

# Comprehensive Scikit-learn Guide for Machine Learning

Machine Learning Team

December 23, 2025

## Contents

<b>1</b>	<b>1. Cheatsheet - Quick Reference</b>	<b>2</b>
1.1	1.1 Basic Import and Setup . . . . .	2
1.2	1.2 Data Preprocessing . . . . .	2
1.3	1.3 Model Training and Evaluation . . . . .	2
1.4	1.4 Model Selection and Validation . . . . .	3
1.5	1.5 Performance Metrics . . . . .	3
1.6	1.6 Feature Selection and Engineering . . . . .	4
1.7	1.7 Pipeline Construction . . . . .	4
1.8	1.8 Model Persistence . . . . .	4
<b>2</b>	<b>2. Glossary - Terms Defined</b>	<b>5</b>
2.1	2.1 A . . . . .	5
2.2	2.2 B . . . . .	5
2.3	2.3 C . . . . .	5
2.4	2.4 D . . . . .	6
2.5	2.5 E . . . . .	6
2.6	2.6 F . . . . .	7
2.7	2.7 G . . . . .	7
2.8	2.8 H . . . . .	7
2.9	2.9 I . . . . .	8
2.10	2.10 J . . . . .	8
2.11	2.11 K . . . . .	8
2.12	2.12 L . . . . .	8
2.13	2.13 M . . . . .	9
2.14	2.14 N . . . . .	9
2.15	2.15 O . . . . .	9
2.16	2.16 P . . . . .	10
2.17	2.17 Q . . . . .	10
2.18	2.18 R . . . . .	10
2.19	2.19 S . . . . .	11
2.20	2.20 T . . . . .	11
2.21	2.21 U . . . . .	11
2.22	2.22 V . . . . .	12
2.23	2.23 W . . . . .	12
2.24	2.24 X . . . . .	12

2.25	2.25 Y . . . . .	12
2.26	2.26 Z . . . . .	12
<b>3</b>	<b>3. Code Snippets</b>	<b>12</b>
3.1	3.1 Data Loading and Preprocessing . . . . .	12
3.1.1	3.1.1 Loading Data . . . . .	12
3.1.2	3.1.2 Data Exploration . . . . .	13
3.1.3	3.1.3 Handling Missing Values . . . . .	13
3.1.4	3.1.4 Feature Encoding . . . . .	14
3.1.5	3.1.5 Feature Scaling . . . . .	15
3.1.6	3.1.6 Feature Engineering . . . . .	16
3.2	3.2 Model Training and Evaluation . . . . .	17
3.2.1	3.2.1 Classification Models . . . . .	17
3.2.2	3.2.2 Regression Models . . . . .	18
3.2.3	3.2.3 Cross-Validation . . . . .	20
3.2.4	3.2.4 Hyperparameter Tuning . . . . .	21
3.2.5	3.2.5 Model Evaluation Metrics . . . . .	23
3.3	3.3 Feature Selection and Engineering . . . . .	25
3.3.1	3.3.1 Feature Selection Methods . . . . .	25
3.3.2	3.3.2 Dimensionality Reduction . . . . .	27
3.4	3.4 Model Selection and Pipeline . . . . .	29
3.4.1	3.4.1 Pipeline Construction . . . . .	29
3.4.2	3.4.2 Model Comparison and Selection . . . . .	32
3.5	3.5 Advanced Topics . . . . .	35
3.5.1	3.5.1 Imbalanced Dataset Handling . . . . .	35
3.5.2	3.5.2 Time Series Analysis . . . . .	38
<b>4</b>	<b>4. Best Practices</b>	<b>41</b>
4.1	4.1 Data Preprocessing Best Practices . . . . .	41
4.2	4.2 Model Training Best Practices . . . . .	45
4.3	4.3 Model Evaluation Best Practices . . . . .	51
4.4	4.4 Deployment and Production Best Practices . . . . .	62
<b>5</b>	<b>5. Learning from Errors</b>	<b>81</b>
5.1	5.1 Common Data Preprocessing Errors . . . . .	81
5.2	5.2 Model Training and Evaluation Errors . . . . .	98
5.3	5.3 Deployment and Production Errors . . . . .	136
<b>6</b>	<b>6. Counter-Intuitive Behaviors</b>	<b>194</b>
6.1	6.1 Model Behavior Surprises . . . . .	194

# 1 1. Cheatsheet - Quick Reference

## 1.1 1.1 Basic Import and Setup

```
1 # Essential imports
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler, LabelEncoder
6 from sklearn.metrics import accuracy_score, classification_report
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn.svm import SVC
10
11 # Basic setup
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
13     random_state=42)
14 scaler = StandardScaler()
15 X_train_scaled = scaler.fit_transform(X_train)
16 X_test_scaled = scaler.transform(X_test)
```

## 1.2 1.2 Data Preprocessing

```
1 # Handling missing values
2 from sklearn.impute import SimpleImputer
3 imputer = SimpleImputer(strategy='mean') # or 'median', 'most_frequent'
4 X_imputed = imputer.fit_transform(X)
5
6 # Encoding categorical variables
7 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
8 le = LabelEncoder()
9 y_encoded = le.fit_transform(y)
10
11 # One-hot encoding
12 from sklearn.preprocessing import OneHotEncoder
13 ohe = OneHotEncoder(sparse=False)
14 X_encoded = ohe.fit_transform(X_categorical)
15
16 # Scaling features
17 from sklearn.preprocessing import StandardScaler, MinMaxScaler
18 scaler = StandardScaler() # or MinMaxScaler()
19 X_scaled = scaler.fit_transform(X)
```

## 1.3 1.3 Model Training and Evaluation

```
1 # Basic model training
2 from sklearn.linear_model import LogisticRegression
3 model = LogisticRegression()
4 model.fit(X_train, y_train)
5 predictions = model.predict(X_test)
6
```

```
7 # Cross-validation
8 from sklearn.model_selection import cross_val_score
9 scores = cross_val_score(model, X, y, cv=5)
10
11 # Grid search
12 from sklearn.model_selection import GridSearchCV
13 param_grid = {'C': [0.1, 1, 10], 'penalty': ['l1', 'l2']}
14 grid_search = GridSearchCV(LogisticRegression(), param_grid, cv=5)
15 grid_search.fit(X_train, y_train)
```

## 1.4 1.4 Model Selection and Validation

```
1 # Train-test split
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
4     random_state=42)
5
6 # Cross-validation
7 from sklearn.model_selection import cross_val_score, StratifiedKFold
8 cv_scores = cross_val_score(model, X, y, cv=5)
9 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
10
11 # Validation curve
12 from sklearn.model_selection import validation_curve
13 train_scores, val_scores = validation_curve(model, X, y, param_name='param
14     ', param_range=param_range)
```

## 1.5 1.5 Performance Metrics

```
1 # Classification metrics
2 from sklearn.metrics import accuracy_score, precision_score, recall_score,
3     f1_score, roc_auc_score
4 accuracy = accuracy_score(y_test, predictions)
5 precision = precision_score(y_test, predictions, average='weighted')
6 recall = recall_score(y_test, predictions, average='weighted')
7 f1 = f1_score(y_test, predictions, average='weighted')
8 roc_auc = roc_auc_score(y_test, predictions_proba)
9
10 # Regression metrics
11 from sklearn.metrics import mean_squared_error, mean_absolute_error,
12     r2_score
13 mse = mean_squared_error(y_test, predictions)
14 mae = mean_absolute_error(y_test, predictions)
15 r2 = r2_score(y_test, predictions)
16
17 # Confusion matrix
18 from sklearn.metrics import confusion_matrix, classification_report
19 cm = confusion_matrix(y_test, predictions)
20 report = classification_report(y_test, predictions)
```

## 1.6 1.6 Feature Selection and Engineering

```
1 # Feature selection
2 from sklearn.feature_selection import SelectKBest, f_classif, RFE
3 selector = SelectKBest(score_func=f_classif, k=10)
4 X_selected = selector.fit_transform(X, y)
5
6 # Recursive Feature Elimination
7 from sklearn.feature_selection import RFE
8 rfe = RFE(estimator=LogisticRegression(), n_features_to_select=5)
9 X_rfe = rfe.fit_transform(X, y)
10
11 # Principal Component Analysis
12 from sklearn.decomposition import PCA
13 pca = PCA(n_components=0.95) # Retain 95% variance
14 X_pca = pca.fit_transform(X)
```

## 1.7 1.7 Pipeline Construction

```
1 # Basic pipeline
2 from sklearn.pipeline import Pipeline
3 pipeline = Pipeline([
4     ('scaler', StandardScaler()),
5     ('classifier', LogisticRegression())
6 ])
7 pipeline.fit(X_train, y_train)
8 predictions = pipeline.predict(X_test)
9
10 # Column transformer pipeline
11 from sklearn.compose import ColumnTransformer
12 from sklearn.preprocessing import StandardScaler, OneHotEncoder
13 preprocessor = ColumnTransformer(
14     transformers=[
15         ('num', StandardScaler(), numeric_features),
16         ('cat', OneHotEncoder(), categorical_features)
17     ])
18 full_pipeline = Pipeline([
19     ('preprocessor', preprocessor),
20     ('classifier', RandomForestClassifier())
21 ])
```

## 1.8 1.8 Model Persistence

```
1 # Save and load model
2 import joblib
3 joblib.dump(model, 'model.pkl')
4 loaded_model = joblib.load('model.pkl')
5
6 # Using pickle
7 import pickle
8 with open('model.pkl', 'wb') as f:
9     pickle.dump(model, f)
```

```
10 with open('model.pkl', 'rb') as f:  
11     loaded_model = pickle.load(f)
```

## 2 2. Glossary - Terms Defined

### 2.1 2.1 A

**Accuracy** - The proportion of correctly predicted instances among all predictions.

**Activation Function** - A function that transforms the weighted sum of inputs in neural network nodes.

**AdaBoost** - Adaptive Boosting, an ensemble method that combines weak learners into a strong learner.

**Agglomerative Clustering** - A hierarchical clustering method that starts with individual points and merges them.

**Algorithm** - A step-by-step procedure for solving a problem or completing a task.

**Anomaly Detection** - The identification of rare items, events or observations which raise suspicions.

**API (Application Programming Interface)** - A set of protocols and tools for building software applications.

**Area Under Curve (AUC)** - A performance measurement for classification problems at various threshold settings.

**Artificial Intelligence** - The simulation of human intelligence processes by machines, especially computer systems.

**Association Rules** - If-then statements that help to show the probability of relationships between data items.

### 2.2 2.2 B

**Bagging** - Bootstrap Aggregating, an ensemble method that trains multiple models on different subsets of data.

**Bias** - Error introduced by approximating a real-world problem, which may be complex, by a simplified model.

**Boosting** - An ensemble technique which converts weak learners to strong learners sequentially.

**Bootstrap** - A statistical method for estimating the sampling distribution of an estimator.

**Box Plot** - A standardized way of displaying the distribution of data based on a five number summary.

### 2.3 2.3 C

**Calibration** - The process of adjusting model outputs to better match observed frequencies.

**CART (Classification and Regression Trees)** - A decision tree learning technique.

**Categorical Data** - Data that represents categories or groups rather than numerical values.

**Classification** - A supervised learning task where the goal is to predict discrete class labels.

**Classifier** - A machine learning model that performs classification tasks.

**Clustering** - An unsupervised learning task that groups similar data points together.

**Coefficient of Determination ( $R^2$ )** - A statistical measure of how well regression predictions approximate real data.

**Collinearity** - A phenomenon where two or more predictor variables are highly correlated.

**Confidence Interval** - A range of values that's likely to contain a population parameter with a certain probability.

**Confusion Matrix** - A table used to describe the performance of a classification model.

**Convergence** - The property of a sequence or algorithm approaching a limit or stable state.

**Correlation** - A statistical measure that describes the extent to which two variables change together.

**Cross-Entropy** - A measure of the difference between two probability distributions.

**Cross-Validation** - A model validation technique for assessing how results will generalize to an independent data set.

**Curse of Dimensionality** - Various phenomena that arise when analyzing data in high-dimensional spaces.

## 2.4 2.4 D

**Data Cleaning** - The process of detecting and correcting corrupt or inaccurate records in a dataset.

**Data Mining** - The process of discovering patterns in large data sets involving methods at the intersection of machine learning, statistics, and database systems.

**Data Preprocessing** - Techniques used to transform raw data into a format suitable for modeling.

**Data Wrangling** - The process of cleaning, structuring and enriching raw data into a desired format.

**Decision Boundary** - A hypersurface that partitions the underlying vector space into sets.

**Decision Tree** - A tree-like model of decisions and their possible consequences.

**Deep Learning** - A subset of machine learning based on artificial neural networks with multiple layers.

**Density Estimation** - The construction of an estimate of an unknown probability density function.

**Dependent Variable** - The variable being predicted or explained in a statistical model.

**Dimensionality Reduction** - The process of reducing the number of random variables under consideration.

**Discriminant Analysis** - A statistical technique used to classify observations into predefined groups.

**Distance Metrics** - Mathematical functions that define the distance between two data points.

**Distribution** - A function that shows the possible values a variable can take and how often.

## 2.5 2.5 E

**Ensemble Methods** - Techniques that combine multiple learning algorithms to obtain better predictive performance.

**Entropy** - A measure of uncertainty or disorder in a dataset used in decision tree algorithms.

**Epoch** - One complete pass through the entire training dataset in machine learning.

**Error Rate** - The proportion of incorrect predictions made by a model.

**Estimator** - An object that fits a statistical model to data and provides predictions.

**Euclidean Distance** - The straight-line distance between two points in Euclidean space.

**Evaluation Metrics** - Quantitative measures used to assess model performance.

**Explained Variance** - The proportion of the dataset's variance that is accounted for by the model.

**Exploratory Data Analysis (EDA)** - An approach to analyzing data sets to summarize their main characteristics.

**Exponential Smoothing** - A time series forecasting method for univariate data.

## 2.6 2.6 F

**F1-Score** - The harmonic mean of precision and recall used in classification evaluation.

**False Negative** - An outcome where the model incorrectly predicts the negative class.

**False Positive** - An outcome where the model incorrectly predicts the positive class.

**Feature Engineering** - The process of using domain knowledge to create new features from raw data.

**Feature Extraction** - The process of reducing the dimensionality of data by transforming it.

**Feature Importance** - A technique to determine which features contribute most to model predictions.

**Feature Scaling** - Methods used to standardize the range of independent variables or features.

**Feature Selection** - The process of selecting a subset of relevant features for model construction.

**Fitting** - The process of training a model on data to learn the underlying patterns.

**Forest** - A collection of decision trees used in ensemble methods like Random Forest.

**Forward Selection** - A feature selection method that starts with no features and adds them one by one.

**Fuzzy Logic** - A form of many-valued logic in which the truth values of variables may be any real number between 0 and 1.

## 2.7 2.7 G

**Gaussian Distribution** - A probability distribution that is symmetric about the mean, showing that data near the mean are more frequent.

**Generalization** - The ability of a model to perform well on unseen data.

**Gradient Descent** - An optimization algorithm used to minimize a function by iteratively moving in the direction of steepest descent.

**Grid Search** - A hyperparameter tuning technique that exhaustively considers all parameter combinations.

**Gini Impurity** - A measure of how often a randomly chosen element would be incorrectly labeled.

**Goodness of Fit** - A statistical test describing how well a regression line fits a set of observations.

## 2.8 2.8 H

**Hierarchical Clustering** - A clustering method that builds a hierarchy of clusters.

**Holdout Validation** - A validation technique that splits data into training and testing sets.

**Homoscedasticity** - A condition where the variance of the error term is constant across all levels of the independent variables.

**Hyperparameter** - A parameter whose value is set before the learning process begins.

**Hypothesis Testing** - A statistical method that uses sample data to evaluate a hypothesis about a population parameter.



## 2.9 2.9 I

**Imputation** - The process of replacing missing data with substituted values.

**Independent Variable** - A variable that is manipulated or controlled in an experiment.

**Information Gain** - The reduction in entropy or surprise by transforming a dataset.

**Initialization** - The process of setting initial values for model parameters before training.

**Instance-Based Learning** - A family of learning algorithms that compare new problem instances with instances seen in training.

**Interpolation** - The estimation of unknown values that fall between known values.

**Iteration** - A single execution of the training process on a batch of data.

**Iris Dataset** - A classic dataset in machine learning containing measurements of iris flowers.

## 2.10 2.10 J

**Jaccard Index** - A statistic used for comparing the similarity and diversity of sample sets.

**Joint Probability** - The probability of two events occurring simultaneously.

## 2.11 2.11 K

**K-Means Clustering** - An unsupervised learning algorithm that partitions data into k clusters.

**K-Nearest Neighbors (KNN)** - A non-parametric method used for classification and regression.

**Kernel** - A function used in support vector machines to transform data into higher-dimensional space.

**K-Fold Cross-Validation** - A resampling procedure used to evaluate machine learning models.

## 2.12 2.12 L

**Label Encoding** - The process of converting categorical labels into numeric form.

**Lasso Regression** - A regression analysis method that performs both variable selection and regularization.

**Latent Variables** - Variables that are not directly observed but are inferred from other variables.

**Learning Rate** - A hyperparameter that controls how much to change the model in response to the estimated error.

**Leave-One-Out Cross-Validation** - A special case of cross-validation where each observation is used as a test set once.

**Linear Discriminant Analysis (LDA)** - A method used in statistics and machine learning to find a linear combination of features.

**Linear Regression** - A linear approach to modeling the relationship between a scalar response and explanatory variables.

**Logistic Regression** - A statistical model that uses a logistic function to model a binary dependent variable.

**Loss Function** - A function that maps an event or values of one or more variables onto a real number representing cost.

## 2.13 2.13 M

**Machine Learning** - A method of data analysis that automates analytical model building.

**Manhattan Distance** - The distance between two points measured along axes at right angles.

**Margin** - The distance between the hyperplane and the observations closest to the hyperplane in SVM.

**Matrix Factorization** - A class of techniques used to reduce the dimensionality of data.

**Maximum Likelihood Estimation** - A method of estimating the parameters of a probability distribution.

**Mean Absolute Error (MAE)** - The average of the absolute differences between predicted and actual values.

**Mean Squared Error (MSE)** - The average of the squares of the differences between predicted and actual values.

**Median** - The middle value in a sorted list of numbers.

**Memory-Based Learning** - Learning algorithms that store training data and use it directly for prediction.

**Metric Learning** - The process of learning a distance function over objects.

**Mini-Batch** - A small subset of the training data used in one iteration of model training.

**Mode** - The value that appears most frequently in a data set.

**Model Selection** - The process of selecting the best model from a set of candidate models.

**Multicollinearity** - The occurrence of high correlation among two or more independent variables.

**Multinomial Distribution** - A generalization of the binomial distribution for experiments with more than two outcomes.

## 2.14 2.14 N

**Naive Bayes** - A classification technique based on Bayes' theorem with an assumption of independence.

**Neural Network** - A computing system inspired by the biological neural networks that constitute animal brains.

**Noise** - Random variation in data that does not carry meaningful information.

**Normalization** - The process of scaling individual samples to have unit norm.

**Numpy** - A fundamental package for scientific computing in Python.

**Nystroem Method** - A technique for approximating the feature map of an RBF kernel.

## 2.15 2.15 O

**One-Hot Encoding** - A process of converting categorical data into a binary vector representation.

**One-Class SVM** - A variant of SVM used for novelty detection.

**One-Vs-All** - A strategy for multi-class classification using binary classifiers.

**One-Vs-One** - A strategy for multi-class classification that trains a binary classifier for each pair of classes.

**Outlier** - An observation point that is distant from other observations.

**Overfitting** - When a model learns the detail and noise in the training data to the extent that it negatively impacts performance on new data.

## 2.16 2.16 P

**Pandas** - A software library for data manipulation and analysis in Python.

**Parameter** - A configuration variable that is internal to the model and whose value can be estimated from data.

**Partial Fit** - A method that allows incremental learning on chunks of data.

**Perceptron** - An algorithm for supervised learning of binary classifiers.

**Pipeline** - A sequence of data processing components chained together.

**Polynomial Features** - Features created by raising existing features to various powers.

**Pooling** - An operation in neural networks that reduces the spatial dimensions of the input.

**Precision** - The ratio of correctly predicted positive observations to the total predicted positives.

**Principal Component Analysis (PCA)** - A statistical procedure that uses an orthogonal transformation to convert correlated variables into linearly uncorrelated variables.

**Probability** - A measure of the likelihood that an event will occur.

**Probability Calibration** - The process of transforming classifier scores into calibrated probabilities.

**Python** - A high-level programming language widely used in data science and machine learning.

## 2.17 2.17 Q

**Quantile** - Values that divide a probability distribution into continuous intervals with equal probabilities.

**Quantization** - The process of mapping input values from a large set to output values in a smaller set.

## 2.18 2.18 R

**Random Forest** - An ensemble learning method that operates by constructing multiple decision trees.

**Random Sampling** - A sampling method where each member of the population has an equal chance of being selected.

**Random State** - A parameter that controls the shuffling applied to the data before applying the split.

**Range** - The difference between the highest and lowest values in a dataset.

**Recall** - The ratio of correctly predicted positive observations to all observations in actual class.

**Regression** - A supervised learning task where the goal is to predict continuous values.

**Regularization** - A technique used to prevent overfitting by adding a penalty term to the loss function.

**Reinforcement Learning** - An area of machine learning concerned with how intelligent agents ought to take actions.

**Resampling** - A statistical method of repeatedly drawing samples from a dataset.

**Ridge Regression** - A technique for analyzing multiple regression data that suffer from multicollinearity.

**ROC Curve** - A graphical plot that illustrates the diagnostic ability of a binary classifier system.

**Root Mean Square Error (RMSE)** - The square root of the average of squared differences between prediction and actual observation.

## 2.19 2.19 S

**Sample** - A subset of a statistical population in which a survey or census is conducted.

**Sampling** - The process of selecting a subset of individuals from a population to estimate characteristics of the whole population.

**Scikit-learn** - A free software machine learning library for Python.

**Score** - A function that evaluates the quality of a model's predictions.

**Scoring** - The process of applying a trained model to new data to make predictions.

**Sensitivity** - The proportion of actual positives that are correctly identified.

**Separability** - The ability to distinguish between different classes in a classification problem.

**Shuffle** - The process of randomly reordering data.

**Silhouette Score** - A measure of how similar an object is to its own cluster compared to other clusters.

**Skewness** - A measure of the asymmetry of the probability distribution of a real-valued random variable.

**Smoothing** - Techniques used to reduce noise in data.

**Specificity** - The proportion of actual negatives that are correctly identified.

**Spline** - A special function defined piecewise by polynomials.

**Standard Deviation** - A measure of the amount of variation or dispersion of a set of values.

**Standardization** - The process of rescaling data to have a mean of 0 and standard deviation of 1.

**Statistical Learning** - The framework for studying and understanding the performance of machine learning algorithms.

**Stratified Sampling** - A method of sampling that involves dividing a population into homogeneous subgroups.

**Support Vector Machine (SVM)** - A supervised learning model that analyzes data for classification and regression analysis.

## 2.20 2.20 T

**Target Variable** - The variable that is being predicted in a supervised learning problem.

**Test Set** - A dataset used to provide an unbiased evaluation of a final model fit on the training dataset.

**Time Series** - A series of data points indexed in time order.

**Training Set** - A dataset used to train a machine learning model.

**Tree Pruning** - The process of reducing the size of decision trees to prevent overfitting.

**True Negative** - An outcome where the model correctly predicts the negative class.

**True Positive** - An outcome where the model correctly predicts the positive class.

**Type I Error** - The incorrect rejection of a true null hypothesis.

**Type II Error** - The failure to reject a false null hypothesis.

## 2.21 2.21 U

**Underfitting** - When a model is too simple to capture the underlying structure of the data.

**Unsupervised Learning** - A type of machine learning that looks for previously undetected patterns in a data set.

**Upsampling** - A technique used to increase the frequency of a signal or to balance imbalanced datasets.

## 2.22 2.22 V

**Validation Curve** - A plot that shows the training and validation scores for different parameter values.

**Validation Set** - A dataset used to tune hyperparameters and prevent overfitting.

**Variance** - A measure of how far a set of numbers are spread out from their average value.

**Variance Inflation Factor (VIF)** - A measure of multicollinearity in regression analysis.

**Vector** - A mathematical object that has both magnitude and direction.

**Voting Classifier** - An ensemble method that combines several classifiers and uses majority voting.

## 2.23 2.23 W

**Weight** - A parameter in a neural network that transforms input data within the network's hidden layers.

**Weight Decay** - A regularization technique that adds a penalty term to the loss function.

**Within-Cluster Sum of Squares** - A measure used in clustering to evaluate the compactness of clusters.

## 2.24 2.24 X

**XGBoost** - An optimized distributed gradient boosting library designed to be highly efficient, flexible and portable.

## 2.25 2.25 Y

**Yield** - A keyword in Python used to create generator functions.

## 2.26 2.26 Z

**Z-Score** - A statistical measurement that describes a value's relationship to the mean of a group of values.

# 3 3. Code Snippets

## 3.1 3.1 Data Loading and Preprocessing

### 3.1.1 3.1.1 Loading Data

```
1 # Load from CSV
2 import pandas as pd
3 data = pd.read_csv('data.csv')
4
5 # Load from Excel
6 data = pd.read_excel('data.xlsx', sheet_name='Sheet1')
7
8 # Load from JSON
9 data = pd.read_json('data.json')
10
11 # Load from database
12 import sqlite3
```

```
13 conn = sqlite3.connect('database.db')
14 data = pd.read_sql_query('SELECT * FROM table_name', conn)
15
16 # Load from URL
17 data = pd.read_csv('https://example.com/data.csv')
18
19 # Load built-in datasets
20 from sklearn.datasets import load_iris, load_boston, load_breast_cancer
21 iris = load_iris()
22 boston = load_boston() # Note: deprecated in newer versions
23 cancer = load_breast_cancer()
```

Listing 1: Load data from various sources

### 3.1.2 3.1.2 Data Exploration

```
1 # Basic dataset information
2 print(data.shape)
3 print(data.info())
4 print(data.describe())
5 print(data.head())
6 print(data.tail())
7
8 # Check for missing values
9 print(data.isnull().sum())
10 print(data.isnull().mean() * 100) # Percentage of missing values
11
12 # Check data types
13 print(data.dtypes)
14
15 # Check unique values in categorical columns
16 for col in data.select_dtypes(include=['object']).columns:
17     print(f {col}: {data[col].nunique()} unique values )
18
19 # Correlation matrix
20 import seaborn as sns
21 import matplotlib.pyplot as plt
22 correlation_matrix = data.corr()
23 plt.figure(figsize=(10, 8))
24 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm')
25 plt.show()
```

Listing 2: Explore dataset characteristics

### 3.1.3 3.1.3 Handling Missing Values

```
1 from sklearn.impute import SimpleImputer, KNNImputer, IterativeImputer
2 import pandas as pd
3 import numpy as np
4
5 # Simple imputation strategies
6 # Mean imputation
7 mean_imputer = SimpleImputer(strategy='mean')
```

```

8 data_mean_imputed = mean_imputer.fit_transform(data_numeric)
9
10 # Median imputation
11 median_imputer = SimpleImputer(strategy='median')
12 data_median_imputed = median_imputer.fit_transform(data_numeric)
13
14 # Mode imputation for categorical
15 mode_imputer = SimpleImputer(strategy='most_frequent')
16 data_mode_imputed = mode_imputer.fit_transform(data_categorical)
17
18 # Forward fill for time series
19 data_ffill = data.fillna(method='ffill')
20
21 # Backward fill for time series
22 data_bfill = data.fillna(method='bfill')
23
24 # KNN imputation
25 knn_imputer = KNNImputer(n_neighbors=5)
26 data_knn_imputed = knn_imputer.fit_transform(data_numeric)
27
28 # Advanced imputation with IterativeImputer
29 from sklearn.experimental import enable_iterative_imputer
30 iterative_imputer = IterativeImputer(random_state=42)
31 data_iterative_imputed = iterative_imputer.fit_transform(data_numeric)
32
33 # Custom imputation function
34 def custom_impute(df):
35     df_imputed = df.copy()
36     for column in df.columns:
37         if df[column].dtype in ['int64', 'float64']:
38             # For numeric columns, use median
39             df_imputed[column].fillna(df[column].median(), inplace=True)
40         else:
41             # For categorical columns, use mode
42             mode_value = df[column].mode()
43             if len(mode_value) > 0:
44                 df_imputed[column].fillna(mode_value[0], inplace=True)
45     return df_imputed
46
47 data_custom_imputed = custom_impute(data)

```

Listing 3: Handle missing data comprehensively

### 3.1.4 3.1.4 Feature Encoding

```

1 from sklearn.preprocessing import LabelEncoder, OneHotEncoder,
   OrdinalEncoder
2 from sklearn.compose import ColumnTransformer
3 import pandas as pd
4
5 # Label encoding for target variable
6 le = LabelEncoder()
7 y_encoded = le.fit_transform(y)

```

```

8
9 # One-hot encoding for categorical features
10 ohe = OneHotEncoder(sparse=False, drop='first') # drop='first' to avoid
    multicollinearity
11 X_encoded = ohe.fit_transform(X_categorical)
12
13 # Ordinal encoding for ordered categorical variables
14 oe = OrdinalEncoder(categories=[['Low', 'Medium', 'High']])
15 X_ordinal = oe.fit_transform(X_ordinal_categorical)
16
17 # Column transformer for mixed data types
18 numeric_features = ['age', 'income', 'score']
19 categorical_features = ['gender', 'education', 'city']
20
21 preprocessor = ColumnTransformer(
22     transformers=[
23         ('num', 'passthrough', numeric_features),
24         ('cat', OneHotEncoder(drop='first'), categorical_features)
25     ])
26
27 X_processed = preprocessor.fit_transform(X)
28
29 # Pandas get_dummies (alternative to OneHotEncoder)
30 X_dummies = pd.get_dummies(X, columns=categorical_features, drop_first=
    True)
31
32 # Target encoding (mean encoding)
33 def target_encode(df, categorical_col, target_col):
34     encoding_map = df.groupby(categorical_col)[target_col].mean().to_dict
    ()
35     return df[categorical_col].map(encoding_map)
36
37 X_target_encoded = target_encode(data, 'category', 'target')

```

Listing 4: Encode categorical variables

### 3.1.5 3.1.5 Feature Scaling

```

1 from sklearn.preprocessing import StandardScaler, MinMaxScaler,
    RobustScaler, Normalizer
2 import numpy as np
3
4 # Standardization (Z-score normalization)
5 scaler = StandardScaler()
6 X_standardized = scaler.fit_transform(X)
7
8 # Min-Max scaling (0-1 range)
9 minmax_scaler = MinMaxScaler()
10 X_minmax = minmax_scaler.fit_transform(X)
11
12 # Robust scaling (using median and IQR)
13 robust_scaler = RobustScaler()
14 X_robust = robust_scaler.fit_transform(X)

```



```

15
16 # Normalization (unit vector)
17 normalizer = Normalizer()
18 X_normalized = normalizer.fit_transform(X)
19
20 # Custom scaling function
21 def custom_scale(X, method='standard'):
22     if method == 'standard':
23         return (X - np.mean(X, axis=0)) / np.std(X, axis=0)
24     elif method == 'minmax':
25         return (X - np.min(X, axis=0)) / (np.max(X, axis=0) - np.min(X,
axis=0))
26     elif method == 'robust':
27         median = np.median(X, axis=0)
28         q75, q25 = np.percentile(X, [75, 25], axis=0)
29         iqr = q75 - q25
30         return (X - median) / iqr
31
32 # Apply inverse transform to get original scale
33 X_original = scaler.inverse_transform(X_standardized)

```

Listing 5: Scale numerical features

### 3.1.6 3.1.6 Feature Engineering

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.preprocessing import PolynomialFeatures
4
5 # Create polynomial features
6 poly = PolynomialFeatures(degree=2, include_bias=False)
7 X_poly = poly.fit_transform(X)
8
9 # Create interaction features
10 def create_interactions(df, feature_pairs):
11     df_new = df.copy()
12     for feat1, feat2 in feature_pairs:
13         df_new[f'{feat1}_{feat2}_interaction'] = df[feat1] * df[feat2]
14     return df_new
15
16 interaction_pairs = [('feature1', 'feature2'), ('feature1', 'feature3')]
17 data_with_interactions = create_interactions(data, interaction_pairs)
18
19 # Create binned features
20 data['age_group'] = pd.cut(data['age'], bins=[0, 18, 35, 50, 100],
21                             labels=['Child', 'Young', 'Adult', 'Senior'])
22
23 # Create date-based features
24 data['date'] = pd.to_datetime(data['date'])
25 data['year'] = data['date'].dt.year
26 data['month'] = data['date'].dt.month
27 data['day_of_week'] = data['date'].dt.dayofweek
28 data['is_weekend'] = data['day_of_week'].isin([5, 6]).astype(int)

```

```

29
30 # Create ratio features
31 data['income_to_age_ratio'] = data['income'] / (data['age'] + 1) # Add 1
    to avoid division by zero
32
33 # Create binary flags
34 data['high_income'] = (data['income'] > data['income'].median()).astype(
    int)
35
36 # Create rolling window features (for time series)
37 data['rolling_mean_7'] = data['value'].rolling(window=7).mean()
38 data['rolling_std_7'] = data['value'].rolling(window=7).std()
39
40 # Create lag features
41 data['lag_1'] = data['value'].shift(1)
42 data['lag_7'] = data['value'].shift(7)
43
44 # Create cyclical features (for periodic data)
45 data['month_sin'] = np.sin(2 * np.pi * data['month'] / 12)
46 data['month_cos'] = np.cos(2 * np.pi * data['month'] / 12)

```

Listing 6: Create new features from existing data

## 3.2 Model Training and Evaluation

### 3.2.1 Classification Models

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.ensemble import RandomForestClassifier,
    GradientBoostingClassifier
4 from sklearn.svm import SVC
5 from sklearn.naive_bayes import GaussianNB
6 from sklearn.neighbors import KNeighborsClassifier
7 from sklearn.tree import DecisionTreeClassifier
8 from sklearn.metrics import accuracy_score, classification_report,
    confusion_matrix
9
10 # Split data
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)
12
13 # Logistic Regression
14 lr = LogisticRegression(random_state=42)
15 lr.fit(X_train, y_train)
16 lr_pred = lr.predict(X_test)
17 lr_proba = lr.predict_proba(X_test)
18
19 # Random Forest
20 rf = RandomForestClassifier(n_estimators=100, random_state=42)
21 rf.fit(X_train, y_train)
22 rf_pred = rf.predict(X_test)
23 rf_proba = rf.predict_proba(X_test)
24

```

```

25 # Support Vector Machine
26 svm = SVC(probability=True, random_state=42)
27 svm.fit(X_train, y_train)
28 svm_pred = svm.predict(X_test)
29 svm_proba = svm.predict_proba(X_test)
30
31 # Gradient Boosting
32 gb = GradientBoostingClassifier(random_state=42)
33 gb.fit(X_train, y_train)
34 gb_pred = gb.predict(X_test)
35 gb_proba = gb.predict_proba(X_test)
36
37 # Naive Bayes
38 nb = GaussianNB()
39 nb.fit(X_train, y_train)
40 nb_pred = nb.predict(X_test)
41 nb_proba = nb.predict_proba(X_test)
42
43 # K-Nearest Neighbors
44 knn = KNeighborsClassifier(n_neighbors=5)
45 knn.fit(X_train, y_train)
46 knn_pred = knn.predict(X_test)
47 knn_proba = knn.predict_proba(X_test)
48
49 # Decision Tree
50 dt = DecisionTreeClassifier(random_state=42)
51 dt.fit(X_train, y_train)
52 dt_pred = dt.predict(X_test)
53 dt_proba = dt.predict_proba(X_test)
54
55 # Evaluate models
56 models = {
57     'Logistic Regression': lr_pred,
58     'Random Forest': rf_pred,
59     'SVM': svm_pred,
60     'Gradient Boosting': gb_pred,
61     'Naive Bayes': nb_pred,
62     'KNN': knn_pred,
63     'Decision Tree': dt_pred
64 }
65
66 for name, predictions in models.items():
67     accuracy = accuracy_score(y_test, predictions)
68     print(f {name} Accuracy: {accuracy:.4f} )
69     print(f {name} Classification Report: )
70     print(classification_report(y_test, predictions))
71     print( - * 50)

```

Listing 7: Train various classification models

### 3.2.2 3.2.2 Regression Models

```

1 from sklearn.linear_model import LinearRegression, Ridge, Lasso,
   ElasticNet

```

```
2 from sklearn.ensemble import RandomForestRegressor,
    GradientBoostingRegressor
3 from sklearn.svm import SVR
4 from sklearn.neighbors import KNeighborsRegressor
5 from sklearn.tree import DecisionTreeRegressor
6 from sklearn.metrics import mean_squared_error, mean_absolute_error,
    r2_score
7
8 # Linear Regression
9 lr = LinearRegression()
10 lr.fit(X_train, y_train)
11 lr_pred = lr.predict(X_test)
12
13 # Ridge Regression
14 ridge = Ridge(alpha=1.0)
15 ridge.fit(X_train, y_train)
16 ridge_pred = ridge.predict(X_test)
17
18 # Lasso Regression
19 lasso = Lasso(alpha=1.0)
20 lasso.fit(X_train, y_train)
21 lasso_pred = lasso.predict(X_test)
22
23 # Elastic Net
24 elastic = ElasticNet(alpha=1.0, l1_ratio=0.5)
25 elastic.fit(X_train, y_train)
26 elastic_pred = elastic.predict(X_test)
27
28 # Random Forest Regressor
29 rf = RandomForestRegressor(n_estimators=100, random_state=42)
30 rf.fit(X_train, y_train)
31 rf_pred = rf.predict(X_test)
32
33 # Support Vector Regression
34 svr = SVR(kernel='rbf')
35 svr.fit(X_train, y_train)
36 svr_pred = svr.predict(X_test)
37
38 # Gradient Boosting Regressor
39 gb = GradientBoostingRegressor(random_state=42)
40 gb.fit(X_train, y_train)
41 gb_pred = gb.predict(X_test)
42
43 # K-Nearest Neighbors Regressor
44 knn = KNeighborsRegressor(n_neighbors=5)
45 knn.fit(X_train, y_train)
46 knn_pred = knn.predict(X_test)
47
48 # Decision Tree Regressor
49 dt = DecisionTreeRegressor(random_state=42)
50 dt.fit(X_train, y_train)
51 dt_pred = dt.predict(X_test)
52
53 # Evaluate regression models
```

```

54 def evaluate_regression(y_true, y_pred, model_name):
55     mse = mean_squared_error(y_true, y_pred)
56     rmse = np.sqrt(mse)
57     mae = mean_absolute_error(y_true, y_pred)
58     r2 = r2_score(y_true, y_pred)
59
60     print(f {model_name}: )
61     print(f     RMSE: {rmse:.4f} )
62     print(f     MAE: {mae:.4f} )
63     print(f     R   : {r2:.4f} )
64     print(f -   * 30)
65
66 # Evaluate all models
67 models_regression = {
68     'Linear Regression': lr_pred,
69     'Ridge': ridge_pred,
70     'Lasso': lasso_pred,
71     'Elastic Net': elastic_pred,
72     'Random Forest': rf_pred,
73     'SVR': svr_pred,
74     'Gradient Boosting': gb_pred,
75     'KNN': knn_pred,
76     'Decision Tree': dt_pred
77 }
78
79 for name, predictions in models_regression.items():
80     evaluate_regression(y_test, predictions, name)

```

Listing 8: Train various regression models

### 3.2.3 3.2.3 Cross-Validation

```

1 from sklearn.model_selection import cross_val_score, StratifiedKFold,
  KFold
2 from sklearn.model_selection import LeaveOneOut, ShuffleSplit
3 import numpy as np
4
5 # Basic cross-validation
6 cv_scores = cross_val_score(model, X, y, cv=5)
7 print(f CV Scores: {cv_scores} )
8 print(f Mean CV Score: {cv_scores.mean():.4f} (+/- {cv_scores.std() * 2:.4
  f}) )
9
10 # Stratified K-Fold (for classification)
11 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
12 skf_scores = cross_val_score(model, X, y, cv=skf)
13 print(f Stratified CV Scores: {skf_scores} )
14
15 # K-Fold (for regression)
16 kf = KFold(n_splits=5, shuffle=True, random_state=42)
17 kf_scores = cross_val_score(model, X, y, cv=kf, scoring='
  neg_mean_squared_error')
18 print(f KFold CV Scores: {-kf_scores} )

```

```

19
20 # Leave-One-Out Cross-Validation
21 loo = LeaveOneOut()
22 loo_scores = cross_val_score(model, X_small, y_small, cv=loo)
23 print(f L00 CV Score: {loo_scores.mean():.4f} )
24
25 # Shuffle Split
26 ss = ShuffleSplit(n_splits=10, test_size=0.2, random_state=42)
27 ss_scores = cross_val_score(model, X, y, cv=ss)
28 print(f ShuffleSplit CV Score: {ss_scores.mean():.4f} )
29
30 # Custom cross-validation with multiple metrics
31 from sklearn.model_selection import cross_validate
32 scoring = ['accuracy', 'precision_macro', 'recall_macro']
33 cv_results = cross_validate(model, X, y, cv=5, scoring=scoring)
34 for metric in scoring:
35     scores = cv_results[f'test_{metric}']
36     print(f {metric}: {scores.mean():.4f} (+/- {scores.std() * 2:.4f}) )
37
38 # Time Series Cross-Validation
39 from sklearn.model_selection import TimeSeriesSplit
40 tscv = TimeSeriesSplit(n_splits=5)
41 tscv_scores = cross_val_score(model, X_time_series, y_time_series, cv=tscv
42                               )
43 print(f Time Series CV Score: {tscv_scores.mean():.4f} )

```

Listing 9: Implement cross-validation strategies

### 3.2.4 3.2.4 Hyperparameter Tuning

```

1 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
2 from sklearn.model_selection import validation_curve, learning_curve
3 from scipy.stats import randint, uniform
4 import numpy as np
5
6 # Grid Search
7 param_grid = {
8     'n_estimators': [50, 100, 200],
9     'max_depth': [3, 5, 7, None],
10    'min_samples_split': [2, 5, 10],
11    'min_samples_leaf': [1, 2, 4]
12 }
13
14 grid_search = GridSearchCV(
15     estimator=RandomForestClassifier(random_state=42),
16     param_grid=param_grid,
17     cv=5,
18     scoring='accuracy',
19     n_jobs=-1,
20     verbose=1
21 )
22
23 grid_search.fit(X_train, y_train)

```

```

24 print(f Best parameters: {grid_search.best_params_} )
25 print(f Best cross-validation score: {grid_search.best_score_:.4f} )
26
27 # Randomized Search
28 param_dist = {
29     'C': uniform(0.1, 10),
30     'gamma': uniform(0.001, 1),
31     'kernel': ['rbf', 'poly', 'sigmoid']
32 }
33
34 random_search = RandomizedSearchCV(
35     estimator=SVC(random_state=42),
36     param_distributions=param_dist,
37     n_iter=100,
38     cv=5,
39     scoring='accuracy',
40     n_jobs=-1,
41     random_state=42
42 )
43
44 random_search.fit(X_train, y_train)
45 print(f Best parameters: {random_search.best_params_} )
46 print(f Best cross-validation score: {random_search.best_score_:.4f} )
47
48 # Validation Curve
49 from sklearn.model_selection import validation_curve
50 param_range = [1, 10, 100, 1000]
51 train_scores, val_scores = validation_curve(
52     SVC(), X, y, param_name='C', param_range=param_range, cv=5
53 )
54
55 train_mean = np.mean(train_scores, axis=1)
56 train_std = np.std(train_scores, axis=1)
57 val_mean = np.mean(val_scores, axis=1)
58 val_std = np.std(val_scores, axis=1)
59
60 # Nested Cross-Validation
61 from sklearn.model_selection import cross_val_score
62 nested_scores = cross_val_score(grid_search, X, y, cv=5)
63 print(f Nested CV Score: {nested_scores.mean():.4f} (+/- {nested_scores.
64     std() * 2:.4f}) )
65
66 # Bayesian Optimization (using scikit-optimize)
67 # pip install scikit-optimize
68 from skopt import BayesSearchCV
69 from skopt.space import Real, Categorical, Integer
70
71 search_space = {
72     'n_estimators': Integer(10, 1000),
73     'max_depth': Integer(1, 50),
74     'min_samples_split': Integer(2, 20),
75     'min_samples_leaf': Integer(1, 10)
76 }

```

```

77 bayes_search = BayesSearchCV(
78     estimator=RandomForestClassifier(random_state=42),
79     search_spaces=search_space,
80     n_iter=50,
81     cv=5,
82     scoring='accuracy',
83     n_jobs=-1,
84     random_state=42
85 )
86
87 bayes_search.fit(X_train, y_train)
88 print(f Bayesian Optimization Best Score: {bayes_search.best_score_:.4f} )

```

Listing 10: Optimize model hyperparameters

### 3.2.5 3.2.5 Model Evaluation Metrics

```

1 from sklearn.metrics import (accuracy_score, precision_score, recall_score
  , f1_score,
2
3                               roc_auc_score, roc_curve,
4                               precision_recall_curve,
5                               confusion_matrix, classification_report,
6                               mean_squared_error, mean_absolute_error,
7                               r2_score,
8                               silhouette_score, adjusted_rand_score)
9
10 # Classification Metrics
11 def evaluate_classification(y_true, y_pred, y_proba=None):
12     # Basic metrics
13     accuracy = accuracy_score(y_true, y_pred)
14     precision = precision_score(y_true, y_pred, average='weighted')
15     recall = recall_score(y_true, y_pred, average='weighted')
16     f1 = f1_score(y_true, y_pred, average='weighted')
17
18     print(f Accuracy: {accuracy:.4f} )
19     print(f Precision: {precision:.4f} )
20     print(f Recall: {recall:.4f} )
21     print(f F1-Score: {f1:.4f} )
22
23     # Confusion Matrix
24     cm = confusion_matrix(y_true, y_pred)
25     plt.figure(figsize=(8, 6))
26     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
27     plt.title('Confusion Matrix')
28     plt.ylabel('Actual')
29     plt.xlabel('Predicted')
30     plt.show()
31
32     # Classification Report
33     print( \nClassification Report: )

```



```

34     print(classification_report(y_true, y_pred))
35
36     # ROC Curve and AUC (for binary classification)
37     if y_proba is not None and len(np.unique(y_true)) == 2:
38         fpr, tpr, _ = roc_curve(y_true, y_proba[:, 1])
39         auc = roc_auc_score(y_true, y_proba[:, 1])
40
41         plt.figure(figsize=(8, 6))
42         plt.plot(fpr, tpr, label=f'ROC Curve (AUC = {auc:.4f})')
43         plt.plot([0, 1], [0, 1], 'k--', label='Random Classifier')
44         plt.xlabel('False Positive Rate')
45         plt.ylabel('True Positive Rate')
46         plt.title('ROC Curve')
47         plt.legend()
48         plt.show()
49
50         print(f AUC Score: {auc:.4f} )
51
52     # Precision-Recall Curve
53     if y_proba is not None and len(np.unique(y_true)) == 2:
54         precision_vals, recall_vals, _ = precision_recall_curve(y_true,
55 y_proba[:, 1])
56
57         plt.figure(figsize=(8, 6))
58         plt.plot(recall_vals, precision_vals)
59         plt.xlabel('Recall')
60         plt.ylabel('Precision')
61         plt.title('Precision-Recall Curve')
62         plt.show()
63
64     # Regression Metrics
65     def evaluate_regression(y_true, y_pred):
66         mse = mean_squared_error(y_true, y_pred)
67         rmse = np.sqrt(mse)
68         mae = mean_absolute_error(y_true, y_pred)
69         r2 = r2_score(y_true, y_pred)
70
71         print(f Mean Squared Error: {mse:.4f} )
72         print(f Root Mean Squared Error: {rmse:.4f} )
73         print(f Mean Absolute Error: {mae:.4f} )
74         print(f R Score: {r2:.4f} )
75
76     # Residuals plot
77     residuals = y_true - y_pred
78     plt.figure(figsize=(10, 4))
79
80     plt.subplot(1, 2, 1)
81     plt.scatter(y_pred, residuals)
82     plt.xlabel('Predicted Values')
83     plt.ylabel('Residuals')
84     plt.title('Residuals vs Predicted')
85     plt.axhline(y=0, color='r', linestyle='--')
86
87     plt.subplot(1, 2, 2)

```

```

87     plt.hist(residuals, bins=30)
88     plt.xlabel('Residuals')
89     plt.ylabel('Frequency')
90     plt.title('Residuals Distribution')
91
92     plt.tight_layout()
93     plt.show()
94
95 # Clustering Metrics
96 def evaluate_clustering(X, labels):
97     if len(np.unique(labels)) > 1:
98         silhouette = silhouette_score(X, labels)
99         print(f Silhouette Score: {silhouette:.4f} )
100
101 # Plot clusters
102 if X.shape[1] >= 2:
103     plt.figure(figsize=(8, 6))
104     plt.scatter(X[:, 0], X[:, 1], c=labels, cmap='viridis')
105     plt.title('Clustering Results')
106     plt.colorbar()
107     plt.show()
108
109 # Custom scoring function
110 from sklearn.metrics import make_scorer
111
112 def custom_f1_score(y_true, y_pred):
113     return f1_score(y_true, y_pred, average='macro')
114
115 custom_scorer = make_scorer(custom_f1_score)
116
117 # Multi-class specific metrics
118 def multiclass_evaluation(y_true, y_pred):
119     # Per-class metrics
120     precision_per_class = precision_score(y_true, y_pred, average=None)
121     recall_per_class = recall_score(y_true, y_pred, average=None)
122     f1_per_class = f1_score(y_true, y_pred, average=None)
123
124     classes = np.unique(y_true)
125     for i, class_label in enumerate(classes):
126         print(f Class {class_label}: )
127         print(f     Precision: {precision_per_class[i]:.4f} )
128         print(f     Recall: {recall_per_class[i]:.4f} )
129         print(f     F1-Score: {f1_per_class[i]:.4f} )
130     print()

```

Listing 11: Comprehensive model evaluation

### 3.3 Feature Selection and Engineering

#### 3.3.1 Feature Selection Methods

```

1 from sklearn.feature_selection import (SelectKBest, SelectPercentile,
    f_classif, f_regression,

```

```

2         chi2, mutual_info_classif,
        mutual_info_regression,
3         RFE, RFECV, SelectFromModel)
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.linear_model import LassoCV
6 import pandas as pd
7 import numpy as np
8
9 # Univariate Feature Selection
10 # For classification
11 selector_f_classif = SelectKBest(score_func=f_classif, k=10)
12 X_selected_f_classif = selector_f_classif.fit_transform(X, y)
13
14 # For regression
15 selector_f_regression = SelectKBest(score_func=f_regression, k=10)
16 X_selected_f_regression = selector_f_regression.fit_transform(X, y)
17
18 # Using Chi-square test (for non-negative features)
19 selector_chi2 = SelectKBest(score_func=chi2, k=10)
20 X_selected_chi2 = selector_chi2.fit_transform(X_non_negative, y)
21
22 # Using mutual information
23 selector_mi_classif = SelectKBest(score_func=mutual_info_classif, k=10)
24 X_selected_mi_classif = selector_mi_classif.fit_transform(X, y)
25
26 # Select Percentile
27 selector_percentile = SelectPercentile(score_func=f_classif, percentile
    =20)
28 X_selected_percentile = selector_percentile.fit_transform(X, y)
29
30 # Recursive Feature Elimination (RFE)
31 estimator = RandomForestClassifier(n_estimators=100, random_state=42)
32 rfe = RFE(estimator=estimator, n_features_to_select=10)
33 X_selected_rfe = rfe.fit_transform(X, y)
34
35 # RFE with Cross-Validation
36 rfecv = RFECV(estimator=estimator, step=1, cv=5, scoring='accuracy')
37 X_selected_rfecv = rfecv.fit_transform(X, y)
38 print(f Optimal number of features: {rfecv.n_features_} )
39
40 # Feature selection based on model coefficients
41 lasso = LassoCV(cv=5, random_state=42)
42 lasso.fit(X, y)
43 selector_lasso = SelectFromModel(lasso, prefit=True)
44 X_selected_lasso = selector_lasso.transform(X)
45
46 # Feature selection based on feature importance
47 rf = RandomForestClassifier(n_estimators=100, random_state=42)
48 rf.fit(X, y)
49 selector_rf = SelectFromModel(rf, prefit=True)
50 X_selected_rf = selector_rf.transform(X)
51
52 # Variance Threshold (remove low-variance features)
53 from sklearn.feature_selection import VarianceThreshold

```

```

54 selector_variance = VarianceThreshold(threshold=0.01)
55 X_selected_variance = selector_variance.fit_transform(X)
56
57 # Correlation-based feature selection
58 def correlation_filter(X, threshold=0.95):
59     corr_matrix = pd.DataFrame(X).corr().abs()
60     upper_triangle = corr_matrix.where(
61         np.triu(np.ones(corr_matrix.shape), k=1).astype(bool)
62     )
63     to_drop = [column for column in upper_triangle.columns
64                 if any(upper_triangle[column] > threshold)]
65     return np.delete(X, to_drop, axis=1)
66
67 X_selected_corr = correlation_filter(X, threshold=0.95)
68
69 # Wrapper methods
70 from mlxtend.feature_selection import SequentialFeatureSelector
71
72 # Forward selection
73 sfs_forward = SequentialFeatureSelector(
74     RandomForestClassifier(n_estimators=100, random_state=42),
75     k_features=10,
76     forward=True,
77     floating=False,
78     scoring='accuracy',
79     cv=5
80 )
81 sfs_forward.fit(X, y)
82 X_selected_sfs = sfs_forward.transform(X)
83
84 # Backward elimination
85 sfs_backward = SequentialFeatureSelector(
86     RandomForestClassifier(n_estimators=100, random_state=42),
87     k_features=10,
88     forward=False,
89     floating=False,
90     scoring='accuracy',
91     cv=5
92 )
93 sfs_backward.fit(X, y)
94 X_selected_sfs_backward = sfs_backward.transform(X)

```

Listing 12: Implement various feature selection techniques

### 3.3.2 3.3.2 Dimensionality Reduction

```

1 from sklearn.decomposition import PCA, TruncatedSVD, FastICA, NMF
2 from sklearn.manifold import TSNE, MDS, Isomap, LocallyLinearEmbedding
3 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
4 import matplotlib.pyplot as plt
5 import numpy as np
6
7 # Principal Component Analysis (PCA)

```

```
8 pca = PCA(n_components=0.95) # Retain 95% of variance
9 X_pca = pca.fit_transform(X)
10 print(f Original dimensions: {X.shape} )
11 print(f PCA dimensions: {X_pca.shape} )
12 print(f Explained variance ratio: {pca.explained_variance_ratio_} )
13
14 # Plot PCA components
15 plt.figure(figsize=(10, 6))
16 plt.plot(np.cumsum(pca.explained_variance_ratio_))
17 plt.xlabel('Number of Components')
18 plt.ylabel('Cumulative Explained Variance')
19 plt.title('PCA Explained Variance')
20 plt.show()
21
22 # Linear Discriminant Analysis (LDA)
23 lda = LinearDiscriminantAnalysis(n_components=2)
24 X_lda = lda.fit_transform(X, y)
25 print(f LDA dimensions: {X_lda.shape} )
26
27 # Truncated SVD (for sparse matrices)
28 svd = TruncatedSVD(n_components=50)
29 X_svd = svd.fit_transform(X_sparse)
30
31 # Independent Component Analysis (ICA)
32 ica = FastICA(n_components=10, random_state=42)
33 X_ica = ica.fit_transform(X)
34
35 # Non-negative Matrix Factorization (NMF)
36 nmf = NMF(n_components=10, random_state=42)
37 X_nmf = nmf.fit_transform(X_non_negative)
38
39 # t-SNE (for visualization)
40 tsne = TSNE(n_components=2, random_state=42, perplexity=30)
41 X_tsne = tsne.fit_transform(X_sample) # Use sample for large datasets
42
43 # Plot t-SNE
44 plt.figure(figsize=(8, 6))
45 plt.scatter(X_tsne[:, 0], X_tsne[:, 1], c=y_sample, cmap='viridis')
46 plt.title('t-SNE Visualization')
47 plt.colorbar()
48 plt.show()
49
50 # Multidimensional Scaling (MDS)
51 mds = MDS(n_components=2, random_state=42)
52 X_mds = mds.fit_transform(X_sample)
53
54 # Isomap
55 isomap = Isomap(n_components=2, n_neighbors=5)
56 X_isomap = isomap.fit_transform(X)
57
58 # Locally Linear Embedding (LLE)
59 lle = LocallyLinearEmbedding(n_components=2, n_neighbors=5)
60 X_lle = lle.fit_transform(X)
61
```

```

62 # Kernel PCA
63 from sklearn.decomposition import KernelPCA
64 kpca = KernelPCA(n_components=10, kernel='rbf', gamma=0.1)
65 X_kpca = kpca.fit_transform(X)
66
67 # Autoencoders (using TensorFlow/Keras)
68 import tensorflow as tf
69 from tensorflow.keras.layers import Input, Dense
70 from tensorflow.keras.models import Model
71
72 def create_autoencoder(input_dim, encoding_dim):
73     input_layer = Input(shape=(input_dim,))
74     encoded = Dense(encoding_dim, activation='relu')(input_layer)
75     decoded = Dense(input_dim, activation='sigmoid')(encoded)
76
77     autoencoder = Model(input_layer, decoded)
78     encoder = Model(input_layer, encoded)
79
80     autoencoder.compile(optimizer='adam', loss='mse')
81     return autoencoder, encoder
82
83 # Use autoencoder for dimensionality reduction
84 input_dim = X.shape[1]
85 encoding_dim = 50
86 autoencoder, encoder = create_autoencoder(input_dim, encoding_dim)
87 autoencoder.fit(X, X, epochs=100, batch_size=256, shuffle=True, verbose=0)
88 X_autoencoded = encoder.predict(X)

```

Listing 13: Apply dimensionality reduction techniques

## 3.4 3.4 Model Selection and Pipeline

### 3.4.1 3.4.1 Pipeline Construction

```

1 from sklearn.pipeline import Pipeline, FeatureUnion
2 from sklearn.compose import ColumnTransformer
3 from sklearn.preprocessing import StandardScaler, OneHotEncoder,
  FunctionTransformer
4 from sklearn.impute import SimpleImputer
5 from sklearn.feature_selection import SelectKBest, f_classif
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.model_selection import train_test_split
8 import pandas as pd
9 import numpy as np
10
11 # Basic Pipeline
12 pipeline = Pipeline([
13     ('scaler', StandardScaler()),
14     ('selector', SelectKBest(f_classif, k=10)),
15     ('classifier', RandomForestClassifier(random_state=42))
16 ])
17
18 pipeline.fit(X_train, y_train)
19 predictions = pipeline.predict(X_test)

```

```

20
21 # Column Transformer Pipeline (for mixed data types)
22 numeric_features = ['age', 'income', 'score']
23 categorical_features = ['gender', 'education', 'city']
24
25 preprocessor = ColumnTransformer(
26     transformers=[
27         ('num', Pipeline([
28             ('imputer', SimpleImputer(strategy='median')),
29             ('scaler', StandardScaler())
30         ]), numeric_features),
31         ('cat', Pipeline([
32             ('imputer', SimpleImputer(strategy='constant', fill_value='
missing')),
33             ('onehot', OneHotEncoder(drop='first', handle_unknown='ignore'
))
34         ]), categorical_features)
35     ])
36
37 # Full pipeline with preprocessing
38 full_pipeline = Pipeline([
39     ('preprocessor', preprocessor),
40     ('feature_selection', SelectKBest(f_classif, k=20)),
41     ('classifier', RandomForestClassifier(n_estimators=100, random_state
=42))
42 ])
43
44 full_pipeline.fit(X_train, y_train)
45 predictions = full_pipeline.predict(X_test)
46
47 # Feature Union Pipeline (combine multiple feature extraction methods)
48 from sklearn.base import BaseEstimator, TransformerMixin
49
50 class CustomFeatureExtractor(BaseEstimator, TransformerMixin):
51     def __init__(self, feature_indices):
52         self.feature_indices = feature_indices
53
54     def fit(self, X, y=None):
55         return self
56
57     def transform(self, X):
58         return X[:, self.feature_indices]
59
60 # Feature union combining different transformations
61 feature_union = FeatureUnion([
62     ('pca_features', Pipeline([
63         ('scaler', StandardScaler()),
64         ('pca', PCA(n_components=10))
65     ])),
66     ('statistical_features', Pipeline([
67         ('scaler', StandardScaler()),
68         ('selector', SelectKBest(f_classif, k=5))
69     ]))
70 ])

