

Lecture 13 Handout

Advanced NumPy Operations

Linear Algebra, Random Numbers, and Applications

INF 605 - Introduction to Programming - Python

Prof. Rongyu Lin

Quinnipiac University

School of Computing and Engineering

Fall 2024

Required Reading

Textbook: Chapter 4, Sections 4.1.11, 4.2, 4.4, 4.8, 4.9, 4.10

Reference Notebooks: ch04/04_NumPy_Advanced.ipynb for advanced operations

Learning Objectives

By the end of this lecture, you will be able to:

1. **Master matrix operations and transposing** - Apply transpose, `swapaxes()`, and matrix multiplication with `@` operator
2. **Apply broadcasting principles** - Understand broadcasting rules and use them for efficient array operations
3. **Generate and control random numbers** - Use `np.random.default_rng()` for reproducible random data generation
4. **Use array-oriented programming** - Apply `np.where()` and `np.meshgrid()` for vectorized logic
5. **Perform linear algebra operations** - Use matrix inverse, solve linear systems, **fit** polynomials with `np.polyfit()`, and **calculate correlations with `np.corrcoef()`**
6. **Save and load NumPy arrays** - Use `np.save()`, `np.load()`, and `np.savez()` for data persistence
7. **Synthesize skills in complex applications** - Implement complete random walk simulations

Prerequisites Review

Building on Your NumPy Foundation:

From Lectures 12-13, you have mastered NumPy fundamentals including array creation (`np.array()`, `np.zeros()`, `np.ones()`, `np.arange()`, `np.linspace()`), array attributes (`shape`,

dtype, ndim, size), basic indexing and slicing, boolean and fancy indexing, element-wise operations, reshaping operations (`reshape()`, `flatten()`, `ravel()`), universal functions (`np.sqrt()`, `np.exp()`, `np.maximum()`), and statistical methods with the axis parameter (`mean()`, `sum()`, `std()`, `cumsum()`).

This lecture builds on that foundation by introducing advanced NumPy operations that are essential for scientific computing, data analysis, and machine learning. You will learn powerful techniques including transpose operations (building on your understanding of array dimensions), broadcasting rules (extending element-wise operations), random number generation (more powerful than Python's random module), vectorized conditional logic (replacing loops with `np.where()`), linear algebra operations (solving systems and fitting models), and comprehensive applications that synthesize all your NumPy skills.

Connection to Previous Knowledge: Transpose builds on Lecture 13's axis understanding (axis=0 goes DOWN, axis=1 goes ACROSS). Broadcasting extends Lecture 12's element-wise operations. `np.where()` extends Lecture 13's boolean indexing to conditional assignment. Linear algebra operations solve mathematical problems using all your array manipulation skills.

1 Section 1: Transposing and Matrix Operations

1.1 Understanding the Transpose Operation

Transpose is a fundamental matrix operation that flips rows and columns - imagine rotating a spreadsheet 90 degrees so that what were rows become columns and vice versa. In NumPy, you use the `.T` attribute (very simple!) to transpose any array, which is essential for many mathematical operations. Think of a class roster where rows are students and columns are test scores - transposing it makes rows become different tests and columns become different students, letting you analyze from a different perspective. The transpose operation is crucial for matrix multiplication because it lets you align dimensions correctly. For a 2D array with shape (3, 5), the transpose has shape (5, 3) - the dimensions swap positions.

```
1 # Basic transpose with .T attribute
2 import numpy as np
3
4 # Create a 3x5 array
5 arr = np.arange(15).reshape((3, 5))
6 print("Original array (3 rows, 5 columns):")
7 print(arr)
8 print(f"Shape: {arr.shape}")
9
10 # Transpose - swap rows and columns
11 transposed = arr.T
12 print("\nTransposed array (5 rows, 3 columns):")
13 print(transposed)
14 print(f"Shape: {transposed.shape}")
```

The transpose operation flipped the array - what were 3 rows became 3 columns, and 5 columns became 5 rows. Notice how the shape changed from (3, 5) to (5, 3). The first row [0, 1, 2, 3, 4] became the first column in the transposed version. This operation is extremely fast because NumPy doesn't actually move data in memory - it just changes how the array is viewed.

```
1 # Real application - analyzing test scores by subject
2 # Student scores: rows are students, columns are tests
3 scores = np.array([[85, 92, 78],
4                    [90, 88, 82],
5                    [75, 80, 85],
6                    [88, 92, 90]])
7
```

```

8 print(f"Original shape (students x tests): {scores.shape}")
9
10 # Transpose: now rows are tests, columns are students
11 scores_by_test = scores.T
12 print(f"Transposed shape (tests x students): {scores_by_test.shape}")
13
14 # Now we can easily calculate statistics per test
15 test_averages = scores_by_test.mean(axis=1)
16 print(f"\nAverage score per test: {test_averages}")

