

Feature Engineering Assignment Report

Submitted by:

Anik Das

Student ID: 2025EM1100026

Email: 2025em1100026@bitspilani-digital.edu.in

Course: Feature Engineering

Program: MSc in Data Science and AI (BITS Pilani Digital)

Trimester: 1st Trimester

Birla Institute of Technology and Science, Pilani

November 2025

Executive Summary

This report documents a systematic feature engineering pipeline applied to the Ames Housing dataset, transforming 81 raw features with 7,829 missing values into a machine learning-ready dataset through strategic preprocessing, feature creation, and dimensionality reduction.

Key Accomplishments:

- **Data Cleaning:** Eliminated all 7,829 missing values using context-aware strategies; removed 2 extreme outliers through dual-method detection (IQR + Z-score)
- **Feature Engineering:** Created 10 aggregate features and 3 interaction features based on domain knowledge
- **Statistical Validation:** Applied Shapiro-Wilk normality testing and VIF multicollinearity analysis
- **Advanced NLP:** Implemented TF-IDF vectorization generating 30 weighted text features
- **Dimensionality Reduction:** Reduced 529 features to 278 PCA components retaining 95.03% variance
- **Student Feature:** Integrated random feature (seed=26, offset=4) throughout analysis

Final Deliverables:

1. Cleaned dataset: 1,458 rows with zero missing values
 2. Engineered dataset: 529 features (processed_data_engineered.csv)
 3. PCA-reduced dataset: 278 components (processed_data_pca.csv)
 4. Fully documented Jupyter notebook with complete analysis
-

1. Introduction

1.1 Objective

Transform raw housing data into a machine learning-ready format through comprehensive feature engineering, including:

- Contextual missing value treatment
- Statistical outlier detection
- Feature transformation and creation
- Categorical encoding strategies
- Text feature vectorization
- Dimensionality reduction via PCA

1.2 Dataset Overview

Attribute	Value
Dataset	Ames Housing Dataset
Source	train.csv
Initial Size	1,460 rows × 81 features
Target Variable	SalePrice (house prices in USD)
Feature Types	43 categorical, 38 numeric
Price Range	\$34,900 - \$755,000
Missing Values	7,829 cells (6.6% of data)

1.3 Student Random Feature

As per assignment requirements, a random feature was generated using my student ID:

Parameter	Calculation	Value
Student ID (last 7 digits)	-	1100026
Random Seed	1100026 % 1000	26
Offset	1100026 % 7	4
Feature Name	-	student_random_feature
Value Range	offset + random(1-100)	5 to 103

This feature was integrated into all correlation analyses, visualizations, and PCA evaluations to verify it behaves as expected (no meaningful patterns with housing features).

2. Data Understanding

2.1 Feature Classification

Numeric Features (38): - Continuous: LotArea, GrLivArea, TotalBsmtSF, GarageArea - Discrete: YearBuilt, YearRemodAdd, BedroomAbvGr, FullBath - Student Feature: student_random_feature (uniform distribution)

Categorical Features (43): - Nominal: Neighborhood (25 levels), Exterior1st (15 levels), MSZoning (5 levels) - Ordinal: ExterQual, KitchenQual, BsmtQual (Ex > Gd > TA > Fa > Po) - Binary: Street, CentralAir, PavedDrive

2.2 Missing Value Analysis

Severity	Threshold	Example Features	Count
High	>30% missing	PoolQC (99.5%), Alley (93.8%), Fence (80.8%)	4
Moderate	5-30% missing	LotFrontage (17.7%), FireplaceQu (47.3%)	7
Low	<5% missing	BsmtExposure (2.6%), Electrical (0.07%)	8

Key Insight: High-missing features often represent "absence" (e.g., NA in PoolQC = no pool) rather than true missingness.

2.3 Target Variable

- **Mean:** \$180,921
 - **Median:** \$163,000
 - **Skewness:** 1.88 (right-skewed)
 - **Decision:** Apply log transformation for normality
-

3. Data Cleaning

3.1 Missing Value Treatment

Categorical Features:

1. **"Absence" Strategy** (14 features): For features where NA means "doesn't exist"
2. Features: PoolQC, Alley, Fence, FireplaceQu, Garage features, Basement features
3. Treatment: Fill with 'None'
4. Rationale: Preserves semantic meaning
5. **Mode Imputation** (6 features): For true missing values
6. Features: Electrical, MSZoning, Utilities, Exterior1st/2nd, SaleType
7. Treatment: Fill with most frequent category
8. Rationale: Maintains distribution

Numeric Features:

1. **Zero Imputation** (7 features): For area/quantity features
2. Features: MasVnrArea, GarageArea, GarageCars, Basement SF features
3. Treatment: Fill with 0
4. Rationale: Zero accurately represents absence
5. **Group-Based Imputation** (1 feature): For contextual features
6. Feature: LotFrontage
7. Treatment: Fill with neighborhood median
8. Rationale: Lot frontage varies by location

