

# NumPy Handbook for Machine Learning

Explore NumPy with beginner – friendly projects and practical applications



First Edition

Samriddha Pathak

# NumPy Handbook for Machine Learning

By Samriddha Pathak

---

## Preface

Welcome to the "NumPy Handbook for Machine Learning." This book is intended to be a comprehensive handbook to learning and utilizing NumPy, the basic package for scientific computing in Python, and how it is applied in machine learning.

As we enter an age of artificial intelligence and information dominance, the skill of efficiently manipulating and analyzing numerical data becomes ever more important. The centerpiece of such a skill is NumPy, which provides efficient tools that make tedious numerical computations simpler and improve computational efficiency greatly.

This book is intended for students, researchers, and practitioners who aspire to utilize the complete capabilities of NumPy on their machine learning endeavors. This book is the necessary background and implementation examples to start with, if you are beginning your data science journey or seeking to know more about numerical computation in the machine learning algorithm.

In this book, not only will we be discussing the basics and operations of NumPy but also real-world usage by implementing detailed projects. By the time you finish this book, you will have a strong understanding of NumPy that will enable you to solve real-world machine learning problems with confidence.

I sincerely hope this handbook will be a useful companion on your machine learning journey.

# Acknowledgement

---

The release of this book, NumPy Handbook for Machine Learning, has not been made possible without the constant encouragement, guidance, and assistance of some individuals. I would like to take this opportunity to thank the individuals who played a significant role in this endeavor.

First and foremost, I would like to thank Mr. Rasum Subedi for his excellent editorial service. His keen eye for detail and commitment to editing my manuscript into a professionally written and well-structured book have been an invaluable help in assisting me in turning my manuscript into a readable and well-structured book. His commitment and efforts have significantly enhanced the readability and comprehensibility of this book.

I am highly grateful to my parents, Kumar Pathak and Radhika Pathak, whose belief, encouragement, and confidence in my potential have been the driving force of my determination. Their encouragement and constant assurance kept me firm and resolute during this difficult process.

I would also like to express my gratitude to my close friends Shrijan Shrestha, Kriti Dahal, Amnisha Luitel, Rasum Subedi, Amit Rai, and Binod Basnet for their constant support, stimulating discussions, and inspiring words. Their support of my idea and their constructive criticism have been vital in the creation of this book.

A genuine word of thanks to my Sister Jyoti Pokhrel and Saroj Dahal Brother, whose inspiration, support, and encouragement have been the driving forces behind me in all this. Their faith in me and support have encouraged me to keep going even when everything seemed to be against me. In addition, I also want to express my gratitude to the whole data science and machine learning community whose collective research, expertise, and shared resources have educated and inspired much of what is presented in this book. Their open-source learning has been invaluable in my own growth and in the production of this book.

Lastly, I would also appreciate it if all the readers and students interested in expanding their knowledge on NumPy and its applications in machine learning would thank me. My sincere hope is that this book is a helpful guide to your studies and assists you in growing in this ever-evolving field.

With sincerest thanks,

Samriddha Pathak

# Table of Contents

---

1. Introduction to NumPy
  - What is NumPy?
  - Why NumPy for Machine Learning?
  - Installation and Setup
2. NumPy Fundamentals
  - Creating NumPy Arrays
  - Array Attributes
  - Array Indexing and Slicing
  - Array Reshaping
  - Array Iteration
3. Mathematical Operations in NumPy
  - Statistical Functions
  - Linear Algebra Operations
4. Advanced NumPy Operations
  - Broadcasting
  - Array Manipulation
  - Sorting and Searching
  - File I/O with NumPy
5. Project 1: K-Nearest Neighbors from Scratch
  - Theory and Concepts
  - Implementation
  - Testing and Evaluation
6. Project 2: Principal Component Analysis (PCA) Implementation
  - Theory and Background
  - Implementation with NumPy
  - Visualization and Application
7. NumPy Best Practices
  - Performance Optimization
  - Memory Management
  - Common Pitfalls and Solutions

8. NumPy and the Machine Learning Ecosystem

- NumPy with Pandas
- NumPy with Scikit-learn
- NumPy with TensorFlow and PyTorch

9. Appendix

- Quick Reference
  - Further Resources
  - Glossary of Terms
-

# Chapter 1: Introduction to NumPy

## What is NumPy?

NumPy (Numerical Python) is a high-performance library for scientific computing in Python. It provides support for large, multi-dimensional arrays and matrices, along with a vast collection of mathematical and numerical functions to efficiently manipulate these arrays.

