

Data Science using Python

Manaranjan Pradhan

A Data Enthusiast

LinkedIn: <https://www.linkedin.com/in/manaranjanpradhan>

He writes blogs at www.awesomestats.in

Complete the following steps...

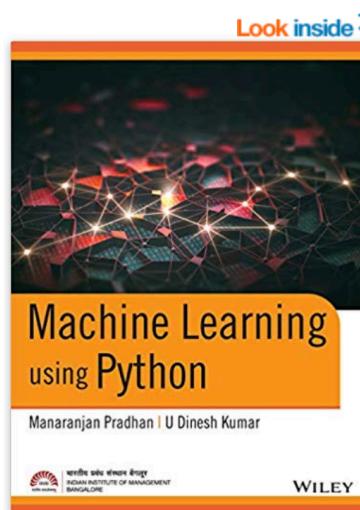
- Create a working directory on your desktop / laptop, where you can store all data and the programs of this workshop
- Download datasets from the link below and store in the above directory
 - https://github.com/manaranjanp/ML_Python
 - Download the complete repository as zip file and unarchive it
- Install latest Anaconda Distribution for Python (3.6+) on your desktop / laptop
 - <https://www.anaconda.com/download/>

About Me

- Worked for HP and iGate for 12+ years
 - Working as freelancer for last 5 years
 - Mostly conduct workshops on Spark, Databricks, Machine Learning and Deep Learning
- <https://www.linkedin.com/in/manaranjanpradhan>



Manaranjan Pradhan (Manu)



<https://www.amazon.in/Machine-Learning-Python-Manaranjan-Pradhan-ebook/dp/B07RLQPNRX/>

Machine Learning using Python Paperback – 2019

by [U Dinesh Kumar Manaranjan Pradhan](#) (Author)

7 customer reviews

[See all 2 formats and editions](#)

Kindle Edition
₹ 423.20

Paperback
₹ 529.00

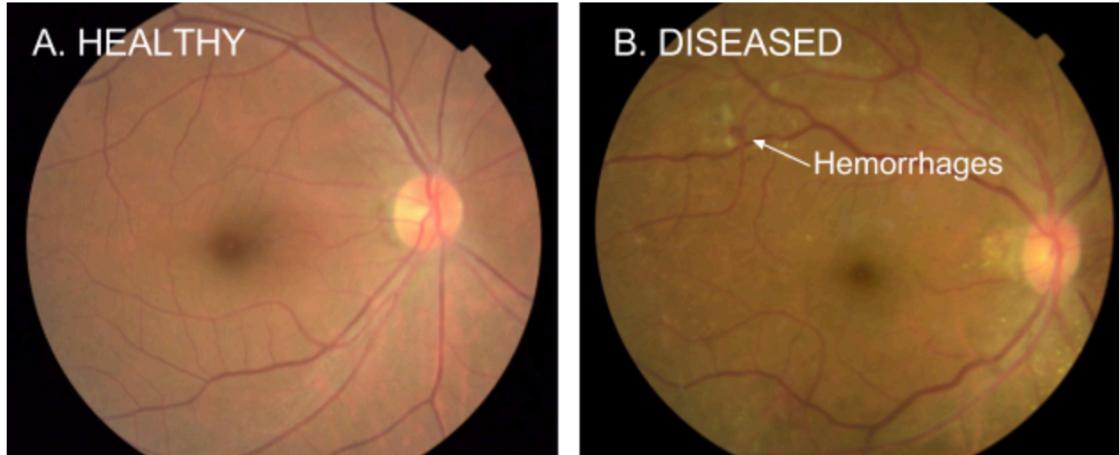
[Read with Our Free App](#) 2 New from ₹ 529.00

This book is written to provide a strong foundation in machine learning using Python libraries by providing real-life case studies and examples. It covers topics such as foundations of machine learning, introduction to Python, descriptive analytics and predictive analytics. Advanced machine learning concepts such as decision tree learning, random forest, boosting, recommended systems, and text analytics are covered. The book takes a balanced approach between theoretical understanding and practical applications. All the topics include real-world examples and provide step-by-step approach on how to explore, build, evaluate, and optimize machine learning models.