```

```

71
72 union_pipeline = Pipeline([
73     ('features', feature_union),
74     ('classifier', RandomForestClassifier(random_state=42))
75 ])
76
77 # Custom Pipeline with Custom Transformers
78 class OutlierRemover(BaseEstimator, TransformerMixin):
79     def __init__(self, columns=None, threshold=3):
80         self.columns = columns
81         self.threshold = threshold
82
83     def fit(self, X, y=None):
84         return self
85
86     def transform(self, X):
87         if isinstance(X, pd.DataFrame):
88             X_copy = X.copy()
89             if self.columns is None:
90                 self.columns = X_copy.select_dtypes(include=[np.number]).
columns
91
92                 for col in self.columns:
93                     z_scores = np.abs((X_copy[col] - X_copy[col].mean()) /
X_copy[col].std())
94                     X_copy = X_copy[z_scores < self.threshold]
95                 return X_copy
96             return X
97
98 # Pipeline with outlier removal
99 outlier_pipeline = Pipeline([
100     ('outlier_remover', OutlierRemover(threshold=2.5)),
101     ('scaler', StandardScaler()),
102     ('classifier', RandomForestClassifier(random_state=42))
103 ])
104
105 # Pipeline with custom feature engineering
106 def create_interaction_features(X):
107     if isinstance(X, pd.DataFrame):
108         X_new = X.copy()
109         numeric_cols = X_new.select_dtypes(include=[np.number]).columns
110         for i in range(len(numeric_cols)):
111             for j in range(i+1, len(numeric_cols)):
112                 col1, col2 = numeric_cols[i], numeric_cols[j]
113                 X_new[f'{col1}_{col2}_interaction'] = X_new[col1] * X_new[
col2]
114             return X_new
115         return X
116
117 interaction_transformer = FunctionTransformer(create_interaction_features,
validate=False)
118
119 feature_engineering_pipeline = Pipeline([
120     ('interactions', interaction_transformer),

```



```

121     ('scaler', StandardScaler()),
122     ('selector', SelectKBest(f_classif, k=15)),
123     ('classifier', RandomForestClassifier(random_state=42))
124 ])
125
126 # Pipeline with parameter tuning
127 from sklearn.model_selection import GridSearchCV
128
129 param_grid = {
130     'classifier__n_estimators': [50, 100, 200],
131     'classifier__max_depth': [3, 5, 7, None],
132     'selector__k': [10, 15, 20]
133 }
134
135 pipeline_grid = GridSearchCV(
136     full_pipeline,
137     param_grid,
138     cv=5,
139     scoring='accuracy',
140     n_jobs=-1
141 )
142
143 pipeline_grid.fit(X_train, y_train)
144 print(f Best parameters: {pipeline_grid.best_params_} )
145 print(f Best score: {pipeline_grid.best_score_:.4f} )
146
147 # Save and load pipeline
148 import joblib
149 joblib.dump(full_pipeline, 'ml_pipeline.pkl')
150 loaded_pipeline = joblib.load('ml_pipeline.pkl')

```

Listing 14: Build comprehensive machine learning pipelines

### 3.4.2 Model Comparison and Selection

```

1 from sklearn.model_selection import cross_validate
2 from sklearn.metrics import accuracy_score, make_scorer
3 import pandas as pd
4 import numpy as np
5
6 # Define models to compare
7 models = {
8     'Logistic Regression': LogisticRegression(random_state=42, max_iter
9         =1000),
10    'Random Forest': RandomForestClassifier(n_estimators=100, random_state
11        =42),
12    'SVM': SVC(random_state=42, probability=True),
13    'Gradient Boosting': GradientBoostingClassifier(random_state=42),
14    'KNN': KNeighborsClassifier(n_neighbors=5),
15    'Naive Bayes': GaussianNB(),
16    'Decision Tree': DecisionTreeClassifier(random_state=42)
17 }

```

```

17 # Define scoring metrics
18 scoring = {
19     'accuracy': 'accuracy',
20     'precision': 'precision_macro',
21     'recall': 'recall_macro',
22     'f1': 'f1_macro',
23     'roc_auc': 'roc_auc_ovr' # One-vs-Rest for multiclass
24 }
25
26 # Compare models using cross-validation
27 results = {}
28 for name, model in models.items():
29     print(f Evaluating {name}... )
30     cv_results = cross_validate(model, X, y, cv=5, scoring=scoring, n_jobs
31                               =-1)
32     results[name] = {
33         'accuracy_mean': cv_results['test_accuracy'].mean(),
34         'accuracy_std': cv_results['test_accuracy'].std(),
35         'precision_mean': cv_results['test_precision'].mean(),
36         'precision_std': cv_results['test_precision'].std(),
37         'recall_mean': cv_results['test_recall'].mean(),
38         'recall_std': cv_results['test_recall'].std(),
39         'f1_mean': cv_results['test_f1'].mean(),
40         'f1_std': cv_results['test_f1'].std(),
41         'roc_auc_mean': cv_results['test_roc_auc'].mean(),
42         'roc_auc_std': cv_results['test_roc_auc'].std()
43     }
44
45 # Convert results to DataFrame for easy comparison
46 results_df = pd.DataFrame(results).T
47 print(results_df.round(4))
48
49 # Visualize model comparison
50 import matplotlib.pyplot as plt
51
52 metrics = ['accuracy', 'precision', 'recall', 'f1', 'roc_auc']
53 fig, axes = plt.subplots(2, 3, figsize=(15, 10))
54 axes = axes.ravel()
55
56 for i, metric in enumerate(metrics):
57     means = results_df[f'{metric}_mean']
58     stds = results_df[f'{metric}_std']
59
60     axes[i].bar(range(len(means)), means, yerr=stds, capsize=5)
61     axes[i].set_xticks(range(len(means)))
62     axes[i].set_xticklabels(means.index, rotation=45, ha='right')
63     axes[i].set_title(f'{metric.capitalize()} Comparison')
64     axes[i].set_ylabel(metric.capitalize())
65
66 # Remove empty subplot
67 axes[5].remove()
68 plt.tight_layout()
69 plt.show()

```

```

70
71 # Statistical significance testing
72 from scipy import stats
73
74 def compare_models_statistically(model1_name, model2_name, cv_scores1,
    cv_scores2):
75     Compare two models using paired t-test
76     t_stat, p_value = stats.ttest_rel(cv_scores1, cv_scores2)
77
78     print(f {model1_name} vs {model2_name}: )
79     print(f    t-statistic: {t_stat:.4f} )
80     print(f    p-value: {p_value:.4f} )
81
82     if p_value < 0.05:
83         print(    Result: Statistically significant difference )
84         if np.mean(cv_scores1) > np.mean(cv_scores2):
85             print(f    {model1_name} is significantly better )
86         else:
87             print(f    {model2_name} is significantly better )
88     else:
89         print(    Result: No statistically significant difference )
90     print()
91
92 # Compare top two models
93 top_models = results_df['accuracy_mean'].nlargest(2).index.tolist()
94 if len(top_models) >= 2:
95     model1, model2 = top_models[0], top_models[1]
96
97     # Get CV scores for top models
98     cv_scores1 = cross_validate(models[model1], X, y, cv=5, scoring='
accuracy')['test_score']
99     cv_scores2 = cross_validate(models[model2], X, y, cv=5, scoring='
accuracy')['test_score']
100
101     compare_models_statistically(model1, model2, cv_scores1, cv_scores2)
102
103 # Ensemble methods for model combination
104 from sklearn.ensemble import VotingClassifier, StackingClassifier
105 from sklearn.linear_model import LogisticRegression
106
107 # Voting Classifier (Hard Voting)
108 voting_clf = VotingClassifier(
109     estimators=[(name, model) for name, model in models.items()],
110     voting='hard'
111 )
112
113 # Voting Classifier (Soft Voting)
114 voting_clf_soft = VotingClassifier(
115     estimators=[(name, model) for name, model in models.items()
116                 if hasattr(model, 'predict_proba')],
117     voting='soft'
118 )
119
120 # Stacking Classifier

```

```

121 stacking_clf = StackingClassifier(
122     estimators=[(name, model) for name, model in models.items()],
123     final_estimator=LogisticRegression(random_state=42),
124     cv=5
125 )
126
127 # Compare ensemble methods
128 ensemble_models = {
129     'Hard Voting': voting_clf,
130     'Soft Voting': voting_clf_soft,
131     'Stacking': stacking_clf
132 }
133
134 ensemble_results = {}
135 for name, model in ensemble_models.items():
136     cv_scores = cross_validate(model, X, y, cv=5, scoring='accuracy',
137                               n_jobs=-1)
138     ensemble_results[name] = {
139         'mean': cv_scores['test_score'].mean(),
140         'std': cv_scores['test_score'].std()
141     }
142
143 ensemble_df = pd.DataFrame(ensemble_results).T
144 print(Ensemble Methods Comparison: )
145 print(ensemble_df.round(4))

```

Listing 15: Compare multiple models systematically

## 3.5 3.5 Advanced Topics

### 3.5.1 3.5.1 Imbalanced Dataset Handling

```

1 from sklearn.utils.class_weight import compute_class_weight
2 from imblearn.over_sampling import SMOTE, ADASYN, RandomOverSampler
3 from imblearn.under_sampling import RandomUnderSampler, TomekLinks
4 from imblearn.combine import SMOTETomek, SMOTEENN
5 from imblearn.pipeline import Pipeline as ImbPipeline
6 from collections import Counter
7 import numpy as np
8
9 # Analyze class distribution
10 class_distribution = Counter(y)
11 print(Original class distribution: , class_distribution)
12
13 # Compute class weights
14 class_weights = compute_class_weight('balanced', classes=np.unique(y), y=y)
15
16 class_weight_dict = dict(zip(np.unique(y), class_weights))
17 print(Class weights: , class_weight_dict)
18
19 # Apply to model
20 weighted_model = RandomForestClassifier(class_weight='balanced',

```

```

21 # Oversampling techniques
22 # SMOTE (Synthetic Minority Oversampling Technique)
23 smote = SMOTE(random_state=42)
24 X_smote, y_smote = smote.fit_resample(X, y)
25 print( SMOTE class distribution: , Counter(y_smote))
26
27 # ADASYN (Adaptive Synthetic Sampling)
28 adasyn = ADASYN(random_state=42)
29 X_adasyn, y_adasyn = adasyn.fit_resample(X, y)
30 print( ADASYN class distribution: , Counter(y_adasyn))
31
32 # Random Oversampling
33 random_oversampler = RandomOverSampler(random_state=42)
34 X_random_over, y_random_over = random_oversampler.fit_resample(X, y)
35 print( Random oversampling distribution: , Counter(y_random_over))
36
37 # Undersampling techniques
38 # Random Undersampling
39 random_undersampler = RandomUnderSampler(random_state=42)
40 X_random_under, y_random_under = random_undersampler.fit_resample(X, y)
41 print( Random undersampling distribution: , Counter(y_random_under))
42
43 # Tomek Links
44 tomek = TomekLinks()
45 X_tomek, y_tomek = tomek.fit_resample(X, y)
46 print( Tomek Links distribution: , Counter(y_tomek))
47
48 # Combined techniques
49 # SMOTE + Tomek Links
50 smote_tomek = SMOTETomek(random_state=42)
51 X_smote_tomek, y_smote_tomek = smote_tomek.fit_resample(X, y)
52 print( SMOTE + Tomek Links distribution: , Counter(y_smote_tomek))
53
54 # SMOTE + Edited Nearest Neighbors
55 from imblearn.combine import SMOTEENN
56 smote_enn = SMOTEENN(random_state=42)
57 X_smote_enn, y_smote_enn = smote_enn.fit_resample(X, y)
58 print( SMOTE + ENN distribution: , Counter(y_smote_enn))
59
60 # Pipeline with resampling
61 resampling_pipeline = ImbPipeline([
62     ('smote', SMOTE(random_state=42)),
63     ('classifier', RandomForestClassifier(random_state=42))
64 ])
65
66 # Advanced resampling with custom sampling strategy
67 from imblearn.over_sampling import SMOTE
68
69 # Custom sampling strategy
70 sampling_strategy = {0: 1000, 1: 1000, 2: 1000} # Equalize all classes to
1000 samples
71 custom_smote = SMOTE(sampling_strategy=sampling_strategy, random_state=42)
72 X_custom, y_custom = custom_smote.fit_resample(X, y)
73

```

```

74 # Evaluate different resampling techniques
75 from sklearn.model_selection import cross_val_score
76
77 techniques = {
78     'Original': (X, y),
79     'SMOTE': (X_smote, y_smote),
80     'ADASYN': (X_adasyn, y_adasyn),
81     'Random Oversample': (X_random_over, y_random_over),
82     'Random Undersample': (X_random_under, y_random_under)
83 }
84
85 for name, (X_data, y_data) in techniques.items():
86     scores = cross_val_score(RandomForestClassifier(random_state=42),
87                             X_data, y_data, cv=5, scoring='f1_macro')
88     print(f {name} - F1 Score: {scores.mean():.4f} (+/- {scores.std() *
89         2:.4f}) )
90
91 # Cost-sensitive learning with custom weights
92 def custom_class_weight(y, multiplier=2):
93     Create custom class weights based on class frequency
94     unique, counts = np.unique(y, return_counts=True)
95     total = len(y)
96     weights = {}
97     for i, (cls, count) in enumerate(zip(unique, counts)):
98         # Higher weight for minority classes
99         weight = (total / (len(unique) * count)) * multiplier
100         weights[cls] = weight
101     return weights
102
103 custom_weights = custom_class_weight(y, multiplier=3)
104 weighted_model_custom = RandomForestClassifier(class_weight=custom_weights
105     , random_state=42)
106
107 # Threshold tuning for imbalanced datasets
108 from sklearn.metrics import precision_recall_curve
109
110 def find_optimal_threshold(y_true, y_proba):
111     Find optimal threshold to maximize F1 score
112     precisions, recalls, thresholds = precision_recall_curve(y_true,
113         y_proba)
114     f1_scores = 2 * (precisions * recalls) / (precisions + recalls)
115     optimal_idx = np.argmax(f1_scores)
116     optimal_threshold = thresholds[optimal_idx]
117     return optimal_threshold
118
119 # Example usage
120 model = RandomForestClassifier(random_state=42)
121 model.fit(X_train, y_train)
122 y_proba = model.predict_proba(X_test)[: , 1]
123 optimal_threshold = find_optimal_threshold(y_test, y_proba)
124 print(f Optimal threshold: {optimal_threshold:.4f} )
125
126 # Apply optimal threshold

```

```
124 y_pred_optimal = (y_proba >= optimal_threshold).astype(int)
```

Listing 16: Handle imbalanced datasets effectively

### 3.5.2 3.5.2 Time Series Analysis

```
1 from sklearn.model_selection import TimeSeriesSplit
2 from sklearn.metrics import mean_squared_error
3 import pandas as pd
4 import numpy as np
5
6 # Time series cross-validation
7 def time_series_cv(model, X, y, n_splits=5):
8     Perform time series cross-validation
9     tscv = TimeSeriesSplit(n_splits=n_splits)
10    scores = []
11
12    for train_index, test_index in tscv.split(X):
13        X_train, X_test = X[train_index], X[test_index]
14        y_train, y_test = y[train_index], y[test_index]
15
16        model.fit(X_train, y_train)
17        y_pred = model.predict(X_test)
18        score = mean_squared_error(y_test, y_pred)
19        scores.append(score)
20
21    return np.array(scores)
22
23 # Create time series features
24 def create_time_series_features(df, target_col, lags=[1, 2, 3, 7],
25                               window_sizes=[3, 7, 14]):
26     Create lag and rolling window features for time series
27     df_features = df.copy()
28
29     # Lag features
30     for lag in lags:
31         df_features[f'{target_col}_lag_{lag}'] = df_features[target_col].
32         shift(lag)
33
34     # Rolling window features
35     for window in window_sizes:
36         df_features[f'{target_col}_rolling_mean_{window}'] = df_features[
37         target_col].rolling(window=window).mean()
38         df_features[f'{target_col}_rolling_std_{window}'] = df_features[
39         target_col].rolling(window=window).std()
40         df_features[f'{target_col}_rolling_min_{window}'] = df_features[
41         target_col].rolling(window=window).min()
42         df_features[f'{target_col}_rolling_max_{window}'] = df_features[
43         target_col].rolling(window=window).max()
44
45     # Time-based features
46     if 'date' in df_features.columns:
47         df_features['date'] = pd.to_datetime(df_features['date'])
```

```

42     df_features['year'] = df_features['date'].dt.year
43     df_features['month'] = df_features['date'].dt.month
44     df_features['day'] = df_features['date'].dt.day
45     df_features['dayofweek'] = df_features['date'].dt.dayofweek
46     df_features['quarter'] = df_features['date'].dt.quarter
47     df_features['is_weekend'] = df_features['dayofweek'].isin([5, 6]).
    astype(int)
48
49     return df_features
50
51 # Example usage
52 # Assuming you have a time series dataframe with 'date' and 'value'
    columns
53 # ts_features = create_time_series_features(time_series_df, 'value')
54
55 # Time series forecasting pipeline
56 from sklearn.preprocessing import StandardScaler
57 from sklearn.linear_model import LinearRegression
58 from sklearn.ensemble import RandomForestRegressor
59
60 def time_series_forecast_pipeline(X_train, y_train, X_test, model_type='rf
    '):
61     Pipeline for time series forecasting
62
63     # Scale features
64     scaler = StandardScaler()
65     X_train_scaled = scaler.fit_transform(X_train)
66     X_test_scaled = scaler.transform(X_test)
67
68     # Choose model
69     if model_type == 'lr':
70         model = LinearRegression()
71     elif model_type == 'rf':
72         model = RandomForestRegressor(n_estimators=100, random_state=42)
73     else:
74         raise ValueError( model_type must be 'lr' or 'rf' )
75
76     # Train model
77     model.fit(X_train_scaled, y_train)
78
79     # Make predictions
80     y_pred = model.predict(X_test_scaled)
81
82     return y_pred, model, scaler
83
84 # Multi-step forecasting
85 def multi_step_forecast(model, scaler, X_last, steps=5):
86     Multi-step forecasting for time series
87     predictions = []
88     X_current = X_last.copy()
89
90     for _ in range(steps):
91         X_scaled = scaler.transform(X_current.reshape(1, -1))
92         pred = model.predict(X_scaled)[0]

```



```

93     predictions.append(pred)
94
95     # Update X_current for next prediction (simplified approach)
96     # In practice, you'd need to update the lag features properly
97     X_current = np.roll(X_current, -1)
98     X_current[-1] = pred
99
100     return np.array(predictions)
101
102 # Seasonal decomposition features
103 from statsmodels.tsa.seasonal import seasonal_decompose
104
105 def add_seasonal_features(df, target_col, period=12):
106     Add seasonal decomposition features
107     decomposition = seasonal_decompose(df[target_col], model='additive',
108     period=period)
109
110     df['trend'] = decomposition.trend
111     df['seasonal'] = decomposition.seasonal
112     df['residual'] = decomposition.resid
113
114     return df
115
116 # Anomaly detection in time series
117 from sklearn.ensemble import IsolationForest
118 from sklearn.svm import OneClassSVM
119
120 def detect_time_series_anomalies(X, method='isolation_forest'):
121     Detect anomalies in time series data
122
123     if method == 'isolation_forest':
124         detector = IsolationForest(contamination=0.1, random_state=42)
125     elif method == 'one_class_svm':
126         detector = OneClassSVM(nu=0.1)
127     else:
128         raise ValueError( method must be 'isolation_forest' or '
129         one_class_svm' )
130
131     # Fit and predict
132     anomalies = detector.fit_predict(X)
133
134     # Convert to boolean (True for anomalies)
135     return anomalies == -1
136
137 # Time series clustering
138 from sklearn.cluster import KMeans
139 from sklearn.metrics import silhouette_score
140
141 def cluster_time_series_patterns(X_sequences, n_clusters=3):
142     Cluster time series patterns
143
144     # Flatten sequences for clustering
145     X_flat = X_sequences.reshape(X_sequences.shape[0], -1)

```

```

145     # Apply K-means clustering
146     kmeans = KMeans(n_clusters=n_clusters, random_state=42)
147     cluster_labels = kmeans.fit_predict(X_flat)
148
149     # Calculate silhouette score
150     sil_score = silhouette_score(X_flat, cluster_labels)
151
152     return cluster_labels, sil_score, kmeans
153
154 # Rolling window cross-validation
155 def rolling_window_cv(model, X, y, window_size=0.8, step_size=0.1):
156     Rolling window cross-validation for time series
157     n_samples = len(X)
158     window_length = int(window_size * n_samples)
159     step_length = int(step_size * n_samples)
160
161     scores = []
162     start = 0
163
164     while start + window_length < n_samples:
165         end = start + window_length
166
167         # Train on window
168         X_train, y_train = X[start:end], y[start:end]
169         X_test, y_test = X[end:end+step_length], y[end:end+step_length]
170
171         if len(X_test) > 0:
172             model.fit(X_train, y_train)
173             y_pred = model.predict(X_test)
174             score = mean_squared_error(y_test, y_pred)
175             scores.append(score)
176
177         start += step_length
178
179     return np.array(scores)

```

Listing 17: Time series machine learning approaches

## 4 4. Best Practices

### 4.1 4.1 Data Preprocessing Best Practices

#### Best Practice 4.1.1: Always Split Data Before Preprocessing

**Why:** Prevents data leakage by ensuring test data doesn't influence training preprocessing.

**How:** Split data first, then apply preprocessing separately to train and test sets.

```

1 # WRONG WAY - Data leakage risk
2 scaler = StandardScaler()
3 X_scaled = scaler.fit_transform(X) # Using entire dataset
4 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size
    =0.2)

```

```

5
6 # CORRECT WAY - No data leakage
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)
8
9 # Fit preprocessing on training data only
10 scaler = StandardScaler()
11 X_train_scaled = scaler.fit_transform(X_train)
12 X_test_scaled = scaler.transform(X_test) # Only transform, don't fit

```

Listing 18: Correct data splitting and preprocessing

### Best Practice 4.1.2: Handle Missing Values Strategically

**Why:** Different missing data patterns require different handling approaches.

**How:** Analyze missing patterns and choose appropriate imputation strategies.

```

1 # Analyze missing data patterns
2 def analyze_missing_data(df):
3     Analyze missing data patterns in dataframe
4     missing_stats = df.isnull().sum()
5     missing_percent = 100 * missing_stats / len(df)
6
7     missing_table = pd.DataFrame({
8         'Missing Count': missing_stats,
9         'Missing Percentage': missing_percent
10    }).sort_values('Missing Percentage', ascending=False)
11
12    return missing_table[missing_table['Missing Count'] > 0]
13
14 # Choose imputation strategy based on data characteristics
15 def smart_impute(df):
16     Smart imputation based on data types and missing patterns
17     df_imputed = df.copy()
18
19     for column in df.columns:
20         missing_percent = df[column].isnull().sum() / len(df) * 100
21
22         if missing_percent > 50:
23             print(f Warning: {column} has {missing_percent:.1f}% missing
24 values )
25             # Consider dropping column
26         elif missing_percent > 5:
27             print(f High missing rate in {column}: {missing_percent:.1f}%
28 )
29             # Use advanced imputation
30         else:
31             # Standard imputation
32             if df[column].dtype in ['int64', 'float64']:
33                 df_imputed[column].fillna(df[column].median(), inplace=
34 True)
35             else:
36                 mode_value = df[column].mode()

```

```

34         if len(mode_value) > 0:
35             df_imputed[column].fillna(mode_value[0], inplace=True)
36
37     return df_imputed
38
39 # Advanced imputation with domain knowledge
40 def domain_aware_impute(df, domain_rules):
41     Impute missing values using domain-specific rules
42     df_imputed = df.copy()
43
44     for column, rule in domain_rules.items():
45         if column in df.columns:
46             if rule['method'] == 'constant':
47                 df_imputed[column].fillna(rule['value'], inplace=True)
48             elif rule['method'] == 'conditional':
49                 condition_col = rule['condition_column']
50                 mapping = rule['mapping']
51                 for cond_value, fill_value in mapping.items():
52                     mask = (df_imputed[condition_col] == cond_value) &
df_imputed[column].isnull()
53                     df_imputed.loc[mask, column] = fill_value
54
55     return df_imputed

```

Listing 19: Comprehensive missing data handling

## Best Practice 4.1.3: Validate Data Quality Before Modeling

**Why:** Poor data quality leads to poor model performance and unreliable results.

**How:** Implement comprehensive data validation checks before training models.

```

1 def validate_data_quality(df, validation_rules=None):
2     Comprehensive data quality validation
3
4     validation_report = {
5         'timestamp': pd.Timestamp.now(),
6         'total_rows': len(df),
7         'total_columns': len(df.columns),
8         'issues': []
9     }
10
11     # Basic structure validation
12     if not isinstance(df, pd.DataFrame):
13         validation_report['issues'].append( Data is not a pandas DataFrame
14     )
15     return validation_report
16
17     # Check for duplicate rows
18     duplicates = df.duplicated().sum()
19     if duplicates > 0:
20         validation_report['issues'].append(f Found {duplicates} duplicate
rows )

```

```

21     # Check for completely empty columns
22     empty_cols = df.columns[df.isnull().all()].tolist()
23     if empty_cols:
24         validation_report['issues'].append(f Empty columns: {empty_cols} )
25
26     # Data type validation
27     for column in df.columns:
28         # Check for mixed data types
29         unique_types = df[column].apply(type).unique()
30         if len(unique_types) > 1:
31             validation_report['issues'].append(f Mixed types in column {
column}: {unique_types} )
32
33         # Check for extreme values in numeric columns
34         if df[column].dtype in ['int64', 'float64']:
35             Q1 = df[column].quantile(0.25)
36             Q3 = df[column].quantile(0.75)
37             IQR = Q3 - Q1
38             outliers = df[(df[column] < (Q1 - 1.5 * IQR)) | (df[column] >
(Q3 + 1.5 * IQR))]
39             if len(outliers) > 0:
40                 outlier_percent = len(outliers) / len(df) * 100
41                 if outlier_percent > 5: # More than 5% outliers
42                     validation_report['issues'].append(
43                         f High outlier rate in {column}: {outlier_percent
:.1f}%
44                     )
45
46     # Custom validation rules
47     if validation_rules:
48         for rule in validation_rules:
49             try:
50                 result = df.eval(rule['condition'])
51                 failed_count = (~result).sum()
52                 if failed_count > 0:
53                     validation_report['issues'].append(
54                         f Rule '{rule['name']}' failed for {failed_count}
rows
55                     )
56             except Exception as e:
57                 validation_report['issues'].append(f Rule '{rule['name']}'
error: {str(e)} )
58
59     return validation_report
60
61 # Example usage
62 validation_rules = [
63     {'name': 'age_positive', 'condition': 'age > 0'},
64     {'name': 'income_non_negative', 'condition': 'income >= 0'},
65     {'name': 'valid_category', 'condition': 'category.isin(['A', 'B', 'C
']) }
66 ]
67
68 quality_report = validate_data_quality(data, validation_rules)

```

```

69 if quality_report['issues']:
70     print( Data Quality Issues Found: )
71     for issue in quality_report['issues']:
72         print(f      - {issue} )
73 else:
74     print( Data quality validation passed! )

```

Listing 20: Data quality validation framework

## 4.2 Model Training Best Practices

### Best Practice 4.2.1: Use Proper Cross-Validation Strategies

**Why:** Proper validation prevents overfitting and provides reliable performance estimates.

**How:** Choose appropriate cross-validation methods based on data characteristics.

```

1 from sklearn.model_selection import (StratifiedKFold, TimeSeriesSplit,
2                                     LeaveOneOut, ShuffleSplit, GroupKFold)
3
4 def choose_cv_strategy(data_type, n_samples, groups=None):
5     Choose appropriate cross-validation strategy based on data
6     characteristics
7
8     if data_type == 'time_series':
9         return TimeSeriesSplit(n_splits=5)
10    elif data_type == 'grouped':
11        if groups is None:
12            raise ValueError( Groups required for grouped data )
13        return GroupKFold(n_splits=5)
14    elif data_type == 'classification':
15        if n_samples < 1000:
16            return StratifiedKFold(n_splits=5, shuffle=True, random_state
17                                  =42)
18        else:
19            return StratifiedKFold(n_splits=10, shuffle=True, random_state
20                                  =42)
21    elif data_type == 'regression':
22        if n_samples < 1000:
23            return KFold(n_splits=5, shuffle=True, random_state=42)
24        else:
25            return KFold(n_splits=10, shuffle=True, random_state=42)
26    else:
27        return ShuffleSplit(n_splits=5, test_size=0.2, random_state=42)
28
29 # Nested cross-validation for unbiased model evaluation
30 def nested_cross_validation(model, param_grid, X, y, outer_cv=None,
31                             inner_cv=None):
32     Perform nested cross-validation for unbiased model evaluation
33
34     if outer_cv is None:
35         outer_cv = KFold(n_splits=5, shuffle=True, random_state=42)
36     if inner_cv is None:
37         inner_cv = KFold(n_splits=3, shuffle=True, random_state=42)

```

```

34
35     # Outer loop for performance evaluation
36     outer_scores = []
37
38     for train_idx, test_idx in outer_cv.split(X):
39         X_train_outer, X_test_outer = X[train_idx], X[test_idx]
40         y_train_outer, y_test_outer = y[train_idx], y[test_idx]
41
42         # Inner loop for hyperparameter tuning
43         grid_search = GridSearchCV(
44             model, param_grid, cv=inner_cv, scoring='accuracy', n_jobs=-1
45         )
46         grid_search.fit(X_train_outer, y_train_outer)
47
48         # Evaluate best model on outer test set
49         best_model = grid_search.best_estimator_
50         score = best_model.score(X_test_outer, y_test_outer)
51         outer_scores.append(score)
52
53     return np.array(outer_scores)
54
55 # Example usage
56 param_grid = {
57     'n_estimators': [50, 100, 200],
58     'max_depth': [3, 5, 7, None]
59 }
60
61 nested_scores = nested_cross_validation(
62     RandomForestClassifier(random_state=42),
63     param_grid,
64     X,
65     y
66 )
67
68 print(f Nested CV Score: {nested_scores.mean():.4f} (+/- {nested_scores.
    std() * 2:.4f}) )

```

Listing 21: Advanced cross-validation strategies

#### Best Practice 4.2.2: Implement Early Stopping and Regularization

**Why:** Prevents overfitting and reduces training time for iterative algorithms.

**How:** Use early stopping criteria and regularization parameters appropriately.

```

1 from sklearn.ensemble import GradientBoostingClassifier
2 from sklearn.linear_model import Ridge, Lasso, ElasticNet
3 from sklearn.neural_network import MLPClassifier
4
5 # Early stopping for gradient boosting
6 gb_model = GradientBoostingClassifier(
7     n_estimators=1000,
8     learning_rate=0.1,
9     max_depth=3,

```

```

10     random_state=42,
11     validation_fraction=0.1, # Use 10% of training data for validation
12     n_iter_no_change=10,      # Stop if no improvement for 10 iterations
13     tol=1e-4                  # Tolerance for early stopping
14 )
15
16 # Regularization for linear models
17 ridge_model = Ridge(alpha=1.0)      # L2 regularization
18 lasso_model = Lasso(alpha=1.0)      # L1 regularization
19 elastic_model = ElasticNet(alpha=1.0, l1_ratio=0.5) # Combined L1/L2
20
21 # Early stopping for neural networks
22 nn_model = MLPClassifier(
23     hidden_layer_sizes=(100, 50),
24     max_iter=1000,
25     early_stopping=True,           # Enable early stopping
26     validation_fraction=0.1,      # Validation split
27     n_iter_no_change=10,          # Patience
28     tol=1e-4,
29     random_state=42
30 )
31
32 # Custom early stopping callback
33 class EarlyStoppingCallback:
34     def __init__(self, patience=10, min_delta=1e-4):
35         self.patience = patience
36         self.min_delta = min_delta
37         self.best_score = None
38         self.wait = 0
39
40     def __call__(self, current_score):
41         if self.best_score is None:
42             self.best_score = current_score
43         elif current_score - self.best_score < self.min_delta:
44             self.wait += 1
45             if self.wait >= self.patience:
46                 return True # Stop training
47         else:
48             self.best_score = current_score
49             self.wait = 0
50         return False # Continue training
51
52 # Adaptive regularization strength selection
53 def find_optimal_regularization(X, y, model_type='ridge'):
54     Find optimal regularization strength using cross-validation
55
56     if model_type == 'ridge':
57         alphas = np.logspace(-4, 4, 50)
58         scores = []
59
60         for alpha in alphas:
61             model = Ridge(alpha=alpha)
62             cv_scores = cross_val_score(model, X, y, cv=5, scoring='
neg_mean_squared_error')

```



```

63         scores.append(cv_scores.mean())
64
65     optimal_alpha = alphas[np.argmax(scores)]
66     return optimal_alpha
67
68     elif model_type == 'lasso':
69         # Use LassoCV for automatic selection
70         lasso_cv = LassoCV(cv=5, random_state=42)
71         lasso_cv.fit(X, y)
72         return lasso_cv.alpha_
73
74 # Example usage
75 optimal_alpha = find_optimal_regularization(X_train, y_train, 'ridge')
76 print(f Optimal regularization strength: {optimal_alpha:.6f} )

```

Listing 22: Early stopping and regularization implementation

### Best Practice 4.2.3: Monitor Training Progress and Diagnose Issues

**Why:** Helps identify problems early and optimize model performance.

**How:** Track metrics during training and implement diagnostic tools.

```

1 import matplotlib.pyplot as plt
2 from sklearn.model_selection import learning_curve, validation_curve
3 import time
4
5 # Learning curves to diagnose bias/variance
6 def plot_learning_curve(estimator, X, y, cv=5, train_sizes=np.linspace
7     (0.1, 1.0, 10)):
8     Plot learning curves to diagnose bias/variance problems
9
10    train_sizes, train_scores, val_scores = learning_curve(
11        estimator, X, y, cv=cv, train_sizes=train_sizes, n_jobs=-1
12    )
13
14    train_mean = np.mean(train_scores, axis=1)
15    train_std = np.std(train_scores, axis=1)
16    val_mean = np.mean(val_scores, axis=1)
17    val_std = np.std(val_scores, axis=1)
18
19    plt.figure(figsize=(10, 6))
20    plt.plot(train_sizes, train_mean, 'o-', color='blue', label='Training
21        score')
22    plt.fill_between(train_sizes, train_mean - train_std, train_mean +
23        train_std, alpha=0.1, color='blue')
24
25    plt.plot(train_sizes, val_mean, 'o-', color='red', label='Cross-
26        validation score')
27    plt.fill_between(train_sizes, val_mean - val_std, val_mean + val_std,
28        alpha=0.1, color='red')
29
30    plt.xlabel('Training Set Size')
31    plt.ylabel('Score')

```

```

27 plt.title('Learning Curve')
28 plt.legend(loc='best')
29 plt.grid(True)
30 plt.show()
31
32 # Diagnose problems
33 final_train_score = train_mean[-1]
34 final_val_score = val_mean[-1]
35
36 if final_train_score < 0.8:
37     print( High bias (underfitting) - Model is too simple )
38 elif final_train_score - final_val_score > 0.1:
39     print( High variance (overfitting) - Model is too complex )
40 else:
41     print( Good fit - Balanced bias and variance )
42
43 # Validation curves to optimize hyperparameters
44 def plot_validation_curve(estimator, X, y, param_name, param_range, cv=5):
45     Plot validation curves for hyperparameter tuning
46
47     train_scores, val_scores = validation_curve(
48         estimator, X, y, param_name=param_name, param_range=param_range,
49         cv=cv, n_jobs=-1
50     )
51
52     train_mean = np.mean(train_scores, axis=1)
53     train_std = np.std(train_scores, axis=1)
54     val_mean = np.mean(val_scores, axis=1)
55     val_std = np.std(val_scores, axis=1)
56
57     plt.figure(figsize=(10, 6))
58     plt.semilogx(param_range, train_mean, 'o-', color='blue', label='
Training score')
59     plt.fill_between(param_range, train_mean - train_std, train_mean +
60 train_std, alpha=0.1, color='blue')
61
62     plt.semilogx(param_range, val_mean, 'o-', color='red', label='Cross-
63 validation score')
64     plt.fill_between(param_range, val_mean - val_std, val_mean + val_std,
65 alpha=0.1, color='red')
66
67     plt.xlabel(param_name)
68     plt.ylabel('Score')
69     plt.title(f'Validation Curve for {param_name}')
70     plt.legend(loc='best')
71     plt.grid(True)
72     plt.show()
73
74 # Find optimal parameter
75 optimal_idx = np.argmax(val_mean)
76 optimal_param = param_range[optimal_idx]
77 print(f Optimal {param_name}: {optimal_param} )
78
79 return optimal_param

```

```

76
77 # Training time monitoring
78 def time_model_training(model, X, y, n_repeats=3):
79     Measure model training time
80
81     times = []
82     for _ in range(n_repeats):
83         start_time = time.time()
84         model.fit(X, y)
85         end_time = time.time()
86         times.append(end_time - start_time)
87
88     avg_time = np.mean(times)
89     std_time = np.std(times)
90
91     print(f Average training time: {avg_time:.4f} seconds (+/- {std_time
92         :.4f}) )
93     return avg_time, std_time
94
95 # Model complexity analysis
96 def analyze_model_complexity(X, y, model_class, param_name, param_values):
97     Analyze how model complexity affects performance
98
99     train_scores = []
100     val_scores = []
101     complexities = []
102
103     for param_value in param_values:
104         model = model_class(**{param_name: param_value})
105
106         # Training score
107         model.fit(X, y)
108         train_score = model.score(X, y)
109         train_scores.append(train_score)
110
111         # Validation score
112         val_score = cross_val_score(model, X, y, cv=5).mean()
113         val_scores.append(val_score)
114
115         # Model complexity (example: number of parameters for tree-based
116         models)
117         if hasattr(model, 'n_estimators'):
118             complexity = model.n_estimators
119         elif hasattr(model, 'max_depth'):
120             complexity = model.max_depth if model.max_depth else float('
121             inf')
122         else:
123             complexity = param_value
124
125         complexities.append(complexity)
126
127     # Plot complexity vs performance
128     plt.figure(figsize=(12, 4))

```

```

127 plt.subplot(1, 2, 1)
128 plt.plot(complexities, train_scores, 'o-', label='Training Score')
129 plt.plot(complexities, val_scores, 'o-', label='Validation Score')
130 plt.xlabel('Model Complexity')
131 plt.ylabel('Score')
132 plt.title('Model Complexity vs Performance')
133 plt.legend()
134 plt.grid(True)
135
136 plt.subplot(1, 2, 2)
137 plt.plot(complexities, np.array(train_scores) - np.array(val_scores),
138 'o-')
139 plt.xlabel('Model Complexity')
140 plt.ylabel('Train-Validation Gap')
141 plt.title('Overfitting Analysis')
142 plt.grid(True)
143
144 plt.tight_layout()
145 plt.show()
146
147 # Example usage
148 plot_learning_curve(RandomForestClassifier(random_state=42), X, y)
149 optimal_n_estimators = plot_validation_curve(
150     RandomForestClassifier(random_state=42),
151     X, y,
152     'n_estimators',
153     [10, 50, 100, 200, 500]
154 )
155 time_model_training(RandomForestClassifier(n_estimators=100), X, y)

```

Listing 23: Training monitoring and diagnostics

### 4.3 Model Evaluation Best Practices

#### Best Practice 4.3.1: Use Multiple Evaluation Metrics

**Why:** Single metrics can be misleading; multiple metrics provide comprehensive model assessment.

**How:** Select appropriate metrics based on problem type and business requirements.

```

1 from sklearn.metrics import (accuracy_score, precision_score, recall_score,
2                               f1_score,
3                               roc_auc_score, precision_recall_curve,
4                               roc_curve,
5                               mean_squared_error, mean_absolute_error,
6                               r2_score,
7                               confusion_matrix, classification_report)
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10
11 def comprehensive_classification_evaluation(y_true, y_pred, y_proba=None,
12                                           class_names=None):
13     """Comprehensive evaluation for classification models"""

```

```

10
11     results = {}
12
13     # Basic metrics
14     results['accuracy'] = accuracy_score(y_true, y_pred)
15     results['precision_macro'] = precision_score(y_true, y_pred, average='
macro')
16     results['recall_macro'] = recall_score(y_true, y_pred, average='macro'
)
17     results['f1_macro'] = f1_score(y_true, y_pred, average='macro')
18     results['precision_weighted'] = precision_score(y_true, y_pred,
average='weighted')
19     results['recall_weighted'] = recall_score(y_true, y_pred, average='
weighted')
20     results['f1_weighted'] = f1_score(y_true, y_pred, average='weighted')
21
22     # Per-class metrics
23     if class_names is None:
24         class_names = [f'Class_{i}' for i in range(len(np.unique(y_true)))
]
25
26     precision_per_class = precision_score(y_true, y_pred, average=None)
27     recall_per_class = recall_score(y_true, y_pred, average=None)
28     f1_per_class = f1_score(y_true, y_pred, average=None)
29
30     for i, class_name in enumerate(class_names):
31         results[f'precision_{class_name}'] = precision_per_class[i]
32         results[f'recall_{class_name}'] = recall_per_class[i]
33         results[f'f1_{class_name}'] = f1_per_class[i]
34
35     # Probability-based metrics (if probabilities provided)
36     if y_proba is not None:
37         if len(np.unique(y_true)) == 2:
38             # Binary classification
39             results['roc_auc'] = roc_auc_score(y_true, y_proba[:, 1])
40
41             # Find optimal threshold
42             precision_vals, recall_vals, thresholds =
precision_recall_curve(y_true, y_proba[:, 1])
43             f1_scores = 2 * (precision_vals * recall_vals) / (
precision_vals + recall_vals)
44             optimal_idx = np.argmax(f1_scores)
45             results['optimal_threshold'] = thresholds[optimal_idx]
46             results['optimal_f1'] = f1_scores[optimal_idx]
47         else:
48             # Multi-class
49             results['roc_auc_ovr'] = roc_auc_score(y_true, y_proba,
multi_class='ovr')
50             results['roc_auc_ovo'] = roc_auc_score(y_true, y_proba,
multi_class='ovo')
51
52     # Print results
53     print( == CLASSIFICATION EVALUATION RESULTS == )
54     print(f Accuracy: {results['accuracy']:.4f} )

```

```

55     print(f Macro Precision: {results['precision_macro']:.4f} )
56     print(f Macro Recall: {results['recall_macro']:.4f} )
57     print(f Macro F1-Score: {results['f1_macro']:.4f} )
58     print(f Weighted Precision: {results['precision_weighted']:.4f} )
59     print(f Weighted Recall: {results['recall_weighted']:.4f} )
60     print(f Weighted F1-Score: {results['f1_weighted']:.4f} )
61
62     if 'roc_auc' in results:
63         print(f AUC-ROC: {results['roc_auc']:.4f} )
64         print(f Optimal Threshold: {results['optimal_threshold']:.4f} )
65         print(f Optimal F1-Score: {results['optimal_f1']:.4f} )
66
67     # Per-class results
68     print( \nPer-Class Results: )
69     for class_name in class_names:
70         print(f {class_name}: )
71         print(f Precision: {results[f'precision_{class_name}']:.4f} )
72         print(f Recall: {results[f'recall_{class_name}']:.4f} )
73         print(f F1-Score: {results[f'f1_{class_name}']:.4f} )
74
75     return results
76
77 def comprehensive_regression_evaluation(y_true, y_pred):
78     Comprehensive evaluation for regression models
79
80     results = {}
81
82     # Error metrics
83     results['mse'] = mean_squared_error(y_true, y_pred)
84     results['rmse'] = np.sqrt(results['mse'])
85     results['mae'] = mean_absolute_error(y_true, y_pred)
86     results['r2'] = r2_score(y_true, y_pred)
87
88     # Relative errors
89     results['mape'] = np.mean(np.abs((y_true - y_pred) / y_true)) * 100 #
90     Mean Absolute Percentage Error
91
92     # Residual analysis
93     residuals = y_true - y_pred
94     results['residual_std'] = np.std(residuals)
95     results['residual_mean'] = np.mean(residuals)
96
97     # Print results
98     print( === REGRESSION EVALUATION RESULTS === )
99     print(f Mean Squared Error: {results['mse']:.4f} )
100    print(f Root Mean Squared Error: {results['rmse']:.4f} )
101    print(f Mean Absolute Error: {results['mae']:.4f} )
102    print(f R Score: {results['r2']:.4f} )
103    print(f Mean Absolute Percentage Error: {results['mape']:.2f}% )
104    print(f Residual Standard Deviation: {results['residual_std']:.4f} )
105    print(f Residual Mean: {results['residual_mean']:.4f} )
106
107    # Visualization
108    fig, axes = plt.subplots(2, 2, figsize=(12, 10))

```

```

108
109     # Predicted vs Actual
110     axes[0, 0].scatter(y_true, y_pred, alpha=0.6)
111     axes[0, 0].plot([y_true.min(), y_true.max()], [y_true.min(), y_true.
112 max()], 'r--', lw=2)
113     axes[0, 0].set_xlabel('Actual Values')
114     axes[0, 0].set_ylabel('Predicted Values')
115     axes[0, 0].set_title('Predicted vs Actual')
116
117     # Residuals vs Predicted
118     axes[0, 1].scatter(y_pred, residuals, alpha=0.6)
119     axes[0, 1].axhline(y=0, color='r', linestyle='--')
120     axes[0, 1].set_xlabel('Predicted Values')
121     axes[0, 1].set_ylabel('Residuals')
122     axes[0, 1].set_title('Residuals vs Predicted')
123
124     # Residuals histogram
125     axes[1, 0].hist(residuals, bins=30, edgecolor='black')
126     axes[1, 0].set_xlabel('Residuals')
127     axes[1, 0].set_ylabel('Frequency')
128     axes[1, 0].set_title('Residuals Distribution')
129
130     # Q-Q plot
131     from scipy import stats
132     stats.probplot(residuals, dist= norm , plot=axes[1, 1])
133     axes[1, 1].set_title('Q-Q Plot of Residuals')
134
135     plt.tight_layout()
136     plt.show()
137
138     return results
139
140 def plot_confusion_matrix(y_true, y_pred, class_names=None, normalize=
False):
141     """
142     Plot confusion matrix with optional normalization
143
144     cm = confusion_matrix(y_true, y_pred)
145
146     if normalize:
147         cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
148         fmt = '.2f'
149         title = 'Normalized Confusion Matrix'
150     else:
151         fmt = 'd'
152         title = 'Confusion Matrix'
153
154     if class_names is None:
155         class_names = [f'Class_{i}' for i in range(len(cm))]
156
157     plt.figure(figsize=(8, 6))
158     sns.heatmap(cm, annot=True, fmt=fmt, cmap='Blues',
159                 xticklabels=class_names, yticklabels=class_names)
160     plt.title(title)
161     plt.ylabel('True Label')

```

```

160 plt.xlabel('Predicted Label')
161 plt.show()
162
163 def plot_roc_curves(y_true, y_proba, class_names=None):
164     Plot ROC curves for multi-class classification
165
166     from sklearn.preprocessing import label_binarize
167     from itertools import cycle
168
169     n_classes = len(np.unique(y_true))
170
171     if n_classes == 2:
172         # Binary classification
173         fpr, tpr, _ = roc_curve(y_true, y_proba[:, 1])
174         roc_auc = roc_auc_score(y_true, y_proba[:, 1])
175
176         plt.figure(figsize=(8, 6))
177         plt.plot(fpr, tpr, color='darkorange', lw=2,
178                 label=f'ROC curve (AUC = {roc_auc:.2f})')
179         plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
180         plt.xlim([0.0, 1.0])
181         plt.ylim([0.0, 1.05])
182         plt.xlabel('False Positive Rate')
183         plt.ylabel('True Positive Rate')
184         plt.title('Receiver Operating Characteristic')
185         plt.legend(loc= lower right )
186         plt.show()
187
188     else:
189         # Multi-class classification
190         y_bin = label_binarize(y_true, classes=np.unique(y_true))
191
192         if class_names is None:
193             class_names = [f'Class_{i}' for i in range(n_classes)]
194
195         fpr = dict()
196         tpr = dict()
197         roc_auc = dict()
198
199         for i in range(n_classes):
200             fpr[i], tpr[i], _ = roc_curve(y_bin[:, i], y_proba[:, i])
201             roc_auc[i] = roc_auc_score(y_bin[:, i], y_proba[:, i])
202
203         # Plot ROC curves
204         plt.figure(figsize=(10, 8))
205         colors = cycle(['aqua', 'darkorange', 'cornflowerblue', 'red', '
206 green'])
207
208         for i, color in zip(range(n_classes), colors):
209             plt.plot(fpr[i], tpr[i], color=color, lw=2,
210                     label=f'ROC curve of class {class_names[i]} (AUC = {
211 roc_auc[i]:.2f})')
210
211         plt.plot([0, 1], [0, 1], 'k--', lw=2)

```



```

212     plt.xlim([0.0, 1.0])
213     plt.ylim([0.0, 1.05])
214     plt.xlabel('False Positive Rate')
215     plt.ylabel('True Positive Rate')
216     plt.title('Multi-class ROC Curves')
217     plt.legend(loc= lower right )
218     plt.show()
219
220 # Example usage
221 # For classification
222 # results = comprehensive_classification_evaluation(y_test, y_pred,
223 #         y_proba, class_names=['Class_A', 'Class_B'])
224
225 # For regression
226 # results = comprehensive_regression_evaluation(y_test, y_pred)
227
228 # Additional visualizations
229 # plot_confusion_matrix(y_test, y_pred, class_names=['Class_A', 'Class_B',
230 #         ''])
231 # plot_roc_curves(y_test, y_proba, class_names=['Class_A', 'Class_B'])

```

Listing 24: Comprehensive model evaluation with multiple metrics

#### Best Practice 4.3.2: Validate Model Assumptions and Check Robustness

**Why:** Ensures model reliability and identifies potential failure modes.

**How:** Test model assumptions and evaluate performance under various conditions.

```

1  import numpy as np
2  import pandas as pd
3  from sklearn.model_selection import train_test_split
4  from sklearn.linear_model import LinearRegression, LogisticRegression
5  from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier
6  from scipy import stats
7  import matplotlib.pyplot as plt
8
9  def validate_linear_regression_assumptions(X, y, model):
10     Validate assumptions for linear regression
11
12     # Fit model
13     model.fit(X, y)
14     y_pred = model.predict(X)
15     residuals = y - y_pred
16
17     assumptions_validated = {}
18
19     # 1. Linearity
20     # Check if relationship between features and target is linear
21     plt.figure(figsize=(15, 5))
22
23     plt.subplot(1, 3, 1)
24     plt.scatter(y_pred, y, alpha=0.6)
25     plt.plot([y_pred.min(), y_pred.max()], [y_pred.min(), y_pred.max()], '

```

```

r--', lw=2)
26 plt.xlabel('Predicted Values')
27 plt.ylabel('Actual Values')
28 plt.title('Linearity Check')
29
30 # 2. Homoscedasticity (Constant Variance)
31 plt.subplot(1, 3, 2)
32 plt.scatter(y_pred, residuals, alpha=0.6)
33 plt.axhline(y=0, color='r', linestyle='--')
34 plt.xlabel('Predicted Values')
35 plt.ylabel('Residuals')
36 plt.title('Homoscedasticity Check')
37
38 # 3. Normality of Residuals
39 plt.subplot(1, 3, 3)
40 stats.probplot(residuals, dist= norm , plot=plt)
41 plt.title('Normality Check (Q-Q Plot)')
42
43 plt.tight_layout()
44 plt.show()
45
46 # Statistical tests
47 # Shapiro-Wilk test for normality
48 shapiro_stat, shapiro_p = stats.shapiro(residuals)
49 assumptions_validated['normality'] = {
50     'statistic': shapiro_stat,
51     'p_value': shapiro_p,
52     'valid': shapiro_p > 0.05 # Fail to reject normality
53 }
54
55 # Breusch-Pagan test for homoscedasticity (simplified)
56 # Correlation between squared residuals and predicted values
57 bp_stat = np.corrcoef(residuals**2, y_pred)[0, 1]
58 assumptions_validated['homoscedasticity'] = {
59     'correlation': bp_stat,
60     'valid': abs(bp_stat) < 0.3 # Low correlation indicates
homoscedasticity
61 }
62
63 # Linearity check using correlation
64 linearity_corr = np.corrcoef(y_pred, y)[0, 1]
65 assumptions_validated['linearity'] = {
66     'correlation': linearity_corr,
67     'valid': linearity_corr > 0.7 # High correlation indicates good
fit
68 }
69
70 # Print results
71 print( '=== LINEAR REGRESSION ASSUMPTIONS VALIDATION === ' )
72 for assumption, result in assumptions_validated.items():
73     status = VALID if result['valid'] else INVALID
74     print(f {assumption.capitalize()}: {status} )
75     for key, value in result.items():
76         if key != 'valid':

```

```

77         print(f'    {key}: {value} ')
78     print()
79
80     return assumptions_validated
81
82 def validate_logistic_regression_assumptions(X, y, model):
83     """Validate assumptions for logistic regression"""
84
85     # Fit model
86     model.fit(X, y)
87     y_pred_proba = model.predict_proba(X)[:, 1]
88
89     assumptions_validated = {}
90
91     # 1. Linearity of logit
92     # Check if log-odds are linearly related to predictors
93     log_odds = np.log(y_pred_proba / (1 - y_pred_proba + 1e-15)) # Add
94     small value to avoid log(0)
95
96     plt.figure(figsize=(12, 4))
97
98     # Plot log-odds vs each feature
99     n_features = min(X.shape[1], 6) # Limit to 6 features for
100     visualization
101     for i in range(n_features):
102         plt.subplot(1, n_features, i + 1)
103         plt.scatter(X[:, i], log_odds, alpha=0.6)
104         plt.xlabel(f'Feature {i}')
105         plt.ylabel('Log-Odds')
106         plt.title(f'Log-Odds vs Feature {i}')
107
108     plt.tight_layout()
109     plt.show()
110
111     # 2. Absence of multicollinearity
112     correlation_matrix = np.corrcoef(X.T)
113     max_corr = np.max(np.abs(correlation_matrix - np.eye(
114     correlation_matrix.shape[0])))
115     assumptions_validated['multicollinearity'] = {
116         'max_correlation': max_corr,
117         'valid': max_corr < 0.8 # Low multicollinearity
118     }
119
120     # 3. Large sample size
121     sample_size_ratio = len(y) / X.shape[1]
122     assumptions_validated['sample_size'] = {
123         'ratio': sample_size_ratio,
124         'valid': sample_size_ratio > 10 # At least 10 observations per
125     feature
126     }
127
128     # Print results
129     print( '=== LOGISTIC REGRESSION ASSUMPTIONS VALIDATION === ')
130     for assumption, result in assumptions_validated.items():

```

```

127     status = VALID if result['valid'] else INVALID
128     print(f {assumption.replace('_', ' ').title()}: {status} )
129     for key, value in result.items():
130         if key != 'valid':
131             print(f {key}: {value} )
132     print()
133
134     return assumptions_validated
135
136 def test_model_robustness(model, X_train, y_train, X_test, y_test,
137     noise_levels=[0.01, 0.05, 0.1]):
138     Test model robustness to input noise
139
140     # Fit model on original data
141     model.fit(X_train, y_train)
142     original_score = model.score(X_test, y_test)
143
144     robustness_results = {
145         'original_score': original_score,
146         'noise_levels': noise_levels,
147         'noisy_scores': [],
148         'score_degradation': []
149     }
150
151     print( === MODEL ROBUSTNESS TEST === )
152     print(f Original performance: {original_score:.4f} )
153     print()
154
155     for noise_level in noise_levels:
156         # Add noise to test data
157         noise = np.random.normal(0, noise_level, X_test.shape)
158         X_test_noisy = X_test + noise
159
160         # Evaluate on noisy data
161         noisy_score = model.score(X_test_noisy, y_test)
162         degradation = (original_score - noisy_score) / original_score *
100
163
164         robustness_results['noisy_scores'].append(noisy_score)
165         robustness_results['score_degradation'].append(degradation)
166
167         print(f Noise level {noise_level}: )
168         print(f Performance: {noisy_score:.4f} )
169         print(f Degradation: {degradation:.2f}% )
170         print()
171
172     # Plot results
173     plt.figure(figsize=(10, 6))
174     plt.plot([0] + noise_levels, [original_score] + robustness_results['
noisy_scores'], 'o-')
175     plt.xlabel('Noise Level')
176     plt.ylabel('Model Performance')
177     plt.title('Model Robustness to Input Noise')
178     plt.grid(True)

```

```

178     plt.show()
179
180     # Robustness assessment
181     avg_degradation = np.mean(robustness_results['score_degradation'])
182     if avg_degradation < 10:
183         robustness_level = High
184     elif avg_degradation < 25:
185         robustness_level = Medium
186     else:
187         robustness_level = Low
188
189     print(f Overall Robustness Level: {robustness_level} )
190     print(f Average Performance Degradation: {avg_degradation:.2f}% )
191
192     return robustness_results
193
194 def test_model_stability(model, X, y, n_bootstrap=100):
195     Test model stability using bootstrap sampling
196
197     original_model = model.__class__(**model.get_params())
198     original_model.fit(X, y)
199     original_coef = getattr(original_model, 'coef_', None)
200
201     bootstrap_scores = []
202     bootstrap_coefs = []
203
204     print( === MODEL STABILITY TEST (BOOTSTRAP) === )
205
206     for i in range(n_bootstrap):
207         # Bootstrap sample
208         indices = np.random.choice(len(X), len(X), replace=True)
209         X_boot = X[indices]
210         y_boot = y[indices]
211
212         # Fit model on bootstrap sample
213         boot_model = model.__class__(**model.get_params())
214         boot_model.fit(X_boot, y_boot)
215         bootstrap_scores.append(boot_model.score(X, y))
216
217         # Store coefficients if available
218         if hasattr(boot_model, 'coef_'):
219             bootstrap_coefs.append(boot_model.coef_)
220
221     # Calculate statistics
222     mean_score = np.mean(bootstrap_scores)
223     std_score = np.std(bootstrap_scores)
224
225     print(f Original performance: {original_model.score(X, y):.4f} )
226     print(f Bootstrap mean performance: {mean_score:.4f} (+/- {std_score
227 :.4f}) )
228     print(f Performance stability: {((1 - std_score/mean_score) * 100):.1f
229 }% )
230
231     # Coefficient stability (if applicable)

```

```

230     if bootstrap_coefs and original_coef is not None:
231         coef_matrix = np.array(bootstrap_coefs)
232         coef_std = np.std(coef_matrix, axis=0)
233         coef_cv = np.mean(coef_std / np.abs(np.mean(coef_matrix, axis=0) +
1e-10)) # Coefficient of variation
234
235     print(f Coefficient stability (CV): {coef_cv:.4f} )
236     if coef_cv < 0.1:
237         print( Coefficients are stable )
238     elif coef_cv < 0.3:
239         print( Coefficients are moderately stable )
240     else:
241         print( Coefficients are unstable )
242
243     # Plot score distribution
244     plt.figure(figsize=(10, 6))
245     plt.hist(bootstrap_scores, bins=30, alpha=0.7, edgecolor='black')
246     plt.axvline(original_model.score(X, y), color='red', linestyle='--',
247                 label=f'Original Score: {original_model.score(X, y):.4f}')
248     plt.axvline(mean_score, color='blue', linestyle='--',
249                 label=f'Mean Bootstrap Score: {mean_score:.4f}')
250     plt.xlabel('Model Performance')
251     plt.ylabel('Frequency')
252     plt.title('Bootstrap Performance Distribution')
253     plt.legend()
254     plt.grid(True)
255     plt.show()
256
257     return {
258         'bootstrap_scores': bootstrap_scores,
259         'bootstrap_coefs': bootstrap_coefs,
260         'mean_score': mean_score,
261         'std_score': std_score
262     }
263
264 # Example usage
265 # For linear regression
266 # lr_model = LinearRegression()
267 # assumptions = validate_linear_regression_assumptions(X, y, lr_model)
268
269 # For logistic regression
270 # log_model = LogisticRegression()
271 # assumptions = validate_logistic_regression_assumptions(X, y, log_model)
272
273 # Robustness testing
274 # rf_model = RandomForestClassifier(random_state=42)
275 # robustness_results = test_model_robustness(rf_model, X_train, y_train,
276     X_test, y_test)
277
277 # Stability testing
278 # stability_results = test_model_stability(rf_model, X, y)

```

Listing 25: Model assumption validation and robustness testing

## 4.4 4.4 Deployment and Production Best Practices

### Best Practice 4.4.1: Implement Proper Model Versioning and Tracking

**Why:** Enables reproducibility, rollback capabilities, and experiment tracking.

**How:** Use version control for code and models, and implement experiment tracking.

```

1 import joblib
2 import json
3 import os
4 from datetime import datetime
5 import git
6 import hashlib
7
8 class ModelVersionTracker:
9     Track model versions, parameters, and performance metrics
10
11     def __init__(self, tracking_dir='model_tracking'):
12         self.tracking_dir = tracking_dir
13         os.makedirs(tracking_dir, exist_ok=True)
14         self.experiments_file = os.path.join(tracking_dir, 'experiments.
15         json')
16
17         # Initialize experiments file if it doesn't exist
18         if not os.path.exists(self.experiments_file):
19             with open(self.experiments_file, 'w') as f:
20                 json.dump([], f)
21
22     def get_git_info(self):
23         Get current git commit information
24         try:
25             repo = git.Repo(search_parent_directories=True)
26             commit_hash = repo.head.object.hexsha
27             commit_message = repo.head.object.message.strip()
28             branch = repo.active_branch.name
29             return {
30                 'commit_hash': commit_hash,
31                 'commit_message': commit_message,
32                 'branch': branch
33             }
34         except:
35             return {
36                 'commit_hash': 'unknown',
37                 'commit_message': 'unknown',
38                 'branch': 'unknown'
39             }
40
41     def calculate_data_hash(self, X, y):
42         Calculate hash of training data for version tracking
43         # Convert to numpy arrays and flatten
44         X_flat = X.flatten() if hasattr(X, 'flatten') else np.array(X).
45         flatten()
46         y_flat = y.flatten() if hasattr(y, 'flatten') else np.array(y).

