

# ML Project Report: Predicting Startup Founder Retention

Nipun Verma (IMT2023591)

Sahil Singh (IMT2023521)

Mohammed Ibrahim (IMT2023112)

[GitHub Repository](#)

## Background and Problem Statement

Startup ecosystems are built around founders. A founder who stays with the company provides continuity of vision, better investor confidence and more stable team dynamics. A founder who exits early can trigger a cascade of problems: stalled roadmaps, new leadership struggles and renegotiated funding terms.

The dataset used in this project comes from the Kaggle competition “*Startup Founder Retention Prediction*”. The task is to predict whether a founder will **Stay** with the venture or **Leave**, based on demographic, financial and organizational information.

Unlike the earlier transport–cost regression task (where the target was a continuous cost value and the modelling focus was linear vs tree vs boosting regressors, this project is a **binary classification** problem. The primary goal is not only to get a high score but also to understand which factors matter for retention and why certain model families behave the way they do on this kind of structured data.

Three model families were used:

1. A **Logistic Regression** pipeline as a fast, interpretable baseline.
2. An **SVM with RBF kernel** as a strong non–linear classical model.
3. A **Neural Network (Tabular MLP)** as a flexible high–capacity model.

For fairness, all three models share a common preprocessing philosophy, but they deliberately use slightly different engineering details to suit their strengths. The report mirrors

the structure of the earlier medical transport project: starting from data description and cleaning, moving through feature engineering and finally giving an in-depth analysis of each model and the overall comparison.

## Dataset Description

The training file contains 59,611 rows and 24 columns. The test file contains 14,900 rows. The target variable is:

- `retention_status` ∈ {Stayed, Left}

The target distribution in the training set is:

`Stayed` : 31,265,   `Left` : 28,346

which is reasonably balanced but not perfectly symmetric. This mild imbalance motivates the use of **macro F1** rather than accuracy.

The remaining features can be grouped naturally into numeric, nominal and ordinal types:

**Numeric features** `founder_age`, `years_with_startup`, `monthly_revenue_generated`, `funding_rounds`, `distance_from_investor_hub`, `num_dependents`, `years_since_founding`.

**Nominal features** `founder_gender`, `founder_role`, `working_overtime`, `personal_status`, `remote_operations`, `leadership_scope`, `innovation_support`.

**Ordinal features** These started life as strings like “Excellent”, “Poor”, “Very High” and had to be carefully mapped to integers:

`work_life_balance_rating`, `venture_satisfaction`, `startup_performance_rating`, `startup_stage`, `team_size_category`, `startup_reputation`, `founder_visibility`, `education_background`.

A recurring theme in this project is that most modelling difficulties come not from the algorithms, but from making these heterogeneous columns numerically meaningful without destroying information.

# EDA and Preprocessing

## Handling Missing Values and Data Types

The raw CSV contained missing values in both numeric and categorical columns, and multiple columns that looked numeric but were actually stored as strings. The cleaning strategy was:

- For **numeric and ordinal** columns, missing values were imputed using the median of each column. Median is more robust than the mean in the presence of skewed distributions, which is important because revenue and distance variables have long right tails.
- For **nominal** columns, missing values were filled with the most frequent category (mode). This is defensible because these features represent membership in discrete groups, and mode imputation avoids inventing artificial categories.

For all models built directly with scikit-learn pipelines (the Logistic Regression and SVM pipelines), imputation was handled via `SimpleImputer`. For the neural-network-oriented code, imputation was performed manually during preprocessing but followed the same median/mode philosophy.

The final sanity check after preprocessing printed the remaining number of NaNs in train and test; both were zero.

## Boolean Normalization

Columns such as `working_overtime`, `remote_operations`, and `leadership_scope` contained values like “Yes”, “No”, “True”, “False” and sometimes even numeric flags like “1” and “0”.

Instead of treating them as generic nominal categories, they were normalised to integer indicators:

$$\text{Yes, True, } 1 \rightarrow 1, \quad \text{No, False, } 0 \rightarrow 0$$

This choice has two benefits. First, it prevents a simple binary decision from exploding into multiple one-hot columns. Second, it conveys the natural ordering (0 vs 1) directly to models like Logistic Regression and SVM that expect numeric features.

