# Sales & Demand Analysis Technical Report

**Disclaimer:** To protect client confidentiality, the data presented in this report has been intentionally adjusted. All numerical values have been modified to use fictional figures, and any geographic, product-specific, or sales location information has been generalised. These changes preserve the analytical value of the work while ensuring anonymity and privacy.

# Business Context

A data science consultancy partnered with a large retailer and consumer goods manufacturer to address supply and demand challenges driven by climate change. In 2022, a severe European heatwave led to a widespread shortage of temperature-sensitive goods, leading to significant missed sales opportunities and operational disruption. This project analyses sales, temperature, and event data from 2021–2022 across retail locations within a major metropolitan area in Europe to uncover patterns influencing demand.

Key questions guiding the analysis include:

**Temperature & Weather**

- What temperature thresholds correlate with spikes in sales?
- How do weather conditions affect purchasing behaviour?

**Time & Location**

- How do sales vary by day of week, month, season?
- How do holidays and events influence demand?

**Demand Forecasting**

- Can historical patterns predict future demand for better stock planning?

**Operational Efficiency**

- What is the financial impact of stockouts versus overstocking?

**Key Questions Summary**

- How can we balance supply with demand forecasting?
- What external factors should drive inventory decisions?
- How can profitability be optimised while ensuring product availability?
- The findings will support more resilient, data-driven supply chain strategies.

# Project Development

## Tools and Libraries

Python was used in Colab for analysis and modeling; Tableau supported visuals.

Core libraries included pandas, numpy, seaborn, matplotlib, with forecasting via HistGradientBoosting, and time series split using scikit-learn. Feature engineering and evaluation used scikit-learn, statsmodels, and meteostat[1] weather data.

## Data Cleaning

**Sales Data (2021 & 2022):**
- Dates standardised (YYYY-MM-DD) with added features (e.g. weekday, season).
- No duplicates/missing values; negative sales confirmed as refunds.
- Columns renamed for consistency; locations cleaned and standardized.
- Filtered to stations within the major metropolitan area (Appendix 1).
- 2022 has fewer days (348 vs. 365 in 2021).

**Temperature Data (2021 & 2022):**
- Dates cleaned with derived features; columns renamed for clarity.
- Dropped empty columns; no duplicates.
- Retained missing values (e.g. avg_temp, pressure) to avoid bias.
- Outliers kept, as they reflect real weather (Appendix 2).

## Final Preparation:

2021-2022 sales and temperature data were cleaned and merged to analyse climate impact on sales.

External data on public holidays and major sporting events was also added.

---

[1]

**Assumptions:**
- Products sold in 2021 and 2022 are assumed to be standard product.
- Standard products are sold at ~ €1.50 with a profit margin of ~ €0.20.
- Premium products are sold at ~ €1.90 with a margin of ~ €0.40.
- Missing daily data for certain sales locations indicates zero sales on those days.

**Limitations:**
- Seasonality is difficult to assess due to missing sales data from 15–31 December 2022.
- Sales limitations from stock shortages are not documented.
- No data is available on other sales points, limiting context on demand fluctuations at sales locations.

# Visualisations & Insights

## Exploratory Analysis

We began by analysing sales and temperature trends across country and metropolitan area, identifying top-selling locations and evaluating the impact of events and holidays. This guided deeper dives into district-level sales, North-South area differences, and the influence of local events.
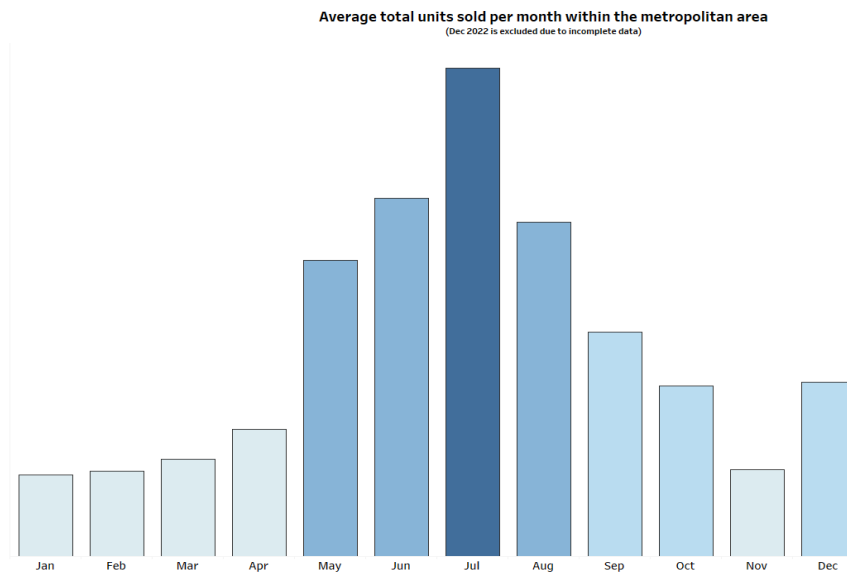
**Visualisations:**

