

6.1 Pandas

6.2 Matplotlib

6.3 Numpy

6.4 Scikit-learn

■ Question Bank

When it comes to machine learning and deep learning projects written in Python, there are thousands of libraries to pick and choose from. However, they're not all on the same level of code quality, diversity, or size. To help you choose, here are the best Python libraries for machine learning. Here we are going to discuss about Pandas, Numpy, Matplotlib and Scikit-learn.

6.1 Pandas

Features of Pandas:

1. Data structures: Pandas provides two main data structures - Series (1D labeled array) and DataFrame (2D labeled data structure) - that allow for easy handling and manipulation of data. These data structures provide powerful and flexible ways to represent and analyze data.
2. Data cleaning: Pandas offers various functionalities for data cleaning, such as handling missing values, removing duplicates, and data type conversion. These operations are crucial for preparing data for machine learning models, as clean and well-structured data is essential for accurate model training.

3. Data transformation: Pandas provides functions for transforming data, such as filtering, sorting, and grouping. These operations are commonly used for data preprocessing and feature engineering in machine learning workflows.
4. Data visualization: Pandas integrates well with other popular data visualization libraries in Python, such as Matplotlib and Seaborn, allowing for easy visualization of data to gain insights and perform exploratory data analysis (EDA).
5. Data I/O: Pandas offers functions to read and write data in various formats, such as CSV, Excel, and SQL, which are commonly used for loading and storing data in machine learning workflows.
6. Flexible indexing and slicing: Pandas allows for flexible indexing and slicing of data, making it easy to extract, filter, and manipulate subsets of data based on various conditions and criteria.
7. Data alignment: Pandas automatically aligns data based on labeled indexes, making it easy to perform operations on data with different shapes and

- sizes, and handle missing data gracefully.
8. Handling time-series data: Pandas provides built-in support for handling time-series data, making it easy to work with data that has a time component, such as financial or sensor data.
 9. Memory-efficient operations: Pandas provides efficient operations for data manipulation, such as vectorized operations and optimized algorithms, making it suitable for handling large datasets.
 10. Interoperability: Pandas seamlessly integrates with other popular Python libraries such as NumPy, Scikit-learn, and more, making it a versatile tool for data analysis and machine learning workflows.

Application of Pandas:

1. Data preprocessing: Pandas provides functions for handling missing values, converting data types, removing duplicates, and handling outliers, which are crucial steps in preparing data for machine learning models.

For example:

```
import pandas as pd
# Load data into a DataFrame
data = pd.read_csv('data.csv')
# Handling missing values
data.dropna() # Drop rows with missing values
data.fillna(value) # Fill missing values with a specified value
# Data type conversion
data['column_name'] =
data['column_name'].astype('data_type')
# Removing duplicates
```

```
data.drop_duplicates() # Drop duplicate rows
```

2. Feature engineering: Pandas provides powerful functions for data transformation, such as filtering, sorting, grouping, and merging, which are commonly used for feature engineering.

For example:

```
# Filtering data
filtered_data = data[data['column_name'] > threshold]
# Grouping data
grouped_data =
    data.groupby('column_name').mean()
```

```
# Merging data
merged_data = pd.merge(data1, data2
    on='common_column', how='inner')
```

3. Data analysis: Pandas makes it easy to perform exploratory data analysis (EDA) and gain insights from data before building machine learning models. For example:

```
# Descriptive statistics
data.describe() # Get summary statistics of data
# Data visualization
import matplotlib.pyplot as plt
data.plot(x='column_name',
y='column_name', kind='scatter')
# Create scatter plot
plt.show()
```

```
# Pivot tables
pivot_table = data.pivot_table(
values='value_column',
index='index_column', columns=
column_name)
```

4. Data preparation for model training: Pandas allows for easy data slicing,

indexing, and splitting, which are essential for preparing data for model training. For example:

```
# Data slicing and indexing
X = data[['feature1', 'feature2']] # Select specific features as input (X)
y = data['target'] # Select target variable (y)

# Data splitting
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state)
```

