

# Introduction to ML for Business

## Case Study Worksheet: Startup Profiling from Investments

**Group:**

**Date:**

**Course:** Machine Learning for Business

### Problem Representation

*How can you profile a startup using investment data?*

Identify and streamline the types of data that are critical for a thorough analysis to understand the factors that influence startup success and investment attractiveness. Consider the impact of different data types and how they can reveal insights into startup viability, market trends and investor behavior.

Consider a situation where investors are evaluating startups for potential investment. How should they approach this data-driven analysis? Consider variables such as industry, funding history, and growth metrics in your evaluation.

### Data Source

**Dataset Description:** The [StartUp Investments \(Crunchbase\) dataset](#) is a comprehensive collection of data on startups and their investments, sourced from Crunchbase.

**Dataset Features:** It includes detailed information such as startup names, sectors, funding rounds, investment amounts, and investor details. This dataset is useful for analyzing investment patterns, identifying characteristics of successful startups, and understanding investor behavior.

**Access the dataset:** To access this dataset, visit the [Kaggle link](#) provided.

Here is an example of the dataset records.

| name          | homepage_url                                    | funding_total_usd | status | country_code |
|---------------|---|-------------------|--------|--------------|
| .Club Domains | <a href="http://nic.club/">http://nic.club/</a> | 70,00,000         | nan    | USA          |

| name          | homepage_url  | funding_total_usd | status    | country_code |
|---------------|---|-------------------|-----------|--------------|
| .Fox Networks | <a href="http://www.dotfox.com">http://www.dotfox.com</a> | 49,12,393         | closed    | ARG          |
| 0-6.com       | <a href="http://www.0-6.com">http://www.0-6.com</a>       | 20,00,000         | operating | nan          |

## Task Identification

With the goal of profiling startups based on comprehensive investment data, consider the following areas of focus for your analysis:

- **Financial pattern analysis:** Drill down into `funding_total_usd` and different funding types (e.g., seed, venture) to uncover patterns in financial support and their correlation with startup success.
- **Evaluate funding timelines:** Examine `'first_funding_at'`, `'last_funding_at'` and funding round data to understand the impact of funding timelines on startup growth and sustainability.
- **Assess startup viability:** Use `'status'` and `'founded_year'` to assess how company age and current status relate to investment attractiveness and success rates.
- **Market Analysis:** Use `'market'` and `'category_list'` to determine which sectors are attracting more investment and why.

## Critical questions

- *What financial benchmarks indicate promising investment opportunities?*
- *How do funding patterns differ across markets, and what does this say about sector viability?*
- *Can startup longevity and status be used as reliable predictors of investment success?*

## Data Exploration

Critically analyze the data set to identify any limitations that may affect our analysis.

Key areas of focus include:

- **Missing value assessment:** Examine the extent of missing data. Are there significant gaps, especially in critical variables essential to our analysis?
- **Relevance of missing data:** Determine whether the missing values relate to critical factors that could skew our understanding of startup investment.
- **Data consistency check:** Examine the data set for anomalies or deviations from expected patterns that could affect the accuracy of our conclusions.

## Data Preparation for Machine Learning

To effectively prepare a machine learning model, consider these critical steps:

- **Handling missing data:** Discuss a strategy for handling missing data points. Options include data imputation, removing rows/columns with excessive missing values, or using models that can inherently handle missing data.
- **Feature Engineering Tactics:** Explore the development of derived metrics or new variables that improve model predictability. For example, create aggregate variables (such as average funding amount per sector) or interaction terms (such as the ratio of seed funding to venture funding) that can provide deeper insights into investment patterns.

## Model selection

In case you adopt a supervised approach by using **regression**, address the following:

- **Feature Inclusion:** Identify key features to include in your regression model.
- **Algorithm Selection:** Select an appropriate regression algorithm. Options include linear regression for simplicity or more complex methods such as ridge or lasso regression to handle multicollinearity. Which best fit the problem?
- **Prediction Target:** Define the specific target you want to predict.

## Model Evaluation

**Train and test the model:** When profiling startups using investment data, provide a proper *dataset split* strategy, ensuring a representative balance. Chose a proper *hyperparameter tuning* techniques. Address data imbalance.

**Interpret the results:** Evaluate the model using metrics relevant to investment profiling, such as *precision*, *recall*, and *F1 score*, to measure both the accuracy and reliability of the predictions. *Which best fit the problem at hand?*

**Model limitations:** Recognize any assumptions inherent in your model and potential **biases** in the data set. For example, the model may assume linear relationships or overlook the impact of external economic factors on startup success.