Course Outline

| | |
|---|---|
| Introduction to Data Science and Setting up data analysis environment | Introduction to Data Science Setting up Python Environment for Data Analysis Overview of Data Analysis Stack - Numpy, Pandas, Matplotlib, scipy and Scikit-learn |
| Accessing and preparing data with Pandas | Loading data from Different Sources Data manipulation - Filtering, Grouping, Ordering, Joining Dealing with missing Data |
| Data Exploration, Visualizations & Statistical Analysis | Histograms, Bar charts Density Plots, Box Plots, Scatter Plots, Heat Maps Understanding Basic Statistics, Distributions, Correlations |
| Algorithms for Regression and Classification Problems | Understanding loss function and gradient descent approach for loss minimization Linear Regression, Logistic Regression, Decision Trees, KNN Model Optimization & Parameter Tuning |
| Clustering | K-means clustering Finding optimal number of clusters |
| Model Evaluation | Creating Training, validation and Test Data Sets Cross validations Understanding Evaluation Metrics: RMSE, R-square, ROC, Confusion Matrix, Precision, Recall, Accuracy etc. |

Detection of Diabetic Eye Disease



haemorrhage

/'hemərɪdʒ/ 🔊

noun

plural noun: **hemorrhages**

1. an escape of blood from a ruptured blood vessel.
"a massive haemorrhage of the brain"

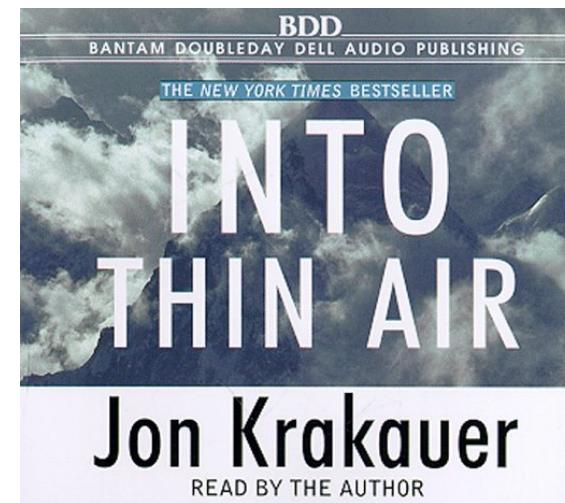
<https://research.googleblog.com/2016/11/deep-learning-for-detection-of-diabetic.html>

- Google working closely with doctors both in India and the US, created a development dataset of 128,000 images which were each evaluated by 3-7 ophthalmologists from a panel of 54 ophthalmologists.
- Trained a deep neural network to detect referable diabetic retinopathy.
- Then tested the algorithm's performance on two separate clinical validation sets totaling ~12,000 images, with the majority decision of a panel 7 or 8 U.S. board-certified ophthalmologists serving as the reference standard.
- the algorithm has a F-score (combined sensitivity and specificity metric, with max=1) of 0.95, which is slightly better than the median F-score of the 8 ophthalmologists we consulted (measured at 0.91).

Into Thin Air

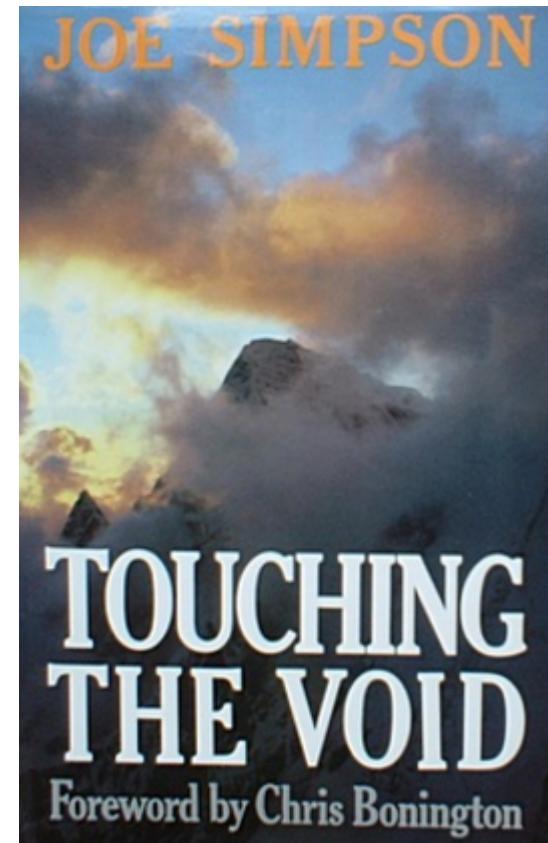


It details the author's presence at [Mount Everest](#) during the [1996 Mount Everest disaster](#), when eight climbers were killed and several others were stranded by a "rogue storm"



Touching the Void

Touching the Void is a 1988 book by [Joe Simpson](#), recounting his and [Simon Yates'](#) successful but disastrous and nearly fatal climb of the 6,344-metre (20,813 foot) [Siula Grande](#) in the [Peruvian Andes](#) in 1985.