6.2 Matplotlib

Matplotlib is a popular data visualization library in Python that provides a wide range of functionalities for creating static, animated, and interactive plots. Some of the key features of Matplotlib include:

Features of Matplotlib:

1. Comprehensive Plotting: Matplotlib provides a vast collection of plot types, including line plots, scatter plots, bar plots, histogram, pie charts, 3D plots, and more. It offers extensive customization options for colors, markers, linestyles, labels, titles, legends, and axes.
2. Multi-platform Support: Matplotlib supports various output formats, including PNG, JPEG, PDF, SVG, and more. It also provides an interactive mode for embedded applications and supports various GUI backends, such as Tkinter, Qt, GTK, and more, making it suitable for different platforms and environments.
3. Data Visualization: Matplotlib allows for visualizing complex data through

different plotting techniques, such as subplots, grids, multiple axes, shared axes, and dual axes. It also provides support for adding text, annotations, arrows, images, and other visual elements to enhance the plot's visual appeal and readability.

4. Customization: Matplotlib provides extensive customization options, allowing users to customize the appearance of plots based on their specific needs. Users can customize plot properties, including figure size, axes limits, ticks, labels, fonts, colors, and more, to create visually appealing and informative plots.
5. Plotting Styles: Matplotlib provides support for various predefined plotting styles, including classic, ggplot, seaborn, and more, allowing users to easily switch between different styles and customize the visual appearance of their plots.
6. Object-oriented and Pyplot Interfaces: Matplotlib offers two interfaces for creating plots: the object-oriented interface and the pyplot interface. The object-oriented interface provides more flexibility and control over plot creation and customization, while the pyplot interface provides a simple and intuitive way to create plots using a MATLAB-like syntax.
7. Animation and Interactivity: Matplotlib provides support for creating animated plots, allowing users to visualize data changes over time. It also provides support for adding interactivity to plots, allowing users to zoom, pan, and interact with plots to explore the data in a more dynamic way.

8. Integration with Other Libraries: Matplotlib integrates well with other popular Python libraries such as NumPy, Pandas, and Scikit-learn, making it a powerful tool for visualizing data in machine learning workflows and data analysis tasks.

Applications of Matplotlib:

1. Exploratory Data Analysis (EDA): Matplotlib can be used to visualize and explore the data before building machine learning models. This can include creating scatter plots, histograms, box plots, bar plots, and other types of plots to gain insights into the distribution, trends, and relationships of the data. EDA with Matplotlib helps in identifying patterns, outliers, and potential issues in the data, which can inform feature engineering and model selection decisions.
2. Model Evaluation: Matplotlib can be used to visualize the performance of machine learning models by plotting various evaluation metrics, such as accuracy, precision, recall, F1-score, and ROC curves. These visualizations can help in comparing the performance of different models, selecting the best-performing model, and identifying potential issues, such as overfitting or underfitting.
3. Feature Visualization: Matplotlib can be used to visualize the importance or contribution of features in a machine learning model. This can include creating bar plots or heatmaps to display feature importances, coefficients, or weights from models such as linear regression or logistic regression. These visualizations can help in understanding the importance of different features in the model's prediction and interpreting the model's behavior.
4. Model Interpretability: Matplotlib can be used to create visualizations that help interpret complex machine learning models, such as decision trees, random forests, or support vector machines. These visualizations can help in understanding the decision-making process of the model, visualizing decision boundaries, and identifying critical features or instances that influence the model's prediction.
5. Model Explainability: Matplotlib can be used to create visualizations that explain the prediction of machine learning models, such as partial dependence plots, individual conditional expectation plots, or LIME plots. These visualizations can help in explaining how the model's prediction changes with respect to specific features or instances and provide insights into the model's behavior and predictions for interpretability and transparency.
6. Model Deployment: Matplotlib can be used to create visualizations for model deployment in production systems, such as creating interactive plots, dashboards, or reports to monitor model performance, visualize real-time predictions, or provide insights to stakeholders.

6.3 Numpy

Numpy is a library for the Python programming language, adding support for large, multi-dimensional arrays and matrices, along with a large collection of

high-level mathematical functions to operate on these arrays. Moreover Numpy forms the foundation of the Machine Learning stack.

