# BI2025 Experiment Report - Group 16

Muhammad Ibrar[*]
TU Wien
Austria

Ahmad Ibrahim[†]
TU Wien
Austria

## Abstract

This report documents a machine learning experiment conducted by Group 16 following the CRISP-DM methodology. The goal is to predict the price range of mobile phones based on their technical specifications using supervised learning techniques. The experiment focuses on data understanding, preparation, modeling, evaluation, and deployment considerations, while emphasizing transparency and reproducibility.

## CCS Concepts

• **Computing methodologies → Machine learning**.

## Keywords

CRISP-DM, Machine Learning, Random Forest, Mobile Price Prediction

## 1 Business Understanding

### 1.1 Business Objectives

The primary business objectives of this project are as follows:

- Support pricing decisions by predicting the most suitable price range for new phone models based on technical specifications.
- Reduce manual effort and time required for estimating price categories during early product planning stages.
- Improve product positioning across budget, mid-range, high-end, and flagship market segments.
- Increase transparency and understanding of how technical features influence pricing decisions.

### 1.2 Business Success Criteria

The success of the project from a business perspective is measured using the following criteria:

- The machine learning system is regularly used by product management and pricing teams.

[*]Student A, Matr.Nr.: 12350094
[†]Student B, Matr.Nr.: 11826186

- A reduction of at least 30% in the time required for initial price-range estimation.
- More than 80% of newly launched smartphone models remain within a ±10% deviation from their initially selected price band after market entry.
- Pricing decisions become more consistent and data-driven across all phone segments.

### 1.3 Data Mining Goals

The data mining goals translate the business objectives into concrete analytical tasks:

- Build a multi-class classification model that predicts the `price_range` (0–3) from 20 phone features.
- Achieve robust predictive performance on unseen data and generalize well to new device configurations.
- Identify and analyze the most influential features (e.g., RAM, pixel resolution, battery power) affecting price range.
- Provide probabilistic outputs to support uncertainty-aware pricing decisions.

### 1.4 Data Mining Success Criteria

The technical success of the data mining task is evaluated using the following criteria:

- Achieve at least 90% classification accuracy on the validation or test dataset.
- Reach a macro F1-score of at least 0.88, with no individual class having recall below 0.80.
- Demonstrate stable model performance across different random train-test splits.
- Produce reasonably calibrated probability estimates suitable for business decision-making.

### 1.5 AI Risk Aspects

Several potential risks related to the use of machine learning in pricing decisions were identified:

- Misclassification may lead to incorrect pricing decisions, negatively affecting revenue or sales volume.
- The dataset may not represent future device trends, introducing the risk of model drift over time.
- Over-reliance on automated predictions could result in poor decisions if expert judgment is excluded.
- Although no personal data is involved, systematic bias across device categories may still occur.
- Exposure of the pricing logic could pose a business security risk if the model is deployed externally.

## 2 Data Understanding

**Dataset Description:** The dataset consists of 2,000 instances with balanced class distribution across four price ranges. All features are

numerical, including binary indicators and continuous measurements.

**Table 1: Raw Data Features**

| Feature Name | Data Type | Description |
|---|---|---|
| battery_power | integer | Battery capacity (mAh). |
| blue | integer | Bluetooth support (0/1). |
| clock_speed | double | Processor clock speed (GHz). |
| dual_sim | integer | Dual SIM support (0/1). |
| fc | integer | Front camera (MP). |
| four_g | integer | 4G support (0/1). |
| int_memory | integer | Internal memory (GB). |
| m_dep | double | Phone thickness (cm). |
| mobile_wt | integer | Weight (grams). |
| n_cores | integer | Number of processor cores. |
| pc | integer | Primary camera (MP). |
| price_range | integer | Target class (0–3). |
| px_height | integer | Screen pixel height. |
| px_width | integer | Screen pixel width. |
| ram | integer | RAM (MB). |
| sc_h | integer | Screen height (cm). |
| sc_w | integer | Screen width (cm). |
| talk_time | integer | Battery talk time (hours). |
| three_g | integer | 3G support (0/1). |
| touch_screen | integer | Touchscreen support (0/1). |
| wifi | integer | WiFi support (0/1). |

## 2.1 Categorical Feature Distributions

Categorical variables, primarily binary indicators (e.g., `blue`, `dual_sim`, `four_g`, `three_g`, `wifi`, `touch_screen`), were analyzed using pie-chart plots (Figure 1).

