# A Comparative Analysis of Data Science Methodologies: Implementation of CRISP-DM, KDD, and SEMMA in Practical Applications

Rishi Patel

## Abstract

This research paper presents a comprehensive analysis of three prominent data science methodologies: CRISP-DM (Cross-Industry Standard Process for Data Mining), KDD (Knowledge Discovery in Databases), and SEMMA (Sample, Explore, Modify, Model, Assess). Through practical implementation across different use cases, we evaluate the effectiveness and applicability of each methodology in real-world scenarios. The study demonstrates how these frameworks can be applied to various prediction tasks, including house price prediction, survival analysis, and medical diagnosis.

## 1. Introduction

Data science methodologies provide structured approaches to solving complex analytical problems. While each methodology has its unique characteristics, they all aim to transform raw data into actionable insights. This research examines the implementation and effectiveness of CRISP-DM, KDD, and SEMMA through three distinct case studies:

- House price prediction using CRISP-DM
- Titanic survival prediction using KDD
- Stroke prediction using SEMMA

## 2. Methodology Overview

### 2.1 CRISP-DM Framework

CRISP-DM consists of six phases:

- Business Understanding
- Data Understanding
- Data Preparation
- Modeling
- Evaluation
- Deployment

### 2.2 KDD Process

KDD focuses on five key stages:

- Selection
- Preprocessing
- Transformation
- Data Mining
- Interpretation/Evaluation

### 2.3 SEMMA Approach

SEMMA encompasses five steps:

- Sample
- Explore
- Modify
- Model
- Assess

# 3. Case Studies

### 3.1 House Price Prediction Using CRISP-DM

# 3.1.1 Implementation

The CRISP-DM methodology was applied to the "House Prices: Advanced Regression Techniques" dataset from Kaggle. The implementation followed all six phases:

**Business Understanding:**

- Primary goal: Predict house sale prices
- Key stakeholders: Real estate agents, buyers, and sellers
- Success criteria: Accurate price predictions based on historical data

**Data Understanding:**

- Dataset features: 79 variables describing residential properties
- Key variables identified: LotArea, YearBuilt, SalePrice
- Initial analysis revealed correlations between features and target variable

**Data Preparation:**

- Missing value imputation using median for numerical features
- Categorical variable encoding through one-hot encoding
- Feature scaling implementation
- Feature selection based on correlation analysis

**Modeling:**

- Models tested: Linear Regression, Random Forest, Gradient Boosting
- Best performer: Gradient Boosting Regressor
- Performance metrics: MAE of 17,417, R² score of 0.87

**Evaluation and Deployment:**

- Model validation through cross-validation
- Deployment using joblib for model preservation

## 3.2 Titanic Survival Prediction Using KDD

# 3.2.1 Implementation

The KDD process was applied to the Titanic dataset, focusing on survival prediction:

**Selection and Preprocessing:**

- Feature selection: Pclass, Sex, Age, Fare, SibSp, Parch, Embarked
- Missing value treatment for Age, Cabin, and Embarked columns

**Transformation:**

- Categorical variable encoding
- Feature scaling for numerical variables
- Data standardization

**Data Mining:**

- Primary model: Logistic Regression
- Secondary model: Random Forest
- Performance metrics:
    - Logistic Regression accuracy: 81.01%
    - Random Forest accuracy: 81.56%

## 3.3 Stroke Prediction Using SEMMA

# 3.3.1 Implementation

The SEMMA methodology was applied to the Stroke Prediction Dataset:

**Sample and Explore:**

- Complete dataset utilization
- Exploratory data analysis
- Missing value identification in BMI column

**Modify:**

- Missing value imputation
- Feature encoding: Label encoding for binary variables
- One-hot encoding for multi-class variables
- Data scaling implementation

**Model and Assess:**

- Models implemented: Logistic Regression, Random Forest, SVM
- Evaluation metrics: Accuracy, Precision, Recall, F1 score, AUC
- Random Forest emerged as the best performer

# 4. Comparative Analysis

## 4.1 Methodological Strengths

**CRISP-DM:**

- Strong business focus
- Structured approach to problem-solving
- Clear deployment guidelines

**KDD:**

- Emphasis on knowledge discovery
- Detailed data transformation process
- Strong pattern recognition focus

**SEMMA:**

- Streamlined modeling process
- Strong emphasis on statistical analysis
- Clear assessment criteria

## 4.2 Application Scenarios

Each methodology showed particular strengths in different scenarios:

- CRISP-DM: Ideal for projects with clear business objectives
- KDD: Excellent for complex pattern discovery
- SEMMA: Optimal for model-centric projects

# 5. Conclusion

This research demonstrates that while each methodology has its unique approach, they all provide effective frameworks for data science projects. CRISP-DM excels in business-oriented projects, KDD in knowledge discovery, and SEMMA in model development and evaluation. The choice of methodology should be based on project requirements, complexity, and objectives.

# Acknowledgments