```

```

45     flatten()
46     # Combine and hash
47     combined = np.concatenate([X_flat, y_flat])
48     return hashlib.md5(combined.tobytes()).hexdigest()
49
50     def log_experiment(self, model, X, y, params, metrics, model_name,
51 notes= ):
52         Log experiment details
53
54         experiment = {
55             'experiment_id': datetime.now().strftime('%Y%m%d_%H%M%S'),
56             'timestamp': datetime.now().isoformat(),
57             'model_name': model_name,
58             'parameters': params,
59             'metrics': metrics,
60             'data_hash': self.calculate_data_hash(X, y),
61             'git_info': self.get_git_info(),
62             'notes': notes
63         }
64
65         # Load existing experiments
66         with open(self.experiments_file, 'r') as f:
67             experiments = json.load(f)
68
69         # Add new experiment
70         experiments.append(experiment)
71
72         # Save updated experiments
73         with open(self.experiments_file, 'w') as f:
74             json.dump(experiments, f, indent=2)
75
76         print(f Experiment logged: {experiment['experiment_id']} )
77         return experiment['experiment_id']
78
79     def save_model(self, model, experiment_id, model_name):
80         Save model with version tracking
81         model_dir = os.path.join(self.tracking_dir, 'models')
82         os.makedirs(model_dir, exist_ok=True)
83
84         model_filename = f {model_name}_{experiment_id}.pkl
85         model_path = os.path.join(model_dir, model_filename)
86
87         joblib.dump(model, model_path)
88         print(f Model saved: {model_path} )
89         return model_path
90
91     def load_model(self, model_name, experiment_id=None):
92         Load specific model version
93         model_dir = os.path.join(self.tracking_dir, 'models')
94
95         if experiment_id:
96             model_filename = f {model_name}_{experiment_id}.pkl
97             model_path = os.path.join(model_dir, model_filename)

```



```

97         if os.path.exists(model_path):
98             return joblib.load(model_path)
99         else:
100             raise FileNotFoundError(f Model {model_filename} not found
101     )
102     else:
103         # Load latest model
104         model_files = [f for f in os.listdir(model_dir) if f.
105 startswith(model_name)]
106         if not model_files:
107             raise FileNotFoundError(f No models found for {model_name}
108         )
109
110         # Sort by timestamp and load latest
111         model_files.sort(reverse=True)
112         latest_model = model_files[0]
113         model_path = os.path.join(model_dir, latest_model)
114         return joblib.load(model_path)
115
116 def get_experiment_history(self, model_name=None):
117     Get experiment history
118     with open(self.experiments_file, 'r') as f:
119         experiments = json.load(f)
120
121     if model_name:
122         experiments = [exp for exp in experiments if exp['model_name']
123 == model_name]
124
125     return experiments
126
127 def compare_experiments(self, experiment_ids):
128     Compare multiple experiments
129     with open(self.experiments_file, 'r') as f:
130         all_experiments = json.load(f)
131
132     experiments = [exp for exp in all_experiments if exp['
133 experiment_id'] in experiment_ids]
134
135     comparison = {}
136     for exp in experiments:
137         exp_id = exp['experiment_id']
138         comparison[exp_id] = {
139             'timestamp': exp['timestamp'],
140             'model_name': exp['model_name'],
141             'metrics': exp['metrics'],
142             'parameters': exp['parameters']
143         }
144
145     return comparison
146
147 # Advanced model versioning with MLflow-like functionality
148 class SimpleMLFlow:
149     Simple experiment tracking inspired by MLflow
150

```

```

146     def __init__(self, experiment_name= default ):
147         self.experiment_name = experiment_name
148         self.runs_dir = f mlruns/{experiment_name}
149         os.makedirs(self.runs_dir, exist_ok=True)
150         self.current_run = None
151
152     def start_run(self, run_name=None):
153         Start a new experiment run
154         run_id = datetime.now().strftime('%Y%m%d_%H%M%S_%f')
155         run_dir = os.path.join(self.runs_dir, run_id)
156         os.makedirs(run_dir, exist_ok=True)
157
158         self.current_run = {
159             'run_id': run_id,
160             'run_name': run_name or f run_{run_id} ,
161             'start_time': datetime.now().isoformat(),
162             'params': {},
163             'metrics': {},
164             'artifacts': {}
165         }
166
167         print(f Started run: {run_name or run_id} )
168         return self.current_run
169
170     def log_param(self, key, value):
171         Log a parameter
172         if self.current_run:
173             self.current_run['params'][key] = value
174             print(f Logged parameter: {key} = {value} )
175
176     def log_metric(self, key, value):
177         Log a metric
178         if self.current_run:
179             self.current_run['metrics'][key] = value
180             print(f Logged metric: {key} = {value} )
181
182     def log_artifact(self, local_path, artifact_path=None):
183         Log an artifact (file)
184         if self.current_run:
185             import shutil
186             artifact_name = artifact_path or os.path.basename(local_path)
187             run_artifacts_dir = os.path.join(self.runs_dir, self.
current_run['run_id'], 'artifacts')
188             os.makedirs(run_artifacts_dir, exist_ok=True)
189
190             artifact_dest = os.path.join(run_artifacts_dir, artifact_name)
191             shutil.copy2(local_path, artifact_dest)
192
193             self.current_run['artifacts'][artifact_name] = artifact_dest
194             print(f Logged artifact: {artifact_name} )
195
196     def end_run(self):
197         End current run and save metadata
198         if self.current_run:

```

```

199         # Save run metadata
200         run_metadata_path = os.path.join(self.runs_dir, self.
current_run['run_id'], 'meta.json')
201         self.current_run['end_time'] = datetime.now().isoformat()
202
203         with open(run_metadata_path, 'w') as f:
204             json.dump(self.current_run, f, indent=2)
205
206         print(f Ended run: {self.current_run['run_name']} )
207         self.current_run = None
208
209     def list_runs(self):
210         List all runs in experiment
211         runs = []
212         for run_id in os.listdir(self.runs_dir):
213             run_dir = os.path.join(self.runs_dir, run_id)
214             if os.path.isdir(run_dir):
215                 meta_file = os.path.join(run_dir, 'meta.json')
216                 if os.path.exists(meta_file):
217                     with open(meta_file, 'r') as f:
218                         meta = json.load(f)
219                         runs.append(meta)
220         return sorted(runs, key=lambda x: x['start_time'], reverse=True)
221
222 # Production model deployment utilities
223 class ModelDeployer:
224     Utilities for model deployment
225
226     def __init__(self, model_path, version_tracker=None):
227         self.model = joblib.load(model_path)
228         self.model_path = model_path
229         self.version_tracker = version_tracker
230         self.load_time = datetime.now()
231
232     def predict(self, X):
233         Make predictions with error handling
234         try:
235             predictions = self.model.predict(X)
236             return predictions
237         except Exception as e:
238             print(f Prediction error: {str(e)} )
239             raise
240
241     def predict_proba(self, X):
242         Make probability predictions if supported
243         if hasattr(self.model, 'predict_proba'):
244             try:
245                 probabilities = self.model.predict_proba(X)
246                 return probabilities
247             except Exception as e:
248                 print(f Probability prediction error: {str(e)} )
249                 raise
250         else:
251             raise AttributeError( Model does not support predict_proba )

```

```

252
253     def get_model_info(self):
254         Get model information
255         info = {
256             'model_type': type(self.model).__name__,
257             'load_time': self.load_time.isoformat(),
258             'model_path': self.model_path
259         }
260
261         # Add model-specific information
262         if hasattr(self.model, 'get_params'):
263             info['parameters'] = self.model.get_params()
264
265         if hasattr(self.model, 'feature_importances_'):
266             info['feature_importance_available'] = True
267
268         return info
269
270     def health_check(self, X_sample, y_sample=None):
271         Perform health check on loaded model
272         try:
273             # Basic prediction test
274             predictions = self.predict(X_sample[:5]) # Test with first 5
samples
275
276             health_status = {
277                 'status': 'healthy',
278                 'timestamp': datetime.now().isoformat(),
279                 'prediction_shape': predictions.shape,
280                 'prediction_types': [type(pred).__name__ for pred in
predictions[:3]]
281             }
282
283             # If sample labels provided, check accuracy
284             if y_sample is not None:
285                 sample_accuracy = np.mean(predictions[:len(y_sample[:5])])
== y_sample[:5])
286                 health_status['sample_accuracy'] = sample_accuracy
287
288             print( Health check passed )
289             return health_status
290
291         except Exception as e:
292             health_status = {
293                 'status': 'unhealthy',
294                 'timestamp': datetime.now().isoformat(),
295                 'error': str(e)
296             }
297             print(f Health check failed: {str(e)} )
298             return health_status
299
300 # Example usage
301 def example_model_tracking():
302     Example of comprehensive model tracking

```

```
303
304     # Initialize trackers
305     version_tracker = ModelVersionTracker()
306     mlflow = SimpleMLFlow( customer_churn_prediction )
307
308     # Start experiment
309     mlflow.start_run( random_forest_baseline )
310
311     # Log parameters
312     params = {
313         'n_estimators': 100,
314         'max_depth': 10,
315         'random_state': 42
316     }
317
318     for key, value in params.items():
319         mlflow.log_param(key, value)
320
321     # Train model
322     model = RandomForestClassifier(**params)
323     model.fit(X_train, y_train)
324
325     # Evaluate model
326     y_pred = model.predict(X_test)
327     metrics = {
328         'accuracy': accuracy_score(y_test, y_pred),
329         'precision': precision_score(y_test, y_pred),
330         'recall': recall_score(y_test, y_pred),
331         'f1': f1_score(y_test, y_pred)
332     }
333
334     for key, value in metrics.items():
335         mlflow.log_metric(key, value)
336
337     # Save model
338     model_path = version_tracker.save_model(model, mlflow.current_run['
run_id'], 'churn_model')
339     mlflow.log_artifact(model_path)
340
341     # Log experiment
342     experiment_id = version_tracker.log_experiment(
343         model, X_train, y_train, params, metrics,
344         'churn_model', 'Baseline Random Forest model'
345     )
346
347     # End run
348     mlflow.end_run()
349
350     # Compare experiments
351     recent_experiments = version_tracker.get_experiment_history('
churn_model')[:3]
352     experiment_ids = [exp['experiment_id'] for exp in recent_experiments]
353     comparison = version_tracker.compare_experiments(experiment_ids)
354
```

```

355     print( Recent experiments comparison: )
356     for exp_id, exp_data in comparison.items():
357         print(f    {exp_id}: Accuracy = {exp_data['metrics']['accuracy']:.4
f} )
358
359 # Model monitoring in production
360 class ModelMonitor:
361     Monitor model performance in production
362
363     def __init__(self, model_deployer, alert_threshold=0.1):
364         self.model_deployer = model_deployer
365         self.alert_threshold = alert_threshold
366         self.performance_history = []
367         self.prediction_history = []
368
369     def log_prediction(self, X, y_true=None, prediction_time=None):
370         Log prediction for monitoring
371         prediction_record = {
372             'timestamp': prediction_time or datetime.now().isoformat(),
373             'n_samples': len(X),
374             'prediction_time': datetime.now().isoformat()
375         }
376
377         # Store prediction data (be careful with data privacy)
378         self.prediction_history.append(prediction_record)
379
380         # If true labels provided, evaluate performance
381         if y_true is not None:
382             try:
383                 y_pred = self.model_deployer.predict(X)
384                 accuracy = accuracy_score(y_true, y_pred)
385
386                 performance_record = {
387                     'timestamp': prediction_time or datetime.now().
isoformat(),
388                     'accuracy': accuracy,
389                     'n_samples': len(y_true)
390                 }
391
392                 self.performance_history.append(performance_record)
393
394                 # Check for performance degradation
395                 if len(self.performance_history) > 5:
396                     recent_accuracy = np.mean([r['accuracy'] for r in self
.performance_history[-5:]])
397                     historical_accuracy = np.mean([r['accuracy'] for r in
self.performance_history[:-5]])
398
399                     if historical_accuracy - recent_accuracy > self.
alert_threshold:
400                         self.send_alert(
401                             f Performance degradation detected!
402                             f Accuracy dropped from {historical_accuracy
:.4f} to {recent_accuracy:.4f}

```

```

403         )
404
405         return accuracy
406     except Exception as e:
407         print(f Performance monitoring error: {str(e)} )
408
409     def send_alert(self, message):
410         Send alert (simplified - in practice, integrate with alerting
411         systems)
412         print(f ALERT: {message} )
413         # In practice: send email, Slack message, PagerDuty alert, etc.
414
415     def get_performance_report(self, window_hours=24):
416         Get performance report for specified time window
417         cutoff_time = datetime.now() - pd.Timedelta(hours=window_hours)
418         cutoff_str = cutoff_time.isoformat()
419
420         recent_performance = [
421             record for record in self.performance_history
422             if record['timestamp'] >= cutoff_str
423         ]
424
425         if recent_performance:
426             avg_accuracy = np.mean([r['accuracy'] for r in
427             recent_performance])
428             total_predictions = sum([r['n_samples'] for r in
429             recent_performance])
430
431             report = {
432                 'period': f Last {window_hours} hours ,
433                 'average_accuracy': avg_accuracy,
434                 'total_predictions': total_predictions,
435                 'n_evaluation_batches': len(recent_performance)
436             }
437
438             return report
439         else:
440             return {
441                 'period': f Last {window_hours} hours ,
442                 'message': 'No performance data available'
443             }
444
445     # Example usage of production deployment
446     def production_deployment_example():
447         Example of production model deployment
448
449         # Load model
450         deployer = ModelDeployer('models/churn_model_latest.pkl')
451
452         # Health check
453         health_status = deployer.health_check(X_test[:5], y_test[:5])
454         print( Model health: , health_status)
455
456         # Initialize monitoring

```

```

454     monitor = ModelMonitor(deployer, alert_threshold=0.05)
455
456     # Simulate production predictions
457     for i in range(10):
458         # Get batch of data
459         batch_size = np.random.randint(10, 100)
460         X_batch = X_test[i*batch_size:(i+1)*batch_size]
461
462         # Make predictions
463         predictions = deployer.predict(X_batch)
464
465         # Log predictions for monitoring
466         monitor.log_prediction(X_batch)
467
468         print(f Batch {i+1}: Processed {len(predictions)} predictions )
469
470     # Get performance report
471     report = monitor.get_performance_report(window_hours=1)
472     print( Performance report: , report)
473
474 # Run examples
475 # example_model_tracking()
476 # production_deployment_example()

```

Listing 26: Model versioning and experiment tracking

#### Best Practice 4.4.2: Implement Proper Error Handling and Logging

**Why:** Ensures system reliability and enables debugging in production environments.

**How:** Use comprehensive logging, error handling, and monitoring throughout the ML pipeline.

```

1  import logging
2  import traceback
3  from functools import wraps
4  import time
5  from datetime import datetime
6  import json
7  import os
8
9  # Configure logging
10 def setup_logging(log_level=logging.INFO, log_file='ml_pipeline.log'):
11     Setup comprehensive logging for ML pipeline
12
13     # Create logger
14     logger = logging.getLogger('ml_pipeline')
15     logger.setLevel(log_level)
16
17     # Create formatters
18     detailed_formatter = logging.Formatter(
19         '%(asctime)s - %(name)s - %(levelname)s - %(funcName)s:%(lineno)d
20         - %(message)s'

```



```

21     simple_formatter = logging.Formatter(
22         '%(asctime)s - %(levelname)s - %(message)s'
23     )
24
25     # Console handler
26     console_handler = logging.StreamHandler()
27     console_handler.setLevel(logging.INFO)
28     console_handler.setFormatter(simple_formatter)
29
30     # File handler
31     file_handler = logging.FileHandler(log_file)
32     file_handler.setLevel(log_level)
33     file_handler.setFormatter(detailed_formatter)
34
35     # Add handlers to logger
36     logger.addHandler(console_handler)
37     logger.addHandler(file_handler)
38
39     return logger
40
41 # Global logger
42 logger = setup_logging()
43
44 # Decorator for function timing and error handling
45 def monitor_function(func):
46     Decorator to monitor function execution time and handle errors
47
48     @wraps(func)
49     def wrapper(*args, **kwargs):
50         start_time = time.time()
51         function_name = func.__name__
52
53         logger.info(f Starting execution of {function_name} )
54
55         try:
56             result = func(*args, **kwargs)
57             execution_time = time.time() - start_time
58
59             logger.info(f Successfully completed {function_name} in {
execution_time:.2f} seconds )
60             return result
61
62         except Exception as e:
63             execution_time = time.time() - start_time
64             error_msg = f Error in {function_name} after {execution_time
:.2f} seconds: {str(e)}
65             logger.error(error_msg)
66             logger.error(f Traceback: {traceback.format_exc()} )
67
68             # Re-raise the exception
69             raise
70
71     return wrapper
72

```

```

73 # Custom exception classes
74 class DataValidationError(Exception):
75     Exception raised for data validation errors
76     def __init__(self, message, details=None):
77         self.message = message
78         self.details = details
79         super().__init__(self.message)
80
81 class ModelTrainingError(Exception):
82     Exception raised for model training errors
83     def __init__(self, message, model_params=None):
84         self.message = message
85         self.model_params = model_params
86         super().__init__(self.message)
87
88 class PredictionError(Exception):
89     Exception raised for prediction errors
90     def __init__(self, message, input_shape=None):
91         self.message = message
92         self.input_shape = input_shape
93         super().__init__(self.message)
94
95 # Data validation with comprehensive error handling
96 @monitor_function
97 def validate_input_data(X, y=None, expected_features=None):
98     Validate input data with comprehensive error checking
99
100     errors = []
101
102     # Check if X is None or empty
103     if X is None:
104         errors.append( Input data X is None )
105     elif len(X) == 0:
106         errors.append( Input data X is empty )
107
108     # Check data types
109     if not isinstance(X, (list, np.ndarray, pd.DataFrame)):
110         errors.append(f X must be array-like, got {type(X)} )
111
112     # Convert to numpy array for consistency
113     try:
114         X_array = np.array(X)
115     except Exception as e:
116         errors.append(f Cannot convert X to numpy array: {str(e)} )
117         X_array = None
118
119     # Check dimensions
120     if X_array is not None:
121         if X_array.ndim != 2:
122             errors.append(f X must be 2-dimensional, got {X_array.ndim}
123                             dimensions )
124
125     # Check for expected number of features
126     if expected_features and X_array.shape[1] != expected_features:

```

```

126         errors.append(f Expected {expected_features} features, got {
X_array.shape[1]} )
127
128     # Check for infinite values
129     if np.any(np.isinf(X_array)):
130         inf_count = np.sum(np.isinf(X_array))
131         errors.append(f Found {inf_count} infinite values in X )
132
133     # Check for NaN values
134     if np.any(np.isnan(X_array)):
135         nan_count = np.sum(np.isnan(X_array))
136         errors.append(f Found {nan_count} NaN values in X )
137
138     # Validate target variable if provided
139     if y is not None:
140         if len(y) != len(X):
141             errors.append(f Length mismatch: X has {len(X)} samples, y has
{len(y)} samples )
142
143     # Check for NaN in target
144     y_array = np.array(y)
145     if np.any(np.isnan(y_array)):
146         nan_count = np.sum(np.isnan(y_array))
147         errors.append(f Found {nan_count} NaN values in y )
148
149     # Raise exception if errors found
150     if errors:
151         error_details = {
152             'timestamp': datetime.now().isoformat(),
153             'input_shape': X_array.shape if X_array is not None else None,
154             'error_count': len(errors)
155         }
156         raise DataValidationError(
157             Data validation failed:  +  ; .join(errors),
158             details=error_details
159         )
160
161     logger.info( Data validation passed )
162     return X_array, y_array if y is not None else None
163
164 # Robust model training with error handling
165 @monitor_function
166 def robust_model_training(X_train, y_train, model_class, model_params=None
):
167     Train model with comprehensive error handling
168
169     if model_params is None:
170         model_params = {}
171
172     try:
173         # Validate input data
174         X_train_valid, y_train_valid = validate_input_data(X_train,
y_train)
175

```

```

176     # Create model instance
177     logger.info(f Creating model: {model_class.__name__} )
178     model = model_class(**model_params)
179
180     # Train model
181     logger.info( Starting model training )
182     model.fit(X_train_valid, y_train_valid)
183
184     # Validate trained model
185     if not hasattr(model, 'predict'):
186         raise ModelTrainingError(
187             Trained model does not have predict method ,
188             model_params=model_params
189         )
190
191     logger.info( Model training completed successfully )
192     return model
193
194 except DataValidationError as e:
195     logger.error(f Data validation error during training: {e.message}
196 )
197
198     raise
199 except Exception as e:
200     error_details = {
201         'timestamp': datetime.now().isoformat(),
202         'model_class': model_class.__name__,
203         'model_params': model_params,
204         'error_type': type(e).__name__
205     }
206     raise ModelTrainingError(
207         f Model training failed: {str(e)} ,
208         model_params=model_params
209     ) from e
210
211 # Safe prediction with error handling
212 @monitor_function
213 def safe_predict(model, X, return_proba=False):
214     Make predictions with comprehensive error handling
215
216     try:
217         # Validate input data
218         X_valid, _ = validate_input_data(X)
219
220         # Check if model has required methods
221         if not hasattr(model, 'predict'):
222             raise PredictionError(
223                 Model does not have predict method ,
224                 input_shape=X_valid.shape
225             )
226
227         # Make predictions
228         logger.info(f Making predictions for {X_valid.shape[0]} samples )
229         predictions = model.predict(X_valid)

```

```

229     # Return probabilities if requested
230     if return_proba:
231         if not hasattr(model, 'predict_proba'):
232             logger.warning( Model does not support predict_proba,
returning only predictions )
233             return predictions, None
234         else:
235             probabilities = model.predict_proba(X_valid)
236             return predictions, probabilities
237
238     return predictions
239
240     except DataValidationError as e:
241         logger.error(f Data validation error during prediction: {e.message
} )
242         raise
243     except Exception as e:
244         error_details = {
245             'timestamp': datetime.now().isoformat(),
246             'input_shape': X.shape if hasattr(X, 'shape') else None,
247             'model_type': type(model).__name__,
248             'error_type': type(e).__name__
249         }
250         raise PredictionError(
251             f Prediction failed: {str(e)} ,
252             input_shape=X.shape if hasattr(X, 'shape') else None
253         ) from e
254
255 # Model performance monitoring with error handling
256 @monitor_function
257 def monitor_model_performance(y_true, y_pred, metrics=None):
258     Monitor model performance with error handling
259
260     if metrics is None:
261         metrics = ['accuracy', 'precision', 'recall', 'f1']
262
263     results = {}
264     errors = []
265
266     try:
267         # Validate inputs
268         if len(y_true) != len(y_pred):
269             raise ValueError(f Length mismatch: y_true={len(y_true)},
y_pred={len(y_pred)} )
270
271         # Calculate metrics
272         for metric in metrics:
273             try:
274                 if metric == 'accuracy':
275                     results[metric] = accuracy_score(y_true, y_pred)
276                 elif metric == 'precision':
277                     results[metric] = precision_score(y_true, y_pred,
average='weighted')
278                 elif metric == 'recall':

```

```

279         results[metric] = recall_score(y_true, y_pred, average
='weighted')
280     elif metric == 'f1':
281         results[metric] = f1_score(y_true, y_pred, average='
weighted')
282     elif metric == 'roc_auc':
283         if len(np.unique(y_true)) == 2:
284             results[metric] = roc_auc_score(y_true, y_pred)
285         else:
286             results[metric] = 'Not applicable for multiclass'
287     else:
288         logger.warning(f Unknown metric: {metric} )
289     except Exception as e:
290         errors.append(f Error calculating {metric}: {str(e)} )
291         results[metric] = None
292
293     # Log results
294     logger.info( Model performance metrics: )
295     for metric, value in results.items():
296         if value is not None:
297             logger.info(f {metric}: {value:.4f} )
298
299     if errors:
300         logger.warning(f Errors in metric calculation: {errors} )
301
302     return results
303
304     except Exception as e:
305         logger.error(f Error in performance monitoring: {str(e)} )
306         raise
307
308 # Context manager for safe model operations
309 class SafeModelOperation:
310     Context manager for safe model operations with automatic cleanup
311
312     def __init__(self, operation_name, timeout_seconds=300):
313         self.operation_name = operation_name
314         self.timeout_seconds = timeout_seconds
315         self.start_time = None
316
317     def __enter__(self):
318         self.start_time = time.time()
319         logger.info(f Starting {self.operation_name} )
320         return self
321
322     def __exit__(self, exc_type, exc_value, traceback):
323         elapsed_time = time.time() - self.start_time
324
325         if exc_type is not None:
326             # An exception occurred
327             logger.error(f Error in {self.operation_name} after {
elapsed_time:.2f} seconds )
328             logger.error(f Exception type: {exc_type.__name__} )
329             logger.error(f Exception message: {str(exc_value)} )

```

```

330         if traceback:
331             logger.error(f Traceback: {traceback.format_exc()} )
332             return False # Don't suppress the exception
333         else:
334             # Operation completed successfully
335             logger.info(f Completed {self.operation_name} in {elapsed_time
:.2f} seconds )
336             return True
337
338 # Retry mechanism for unreliable operations
339 def retry_on_failure(max_retries=3, delay_seconds=1, backoff_factor=2):
340     Decorator to retry function execution on failure
341
342     def decorator(func):
343         @wraps(func)
344         def wrapper(*args, **kwargs):
345             retries = 0
346             current_delay = delay_seconds
347
348             while retries < max_retries:
349                 try:
350                     return func(*args, **kwargs)
351                 except Exception as e:
352                     retries += 1
353                     if retries >= max_retries:
354                         logger.error(f Function {func.__name__} failed
after {max_retries} retries )
355                         raise
356
357                         logger.warning(
358                             f Function {func.__name__} failed (attempt {
retries}/{max_retries}): {str(e)}
359                         )
360                         logger.info(f Retrying in {current_delay} seconds... )
361                         time.sleep(current_delay)
362                         current_delay *= backoff_factor
363
364             return wrapper
365         return decorator
366
367 # Example usage of error handling and logging
368 @retry_on_failure(max_retries=3, delay_seconds=2)
369 @monitor_function
370 def example_ml_pipeline(X_train, y_train, X_test, y_test):
371     Example ML pipeline with comprehensive error handling
372
373     with SafeModelOperation( Data Validation ) as op:
374         # Validate training data
375         X_train_valid, y_train_valid = validate_input_data(X_train,
y_train)
376         X_test_valid, y_test_valid = validate_input_data(X_test, y_test)
377
378     with SafeModelOperation( Model Training ) as op:
379         # Train model

```

```

380     model = robust_model_training(
381         X_train_valid, y_train_valid,
382         RandomForestClassifier,
383         {'n_estimators': 100, 'random_state': 42}
384     )
385
386     with SafeModelOperation( Model Evaluation ) as op:
387         # Make predictions
388         y_pred = safe_predict(model, X_test_valid)
389
390         # Evaluate performance
391         metrics = monitor_model_performance(y_test_valid, y_pred)
392
393     return model, metrics
394
395 # Error handling for data loading
396 @monitor_function
397 def safe_data_loading(file_path, file_type='csv'):
398     Safely load data with comprehensive error handling
399
400     try:
401         if not os.path.exists(file_path):
402             raise FileNotFoundError(f File not found: {file_path} )
403
404         if file_type == 'csv':
405             data = pd.read_csv(file_path)
406         elif file_type == 'json':
407             data = pd.read_json(file_path)
408         elif file_type == 'excel':
409             data = pd.read_excel(file_path)
410         else:
411             raise ValueError(f Unsupported file type: {file_type} )
412
413         logger.info(f Successfully loaded data from {file_path} )
414         logger.info(f Data shape: {data.shape} )
415
416     return data
417
418     except FileNotFoundError as e:
419         logger.error(f File not found error: {str(e)} )
420         raise
421     except pd.errors.EmptyDataError:
422         logger.error( File is empty )
423         raise DataValidationError( Data file is empty )
424     except pd.errors.ParserError as e:
425         logger.error(f Data parsing error: {str(e)} )
426         raise DataValidationError(f Data parsing failed: {str(e)} )
427     except Exception as e:
428         logger.error(f Unexpected error loading data: {str(e)} )
429         raise DataValidationError(f Failed to load data: {str(e)} )
430
431 # Graceful degradation for model services
432 class GracefulModelService:
433     Model service with graceful degradation

```



```

434
435     def __init__(self, primary_model, fallback_model=None):
436         self.primary_model = primary_model
437         self.fallback_model = fallback_model
438         self.primary_model_failed = False
439
440     def predict(self, X, use_fallback_if_needed=True):
441         Make predictions with fallback capability
442
443         # Try primary model first
444         if not self.primary_model_failed:
445             try:
446                 logger.info( Using primary model for prediction )
447                 return safe_predict(self.primary_model, X)
448             except Exception as e:
449                 logger.warning(f Primary model failed: {str(e)} )
450                 self.primary_model_failed = True
451
452                 # If fallback available and allowed, use it
453                 if self.fallback_model and use_fallback_if_needed:
454                     logger.info( Switching to fallback model )
455                     return safe_predict(self.fallback_model, X)
456                 else:
457                     raise
458         else:
459             # Primary model already failed, use fallback if available
460             if self.fallback_model and use_fallback_if_needed:
461                 logger.info( Using fallback model for prediction )
462                 return safe_predict(self.fallback_model, X)
463             else:
464                 raise PredictionError( Primary model failed and no
465                 fallback available )
466
467 # Example usage
468 def example_error_handling():
469     Example demonstrating comprehensive error handling
470
471     try:
472         # Load data safely
473         data = safe_data_loading('data.csv')
474
475         # Prepare data
476         X = data.drop('target', axis=1)
477         y = data['target']
478
479         # Split data
480         X_train, X_test, y_train, y_test = train_test_split(X, y,
481         test_size=0.2, random_state=42)
482
483         # Run pipeline with error handling
484         model, metrics = example_ml_pipeline(X_train, y_train, X_test,
485         y_test)
486
487         logger.info( ML pipeline completed successfully )

```

```

485         return model, metrics
486
487     except DataValidationError as e:
488         logger.error(f Data validation error: {e.message} )
489         if e.details:
490             logger.error(f Details: {e.details} )
491     except ModelTrainingError as e:
492         logger.error(f Model training error: {e.message} )
493         if e.model_params:
494             logger.error(f Model parameters: {e.model_params} )
495     except PredictionError as e:
496         logger.error(f Prediction error: {e.message} )
497         if e.input_shape:
498             logger.error(f Input shape: {e.input_shape} )
499     except Exception as e:
500         logger.error(f Unexpected error: {str(e)} )
501         logger.error(f Traceback: {traceback.format_exc()} )
502
503 # Run example
504 # example_error_handling()

```

Listing 27: Production-ready error handling and logging

## 5 5. Learning from Errors

### 5.1 5.1 Common Data Preprocessing Errors

```

1 # Buggy: Data leakage during preprocessing
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.model_selection import train_test_split
4
5 # WRONG WAY - Data leakage
6 scaler = StandardScaler()
7 X_scaled = scaler.fit_transform(X) # Using entire dataset!
8 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size
    =0.2)
9
10 # This causes data leakage because test data influences the scaling
    parameters
11
12 # CORRECT WAY - No data leakage
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
    random_state=42)
14
15 # Fit preprocessing on training data only
16 scaler = StandardScaler()
17 X_train_scaled = scaler.fit_transform(X_train) # Fit only on training
    data
18 X_test_scaled = scaler.transform(X_test) # Transform test data
    using training parameters
19
20 print( Correct preprocessing completed without data leakage )

```

Listing 28: Error 5.1.1: Data Leakage During Preprocessing - Buggy Code

```

1 # Buggy: Incorrect categorical variable handling
2 import pandas as pd
3 from sklearn.preprocessing import LabelEncoder
4
5 # WRONG WAY - Problems with LabelEncoder for features
6 data = pd.DataFrame({
7     'category': ['Low', 'Medium', 'High', 'Low', 'Medium'],
8     'target': [0, 1, 1, 0, 1]
9 })
10
11 # This creates artificial ordering: Low=0, Medium=1, High=2
12 le = LabelEncoder()
13 data['category_encoded'] = le.fit_transform(data['category'])
14
15 print( Buggy encoding result: )
16 print(data)
17
18 # PROBLEMS:
19 # 1. Artificial ordinal relationship (Low < Medium < High)
20 # 2. Model might interpret this as numerical relationship
21 # 3. Can lead to incorrect predictions
22
23 # CORRECT WAY - Use OneHotEncoder for features
24 from sklearn.preprocessing import OneHotEncoder
25
26 ohe = OneHotEncoder(sparse=False, drop='first') # drop='first' to avoid
        multicollinearity
27 category_encoded = ohe.fit_transform(data[['category']])
28
29 print( \nCorrect encoding result: )
30 print(category_encoded)
31 print( Categories: , ohe.categories_)
32
33 # For target variables, LabelEncoder is appropriate
34 target_le = LabelEncoder()
35 data['target_encoded'] = target_le.fit_transform(data['target'])
36 print( \nTarget encoding: , data['target_encoded'])

```

Listing 29: Error 5.1.2: Incorrect Handling of Categorical Variables - Buggy Code

```

1 # Buggy: Improper missing value handling
2 import pandas as pd
3 import numpy as np
4 from sklearn.impute import SimpleImputer
5
6 # WRONG WAY - Dropping important data
7 data = pd.DataFrame({
8     'feature1': [1, 2, np.nan, 4, 5],
9     'feature2': [10, np.nan, 30, 40, 50],
10    'target': [0, 1, 1, 0, 1]
11 })
12
13 print( Original data: )
14 print(data)
15 print(f Missing values: {data.isnull().sum().sum()} )

```

```

16
17 # Buggy approach - Dropping rows with missing values
18 data_dropped = data.dropna()
19 print(f \nAfter dropping rows: {len(data_dropped)} rows remaining )
20 print( This loses 40% of data! )
21
22 # Another buggy approach - Filling with zeros
23 data_zero_filled = data.fillna(0)
24 print( \nZero-filled data: )
25 print(data_zero_filled)
26 # Problem: Zeros might not be meaningful and can skew results
27
28 # CORRECT WAY - Strategic imputation
29 # For numerical features, use median or mean
30 num_imputer = SimpleImputer(strategy='median')
31 data['feature1'] = num_imputer.fit_transform(data[['feature1']])
32
33 # For features with specific patterns, use domain knowledge
34 # Or use advanced imputation methods
35 from sklearn.impute import KNNImputer
36 knn_imputer = KNNImputer(n_neighbors=2)
37 data_imputed = pd.DataFrame(
38     knn_imputer.fit_transform(data[['feature1', 'feature2']]),
39     columns=['feature1', 'feature2']
40 )
41 data_imputed['target'] = data['target']
42
43 print( \nCorrectly imputed data: )
44 print(data_imputed)

```

Listing 30: Error 5.1.3: Improper Missing Value Handling - Buggy Code

```

1 # Buggy: Incorrect feature scaling on sparse data
2 import numpy as np
3 from scipy.sparse import csr_matrix
4 from sklearn.preprocessing import StandardScaler
5
6 # Create sparse data
7 sparse_data = np.array([
8     [1, 0, 0, 0, 5],
9     [0, 2, 0, 0, 0],
10    [0, 0, 3, 0, 0],
11    [0, 0, 0, 4, 0]
12 ])
13
14 sparse_matrix = csr_matrix(sparse_data)
15 print( Original sparse matrix: )
16 print(sparse_matrix.toarray())
17 print(f Sparsity: {1 - sparse_matrix.nnz / (sparse_matrix.shape[0] *
18     sparse_matrix.shape[1]):.2%} )
19
20 # WRONG WAY - Using StandardScaler on sparse data
21 scaler = StandardScaler()
22 try:
23     scaled_data_wrong = scaler.fit_transform(sparse_matrix)

```

```

23     print( \nAfter StandardScaler (WRONG): )
24     print(scaled_data_wrong)
25     print( Result is dense matrix - lost sparsity! )
26 except Exception as e:
27     print(f Error: {e} )
28
29 # CORRECT WAY - Use appropriate scaling for sparse data
30 from sklearn.preprocessing import MaxAbsScaler
31
32 # MaxAbsScaler preserves sparsity
33 maxabs_scaler = MaxAbsScaler()
34 scaled_data_correct = maxabs_scaler.fit_transform(sparse_matrix)
35
36 print( \nAfter MaxAbsScaler (CORRECT): )
37 print(scaled_data_correct.toarray())
38 print(f Sparsity preserved: {1 - scaled_data_correct.nnz / (
    scaled_data_correct.shape[0] * scaled_data_correct.shape[1]):.2%} )
39
40 # Alternative: Use MinMaxScaler with feature_range=(-1, 1)
41 from sklearn.preprocessing import MinMaxScaler
42
43 minmax_scaler = MinMaxScaler(feature_range=(-1, 1))
44 scaled_data_minmax = minmax_scaler.fit_transform(sparse_matrix)
45
46 print( \nAfter MinMaxScaler (Alternative): )
47 print(scaled_data_minmax.toarray())

```

Listing 31: Error 5.1.4: Feature Scaling on Sparse Data - Buggy Code

```

1 # Buggy: Incorrect train-test split approaches
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5
6 # Create time series data
7 dates = pd.date_range('2020-01-01', periods=1000, freq='D')
8 data = pd.DataFrame({
9     'date': dates,
10    'value': np.random.randn(1000).cumsum() + 100,
11    'feature1': np.random.randn(1000),
12    'feature2': np.random.randn(1000)
13 })
14
15 print( Data shape: , data.shape)
16 print( Date range: , data['date'].min(), to , data['date'].max())
17
18 # WRONG WAY 1 - Random split on time series data
19 X = data[['feature1', 'feature2']]
20 y = data['value']
21
22 X_train_wrong, X_test_wrong, y_train_wrong, y_test_wrong =
    train_test_split(
23     X, y, test_size=0.2, random_state=42
24 )
25

```

```

26 print( \nWRONG WAY 1 - Random split: )
27 print(f Train date range: {data.iloc[X_train_wrong.index]['date'].min()}
    to {data.iloc[X_train_wrong.index]['date'].max()} )
28 print(f Test date range: {data.iloc[X_test_wrong.index]['date'].min()} to
    {data.iloc[X_test_wrong.index]['date'].max()} )
29 print( Problem: Future data leaks into training! )
30
31 # WRONG WAY 2 - Incorrect stratification
32 categorical_data = pd.DataFrame({
33     'feature': np.random.randn(1000),
34     'category': np.random.choice(['A', 'B', 'C'], 1000, p=[0.95, 0.03,
    0.02])
35 })
36
37 X_cat = categorical_data[['feature']]
38 y_cat = categorical_data['category']
39
40 # Not using stratification for imbalanced data
41 X_train_no_strat, X_test_no_strat, y_train_no_strat, y_test_no_strat =
    train_test_split(
42     X_cat, y_cat, test_size=0.2, random_state=42
43 )
44
45 print(f \nWRONG WAY 2 - No stratification: )
46 print(f Original distribution: A={sum(y_cat=='A')}, B={sum(y_cat=='B')}, C
    ={sum(y_cat=='C')} )
47 print(f Train distribution: A={sum(y_train_no_strat=='A')}, B={sum(
    y_train_no_strat=='B')}, C={sum(y_train_no_strat=='C')} )
48 print( Problem: Rare classes might be underrepresented in train/test sets!
    )
49
50 # CORRECT WAY 1 - Time series split
51 # Use temporal order for time series
52 split_date = data['date'].quantile(0.8) # 80th percentile
53 train_mask = data['date'] <= split_date
54 test_mask = data['date'] > split_date
55
56 X_train_correct = X[train_mask]
57 X_test_correct = X[test_mask]
58 y_train_correct = y[train_mask]
59 y_test_correct = y[test_mask]
60
61 print( \nCORRECT WAY 1 - Temporal split: )
62 print(f Train date range: {data[train_mask]['date'].min()} to {data[
    train_mask]['date'].max()} )
63 print(f Test date range: {data[test_mask]['date'].min()} to {data[
    test_mask]['date'].max()} )
64
65 # CORRECT WAY 2 - Stratified split
66 X_train_strat, X_test_strat, y_train_strat, y_test_strat =
    train_test_split(
67     X_cat, y_cat, test_size=0.2, random_state=42, stratify=y_cat
68 )
69

```

```

70 print(f \nCORRECT WAY 2 - Stratified split: )
71 print(f Train distribution: A={sum(y_train_strat=='A')}, B={sum(
    y_train_strat=='B')}, C={sum(y_train_strat=='C')} )
72 print(f Test distribution: A={sum(y_test_strat=='A')}, B={sum(y_test_strat
    =='B')}, C={sum(y_test_strat=='C')} )

```

Listing 32: Error 5.1.5: Incorrect Train-Test Split - Buggy Code

```

1 # Buggy: Data type conversion problems
2 import pandas as pd
3 import numpy as np
4 from sklearn.preprocessing import LabelEncoder
5
6 # Create problematic data
7 data = pd.DataFrame({
8     'numeric_string': ['1', '2', '3', '4', '5'],
9     'mixed_numeric': [1, 2.5, 3, 4.7, 5],
10    'categorical_numeric': [1, 2, 3, 1, 2], # Actually categories
11    'boolean_like': ['True', 'False', 'True', 'False', 'True'],
12    'date_string': ['2020-01-01', '2020-01-02', '2020-01-03', '2020-01-04'
13    , '2020-01-05']
14 })
15 print( Original data types: )
16 print(data.dtypes)
17 print( \nOriginal data: )
18 print(data)
19
20 # WRONG WAY 1 - Incorrect numeric conversion
21 try:
22     # This might work but can cause issues
23     data['numeric_string_converted'] = pd.to_numeric(data['numeric_string'
24     ])
25     print( \nNumeric string conversion: )
26     print(data['numeric_string_converted'].dtype)
27 except Exception as e:
28     print(f Error in numeric conversion: {e} )
29
30 # WRONG WAY 2 - Treating categorical as numeric
31 # Problem: categorical_numeric should be treated as categories, not
32 # numbers
33 le = LabelEncoder()
34 data['categorical_numeric_encoded'] = le.fit_transform(data['
35 categorical_numeric'])
36 print(f \nCategorical numeric encoding: {data['categorical_numeric_encoded'
37     ''].tolist()} )
38 print( Problem: Model might treat these as ordinal relationships! )
39
40 # WRONG WAY 3 - Boolean conversion issues
41 try:
42     data['boolean_converted'] = data['boolean_like'].astype(bool)
43     print(f \nBoolean conversion: {data['boolean_converted'].tolist()} )
44     print( Problem: All strings convert to True! )
45 except Exception as e:
46     print(f Error in boolean conversion: {e} )

```

```

43
44 # WRONG WAY 4 - Date parsing without proper format
45 try:
46     data['date_parsed_wrong'] = pd.to_datetime(data['date_string'])
47     print(f \nDate parsing (might work): {data['date_parsed_wrong'].dtype}
48         )
49 except Exception as e:
50     print(f Error in date parsing: {e} )
51
52 # CORRECT WAY - Proper data type handling
53 # 1. Explicit numeric conversion with error handling
54 data['numeric_string_correct'] = pd.to_numeric(data['numeric_string'],
55     errors='coerce')
56
57 # 2. Proper categorical handling
58 data['categorical_numeric_correct'] = data['categorical_numeric'].astype('
59     category')
60
61 # 3. Boolean conversion with proper mapping
62 bool_mapping = {'True': True, 'False': False}
63 data['boolean_correct'] = data['boolean_like'].map(bool_mapping)
64
65 # 4. Explicit date parsing with format
66 data['date_correct'] = pd.to_datetime(data['date_string'], format='%Y-%m-%
67     d')
68
69 print( \nCorrect data types: )
70 print(data.dtypes)
71
72 print( \nCorrectly processed data: )
73 print(data[['numeric_string_correct', 'categorical_numeric_correct',
74     'boolean_correct', 'date_correct']].head())
75
76 # Additional safety checks
77 def safe_convert_types(df):
78     Safely convert data types with proper error handling
79     df_converted = df.copy()
80
81     for column in df.columns:
82         # Try to infer and convert appropriate types
83         if df[column].dtype == 'object':
84             # Check if it's numeric
85             numeric_converted = pd.to_numeric(df[column], errors='coerce')
86             if not numeric_converted.isna().all():
87                 # Check if conversion lost too much data
88                 if numeric_converted.isna().sum() / len(df) < 0.1: # Less
89                     than 10% loss
90                     df_converted[column] = numeric_converted
91                     continue
92
93             # Check if it's boolean-like
94             unique_values = df[column].unique()
95             if set(unique_values).issubset({True, False, 'True', 'False',
96                 'true', 'false', 1, 0}):

```



```

91         bool_mapping = {True: True, False: False, 'True': True, '
False': False,
92                        'true': True, 'false': False, 1: True, 0:
False}
93         df_converted[column] = df[column].map(bool_mapping)
94         continue
95
96     # Check if it's date-like
97     try:
98         pd.to_datetime(df[column].iloc[:5]) # Test first 5 values
99         df_converted[column] = pd.to_datetime(df[column])
100        continue
101    except:
102        pass
103
104    # Otherwise, keep as categorical if few unique values
105    if df[column].nunique() / len(df) < 0.05: # Less than 5%
unique values
106        df_converted[column] = df[column].astype('category')
107
108    return df_converted
109
110 # Apply safe conversion
111 data_safe = safe_convert_types(data)
112 print( \nSafely converted data types: )
113 print(data_safe.dtypes)

```

Listing 33: Error 5.1.6: Data Type Conversion Issues - Buggy Code

```

1 # Buggy: Common feature engineering mistakes
2 import pandas as pd
3 import numpy as np
4 from sklearn.preprocessing import StandardScaler
5
6 # Create sample data
7 np.random.seed(42)
8 data = pd.DataFrame({
9     'feature1': np.random.randn(1000),
10    'feature2': np.random.randn(1000),
11    'category': np.random.choice(['A', 'B', 'C'], 1000),
12    'date': pd.date_range('2020-01-01', periods=1000, freq='D')
13 })
14
15 print( Original data shape: , data.shape)
16
17 # WRONG WAY 1 - Creating data leakage features
18 # Buggy: Using future information to create features
19 def create_leaky_features_buggy(df):
20     Buggy feature engineering that creates data leakage
21     df_new = df.copy()
22
23     # WRONG: Using mean of entire dataset (data leakage!)
24     df_new['feature1_normalized'] = df_new['feature1'] / df_new['feature1'
].mean()
25

```

```

26     # WRONG: Using global statistics instead of training set statistics
27     df_new['feature2_scaled'] = (df_new['feature2'] - df_new['feature2'].
    min()) / (df_new['feature2'].max() - df_new['feature2'].min())
28
29     # WRONG: Creating features that wouldn't be available at prediction
    time
30     df_new['future_trend'] = df_new['feature1'].rolling(window=10).mean().
    shift(-5) # Future data!
31
32     return df_new
33
34 # Apply buggy feature engineering
35 data_leaky = create_leaky_features_buggy(data)
36 print( Buggy feature engineering completed )
37 print( Problems: Data leakage, future data usage )
38
39 # WRONG WAY 2 - Overfitting through excessive feature creation
40 def create_overfitting_features_buggy(df):
41     Buggy: Creating too many features leading to overfitting
42     df_new = df.copy()
43
44     # WRONG: Creating polynomial features for all combinations
45     for i in range(10): # Too many features!
46         for j in range(10):
47             df_new[f'feature1_poly_{i}_feature2_poly_{j}'] = (df_new['
    feature1'] ** i) * (df_new['feature2'] ** j)
48
49     # WRONG: Creating features for every possible combination
50     categories = df_new['category'].unique()
51     for cat1 in categories:
52         for cat2 in categories:
53             df_new[f'category_interaction_{cat1}_{cat2}'] = ((df_new['
    category'] == cat1) & (df_new['category'] == cat2)).astype(int)
54
55     return df_new
56
57 # WRONG WAY 3 - Inconsistent feature engineering
58 def inconsistent_feature_engineering(df_train, df_test):
59     Buggy: Inconsistent feature engineering between train and test
60     # WRONG: Different transformations on train vs test
61     df_train_new = df_train.copy()
62     df_test_new = df_test.copy()
63
64     # Train gets log transformation
65     df_train_new['feature1_transformed'] = np.log(df_train['feature1'] +
    1)
66
67     # Test gets square root transformation (INCONSISTENT!)
68     df_test_new['feature1_transformed'] = np.sqrt(df_test['feature1'])
69
70     return df_train_new, df_test_new
71
72 # CORRECT WAY - Proper feature engineering practices
73 def create_proper_features(df_train, df_test=None):

```

```

74     Correct feature engineering with proper practices
75
76     # 1. Feature engineering on training data only
77     df_train_new = df_train.copy()
78
79     # Safe feature creation using training statistics only
80     feature1_mean = df_train['feature1'].mean()
81     feature1_std = df_train['feature1'].std()
82     feature2_min = df_train['feature2'].min()
83     feature2_max = df_train['feature2'].max()
84
85     df_train_new['feature1_normalized'] = (df_train['feature1'] -
86     feature1_mean) / feature1_std
87     df_train_new['feature2_scaled'] = (df_train['feature2'] - feature2_min
88     ) / (feature2_max - feature2_min)
89
90     # 2. Safe date features (no future data)
91     df_train_new['year'] = df_train['date'].dt.year
92     df_train_new['month'] = df_train['date'].dt.month
93     df_train_new['day_of_week'] = df_train['date'].dt.dayofweek
94     df_train_new['is_weekend'] = df_train['date'].dt.dayofweek.isin([5,
95     6]).astype(int)
96
97     # 3. Safe categorical encoding
98     df_train_new = pd.get_dummies(df_train_new, columns=['category'],
99     prefix='cat')
100
101     # 4. Lag features (past information only)
102     df_train_new['feature1_lag_1'] = df_train['feature1'].shift(1)
103     df_train_new['feature1_rolling_mean_3'] = df_train['feature1'].rolling
104     (window=3).mean()
105
106     # 5. If test data provided, apply same transformations
107     if df_test is not None:
108         df_test_new = df_test.copy()
109
110         # Apply same transformations using training statistics
111         df_test_new['feature1_normalized'] = (df_test['feature1'] -
112         feature1_mean) / feature1_std
113         df_test_new['feature2_scaled'] = (df_test['feature2'] -
114         feature2_min) / (feature2_max - feature2_min)
115
116         # Date features
117         df_test_new['year'] = df_test['date'].dt.year
118         df_test_new['month'] = df_test['date'].dt.month
119         df_test_new['day_of_week'] = df_test['date'].dt.dayofweek
120         df_test_new['is_weekend'] = df_test['date'].dt.dayofweek.isin([5,
121         6]).astype(int)
122
123         # Categorical encoding (ensure same columns)
124         df_test_new = pd.get_dummies(df_test_new, columns=['category'],
125         prefix='cat')
126
127         # Ensure same columns in both datasets

```

```

119     train_cols = set(df_train_new.columns) - {'feature1', 'feature2',
120     'date', 'category'}
121     test_cols = set(df_test_new.columns) - {'feature1', 'feature2', '
122     date', 'category'}
123
124     # Add missing columns with zeros
125     for col in train_cols - test_cols:
126         df_test_new[col] = 0
127     for col in test_cols - train_cols:
128         df_train_new[col] = 0
129
130     # Reorder columns to match
131     common_cols = sorted(list(train_cols & test_cols))
132     df_train_new = df_train_new[common_cols]
133     df_test_new = df_test_new[common_cols]
134
135     return df_train_new, df_test_new
136
137 return df_train_new
138
139 # Demonstrate correct approach
140 train_data = data.iloc[:800] # First 800 rows for training
141 test_data = data.iloc[800:] # Last 200 rows for testing
142
143 train_features, test_features = create_proper_features(train_data,
144     test_data)
145
146 print(f \nCorrect feature engineering: )
147 print(f Training features shape: {train_features.shape} )
148 print(f Testing features shape: {test_features.shape} )
149 print(f Feature columns: {list(train_features.columns)} )
150
151 # Additional best practices for feature engineering
152 class FeatureEngineer:
153     Robust feature engineering class
154
155     def __init__(self):
156         self.feature_stats = {}
157         self.categorical_mappings = {}
158
159     def fit_transform(self, df, target_col=None):
160         Fit feature engineering parameters and transform training data
161
162         df_transformed = df.copy()
163
164         # Store statistics for numerical features
165         numerical_cols = df.select_dtypes(include=[np.number]).columns
166         if target_col in numerical_cols:
167             numerical_cols = numerical_cols.drop(target_col)
168
169         for col in numerical_cols:
170             self.feature_stats[col] = {
171                 'mean': df[col].mean(),
172                 'std': df[col].std(),

```

```

169         'min': df[col].min(),
170         'max': df[col].max()
171     }
172
173     # Apply standardization
174     mean_val = self.feature_stats[col]['mean']
175     std_val = self.feature_stats[col]['std']
176     if std_val > 0:
177         df_transformed[f'{col}_standardized'] = (df[col] -
mean_val) / std_val
178
179     # Handle categorical features
180     categorical_cols = df.select_dtypes(include=['object', 'category'
]).columns
181     for col in categorical_cols:
182         # Store value counts for reference
183         self.categorical_mappings[col] = df[col].value_counts().
to_dict()
184
185     # One-hot encoding
186     dummies = pd.get_dummies(df[col], prefix=col)
187     df_transformed = pd.concat([df_transformed, dummies], axis=1)
188
189     return df_transformed
190
191     def transform(self, df):
192         Transform new data using fitted parameters
193         df_transformed = df.copy()
194
195         # Apply same transformations using stored parameters
196         for col, stats in self.feature_stats.items():
197             if col in df.columns:
198                 mean_val = stats['mean']
199                 std_val = stats['std']
200                 if std_val > 0:
201                     df_transformed[f'{col}_standardized'] = (df[col] -
mean_val) / std_val
202
203         # Handle categorical features consistently
204         for col, mappings in self.categorical_mappings.items():
205             if col in df.columns:
206                 dummies = pd.get_dummies(df[col], prefix=col)
207                 # Ensure same columns as training
208                 for dummy_col in [f'{col}_{val}' for val in mappings.keys
()]:
209                     if dummy_col not in dummies.columns:
210                         dummies[dummy_col] = 0
211                 df_transformed = pd.concat([df_transformed, dummies], axis
=1)
212
213         return df_transformed
214
215     # Example usage
216     engineer = FeatureEngineer()

```

```

217 train_transformed = engineer.fit_transform(train_data)
218 test_transformed = engineer.transform(test_data)
219
220 print(f \nUsing FeatureEngineer class: )
221 print(f Training transformed shape: {train_transformed.shape} )
222 print(f Testing transformed shape: {test_transformed.shape} )

```

Listing 34: Error 5.1.7: Feature Engineering Mistakes - Buggy Code

```

1 # Buggy: Ignoring data quality issues
2 import pandas as pd
3 import numpy as np
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6
7 # Create problematic data with quality issues
8 np.random.seed(42)
9 data = pd.DataFrame({
10     'feature1': np.random.randn(1000),
11     'feature2': np.random.randn(1000),
12     'feature3': np.random.choice([1, 2, 3, 999], 1000, p=[0.3, 0.3, 0.3,
13         0.1]), # 999 is outlier
14     'category': np.random.choice(['A', 'B', 'C', ''], 1000, p=[0.4, 0.4,
15         0.1, 0.1]), # Empty strings
16     'target': np.random.choice([0, 1], 1000, p=[0.7, 0.3])
17 })
18
19 # Add some missing values and duplicates
20 data.loc[100:105, 'feature1'] = np.nan
21 data.loc[200:202] = data.loc[100:102].values # Duplicate rows
22
23 print( Original data shape: , data.shape)
24 print( Missing values: )
25 print(data.isnull().sum())
26 print( \nDuplicate rows: , data.duplicated().sum())
27 print( \nValue counts for category: )
28 print(data['category'].value_counts())
29
30 # WRONG WAY - Ignoring data quality issues
31 def train_model_buggy(X, y):
32     Buggy model training ignoring data quality
33     # No data quality checks!
34     model = RandomForestClassifier(random_state=42)
35     model.fit(X, y)
36     return model
37
38 # Prepare data the buggy way
39 X_buggy = data.drop('target', axis=1)
40 y_buggy = data['target']
41
42 # This will cause issues:
43 # 1. Missing values in X
44 # 2. Categorical variables not encoded
45 # 3. Outliers not handled
46 # 4. Duplicates not removed

```

```

45
46 try:
47     model_buggy = train_model_buggy(X_buggy, y_buggy)
48     print( Buggy model training completed (but with issues!) )
49 except Exception as e:
50     print(f Buggy approach failed: {e} )
51
52 # CORRECT WAY - Comprehensive data quality handling
53 def comprehensive_data_quality_check(df):
54     Comprehensive data quality assessment
55     quality_report = {
56         'shape': df.shape,
57         'missing_values': df.isnull().sum().to_dict(),
58         'duplicate_rows': df.duplicated().sum(),
59         'data_types': df.dtypes.to_dict(),
60         'issues': []
61     }
62
63     # Check for missing values
64     missing_cols = df.isnull().sum()
65     for col, count in missing_cols.items():
66         if count > 0:
67             percentage = (count / len(df)) * 100
68             quality_report['issues'].append(f Column '{col}' has {count}
({percentage:.1f}%) missing values )
69
70     # Check for duplicates
71     if quality_report['duplicate_rows'] > 0:
72         quality_report['issues'].append(f Found {quality_report['
duplicate_rows']} duplicate rows )
73
74     # Check for suspicious values (outliers, invalid categories)
75     for col in df.select_dtypes(include=[np.number]).columns:
76         Q1 = df[col].quantile(0.25)
77         Q3 = df[col].quantile(0.75)
78         IQR = Q3 - Q1
79         outliers = df[(df[col] < (Q1 - 1.5 * IQR)) | (df[col] > (Q3 + 1.5
* IQR))]
80         if len(outliers) > 0:
81             outlier_percentage = (len(outliers) / len(df)) * 100
82             if outlier_percentage > 5: # More than 5% outliers
83                 quality_report['issues'].append(f Column '{col}' has {len(
outliers)} ({outlier_percentage:.1f}%) outliers )
84
85     for col in df.select_dtypes(include=['object']).columns:
86         empty_values = df[col].isin(['', ' ', 'None', 'null']).sum()
87         if empty_values > 0:
88             percentage = (empty_values / len(df)) * 100
89             quality_report['issues'].append(f Column '{col}' has {
empty_values} ({percentage:.1f}%) empty/invalid values )
90
91     return quality_report
92
93 def clean_data_properly(df, target_col):

```

```

94     Properly clean data with quality checks
95     df_clean = df.copy()
96
97     # Remove duplicates
98     duplicates_before = df_clean.duplicated().sum()
99     df_clean = df_clean.drop_duplicates()
100    print(f Removed {duplicates_before - df_clean.duplicated().sum()}
duplicate rows )
101
102    # Handle missing values
103    missing_before = df_clean.isnull().sum().sum()
104
105    # For numerical columns, fill with median
106    numerical_cols = df_clean.select_dtypes(include=[np.number]).columns
107    if target_col in numerical_cols:
108        numerical_cols = numerical_cols.drop(target_col)
109
110    for col in numerical_cols:
111        if df_clean[col].isnull().sum() > 0:
112            median_val = df_clean[col].median()
113            df_clean[col].fillna(median_val, inplace=True)
114
115    # For categorical columns, fill with mode or create 'Unknown' category
116    categorical_cols = df_clean.select_dtypes(include=['object']).columns
117    for col in categorical_cols:
118        if df_clean[col].isnull().sum() > 0:
119            mode_val = df_clean[col].mode()
120            if len(mode_val) > 0:
121                df_clean[col].fillna(mode_val[0], inplace=True)
122            else:
123                df_clean[col].fillna('Unknown', inplace=True)
124
125    # Handle empty strings
126    for col in categorical_cols:
127        df_clean[col] = df_clean[col].replace(['', ' ', 'None', 'null'], '
Unknown')
128
129    # Handle outliers (capping approach)
130    for col in numerical_cols:
131        Q1 = df_clean[col].quantile(0.25)
132        Q3 = df_clean[col].quantile(0.75)
133        IQR = Q3 - Q1
134        lower_bound = Q1 - 1.5 * IQR
135        upper_bound = Q3 + 1.5 * IQR
136
137        # Cap outliers instead of removing them
138        df_clean[col] = df_clean[col].clip(lower=lower_bound, upper=
upper_bound)
139
140    print(f Handled {missing_before - df_clean.isnull().sum().sum()}
missing values )
141
142    return df_clean
143

