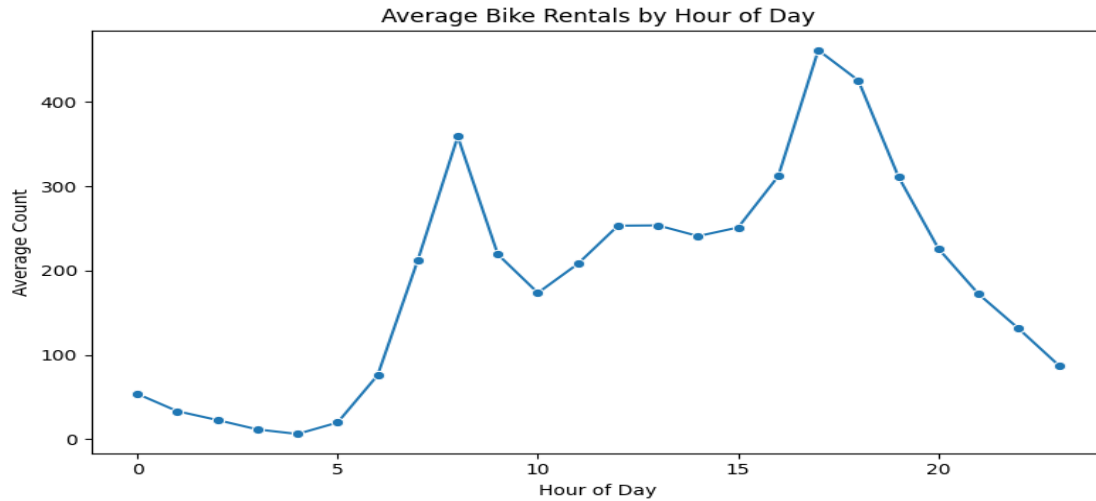


Figure 2: hourly_avg.png

3. **cnt_by_weathersit.png:**

Here, a box plot is used to show the distribution of bike rental counts across different weather situations (indicated by the "weathersit" variable). The box plot displays the median, quartiles, and potential outliers for rental counts under each weather condition, allowing an assessment of how weather impacts bike usage.

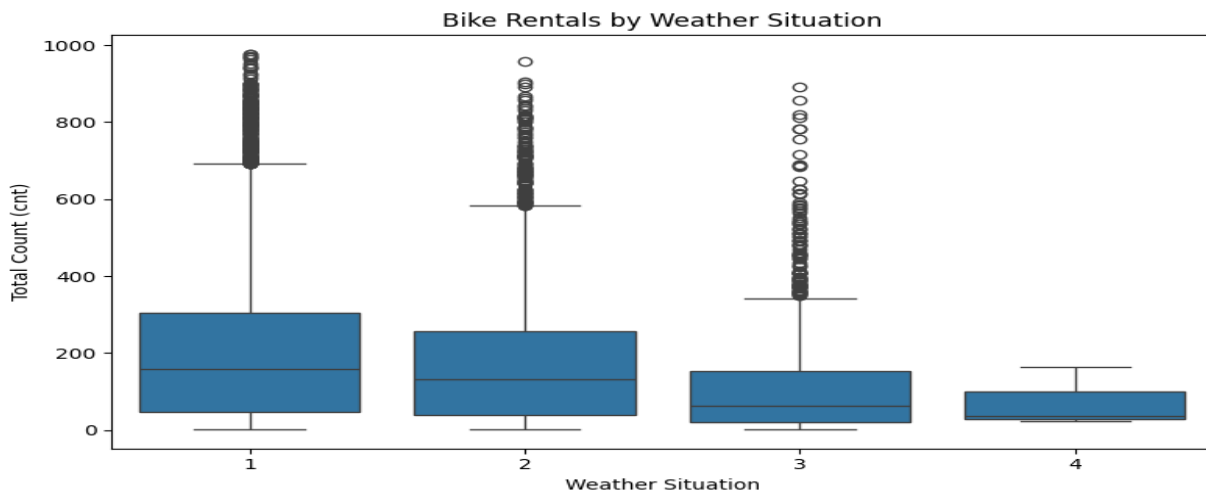


Figure 3: cnt_by_weathersit.png

4. **correlation_heatmap.png:**

This heatmap represents the correlation matrix among selected variables: temperature (temp), “feels like” temperature (atemp), humidity (hum), windspeed, and the total rental count (cnt). By visualizing these correlations, one can quickly identify which factors have stronger or weaker linear relationships with bike rental counts, providing insight into potential predictors.

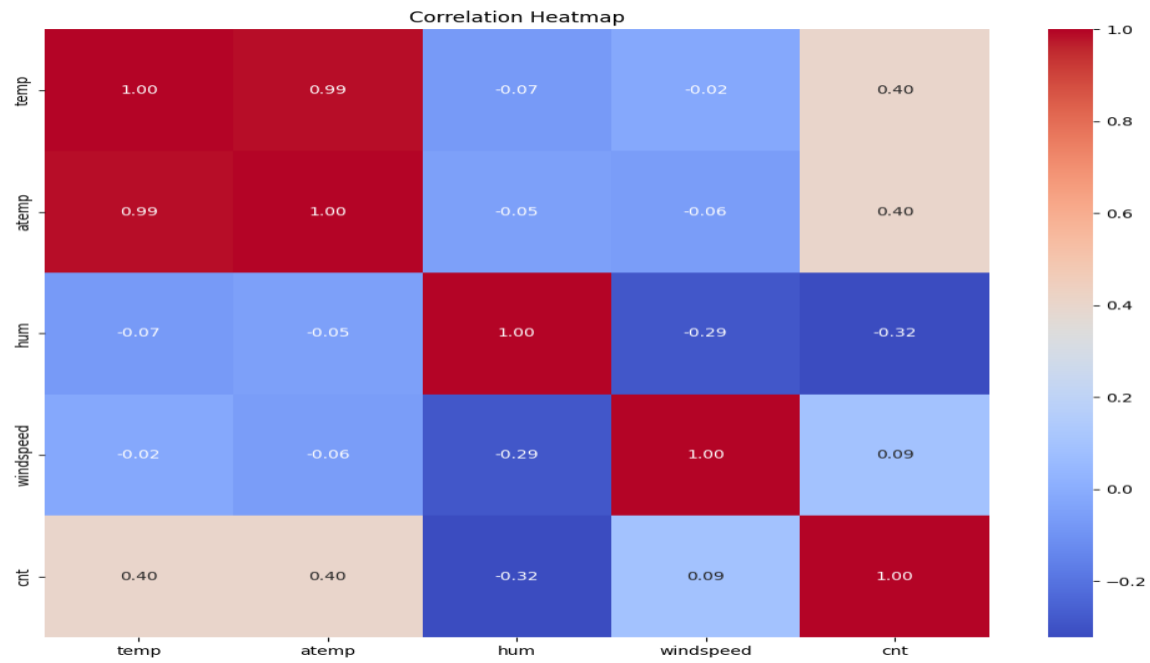


Figure 4: correlation_heatmap.png

Each visualization offers a different perspective on the dataset, enabling a comprehensive exploratory data analysis that informs the subsequent model selection and development.