NumPy (Numerical Python) is a powerful numerical computing library in Python that provides essential functionalities for performing numerical operations efficiently.

Features of NumPy:

Some of the key features of NumPy include:

1. N-dimensional array: NumPy provides a powerful N-dimensional array object called **ndarray**, which is a homogeneous collection of elements with fixed-size and enables efficient operations on large arrays of data. It allows for easy manipulation and computation of multi-dimensional data, such as matrices, images, and time series.
2. Broadcasting: NumPy allows for element-wise operations on arrays of different shapes and sizes through broadcasting, which automatically handles shape and size compatibility. This enables efficient and concise operations on arrays with different dimensions, eliminating the need for explicit loops.
3. Mathematical Functions: NumPy provides a vast collection of mathematical functions for performing operations such as arithmetic, trigonometry, logarithm, exponential, linear algebra, statistical analysis, and more. These functions are optimized for performance and provide efficient computations for numerical operations.
4. Array Manipulation: NumPy provides

a wide range of functions for manipulating arrays, including reshaping, resizing, stacking, splitting, transposing, and more. These functions allow for efficient manipulation of array data, making it easy to preprocess and transform data for machine learning and other numerical computations.

5. Broadcasting: NumPy allows for element-wise operations on arrays of different shapes and sizes through broadcasting, which automatically handles shape and size compatibility. This enables efficient and concise operations on arrays with different dimensions, eliminating the need for explicit loops.
6. Vectorized Operations: NumPy allows for vectorized operations, which enable efficient computation on entire arrays without the need for explicit loops. This makes NumPy code more concise, readable, and optimized for performance.
7. Interoperability: NumPy provides support for interoperability with other libraries and data formats. It can seamlessly integrate with other Python libraries, such as pandas, Matplotlib, and scikit-learn, making it a fundamental building block for scientific computing and data analysis in Python.
8. Performance Optimization: NumPy is designed for efficient numerical computing and provides optimized implementations of mathematical functions that are written in C or Fortran. It also allows for vectorized operations and provides support for parallel computing, making it highly

efficient for processing large datasets and performing computationally-intensive tasks.

9. Random Number Generation: NumPy provides functions for generating random numbers, including various probability distributions, random sampling, and random data generation. This is useful for simulations, statistical analysis, and machine learning applications that require randomization.

Applications of Numpy:

1. Scientific Computing: NumPy provides a powerful array object, `ndarray`, which allows for efficient handling of large numerical data, such as matrices, vectors, and tensors. It provides essential functions for performing mathematical operations on arrays, such as linear algebra, Fourier analysis, numerical integration, interpolation, and more, making it widely used in scientific computing for tasks such as simulations, data analysis, and modeling.
2. Data Analysis: NumPy is often used in combination with other data analysis libraries, such as pandas, for data manipulation, preprocessing, and analysis tasks. NumPy provides efficient array operations for handling large datasets, numerical computations, and mathematical operations, making it a foundational library for data analysis tasks, including data cleaning, feature engineering, statistical analysis, and more.
3. Machine Learning: NumPy is a fundamental library used in machine learning tasks for data preprocessing, feature engineering, and model

evaluation. It provides efficient operations for manipulating and transforming data in arrays, such as reshaping, slicing, and stacking, which are commonly used in machine learning workflows. NumPy is also used in implementing numerical algorithms for machine learning models, such as linear regression, logistic regression, support vector machines, and deep learning models.

4. Image and Signal Processing: NumPy provides functions for image and signal processing tasks, such as image manipulation, filtering, Fourier analysis, and convolution operations. It is widely used in computer vision, image processing, and signal processing applications for tasks such as image filtering, feature extraction, and image reconstruction.
5. Computational Mathematics: NumPy is used in various computational mathematics applications, including numerical optimization, numerical integration, numerical solving of equations, and more. It provides efficient array operations, mathematical functions, and linear algebra operations, which are essential in solving numerical problems in various fields, such as physics, engineering, and finance.
6. Data Visualization: NumPy is often used in conjunction with data visualization libraries, such as Matplotlib and Seaborn, for creating visualizations of numerical data. NumPy provides efficient operations for data manipulation, aggregation, and computation, which can be used to prepare data for visualization and

perform calculations for generating plots, charts, and graphs.