## Ordinal Encoding

The more interesting part of preprocessing was the manual mapping of ordinal categories. For example:

`work_life_balance`: Poor → 0, Fair → 1, Good → 2, Excellent → 3.

Similarly:

`venture_satisfaction`: Low = 0, Medium = 1, High = 2, Very High = 3.

`startup_stage`: Entry = 0, Mid = 1, Senior = 2.

`education_background`: High School = 0, Associate Degree = 1, Bachelor's Degree = 2, Master's Degree = 3.

The key design decision here was to treat these as genuinely ordered scales, not just categories. One-hot encoding them would have hidden the natural progression (for example, between “Low” and “Very High”) and would have created unnecessary dimensionality. Instead, the models can now use the numeric gap between levels as signal.

As a safety check, the preprocessing function also printed out any category values that did not appear in the mapping, so that unexpected labels (due to typos or hidden categories) could be identified early.

## Outlier Treatment via IQR Capping

During EDA, boxplots of numeric features, especially `monthly_revenue_generated`, revealed extreme outliers, with some entries far above the bulk of the data. Leaving these untouched would make scaling unstable and could bias models that rely on distances or linear combinations.

Rather than deleting outlier rows—which might remove genuinely informative edge cases—we applied IQR-based capping. For each numeric column:

$$Q_1 = \text{25th percentile}, \quad Q_3 = \text{75th percentile}, \quad \text{IQR} = Q_3 - Q_1,$$

$$\text{lower cap} = Q_1 - 1.5 \cdot \text{IQR}, \quad \text{upper cap} = Q_3 + 0.7 \cdot \text{IQR}.$$

Values outside this band are clipped. The asymmetric choice (1.5 below, 0.7 above) reflects the observation that extreme high values are more common than extreme low values, especially for revenue and distance. The cap is therefore more conservative on the upper side to avoid distorting genuine high-performing founders.

The motivation is similar to the earlier transport–cost project, where heavy-tailed cost variables motivated robust handling rather than aggressive trimming.

## Scaling and Encoding for Modelling

After imputation and encoding, different models needed slightly different final feature transformations:

- For the **Logistic Regression** and **SVM** pipelines written with scikit-learn’s `ColumnTransformer`, numeric features were log-transformed (for highly skewed variables like revenue and distance) and then standardised with `StandardScaler`. Categorical features were one-hot encoded.
- For the **Neural Network** code path, after one-hot encoding of nominal features, a `RobustScaler` was applied to all features at once. This matches the NN’s preference for roughly standardised, but outlier-robust, inputs.

A deliberate design choice was to keep the NN preprocessing relatively simple and uniform (one big numeric matrix), while the classical pipelines used a more explicit per-column treatment using transformers.

## Feature Importance and Mutual Information

To understand which features actually carry information about retention, we computed mutual information between each input feature and the encoded target.

The top-ranked features by mutual information were:

- `startup_stage`
- `personal_status`
- `remote_operations`
- `work_life_balance_rating`
- `startup_reputation`
- `funding_rounds_led`

Lower-ranked but still nonzero contributions came from revenue, dependents, age and the various leadership and innovation indicators.

The fact that `startup_stage` and `personal_status` appear at the top is intuitive: founders at senior stages and with stable personal circumstances are more likely to stay. Interestingly, raw financial numbers such as monthly revenue are not as dominant as one might expect, suggesting that retention is driven by structural and lifestyle factors at least as much as short-term income.

Even features with small individual mutual information were not dropped immediately. In the earlier regression project, low-contributing features were pruned aggressively using Lasso and then removed for linear and ensemble models. Here, the task is classification, and certain models (particularly NNs and RBF-SVMs) can benefit from interactions between many weak features. Therefore, manual feature removal was kept to a minimum.

## Modelling Strategy

### Train–Validation Strategy

All three models share the same basic data split:

Training set : 80% (47,688 rows),      Validation set : 20% (11,923 rows),

with **stratification** on the target label to preserve the Stayed/Left proportions in both splits.

For some models, an additional level of cross-validation was used:

- Logistic Regression: hyperparameters tuned with 3-fold cross-validation on the training split.
- SVM: performance estimated with 5-fold StratifiedKFold on the held-out validation set.

This “outer” validation set is used for the final fair comparison across models, mirroring the structure of the previous project where a fixed public test set was evaluated after tuning on the training folds.