- Colour was used to highlight contrasts (e.g., sales vs. temperature).
- Labels aided quick insights; axes/grid lines were removed if unhelpful.
- Maps were used to show geospatial aspects, with colour and size indicating sales volume. Filters allowed detailed data exploration, and a video illustrated how locations disappeared as sales grew.

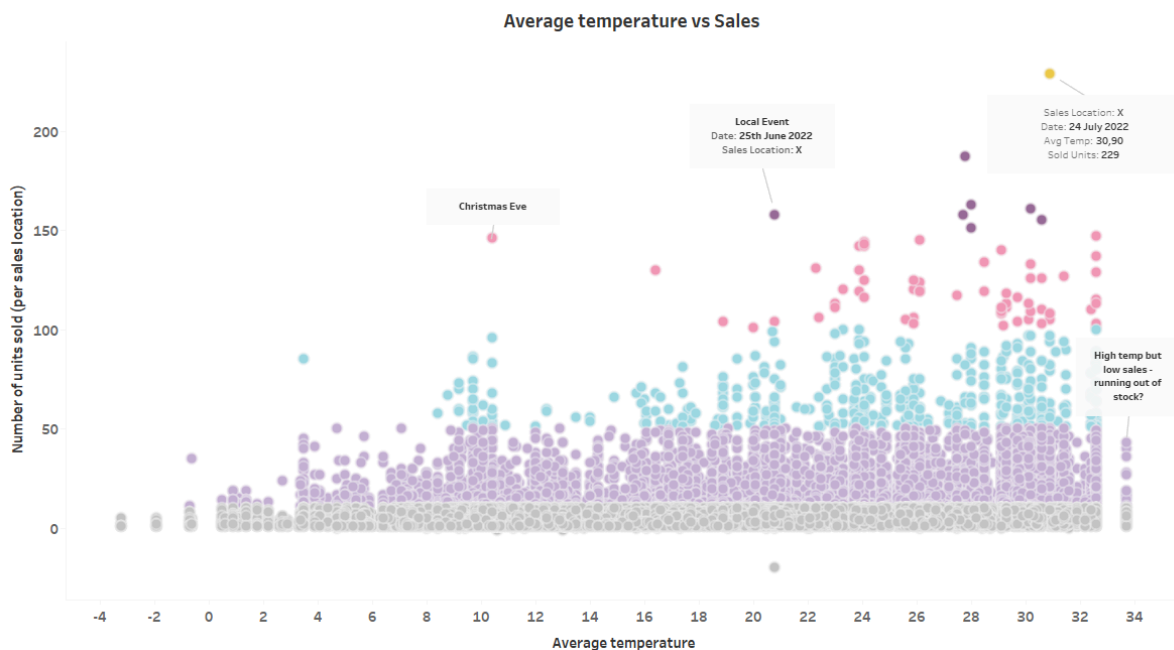Sales Sensitivity to Temperature Trends:

1. **Seasonal temperature patterns and peak sales**:
   A monthly bar chart shows sales peaking in July, confirming summer heat drives strong seasonal demand.

**Average total units sold per month within the metropolitan area**
(Dec 2022 is excluded due to incomplete data)



2. **Temperature thresholds and sales density**:
Sales rise above 28°C, suggesting a demand tipping point, but drop beyond 34°C - likely due to stockouts. High heat does not guarantee high sales, as many locations still show low volumes.

