

Machine Learning – Deepdive

Raghu Prasad K S



Introduction

- Raghu Prasad BE, MS
- Total of 29 years of experience
- 7 years as a lecturer in Engineering College
- 22 Years into IT
- Worked with companies like CISCO, CSC, ICICI, First Apex NTT Data
- Currently into Corporate training and consultancy
- Worked with corporates and public sector
- Technologies Java, Python, Data Sciences, Web technologies, Java Script technologies (MEAN stack), IOT, Test Automation – Selenium, JMeter



Topics

- Introduction to Data science
- Python for Data Sciences and project lifecycle
- Machine learning examples
- How machine learning works
- Types of Machine Learning
- Numpy
- Pandas
- Matlablib
- Machine Learning Algorithms
- Introduction to Scikit-learn

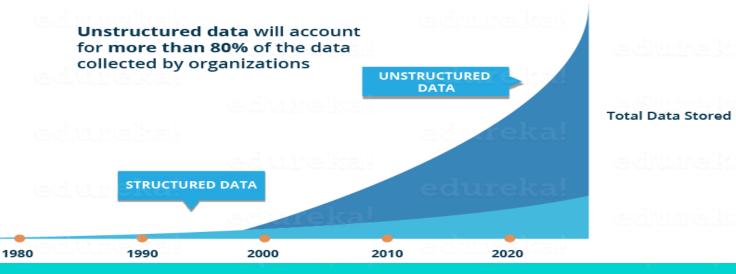


 Data science is a multi-disciplinary field that uses scientific methods, processes, algorithms and systems to extract knowledge and insights from structured and unstructured data. [1][2] Data science is the same concept as data mining and big data: "use the most powerful hardware, the most powerful programming systems, and the most efficient algorithms to solve problems"



- Need for data science
- https://jeremyronk.wordpress.com/2014/09/01/s tructured-semi-structured-and-unstructured-

data/

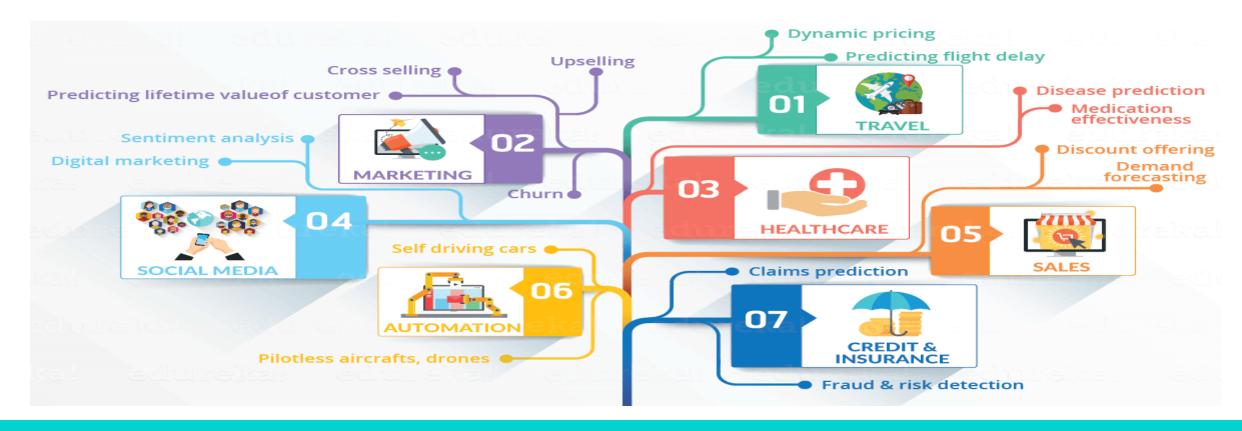




- Need for data science
- This data is generated from different sources like financial logs, text files, multimedia forms, sensors, and instruments.
- Source https://www.edureka.co/blog/what-is-data-science/
- Simple BI tools are not capable of processing this huge volume and variety of data.
- This is why we need more complex and advanced analytical tools and algorithms for processing, analyzing and drawing meaningful insights out of it.
- https://www.geeksforgeeks.org/difference-between-structured-semi-structured-and-unstructured-data/
- https://www.w3trainingschool.com/structured-semi-structured-unstructured-data