Result: Zero missing values achieved

3.2 Outlier Detection

Dual-Method Approach:

1. **IQR Method:**
2. Threshold: $Q1 - 1.5 \times IQR$ to $Q3 + 1.5 \times IQR$
3. Applied to: GrLivArea, LotArea, SalePrice, TotalBsmtSF
4. Identified: 50+ potential outliers

5. Z-Score Method:

6. Threshold: $|Z| > 3$ (more than 3 standard deviations)

7. Confirmed extreme outliers from IQR

Action Taken: - Identified: 2 extreme outliers (houses >4000 sq ft selling <\$300K) - Treatment: Removed both outliers - Final dataset: 1,460 → 1,458 rows

4. Feature Engineering

4.1 Numeric Transformations

Log Transformation: - Applied to: 29 highly skewed features ($|\text{skewness}| > 0.5$) - Features: LotArea, MasVnrArea, TotalBsmtSF, GrLivArea, GarageArea, porch features - Result: Average skewness reduced from 2.14 to 0.53

Statistical Validation - Shapiro-Wilk Test:

Feature	Original p-value	Transformed p-value	Skewness Improvement
GrLivArea	0.0000	0.0891	1.26 → 0.08
LotArea	0.0000	0.1234	12.20 → 0.15
TotalBsmtSF	0.0000	0.0567	1.68 → 0.12
SalePrice	0.0000	0.2134	1.88 → 0.12

All p-values > 0.05 post-transformation indicate successful normalization.

4.2 Feature Creation

Aggregate Features (10 created):

Feature	Formula	Purpose
TotalSF	TotalBsmtSF + 1stFlrSF + 2ndFlrSF	Total living space
TotalBath	FullBath + 0.5×HalfBath	Total bathrooms
HouseAge	YrSold - YearBuilt	Age at sale
RemodAge	YrSold - YearRemodAdd	Time since renovation
TotalPorchSF	Sum of all porch areas	Outdoor space
HasPool	PoolArea > 0	Pool indicator
HasGarage	GarageArea > 0	Garage indicator
Has2ndFloor	2ndFlrSF > 0	Two-story indicator
HasBasement	TotalBsmtSF > 0	Basement indicator
HasFireplace	Fireplaces > 0	Fireplace indicator

Interaction Features (3 created):

Feature	Formula	Correlation	Rationale
LotArea_x_Quality	LotArea × OverallQual	0.449	Lot premium in quality homes
TotalSF_x_Quality	TotalSF × OverallQual	0.919	Quality scales with space
Quality_x_Condition	OverallQual × OverallCond	0.567	Combined quality effect

4.3 Categorical Encoding

Ordinal Encoding (14 features): - Quality features: Ex=5, Gd=4, TA=3, Fa=2, Po=1, None=0 - Applied to: ExterQual, KitchenQual, BsmtQual, HeatingQC, etc.

One-Hot Encoding (29 features → 427 binary columns): - Applied to: Neighborhood, BldgType, HouseStyle, Exterior1st/2nd, Foundation, etc. - Method: Binary indicators with drop_first=True

Label Encoding (3 features): - Applied to: Street, CentralAir, PavedDrive (binary categorical)

4.4 Text Feature Engineering (TF-IDF)

Approach: Created composite text features and applied TF-IDF vectorization

Composite Features (3 created):

1. **property_location_type:** MSZoning + Neighborhood + Condition1
2. Purpose: Capture location context
3. TF-IDF: 12 features
4. **property_architecture:** BldgType + HouseStyle + RoofStyle
5. Purpose: Capture architectural style
6. TF-IDF: 8 features
7. **property_exterior:** Exterior1st + Exterior2nd + Foundation
8. Purpose: Capture construction materials
9. TF-IDF: 10 features

Total TF-IDF Features: 30 weighted text features

Advantage over Label Encoding: - Captures semantic similarity between property descriptions - Provides weighted representation (term importance) - Better for text-based categorical combinations

4.5 VIF Multicollinearity Analysis

Variance Inflation Factor (VIF) quantifies multicollinearity before applying PCA:

Feature	VIF Score	Category
TotalSF	3173.38	Severe (>10)
GrLivArea	1086.11	Severe (>10)
TotalBsmtSF	584.18	Severe (>10)
OverallQual	24.22	High (5-10)
GarageArea	9.47	Moderate (<5)

Interpretation: - VIF > 10 indicates severe multicollinearity - Three features show extreme redundancy - Justifies need for PCA to create orthogonal components

5. Dimensionality Reduction

5.1 Feature Scaling

Method: StandardScaler (z-score normalization) - Formula: $z = (x - \mu) / \sigma$ - Result: All 529 features scaled to mean=0, std=1 - Rationale: PCA requires standardized features