NumPy serves as the foundation for many scientific computing and machine learning libraries, including Pandas, SciPy, TensorFlow, and PyTorch. It offers significant advantages over Python's built-in lists due to its faster computation, lower memory consumption, and efficient vectorized operations.

## Why NumPy for Machine Learning?

Machine learning and artificial intelligence involve handling and processing large datasets efficiently. NumPy plays a crucial role in optimizing this process by providing:

- High-performance computations: NumPy is implemented in C and Fortran, making it significantly faster than using Python lists for numerical operations.
- Memory efficiency: NumPy arrays are stored in contiguous memory locations, reducing the overhead of Python objects and enabling faster access.
- Vectorized operations: Unlike traditional Python loops, NumPy supports vectorization, which allows operations to be performed on entire arrays at once. This eliminates explicit loops and speeds up execution.
- Broadcasting: Enables element-wise operations on arrays of different shapes without the need for explicit reshaping.
- Interoperability: NumPy integrates seamlessly with other libraries, such as:
  - Pandas (data manipulation)
  - Matplotlib (data visualization)
  - Scikit-learn (machine learning)
  - TensorFlow & PyTorch (deep learning)

## Key Features of NumPy:

Here are some key features that make NumPy an indispensable tool for machine learning:

1. N-dimensional array (`ndarray`): A powerful object that allows for handling of large datasets efficiently.
2. Mathematical and statistical functions: Includes operations like mean, median, standard deviation, and linear algebra functions for numerical analysis.

3. Random number generation: Used in machine learning for initializing weights, bootstrapping, and creating synthetic data.
4. Broadcasting: Allows operations on arrays of different sizes without the need for manual reshaping.
5. File handling: NumPy allows reading and writing from CSV, binary, and other formats, making it easy to store and manipulate datasets.
6. Memory efficiency: Uses less memory than Python lists due to optimized storage.

## Installation and Setup

Before using NumPy, you need to install it in your Python environment.

Installing NumPy using pip (Recommended Method)

You can install NumPy using the Python package manager pip by running:

```
pip install numpy
```

# Chapter 2: NumPy Fundamentals

## Creating NumPy Arrays

A NumPy array is the fundamental building block of the NumPy library. Unlike Python lists, NumPy arrays provide a more efficient, faster, and memory-optimized way to store and manipulate numerical data.

### 1. Creating a NumPy Array from a List

```
import numpy as np

# Creating a 1D array from a list
arr_1d = np.array([1, 2, 3, 4, 5])

print("1D Array:", arr_1d)
print("Type of the array:", type(arr_1d)) # Verify it's a NumPy array
```

 **Note:**

*Unlike Python lists, NumPy arrays provide faster numerical operations and take up less memory.*

### 2. Creating Multi-Dimensional Arrays

You can also create 2D and 3D arrays in NumPy:

```
# Creating a 2D NumPy array (Matrix)
arr_2d = np.array([[1, 2, 3], [4, 5, 6]])
print("2D Array:\n", arr_2d)

# Creating a 3D NumPy array (Tensor)
arr_3d = np.array([[[1, 2, 3], [4, 5, 6]], [[7, 8, 9], [10, 11, 12]]])
print("3D Array:\n", arr_3d)
```

 **Key Differences:**

- 1D Array: A simple list of numbers
- 2D Array: A matrix (rows & columns)
- 3D Array: A collection of 2D arrays

### 3. Creating Special NumPy Arrays

NumPy provides several built-in functions to create arrays quickly:

#### a) Creating an Array of Zeros

```
zeros_arr = np.zeros((3, 3)) # Creates a 3x3 matrix filled with zeros
print("Zeros Array:\n", zeros_arr)
```

### b) Creating an Array of Ones

```
ones_arr = np.ones((2, 5)) # Creates a 2x5 matrix filled with ones
print("Ones Array:\n", ones_arr)
```

### c) Creating an Identity Matrix

```
identity_matrix = np.eye(4) # Creates a 4x4 identity matrix
print("Identity Matrix:\n", identity_matrix)
```

### d) Creating an Array with a Range of Values

```
range_arr = np.arange(1, 11, 2) # Creates an array with values from 1 to 10
with a step of 2
print("Range Array:", range_arr)
```

### e) Creating an Array with Evenly Spaced Values

```
linspace_arr = np.linspace(0, 1, 5) # Creates 5 evenly spaced values between
0 and 1
print("Linspace Array:", linspace_arr)
```

### f) Creating a Random Array

```
random_arr = np.random.rand(3, 3) # Generates a 3x3 array of random values
between 0 and 1
print("Random Array:\n", random_arr)
```

#### Tip:

These functions are useful for machine learning, especially when initializing weights and datasets.

