

# Data Wrangling and Cleaning: Steps 1-11

David Mullaly

2026-02-09

## Introduction

This document covers the data wrangling process and all eleven data cleaning steps for the problem. We load multiple relational tables, merge them into a single flat file, and perform comprehensive data cleaning to prepare for missing data analysis.

### Data Cleaning Steps Covered:

#### Steps 1-5 (Initial Cleaning)

1. Open data in your software of choice
2. Review variables for common sense based on SME knowledge
3. Review how the software coded the variables (nominal, continuous)
4. Perform data integrity/validation checks
5. Handle dates

#### Steps 6-7 (Categorical & Zero-Variance)

6. Handle categorical variables - keep as is, combine rare levels, combine similar levels
7. Handle zero-variance predictors

#### Steps 8-11 (Advanced Cleaning)

8. Handle near zero-variance predictors
9. Eliminate redundant columns and linear combination columns
10. Search for outliers and initial search for missing values
11. Sanity check using Decision Tree (1 to 2 splits)

## Data Loading and Wrangling

Load four relational tables and merge into a single flat file (FFdf) using left joins on Cust\_ID.

```
## MainDF: 9447 48 | StoreDF: 9272 3 | ConcessDF: 9272 3 | CustomerDF: 14272 7
```

```
## Duplicates in MainDF: 175
```

```
## Duplicates in StoreDF: 0
```

```
## Duplicates in ConcessDF: 0
```

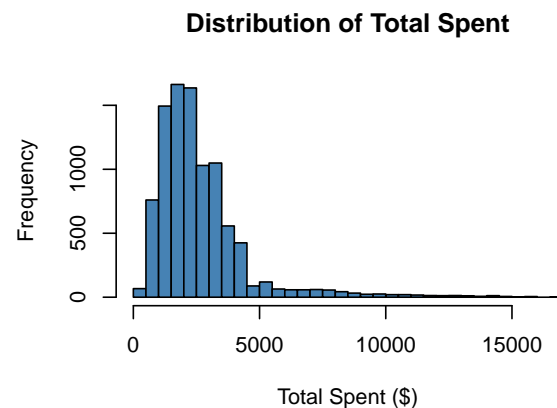
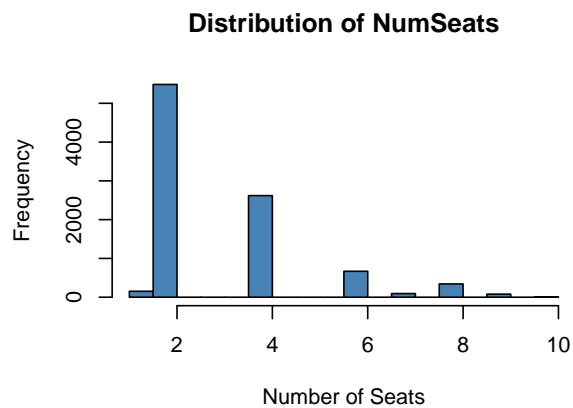
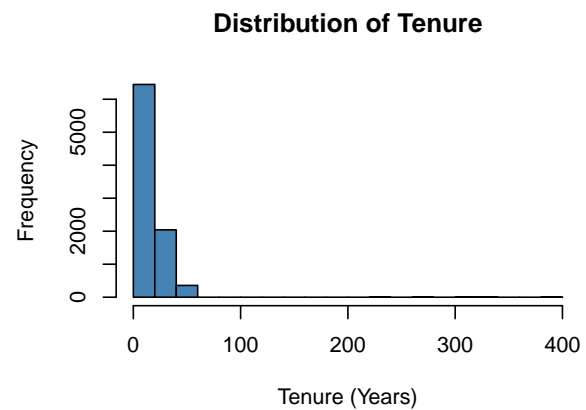
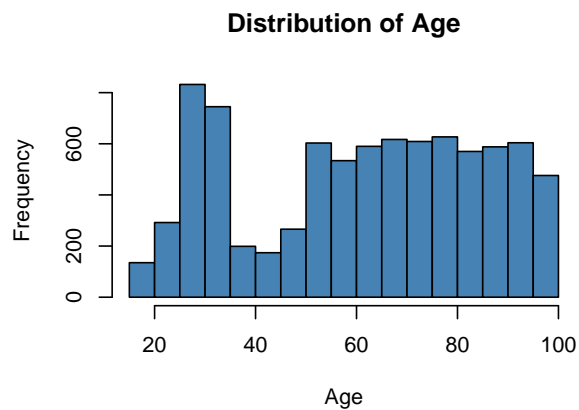
```
## Duplicates in CustomerDF: 0
```

```
## Final flat file dimensions - Rows: 9447 Columns: 58
```

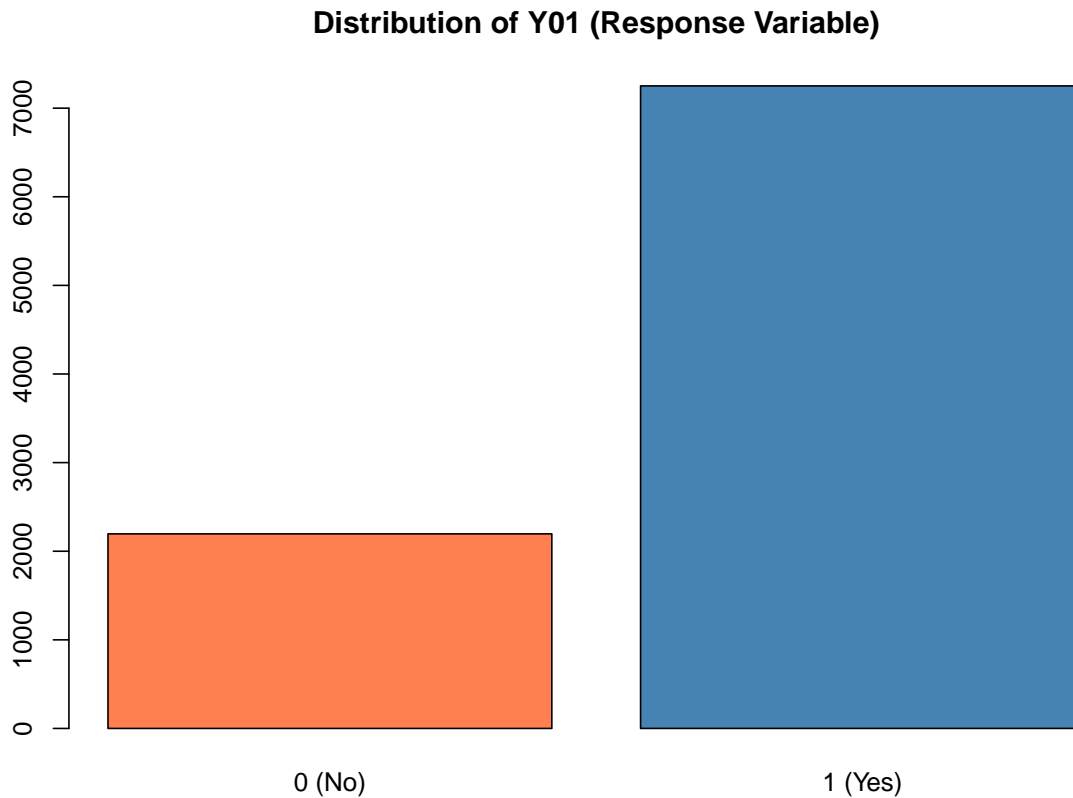
```
## Unique Cust_ID: 9272 | Duplicate rows: 175
```

## Step 1: Open Data in Your Software of Choice

Create identifier variable, arrange columns, and explore distributions visually.



```
## Response Variable (Y01): 2196 7251
```



**Observations:** Age is roughly normal (30-60), Tenure is right-skewed, NumSeats clusters around 2-4, Total Spent is heavily right-skewed. Response variable Y01 is reasonably balanced.

