# PYTHON

## DATA

### Preparation & Visualization

# Foundations of Data Preparation for Machine Learning

**Lecturer:** **Dr. Nguyen Tuan Long**

Email: ntlong@neu.edu.vn

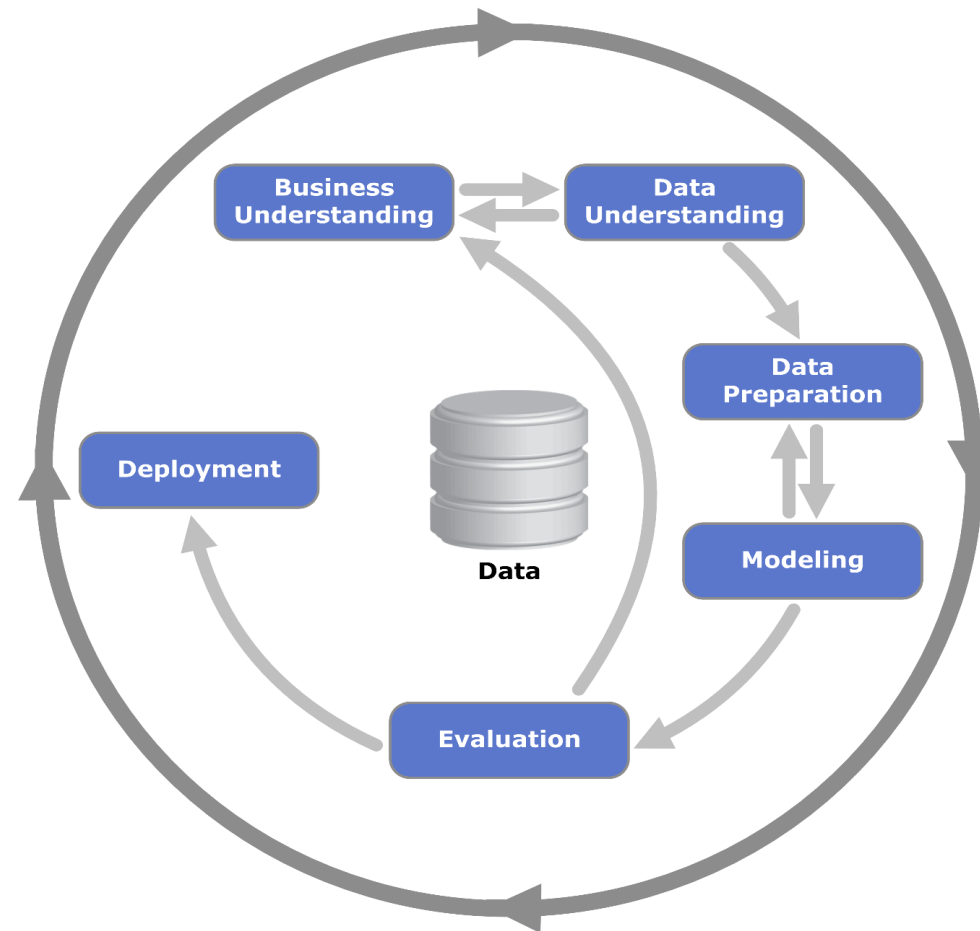Mobile: 0982 746 235

**Main Content**

1.  **The Role** of Data Preparation in a Machine Learning Project.

2.  The Data Preparation **Process** Overview.

3.  **Core Concept:** Data Leakage and How to Avoid It.

# A Machine Learning project is not a straight line, but an iterative cycle.

**Step 1: Define Problem**

- Understand the problem, collect data, choose a metric.

**Step 2: Prepare Data**

- Clean, transform, and select features.

- The **foundation** for the next steps.

**Step 3: Evaluate Algorithms**

- Test multiple models, use cross-validation.