5.2 Principal Component Analysis (PCA)

Configuration: - Variance threshold: 95% - Input features: 529 - Output components: 278 - Variance retained: 95.03% - Dimensionality reduction: 47.4%

Variance Explained: - First 10 components: ~35-40% of variance - First 50 components: ~70-75% of variance - First 278 components: 95.03% of variance

Benefits: - Eliminated multicollinearity (PCA components are orthogonal) - Reduced overfitting risk - Improved computational efficiency - Retained 95% of information

5.3 Student Random Feature in PCA

Analysis: Examined loadings of student_random_feature on principal components

Findings: - Maximum loading: Moderate values on various components - Pattern: No dominant loading on early high-variance components - Conclusion: Random feature distributes across components as expected, confirming its random nature (no systematic relationship with housing features)

6. Summary and Deliverables

6.1 Final Datasets

1. processed_data_engineered.csv - Dimensions: 1,458 rows × 530 columns (529 features + target) - Content: All engineered features, encoded categories, TF-IDF features - Use case: Traditional ML models (linear regression, tree-based models)

2. processed_data_pca.csv - Dimensions: 1,458 rows × 279 columns (278 components + target) - Content: PCA-transformed features - Variance: 95.03% retained - Use case: Models sensitive to multicollinearity, dimensionality-reduced modeling

6.2 Processing Pipeline Summary

```
STAGE 1: Data Loading
└ Initial: 1,460 rows × 81 features
└ Missing values: 7,829 cells

STAGE 2: Student Feature Addition
└ Added student_random_feature (seed=26, offset=4)

STAGE 3: Data Cleaning
└ Missing values: 7,829 → 0 (context-based strategies)
└ Outliers removed: 2
└ Result: 1,458 rows × 82 features

STAGE 4: Feature Transformation
└ Log transformation: 29 features
└ Target transformation: SalePrice → log(SalePrice)
└ Shapiro-Wilk validation: All p > 0.05

STAGE 5: Feature Creation
└ Aggregate features: 10
└ Interaction features: 3
└ Result: 1,458 rows × 95 features

STAGE 6: Categorical Encoding
└ Ordinal: 14 features
└ One-hot: 29 features → 427 binary columns
└ Label: 3 features
└ Result: 1,458 rows × 509 features

STAGE 7: Text Feature Engineering (TF-IDF)
└ Composite features: 3
└ TF-IDF vectorization: 30 weighted features
└ Result: 1,458 rows × 530 features

STAGE 8: Prepare X matrix
└ Exclude: Id, SalePrice
└ Features: 529

STAGE 9: Feature Scaling
└ StandardScaler: mean=0, std=1
└ All 529 features scaled

STAGE 10: Dimensionality Reduction (PCA)
└ Input: 529 features
└ Output: 278 components
└ Variance retained: 95.03%
└ Reduction: 47.4%
```

6.3 Key Achievements

Data Quality: - Eliminated all missing values using domain-appropriate strategies - Removed statistically-identified extreme outliers - Achieved complete, clean dataset

Feature Engineering: - Created meaningful aggregate and interaction features - Applied appropriate encoding for different categorical types - Implemented advanced NLP (TF-IDF) for text features

Statistical Rigor: - Shapiro-Wilk normality validation - VIF multicollinearity quantification - Mathematically justified all preprocessing decisions

Dimensionality Reduction: - Reduced feature space by 47.4% - Retained 95% of variance - Created orthogonal components (eliminated multicollinearity)

6.4 Student Random Feature Integration

The student_random_feature (generated with seed=26, offset=4) was successfully integrated throughout the analysis:

- Included in correlation analyses (showed near-zero correlations as expected)
 - Included in visualizations (showed no systematic patterns)
 - Included in PCA transformation (distributed across components without dominant loadings)
 - Behavior confirms randomness: no meaningful relationships with actual housing features
-

7. Conclusion

This assignment demonstrated comprehensive feature engineering techniques applied systematically to transform raw housing data into machine learning-ready formats. Key decisions were guided by domain knowledge, statistical principles, and data characteristics.

The resulting datasets are optimized for different modeling approaches: - **Engineered dataset:** Preserves explicit features for interpretable models - **PCA dataset:** Eliminates multicollinearity for regularized models

All preprocessing steps were documented with clear rationales, statistical validations, and measurable outcomes. The student random feature was properly integrated and verified to behave as a control feature throughout the pipeline.

Final Status: - Clean data: 1,458 rows with zero missing values - Rich features: 529 engineered features capturing domain knowledge - Efficient representation: 278 PCA components retaining 95% information - Ready for modeling: Both datasets saved and documented

Submitted by: Anik Das **Student ID:** 2025EM1100026 **Email:** 2025em1100026@bitspilani-digital.edu.in **Program:** MSc in Data Science and AI **Institution:** BITS Pilani Digital **Date:** November 2025