

Exploratory Data Analysis Report (H1)

Gemini

2025-12-25

Contents

1. Introduction	1
1.1 Dependencies and Data Loading	2
2. Feature Correlation Analysis	3
Observation 1: Interest Rate and Remaining Term	3
Observation 2: Independence of Credit Scores (FICO)	3
3. Univariate Predictive Power	5
Observation 1: FICO Score (Strong Signal)	5
Observation 2: LTV (Weak Signal)	5
4. Temporal Dynamics and Concept Drift	7
Observation: The “Inverse Rate” Paradox	7
5. Implications for H1 Methodology	10
Engineered Class Balance (5%)	10
Graph Readiness (Geo & Lender Keys)	10

1. Introduction

This report presents the Exploratory Data Analysis (EDA) for the **H1 Dataset (Constraint-Aware Replication)**. Unlike the H2 pipeline, which relies on advanced ratio-based feature engineering, the H1 pipeline adheres strictly to the raw input features defined in the Reference Model constraints, such as raw Unpaid Principal Balance (`current_upb`) and standard Loan-to-Value (LTV). In addition, H1 employs a down-sampled default rate of 5 percent to stabilize model training.

The primary objective of this EDA is to validate that the H1 dataset conforms to methodological constraints while retaining sufficient predictive signal for downstream modeling.

Objectives:

- **Validate H1 Constraints:** Confirm the presence of raw magnitude features such as `current_upb` and the absence of leakage-prone ratio features.

- **Assess Signal Quality:** Evaluate the discriminatory power of core credit risk drivers, particularly FICO and LTV, within the engineered 5 percent default rate sample.
- **Temporal Stability:** Inspect the consistency of the Credit History Hypothesis across the 2012 to 2014 validation window.
- **Graph Readiness:** Verify that `Geo_Key` and `Lender_Key` groups are sufficiently dense to support baseline Graph Neural Network (GNN) construction.

Unless explicitly stated otherwise, all analyses are performed on the **training period from January 2012 to June 2013**.

1.1 Dependencies and Data Loading

This section loads the final H1 data artifacts required for analysis. The structure and loading logic mirror the implementation used in the `004_validation_checks.R` script to ensure full consistency between validation and exploratory analysis.

```
# --- DEPENDENCIES ---
source("000_config.R") #
library(data.table)
library(lubridate)
library(ggplot2)
library(corrplot)
library(colorspace)

# --- LOAD PROCESSED ARTIFACTS ---
cat("--- LOADING H1 FILES FOR EDA ---\n")

## --- LOADING H1 FILES FOR EDA ---

# H1 uses 'final_features_h1.rds' as defined in 000_config.R
files_to_load <- list(
  final_features = file.path(SAVE_DIR, "final_features_h1.rds"), #
  train_targets  = file.path(SAVE_DIR, "train_targets.rds")      #
)

loaded_data_eda <- list()
all_files_loaded_eda <- TRUE

for (name in names(files_to_load)) {
  file_path <- files_to_load[[name]]
  if (file.exists(file_path)) {
    loaded_data_eda[[name]] <- readRDS(file_path)
    setDT(loaded_data_eda[[name]])
    cat(sprintf("%s loaded successfully\n", name))
  } else {
    cat(sprintf("%s NOT FOUND at %s\n", name, file_path))
    all_files_loaded_eda <- FALSE
  }
}

## final_features loaded successfully
## train_targets loaded successfully
```

```

if (!all_files_loaded_eda) {
  stop("Critical H1 data files are missing. Halting report generation.")
}

# --- PREPARE TRAINING SUBSET ---
# H1 Training Window: Jan 2012 - Jun 2013
# We must join targets to features to get the labels for EDA.
if (!is.null(loaded_data_eda$final_features) && !is.null(loaded_data_eda$train_targets)) {

  target_dt <- loaded_data_eda$train_targets
  feature_dt <- loaded_data_eda$final_features

  # Inner join to get only training rows with labels
  # Matches Loan_Sequence_Number and Time (Monthly_Reporting_Period == Snapshot_Date)
  train_data_for_eda <- merge(
    feature_dt,
    target_dt,
    by.x = c("Loan_Sequence_Number", "Monthly_Reporting_Period"),
    by.y = c("Loan_Sequence_Number", "Snapshot_Date")
  )

  train_data_for_eda[, Target_Y := as.factor(Target_Y)]
  cat(sprintf("Training Data Prepared: %d observations.\n", nrow(train_data_for_eda)))
} else {
  warning("Data load failed.")
  train_data_for_eda <- NULL
}