```

1.2 Matrix Multiplication and Dot Products

Matrix multiplication is a fundamental operation in linear algebra where you combine two matrices following specific rules to produce a new matrix - it is not just multiplying corresponding elements! For two matrices to multiply, the number of columns in the first must equal the number of rows in the second - if A is (2,3) and B is (3,4), then A @ B produces a (2,4) result. In Python, you use the @ operator (the modern way) or np.dot() function (the traditional way) to multiply matrices. This operation is crucial for machine learning, computer graphics, solving systems of equations, and analyzing relationships in data. Do not confuse matrix multiplication (@) with element-wise multiplication (*) - they are completely different operations!

```

1 # Matrix multiplication with @ operator
2 arr = np.array([[0, 1, 0],
3                 [1, 2, -2],
4                 [6, 3, 2],
5                 [-1, 0, -1],
6                 [1, 0, 1]])
7 # arr is 5x3
8
9 print(f"Array shape: {arr.shape}")
10
11 # Compute arr.T @ arr (transpose times original)
12 result = arr.T @ arr
13 print("Result of arr.T @ arr:")
14 print(result)
15 print(f"Result shape: {result.shape}") # (3, 3)

```

When we multiply arr.T (3x5) by arr (5x3), we get a (3x3) square matrix. This operation (transpose times original) is very common in statistics and machine learning - it computes how each column relates to every other column. The resulting matrix is symmetric, meaning it equals its own transpose.

```

1 # Practical example: calculating revenue
2 # Prices: 2 products in 3 stores
3 # Rows = products, Columns = stores
4 prices = np.array([[2.50, 2.75, 2.60], # Product 1
5                   [3.20, 3.00, 3.10]]) # Product 2
6
7 # Quantities sold: 3 stores, 2 days
8 # Rows = stores, Columns = days
9 quantities = np.array([[100, 120], # Store 1
10                       [150, 140], # Store 2
11                       [80, 90]]) # Store 3
12
13 # Calculate revenue: prices @ quantities
14 revenue = prices @ quantities
15 print(f"Revenue (products x days):\n{revenue}")
16 # Result shows total revenue for each product on each day
17 # First row: Product 1 revenue on day 1 and day 2
18 # Second row: Product 2 revenue on day 1 and day 2

```

1.3 Swapping Axes in Multidimensional Arrays

For arrays with more than 2 dimensions, `swapaxes()` gives you precise control over which dimensions to swap - it is like transpose but for any two axes you choose. While `.T` reverses all axes (in 2D, swaps rows and columns), `swapaxes()` lets you swap specific axes by number: axis 0, axis 1, axis 2, etc. This is useful for reorganizing 3D data like images (height x width x color channels) or scientific datasets (time x location x measurement). Understanding axis swapping helps you reshape data to match what different functions expect.

```
1 # Swapping specific axes in 3D array
2 arr = np.arange(16).reshape((2, 2, 4))
3 print("Original array (2x2x4):")
4 print(arr)
5 print(f"Shape: {arr.shape}")
6
7 # Swap axes 0 and 1
8 swapped = arr.swapaxes(0, 1)
9 print("\nAfter swapping axes 0 and 1:")
10 print(swapped)
11 print(f"Shape: {swapped.shape}")
12 # Shape stays (2,2,4) but data is reorganized
```

Key takeaway: Transpose (`.T`) reverses all axes, but `swapaxes(i, j)` swaps only the specified axes `i` and `j`. For most 2D work, use `.T`. For complex 3D+ arrays, `swapaxes()` gives precise control.

2 Section 2: Broadcasting Principles

2.1 Understanding Broadcasting Rules

Broadcasting is NumPy's powerful mechanism for performing operations on arrays of different shapes without creating copies - it is what makes NumPy code so concise and fast! Think of it like filling a multiplication table: when you multiply a single number by a list of numbers, the single number "broadcasts" to match the list's shape. NumPy automatically stretches smaller arrays to match larger ones during operations, following specific rules: (1) Arrays must have compatible shapes (dimensions match or one is 1), (2) Broadcasting happens along missing dimensions by repeating values. For example, adding a (3,) array to a (4,3) array broadcasts the (3,) array across all 4 rows automatically - no loops needed! Understanding broadcasting is essential for writing efficient NumPy code and avoiding unnecessary memory usage. Most shape errors you encounter are broadcasting compatibility issues!

```
1 # Broadcasting Compatibility Rules:
2 #
3 # Rule 1: Compare dimensions from right to left
4 # Rule 2: Dimensions are compatible if they are:
5 #         - Equal, OR
6 #         - One of them is 1
7 #
8 # Examples:
9 # (3, 4) + (4,)    Compatible! (4,) becomes (1, 4) then broadcasts to (3, 4)
10 # (3, 4) + (3, 1)  Compatible! (3, 1) broadcasts to (3, 4)
11 # (3, 4) + (3,)    ERROR! Cannot broadcast (3,) to match (3, 4)
```

```
1 # Basic broadcasting patterns
2 # Pattern 1: Scalar broadcasting
3 arr = np.array([[1, 2, 3],
4                 [4, 5, 6]])
5 result = arr + 10 # 10 broadcasts to (2,3)
6 print("Array + scalar:")
```

```

7 print(result)
8
9 # Pattern 2: 1D array broadcasting to 2D
10 row = np.array([10, 20, 30])
11 result = arr + row # row broadcasts across all rows
12 print("\nArray + row:")
13 print(result)

