



DATA SCIENCE DEVELOPMENT TOOLS

PART 1: PYTHON, NUMPY

TOOLS FOR DATA MINING

- Start with Pickaxe and Shovel
- Move to jack hammer and other automated tools
- Python, SQL, R most common languages
- NumPy, Pandas, Matplotlib, Seaborn, Scikit-Learn
- Tableau, PowerBI
- AWS SageMaker, Dataiku, DataRobot, DataBricks



COMMON DATA SCIENCE PACKAGES FOR PYTHON

CS 180

Import name	Common alias	Description
<code>numpy</code>	<code>np</code>	NumPy includes functions and classes that aid in numerical computation. NumPy is used in many other data science packages.
<code>pandas</code>	<code>pd</code>	pandas provides methods and classes for tabular and time-series data.
<code>sklearn</code>	<code>sk</code>	scikit-learn provides implementations of many machine learning algorithms with a uniform syntax for preprocessing data, specifying models, fitting models with cross-validation, and assessing models.
<code>matplotlib.pyplot</code>	<code>plt</code>	matplotlib allows the creation of data visualizations in Python. The functions mostly expect NumPy arrays.
<code>seaborn</code>	<code>sns</code>	seaborn also allows the creation of data visualizations but works better with pandas DataFrame.
<code>scipy.stats</code>	<code>sp.stats</code>	SciPy provides algorithms and functions for computing problems that arise in science, engineering and statistics. <code>scipy.stats</code> provides the functions for statistics.
<code>statsmodels</code>	<code>sm</code>	statsmodels adds functionality to Python to estimate many different kinds of statistical models, make inferences from those models, and explore data.



NUMPY

- **Spelled:** NumPy, Pronounced “num-pie”
- **What is NumPy?** NumPy (Numerical Python) is the fundamental package for scientific computing with Python. It provides a high-performance multidimensional array object and tools for working with these arrays.
- **Why NumPy?**
 - **Speed:** NumPy arrays are more efficient and faster than Python lists for numerical operations, as they are implemented in C.
 - **Functionality:** It provides a rich set of functions for linear algebra, Fourier transforms, and random number generation.
 - **Foundation:** Many other data science libraries like Pandas, SciPy, and Scikit-learn are built on top of NumPy.

NUMPY ARRAY FUNCTIONS

Function	Parameters	Description
<code>array()</code>	<code>object</code> <code>dtype=None</code> <code>ndim=0</code>	Returns an array constructed from <code>object</code> . <code>object</code> must be a scalar or an ordered container, such as tuple or list. The array element type is inferred from <code>object</code> unless a <code>dtype</code> is specified. <code>ndim</code> is the minimum number of array dimensions.
<code>delete()</code>	<code>arr</code> <code>obj</code> <code>axis=None</code>	Deletes a slice of input array <code>arr</code> . <code>axis</code> is the axis along which to remove a slice. <code>obj</code> is the index of the slice along the axis.
<code>full()</code>	<code>shape</code> <code>fill_value</code> <code>dtype=None</code>	Returns an array filled with <code>fill_value</code> . The <code>shape</code> tuple specifies array shape. <code>dtype</code> specifies the array type. If <code>dtype=None</code> , the type is inferred from <code>fill_value</code> .
<code>insert()</code>	<code>arr</code> <code>obj</code> <code>values</code> <code>axis=None</code>	Inserts array <code>values</code> to input array <code>arr</code> . <code>axis</code> is the axis along which to insert. <code>obj</code> is the index before which <code>values</code> is inserted.
<code>zeros()</code>	<code>shape</code> <code>dtype=float</code>	Returns an array filled with zeros. The <code>shape</code> tuple specifies array shape. <code>dtype</code> specifies the array type.
<code>ones()</code>	<code>shape</code> <code>dtype=None</code>	Returns an array filled with ones. The <code>shape</code> tuple specifies array shape. <code>dtype</code> specifies the array type. If <code>dtype=None</code> , the type is float64.
<code>sort()</code>	<code>a</code> <code>axis=-1</code>	Sorts array <code>a</code> along <code>axis</code> . The default <code>axis=-1</code> sorts along the last axis in <code>a</code> . <code>axis=None</code> flattens <code>a</code> before sorting.