Customers who bought this item also bought...

[Look inside](#) ↴

#1 NATIONAL BESTSELLER
A Personal Account of the Mt. Everest Disaster

INTO THIN AIR

"Ranks among the great adventure books of all time." — *THE WALL STREET JOURNAL*

Jon Krakauer
AUTHOR OF *INTO THE WILD* AND *TIGER DREAMS*

Paperback – October 19, 1999
by Jon Krakauer (Author, Photographer), Randy Rackliff (Illustrator), Daniel Rembert (Contributor), & 2 more
★★★★★ 2,414 customer reviews
#1 Best Seller in Mountain Climbing

See all 65 formats and editions

| | | | |
|------------------|----------------------|----------------------|--------------------------------------|
| Kindle \$3.40 | Hardcover \$18.66 | Paperback \$10.36 | Mass Market Paperback from \$0.01 |
|------------------|----------------------|----------------------|--------------------------------------|

Read with our [free app](#)

667 Used from \$0.01
82 New from \$4.49
51 Collectible from \$6.37

483 Used from \$0.01
117 New from \$6.22
11 Collectible from \$9.70

487 Used from \$0.01
18 New from \$4.59
19 Collectible from \$3.00

National Bestseller

A bank of clouds was assembling on the not-so-distant horizon, but journalist-mountaineer Jon

Customers Who Bought This Item Also Bought

| | | | | |
|---|---|--|---|---|
| LOOK INSIDE! | LOOK INSIDE! | LOOK INSIDE! | LOOK INSIDE! | LOOK INSIDE! |
| Into the Wild » Jon Krakauer ★★★★★ 2,304 #1 Best Seller in Travelogues & Travel Essays Paperback \$7.34 | Under the Banner of... » Jon Krakauer ★★★★★ 1,361 Paperback \$10.03 Get it by Tomorrow | Missoula: Rape and the... Jon Krakauer ★★★★★ 361 Hardcover \$18.09 Get it by Tomorrow | Touching the Void: The... » Joe Simpson ★★★★★ 315 Paperback \$11.22 Get it by Tomorrow | Buried in the Sky: The... Peter Zuckerman and Amanda Paddan ★★★★★ 225 Paperback \$10.63 Get it by Tomorrow |

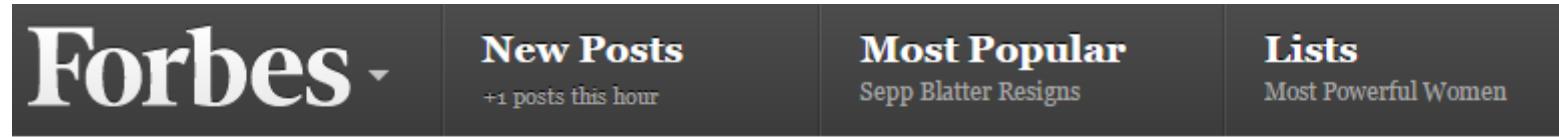
Recommendations are key to personalization

Amazon's recommendation secret

McKinsey estimated that 35 percent of consumer purchases on Amazon come from product recommendations, although the e-commerce giant itself has never revealed its own estimates. In 2016, it offered its open-source artificial intelligence (AI) framework called, DSSTNE (pronounced as "destiny"), for free to encourage the development of artificial intelligence apps.

<https://martechtoday.com/roi-recommendation-engines-marketing-205787>

What customers bought together



4/06/1998 @ 12:00AM

f Share

Diaper-beer syndrome

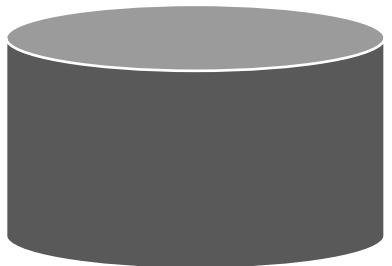
IT'S PART OF the folklore of data processing. A retail chain put all its checkout-counter data into a giant digital warehouse and set the disk drives spinning.