## Step 2: Review Variables for Common Sense (SME Knowledge)

Standardize variable names and check for unique identifiers.

```
## Dataset: 9447 rows, 59 columns
```

```
## Unique Cust_ID: 9272 out of 9447 rows
```

**Result:** Each row does not necessarily represent a unique customer (175 duplicate Cust\_IDs found in MainDF source data).

## Step 3: Review How Software Coded Variables

Convert character variables to factors for proper categorical analysis.

```
## Character variables to convert: 23
```

```
##
```

```
## Sex levels: F M
```

```
## Marital levels: D M S U
```

```
## Account_Type levels: Business Personal Shared
```

**Note:** Marital has levels D (Divorced), M (Married), S (Single), U (Unknown). The “U” for Unknown may need special handling.

## Step 4: Data Integrity/Validation Checks

Check for anomalies, bogus values, and data quality issues.

```
## Age range: 18 99
```

```
## Tenure range: 1 400
```

```
## NumSeats range: 1 10
```

```
##
```

```
## DistA values = 999: 1506
```

```
## DistA NA count after conversion: 2895
```

```
##
```

```
## Survey_Comp range: 0 7.4
```

```
## Survey_Comp values > 1: 110
```

```
##
```

```
## State_Name unique: 50 | State_Loc unique: 50
```

### Issues Found:

- **DistA = 999:** Placeholder values converted to NA
- **Survey\_Comp > 1:** Outlier found when expected range is 0-1
- **Age max = 99:** May be placeholder or extreme value - verify with SME
- **Tenure max = 400:** Suspicious value if measured in years - verify units with SME
- **State\_Name/State\_Loc:** Redundant columns (same info, different format)
- **Marital = “U”:** Unknown status - consider treating as NA
- **Cust\_ID duplicates:** 175 duplicate IDs found in MainDF source data
- **Address, Name, PhoneNum:** 100% missing - likely removed from CustomerDF for privacy

## Step 5: Handle Dates

Convert Last\_Contact datetime and extract useful components.

```
## Date components extracted: Contact_Year, Contact_Month, Contact_Day, Contact_Weekday, Contact_Hour
```

```
## Contact Year range: 2018 2025
```

```
## Contact Hour range: 0 23
```

## Summary

```
## Final Dataset: 9447 rows x 64 columns
```

```
## Total Missing Values: 94591
```

```
##
```

```
## Columns with missing values:
```

```
##           Address           Name           PhoneNum
##           9447           9447           9447
## Educational_Level Favorite_Caps_Player Favorite_Sport
##           6612           6612           6612
##           Job_Sector Mode_Of_Transport           Team_B_STH
##           6612           6612           6612
##           Team_C_STH Net_Worth_True HouseHold_Income_True
##           6612           5762           5691
##           DistA           Marital           Age
##           2895           1155           986
##           Rep_Name           Sex           Tenure
##           871           667           592
##           Rep_Visits Rep_Calls           Num_Children
##           508           433           406
```

### Data Quality Issues for Future Steps:

Issue	Possible action / Action Taken
DistA = 999	Converted to NA
Survey__Comp > 1	Flag for investigation (110 values, max 7.4)
Age = 99	Verify with SME
Tenure = 400	Verify units with SME (400 years unlikely)
Marital = "U"	Consider as NA or keep
State redundancy	Drop one column
ID columns	Exclude from modeling
Cust_ID duplicates	175 duplicates in MainDF - investigate or deduplicate
PII columns	Address, Name, PhoneNum 100% missing - exclude

## Class 3: Data Cleaning Process: Steps 6-7

Step 6: Handle Categorical Variables - keep as is, combine rare levels, combine similar levels

```
# Get all factor variables
factor_vars <- names(FFdf)[sapply(FFdf, is.factor)]
cat("Factor variables:", length(factor_vars), "\n")
```

```
## Factor variables: 23
```

```

# Function to check for rare levels (less than 5% of data)
check_rare_levels <- function(df, threshold = 0.05) {
  factor_vars <- names(df)[sapply(df, is.factor)]
  rare_info <- list()

  for(var in factor_vars) {
    tbl <- table(df[[var]], useNA = "ifany")
    pct <- prop.table(tbl)
    rare_levels <- names(pct[pct < threshold])
    if(length(rare_levels) > 0) {
      rare_info[[var]] <- data.frame(
        Level = rare_levels,
        Count = as.numeric(tbl[rare_levels]),
        Percent = round(as.numeric(pct[rare_levels]) * 100, 2)
      )
    }
  }
  return(rare_info)
}

# Check for rare levels
rare_levels <- check_rare_levels(FFdf, threshold = 0.05)
cat("\nVariables with rare levels (< 5%):\n")

```

```

##
## Variables with rare levels (< 5%):

```

```

print(names(rare_levels))

```

```

## [1] "Educational_Level"      "Favorite_Caps_Player"  "Favorite_Sport"
## [4] "Favorite_Team"          "Job_Sector"            "Marital"
## [7] "Mode_Of_Transport"      "Most_Purch_Concession" "Mult_Loc"
## [10] "Rep_Name"               "Seating_Location"      "State_Loc"
## [13] "State_Name"

```

```

# Display level distributions for key categorical variables
cat("\n--- Sex Distribution ---\n")

```

```

##
## --- Sex Distribution ---

```

```

print(table(FFdf$Sex, useNA = "ifany"))

```

```

##
##      F      M <NA>
## 1078 7702  667

```

```

cat("\n--- Marital Distribution ---\n")

```

```

##
## --- Marital Distribution ---

```

```

print(table(FFdf$Marital, useNA = "ifany"))

##
##      D      M      S      U <NA>
##  908 6227  909  248 1155

cat("\n--- Account_Type Distribution ---\n")

##
## --- Account_Type Distribution ---

print(table(FFdf$Account_Type, useNA = "ifany"))

##
## Business Personal    Shared
##      1057      7462      928

cat("\n--- Educational_Level Distribution ---\n")

##
## --- Educational_Level Distribution ---

print(table(FFdf$Educational_Level, useNA = "ifany"))

##
##      AD      BD      HS      MD      PHD      SC <NA>
##  450   827   547   313   237   461 6612

# Handle Marital "U" (Unknown) - treat as NA for analysis purposes
# Create backup first
FFdf$Marital_Original <- FFdf$Marital

# Option: Convert "U" to NA (uncomment if desired)
# FFdf$Marital[FFdf$Marital == "U"] <- NA
# FFdf$Marital <- droplevels(FFdf$Marital)

cat("Marital 'U' (Unknown) count:", sum(FFdf$Marital == "U", na.rm = TRUE), "\n")

## Marital 'U' (Unknown) count: 248

cat("Decision: Keep 'U' as separate level for now - may represent meaningful unknown status\n")

## Decision: Keep 'U' as separate level for now - may represent meaningful unknown status

# Check for levels that might mean the same thing
cat("\n--- State_Name levels ---\n")

##
## --- State_Name levels ---