The target variable `price_range` is perfectly balanced, with exactly 500 samples per class, which is advantageous for multiclass classification. Most binary features show near-even splits between 0 and 1. An exception is `three_g`, where approximately 76% of devices support 3G. No unexpected or degenerate category distributions were observed.

Histograms with kernel density estimates were generated for all numerical attributes. Most numerical features exhibit unimodal distributions with realistic value ranges. However, some variables such as `px_height` and `sc_w` show a notable concentration of values near zero, suggesting potential noise or atypical records.

Skewness was computed for all numerical variables to assess distribution asymmetry. While several features exhibit mild skewness, no extreme distortions were observed that would immediately invalidate modeling. These findings inform later decisions regarding scaling or transformation during data preparation.

## 2.2 Outlier Detection

Outlier detection was performed using a z-score threshold of 3.0 across all numerical attributes. The analysis revealed that nearly
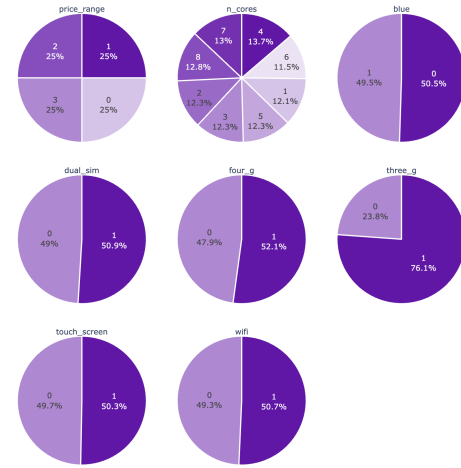


**Figure 1: Distribution of categorical features and target variable `price_range`.**



**Figure 2: Distributions of numerical features with histograms and kernel density estimates.**

all features contain no statistically significant outliers. The only exception is the `fc` (front camera megapixels) attribute, for which 12 observations were flagged as potential outliers.

Overall, the dataset appears largely clean, with only a small number of extreme values that may require attention in the data preparation phase.

## 2.3 Correlation Analysis

A correlation heatmap was computed for all numerical variables, including the target `price_range`. Most feature pairs show weak linear correlations, indicating low multicollinearity and suggesting that features provide largely complementary information.
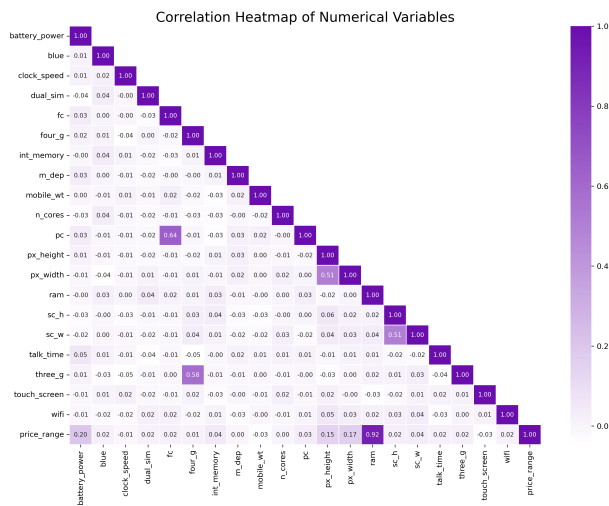


**Figure 3: Correlation heatmap of numerical features, including the target variable `price_range`.**

Moderate positive correlations were observed between related attributes, such as screen height and width, pixel height and width, and front versus primary camera resolution. Notably, a strong positive correlation of 0.92 was observed between `ram` and the target variable `price_range`, indicating that memory capacity is a key determinant of price category. Correlations between most other individual features and the target variable are moderate, supporting the suitability of multivariate modeling approaches.

## 2.4 Data Quality Assessment

A systematic data quality assessment showed that the dataset contains no missing values in any of the 21 columns and no duplicate rows. All data types are consistent with the intended semantics of the attributes.