```

In the first example, the scalar 10 broadcasts to match the array's (2,3) shape, adding 10 to every element. In the second, the (3,) row array broadcasts across both rows of the (2,3) array, adding corresponding values to each column. This is equivalent to manually copying the row twice and then adding, but NumPy does it automatically and efficiently.

```

1 # Broadcasting with columns
2 arr = np.array([[1, 2, 3],
3                 [4, 5, 6],
4                 [7, 8, 9]])
5
6 # Add different value to each row (broadcast column)
7 col = np.array([[10],
8                 [20],
9                 [30]])
10 result = arr + col
11 print("Array + column:")
12 print(result)
13 # [[11, 12, 13],
14 #  [24, 25, 26],
15 #  [37, 38, 39]]

```

The column array with shape (3,1) broadcasts across all columns, adding 10 to the first row, 20 to the second row, and 30 to the third row. Notice how the shape (3,1) is compatible with (3,3) because the second dimension is 1.

2.2 Practical Broadcasting - Data Normalization

One of the most common uses of broadcasting is data normalization - transforming data to have mean of 0 and standard deviation of 1 (called "standardization"). This is crucial in machine learning and statistics because it puts all features on the same scale, making comparisons fair. Using broadcasting, you can calculate column means and standard deviations, then subtract means and divide by stds in just two lines - no loops! This is exactly the kind of operation where broadcasting shines: clean, readable, fast code that operates on entire datasets at once. Understanding this pattern will help you with Assignment 4 Problem 6 and many real data analysis tasks.

```

1 # Normalizing data with broadcasting
2 # Sample data: rows = observations, columns = features
3 data = np.array([[1, 200, 3],
4                  [4, 250, 6],
5                  [7, 300, 9],
6                  [10, 350, 12]])
7
8 # Calculate mean and std for each column
9 mean = data.mean(axis=0) # Shape: (3,)
10 std = data.std(axis=0) # Shape: (3,)
11
12 # Normalize: broadcasting makes this clean!
13 normalized = (data - mean) / std
14
15 print("Original data:")
16 print(data)
17 print(f"\nMeans: {mean}")

```

```

18 print(f"Stds: {std}")
19 print("\nNormalized data:")
20 print(normalized)

```

Broadcasting happened twice: first when subtracting mean (3,) from data (4,3), then dividing by std (3,). Each column now has mean approximately 0 and std approximately 1, making all features comparable despite originally different scales (1-10 vs 200-350). This two-line pattern - `(data - mean) / std` - is fundamental in data science.

3 Section 3: Random Number Generation

3.1 Modern Random Number API - `np.random.default_rng()`

NumPy's random number generation is dramatically faster than Python's built-in random module - we are talking 10-100 times faster for large arrays! The modern approach uses `np.random.default_rng()` to create a random number generator object, which then has methods for generating different types of random numbers. Think of it like creating a specialized random number machine that you can configure and use repeatedly - you create the generator once, then call its methods many times. This API replaced the old `np.random` functions (which still work but are discouraged) to give better control and reproducibility. Whether you need uniform random numbers, normally distributed data, random integers, or samples from datasets, the Generator object has optimized methods for all of these.

```

1  # Creating and using the random number generator
2  # Create random number generator with seed
3  rng = np.random.default_rng(seed=12345)
4
5  # Generate uniform random numbers [0, 1)
6  uniform = rng.random(5)
7  print("Uniform [0,1):", uniform)
8
9  # Generate standard normal (mean=0, std=1)
10 normal = rng.standard_normal(10)
11 print("\nStandard normal:", normal)
12
13 # Generate random integers
14 dice_rolls = rng.integers(1, 7, size=20) # 20 dice rolls
15 print("\nDice rolls:", dice_rolls)

```

We created one generator with a seed (12345), then used it multiple times to generate different types of random numbers. The seed makes results reproducible - running this code again with same seed produces identical "random" numbers. This is essential for debugging and sharing results.

3.2 Distributions and Seeding for Reproducibility

Random number generators can produce values from many different probability distributions - uniform (all values equally likely), normal/Gaussian (bell curve), integers (discrete values), and many more. In real work, you almost always want seeded random generation for reproducibility - if you discover interesting results from a simulation, you need to be able to reproduce them exactly! Setting a seed (any number) makes the "random" sequence deterministic and repeatable, which is essential for debugging, sharing results, and scientific reproducibility. Think of the seed as a starting point for the random sequence - same seed always produces same sequence.

```

1 # Different distributions with seeding
2 rng = np.random.default_rng(seed=42)
3