Domains currently using data science

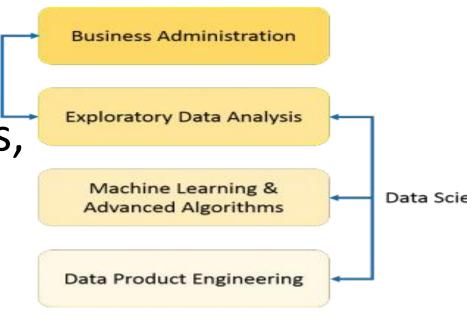




Analyst

What is data science

Data Science is a blend of various tools, algorithms, and machine learning principles with the goal to discover hidden patterns from the raw data, make sense of the data, predicting the future and take business decisions.





Predictive causal analytics – If you want a model which can predict the possibilities of a particular event in the future, you need to apply predictive causal analytics.

Example: If you are providing money on credit, then the probability of customers making future credit payments on time is a matter of concern for you. Here, you can build a model which can perform predictive analytics on the payment history of the customer to predict if the future payments will be on time or not.

Prescriptive analytics: If you want a model which has the intelligence of taking its own decisions and the ability to modify it with dynamic parameters, you certainly need prescriptive analytics for it. This relatively new field is all about providing advice. In other terms, it not only predicts but suggests a range of prescribed actions and associated outcomes.

Example for this is Google's self-driving car which I had discussed earlier too. The data gathered by vehicles can be used to train self-driving cars. You can run algorithms on this data to bring intelligence to it. This will enable your car to take decisions like when to turn, which path to take, when to slow down or speed up.



Machine learning for making predictions — If you have transactional data of a finance company and need to build a model to determine the future trend, then machine learning algorithms are the best bet.

This falls under the paradigm of supervised learning. It is called supervised because you already have the data based on which you can train your machines.

For example, a fraud detection model can be trained using a historical record of fraudulent purchases.

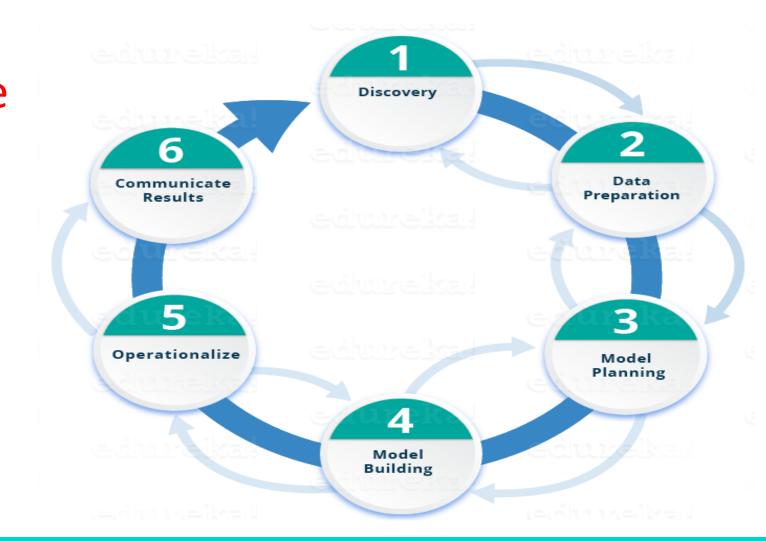


Machine learning for pattern discovery — If you don't have the parameters based on which you can make predictions, then you need to find out the hidden patterns within the dataset to be able to make meaningful predictions. This is nothing but the unsupervised model as you don't have any predefined labels for grouping.