A potential issue is that `px_height` contains 2 zero values and `sc_w` contains 180 zero values. Zero screen height or width is implausible for real mobile devices and may indicate noisy, special, or incorrectly recorded cases. These findings should be considered in the data preparation phase (e.g., deciding whether to filter, impute, or keep these records).

## 2.5 Ethical and Bias Aspects

The dataset contains no personal or sensitive attributes; all features describe technical phone specifications only (e.g., battery capacity, memory, screen resolution, connectivity options). The target variable `price_range` is perfectly balanced, with 500 instances per class, which reduces the risk of systematic bias in terms of underrepresented target categories.

Some device types and feature combinations are less frequent (for example, rare combinations of extreme screen or camera configurations). However, these imbalances primarily affect model performance and generalization rather than human fairness, because they do not correspond to demographic or protected population groups. Based on the available information, no ethical issues related to discrimination or sensitive attributes are apparent, and no specific demographic bias risks are present in this dataset.

## 2.6 Risks and Expert Questions

Despite the absence of personal data, several data- and model-related risks were identified. First, the representativeness of the dataset with respect to the real smartphone market is uncertain. Some feature combinations—such as very high camera megapixel values or screen dimensions close to zero—may not reflect realistic devices and could introduce noise or distort model behaviour. Second, the dataset may not include newer technologies or recent device trends, which increases the risk of model drift if the model is applied to future products.

These uncertainties motivate consultation with a domain expert. In particular, the following questions should be clarified:

- Are zero values for `px_height` or `sc_w` technically valid (e.g., encoding a special case), or are they more likely to be measurement or recording artifacts?
- Are unusually high values in `fc` (front camera megapixels) realistic device specifications, or should they be treated as outliers?
- Does the dataset represent a realistic mix of budget, midrange, and high-end devices, or are certain market segments over- or underrepresented?
- Are any important device characteristics missing that strongly influence real-world pricing (e.g., brand, release year, or build quality)?

Answers to these questions would help to better assess potential data quality issues, refine preprocessing decisions, and estimate the robustness of the model in practical deployment scenarios.

## 2.7 Actions Required for Data Preparation

Based on the data understanding analysis, several concrete actions are recommended for the data preparation phase:

- Investigate and handle zero values in `px_height` (2 cases) and `sc_w` (180 cases), as such values are unlikely for real devices. Depending on domain expert feedback, these records may be filtered out or adjusted via imputation.
- Re-assess the 12 outliers detected in `fc` (front camera megapixels). Depending on whether they are confirmed as valid device specifications or artifacts, decide whether to retain, cap, or remove them.

- Standardize or scale numerical features (e.g., `ram`, `battery_power`, pixel resolution attributes) to account for differing value ranges and to support algorithms that are sensitive to feature scaling.
- Ensure that categorical binary features (e.g., `blue`, `dual_sim`, `four_g`, `three_g`, `wifi`, `touch_screen`) are stored in consistent numeric types and encodings. No additional encoding is required, as they are already represented as 0/1 indicators.
- Review potential skewness in selected numerical variables and consider applying transformations (such as log or power transforms) for algorithms that assume more symmetric or approximately normal feature distributions.
- Ensure proper train–validation–test splitting to maintain the balanced distribution of price_range classes.

These actions provide a structured basis for subsequent preprocessing and help to increase the robustness, interpretability, and reliability of the resulting machine learning models.

## 3 Data Preparation

### 3.1 Initial Preprocessing Actions

Basic preprocessing checks were performed based on the Data Understanding phase. No duplicate rows were detected, and no missing values were found across any of the 21 attributes.

Potential noise-like values were identified for selected features. Screen width values below 2 cm (`sc_w < 2`) occurred in 390 records, and pixel height values below 5 pixels (`px_height < 5`) occurred in 9 records. These records were tagged but not removed, as their domain validity is unclear and they may represent atypical devices. Their handling was deferred to the modeling phase.