```



```

144 # Apply proper data quality handling
145 print( === DATA QUALITY ASSESSMENT === )
146 quality_report = comprehensive_data_quality_check(data)
147 for issue in quality_report['issues']:
148     print(f    - {issue} )
149
150 print( \n=== CLEANING DATA === )
151 data_clean = clean_data_properly(data, 'target')
152
153 print(f \nAfter cleaning: )
154 print(f    Shape: {data_clean.shape} )
155 print(f    Missing values: {data_clean.isnull().sum().sum()} )
156 print(f    Duplicates: {data_clean.duplicated().sum()} )
157
158 # Proper model training with clean data
159 def train_model_properly(df, target_col):
160     Proper model training with data quality checks
161
162     # Separate features and target
163     X = df.drop(target_col, axis=1)
164     y = df[target_col]
165
166     # Encode categorical variables
167     X_encoded = pd.get_dummies(X, drop_first=True)
168
169     # Split data
170     X_train, X_test, y_train, y_test = train_test_split(
171         X_encoded, y, test_size=0.2, random_state=42, stratify=y
172     )
173
174     # Train model
175     model = RandomForestClassifier(random_state=42, n_estimators=100)
176     model.fit(X_train, y_train)
177
178     # Evaluate
179     train_score = model.score(X_train, y_train)
180     test_score = model.score(X_test, y_test)
181
182     print(f Training accuracy: {train_score:.4f} )
183     print(f Testing accuracy: {test_score:.4f} )
184
185     return model, X_train, X_test, y_train, y_test
186
187 # Train model with clean data
188 print( \n=== TRAINING MODEL WITH CLEAN DATA === )
189 model_proper, X_train, X_test, y_train, y_test = train_model_properly(
190     data_clean, 'target')
191
192 # Advanced data quality monitoring
193 class DataQualityMonitor:
194     Monitor data quality in production
195
196     def __init__(self):
197         self.baseline_stats = {}

```

```

197
198     def set_baseline(self, df):
199         Set baseline statistics for monitoring
200         self.baseline_stats = {}
201
202         for col in df.columns:
203             if df[col].dtype in ['int64', 'float64']:
204                 self.baseline_stats[col] = {
205                     'mean': df[col].mean(),
206                     'std': df[col].std(),
207                     'min': df[col].min(),
208                     'max': df[col].max()
209                 }
210             else:
211                 self.baseline_stats[col] = {
212                     'value_counts': df[col].value_counts().to_dict()
213                 }
214
215     def check_drift(self, new_df, threshold=0.1):
216         Check for data drift
217         drift_alerts = []
218
219         for col in new_df.columns:
220             if col in self.baseline_stats and new_df[col].dtype in ['int64', 'float64']:
221                 current_mean = new_df[col].mean()
222                 baseline_mean = self.baseline_stats[col]['mean']
223
224                 if baseline_mean != 0:
225                     drift = abs(current_mean - baseline_mean) / abs(
baseline_mean)
226                     if drift > threshold:
227                         drift_alerts.append(f Column '{col}' drift: {drift
:.4f} )
228
229         return drift_alerts
230
231 # Example usage of data quality monitoring
232 monitor = DataQualityMonitor()
233 monitor.set_baseline(data_clean)
234
235 # Simulate new data
236 new_data = data_clean.sample(100)
237 drift_alerts = monitor.check_drift(new_data)
238
239 if drift_alerts:
240     print( \n=== DATA DRIFT ALERTS === )
241     for alert in drift_alerts:
242         print(f - {alert} )
243 else:
244     print( \nNo significant data drift detected )

```

Listing 35: Error 5.1.8: Data Quality Issues - Buggy Code

## 5.2 Model Training and Evaluation Errors

```

1 # Buggy: Incorrect cross-validation usage
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import cross_val_score, KFold,
   train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.metrics import accuracy_score
8
9 # Create time series data
10 np.random.seed(42)
11 dates = pd.date_range('2020-01-01', periods=1000, freq='D')
12 data = pd.DataFrame({
13     'date': dates,
14     'feature1': np.random.randn(1000).cumsum(),
15     'feature2': np.random.randn(1000),
16     'target': np.random.choice([0, 1], 1000, p=[0.7, 0.3])
17 })
18
19 X = data[['feature1', 'feature2']]
20 y = data['target']
21
22 print( Data shape: , X.shape)
23 print( Date range: , data['date'].min(), to , data['date'].max())
24
25 # WRONG WAY 1 - Using regular CV on time series data
26 print( \n=== WRONG WAY 1: Regular CV on Time Series === )
27
28 # This is wrong because it mixes future and past data
29 model = RandomForestClassifier(random_state=42)
30 cv_scores_wrong = cross_val_score(model, X, y, cv=5)
31
32 print( Cross-validation scores (WRONG for time series): )
33 print(f Mean: {cv_scores_wrong.mean():.4f} (+/- {cv_scores_wrong.std() *
   2:.4f}) )
34 print( Problem: Future data leaks into training! )
35
36 # WRONG WAY 2 - Data leakage in CV
37 print( \n=== WRONG WAY 2: Data Leakage in CV === )
38
39 # Split data first, then do CV on entire dataset (data leakage!)
40 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
   random_state=42)
41
42 # WRONG: Doing CV on entire dataset after splitting
43 cv_scores_leakage = cross_val_score(model, X, y, cv=5) # Uses test data!
44
45 print( CV scores with data leakage: )
46 print(f Mean: {cv_scores_leakage.mean():.4f} )
47 print( Problem: Test data used in cross-validation! )
48
49 # WRONG WAY 3 - Inconsistent preprocessing in CV

```

```

50 print( \n=== WRONG WAY 3: Inconsistent Preprocessing === )
51
52 from sklearn.preprocessing import StandardScaler
53
54 # WRONG: Fitting scaler on entire dataset before CV
55 scaler = StandardScaler()
56 X_scaled = scaler.fit_transform(X) # Data leakage!
57
58 cv_scores_inconsistent = cross_val_score(model, X_scaled, y, cv=5)
59
60 print( CV scores with inconsistent preprocessing: )
61 print(f Mean: {cv_scores_inconsistent.mean():.4f} )
62 print( Problem: Scaler fitted on entire dataset! )
63
64 # WRONG WAY 4 - Wrong CV for imbalanced data
65 print( \n=== WRONG WAY 4: Wrong CV for Imbalanced Data === )
66
67 # Create imbalanced dataset
68 imbalanced_data = pd.DataFrame({
69     'feature1': np.random.randn(1000),
70     'feature2': np.random.randn(1000),
71     'target': np.random.choice([0, 1], 1000, p=[0.95, 0.05]) # 95% class
72     0
73 })
74 X_imb = imbalanced_data[['feature1', 'feature2']]
75 y_imb = imbalanced_data['target']
76
77 print(f Class distribution: {np.bincount(y_imb)} )
78
79 # WRONG: Regular CV on imbalanced data
80 cv_scores_imbalanced = cross_val_score(model, X_imb, y_imb, cv=5)
81
82 print( CV scores on imbalanced data (regular CV): )
83 print(f Mean: {cv_scores_imbalanced.mean():.4f} )
84 print( Problem: May not represent minority class performance! )
85
86 # CORRECT WAY 1 - Time Series Cross-Validation
87 print( \n=== CORRECT WAY 1: Time Series CV === )
88
89 from sklearn.model_selection import TimeSeriesSplit
90
91 # Use TimeSeriesSplit for time series data
92 tscv = TimeSeriesSplit(n_splits=5)
93 cv_scores_correct_ts = cross_val_score(model, X, y, cv=tscv)
94
95 print( Time series cross-validation scores: )
96 print(f Mean: {cv_scores_correct_ts.mean():.4f} (+/- {
97     cv_scores_correct_ts.std() * 2:.4f}) )
98
99 # CORRECT WAY 2 - Proper data splitting
100 print( \n=== CORRECT WAY 2: Proper Data Splitting === )
101
102 # Split first, then do CV only on training data

```

```

102 X_train_correct, X_test_correct, y_train_correct, y_test_correct =
    train_test_split(
103     X, y, test_size=0.2, random_state=42
104 )
105
106 # CV only on training data
107 cv_scores_correct_split = cross_val_score(model, X_train_correct,
    y_train_correct, cv=5)
108
109 print( CV scores on training data only: )
110 print(f    Mean: {cv_scores_correct_split.mean():.4f} (+/- {
    cv_scores_correct_split.std() * 2:.4f}) )
111
112 # Final evaluation on test set
113 model_final = RandomForestClassifier(random_state=42)
114 model_final.fit(X_train_correct, y_train_correct)
115 test_score = model_final.score(X_test_correct, y_test_correct)
116 print(f Final test score: {test_score:.4f} )
117
118 # CORRECT WAY 3 - Pipeline with proper preprocessing
119 print( \n=== CORRECT WAY 3: Pipeline with Preprocessing === )
120
121 from sklearn.pipeline import Pipeline
122 from sklearn.model_selection import cross_validate
123
124 # Create pipeline that handles preprocessing within CV
125 pipeline = Pipeline([
126     ('scaler', StandardScaler()),
127     ('classifier', RandomForestClassifier(random_state=42))
128 ])
129
130 # CV will properly handle preprocessing within each fold
131 cv_results = cross_validate(pipeline, X, y, cv=5, return_train_score=True)
132
133 print( Pipeline CV scores: )
134 print(f    Train mean: {cv_results['train_score'].mean():.4f} )
135 print(f    Test mean: {cv_results['test_score'].mean():.4f} )
136
137 # CORRECT WAY 4 - Stratified CV for imbalanced data
138 print( \n=== CORRECT WAY 4: Stratified CV for Imbalanced Data === )
139
140 from sklearn.model_selection import StratifiedKFold
141
142 # Use StratifiedKFold to maintain class distribution
143 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
144 cv_scores_stratified = cross_val_score(model, X_imb, y_imb, cv=skf)
145
146 print( Stratified CV scores on imbalanced data: )
147 print(f    Mean: {cv_scores_stratified.mean():.4f} (+/- {
    cv_scores_stratified.std() * 2:.4f}) )
148
149 # Compare with regular CV
150 cv_scores_regular = cross_val_score(model, X_imb, y_imb, cv=5)
151 print(f Regular CV scores: {cv_scores_regular.mean():.4f} )

```

```

152
153 print( Stratified CV better represents minority class performance! )
154
155 # Advanced CV techniques
156 class AdvancedCrossValidation:
157     Advanced cross-validation techniques
158
159     @staticmethod
160     def nested_cv(model, param_grid, X, y, outer_cv=5, inner_cv=3):
161         Nested cross-validation for unbiased model evaluation
162         from sklearn.model_selection import GridSearchCV
163
164         outer_scores = []
165
166         if isinstance(outer_cv, int):
167             outer_cv = KFold(n_splits=outer_cv, shuffle=True, random_state
=42)
168
169         for train_idx, test_idx in outer_cv.split(X):
170             X_train_outer, X_test_outer = X.iloc[train_idx], X.iloc[
test_idx]
171             y_train_outer, y_test_outer = y.iloc[train_idx], y.iloc[
test_idx]
172
173             # Inner CV for hyperparameter tuning
174             grid_search = GridSearchCV(
175                 model, param_grid, cv=inner_cv, scoring='accuracy'
176             )
177             grid_search.fit(X_train_outer, y_train_outer)
178
179             # Evaluate best model on outer test set
180             best_model = grid_search.best_estimator_
181             score = best_model.score(X_test_outer, y_test_outer)
182             outer_scores.append(score)
183
184         return np.array(outer_scores)
185
186     @staticmethod
187     def group_cv(model, X, y, groups, cv=5):
188         Cross-validation with group constraints
189         from sklearn.model_selection import GroupKFold
190
191         group_kfold = GroupKFold(n_splits=cv)
192         scores = cross_val_score(model, X, y, cv=group_kfold, groups=
groups)
193         return scores
194
195 # Example of nested CV
196 print( \n=== ADVANCED: Nested Cross-Validation === )
197
198 param_grid = {
199     'n_estimators': [50, 100, 200],
200     'max_depth': [3, 5, 7]
201 }

```

```

202
203 nested_scores = AdvancedCrossValidation.nested_cv(
204     RandomForestClassifier(random_state=42),
205     param_grid,
206     X_train_correct,
207     y_train_correct
208 )
209
210 print( Nested CV scores: )
211 print(f    Mean: {nested_scores.mean():.4f} (+/- {nested_scores.std() *
212     2:.4f}) )
213 print( Unbiased estimate of model performance! )

```

Listing 36: Error 5.2.1: Incorrect Cross-Validation Usage - Buggy Code

```

1 # Buggy: Common hyperparameter tuning mistakes
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV,
   train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.svm import SVC
7 from sklearn.metrics import accuracy_score
8 from sklearn.preprocessing import StandardScaler
9
10 # Create sample data
11 np.random.seed(42)
12 X = np.random.randn(1000, 10)
13 y = np.random.choice([0, 1], 1000, p=[0.7, 0.3])
14
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
   random_state=42)
16
17 print( Data shape: , X.shape)
18 print( Train/Test split: , X_train.shape, X_test.shape)
19
20 # WRONG WAY 1 - Tuning on test set (data leakage)
21 print( \n=== WRONG WAY 1: Tuning on Test Set === )
22
23 # This is completely wrong - using test set for hyperparameter tuning!
24 param_grid_wrong = {
25     'n_estimators': [50, 100, 200],
26     'max_depth': [3, 5, 7]
27 }
28
29 grid_search_wrong = GridSearchCV(
30     RandomForestClassifier(random_state=42),
31     param_grid_wrong,
32     cv=5,
33     scoring='accuracy'
34 )
35
36 # WRONG: Fitting on entire dataset including test data
37 grid_search_wrong.fit(np.vstack([X_train, X_test]), np.hstack([y_train,
   y_test]))

```

```

38
39 print( Best parameters (WRONG - data leakage): , grid_search_wrong.
    best_params_)
40 print( Best score: , grid_search_wrong.best_score_)
41
42 # Test on same data - overly optimistic!
43 test_score_wrong = grid_search_wrong.score(X_test, y_test)
44 print( Test score (overly optimistic): , test_score_wrong)
45 print( Problem: Test set used in hyperparameter tuning! )
46
47 # WRONG WAY 2 - No proper validation set
48 print( \n=== WRONG WAY 2: No Proper Validation Set === )
49
50 # WRONG: Using GridSearchCV without proper separation
51 grid_search_no_val = GridSearchCV(
52     RandomForestClassifier(random_state=42),
53     param_grid_wrong,
54     cv=3, # Too few folds
55     scoring='accuracy'
56 )
57
58 # Fitting on training data is correct, but...
59 grid_search_no_val.fit(X_train, y_train)
60
61 # But then evaluating on test set without considering overfitting to
    validation folds
62 best_model_no_val = grid_search_no_val.best_estimator_
63 test_score_no_val = best_model_no_val.score(X_test, y_test)
64
65 print( Best parameters: , grid_search_no_val.best_params_)
66 print( CV score: , grid_search_no_val.best_score_)
67 print( Test score: , test_score_no_val)
68 print( Problem: May be overfitting to validation folds! )
69
70 # WRONG WAY 3 - Exhaustive search with too many parameters
71 print( \n=== WRONG WAY 3: Exhaustive Search with Too Many Parameters === )
72
73 # WRONG: Too many parameter combinations
74 param_grid_exhaustive = {
75     'C': [0.001, 0.01, 0.1, 1, 10, 100, 1000], # 7 values
76     'gamma': [0.001, 0.01, 0.1, 1, 10, 100], # 6 values
77     'kernel': ['rbf', 'poly', 'sigmoid'] # 3 values
78 }
79
80 # Total combinations: 7      6      3 = 126 combinations!
81 # With 5-fold CV: 126      5 = 630 model fits!
82
83 print(f Total combinations: {7 * 6 * 3} )
84 print(f Total model fits with 5-fold CV: {7 * 6 * 3 * 5} )
85 print( Problem: Computationally expensive, may not be necessary! )
86
87 # WRONG WAY 4 - Inappropriate scoring metric
88 print( \n=== WRONG WAY 4: Inappropriate Scoring Metric === )
89

```



```

90 # For imbalanced dataset, accuracy might not be the best metric
91 imbalanced_y = np.random.choice([0, 1], 1000, p=[0.95, 0.05])
92
93 param_grid_metric = {
94     'n_estimators': [50, 100],
95     'max_depth': [3, 5]
96 }
97
98 # WRONG: Using accuracy for imbalanced data
99 grid_search_accuracy = GridSearchCV(
100     RandomForestClassifier(random_state=42),
101     param_grid_metric,
102     cv=5,
103     scoring='accuracy' # Not appropriate for imbalanced data
104 )
105
106 grid_search_accuracy.fit(X_train, imbalanced_y[:800])
107
108 print( Best parameters with accuracy scoring: , grid_search_accuracy.
109       best_params_)
109 print( Best accuracy score: , grid_search_accuracy.best_score_)
110 print( Problem: May favor models that predict majority class! )
111
112 # CORRECT WAY 1 - Proper train/validation/test split
113 print( \n=== CORRECT WAY 1: Proper Data Splitting === )
114
115 # Split into train/validation/test
116 X_temp, X_test_final, y_temp, y_test_final = train_test_split(
117     X, y, test_size=0.2, random_state=42
118 )
119
120 X_train_final, X_val_final, y_train_final, y_val_final = train_test_split(
121     X_temp, y_temp, test_size=0.25, random_state=42 # 0.25 * 0.8 = 0.2 of
122     total
123 )
124
125 print(f Train: {X_train_final.shape}, Val: {X_val_final.shape}, Test: {
126     X_test_final.shape} )
127
128 # Hyperparameter tuning on validation set
129 param_grid_correct = {
130     'n_estimators': [50, 100, 200],
131     'max_depth': [3, 5, 7, None]
132 }
133
134 # Grid search on training data with cross-validation
135 grid_search_correct = GridSearchCV(
136     RandomForestClassifier(random_state=42),
137     param_grid_correct,
138     cv=5,
139     scoring='accuracy',
140     n_jobs=-1 # Use all processors

```

```

141 grid_search_correct.fit(X_train_final, y_train_final)
142
143 # Evaluate on validation set
144 val_score = grid_search_correct.score(X_val_final, y_val_final)
145 print( Best parameters: , grid_search_correct.best_params_)
146 print( CV score: , grid_search_correct.best_score_)
147 print( Validation score: , val_score)
148
149 # Final evaluation on test set
150 best_model = grid_search_correct.best_estimator_
151 test_score_final = best_model.score(X_test_final, y_test_final)
152 print( Final test score: , test_score_final)
153
154 # CORRECT WAY 2 - Randomized search for efficiency
155 print( \n=== CORRECT WAY 2: Randomized Search === )
156
157 from scipy.stats import randint, uniform
158
159 # Use RandomizedSearchCV for efficiency
160 param_dist = {
161     'n_estimators': randint(50, 300),
162     'max_depth': randint(3, 20),
163     'min_samples_split': randint(2, 20),
164     'min_samples_leaf': randint(1, 10)
165 }
166
167 random_search = RandomizedSearchCV(
168     RandomForestClassifier(random_state=42),
169     param_distributions=param_dist,
170     n_iter=50, # Try 50 random combinations instead of all
171     cv=5,
172     scoring='accuracy',
173     n_jobs=-1,
174     random_state=42
175 )
176
177 random_search.fit(X_train, y_train)
178
179 print( Random search best parameters: , random_search.best_params_)
180 print( Random search best score: , random_search.best_score_)
181
182 # CORRECT WAY 3 - Appropriate scoring for imbalanced data
183 print( \n=== CORRECT WAY 3: Appropriate Scoring === )
184
185 # For imbalanced data, use F1-score or AUC
186 from sklearn.metrics import make_scorer, f1_score
187
188 # Custom scorer for F1-score
189 f1_scorer = make_scorer(f1_score, average='weighted')
190
191 grid_search_f1 = GridSearchCV(
192     RandomForestClassifier(random_state=42),
193     param_grid_correct,
194     cv=5,

```

```

195     scoring=f1_scorer, # Appropriate for imbalanced data
196     n_jobs=-1
197 )
198
199 grid_search_f1.fit(X_train, imbalanced_y[:800])
200
201 print( Best parameters with F1 scoring: , grid_search_f1.best_params_)
202 print( Best F1 score: , grid_search_f1.best_score_)
203
204 # CORRECT WAY 4 - Pipeline with preprocessing
205 print( \n=== CORRECT WAY 4: Pipeline with Preprocessing === )
206
207 from sklearn.pipeline import Pipeline
208
209 # Create pipeline that includes preprocessing
210 pipeline = Pipeline([
211     ('scaler', StandardScaler()),
212     ('classifier', SVC(random_state=42))
213 ])
214
215 param_grid_pipeline = {
216     'classifier__C': [0.1, 1, 10],
217     'classifier__gamma': [0.001, 0.01, 0.1],
218     'classifier__kernel': ['rbf', 'linear']
219 }
220
221 grid_search_pipeline = GridSearchCV(
222     pipeline,
223     param_grid_pipeline,
224     cv=5,
225     scoring='accuracy',
226     n_jobs=-1
227 )
228
229 grid_search_pipeline.fit(X_train, y_train)
230
231 print( Pipeline best parameters: , grid_search_pipeline.best_params_)
232 print( Pipeline best score: , grid_search_pipeline.best_score_)
233
234 # Advanced hyperparameter tuning techniques
235 class AdvancedHyperparameterTuning:
236     Advanced hyperparameter tuning techniques
237
238     @staticmethod
239     def bayesian_optimization(X, y, model_class, param_space, n_iter=50):
240         Bayesian optimization using scikit-optimize
241         try:
242             from skopt import BayesSearchCV
243             from skopt.space import Real, Integer, Categorical
244
245             # Convert parameter space
246             skopt_space = {}
247             for param, values in param_space.items():
248                 if isinstance(values, list) and all(isinstance(v, (int, np

```

```

249         .integer)) for v in values):
250             skopt_space[param] = Integer(min(values), max(values))
251             elif isinstance(values, list) and all(isinstance(v, (float
, np.floating)) for v in values):
252                 skopt_space[param] = Real(min(values), max(values))
253                 elif isinstance(values, list):
254                     skopt_space[param] = Categorical(values)
255                 else:
256                     skopt_space[param] = values
257
258     bayes_search = BayesSearchCV(
259         estimator=model_class(random_state=42),
260         search_spaces=skopt_space,
261         n_iter=n_iter,
262         cv=5,
263         scoring='accuracy',
264         n_jobs=-1,
265         random_state=42
266     )
267
268     bayes_search.fit(X, y)
269     return bayes_search
270
271 except ImportError:
272     print( scikit-optimize not installed. Using RandomizedSearchCV
instead. )
273     from scipy.stats import randint, uniform
274
275     # Convert to scipy distributions
276     param_dist = {}
277     for param, values in param_space.items():
278         if isinstance(values, list) and all(isinstance(v, (int, np
.integer)) for v in values):
279             param_dist[param] = randint(min(values), max(values) +
1)
280             elif isinstance(values, list) and all(isinstance(v, (float
, np.floating)) for v in values):
281                 param_dist[param] = uniform(min(values), max(values) -
min(values))
282             else:
283                 param_dist[param] = values
284
285     random_search = RandomizedSearchCV(
286         model_class(random_state=42),
287         param_distributions=param_dist,
288         n_iter=n_iter,
289         cv=5,
290         scoring='accuracy',
291         n_jobs=-1,
292         random_state=42
293     )
294
295     random_search.fit(X, y)
296     return random_search

```

```

296
297     @staticmethod
298     def multi_metric_tuning(X, y, model, param_grid, scoring_metrics):
299         Tune hyperparameters using multiple metrics
300         from sklearn.model_selection import GridSearchCV
301
302         multi_search = GridSearchCV(
303             model,
304             param_grid,
305             cv=5,
306             scoring=scoring_metrics,
307             refit='f1_weighted', # Which metric to use for final model
308             n_jobs=-1
309         )
310
311         multi_search.fit(X, y)
312         return multi_search
313
314 # Example of advanced tuning
315 print( \n=== ADVANCED: Bayesian Optimization === )
316
317 param_space = {
318     'n_estimators': [50, 100, 200, 300],
319     'max_depth': [3, 5, 7, 10, None],
320     'min_samples_split': [2, 5, 10, 20],
321     'min_samples_leaf': [1, 2, 4, 8]
322 }
323
324 # This would use Bayesian optimization if scikit-optimize is installed
325 # advanced_search = AdvancedHyperparameterTuning.bayesian_optimization(
326 #     X_train, y_train, RandomForestClassifier, param_space, n_iter=30
327 # )
328
329 print( Advanced tuning techniques available for more efficient
        hyperparameter optimization! )

```

Listing 37: Error 5.2.2: Hyperparameter Tuning Mistakes - Buggy Code

```

1 # Buggy: Common model selection mistakes
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split, cross_val_score
5 from sklearn.ensemble import RandomForestClassifier,
6     GradientBoostingClassifier
7 from sklearn.svm import SVC
8 from sklearn.linear_model import LogisticRegression
9 from sklearn.naive_bayes import GaussianNB
10 from sklearn.metrics import accuracy_score, classification_report
11 from sklearn.preprocessing import StandardScaler
12
13 # Create different types of datasets to demonstrate various issues
14 np.random.seed(42)
15
16 # 1. Linearly separable data
17 X_linear = np.random.randn(1000, 2)

```

```

17 y_linear = (X_linear[:, 0] + X_linear[:, 1] > 0).astype(int)
18
19 # 2. Non-linear data (XOR pattern)
20 X_nonlinear = np.random.randn(1000, 2)
21 y_nonlinear = ((X_nonlinear[:, 0] > 0) ^ (X_nonlinear[:, 1] > 0)).astype(
    int)
22
23 # 3. High-dimensional sparse data
24 X_sparse = np.random.choice([0, 1], size=(1000, 1000), p=[0.95, 0.05])
25 y_sparse = np.random.choice([0, 1], 1000, p=[0.7, 0.3])
26
27 # 4. Imbalanced data
28 X_imbalanced = np.random.randn(1000, 5)
29 y_imbalanced = np.random.choice([0, 1], 1000, p=[0.95, 0.05])
30
31 print( Dataset shapes: )
32 print(f Linear: {X_linear.shape} )
33 print(f Non-linear: {X_nonlinear.shape} )
34 print(f Sparse: {X_sparse.shape} )
35 print(f Imbalanced: {X_imbalanced.shape} )
36
37 # WRONG WAY 1 - Always using the same model
38 print( \n=== WRONG WAY 1: Always Using Same Model === )
39
40 # WRONG: Always using Random Forest regardless of data characteristics
41 def wrong_model_selection(X, y):
42     Always use Random Forest - not considering data characteristics
43     model = RandomForestClassifier(random_state=42)
44     model.fit(X, y)
45     return model
46
47 # This works but might not be optimal
48 model_linear_wrong = wrong_model_selection(X_linear, y_linear)
49 model_nonlinear_wrong = wrong_model_selection(X_nonlinear, y_nonlinear)
50
51 print( Random Forest accuracy on linear data: , model_linear_wrong.score(
    X_linear, y_linear))
52 print( Random Forest accuracy on non-linear data: , model_nonlinear_wrong.
    score(X_nonlinear, y_nonlinear))
53 print( Problem: Not considering that linear models might be better for
    linear data! )
54
55 # WRONG WAY 2 - No proper model comparison
56 print( \n=== WRONG WAY 2: No Proper Model Comparison === )
57
58 # WRONG: Just trying a few models without systematic comparison
59 models_quick = {
60     'RF': RandomForestClassifier(random_state=42),
61     'SVM': SVC(random_state=42),
62     'LR': LogisticRegression(random_state=42)
63 }
64
65 def quick_model_comparison(X, y, models):
66     Quick comparison without proper validation

```

```

67     results = {}
68     for name, model in models.items():
69         model.fit(X, y) # WRONG: Fitting on entire dataset
70         score = model.score(X, y) # WRONG: Evaluating on training data
71         results[name] = score
72     return results
73
74 # This gives overly optimistic results
75 quick_results = quick_model_comparison(X_linear, y_linear, models_quick)
76 print( Quick comparison results (OVERLY OPTIMISTIC): )
77 for name, score in quick_results.items():
78     print(f {name}: {score:.4f} )
79 print( Problem: No proper validation, results are overly optimistic! )
80
81 # WRONG WAY 3 - Ignoring data preprocessing requirements
82 print( \n=== WRONG WAY 3: Ignoring Preprocessing === )
83
84 # WRONG: Using SVM without scaling on data that needs it
85 X_scaled_train, X_scaled_test, y_scaled_train, y_scaled_test =
86     train_test_split(
87         X_linear, y_linear, test_size=0.2, random_state=42
88     )
89 # Without scaling (WRONG for SVM)
90 svm_no_scaling = SVC(random_state=42)
91 svm_no_scaling.fit(X_scaled_train, y_scaled_train)
92 score_no_scaling = svm_no_scaling.score(X_scaled_test, y_scaled_test)
93
94 # With proper scaling
95 scaler = StandardScaler()
96 X_scaled_proper = scaler.fit_transform(X_linear)
97 X_scaled_train_proper, X_scaled_test_proper, y_scaled_train_proper,
98     y_scaled_test_proper = train_test_split(
99         X_scaled_proper, y_linear, test_size=0.2, random_state=42
100     )
101 svm_with_scaling = SVC(random_state=42)
102 svm_with_scaling.fit(X_scaled_train_proper, y_scaled_train_proper)
103 score_with_scaling = svm_with_scaling.score(X_scaled_test_proper,
104     y_scaled_test_proper)
105 print( SVM without scaling: , score_no_scaling)
106 print( SVM with scaling: , score_with_scaling)
107 print( Problem: SVM performance significantly affected by scaling! )
108
109 # WRONG WAY 4 - Inappropriate models for data types
110 print( \n=== WRONG WAY 4: Inappropriate Models === )
111
112 # WRONG: Using tree-based models on high-dimensional sparse data
113 print( High-dimensional sparse data shape: , X_sparse.shape)
114
115 # Tree-based models don't work well with sparse high-dimensional data
116 rf_sparse_wrong = RandomForestClassifier(random_state=42)
117 rf_sparse_wrong.fit(X_sparse, y_sparse)

```

```

118 sparse_score = rf_sparse_wrong.score(X_sparse, y_sparse)
119
120 print( Random Forest on sparse data: , sparse_score)
121 print( Problem: Tree-based models often perform poorly on high-dimensional
      sparse data! )
122
123 # CORRECT WAY 1 - Systematic model selection based on data characteristics
124 print( \n== CORRECT WAY 1: Systematic Model Selection == )
125
126 def analyze_data_characteristics(X, y):
127     Analyze data characteristics for model selection
128     characteristics = {
129         'n_samples': X.shape[0],
130         'n_features': X.shape[1],
131         'feature_to_sample_ratio': X.shape[1] / X.shape[0],
132         'sparsity': np.mean(X == 0) if X.shape[1] > 100 else 0,
133         'class_balance': min(np.bincount(y)) / max(np.bincount(y)),
134         'linearity_indicator': None # Would require more complex analysis
135     }
136
137     return characteristics
138
139 def recommend_models(characteristics):
140     Recommend appropriate models based on data characteristics
141     recommendations = []
142
143     # Feature-to-sample ratio analysis
144     if characteristics['feature_to_sample_ratio'] > 0.5:
145         recommendations.append( Use regularization (Ridge, Lasso) or
dimensionality reduction )
146         recommendations.append( Avoid tree-based models )
147
148     # Sparsity analysis
149     if characteristics['sparsity'] > 0.8:
150         recommendations.append( Use linear models with L1 regularization )
151         recommendations.append( Consider Naive Bayes for text-like data )
152
153     # Class balance analysis
154     if characteristics['class_balance'] < 0.1:
155         recommendations.append( Use appropriate metrics (F1, AUC) )
156         recommendations.append( Consider resampling techniques )
157
158     # General recommendations
159     recommendations.append( Use cross-validation for model comparison )
160     recommendations.append( Consider preprocessing requirements )
161
162     return recommendations
163
164 # Analyze datasets
165 char_linear = analyze_data_characteristics(X_linear, y_linear)
166 char_sparse = analyze_data_characteristics(X_sparse, y_sparse)
167
168 print( Linear data characteristics: )
169 for key, value in char_linear.items():

```



```

170     print(f    {key}: {value} )
171
172 print( \nSparse data characteristics: )
173 for key, value in char_sparse.items():
174     print(f    {key}: {value} )
175
176 print( \nModel recommendations for linear data: )
177 for rec in recommend_models(char_linear):
178     print(f    - {rec} )
179
180 print( \nModel recommendations for sparse data: )
181 for rec in recommend_models(char_sparse):
182     print(f    - {rec} )
183
184 # CORRECT WAY 2 - Proper model comparison with cross-validation
185 print( \n=== CORRECT WAY 2: Proper Model Comparison === )
186
187 def proper_model_comparison(X, y, test_size=0.2):
188     Proper model comparison with cross-validation
189
190     # Split data
191     X_train, X_test, y_train, y_test = train_test_split(
192         X, y, test_size=test_size, random_state=42
193     )
194
195     # Define models
196     models = {
197         'Logistic Regression': LogisticRegression(random_state=42,
198 max_iter=1000),
199         'Random Forest': RandomForestClassifier(random_state=42,
200 n_estimators=100),
201         'SVM': SVC(random_state=42, probability=True),
202         'Naive Bayes': GaussianNB()
203     }
204
205     # Preprocessing pipeline
206     scaler = StandardScaler()
207     X_train_scaled = scaler.fit_transform(X_train)
208     X_test_scaled = scaler.transform(X_test)
209
210     results = {}
211
212     for name, model in models.items():
213         # Use appropriate data (scaled for SVM and LR)
214         if name in ['Logistic Regression', 'SVM']:
215             X_train_use = X_train_scaled
216             X_test_use = X_test_scaled
217         else:
218             X_train_use = X_train
219             X_test_use = X_test
220
221         # Cross-validation
222         cv_scores = cross_val_score(model, X_train_use, y_train, cv=5)

```

```

222     # Fit and test
223     model.fit(X_train_use, y_train)
224     test_score = model.score(X_test_use, y_test)
225
226     results[name] = {
227         'cv_mean': cv_scores.mean(),
228         'cv_std': cv_scores.std(),
229         'test_score': test_score
230     }
231
232     return results
233
234 # Compare models properly
235 print( Proper model comparison on linear data: )
236 linear_results = proper_model_comparison(X_linear, y_linear)
237
238 for name, scores in linear_results.items():
239     print(f    {name}: )
240     print(f    CV: {scores['cv_mean']:.4f} (+/- {scores['cv_std'] * 2:.4f
241         }) )
242     print(f    Test: {scores['test_score']:.4f} )
243
244 # CORRECT WAY 3 - Model selection for specific data types
245 print( \n=== CORRECT WAY 3: Data-Type Specific Selection === )
246
247 def select_model_for_data_type(X, y):
248     Select appropriate model based on data type
249
250     n_samples, n_features = X.shape
251
252     # High-dimensional data
253     if n_features > n_samples or n_features > 100:
254         print( High-dimensional data detected )
255         if np.mean(X == 0) > 0.8: # Sparse data
256             print( Using Naive Bayes (good for sparse data) )
257             return GaussianNB()
258         else:
259             print( Using Logistic Regression with L1 regularization )
260             return LogisticRegression(penalty='l1', solver='liblinear',
261                                     random_state=42)
262
263     # Low-dimensional data
264     elif n_features < 20:
265         print( Low-dimensional data detected )
266         print( Using Random Forest (good general-purpose model) )
267         return RandomForestClassifier(random_state=42, n_estimators=100)
268
269     # Medium-dimensional data
270     else:
271         print( Medium-dimensional data detected )
272         print( Using Gradient Boosting (good for structured data) )
273         return GradientBoostingClassifier(random_state=42)
274
275 # Test model selection

```

```

274 print( Model selection for linear data: )
275 model_linear_correct = select_model_for_data_type(X_linear, y_linear)
276
277 print( \nModel selection for sparse data: )
278 model_sparse_correct = select_model_for_data_type(X_sparse, y_sparse)
279
280 # CORRECT WAY 4 - Ensemble methods for robust selection
281 print( \n== CORRECT WAY 4: Ensemble Methods == )
282
283 from sklearn.ensemble import VotingClassifier
284
285 def create_ensemble_model(X, y):
286     Create ensemble model combining multiple approaches
287
288     # Split data
289     X_train, X_test, y_train, y_test = train_test_split(
290         X, y, test_size=0.2, random_state=42
291     )
292
293     # Scale data for models that need it
294     scaler = StandardScaler()
295     X_train_scaled = scaler.fit_transform(X_train)
296     X_test_scaled = scaler.transform(X_test)
297
298     # Create diverse set of models
299     models = [
300         ('lr', LogisticRegression(random_state=42, max_iter=1000)),
301         ('rf', RandomForestClassifier(random_state=42, n_estimators=100)),
302         ('svm', SVC(probability=True, random_state=42)),
303         ('gb', GradientBoostingClassifier(random_state=42))
304     ]
305
306     # Create voting classifier
307     ensemble = VotingClassifier(
308         estimators=models,
309         voting='soft' # Use probability averaging
310     )
311
312     # Fit ensemble
313     ensemble.fit(X_train_scaled, y_train) # Use scaled data for all
314
315     return ensemble, scaler
316
317 # Create ensemble for linear data
318 ensemble_linear, scaler_linear = create_ensemble_model(X_linear, y_linear)
319 ensemble_score = ensemble_linear.score(scaler_linear.transform(X_linear),
320 y_linear)
321 print( Ensemble model accuracy on linear data: , ensemble_score)
322
323 # Advanced model selection techniques
324 class AdvancedModelSelector:
325     Advanced model selection techniques
326
327     @staticmethod

```

```

327     def automated_model_selection(X, y, time_limit=300):
328         Automated model selection with time constraints
329         import time
330         from sklearn.model_selection import cross_validate
331
332         models = {
333             'LogisticRegression': LogisticRegression(random_state=42,
max_iter=1000),
334             'RandomForest': RandomForestClassifier(random_state=42,
n_estimators=100),
335             'SVM': SVC(random_state=42, probability=True),
336             'GradientBoosting': GradientBoostingClassifier(random_state
=42),
337             'NaiveBayes': GaussianNB()
338         }
339
340         results = {}
341         start_time = time.time()
342
343         for name, model in models.items():
344             # Check time limit
345             if time.time() - start_time > time_limit:
346                 print(f Time limit reached. Skipping {name} )
347                 continue
348
349             try:
350                 # Preprocessing if needed
351                 if name in ['LogisticRegression', 'SVM']:
352                     scaler = StandardScaler()
353                     X_processed = scaler.fit_transform(X)
354                 else:
355                     X_processed = X
356
357                 # Quick cross-validation
358                 cv_results = cross_validate(
359                     model, X_processed, y,
360                     cv=3, # Quick CV
361                     scoring='accuracy',
362                     return_train_score=True
363                 )
364
365                 results[name] = {
366                     'mean_cv_score': cv_results['test_score'].mean(),
367                     'std_cv_score': cv_results['test_score'].std(),
368                     'training_time': time.time() - start_time
369                 }
370
371             except Exception as e:
372                 print(f Error with {name}: {e} )
373                 continue
374
375         # Select best model
376         if results:
377             best_model_name = max(results.keys(), key=lambda x: results[x

```

```

    [['mean_cv_score']]
378         print(f Best model: {best_model_name} )
379         print(f CV Score: {results[best_model_name]['mean_cv_score
    ']:.4f} )
380         return best_model_name, results[best_model_name]
381     else:
382         return None, None
383
384     @staticmethod
385     def model_complexity_analysis(X, y, model_class, complexity_param,
    param_values):
386         Analyze model complexity vs performance
387         from sklearn.model_selection import validation_curve
388
389         train_scores, val_scores = validation_curve(
390             model_class(random_state=42), X, y,
391             param_name=complexity_param,
392             param_range=param_values,
393             cv=5, scoring='accuracy'
394         )
395
396         train_mean = np.mean(train_scores, axis=1)
397         train_std = np.std(train_scores, axis=1)
398         val_mean = np.mean(val_scores, axis=1)
399         val_std = np.std(val_scores, axis=1)
400
401         # Find optimal complexity
402         optimal_idx = np.argmax(val_mean)
403         optimal_value = param_values[optimal_idx]
404
405         return {
406             'train_mean': train_mean,
407             'train_std': train_std,
408             'val_mean': val_mean,
409             'val_std': val_std,
410             'optimal_value': optimal_value,
411             'optimal_score': val_mean[optimal_idx]
412         }
413
414     # Example of advanced model selection
415     print( \n=== ADVANCED: Automated Model Selection === )
416
417     # This would run automated selection
418     # best_model, best_results = AdvancedModelSelector.
    automated_model_selection(X_linear, y_linear)
419
420     print( Advanced model selection techniques provide automated and robust
    model selection! )
421
422     # Model selection based on problem requirements
423     def select_model_by_requirements(X, y, requirements):
424         Select model based on specific requirements
425
426         recommended_models = []

```

```

427
428     if requirements.get('interpretability', False):
429         recommended_models.extend(['LogisticRegression', 'DecisionTree'])
430
431     if requirements.get('speed', False):
432         recommended_models.extend(['LogisticRegression', 'NaiveBayes'])
433
434     if requirements.get('accuracy', False):
435         recommended_models.extend(['RandomForest', 'GradientBoosting'])
436
437     if requirements.get('scalability', False):
438         recommended_models.extend(['LogisticRegression', 'SGDClassifier'])
439
440     return list(set(recommended_models)) # Remove duplicates
441
442 # Example usage
443 requirements = {
444     'interpretability': True,
445     'speed': True,
446     'accuracy': True
447 }
448
449 recommended = select_model_by_requirements(X_linear, y_linear,
450                                           requirements)
451 print(f \nModels recommended for requirements {requirements}: {recommended}
452       )

```

Listing 38: Error 5.2.3: Model Selection Mistakes - Buggy Code

```

1 # Buggy: Common evaluation metric mistakes
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.metrics import (accuracy_score, precision_score, recall_score
8                               , f1_score,
9                               roc_auc_score, confusion_matrix,
10                               classification_report)
11
12 # Create different types of datasets to demonstrate issues
13 np.random.seed(42)
14
15 # 1. Balanced dataset
16 X_balanced, y_balanced = make_classification(
17     n_samples=1000, n_features=10, n_classes=2,
18     weights=[0.5, 0.5], random_state=42
19 )
20
21 # 2. Imbalanced dataset (95% class 0, 5% class 1)
22 X_imbalanced, y_imbalanced = make_classification(
23     n_samples=1000, n_features=10, n_classes=2,
24     weights=[0.95, 0.05], random_state=42

```

```

25 )
26
27 # 3. Multi-class dataset
28 X_multiclass, y_multiclass = make_classification(
29     n_samples=1000, n_features=10, n_classes=3,
30     n_informative=8, random_state=42
31 )
32
33 print( Dataset shapes: )
34 print(f Balanced: {X_balanced.shape}, classes: {np.bincount(y_balanced)} )
35 print(f Imbalanced: {X_imbalanced.shape}, classes: {np.bincount(
36     y_imbalanced)} )
37
38 # WRONG WAY 1 - Using accuracy for imbalanced data
39 print( \n=== WRONG WAY 1: Accuracy for Imbalanced Data === )
40
41 # Train models
42 rf_balanced = RandomForestClassifier(random_state=42)
43 rf_imbalanced = RandomForestClassifier(random_state=42)
44
45 rf_balanced.fit(X_balanced, y_balanced)
46 rf_imbalanced.fit(X_imbalanced, y_imbalanced)
47
48 # WRONG: Using accuracy for imbalanced dataset
49 y_pred_imbalanced = rf_imbalanced.predict(X_imbalanced)
50 accuracy_wrong = accuracy_score(y_imbalanced, y_pred_imbalanced)
51
52 print( Imbalanced dataset results: )
53 print(f Class distribution: {np.bincount(y_imbalanced)} )
54 print(f Accuracy (MISLEADING): {accuracy_wrong:.4f} )
55
56 # The model might just predict the majority class all the time!
57 majority_class = np.bincount(y_imbalanced).argmax()
58 naive_accuracy = np.mean(y_imbalanced == majority_class)
59 print(f Naive classifier accuracy: {naive_accuracy:.4f} )
60 print( Problem: High accuracy doesn't mean good performance on minority
61     class! )
62
63 # WRONG WAY 2 - Inappropriate metrics for multi-class
64 print( \n=== WRONG WAY 2: Inappropriate Multi-class Metrics === )
65
66 # Train on multi-class data
67 rf_multiclass = RandomForestClassifier(random_state=42)
68 rf_multiclass.fit(X_multiclass, y_multiclass)
69
70 y_pred_multiclass = rf_multiclass.predict(X_multiclass)
71
72 # WRONG: Using binary metrics without proper averaging
73 try:
74     precision_wrong = precision_score(y_multiclass, y_pred_multiclass) #
75     Error!
76 except ValueError as e:

```

```

75     print(f Error with precision_score: {e} )
76
77 # WRONG: Not specifying averaging method
78 try:
79     recall_wrong = recall_score(y_multiclass, y_pred_multiclass) # Error!
80 except ValueError as e:
81     print(f Error with recall_score: {e} )
82
83 print( Problem: Multi-class requires averaging method specification! )
84
85 # WRONG WAY 3 - ROC AUC without probability scores
86 print( \n=== WRONG WAY 3: ROC AUC without Probabilities === )
87
88 # WRONG: Using predictions instead of probabilities for ROC AUC
89 try:
90     auc_wrong = roc_auc_score(y_balanced, y_pred_balanced) # Using
    predictions!
91     print(f AUC with predictions (WRONG): {auc_wrong:.4f} )
92 except Exception as e:
93     print(f Error with AUC calculation: {e} )
94
95 print( Problem: ROC AUC needs probability scores, not class predictions! )
96
97 # WRONG WAY 4 - Single metric evaluation
98 print( \n=== WRONG WAY 4: Single Metric Evaluation === )
99
100 # WRONG: Only looking at one metric
101 model1 = RandomForestClassifier(random_state=42)
102 model2 = LogisticRegression(random_state=42)
103
104 model1.fit(X_balanced, y_balanced)
105 model2.fit(X_balanced, y_balanced)
106
107 y_pred1 = model1.predict(X_balanced)
108 y_pred2 = model2.predict(X_balanced)
109
110 acc1 = accuracy_score(y_balanced, y_pred1)
111 acc2 = accuracy_score(y_balanced, y_pred2)
112
113 print( Single metric comparison: )
114 print(f Model 1 accuracy: {acc1:.4f} )
115 print(f Model 2 accuracy: {acc2:.4f} )
116 print( Problem: Single metric doesn't tell the full story! )
117
118 # CORRECT WAY 1 - Appropriate metrics for imbalanced data
119 print( \n=== CORRECT WAY 1: Appropriate Imbalanced Metrics === )
120
121 # Use multiple metrics for imbalanced data
122 def evaluate_imbalanced_model(y_true, y_pred, y_proba=None):
123     Proper evaluation for imbalanced datasets
124
125     results = {
126         'accuracy': accuracy_score(y_true, y_pred),
127         'precision': precision_score(y_true, y_pred, average='weighted'),

```



```

128         'recall': recall_score(y_true, y_pred, average='weighted'),
129         'f1': f1_score(y_true, y_pred, average='weighted')
130     }
131
132     # Per-class metrics for minority class
133     precision_per_class = precision_score(y_true, y_pred, average=None)
134     recall_per_class = recall_score(y_true, y_pred, average=None)
135
136     # Focus on minority class (class 1 in our case)
137     results['minority_precision'] = precision_per_class[1]
138     results['minority_recall'] = recall_per_class[1]
139
140     # ROC AUC if probabilities provided
141     if y_proba is not None:
142         results['auc'] = roc_auc_score(y_true, y_proba[:, 1])
143
144     return results
145
146 # Proper evaluation
147 y_pred_imb_correct = rf_imbalanced.predict(X_imbalanced)
148 y_proba_imb_correct = rf_imbalanced.predict_proba(X_imbalanced)
149
150 results_imb = evaluate_imbalanced_model(
151     y_imbalanced, y_pred_imb_correct, y_proba_imb_correct
152 )
153
154 print( Proper imbalanced evaluation: )
155 for metric, value in results_imb.items():
156     print(f    {metric}: {value:.4f} )
157
158 print(f Majority class accuracy: {np.mean(y_imbalanced == 0):.4f} )
159 print( Notice the difference between overall accuracy and minority class
160       performance! )
161
162 # CORRECT WAY 2 - Proper multi-class metrics
163 print( \n=== CORRECT WAY 2: Proper Multi-class Metrics === )
164
165 def evaluate_multiclass_model(y_true, y_pred, y_proba=None):
166     Proper evaluation for multi-class datasets
167
168     results = {
169         'accuracy': accuracy_score(y_true, y_pred),
170         'precision_macro': precision_score(y_true, y_pred, average='macro'
171         ),
172         'precision_weighted': precision_score(y_true, y_pred, average='
173         weighted'),
174         'recall_macro': recall_score(y_true, y_pred, average='macro'),
175         'recall_weighted': recall_score(y_true, y_pred, average='weighted'
176         ),
177         'f1_macro': f1_score(y_true, y_pred, average='macro'),
178         'f1_weighted': f1_score(y_true, y_pred, average='weighted')
179     }
180
181     # Per-class metrics

```

```

178     precision_per_class = precision_score(y_true, y_pred, average=None)
179     recall_per_class = recall_score(y_true, y_pred, average=None)
180
181     for i in range(len(precision_per_class)):
182         results[f'class_{i}_precision'] = precision_per_class[i]
183         results[f'class_{i}_recall'] = recall_per_class[i]
184
185     # ROC AUC for multi-class
186     if y_proba is not None:
187         results['auc_ovr'] = roc_auc_score(y_true, y_proba, multi_class='
188         results['auc_ovo'] = roc_auc_score(y_true, y_proba, multi_class='
189
190     return results
191
192 # Proper multi-class evaluation
193 y_pred_mc_correct = rf_multiclass.predict(X_multiclass)
194 y_proba_mc_correct = rf_multiclass.predict_proba(X_multiclass)
195
196 results_mc = evaluate_multiclass_model(
197     y_multiclass, y_pred_mc_correct, y_proba_mc_correct
198 )
199
200 print( Proper multi-class evaluation: )
201 for metric, value in results_mc.items():
202     if not metric.startswith('class_'): # Skip per-class for brevity
203         print(f {metric}: {value:.4f} )
204
205 print( \nPer-class performance: )
206 for i in range(3):
207     print(f Class {i}: Precision={results_mc[f'class_{i}_precision']:.4f
208           f Recall={results_mc[f'class_{i}_recall']:.4f} )
209
210 # CORRECT WAY 3 - Comprehensive evaluation with proper ROC AUC
211 print( \n=== CORRECT WAY 3: Proper ROC AUC === )
212
213 # Correct way to calculate ROC AUC
214 def proper_auc_calculation(model, X, y):
215     Proper ROC AUC calculation
216
217     # Get probability predictions
218     if hasattr(model, 'predict_proba'):
219         y_proba = model.predict_proba(X)
220         if y_proba.shape[1] == 2: # Binary classification
221             auc = roc_auc_score(y, y_proba[:, 1])
222         else: # Multi-class
223             auc = roc_auc_score(y, y_proba, multi_class='ovr')
224     else:
225         # For models without predict_proba, use decision function
226         if hasattr(model, 'decision_function'):
227             y_scores = model.decision_function(X)
228             auc = roc_auc_score(y, y_scores)

```

```

229         else:
230             raise ValueError( Model doesn't support probability or
                decision function )
231
232         return auc
233
234 # Calculate AUC properly
235 auc_balanced = proper_auc_calculation(rf_balanced, X_balanced, y_balanced)
236 auc_imbalanced = proper_auc_calculation(rf_imbalanced, X_imbalanced,
                y_imbalanced)
237
238 print( Proper AUC calculation: )
239 print(f Balanced dataset AUC: {auc_balanced:.4f} )
240 print(f Imbalanced dataset AUC: {auc_imbalanced:.4f} )
241
242 # CORRECT WAY 4 - Multi-metric evaluation with context
243 print( \n=== CORRECT WAY 4: Multi-metric Evaluation === )
244
245 def comprehensive_model_evaluation(models, X_train, y_train, X_test,
                y_test):
246     Comprehensive model evaluation with multiple metrics
247
248     results = {}
249
250     for name, model in models.items():
251         # Train model
252         model.fit(X_train, y_train)
253
254         # Make predictions
255         y_pred = model.predict(X_test)
256         if hasattr(model, 'predict_proba'):
257             y_proba = model.predict_proba(X_test)
258         else:
259             y_proba = None
260
261         # Calculate multiple metrics
262         metrics = {
263             'accuracy': accuracy_score(y_test, y_pred),
264             'precision_weighted': precision_score(y_test, y_pred, average=
'weighted'),
265             'recall_weighted': recall_score(y_test, y_pred, average='
weighted'),
266             'f1_weighted': f1_score(y_test, y_pred, average='weighted'),
267             'confusion_matrix': confusion_matrix(y_test, y_pred)
268         }
269
270         # Add AUC if probabilities available
271         if y_proba is not None:
272             if y_proba.shape[1] == 2: # Binary
273                 metrics['auc'] = roc_auc_score(y_test, y_proba[:, 1])
274             else: # Multi-class
275                 metrics['auc'] = roc_auc_score(y_test, y_proba,
multi_class='ovr')
276

```

```

277         results[name] = metrics
278
279     return results
280
281 # Compare models properly
282 X_train, X_test, y_train, y_test = train_test_split(
283     X_balanced, y_balanced, test_size=0.2, random_state=42
284 )
285
286 models_to_compare = {
287     'RandomForest': RandomForestClassifier(random_state=42),
288     'LogisticRegression': LogisticRegression(random_state=42, max_iter
289 =1000)
290 }
291 comprehensive_results = comprehensive_model_evaluation(
292     models_to_compare, X_train, y_train, X_test, y_test
293 )
294
295 print( Comprehensive model comparison: )
296 for model_name, metrics in comprehensive_results.items():
297     print(f \n{model_name}: )
298     print(f     Accuracy: {metrics['accuracy']:.4f} )
299     print(f     Precision: {metrics['precision_weighted']:.4f} )
300     print(f     Recall: {metrics['recall_weighted']:.4f} )
301     print(f     F1-Score: {metrics['f1_weighted']:.4f} )
302     if 'auc' in metrics:
303         print(f     AUC: {metrics['auc']:.4f} )
304     print(f     Confusion Matrix:\n{metrics['confusion_matrix']} )
305
306 # Advanced evaluation techniques
307 class AdvancedEvaluator:
308     Advanced model evaluation techniques
309
310     @staticmethod
311     def cost_sensitive_evaluation(y_true, y_pred, cost_matrix):
312         Evaluate model with custom cost matrix
313         cm = confusion_matrix(y_true, y_pred)
314         total_cost = np.sum(cm * cost_matrix)
315         return total_cost
316
317     @staticmethod
318     def bootstrap_confidence_intervals(y_true, y_pred, metric_func,
319 n_bootstrap=1000, confidence=0.95):
320         Calculate confidence intervals using bootstrap
321         bootstrap_scores = []
322
323         for _ in range(n_bootstrap):
324             indices = np.random.choice(len(y_true), len(y_true), replace=
325 True)
326             score = metric_func(y_true[indices], y_pred[indices])
327             bootstrap_scores.append(score)
328
329         alpha = 1 - confidence

```

```

328     lower_percentile = (alpha/2) * 100
329     upper_percentile = (1 - alpha/2) * 100
330
331     lower_bound = np.percentile(bootstrap_scores, lower_percentile)
332     upper_bound = np.percentile(bootstrap_scores, upper_percentile)
333
334     return {
335         'mean': np.mean(bootstrap_scores),
336         'std': np.std(bootstrap_scores),
337         'confidence_interval': (lower_bound, upper_bound)
338     }
339
340     @staticmethod
341     def evaluate_model_stability(model, X, y, n_bootstrap=100):
342         Evaluate model stability using bootstrap
343         scores = []
344
345         for _ in range(n_bootstrap):
346             indices = np.random.choice(len(X), len(X), replace=True)
347             X_boot, y_boot = X[indices], y[indices]
348
349             model.fit(X_boot, y_boot)
350             score = model.score(X, y)
351             scores.append(score)
352
353         return {
354             'mean_score': np.mean(scores),
355             'std_score': np.std(scores),
356             'stability': 1 - np.std(scores) / np.mean(scores) # Higher is
357             better
358         }
359
360 # Example of advanced evaluation
361 print( \n=== ADVANCED: Bootstrap Confidence Intervals === )
362
363 # Calculate confidence intervals for accuracy
364 def accuracy_metric(y_true, y_pred):
365     return accuracy_score(y_true, y_pred)
366
367 bootstrap_results = AdvancedEvaluator.bootstrap_confidence_intervals(
368     y_test, rf_balanced.predict(X_test), accuracy_metric
369 )
370
371 print( Bootstrap confidence intervals for accuracy: )
372 print(f Mean: {bootstrap_results['mean']:.4f} )
373 print(f Std: {bootstrap_results['std']:.4f} )
374 print(f 95% CI: ({bootstrap_results['confidence_interval'][0]:.4f},
375     f {bootstrap_results['confidence_interval'][1]:.4f}) )
376
377 # Example of cost-sensitive evaluation
378 print( \n=== ADVANCED: Cost-Sensitive Evaluation === )
379
380 # Define cost matrix (higher cost for false negatives in medical diagnosis
381 )

```

```

380 # [TN, FP]
381 # [FN, TP]
382 cost_matrix = np.array([
383     [0, 1],    # TN cost = 0, FP cost = 1
384     [10, 0]    # FN cost = 10, TP cost = 0
385 ])
386
387 cost_sensitive_score = AdvancedEvaluator.cost_sensitive_evaluation(
388     y_test, rf_balanced.predict(X_test), cost_matrix
389 )
390
391 print(f Cost-sensitive evaluation score: {cost_sensitive_score} )
392
393 # Business impact metrics
394 def calculate_business_metrics(y_true, y_pred, value_per_tp=100,
395     cost_per_fp=10):
396     Calculate business impact metrics
397     cm = confusion_matrix(y_true, y_pred)
398     tn, fp, fn, tp = cm.ravel()
399
400     revenue = tp * value_per_tp
401     cost = fp * cost_per_fp
402     net_benefit = revenue - cost
403
404     return {
405         'true_positives': tp,
406         'false_positives': fp,
407         'revenue': revenue,
408         'cost': cost,
409         'net_benefit': net_benefit
410     }
411
412 business_metrics = calculate_business_metrics(y_test, rf_balanced.predict(
413     X_test))
414 print( \nBusiness impact metrics: )
415 for metric, value in business_metrics.items():
416     print(f {metric}: {value} )

```

Listing 39: Error 5.2.4: Evaluation Metric Mistakes - Buggy Code

```

1 # Buggy: Common overfitting and underfitting issues
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split, validation_curve,
5     learning_curve
6 from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor
7 from sklearn.linear_model import LinearRegression, Ridge, Lasso
8 from sklearn.svm import SVR
9 from sklearn.preprocessing import PolynomialFeatures, StandardScaler
10 from sklearn.metrics import mean_squared_error, accuracy_score
11 from sklearn.datasets import make_classification, make_regression
12 import matplotlib.pyplot as plt
13
14 # Create datasets to demonstrate overfitting/underfitting
15 np.random.seed(42)