## Array Attributes

NumPy arrays have several attributes that provide important information about their properties.

```
# Creating a sample NumPy array
arr = np.array([[10, 20, 30], [40, 50, 60]])

print("Array:\n", arr)
print("Shape of array:", arr.shape) # Returns (rows, columns)
print("Size of array:", arr.size) # Returns total number of elements
print("Number of dimensions:", arr.ndim) # Returns 2 (for 2D array)
print("Data type of elements:", arr.dtype) # Data type of elements
print("Item size in bytes:", arr.itemsize) # Size of each element in bytes
print("Total memory consumed:", arr.nbytes, "bytes") # Total memory usage
```

### Key Takeaways:

- .shape tells the (rows, columns) of an array
- .size gives the total number of elements
- .dtype shows the data type of elements
- .ndim returns the number of dimensions

## Array Indexing and Slicing

Just like Python lists, NumPy arrays support indexing and slicing.

### 1. Accessing Elements in a 1D Array

```
arr_1d = np.array([10, 20, 30, 40, 50])
print("First Element:", arr_1d[0])      # Access first element
print("Last Element:", arr_1d[-1])     # Access last element
print("Elements from index 1 to 3:", arr_1d[1:4]) # Slicing
print("Every second element:", arr_1d[::-2]) # Step slicing
```

### 2. Accessing Elements in a 2D Array

```
arr_2d = np.array([[10, 20, 30], [40, 50, 60], [70, 80, 90]])

print("Element at row 1, column 2:", arr_2d[1, 2]) # Accessing element
print("First row:", arr_2d[0, :])    # Accessing entire row
print("First column:", arr_2d[:, 0])  # Accessing entire column
```

## Array Reshaping

Reshaping allows you to change the shape of an array without altering the data.

```
arr = np.arange(12) # Create an array with 12 elements
reshaped_arr = arr.reshape(3, 4) # Reshape it into a 3x4 matrix
print("Original Array:", arr)
print("Reshaped Array:\n", reshaped_arr)
```

### Tip:

Use -1 to automatically infer one dimension:

```
auto_reshaped = arr.reshape(4, -1)    # NumPy automatically calculates the
                                         # missing dimension
print(auto_reshaped.shape) # Output: (4, 3)
```

## Array Iteration

### 1. Iterating Over a 1D Array

```
arr_1d = np.array([10, 20, 30, 40])  
  
for num in arr_1d:  
    print(num)
```

### 2. Iterating Over a 2D Array

```
arr_2d = np.array([[1, 2, 3], [4, 5, 6]])  
  
# Iterating row by row  
for row in arr_2d:  
    print("Row:", row)  
  
# Iterating element by element  
for element in np.nditer(arr_2d):  
    print("Element:", element)
```

# Chapter 3: Mathematical Operations in NumPy

NumPy provides an extensive library of numerical computation to aid efficient computation in matrices, arrays, and numerical arrays. NumPy is most efficient in these computations than loops in ordinary Python because NumPy is most efficiently efficient in these computations to aid efficient computation in matrices, arrays, and numerical arrays.