**Step 4: Finalize Model**
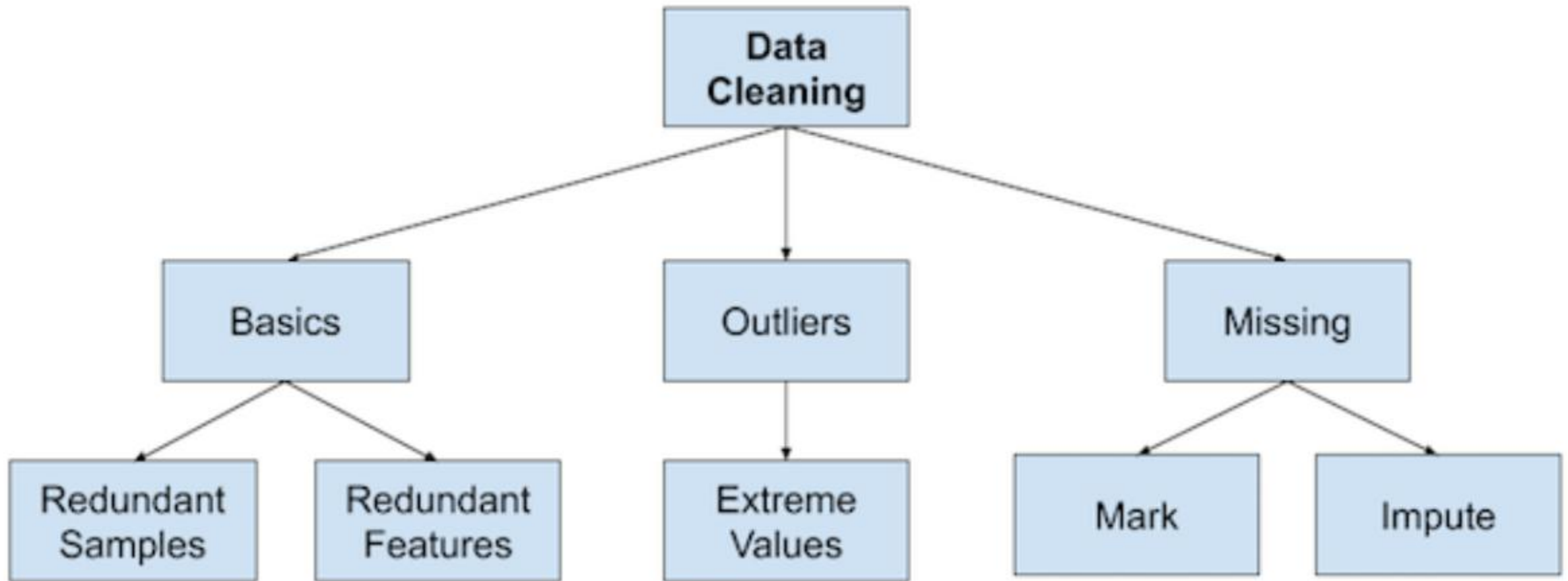
- Train the final model and deploy it.

**5 Main Task Groups**

1. Data Cleaning

2. Feature Selection

3. Data Transforms

4. Feature Engineering
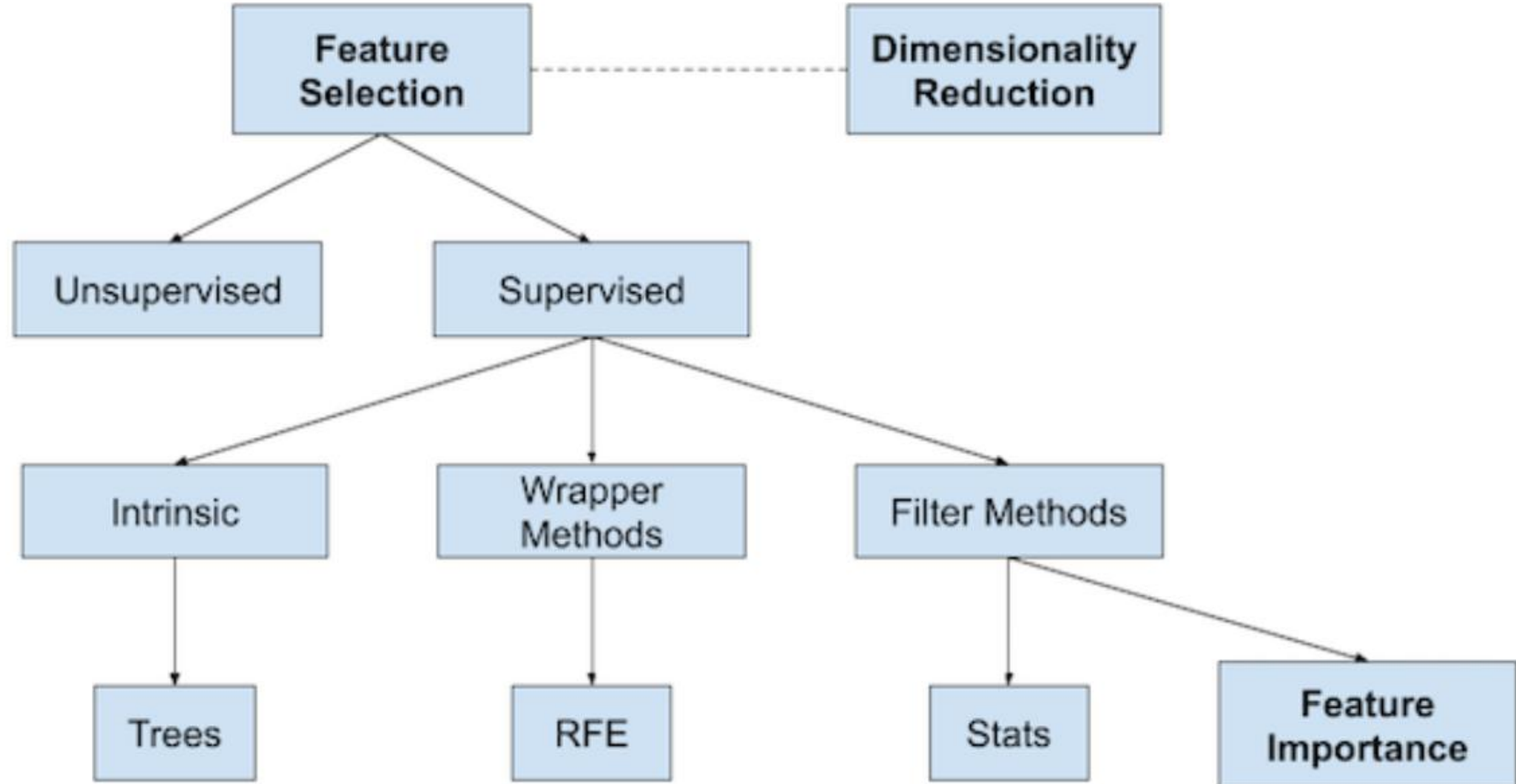
5. Dimensionality Reduction

## Overview of Data Cleaning

## Overview of Feature Selection Techniques
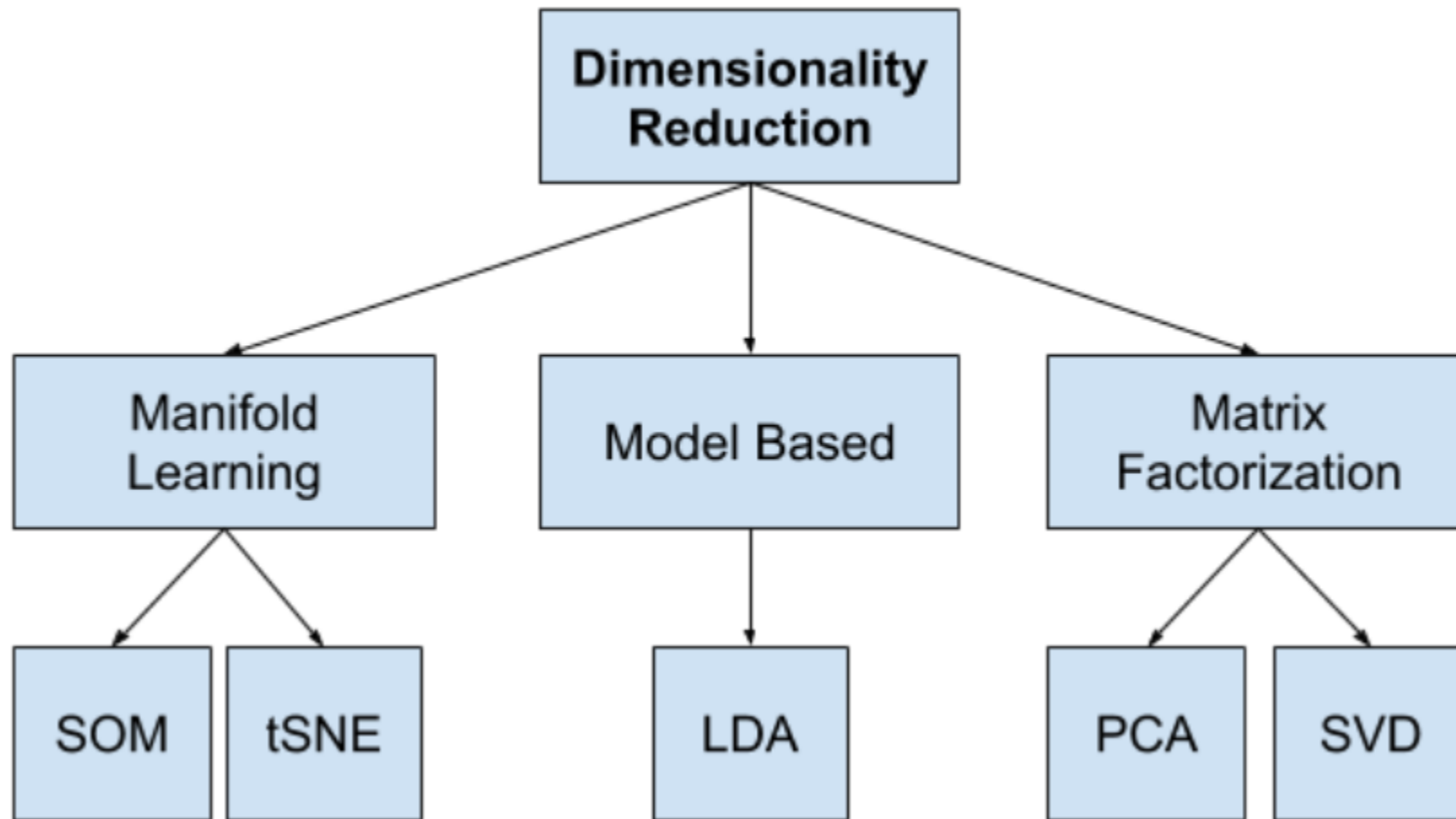
**Overview of Data Transforms**

- **Purpose:** To create new features from existing data.

- **Examples:**

  - Split a date-month-year column into day, month, year.

  - Combine population and area to create population_density.

Overview of Dimensionality Reduction Techniques

**The MOST IMPORTANT Concept**

- Misunderstanding this can make your model **completely useless** in practice.

- **Data Leakage:** Occurs when information from the **test set** is accidentally "leaked" into the model training process.

**The Golden Rule: SPLIT FIRST, PREPARE LATER**

1. **Split:** Split raw data into train and test sets.

2. **Fit:** Learn the transform parameters **ONLY on the TRAIN set**.

3. **Transform:** Apply the transform to both the train and test sets.

4. Train the model.

1. Take the **ENTIRE** dataset.

2. **Prepare the data** (e.g., Scaling).

3. **Split** into train/test sets.

4. Train the model.

=> **PROBLEM:** Information (e.g., min, max) from the test set has "leaked" into the preparation step.

**Key Takeaways**

1. Data preparation is an **iterative process** and the foundation of a project.

2. Master the 5 main task groups: Cleaning, Selection, Transforms, Engineering, & Reduction.

3. **Always split data before preparation** to avoid data leakage.

4. Using a **Pipeline** is a best practice.