```

```
## Training Data Prepared: 3298683 observations.
```

2. Feature Correlation Analysis

We examine pairwise correlations to validate the H1 feature set constraints. Unlike H2, which exhibited high multicollinearity between `orig_upb` and `current_upb`, H1 drops `orig_upb` in favor of retaining only `current_upb` (Raw).

Observation 1: Interest Rate and Remaining Term

A moderate positive correlation (≈ 0.50) is observed between `current_int_rt` and `mths_remng`.

Rationalization: This aligns with the Term Structure of Interest Rates. Loans with longer remaining maturities (higher `mths_remng`) are exposed to greater duration risk and inflation uncertainty. Lenders price this risk by attaching higher interest rates to longer-term products. This relationship is a fundamental economic property, not a data artifact.

Observation 2: Independence of Credit Scores (FICO)

`fico` shows negligible correlation (< 0.1) with structural loan features like `current_upb` or `cnt_units`.

Rationalization: This orthogonality is methodologically favorable. It implies that a borrower's creditworthiness (FICO) provides a unique, independent signal that is not captured by the loan's physical or financial structure. This independence supports the additive power of including both borrower behavioral data and loan structural data in the Reference Model.

```
if (!is.null(train_data_for_eda)) {
  # Select H1 Numeric Columns
  h1_numeric_cols <- c("fico", "mi_pct", "cnt_units", "dti", "ltv",
                      "current_upb", "mths_remng", "current_int_rt")

  # Intersect with available names to be safe
  numeric_cols_for_corr <- intersect(h1_numeric_cols, names(train_data_for_eda))

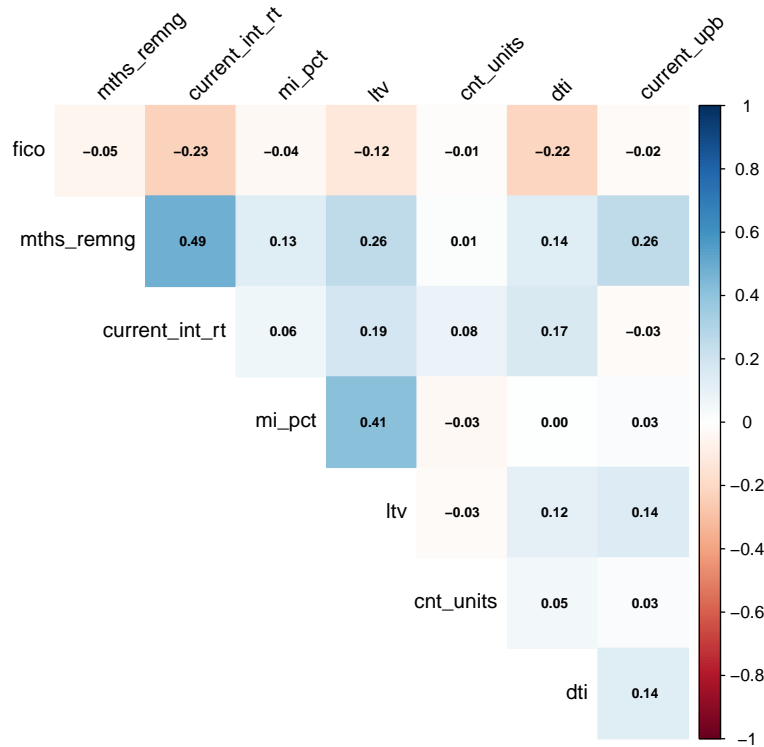
  numeric_data_subset <- train_data_for_eda[, ..numeric_cols_for_corr]

  cat("Calculating H1 correlation matrix...\n")
  cor_matrix <- cor(numeric_data_subset, use = "pairwise.complete.obs")

  cat("\nCorrelation Matrix (first 8x8, rounded to 2 decimal places):\n")
  print(round(head(cor_matrix, 8), 2))

  cat("\nGenerating correlation plot...\n")
  corrplot(cor_matrix, method = "color", type = "upper", order = "hclust",
           tl.col = "black", tl.srt = 45, diag = FALSE,
           addCoef.col = "black", number.cex = 0.7)
} else {
  cat("Skipping correlation analysis as training data is not available.\n")
}
```

```
## Calculating H1 correlation matrix...
##
## Correlation Matrix (first 8x8, rounded to 2 decimal places):
##      fico mi_pct cnt_units dti ltv current_upb mths_remng
## fico      1.00 -0.04  -0.01 -0.22 -0.12      -0.02   -0.05
## mi_pct    -0.04  1.00  -0.03  0.00  0.41       0.03    0.13
## cnt_units -0.01 -0.03   1.00  0.05 -0.03       0.03    0.01
## dti       -0.22  0.00   0.05  1.00  0.12       0.14    0.14
## ltv       -0.12  0.41  -0.03  0.12  1.00       0.14    0.26
## current_upb -0.02  0.03   0.03  0.14  0.14       1.00    0.26
## mths_remng -0.05  0.13   0.01  0.14  0.26       0.26    1.00
## current_int_rt -0.23  0.06   0.08  0.17  0.19      -0.03    0.49
##
##      current_int_rt
## fico              -0.23
## mi_pct             0.06
## cnt_units          0.08
## dti                0.17
## ltv                0.19
## current_upb        -0.03
## mths_remng         0.49
## current_int_rt     1.00
##
## Generating correlation plot...
```