This chapter covers:

- Element-wise operations
- Universal functions (ufuncs)
- Statistical computations
- Linear algebra operations
- Advanced mathematical functions

## 1. Element-wise Arithmetic Operations

NumPy allows fast and efficient element-wise operations on arrays.

```
import numpy as np

# Creating two NumPy arrays
a = np.array([1, 2, 3, 4])
b = np.array([5, 6, 7, 8])

# Element-wise arithmetic operations
print("Addition:", a + b) # [6 8 10 12]
print("Subtraction:", a - b) # [-4 -4 -4 -4]
print("Multiplication:", a * b) # [5 12 21 32]
print("Division:", a / b) # [0.2 0.333 0.428 0.5]
print("Exponentiation:", a ** 2) # [1 4 9 16] (Square each element)
print("Modulus:", b % a) # [0 0 1 0] (Remainder of division)
```

### ❖ Key Takeaways:

*Element-wise operations are performed without using explicit loops.  
Efficient and faster than regular Python lists.*

## 2. Universal Functions (ufuncs)

NumPy provides vectorized functions that perform element-wise operations efficiently. These are called Universal Functions (ufuncs).

```

# Creating a sample array
arr = np.array([1, 2, 3, 4, 5])

# Applying universal functions
print("Square root:", np.sqrt(arr))
print("Exponential (e^x):", np.exp(arr))
print("Natural logarithm (ln):", np.log(arr))
print("Base-10 logarithm (log10):", np.log10(arr))
print("Sine:", np.sin(arr))
print("Cosine:", np.cos(arr))
print("Tangent:", np.tan(arr))

```

### Advantages of ufuncs:

- Faster execution using vectorized operations
- Automatically handles broadcasting for arrays of different shapes
- Supports parallel execution on multi-core processors

## 3. Statistical Functions

NumPy provides various functions to calculate statistical metrics for datasets.

```

# Creating a dataset
data = np.array([10, 20, 30, 40, 50])

# Calculating key statistics
print("Mean (Average):", np.mean(data)) # 30.0
print("Median (Middle value):", np.median(data)) # 30.0
print("Standard Deviation:", np.std(data)) # Measures data spread
print("Variance:", np.var(data)) # Square of standard deviation
print("Minimum value:", np.min(data)) # Smallest value
print("Maximum value:", np.max(data)) # Largest value
print("Index of Minimum Value:", np.argmin(data)) # Index of smallest value
print("Index of Maximum Value:", np.argmax(data)) # Index of largest value

```

### Applications in Machine Learning:

Used to normalize data for machine learning models.  
 Helps in data preprocessing and feature engineering.  
 Essential for data analysis and insights.

## 4. Linear Algebra Operations

NumPy provides a powerful module called numpy.linalg for linear algebra operations, which are essential in machine learning, physics, and engineering.

## Matrix Creation

```
# Creating two matrices
A = np.array([[1, 2], [3, 4]])
B = np.array([[5, 6], [7, 8]])

print("Matrix A:\n", A)
print("Matrix B:\n", B)
```

## Matrix Multiplication (Dot Product)

```
# Matrix multiplication
dot_product = np.dot(A, B)
print("Dot Product:\n", dot_product)
```



*Alternative:*

*You can also use the @ operator for matrix multiplication.*

```
print("Dot Product using @ operator:\n", A @ B)
```

## Matrix Transpose

```
# Swaps rows and columns
print("Transpose of A:\n", A.T)
```

## Determinant of a Matrix

```
# Computing determinant
print("Determinant of A:", np.linalg.det(A))
```

## Inverse of a Matrix

```
# Computing inverse
print("Inverse of A:\n", np.linalg.inv(A))
```

## Eigenvalues and Eigenvectors

```
# Computing eigenvalues and eigenvectors
eigenvalues, eigenvectors = np.linalg.eig(A)
print("Eigenvalues:", eigenvalues)
print("Eigenvectors:\n", eigenvectors)
```

### Why Linear Algebra is Important?

Used in machine learning models, especially Principal Component Analysis (PCA).  
Essential for neural networks and deep learning.  
Helps in solving systems of linear equations.

## 5. Advanced Mathematical Functions

### Summation & Cumulative Sum

```
arr = np.array([1, 2, 3, 4])  
  
print("Sum of elements:", np.sum(arr)) # 10  
print("Cumulative sum:", np.cumsum(arr)) # [1 3 6 10]
```

### Finding Unique Elements & Counting Occurrences

```
arr = np.array([1, 2, 3, 1, 2, 3, 4, 4, 4])  
  
unique_elements, counts = np.unique(arr, return_counts=True)  
print("Unique elements:", unique_elements)  
print("Counts:", counts)
```

### Sorting an Array

```
arr = np.array([3, 1, 4, 1, 5, 9])  
  
print("Sorted array:", np.sort(arr))  
print("Indices of sorted elements:", np.argsort(arr))
```

### Finding Percentiles & Quantiles

```
data = np.array([10, 20, 30, 40, 50])  
  
print("25th percentile:", np.percentile(data, 25))  
print("50th percentile (median):", np.percentile(data, 50))  
print("75th percentile:", np.percentile(data, 75))
```

# Chapter 4: Advanced NumPy Operations

NumPy provides several advanced operations that optimize numerical computing, including:

- Broadcasting for efficient arithmetic operations on different-sized arrays.
- Sorting and searching for handling large datasets.
- File I/O operations to store and retrieve data efficiently.