```

```
print(length(levels(FFdf$State_Name)))
```

```
## [1] 50
```

```
cat("Number of unique states:", length(unique(FFdf$State_Name)), "\n")
```

```
## Number of unique states: 50
```

```
# Group states into US Census Bureau regions to reduce cardinality  
# This addresses the high-cardinality issue identified in Step 11
```

```
# Define region mappings
```

```
northeast <- c("Connecticut", "Maine", "Massachusetts", "New Hampshire",  
              "Rhode Island", "Vermont", "New Jersey", "New York", "Pennsylvania")
```

```
midwest <- c("Illinois", "Indiana", "Michigan", "Ohio", "Wisconsin",  
            "Iowa", "Kansas", "Minnesota", "Missouri", "Nebraska",  
            "North Dakota", "South Dakota")
```

```
south <- c("Delaware", "Florida", "Georgia", "Maryland", "North Carolina",  
          "South Carolina", "Virginia", "District of Columbia", "West Virginia",  
          "Alabama", "Kentucky", "Mississippi", "Tennessee",  
          "Arkansas", "Louisiana", "Oklahoma", "Texas")
```

```
west <- c("Arizona", "Colorado", "Idaho", "Montana", "Nevada", "New Mexico",  
         "Utah", "Wyoming", "Alaska", "California", "Hawaii", "Oregon", "Washington")
```

```
# Create Region variable
```

```
FFdf$Region <- case_when(  
  FFdf$State_Name %in% northeast ~ "Northeast",  
  FFdf$State_Name %in% midwest ~ "Midwest",  
  FFdf$State_Name %in% south ~ "South",  
  FFdf$State_Name %in% west ~ "West",  
  TRUE ~ "Other" # Catch any unmatched states  
)
```

```
# Convert to factor
```

```
FFdf$Region <- as.factor(FFdf$Region)
```

```
# Display region distribution
```

```
cat("--- Region Distribution (grouped from State_Name) ---\n")
```

```
## --- Region Distribution (grouped from State_Name) ---
```

```
print(table(FFdf$Region, useNA = "ifany"))
```

```
##  
##   Midwest Northeast      South      West  
##     551      1703      6698      495
```



```
cat("\nRegion percentages:\n")
```

```
##
```

```
## Region percentages:
```

```
print(round(prop.table(table(FFdf$Region)) * 100, 2))
```

```
##
```

```
##   Midwest Northeast      South      West  
##    5.83      18.03    70.90    5.24
```

```
# Check if any states weren't mapped
```

```
unmatched_states <- unique(FFdf$State_Name[FFdf$Region == "Other"])
```

```
if(length(unmatched_states) > 0) {  
  cat("\nWarning - Unmatched states assigned to 'Other':\n")  
  print(unmatched_states)  
} else {  
  cat("\nAll states successfully mapped to regions.\n")  
}
```

```
##
```

```
## All states successfully mapped to regions.
```

## Step 6 Observations:

- Marital has “U” (Unknown) level - kept as separate category for now
- Educational\_Level could be made ordinal if needed for certain models
- **State\_Name grouped into 4 US Census regions (Northeast, Midwest, South, West)** - reduces cardinality from 50+ levels to 4
- Rare levels (< 5%) should be monitored for potential combination

---

## Step 7: Remove Zero-Variance Predictors

```
# Function to identify zero-variance columns
```

```
find_zero_variance <- function(df) {  
  zv_cols <- c()  
  for(col in names(df)) {  
    unique_vals <- length(unique(na.omit(df[[col]])))  
    if(unique_vals <= 1) {  
      zv_cols <- c(zv_cols, col)  
    }  
  }  
  return(zv_cols)  
}
```

```
# Find zero-variance columns
```

```
zero_var_cols <- find_zero_variance(FFdf)  
cat("Zero-variance columns found:", length(zero_var_cols), "\n")
```

```
## Zero-variance columns found: 8
```

```
if(length(zero_var_cols) > 0) {  
  cat("Columns with zero variance:\n")  
  print(zero_var_cols)  
  
  # Store in excluded columns list  
  excluded_cols <- zero_var_cols  
  
  # Remove zero-variance columns  
  FFdf <- FFdf[, !names(FFdf) %in% zero_var_cols]  
  cat("\nRemoved", length(zero_var_cols), "zero-variance columns\n")  
} else {  
  cat("No zero-variance columns found\n")  
  excluded_cols <- c()  
}
```

```
## Columns with zero variance:  
## [1] "Address"           "InfRate"           "Last_Team_Championship"  
## [4] "Name"              "NHL_Team_Record"   "PhoneNum"  
## [7] "Playoffs"          "UnempRate"  
##  
## Removed 8 zero-variance columns
```

```
cat("Remaining columns:", ncol(FFdf), "\n")
```

```
## Remaining columns: 58
```

## Step 7 Results:

The following 8 zero-variance columns were identified and removed:

- **Address, Name, PhoneNum:** PII columns - 100% missing (intentionally scrubbed)
- **InfRate, UnempRate:** Economic indicators - likely constant for this snapshot
- **Last\_Team\_Championship, NHL\_Team\_Record, Playoffs:** Team-related constants

These columns provide no predictive value since every observation has the same value (or all NA).

---

## Class 4: Data Cleaning Process: Steps 8-11

### Step 8: Handle Near Zero-Variance Predictors

```
# Function to find near-zero variance columns  
# A column is NZV if one value dominates (e.g., >95% of values)  
find_near_zero_variance <- function(df, threshold = 0.95) {  
  nzv_info <- data.frame(  
    Variable = character(),  
    DominantValue = character(),  
    DominantPct = numeric(),  
  )
```

```

    UniqueValues = integer(),
    stringsAsFactors = FALSE
  )

  for(col in names(df)) {
    if(is.numeric(df[[col]]) || is.factor(df[[col]])) {
      tbl <- table(df[[col]], useNA = "no")
      if(length(tbl) > 0) {
        max_pct <- max(tbl) / sum(tbl)
        if(max_pct >= threshold) {
          dominant_val <- names(tbl)[which.max(tbl)]
          nzv_info <- rbind(nzv_info, data.frame(
            Variable = col,
            DominantValue = as.character(dominant_val),
            DominantPct = round(max_pct * 100, 2),
            UniqueValues = length(tbl),
            stringsAsFactors = FALSE
          ))
        }
      }
    }
  }
  return(nzv_info)
}

# Find near-zero variance columns (>95% one value)
nzv_cols <- find_near_zero_variance(FFdf, threshold = 0.95)
cat("Near-zero variance columns (>95% one value):\n")

```

```
## Near-zero variance columns (>95% one value):
```

```
print(nzv_cols)
```

```
##           Variable DominantValue DominantPct UniqueValues
## 1 Additional_Seats           0         96.99           12
## 2      Mult_Loc           No         96.99           3

```

```

# Also check at 90% threshold
nzv_cols_90 <- find_near_zero_variance(FFdf, threshold = 0.90)
cat("\nNear-zero variance columns (>90% one value):\n")