```

```

15
16 # 1. Simple linear relationship (for underfitting demo)
17 X_simple, y_simple = make_regression(n_samples=100, n_features=1, noise
    =10, random_state=42)
18
19 # 2. Complex non-linear relationship (for overfitting demo)
20 X_complex = np.linspace(0, 2*np.pi, 100).reshape(-1, 1)
21 y_complex = np.sin(X_complex).ravel() + np.random.normal(0, 0.1, 100)
22
23 # 3. High-dimensional dataset (for overfitting demo)
24 X_high_dim, y_high_dim = make_classification(
25     n_samples=100, n_features=50, n_informative=5,
26     n_redundant=10, random_state=42
27 )
28
29 print( Dataset shapes: )
30 print(f Simple: {X_simple.shape} )
31 print(f Complex: {X_complex.shape} )
32 print(f High-dimensional: {X_high_dim.shape} )
33
34 # WRONG WAY 1 - Underfitting with too simple model
35 print( \n=== WRONG WAY 1: Underfitting with Too Simple Model === )
36
37 # WRONG: Using linear model for non-linear data
38 X_train_complex, X_test_complex, y_train_complex, y_test_complex =
    train_test_split(
39     X_complex, y_complex, test_size=0.2, random_state=42
40 )
41
42 # Linear model on non-linear data (underfitting)
43 linear_model = LinearRegression()
44 linear_model.fit(X_train_complex, y_train_complex)
45
46 train_score_linear = linear_model.score(X_train_complex, y_train_complex)
47 test_score_linear = linear_model.score(X_test_complex, y_test_complex)
48
49 print( Linear model on non-linear data: )
50 print(f Training score: {train_score_linear:.4f} )
51 print(f Test score: {test_score_linear:.4f} )
52 print( Problem: High bias, underfitting - model too simple! )
53
54 # WRONG WAY 2 - Overfitting with too complex model
55 print( \n=== WRONG WAY 2: Overfitting with Too Complex Model === )
56
57 # WRONG: High-degree polynomial on small dataset
58 poly_features_overfit = PolynomialFeatures(degree=15) # Too high!
59 X_poly_overfit = poly_features_overfit.fit_transform(X_train_complex)
60
61 # Complex model on small dataset (overfitting)
62 complex_model = LinearRegression()
63 complex_model.fit(X_poly_overfit, y_train_complex)
64
65 # Training score will be very high
66 train_pred_overfit = complex_model.predict(X_poly_overfit)

```

```

67 train_mse_overfit = mean_squared_error(y_train_complex, train_pred_overfit
68 )
69 # Test score will be poor
70 X_test_poly_overfit = poly_features_overfit.transform(X_test_complex)
71 test_pred_overfit = complex_model.predict(X_test_poly_overfit)
72 test_mse_overfit = mean_squared_error(y_test_complex, test_pred_overfit)
73
74 print( High-degree polynomial on small dataset: )
75 print(f Training MSE: {train_mse_overfit:.4f} )
76 print(f Test MSE: {test_mse_overfit:.4f} )
77 print( Problem: High variance, overfitting - model too complex! )
78
79 # WRONG WAY 3 - No regularization on high-dimensional data
80 print( \n=== WRONG WAY 3: No Regularization on High-dimensional Data === )
81
82 X_train_hd, X_test_hd, y_train_hd, y_test_hd = train_test_split(
83     X_high_dim, y_high_dim, test_size=0.2, random_state=42
84 )
85
86 # WRONG: No regularization on high-dimensional data
87 rf_no_reg = RandomForestClassifier(random_state=42, min_samples_split=2,
88     min_samples_leaf=1)
89 rf_no_reg.fit(X_train_hd, y_train_hd)
90
91 train_score_no_reg = rf_no_reg.score(X_train_hd, y_train_hd)
92 test_score_no_reg = rf_no_reg.score(X_test_hd, y_test_hd)
93
94 print( Random Forest without regularization on high-dimensional data: )
95 print(f Training accuracy: {train_score_no_reg:.4f} )
96 print(f Test accuracy: {test_score_no_reg:.4f} )
97 print( Problem: Overfitting due to high variance in high-dimensional space
98 ! )
99
100 # WRONG WAY 4 - Ignoring learning curves
101 print( \n=== WRONG WAY 4: Ignoring Learning Curves === )
102
103 # WRONG: Not checking if more data would help
104 def quick_model_evaluation(X, y):
105     Quick evaluation without learning curve analysis
106     model = RandomForestClassifier(random_state=42)
107     model.fit(X, y)
108     return model.score(X, y) # WRONG: Only training score
109
110 quick_score = quick_model_evaluation(X_high_dim, y_high_dim)
111 print(f Quick evaluation score: {quick_score:.4f} )
112 print( Problem: No insight into bias-variance tradeoff! )
113
114 # CORRECT WAY 1 - Diagnosing bias-variance with validation curves
115 print( \n=== CORRECT WAY 1: Validation Curves for Diagnosis === )
116
117 def plot_validation_curve_analysis(X, y, model, param_name, param_range,
118     cv=5):
119     Plot validation curves to diagnose bias-variance

```



```

117
118     train_scores, val_scores = validation_curve(
119         model, X, y, param_name=param_name, param_range=param_range, cv=cv
120     )
121
122     train_mean = np.mean(train_scores, axis=1)
123     train_std = np.std(train_scores, axis=1)
124     val_mean = np.mean(val_scores, axis=1)
125     val_std = np.std(val_scores, axis=1)
126
127     plt.figure(figsize=(10, 6))
128     plt.plot(param_range, train_mean, 'o-', color='blue', label='Training
129         score')
129     plt.fill_between(param_range, train_mean - train_std, train_mean +
130         train_std, alpha=0.1, color='blue')
131
131     plt.plot(param_range, val_mean, 'o-', color='red', label='Cross-
132         validation score')
132     plt.fill_between(param_range, val_mean - val_std, val_mean + val_std,
133         alpha=0.1, color='red')
134
134     plt.xlabel(param_name)
135     plt.ylabel('Score')
136     plt.title(f'Validation Curve for {param_name}')
137     plt.legend()
138     plt.grid(True)
139     plt.show()
140
141     # Diagnosis
142     final_train = train_mean[-1]
143     final_val = val_mean[-1]
144
145     if final_train < 0.8:
146         print( Diagnosis: High bias (underfitting) - Model is too simple )
147     elif final_train - final_val > 0.1:
148         print( Diagnosis: High variance (overfitting) - Model is too
149             complex )
150     else:
151         print( Diagnosis: Good fit - Balanced bias and variance )
152
152     return final_train, final_val
153
154 # Analyze polynomial degree complexity
155 print( Analyzing polynomial degree complexity: )
156 poly_model = LinearRegression()
157 degrees = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
158 train_scores_poly, val_scores_poly = [], []
159
160 for degree in degrees:
161     poly_features = PolynomialFeatures(degree=degree)
162     X_poly = poly_features.fit_transform(X_train_complex)
163
164     # Cross-validation score
165     from sklearn.model_selection import cross_val_score

```

```

166     cv_scores = cross_val_score(LinearRegression(), X_poly,
167     y_train_complex, cv=5)
168     val_scores_poly.append(cv_scores.mean())
169
170     # Training score
171     model = LinearRegression()
172     model.fit(X_poly, y_train_complex)
173     train_scores_poly.append(model.score(X_poly, y_train_complex))
174
175 # Plot results
176 plt.figure(figsize=(10, 6))
177 plt.plot(degrees, train_scores_poly, 'o-', label='Training Score')
178 plt.plot(degrees, val_scores_poly, 'o-', label='Validation Score')
179 plt.xlabel('Polynomial Degree')
180 plt.ylabel('Score')
181 plt.title('Polynomial Degree vs Model Performance')
182 plt.legend()
183 plt.grid(True)
184 plt.show()
185
186 # CORRECT WAY 2 - Proper regularization
187 print( \n=== CORRECT WAY 2: Proper Regularization === )
188
189 # Ridge regression with proper regularization
190 def find_optimal_regularization(X, y, alpha_range=np.logspace(-4, 4, 50)):
191     Find optimal regularization strength
192
193     train_scores = []
194     val_scores = []
195
196     for alpha in alpha_range:
197         ridge = Ridge(alpha=alpha)
198
199         # Cross-validation
200         cv_scores = cross_val_score(ridge, X, y, cv=5)
201         val_scores.append(cv_scores.mean())
202
203         # Training score
204         ridge.fit(X, y)
205         train_scores.append(ridge.score(X, y))
206
207     optimal_idx = np.argmax(val_scores)
208     optimal_alpha = alpha_range[optimal_idx]
209
210     return optimal_alpha, train_scores, val_scores, alpha_range
211
212 # Find optimal regularization for high-dimensional data
213 optimal_alpha, train_reg, val_reg, alphas = find_optimal_regularization(
214     X_train_hd, y_train_hd)
215
216 print(f Optimal regularization strength: {optimal_alpha:.6f} )
217
218 # Compare with no regularization
219 ridge_optimal = Ridge(alpha=optimal_alpha)

```

```

218 ridge_optimal.fit(X_train_hd, y_train_hd)
219
220 rf_regularized = RandomForestClassifier(
221     random_state=42,
222     min_samples_split=10, # Increase minimum samples
223     min_samples_leaf=5,   # Increase minimum leaf samples
224     max_depth=10          # Limit depth
225 )
226 rf_regularized.fit(X_train_hd, y_train_hd)
227
228 print( Regularized models performance: )
229 print(f    Ridge (   {ridge_optimal.alpha:.4f}) - Train: {ridge_optimal.score(
230     X_train_hd, y_train_hd):.4f},
231     f Test: {ridge_optimal.score(X_test_hd, y_test_hd):.4f} )
232
233 print(f    Random Forest (regularized) - Train: {rf_regularized.score(
234     X_train_hd, y_train_hd):.4f},
235     f Test: {rf_regularized.score(X_test_hd, y_test_hd):.4f} )
236
237 # CORRECT WAY 3 - Learning curve analysis
238 print( \n=== CORRECT WAY 3: Learning Curve Analysis === )
239
240 def plot_learning_curve_analysis(estimator, X, y, cv=5, train_sizes=np.
241     linspace(0.1, 1.0, 10)):
242     Plot learning curves to understand bias-variance tradeoff
243
244     train_sizes, train_scores, val_scores = learning_curve(
245         estimator, X, y, cv=cv, train_sizes=train_sizes, n_jobs=-1
246     )
247
248     train_mean = np.mean(train_scores, axis=1)
249     train_std = np.std(train_scores, axis=1)
250     val_mean = np.mean(val_scores, axis=1)
251     val_std = np.std(val_scores, axis=1)
252
253     plt.figure(figsize=(10, 6))
254     plt.plot(train_sizes, train_mean, 'o-', color='blue', label='Training
255     score')
256     plt.fill_between(train_sizes, train_mean - train_std, train_mean +
257     train_std, alpha=0.1, color='blue')
258
259     plt.plot(train_sizes, val_mean, 'o-', color='red', label='Cross-
260     validation score')
261     plt.fill_between(train_sizes, val_mean - val_std, val_mean + val_std,
262     alpha=0.1, color='red')
263
264     plt.xlabel('Training Set Size')
265     plt.ylabel('Score')
266     plt.title('Learning Curve')
267     plt.legend()
268     plt.grid(True)
269     plt.show()
270
271 # Diagnosis based on learning curves

```

```

265     final_train = train_mean[-1]
266     final_val = val_mean[-1]
267     gap = final_train - final_val
268
269     print( Learning Curve Analysis: )
270     print(f    Final training score: {final_train:.4f} )
271     print(f    Final validation score: {final_val:.4f} )
272     print(f    Gap: {gap:.4f} )
273
274     if final_train < 0.8:
275         print(        High bias problem (underfitting) )
276         print(        Consider more complex model or feature engineering )
277     elif gap > 0.1:
278         print(        High variance problem (overfitting) )
279         print(        Consider more data, regularization, or simpler model
280     )
281     else:
282         print(        Good bias-variance balance )
283
284     # Data sufficiency analysis
285     if final_val < 0.8 and val_mean[-1] - val_mean[-3] > 0.01:
286         print(        More data might improve performance )
287     elif final_val < 0.8:
288         print(        More data unlikely to help significantly )
289
290     return train_sizes, train_mean, val_mean
291
292 # Analyze learning curves for different models
293 print( Learning curve analysis for Random Forest: )
294 plot_learning_curve_analysis(
295     RandomForestClassifier(random_state=42, n_estimators=100),
296     X_high_dim, y_high_dim
297 )
298
299 # CORRECT WAY 4 - Early stopping and model complexity control
300 print( \n=== CORRECT WAY 4: Early Stopping and Complexity Control === )
301
302 class EarlyStoppingModel:
303     Model with early stopping capability
304
305     def __init__(self, base_model, patience=10, min_delta=1e-4):
306         self.base_model = base_model
307         self.patience = patience
308         self.min_delta = min_delta
309         self.best_score = None
310         self.wait = 0
311         self.stopped_epoch = 0
312
313     def fit_with_early_stopping(self, X, y, X_val, y_val, max_epochs=1000)
314     :
315         Fit model with early stopping
316
317         scores = []

```

```

317         for epoch in range(max_epochs):
318             # Fit model for this epoch (simplified for demonstration)
319             self.base_model.fit(X, y)
320
321             # Evaluate on validation set
322             score = self.base_model.score(X_val, y_val)
323             scores.append(score)
324
325             # Early stopping logic
326             if self.best_score is None:
327                 self.best_score = score
328             elif score - self.best_score < self.min_delta:
329                 self.wait += 1
330                 if self.wait >= self.patience:
331                     self.stopped_epoch = epoch
332                     print(f Early stopping at epoch {epoch} )
333                     break
334             else:
335                 self.best_score = score
336                 self.wait = 0
337
338         return scores
339
340 # Demonstrate early stopping concept
341 from sklearn.ensemble import GradientBoostingClassifier
342
343 # Split data for validation
344 X_train_es, X_val_es, y_train_es, y_val_es = train_test_split(
345     X_high_dim, y_high_dim, test_size=0.2, random_state=42
346 )
347
348 # Model with potential overfitting
349 gb_model = GradientBoostingClassifier(
350     n_estimators=1000, # Many estimators
351     learning_rate=0.1,
352     max_depth=3,
353     random_state=42
354 )
355
356 # Early stopping approach
357 gb_model_early = GradientBoostingClassifier(
358     n_estimators=1000,
359     learning_rate=0.1,
360     max_depth=3,
361     validation_fraction=0.1, # Use 10% for validation
362     n_iter_no_change=10,     # Patience
363     tol=1e-4,                # Minimum improvement
364     random_state=42
365 )
366
367 # Fit both models
368 gb_model.fit(X_train_es, y_train_es)
369 gb_model_early.fit(X_train_es, y_train_es)
370

```

```

371 print( Early stopping comparison: )
372 print(f    Regular GB - Train: {gb_model.score(X_train_es, y_train_es):.4f
    },
373       f Test: {gb_model.score(X_val_es, y_val_es):.4f} )
374
375 print(f    Early stopping GB - Train: {gb_model_early.score(X_train_es,
    y_train_es):.4f},
376       f Test: {gb_model_early.score(X_val_es, y_val_es):.4f} )
377
378 print(f    Regular GB estimators used: {gb_model.n_estimators} )
379 print(f    Early stopping GB estimators used: {gb_model_early.n_estimators_
    } )
380
381 # Advanced bias-variance analysis
382 class BiasVarianceAnalyzer:
383     Advanced bias-variance decomposition analysis
384
385     @staticmethod
386     def bias_variance_decomposition(model, X, y, n_bootstrap=100):
387         Decompose error into bias, variance, and noise components
388
389         predictions = []
390
391         # Bootstrap sampling
392         for _ in range(n_bootstrap):
393             indices = np.random.choice(len(X), len(X), replace=True)
394             X_boot, y_boot = X[indices], y[indices]
395
396             model.fit(X_boot, y_boot)
397             pred = model.predict(X)
398             predictions.append(pred)
399
400         predictions = np.array(predictions)
401
402         # Calculate components
403         avg_prediction = np.mean(predictions, axis=0)
404         bias_squared = np.mean((avg_prediction - y) ** 2)
405         variance = np.mean(np.var(predictions, axis=0))
406         noise = np.mean(np.var(predictions, axis=0)) # Simplified
407
408         return {
409             'bias_squared': bias_squared,
410             'variance': variance,
411             'noise': noise,
412             'total_error': bias_squared + variance + noise
413         }
414
415     @staticmethod
416     def model_complexity_analysis(X, y, model_class, complexity_param,
    param_values, cv=5):
417         Analyze how model complexity affects bias-variance tradeoff
418
419         results = {
420             'param_values': param_values,

```

```

421         'bias_squared': [],
422         'variance': [],
423         'total_error': []
424     }
425
426     for param_value in param_values:
427         model = model_class(**{complexity_param: param_value})
428
429         # Cross-validation for variance estimation
430         from sklearn.model_selection import cross_val_predict
431         cv_predictions = cross_val_predict(model, X, y, cv=cv)
432
433         # Bias calculation
434         model.fit(X, y)
435         train_predictions = model.predict(X)
436
437         bias_sq = np.mean((np.mean(cv_predictions) - y) ** 2)
438         var = np.var(cv_predictions)
439
440         results['bias_squared'].append(bias_sq)
441         results['variance'].append(var)
442         results['total_error'].append(bias_sq + var)
443
444     return results
445
446 # Example of bias-variance analysis
447 print( \n=== ADVANCED: Bias-Variance Analysis === )
448
449 # This would perform detailed bias-variance decomposition
450 # bias_var_results = BiasVarianceAnalyzer.bias_variance_decomposition(
451 #     RandomForestClassifier(random_state=42), X_high_dim, y_high_dim
452 # )
453
454 print( Bias-variance analysis helps understand the sources of prediction
         error! )
455
456 # Practical overfitting prevention strategies
457 def prevent_overfitting_checklist():
458     Checklist for preventing overfitting
459
460     checklist = [
461         1. Use cross-validation for model evaluation ,
462         2. Split data properly (train/validation/test) ,
463         3. Apply appropriate regularization ,
464         4. Use early stopping for iterative algorithms ,
465         5. Limit model complexity (depth, features, etc.) ,
466         6. Increase minimum samples for splits/leaves ,
467         7. Use ensemble methods to reduce variance ,
468         8. Analyze learning curves to understand data needs ,
469         9. Check validation curves for optimal hyperparameters ,
470         10. Monitor training vs validation performance gap
471     ]
472
473     print( Overfitting Prevention Checklist: )

```

```

474     for item in checklist:
475         print(f    {item} )
476
477 prevent_overfitting_checklist()
478
479 # Model selection based on bias-variance tradeoff
480 def select_model_by_bias_variance(X, y, models_dict):
481     Select model based on bias-variance analysis
482
483     results = {}
484
485     for name, model in models_dict.items():
486         # Cross-validation scores
487         cv_scores = cross_val_score(model, X, y, cv=5)
488
489         # Training score
490         model.fit(X, y)
491         train_score = model.score(X, y)
492
493         # Bias-variance indicators
494         bias_indicator = 1 - train_score # Higher = more bias
495         variance_indicator = train_score - cv_scores.mean() # Higher =
more variance
496
497         results[name] = {
498             'train_score': train_score,
499             'cv_mean': cv_scores.mean(),
500             'cv_std': cv_scores.std(),
501             'bias_indicator': bias_indicator,
502             'variance_indicator': variance_indicator
503         }
504
505     return results
506
507 # Example model selection
508 models_for_analysis = {
509     'Simple Model': LogisticRegression(random_state=42),
510     'Moderate Model': RandomForestClassifier(n_estimators=100,
random_state=42),
511     'Complex Model': RandomForestClassifier(n_estimators=500, max_depth
=20, random_state=42)
512 }
513
514 analysis_results = select_model_by_bias_variance(X_train_hd, y_train_hd,
models_for_analysis)
515
516 print( \nModel Selection Based on Bias-Variance Tradeoff: )
517 for name, metrics in analysis_results.items():
518     print(f \n{name}: )
519     print(f    Train Score: {metrics['train_score']:.4f} )
520     print(f    CV Score: {metrics['cv_mean']:.4f} (+/- {metrics['cv_std'] *
2:.4f}) )
521     print(f    Bias Indicator: {metrics['bias_indicator']:.4f} )

```



```
522 print(f Variance Indicator: {metrics['variance_indicator']:.4f} )
```

Listing 40: Error 5.2.5: Overfitting and Underfitting Issues - Buggy Code

### 5.3 5.3 Deployment and Production Errors

```
1 # Buggy: Common model versioning issues
2 import joblib
3 import pickle
4 import os
5 import numpy as np
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.model_selection import train_test_split
8 from sklearn.datasets import make_classification
9
10 # Create sample data and model
11 X, y = make_classification(n_samples=1000, n_features=10, random_state=42)
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
13     random_state=42)
14
15 # Train a model
16 model_v1 = RandomForestClassifier(n_estimators=100, random_state=42)
17 model_v1.fit(X_train, y_train)
18
19 print( Model V1 trained )
20 print(f Model V1 accuracy: {model_v1.score(X_test, y_test):.4f} )
21
22 # WRONG WAY 1 - No versioning at all
23 print( \n=== WRONG WAY 1: No Versioning === )
24
25 # WRONG: Overwriting model file without version control
26 joblib.dump(model_v1, 'model.pkl') # No version information
27 print( Model saved without versioning )
28
29 # Train improved model
30 model_v2 = RandomForestClassifier(n_estimators=200, random_state=42) #
31     Better parameters
32 model_v2.fit(X_train, y_train)
33
34 # WRONG: Overwrite previous model
35 joblib.dump(model_v2, 'model.pkl') # Previous model lost!
36 print( Model V2 saved, V1 overwritten )
37
38 # Problem: No way to rollback or compare versions
39 print( Problem: Previous model version lost forever! )
40
41 # WRONG WAY 2 - Inconsistent naming
42 print( \n=== WRONG WAY 2: Inconsistent Naming === )
43
44 # WRONG: Inconsistent version naming
45 joblib.dump(model_v1, 'my_model_v1.pkl')
46 joblib.dump(model_v2, 'model_version_2.pkl') # Different naming
47     convention
```

```

45 joblib.dump(model_v1, 'final_model.pkl')          # No version info
46
47 print( Models saved with inconsistent naming )
48 print( Problem: Hard to track and manage versions! )
49
50 # WRONG WAY 3 - No metadata tracking
51 print( \n=== WRONG WAY 3: No Metadata Tracking === )
52
53 # WRONG: Only saving model, no metadata
54 joblib.dump(model_v2, 'model_v2.pkl')
55
56 # Later, when loading...
57 loaded_model = joblib.load('model_v2.pkl')
58 print( Model loaded, but no information about: )
59 print( - Training parameters )
60 print( - Performance metrics )
61 print( - Training data characteristics )
62 print( - Creation date )
63 print( Problem: No context about the model! )
64
65 # WRONG WAY 4 - Manual version management
66 print( \n=== WRONG WAY 4: Manual Version Management === )
67
68 # WRONG: Manual file management
69 def save_model_manual(model, version):
70     Manual model saving - error prone
71     filename = f'model_v{version}.pkl'
72     joblib.dump(model, filename)
73     print(f Model saved as {filename} )
74
75 save_model_manual(model_v1, 1)
76 save_model_manual(model_v2, 2)
77
78 # Manual rollback - error prone
79 def load_model_manual(version):
80     Manual model loading - error prone
81     filename = f'model_v{version}.pkl'
82     if os.path.exists(filename):
83         return joblib.load(filename)
84     else:
85         print(f Model version {version} not found )
86         return None
87
88 # This approach is fragile and error-prone
89 print( Problem: Manual management is error-prone and inconsistent! )
90
91 # CORRECT WAY 1 - Systematic versioning
92 print( \n=== CORRECT WAY 1: Systematic Versioning === )
93
94 import datetime
95 import json
96
97 class ModelVersionManager:
98     Systematic model version management

```

```

99
100     def __init__(self, model_dir='models'):
101         self.model_dir = model_dir
102         os.makedirs(model_dir, exist_ok=True)
103         self.versions_file = os.path.join(model_dir, 'versions.json')
104
105         # Initialize versions file if it doesn't exist
106         if not os.path.exists(self.versions_file):
107             with open(self.versions_file, 'w') as f:
108                 json.dump([], f)
109
110     def save_model(self, model, version_name, metadata=None):
111         Save model with version tracking
112
113         # Create version identifier
114         timestamp = datetime.datetime.now().strftime('%Y%m%d_%H%M%S')
115         version_id = f {version_name}_{timestamp}
116
117         # Save model
118         model_filename = f {version_id}.pkl
119         model_path = os.path.join(self.model_dir, model_filename)
120         joblib.dump(model, model_path)
121
122         # Create metadata
123         version_info = {
124             'version_id': version_id,
125             'version_name': version_name,
126             'timestamp': timestamp,
127             'model_path': model_filename,
128             'metadata': metadata or {}
129         }
130
131         # Update versions file
132         with open(self.versions_file, 'r') as f:
133             versions = json.load(f)
134
135         versions.append(version_info)
136
137         with open(self.versions_file, 'w') as f:
138             json.dump(versions, f, indent=2)
139
140         print(f Model saved: {version_id} )
141         return version_id
142
143     def load_model(self, version_id=None, version_name=None):
144         Load specific model version
145
146         with open(self.versions_file, 'r') as f:
147             versions = json.load(f)
148
149         # Find version
150         if version_id:
151             version_info = next((v for v in versions if v['version_id'] ==
version_id), None)

```

```

152         elif version_name:
153             # Get latest version with this name
154             matching_versions = [v for v in versions if v['version_name']
== version_name]
155             version_info = max(matching_versions, key=lambda x: x['
timestamp']) if matching_versions else None
156         else:
157             # Get latest version
158             version_info = max(versions, key=lambda x: x['timestamp']) if
versions else None
159
160         if not version_info:
161             raise ValueError( Model version not found )
162
163         model_path = os.path.join(self.model_dir, version_info['model_path
'])
164         model = joblib.load(model_path)
165
166         return model, version_info
167
168     def list_versions(self):
169         List all model versions
170         with open(self.versions_file, 'r') as f:
171             versions = json.load(f)
172         return versions
173
174 # Use proper version management
175 version_manager = ModelVersionManager()
176
177 # Save models with proper versioning
178 metadata_v1 = {
179     'n_estimators': 100,
180     'accuracy': model_v1.score(X_test, y_test),
181     'training_samples': len(X_train)
182 }
183
184 metadata_v2 = {
185     'n_estimators': 200,
186     'accuracy': model_v2.score(X_test, y_test),
187     'training_samples': len(X_train)
188 }
189
190 version_id_1 = version_manager.save_model(model_v1, 'random_forest',
metadata_v1)
191 version_id_2 = version_manager.save_model(model_v2, 'random_forest',
metadata_v2)
192
193 print(f Saved versions: {version_id_1}, {version_id_2} )
194
195 # List all versions
196 print( \nAll model versions: )
197 versions = version_manager.list_versions()
198 for version in versions:
199     print(f {version['version_id']}: {version['metadata']} )

```

```

200
201 # Load specific version
202 loaded_model_v1, info_v1 = version_manager.load_model(version_id=
    version_id_1)
203 print(f \nLoaded model {info_v1['version_id']} )
204 print(f Accuracy: {loaded_model_v1.score(X_test, y_test):.4f} )
205
206 # CORRECT WAY 2 - Semantic versioning
207 print( \n=== CORRECT WAY 2: Semantic Versioning === )
208
209 class SemanticModelVersionManager:
210     Model version management with semantic versioning
211
212     def __init__(self, model_dir='models_semantic'):
213         self.model_dir = model_dir
214         os.makedirs(model_dir, exist_ok=True)
215         self.versions_file = os.path.join(model_dir, 'semantic_versions.
    json')
216
217         if not os.path.exists(self.versions_file):
218             # Initialize with version 1.0.0
219             initial_info = {
220                 'latest_major': 1,
221                 'latest_minor': 0,
222                 'latest_patch': 0,
223                 'versions': []
224             }
225             with open(self.versions_file, 'w') as f:
226                 json.dump(initial_info, f)
227
228         def _get_next_version(self, version_type='patch'):
229             Get next version number
230             with open(self.versions_file, 'r') as f:
231                 version_info = json.load(f)
232
233             major = version_info['latest_major']
234             minor = version_info['latest_minor']
235             patch = version_info['latest_patch']
236
237             if version_type == 'major':
238                 major += 1
239                 minor = 0
240                 patch = 0
241             elif version_type == 'minor':
242                 minor += 1
243                 patch = 0
244             elif version_type == 'patch':
245                 patch += 1
246
247             version_info['latest_major'] = major
248             version_info['latest_minor'] = minor
249             version_info['latest_patch'] = patch
250
251             return f {major}.{minor}.{patch} , version_info

```

```

252
253     def save_model_semantic(self, model, change_type='patch', description=
    ''):
254         Save model with semantic versioning
255
256         version, version_info = self._get_next_version(change_type)
257         timestamp = datetime.datetime.now().strftime('%Y%m%d_%H%M%S')
258         version_id = f v{version}_{timestamp}
259
260         # Save model
261         model_filename = f {version_id}.pkl
262         model_path = os.path.join(self.model_dir, model_filename)
263         joblib.dump(model, model_path)
264
265         # Create version record
266         version_record = {
267             'version': version,
268             'version_id': version_id,
269             'timestamp': timestamp,
270             'model_path': model_filename,
271             'change_type': change_type,
272             'description': description
273         }
274
275         version_info['versions'].append(version_record)
276
277         with open(self.versions_file, 'w') as f:
278             json.dump(version_info, f, indent=2)
279
280         print(f Model saved with semantic version: {version} )
281         return version_id, version
282
283     def get_version_history(self):
284         Get version history
285         with open(self.versions_file, 'r') as f:
286             version_info = json.load(f)
287             return version_info['versions']
288
289 # Use semantic versioning
290 semantic_manager = SemanticModelVersionManager()
291
292 # Save models with semantic versioning
293 version_id_1, sem_version_1 = semantic_manager.save_model_semantic(
294     model_v1, 'major', 'Initial release with basic Random Forest'
295 )
296
297 version_id_2, sem_version_2 = semantic_manager.save_model_semantic(
298     model_v2, 'minor', 'Improved with more estimators'
299 )
300
301 print(f Semantic versions: {sem_version_1}, {sem_version_2} )
302
303 # View version history
304 print( \nSemantic version history: )

```

```

305 history = semantic_manager.get_version_history()
306 for record in history:
307     print(f    v{record['version']} ({record['change_type']}): {record['
description']} )
308
309 # CORRECT WAY 3 - Model registry with comprehensive metadata
310 print( \n=== CORRECT WAY 3: Comprehensive Model Registry === )
311
312 class ModelRegistry:
313     Comprehensive model registry with metadata
314
315     def __init__(self, registry_dir='model_registry'):
316         self.registry_dir = registry_dir
317         os.makedirs(registry_dir, exist_ok=True)
318         self.registry_file = os.path.join(registry_dir, 'registry.json')
319
320         if not os.path.exists(self.registry_file):
321             with open(self.registry_file, 'w') as f:
322                 json.dump({'models': {}}, f)
323
324     def register_model(self, model, model_name, version, metrics=None,
325                       training_data_info=None, hyperparameters=None,
326                       feature_names=None, target_name=None):
327         Register model with comprehensive metadata
328
329         timestamp = datetime.datetime.now().strftime('%Y%m%d_%H%M%S')
330         model_id = f {model_name}_v{version}_{timestamp}
331
332         # Save model
333         model_filename = f {model_id}.pkl
334         model_path = os.path.join(self.registry_dir, model_filename)
335         joblib.dump(model, model_path)
336
337         # Create comprehensive metadata
338         model_info = {
339             'model_id': model_id,
340             'model_name': model_name,
341             'version': version,
342             'timestamp': timestamp,
343             'model_path': model_filename,
344             'model_type': type(model).__name__,
345             'hyperparameters': hyperparameters or {},
346             'metrics': metrics or {},
347             'training_data_info': training_data_info or {},
348             'feature_names': feature_names or [],
349             'target_name': target_name,
350             'python_version': '.'.join(map(str, [x for x in __import__(
sys').version_info[:3]])),
351             'scikit_learn_version': __import__('sklearn').__version__
352         }
353
354         # Update registry
355         with open(self.registry_file, 'r') as f:
356             registry = json.load(f)

```

```

357
358     if model_name not in registry['models']:
359         registry['models'][model_name] = []
360
361     registry['models'][model_name].append(model_info)
362
363     with open(self.registry_file, 'w') as f:
364         json.dump(registry, f, indent=2)
365
366     print(f Model registered: {model_id} )
367     return model_id
368
369     def get_model(self, model_name, version=None):
370         Retrieve model by name and optional version
371         with open(self.registry_file, 'r') as f:
372             registry = json.load(f)
373
374         if model_name not in registry['models']:
375             raise ValueError(f Model {model_name} not found )
376
377         model_versions = registry['models'][model_name]
378
379         if version:
380             model_info = next((m for m in model_versions if m['version']
== version), None)
381         else:
382             # Get latest version
383             model_info = max(model_versions, key=lambda x: x['timestamp'])
384
385         if not model_info:
386             raise ValueError(f Model version not found )
387
388         model_path = os.path.join(self.registry_dir, model_info['
model_path'])
389         model = joblib.load(model_path)
390
391         return model, model_info
392
393     def compare_models(self, model_name):
394         Compare different versions of a model
395         with open(self.registry_file, 'r') as f:
396             registry = json.load(f)
397
398         if model_name not in registry['models']:
399             return None
400
401         versions = registry['models'][model_name]
402         return sorted(versions, key=lambda x: x['timestamp'])
403
404 # Use comprehensive model registry
405 registry = ModelRegistry()
406
407 # Register models with full metadata
408 model_v1_id = registry.register_model(

```



```

409     model_v1,
410     model_name='customer_churn_rf',
411     version='1.0.0',
412     metrics={'accuracy': model_v1.score(X_test, y_test)},
413     training_data_info={'samples': len(X_train), 'features': X_train.shape
414                          [1]},
415     hyperparameters={'n_estimators': 100, 'random_state': 42},
416     feature_names=[f'feature_{i}' for i in range(X.shape[1])],
417     target_name='churn'
418 )
419 model_v2_id = registry.register_model(
420     model_v2,
421     model_name='customer_churn_rf',
422     version='1.1.0',
423     metrics={'accuracy': model_v2.score(X_test, y_test)},
424     training_data_info={'samples': len(X_train), 'features': X_train.shape
425                          [1]},
426     hyperparameters={'n_estimators': 200, 'random_state': 42},
427     feature_names=[f'feature_{i}' for i in range(X.shape[1])],
428     target_name='churn'
429 )
430 print(f Registered models: {model_v1_id}, {model_v2_id} )
431
432 # Compare model versions
433 print( \nModel comparison: )
434 comparison = registry.compare_models('customer_churn_rf')
435 for model_info in comparison:
436     print(f v{model_info['version']}: Accuracy={model_info['metrics'].
437           get('accuracy', 'N/A')} )
438
439 # Retrieve specific model version
440 retrieved_model, model_info = registry.get_model('customer_churn_rf', '
441 1.0.0')
442 print(f \nRetrieved model: {model_info['model_id']} )
443 print(f Model type: {model_info['model_type']} )
444 print(f Hyperparameters: {model_info['hyperparameters']} )
445
446 # Advanced versioning with Git integration
447 class GitModelVersioning:
448     Model versioning with Git integration
449
450     def __init__(self, repo_path='.'):
451         try:
452             import git
453             self.repo = git.Repo(repo_path)
454         except ImportError:
455             print( GitPython not installed. Git integration disabled. )
456             self.repo = None
457         except Exception as e:
458             print(f Git repository not found: {e} )
459             self.repo = None

```

```

459     def get_git_info(self):
460         Get current Git information
461         if not self.repo:
462             return {}
463
464         try:
465             return {
466                 'commit_hash': self.repo.head.object.hexsha,
467                 'commit_message': self.repo.head.object.message.strip(),
468                 'branch': self.repo.active_branch.name,
469                 'timestamp': datetime.datetime.now().isoformat()
470             }
471         except Exception as e:
472             print(f Error getting Git info: {e} )
473             return {}
474
475     def save_model_with_git(self, model, model_name, description=''):
476         Save model with Git integration
477
478         git_info = self.get_git_info()
479         timestamp = datetime.datetime.now().strftime('%Y%m%d_%H%M%S')
480         model_id = f {model_name}_{timestamp}
481
482         # Save model
483         model_filename = f {model_id}.pkl
484         joblib.dump(model, model_filename)
485
486         # Create metadata with Git info
487         metadata = {
488             'model_id': model_id,
489             'timestamp': timestamp,
490             'model_name': model_name,
491             'description': description,
492             'git_info': git_info
493         }
494
495         metadata_filename = f {model_id}_metadata.json
496         with open(metadata_filename, 'w') as f:
497             json.dump(metadata, f, indent=2)
498
499         print(f Model saved with Git integration: {model_id} )
500         return model_id
501
502     # Example of Git-integrated versioning
503     print( \n=== ADVANCED: Git-Integrated Versioning === )
504
505     # This would work if GitPython is installed and in a Git repository
506     # git_versioning = GitModelVersioning()
507     # git_model_id = git_versioning.save_model_with_git(
508     #     model_v1, 'git_tracked_model', 'Model saved with Git integration'
509     # )
510

```

```
511 print( Git-integrated versioning provides commit-level tracking! )
```

Listing 41: Error 5.3.1: Model Versioning Issues - Buggy Code

```
1 # Buggy: Common data drift and model monitoring issues
2 import numpy as np
3 import pandas as pd
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score
7 from sklearn.datasets import make_classification
8 import matplotlib.pyplot as plt
9 from datetime import datetime, timedelta
10
11 # Create initial training data
12 np.random.seed(42)
13 X_train_orig, y_train_orig = make_classification(
14     n_samples=1000, n_features=10, n_classes=2, random_state=42
15 )
16
17 # Train initial model
18 model = RandomForestClassifier(random_state=42)
19 model.fit(X_train_orig, y_train_orig)
20
21 print( Initial model trained )
22 print(f Training accuracy: {model.score(X_train_orig, y_train_orig):.4f} )
23
24 # Simulate production data over time
25 def generate_production_data(base_date, n_samples=200, drift_factor=0):
26     Generate production data with potential drift
27     # Add some drift to feature distributions
28     X = np.random.randn(n_samples, 10)
29     X[:, 0] += drift_factor * 0.5 # Drift in first feature
30     X[:, 1] += drift_factor * 0.3 # Drift in second feature
31
32     # Change relationship between features and target (concept drift)
33     if drift_factor > 2:
34         y = (X[:, 0] + X[:, 2] > 0).astype(int) # Different relationship
35     else:
36         y = (X[:, 0] + X[:, 1] > 0).astype(int) # Original relationship
37
38     dates = [base_date + timedelta(days=i) for i in range(n_samples)]
39     return X, y, dates
40
41 # WRONG WAY 1 - No monitoring at all
42 print( \n=== WRONG WAY 1: No Monitoring === )
43
44 # WRONG: Deploy model and ignore monitoring
45 production_dates = [datetime(2024, 1, 1) + timedelta(days=i) for i in
46     range(30)]
47 X_prod, y_prod, dates_prod = generate_production_data(datetime(2024, 1, 1),
48     , 300, drift_factor=3)
49
49 # Make predictions without monitoring
50 predictions = model.predict(X_prod)
```

```

50 accuracy = accuracy_score(y_prod, predictions)
51
52 print(f Production accuracy: {accuracy:.4f} )
53 print( Problem: No detection of performance degradation due to data drift!
54       )
55
56 # WRONG WAY 2 - Manual monitoring with static thresholds
57 print( \n=== WRONG WAY 2: Static Thresholds === )
58
59 # WRONG: Using fixed thresholds that don't adapt
60 def manual_monitoring(predictions, actuals, threshold=0.8):
61     Manual monitoring with static threshold
62     current_accuracy = accuracy_score(actuals, predictions)
63
64     if current_accuracy < threshold:
65         print(f ALERT: Accuracy {current_accuracy:.4f} below threshold {
66             threshold} )
67         return True
68     else:
69         print(f Normal: Accuracy {current_accuracy:.4f} above threshold )
70         return False
71
72 # This approach fails when data distribution changes
73 alert_triggered = manual_monitoring(predictions, y_prod, threshold=0.8)
74 print(f Alert triggered: {alert_triggered} )
75
76 # WRONG WAY 3 - No concept drift detection
77 print( \n=== WRONG WAY 3: No Concept Drift Detection === )
78
79 # WRONG: Only monitoring input features, not target relationships
80 def simple_drift_detection(X_current, X_reference):
81     Simple drift detection - only looks at input features
82     # Compare means
83     current_means = np.mean(X_current, axis=0)
84     reference_means = np.mean(X_reference, axis=0)
85
86     drift_scores = np.abs(current_means - reference_means)
87     total_drift = np.mean(drift_scores)
88
89     return total_drift
90
91 # This misses concept drift (changes in P(Y|X))
92 drift_score = simple_drift_detection(X_prod, X_train_orig)
93 print(f Feature drift score: {drift_score:.4f} )
94 print( Problem: Feature drift detected, but concept drift (P(Y|X)) not
95       monitored! )
96
97 # WRONG WAY 4 - No statistical significance testing
98 print( \n=== WRONG WAY 4: No Statistical Testing === )
99
100 # WRONG: Not using statistical tests for significant changes

```

```

100 def naive_performance_comparison(old_perf, new_perf):
101     Naive performance comparison without statistical testing
102     if new_perf < old_perf:
103         print( Performance decreased! )
104         return True
105     else:
106         print( Performance stable )
107         return False
108
109 # This can lead to false alarms due to random variation
110 old_performance = 0.85
111 new_performance = 0.82
112 decrease_detected = naive_performance_comparison(old_performance,
113         new_performance)
113 print(f Performance decrease detected: {decrease_detected} )
114 print( Problem: No statistical significance testing for performance
115         changes! )
116
116 # CORRECT WAY 1 - Statistical drift detection
117 print( \n=== CORRECT WAY 1: Statistical Drift Detection === )
118
119 from scipy import stats
120
121 class StatisticalDriftDetector:
122     Statistical approach to drift detection
123
124     def __init__(self, reference_data, alpha=0.05):
125         self.reference_data = reference_data
126         self.alpha = alpha
127         self.reference_stats = self._calculate_reference_stats()
128
129     def _calculate_reference_stats(self):
130         Calculate reference statistics
131         return {
132             'means': np.mean(self.reference_data, axis=0),
133             'stds': np.std(self.reference_data, axis=0),
134             'medians': np.median(self.reference_data, axis=0)
135         }
136
137     def detect_feature_drift(self, current_data):
138         Detect drift in input features using statistical tests
139         drift_results = {}
140
141         for i in range(current_data.shape[1]):
142             current_feature = current_data[:, i]
143             reference_feature = self.reference_data[:, i]
144
145             # Two-sample t-test for mean comparison
146             t_stat, p_value = stats.ttest_ind(current_feature,
147         reference_feature)
148
148             # Kolmogorov-Smirnov test for distribution comparison
149             ks_stat, ks_p_value = stats.ks_2samp(current_feature,
150         reference_feature)

```

```

150
151         drift_results[f'feature_{i}'] = {
152             't_test_p_value': p_value,
153             'ks_test_p_value': ks_p_value,
154             'drift_detected': p_value < self.alpha or ks_p_value <
self.alpha,
155             'mean_change': np.mean(current_feature) - self.
reference_stats['means'][i]
156         }
157
158     return drift_results
159
160     def detect_multivariate_drift(self, current_data):
161         Detect multivariate drift using Hotelling's T-squared
162         # Simplified approach - in practice, use more sophisticated
methods
163         reference_mean = np.mean(self.reference_data, axis=0)
164         current_mean = np.mean(current_data, axis=0)
165
166         # This is a simplified version - real implementation would be more
complex
167         mean_diff = np.linalg.norm(current_mean - reference_mean)
168
169         return {
170             'mean_difference': mean_diff,
171             'drift_magnitude': mean_diff / np.linalg.norm(reference_mean)
172         }
173
174 # Use statistical drift detection
175 drift_detector = StatisticalDriftDetector(X_train_orig)
176
177 # Generate data with different drift levels
178 X_no_drift, y_no_drift, _ = generate_production_data(datetime(2024, 1, 1),
200, drift_factor=0)
179 X_moderate_drift, y_moderate_drift, _ = generate_production_data(datetime
(2024, 1, 1), 200, drift_factor=1)
180 X_heavy_drift, y_heavy_drift, _ = generate_production_data(datetime(2024,
1, 1), 200, drift_factor=3)
181
182 print( Statistical drift detection results: )
183
184 # No drift case
185 no_drift_results = drift_detector.detect_feature_drift(X_no_drift)
186 drift_count_no = sum(1 for r in no_drift_results.values() if r['
drift_detected'])
187 print(f    No drift: {drift_count_no} features showing drift )
188
189 # Moderate drift case
190 moderate_drift_results = drift_detector.detect_feature_drift(
X_moderate_drift)
191 drift_count_moderate = sum(1 for r in moderate_drift_results.values() if r
['drift_detected'])
192 print(f    Moderate drift: {drift_count_moderate} features showing drift )
193

```

```

194 # Heavy drift case
195 heavy_drift_results = drift_detector.detect_feature_drift(X_heavy_drift)
196 drift_count_heavy = sum(1 for r in heavy_drift_results.values() if r['
    drift_detected'])
197 print(f    Heavy drift: {drift_count_heavy} features showing drift )
198
199 # CORRECT WAY 2 - Concept drift detection
200 print( \n== CORRECT WAY 2: Concept Drift Detection == )
201
202 class ConceptDriftDetector:
203     Detect changes in the relationship between features and target
204
205     def __init__(self, window_size=100):
206         self.window_size = window_size
207         self.performance_history = []
208         self.prediction_history = []
209
210     def add_predictions(self, y_true, y_pred, timestamps=None):
211         Add prediction results to history
212         if timestamps is None:
213             timestamps = [datetime.now()] * len(y_true)
214
215         for i in range(len(y_true)):
216             self.prediction_history.append({
217                 'timestamp': timestamps[i],
218                 'y_true': y_true[i],
219                 'y_pred': y_pred[i],
220                 'correct': y_true[i] == y_pred[i]
221             })
222
223     # Calculate current window accuracy
224     if len(self.prediction_history) >= self.window_size:
225         recent_predictions = self.prediction_history[-self.window_size
:]
226         current_accuracy = np.mean([p['correct'] for p in
recent_predictions])
227         self.performance_history.append(current_accuracy)
228
229     def detect_concept_drift(self, significance_level=0.05):
230         Detect concept drift using performance degradation
231         if len(self.performance_history) < 2:
232             return False, 0.0
233
234     # Compare recent performance with historical baseline
235     if len(self.performance_history) > 5:
236         baseline_performance = np.mean(self.performance_history[:-3])
237         recent_performance = np.mean(self.performance_history[-3:])
238
239     # Statistical test for performance degradation
240     if len(self.performance_history) > 6:
241         # Use t-test to check if difference is significant
242         t_stat, p_value = stats.ttest_rel(
243             self.performance_history[-6:-3], # Previous period
244             self.performance_history[-3:]    # Current period

```

```

245         )
246
247         performance_degradation = baseline_performance -
recent_performance
248         significant_degradation = p_value < significance_level and
performance_degradation > 0.05
249
250         return significant_degradation, performance_degradation
251
252     return False, 0.0
253
254     def get_performance_trend(self):
255         Get performance trend analysis
256         if len(self.performance_history) < 2:
257             return Insufficient data
258
259         recent_avg = np.mean(self.performance_history[-3:])
260         overall_avg = np.mean(self.performance_history)
261
262         if recent_avg < overall_avg - 0.05:
263             return Degrading
264         elif recent_avg > overall_avg + 0.05:
265             return Improving
266         else:
267             return Stable
268
269 # Use concept drift detection
270 concept_drift_detector = ConceptDriftDetector(window_size=50)
271
272 # Simulate predictions over time with concept drift
273 base_date = datetime(2024, 1, 1)
274 all_predictions = []
275 all_actuals = []
276 all_timestamps = []
277
278 # Period 1: No drift
279 X_period1, y_period1, dates_period1 = generate_production_data(base_date,
100, drift_factor=0)
280 pred_period1 = model.predict(X_period1)
281 concept_drift_detector.add_predictions(y_period1, pred_period1,
dates_period1)
282
283 # Period 2: Concept drift
284 X_period2, y_period2, dates_period2 = generate_production_data(base_date +
timedelta(days=100), 100, drift_factor=3)
285 pred_period2 = model.predict(X_period2)
286 concept_drift_detector.add_predictions(y_period2, pred_period2,
dates_period2)
287
288 print( Concept drift detection results: )
289 drift_detected, degradation = concept_drift_detector.detect_concept_drift
()
290 trend = concept_drift_detector.get_performance_trend()
291

```



```

292 print(f    Concept drift detected: {drift_detected} )
293 print(f    Performance degradation: {degradation:.4f} )
294 print(f    Performance trend: {trend} )
295
296 # CORRECT WAY 3 - Adaptive monitoring thresholds
297 print( \n=== CORRECT WAY 3: Adaptive Monitoring === )
298
299 class AdaptivePerformanceMonitor:
300     Adaptive performance monitoring with dynamic thresholds
301
302     def __init__(self, initial_window_size=100, confidence_level=0.95):
303         self.initial_window_size = initial_window_size
304         self.confidence_level = confidence_level
305         self.performance_history = []
306         self.baseline_stats = None
307         self.alert_threshold = None
308
309     def establish_baseline(self, performances):
310         Establish baseline performance statistics
311         self.performance_history.extend(performances)
312
313         # Calculate baseline statistics
314         mean_perf = np.mean(performances)
315         std_perf = np.std(performances)
316
317         # Calculate adaptive threshold using confidence intervals
318         alpha = 1 - self.confidence_level
319         margin_of_error = stats.norm.ppf(1 - alpha/2) * (std_perf / np.
sqrt(len(performances)))
320
321         self.baseline_stats = {
322             'mean': mean_perf,
323             'std': std_perf,
324             'margin_of_error': margin_of_error
325         }
326
327         # Set alert threshold (mean - 2 * margin of error)
328         self.alert_threshold = mean_perf - 2 * margin_of_error
329
330         print(f Baseline established: Mean={mean_perf:.4f}, Std={std_perf
:.4f} )
331         print(f Alert threshold: {self.alert_threshold:.4f} )
332
333     def monitor_performance(self, current_performance):
334         Monitor current performance against adaptive thresholds
335         self.performance_history.append(current_performance)
336
337         if self.baseline_stats is None:
338             print( Warning: Baseline not established yet )
339             return False
340
341         # Check if performance is below threshold
342         alert_triggered = current_performance < self.alert_threshold
343

```

```

344     # Update threshold adaptively (moving average approach)
345     if len(self.performance_history) > self.initial_window_size:
346         recent_performances = self.performance_history[-self.
initial_window_size:]
347         recent_mean = np.mean(recent_performances)
348         recent_std = np.std(recent_performances)
349
350     # Update threshold based on recent performance
351     alpha = 1 - self.confidence_level
352     margin_of_error = stats.norm.ppf(1 - alpha/2) * (recent_std /
np.sqrt(len(recent_performances)))
353     self.alert_threshold = recent_mean - 2 * margin_of_error
354
355     return alert_triggered
356
357     def get_performance_summary(self):
358         Get performance summary statistics
359         if not self.performance_history:
360             return {}
361
362         return {
363             'current_performance': self.performance_history[-1],
364             'historical_mean': np.mean(self.performance_history),
365             'historical_std': np.std(self.performance_history),
366             'alert_threshold': self.alert_threshold,
367             'samples_monitored': len(self.performance_history)
368         }
369
370 # Use adaptive monitoring
371 adaptive_monitor = AdaptivePerformanceMonitor()
372
373 # Establish baseline with initial good performance
374 baseline_performances = [0.85, 0.87, 0.86, 0.88, 0.84] # Good performance
375 adaptive_monitor.establish_baseline(baseline_performances)
376
377 # Monitor degrading performance
378 degrading_performances = [0.83, 0.81, 0.79, 0.75, 0.72, 0.68] # Degrading
379
380 print( Adaptive monitoring results: )
381 for i, perf in enumerate(degrading_performances):
382     alert = adaptive_monitor.monitor_performance(perf)
383     summary = adaptive_monitor.get_performance_summary()
384
385     print(f    Sample {i+1}: Performance={perf:.4f},
386           f    Threshold={summary['alert_threshold']:.4f},
387           f    Alert={alert} )
388
389 # CORRECT WAY 4 - Comprehensive monitoring dashboard
390 print( \n=== CORRECT WAY 4: Comprehensive Monitoring Dashboard === )
391
392 class ModelMonitoringDashboard:
393     Comprehensive model monitoring dashboard
394
395     def __init__(self):

```

```

396     self.drift_detector = None
397     self.concept_drift_detector = ConceptDriftDetector()
398     self.performance_monitor = AdaptivePerformanceMonitor()
399     self.data_quality_monitor = None
400     self.alerts = []
401
402     def initialize_with_training_data(self, X_train, y_train):
403         Initialize monitoring with training data
404         self.drift_detector = StatisticalDriftDetector(X_train)
405         print( Monitoring dashboard initialized with training data )
406
407     def process_batch(self, X_batch, y_batch, timestamps=None, model=None)
408     :
409         Process a batch of predictions
410         if model is not None:
411             y_pred = model.predict(X_batch)
412         else:
413             # Assume predictions are provided
414             y_pred = y_batch # Simplified for demo
415
416         # Feature drift detection
417         if self.drift_detector:
418             drift_results = self.drift_detector.detect_feature_drift(
419             X_batch)
420             drifted_features = sum(1 for r in drift_results.values() if r[
421             'drift_detected'])
422
423             if drifted_features > 0:
424                 self._add_alert(f Feature drift detected in {
425                 drifted_features} features )
426
427         # Concept drift detection
428         self.concept_drift_detector.add_predictions(y_batch, y_pred,
429         timestamps)
430         drift_detected, degradation = self.concept_drift_detector.
431         detect_concept_drift()
432
433         if drift_detected:
434             self._add_alert(f Concept drift detected (degradation: {
435             degradation:.4f}) )
436
437         # Performance monitoring
438         current_accuracy = accuracy_score(y_batch, y_pred)
439         alert_triggered = self.performance_monitor.monitor_performance(
440         current_accuracy)
441
442         if alert_triggered:
443             self._add_alert(f Performance below threshold: {
444             current_accuracy:.4f} )
445
446         return {
447             'accuracy': current_accuracy,
448             'feature_drift_count': drifted_features if 'drifted_features'
449             in locals() else 0,

```

```

440         'concept_drift': drift_detected,
441         'performance_alert': alert_triggered
442     }
443
444     def _add_alert(self, message):
445         Add alert to monitoring system
446         alert = {
447             'timestamp': datetime.now(),
448             'message': message,
449             'severity': 'HIGH' if 'drift' in message.lower() or '
performance' in message.lower() else 'MEDIUM'
450         }
451         self.alerts.append(alert)
452         print(f ALERT ({alert['severity']}): {message} )
453
454     def get_dashboard_summary(self):
455         Get comprehensive monitoring summary
456         summary = {
457             'total_alerts': len(self.alerts),
458             'recent_alerts': self.alerts[-5:] if self.alerts else [],
459             'performance_trend': self.concept_drift_detector.
get_performance_trend(),
460             'samples_processed': len(self.concept_drift_detector.
prediction_history)
461         }
462
463         if self.alerts:
464             severity_counts = {}
465             for alert in self.alerts:
466                 severity = alert['severity']
467                 severity_counts[severity] = severity_counts.get(severity,
0) + 1
468             summary['alert_severity_distribution'] = severity_counts
469
470         return summary
471
472     def generate_report(self):
473         Generate monitoring report
474         summary = self.get_dashboard_summary()
475
476         print( \n=== MODEL MONITORING DASHBOARD === )
477         print(f Samples Processed: {summary['samples_processed']} )
478         print(f Performance Trend: {summary['performance_trend']} )
479         print(f Total Alerts: {summary['total_alerts']} )
480
481         if 'alert_severity_distribution' in summary:
482             print( Alert Severity Distribution: )
483             for severity, count in summary['alert_severity_distribution'].
items():
484                 print(f {severity}: {count} )
485
486         if summary['recent_alerts']:
487             print( \nRecent Alerts: )
488             for alert in summary['recent_alerts']:

```

```

489         print(f      {alert['timestamp']}.strftime('%Y-%m-%d %H:%M:%S
        ')} -
490             f {alert['severity']} - {alert['message']} )
491
492 # Use comprehensive monitoring dashboard
493 dashboard = ModelMonitoringDashboard()
494 dashboard.initialize_with_training_data(X_train_orig, y_train_orig)
495
496 # Simulate production monitoring
497 print( Simulating production monitoring... )
498
499 # Batch 1: Normal performance
500 X_batch1, y_batch1, dates_batch1 = generate_production_data(datetime(2024,
        1, 1), 50, drift_factor=0)
501 results_batch1 = dashboard.process_batch(X_batch1, y_batch1, dates_batch1,
        model)
502
503 # Batch 2: Feature drift
504 X_batch2, y_batch2, dates_batch2 = generate_production_data(datetime(2024,
        1, 2), 50, drift_factor=1)
505 results_batch2 = dashboard.process_batch(X_batch2, y_batch2, dates_batch2,
        model)
506
507 # Batch 3: Concept drift + performance degradation
508 X_batch3, y_batch3, dates_batch3 = generate_production_data(datetime(2024,
        1, 3), 50, drift_factor=3)
509 results_batch3 = dashboard.process_batch(X_batch3, y_batch3, dates_batch3,
        model)
510
511 print( \nBatch Results: )
512 print(f Batch 1: {results_batch1} )
513 print(f Batch 2: {results_batch2} )
514 print(f Batch 3: {results_batch3} )
515
516 # Generate final report
517 dashboard.generate_report()
518
519 # Advanced monitoring techniques
520 class AdvancedMonitoring:
521     Advanced monitoring techniques
522
523     @staticmethod
524     def cusum_monitoring(baseline_mean, current_values, threshold=5, drift
        =1):
525         CUSUM (Cumulative Sum) control chart for drift detection
526         # CUSUM algorithm for detecting small shifts
527         cusum_pos = 0
528         cusum_neg = 0
529         alarms = []
530
531         for i, value in enumerate(current_values):
532             # Calculate deviation from baseline
533             deviation = value - baseline_mean

```

```

535         # Update CUSUM statistics
536         cusum_pos = max(0, cusum_pos + (deviation - drift/2))
537         cusum_neg = max(0, cusum_neg - (deviation + drift/2))
538
539         # Check for alarms
540         if cusum_pos > threshold or cusum_neg > threshold:
541             alarms.append({
542                 'index': i,
543                 'value': value,
544                 'cusum_pos': cusum_pos,
545                 'cusum_neg': cusum_neg
546             })
547
548     return alarms
549
550     @staticmethod
551     def ewma_monitoring(baseline_mean, current_values, lambda_param=0.2,
552                         threshold=2):
553         EWMA (Exponentially Weighted Moving Average) for monitoring
554         ewma_values = []
555         alarms = []
556
557         # Initialize EWMA
558         ewma = baseline_mean
559
560         for i, value in enumerate(current_values):
561             # Update EWMA
562             ewma = lambda_param * value + (1 - lambda_param) * ewma
563             ewma_values.append(ewma)
564
565             # Calculate standardized EWMA (simplified)
566             std_ewma = abs(ewma - baseline_mean)
567
568             # Check for alarms
569             if std_ewma > threshold:
570                 alarms.append({
571                     'index': i,
572                     'value': value,
573                     'ewma': ewma,
574                     'std_ewma': std_ewma
575                 })
576
577         return ewma_values, alarms
578
579 # Example of advanced monitoring techniques
580 print( \n=== ADVANCED: Statistical Process Control === )
581
582 # Simulate performance data with gradual degradation
583 baseline_performance = 0.85
584 performance_data = [baseline_performance - 0.01 * i + np.random.normal(0,
585                             0.02)
586                     for i in range(100)]
587
588 # CUSUM monitoring

```

```

587 cusum_alarms = AdvancedMonitoring.cusum_monitoring(baseline_performance,
    performance_data)
588 print(f CUSUM alarms detected: {len(cusum_alarms)} )
589
590 # EWMA monitoring
591 ewma_values, ewma_alarms = AdvancedMonitoring.ewma_monitoring(
    baseline_performance, performance_data)
592 print(f EWMA alarms detected: {len(ewma_alarms)} )
593
594 # Real-time alerting system
595 class RealTimeAlerting:
596     Real-time alerting system with multiple channels
597
598     def __init__(self):
599         self.alert_rules = []
600         self.alert_channels = []
601
602     def add_alert_rule(self, name, condition_func, severity='MEDIUM'):
603         Add alert rule
604         self.alert_rules.append({
605             'name': name,
606             'condition': condition_func,
607             'severity': severity
608         })
609
610     def add_alert_channel(self, channel_func):
611         Add alert channel (email, Slack, etc.)
612         self.alert_channels.append(channel_func)
613
614     def check_alerts(self, monitoring_data):
615         Check all alert rules
616         triggered_alerts = []
617
618         for rule in self.alert_rules:
619             if rule['condition'](monitoring_data):
620                 alert = {
621                     'name': rule['name'],
622                     'severity': rule['severity'],
623                     'timestamp': datetime.now(),
624                     'data': monitoring_data
625                 }
626                 triggered_alerts.append(alert)
627
628             # Send alerts through all channels
629             for channel in self.alert_channels:
630                 try:
631                     channel(alert)
632                 except Exception as e:
633                     print(f Failed to send alert through channel: {e}
634
635             )
636
637         return triggered_alerts
638
639 # Example alert channels