**Average temperature vs Sales**



3. **Temporal interaction and potential stockouts:**
The chart shows sales rising with temperature, but two sharp drops despite continued heat suggest missed opportunities or stockouts. A smaller cold-weather spike indicates temperature is not the only sales driver.

4. **Multivariate Relationships Driving Sales:**
Correlation analysis links sales to temperature, while Random Forest's feature importance highlights timing and location as key drivers, showing sales depend on more than just weather.

What Drives Daily Sales Most? (Including Day of Week)

Sales, location, and event-driven demand:

1. **Impact of location on sales**: A map was created in Tableau to display location of sales locations and volume of sales at each site (using size and colour density). There is a higher proportion of sales in the North, suggesting that it should be a higher priority.

(Maps were removed to ensure the anonymity of our client)

2. **Impact of events and holidays on sales**: Holiday-only days showed the highest average sales, with a 27% uplift over non-event days. A bar chart with error bars was used to clearly compare group averages while visualising variability across categories.



Average daily sales by event category

● Some specific events have a significant impact on sales.



Weekly number of sold units and average weekly maximum temperature

3. **Calendar-based sales trends**: A day-by-month heatmap revealed peak sales during summer and major holidays, like Christmas and New Year, reflecting alignment with leisure-driven consumer behaviours. The heatmap enabled visualisation of temporal trends across the calendar year.
   ● Calendar-style heatmaps for 2022 and 2021 (see Appendix 3) confirmed seasonal and weekly consistency but were excluded from the presentation as it did not highlight individual holiday impact - the focus of our heatmap analysis.
   ● Day-of-week chart reveals that ~60% of sales are from Fridays to Sundays, with ⅓ from Saturdays (Appendix 3).

Total Daily Sales by Day and Month (Including Estimated Data for 15–31 Dec 2022)

# Modeling

**Can historical patterns predict future demand for better stock planning?**

We tested a number of progressively more powerful and sophisticated machine learning models, each addressing the limitations of the previous. We selected the Histogram-based Gradient Boosting (HGB) paired with a Monte Carlo simulation as our final forecasting approach.

**Data Preparation**

To ensure optimal conditions for machine learning, preparation included:
- Insert row per date for each location where gaps, with assumed 0 sales (26,000 rows)
- Applying cyclical encoding to preserve continuing of temporal variables
- Encoding features with rolling mean to reduce noise and highlight seasonality
- Creating lag variables to capture temporal dependencies
- Retrieving temperature normals from Meteostat and filling gaps using KNNImputer

Please refer to Appendix 4 for the findings from the Decision Tree, Random Forest and the Extreme Gradient Boosting model (XGBoost).

**HGBoost**

The final HistGradientBoosting Model showed clear improvement after aggregating sales data by town/city, reducing noise from sales locations with frequent zero-sales days. Across cross-validation folds, it achieved an average R² of 0.66, RMSE of

16.05, and MAE of 6.38, indicating stronger predictive accuracy than previous versions. Overfitting persists, as evidenced by a training R² of 0.89 versus a test R² of 0.62 in the final fold, likely stemming from the limited data range (2021–2022).

Key features included town/city and day-of-week patterns, while weather variables had minimal impact. Binary flags for holidays and recurring weekly/monthly events helped capture predictable demand shifts. Residuals were generally well distributed with a slight leftward skew (see image below), suggesting minor underprediction but overall improved generalization and reduced bias compared to prior models.

While residuals indicate a generally good fit, the model's high Mean Absolute Percentage Error (MAPE) of 169.66% reflects considerable variability or scale issues in the target data. This limitation could be mitigated by adding more historical data to improve generalization and reduce overfitting.