## Evaluation Metrics

The primary metric throughout is the **macro F1 score**. Unlike accuracy, macro F1 gives equal weight to both classes, even though “Stayed” is slightly more frequent than “Left”.

This discourages trivial strategies that over-predict the majority class.

For the Logistic Regression and SVM pipelines we also monitored ROC–AUC across folds, mainly to check calibration and ranking performance, but the leaderboard score and main discussion focus on macro F1.

Confusion matrices are used to analyse which class each model struggles with, and whether misclassifications are symmetric.

## Model 1: Logistic Regression

### Motivation

Logistic Regression is the simplest baseline capable of handling thousands of one-hot encoded features while remaining interpretable. In the earlier regression project, linear models (Lasso, Ridge, ElasticNet) performed surprisingly well because the underlying relationship between features and cost was largely linear.

For founder retention, it is not clear a priori whether the decision boundary is linear. Running Logistic Regression first provides a sanity check: if it performs close to more complex models, the problem may be mostly linear; if it lags significantly, we need non-linear classifiers.

### Pipeline and Hyperparameters

The Logistic Regression model was implemented as a scikit-learn **Pipeline**:

1. **Preprocessing** with a `ColumnTransformer` that:

- imputes and scales numeric columns (with log transform for skewed variables),
- imputes and one-hot encodes categorical columns,
- passes through binary columns as numeric indicators.

2. A **LogisticRegression** classifier with:

- penalty: L2,
- solvers considered: `lbfgs`, `saga`,
- `class_weight`: `balanced`,
- `max_iter`: 500.

The regularisation strength  $C$  was tuned over  $\{0.01, 0.1, 1, 3, 10\}$  using 3-fold cross-validation with macro F1 as the scoring metric.

The best configuration was:

$$C = 1, \quad \text{solver} = \text{saga}.$$

## Validation Performance

On the held-out validation set of 11,923 founders, the tuned Logistic Regression achieved:

$$\text{Macro F1} = 0.7364, \quad \text{Accuracy} \approx 0.74.$$

The confusion matrix was:

$$\begin{bmatrix} 4199 & 1471 \\ 1669 & 4584 \end{bmatrix}$$

where rows correspond to true labels (`Left`, `Stayed`) and columns to predicted labels.

The model correctly identifies most of the “Stayed” founders, but it is slightly more inclined to predict “Stayed” than “Left”, as shown by the 1,471 false “Stayed” predictions for founders who actually left. This mild bias is expected given the near-balanced but slightly skewed target and the use of class weights rather than explicit threshold tuning.

## Model 2: Support Vector Machine (RBF Kernel)

### Motivation

SVMs with RBF kernels are classic strong performers on medium-sized tabular datasets. They are particularly appealing when:

- The non-linear boundary is important.
- The number of features is moderate after encoding.
- The dataset is not so large that training becomes impossible.

In this project, training almost 60k samples with an RBF kernel is computationally heavy but still feasible on Kaggle hardware. The expectation was that SVM would outperform Logistic Regression if the retention decision truly depends on non-linear combinations of features.

## Preprocessing and Model Setup

The SVM model used a similar preprocessing pipeline as Logistic Regression, with a few differences:

- Numeric columns were log-transformed where needed and then scaled.
- Ordinal columns were scaled as numeric.
- Nominal columns were one-hot encoded.

The classifier itself was an `SVC` with:

```
kernel = rbf, C = 1.5, gamma = scale, probability = True, class_weight = balanced.
```

Instead of a large grid-search (which would be prohibitively slow), a good configuration was chosen and evaluated using 5-fold StratifiedKFold on the validation split.

## Performance and Behaviour

Across 5 folds on the validation set, the SVM achieved:

Average Macro F1  $\approx 0.7354$ ,

The combined confusion matrix over all folds was:

$$\begin{bmatrix} 4117 & 1553 \\ 1562 & 4691 \end{bmatrix}$$

compared to Logistic Regression, SVM:

- slightly improves both true positives for “Stayed” and true negatives for “Left”,
- reduces the asymmetry between the two classes,
- produces a smoother probability distribution for AUC.

The downside is training time: the SVM model took tens of minutes to fit, whereas Logistic Regression completed in a fraction of that time. In a production setting, this cost must be weighed against the roughly 1–1.5 percentage point improvement in macro F1.