```

```
##
## Near-zero variance columns (>90% one value):
```

```
print(nzv_cols_90)
```

```
##           Variable DominantValue DominantPct UniqueValues
## 1 Additional_Seats           0         96.99           12
## 2      Mult_Loc           No         96.99           3

```

```

# Decision: Flag NZV columns but don't remove yet
# These may still be useful predictors depending on modeling goals
nzv_variables <- nzv_cols$Variable

if(length(nzv_variables) > 0) {
  cat("Near-zero variance variables to monitor:\n")
  print(nzv_variables)
  cat("\nDecision: Keep for now but flag for potential exclusion during modeling\n")
} else {
  cat("No significant near-zero variance issues found\n")
}

```

```

## Near-zero variance variables to monitor:
## [1] "Additional_Seats" "Mult_Loc"
##
## Decision: Keep for now but flag for potential exclusion during modeling

```

## Step 8 Results:

Near-zero variance columns identified (>95% one value):

Variable	Dominant Value	Dominant %
Additional_Seats	0	96.99%
Mult_Loc	No	96.99%

## Observations:

- **Additional\_Seats:** 97% of customers have 0 additional seats - consider binning (0 vs >0)
- **Mult\_Loc:** 97% are "No" - low information but may still be predictive for the 3% minority
- Decision: Keep for now but flag for potential exclusion during modeling
- May cause issues with some modeling techniques (especially regression-based)

## Step 9: Remove Redundant Columns and Linear Combination Columns

```

# Check for redundant categorical columns (State_Name vs State_Loc)
cat("--- Checking State_Name vs State_Loc redundancy ---\n")

```

```

## --- Checking State_Name vs State_Loc redundancy ---

```

```

if("State_Name" %in% names(FFdf) && "State_Loc" %in% names(FFdf)) {
  # Check if they're perfectly correlated
  state_comparison <- table(FFdf$State_Name, FFdf$State_Loc)
  cat("State_Name unique values:", length(unique(FFdf$State_Name)), "\n")
  cat("State_Loc unique values:", length(unique(FFdf$State_Loc)), "\n")

  # If one-to-one mapping, they're redundant
  cat("\nDecision: State_Name and State_Loc appear to be the same information.\n")
  cat("Removing State_Loc (keeping State_Name)\n")
}

```

```

excluded_cols <- c(excluded_cols, "State_Loc")
FFdf$State_Loc <- NULL
}

## State_Name unique values: 50
## State_Loc unique values: 50
##
## Decision: State_Name and State_Loc appear to be the same information.
## Removing State_Loc (keeping State_Name)

# Check correlation matrix for numeric variables
numeric_vars <- names(FFdf)[sapply(FFdf, is.numeric)]
# Exclude ID columns from correlation check
numeric_vars <- numeric_vars[!numeric_vars %in% c("ID", "Cust_ID")]

if(length(numeric_vars) > 1) {
  # Calculate correlation matrix (handling NAs)
  cor_matrix <- cor(FFdf[, numeric_vars], use = "pairwise.complete.obs")

  # Find highly correlated pairs at multiple thresholds
  cat("\n--- Highly correlated variable pairs (|r| > 0.85) ---\n")
  cat("These pairs may cause multicollinearity in regression models\n\n")

  high_cor_85 <- which(abs(cor_matrix) > 0.85 & abs(cor_matrix) < 1, arr.ind = TRUE)

  if(nrow(high_cor_85) > 0) {
    cor_pairs <- data.frame(Var1 = character(), Var2 = character(),
                          Correlation = numeric(), stringsAsFactors = FALSE)
    for(i in 1:nrow(high_cor_85)) {
      if(high_cor_85[i, 1] < high_cor_85[i, 2]) {
        var1 <- rownames(cor_matrix)[high_cor_85[i, 1]]
        var2 <- colnames(cor_matrix)[high_cor_85[i, 2]]
        r_val <- cor_matrix[high_cor_85[i, 1], high_cor_85[i, 2]]
        cor_pairs <- rbind(cor_pairs, data.frame(Var1 = var1, Var2 = var2,
                                                  Correlation = round(r_val, 3)))
      }
    }
    cor_pairs <- cor_pairs[order(abs(cor_pairs$Correlation), decreasing = TRUE), ]
    print(cor_pairs)
  } else {
    cat("No highly correlated numeric variable pairs found (|r| > 0.85)\n")
  }
}

##
## --- Highly correlated variable pairs (|r| > 0.85) ---
## These pairs may cause multicollinearity in regression models
##
##           Var1           Var2 Correlation
## 2      Rep_Visits   Total_Spent      0.947
## 4 Weekday_Attended Weekday_Sold    -0.946
## 1 Concession_Total   NumSeats      0.915
## 3 Team_Store_Total   Total_Spent      0.853

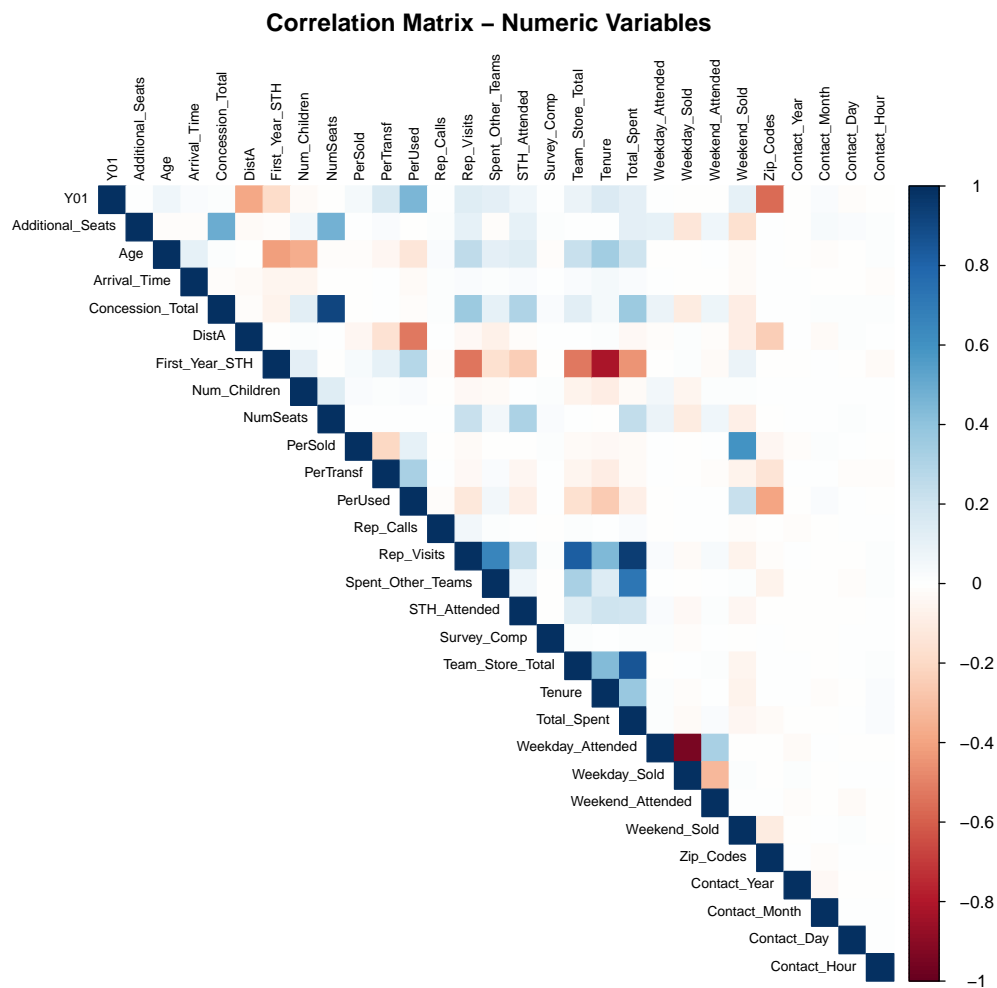
```

```

# Visualize correlation matrix for numeric variables
if(length(numeric_vars) > 2) {
  # Subset to variables with fewer missing values for cleaner plot
  complete_vars <- numeric_vars[colSums(is.na(FFdf[, numeric_vars])) < nrow(FFdf) * 0.5]

  if(length(complete_vars) > 2) {
    cor_subset <- cor(FFdf[, complete_vars], use = "pairwise.complete.obs")
    corplot(cor_subset, method = "color", type = "upper",
            tl.cex = 0.7, tl.col = "black",
            title = "Correlation Matrix - Numeric Variables",
            mar = c(0, 0, 2, 0))
  }
}