```

```

638 def console_alert_channel(alert):
639     Console alert channel
640     print(f      ALERT: {alert['name']} ({alert['severity']}) at {alert['
timestamp']}] )
641
642 def email_alert_channel(alert):
643     Email alert channel (simplified)
644     print(f      Email alert sent: {alert['name']} )
645
646 # Setup real-time alerting
647 alerting_system = RealTimeAlerting()
648 alerting_system.add_alert_channel(console_alert_channel)
649
650 # Add alert rules
651 alerting_system.add_alert_rule(
652     'High Drift',
653     lambda data: data.get('feature_drift_count', 0) > 3,
654     'HIGH'
655 )
656
657 alerting_system.add_alert_rule(
658     'Performance Drop',
659     lambda data: data.get('accuracy', 1.0) < 0.7,
660     'CRITICAL'
661 )
662
663 # Test alerting system
664 test_data = {'feature_drift_count': 5, 'accuracy': 0.65}
665 alerts = alerting_system.check_alerts(test_data)
666 print(f Real-time alerts triggered: {len(alerts)} )

```

Listing 42: Error 5.3.2: Data Drift and Model Monitoring Issues - Buggy Code

```

1 # Buggy: Common feature pipeline inconsistencies
2 import numpy as np
3 import pandas as pd
4 from sklearn.preprocessing import StandardScaler, LabelEncoder,
OneHotEncoder
5 from sklearn.impute import SimpleImputer
6 from sklearn.model_selection import train_test_split
7 from sklearn.ensemble import RandomForestClassifier
8 from sklearn.metrics import accuracy_score
9 from sklearn.datasets import make_classification
10 import warnings
11 warnings.filterwarnings('ignore')
12
13 # Create sample data with various data types and issues
14 np.random.seed(42)
15
16 # Create realistic dataset with mixed data types
17 data = pd.DataFrame({
18     'numeric_feature1': np.random.randn(1000),
19     'numeric_feature2': np.random.randn(1000),
20     'categorical_feature1': np.random.choice(['A', 'B', 'C'], 1000),
21     'categorical_feature2': np.random.choice(['X', 'Y'], 1000),

```



```

22     'binary_feature': np.random.choice([0, 1], 1000),
23     'date_feature': pd.date_range('2020-01-01', periods=1000, freq='D'),
24     'target': np.random.choice([0, 1], 1000)
25 })
26
27 # Add some missing values and data issues
28 data.loc[100:105, 'numeric_feature1'] = np.nan
29 data.loc[200:205, 'categorical_feature1'] = np.nan
30 data.loc[300:305, 'date_feature'] = pd.NaT
31
32 print( Original data shape: , data.shape)
33 print( Missing values: )
34 print(data.isnull().sum())
35
36 # WRONG WAY 1 - Inconsistent preprocessing between train/test
37 print( \n=== WRONG WAY 1: Inconsistent Train/Test Preprocessing === )
38
39 # Split data
40 X = data.drop('target', axis=1)
41 y = data['target']
42 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
43     random_state=42)
44
45 # WRONG: Different preprocessing for train and test
46 # Training preprocessing
47 numeric_features = ['numeric_feature1', 'numeric_feature2']
48 categorical_features = ['categorical_feature1', 'categorical_feature2']
49
50 # WRONG: Fit scaler on training data
51 scaler_train = StandardScaler()
52 X_train_numeric_scaled = scaler_train.fit_transform(X_train[
53     numeric_features])
54
55 # WRONG: Fit scaler on test data (INCONSISTENT!)
56 scaler_test = StandardScaler() # Different scaler!
57 X_test_numeric_scaled = scaler_test.fit_transform(X_test[numeric_features
58     ]) # WRONG: fit on test!
59
60 print( Training scaler mean: , scaler_train.mean_)
61 print( Test scaler mean: , scaler_test.mean_)
62 print( Problem: Different scaling parameters - data leakage and
63     inconsistency! )
64
65 # WRONG WAY 2 - Fitting preprocessing on entire dataset
66 print( \n=== WRONG WAY 2: Fitting on Entire Dataset === )
67
68 # WRONG: Fitting preprocessing on entire dataset (data leakage)
69 scaler_leakage = StandardScaler()
70 X_numeric_scaled_leakage = scaler_leakage.fit_transform(X[numeric_features
71     ]) # WRONG!
72
73 # Split after preprocessing
74 X_train_leak, X_test_leak, y_train_leak, y_test_leak = train_test_split(
75     X_numeric_scaled_leakage, y, test_size=0.2, random_state=42

```

```

71 )
72
73 print( Scaler fitted on entire dataset )
74 print( Problem: Test data influences training preprocessing - data leakage
       ! )
75
76 # WRONG WAY 3 - Inconsistent categorical encoding
77 print( \n=== WRONG WAY 3: Inconsistent Categorical Encoding === )
78
79 # WRONG: Different categories in train vs test
80 X_train_cat = X_train.copy()
81 X_test_cat = X_test.copy()
82
83 # Add new category to test set that doesn't exist in training
84 X_test_cat.loc[X_test_cat.index[0], 'categorical_feature1'] = 'D' # New
    category!
85
86 # WRONG: Fit encoder on training data
87 ohe_train = OneHotEncoder(sparse=False, handle_unknown='error')
88 X_train_encoded = ohe_train.fit_transform(X_train_cat[categorical_features
    ])
89
90 # WRONG: Transform test data with new category
91 try:
92     X_test_encoded = ohe_train.transform(X_test_cat[categorical_features])
93 except ValueError as e:
94     print(f Error with inconsistent categories: {e} )
95     print( Problem: New categories in test set not seen during training! )
96
97 # WRONG WAY 4 - Manual feature engineering inconsistencies
98 print( \n=== WRONG WAY 4: Manual Feature Engineering Inconsistencies === )
99
100 def manual_feature_engineering_train(df):
101     Manual feature engineering for training data
102     df_new = df.copy()
103
104     # WRONG: Hard-coded values based on training data
105     df_new['numeric_feature1_normalized'] = (df_new['numeric_feature1'] -
df_new['numeric_feature1'].mean()) / df_new['numeric_feature1'].std()
106
107     # WRONG: Date feature engineering with training-specific logic
108     df_new['month'] = df_new['date_feature'].dt.month
109     df_new['day_of_year'] = df_new['date_feature'].dt.dayofyear
110
111     # WRONG: Categorical encoding with training-specific mappings
112     category_mapping = {'A': 1, 'B': 2, 'C': 3} # Hard-coded
113     df_new['categorical_feature1_encoded'] = df_new['categorical_feature1'
].map(category_mapping)
114
115     return df_new
116
117 def manual_feature_engineering_test(df):
118     Manual feature engineering for test data - INCONSISTENT!
119     df_new = df.copy()

```

```

120
121     # WRONG: Different normalization using test data statistics
122     df_new['numeric_feature1_normalized'] = (df_new['numeric_feature1'] -
df_new['numeric_feature1'].mean()) / df_new['numeric_feature1'].std()
123
124     # WRONG: Different date processing
125     df_new['month'] = df_new['date_feature'].dt.month
126     # Missing day_of_year feature!
127
128     # WRONG: Different categorical encoding
129     category_mapping = {'A': 10, 'B': 20, 'C': 30} # Different mapping!
130     df_new['categorical_feature1_encoded'] = df_new['categorical_feature1'
].map(category_mapping)
131
132     return df_new
133
134 # Apply inconsistent feature engineering
135 X_train_manual = manual_feature_engineering_train(X_train)
136 X_test_manual = manual_feature_engineering_test(X_test)
137
138 print( Manual feature engineering applied )
139 print( Training features: , list(X_train_manual.columns))
140 print( Test features: , list(X_test_manual.columns))
141 print( Problem: Inconsistent feature engineering between train and test! )
142
143 # CORRECT WAY 1 - Proper pipeline with consistent preprocessing
144 print( \n=== CORRECT WAY 1: Proper Pipeline with Consistent Preprocessing
=== )
145
146 from sklearn.pipeline import Pipeline
147 from sklearn.compose import ColumnTransformer
148 from sklearn.base import BaseEstimator, TransformerMixin
149
150 class DateFeatureExtractor(BaseEstimator, TransformerMixin):
151     Custom transformer for date feature extraction
152
153     def __init__(self, date_columns):
154         self.date_columns = date_columns
155
156     def fit(self, X, y=None):
157         return self
158
159     def transform(self, X):
160         X_transformed = X.copy()
161         for col in self.date_columns:
162             if col in X_transformed.columns:
163                 X_transformed[f'{col}_year'] = X_transformed[col].dt.year
164                 X_transformed[f'{col}_month'] = X_transformed[col].dt.
month
165                 X_transformed[f'{col}_day'] = X_transformed[col].dt.day
166                 X_transformed[f'{col}_dayofweek'] = X_transformed[col].dt.
dayofweek
167         return X_transformed
168

```

```

169 # Define feature types
170 numeric_features = ['numeric_feature1', 'numeric_feature2']
171 categorical_features = ['categorical_feature1', 'categorical_feature2']
172 binary_features = ['binary_feature']
173 date_features = ['date_feature']
174
175 # Create preprocessing pipeline
176 preprocessor = ColumnTransformer(
177     transformers=[
178         ('num', Pipeline([
179             ('imputer', SimpleImputer(strategy='median')),
180             ('scaler', StandardScaler())
181         ]), numeric_features),
182
183         ('cat', Pipeline([
184             ('imputer', SimpleImputer(strategy='constant', fill_value='
missing'))),
185             ('encoder', OneHotEncoder(sparse=False, handle_unknown='ignore
'))
186         ]), categorical_features),
187
188         ('binary', 'passthrough', binary_features)
189     ],
190     remainder='drop'
191 )
192
193 # Full pipeline including date feature extraction
194 full_pipeline = Pipeline([
195     ('date_features', DateFeatureExtractor(date_features)),
196     ('preprocessor', preprocessor),
197     ('classifier', RandomForestClassifier(random_state=42))
198 ])
199
200 # Split data properly
201 X_train_correct, X_test_correct, y_train_correct, y_test_correct =
    train_test_split(
202     X, y, test_size=0.2, random_state=42
203 )
204
205 # Fit pipeline on training data only
206 full_pipeline.fit(X_train_correct, y_train_correct)
207
208 # Transform test data using fitted pipeline
209 predictions_correct = full_pipeline.predict(X_test_correct)
210 accuracy_correct = accuracy_score(y_test_correct, predictions_correct)
211
212 print( Proper pipeline accuracy: , accuracy_correct)
213 print( Problem solved: Consistent preprocessing applied to both train and
    test! )
214
215 # CORRECT WAY 2 - Custom feature engineering pipeline
216 print( \n=== CORRECT WAY 2: Custom Feature Engineering Pipeline === )
217
218 class FeatureEngineer(BaseEstimator, TransformerMixin):

```

```

219     Custom feature engineering transformer
220
221     def __init__(self):
222         self.numeric_stats = {}
223         self.category_mappings = {}
224         self.date_stats = {}
225
226     def fit(self, X, y=None):
227         Fit feature engineering parameters on training data
228         # Store numeric statistics
229         for col in ['numeric_feature1', 'numeric_feature2']:
230             if col in X.columns:
231                 self.numeric_stats[col] = {
232                     'mean': X[col].mean(),
233                     'std': X[col].std(),
234                     'min': X[col].min(),
235                     'max': X[col].max()
236                 }
237
238         # Store categorical mappings
239         for col in ['categorical_feature1', 'categorical_feature2']:
240             if col in X.columns:
241                 unique_values = X[col].dropna().unique()
242                 self.category_mappings[col] = {val: idx for idx, val in
enumerate(unique_values)}
243
244         # Store date statistics
245         for col in ['date_feature']:
246             if col in X.columns and not X[col].isna().all():
247                 self.date_stats[col] = {
248                     'min_date': X[col].min(),
249                     'max_date': X[col].max()
250                 }
251
252         return self
253
254     def transform(self, X):
255         Apply feature engineering using fitted parameters
256         X_transformed = X.copy()
257
258         # Apply numeric transformations
259         for col, stats in self.numeric_stats.items():
260             if col in X_transformed.columns:
261                 # Standardization using training statistics
262                 X_transformed[f'{col}_standardized'] = (
263                     X_transformed[col] - stats['mean']
264                 ) / stats['std']
265
266                 # Min-max scaling
267                 range_val = stats['max'] - stats['min']
268                 if range_val > 0:
269                     X_transformed[f'{col}_normalized'] = (
270                         X_transformed[col] - stats['min']
271                     ) / range_val

```

```

272
273     # Apply categorical transformations
274     for col, mapping in self.category_mappings.items():
275         if col in X_transformed.columns:
276             # Use training mappings, handle unknown categories
277             X_transformed[f'{col}_encoded'] = X_transformed[col].map(
mapping).fillna(-1)
278
279     # Apply date transformations
280     for col, stats in self.date_stats.items():
281         if col in X_transformed.columns and not X_transformed[col].
isna().all():
282             # Extract date components
283             X_transformed[f'{col}_year'] = X_transformed[col].dt.year
284             X_transformed[f'{col}_month'] = X_transformed[col].dt.
month
285             X_transformed[f'{col}_day'] = X_transformed[col].dt.day
286             X_transformed[f'{col}_dayofweek'] = X_transformed[col].dt.
dayofweek
287
288             # Days since reference date
289             if 'min_date' in stats:
290                 X_transformed[f'{col}_days_since_start'] = (
291                     X_transformed[col] - stats['min_date']
292                 ).dt.days
293
294     return X_transformed
295
296 # Use custom feature engineering pipeline
297 feature_engineer = FeatureEngineer()
298 feature_engineer.fit(X_train)
299
300 # Transform both training and test data consistently
301 X_train_engineered = feature_engineer.transform(X_train)
302 X_test_engineered = feature_engineer.transform(X_test)
303
304 print( Custom feature engineering applied consistently )
305 print( Training engineered features: , list(X_train_engineered.columns))
306 print( Test engineered features: , list(X_test_engineered.columns))
307
308 # CORRECT WAY 3 - Pipeline with proper handling of missing categories
309 print( \n=== CORRECT WAY 3: Pipeline with Missing Category Handling === )
310
311 # Create pipeline that properly handles unseen categories
312 robust_preprocessor = ColumnTransformer(
313     transformers=[
314         ('num', Pipeline([
315             ('imputer', SimpleImputer(strategy='median')),
316             ('scaler', StandardScaler())
317         ]), numeric_features),
318
319         ('cat', Pipeline([
320             ('imputer', SimpleImputer(strategy='constant', fill_value='
missing'))],

```

```

321         ('encoder', OneHotEncoder(sparse=False, handle_unknown='ignore
    '))
322     ]), categorical_features),
323
324     ('binary', 'passthrough', binary_features)
325 ],
326     remainder='drop'
327 )
328
329 # Test with unseen categories
330 X_test_with_new_cat = X_test.copy()
331 X_test_with_new_cat.loc[X_test_with_new_cat.index[0], '
    categorical_feature1'] = 'NEW_CATEGORY'
332
333 # Fit on training data
334 robust_preprocessor.fit(X_train)
335
336 # Transform test data with new category - should handle gracefully
337 X_test_transformed = robust_preprocessor.transform(X_test_with_new_cat)
338 print( Pipeline handles unseen categories gracefully )
339 print( Transformed shape: , X_test_transformed.shape)
340
341 # CORRECT WAY 4 - Feature validation and consistency checking
342 print( \n=== CORRECT WAY 4: Feature Validation and Consistency === )
343
344 class FeatureValidator(BaseEstimator, TransformerMixin):
345     Feature validator to ensure consistency
346
347     def __init__(self, expected_features=None):
348         self.expected_features = expected_features
349         self.feature_stats = {}
350
351     def fit(self, X, y=None):
352         Store feature statistics for validation
353         self.feature_stats = {}
354         for col in X.columns:
355             if X[col].dtype in ['int64', 'float64']:
356                 self.feature_stats[col] = {
357                     'mean': X[col].mean(),
358                     'std': X[col].std(),
359                     'min': X[col].min(),
360                     'max': X[col].max()
361                 }
362         return self
363
364     def transform(self, X):
365         Validate and transform features
366         # Check feature consistency
367         if self.expected_features and set(X.columns) != set(self.
    expected_features):
368             missing_features = set(self.expected_features) - set(X.columns
    )
369             extra_features = set(X.columns) - set(self.expected_features)
370

```

```

371         if missing_features:
372             print(f Warning: Missing features: {missing_features} )
373         if extra_features:
374             print(f Warning: Extra features: {extra_features} )
375
376         # Check for extreme values (data drift detection)
377         for col in X.columns:
378             if col in self.feature_stats and X[col].dtype in ['int64', '
float64']:
379                 stats = self.feature_stats[col]
380                 outliers = X[(X[col] < stats['mean'] - 3*stats['std']) |
381                             (X[col] > stats['mean'] + 3*stats['std'])]
382                 if len(outliers) > 0:
383                     print(f Warning: {len(outliers)} outliers detected in
{col} )
384
385         return X
386
387 # Use feature validator in pipeline
388 validated_pipeline = Pipeline([
389     ('feature_engineer', FeatureEngineer()),
390     ('validator', FeatureValidator()),
391     ('classifier', RandomForestClassifier(random_state=42))
392 ])
393
394 # Fit and validate
395 validated_pipeline.fit(X_train, y_train)
396 predictions_validated = validated_pipeline.predict(X_test)
397 accuracy_validated = accuracy_score(y_test, predictions_validated)
398
399 print( Validated pipeline accuracy: , accuracy_validated)
400
401 # Advanced feature pipeline with error handling
402 class RobustFeaturePipeline:
403     Robust feature pipeline with comprehensive error handling
404
405     def __init__(self):
406         self.pipeline_steps = []
407         self.fitted_params = {}
408
409     def add_step(self, name, transformer):
410         Add processing step to pipeline
411         self.pipeline_steps.append((name, transformer))
412
413     def fit_transform(self, X, y=None):
414         Fit and transform data through pipeline
415         X_processed = X.copy()
416
417         for name, transformer in self.pipeline_steps:
418             try:
419                 if hasattr(transformer, 'fit_transform'):
420                     X_processed = transformer.fit_transform(X_processed, y
)
421
                     elif hasattr(transformer, 'fit') and hasattr(transformer,

```



```

    'transform'):
422         transformer.fit(X_processed, y)
423         X_processed = transformer.transform(X_processed)
424     else:
425         X_processed = transformer(X_processed)
426
427     print(f Step '{name}' completed successfully )
428
429     except Exception as e:
430         print(f Error in step '{name}': {e} )
431         # Continue with previous state or apply fallback
432         continue
433
434     return X_processed
435
436 def transform(self, X):
437     Transform new data using fitted pipeline
438     X_processed = X.copy()
439
440     for name, transformer in self.pipeline_steps:
441         try:
442             if hasattr(transformer, 'transform'):
443                 X_processed = transformer.transform(X_processed)
444             else:
445                 X_processed = transformer(X_processed)
446
447         except Exception as e:
448             print(f Error transforming with step '{name}': {e} )
449             continue
450
451     return X_processed
452
453 # Example of robust pipeline
454 print( \n=== ADVANCED: Robust Feature Pipeline === )
455
456 robust_pipeline = RobustFeaturePipeline()
457 robust_pipeline.add_step('imputer', SimpleImputer(strategy='median'))
458 robust_pipeline.add_step('scaler', StandardScaler())
459 robust_pipeline.add_step('feature_selection', SelectKBest(k=5)) # This
    would need to be defined
460
461 print( Robust pipeline framework ready for production use! )
462
463 # Feature pipeline versioning and rollback
464 class VersionedFeaturePipeline:
465     Feature pipeline with versioning and rollback capabilities
466
467     def __init__(self):
468         self.versions = {}
469         self.current_version = None
470
471     def save_version(self, version_name, pipeline):
472         Save pipeline version
473         import pickle

```

```

474     import datetime
475
476     version_info = {
477         'pipeline': pickle.dumps(pipeline),
478         'timestamp': datetime.datetime.now(),
479         'version_name': version_name
480     }
481
482     self.versions[version_name] = version_info
483     self.current_version = version_name
484     print(f Pipeline version '{version_name}' saved )
485
486     def load_version(self, version_name):
487         Load specific pipeline version
488         import pickle
489
490         if version_name not in self.versions:
491             raise ValueError(f Version {version_name} not found )
492
493         version_info = self.versions[version_name]
494         pipeline = pickle.loads(version_info['pipeline'])
495         self.current_version = version_name
496
497         print(f Pipeline version '{version_name}' loaded )
498         return pipeline
499
500     def rollback(self, version_name=None):
501         Rollback to previous version
502         if not version_name and len(self.versions) > 1:
503             # Get previous version
504             versions_list = list(self.versions.keys())
505             current_idx = versions_list.index(self.current_version)
506             if current_idx > 0:
507                 version_name = versions_list[current_idx - 1]
508
509         if version_name:
510             return self.load_version(version_name)
511         else:
512             raise ValueError( No previous version to rollback to )
513
514     # Example usage would require proper versioning setup
515     print( \nFeature pipeline versioning system ready for production
           deployment! )

```

Listing 43: Error 5.3.3: Feature Pipeline Inconsistencies - Buggy Code

```

1 # Buggy: Common error handling and graceful degradation issues
2 import numpy as np
3 import pandas as pd
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score
7 from sklearn.datasets import make_classification
8 import time
9 import logging

```

```

10 from datetime import datetime
11
12 # Setup logging
13 logging.basicConfig(level=logging.INFO)
14 logger = logging.getLogger(__name__)
15
16 # Create sample data
17 np.random.seed(42)
18 X, y = make_classification(n_samples=1000, n_features=10, random_state=42)
19 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
20                                                     random_state=42)
21
22 # Train model
23 model = RandomForestClassifier(random_state=42)
24 model.fit(X_train, y_train)
25
26 print( Model trained successfully )
27 print(f Training accuracy: {model.score(X_train, y_train):.4f} )
28
29 # WRONG WAY 1 - No error handling at all
30 print( \n=== WRONG WAY 1: No Error Handling === )
31
32 # WRONG: No error handling - will crash on any issue
33 def predict_no_error_handling(model, X):
34     Prediction without error handling
35     predictions = model.predict(X) # No validation, no error handling
36     return predictions
37
38 # This works fine with good data
39 try:
40     good_predictions = predict_no_error_handling(model, X_test)
41     print(f Predictions successful: {len(good_predictions)} predictions )
42 except Exception as e:
43     print(f Error occurred: {e} )
44
45 # But what if there's an issue?
46 print( Problem: No protection against data issues, model failures, or
47       system problems! )
48
49 # WRONG WAY 2 - Basic try-except without proper handling
50 print( \n=== WRONG WAY 2: Basic Error Handling === )
51
52 # WRONG: Generic exception handling that doesn't address specific issues
53 def predict_basic_error_handling(model, X):
54     Basic error handling - too generic
55     try:
56         predictions = model.predict(X)
57         return predictions
58     except Exception as e:
59         print(f Error: {e} ) # Just print and continue
60         return None # Return None - caller must handle this!
61
62 # This can lead to cascading failures
63 result = predict_basic_error_handling(model, X_test)

```

```

62 if result is None:
63     print( Prediction failed - need to handle None result everywhere! )
64 print( Problem: Generic error handling doesn't solve the underlying issues
    ! )
65
66 # WRONG WAY 3 - Silent failures
67 print( \n=== WRONG WAY 3: Silent Failures === )
68
69 # WRONG: Silently handling errors without notification
70 def predict_silent_failure(model, X):
71     Silent failure handling
72     try:
73         predictions = model.predict(X)
74         return predictions
75     except:
76         return np.zeros(len(X)) # Return zeros silently!
77
78 # This can mask serious problems
79 silent_result = predict_silent_failure(model, X_test)
80 print(f Silent failure result shape: {silent_result.shape} )
81 print( Problem: Serious errors masked as successful predictions! )
82
83 # WRONG WAY 4 - No fallback mechanisms
84 print( \n=== WRONG WAY 4: No Fallback Mechanisms === )
85
86 # WRONG: No backup plan when primary model fails
87 def predict_no_fallback(model, X):
88     No fallback when model fails
89     predictions = model.predict(X)
90     return predictions
91
92 # What if the model file is corrupted or missing?
93 # In production, this would crash the entire system
94 print( Problem: Single point of failure - no redundancy or fallback! )
95
96 # CORRECT WAY 1 - Comprehensive error handling with logging
97 print( \n=== CORRECT WAY 1: Comprehensive Error Handling === )
98
99 import traceback
100
101 def predict_with_comprehensive_error_handling(model, X, model_name=
    default_model ):
102     Comprehensive error handling with logging and context
103
104     try:
105         # Input validation
106         if X is None:
107             raise ValueError( Input data X cannot be None )
108
109         if len(X) == 0:
110             raise ValueError( Input data X is empty )
111
112         # Check if model exists and is fitted
113         if model is None:

```

```

114         raise ValueError( Model is None )
115
116     if not hasattr(model, 'predict'):
117         raise AttributeError( Model does not have predict method )
118
119     # Make predictions
120     logger.info(f Starting prediction with {model_name} )
121     start_time = time.time()
122
123     predictions = model.predict(X)
124
125     end_time = time.time()
126     logger.info(f Prediction completed in {end_time - start_time:.4f}
seconds )
127     logger.info(f Generated {len(predictions)} predictions )
128
129     return predictions
130
131     except ValueError as e:
132         logger.error(f ValueError in {model_name} prediction: {e} )
133         logger.error(f Input shape: {X.shape if hasattr(X, 'shape') else '
Unknown'} )
134         raise
135     except AttributeError as e:
136         logger.error(f AttributeError in {model_name} prediction: {e} )
137         raise
138     except Exception as e:
139         logger.error(f Unexpected error in {model_name} prediction: {e} )
140         logger.error(f Traceback: {traceback.format_exc()} )
141         raise
142
143 # Test comprehensive error handling
144 try:
145     comprehensive_result = predict_with_comprehensive_error_handling(model
, X_test, random_forest )
146     print(f Comprehensive error handling successful: {len(
comprehensive_result)} predictions )
147 except Exception as e:
148     print(f Comprehensive error handling caught: {e} )
149
150 # CORRECT WAY 2 - Graceful degradation with fallback
151 print( \n=== CORRECT WAY 2: Graceful Degradation === )
152
153 from sklearn.linear_model import LogisticRegression
154 from sklearn.dummy import DummyClassifier
155
156 class GracefulPredictionService:
157     Prediction service with graceful degradation
158
159     def __init__(self):
160         self.primary_model = None
161         self.backup_model = None
162         self.fallback_model = None
163         self.model_status = {

```

```

164         'primary': 'unknown',
165         'backup': 'unknown',
166         'fallback': 'unknown'
167     }
168
169     def initialize_models(self, X_train, y_train):
170         Initialize models with different complexity levels
171         try:
172             # Primary model - complex
173             self.primary_model = RandomForestClassifier(n_estimators=100,
174 random_state=42)
175             self.primary_model.fit(X_train, y_train)
176             self.model_status['primary'] = 'ready'
177             logger.info( Primary model initialized successfully )
178         except Exception as e:
179             logger.error(f Failed to initialize primary model: {e} )
180             self.model_status['primary'] = 'failed'
181
182         try:
183             # Backup model - simpler
184             self.backup_model = LogisticRegression(random_state=42,
185 max_iter=1000)
186             self.backup_model.fit(X_train, y_train)
187             self.model_status['backup'] = 'ready'
188             logger.info( Backup model initialized successfully )
189         except Exception as e:
190             logger.error(f Failed to initialize backup model: {e} )
191             self.model_status['backup'] = 'failed'
192
193         try:
194             # Fallback model - very simple
195             self.fallback_model = DummyClassifier(strategy='most_frequent'
196 )
197             self.fallback_model.fit(X_train, y_train)
198             self.model_status['fallback'] = 'ready'
199             logger.info( Fallback model initialized successfully )
200         except Exception as e:
201             logger.error(f Failed to initialize fallback model: {e} )
202             self.model_status['fallback'] = 'failed'
203
204     def predict_with_fallback(self, X, timeout=30):
205         Predict with automatic fallback on failure
206         start_time = time.time()
207
208         # Try primary model first
209         if self.model_status['primary'] == 'ready':
210             try:
211                 logger.info( Attempting prediction with primary model )
212                 predictions = self._safe_predict(self.primary_model, X,
213 timeout)
214                 logger.info( Primary model prediction successful )
215                 return predictions, 'primary'
216             except Exception as e:
217                 logger.warning(f Primary model failed: {e} )

```

```

214         self.model_status['primary'] = 'failed'
215
216     # Try backup model
217     if self.model_status['backup'] == 'ready':
218         try:
219             logger.info( Attempting prediction with backup model )
220             predictions = self._safe_predict(self.backup_model, X,
timeout)
221             logger.info( Backup model prediction successful )
222             return predictions, 'backup'
223         except Exception as e:
224             logger.warning(f Backup model failed: {e} )
225             self.model_status['backup'] = 'failed'
226
227     # Try fallback model
228     if self.model_status['fallback'] == 'ready':
229         try:
230             logger.info( Attempting prediction with fallback model )
231             predictions = self._safe_predict(self.fallback_model, X,
timeout)
232             logger.info( Fallback model prediction successful )
233             return predictions, 'fallback'
234         except Exception as e:
235             logger.error(f Fallback model failed: {e} )
236             self.model_status['fallback'] = 'failed'
237
238     # All models failed
239     raise RuntimeError( All prediction models failed )
240
241     def _safe_predict(self, model, X, timeout):
242         Safe prediction with timeout
243         # In practice, you might use threading or other timeout mechanisms
244         # This is a simplified version
245         return model.predict(X)
246
247 # Test graceful degradation
248 graceful_service = GracefulPredictionService()
249 graceful_service.initialize_models(X_train, y_train)
250
251 try:
252     predictions, model_used = graceful_service.predict_with_fallback(
X_test)
253     print(f Graceful prediction successful using {model_used} model )
254     print(f Predictions: {len(predictions)} )
255 except Exception as e:
256     print(f Graceful prediction failed: {e} )
257
258 # CORRECT WAY 3 - Circuit breaker pattern
259 print( \n=== CORRECT WAY 3: Circuit Breaker Pattern === )
260
261 import time
262 from enum import Enum
263
264 class CircuitState(Enum):

```

```

265     CLOSED =    closed
266     OPEN =     open
267     HALF_OPEN = half_open
268
269 class CircuitBreaker:
270     Circuit breaker for preventing cascading failures
271
272     def __init__(self, failure_threshold=5, timeout=60):
273         self.failure_threshold = failure_threshold
274         self.timeout = timeout
275         self.failure_count = 0
276         self.last_failure_time = None
277         self.state = CircuitState.CLOSED
278
279     def call(self, func, *args, **kwargs):
280         Call function with circuit breaker protection
281
282         if self.state == CircuitState.OPEN:
283             if self.last_failure_time and time.time() - self.
last_failure_time > self.timeout:
284                 self.state = CircuitState.HALF_OPEN
285                 logger.info( Circuit breaker half-open - testing service )
286             else:
287                 raise Exception( Circuit breaker is OPEN - service
unavailable )
288
289         try:
290             result = func(*args, **kwargs)
291             self._on_success()
292             return result
293         except Exception as e:
294             self._on_failure()
295             raise e
296
297     def _on_success(self):
298         Handle successful call
299         self.failure_count = 0
300         self.state = CircuitState.CLOSED
301         logger.info( Circuit breaker closed - service restored )
302
303     def _on_failure(self):
304         Handle failed call
305         self.failure_count += 1
306         self.last_failure_time = time.time()
307
308         if self.failure_count >= self.failure_threshold:
309             self.state = CircuitState.OPEN
310             logger.warning(f Circuit breaker OPEN - {self.failure_count}
failures )
311         else:
312             logger.warning(f Circuit breaker failure count: {self.
failure_count} )
313
314 # Test circuit breaker

```



```

315 circuit_breaker = CircuitBreaker(failure_threshold=3, timeout=30)
316
317 def unreliable_service(X):
318     Simulate unreliable service
319     # Randomly fail 30% of the time
320     if np.random.random() < 0.3:
321         raise Exception( Service temporarily unavailable )
322     return model.predict(X)
323
324 print( Testing circuit breaker with unreliable service: )
325 for i in range(10):
326     try:
327         result = circuit_breaker.call(unreliable_service, X_test[:10])
328         print(f Call {i+1}: Success )
329     except Exception as e:
330         print(f Call {i+1}: Failed - {e} )
331         print(f Circuit state: {circuit_breaker.state} )
332         time.sleep(1) # Small delay between calls
333
334 # CORRECT WAY 4 - Health checks and monitoring
335 print( \n=== CORRECT WAY 4: Health Checks and Monitoring === )
336
337 class ModelHealthMonitor:
338     Comprehensive model health monitoring
339
340     def __init__(self):
341         self.health_metrics = {
342             'prediction_count': 0,
343             'error_count': 0,
344             'last_prediction_time': None,
345             'average_prediction_time': 0,
346             'prediction_success_rate': 1.0
347         }
348         self.error_log = []
349
350     def record_prediction(self, success, prediction_time, error=None):
351         Record prediction outcome for monitoring
352         self.health_metrics['prediction_count'] += 1
353
354         if success:
355             # Update timing metrics
356             if self.health_metrics['last_prediction_time']:
357                 self.health_metrics['average_prediction_time'] = (
358                     self.health_metrics['average_prediction_time'] * 0.9 +
359                     prediction_time * 0.1
360                 )
361             else:
362                 self.health_metrics['average_prediction_time'] =
prediction_time
363         else:
364             self.health_metrics['error_count'] += 1
365             if error:
366                 self.error_log.append({
367                     'timestamp': datetime.now(),

```

```

368         'error': str(error)
369     })
370     # Keep only last 100 errors
371     if len(self.error_log) > 100:
372         self.error_log = self.error_log[-100:]
373
374     # Update success rate
375     total_predictions = self.health_metrics['prediction_count']
376     successful_predictions = total_predictions - self.health_metrics['
error_count']
377     self.health_metrics['prediction_success_rate'] =
successful_predictions / total_predictions
378
379     self.health_metrics['last_prediction_time'] = datetime.now()
380
381     def get_health_status(self):
382         Get current health status
383         return {
384             'status': self._calculate_status(),
385             'metrics': self.health_metrics.copy(),
386             'recent_errors': self.error_log[-5:] if self.error_log else []
387         }
388
389     def _calculate_status(self):
390         Calculate overall health status
391         success_rate = self.health_metrics['prediction_success_rate']
392         error_count = self.health_metrics['error_count']
393
394         if success_rate >= 0.95 and error_count < 10:
395             return HEALTHY
396         elif success_rate >= 0.80 and error_count < 50:
397             return DEGRADED
398         else:
399             return UNHEALTHY
400
401     def is_healthy(self):
402         Check if service is healthy
403         return self._calculate_status() == HEALTHY
404
405 # Test health monitoring
406 health_monitor = ModelHealthMonitor()
407
408 def monitored_prediction(model, X, health_monitor):
409     Prediction with health monitoring
410     start_time = time.time()
411
412     try:
413         predictions = model.predict(X)
414         prediction_time = time.time() - start_time
415
416         health_monitor.record_prediction(True, prediction_time)
417         return predictions
418     except Exception as e:
419         prediction_time = time.time() - start_time

```

```

420     health_monitor.record_prediction(False, prediction_time, e)
421     raise
422
423 # Test monitored predictions
424 print( Testing monitored predictions: )
425 for i in range(20):
426     try:
427         # Simulate some failures
428         if i == 5 or i == 15: # Force some failures
429             raise Exception( Simulated prediction error )
430
431         result = monitored_prediction(model, X_test[:5], health_monitor)
432         print(f Prediction {i+1}: Success )
433     except Exception as e:
434         print(f Prediction {i+1}: Failed - {e} )
435
436 # Check health status periodically
437 if (i + 1) % 5 == 0:
438     health_status = health_monitor.get_health_status()
439     print(f Health status: {health_status['status']})
440     print(f Success rate: {health_status['metrics']['prediction_success_rate']:.2%} )
441
442 # Advanced error handling patterns
443 class AdvancedErrorHandling:
444     Advanced error handling patterns for production systems
445
446     @staticmethod
447     def retry_with_exponential_backoff(func, max_retries=3, base_delay=1):
448         Retry function with exponential backoff
449         import random
450
451         for attempt in range(max_retries + 1):
452             try:
453                 return func()
454             except Exception as e:
455                 if attempt == max_retries:
456                     raise e
457
458                 delay = base_delay * (2 ** attempt) + random.uniform(0, 1)
459                 logger.warning(f Attempt {attempt + 1} failed: {e}.
460 Retrying in {delay:.2f} seconds... )
461                 time.sleep(delay)
462
463     @staticmethod
464     def bulk_prediction_with_error_handling(model, X_batches, batch_size
465 =100):
466         Bulk prediction with individual batch error handling
467         all_predictions = []
468         batch_results = []
469
470         for i, batch in enumerate(X_batches):
471             try:
472                 batch_predictions = model.predict(batch)

```

```

471         all_predictions.extend(batch_predictions)
472         batch_results.append({
473             'batch': i,
474             'status': 'success',
475             'count': len(batch_predictions)
476         })
477         logger.info(f Batch {i}: Success - {len(batch_predictions)
} predictions )
478     except Exception as e:
479         logger.error(f Batch {i}: Failed - {e} )
480         batch_results.append({
481             'batch': i,
482             'status': 'failed',
483             'error': str(e)
484         })
485         # Use fallback predictions for this batch
486         fallback_predictions = np.full(len(batch), 0) # Default
to class 0
487         all_predictions.extend(fallback_predictions)
488
489     return np.array(all_predictions), batch_results
490
491 @staticmethod
492 def prediction_with_timeout(func, timeout_seconds=30):
493     Prediction with timeout protection
494     import signal
495
496     def timeout_handler(signum, frame):
497         raise TimeoutError(f Prediction timed out after {
timeout_seconds} seconds )
498
499     # Set up timeout
500     old_handler = signal.signal(signal.SIGALRM, timeout_handler)
501     signal.alarm(timeout_seconds)
502
503     try:
504         result = func()
505         signal.alarm(0) # Cancel timeout
506         return result
507     except TimeoutError:
508         raise
509     except Exception as e:
510         signal.alarm(0) # Cancel timeout
511         raise e
512     finally:
513         signal.signal(signal.SIGALRM, old_handler) # Restore old
handler
514
515 # Example usage of advanced error handling
516 print( \n=== ADVANCED: Production-Ready Error Handling === )
517
518 # Retry with exponential backoff
519 def unreliable_api_call():
520     Simulate unreliable API call

```

```

521     if np.random.random() < 0.7: # 70% failure rate
522         raise Exception( API temporarily unavailable )
523     return API call successful
524
525 try:
526     result = AdvancedErrorHandling.retry_with_exponential_backoff(
527         unreliable_api_call, max_retries=3, base_delay=0.5
528     )
529     print(f Retry successful: {result} )
530 except Exception as e:
531     print(f Retry failed after all attempts: {e} )
532
533 # Bulk prediction with error handling
534 X_batches = [X_test[i:i+50] for i in range(0, len(X_test), 50)]
535 bulk_predictions, batch_results = AdvancedErrorHandling.
536     bulk_prediction_with_error_handling(
537     model, X_batches
538 )
539 successful_batches = sum(1 for r in batch_results if r['status'] == '
540     success')
541 failed_batches = sum(1 for r in batch_results if r['status'] == 'failed')
542 print(f Bulk prediction results: )
543 print(f     Total batches: {len(batch_results)} )
544 print(f     Successful batches: {successful_batches} )
545 print(f     Failed batches: {failed_batches} )
546
547 # Error handling best practices summary
548 def error_handling_best_practices():
549     Summary of error handling best practices
550
551     practices = [
552         1. Always validate inputs before processing ,
553         2. Use specific exception handling, not generic except clauses ,
554         3. Log errors with context and stack traces ,
555         4. Implement graceful degradation with fallback mechanisms ,
556         5. Use circuit breakers to prevent cascading failures ,
557         6. Monitor service health and performance metrics ,
558         7. Implement retry logic with exponential backoff ,
559         8. Set timeouts to prevent hanging operations ,
560         9. Use bulk processing with individual error handling ,
561         10. Maintain error logs for debugging and monitoring
562     ]
563
564     print( \nError Handling Best Practices: )
565     for practice in practices:
566         print(f     {practice} )
567
568 error_handling_best_practices()
569
570 # Production deployment considerations
571 class ProductionDeploymentChecklist:
572     Checklist for production deployment error handling

```

```

573
574     def __init__(self):
575         self.checks = {
576             'input_validation': False,
577             'error_logging': False,
578             'fallback_mechanisms': False,
579             'circuit_breaker': False,
580             'health_monitoring': False,
581             'timeout_handling': False,
582             'retry_logic': False,
583             'bulk_processing': False
584         }
585
586     def run_checklist(self):
587         Run production deployment checklist
588         print( \n=== PRODUCTION DEPLOYMENT CHECKLIST === )
589
590         for check, completed in self.checks.items():
591             status = COMPLETED if completed else PENDING
592             print(f {check.replace('_', ' ').title()}: {status} )
593
594         completed_checks = sum(self.checks.values())
595         total_checks = len(self.checks)
596         print(f \nProgress: {completed_checks}/{total_checks} checks
597         completed )
598
599         if completed_checks == total_checks:
600             print( All production readiness checks passed! )
601         else:
602             print( Some checks still pending - review before
603             production deployment )
604
605 # Run production checklist
606 deployment_checklist = ProductionDeploymentChecklist()
607 deployment_checklist.run_checklist()

```

Listing 44: Error 5.3.4: Error Handling and Graceful Degradation - Buggy Code

```

1 # Buggy: Common performance and scalability issues
2 import numpy as np
3 import pandas as pd
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score
7 from sklearn.datasets import make_classification
8 import time
9 import memory_profiler
10 from functools import wraps
11
12 # Create large dataset to demonstrate performance issues
13 np.random.seed(42)
14 print( Creating large dataset for performance testing... )
15
16 # Large dataset (100K samples, 100 features)

```

```

17 X_large, y_large = make_classification(n_samples=100000, n_features=100,
    random_state=42)
18 X_train_large, X_test_large, y_train_large, y_test_large =
    train_test_split(
19     X_large, y_large, test_size=0.2, random_state=42
20 )
21
22 print(f Large dataset created: {X_large.shape} )
23
24 # WRONG WAY 1 - No performance optimization
25 print( \n=== WRONG WAY 1: No Performance Optimization === )
26
27 # WRONG: Using default parameters without optimization
28 def train_model_no_optimization(X, y):
29     Train model without performance considerations
30     model = RandomForestClassifier(random_state=42) # Default parameters
31     model.fit(X, y)
32     return model
33
34 # This will be slow and memory-intensive
35 start_time = time.time()
36 model_slow = train_model_no_optimization(X_train_large, y_train_large)
37 training_time_slow = time.time() - start_time
38
39 print(f Slow training time: {training_time_slow:.2f} seconds )
40 print( Problem: Default parameters may not be optimal for large datasets!
    )
41
42 # WRONG WAY 2 - Loading entire dataset into memory
43 print( \n=== WRONG WAY 2: Loading Entire Dataset === )
44
45 # WRONG: Loading all data at once without batching
46 def process_entire_dataset_at_once(X, y):
47     Process entire dataset at once - memory inefficient
48     # This approach loads everything into memory
49     processed_data = []
50     for i in range(len(X)):
51         # Expensive processing for each sample
52         processed_sample = X[i] * 2 + np.random.normal(0, 0.1, len(X[i]))
53         processed_data.append(processed_sample)
54     return np.array(processed_data)
55
56 # This can cause memory issues with large datasets
57 print( Processing large dataset without batching... )
58 start_time = time.time()
59 try:
60     processed_data = process_entire_dataset_at_once(X_large[:10000],
        y_large[:10000])
61     processing_time = time.time() - start_time
62     print(f Processing time: {processing_time:.2f} seconds )
63 except MemoryError:
64     print( Memory error occurred - dataset too large for available memory!
        )
65 print( Problem: Loading entire dataset can cause memory issues! )

```

```

66
67 # WRONG WAY 3 - No parallel processing
68 print( \n=== WRONG WAY 3: No Parallel Processing === )
69
70 # WRONG: Sequential processing without parallelization
71 def sequential_feature_engineering(X):
72     Sequential feature engineering - slow
73     engineered_features = np.zeros((X.shape[0], X.shape[1] * 3))
74
75     for i in range(X.shape[0]): # Process each sample sequentially
76         for j in range(X.shape[1]):
77             engineered_features[i, j*3] = X[i, j] ** 2
78             engineered_features[i, j*3 + 1] = np.sqrt(np.abs(X[i, j]))
79             engineered_features[i, j*3 + 2] = np.log(np.abs(X[i, j]) + 1)
80
81     return engineered_features
82
83 # This is very slow for large datasets
84 print( Sequential feature engineering on small subset... )
85 start_time = time.time()
86 sequential_result = sequential_feature_engineering(X_large[:1000]) #
    Small subset
87 sequential_time = time.time() - start_time
88 print(f Sequential processing time: {sequential_time:.2f} seconds )
89 print( Problem: Sequential processing doesn't utilize multiple CPU cores!
    )
90
91 # WRONG WAY 4 - No caching or memoization
92 print( \n=== WRONG WAY 4: No Caching === )
93
94 # WRONG: Recomputing expensive operations
95 def expensive_computation_no_cache(data):
96     Expensive computation without caching
97     # Simulate expensive computation
98     time.sleep(0.01) # Artificial delay
99     return np.sum(data ** 2)
100
101 def process_data_no_cache(X):
102     Process data without caching expensive operations
103     results = []
104     for i in range(len(X)):
105         # Same computation repeated multiple times
106         result1 = expensive_computation_no_cache(X[i])
107         result2 = expensive_computation_no_cache(X[i]) # Same computation
            again!
108         result3 = expensive_computation_no_cache(X[i]) # And again!
109         results.append(result1 + result2 + result3)
110     return results
111
112 print( Processing without caching... )
113 start_time = time.time()
114 no_cache_results = process_data_no_cache(X_large[:100]) # Small subset
115 no_cache_time = time.time() - start_time
116 print(f No caching processing time: {no_cache_time:.2f} seconds )

```



```

117 print( Problem: Redundant computations waste time and resources! )
118
119 # CORRECT WAY 1 - Optimized model parameters
120 print( \n=== CORRECT WAY 1: Optimized Model Parameters === )
121
122 # CORRECT: Using optimized parameters for large datasets
123 def train_model_optimized(X, y):
124     Train model with performance-optimized parameters
125     model = RandomForestClassifier(
126         n_estimators=100,          # Reasonable number
127         max_depth=10,             # Limit depth to prevent overfitting
128         min_samples_split=100,    # Increase to reduce overfitting on large
data
129         min_samples_leaf=50,      # Increase to reduce overfitting
130         max_features='sqrt',      # Use subset of features
131         n_jobs=-1,                # Use all available cores
132         random_state=42
133     )
134     model.fit(X, y)
135     return model
136
137 # Compare performance
138 start_time = time.time()
139 model_optimized = train_model_optimized(X_train_large, y_train_large)
140 training_time_optimized = time.time() - start_time
141
142 print(f Optimized training time: {training_time_optimized:.2f} seconds )
143 print(f Performance improvement: {training_time_slow/
training_time_optimized:.2f}x faster )
144
145 # CORRECT WAY 2 - Batch processing for large datasets
146 print( \n=== CORRECT WAY 2: Batch Processing === )
147
148 def process_data_in_batches(X, y, batch_size=1000):
149     Process data in batches to manage memory
150     n_samples = X.shape[0]
151     results = []
152
153     for i in range(0, n_samples, batch_size):
154         batch_end = min(i + batch_size, n_samples)
155         X_batch = X[i:batch_end]
156         y_batch = y[i:batch_end]
157
158         # Process batch
159         batch_result = {
160             'batch_start': i,
161             'batch_end': batch_end,
162             'batch_size': batch_end - i,
163             'processing_time': None
164         }
165
166         # Simulate batch processing
167         start_time = time.time()
168         # Actual processing would go here

```

```

169     batch_processing_time = time.time() - start_time
170     batch_result['processing_time'] = batch_processing_time
171
172     results.append(batch_result)
173
174     # Progress indicator
175     if (i // batch_size) % 10 == 0:
176         print(f Processed batch {(i // batch_size) + 1}/{(n_samples +
177             batch_size - 1) // batch_size} )
178
179     return results
180
181 # Process large dataset in batches
182 print( Processing large dataset in batches... )
183 batch_results = process_data_in_batches(X_large, y_large, batch_size=5000)
184 print(f Processed {len(batch_results)} batches successfully )
185 print( Problem solved: Memory usage controlled through batching! )
186
187 # CORRECT WAY 3 - Parallel processing
188 print( \n=== CORRECT WAY 3: Parallel Processing === )
189
190 from multiprocessing import Pool
191 import multiprocessing as mp
192
193 def parallel_feature_computation(sample):
194     Compute features for a single sample
195     features = []
196     for value in sample:
197         features.extend([value ** 2, np.sqrt(np.abs(value)), np.log(np.abs
198             (value) + 1)])
199     return features
200
201 def parallel_feature_engineering(X, n_processes=None):
202     Parallel feature engineering using multiprocessing
203     if n_processes is None:
204         n_processes = mp.cpu_count()
205
206     print(f Using {n_processes} processes for parallel processing )
207
208     with Pool(processes=n_processes) as pool:
209         results = pool.map(parallel_feature_computation, X)
210
211     return np.array(results)
212
213 # Compare sequential vs parallel processing
214 print( Parallel feature engineering on subset... )
215 start_time = time.time()
216 parallel_result = parallel_feature_engineering(X_large[:1000], n_processes
217     =4)
218 parallel_time = time.time() - start_time
219
220 print(f Parallel processing time: {parallel_time:.2f} seconds )
221 print(f Speedup: {sequential_time/parallel_time:.2f}x faster )

```

```

220 # CORRECT WAY 4 - Caching and memoization
221 print( \n=== CORRECT WAY 4: Caching and Memoization === )
222
223 from functools import lru_cache
224
225 # Cache expensive computations
226 @lru_cache(maxsize=1000)
227 def expensive_computation_cached(data_tuple):
228     Expensive computation with caching
229     # Convert tuple back to array for computation
230     data = np.array(data_tuple)
231     time.sleep(0.001) # Artificial delay
232     return np.sum(data ** 2)
233
234 def process_data_with_cache(X):
235     Process data with caching
236     results = []
237     for i in range(len(X)):
238         # Convert to tuple for hashing (cache key)
239         data_tuple = tuple(X[i])
240
241         # These calls will be cached after first computation
242         result1 = expensive_computation_cached(data_tuple)
243         result2 = expensive_computation_cached(data_tuple) # Cached!
244         result3 = expensive_computation_cached(data_tuple) # Cached!
245
246         results.append(result1 + result2 + result3)
247     return results
248
249 # Compare caching vs no caching
250 print( Processing with caching... )
251 start_time = time.time()
252 cache_results = process_data_with_cache(X_large[:100])
253 cache_time = time.time() - start_time
254
255 print(f Caching processing time: {cache_time:.2f} seconds )
256 print(f Caching speedup: {no_cache_time/cache_time:.2f}x faster )
257
258 # Advanced performance optimization techniques
259 class PerformanceOptimizer:
260     Advanced performance optimization techniques
261
262     @staticmethod
263     def memory_efficient_data_loading(file_path, chunk_size=10000):
264         Memory-efficient data loading with chunking
265         # This would be implemented for actual file loading
266         # For demo, simulate chunked processing
267         print(f Loading data in chunks of {chunk_size} )
268
269         # Simulate chunks
270         total_samples = 100000
271         n_chunks = (total_samples + chunk_size - 1) // chunk_size
272
273         for chunk_idx in range(n_chunks):

```

```

274         start_idx = chunk_idx * chunk_size
275         end_idx = min((chunk_idx + 1) * chunk_size, total_samples)
276         chunk_size_actual = end_idx - start_idx
277
278         # Process chunk
279         yield {
280             'chunk_index': chunk_idx,
281             'start_index': start_idx,
282             'end_index': end_idx,
283             'chunk_size': chunk_size_actual
284         }
285
286     @staticmethod
287     def optimized_model_training(X, y, model_type='random_forest'):
288         Optimized model training with performance considerations
289
290         if model_type == 'random_forest':
291             # Use optimized parameters for large datasets
292             from sklearn.ensemble import RandomForestClassifier
293
294             model = RandomForestClassifier(
295                 n_estimators=100,
296                 max_depth=15,
297                 min_samples_split=100,
298                 min_samples_leaf=50,
299                 max_features='sqrt',
300                 n_jobs=-1, # Use all cores
301                 random_state=42,
302                 verbose=1 # Show progress
303             )
304
305         elif model_type == 'gradient_boosting':
306             from sklearn.ensemble import HistGradientBoostingClassifier
307
308             # HistGradientBoosting is more memory efficient for large
309             datasets
310             model = HistGradientBoostingClassifier(
311                 max_iter=100,
312                 random_state=42,
313                 verbose=1
314             )
315
316             # Train model
317             start_time = time.time()
318             model.fit(X, y)
319             training_time = time.time() - start_time
320
321             return model, training_time
322
323     @staticmethod
324     def model_prediction_optimization(model, X, batch_size=1000):
325         Optimized model prediction with batching
326         n_samples = X.shape[0]
327         predictions = []

```

```

327
328     for i in range(0, n_samples, batch_size):
329         batch_end = min(i + batch_size, n_samples)
330         X_batch = X[i:batch_end]
331
332         batch_predictions = model.predict(X_batch)
333         predictions.extend(batch_predictions)
334
335         # Progress indicator for large datasets
336         if n_samples > 10000 and (i // batch_size) % 10 == 0:
337             progress = (i + batch_size) / n_samples * 100
338             print(f Prediction progress: {progress:.1f}% )
339
340     return np.array(predictions)
341
342 # Test advanced optimization techniques
343 print( \n=== ADVANCED: Performance Optimization Techniques === )
344
345 # Memory-efficient data loading
346 print( Memory-efficient data loading: )
347 chunks = list(PerformanceOptimizer.memory_efficient_data_loading('
348     dummy_path.csv', chunk_size=25000))
349 print(f Generated {len(chunks)} chunks )
350
351 # Optimized model training comparison
352 print( \nComparing model training optimization: )
353
354 # Standard Random Forest
355 rf_standard = RandomForestClassifier(n_estimators=100, random_state=42,
356     n_jobs=1)
357 start_time = time.time()
358 rf_standard.fit(X_train_large[:10000], y_train_large[:10000]) # Smaller
359     subset for demo
360 standard_time = time.time() - start_time
361
362 # Optimized Random Forest
363 rf_optimized, optimized_time = PerformanceOptimizer.
364     optimized_model_training(
365     X_train_large[:10000], y_train_large[:10000], 'random_forest'
366 )
367
368 print(f Standard RF training time: {standard_time:.2f} seconds )
369 print(f Optimized RF training time: {optimized_time:.2f} seconds )
370 if optimized_time > 0:
371     print(f Optimization speedup: {standard_time/optimized_time:.2f}x )
372
373 # Batch prediction optimization
374 print( \nBatch prediction optimization: )
375 batch_predictions = PerformanceOptimizer.model_prediction_optimization(
376     rf_optimized, X_test_large[:5000], batch_size=1000
377 )
378 print(f Batch predictions completed: {len(batch_predictions)} predictions
379     )
380

```

```

376 # Scalability testing framework
377 class ScalabilityTester:
378     Framework for testing model scalability
379
380     def __init__(self):
381         self.results = []
382
383     def test_scalability(self, model, X_train, y_train, X_test, y_test,
384                         sample_sizes=[1000, 5000, 10000, 50000]):
385         Test model scalability with different sample sizes
386
387         for size in sample_sizes:
388             if size > len(X_train):
389                 continue
390
391             print(f \nTesting with {size} samples... )
392
393             # Training scalability
394             X_train_subset = X_train[:size]
395             y_train_subset = y_train[:size]
396
397             start_time = time.time()
398             model.fit(X_train_subset, y_train_subset)
399             train_time = time.time() - start_time
400
401             # Prediction scalability
402             start_time = time.time()
403             predictions = model.predict(X_test[:min(1000, len(X_test))])
404             pred_time = time.time() - start_time
405
406             # Accuracy
407             accuracy = accuracy_score(y_test[:len(predictions)],
408 predictions)
409
410             result = {
411                 'sample_size': size,
412                 'training_time': train_time,
413                 'prediction_time': pred_time,
414                 'accuracy': accuracy
415             }
416
417             self.results.append(result)
418             print(f Training time: {train_time:.2f}s )
419             print(f Prediction time: {pred_time:.4f}s )
420             print(f Accuracy: {accuracy:.4f} )
421
422         return self.results
423
424     def plot_scalability_results(self):
425         Plot scalability results
426         if not self.results:
427             print( No results to plot )
428             return

```

```

429     import matplotlib.pyplot as plt
430
431     sample_sizes = [r['sample_size'] for r in self.results]
432     train_times = [r['training_time'] for r in self.results]
433     pred_times = [r['prediction_time'] for r in self.results]
434     accuracies = [r['accuracy'] for r in self.results]
435
436     fig, axes = plt.subplots(1, 3, figsize=(15, 5))
437
438     # Training time vs sample size
439     axes[0].plot(sample_sizes, train_times, 'o-')
440     axes[0].set_xlabel('Sample Size')
441     axes[0].set_ylabel('Training Time (s)')
442     axes[0].set_title('Training Time vs Sample Size')
443     axes[0].grid(True)
444
445     # Prediction time vs sample size
446     axes[1].plot(sample_sizes, pred_times, 'o-')
447     axes[1].set_xlabel('Sample Size')
448     axes[1].set_ylabel('Prediction Time (s)')
449     axes[1].set_title('Prediction Time vs Sample Size')
450     axes[1].grid(True)
451
452     # Accuracy vs sample size
453     axes[2].plot(sample_sizes, accuracies, 'o-')
454     axes[2].set_xlabel('Sample Size')
455     axes[2].set_ylabel('Accuracy')
456     axes[2].set_title('Accuracy vs Sample Size')
457     axes[2].grid(True)
458
459     plt.tight_layout()
460     plt.show()
461
462     # Test scalability
463     print( '\n=== SCALABILITY TESTING === ')
464
465     scaler = ScalabilityTester()
466     test_model = RandomForestClassifier(n_estimators=50, random_state=42,
467                                       n_jobs=-1)
468
469     # Test with smaller dataset for demo
470     scalability_results = scaler.test_scalability(
471         test_model, X_large, y_large, X_large, y_large,
472         sample_sizes=[1000, 5000, 10000, 20000]
473     )
474
475     print( '\nScalability test completed! ')
476     print( 'Performance scales reasonably with dataset size ')
477
478     # Resource monitoring and optimization
479     class ResourceMonitor:
480         Monitor system resources during model operations
481
482         def __init__(self):

```

```

482         self.metrics = []
483
484     def monitor_resources(self, func, *args, **kwargs):
485         Monitor resources during function execution
486         import psutil
487         import os
488
489         # Get initial resource usage
490         process = psutil.Process(os.getpid())
491         initial_memory = process.memory_info().rss / 1024 / 1024 # MB
492         initial_cpu = process.cpu_percent()
493
494         start_time = time.time()
495
496         # Execute function
497         result = func(*args, **kwargs)
498
499         end_time = time.time()
500         final_memory = process.memory_info().rss / 1024 / 1024 # MB
501         final_cpu = process.cpu_percent()
502
503         metrics = {
504             'execution_time': end_time - start_time,
505             'memory_used': final_memory - initial_memory,
506             'cpu_usage': final_cpu - initial_cpu,
507             'peak_memory': final_memory
508         }
509
510         self.metrics.append(metrics)
511         return result, metrics
512
513     def get_resource_summary(self):
514         Get resource usage summary
515         if not self.metrics:
516             return {}
517
518         return {
519             'avg_execution_time': np.mean([m['execution_time'] for m in
520 self.metrics]),
521             'max_memory_usage': np.max([m['peak_memory'] for m in self.
522 metrics]),
523             'avg_memory_increase': np.mean([m['memory_used'] for m in self
524 .metrics])
525         }
526
527 # Example resource monitoring
528 print( \n=== RESOURCE MONITORING === )
529
530 monitor = ResourceMonitor()
531
532 # Monitor a function execution
533 def sample_processing(X, y):
534     Sample processing function
535     model = RandomForestClassifier(n_estimators=50, random_state=42)

