

# MoDeVa Guide to Data Operations

July 5, 2025

# Tutorial Outline

- 1 Introduction to MoDeVa
- 2 Data Loading
- 3 Data Summary
- 4 Data Preprocessing
- 5 Data Preparation
- 6 Data Registration
- 7 Complete Workflow Example

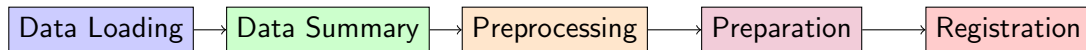
# MoDeVa Data Preparation

- Python library for data operations for MoDeVa
- Built around the central DataSet class
- Supports end-to-end ML data workflows

## Key Features

- Data loading (built-in and external datasets)
- Data summarization
- Preprocessing pipelines
- Dataset registration and management
- MLflow integration for experiment tracking

# MoDeVa Data Workflow Overview



## Central Concept

All operations are performed through the DataSet class, which maintains data state and provides method chaining capabilities.

# Data Loading Overview

## Built-in Datasets:

- BikeSharing (Regression)
- CaliforniaHousing (Regression)
- SimuCredit (Classification)
- TaiwanCredit (Classification)

## External Data Sources:

- CSV files (`load_csv`)
- Pandas DataFrames (`load_dataframe`)
- Spark DataFrames (`load_spark`)

## Usage Pattern

- 1 Create DataSet instance
- 2 Load data using appropriate method
- 3 Data becomes available in `ds.data`

# Loading Built-in Datasets

```
1 # Create DataSet instance
2 from modeva import DataSet
3 ds = DataSet()
4
5 # Load built-in dataset
6 ds.load("SimuCredit")
7
8 # Explore the data
9 ds.data.head(5)
```

## Available Built-in Datasets

- "BikeSharing" - Bike rental prediction (regression)
- "CaliforniaHousing" - Housing prices (regression)
- "SimuCredit" - Credit risk simulation (classification)
- "TaiwanCredit" - Credit default prediction (classification)

# Loading External Data

```
1 # Loading from external sources
2 import pandas as pd
3 from sklearn.datasets import load_iris
4 from modeva import DataSet
5 # Load and prepare external data
6 iris = load_iris()
7 df = pd.DataFrame(data=iris.data, columns=iris.feature_names)
8 df['species'] = pd.Categorical.from_codes(iris.target, iris.target_names)
9
10 # Create named DataSet and load DataFrame
11 ds = DataSet(name="IrisData")
12 ds.load_dataframe(df)
```

## Loading Methods

- `load_csv(filepath)` - Direct CSV loading
- `load_dataframe(df)` - From pandas DataFrame
- `load_spark(spark_df)` - From Spark DataFrame

# Data Summary Components

## Summary Structure

- ① **Overall Summary** (`res.table["summary"]`)
  - Sample count, feature types, duplicates
  - Missing and infinite value statistics
  
- ② **Categorical Variables** (`res.table["categorical"]`)
  - Missing values, unique counts
  - Top 2 value frequencies
  
- ③ **Numerical Variables** (`res.table["numerical"]`)
  - Descriptive statistics (mean, std, percentiles)
  - Missing/infinite value counts



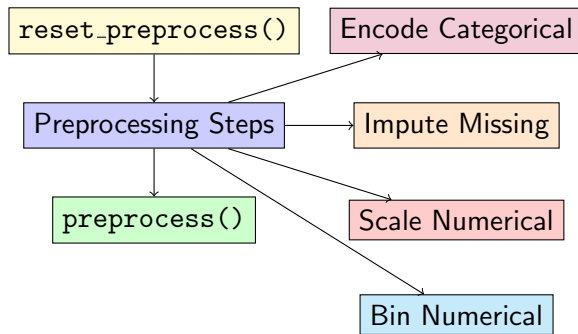
# Generating Data Summary

```
1 # Generate comprehensive summary
2 res = ds.summary()
3 # Access overall dataset summary
4 overall_summary = res.table["summary"]
5 print(overall_summary)
6 # Access categorical variable statistics
7 categorical_stats = res.table["categorical"]
8 print(categorical_stats)
9 # Access numerical variable statistics
10 numerical_stats = res.table["numerical"]
11 print(numerical_stats)
```

## Key Insight

The summary provides a comprehensive view of data quality issues, distribution characteristics, and variable types before preprocessing.

# Preprocessing Pipeline Architecture



## Preprocessing Workflow

- 1 Initialize preprocessing pipeline
- 2 Define preprocessing steps
- 3 Execute all steps in sequence

# Handling Missing Values

```
1 # Initialize preprocessing
2 ds.reset_preprocess()
3 # Impute numerical features with mean and add indicators
4 ds.impute_missing(
5     features=ds.feature_names_numerical,
6     method='mean', add_indicators=True
7 )
8 # Impute categorical features with most frequent value
9 ds.impute_missing(
10     features=ds.feature_names_categorical,
11     method='most_frequent', add_indicators=True
12 )
13 # Handle mixed features with special values
14 ds.impute_missing(
15     features=ds.feature_names_mixed, method='median',
16     add_indicators=True, special_values=["SV1", "SV2"]
17 )
```

# Categorical Variable Encoding

```
1 # One-hot encoding for categorical features
2 ds.encode_categorical(
3     features=("Gender", "Race"),
4     method="onehot"
5 )
6
7 # Ordinal encoding (alternative)
8 ds.encode_categorical(
9     features=("Education_Level",),
10    method="ordinal"
11 )
```