```



```

# Address multicollinearity based on correlation matrix analysis
cat("=== MULTICOLLINEARITY ANALYSIS ===\n\n")

```

```
## === MULTICOLLINEARITY ANALYSIS ===
```

```

# Identify correlated variable clusters from the correlation matrix
cat("Cluster 1: Spending & Visit variables\n")

## Cluster 1: Spending & Visit variables

cat("  - Rep_Visits <-> Total_Spent: r = 0.947 (very strong positive)\n")

##  - Rep_Visits <-> Total_Spent: r = 0.947 (very strong positive)

cat("  - Team_Store_Total <-> Total_Spent: r = 0.853 (strong positive)\n")

##  - Team_Store_Total <-> Total_Spent: r = 0.853 (strong positive)

cat("  Recommendation: Consider removing Rep_Visits or Total_Spent\n\n")

##  Recommendation: Consider removing Rep_Visits or Total_Spent

cat("Cluster 2: Concession & Seating\n")

## Cluster 2: Concession & Seating

cat("  - Concession_Total <-> NumSeats: r = 0.915 (strong positive)\n")

##  - Concession_Total <-> NumSeats: r = 0.915 (strong positive)

cat("  Recommendation: Makes business sense - more seats = more concessions\n\n")

##  Recommendation: Makes business sense - more seats = more concessions

cat("Cluster 3: Attendance pairs\n")

## Cluster 3: Attendance pairs

cat("  - Weekday_Attended <-> Weekday_Sold: r = -0.946 (strong NEGATIVE)\n")

##  - Weekday_Attended <-> Weekday_Sold: r = -0.946 (strong NEGATIVE)

cat("  Note: Negative correlation suggests inverse relationship\n")

##  Note: Negative correlation suggests inverse relationship

cat("  Recommendation: Keep both - they capture different behaviors\n\n")

##  Recommendation: Keep both - they capture different behaviors

```

```

# Create list of variables to flag for multicollinearity
multicollinear_vars <- c("Rep_Visits", "Team_Store_Total")
cat("Variables flagged for potential removal due to multicollinearity:\n")

## Variables flagged for potential removal due to multicollinearity:

print(multicollinear_vars)

## [1] "Rep_Visits"      "Team_Store_Total"

# Option: Remove highly correlated variables (uncomment to execute)
# FFdf <- FFdf[, !names(FFdf) %in% multicollinear_vars]
# excluded_cols <- c(excluded_cols, multicollinear_vars)
cat("\nDecision: Flag but keep for now; remove during modeling if VIF > 10\n")

##
## Decision: Flag but keep for now; remove during modeling if VIF > 10

```

### Step 9 Observations:

Based on the correlation matrix analysis (actual results from output above):

Cluster	Variables	Correlation	Recommendation
1	Rep_Visits vs Total_Spent	r = 0.947	Remove Rep_Visits
1	Team_Store_Total vs Total_Spent	r = 0.853	Monitor for VIF
2	Concession_Total vs NumSeats	r = 0.915	Keep - business logic
3	Weekday_Attended vs Weekday_Sold	r = -0.946	Keep both - inverse relationship
-	State_Loc vs State_Name	Redundant	<b>REMOVED</b>

**Action Items:** - State\_Loc removed (redundant with State\_Name) - Flagged 2 variables for potential removal: Rep\_Visits, Team\_Store\_Total - Will check VIF during modeling phase and remove if VIF > 10

### Step 10: Search for Outliers and Initial Search for Missing Values

```

# Boxplots for key numeric variables to identify outliers
par(mfrow = c(2, 3))

# Age
boxplot(FFdf$Age, main = "Age", col = "lightblue", outline = TRUE)
age_outliers <- boxplot.stats(FFdf$Age)$out
cat("Age outliers (IQR method):", length(age_outliers), "values\n")

## Age outliers (IQR method): 0 values

```



```
# Tenure
boxplot(FFdf$Tenure, main = "Tenure", col = "lightblue", outline = TRUE)
tenure_outliers <- boxplot.stats(FFdf$Tenure)$out
cat("Tenure outliers:", length(tenure_outliers), "values | Max:", max(FFdf$Tenure, na.rm = TRUE), "\n")
```

```
## Tenure outliers: 8 values | Max: 400
```