The most common algorithm used for pattern discovery is Clustering. Let's say you are working in a telephone company and you need to establish a network by putting towers in a region. Then, you can use the clustering technique to find those tower locations which will ensure that all the users receive optimum signal strength.



Life cycle of data science





Python Data Science

- Why Learn Python For Data Science?
- Python is no-doubt the best-suited language for a Data Scientist. I
 have listed down a few points which will help you understand why
 people go with Python for Data Science:
 - Python is a free, flexible and powerful open source language
 - Python cuts development time in half with its simple and easy to read syntax
 - With Python, you can perform data manipulation, analysis, and visualization
 - Python provides powerful libraries for Machine learning applications and other scientific computations

-

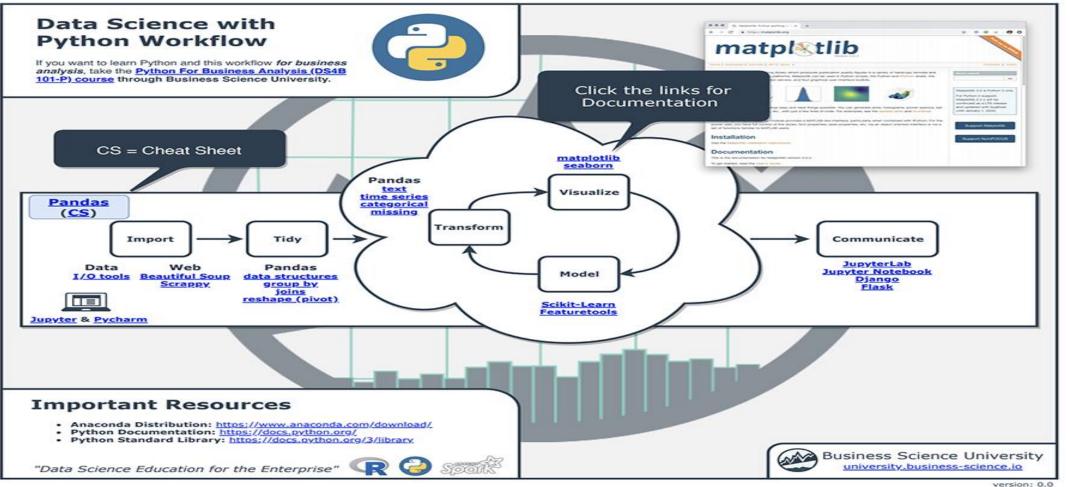


Python Data Science

- Python Libraries For Data Science
- This is the part where the actual power of Python with data science comes into the picture.
 Python comes with numerous libraries for scientific computing, analysis, visualization etc.
 Some of them are listed below:
- Numpy NumPy is a core library of Python for Data Science which stands for 'Numerical Python'. It is used for scientific computing, which contains a powerful n-dimensional array object and provide tools for integrating C, C++ etc. It can also be used as multi-dimensional container for generic data where you can perform various Numpy Operations and special functions.
- <u>Pandas</u> Pandas is an important library in Python for data science. It is used for data manipulation and analysis. It is well suited for different data such as tabular, ordered and unordered <u>time series</u>, matrix data etc.
- <u>Matplotlib</u> Matplotlib is a powerful library for visualization in Python. It can be used in Python scripts, shell, web application servers and other GUI toolkits. You can use different <u>types of plots</u> and how <u>multiple</u> <u>plots</u> work using Matplotlib
- <u>Scikit-learn</u> Scikit learn is one of the main attractions, where in you can implement machine learning using Python. It is a free library which contains simple and efficient tools for data analysis and mining purposes. You can implement various algorithm, such as <u>logistic regression</u>, <u>time series algorithm</u> using scikit-learn.