6.4 Scikit-learn

Scikit-learn is an open-source Python library that implements a range of machine learning, pre-processing, cross-validation, and visualization algorithms using a unified interface.

Scikit-learn, a popular machine learning library in Python, offers a wide range of features that make it a powerful tool for various machine learning tasks.

Features of Scikit-learn :

Some of the key features of scikit-learn are:

1. Broad Range of Machine Learning Algorithms: Scikit-learn provides a diverse collection of machine learning algorithms for classification, regression, clustering, dimensionality reduction, ensemble methods, and more. It includes popular algorithms such as linear regression, logistic regression, decision trees, support vector machines, k-nearest neighbors, random forests, gradient boosting, and many more.
2. Consistent API: Scikit-learn follows a consistent API for all its algorithms, making it easy to learn and use. The “fit-transform-predict” paradigm provides a consistent approach for training models, transforming data, and making predictions. This consistent API allows for easy switching between different algorithms and models.
3. Data Preprocessing: Scikit-learn provides tools for data preprocessing, including handling missing values, feature scaling, feature encoding, and feature selection. These tools help in preparing data for machine learning algorithms by handling common data-related issues.
4. Model Selection and Evaluation: Scikit-learn offers tools for model selection and evaluation, including techniques such as cross-validation, grid search, and randomized search for hyperparameter tuning. It also includes tools for evaluating model performance using various metrics, such as accuracy, precision, recall, F1-score, and more.
5. Feature Extraction and Dimensionality Reduction: Scikit-learn provides methods for feature extraction, such as Principal Component Analysis (PCA) and feature selection techniques like SelectKBest, SelectPercentile, and Recursive Feature Elimination (RFE). These techniques help in reducing the dimensionality of the data and selecting relevant features for machine learning models.
6. Model Persistence: Scikit-learn allows for saving trained models to disk and reloading them later for making predictions on new data. This makes it easy to reuse models in production environments without the need to retrain them every time.
7. Integrated Data Visualization: Scikit-learn integrates well with other popular Python libraries for data visualization, such as Matplotlib and Seaborn, allowing for visual exploration and analysis of data.
8. Robust Documentation and Community Support: Scikit-learn has extensive documentation and a large community of users and developers, providing resources and support for beginners.

and advanced users alike. The library is actively maintained and regularly updated with new features and bug fixes.

9. Scalability: Scikit-learn is designed to be scalable and efficient, making it suitable for handling large datasets and running on various computing platforms, including desktops, servers, and distributed clusters.

Applications of Scikit-learn:

1. Classification: Scikit-learn provides a wide range of algorithms for classification tasks, such as logistic regression, support vector machines, decision trees, random forests, and naive Bayes, among others. These algorithms can be used for tasks like spam detection, fraud detection, sentiment analysis, and image classification.
2. Regression : Scikit-learn offers algorithms for regression tasks, such as linear regression, ridge regression, and support vector regression. These algorithms can be used for tasks like predicting housing prices, stock prices, and other continuous variables.
3. Clustering : Scikit-learn provides algorithms for unsupervised learning tasks, such as clustering, including k-means, hierarchical clustering, and DBSCAN. These algorithms can be used for tasks like customer segmentation, anomaly detection, and image segmentation.
4. Dimensionality Reduction: Scikit-learn includes methods for dimensionality

reduction, such as Principal Component Analysis (PCA), which can be used for tasks like reducing the dimensionality of high-dimensional data for visualization, feature extraction, and model building.