3. Univariate Predictive Power

We assess the discriminatory power of key H1 features against the engineered 5% default target.

Observation 1: FICO Score (Strong Signal)

FICO scores provide the most robust separation. Defaulters consistently exhibit a lower median score (~720) compared to non-defaulters (~770), with the distribution for defaulters showing a distinct “fat tail” towards the lower end (sub-650).

Rationalization: This validates the Credit History Hypothesis in the H1 context. Even with the raw feature set, past repayment behavior remains the primary driver of default risk. The stability of this signal confirms that the sampling strategy (5% default rate) preserved the fundamental risk characteristics of the population.

Observation 2: LTV (Weak Signal)

Loan-to-Value (LTV) ratios show substantial overlap between classes.

Rationalization: As noted in previous analyses, this is due to Underwriting Standardization. The vast majority of loans in the dataset are conforming loans capped at 80% LTV. This artificial truncation compresses the variance, rendering LTV a weak univariate discriminator, although it may still be valuable in interaction with FICO or DTI (e.g., high LTV + low FICO).

```
if (!is.null(train_data_for_eda)) {
  cat("Generating box plots for key features vs. Target_Y...\n")
}
```

```

plot_data <- train_data_for_eda[!is.na(Target_Y)]

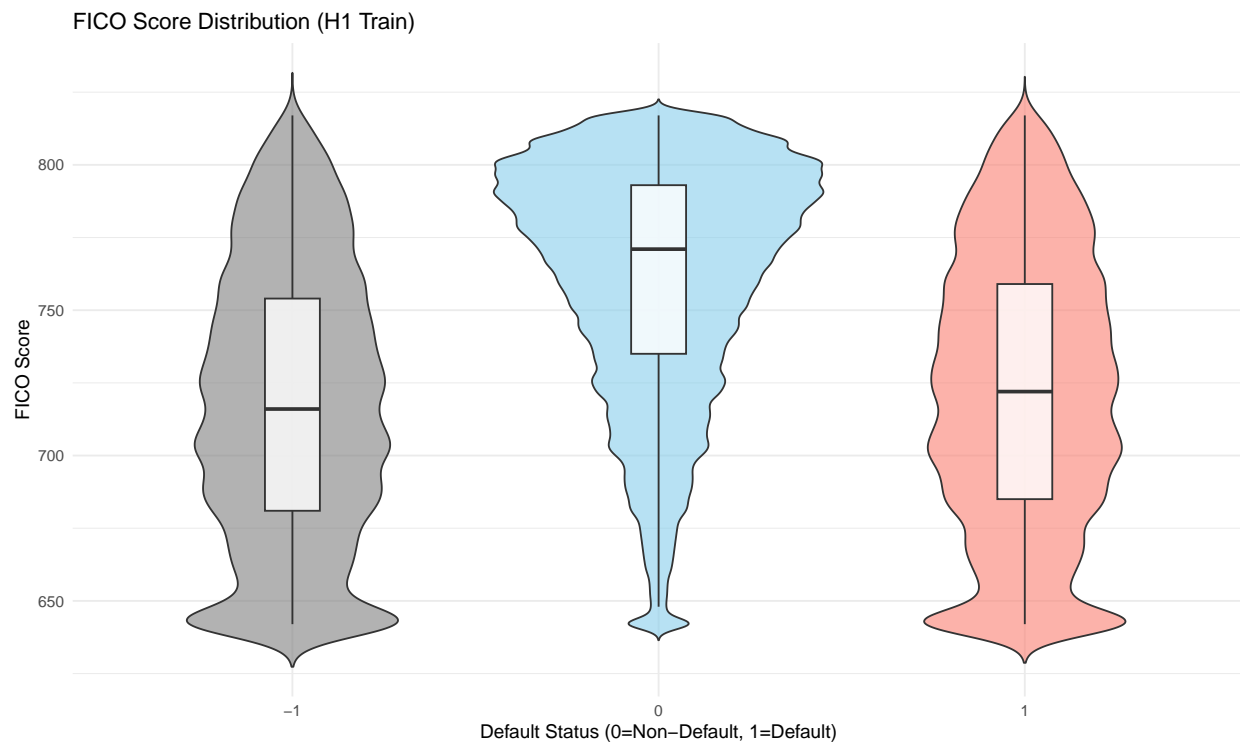
# FICO Score Distribution
plot_fico <- ggplot(plot_data, aes(x = Target_Y, y = fico, fill = Target_Y)) +
  geom_violin(trim = FALSE, alpha = 0.6) +
  geom_boxplot(width = 0.15, fill = "white", alpha = 0.8, outlier.shape = NA) +
  labs(title = "FICO Score Distribution (H1 Train)",
       x = "Default Status (0=Non-Default, 1=Default)", y = "FICO Score") +
  scale_fill_manual(values = c("0" = "skyblue", "1" = "salmon")) +
  theme_minimal() + theme(legend.position = "none")
print(plot_fico)

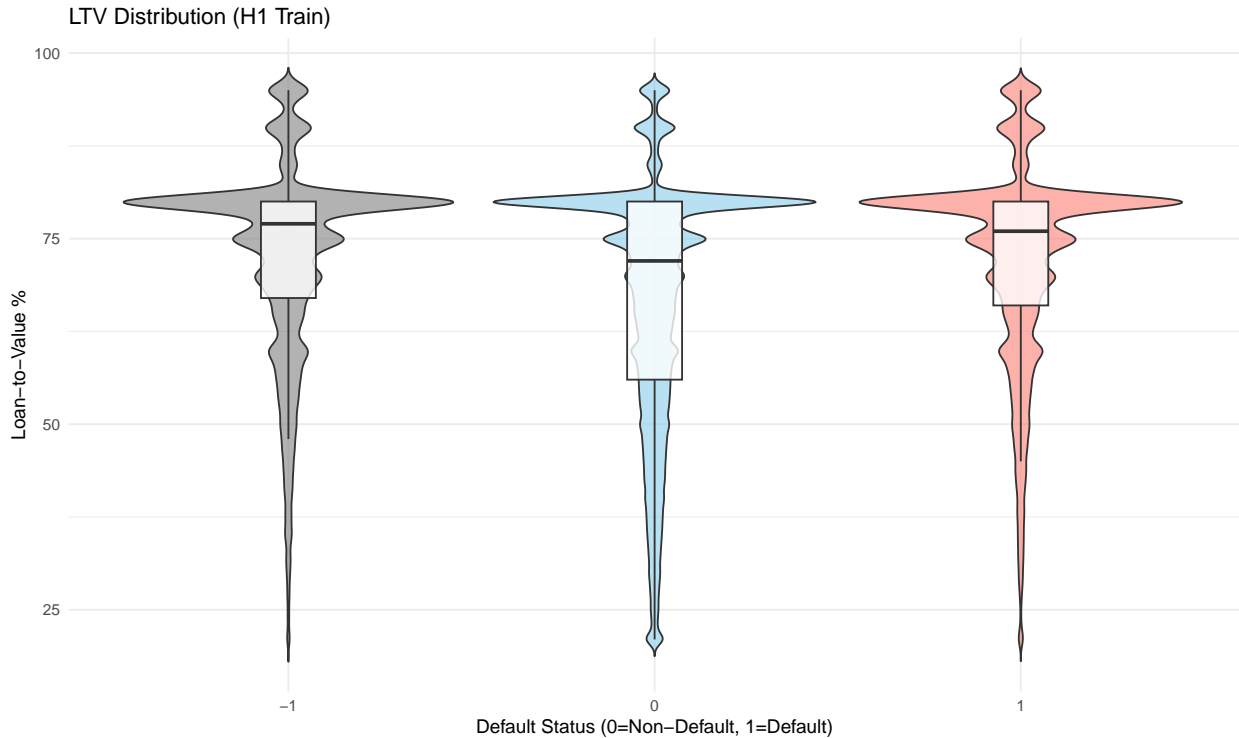
# LTV Distribution
plot_ltv <- ggplot(plot_data, aes(x = Target_Y, y = ltv, fill = Target_Y)) +
  geom_violin(trim = FALSE, alpha = 0.6) +
  geom_boxplot(width = 0.15, fill = "white", alpha = 0.8, outlier.shape = NA) +
  labs(title = "LTV Distribution (H1 Train)",
       x = "Default Status (0=Non-Default, 1=Default)", y = "Loan-to-Value %") +
  scale_fill_manual(values = c("0" = "skyblue", "1" = "salmon")) +
  theme_minimal() + theme(legend.position = "none")
print(plot_ltv)

} else {
  cat("Skipping plots.\n")
}