### 3.2 Preprocessing Steps Considered but Not Applied

During data preparation, several preprocessing steps were considered but not applied at this stage. The rationale for each is summarized below:

- **Outlier removal:** Outliers in the front camera feature (`fc`) were retained, as they may represent legitimate device variations. Their effect will be evaluated during modeling, where extreme values can be handled if necessary.
- **Noise cleaning:** Values such as `sc_w < 2` and `px_height < 5` were kept. These may correspond to early or atypical devices. Noise handling will be revisited if model performance is affected.
- **Feature removal:** Some features showed weak linear correlation with price_range, but correlation alone is insufficient for elimination. Features were retained, and model-based selection or regularization will guide removal if needed.
- **Scaling and normalization:** Standardization and normalization were considered but deferred. They will be applied within model-specific pipelines if required, particularly for algorithms sensitive to feature magnitudes (e.g., SVM, KNN).
- **Encoding:** Binary features (`blue`, `dual_sim`, `three_g`, `four_g`, `touch_screen`, `wifi`) are already numeric (0/1). One-hot or additional encoding was deemed unnecessary.

- **Zero or missing-like values:** Zero entries in `px_height` and `sc_w` were retained for now, as they may reflect unusual but valid devices. Adjustments will depend on domain feedback or model sensitivity.

These decisions preserve the integrity of the dataset while leaving flexibility for model-specific preprocessing and evaluation.

### 3.3 Derived Attributes and External Data

During the Data Preparation phase, the potential creation of derived attributes and the inclusion of external data were considered as ways to enhance model performance and interpretability.

### 3.4 Derived Attributes

Several derived features could potentially improve model performance or interpretability:

- *Screen-related features:*
  - `pixel_area = px_height * px_width` (proxy for screen resolution)
  - `screen_area = sc_h * sc_w` (approximate physical display size)
  - `pixel_density_ratio = pixel_area / screen_area` (requires reliable screen dimensions)
- *Performance and capacity ratios:*
  - `ram_per_internal_memory = ram / int_memory` (relative memory configuration)
  - `battery_per_weight = battery_power / mobile_wt` (capacity relative to device weight)
  - `camera_total_mp = fc + pc` (overall camera capability)
- *Connectivity and feature counts:*
  - `connectivity_score = blue + three_g + four_g + wifi` (simple connectivity count)
  - `feature_richness = connectivity_score + touch_screen + dual_sim` (overall feature richness)

These derived attributes could help models distinguish between devices within the same price_range. Their creation will be decided based on modeling needs and complexity trade-offs.

### 3.5 External Data Sources

In addition to derived attributes, external data sources could provide further context and improve alignment with real-world pricing:

- *Real retail prices:* Linking devices to historical market prices from online shops or price comparison portals would allow training regression models for actual prices and validating the price_range labels.
- *Brand and model metadata:* Adding manufacturer, model family, and release year could capture brand and generation effects that influence perceived value.
- *Market segment and region:* Including information about target segment (budget, mid-range, flagship) or region (EU, US, Asia) enables finer-grained analysis of pricing expectations.
- *User or expert ratings:* Aggregated review scores (camera, battery, display) can help connect technical specifications to perceived quality and justify differences within the same price_range.

While these external sources are not integrated in this project, documenting them clarifies potential extensions for richer pricing and marketing analysis.

## 3.6 Summary

Data preparation focused primarily on validation rather than transformation. The dataset was found to be complete, clean, and well-structured, and all identified noise or outliers were retained. More complex preprocessing steps were deliberately deferred to later modeling stages.

Additional decisions include:

- **Scaling:** Not applied globally; only model-specific pipelines (e.g., SVM) will scale features as needed.
- **Binning:** Not performed, as numerical attributes have meaningful continuous ranges and the `price_range` target is balanced.
- **Outlier removal:** Front camera outliers (`fc`) were retained due to uncertain domain validity; removal will be considered only if models show sensitivity.
- **Noise values:** Noise-like entries in `sc_w` and `px_height` were kept, as they may represent early-generation devices; planned models are robust to such noise.

- **Encoding:** No additional encoding required; categorical features are already numeric (0/1 or ordinal).
- **Feature removal:** Features were not removed solely based on weak linear correlation; non-linear effects will be captured via model-based importance.

In conclusion, only essential checks—duplicates, missing values, and noise tagging—were performed. All other transformations were deferred to the modeling phase or deemed unnecessary given the dataset's clean and structured nature.

## 4 Modeling

...

## 5 Evaluation

...

## 6 Deployment

...

## 7 Conclusion

...