```

```

4 # Uniform distribution
5 uniform_data = rng.random(5)
6 print("Uniform [0,1):", uniform_data)
7
8 # Custom normal distribution (mean=100, std=15)
9 iq_scores = rng.normal(loc=100, scale=15, size=10)
10 print("\nSimulated IQ scores:", iq_scores.astype(int))
11
12 # Random integers from range
13 lottery = rng.integers(1, 50, size=6) # 6 numbers between 1-49
14 print("\nLottery numbers:", np.sort(lottery))

```

```

1 # Demonstrating reproducibility
2 # Two generators with same seed produce identical results
3 rng1 = np.random.default_rng(seed=999)
4 rng2 = np.random.default_rng(seed=999)
5
6 sample1 = rng1.random(5)
7 sample2 = rng2.random(5)
8
9 print("Sample 1:", sample1)
10 print("Sample 2:", sample2)
11 print("Identical?", np.array_equal(sample1, sample2)) # True!

```

The two generators with identical seeds produced exactly the same "random" numbers. This is crucial for scientific reproducibility - anyone running your code with the same seed will get identical results.

3.3 Sampling and Permutations

Beyond generating random numbers from distributions, you often need to sample from existing data (picking random elements) or shuffle data randomly. The generator's `choice()` method lets you sample elements from an array, with or without replacement (putting them back), which is perfect for creating random train/test splits or selecting random subsets. The `shuffle()` method randomly reorders an array in place, useful for randomizing data order before processing. These operations are fundamental for machine learning workflows, statistical sampling, and any situation where you need random subsets of your data.

```

1 # Sampling and shuffling operations
2 rng = np.random.default_rng(seed=42)
3
4 # Sample without replacement (no duplicates)
5 names = np.array(["Alice", "Bob", "Charlie", "David", "Eve"])
6 selected = rng.choice(names, size=3, replace=False)
7 print("Random sample (3 names):", selected)
8
9 # Shuffle array in place
10 arr = np.arange(10)
11 print("\nBefore shuffle:", arr)
12 rng.shuffle(arr)
13 print("After shuffle:", arr)
14 # Note: shuffle modifies the original array!

```

Use `choice()` to sample elements (with or without replacement), and `shuffle()` to randomize order. Both are essential for data science workflows like cross-validation and random sampling.

4 Section 4: Array-Oriented Programming

4.1 Vectorized Conditional Logic with np.where()

The `np.where()` function is NumPy's answer to vectorized if-else logic - instead of looping through elements one by one asking "if this, then that", `np.where()` evaluates the entire array at once and assigns values based on conditions. Think of it like Excel's IF function but applied to thousands of cells simultaneously: `IF(condition, value_if_true, value_if_false)`. This is dramatically faster than Python loops or list comprehensions for large arrays, and the code is often more readable too. You can use `np.where()` for simple substitution (replace negatives with zero) or complex conditional transformations. This function is essential for any time you need different operations based on array values - data cleaning, categorization, or conditional calculations.

```
1 # Basic np.where() usage
2 xarr = np.array([1.1, 1.2, 1.3, 1.4, 1.5])
3 yarr = np.array([2.1, 2.2, 2.3, 2.4, 2.5])
4 cond = np.array([True, False, True, True, False])
5
6 # Without np.where() - SLOW with loops
7 result_slow = [(x if c else y)
8                 for x, y, c in zip(xarr, yarr, cond)]
9 print("Loop result:", result_slow)
10
11 # With np.where() - FAST and clean!
12 result_fast = np.where(cond, xarr, yarr)
13 print("np.where result:", result_fast)
14 # Same result, much faster!
```

`np.where(cond, xarr, yarr)` says "where condition is True, take from xarr; where False, take from yarr". The result has values from xarr at positions 0, 2, 3 (where cond is True) and values from yarr at positions 1, 4 (where cond is False).

```
1 # Common np.where() patterns
2 rng = np.random.default_rng(seed=42)
3 arr = rng.standard_normal((4, 4))
4
5 print("Original array:")
6 print(arr)
7
8 # Pattern 1: Replace positive values with 2, keep negative
9 result1 = np.where(arr > 0, 2, arr)
10 print("\nPositive to 2, negative unchanged:")
11 print(result1)
12
13 # Pattern 2: Replace positive with 2, negative with -2
14 result2 = np.where(arr > 0, 2, -2)
15 print("\nPositive to 2, negative to -2:")
16 print(result2)
```

These patterns show how `np.where()` can transform arrays based on conditions. The first pattern keeps negative values unchanged (using `arr` as the false value), while the second pattern assigns a constant -2 to all negative values.

4.2 Creating Coordinate Grids with np.meshgrid()

The `np.meshgrid()` function creates coordinate matrices from coordinate vectors, which is essential for evaluating functions over 2D grids - imagine creating a checkerboard where you want to calculate something at every square. You provide 1D arrays for x-coordinates and y-coordinates, and `meshgrid` gives you 2D arrays where every combination of (x,y) is represented.

