

Lesson: The Machine Learning Workflow

What is the Machine Learning Workflow?

"The Machine Learning workflow is like a step-by-step process to teach computers how to learn and solve problems. Just like we follow steps to solve math problems or bake a cake, machine learning has its own set of steps that we need to follow.

These steps help us go from a problem we want to solve to a working machine learning model that can make predictions or decisions."

Steps in the Machine Learning Workflow

1. Problem Definition

"Before starting, we need to understand the problem we're trying to solve.

For example:

- If we want to predict house prices, our problem is: 'How can we predict the price of a house based on its size, location, and other details?'
- If we want to classify emails as spam or not spam, the problem is: 'How can we teach the computer to identify spam emails?'

Defining the problem clearly is the first and most important step."

2. Data Collection

"Next, we need data. Data is what the computer uses to learn. It's like practice problems for the computer.

For example:

- If we're predicting house prices, our data might include information about houses: their sizes, locations, and prices.
- If we're identifying spam emails, our data might include lots of emails labeled as 'spam' or 'not spam.'

We collect as much data as possible to help the computer learn better."

3. Data Preprocessing

"Raw data isn't always perfect—it can have missing values, incorrect entries, or unnecessary details. In this step, we clean the data and prepare it for the computer.

For example:

- If some house prices are missing, we can fill them with the average price.
- If the data includes words (like 'big' or 'small'), we turn them into numbers because computers understand numbers better."

This step ensures that the data is accurate and ready to be used for training."

4. Model Training

"In this step, we teach the computer using the data.

Here's how it works:

- We give the computer the training data (examples with answers).
- The computer studies the patterns in the data.
- It learns how to predict the answers for new data."

For example:

- If we're teaching the computer to predict house prices, we show it many houses along with their prices.
-

5. Model Evaluation

"After training, we test the computer to see how well it has learned.

We use testing data (examples that the computer hasn't seen before) to check if the computer can predict the answers correctly. This step helps us know if the computer is ready to solve real-world problems or if it needs more training."

6. Deployment

"Finally, once we're happy with how well the computer has learned, we use the trained model to make predictions or decisions.

For example:

- We can use the model to predict house prices for new houses.
- We can use the model to automatically filter spam emails from your inbox.

This step brings the machine learning model into action!"

Overview of Data Pipelines

"A data pipeline is like a conveyor belt that moves data through all the steps of the machine learning workflow.

1. The data starts as raw and messy.
2. It gets cleaned and organized during preprocessing.
3. Then it's used for training and evaluation.

Think of it as a smooth process that makes sure data flows from one step to the next, just like ingredients in a recipe!"

Activity: Splitting Data into Training and Testing Sets

"Before training a machine learning model, we need to divide our data into two parts:

1. **Training Data:** Used to teach the computer.
 2. **Testing Data:** Used to check how well the computer has learned."
-

Why Do We Split Data?

"If we use the same data for both training and testing, the computer might simply memorize the answers instead of actually learning. By using different data for testing, we can see how well the computer performs on new examples."

Example: Splitting Data Using Python

"Let's say we have a dataset of house prices, and we want to split it into training and testing sets."

```
# Importing necessary library
from sklearn.model_selection import train_test_split
```

```
# Sample dataset
data = [[2000, 3, 300000], [1500, 2, 200000], [2500, 4, 400000], [1800, 3,
250000]]
features = [row[:2] for row in data] # Size and number of rooms
labels = [row[2] for row in data] # Prices

# Splitting the data
X_train, X_test, y_train, y_test = train_test_split(features, labels,
test_size=0.25, random_state=42)

print("Training Features:", X_train)
print("Testing Features:", X_test)
print("Training Labels:", y_train)
print("Testing Labels:", y_test)
```

Introduction to the Boston Housing Dataset

"The Boston Housing dataset is a famous dataset used to predict house prices. It contains information about houses in Boston, like:

- The size of the house.
- The number of rooms.
- The location.
- The distance to schools and work areas.

For our upcoming lessons, we'll use this dataset to train and test machine learning models.