These advanced features make NumPy an essential tool for data analysis, machine learning, and scientific computing.

## 1. Broadcasting

### What is Broadcasting?

Broadcasting allows NumPy to perform element-wise operations on arrays of different shapes without the need for manual reshaping or looping.

Instead of repeating values, NumPy automatically expands smaller arrays to match the dimensions of larger ones.

### Example: Broadcasting a 1D and 2D Array

```
import numpy as np

# Creating two arrays of different shapes
a = np.array([1, 2, 3])
b = np.array([[1], [2], [3]])

# Broadcasting: NumPy automatically expands 'b' to match 'a'
result = a + b

print("Array a:\n", a)
print("Array b:\n", b)
print("Broadcasted Addition Result:\n", result)
```

- How does this work?*

*a has a shape of (3,) → [1, 2, 3]*

*b has a shape of (3,1) → [[1], [2], [3]]*

*NumPy broadcasts b to (3,3) and adds element-wise with a.*

## More Examples of Broadcasting

### Broadcasting a Scalar to a Matrix

```
matrix = np.array([[1, 2, 3], [4, 5, 6]])
scalar = 10

result = matrix + scalar # The scalar gets "broadcasted" to each element
print(result)
```

#### 💡 Real-World Use Case:

*Broadcasting is widely used in machine learning for normalization, feature scaling, and batch operations.*

## 2. Sorting and Searching

### Sorting Arrays

Sorting is essential for data analysis, ranking, and organizing datasets.

#### Sorting a 1D Array

```
a = np.array([3, 1, 4, 1, 5])

print("Original Array:", a)
print("Sorted Array:", np.sort(a)) # Sorts in ascending order
```

#### Sorting a 2D Array

```
matrix = np.array([[3, 2, 1], [6, 5, 4]])

print("Original Matrix:\n", matrix)
print("Sorted along each row:\n", np.sort(matrix, axis=1)) # Row-wise sorting
print("Sorted along each column:\n", np.sort(matrix, axis=0)) # Column-wise sorting
```

### Searching for Elements in an Array

NumPy provides fast searching operations to find values in an array.

#### Finding the Index of the Maximum and Minimum Values

```
a = np.array([3, 1, 4, 1, 5, 9, 2])

print("Index of Maximum Value:", np.argmax(a)) # Returns index of max value
print("Index of Minimum Value:", np.argmin(a)) # Returns index of min value
```

## Finding Elements that Satisfy a Condition

```
a = np.array([10, 20, 30, 40, 50])

# Find elements greater than 25
filtered_elements = a[a > 25]
print("Elements greater than 25:", filtered_elements)
```

### Why is Searching Important?

Used in data preprocessing to filter values in large datasets.  
Helps in feature selection and anomaly detection in machine learning.

## 3. File I/O Operations with NumPy

### Saving and Loading Arrays

NumPy allows us to save and retrieve arrays efficiently in binary or text format.

#### a) Saving and Loading a NumPy Array (Binary Format)

```
# Create an array
a = np.array([1, 2, 3, 4, 5])

# Save array to a binary file
np.save('data.npy', a)

# Load the array from the file
b = np.load('data.npy')
print("Loaded Array:", b)
```

### Why Use .npy Format?

Stores data in a compact binary format.  
Faster loading and saving compared to CSV or text files.  
Retains data type and structure of arrays.

#### b) Saving and Loading Multiple Arrays

We can save multiple arrays in a single file using np.savez().

```
# Create multiple arrays
arr1 = np.array([1, 2, 3])
arr2 = np.array([10, 20, 30])

# Save both arrays in a single file
np.savez('multi_data.npz', array1=arr1, array2=arr2)

# Load the arrays from file
loaded_data = np.load('multi_data.npz')
print("Array 1:", loaded_data['array1'])
print("Array 2:", loaded_data['array2'])
```

### When to Use .npz Format?

*When working with large datasets where multiple arrays need to be stored efficiently.  
Reduces disk space usage while maintaining fast I/O speeds.*

## c) Saving and Loading CSV Files

Many real-world datasets are stored in CSV format. NumPy makes it easy to save and load CSV files.