This is fundamental for plotting 3D surfaces, creating grids for image processing, or evaluating mathematical functions over a region. For example, to evaluate the function $z = \sqrt{x^2 + y^2}$ over a grid, you first create coordinate grids with meshgrid, then apply the function using broadcasting.

```

1 # Creating and using meshgrid
2 # Create 1D arrays for coordinates
3 points = np.arange(-5, 5, 0.1) # 100 points from -5 to 5
4 print(f"Points shape: {points.shape}") # (100,)
5
6 # Create 2D coordinate grids
7 xs, ys = np.meshgrid(points, points)
8 print(f"xs shape: {xs.shape}") # (100, 100)
9 print(f"ys shape: {ys.shape}") # (100, 100)
10
11 # Calculate function over grid: distance from origin
12 z = np.sqrt(xs**2 + ys**2)
13 print(f"z shape: {z.shape}") # (100, 100)
14 # Now z contains the distance from origin at every (x,y) point

```

Meshgrid converted two 1D arrays (100,) into two 2D arrays (100, 100) where every (x,y) coordinate pair is represented. We then calculated $z = \sqrt{x^2 + y^2}$ over the entire grid using broadcasting - this evaluated the function at 10,000 points without any loops!

5 Section 5: Linear Algebra Operations

5.1 Matrix Inverse and Solving Linear Systems

Linear algebra operations in NumPy solve systems of equations and perform matrix calculations you learned in algebra class - but now working with large matrices automatically! The matrix inverse is like division for matrices: just as dividing by x is the same as multiplying by 1/x, multiplying by a matrix inverse "undoes" the original matrix. NumPy's `np.linalg.inv()` computes matrix inverses, and `np.linalg.solve()` efficiently solves systems of linear equations ($Ax = b$) without explicitly computing the inverse. These operations are fundamental in many fields: solving circuit equations in engineering, fitting models in statistics, transforming coordinates in graphics, and optimization in machine learning.

```

1 # Matrix inverse
2 rng = np.random.default_rng(seed=42)
3
4 # Create symmetric positive definite matrix
5 X = rng.standard_normal((5, 5))
6 mat = X.T @ X # This creates a matrix we can invert
7
8 # Calculate inverse
9 inv_mat = np.linalg.inv(mat)
10
11 # Verify: matrix times inverse should be identity
12 identity = mat @ inv_mat
13 print("mat @ inv(mat) should be identity:")
14 print(identity)
15 # Should be very close to identity matrix (1s on diagonal, 0s elsewhere)

```

The matrix multiplied by its inverse produces an identity matrix (like multiplying a number by its reciprocal gives 1). Small rounding errors cause diagonal to be approximately 1.0 instead of exactly 1.0, and off-diagonal approximately 0.0 instead of exactly 0.0.

```

1 # Solving linear systems
2 # Solve system of equations:
3 # 3x + y = 9

```

```

4 # x + 2y = 8
5
6 # Matrix form: Ax = b
7 A = np.array([[3, 1],
8               [1, 2]])
9 b = np.array([9, 8])
10
11 # Solve for x
12 x = np.linalg.solve(A, b)
13 print(f"Solution: x = {x}")
14
15 # Verify solution
16 print(f"Check: A @ x = {A @ x}") # Should equal b

```

The solution $x = [2, 3]$ satisfies both equations: $3(2) + 3 = 9$ and $2 + 2(3) = 8$. Always use `np.linalg.solve()` instead of computing inverse explicitly - it is more efficient and numerically stable.

5.2 Polynomial Fitting with `np.polyfit()` - CRITICAL FOR ASSIGNMENT 4

Polynomial fitting is one of the most common data analysis tasks - finding the best polynomial curve that matches your data points. NumPy's `np.polyfit()` does this automatically using least squares fitting (minimizing the squared errors between the curve and data points). You provide x-values, y-values, and the polynomial degree (1 for linear, 2 for quadratic, 3 for cubic, etc.), and it returns the coefficients of the polynomial that best fits your data. This is fundamental for trend analysis, making predictions, and understanding relationships in data - is the relationship linear? curved? exponential-looking? Polynomial fitting is used everywhere: economics (trend forecasting), physics (fitting experimental data), machine learning (polynomial regression), and Assignment 4 Problem 7 explicitly requires using `np.polyfit()`!

```

1 # Linear fit (degree 1)
2 # Sample data with linear trend
3 x_data = np.array([1, 2, 3, 4, 5])
4 y_data = np.array([2.1, 3.9, 6.2, 8.1, 9.8])
5
6 # Fit polynomial of degree 1 (linear: y = mx + b)
7 coefficients = np.polyfit(x_data, y_data, deg=1)
8 slope, intercept = coefficients
9
10 print(f"Best fit line: y = {slope:.2f}x + {intercept:.2f}")
11
12 # Make predictions
13 x_new = 6
14 y_pred = slope * x_new + intercept
15 print(f"Prediction at x=6: {y_pred:.2f}")