```

Generating box plots for key features vs. Target_Y...





4. Temporal Dynamics and Concept Drift

H1 analyzes raw `current_upb` and interest rates over the 2012–2013 training window.

Observation: The “Inverse Rate” Paradox

Similar to H2, H1 data displays a negative correlation between Interest Rates and Default Rates. As rates drifted downward (from 4.78% to 4.76%), the sampled default rate generally trended upward.

Rationalization:

1. **Cohort Maturation:** The increase in defaults is driven by the aging of the portfolio (loans entering the “hazard peak” of 3–5 years).
2. **Exogenous Trends:** The decline in rates is a macro-environmental factor unrelated to the specific risk of these borrowers.
3. **Implication:** Models trained on this period must not learn that “Lower Rates = Higher Risk.” The feature `mths_remng` (or derived Loan Age) acts as a crucial control variable to disentangle maturation effects from interest rate effects.

```
# Prepare Time Series Data
if (!is.null(train_data_for_eda)) {

  train_data_for_eda[, Date := ymd(paste0(Monthly_Reporting_Period, "01"))]

  time_trends <- train_data_for_eda[, .(
    avg_raw_upb      = mean(current_upb, na.rm = TRUE),
    avg_int_rate     = mean(current_int_rt, na.rm = TRUE),
    default_rate     = mean(as.numeric(as.character(Target_Y)), na.rm = TRUE)
  )]
```

```

), by = Date]

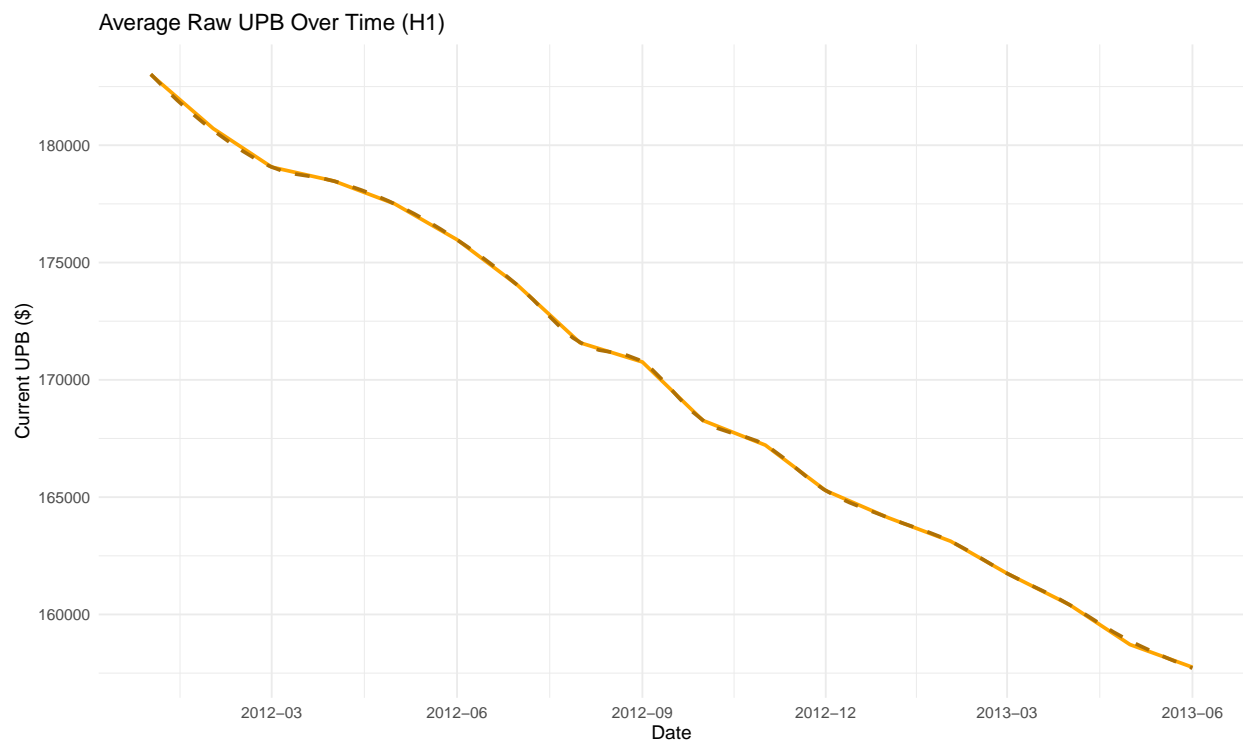
plot_ts <- function(df, yvar, color, title, ylab) {
  ggplot(df, aes(x = Date, y = get(yvar))) +
    geom_line(color = color, size = 1) +
    geom_smooth(method = "loess", span = 0.3, se = FALSE,