**Monte Carlo simulation**

To project demand under future weather scenarios, we used Monte Carlo to simulate 1000 weather outcomes based on 30 years of historical Meteostat data—balancing computational efficiency with prediction reliability. Applying our Boost model to these simulations, we generated a 2023 demand forecast with 80% and 95% confidence intervals, capturing likely variability in sales:

2023 Forecast with Confidence Bands



Monthly Sales: 2021–2022 Actuals & 2023 Predictions

## Revenue Optimisation

Building on the model forecasts, we aimed to optimise revenue generation by applying a further layer of logic to dynamically adjust the sales strategy, accounting for event-driven demand and weather forecast data. This allows us to also compute corresponding revenue and expected profit.

We assumed a base set of product parameters based on typical industry assumptions.

The key steps of optimisation include:

1. Demand Estimation: Uses the 95% confidence interval for peak event days otherwise the 80% confidence interval is used to estimate total units sold.

2. Heat Signal Calculation: Combines the probability of heat events with a boost based on expected temperature. Applies a dampening factor of 30% for highly uncertain temperature forecasts which deviate from the typical temperature ($z\_score > 3$).

3. Product Ratio Adjustment: Gradually increases the premium product ratio above a defined probability and temperature threshold (capped at 10% of all units).

4. Revenue & Profit Estimation: Computes expected revenue and profit based on adjusted product mix and pricing/margin assumptions.

The final optimised forecast was visualised in a dashboard via Tableau to demonstrate the effect of events and temperature variables on the final forecast via the implementation of filters.

**Forecasted Units for 2023**

2023

| Q1 | Q2 | Q3 | Q4 |

Total Units

January · February · March · April · May · June · July · August · September · October · November · December

# Strategic Recommendations

**Forecasting:**
- Use forecasting model to plan production volumes for peak times (weekends, summer, Christmas, New Year and major events)
  - Use the upper 80% confidence band as standard, increase to 95% for peak events.
- Use revenue optimiser to adjust and plan distribution of premium products
  - Recommend trial at 10% cap for first year.

**Geospatial:**
- Align distribution schedules with weekly and seasonal cycles ensuring stock for peak periods such as weekends and events
- Prioritise stock and delivery to sales locations in the north which have nearly 75% of all sales
- Prioritise delivery to the locations with sales on over 75% of days
- Stop distributing to the locations with sales on less than 25% of days

**Technical Recommendations:**
- Add additional historical sales data to the model to improve the predictive power and accuracy
- Differentiate between premium and standard products in data collection
- Include stock-flow data to test assumptions on missed sales opportunities

# Appendix 1 - Geographic Filtering

We used postcode-based radius filtering. Minor margin-of-error acknowledged due to measurement method.

# Appendix 2 - Outlier Analysis

- Precipitation: Max value identified as outlier but retained as plausible during storm events.
- Wind Speed & Pressure: High-end values flagged as outliers but consistent with known weather extremes for the metropolitan area.
- Temperature: All deemed realistic.

# Appendix 3 - Other Visualisations

Calendar-style heatmaps:

**January**

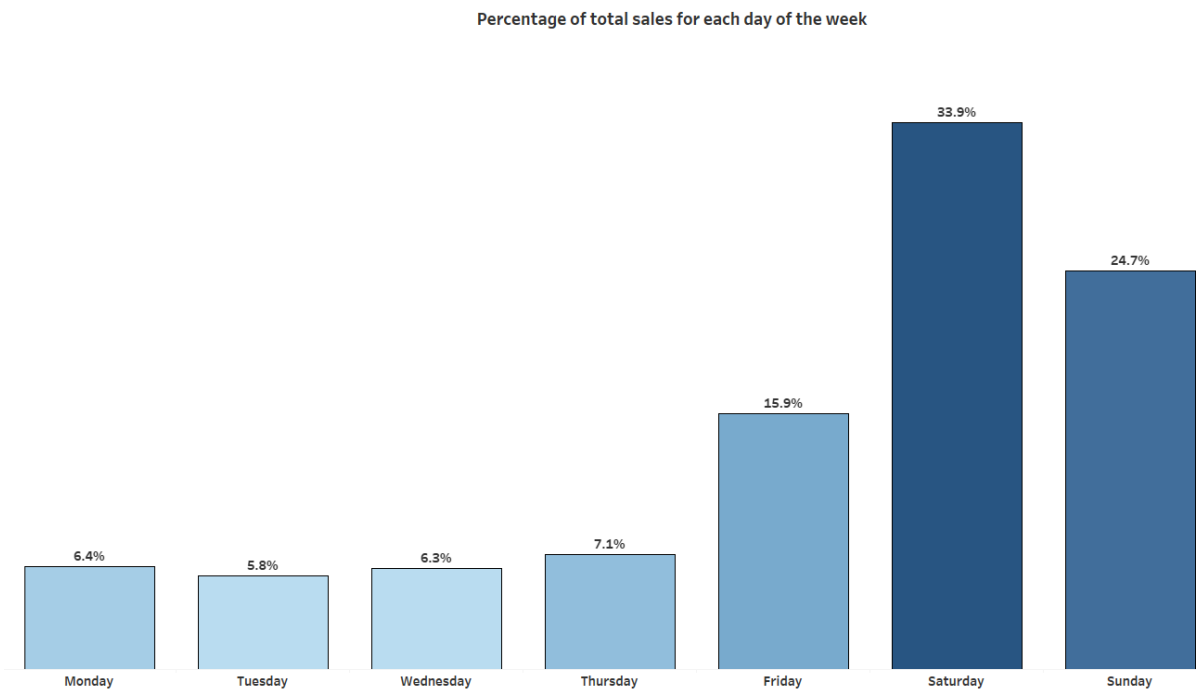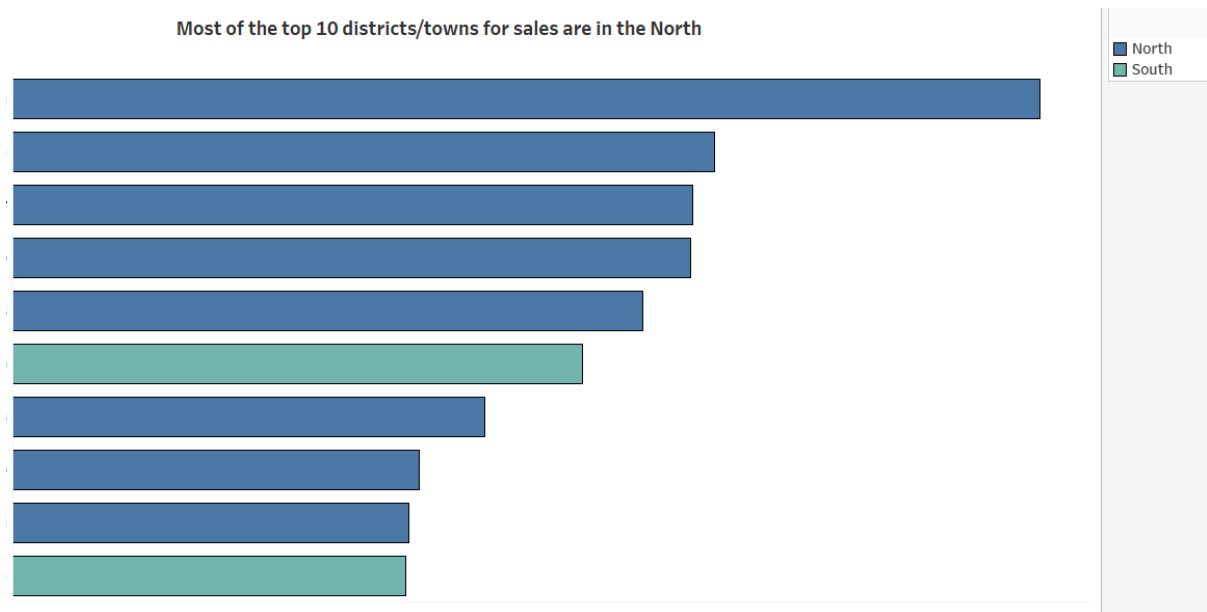| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 54 | | | | | | |
| 43 | 22 | 56 | 73 | 134 | 376 | 261 |
| 30 | 37 | 41 | 68 | 176 | 364 | 171 |
| 42 | 63 | 74 | 54 | 117 | 342 | 211 |
| 76 | 59 | 147 | 67 | 226 | 531 | 332 |
| | 57 | 38 | | | 1224 | 199 |