```

`polyfit()` found the line that best fits the data points. The slope approximately 2 and intercept approximately 0 mean the relationship is approximately $y = 2x$. We can use these coefficients to predict y-values for new x-values not in the original data.

```

1 # Quadratic fit (degree 2)
2 # Data with curved relationship
3 x_data = np.array([0, 1, 2, 3, 4, 5])
4 y_data = np.array([1, 2.5, 5.8, 10.9, 17.8, 26.5])
5
6 # Fit polynomial of degree 2 (quadratic: y = ax^2 + bx + c)
7 coefficients = np.polyfit(x_data, y_data, deg=2)
8 a, b, c = coefficients
9

```

```

10 print(f"Best fit curve: y = {a:.2f}x^2 + {b:.2f}x + {c:.2f}")
11
12 # Evaluate polynomial at new points
13 x_new = np.array([6, 7, 8])
14 y_pred = a*x_new**2 + b*x_new + c
15 print(f"Predictions: {y_pred}")

```

For curved data, a quadratic polynomial (degree 2) provides a better fit than a straight line. The three coefficients (a, b, c) define the parabola that best matches the data points. This pattern - fitting, extracting coefficients, making predictions - is exactly what Assignment 4 Problem 7 requires.

5.3 Correlation Analysis with np.corrcoef() - CRITICAL FOR ASSIGNMENT 4

Correlation measures how strongly two variables are related - do they increase together? decrease together? or are they unrelated? NumPy's `np.corrcoef()` calculates correlation coefficients between variables, producing a correlation matrix showing all pairwise relationships. Correlation values range from -1 (perfect negative correlation) through 0 (no correlation) to +1 (perfect positive correlation). For example, height and weight have positive correlation (taller people tend to weigh more), while temperature and heating bills have negative correlation (higher temperature means lower bills). Understanding correlation is fundamental for data analysis - identifying related variables, finding redundant features, or discovering hidden relationships. Assignment 4 Problem 7 requires using `np.corrcoef()` to analyze relationships in data!

```

1 # Basic correlation analysis
2 # Two variables: study hours and test scores
3 study_hours = np.array([2, 3, 4, 5, 6, 7, 8])
4 test_scores = np.array([65, 70, 75, 80, 85, 88, 92])
5
6 # Calculate correlation matrix
7 correlation_matrix = np.corrcoef(study_hours, test_scores)
8 print("Correlation matrix:")
9 print(correlation_matrix)
10
11 # Extract correlation coefficient (off-diagonal)
12 correlation = correlation_matrix[0, 1]
13 print(f"\nCorrelation between study hours and scores: {correlation:.3f}")
14 # Interpretation: value near 1 means strong positive relationship

```

The correlation matrix is 2x2: diagonal elements are 1.0 (each variable correlates perfectly with itself), and off-diagonal elements show correlation between the two variables. A value near 1.0 indicates strong positive correlation - more study hours are associated with higher scores.

```

1 # Multiple variables correlation
2 # Three variables: temperature, ice cream sales, sunglasses sales
3 temp = np.array([70, 75, 80, 85, 90, 95])
4 ice_cream = np.array([100, 120, 150, 180, 220, 250])
5 sunglasses = np.array([50, 55, 65, 75, 90, 100])
6
7 # Stack variables as rows
8 data = np.array([temp, ice_cream, sunglasses])
9
10 # Calculate all pairwise correlations
11 corr_matrix = np.corrcoef(data)
12 print("Correlation matrix:")
13 print(corr_matrix)
14 # All three variables positively correlated (all increase together)

```

With three variables, the correlation matrix is 3x3, showing all nine pairwise correlations. The diagonal is still 1.0, and off-diagonal elements show how each pair of variables relates. All three variables show strong positive correlation because they all tend to increase together with warmer weather.

5.4 Other Linear Algebra Functions

NumPy's `linalg` module includes many other linear algebra operations you might encounter: determinant (`np.linalg.det`), eigenvalues/eigenvectors (`np.linalg.eig`), QR decomposition (`np.linalg.qr`), and SVD (`np.linalg.svd`). These advanced operations are crucial in machine learning and scientific computing but beyond our current scope - just know they exist for future reference. For this course, focus on matrix multiplication (`@`), inverse, solve, polyfit, and `corrcoef` - these cover 90 percent of practical linear algebra needs.

```
1 # Other operations available (awareness only):
2 # np.linalg.det(A)           # Determinant
3 # np.linalg.eig(A)           # Eigenvalues and eigenvectors
4 # np.linalg.qr(A)            # QR decomposition
5 # np.linalg.svd(A)           # Singular Value Decomposition
6 # np.linalg.matrix_rank(A)   # Matrix rank
7
8 # For this course, focus on:
9 # - Matrix multiplication: @
10 # - Solving systems: np.linalg.solve()
11 # - Polynomial fitting: np.polyfit()
12 # - Correlation: np.corrcoef()
```