Python Data Science



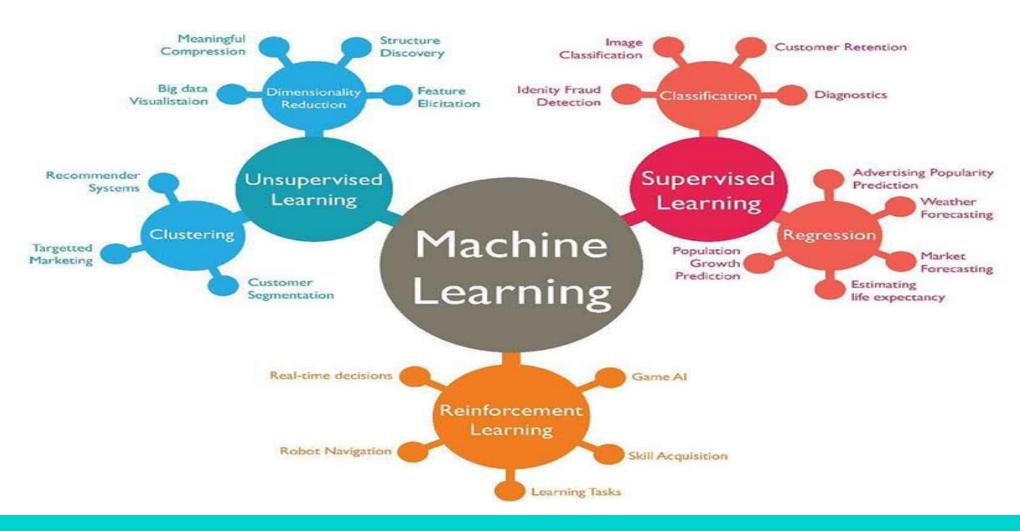


How does machine learning works?

How does Machine Learning Work? Training Data Train The ML Algorithm Successfull Model Model Input Data Predication **New Input** Data ML Algorithm



Types of machine learning





Classification of machine learning

Supervised Learning

Supervised learning is a type of machine learning algorithm that uses a known dataset (called the training dataset) to make predictions. The training dataset includes input data and response values. From it, the supervised learning algorithm seeks to build a model that can make predictions of the response values for a new dataset.

Supervised Learning

Let's take an example here. Say you are a teacher, and your way of teaching is,

To teach by example, i.e for every problem in their life you are providing solutions to them, this type of learning is called **supervised learning**.

Let's take the same example forward:



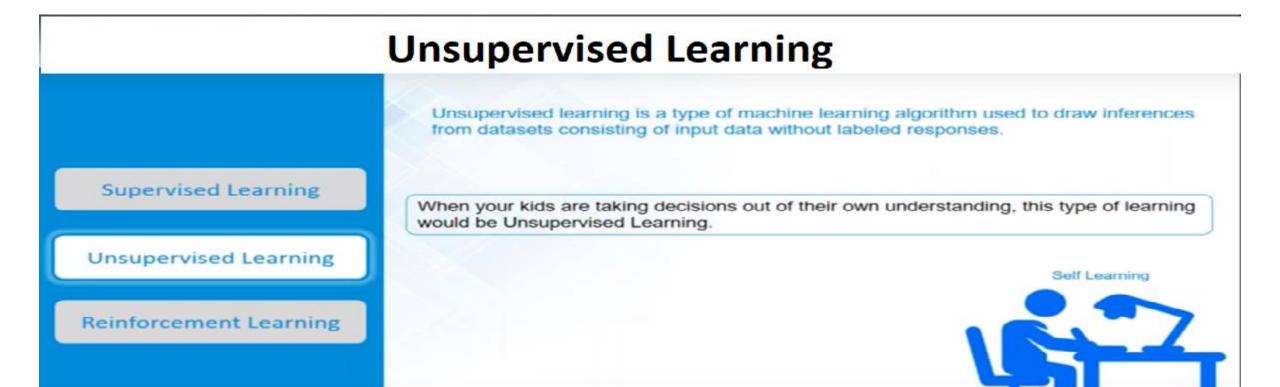
gisul.co.in

Unsupervised Learning

Reinforcement Learning



Classification of machine learning





Classification of machine learning

Reinforcement Learning

Supervised Learning

Unsupervised Learning

Reinforcement Learning

Reinforcement learning is an area of machine learning inspired by behaviorist psychology, concerned with how software agents ought to take actions in an environment so as to maximize some notion of cumulative reward.