5. Model Selection and Evaluation: Scikit-learn offers tools for model selection and evaluation, such as cross-validation, grid search, and randomized search for hyperparameter tuning. These tools help in selecting the best model and optimizing its performance.
6. Ensemble Methods: Scikit-learn includes ensemble methods, such as Random Forests, Gradient Boosting, and Bagging, which can be used for tasks like improving the accuracy and robustness of models.
7. Preprocessing: Scikit-learn provides tools for data preprocessing, such as handling missing values, feature scaling, and feature encoding. These tools help in preparing data for machine learning algorithms.
8. Text Mining: Scikit-learn includes text mining tools, such as text feature extraction techniques, vectorizers, and natural language processing (NLP) tools, which can be used for tasks like text classification, sentiment analysis, and topic modeling.
9. Model Deployment: Scikit-learn models can be easily deployed in production environments for making real-time predictions on new data, making it suitable for building machine learning applications.

Difference between numpy, pandas, matplotlib and scikit-learn

Library	Purpose	Main Features
Numpy	Numerical computing	<ul style="list-style-type: none"> - Provides support for handling large, multi-dimensional arrays and matrices. - Offers mathematical functions for operations on arrays, such as linear algebra and statistics.
Pandas	Data manipulation and analysis	<ul style="list-style-type: none"> - Provides data structures like Series and DataFrame for efficient data handling. - Offers powerful tools for data cleaning, aggregation, filtering, and transformation.
Matplotlib	Data visualization	<ul style="list-style-type: none"> - Provides a wide range of plotting and visualization options for creating static visualizations. - Supports various types of plots, including line plots, scatter plots, histograms, and more.
Scikit-learn	Machine learning and data modeling	<ul style="list-style-type: none"> - Offers a wide range of machine learning algorithms for classification, regression, clustering, and dimensionality reduction tasks. - Provides tools for model selection, evaluation, and preprocessing of data. - Includes utilities for model evaluation, metrics, and model deployment.

Example of classification using numpy, pandas, Matplotlib, Scikit-learn.

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score, confusion_matrix
# Load data into a Pandas DataFrame
data = pd.read_csv('data.csv')

# Data preprocessing
X = data[['feature1', 'feature2']] # Select input features (X)
y = data['target'] # Select target variable (y)

```

```
# Split 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)

# Model training
clf = LogisticRegression() # Instantiate a logistic regression model
clf.fit(X_train, y_train) # Fit the model to the training data

# Model prediction
y_pred = clf.predict(X_test) # Predict the target variable for the test data

# Model evaluation
accuracy = accuracy_score(y_test, y_pred) # Calculate accuracy
confusion = confusion_matrix(y_test, y_pred) # Calculate confusion matrix
# Data visualization
plt.scatter(X_test['feature1'], X_test['feature2'], c=y_pred) # Scatter plot of predicted data
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
plt.title('Classification Results')
plt.show( )
```

Question Bank

■ Short-Answer Questions

(3 or 4 marks) :

- Give the purpose of Numpy in Machine Learning.
- Compare and contrast Pandas and Numpy.
- What is the main purpose of matplotlib in Machine Learning ?
- What are the advantages of using Scikit-learn?
- How to load dataset using Numpy ? Explain.
- How to load dataset using Panda ? Explain.
- How to plot a vertical line and a horizontal line using Matplotlib ?

■ Long-Answer Questions

(7 marks):

- Explain features and applications of Pandas.
- Explain features and applications of Numpy.
- Explain features and applications of Matplotlib.
- Explain features and applications of Scikit-Learn.
- Write a Python program to load the iris data from a given csv file into a dataframe and print the shape of the data, type of the data and first 3 rows using Scikit-Learn.

■ Check your knowledge :

- What is a DataFrame in pandas ?
 - A type of function
 - A two-dimensional table-like data structure with rows and columns
 - A way to iterate through a pandas object
 - None of above
- How can you create a NumPy array ?
 - Using the array() function
 - Using the list() function
 - All of above
- How can you add a label to the x-axis of a plot in Matplotlib ?
 - Using the xlabel() function
 - Using the ylabel() function
 - Using the xplot() function
- How can you split a dataset into training and testing sets in Scikit-learn ?
 - Using the split() function
 - Using the train_test_split() function
 - Using the shuffle() function
 - Using the train_test_shuffle() function
- What is pandas?
 - A type of bear
 - A Python library used for data manipulation and analysis
 - A programming language
 - All of above