6 Section 6: File I/O for Arrays

6.1 Saving and Loading Single Arrays

After spending hours processing data and creating NumPy arrays, you do not want to lose your work or recompute everything next time! NumPy provides efficient binary file formats (`.npy`) for saving and loading arrays - much faster and more space-efficient than text formats like CSV. Use `np.save()` to save a single array to disk, and `np.load()` to load it back - NumPy automatically preserves the exact shape, data type, and values. Think of it like saving a Word document - you can close your program, come back later, and pick up exactly where you left off.

```
1 # Save and load single array
2 # Create array
3 arr = np.arange(10)
4 print("Original array:", arr)
5
6 # Save to file (automatically adds .npy extension)
7 np.save("my_array", arr)
8 print("Array saved to 'my_array.npy'")
9
10 # Load from file
11 loaded = np.load("my_array.npy")
12 print("Loaded array:", loaded)
13 print("Arrays equal?", np.array_equal(arr, loaded)) # True
```

`np.save()` writes the array in NumPy's efficient binary format. `np.load()` reads it back perfectly - same values, same shape, same dtype. The `.npy` file is much smaller and faster to read than text formats.

6.2 Saving Multiple Arrays with np.savez()

Often you need to save multiple related arrays together - think of input data, processed results, and calculated statistics all belonging to one analysis. NumPy's `np.savez()` creates compressed archives (like zip files) that store multiple arrays with named keys, and `np.savez_compressed()` adds extra compression for large files. You load the archive with `np.load()` and access individual arrays by their names like a dictionary.

```
1 # Save multiple arrays together
2 # Create multiple arrays
3 arr1 = np.arange(10)
4 arr2 = np.random.randn(5, 3)
5
6 # Save multiple arrays with names
7 np.savez("array_archive.npz", first=arr1, second=arr2)
8 print("Multiple arrays saved to 'array_archive.npz'")
9
10 # Load archive
11 archive = np.load("array_archive.npz")
12
13 # Access by name
14 loaded_arr1 = archive["first"]
15 loaded_arr2 = archive["second"]
16
17 print("First array:", loaded_arr1)
18 print("Second array shape:", loaded_arr2.shape)
19
20 # For large arrays, use compressed version:
21 np.savez_compressed("compressed_archive.npz", first=arr1, second=arr2)
```

Use `np.save()` for single arrays, `np.savez()` for multiple arrays, and `np.savez_compressed()` when file size matters. This preserves your work and makes analysis reproducible.

7 Section 7: Comprehensive Example - Random Walks

7.1 Understanding Random Walks

Random walks are a fundamental model in mathematics, physics, finance, and computer science - imagine flipping a coin repeatedly and taking a step left or right based on heads or tails. This simple process models everything from stock prices to molecular motion to gamblers' fortunes! What makes this example amazing is how NumPy lets us simulate not just one walk, but THOUSANDS simultaneously with just a few lines of code. We will start with a Python loop version (slow, many lines), then show the NumPy version (fast, few lines), then simulate 5000 walks at once and analyze them all together. This demonstrates everything we have learned: random numbers, cumsum, boolean operations, axis parameter, argmax, broadcasting - all working together in a real, meaningful application!

7.2 Pure Python Implementation

Let us first implement a random walk the "traditional" Python way using loops - this helps us understand what a random walk is before we accelerate it with NumPy. We will flip a coin 1000 times (`random.randint()`), take steps based on results, and track our position over time. This version is clear but slow - it would take forever to simulate thousands of walks.

```
1 # Python loop version
2 import random
3
4 # Start at position 0
5 position = 0
```

```

6 walk = [position]
7 nsteps = 1000
8
9 # Take 1000 random steps
10 for _ in range(nsteps):
11     step = 1 if random.randint(0, 1) else -1
12     position += step
13     walk.append(position)
14
15 print(f"Final position: {walk[-1]}")
16 print(f"Maximum position: {max(walk)}")
17 print(f"Minimum position: {min(walk)}")

```

After 1000 random steps, we ended at some random position. The walk randomly drifts away from origin - sometimes reaching high values, sometimes low. This took Python many loop iterations to compute.

7.3 NumPy Implementation - Single Walk

Now let us implement the same random walk using NumPy - watch how much shorter and faster this is! Instead of a loop, we generate all 1000 random numbers at once, convert them to steps with `np.where()`, then use `cumsum()` to calculate positions. This is 10-100 times faster than the loop version and demonstrates the power of vectorized thinking: generate all random data, transform it, accumulate it - no loops needed!

```

1 # NumPy version - single walk
2 rng = np.random.default_rng(seed=12345)
3 nsteps = 1000
4
5 # Generate all random values at once
6 draws = rng.integers(0, 2, size=nsteps) # 1000 random 0s and 1s
7
8 # Convert to steps: 0 -> +1, 1 -> -1
9 steps = np.where(draws == 0, 1, -1)
10
11 # Calculate cumulative position (this is the walk!)
12 walk = steps.cumsum()
13
14 print(f"Final position: {walk[-1]}")
15 print(f"Maximum position: {walk.max()}")
16 print(f"Minimum position: {walk.min()}")
17
18 # Find first crossing time for plus/minus 10
19 crossing_time = (np.abs(walk) >= 10).argmax()
20 print(f"First crossed +/-10 at step: {crossing_time}")

```

We generated the entire walk in just 4 lines of NumPy code! `cumsum()` calculated all 1000 positions at once. `argmax()` on boolean array finds the first True value - this tells us when the walk first crossed the plus/minus 10 threshold.