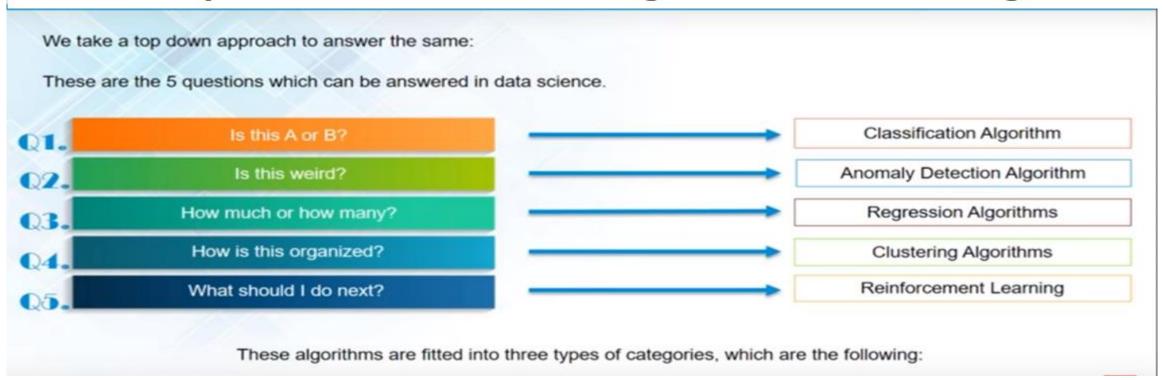
If a new situation comes up, the kid will take actions on his own i.e from his past experiences, but as a parent towards the end of an action you can tell him whether he did good or not.





How to solve a problem

How a problem is solved using Machine Learning





How to solve a problem

Classification Algorithms

Classification Algorithms are used to classify a record.

It is used for questions which can have only a limited number of answers.

For Example:

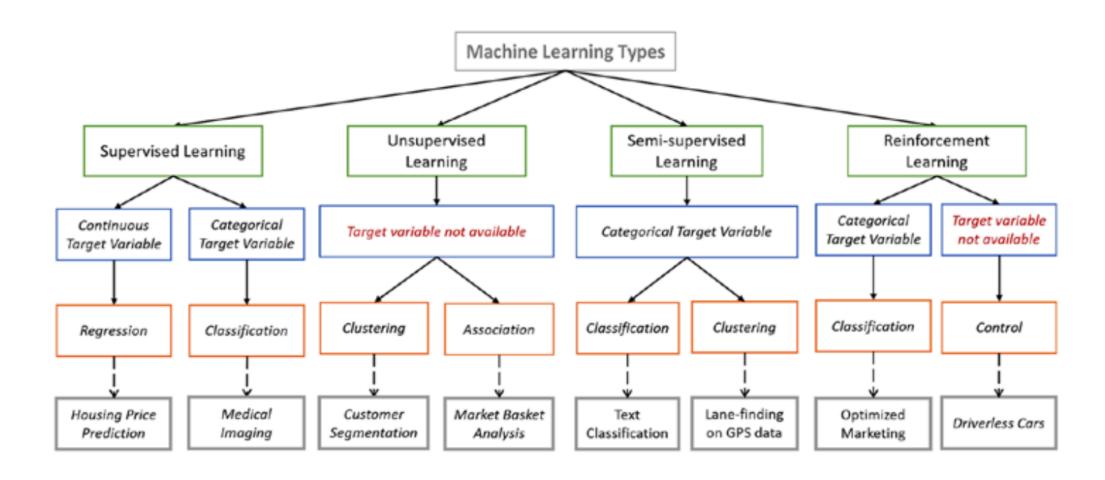
s or No
o or Maybe



When you have only two choices, its called 2 class Classification, if you have more than 2 choices its called Multi Class Classification