## Encoding Methods

- **One-hot:** Creates binary columns for each category
- **Ordinal:** Maps categories to ordered integers

# Numerical Variable Scaling

```
1 # Log transformation for skewed features
2 ds.scale_numerical(
3     features=("Mortgage", "Balance"), method="log1p"
4 )
5
6 # Min-max scaling for bounded features
7 ds.scale_numerical(
8     features=("Delinquency",), method="minmax"
9 )
10
11 # Quantile transformation for robust scaling
12 ds.scale_numerical(
13     features=("Inquiry",), method="quantile"
14 )
15
16 # Standardization (alternative)
17 ds.scale_numerical(
18     features=("Income",), method="standardize"
19 )
```

# Numerical Variable Binning

```
1 # Uniform binning - equal-width intervals
2 ds.bin_numerical(
3     features=("Utilization",), bins=10, method="uniform"
4 )
5 # Quantile binning - equal-frequency intervals
6 ds.bin_numerical(
7     features=("Mortgage", "Balance", "Amount_Past_Due"),
8     bins=10, method="quantile"
9 )
10 # Execute all preprocessing steps
11 ds.preprocess()
```

## Binning Methods

- **Uniform:** Equal-width intervals
- **Quantile:** Equal-frequency intervals
- **Precompute:** User-defined bin edges

# Data Preparation for Modeling

## Essential Configuration Steps

- 1 **Train-Test Splitting:** Define data splits for validation
- 2 **Target Variable:** Specify the prediction target
- 3 **Task Type:** Set regression or classification mode
- 4 **Feature Selection:** Choose active/inactive features
- 5 **Sample Weighting:** Handle imbalanced datasets

## Modeling Readiness

Data preparation ensures the dataset is properly configured for machine learning algorithms with clear targets, appropriate splits, and selected features.

# Configuring Dataset for Modeling

```
1 # Split data into training and testing sets
2 ds.set_random_split()
3 # Set target variable for prediction
4 ds.set_target("Status")
5 # Set task type (Classification or Regression)
6 ds.set_task_type("Classification")
7 # Exclude features from modeling
8 ds.set_inactive_features(features=('Gender', 'Race'))
9 # Set sample weights (optional)
10 ds.set_sample_weight("sample_weight_column")
11 # Override active features (optional)
12 ds.set_active_features(features=('Income', 'Age', 'Balance'))
```

## Key Configuration Methods

Feature selection can be done through inclusion (`set_active_features`) or exclusion (`set_inactive_features`) approaches.



# Dataset Registration and Management

## MLflow Integration Benefits

- **Version Control:** Track dataset changes over time
- **Reproducibility:** Ensure consistent data across experiments
- **Collaboration:** Share datasets across team members
- **Metadata Tracking:** Store dataset properties and lineage

## Registration Operations:

- Register datasets
- List registered datasets
- Delete datasets

## Use Cases:

- Experiment tracking
- Dataset versioning
- Team collaboration

# Dataset Registration Operations

```
1 # Register the processed dataset
2 ds.register(name="A0-SimuCredit", override=True)
3
4 # List all registered datasets
5 registered_datasets = ds.list_registered_data()
6 print(registered_datasets)
7
8 # Delete a registered dataset (if needed)
9 ds.delete_registered_data(name="old_dataset_name")
```

## Registration Best Practices

- Use descriptive names with version indicators
- Include preprocessing information in dataset names
- Use `override=True` carefully to avoid data loss
- Regularly clean up unused registered datasets

# End-to-End MoDeVa Workflow

```
1 from modeva import DataSet
2 # 1. Data Loading
3 ds = DataSet()
4 ds.load("SimuCredit")
5
6 # 2. Data Summary
7 summary_results = ds.summary()
8 print(summary_results.table["summary"])
9
10 # 3. Data Preprocessing
11 ds.reset_preprocess()
12 ds.impute_missing(features=ds.feature_names_numerical, method='mean')
13 ds.encode_categorical(features=ds.feature_names_categorical, method="onehot")
14 ds.scale_numerical(features=("Income", "Balance"), method="standardize")
15 ds.preprocess()
16
17 # 4. Data Preparation
18 ds.set_random_split()
19 ds.set_target("Status")
20 ds.set_task_type("Classification")
21 ds.set_inactive_features(features=('ID',))
22
23 # 5. Data Registration
24 ds.register(name="processed_simucredit_v1", override=True)
```

# Best Practices and Tips

## Data Quality

- Always run `summary()` before preprocessing
- Check for data leakage in feature selection
- Validate preprocessing results

## Preprocessing Strategy

- Handle missing values before encoding
- Scale features after encoding categorical variables
- Consider domain knowledge in binning decisions

## Experiment Management

- Use consistent naming conventions for registered datasets
- Document preprocessing steps in dataset names
- Maintain preprocessing pipeline documentation

# Summary

## MoDeVa Capabilities

- **Unified Interface:** Single DataSet class for all operations
- **Comprehensive Processing:** End-to-end data workflow support
- **Built-in Datasets:** Ready-to-use demo datasets
- **Flexible Integration:** Support for various data sources
- **Experiment Tracking:** MLflow integration for reproducibility

## Next Steps

- Practice with built-in datasets
- Experiment with different preprocessing combinations
- Integrate with your ML modeling workflows
- Explore advanced MoDeVa features