

## Tools for Data Science

### Glossary

Term	Definition
<b>Apache MLlib</b>	Language that makes machine learning scalable
<b>Apache Spark</b>	A general-purpose cluster-computing framework allowing you to process data using compute clusters
<b>API</b>	Application programming interface allows communication between two pieces of software
<b>Caffe</b>	A deep learning algorithm repository built with C++ with Python and Matlab bindings
<b>CDLA</b>	Community Data License Agreement
<b>Classification models</b>	Are used to predict whether some information or data belongs to a category (or “class”)
<b>CLI</b>	Command line interface
<b>C++</b>	A general-purpose programming language. It is an extension of the C programming language or C with Classes
<b>Data set</b>	A structured collection of data
<b>Deeplearning4</b>	Language for deep learning
<b>Deep learning</b>	A specialized type of machine learning. It refers to a general set of models and techniques that loosely emulate the way the human brain solves a wide range of problems
<b>ELT</b>	Extract, Load, Transform
<b>ETL</b>	Extract, Transform, and Load
<b>FSF</b>	Free Software Foundation
<b>ggplot2</b>	A popular library for data visualization in R
<b>GPU</b>	Graphics processing units
<b>Git</b>	De facto standard for code asset management, also known as version management or version control. Around Git emerged several services, GitHub, and GitLab
<b>Hadoop</b>	Application of Java which manages data processing and storage for big data applications running in clustered systems
<b>Java</b>	Object-oriented programming language
<b>Java-ML</b>	Language for machine learning
<b>JVM</b>	Java Virtual Machine
<b>JavaScript</b>	A general-purpose language that extended beyond the

	browser with the creation of Node.js and other server-side approaches
<b>Julia</b>	A language for high-performance numerical analysis and computational science
<b>Jupyter Notebook</b>	A browser-based application that allows you to create and share documents containing code, equations, visualizations, narrative text links, and more
<b>Jupyter Lab</b>	A browser-based application that allows you to access multiple Jupyter Notebook files, other code, and data files
<b>Kernel</b>	An execution environment for the different programming languages
<b>Lattice</b>	It is a high-level data visualization library that can handle graphics without customizations
<b>Library</b>	A collection of functions and methods that allow you to perform many actions without writing the code
<b>Leaflet</b>	Used for creating interactive plots
<b>ML</b>	Machine learning uses algorithms – also known as “models” - to identify patterns in the data
<b>Matplotlib</b>	package for data visualization
<b>Model training</b>	The process by which the model learns patterns from data
<b>MNIST</b>	Modified National Institute of Standards and Technology
<b>MongoDB</b>	A NoSQL database for big data management that was built with C++
<b>NLP</b>	Natural Language Processing
<b>NLTK</b>	Natural Language Toolkit
<b>NumPy</b>	Libraries are based on arrays and matrices, allowing you to apply mathematical functions to the arrays
<b>OSI</b>	Open-Source Initiative
<b>PaaS</b>	Platform as a service
<b>Pandas</b>	A library that offers data structures and tools for effective data cleaning, manipulation, and analysis
<b>Plotly</b>	Used for web-based data visualizations that can be displayed or saved as individual HTML files
<b>PMML</b>	Predictive Model Markup Language
<b>Python</b>	A high-level, general-purpose programming language. It has a large, standard library that provides tools suited to many different tasks, including Databases,

	Automation, Web scraping, Text processing, Image processing, Machine learning, and Data analytics
<b>R</b>	A statistical computing language
<b>Regression models</b>	Are used to predict a numeric (or “real”) value
<b>Reinforcement Learning</b>	Loosely based on the way human beings and other organisms learn.
<b>REST</b>	RE stands for Representational; the S stands for State, and the T stands for Transfer
<b>RStudio</b>	Unifies programming, execution, debugging, remote data access, data exploration, and visualization into one tool
<b>SaaS</b>	Software as a service
<b>Scala</b>	Is a combination of scalable and language. A general-purpose programming language that provides support for functional programming and is a strong static type system
<b>Spyder</b>	Integrates code, documentation, and visualizations, among others, into a single canvas
<b>SQL</b>	Structured Query Language that is non-procedural, used for querying and managing data
<b>Supervised Learning</b>	A learning in which a human provides input data and correct outputs
<b>TensorFlow</b>	Deep Learning library for dataflow that was built with C++
<b>Unsupervised Learning</b>	The data is not labeled by a human. Examples are Clustering models used to divide each record of a dataset into one of a similar group
<b>Watson Studio</b>	A fully integrated development environment for data scientists
<b>Weka</b>	Language for data mining