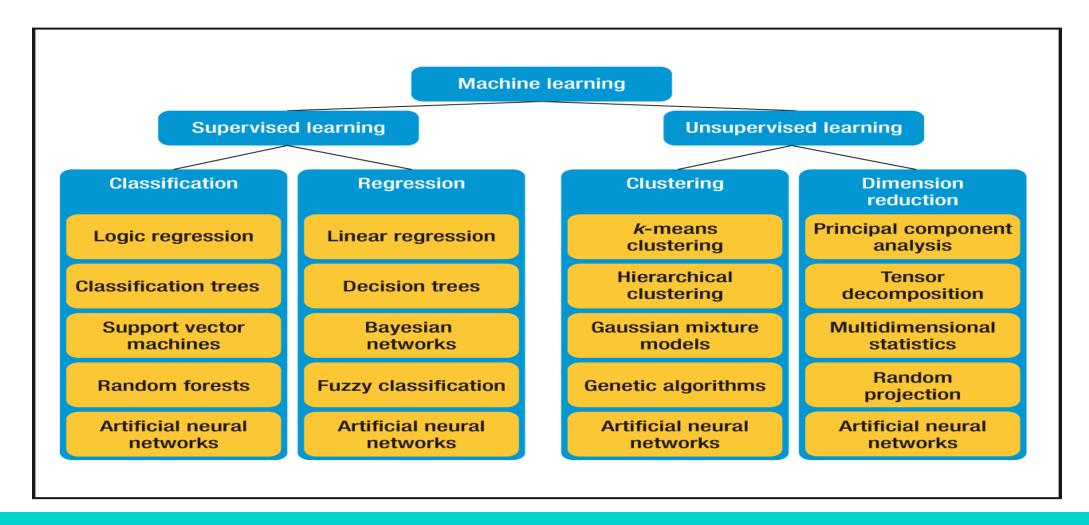


Machine Learning Algorithms



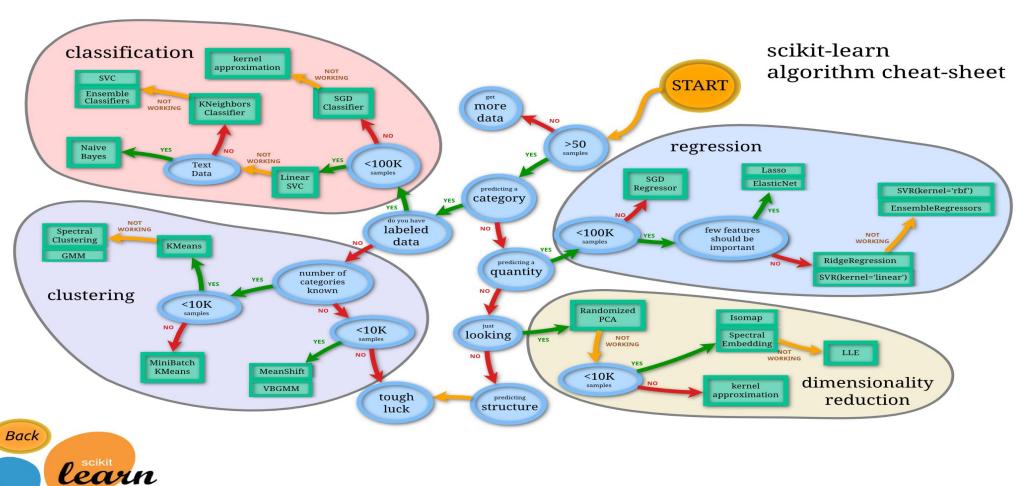


Machine Learning Algorithms





Machine Learning Algorithms





PANDAS

- The <u>Pandas</u> library is one of the most preferred tools for data scientists to do data manipulation and analysis, next to <u>matplotlib</u> for data visualization and <u>NumPy</u>, the fundamental library for scientific computing in Python on which Pandas was built.
- The fast, flexible, and expressive Pandas data structures are designed to make real-world data analysis significantly easier, but this might not be immediately the case for those who are just getting started with it. Exactly because there is so much functionality built into this package that the options are overwhelming.



PANDAS-Cheat Sheet

Data Wrangling

with pandas **Cheat Sheet** http://pandas.pydata.org

Syntax - Creating DataFrames

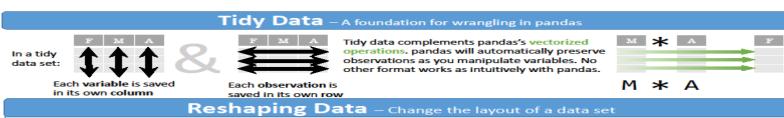
11 df = pd.DataFrame({"a" : [4 ,5, 6], "b" : [7, 8, 9], "c" : [10, 11, 12]}, index = [1, 2, 3])Specify values for each column. df = pd.DataFrame([[4, 7, 10], [5, 8, 11], [6, 9, 12]], index=[1, 2, 3], columns=['a', 'b', 'c'])

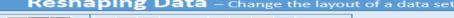
Specify values for each row.

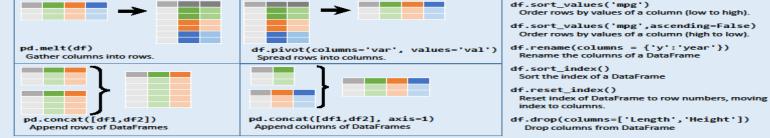
				ь	c	
	n	w				
	_	1	4	7	10	
	•	2	5	8	11	
	e	2	6	9	12	
df = pd.DataFrame(
{"a" : [4 ,5, 6],						
"b" : [7, 8, 9],						
"c" : [10, 11, 12]},						
<pre>index = pd.MultiIndex.from_tuples(</pre>						
[('d',1),('d',2),('e',2)],						
names=['n','v']))						
Create DataFrame with a MultiIndex						

Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code. df = (pd.melt(df)).rename(columns={ 'variable' : 'var', 'value' : 'val'}) .query('val >= 200')