Saving a NumPy Array as a CSV File

```
# Create a sample array
data = np.array([[1, 2, 3], [4, 5, 6]])

# Save to CSV
np.savetxt('data.csv', data, delimiter=',')
```

Loading Data from a CSV File

```
# Load CSV file
loaded_data = np.loadtxt('data.csv', delimiter=',')
print("Loaded CSV Data:\n", loaded_data)
```

### Use Case:

*CSV files are widely used in data science, machine learning, and finance.*

# Chapter 5: Project 1 - K-Nearest Neighbors from Scratch

## Theory and Concepts

The K-Nearest Neighbors (KNN) algorithm is a simple, yet powerful supervised learning algorithm used for both classification and regression tasks. The main idea behind KNN is that similar data points tend to be found near each other in feature space.

1. Calculate Distance: Compute the distance between the test point and all points in the training dataset (commonly using Euclidean distance).
2. Find Neighbors: Select the K closest points based on the computed distances.
3. Make a Decision:
  - o Classification: Assign the most common label among the K nearest neighbors.
  - o Regression: Take the average of the values of the K nearest neighbors.

## Implementation

```
import numpy as np
from collections import Counter

def euclidean_distance(a, b):
    return np.sqrt(np.sum((a - b) ** 2))

class KNN:
    def __init__(self, k=3):
        self.k = k

    def fit(self, X_train, y_train):
        self.X_train = X_train
        self.y_train = y_train

    def predict(self, X_test):
        predictions = [self._predict(x) for x in X_test]
        return np.array(predictions)

    def _predict(self, x):
        distances = [euclidean_distance(x, x_train) for x_train in
self.X_train]
        k_indices = np.argsort(distances)[:self.k]
        k_nearest_labels = [self.y_train[i] for i in k_indices]
        return Counter(k_nearest_labels).most_common(1)[0][0]
```

## Testing and Evaluation

```
from sklearn.model_selection import train_test_split
from sklearn.datasets import load_iris
from sklearn.metrics import accuracy_score

data = load_iris()
X, y = data.data, data.target
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

knn = KNN(k=5)
knn.fit(X_train, y_train)
y_pred = knn.predict(X_test)

print("Accuracy:", accuracy_score(y_test, y_pred))
```

# Chapter 6: Project 2 - Principal Component Analysis (PCA) Implementation

## Theory and Background

Principal Component Analysis (PCA) is a dimensionality reduction technique used to transform high-dimensional data into a lower-dimensional space while retaining as much variance as possible. It is widely used in data preprocessing, visualization, and noise reduction.

1. Compute the Mean: Center the data by subtracting the mean of each feature.
2. Calculate the Covariance Matrix: Measures how different features vary with one another.
3. Compute Eigenvalues and Eigenvectors: Solve for principal components.
4. Select Top Principal Components: Choose the top k eigenvectors based on eigenvalues.
5. Transform the Data: Project data onto the new principal component axes.

## Implementation

```
import numpy as np

def pca(X, n_components):
    # Center the data
    X_meaned = X - np.mean(X, axis=0)

    # Compute the covariance matrix
    cov_matrix = np.cov(X_meaned, rowvar=False)

    # Compute eigenvalues and eigenvectors
    eigenvalues, eigenvectors = np.linalg.eigh(cov_matrix)

    # Sort eigenvectors by eigenvalues in descending order
    sorted_indices = np.argsort(eigenvalues)[::-1]
    eigenvectors = eigenvectors[:, sorted_indices]
    eigenvalues = eigenvalues[sorted_indices]

    # Select top n_components eigenvectors
    selected_eigenvectors = eigenvectors[:, :n_components]

    # Transform data
    X_reduced = np.dot(X_meaned, selected_eigenvectors)
    return X_reduced
```