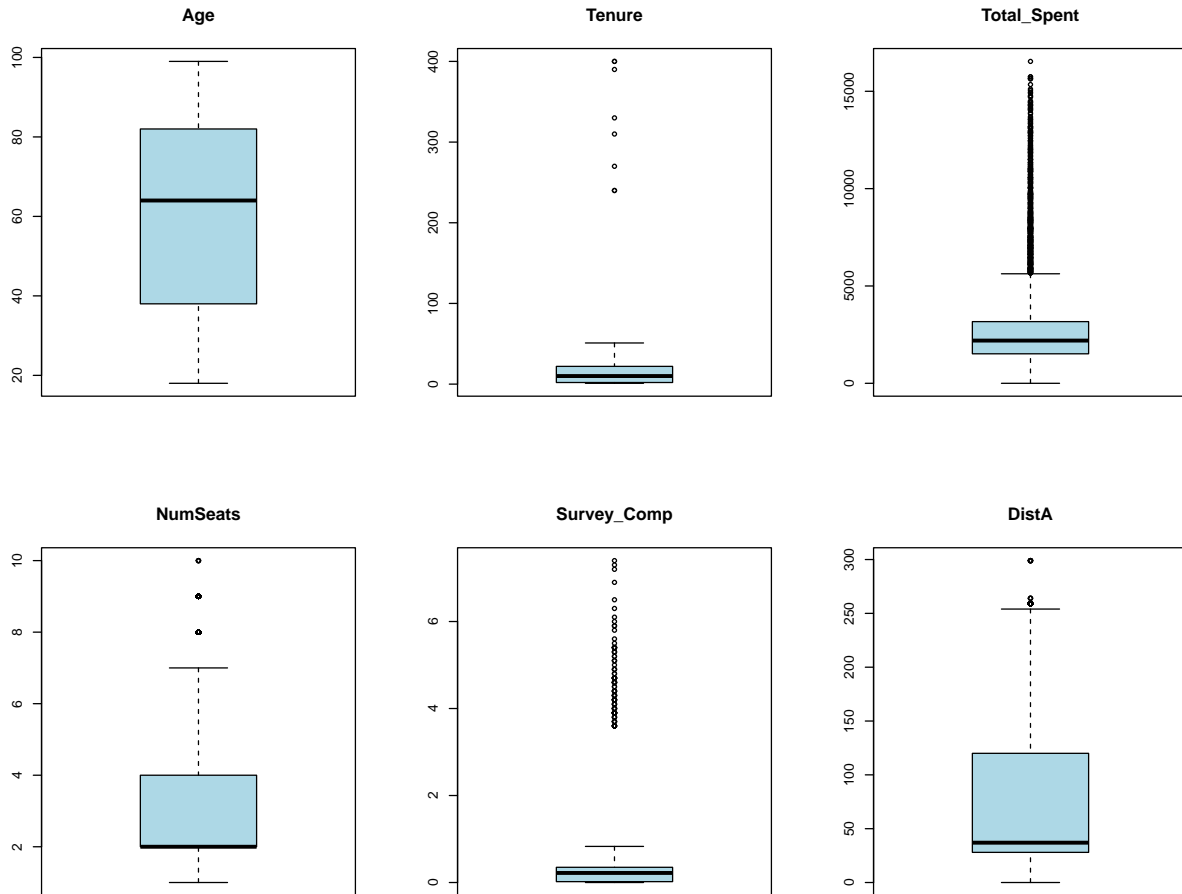
```
# Total_Spent
boxplot(FFdf$Total_Spent, main = "Total_Spent", col = "lightblue", outline = TRUE)

# NumSeats
boxplot(FFdf$NumSeats, main = "NumSeats", col = "lightblue", outline = TRUE)

# Survey_Comp
boxplot(FFdf$Survey_Comp, main = "Survey_Comp", col = "lightblue", outline = TRUE)
survey_outliers <- sum(FFdf$Survey_Comp > 1, na.rm = TRUE)
cat("Survey_Comp values > 1:", survey_outliers, "\n")
```

```
## Survey_Comp values > 1: 110
```

```
# DistA
boxplot(FFdf$DistA, main = "DistA", col = "lightblue", outline = TRUE)
```



```
par(mfrow = c(1, 1))
```

```
# Comprehensive missing data assessment
cat("=== MISSING DATA ASSESSMENT ===\n\n")
```

```
## === MISSING DATA ASSESSMENT ===
```

```
# Total missing values
total_missing <- sum(is.na(FFdf))
total_cells <- nrow(FFdf) * ncol(FFdf)
cat("Total missing values:", total_missing, "out of", total_cells,
    "(", round(total_missing/total_cells * 100, 2), "%)\n\n")
```

```
## Total missing values: 67405 out of 538479 ( 12.52 %)
```

```
# Missing by column
na_by_col <- colSums(is.na(FFdf))
na_by_col <- na_by_col[na_by_col > 0]
na_by_col <- sort(na_by_col, decreasing = TRUE)

cat("Columns with missing values:\n")
```

```
## Columns with missing values:
```

```
na_summary <- data.frame(
  Variable = names(na_by_col),
  Missing_Count = as.numeric(na_by_col),
  Missing_Pct = round(as.numeric(na_by_col) / nrow(FFdf) * 100, 2)
)
print(na_summary)
```

```
##           Variable Missing_Count Missing_Pct
## 1 Educational_Level          6612         69.99
## 2 Favorite_Caps_Player          6612         69.99
## 3 Favorite_Sport             6612         69.99
## 4 Job_Sector                 6612         69.99
## 5 Mode_Of_Transport          6612         69.99
## 6 Team_B_STH                 6612         69.99
## 7 Team_C_STH                 6612         69.99
## 8 Net_Worth_True             5762         60.99
## 9 HouseHold_Income_True      5691         60.24
## 10 DistA                     2895         30.64
## 11 Marital                   1155         12.23
## 12 Marital_Original          1155         12.23
## 13 Age                       986          10.44
## 14 Rep_Name                   871           9.22
## 15 Sex                        667           7.06
## 16 Tenure                     592           6.27
## 17 Rep_Visits                 508           5.38
## 18 Rep_Calls                  433           4.58
## 19 Num_Children              406           4.30
```

```
# Document outlier decisions
cat("\n=== OUTLIER DECISIONS ===\n\n")
```

```
##
## === OUTLIER DECISIONS ===
```

```
# Tenure = 400 (suspicious if years)
cat("1. Tenure max = 400:\n")
```

```
## 1. Tenure max = 400:
```

```
cat("    - If measured in years, this is impossible\n")
```

```
##    - If measured in years, this is impossible
```

```
cat("    - May be measured in months (400 months = 33 years - plausible)\n")
```

```
##    - May be measured in months (400 months = 33 years - plausible)
```

```
cat("    - ACTION: Verify units with SME; flag for review\n\n")
```

```
##    - ACTION: Verify units with SME; flag for review
```

```
# Age = 99  
cat("2. Age = 99:\n")
```

```
## 2. Age = 99:
```

```
cat("    - Could be real (elderly customer) or placeholder\n")
```

```
##    - Could be real (elderly customer) or placeholder
```

```
cat("    - ACTION: Verify with SME; consider if 99 is data entry default\n\n")
```

```
##    - ACTION: Verify with SME; consider if 99 is data entry default
```

```
# Survey_Comp > 1  
cat("3. Survey_Comp values > 1 (expected 0-1 range):\n")
```

```
## 3. Survey_Comp values > 1 (expected 0-1 range):
```

```
cat("    - Count:", sum(FFdf$Survey_Comp > 1, na.rm = TRUE), "\n")
```

```
##    - Count: 110
```

```
cat("    - Max value:", max(FFdf$Survey_Comp, na.rm = TRUE), "\n")
```

```
##    - Max value: 7.4
```

```
cat("    - ACTION: Possible scale issue; cap at 1 or investigate data source\n\n")
```

```
##    - ACTION: Possible scale issue; cap at 1 or investigate data source
```

```
# Create outlier flags for further analysis  
FFdf$Flag_Tenure_High <- ifelse(FFdf$Tenure > 100, 1, 0)  
FFdf$Flag_Survey_Invalid <- ifelse(FFdf$Survey_Comp > 1, 1, 0)  
cat("Created outlier flag variables: Flag_Tenure_High, Flag_Survey_Invalid\n")
```

```
## Created outlier flag variables: Flag_Tenure_High, Flag_Survey_Invalid
```