Subset Observations (Rows)

- df[df.Length > 7]Extract rows that meet logical
- df.drop_duplicates() Remove duplicate rows (only considers columns). df.head(n)
- Select first n rows. df.tail(n) Select last n rows.

- df.sample(frac=0.5) Randomly select fraction of rows.
- df.sample(n=10) Randomly select n rows. df.iloc[10:20]
- Select rows by position. df.nlargest(n, 'value') Select and order top n entries.
- df.nsmallest(n, 'value') Select and order bottom n entries.

	Logic in Python (and pandas)						
<	Less than	!=	Not equal to				
>	Greater than	df.column.isin(values)	Group membership				
==	Equals	pd.isnull(obj)	Is NaN				
<=	Less than or equals	pd.notnull(obj)	Is not NaN				
>=	Greater than or equals	8, ,~,^,df.any(),df.all()	Logical and, or, not, xor, any, all				
the decrease of the second state of the second seco							

Subset Variables (Columns)



- df[['width','length','species']] Select multiple columns with specific names.
- df['width'] or df.width Select single column with specific name. df.filter(regex='regex')
 - Select columns whose name matches regular expression regex. regey (Regular Expressions) Examples

8 (8				
٠,	Matches strings containing a period '.'			
'Length\$'	Matches strings ending with word 'Length'			
'^Sepal'	Matches strings beginning with the word 'Sepal'			
'^x[1-5]\$'	Matches strings beginning with 'x' and ending with 1,2,3,4,5			
'^(?!Species\$).+'	Matches strings except the string 'Species'			

- df.loc[:,'x2':'x4']
- Select all columns between x2 and x4 (inclusive). df.iloc[:,[1,2,5]]
- Select columns in positions 1, 2 and 5 (first column is 0). df.loc[df['a'] > 10, ['a','c']]
- Select rows meeting logical condition, and only the specific columns .