## Visualization and Application

```
import matplotlib.pyplot as plt
from sklearn.datasets import load_iris

# Load dataset
data = load_iris()
X, y = data.data, data.target

# Apply PCA
X_pca = pca(X, 2)

# Plot the transformed data
plt.scatter(X_pca[:, 0], X_pca[:, 1], c=y, cmap='viridis', edgecolor='k')
plt.xlabel('Principal Component 1')
plt.ylabel('Principal Component 2')
plt.title('PCA on Iris Dataset')
plt.colorbar(label='Target Label')
plt.show()
```

# Chapter 7: NumPy Best Practices

## Performance Optimization

NumPy provides several techniques to optimize performance when working with large datasets:

- Vectorization: Use NumPy's built-in functions instead of Python loops.
- Avoiding Copies: Use `view()` instead of `copy()` when possible to save memory.
- Memory Layout: Use `np.ascontiguousarray()` for better performance in C-optimized functions.
- Parallelization: Utilize NumPy's multithreading features for operations like `np.dot()`.

## Memory Management

Efficient memory usage is crucial for handling large datasets:

- Preallocate Arrays: Use `np.empty()` instead of repeatedly appending to lists.
- Use Appropriate Data Types: Convert `float64` arrays to `float32` if precision is not critical.
- Sparse Matrices: For large, sparse datasets, consider using `scipy.sparse` to save memory.

## Common Pitfalls and Solutions

- Floating Point Precision Errors: Due to limited precision, avoid direct equality checks for floating-point numbers.

```
if np.isclose(a, b): # Instead of a == b
```

- Unexpected Shape Changes: Ensure correct array shapes when performing operations to avoid broadcasting issues.
- Indexing Errors: Be cautious when slicing; using `arr[1]` instead of `arr[1, :]` may lead to shape mismatches.

# Chapter 8: NumPy and the Machine Learning Ecosystem

## NumPy with Pandas

Pandas is built on top of NumPy and provides DataFrame and Series objects for handling structured data. Some key integrations include:

### Creating DataFrames from NumPy Arrays:

```
import pandas as pd
import numpy as np
data = np.array([[1, 2], [3, 4], [5, 6]])
df = pd.DataFrame(data, columns=['A', 'B'])
print(df)
```

### Converting Pandas DataFrames to NumPy Arrays:

```
np_array = df.to_numpy()
```

## NumPy with Scikit-learn

Scikit-learn heavily relies on NumPy for numerical computations in machine learning models. Some common applications include:

### Feature Scaling using NumPy:

```
from sklearn.preprocessing import StandardScaler
X = np.array([[1, 2], [3, 4], [5, 6]])
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)
```

### Using NumPy Arrays in Model Training:

```
from sklearn.linear_model import LinearRegression
model = LinearRegression()
model.fit(X, np.array([1, 2, 3]))
```

## NumPy with TensorFlow and PyTorch

Deep learning frameworks like TensorFlow and PyTorch also integrate well with NumPy:

### Converting NumPy Arrays to Tensors in TensorFlow:

```
import tensorflow as tf
np_array = np.array([1, 2, 3, 4])
tensor = tf.convert_to_tensor(np_array)
```

### Converting NumPy Arrays to Tensors in PyTorch:

```
import torch  
torch_tensor = torch.tensor(np_array)
```

# Chapter 9: Appendix

## Quick Reference

- Creating an array: `np.array([1, 2, 3])`
- Generating random numbers: `np.random.rand(3, 3)`
- Reshaping an array: `arr.reshape(2, 3)`
- Computing mean: `np.mean(arr)`
- Matrix multiplication: `np.dot(A, B)`

## Further Resources

- NumPy Documentation: <https://numpy.org/doc/>
- Machine Learning with Python: <https://scikit-learn.org/>
- Deep Learning Frameworks: <https://www.tensorflow.org/> | <https://pytorch.org/>

## Glossary of Terms

- Array: A multi-dimensional container for numerical data.
- Broadcasting: A technique to perform operations on arrays of different shapes.
- Vectorization: The process of replacing explicit loops with array-based operations.
- Tensor: A multi-dimensional array used in deep learning frameworks.

---

## Conclusion

This handbook is an in-depth guide to NumPy both in fundamentals and advanced aspects, including NumPy in context to machine learning libraries. Application of NumPy functionality can speed performance and efficiency in your machine learning pipelines.

Thank you for persevering through NumPy, and hopefully this guide becomes something that is by your side in your quest through machine learning!