Out popped a most unexpected correlation: sales of diapers and beer.

Evidently, young fathers would make a late-night run to the store to pick up Pampers and get some Bud Light while they were there.

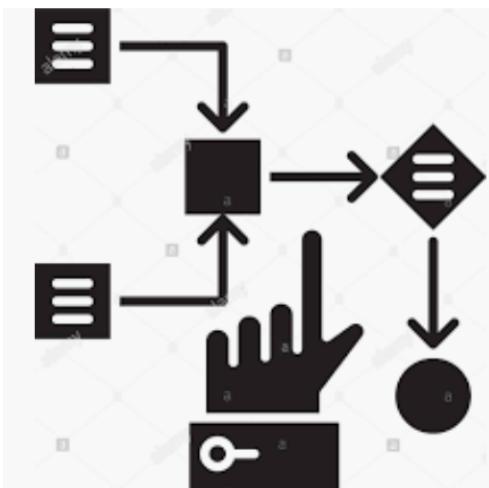
Capitalizing on the discovery, the store placed the disparate items together. Sales zoomed.

Key Components



Data

- Historical Evidences
- Samples representing problem scenario



Learning Algorithm

- Statistical Learning
- Machine Learning or Deep Learning

$$\pi(R-\ell)^2 \leq NS \leq \pi(R+\ell)^2$$
$$(2q + p + 2)\frac{s}{2} = \left(q + \frac{p}{2} + 1\right)s = ns$$

Model

A mathematical expression of the pattern or evidence found in data and can be used to find insights and applied in future to predict.

How AI, ML and DL are related?

ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



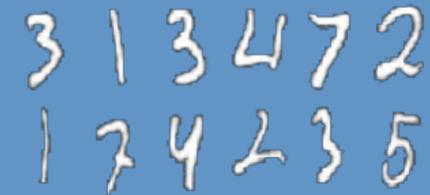
MACHINE LEARNING

Ability to learn without explicitly being programmed

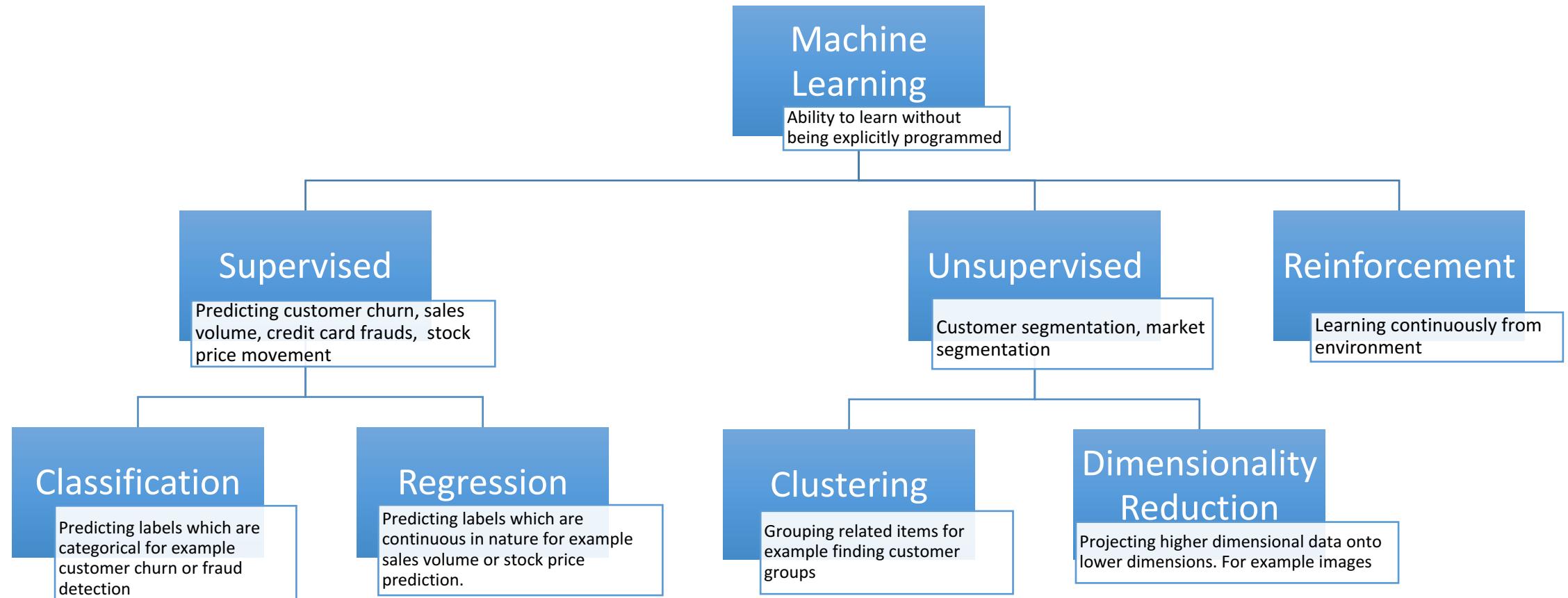


DEEP LEARNING

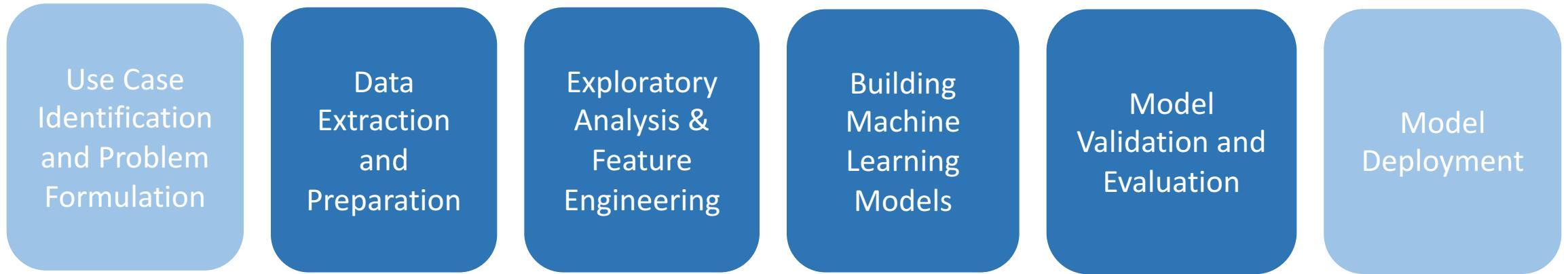
Learn underlying features in data using neural networks



