# Comprehensive Revision Notes for the Loan Tap Machine Learning Case Study

## Introduction

This session covers a machine learning case study based on Loan Tap, an online platform providing customized loan products to millennials. The task involves building an underwriting layer to determine the creditworthiness of MSMEs (Micro, Small, and Medium Enterprises) and involves using a dataset specific to personal loans provided by Loan Tap【4:0†transcript.txt】.

## Problem Statement

The objective of the case study is to predict whether an individual should be extended a credit line based on historical loan data. This involves deciding if a credit line should be extended and determining repayment terms【4:0†transcript.txt】.

## Dataset Overview

1. **Features**: The dataset contains 27 features or columns that include loan-related attributes【4:4†transcript.txt】.
2. **Target Variable**: The target is a binary classification problem with "loan status" being the target column - predicting whether a loan is "fully paid" or "charged off"【4:7†transcript.txt】.
3. **Data Dictionary**: It provides meanings for each feature, which is crucial for understanding what the data encapsulates【4:10†transcript.txt】.

## Data Preprocessing

1. **Handling Missing Values**: Addressing null values using logical strategies instead of outright removal【4:5†transcript.txt】.
2. **Data Encoding**: Using one-hot encoding for categorical variables such as 'purpose', 'zip code', etc.【4:11†transcript.txt】.
3. **Feature Engineering**:

- Dropping irrelevant columns that might introduce bias 【4:11†transcript.txt】.

- Handling skewness in data by transformations like box-cox or log transformations to achieve a normal distribution 【4:17†transcript.txt】.

## Exploratory Data Analysis (EDA)

1. **Univariate and Bivariate Analysis**: Initial analysis to identify important features and their relationships, especially focusing on the loan amount and home ownership 【4:3†transcript.txt】.

2. **Understanding Distribution**: Identifying skewness and variance in columns to decide on further transformation steps 【4:3†transcript.txt】.

3. **Data Visualization**: Visualizing categorical data like grading systems to understand customer distribution and payment likelihood 【4:13†transcript.txt】.

## Model Building and Evaluation

1. **Model Selection**: Starting with logistic regression for this binary classification problem 【4:7†transcript.txt】.

2. **Handling Imbalance**: Techniques like oversampling using SMOTE to balance classes 【4:1†transcript.txt】.

3. **Evaluation Metrics**: Focusing beyond accuracy to metrics like F1-score due to data imbalance issues 【4:12†transcript.txt】.

4. **Experimentation**: Testing various methods such as Grid Search for hyperparameter tuning, Regularization techniques (like L1/L2) to avoid overfitting 【4:16†transcript.txt】.

## Advanced Considerations

1. **Feature Selection**: Using methods to make feature weights zero for irrelevant features 【4:14†transcript.txt】.

2. **Ensembling Techniques**: Combining multiple models to improve prediction performance, showcasing the advantage through ensemble methods 【4:18†transcript.txt】.

3. **Handling Overengineering**: Understanding the limits of model improvement through feature engineering and recognizing when to shift strategies 【4:15†transcript.txt】.

- Real-world data preprocessing involves iterative validation.
- Interpretation of results and method adaptations require consultation with domain experts【4:19†transcript.txt】.
- Experimentation and validation are critical components in machine learning system design【4:14†transcript.txt】.

These concepts lay the foundation for understanding how machine learning approaches can be applied to financial datasets to predict creditworthiness and optimize lending practices.