```



```

533     model.fit(X[:5000], y[:5000])
534     return model.predict(X[:1000])
535
536 result, metrics = monitor.monitor_resources(
537     sample_processing, X_large, y_large
538 )
539
540 print( Resource usage metrics: )
541 for key, value in metrics.items():
542     print(f    {key}: {value:.2f} )
543
544 # Performance optimization checklist
545 def performance_optimization_checklist():
546     Checklist for performance optimization
547
548     checklist = [
549         1. Use appropriate model parameters for dataset size ,
550         2. Implement batch processing for large datasets ,
551         3. Utilize parallel processing when possible ,
552         4. Apply caching for expensive repeated computations ,
553         5. Monitor memory usage and implement garbage collection ,
554         6. Use memory-efficient data structures ,
555         7. Optimize feature engineering pipelines ,
556         8. Consider approximate algorithms for very large datasets ,
557         9. Implement proper data loading strategies ,
558         10. Profile code to identify bottlenecks
559     ]
560
561     print( \nPerformance Optimization Checklist: )
562     for item in checklist:
563         print(f    {item} )
564
565 performance_optimization_checklist()
566
567 # Production scalability considerations
568 class ProductionScalability:
569     Production scalability considerations and best practices
570
571     @staticmethod
572     def scalability_assessment(X, y):
573         Assess scalability requirements
574
575         n_samples, n_features = X.shape
576         data_size_gb = (X.nbytes + y.nbytes) / (1024**3)
577
578         assessment = {
579             'dataset_size': (n_samples, n_features),
580             'data_size_gb': data_size_gb,
581             'scalability_level': None,
582             'recommendations': []
583         }
584
585         if n_samples < 10000:
586             assessment['scalability_level'] = 'Small'

```

```

587         assessment['recommendations'].extend([
588             'Standard scikit-learn models should work fine',
589             'Consider using cross-validation for better estimates'
590         ])
591     elif n_samples < 100000:
592         assessment['scalability_level'] = 'Medium'
593         assessment['recommendations'].extend([
594             'Use n_jobs=-1 for parallel processing',
595             'Consider batch processing for predictions',
596             'Monitor memory usage during training'
597         ])
598     elif n_samples < 1000000:
599         assessment['scalability_level'] = 'Large'
600         assessment['recommendations'].extend([
601             'Use HistGradientBoosting for better memory efficiency',
602             'Implement data streaming/batching',
603             'Consider dimensionality reduction',
604             'Use incremental learning algorithms when possible'
605         ])
606     else:
607         assessment['scalability_level'] = 'Very Large'
608         assessment['recommendations'].extend([
609             'Consider distributed computing frameworks',
610             'Use online/incremental learning algorithms',
611             'Implement data sampling strategies',
612             'Consider approximate algorithms',
613             'Use specialized libraries like Dask or Vaex'
614         ])
615
616     return assessment
617
618     @staticmethod
619     def distributed_computing_recommendations():
620         Recommendations for distributed computing
621
622         recommendations = {
623             'libraries': [
624                 'Dask - Parallel computing library',
625                 'Vaex - Out-of-core dataframes',
626                 'Ray - Distributed computing framework',
627                 'Spark - Big data processing'
628             ],
629             'strategies': [
630                 'Data parallelism - distribute data across nodes',
631                 'Model parallelism - distribute model training',
632                 'Pipeline parallelism - parallelize pipeline steps',
633                 'Ensemble parallelism - train ensemble members in parallel'
634             ],
635             'considerations': [
636                 'Network overhead vs computation time',
637                 'Data serialization costs',
638                 'Fault tolerance and recovery',
639                 'Resource allocation and scheduling'

```

```

640         ]
641     }
642
643     return recommendations
644
645 # Run scalability assessment
646 print( \n=== PRODUCTION SCALABILITY ASSESSMENT === )
647
648 assessment = ProductionScalability.scalability_assessment(X_large, y_large
649 )
650 print(f Dataset size: {assessment['dataset_size']} )
651 print(f Data size: {assessment['data_size_gb']:.2f} GB )
652 print(f Scalability level: {assessment['scalability_level']} )
653
654 print( \nRecommendations: )
655 for rec in assessment['recommendations']:
656     print(f - {rec} )
657
658 # Distributed computing recommendations
659 dist_recommendations = ProductionScalability.
660     distributed_computing_recommendations()
661 print( \nDistributed Computing Recommendations: )
662 print( Libraries: )
663 for lib in dist_recommendations['libraries']:
664     print(f - {lib} )
665
666 print( \nStrategies: )
667 for strategy in dist_recommendations['strategies']:
668     print(f - {strategy} )
669
670 print( \nConsiderations: )
671 for consideration in dist_recommendations['considerations']:
672     print(f - {consideration} )

```

Listing 45: Error 5.3.5: Performance and Scalability Issues - Buggy Code

## 6 Counter-Intuitive Behaviors

### 6.1 Model Behavior Surprises

```

1 # Counter-Intuitive: Random Forest feature importance can be misleading
2 import numpy as np
3 import pandas as pd
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import accuracy_score
7 from sklearn.datasets import make_classification
8 import matplotlib.pyplot as plt
9
10 # Create dataset with correlated features
11 np.random.seed(42)
12
13 # Create correlated features

```

```

14 n_samples = 1000
15 X_base = np.random.randn(n_samples, 3)
16
17 # Create correlated copies
18 X_correlated = np.hstack([
19     X_base,
20     X_base[:, 0:1] + np.random.normal(0, 0.1, (n_samples, 1)), #
21     X_base[:, 1:2] + np.random.normal(0, 0.1, (n_samples, 1)), #
22     X_base[:, 2:3] + np.random.normal(0, 0.1, (n_samples, 1)) #
23 ])
24
25 # Create target with only first original feature being important
26 y = (X_base[:, 0] > 0).astype(int)
27
28 feature_names = [f'feature_{i}' for i in range(X_correlated.shape[1])]
29
30 print( Dataset created with correlated features )
31 print( Only feature_0 is truly important for prediction )
32
33 # Train Random Forest
34 rf = RandomForestClassifier(n_estimators=100, random_state=42)
35 rf.fit(X_correlated, y)
36
37 # Get feature importances
38 importances = rf.feature_importances_
39 feature_importance_df = pd.DataFrame({
40     'feature': feature_names,
41     'importance': importances
42 }).sort_values('importance', ascending=False)
43
44 print( \nRandom Forest Feature Importances: )
45 print(feature_importance_df)
46
47 # Counter-intuitive result: Correlated features split importance
48 print( \n=== COUNTER-INTUITIVE BEHAVIOR === )
49 print( Observation: Features 0 and 3 are highly correlated (both derived
50     from same base) )
51 print( Feature 0 (original): Importance = , importances[0])
52 print( Feature 3 (correlated copy): Importance = , importances[3])
53 print( Problem: Importance is split between correlated features! )
54 print( The true importance of feature 0 is diluted across its correlated
55     copies. )
56
57 # Demonstrate the effect with permutation importance
58 from sklearn.inspection import permutation_importance
59
60 perm_importance = permutation_importance(rf, X_correlated, y, n_repeats
61     =10, random_state=42)
62
63 perm_importance_df = pd.DataFrame({
64     'feature': feature_names,

```

```

62     'perm_importance': perm_importance.importances_mean
63 }).sort_values('perm_importance', ascending=False)
64
65 print( \nPermutation Importances (more reliable): )
66 print(perm_importance_df)
67
68 # Visualize the difference
69 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 5))
70
71 # Feature importances
72 ax1.barh(feature_importance_df['feature'], feature_importance_df['
    importance'])
73 ax1.set_xlabel('Feature Importance')
74 ax1.set_title('Random Forest Feature Importances\n(Split among correlated
    features)')
75
76 # Permutation importances
77 ax2.barh(perm_importance_df['feature'], perm_importance_df['
    perm_importance'])
78 ax2.set_xlabel('Permutation Importance')
79 ax2.set_title('Permutation Importances\n(Shows true predictive power)')
80
81 plt.tight_layout()
82 plt.show()
83
84 # Another counter-intuitive example: High cardinality categorical features
85 print( \n=== ANOTHER COUNTER-INTUITIVE EXAMPLE === )
86
87 # Create dataset with high cardinality categorical feature
88 np.random.seed(42)
89 n_samples = 1000
90
91 # Create meaningful feature
92 meaningful_feature = np.random.randn(n_samples)
93 target = (meaningful_feature > 0).astype(int)
94
95 # Create high cardinality categorical feature (random noise)
96 high_cardinality_feature = np.random.choice(range(100), n_samples) # 100
    categories
97
98 # Combine features
99 X_mixed = np.column_stack([meaningful_feature, high_cardinality_feature])
100
101 # Train model
102 rf_mixed = RandomForestClassifier(n_estimators=100, random_state=42)
103 rf_mixed.fit(X_mixed, target)
104
105 importances_mixed = rf_mixed.feature_importances_
106 print(f Meaningful feature importance: {importances_mixed[0]:.4f} )
107 print(f High cardinality feature importance: {importances_mixed[1]:.4f} )
108 print( Counter-intuitive: Random noise with many categories can appear
    important! )
109 print( Reason: High cardinality features create many splits, appearing
    more important. )

```

```

110
111 # Demonstrate with different random states
112 print( \nTesting with different random states: )
113 for rs in [1, 42, 123, 456]:
114     rf_test = RandomForestClassifier(n_estimators=100, random_state=rs)
115     rf_test.fit(X_mixed, target)
116     imp = rf_test.feature_importances_
117     print(f Random state {rs}: Meaningful={imp[0]:.4f}, High cardinality={
        imp[1]:.4f} )

```

Listing 46: Behavior 6.1.1: Random Forest Feature Importance Misconceptions

```

1 # Counter-Intuitive: Cross-validation can give misleading results
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import (cross_val_score, KFold,
    StratifiedKFold,
5                                     TimeSeriesSplit, LeaveOneOut)
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.datasets import make_classification
9 from sklearn.metrics import accuracy_score
10 import matplotlib.pyplot as plt
11
12 # Create time series data to demonstrate CV issues
13 np.random.seed(42)
14
15 # Generate time series data with trend
16 n_samples = 1000
17 time_points = np.arange(n_samples)
18 trend = 0.01 * time_points
19 seasonal = 2 * np.sin(2 * np.pi * time_points / 50)
20 noise = np.random.normal(0, 0.5, n_samples)
21 y_values = trend + seasonal + noise
22
23 # Create binary classification target
24 y_ts = (y_values > np.median(y_values)).astype(int)
25
26 # Create features with temporal correlation
27 X_ts = np.column_stack([
28     np.roll(y_values, 1), # Lag 1
29     np.roll(y_values, 2), # Lag 2
30     np.roll(y_values, 3), # Lag 3
31 ])
32
33 # Remove first few rows with NaN from rolling
34 X_ts = X_ts[3:]
35 y_ts = y_ts[3:]
36
37 print( Time series dataset created )
38 print( Features are temporally correlated - future values can predict past
    ! )
39
40 # WRONG WAY: Regular K-Fold on time series
41 print( \n=== WRONG WAY: Regular K-Fold on Time Series === )

```

```

42
43 # This violates temporal ordering
44 kfold = KFold(n_splits=5, shuffle=True, random_state=42)
45 rf_model = RandomForestClassifier(n_estimators=50, random_state=42)
46
47 # Cross-validation that ignores temporal order
48 cv_scores_kfold = cross_val_score(rf_model, X_ts, y_ts, cv=kfold)
49
50 print( Regular K-Fold CV scores: , cv_scores_kfold)
51 print(f Mean CV score: {cv_scores_kfold.mean():.4f} (+/- {cv_scores_kfold.
    std() * 2:.4f}) )
52 print( Counter-intuitive: Very high scores due to data leakage! )
53
54 # CORRECT WAY: Time Series Split
55 print( \n=== CORRECT WAY: Time Series Split === )
56
57 # Proper temporal cross-validation
58 tscv = TimeSeriesSplit(n_splits=5)
59 cv_scores_tscv = cross_val_score(rf_model, X_ts, y_ts, cv=tscv)
60
61 print( Time Series CV scores: , cv_scores_tscv)
62 print(f Mean CV score: {cv_scores_tscv.mean():.4f} (+/- {cv_scores_tscv.
    std() * 2:.4f}) )
63 print( Much more realistic performance estimate! )
64
65 # Demonstrate the difference
66 print(f \nDifference in estimates: )
67 print(f K-Fold overestimates by: {cv_scores_kfold.mean() - cv_scores_tscv.
    mean():.4f} )
68
69 # Another counter-intuitive example: Stratified CV with imbalanced small
    datasets
70 print( \n=== ANOTHER COUNTER-INTUITIVE EXAMPLE: Small Imbalanced Datasets
    === )
71
72 # Create small, highly imbalanced dataset
73 np.random.seed(42)
74 X_imb, y_imb = make_classification(
75     n_samples=50,          # Very small
76     n_features=10,
77     n_classes=2,
78     weights=[0.9, 0.1],    # 90% class 0, 10% class 1
79     random_state=42
80 )
81
82 print(f Imbalanced dataset: {len(y_imb)} samples, class distribution: {np.
    bincount(y_imb)} )
83
84 # Regular K-Fold
85 kfold_imb = KFold(n_splits=5, shuffle=True, random_state=42)
86 cv_scores_regular = cross_val_score(LogisticRegression(random_state=42),
87     X_imb, y_imb, cv=kfold_imb)
88
89 # Stratified K-Fold

```

```

90 strat_kfold = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
91 cv_scores_strat = cross_val_score(LogisticRegression(random_state=42),
92                                   X_imb, y_imb, cv=strat_kfold)
93
94 print(f \nRegular K-Fold scores: {cv_scores_regular} )
95 print(f Stratified K-Fold scores: {cv_scores_strat} )
96 print(f Regular mean: {cv_scores_regular.mean():.4f} )
97 print(f Stratified mean: {cv_scores_strat.mean():.4f} )
98
99 # Counter-intuitive: Stratified CV can be more variable with small
   datasets
100 print( Counter-intuitive: Stratified CV has higher variance with small
   datasets! )
101 print( Reason: Forced class balance in small folds can create extreme
   cases. )
102
103 # Demonstrate with multiple random states
104 print( \nTesting variance with different random states: )
105
106 regular_scores_all = []
107 strat_scores_all = []
108
109 for rs in range(10):
110     kfold_temp = KFold(n_splits=5, shuffle=True, random_state=rs)
111     strat_kfold_temp = StratifiedKFold(n_splits=5, shuffle=True,
112                                       random_state=rs)
113
114     scores_reg = cross_val_score(LogisticRegression(random_state=42,
115                                                     max_iter=1000),
116                                 X_imb, y_imb, cv=kfold_temp)
117     scores_strat = cross_val_score(LogisticRegression(random_state=42,
118                                                       max_iter=1000),
119                                   X_imb, y_imb, cv=strat_kfold_temp)
120
121     regular_scores_all.append(scores_reg.mean())
122     strat_scores_all.append(scores_strat.mean())
123
124 print(f Regular K-Fold variance: {np.var(regular_scores_all):.6f} )
125 print(f Stratified K-Fold variance: {np.var(strat_scores_all):.6f} )
126 print( Counter-intuitive: Stratified CV can have higher variance on small
   datasets! )
127
128 # Leave-One-Out CV surprises
129 print( \n=== LEAVE-ONE-OUT CV SURPRISES === )
130
131 # L00 CV on small dataset
132 loo = LeaveOneOut()
133 loo_scores = cross_val_score(LogisticRegression(random_state=42, max_iter
134                                                 =1000),
135                               X_imb[:20], y_imb[:20], cv=loo)
136
137 print(f L00 CV scores (first 10): {loo_scores[:10]} )
138 print(f L00 CV mean: {loo_scores.mean():.4f} )
139 print(f L00 CV std: {loo_scores.std():.4f} )

```



```

136
137 # Counter-intuitive: L00 CV can be overly optimistic
138 print( Counter-intuitive: L00 CV often gives overly optimistic estimates!
139 )
139 print( Reason: Each model is trained on almost the full dataset. )
140
141 # Visualize CV score distributions
142 fig, axes = plt.subplots(1, 3, figsize=(15, 5))
143
144 axes[0].hist(cv_scores_kfold, bins=10, alpha=0.7, label='K-Fold')
145 axes[0].axvline(cv_scores_kfold.mean(), color='red', linestyle='--',
146                label=f'Mean: {cv_scores_kfold.mean():.4f}')
147 axes[0].set_xlabel('CV Score')
148 axes[0].set_ylabel('Frequency')
149 axes[0].set_title('Regular K-Fold CV Scores\n(Overly optimistic)')
150 axes[0].legend()
151
152 axes[1].hist(cv_scores_tscv, bins=10, alpha=0.7, label='Time Series CV')
153 axes[1].axvline(cv_scores_tscv.mean(), color='red', linestyle='--',
154                label=f'Mean: {cv_scores_tscv.mean():.4f}')
155 axes[1].set_xlabel('CV Score')
156 axes[1].set_ylabel('Frequency')
157 axes[1].set_title('Time Series CV Scores\n(Realistic)')
158 axes[1].legend()
159
160 axes[2].hist(loos_scores, bins=10, alpha=0.7, label='L00 CV')
161 axes[2].axvline(loos_scores.mean(), color='red', linestyle='--',
162                label=f'Mean: {loos_scores.mean():.4f}')
163 axes[2].set_xlabel('CV Score')
164 axes[2].set_ylabel('Frequency')
165 axes[2].set_title('Leave-One-Out CV Scores\n(Can be overly optimistic)')
166 axes[2].legend()
167
168 plt.tight_layout()
169 plt.show()
170
171 # Advanced CV behavior demonstration
172 class CVBehaviorAnalyzer:
173     Analyze counter-intuitive CV behaviors
174
175     @staticmethod
176     def analyze_cv_variance(X, y, cv_methods, n_trials=20):
177         Analyze variance across different CV methods
178
179         results = {}
180
181         for name, cv_method in cv_methods.items():
182             trial_scores = []
183
184             for trial in range(n_trials):
185                 # Use different random state for each trial
186                 scores = cross_val_score(
187                     LogisticRegression(random_state=trial, max_iter=1000),
188                     X, y, cv=cv_method

```

```

189         )
190         trial_scores.append(scores.mean())
191
192     results[name] = {
193         'mean': np.mean(trial_scores),
194         'std': np.std(trial_scores),
195         'scores': trial_scores
196     }
197
198     return results
199
200 @staticmethod
201 def demonstrate_data_leakage():
202     Demonstrate how CV can leak data
203
204     # Create dataset where future predicts past
205     np.random.seed(42)
206     n = 100
207     data = np.random.randn(n)
208
209     # Features are future values, target is past values
210     X_leak = data[2:].reshape(-1, 1) # Future values
211     y_leak = (data[:-2] > 0).astype(int) # Past values as target
212
213     # Regular CV will give overly optimistic results
214     kfold = KFold(n_splits=5, shuffle=True)
215     rf = RandomForestClassifier(n_estimators=50, random_state=42)
216
217     scores = cross_val_score(rf, X_leak, y_leak, cv=kfold)
218
219     print( Data leakage example: )
220     print(f CV scores with data leakage: {scores} )
221     print(f Mean: {scores.mean():.4f} )
222     print( Counter-intuitive: High scores due to using future to
predict past! )
223
224 # Run advanced analysis
225 print( \n=== ADVANCED CV BEHAVIOR ANALYSIS === )
226
227 cv_methods = {
228     'K-Fold': KFold(n_splits=5, shuffle=True, random_state=42),
229     'Stratified': StratifiedKFold(n_splits=5, shuffle=True, random_state
=42),
230     'Time Series': TimeSeriesSplit(n_splits=5)
231 }
232
233 analyzer = CVBehaviorAnalyzer()
234 cv_analysis = analyzer.analyze_cv_variance(X_imb, y_imb, cv_methods)
235
236 print( CV Method Comparison (Small Imbalanced Dataset): )
237 for method, results in cv_analysis.items():
238     print(f {method}: Mean={results['mean']:.4f}, Std={results['std']:.4
f} )
239

```

```

240 # Demonstrate data leakage
241 analyzer.demonstrate_data_leakage()
242
243 print( \nKey Counter-Intuitive CV Behaviors: )
244 print( 1. Regular CV on time series data causes data leakage )
245 print( 2. Stratified CV can have higher variance on small datasets )
246 print( 3. L00 CV often gives overly optimistic estimates )
247 print( 4. Feature engineering before CV causes data leakage )
248 print( 5. CV scores can be misleading without proper context )

```

Listing 47: Behavior 6.1.2: Cross-Validation Surprises

```

1 # Counter-Intuitive: Feature scaling can have unexpected effects
2 import numpy as np
3 import pandas as pd
4 from sklearn.preprocessing import StandardScaler, MinMaxScaler,
   RobustScaler
5 from sklearn.model_selection import train_test_split
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.svm import SVC
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn.metrics import accuracy_score
10 from sklearn.datasets import make_classification
11 import matplotlib.pyplot as plt
12
13 # Create dataset with features of different scales
14 np.random.seed(42)
15
16 # Create features with very different scales
17 n_samples = 1000
18
19 # Feature 1: Small scale (0-1 range)
20 feature1 = np.random.uniform(0, 1, n_samples)
21
22 # Feature 2: Medium scale (0-100 range)
23 feature2 = np.random.uniform(0, 100, n_samples)
24
25 # Feature 3: Large scale (0-10000 range)
26 feature3 = np.random.uniform(0, 10000, n_samples)
27
28 # Feature 4: Very large scale with outliers
29 feature4 = np.random.normal(0, 1, n_samples) * 100000
30 feature4[::100] = np.random.normal(0, 1, 10) * 1000000 # Add outliers
31
32 # Create target based on all features
33 target = ((feature1 * 10 + feature2 * 0.1 + feature3 * 0.001 + feature4 *
   0.00001) > 500).astype(int)
34
35 # Combine features
36 X_mixed = np.column_stack([feature1, feature2, feature3, feature4])
37 feature_names = ['small_scale', 'medium_scale', 'large_scale', '
   very_large_with_outliers']
38
39 print( Dataset created with features of different scales )
40 print( Feature statistics: )

```

```

41 for i, name in enumerate(feature_names):
42     print(f    {name}: Mean={X_mixed[:, i].mean():.2f}, Std={X_mixed[:, i].
        std():.2f} )
43
44 # Split data
45 X_train, X_test, y_train, y_test = train_test_split(X_mixed, target,
        test_size=0.2, random_state=42)
46
47 # Test different models with and without scaling
48 models = {
49     'LogisticRegression': LogisticRegression(random_state=42, max_iter
        =1000),
50     'SVM': SVC(random_state=42),
51     'RandomForest': RandomForestClassifier(random_state=42, n_estimators
        =100)
52 }
53
54 scalers = {
55     'No Scaling': None,
56     'StandardScaler': StandardScaler(),
57     'MinMaxScaler': MinMaxScaler(),
58     'RobustScaler': RobustScaler()
59 }
60
61 # Counter-intuitive result 1: Different models react differently to
    scaling
62 print( \n=== COUNTER-INTUITIVE: Models React Differently to Scaling === )
63
64 results = {}
65
66 for model_name, model in models.items():
67     results[model_name] = {}
68
69     for scaler_name, scaler in scalers.items():
70         try:
71             if scaler is not None:
72                 X_train_scaled = scaler.fit_transform(X_train)
73                 X_test_scaled = scaler.transform(X_test)
74             else:
75                 X_train_scaled = X_train
76                 X_test_scaled = X_test
77
78             # Train model
79             model.fit(X_train_scaled, y_train)
80
81             # Predict
82             y_pred = model.predict(X_test_scaled)
83             accuracy = accuracy_score(y_test, y_pred)
84
85             results[model_name][scaler_name] = accuracy
86
87         except Exception as e:
88             results[model_name][scaler_name] = f Error: {e}
89

```

```

90 # Display results
91 print( \nModel Performance with Different Scaling: )
92 for model_name, scaler_results in results.items():
93     print(f \n{model_name}: )
94     for scaler_name, accuracy in scaler_results.items():
95         print(f     {scaler_name}: {accuracy} )
96
97 print( \nCounter-intuitive observations: )
98 print( 1. Random Forest performance is similar with/without scaling )
99 print(     Reason: Tree-based models are scale-invariant )
100 print( 2. Logistic Regression and SVM performance varies significantly
        with scaling )
101 print(     Reason: Distance-based algorithms are sensitive to feature scales
        )
102 print( 3. Different scalers can give different results )
103 print(     StandardScaler: Mean=0, Std=1 )
104 print(     MinMaxScaler: Range=[0,1] )
105 print(     RobustScaler: Uses median and IQR (robust to outliers) )
106
107 # Counter-intuitive result 2: Scaling can hurt performance
108 print( \n=== COUNTER-INTUITIVE: Scaling Can Hurt Performance === )
109
110 # Create dataset where original scale is meaningful
111 np.random.seed(42)
112
113 # Features where scale represents importance
114 X_meaningful = np.array([
115     [1, 100, 1000],    # Feature 3 (large scale) is most important
116     [2, 200, 2000],
117     [3, 300, 3000],
118     [4, 400, 4000],
119     [5, 500, 5000]
120 ])
121
122 y_meaningful = (X_meaningful[:, 2] > 2500).astype(int) # Only feature 3
        matters
123
124 # Train models
125 lr_original = LogisticRegression(random_state=42)
126 lr_scaled = LogisticRegression(random_state=42)
127
128 # Without scaling
129 lr_original.fit(X_meaningful, y_meaningful)
130 coef_original = lr_original.coef_[0]
131
132 # With scaling
133 scaler_meaningful = StandardScaler()
134 X_meaningful_scaled = scaler_meaningful.fit_transform(X_meaningful)
135 lr_scaled.fit(X_meaningful_scaled, y_meaningful)
136 coef_scaled = lr_scaled.coef_[0]
137
138 print( Feature coefficients comparison: )
139 print( Original scale coefficients: , coef_original)
140 print( Scaled coefficients: , coef_scaled)

```

```

141
142 print( \nCounter-intuitive: Scaling can obscure feature importance! )
143 print( In original scale: Feature 3 coefficient is largest (correct) )
144 print( After scaling: All coefficients similar (feature importance hidden)
      )
145
146 # Counter-intuitive result 3: Different scalers handle outliers
      differently
147 print( \n=== COUNTER-INTUITIVE: Different Scaler Outlier Handling === )
148
149 # Create data with outliers
150 np.random.seed(42)
151 data_with_outliers = np.random.normal(0, 1, 100)
152 data_with_outliers[::20] = np.random.normal(0, 10, 5) # Add outliers
153
154 print( Data with outliers statistics: )
155 print(f      Mean: {data_with_outliers.mean():.4f} )
156 print(f      Std: {data_with_outliers.std():.4f} )
157 print(f      Min: {data_with_outliers.min():.4f} )
158 print(f      Max: {data_with_outliers.max():.4f} )
159
160 # Apply different scalers
161 scalers_comparison = {
162     'StandardScaler': StandardScaler(),
163     'MinMaxScaler': MinMaxScaler(),
164     'RobustScaler': RobustScaler()
165 }
166
167 print( \nEffect of different scalers on outliers: )
168 for name, scaler in scalers_comparison.items():
169     scaled_data = scaler.fit_transform(data_with_outliers.reshape(-1, 1)).
        flatten()
170     print(f {name}: )
171     print(f      Mean: {scaled_data.mean():.4f} )
172     print(f      Std: {scaled_data.std():.4f} )
173     print(f      Min: {scaled_data.min():.4f} )
174     print(f      Max: {scaled_data.max():.4f} )
175
176 print( \nCounter-intuitive: Different scalers treat outliers very
      differently! )
177 print( StandardScaler: Outliers affect mean/std significantly )
178 print( MinMaxScaler: Outliers dominate the range )
179 print( RobustScaler: Outliers have minimal effect (uses median/IQR) )
180
181 # Visualize the scaling effects
182 fig, axes = plt.subplots(2, 2, figsize=(12, 10))
183
184 # Original data
185 axes[0, 0].hist(data_with_outliers, bins=30, alpha=0.7)
186 axes[0, 0].set_title('Original Data with Outliers')
187 axes[0, 0].set_xlabel('Value')
188 axes[0, 0].set_ylabel('Frequency')
189
190 # StandardScaler

```

```

191 standard_scaled = StandardScaler().fit_transform(data_with_outliers.
    reshape(-1, 1)).flatten()
192 axes[0, 1].hist(standard_scaled, bins=30, alpha=0.7)
193 axes[0, 1].set_title('StandardScaler (Mean=0, Std=1)')
194 axes[0, 1].set_xlabel('Value')
195 axes[0, 1].set_ylabel('Frequency')
196
197 # MinMaxScaler
198 minmax_scaled = MinMaxScaler().fit_transform(data_with_outliers.reshape
    (-1, 1)).flatten()
199 axes[1, 0].hist(minmax_scaled, bins=30, alpha=0.7)
200 axes[1, 0].set_title('MinMaxScaler (Range=[0,1])')
201 axes[1, 0].set_xlabel('Value')
202 axes[1, 0].set_ylabel('Frequency')
203
204 # RobustScaler
205 robust_scaled = RobustScaler().fit_transform(data_with_outliers.reshape
    (-1, 1)).flatten()
206 axes[1, 1].hist(robust_scaled, bins=30, alpha=0.7)
207 axes[1, 1].set_title('RobustScaler (Robust to Outliers)')
208 axes[1, 1].set_xlabel('Value')
209 axes[1, 1].set_ylabel('Frequency')
210
211 plt.tight_layout()
212 plt.show()
213
214 # Counter-intuitive result 4: Scaling order matters in pipelines
215 print( \n=== COUNTER-INTUITIVE: Scaling Order in Pipelines === )
216
217 from sklearn.pipeline import Pipeline
218 from sklearn.feature_selection import SelectKBest, f_classif
219
220 # Create pipeline where order matters
221 X_pipeline, y_pipeline = make_classification(n_samples=1000, n_features
    =20, random_state=42)
222
223 # Pipeline 1: Scale then select features
224 pipeline1 = Pipeline([
225     ('scaler', StandardScaler()),
226     ('selector', SelectKBest(f_classif, k=10)),
227     ('classifier', LogisticRegression(random_state=42))
228 ])
229
230 # Pipeline 2: Select features then scale
231 pipeline2 = Pipeline([
232     ('selector', SelectKBest(f_classif, k=10)),
233     ('scaler', StandardScaler()),
234     ('classifier', LogisticRegression(random_state=42))
235 ])
236
237 # Split data
238 X_train_p, X_test_p, y_train_p, y_test_p = train_test_split(
239     X_pipeline, y_pipeline, test_size=0.2, random_state=42
240 )

```

```

241
242 # Fit pipelines
243 pipeline1.fit(X_train_p, y_train_p)
244 pipeline2.fit(X_train_p, y_train_p)
245
246 # Predict
247 pred1 = pipeline1.predict(X_test_p)
248 pred2 = pipeline2.predict(X_test_p)
249
250 acc1 = accuracy_score(y_test_p, pred1)
251 acc2 = accuracy_score(y_test_p, pred2)
252
253 print( Pipeline performance comparison: )
254 print(f Scale then select: {acc1:.4f} )
255 print(f Select then scale: {acc2:.4f} )
256
257 print( \nCounter-intuitive: Order of operations affects results! )
258 print( Scale then select: Uses all features for scaling, then selects )
259 print( Select then scale: Selects features first, then scales selected
        features )
260 print( The scaling parameters will be different in each case! )
261
262 # Advanced scaling behavior analysis
263 class ScalingBehaviorAnalyzer:
264     Analyze counter-intuitive scaling behaviors
265
266     @staticmethod
267     def analyze_model_sensitivity(models, scalers, X, y):
268         Analyze how different models respond to scaling
269
270         results = {}
271
272         for model_name, model in models.items():
273             results[model_name] = {}
274
275             for scaler_name, scaler in scalers.items():
276                 try:
277                     if scaler_name == 'No Scaling':
278                         X_processed = X
279                     else:
280                         X_processed = scaler.fit_transform(X)
281
282                     # Cross-validation to get stable estimates
283                     from sklearn.model_selection import cross_val_score
284                     scores = cross_val_score(model, X_processed, y, cv=5)
285
286                     results[model_name][scaler_name] = {
287                         'mean': scores.mean(),
288                         'std': scores.std(),
289                         'scores': scores
290                     }
291
292                 except Exception as e:
293                     results[model_name][scaler_name] = {'error': str(e)}

```



```

294
295     return results
296
297     @staticmethod
298     def demonstrate_scaling_impact_on_coefficients():
299         Demonstrate how scaling affects model coefficients
300
301         # Create dataset where feature scale represents importance
302         np.random.seed(42)
303         n_samples = 1000
304
305         # Features with different scales but same predictive power
306         X_impact = np.column_stack([
307             np.random.normal(0, 1, n_samples),      # Unit scale
308             np.random.normal(0, 10, n_samples),     # 10x scale
309             np.random.normal(0, 100, n_samples)     # 100x scale
310         ])
311
312         # Target depends equally on all features
313         y_impact = (X_impact[:, 0] + X_impact[:, 1]/10 + X_impact[:,
314             2]/100 > 0).astype(int)
315
316         # Models
317         lr_unscaled = LogisticRegression(random_state=42)
318         lr_scaled = LogisticRegression(random_state=42)
319
320         # Fit without scaling
321         lr_unscaled.fit(X_impact, y_impact)
322         coef_unscaled = lr_unscaled.coef_[0]
323
324         # Fit with scaling
325         scaler_impact = StandardScaler()
326         X_impact_scaled = scaler_impact.fit_transform(X_impact)
327         lr_scaled.fit(X_impact_scaled, y_impact)
328         coef_scaled = lr_scaled.coef_[0]
329
330         print( Coefficient analysis: )
331         print( True importance: All features equally important )
332         print(f Unscaled coefficients: {coef_unscaled} )
333         print(f Scaled coefficients: {coef_scaled} )
334
335         print( \nCounter-intuitive: Unscaled coefficients reflect data
336             scale, not importance! )
337         print( After scaling: Coefficients reflect true feature importance
338             )
339
340     # Run advanced analysis
341     print( \n=== ADVANCED SCALING BEHAVIOR ANALYSIS === )
342
343     # Analyze model sensitivity
344     models_analysis = {
345         'LinearSVM': SVC(kernel='linear', random_state=42),
346         'LogisticRegression': LogisticRegression(random_state=42),
347         'RandomForest': RandomForestClassifier(random_state=42, n_estimators

```

```

    =50)
345 }
346
347 scalers_analysis = {
348     'No Scaling': 'No Scaling',
349     'StandardScaler': StandardScaler(),
350     'MinMaxScaler': MinMaxScaler(),
351     'RobustScaler': RobustScaler()
352 }
353
354 # Use a smaller dataset for demonstration
355 X_small, y_small = make_classification(n_samples=500, n_features=10,
    random_state=42)
356
357 analyzer = ScalingBehaviorAnalyzer()
358 sensitivity_results = analyzer.analyze_model_sensitivity(
359     models_analysis, scalers_analysis, X_small, y_small
360 )
361
362 print( Model sensitivity to scaling: )
363 for model_name, scaler_results in sensitivity_results.items():
364     print(f \n{model_name}: )
365     for scaler_name, results in scaler_results.items():
366         if 'error' not in results:
367             print(f    {scaler_name}: {results['mean']:.4f} (+/- {results['
std']:.4f}) )
368         else:
369             print(f    {scaler_name}: Error - {results['error']} )
370
371 # Demonstrate coefficient impact
372 analyzer.demonstrate_scaling_impact_on_coefficients()
373
374 print( \nKey Counter-Intuitive Scaling Behaviors: )
375 print( 1. Tree-based models are scale-invariant )
376 print( 2. Distance-based models are highly scale-sensitive )
377 print( 3. Scaling can obscure true feature importance )
378 print( 4. Different scalers handle outliers differently )
379 print( 5. Order of operations in pipelines matters )
380 print( 6. Scaling parameters should only be fit on training data )

```

Listing 48: Behavior 6.1.3: Feature Scaling Surprises

```

1 # Counter-Intuitive: Ensemble methods can behave unexpectedly
2 import numpy as np
3 import pandas as pd
4 from sklearn.ensemble import (RandomForestClassifier,
    GradientBoostingClassifier,
5                                AdaBoostClassifier, VotingClassifier,
    BaggingClassifier)
6 from sklearn.model_selection import train_test_split, cross_val_score
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.tree import DecisionTreeClassifier
9 from sklearn.metrics import accuracy_score, classification_report
10 from sklearn.datasets import make_classification
11 import matplotlib.pyplot as plt

```

```

12
13 # Create challenging dataset to demonstrate ensemble behaviors
14 np.random.seed(42)
15
16 # Create dataset with different difficulty levels for different models
17 X_complex, y_complex = make_classification(
18     n_samples=2000,
19     n_features=20,
20     n_informative=10,
21     n_redundant=5,
22     n_clusters_per_class=2,
23     class_sep=0.8,
24     flip_y=0.1, # Add some noise
25     random_state=42
26 )
27
28 # Add some highly correlated features to make it more challenging
29 correlated_features = X_complex[:, :5] + np.random.normal(0, 0.1, (2000,
30     5))
31
32 X_complex = np.hstack([X_complex, correlated_features])
33
34 X_train, X_test, y_train, y_test = train_test_split(
35     X_complex, y_complex, test_size=0.3, random_state=42
36 )
37
38 print( Complex dataset created for ensemble testing )
39 print(f Dataset shape: {X_complex.shape} )
40
41 # Individual models for comparison
42 models_individual = {
43     'LogisticRegression': LogisticRegression(random_state=42, max_iter
44         =1000),
45     'DecisionTree': DecisionTreeClassifier(random_state=42),
46     'RandomForest': RandomForestClassifier(random_state=42, n_estimators
47         =100)
48 }
49
50 # Test individual models
51 print( \n=== INDIVIDUAL MODEL PERFORMANCE === )
52 individual_results = {}
53
54 for name, model in models_individual.items():
55     model.fit(X_train, y_train)
56     y_pred = model.predict(X_test)
57     accuracy = accuracy_score(y_test, y_pred)
58     individual_results[name] = accuracy
59     print(f {name}: {accuracy:.4f} )
60
61 # Counter-intuitive result 1: Ensemble not always better than best
62     individual model
63 print( \n=== COUNTER-INTUITIVE: Ensemble Not Always Better === )
64
65 # Create voting ensemble
66 voting_clf = VotingClassifier(

```

```

62     estimators=[
63         ('lr', LogisticRegression(random_state=42, max_iter=1000)),
64         ('dt', DecisionTreeClassifier(random_state=42)),
65         ('rf', RandomForestClassifier(random_state=42, n_estimators=50))
66     ],
67     voting='hard'
68 )
69
70 voting_clf.fit(X_train, y_train)
71 y_pred_voting = voting_clf.predict(X_test)
72 voting_accuracy = accuracy_score(y_test, y_pred_voting)
73
74 print( Voting Ensemble Results: )
75 print(f Voting Classifier: {voting_accuracy:.4f} )
76 print( Individual model accuracies: )
77 for name, acc in individual_results.items():
78     print(f {name}: {acc:.4f} )
79
80 best_individual = max(individual_results.values())
81 print(f \nBest individual model: {best_individual:.4f} )
82 print(f Ensemble vs best individual: {voting_accuracy - best_individual:.4f} )
83
84 print( \nCounter-intuitive: Ensemble can be worse than best individual
      model! )
85 print( Reasons: )
86 print( 1. Poor individual models can drag down ensemble performance )
87 print( 2. Correlated errors among ensemble members )
88 print( 3. Voting can average away the strengths of good models )
89
90 # Counter-intuitive result 2: Bagging can hurt performance on simple
  problems
91 print( \n=== COUNTER-INTUITIVE: Bagging Can Hurt Performance === )
92
93 # Simple, well-behaved dataset
94 X_simple, y_simple = make_classification(
95     n_samples=1000,
96     n_features=5,
97     n_informative=5,
98     n_redundant=0,
99     class_sep=2.0, # Well separated
100     random_state=42
101 )
102
103 X_train_simple, X_test_simple, y_train_simple, y_test_simple =
    train_test_split(
104     X_simple, y_simple, test_size=0.3, random_state=42
105 )
106
107 # Simple model (should work well on simple data)
108 simple_model = LogisticRegression(random_state=42)
109 simple_model.fit(X_train_simple, y_train_simple)
110 simple_accuracy = simple_model.score(X_test_simple, y_test_simple)
111

```

```

112 # Bagged version of simple model
113 bagged_model = BaggingClassifier(
114     base_estimator=LogisticRegression(random_state=42),
115     n_estimators=10,
116     random_state=42
117 )
118 bagged_model.fit(X_train_simple, y_train_simple)
119 bagged_accuracy = bagged_model.score(X_test_simple, y_test_simple)
120
121 print( Simple Dataset Results: )
122 print(f Simple LogisticRegression: {simple_accuracy:.4f} )
123 print(f Bagged LogisticRegression: {bagged_accuracy:.4f} )
124 print(f Difference: {bagged_accuracy - simple_accuracy:.4f} )
125
126 print( \nCounter-intuitive: Bagging can hurt performance on simple
        problems! )
127 print( Reasons: )
128 print( 1. Simple problems don't benefit from variance reduction )
129 print( 2. Bagging introduces additional complexity without benefit )
130 print( 3. Small datasets may not support effective bagging )
131
132 # Counter-intuitive result 3: More estimators doesn't always mean better
133 print( \n=== COUNTER-INTUITIVE: More Estimators Doesn't Always Mean Better
        === )
134
135 # Test Random Forest with different numbers of estimators
136 n_estimators_range = [1, 5, 10, 50, 100, 200, 500]
137 rf_accuracies = []
138
139 print( Testing Random Forest with different n_estimators: )
140 for n_est in n_estimators_range:
141     rf = RandomForestClassifier(n_estimators=n_est, random_state=42)
142     rf.fit(X_train, y_train)
143     accuracy = rf.score(X_test, y_test)
144     rf_accuracies.append(accuracy)
145     print(f n_estimators={n_est}: {accuracy:.4f} )
146
147 # Plot the results
148 plt.figure(figsize=(10, 6))
149 plt.plot(n_estimators_range, rf_accuracies, 'o-')
150 plt.xlabel('Number of Estimators')
151 plt.ylabel('Accuracy')
152 plt.title('Random Forest Performance vs Number of Estimators')
153 plt.grid(True)
154 plt.show()
155
156 # Find optimal point
157 optimal_idx = np.argmax(rf_accuracies)
158 optimal_n_est = n_estimators_range[optimal_idx]
159 optimal_accuracy = rf_accuracies[optimal_idx]
160
161 print(f \nOptimal n_estimators: {optimal_n_est} )
162 print(f Optimal accuracy: {optimal_accuracy:.4f} )
163 print(f Diminishing returns after n_estimators={optimal_n_est} )

```

```

164
165 print( \nCounter-intuitive: More estimators can lead to diminishing
      returns! )
166 print( Reasons: )
167 print( 1. Each additional estimator provides marginal improvement )
168 print( 2. Computational cost increases linearly )
169 print( 3. Overfitting risk with too many estimators )
170 print( 4. Random variation can cause performance to plateau or decrease )
171
172 # Counter-intuitive result 4: Ensemble diversity vs performance trade-off
173 print( \n=== COUNTER-INTUITIVE: Ensemble Diversity Trade-off === )
174
175 from sklearn.metrics import accuracy_score
176
177 # Create diverse ensemble members
178 ensemble_members = [
179     LogisticRegression(random_state=42),
180     DecisionTreeClassifier(random_state=42),
181     RandomForestClassifier(n_estimators=50, random_state=42),
182     GradientBoostingClassifier(random_state=42, n_estimators=50)
183 ]
184
185 # Train individual members
186 individual_predictions = []
187 individual_accuracies = []
188
189 for i, model in enumerate(ensemble_members):
190     model.fit(X_train, y_train)
191     y_pred = model.predict(X_test)
192     accuracy = accuracy_score(y_test, y_pred)
193     individual_predictions.append(y_pred)
194     individual_accuracies.append(accuracy)
195     print(f Model {i+1} ({type(model).__name__}): {accuracy:.4f} )
196
197 # Calculate diversity (disagreement rate)
198 def calculate_diversity(predictions):
199     Calculate average disagreement rate among predictions
200     n_samples = len(predictions[0])
201     disagreements = 0
202
203     for i in range(n_samples):
204         sample_predictions = [pred[i] for pred in predictions]
205         if len(set(sample_predictions)) > 1: # If predictions differ
206             disagreements += 1
207
208     return disagreements / n_samples
209
210 diversity = calculate_diversity(individual_predictions)
211 print(f \nEnsemble diversity (disagreement rate): {diversity:.4f} )
212
213 # Test different ensemble strategies
214 print( \nDifferent Ensemble Strategies: )
215
216 # Simple average

```

```

217 avg_predictions = np.mean(individual_predictions, axis=0)
218 avg_predictions_binary = (avg_predictions > 0.5).astype(int)
219 avg_accuracy = accuracy_score(y_test, avg_predictions_binary)
220 print(f Simple average: {avg_accuracy:.4f} )
221
222 # Majority vote
223 majority_vote = []
224 for i in range(len(y_test)):
225     votes = [pred[i] for pred in individual_predictions]
226     majority_vote.append(1 if sum(votes) > len(votes)/2 else 0)
227 majority_accuracy = accuracy_score(y_test, majority_vote)
228 print(f Majority vote: {majority_accuracy:.4f} )
229
230 # Weighted average (weight by individual performance)
231 weights = np.array(individual_accuracies)
232 weights = weights / np.sum(weights) # Normalize
233 weighted_predictions = np.average(individual_predictions, axis=0, weights=
    weights)
234 weighted_predictions_binary = (weighted_predictions > 0.5).astype(int)
235 weighted_accuracy = accuracy_score(y_test, weighted_predictions_binary)
236 print(f Weighted average: {weighted_accuracy:.4f} )
237
238 print(f \nCounter-intuitive: High diversity doesn't always mean better
    performance! )
239 print( Sometimes diverse models make different types of errors that cancel
    out )
240 print( Sometimes diverse models make correlated errors that compound )
241
242 # Advanced ensemble behavior analysis
243 class EnsembleBehaviorAnalyzer:
244     Analyze counter-intuitive ensemble behaviors
245
246     @staticmethod
247     def analyze_correlation_effects(models, X_train, y_train, X_test,
    y_test):
248         Analyze how model correlation affects ensemble performance
249
250         # Train models
251         predictions = []
252         accuracies = []
253
254         for model in models:
255             model.fit(X_train, y_train)
256             y_pred = model.predict(X_test)
257             predictions.append(y_pred)
258             accuracies.append(accuracy_score(y_test, y_pred))
259
260         # Calculate pairwise correlations
261         correlations = []
262         for i in range(len(predictions)):
263             for j in range(i+1, len(predictions)):
264                 corr = np.corrcoef(predictions[i], predictions[j])[0, 1]
265                 correlations.append(corr)
266

```

```

267     avg_correlation = np.mean(correlations)
268
269     # Ensemble performance
270     ensemble_pred = np.mean(predictions, axis=0)
271     ensemble_pred_binary = (ensemble_pred > 0.5).astype(int)
272     ensemble_accuracy = accuracy_score(y_test, ensemble_pred_binary)
273
274     return {
275         'individual_accuracies': accuracies,
276         'average_correlation': avg_correlation,
277         'ensemble_accuracy': ensemble_accuracy,
278         'correlations': correlations
279     }
280
281     @staticmethod
282     def demonstrate_overfitting_in_ensembles():
283         Demonstrate how ensembles can overfit
284
285         # Create dataset with noise
286         np.random.seed(42)
287         X_noise, y_noise = make_classification(
288             n_samples=500,
289             n_features=100, # High dimensional
290             n_informative=10,
291             n_redundant=20,
292             flip_y=0.3, # High noise
293             random_state=42
294         )
295
296         X_train_n, X_test_n, y_train_n, y_test_n = train_test_split(
297             X_noise, y_noise, test_size=0.3, random_state=42
298         )
299
300         # Test with increasing ensemble size
301         ensemble_sizes = [1, 5, 10, 50, 100, 200]
302         train_scores = []
303         test_scores = []
304
305         for size in ensemble_sizes:
306             rf_overfit = RandomForestClassifier(
307                 n_estimators=size,
308                 max_depth=20, # Deep trees to encourage overfitting
309                 min_samples_split=2,
310                 random_state=42
311             )
312
313             rf_overfit.fit(X_train_n, y_train_n)
314             train_score = rf_overfit.score(X_train_n, y_train_n)
315             test_score = rf_overfit.score(X_test_n, y_test_n)
316
317             train_scores.append(train_score)
318             test_scores.append(test_score)
319
320         print(f Ensemble size {size}: Train={train_score:.4f}, Test={

```



```

    test_score:.4f} )
321
322     # Plot overfitting
323     plt.figure(figsize=(10, 6))
324     plt.plot(ensemble_sizes, train_scores, 'o-', label='Training Score
    ')
325     plt.plot(ensemble_sizes, test_scores, 'o-', label='Test Score')
326     plt.xlabel('Ensemble Size')
327     plt.ylabel('Accuracy')
328     plt.title('Ensemble Overfitting Demonstration')
329     plt.legend()
330     plt.grid(True)
331     plt.show()
332
333     print( \nCounter-intuitive: Larger ensembles can overfit! )
334     print( Reason: More complex models with more parameters )
335
336 # Run advanced analysis
337 print( \n=== ADVANCED ENSEMBLE BEHAVIOR ANALYSIS === )
338
339 # Analyze correlation effects
340 diverse_models = [
341     LogisticRegression(random_state=42),
342     DecisionTreeClassifier(random_state=42),
343     RandomForestClassifier(n_estimators=50, random_state=42),
344     GradientBoostingClassifier(random_state=42, n_estimators=50)
345 ]
346
347 analyzer = EnsembleBehaviorAnalyzer()
348 correlation_analysis = analyzer.analyze_correlation_effects(
349     diverse_models, X_train, y_train, X_test, y_test
350 )
351
352 print( Correlation Analysis Results: )
353 print(f Average model correlation: {correlation_analysis['
    average_correlation']:.4f} )
354 print(f Individual accuracies: {[f'{acc:.4f}' for acc in
    correlation_analysis['individual_accuracies']]})
355 print(f Ensemble accuracy: {correlation_analysis['ensemble_accuracy']:.4f}
    )
356
357 # Demonstrate overfitting
358 print( \nDemonstrating ensemble overfitting: )
359 analyzer.demonstrate_overfitting_in_ensembles()
360
361 print( \nKey Counter-Intuitive Ensemble Behaviors: )
362 print( 1. Ensembles are not always better than the best individual model )
363 print( 2. Bagging can hurt performance on simple problems )
364 print( 3. More estimators lead to diminishing returns )
365 print( 4. High diversity doesn't always improve performance )
366 print( 5. Ensembles can overfit with too many members )
367 print( 6. Correlated errors can reduce ensemble effectiveness )

```

```
368 print( 7. Weighted ensembles may not always outperform simple voting )
```

Listing 49: Behavior 6.1.4: Ensemble Method Surprises

```
1 # Counter-Intuitive: Hyperparameter tuning can have unexpected behaviors
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import (GridSearchCV, RandomizedSearchCV,
5                                     cross_val_score, validation_curve)
6 from sklearn.ensemble import RandomForestClassifier,
7   GradientBoostingClassifier
8 from sklearn.svm import SVC
9 from sklearn.linear_model import LogisticRegression
10 from sklearn.datasets import make_classification
11 from sklearn.model_selection import train_test_split
12 from sklearn.metrics import accuracy_score
13 import matplotlib.pyplot as plt
14 from scipy.stats import randint, uniform
15 import time
16
17 # Create challenging dataset for hyperparameter tuning
18 np.random.seed(42)
19 X, y = make_classification(
20     n_samples=2000,
21     n_features=20,
22     n_informative=15,
23     n_redundant=5,
24     n_clusters_per_class=2,
25     class_sep=1.0,
26     flip_y=0.1,
27     random_state=42
28 )
29 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
30     random_state=42)
31
32 print( Dataset created for hyperparameter tuning experiments )
33 print(f Dataset shape: {X.shape} )
34
35 # Counter-intuitive result 1: Grid search can miss optimal parameters
36 print( \n=== COUNTER-INTUITIVE: Grid Search Can Miss Optimal Parameters
37     === )
38
39 # Define coarse grid
40 param_grid_coarse = {
41     'n_estimators': [10, 50, 100, 200],
42     'max_depth': [3, 5, 7, 10, None],
43     'min_samples_split': [2, 5, 10]
44 }
45
46 # Define fine grid around suspected optimal region
47 param_grid_fine = {
48     'n_estimators': [40, 45, 50, 55, 60],
49     'max_depth': [4, 5, 6],
50     'min_samples_split': [4, 5, 6]
```

```

49 }
50
51 rf_coarse = RandomForestClassifier(random_state=42)
52 rf_fine = RandomForestClassifier(random_state=42)
53
54 # Coarse grid search
55 print( Performing coarse grid search... )
56 start_time = time.time()
57 grid_coarse = GridSearchCV(rf_coarse, param_grid_coarse, cv=5, n_jobs=-1)
58 grid_coarse.fit(X_train, y_train)
59 coarse_time = time.time() - start_time
60
61 print(f Coarse grid search completed in {coarse_time:.2f} seconds )
62 print(f Best parameters (coarse): {grid_coarse.best_params_} )
63 print(f Best score (coarse): {grid_coarse.best_score_:.4f} )
64
65 # Fine grid search around coarse best
66 print( \nPerforming fine grid search... )
67 start_time = time.time()
68 grid_fine = GridSearchCV(rf_fine, param_grid_fine, cv=5, n_jobs=-1)
69 grid_fine.fit(X_train, y_train)
70 fine_time = time.time() - start_time
71
72 print(f Fine grid search completed in {fine_time:.2f} seconds )
73 print(f Best parameters (fine): {grid_fine.best_params_} )
74 print(f Best score (fine): {grid_fine.best_score_:.4f} )
75
76 # Test both on test set
77 coarse_model = grid_coarse.best_estimator_
78 fine_model = grid_fine.best_estimator_
79
80 coarse_test_score = coarse_model.score(X_test, y_test)
81 fine_test_score = fine_model.score(X_test, y_test)
82
83 print(f \nTest set performance: )
84 print(f Coarse grid model: {coarse_test_score:.4f} )
85 print(f Fine grid model: {fine_test_score:.4f} )
86
87 print( \nCounter-intuitive: Fine grid search might not find better
      parameters! )
88 print( Reasons: )
89 print( 1. Coarse grid might have skipped the optimal region )
90 print( 2. Cross-validation variance can mask true performance differences
      )
91 print( 3. The 'optimal' parameters might overfit to validation set )
92
93 # Counter-intuitive result 2: Random search can outperform grid search
94 print( \n=== COUNTER-INTUITIVE: Random Search Can Outperform Grid Search
      === )
95
96 # Parameter distributions for random search
97 param_dist = {
98     'n_estimators': randint(10, 300),
99     'max_depth': randint(3, 20),

```

```

100     'min_samples_split': randint(2, 20),
101     'min_samples_leaf': randint(1, 10)
102 }
103
104 rf_random = RandomForestClassifier(random_state=42)
105
106 # Random search with same budget as grid search
107 print( Performing random search... )
108 start_time = time.time()
109 random_search = RandomizedSearchCV(
110     rf_random, param_dist, n_iter=100, cv=5, n_jobs=-1, random_state=42
111 )
112 random_search.fit(X_train, y_train)
113 random_time = time.time() - start_time
114
115 print(f Random search completed in {random_time:.2f} seconds )
116 print(f Best parameters (random): {random_search.best_params_} )
117 print(f Best score (random): {random_search.best_score_:.4f} )
118
119 # Compare with grid search
120 print(f \nComparison with grid search: )
121 print(f Grid search time: {coarse_time:.2f}s, Score: {grid_coarse.
122     best_score_:.4f} )
123 print(f Random search time: {random_time:.2f}s, Score: {random_search.
124     best_score_:.4f} )
125
126 random_model = random_search.best_estimator_
127 random_test_score = random_model.score(X_test, y_test)
128 print(f Random search test score: {random_test_score:.4f} )
129
130 print( \nCounter-intuitive: Random search can be more efficient than grid
131     search! )
132 print( Reasons: )
133 print( 1. Random search explores parameter space more effectively )
134 print( 2. Grid search can waste evaluations on redundant parameter
135     combinations )
136 print( 3. Random search is less sensitive to parameter scaling )
137
138 # Counter-intuitive result 3: More hyperparameter combinations don't
139     always help
140 print( \n=== COUNTER-INTUITIVE: More Combinations Don't Always Help === )
141
142 # Test with different numbers of parameter combinations
143 param_grids = {
144     'Small': {
145         'n_estimators': [50, 100],
146         'max_depth': [5, 10]
147     },
148     'Medium': {
149         'n_estimators': [25, 50, 100, 200],
150         'max_depth': [3, 5, 7, 10]
151     },
152     'Large': {
153         'n_estimators': [10, 25, 50, 100, 200, 300],

```

```

149         'max_depth': [3, 5, 7, 10, 15, 20, None]
150     }
151 }
152
153 results_complexity = {}
154
155 for grid_name, param_grid in param_grids.items():
156     n_combinations = 1
157     for param_values in param_grid.values():
158         n_combinations *= len(param_values)
159
160     print(f '\nTesting {grid_name} grid ({n_combinations} combinations)...
161         )
162
163     rf_test = RandomForestClassifier(random_state=42)
164     start_time = time.time()
165
166     grid_test = GridSearchCV(rf_test, param_grid, cv=3) # Reduced CV for
167     speed
168     grid_test.fit(X_train[:500], y_train[:500]) # Smaller dataset for
169     demo
170
171     search_time = time.time() - start_time
172     results_complexity[grid_name] = {
173         'combinations': n_combinations,
174         'time': search_time,
175         'best_score': grid_test.best_score_
176     }
177
178     print(f    Combinations: {n_combinations} )
179     print(f    Time: {search_time:.2f}s )
180     print(f    Best score: {grid_test.best_score_:.4f} )
181
182 print( \nComplexity Analysis Results: )
183
184 for name, results in results_complexity.items():
185     print(f {name} Grid: {results['combinations']} combinations,
186           f {results['time']:.2f}s, Score: {results['best_score']:.4f} )
187
188 print( \nCounter-intuitive: More parameter combinations don't always
189       improve results! )
190
191 print( Reasons: )
192 print( 1. Diminishing returns - many combinations are similar )
193 print( 2. Computational cost increases exponentially )
194 print( 3. Overfitting to validation set with too many combinations )
195 print( 4. Some parameters have little impact on performance )
196
197 # Counter-intuitive result 4: Validation curves can be misleading
198 print( \n=== COUNTER-INTUITIVE: Validation Curves Can Be Misleading === )
199
200 # Generate validation curves for different parameters
201 def plot_validation_analysis(model, X, y, param_name, param_range, cv=5):
202     Plot validation curve and analyze behavior
203
204     train_scores, val_scores = validation_curve(

```

```

199     model, X, y, param_name=param_name, param_range=param_range, cv=cv
200 )
201
202     train_mean = np.mean(train_scores, axis=1)
203     train_std = np.std(train_scores, axis=1)
204     val_mean = np.mean(val_scores, axis=1)
205     val_std = np.std(val_scores, axis=1)
206
207     plt.figure(figsize=(10, 6))
208     plt.plot(param_range, train_mean, 'o-', color='blue', label='Training
209 score')
210     plt.fill_between(param_range, train_mean - train_std, train_mean +
211 train_std, alpha=0.1, color='blue')
212
213     plt.plot(param_range, val_mean, 'o-', color='red', label='Cross-
214 validation score')
215     plt.fill_between(param_range, val_mean - val_std, val_mean + val_std,
216 alpha=0.1, color='red')
217
218     plt.xlabel(param_name)
219     plt.ylabel('Score')
220     plt.title(f'Validation Curve for {param_name}')
221     plt.legend()
222     plt.grid(True)
223     plt.show()
224
225     # Find optimal parameter
226     optimal_idx = np.argmax(val_mean)
227     optimal_param = param_range[optimal_idx]
228     optimal_score = val_mean[optimal_idx]
229
230     # Check for overfitting
231     train_at_optimal = train_mean[optimal_idx]
232     gap = train_at_optimal - optimal_score
233
234     print(f Optimal {param_name}: {optimal_param} )
235     print(f Optimal CV score: {optimal_score:.4f} )
236     print(f Training score at optimum: {train_at_optimal:.4f} )
237     print(f Train-CV gap: {gap:.4f} )
238
239     if gap > 0.1:
240         print( Warning: Significant overfitting detected! )
241     elif gap < 0.01:
242         print( Good: Low bias and variance )
243     else:
244         print( Acceptable: Reasonable bias-variance tradeoff )
245
246     return optimal_param, optimal_score
247
248 # Test validation curves
249 print( Analyzing validation curves... )
250
251 # Max depth validation curve
252 print( \n1. Max Depth Analysis: )