**February**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 85 | | | | | | |
| 53 | 56 | 50 | 80 | 216 | 323 | 339 |
| 66 | 92 | 56 | 52 | 208 | 458 | 368 |
| 48 | 93 | 54 | 59 | 186 | 457 | 215 |
| | 53 | 53 | 49 | 186 | 479 | 289 |

**March**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 76 | 36 | 52 | 59 | | | |
| 23 | 36 | 37 | 48 | 121 | 407 | 289 |
| 26 | 50 | 60 | 58 | 192 | 464 | 241 |
| 36 | 59 | 127 | 84 | 154 | 478 | 210 |
| | 53 | 53 | 49 | 214 | 400 | 250 |

**April**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 60 | 88 | 67 | 89 | 247 | 877 | |
| 146 | 38 | 52 | 96 | 171 | 407 | 379 |
| 76 | 59 | 177 | 426 | 735 | 512 | 574 |
| 64 | 70 | 68 | 101 | 141 | 436 | 386 |
| | | | | 214 | 400 | 1301 |

**May**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 196 | 215 | | | | 877 | |
| 127 | 117 | 176 | 201 | 501 | 2375 | 1302 |
| 416 | 173 | 203 | 287 | 724 | 2322 | 1192 |
| 84 | 116 | 136 | 141 | 510 | 1145 | 1268 |
| 297 | 50 | 128 | 250 | 560 | 1501 | 1089 |
| | | 256 | | | | 1301 |

**June**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 185 | 207 | 310 | 343 | | | |
| 272 | 264 | 201 | 237 | 583 | 1388 | 900 |
| 333 | 490 | 547 | 694 | 1313 | 2795 | 1723 |
| 269 | 201 | 217 | 280 | 793 | 2368 | 1830 |
| | | 256 | 250 | 755 | 1556 | 1174 |

**July**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 1574 | 724 | 517 | 584 | 925 | 2466 | 1634 |
| 654 | 619 | 594 | 547 | 987 | 2376 | 2952 |
| 304 | 450 | 543 | 648 | 991 | 2511 | 2367 |
| 281 | 281 | 327 | 303 | 643 | 1635 | 1390 |
| | | | | 755 | 1556 | 1174 |

**August**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 294 | 302 | 386 | | | | |
| 284 | 309 | 315 | 445 | 459 | 899 | 814 |
| 503 | 224 | 246 | 239 | 457 | 892 | 764 |
| 462 | 534 | 615 | 690 | 948 | 806 | 911 |
| 757 | 776 | 793 | 364 | 701 | 1126 | 1061 |

**September**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 96 | 89 | 100 | 96 | 240 | | |
| 107 | 134 | 108 | 170 | 404 | 804 | 353 |
| 185 | 130 | 101 | 108 | 399 | 1006 | 605 |
| 242 | 208 | 180 | 204 | 604 | 1669 | 1175 |
| | | | 364 | 701 | 725 | 594 |

**October**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 287 | | | | | | |
| 59 | 55 | 76 | 98 | 200 | 626 | 416 |
| 61 | 112 | 96 | 83 | 216 | 433 | 274 |
| 100 | 191 | 521 | 138 | 301 | 648 | 446 |
| 139 | 90 | 109 | 89 | 248 | 367 | 282 |
| | 372 | 62 | | | 725 | 594 |

**November**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 25 | 56 | 45 | | | | |
| 78 | 99 | 49 | 43 | 200 | 346 | 311 |
| 52 | 38 | 46 | 52 | 193 | 401 | 191 |
| 40 | 102 | 125 | 69 | 179 | 343 | 289 |
| | 372 | 62 | 68 | 132 | 256 | 243 |

**December**

| Mon | Tue | Wed | Thu | Fri | Sat | Sun |
| --- | --- | --- | --- | --- | --- | --- |
| 0 | 0 | 0 | 0 | 0 | 0 | |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 44 | 52 | 14 | 0 | 0 | 0 | 0 |
| 151 | 158 | 119 | 195 | 104 | 290 | 205 |
| | | | 68 | 132 | 256 | 243 |