Machine Learning Algorithms



Lifecycle

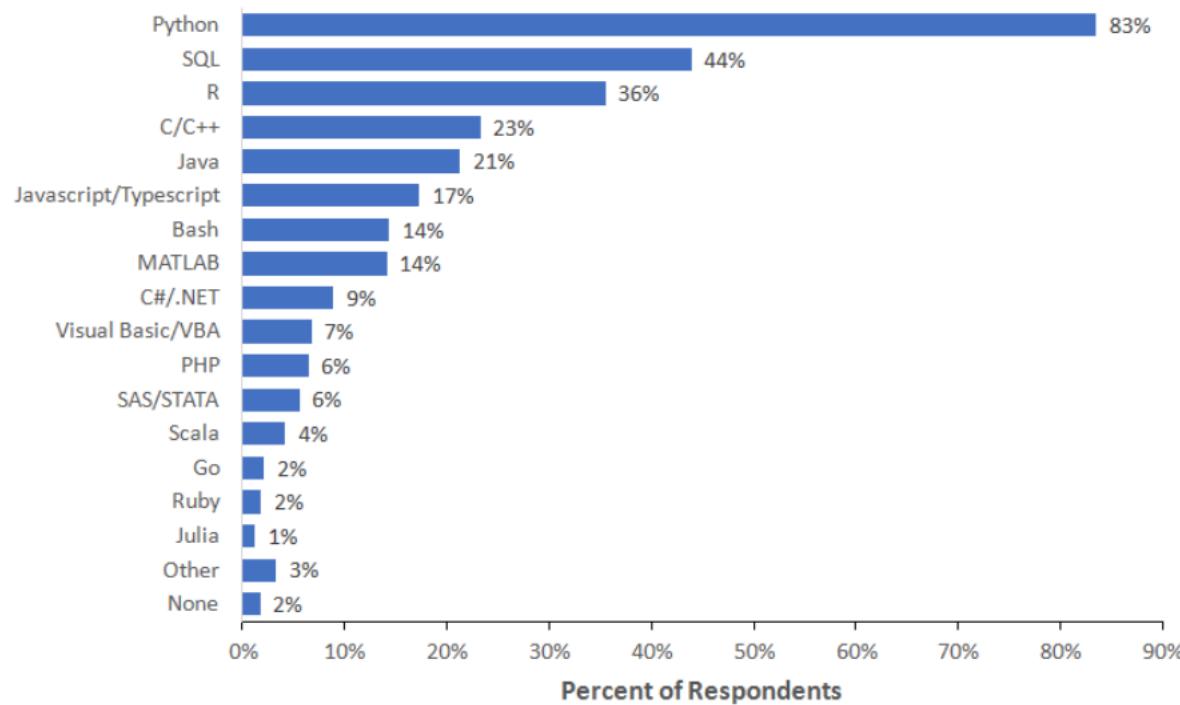


Iterative Steps

What Tools are available?

Language for Machine Learning

What programming language do you use on a regular basis?



Note: Data are from the 2018 Kaggle Machine Learning and Data Science Survey. You can learn more about the study here: <http://www.kaggle.com/kaggle/kaggle-survey-2018>. A total of 18827 respondents answered the question.

<https://businessoverbroadway.com/2019/01/13/programming-languages-most-used-and-recommended-by-data-scientists/>

Python Stack For Data Science

Efficient storage of arrays and matrices. Backbone of all scientific calculations and algorithms.



Library for scientific computing.
Linear algebra, statistical computations, optimization algorithm.



seaborn

Plotting and visualization



IP[y]: IPython
Interactive Computing

pandas



High-performance, easy-to-use data structures for data manipulation and analysis. Pandas provide the features of dataframe, which is very popular in the area of analytics for data munging, cleaning & transformation.

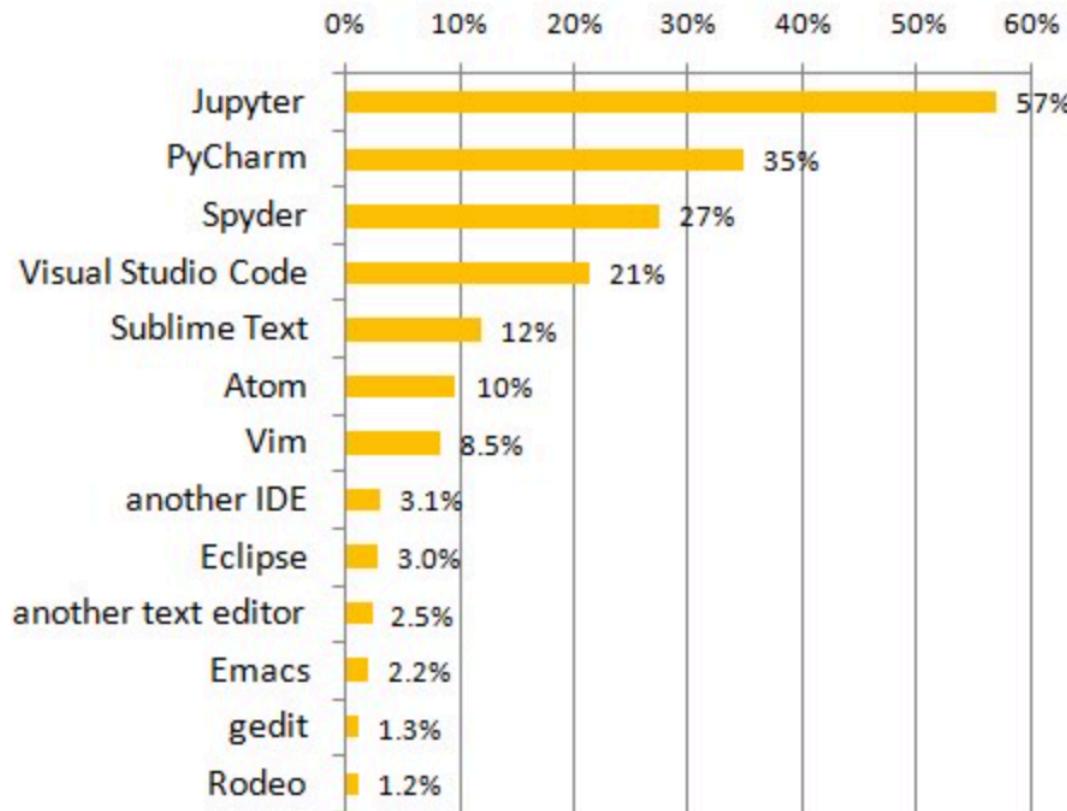
IDE or Development environment for data analysis in python.