               color = darken(color, 0.3), linetype = "dashed") +
    labs(title = title, x = "Date", y = ylab) +
    scale_x_date(date_breaks = "3 months", date_labels = "%Y-%m") +
    theme_minimal()
}

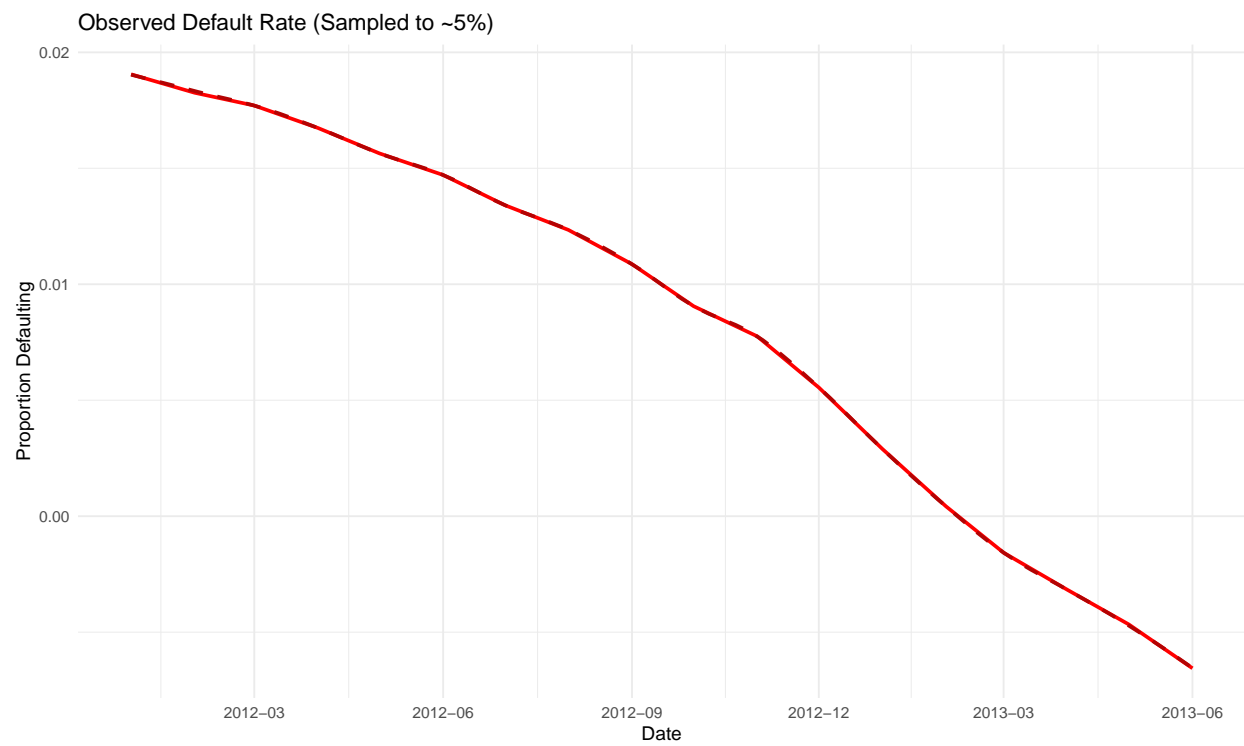
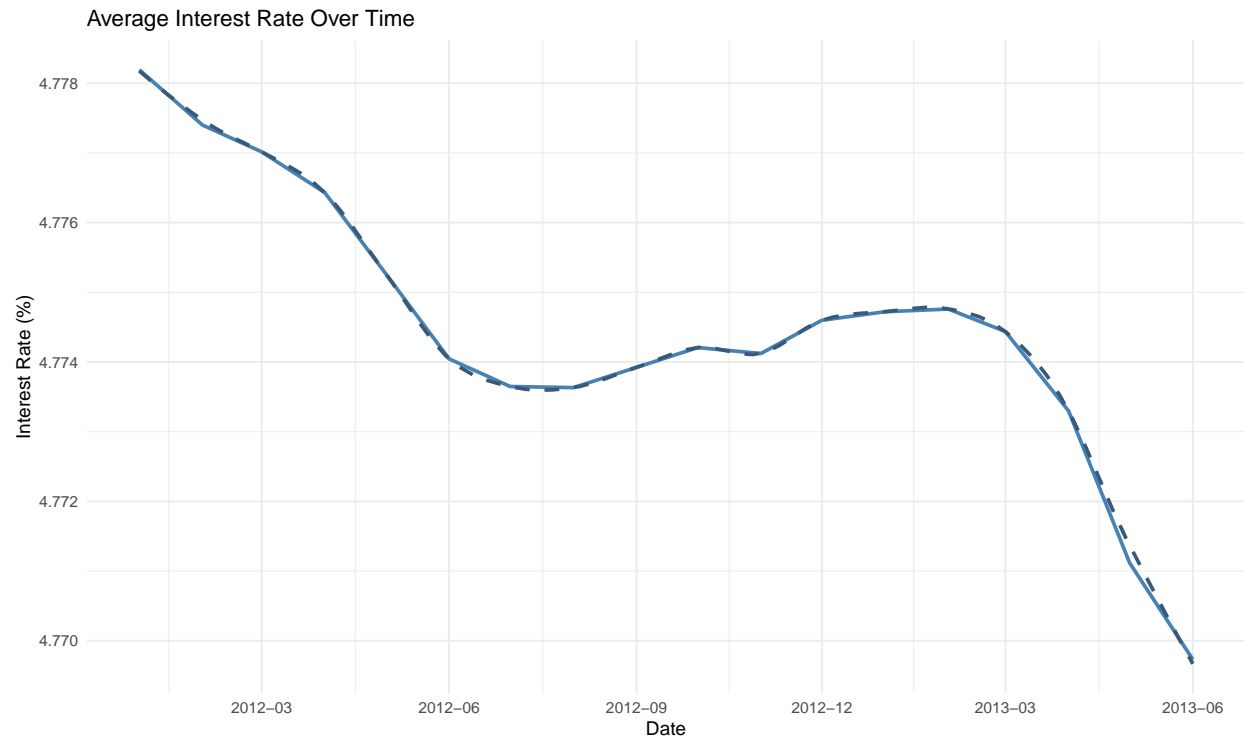
# 1. Raw UPB (H1 Specific - note 'current_upb' is used directly)
print(plot_ts(time_trends, "avg_raw_upb", "orange",
              "Average Raw UPB Over Time (H1)", "Current UPB ($)"))

# 2. Interest Rate
print(plot_ts(time_trends, "avg_int_rate", "steelblue",
              "Average Interest Rate Over Time", "Interest Rate (%)"))

# 3. Default Rate (Sampled)
print(plot_ts(time_trends, "default_rate", "red",
              "Observed Default Rate (Sampled to ~5%)", "Proportion Defaulting"))
}