Data wrangling

Data wrangling, sometimes referred to as data munging, is the process of transforming and mapping data from one "raw" data form into another format with the intent of making it more appropriate and valuable for a variety of downstream purposes such as analytics.



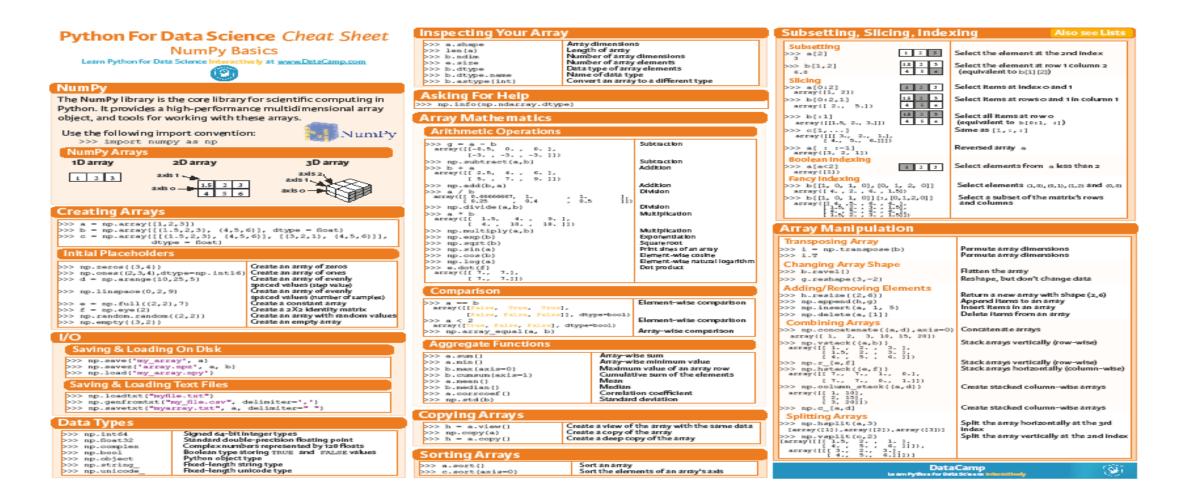


NumPy

- NumPy is the fundamental package for scientific computing with Python.
- It contains among other things:
 - a powerful N-dimensional array object
 - sophisticated (broadcasting) functions
 - tools for integrating C/C++ and Fortran code
 - useful linear algebra, Fourier transform, and random number capabilities
- Besides its obvious scientific uses, NumPy can also be used as an efficient multi-dimensional container of generic data. Arbitrary datatypes can be defined. This allows NumPy to seamlessly and speedily integrate with a wide variety of databases.



NumPy-Cheat Sheet





MatplobLib

- Data visualization and storytelling with your data are essential skills that every data scientist needs to communicate insights gained from analyses effectively to any audience out there.
- Matplotlib is a Python 2D plotting library which produces publication quality figures in a variety of hardcopy formats and interactive environments across platforms. Matplotlib can be used in Python scripts, the Python and <u>IPython</u> shells, the <u>Jupyter</u> notebook, web application servers, and four graphical user interface toolkits.
- Matplotlib tries to make easy things easy and hard things possible.
 You can generate plots, histograms, power spectra, bar charts, errorcharts, scatterplots, etc., with just a few lines of code.



MatplobLib-Cheat Sheet