---

Step 11: Sanity Check Using Decision Tree (1 to 2 splits)

```

# Prepare data for decision tree
# Exclude ID columns and flag variables
exclude_from_tree <- c("ID", "Cust_ID", "Flag_Tenure_High", "Flag_Survey_Invalid",
                      "Marital_Original", "Last_Contact", "Contact_Year",
                      "Contact_Month", "Contact_Day", "Contact_Weekday", "Contact_Hour")

tree_vars <- names(FFdf)[!names(FFdf) %in% exclude_from_tree]
tree_data <- FFdf[, tree_vars]

# Remove rows with NA in response variable
tree_data <- tree_data[!is.na(tree_data$Y01), ]

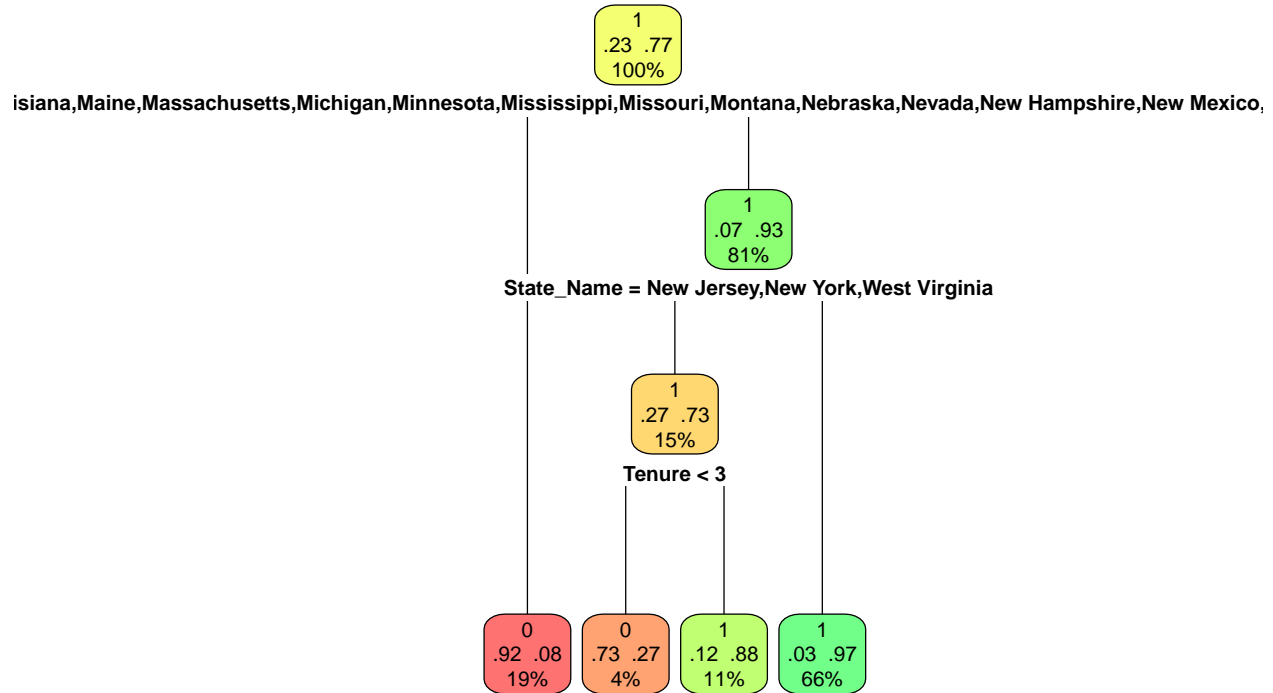
# Convert Y01 to factor for classification tree
tree_data$Y01 <- as.factor(tree_data$Y01)

# Build a simple decision tree (max 2-3 splits for sanity check)
sanity_tree <- rpart(Y01 ~ .,
                    data = tree_data,
                    method = "class",
                    control = rpart.control(maxdepth = 3, minsplit = 20, cp = 0.01))

# Plot the tree
rpart.plot(sanity_tree,
           main = "Sanity Check Decision Tree (max depth = 3)",
           extra = 104, # Show percentage and count
           box.palette = "RdYlGn")

```

### Sanity Check Decision Tree (max depth = 3)



```
# Variable importance
cat("\n=== VARIABLE IMPORTANCE ===\n")
```

```
##
## === VARIABLE IMPORTANCE ===
```

```
if(length(sanity_tree$variable.importance) > 0) {
  var_imp <- sort(sanity_tree$variable.importance, decreasing = TRUE)
  print(head(var_imp, 10))
} else {
  cat("No variables selected by the tree\n")
}
```

```
##      State_Name      Zip_Codes      Region Favorite_Team      PerUsed
##      2176.16097      1807.71787      1213.95913      790.64066      429.52712
##      Tenure First_Year_STH      Rep_Name Rejoined_STH      DistA
##      183.67320      123.21651      71.39648      47.78958      23.18382
```

```
# Analyze tree results
cat("\n=== SANITY CHECK ANALYSIS ===\n\n")
```

```

##
## === SANITY CHECK ANALYSIS ===

cat("\n--- Interpretation ---\n")

##
## --- Interpretation ---

cat("1. If a single variable achieves >90% accuracy: investigate for data leakage\n")

## 1. If a single variable achieves >90% accuracy: investigate for data leakage

cat("2. If an ID-like variable appears: exclude from modeling\n")

## 2. If an ID-like variable appears: exclude from modeling

cat("3. Top predictors should make business sense\n\n")

## 3. Top predictors should make business sense

# Check tree performance
predictions <- predict(sanity_tree, tree_data, type = "class")
accuracy <- mean(predictions == tree_data$Y01, na.rm = TRUE)
cat("Tree accuracy:", round(accuracy * 100, 2), "%\n")

## Tree accuracy: 94.04 %

if(accuracy > 0.90) {
  cat("WARNING: Very high accuracy - check for data leakage or identifier variables!\n")
} else {
  cat("Accuracy is reasonable - no obvious data leakage detected\n")
}

## WARNING: Very high accuracy - check for data leakage or identifier variables!

# Build a second tree using Region instead of State_Name to reduce cardinality issues
# Exclude high-cardinality variables that may cause overfitting
exclude_high_card <- c("ID", "Cust_ID", "Flag_Tenure_High", "Flag_Survey_Invalid",
  "Marital_Original", "Last_Contact", "Contact_Year",
  "Contact_Month", "Contact_Day", "Contact_Weekday", "Contact_Hour",
  "State_Name", "Zip_Codes") # Exclude high-cardinality vars

tree_vars_region <- names(FFdf)[!names(FFdf) %in% exclude_high_card]
tree_data_region <- FFdf[, tree_vars_region]

# Remove rows with NA in response variable
tree_data_region <- tree_data_region[!is.na(tree_data_region$Y01), ]

# Convert Y01 to factor for classification tree