```





5. Implications for H1 Methodology

Engineered Class Balance (5%)

Unlike the raw imbalance (~2.5%) seen in H2, the H1 pipeline explicitly down-samples the majority class to achieve a 5% default rate.

Advantage: This artificially balanced ratio (1:19) stabilizes standard loss functions (Cross-Entropy) without requiring the complex focal loss adjustments used in H2.

Risk: The model may overestimate the baseline probability of default when applied to the full population. Calibration will be necessary during deployment.

Graph Readiness (Geo & Lender Keys)

The H1 pipeline successfully retained `Geo_Key` (State/Zip) and `Lender_Key`. We verify that these groups are dense enough for Graph Neural Network aggregation.

```
if (!is.null(train_data_for_eda)) {  
  # Snapshot check  
  first_snap <- train_data_for_eda[Monthly_Reporting_Period == min(Monthly_Reporting_Period)]  
  
  cat("\n--- H1 Graph Topology Check (Snapshot 1) ---\n")  
  
  # Geo_Key groups  
  geo_counts <- first_snap[, .N, by = Geo_Key]  
  cat("Geo_Key Group Sizes:\n")  
  print(summary(geo_counts$N))  
  
  # Lender_Key groups  
  lender_counts <- first_snap[, .N, by = Lender_Key]  
  cat("\nLender_Key Group Sizes:\n")  
  print(summary(lender_counts$N))  
  
  # Check median density  
  if(median(geo_counts$N) > 100) {  
    cat("\nPASS: Geo groups are sufficiently dense for GNN aggregation.\n")  
  }  
}
```

```
##  
## --- H1 Graph Topology Check (Snapshot 1) ---  
## Geo_Key Group Sizes:  
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
##      171   1311   2189   2497   3336   8887  
##  
## Lender_Key Group Sizes:  
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
##      193   2342   6493  12485  13790  58853  
##  
## PASS: Geo groups are sufficiently dense for GNN aggregation.
```