Here's an example of what the data looks like:

1. Features:
 - Size: 2000 square feet.
 - Rooms: 4.
 - Distance to school: 2 miles.
 2. Label:
 - Price: \$350,000."
-

Wrap-Up

"Today, we learned about the Machine Learning Workflow:

1. Define the problem you want to solve.
2. Collect data for the computer to learn from.
3. Preprocess and clean the data to make it ready.
4. Train the model using the training data.
5. Evaluate the model using the testing data.
6. Deploy the model to solve real-world problems.

Next, we'll explore how to use the Boston Housing dataset to build our own machine learning model!"

Splitting Data with the Boston Housing Dataset

Introduction

"The Boston Housing dataset is a collection of data about houses in Boston. It includes features like the size of the house, the number of rooms, and the location. Our goal is to teach the computer to predict the price of a house based on these features.

Before we train our model, we need to split the dataset into two parts:

1. **Training Data:** The part the computer will learn from.
 2. **Testing Data:** The part we'll use to check if the computer has learned well."
-

Code: Splitting the Boston Housing Dataset

"Let's see how to split the Boston Housing dataset using Python. I'll explain each line of the code so it's easy to understand."

```
# Import necessary libraries
from sklearn.datasets import load_boston
from sklearn.model_selection import train_test_split
import pandas as pd
```

1. **from sklearn.datasets import load_boston:**
"This imports the Boston Housing dataset from a library called scikit-learn, which provides many useful machine learning tools."
 2. **from sklearn.model_selection import train_test_split:**
"This imports a function to split our data into training and testing sets."
 3. **import pandas as pd:**
"This imports pandas, a library that helps us work with data in table format, like an Excel sheet."
-

```
# Load the Boston Housing dataset
boston = load_boston()
```

4. **boston = load_boston():**
"This loads the Boston Housing dataset into a variable called boston. Now, the dataset is ready for us to use."

```
# Convert the dataset into a DataFrame
data = pd.DataFrame(boston.data, columns=boston.feature_names)
data['PRICE'] = boston.target
```

5. **pd.DataFrame(boston.data, columns=boston.feature_names):**

"Here, we convert the dataset into a DataFrame, which is like a table with rows and columns.

- o boston.data contains the actual data (features).
- o boston.feature_names gives the names of each column, like 'RM' (number of rooms) or 'LSTAT' (lower status population)."

6. **data['PRICE'] = boston.target:**

"This adds a new column to the table called PRICE, which is the price of each house. Now the table has both features and labels."

```
# Define features (X) and labels (y)
X = data.drop('PRICE', axis=1)
y = data['PRICE']
```

7. **x = data.drop('PRICE', axis=1):**

"Here, we define x, which contains all the features (like the number of rooms or size of the house). We drop the PRICE column because it's the label, not a feature."

8. **y = data['PRICE']:**

"This defines y, which contains the labels (house prices). This is what we want to predict."

```
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

9. **train_test_split(x, y, test_size=0.2, random_state=42):**

"This splits the data into training and testing sets:

- o X_train and y_train: The training features and labels (80% of the data).
- o X_test and y_test: The testing features and labels (20% of the data).
- o test_size=0.2 means 20% of the data is used for testing, and 80% is used for training.
- o random_state=42 ensures the split is the same every time we run the code."

```
# Print the sizes of each set
print("Training features shape:", X_train.shape)
print("Testing features shape:", X_test.shape)
print("Training labels shape:", y_train.shape)
print("Testing labels shape:", y_test.shape)
```

10. **x_train.shape and x_test.shape:**

"This shows the number of rows and columns in the training and testing features. It helps us check if the split worked correctly."

Summary of What the Code Does

1. Load the Boston Housing dataset.
 2. Organize the data into features (x) and labels (y).
 3. Split the data into training (80%) and testing (20%) sets.
 4. Print the sizes of the training and testing sets to confirm the split.
-

Wrap-Up

"Now the dataset is ready for us to use in upcoming lessons! We'll use the training data to teach our model and the testing data to check how well the model has learned."
