# Data Science Assignment (2)

**Full Name:** Ibrahim Qahtan Adnan

**Theoretical Group:** B

## 1) What did you learn from the assigned reading?

This section of the book introduces the foundational object for all data science in Python: the **NumPy array** (ndarray). The key takeaway is that NumPy provides a high-performance data structure that is far more efficient and powerful than standard Python lists for numerical calculations.

The main topics covered are:

- **Array Creation:** How to create arrays from scratch. This includes:
  - np.array(): Converting existing Python lists or tuples.
  - np.arange(): Creating arrays with regularly spaced values (like Python's range but for arrays).
  - np.linspace(): Creating arrays with a specific number of points between a start and end value.
  - np.zeros(), np.ones(), np.empty(): Creating arrays filled with zeros, ones, or uninitialized data, which is useful for reserving memory.
- **Array Attributes:** Inspecting the properties of an array, such as:
  - .shape: The dimensions of the array (e.g., (3, 4) for a 3x4 matrix).
  - .ndim: The number of dimensions (e.g., 2 for a matrix).
  - .size: The total number of elements in the array.
  - .dtype: The data type of the elements (e.g., float64, int32), which is crucial for memory efficiency.
- **Indexing and Slicing:** This is a major topic. It explains how to access and modify parts of an array. This goes far beyond Python lists and includes:

- o **Basic Slicing:** Accessing rows, columns, or sub-regions (e.g., arr[0:2, 5:10]).
- o **Boolean Indexing:** Using an array of True/False values to filter and select elements that meet a condition (e.g., arr[arr > 50]). This is extremely powerful.
- o **"Fancy" Indexing:** Using an array of integers to pick out specific elements in any order.
- **Vectorization:** This is the most important concept. NumPy allows you to perform batch operations on entire arrays without writing for loops. For example, to add 5 to every number in an array, you just write arr + 5. This "vectorized" code is not only cleaner to read but also hundreds or thousands of times faster, as the loops are executed in highly optimized, pre-compiled C code.

**2) Can your learning help you solve real-life problems?**
   **A) If yes, describe how (Write approximately 200–300 words):**

Yes, this learning is not just helpful—it is absolutely essential for solving almost any real-life data-driven problem.

NumPy is the bedrock of the entire scientific Python ecosystem (which includes Pandas, SciPy, Matplotlib, and scikit-learn). The reason it's so critical is speed and scale. Real-world problems involve datasets that are far too large and complex for standard Python lists and loops.

For example, a high-definition color image is just a 3D NumPy array (Height x Width x 3 channels). To increase its brightness, you don't loop through millions of pixels; you simply multiply the entire array by 1.2. This is a single, vectorized operation.

In data science, you might have a file with 10 million temperature readings. With NumPy, you can find the average, maximum, and minimum temperatures, or identify all readings above a certain threshold, in a fraction of a second.

In finance, analysts model stock market fluctuations using massive arrays of price data to run simulations. In scientific research, complex physical systems, from brain activity to climate patterns, are modeled using multi-dimensional arrays representing space and time.

The "vectorization" and "indexing" techniques from this reading are the high-performance tools that make these large-scale tasks feasible. They allow you to express complex mathematical operations clearly and have them execute with incredible speed.

**B) If you are able to answer question A, provide an example with a clear explanation and include Python code to support your answer.**

**Example: Filtering Sensor Data**

Here is a practical example that directly uses the concepts of array creation, vectorization, and boolean indexing.

**A) The Problem:** Imagine you are monitoring a network of 100 environmental sensors. You receive their temperature readings every minute. Your data comes in as a large list. You need to perform two tasks quickly:

1.  Convert all readings from Celsius to Fahrenheit.

2.  Count how many of those readings are "critical" (e.g., over 90°F).

If you have millions of readings, looping through a Python list is extremely slow.

**B) The NumPy Solution:** Using NumPy, you can perform these tasks almost instantly.

1.  First, create a sample NumPy array (in a real app, you'd load this from a file).

2.  Use a vectorized operation to convert all temperatures at once.

3.  Use a boolean mask to find all critical temperatures.

4.  Use np.sum() on the boolean mask to count them (since True equals 1 and False equals 0).

## C) Python Code:

```python
import random

print("Running standard Python loop...")
celsius_list = [random.uniform(15.0, 40.0) for _ in range(1_000_000)]
fahrenheit_list = []
critical_count = 0

for temp in celsius_list:
    f_temp = (temp * 9/5) + 32
    fahrenheit_list.append(f_temp)
    if f_temp > 90:
        critical_count += 1

print(f"Standard Python: Found {critical_count} critical readings.\n")
import numpy as np

print("Running vectorized NumPy solution...")

celsius_arr = np.array(celsius_list)

fahrenheit_arr = (celsius_arr * 9/5) + 32
critical_mask = fahrenheit_arr > 90

num_critical = np.sum(critical_mask)

print(f"NumPy: Found {num_critical} critical readings.")
```

**Explanation:** The Python for loop has to individually check one million items. The NumPy solution ((celsius_arr * 9/5) + 32 and fahrenheit_arr > 90) runs a single command on the entire dataset at C-language speed. For large datasets, the NumPy version isn't just a little faster; it can be **100x to 10,000x faster**, making it the only practical choice for real-world data analysis.