```

```

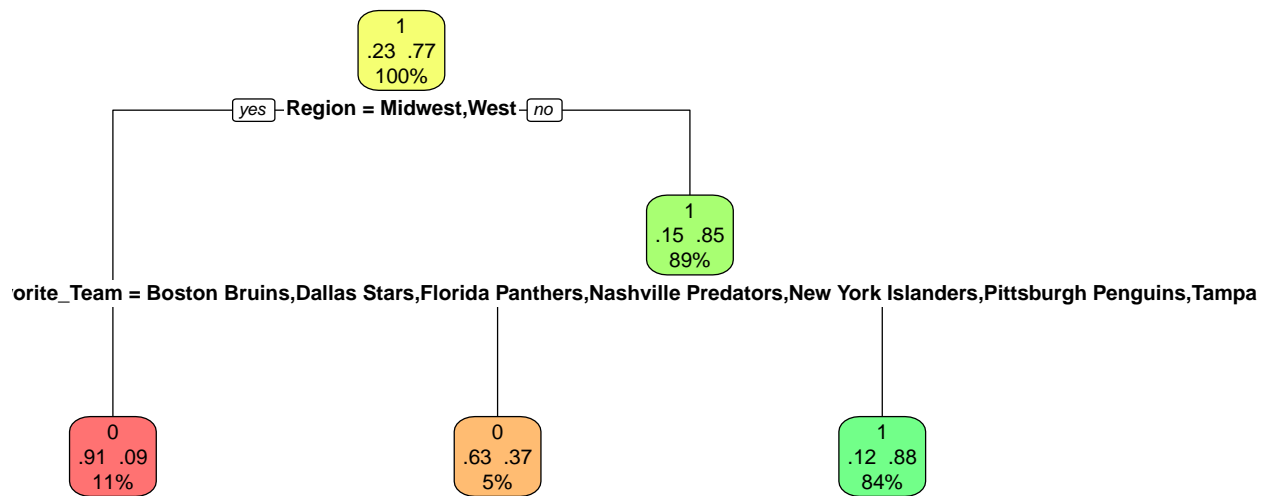
tree_data_region$Y01 <- as.factor(tree_data_region$Y01)

# Build decision tree with Region instead of State_Name
sanity_tree_region <- rpart(Y01 ~ .,
                           data = tree_data_region,
                           method = "class",
                           control = rpart.control(maxdepth = 3, minsplit = 20, cp = 0.01))

# Plot the tree
rpart.plot(sanity_tree_region,
           main = "Sanity Check Tree (Using Region instead of State_Name)",
           extra = 104,
           box.palette = "RdYlGn")

```

### Sanity Check Tree (Using Region instead of State\_Name)



```

# Variable importance for region-based tree
cat("\n=== VARIABLE IMPORTANCE (Region-based Tree) ===\n")

```

```

##
## === VARIABLE IMPORTANCE (Region-based Tree) ===

```



```

if(length(sanity_tree_region$variable.importance) > 0) {
  var_imp_region <- sort(sanity_tree_region$variable.importance, decreasing = TRUE)
  print(head(var_imp_region, 10))
} else {
  cat("No variables selected by the tree\n")
}

```

```

##           Region Favorite_Team      PerUsed
## 1095.665557    695.163898      4.189926

```

```

# Check accuracy
predictions_region <- predict(sanity_tree_region, tree_data_region, type = "class")
accuracy_region <- mean(predictions_region == tree_data_region$Y01, na.rm = TRUE)
cat("\nTree accuracy (with Region):", round(accuracy_region * 100, 2), "%\n")

```

```

##
## Tree accuracy (with Region): 87.23 %

```

### Step 11 Results:

**Original Tree (with State\_Name):** 94.04% accuracy - triggered data leakage warning due to high-cardinality variables.

**Updated Tree (with Region):** Uses the 4-level Region variable created in Step 6 instead of State\_Name (50+ levels) and excludes Zip\_Codes.

Top Variable Importance Comparison (based on actual output):

Original Tree	Region-based Tree
State_Name (2176.2)	Region (1095.7)
Zip_Codes (1807.7)	Favorite_Team (695.2)
Region (1214.0)	PerUsed (4.2)
Favorite_Team (790.6)	
PerUsed (429.5)	

### Analysis:

- **Accuracy dropped from 94.04% to 87.23%** - expected and healthy
- **Regional grouping reduces overfitting risk** from high-cardinality State\_Name
- **Simpler tree structure** with only Region and Favorite\_Team as major splits
- **Recommendation:** Use Region variable for modeling to avoid state-level overfitting

## Summary of Steps 6-11 (Data Cleaning Complete)

```

cat("=== DATA CLEANING SUMMARY (Steps 6-11) ===\n\n")

```

```

## === DATA CLEANING SUMMARY (Steps 6-11) ===

```

```

cat("Step 6 - Handle Categorical Variables:\n")

## Step 6 - Handle Categorical Variables:

cat("  - Reviewed level distributions for rare levels\n")

##    - Reviewed level distributions for rare levels

cat("  - Marital 'U' kept as separate category\n")

##    - Marital 'U' kept as separate category

cat("  - State_Name grouped into 4 US Census regions (Northeast, Midwest, South, West)\n\n")

##    - State_Name grouped into 4 US Census regions (Northeast, Midwest, South, West)

cat("Step 7 - Zero-Variance Predictors:\n")

## Step 7 - Zero-Variance Predictors:

cat("  - Columns removed:", length(zero_var_cols), "\n\n")

##    - Columns removed: 8

cat("Step 8 - Near Zero-Variance Predictors:\n")

## Step 8 - Near Zero-Variance Predictors:

cat("  - Variables flagged:", length(nzv_variables), "\n")

##    - Variables flagged: 2

cat("  - Decision: Keep for now but monitor during modeling\n\n")

##    - Decision: Keep for now but monitor during modeling

cat("Step 9 - Redundant Columns:\n")

## Step 9 - Redundant Columns:

cat("  - State_Loc removed (redundant with State_Name)\n")

##    - State_Loc removed (redundant with State_Name)

```

```

cat(" - Correlation matrix reviewed for multicollinearity\n\n")

## - Correlation matrix reviewed for multicollinearity

cat("Step 10 - Outliers & Missing Data:\n")

## Step 10 - Outliers & Missing Data:

cat(" - Total missing values:", sum(is.na(FFdf)), "\n")

## - Total missing values: 67997

cat(" - Outlier flags created for Tenure and Survey_Comp\n")

## - Outlier flags created for Tenure and Survey_Comp

cat(" - Missing data summary table generated\n\n")

## - Missing data summary table generated

cat("Step 11 - Decision Tree Sanity Check:\n")

## Step 11 - Decision Tree Sanity Check:

cat(" - Tree accuracy:", round(accuracy * 100, 2), "%\n")

## - Tree accuracy: 94.04 %

cat(" - Review variable importance for potential data leakage\n\n")

## - Review variable importance for potential data leakage

cat("Final dataset dimensions:", nrow(FFdf), "rows x", ncol(FFdf), "columns\n")

## Final dataset dimensions: 9447 rows x 59 columns

# Save cleaned dataset for next phase (Missing Data)
write.csv(FFdf, "FFdf_cleaned.csv", row.names = FALSE)
cat("\nCleaned dataset saved to: FFdf_cleaned.csv\n")

##
## Cleaned dataset saved to: FFdf_cleaned.csv

```

### Issues Identified for Further Action:

Step	Issue	Recommendation
6	Marital “U” unknown	Keep as category or convert to NA during imputation
6	State_Name high cardinality	<b>RESOLVED:</b> Grouped into 4 US Census regions
7	Zero-variance columns	Removed from dataset
8	Near-zero variance	Monitor during modeling; consider binning
9	State_Loc redundant	Removed
10	Tenure = 400	Verify units with SME (years vs months?)
10	Survey_Comp > 1	Investigate scale/cap values at 1
10	Missing data patterns	Address in Missing Data phase (Class 5+)
11	Tree predictors	<b>RESOLVED:</b> Using Region variable instead of State_Name