7.4 Simulating 5000 Walks Simultaneously

Here is where NumPy truly shines: we can simulate 5000 different random walks SIMULTANEOUSLY using 2D arrays! Instead of a 1D array of 1000 steps, we create a 2D array with shape (5000, 1000) - 5000 walks, each with 1000 steps. We use `cumsum(axis=1)` to accumulate along the steps dimension (`axis=1`), computing all 5000 walks at once. Then we analyze all walks together: how many crossed plus/minus 30? what was the average crossing time? This would take HOURS with Python loops but runs in seconds with NumPy!

```

1 # Simulate 5000 walks simultaneously

```

```

2 nwalks = 5000
3 nsteps = 1000
4
5 # Generate ALL random draws at once: 5000 x 1000 array!
6 draws = rng.integers(0, 2, size=(nwalks, nsteps))
7
8 # Convert all to steps
9 steps = np.where(draws > 0, 1, -1)
10
11 # Calculate ALL 5000 walks simultaneously!
12 walks = steps.cumsum(axis=1) # cumsum along steps (axis=1)
13
14 print(f"Walks array shape: {walks.shape}") # (5000, 1000)
15 print(f"Maximum any walk reached: {walks.max()}")
16 print(f"Minimum any walk reached: {walks.min()}")

```

```

1 # Analyze the 5000 walks
2 # How many walks crossed +/-30?
3 hits30 = (np.abs(walks) >= 30).any(axis=1)
4 print(f"\nOut of {nwalks} walks:")
5 print(f"    {hits30.sum()} walks crossed +/-30")
6 print(f"    {nwalks - hits30.sum()} walks stayed within +/-30")
7
8 # Average crossing time for walks that crossed
9 crossing_times = (np.abs(walks[hits30]) >= 30).argmax(axis=1)
10 print(f"    Average crossing time: {crossing_times.mean():.1f} steps")
11
12 # Distribution of crossing times
13 print(f"    Fastest crossing: {crossing_times.min()} steps")
14 print(f"    Slowest crossing: {crossing_times.max()} steps")

```

We simulated 5 MILLION individual steps (5000 walks times 1000 steps) in under a second! The analysis shows fascinating statistical patterns: roughly 3400 walks crossed plus/minus 30, taking an average of approximately 500 steps to do so. This demonstrates NumPy's incredible power - what would take hours with loops takes seconds with vectorization.

7.5 Synthesis: What We Demonstrated

Key points from the random walks example:

1. **Random generation:** Created 5 million random numbers instantly
2. **Broadcasting:** Applied `np.where()` to entire 2D array at once
3. **Axis operations:** Used `cumsum(axis=1)` to process each walk independently
4. **Boolean operations:** Found crossing events with `any()` along `axis=1`
5. **Advanced indexing:** Selected only walks that crossed with `walks[hits30]`
6. **Analysis:** Used `argmax()`, `mean()`, `min()`, `max()` to understand results

Real-world applications of random walks:

- **Finance:** Monte Carlo simulation for option pricing
- **Physics:** Brownian motion and diffusion models
- **Biology:** Random mutations and genetic drift
- **Computer Science:** Algorithm analysis and randomized algorithms

- **Statistics:** Understanding random processes and probability

The power of NumPy: Pure Python with loops takes approximately 30-60 seconds for this simulation. NumPy vectorized version takes approximately 0.5-1 seconds. Speed improvement: 30-60 times faster!

This is why NumPy is essential for data science and scientific computing - it lets you work with large datasets efficiently, writing clear, concise code that runs at nearly the speed of compiled languages.

Key Takeaways

Essential Concepts to Remember:

- Transpose is simple: Just use `.T` for 2D arrays
- Broadcasting is powerful: Eliminates loops, makes code clean
- Modern random API: `np.random.default_rng()` is the new standard
- `np.where()` is vectorized if-else: Fundamental for array-oriented thinking
- `np.polyfit()` and `np.corrcoef()`: Essential for Assignment 4 Problem 7!
- Random walks synthesis: Demonstrates everything learned, shows NumPy's power
- Real applications: Connect every technique to practical uses

Connection to Assignment 4

Direct Assignment 4 Support:

- **Problem 4:** Reshaping and transpose operations
- **Problem 5:** Boolean filtering with percentiles
- **Problem 6:** Normalization with broadcasting
- **Problem 7:** Polynomial fitting with `np.polyfit()` and correlation with `np.corrcoef()`
- **Problem 8:** Comprehensive analytics combining all techniques

Preview of Next Lecture

In the next lecture, we will transition from NumPy arrays to pandas DataFrames - pandas builds on NumPy, adding labels, handling mixed types, managing missing data, and providing powerful tools for real-world data analysis. All your NumPy skills transfer directly to pandas!