Day-of-week chart:

**Percentage of total sales for each day of the week**

| Day | Percentage |
|-----|-----------|
| Monday | 6.4% |
| Tuesday | 5.8% |
| Wednesday | 6.3% |
| Thursday | 7.1% |
| Friday | 15.9% |
| Saturday | 33.9% |
| Sunday | 24.7% |

Percentage of days that have sales map:

(Maps were removed to ensure the anonymity of our client)

Number of total sales in the top 10 districts in the metropolitan area or the towns surrounding it:

Most of the top 10 districts/towns for sales are in the North

North
South

## Appendix 4 - Models tested

## Decision Tree

We initially tested a **Decision Tree model** to quickly identify non-linear relationships between non-linear variables, and to determine which are most predictive. The model utilised all provided weather data, location variables, dates plus external data gathered on public holidays and national and major sporting events with regional relevance.

Timeseriessplit using scikit was used, with 5 splits.  To attempt to improve the model various iterations were tested:

Results of first model - no adjustments:

```
Average MSE no lag: 81.18123243451473
Average MAE no lag: 4.555385906225992
Average RMSE no lag: 8.533054503555494
Average R-squared no lag: 0.10794599750411693
```

Results when lag and rolling features were added for sold units over 1, 2, 7, 14 days:

```
Average MSE with lag: 71.21402085142127
Average MAE with lag: 3.9077686583260274
Average RMSE with lag: 7.814870549053825
Average R-squared with lag: 0.2858615095203615
```

Results when model pruned to dept of 5, with only 3 splits:

```
Average MSE pruned: 57.1978521742779
Average MAE pruned: 3.8637096928360903
Average RMSE pruned: 7.511154019007434
Average R-squared pruned: 0.39742565464828
```

Whilst the accuracy and prediction capability improved with each adjustment, the performance is not strong enough to use for the complex forecasting needed in this project and a decision was quickly taken to test with a Random Forest model.

## Random Forest

In preparing the data for the Random Forest Model various engineered features were added, including: lag-features for climate and sold_units, target encoding location date (sales location and company codes and town/city).

The performance was tested with various splits, estimators, depths and leafs and a final model run with the best hyperparameters.

Final results:

```
Average Train Metrics Across All Folds:
R-squared: 0.81
Mean Squared Error: 12.45
Root Mean Squared Error: 3.48
Mean Absolute Error: 1.65
Mean Absolute Percentage Error: 63.37%
RMSE - MAE: 1.83
RMSE - MAE (percentage): 52.69%

Average Test Metrics Across All Folds:
R-squared: 0.42
Mean Squared Error: 45.10
Root Mean Squared Error: 6.32
Mean Absolute Error: 3.36
Mean Absolute Percentage Error: 127.19%
RMSE - MAE: 2.96
RMSE - MAE (percentage): 47.30%
```

From investigations of this model we can see that the generalisation is poor. There were also additional data leakage issues, which we would later address in our boosted models. Dropping additional features and undertaking additional feature engineering could improve this model, but as a random forest overall performance is slow and time consuming. We decided that this is not an appropriate fit for a production model which will need to run in a reasonable amount of time and we progressed with further exploration using a boosted tree model.

# XGBoost

We explored **XGBoost** due to its strong performance on tabular data and its ability to model complex, non-linear relationships. Using a comprehensive feature set of 35 engineered variables—including temporal, location-based, weather, and event-related features—XGBoost achieved the following average performance across three time-based folds:

- **RMSE**: 7.02
- **MAE**: 3.49
- **R²**: 0.47
- **MAPE**: 134.15%

Despite extensive tuning—testing different hyperparameters, feature subsets, and feature eliminations—performance did not improve meaningfully. Analysis of the residuals (see image below) revealed clear patterns, indicating that the model was struggling to capture important dynamics in the data. In addition, signs of **overfitting** suggested limited generalization to unseen time periods.

This led to a shift in approach: we **restarted from scratch** with a smaller, more focused feature set and found that **Histogram-Based Gradient Boosting (HGB)** consistently outperformed XGBoost on this reduced input, offering better generalization and training efficiency.