# Model 3: Neural Network with Nystroem Kernel Features

## Motivation

Support Vector Machines with RBF kernels perform strongly on medium-sized tabular datasets but do not scale well computationally to large samples. To enable nonlinear modelling while remaining GPU-efficient, we use a neural network trained on top of *Nystroem kernel approximations*, which map the preprocessed features into a 512-dimensional nonlinear space. This allows the network to capture interactions similar to an RBF-SVM, but with significantly faster training on large datasets.

## Architecture and Training Setup

After preprocessing (log transforms, scaling, ordinal mappings, one-hot encoding), the features are embedded using a Nystroem approximation with 512 components. The neural network operating on these features consists of:

- **Input:** 512 kernel features.
- **Hidden layers:**  $512 \rightarrow 512 \rightarrow 256 \rightarrow 128$ , each followed by BatchNorm, GELU, and dropout.
- **Output:** Linear layer producing two logits corresponding to *Left* and *Stayed*.

Training configuration:

- Optimiser: AdamW, learning rate  $10^{-3}$
- Weight decay:  $10^{-1}$
- Batch size: 512
- Loss: Cross-entropy
- Early stopping
- Internal validation: 10% of the training partition

## Validation Performance

On the standard 80/20 train-validation split applied to the *full* dataset, the neural network achieved:

$$\text{Macro F1} = 0.7444.$$

The corresponding confusion matrix on the validation data was:

$$\begin{bmatrix} 4157 & 1513 \\ 1528 & 4725 \end{bmatrix}$$

Relative to the SVM, the neural network produces:

- fewer overly confident “Stayed” predictions,
- a more balanced distribution of errors across the two classes,
- and higher recall for the minority class.

## Model Comparison and Discussion

### SVM vs. Neural Network Under Varying Data Sizes

To better understand how model performance depends on data availability, two controlled experiments were conducted:

- **Full Dataset Experiment:** 80/20 split on the entire training set ( $\approx 60,000$  samples).
- **Mini Dataset Experiment:** A randomly selected 20% subset of the data was treated as a “small dataset”, and an 80/20 split was applied within it (training on only  $\approx 9,500$  samples).

The resulting macro F1 scores were:

Model	Full 80/20 F1	Mini 80/20 F1
SVM (RBF Kernel)	0.7354	<b>0.7316</b>
Neural Network (MLP)	<b>0.7444</b>	0.6616

Table 1: Performance of SVM and Neural Network under different dataset sizes.

## Why SVM Outperforms the Neural Network on the Small Dataset

Support Vector Machines, especially with an RBF kernel, have a form of built-in capacity control through the margin maximisation principle. Even with relatively few samples, the model focuses on seeking an optimal separating hyperplane by using only the support vectors—the samples closest to the decision boundary. This makes SVMs highly effective in low-data regimes.

By contrast, this project uses a high-capacity neural network, consisting of multiple wide nonlinear layers with dropout, weight decay, and batch normalisation. Such a model needs a high volume of data to reliably estimate its parameters. When trained on the small dataset, the network underfits and shows unstable validation behavior, which results in significantly lower macro F1 of 0.6616.

Summary: the SVM’s regularised decision boundary with limited data generalizes better, whereas the neural network’s expressive power becomes a liability due to insufficient training signal.

## Why the Neural Network Overtakes SVM on the Full Dataset

When trained on the full dataset, comprising about 60,000 samples in the training split, the high capacity becomes a significant advantage for the neural network. With more data, it can learn richer nonlinear interactions that the RBF-SVM cannot efficiently capture. The combination of Batch Normalisation, and large hidden layers enable the network to model subtler behavioral patterns that result in a macro F1 of 0.7444, outperforming the SVM.

While the SVM starts to reap fewer benefits from more data, kernel methods do not scale their representational power with dataset size. It only accumulates more support vectors and incurs heavier computational cost. As a result, the SVM improves only marginally from 0.7316 to 0.7354 despite a six-fold increase in the training samples.

## Interpretation

These experiments clearly demonstrate a well-known principle in machine learning:

**Low-capacity models dominate when data is scarce; high-capacity models dominate when data is abundant.**

The SVM is stable and strong in the small-data regime because it is essentially a controlled geometric classifier. The neural network performs poorly with limited data but excels on large datasets, where its expressive power can be fully exploited.