```

```

249 optimal_depth, score_depth = plot_validation_analysis(
250     RandomForestClassifier(n_estimators=100, random_state=42),
251     X_train[:1000], y_train[:1000],
252     'max_depth', [3, 5, 7, 10, 15, 20, None]
253 )
254
255 # n_estimators validation curve
256 print( \n2. n_estimators Analysis: )
257 optimal_n_est, score_n_est = plot_validation_analysis(
258     RandomForestClassifier(max_depth=10, random_state=42),
259     X_train[:1000], y_train[:1000],
260     'n_estimators', [10, 50, 100, 200, 300, 500]
261 )
262
263 print( \nCounter-intuitive: Validation curves can suggest suboptimal
      parameters! )
264 print( Reasons: )
265 print( 1. Cross-validation variance can create false peaks )
266 print( 2. Small validation sets may not represent true performance )
267 print( 3. Parameter interactions not captured in single-parameter curves )
268 print( 4. Overfitting to validation set in the curve analysis )
269
270 # Advanced hyperparameter tuning analysis
271 class HyperparameterTuningAnalyzer:
272     Analyze counter-intuitive hyperparameter tuning behaviors
273
274     @staticmethod
275     def analyze_parameter_interactions(model_class, X, y, param_grid, cv
      =5):
276         Analyze interactions between hyperparameters
277
278         # Perform grid search
279         model = model_class()
280         grid_search = GridSearchCV(model, param_grid, cv=cv, n_jobs=-1)
281         grid_search.fit(X, y)
282
283         # Get results
284         results_df = pd.DataFrame(grid_search.cv_results_)
285
286         # Analyze parameter importance
287         param_importance = {}
288         for param in param_grid.keys():
289             param_name = f'param_{param}'
290             if param_name in results_df.columns:
291                 grouped = results_df.groupby(param_name)['mean_test_score'
      ].mean()
292                 param_importance[param] = grouped.std()
293
294         return {
295             'best_params': grid_search.best_params_,
296             'best_score': grid_search.best_score_,
297             'param_importance': param_importance,
298             'results_df': results_df
299         }

```

```

300
301 @staticmethod
302 def demonstrate_local_optima():
303     Demonstrate how local optima can trap optimization
304
305     # Create a parameter landscape with multiple optima
306     np.random.seed(42)
307
308     # Simulate parameter space evaluation
309     n_estimators_range = range(10, 200, 10)
310     max_depth_range = range(3, 20)
311
312     scores = {}
313     true_optimum = (100, 10) # True optimal parameters
314
315     for n_est in n_estimators_range:
316         for max_depth in max_depth_range:
317             # Simulate score with some noise and multiple optima
318             base_score = 0.8 - 0.1 * ((n_est - true_optimum[0])**2 + (
max_depth - true_optimum[1])**2) / 10000
319             noise = np.random.normal(0, 0.02)
320             local_boost = 0.05 * np.exp(-((n_est - 50)**2 + (max_depth
- 5)**2) / 1000) # Local optimum
321             scores[(n_est, max_depth)] = base_score + noise +
local_boost
322
323     # Find local and global optima
324     best_params = max(scores.keys(), key=lambda x: scores[x])
325     best_score = scores[best_params]
326
327     print( Parameter Landscape Analysis: )
328     print(f Global optimum at: n_estimators={true_optimum[0]},
max_depth={true_optimum[1]} )
329     print(f Found optimum at: n_estimators={best_params[0]}, max_depth
={best_params[1]} )
330     print(f Best score: {best_score:.4f} )
331
332     # This demonstrates how optimization can get trapped in local
optima
333     print( \nCounter-intuitive: Optimization can find local rather
than global optima! )
334     print( Reason: Greedy search strategies may miss better regions of
parameter space )
335
336 # Run advanced analysis
337 print( \n=== ADVANCED HYPERPARAMETER TUNING ANALYSIS === )
338
339 # Analyze parameter interactions
340 param_grid_analysis = {
341     'n_estimators': [50, 100, 200],
342     'max_depth': [5, 10, 15],
343     'min_samples_split': [2, 5, 10]
344 }
345

```



```

346 analyzer = HyperparameterTuningAnalyzer()
347 interaction_analysis = analyzer.analyze_parameter_interactions(
348     RandomForestClassifier, X_train[:500], y_train[:500],
349     param_grid_analysis, cv=3
350 )
351 print( Parameter Interaction Analysis: )
352 print(f Best parameters: {interaction_analysis['best_params']} )
353 print(f Best score: {interaction_analysis['best_score']:.4f} )
354 print( Parameter importance (std of CV scores): )
355 for param, importance in interaction_analysis['param_importance'].items():
356     print(f {param}: {importance:.4f} )
357
358 # Demonstrate local optima
359 print( \nDemonstrating local optima in parameter space: )
360 analyzer.demonstrate_local_optima()
361
362 print( \nKey Counter-Intuitive Hyperparameter Tuning Behaviors: )
363 print( 1. Grid search can miss optimal parameters due to coarse
364     discretization )
365 print( 2. Random search can be more efficient than systematic grid search
366     )
367 print( 3. More parameter combinations don't always lead to better results
368     )
369 print( 4. Validation curves can be misleading due to variance or
370     interactions )
371 print( 5. Optimization can get trapped in local optima )
372 print( 6. Cross-validation scores can overfit to validation sets )
373 print( 7. Best parameters on validation may not be best on test set )
374 print( 8. Parameter importance varies by dataset and problem type )

```

Listing 50: Behavior 6.1.5: Hyperparameter Tuning Surprises

```

1 # Counter-Intuitive: Data preprocessing can have unexpected effects
2 import numpy as np
3 import pandas as pd
4 from sklearn.preprocessing import (StandardScaler, MinMaxScaler,
5     RobustScaler,
6     LabelEncoder, OneHotEncoder,
7     PolynomialFeatures)
8 from sklearn.impute import SimpleImputer, KNNImputer
9 from sklearn.model_selection import train_test_split, cross_val_score
10 from sklearn.ensemble import RandomForestClassifier
11 from sklearn.linear_model import LogisticRegression
12 from sklearn.metrics import accuracy_score
13 from sklearn.datasets import make_classification
14 import matplotlib.pyplot as plt
15 from datetime import datetime, timedelta
16
17 # Create complex dataset with various data issues
18 np.random.seed(42)
19
20 # Create dataset with mixed data types and issues
21 n_samples = 1000

```

```

21 # Numerical features with different scales and issues
22 numerical_data = {
23     'small_range': np.random.uniform(0, 1, n_samples),
24     'large_range': np.random.uniform(0, 10000, n_samples),
25     'with_outliers': np.random.normal(0, 1, n_samples),
26     'with_missing': np.random.normal(50, 10, n_samples)
27 }
28
29 # Add outliers to some features
30 outlier_indices = np.random.choice(n_samples, 20, replace=False)
31 numerical_data['with_outliers'][outlier_indices] = np.random.normal(0, 10,
    20)
32
33 # Add missing values
34 missing_indices = np.random.choice(n_samples, 50, replace=False)
35 numerical_data['with_missing'][missing_indices] = np.nan
36
37 # Categorical features
38 categorical_data = {
39     'low_cardinality': np.random.choice(['A', 'B', 'C'], n_samples),
40     'high_cardinality': np.random.choice([f'Cat_{i}' for i in range(50)],
    n_samples),
41     'with_rare_categories': np.random.choice(
42         ['Common1', 'Common2', 'Common3'] + [f'Rare_{i}' for i in range
    (20)],
43         n_samples, p=[0.3, 0.3, 0.3] + [0.1/20]*20
44     )
45 }
46
47 # Date/time features
48 base_date = datetime(2020, 1, 1)
49 date_data = [base_date + timedelta(days=i) for i in range(n_samples)]
50 # Add some missing dates
51 date_missing_indices = np.random.choice(n_samples, 30, replace=False)
52 for idx in date_missing_indices:
53     date_data[idx] = pd.NaT
54
55 # Create target variable based on features
56 target = (
57     numerical_data['small_range'] * 10 +
58     numerical_data['large_range'] / 1000 +
59     (numerical_data['with_outliers'] > 0).astype(int) +
60     (categorical_data['low_cardinality'] == 'A').astype(int)
61 ) > 2
62
63 # Combine into DataFrame
64 data = pd.DataFrame(numerical_data)
65 for col, values in categorical_data.items():
66     data[col] = values
67 data['date_feature'] = date_data
68 data['target'] = target.astype(int)
69
70 print( Complex dataset created with various data issues )
71 print(f Dataset shape: {data.shape} )

```

```

72 print(f Missing values per column: )
73 print(data.isnull().sum())
74
75 # Counter-intuitive result 1: Imputation can introduce bias
76 print( \n=== COUNTER-INTUITIVE: Imputation Can Introduce Bias === )
77
78 # Test different imputation strategies
79 X_impute = data[['small_range', 'large_range', 'with_outliers', '
    with_missing']].copy()
80 y_impute = data['target']
81
82 # Split data
83 X_train_imp, X_test_imp, y_train_imp, y_test_imp = train_test_split(
84     X_impute, y_impute, test_size=0.3, random_state=42
85 )
86
87 imputers = {
88     'Mean': SimpleImputer(strategy='mean'),
89     'Median': SimpleImputer(strategy='median'),
90     'KNN': KNNImputer(n_neighbors=5)
91 }
92
93 results_imputation = {}
94
95 print( Testing different imputation strategies: )
96 for name, imputer in imputers.items():
97     # Fit imputer on training data
98     X_train_imputed = imputer.fit_transform(X_train_imp)
99     X_test_imputed = imputer.transform(X_test_imp)
100
101     # Train model
102     model = LogisticRegression(random_state=42, max_iter=1000)
103     model.fit(X_train_imputed, y_train_imp)
104
105     # Evaluate
106     y_pred = model.predict(X_test_imputed)
107     accuracy = accuracy_score(y_test_imp, y_pred)
108
109     results_imputation[name] = accuracy
110     print(f {name} imputation: {accuracy:.4f} )
111
112 # Show the bias introduced by different imputation methods
113 print(f \nOriginal data statistics for 'with_missing' column: )
114 print(f Mean: {X_train_imp['with_missing'].mean():.4f} )
115 print(f Std: {X_train_imp['with_missing'].std():.4f} )
116
117 # Calculate imputed values
118 for name, imputer in imputers.items():
119     if name == 'Mean':
120         imputed_values = imputer.fit_transform(X_train_imp[:, 3].reshape
(-1, 1))
121         print(f {name} imputed value: {imputed_values[0][0]:.4f} )
122     elif name == 'Median':
123         imputed_values = imputer.fit_transform(X_train_imp[:, 3].reshape

```

```

    (-1, 1))
124     print(f    {name} imputed value: {imputed_values[0][0]:.4f} )
125
126 print( \nCounter-intuitive: Imputation can introduce systematic bias! )
127 print( Reasons: )
128 print( 1. Mean imputation reduces variance )
129 print( 2. Imputed values may not represent true missing data distribution
    )
130 print( 3. Different imputation methods can lead to different model
    performances )
131 print( 4. Imputation assumes data is missing at random (often not true) )
132
133 # Counter-intuitive result 2: One-hot encoding high cardinality features
    can hurt performance
134 print( \n=== COUNTER-INTUITIVE: One-Hot Encoding High Cardinality Features
    === )
135
136 # Test categorical encoding strategies
137 X_cat = data[['low_cardinality', 'high_cardinality', 'with_rare_categories
    ']].copy()
138 y_cat = data['target']
139
140 X_train_cat, X_test_cat, y_train_cat, y_test_cat = train_test_split(
141     X_cat, y_cat, test_size=0.3, random_state=42
142 )
143
144 # Test different encoding approaches
145 print( Testing categorical encoding strategies: )
146
147 # Approach 1: One-hot encode everything
148 ohe_all = OneHotEncoder(sparse=False, handle_unknown='ignore')
149 X_train_ohe = ohe_all.fit_transform(X_train_cat)
150 X_test_ohe = ohe_all.transform(X_test_cat)
151
152 model_ohe = RandomForestClassifier(random_state=42, n_estimators=50)
153 model_ohe.fit(X_train_ohe, y_train_cat)
154 accuracy_ohe = model_ohe.score(X_test_ohe, y_test_cat)
155 print(f    One-hot all features: {accuracy_ohe:.4f} )
156 print(f    Feature dimensions after OHE: {X_train_ohe.shape[1]} )
157
158 # Approach 2: Handle high cardinality differently
159 from sklearn.preprocessing import LabelEncoder
160
161 X_train_smart = X_train_cat.copy()
162 X_test_smart = X_test_cat.copy()
163
164 # Low cardinality: One-hot encode
165 ohe_low = OneHotEncoder(sparse=False, handle_unknown='ignore')
166 low_card_features = ['low_cardinality']
167 X_train_low = ohe_low.fit_transform(X_train_smart[low_card_features])
168 X_test_low = ohe_low.transform(X_test_smart[low_card_features])
169
170 # High cardinality: Label encode or target encode
171 le_high = LabelEncoder()

```

```

172 for col in ['high_cardinality', 'with_rare_categories']:
173     X_train_smart[col] = le_high.fit_transform(X_train_smart[col].astype(
174         str))
175     X_test_smart[col] = le_high.transform(X_test_smart[col].astype(str))
176
177 # Combine features
178 X_train_smart_combined = np.column_stack([
179     X_train_low,
180     X_train_smart[['high_cardinality', 'with_rare_categories']].values
181 ])
182 X_test_smart_combined = np.column_stack([
183     X_test_low,
184     X_test_smart[['high_cardinality', 'with_rare_categories']].values
185 ])
186
187 model_smart = RandomForestClassifier(random_state=42, n_estimators=50)
188 model_smart.fit(X_train_smart_combined, y_train_cat)
189 accuracy_smart = model_smart.score(X_test_smart_combined, y_test_cat)
190
191 print(f Smart encoding: {accuracy_smart:.4f} )
192 print(f Feature dimensions after smart encoding: {X_train_smart_combined
193     .shape[1]} )
194
195 print( \nCounter-intuitive: One-hot encoding high cardinality features can
196     hurt performance! )
197 print( Reasons: )
198 print( 1. Creates too many dimensions (curse of dimensionality) )
199 print( 2. Increases computational complexity )
200 print( 3. Can lead to overfitting )
201 print( 4. Sparse features may not be informative )
202
203 # Counter-intuitive result 3: Feature engineering can decrease performance
204 print( \n=== COUNTER-INTUITIVE: Feature Engineering Can Decrease
205     Performance === )
206
207 # Test feature engineering impact
208 X_fe = data[['small_range', 'large_range', 'with_outliers']].copy()
209 y_fe = data['target']
210
211 X_train_fe, X_test_fe, y_train_fe, y_test_fe = train_test_split(
212     X_fe, y_fe, test_size=0.3, random_state=42
213 )
214
215 # Baseline model
216 model_baseline = LogisticRegression(random_state=42, max_iter=1000)
217 model_baseline.fit(X_train_fe, y_train_fe)
218 baseline_accuracy = model_baseline.score(X_test_fe, y_test_fe)
219 print(f Baseline accuracy: {baseline_accuracy:.4f} )
220
221 # Add polynomial features (can overfit)
222 poly_features = PolynomialFeatures(degree=2, include_bias=False)
223 X_train_poly = poly_features.fit_transform(X_train_fe)
224 X_test_poly = poly_features.transform(X_test_fe)

```

```

222 print(f Feature dimensions increased from {X_train_fe.shape[1]} to {
      X_train_poly.shape[1]} )
223
224 model_poly = LogisticRegression(random_state=42, max_iter=1000)
225 model_poly.fit(X_train_poly, y_train_fe)
226 poly_accuracy = model_poly.score(X_test_poly, y_test_fe)
227 print(f Polynomial features accuracy: {poly_accuracy:.4f} )
228
229 # Add interaction features that don't help
230 X_train_interact = X_train_fe.copy()
231 X_train_interact['useless_interaction'] = (
232     X_train_fe['small_range'] * np.random.normal(0, 1, len(X_train_fe))
233 )
234 X_test_interact = X_test_fe.copy()
235 X_test_interact['useless_interaction'] = (
236     X_test_fe['small_range'] * np.random.normal(0, 1, len(X_test_fe))
237 )
238
239 model_interact = LogisticRegression(random_state=42, max_iter=1000)
240 model_interact.fit(X_train_interact, y_train_fe)
241 interact_accuracy = model_interact.score(X_test_interact, y_test_fe)
242 print(f Useless interaction accuracy: {interact_accuracy:.4f} )
243
244 print( \nCounter-intuitive: Feature engineering can decrease performance!
      )
245 print( Reasons: )
246 print( 1. Overfitting to training data )
247 print( 2. Adding noise instead of signal )
248 print( 3. Increasing dimensionality without adding information )
249 print( 4. Creating correlated features that confuse the model )
250
251 # Counter-intuitive result 4: Scaling can make models worse
252 print( \n=== COUNTER-INTUITIVE: Scaling Can Make Models Worse === )
253
254 # Test scaling impact on different models
255 X_scale = data[['small_range', 'large_range', 'with_outliers']].copy()
256 y_scale = data['target']
257
258 X_train_scale, X_test_scale, y_train_scale, y_test_scale =
      train_test_split(
259     X_scale, y_scale, test_size=0.3, random_state=42
260 )
261
262 models_scale = {
263     'RandomForest': RandomForestClassifier(random_state=42, n_estimators
      =50),
264     'LogisticRegression': LogisticRegression(random_state=42, max_iter
      =1000)
265 }
266
267 scalers_scale = {
268     'No Scaling': None,
269     'StandardScaler': StandardScaler(),
270     'MinMaxScaler': MinMaxScaler(),

```

```

271     'RobustScaler': RobustScaler()
272 }
273
274 results_scaling = {}
275
276 print( Testing scaling impact on different models: )
277 for model_name, model in models_scale.items():
278     results_scaling[model_name] = {}
279
280     for scaler_name, scaler in scalers_scale.items():
281         if scaler is None:
282             X_train_proc = X_train_scale
283             X_test_proc = X_test_scale
284         else:
285             X_train_proc = scaler.fit_transform(X_train_scale)
286             X_test_proc = scaler.transform(X_test_scale)
287
288         model.fit(X_train_proc, y_train_scale)
289         accuracy = model.score(X_test_proc, y_test_scale)
290         results_scaling[model_name][scaler_name] = accuracy
291
292 # Display results
293 for model_name, scaler_results in results_scaling.items():
294     print(f \n{model_name}: )
295     for scaler_name, accuracy in scaler_results.items():
296         print(f     {scaler_name}: {accuracy:.4f} )
297
298 print( \nCounter-intuitive: Scaling can make some models worse! )
299 print( Reasons: )
300 print( 1. Tree-based models are scale-invariant )
301 print( 2. Scaling can obscure feature importance )
302 print( 3. Different scalers handle outliers differently )
303 print( 4. Scaling parameters fitted on wrong data (data leakage) )
304
305 # Advanced preprocessing behavior analysis
306 class PreprocessingBehaviorAnalyzer:
307     Analyze counter-intuitive preprocessing behaviors
308
309     @staticmethod
310     def analyze_missing_data_patterns(X, y):
311         Analyze how different missing data patterns affect imputation
312
313         # Create different missing data patterns
314         patterns = {
315             'MCAR': 'Missing Completely at Random',
316             'MAR': 'Missing at Random',
317             'MNAR': 'Missing Not at Random'
318         }
319
320         results = {}
321
322         for pattern_name, description in patterns.items():
323             print(f \nAnalyzing {pattern_name}: {description} )

```

```

324
325     # Simulate missing data pattern
326     X_missing = X.copy()
327
328     if pattern_name == 'MCAR':
329         # Random missing
330         missing_mask = np.random.random(X.shape) < 0.1
331         X_missing = np.where(missing_mask, np.nan, X_missing)
332
333     elif pattern_name == 'MAR':
334         # Missing related to other features
335         missing_prob = 0.1 * (X[:, 0] > np.median(X[:, 0]))
336         missing_mask = np.random.random(X.shape[0]) < missing_prob
337         X_missing[missing_mask, 1] = np.nan
338
339     elif pattern_name == 'MNAR':
340         # Missing related to missing values themselves
341         missing_mask = X[:, 1] > np.percentile(X[:, 1], 90)
342         X_missing[missing_mask, 1] = np.nan
343
344     # Test imputation
345     imputer = SimpleImputer(strategy='mean')
346     try:
347         X_imputed = imputer.fit_transform(X_missing)
348
349         # Check if imputation introduced bias
350         original_mean = np.nanmean(X[:, 1])
351         imputed_mean = np.mean(X_imputed[:, 1])
352         bias = abs(original_mean - imputed_mean)
353
354         results[pattern_name] = {
355             'bias': bias,
356             'description': description
357         }
358
359         print(f    Original mean: {original_mean:.4f} )
360         print(f    Imputed mean: {imputed_mean:.4f} )
361         print(f    Bias: {bias:.4f} )
362
363     except Exception as e:
364         print(f    Error: {e} )
365         results[pattern_name] = {'error': str(e)}
366
367     return results
368
369     @staticmethod
370     def demonstrate_feature_selection_pitfalls():
371         Demonstrate how feature selection can go wrong
372
373         # Create dataset where feature selection might remove important
374         features
375         np.random.seed(42)
376         n_samples = 1000

```



```

377     # Create features with different importance levels
378     X_fs = np.column_stack([
379         np.random.normal(0, 1, n_samples),      # Low importance
380         np.random.normal(0, 1, n_samples),      # Low importance
381         np.random.normal(0, 0.1, n_samples),    # Very low importance
382         (but necessary interaction)
383         np.random.normal(0, 1, n_samples)      # Medium importance
384     ])
385     # Target depends on interaction that individual feature selection
386     might miss
387     y_fs = ((X_fs[:, 0] * X_fs[:, 2]) + X_fs[:, 3] > 0).astype(int)
388
389     print( Feature Selection Pitfall Demonstration: )
390     print( Target depends on interaction between features 1 and 3 )
391     print( But individually, feature 3 has low variance and might be
392     removed )
393
394     # Individual feature importance
395     from sklearn.feature_selection import f_classif
396     f_scores, p_values = f_classif(X_fs, y_fs)
397
398     print( \nIndividual feature F-scores: )
399     for i, (f_score, p_value) in enumerate(zip(f_scores, p_values)):
400         print(f    Feature {i}: F-score={f_score:.4f}, p-value={p_value
401         :.4f} )
402
403     # This shows how univariate feature selection can miss important
404     interactions
405     print( \nCounter-intuitive: Univariate feature selection can
406     remove important features! )
407     print( Reason: It doesn't consider feature interactions )
408
409 # Run advanced analysis
410 print( \n=== ADVANCED PREPROCESSING BEHAVIOR ANALYSIS === )
411
412 # Analyze missing data patterns
413 X_missing_test, y_missing_test = make_classification(n_samples=500,
414     n_features=5, random_state=42)
415 analyzer = PreprocessingBehaviorAnalyzer()
416 missing_analysis = analyzer.analyze_missing_data_patterns(X_missing_test,
417     y_missing_test)
418
419 print( \nMissing Data Pattern Analysis: )
420 for pattern, results in missing_analysis.items():
421     if 'error' not in results:
422         print(f    {pattern}: Bias={results['bias']:.4f} )
423     else:
424         print(f    {pattern}: Error={results['error']} )
425
426 # Demonstrate feature selection pitfalls
427 analyzer.demonstrate_feature_selection_pitfalls()
428
429 # Demonstrate date feature engineering surprises

```

```

423 print( \n=== DATE FEATURE ENGINEERING SURPRISES === )
424
425 # Create date features and show common pitfalls
426 date_series = pd.date_range('2020-01-01', periods=1000, freq='D')
427
428 # Common date features
429 date_features = pd.DataFrame({
430     'date': date_series,
431     'year': date_series.year,
432     'month': date_series.month,
433     'day': date_series.day,
434     'dayofweek': date_series.dayofweek,
435     'quarter': date_series.quarter,
436     'is_weekend': date_series.dayofweek.isin([5, 6]).astype(int)
437 })
438
439 print( Date feature engineering considerations: )
440 print( 1. Cyclical features (month, dayofweek) should use sine/cosine
         encoding )
441 print( 2. High cardinality (daily dates) can create too many features )
442 print( 3. Seasonal patterns might be better captured with continuous
         features )
443
444 # Show cyclical encoding
445 date_features['month_sin'] = np.sin(2 * np.pi * date_features['month'] /
         12)
446 date_features['month_cos'] = np.cos(2 * np.pi * date_features['month'] /
         12)
447 date_features['dayofweek_sin'] = np.sin(2 * np.pi * date_features['
         dayofweek'] / 7)
448 date_features['dayofweek_cos'] = np.cos(2 * np.pi * date_features['
         dayofweek'] / 7)
449
450 print(f \nCyclical encoding example: )
451 print(f January (1) encoded as: sin={date_features['month_sin'].iloc[0]:.4
         f}, cos={date_features['month_cos'].iloc[0]:.4f} )
452 print(f December (12) encoded as: sin={date_features['month_sin'].iloc
         [11]:.4f}, cos={date_features['month_cos'].iloc[11]:.4f} )
453
454 print( \nKey Counter-Intuitive Preprocessing Behaviors: )
455 print( 1. Imputation can introduce systematic bias )
456 print( 2. One-hot encoding high cardinality features can hurt performance
         )
457 print( 3. Feature engineering can decrease rather than increase
         performance )
458 print( 4. Scaling can make some models worse, not better )
459 print( 5. Different missing data patterns require different imputation
         strategies )
460 print( 6. Univariate feature selection can miss important interactions )
461 print( 7. Cyclical features need special encoding (sine/cosine) )
462 print( 8. Outlier handling methods can significantly impact results )
463 print( 9. Data leakage in preprocessing is easy to accidentally introduce
         )

```

```
464 print( 10. The 'best' preprocessing for one model may hurt another model )
```

Listing 51: Behavior 6.1.6: Data Preprocessing Surprises

```
1 # Counter-Intuitive: Model interpretation can be misleading
2 import numpy as np
3 import pandas as pd
4 from sklearn.ensemble import RandomForestClassifier,
   GradientBoostingClassifier
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import accuracy_score, classification_report
8 from sklearn.datasets import make_classification
9 from sklearn.inspection import permutation_importance
10 import matplotlib.pyplot as plt
11 import shap
12
13 # Create dataset with correlated and interacting features
14 np.random.seed(42)
15
16 # Create base features
17 n_samples = 1000
18 base_feature1 = np.random.normal(0, 1, n_samples)
19 base_feature2 = np.random.normal(0, 1, n_samples)
20
21 # Create correlated features
22 correlated_feature1 = base_feature1 + np.random.normal(0, 0.1, n_samples)
23 correlated_feature2 = base_feature2 + np.random.normal(0, 0.1, n_samples)
24
25 # Create interaction features
26 interaction_feature = base_feature1 * base_feature2
27
28 # Create target with complex relationships
29 target = (
30     base_feature1 * 2 + # Strong effect
31     base_feature2 * 0.5 + # Weak effect
32     interaction_feature * 1.5 + # Strong interaction
33     np.random.normal(0, 0.5, n_samples) # Noise
34 ) > 0
35
36 # Combine into dataset
37 X_interpret = np.column_stack([
38     base_feature1,
39     base_feature2,
40     correlated_feature1,
41     correlated_feature2,
42     interaction_feature
43 ])
44
45 feature_names = [
46     'base_feature1',
47     'base_feature2',
48     'correlated_feature1',
49     'correlated_feature2',
50     'interaction_feature'
```

```

51 ]
52
53 y_interpret = target.astype(int)
54
55 print( Dataset created with correlated and interacting features )
56 print( True importance: base_feature1 > interaction_feature >
        base_feature2 )
57
58 # Split data
59 X_train_int, X_test_int, y_train_int, y_test_int = train_test_split(
60     X_interpret, y_interpret, test_size=0.3, random_state=42
61 )
62
63 # Counter-intuitive result 1: Feature importance can be misleading with
    correlated features
64 print( \n=== COUNTER-INTUITIVE: Feature Importance with Correlated
        Features === )
65
66 # Train Random Forest
67 rf_interpret = RandomForestClassifier(n_estimators=100, random_state=42)
68 rf_interpret.fit(X_train_int, y_train_int)
69
70 # Get feature importances
71 rf_importances = rf_interpret.feature_importances_
72 importance_df = pd.DataFrame({
73     'feature': feature_names,
74     'importance': rf_importances
75 }).sort_values('importance', ascending=False)
76
77 print( Random Forest Feature Importances: )
78 print(importance_df)
79
80 print( \nCounter-intuitive observation: )
81 print( base_feature1 and correlated_feature1 should have similar
        importance )
82 print( But the importance is split between them! )
83 print( The true importance of base_feature1 is diluted. )
84
85 # Compare with permutation importance
86 print( \nPermutation Importances (more reliable): )
87 perm_importance = permutation_importance(
88     rf_interpret, X_test_int, y_test_int, n_repeats=10, random_state=42
89 )
90
91 perm_importance_df = pd.DataFrame({
92     'feature': feature_names,
93     'perm_importance': perm_importance.importances_mean,
94     'perm_std': perm_importance.importances_std
95 }).sort_values('perm_importance', ascending=False)
96
97 print(perm_importance_df)
98
99 # Visualize the difference
100 fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 5))

```

```

101
102 ax1.barh(importance_df['feature'], importance_df['importance'])
103 ax1.set_xlabel('Feature Importance')
104 ax1.set_title('Random Forest Feature Importances\n(Split among correlated
    features)')
105
106 ax2.barh(perm_importance_df['feature'], perm_importance_df['
    perm_importance'])
107 ax2.set_xlabel('Permutation Importance')
108 ax2.set_title('Permutation Importances\n(Shows true predictive power)')
109
110 plt.tight_layout()
111 plt.show()
112
113 # Counter-intuitive result 2: Linear model coefficients can be misleading
114 print( \n=== COUNTER-INTUITIVE: Linear Model Coefficients === )
115
116 # Train Logistic Regression
117 lr_interpret = LogisticRegression(random_state=42, max_iter=1000)
118 lr_interpret.fit(X_train_int, y_train_int)
119
120 # Get coefficients
121 lr_coefficients = lr_interpret.coef_[0]
122 coef_df = pd.DataFrame({
123     'feature': feature_names,
124     'coefficient': lr_coefficients,
125     'abs_coefficient': np.abs(lr_coefficients)
126 }).sort_values('abs_coefficient', ascending=False)
127
128 print( Logistic Regression Coefficients: )
129 print(coef_df)
130
131 print( \nCounter-intuitive observation: )
132 print( Coefficients don't reflect true feature importance due to scale! )
133 print( Features with larger scales appear more important. )
134 print( Correlated features can have opposite signs due to
    multicollinearity. )
135
136 # Compare with standardized coefficients
137 from sklearn.preprocessing import StandardScaler
138
139 scaler_int = StandardScaler()
140 X_train_scaled = scaler_int.fit_transform(X_train_int)
141 X_test_scaled = scaler_int.transform(X_test_int)
142
143 lr_scaled = LogisticRegression(random_state=42, max_iter=1000)
144 lr_scaled.fit(X_train_scaled, y_train_int)
145
146 lr_coefficients_scaled = lr_scaled.coef_[0]
147 coef_scaled_df = pd.DataFrame({
148     'feature': feature_names,
149     'coefficient': lr_coefficients_scaled,
150     'abs_coefficient': np.abs(lr_coefficients_scaled)
151 }).sort_values('abs_coefficient', ascending=False)

```

```

152
153 print( \nStandardized Logistic Regression Coefficients: )
154 print(coef_scaled_df)
155
156 print( \nCounter-intuitive: Even standardized coefficients can be
      misleading! )
157 print( Reason: Multicollinearity still affects coefficient interpretation
      )
158
159 # Counter-intuitive result 3: SHAP values can be counter-intuitive
160 print( \n=== COUNTER-INTUITIVE: SHAP Values === )
161
162 try:
163     # Calculate SHAP values
164     explainer = shap.TreeExplainer(rf_interpret)
165     shap_values = explainer.shap_values(X_test_int[:100]) # Use subset
166     for speed
167
168     # Summary plot
169     shap.summary_plot(shap_values[1], X_test_int[:100], feature_names=
170     feature_names, show=False)
171     plt.title('SHAP Summary Plot')
172     plt.show()
173
174     # Mean absolute SHAP values
175     mean_shap = np.mean(np.abs(shap_values[1]), axis=0)
176     shap_df = pd.DataFrame({
177         'feature': feature_names,
178         'mean_abs_shap': mean_shap
179     }).sort_values('mean_abs_shap', ascending=False)
180
181     print( SHAP Feature Importance: )
182     print(shap_df)
183
184     print( \nCounter-intuitive: SHAP values show individual contributions
185     )
186     print( But global importance might differ from individual explanations
187     )
188
189 except ImportError:
190     print( SHAP not installed. Install with: pip install shap )
191     print( SHAP provides game-theoretic approach to explain model
192     predictions )
193
194 # Counter-intuitive result 4: Partial dependence plots can be misleading
195 print( \n=== COUNTER-INTUITIVE: Partial Dependence Plots === )
196
197 from sklearn.inspection import partial_dependence, plot_partial_dependence
198
199 # Create PDP for features
200 try:
201     features_pdp = [0, 1, 4] # base_feature1, base_feature2,
202     interaction_feature
203     plot_partial_dependence(

```

```

198         rf_interpret, X_train_int, features_pdp,
199         feature_names=feature_names, random_state=42
200     )
201     plt.suptitle('Partial Dependence Plots')
202     plt.show()
203
204     print( Partial Dependence Plots created )
205     print( Counter-intuitive: PDPs show marginal effects )
206     print( But ignore interactions between features )
207     print( For interaction_feature, PDP might show little effect )
208     print( Even though it's actually very important in combination with
209           others )
210
211 except Exception as e:
212     print(f Error creating PDPs: {e} )
213
214 # Counter-intuitive result 5: Model confidence doesn't equal accuracy
215 print( \n=== COUNTER-INTUITIVE: Model Confidence vs Accuracy === )
216
217 # Get prediction probabilities
218 y_proba = rf_interpret.predict_proba(X_test_int)
219
220 # Separate high and low confidence predictions
221 high_confidence_mask = np.max(y_proba, axis=1) > 0.8
222 low_confidence_mask = np.max(y_proba, axis=1) <= 0.8
223
224 y_pred = rf_interpret.predict(X_test_int)
225
226 high_conf_accuracy = accuracy_score(
227     y_test_int[high_confidence_mask],
228     y_pred[high_confidence_mask]
229 )
230
231 low_conf_accuracy = accuracy_score(
232     y_test_int[low_confidence_mask],
233     y_pred[low_confidence_mask]
234 )
235
236 print(f High confidence predictions (>{0.8*100}%): )
237 print(f     Count: {np.sum(high_confidence_mask)} )
238 print(f     Accuracy: {high_conf_accuracy:.4f} )
239
240 print(f Low confidence predictions ( {0.8*100}%): )
241 print(f     Count: {np.sum(low_confidence_mask)} )
242 print(f     Accuracy: {low_conf_accuracy:.4f} )
243
244 print( \nCounter-intuitive: High model confidence doesn't guarantee high
245       accuracy! )
246 print( Reasons: )
247 print( 1. Model can be confidently wrong )
248 print( 2. Calibration issues )
249 print( 3. Distribution shift between training and test data )
250 print( 4. Model overfitting leading to false confidence )
251
252 # Advanced interpretation analysis

```

```

250 class ModelInterpretationAnalyzer:
251     Analyze counter-intuitive model interpretation behaviors
252
253     @staticmethod
254     def analyze_feature_interaction_effects(model, X, y, feature_pairs):
255         Analyze how feature interactions affect interpretation
256
257         # Train model
258         model.fit(X, y)
259
260         # Individual feature importance
261         if hasattr(model, 'feature_importances_'):
262             individual_importance = model.feature_importances_
263         else:
264             # Use permutation importance
265             perm_imp = permutation_importance(model, X, y, n_repeats=5,
266 random_state=42)
267             individual_importance = perm_imp.importances_mean
268
269         # Create interaction features and retrain
270         X_with_interactions = X.copy()
271         interaction_names = []
272
273         for i, (f1, f2) in enumerate(feature_pairs):
274             interaction_name = f'interaction_{f1}_{f2}'
275             X_with_interactions = np.column_stack([
276                 X_with_interactions,
277                 X[:, f1] * X[:, f2]
278             ])
279             interaction_names.append(interaction_name)
280
281         # Retrain with interactions
282         model_interactions = type(model)(**model.get_params())
283         model_interactions.fit(X_with_interactions, y)
284
285         # Compare importances
286         if hasattr(model_interactions, 'feature_importances_'):
287             interaction_importance = model_interactions.
288 feature_importances_
289         else:
290             perm_imp_int = permutation_importance(
291                 model_interactions, X_with_interactions, y, n_repeats=5,
292 random_state=42
293             )
294             interaction_importance = perm_imp_int.importances_mean
295
296         return {
297             'individual_importance': individual_importance,
298             'interaction_importance': interaction_importance,
299             'interaction_names': interaction_names
300         }
301
302     @staticmethod
303     def demonstrate_calibration_issues():

```



```

301         Demonstrate model calibration problems
302
303     # Create dataset where model confidence is misleading
304     np.random.seed(42)
305     X_cal, y_cal = make_classification(
306         n_samples=1000, n_features=10, n_informative=5,
307         n_redundant=2, flip_y=0.3, random_state=42
308     )
309
310     # Train overconfident model
311     rf_overconfident = RandomForestClassifier(
312         n_estimators=50, max_depth=20, random_state=42
313     )
314     rf_overconfident.fit(X_cal[:800], y_cal[:800])
315
316     # Get probabilities
317     y_proba_cal = rf_overconfident.predict_proba(X_cal[800:])
318     y_pred_cal = rf_overconfident.predict(X_cal[800:])
319     y_true_cal = y_cal[800:]
320
321     # Analyze calibration
322     from sklearn.calibration import calibration_curve
323
324     fraction_of_positives, mean_predicted_value = calibration_curve(
325         y_true_cal, y_proba_cal[:, 1], n_bins=10
326     )
327
328     # Plot calibration curve
329     plt.figure(figsize=(8, 6))
330     plt.plot(mean_predicted_value, fraction_of_positives, 's-', label=
Random Forest )
331     plt.plot([0, 1], [0, 1], 'k:', label= Perfectly calibrated )
332     plt.xlabel( Mean Predicted Probability )
333     plt.ylabel( Fraction of Positives )
334     plt.title( Calibration Plot - Overconfident Model )
335     plt.legend()
336     plt.grid(True)
337     plt.show()
338
339     print( Calibration Analysis: )
340     print( Perfect calibration: points should follow diagonal line )
341     print( Overconfident model: points below diagonal (predictions too
high) )
342     print( Underconfident model: points above diagonal (predictions
too low) )
343
344     # Calculate calibration error
345     calibration_error = np.mean(np.abs(fraction_of_positives -
mean_predicted_value))
346     print(f Calibration error: {calibration_error:.4f} )
347
348     print( \nCounter-intuitive: Well-calibrated models are often less
confident! )
349     print( Reason: They honestly represent uncertainty in predictions

```

```

    )
350
351 # Run advanced analysis
352 print( \n=== ADVANCED MODEL INTERPRETATION ANALYSIS === )
353
354 # Analyze feature interactions
355 feature_pairs = [(0, 1)] # base_feature1 and base_feature2
356 analyzer = ModelInterpretationAnalyzer()
357
358 interaction_analysis = analyzer.analyze_feature_interaction_effects(
359     RandomForestClassifier(n_estimators=50, random_state=42),
360     X_train_int, y_train_int, feature_pairs
361 )
362
363 print( Feature Interaction Analysis: )
364 print( Individual feature importance may miss important interactions )
365 print( Adding explicit interaction features can reveal hidden
        relationships )
366
367 # Demonstrate calibration issues
368 print( \nDemonstrating calibration issues: )
369 analyzer.demonstrate_calibration_issues()
370
371 # Demonstrate model-specific interpretation surprises
372 print( \n=== MODEL-SPECIFIC INTERPRETATION SURPRISES === )
373
374 # Tree-based models
375 print( Tree-based models (Random Forest, XGBoost): )
376 print( 1. Feature importance splits importance among correlated features )
377 print( 2. Deep trees can memorize training data, making importance
        misleading )
378 print( 3. Individual trees can have very different feature importance )
379
380 # Linear models
381 print( \nLinear models (Logistic Regression, Linear SVM): )
382 print( 1. Coefficients reflect scale, not necessarily importance )
383 print( 2. Multicollinearity can flip coefficient signs )
384 print( 3. Regularization can shrink important features to zero )
385
386 # Neural networks
387 print( \nNeural networks: )
388 print( 1. Individual weights are not interpretable )
389 print( 2. Feature importance methods are approximations )
390 print( 3. Hidden layer activations are difficult to interpret )
391
392 # Ensemble models
393 print( \nEnsemble models: )
394 print( 1. Ensemble importance is average of individual model importances )
395 print( 2. Different ensemble members may disagree on feature importance )
396 print( 3. Voting can obscure individual model insights )
397
398 print( \nKey Counter-Intuitive Model Interpretation Behaviors: )
399 print( 1. Feature importance splits among correlated features )
400 print( 2. Linear coefficients reflect scale, not importance )

```

```

401 print( 3. SHAP values show individual contributions, not global importance
        )
402 print( 4. Partial dependence ignores feature interactions )
403 print( 5. High confidence doesn't equal high accuracy )
404 print( 6. Models can be confidently wrong )
405 print( 7. Calibration and discrimination are different concepts )
406 print( 8. Individual model insights may not apply to ensembles )
407 print( 9. Feature importance can change with data distribution )
408 print( 10. The most important feature may not be the most actionable )

```

Listing 52: Behavior 6.1.7: Model Interpretation Surprises

```

1 # Counter-Intuitive: Evaluation metrics can be misleading
2 import numpy as np
3 import pandas as pd
4 from sklearn.metrics import (accuracy_score, precision_score, recall_score
    , f1_score,
5                               roc_auc_score, precision_recall_curve,
6                               roc_curve,
7                               confusion_matrix, classification_report)
8 from sklearn.model_selection import train_test_split
9 from sklearn.ensemble import RandomForestClassifier
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.datasets import make_classification
12 import matplotlib.pyplot as plt
13 from sklearn.calibration import calibration_curve
14
15 # Create imbalanced dataset to demonstrate metric surprises
16 np.random.seed(42)
17
18 # Highly imbalanced dataset (95% class 0, 5% class 1)
19 X_imb, y_imb = make_classification(
20     n_samples=2000,
21     n_features=10,
22     n_classes=2,
23     weights=[0.95, 0.05], # 95% class 0, 5% class 1
24     random_state=42
25 )
26
27 X_train_imb, X_test_imb, y_train_imb, y_test_imb = train_test_split(
28     X_imb, y_imb, test_size=0.3, random_state=42
29 )
30
31 print( Highly imbalanced dataset created )
32 print(f Class distribution: {np.bincount(y_imb)} )
33 print(f Class 0: {np.sum(y_imb == 0) / len(y_imb) * 100:.1f}% )
34 print(f Class 1: {np.sum(y_imb == 1) / len(y_imb) * 100:.1f}% )
35
36 # Train models
37 rf_imb = RandomForestClassifier(random_state=42, n_estimators=100)
38 lr_imb = LogisticRegression(random_state=42, max_iter=1000)
39
40 rf_imb.fit(X_train_imb, y_train_imb)
41 lr_imb.fit(X_train_imb, y_train_imb)

```

```

42 # Get predictions
43 y_pred_rf = rf_imb.predict(X_test_imb)
44 y_pred_lr = lr_imb.predict(X_test_imb)
45
46 y_proba_rf = rf_imb.predict_proba(X_test_imb)[: , 1]
47 y_proba_lr = lr_imb.predict_proba(X_test_imb)[: , 1]
48
49 # Counter-intuitive result 1: Accuracy can be misleading for imbalanced
    data
50 print( \n=== COUNTER-INTUITIVE: Accuracy on Imbalanced Data === )
51
52 # Naive classifier that always predicts majority class
53 y_naive = np.zeros(len(y_test_imb)) # Always predict class 0
54
55 # Calculate metrics
56 accuracy_naive = accuracy_score(y_test_imb, y_naive)
57 accuracy_rf = accuracy_score(y_test_imb, y_pred_rf)
58 accuracy_lr = accuracy_score(y_test_imb, y_pred_lr)
59
60 print( Accuracy Comparison: )
61 print(f Naive classifier (always 0): {accuracy_naive:.4f} )
62 print(f Random Forest: {accuracy_rf:.4f} )
63 print(f Logistic Regression: {accuracy_lr:.4f} )
64
65 print( \nCounter-intuitive: High accuracy doesn't mean good performance! )
66 print( Naive classifier achieves 95% accuracy by always predicting
    majority class )
67 print( But it completely fails to identify the minority class (class 1) )
68
69 # Show detailed confusion matrices
70 print( \nConfusion Matrices: )
71
72 cm_naive = confusion_matrix(y_test_imb, y_naive)
73 cm_rf = confusion_matrix(y_test_imb, y_pred_rf)
74 cm_lr = confusion_matrix(y_test_imb, y_pred_lr)
75
76 print(f Naive classifier:\n{cm_naive} )
77 print(f Random Forest:\n{cm_rf} )
78 print(f Logistic Regression:\n{cm_lr} )
79
80 # Counter-intuitive result 2: F1-score can be misleading without context
81 print( \n=== COUNTER-INTUITIVE: F1-Score Without Context === )
82
83 # Calculate detailed metrics
84 def detailed_metrics(y_true, y_pred, y_proba=None, model_name= ):
85     Calculate comprehensive metrics
86
87     metrics = {
88         'accuracy': accuracy_score(y_true, y_pred),
89         'precision': precision_score(y_true, y_pred, zero_division=0),
90         'recall': recall_score(y_true, y_pred, zero_division=0),
91         'f1': f1_score(y_true, y_pred, zero_division=0)
92     }
93

```

```

94     if y_proba is not None:
95         metrics['auc'] = roc_auc_score(y_true, y_proba)
96
97     # Per-class metrics
98     precision_per_class = precision_score(y_true, y_pred, average=None,
99     zero_division=0)
100     recall_per_class = recall_score(y_true, y_pred, average=None,
101     zero_division=0)
102
103     metrics['class_0_precision'] = precision_per_class[0]
104     metrics['class_1_precision'] = precision_per_class[1]
105     metrics['class_0_recall'] = recall_per_class[0]
106     metrics['class_1_recall'] = recall_per_class[1]
107
108     return metrics
109
110 # Calculate metrics for all models
111 naive_metrics = detailed_metrics(y_test_imb, y_naive, model_name= Naive )
112 rf_metrics = detailed_metrics(y_test_imb, y_pred_rf, y_proba_rf,
113     model_name= Random Forest )
114 lr_metrics = detailed_metrics(y_test_imb, y_pred_lr, y_proba_lr,
115     model_name= Logistic Regression )
116
117 print( Detailed Metrics Comparison: )
118 models_metrics = [
119     ( Naive , naive_metrics),
120     ( Random Forest , rf_metrics),
121     ( Logistic Regression , lr_metrics)
122 ]
123
124 for model_name, metrics in models_metrics:
125     print(f \n{model_name}: )
126     print(f     Accuracy: {metrics['accuracy']:.4f} )
127     print(f     F1-Score: {metrics['f1']:.4f} )
128     print(f     Class 1 Precision: {metrics['class_1_precision']:.4f} )
129     print(f     Class 1 Recall: {metrics['class_1_recall']:.4f} )
130     if 'auc' in metrics:
131         print(f     AUC: {metrics['auc']:.4f} )
132
133 print( \nCounter-intuitive: F1-score can hide poor minority class
134     performance! )
135 print( A model with high overall F1 might have terrible recall for
136     minority class )
137 print( Context matters: What's more important - precision or recall? )
138
139 # Counter-intuitive result 3: ROC-AUC can be misleading for imbalanced
140     data
141 print( \n=== COUNTER-INTUITIVE: ROC-AUC for Imbalanced Data === )
142
143 # Create ROC curves
144 fpr_rf, tpr_rf, _ = roc_curve(y_test_imb, y_proba_rf)
145 fpr_lr, tpr_lr, _ = roc_curve(y_test_imb, y_proba_lr)
146 fpr_naive, tpr_naive, _ = roc_curve(y_test_imb, np.zeros_like(y_test_imb))
147

```

```

141 plt.figure(figsize=(10, 6))
142 plt.plot(fpr_rf, tpr_rf, label=f'Random Forest (AUC = {rf_metrics[ auc
    ]:.4f})')
143 plt.plot(fpr_lr, tpr_lr, label=f'Logistic Regression (AUC = {lr_metrics[
    auc ]:.4f})')
144 plt.plot(fpr_naive, tpr_naive, label='Naive Classifier')
145 plt.plot([0, 1], [0, 1], 'k--', label='Random Classifier')
146 plt.xlabel('False Positive Rate')
147 plt.ylabel('True Positive Rate')
148 plt.title('ROC Curves - Imbalanced Dataset')
149 plt.legend()
150 plt.grid(True)
151 plt.show()
152
153 print( ROC-AUC Analysis: )
154 print( High ROC-AUC doesn't mean good practical performance on imbalanced
    data )
155 print( ROC-AUC is insensitive to class distribution )
156 print( A model can have high AUC but poor precision on minority class )
157
158 # Counter-intuitive result 4: Precision-Recall curves tell a different
    story
159 print( \n=== COUNTER-INTUITIVE: Precision-Recall vs ROC === )
160
161 # Create Precision-Recall curves
162 precision_rf, recall_rf, _ = precision_recall_curve(y_test_imb, y_proba_rf
    )
163 precision_lr, recall_lr, _ = precision_recall_curve(y_test_imb, y_proba_lr
    )
164
165 # Calculate average precision
166 from sklearn.metrics import average_precision_score
167 ap_rf = average_precision_score(y_test_imb, y_proba_rf)
168 ap_lr = average_precision_score(y_test_imb, y_proba_lr)
169
170 plt.figure(figsize=(10, 6))
171 plt.plot(recall_rf, precision_rf, label=f'Random Forest (AP = {ap_rf:.4f})
    ')
172 plt.plot(recall_lr, precision_lr, label=f'Logistic Regression (AP = {ap_lr
    :.4f})')
173 plt.xlabel('Recall')
174 plt.ylabel('Precision')
175 plt.title('Precision-Recall Curves - Imbalanced Dataset')
176 plt.legend()
177 plt.grid(True)
178 plt.show()
179
180 print( Precision-Recall Analysis: )
181 print(f Random Forest Average Precision: {ap_rf:.4f} )
182 print(f Logistic Regression Average Precision: {ap_lr:.4f} )
183 print( AP is more informative than AUC for imbalanced datasets )
184 print( PR curves focus on minority class performance )
185
186 # Counter-intuitive result 5: Threshold selection can dramatically change

```

```

results
187 print( \n=== COUNTER-INTUITIVE: Threshold Selection Impact === )
188
189 def evaluate_threshold(y_true, y_proba, threshold):
190     Evaluate model at specific threshold
191     y_pred_thresh = (y_proba >= threshold).astype(int)
192     return {
193         'threshold': threshold,
194         'accuracy': accuracy_score(y_true, y_pred_thresh),
195         'precision': precision_score(y_true, y_pred_thresh, zero_division
=0),
196         'recall': recall_score(y_true, y_pred_thresh, zero_division=0),
197         'f1': f1_score(y_true, y_pred_thresh, zero_division=0)
198     }
199
200 # Test different thresholds
201 thresholds = [0.1, 0.3, 0.5, 0.7, 0.9]
202 threshold_results = []
203
204 print( Impact of Different Classification Thresholds: )
205 print( Random Forest with probability threshold: )
206
207 for thresh in thresholds:
208     results = evaluate_threshold(y_test_imb, y_proba_rf, thresh)
209     threshold_results.append(results)
210     print(f Threshold {thresh}: Acc={results['accuracy']:.4f},
211           f Precision={results['precision']:.4f}, Recall={results['recall
']:.4f},
212           f F1={results['f1']:.4f} )
213
214 # Find optimal thresholds
215 def find_optimal_threshold(y_true, y_proba, metric='f1'):
216     Find threshold that optimizes specific metric
217     thresholds_search = np.arange(0.01, 1.0, 0.01)
218     best_score = 0
219     best_threshold = 0.5
220
221     for thresh in thresholds_search:
222         y_pred_thresh = (y_proba >= thresh).astype(int)
223         if metric == 'f1':
224             score = f1_score(y_true, y_pred_thresh, zero_division=0)
225         elif metric == 'precision':
226             score = precision_score(y_true, y_pred_thresh, zero_division
=0)
227         elif metric == 'recall':
228             score = recall_score(y_true, y_pred_thresh, zero_division=0)
229
230         if score > best_score:
231             best_score = score
232             best_threshold = thresh
233
234     return best_threshold, best_score
235
236 opt_f1_thresh, opt_f1_score = find_optimal_threshold(y_test_imb,

```

```

    y_proba_rf, 'f1')
237 opt_prec_thresh, opt_prec_score = find_optimal_threshold(y_test_imb,
    y_proba_rf, 'precision')
238 opt_rec_thresh, opt_rec_score = find_optimal_threshold(y_test_imb,
    y_proba_rf, 'recall')
239
240 print(f \nOptimal thresholds: )
241 print(f    F1-optimized: {opt_f1_thresh:.3f} (F1 = {opt_f1_score:.4f}) )
242 print(f    Precision-optimized: {opt_prec_thresh:.3f} (Precision = {
    opt_prec_score:.4f}) )
243 print(f    Recall-optimized: {opt_rec_thresh:.3f} (Recall = {opt_rec_score
    :.4f}) )
244
245 print( \nCounter-intuitive: Default 0.5 threshold is often suboptimal! )
246 print( The 'best' threshold depends on your specific requirements )
247 print( Optimizing for F1 might give different threshold than optimizing
    for precision )
248
249 # Counter-intuitive result 6: Calibration affects reliability of
    probabilities
250 print( \n=== COUNTER-INTUITIVE: Model Calibration === )
251
252 # Check calibration
253 fraction_positives_rf, mean_pred_rf = calibration_curve(
254     y_test_imb, y_proba_rf, n_bins=10
255 )
256 fraction_positives_lr, mean_pred_lr = calibration_curve(
257     y_test_imb, y_proba_lr, n_bins=10
258 )
259
260 plt.figure(figsize=(10, 6))
261 plt.plot(mean_pred_rf, fraction_positives_rf, s- , label= Random Forest )
262 plt.plot(mean_pred_lr, fraction_positives_lr, s- , label= Logistic
    Regression )
263 plt.plot([0, 1], [0, 1], k: , label= Perfectly calibrated )
264 plt.xlabel( Mean Predicted Probability )
265 plt.ylabel( Fraction of Positives )
266 plt.title( Calibration Curves )
267 plt.legend()
268 plt.grid(True)
269 plt.show()
270
271 # Calculate calibration error
272 def calibration_error(fraction_positives, mean_predicted):
273     Calculate mean calibration error
274     return np.mean(np.abs(fraction_positives - mean_predicted))
275
276 cal_error_rf = calibration_error(fraction_positives_rf, mean_pred_rf)
277 cal_error_lr = calibration_error(fraction_positives_lr, mean_pred_lr)
278
279 print( Calibration Analysis: )
280 print(f Random Forest calibration error: {cal_error_rf:.4f} )
281 print(f Logistic Regression calibration error: {cal_error_lr:.4f} )
282

```



```

283 print( \nCounter-intuitive: Well-calibrated models may have lower
      discrimination! )
284 print( A perfectly calibrated model might be less confident but more
      reliable )
285 print( Overconfident models can have better metrics but poor real-world
      performance )
286
287 # Advanced evaluation analysis
288 class EvaluationAnalyzer:
289     Analyze counter-intuitive evaluation behaviors
290
291     @staticmethod
292     def analyze_cost_sensitive_metrics(y_true, y_pred_proba, cost_matrix):
293         Analyze cost-sensitive evaluation
294
295         # Cost matrix: [TN, FP, FN, TP costs]
296         tn_cost, fp_cost, fn_cost, tp_cost = cost_matrix
297
298         # Calculate expected costs for different thresholds
299         thresholds = np.arange(0.01, 1.0, 0.01)
300         costs = []
301
302         for thresh in thresholds:
303             y_pred = (y_pred_proba >= thresh).astype(int)
304             cm = confusion_matrix(y_true, y_pred)
305
306             if cm.shape == (2, 2):
307                 tn, fp, fn, tp = cm.ravel()
308             else:
309                 # Handle edge cases
310                 tn, fp, fn, tp = 0, 0, 0, 0
311                 if len(np.unique(y_true)) == 1:
312                     if np.unique(y_true)[0] == 0:
313                         tn = len(y_true)
314                     else:
315                         tp = len(y_true)
316
317             total_cost = tn * tn_cost + fp * fp_cost + fn * fn_cost + tp *
tp_cost
318             costs.append(total_cost)
319
320             optimal_idx = np.argmin(costs)
321             optimal_threshold = thresholds[optimal_idx]
322             optimal_cost = costs[optimal_idx]
323
324         return {
325             'optimal_threshold': optimal_threshold,
326             'optimal_cost': optimal_cost,
327             'costs': list(zip(thresholds, costs))
328         }
329
330     @staticmethod
331     def demonstrate_metric_sensitivity():
332         Demonstrate how small changes can dramatically affect metrics

```

```

333
334     # Create scenario where small changes have big metric impacts
335     np.random.seed(42)
336
337     # Create predictions where one wrong prediction flips F1 score
338     y_true_demo = np.array([0, 0, 0, 0, 1, 1, 1, 1]) # 4 of each
class
339     y_pred_good = np.array([0, 0, 0, 1, 1, 1, 1, 1]) # 1 FP, 0 FN
340     y_pred_bad = np.array([0, 0, 1, 1, 0, 1, 1, 1]) # 2 FP, 1 FN (1
more wrong)
341
342     metrics_good = {
343         'accuracy': accuracy_score(y_true_demo, y_pred_good),
344         'precision': precision_score(y_true_demo, y_pred_good,
zero_division=0),
345         'recall': recall_score(y_true_demo, y_pred_good, zero_division
=0),
346         'f1': f1_score(y_true_demo, y_pred_good, zero_division=0)
347     }
348
349     metrics_bad = {
350         'accuracy': accuracy_score(y_true_demo, y_pred_bad),
351         'precision': precision_score(y_true_demo, y_pred_bad,
zero_division=0),
352         'recall': recall_score(y_true_demo, y_pred_bad, zero_division
=0),
353         'f1': f1_score(y_true_demo, y_pred_bad, zero_division=0)
354     }
355
356     print( Metric Sensitivity Demonstration: )
357     print( Small change (1 additional error) causing large metric
differences: )
358     print(f Good predictions: {y_pred_good} )
359     print(f Bad predictions: {y_pred_bad} )
360     print( \nMetrics comparison: )
361     for metric in metrics_good.keys():
362         diff = metrics_bad[metric] - metrics_good[metric]
363         print(f {metric}: Good={metrics_good[metric]:.4f}, Bad={
metrics_bad[metric]:.4f}, Diff={diff:.4f} )
364
365     print( \nCounter-intuitive: Small absolute changes can cause large
relative metric changes! )
366     print( Especially problematic with small datasets or critical
thresholds )
367
368 # Run advanced analysis
369 print( \n=== ADVANCED EVALUATION ANALYSIS === )
370
371 # Cost-sensitive analysis
372 print( Cost-Sensitive Evaluation: )
373 # Example: Medical diagnosis where false negatives are very costly
374 cost_matrix = [0, 10, 1000, 0] # TN=0, FP=10, FN=1000, TP=0
375 analyzer = EvaluationAnalyzer()

```

```

376
377 cost_analysis = analyzer.analyze_cost_sensitive_metrics(
378     y_test_imb, y_proba_rf, cost_matrix
379 )
380
381 print(f Optimal threshold for cost matrix {cost_matrix}: {cost_analysis['
    optimal_threshold']:.3f} )
382 print(f Minimum expected cost: {cost_analysis['optimal_cost']:.2f} )
383
384 # Demonstrate metric sensitivity
385 print( \nDemonstrating metric sensitivity: )
386 analyzer.demonstrate_metric_sensitivity()
387
388 # Demonstrate different evaluation scenarios
389 print( \n=== DIFFERENT EVALUATION SCENARIOS === )
390
391 # Scenario 1: Medical diagnosis (high cost of false negatives)
392 print( Scenario 1: Medical Diagnosis )
393 print( High cost of missing actual cases (false negatives) )
394 print( Prioritize recall over precision )
395
396 # Scenario 2: Spam detection (high cost of false positives)
397 print( \nScenario 2: Spam Detection )
398 print( High cost of marking legitimate emails as spam (false positives) )
399 print( Prioritize precision over recall )
400
401 # Scenario 3: Recommendation system (balanced approach)
402 print( \nScenario 3: Recommendation System )
403 print( Balanced cost of false positives and false negatives )
404 print( Optimize for F1-score or AUC )
405
406 print( \nKey Counter-Intuitive Evaluation Behaviors: )
407 print( 1. Accuracy is misleading for imbalanced datasets )
408 print( 2. F1-score can hide poor performance on minority class )
409 print( 3. ROC-AUC can be overly optimistic for imbalanced data )
410 print( 4. Precision-Recall curves are more informative for imbalanced data
    )
411 print( 5. Default 0.5 threshold is often suboptimal )
412 print( 6. Model calibration affects reliability of probability estimates )
413 print( 7. Small changes in predictions can cause large metric fluctuations
    )
414 print( 8. The 'best' metric depends on the specific business context )
415 print( 9. Cost-sensitive evaluation can reveal different optimal models )
416 print( 10. Cross-validation scores can be misleading without proper
    stratification )

```

Listing 53: Behavior 6.1.8: Evaluation Metric Surprises