This validates the intuition that the neural network used in this project is a *data-hungry, high-capacity learner* whose strengths only emerge when the training set is sufficiently large.

## Qualitative Behaviour

A few qualitative patterns emerged from confusion matrices and score distributions:

- All models are slightly better at predicting “Stayed” than “Left”. This might reflect true asymmetries in the data: staying behaviour could be more homogeneous, while leaving may be driven by diverse reasons, some of which are not captured in the features.
- Logistic Regression tends to under-predict the “Left” class more than SVM and NN, because a linear decision boundary in the transformed space still cannot carve out certain small pockets of “at-risk” founders.
- SVM and the NN tend to disagree on a small fraction of borderline cases. In those instances, SVM is often more conservative (predicting “Stayed”), whereas the NN sometimes predicts “Left” when certain non-linear combinations of features suggest risk.

## Design Choices and Alternatives

Throughout the pipeline, several fork-in-the-road decisions were made. A few of these are worth emphasising:

### Median vs Mean Imputation

Mean imputation is common but fragile in the presence of skewed numeric distributions. Because we observed long-tailed revenue and distance variables, median imputation was used consistently. This reduces the influence of a few very high-revenue founders on the imputed values.

### Manual Ordinal Mapping vs Automatic Encoders

Libraries such as `category_encoders` offer ordinal encoders, but here the ordering had to be explicitly reasoned about. For example, the distance between “High” and “Very High” is conceptually smaller than between “Low” and “High”.

## Additional Model: CatBoost Classifier

To complement the baseline models (Logistic Regression, SVM and Neural Networks), an additional experiment was conducted using **CatBoost**. Since the founder–retention dataset is tabular and contains several categorical attributes, CatBoost is well-suited due to its native handling of categorical features through ordered target statistics, minimal preprocessing requirements, and strong performance on structured data.

### Training Configuration

The model was trained using 5-fold cross-validation on the 20% validation split. Key settings included:

- Loss function: `Logloss`
- Task: binary classification (`Stayed` vs. `Left`)
- Iterations: 4000 (with early stopping)
- Depth: 6
- Learning rate: adaptive (CatBoost default)

CatBoost was trained using the preprocessed feature matrix with redundant columns removed. The model reached strong separation early in training (typically at 300–600 iterations), after which early stopping prevented overfitting.

### Cross-Validation Results

CatBoost achieved consistently strong performance across all folds. The validation F1 and AUC values for the 5 folds were:

Fold 1:  $F1 = 0.7601$ ,  $AUC = 0.8408$ ,  
Fold 2:  $F1 = 0.7616$ ,  $AUC = 0.8367$ ,  
Fold 3:  $F1 = 0.7667$ ,  $AUC = 0.8471$ ,  
Fold 4:  $F1 = 0.7656$ ,  $AUC = 0.8447$ ,  
Fold 5:  $F1 = 0.7752$ ,  $AUC = 0.8545$ .

The combined metrics across all folds were:

Average F1 = **0.7658**,      Average AUC = **0.8448**.

These values place CatBoost among the top-performing models for this prediction task, achieving both high discrimination (AUC) and balanced performance across classes (F1 score).

## Confusion Matrix

Aggregating predictions from all validation folds gives the combined confusion matrix:

$$\begin{bmatrix} 4178 & 1490 \\ 1448 & 4804 \end{bmatrix}$$

This indicates strong recall for both classes, with well-balanced prediction behaviour—especially notable given the natural imbalance in founder retention outcomes.

## Conclusion

This project shows that startup founder retention can be predicted reasonably well using standard machine learning methods, provided that the preprocessing and encoding steps are carefully designed.

The main conclusions are:

- Clean preprocessing — especially correct ordinal mapping, boolean normalisation and sensible outlier handling — is essential. The models depend critically on the quality of this numeric representation.
- Logistic Regression already provides a strong and interpretable baseline with macro F1 around 0.736 and is suitable when training time or transparency is more important than squeezing out the last few percentage points.
- SVM with an RBF kernel achieves the best performance (macro F1 around 0.751) but at a high computational cost. It is best viewed as a high-accuracy classical reference model.
- A modest two-layer Neural Network matches SVM-level performance within a small margin and opens the door for more advanced deep-learning approaches tailored to tabular data.