Machine learning library. Collection of ML algorithms.

Language for Machine Learning

Most Popular Python IDE, Editors

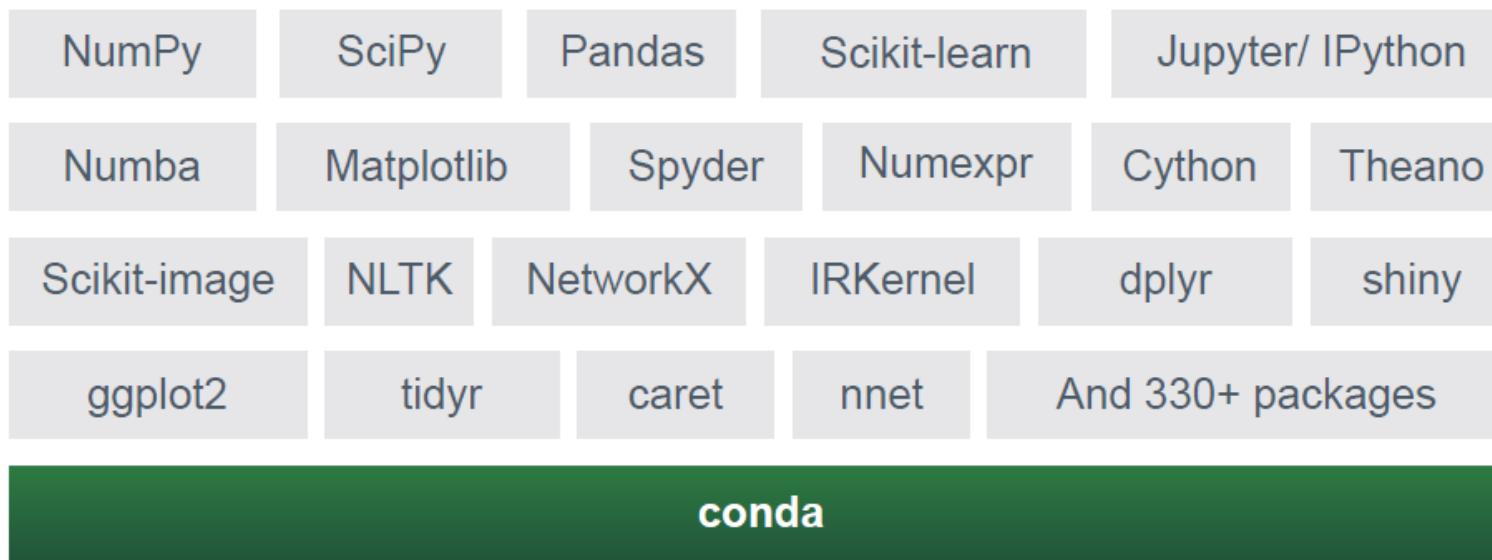


<https://www.anaconda.com/download/>

Python Distribution



ANACONDA®



**Game-Changing
Enterprise Ready
Python Distribution**

- 2 million downloads in last 2 years
- 200k / month and growing
- conda package manager serves up 5 *million* packages per month
- Recommended installer for IPython/Jupyter, Pandas, SciPy, Scikit-learn, etc.

Source: Continuum Analytics

Download link: <https://www.continuum.io/downloads>

Start Jupyter notebook

- For MAC
 - Click on Anaconda Navigator and click on “launch notebook”
 - Or go to command prompt and enter
 - **jupyter notebook**
- For Windows
 - Go to anaconda command prompt and enter
 - **jupyter notebook**

Start a jupyter notebook



Files Running Clusters Conda

Select items to perform actions on them.

Upload New ▾ ⌂

The screenshot shows the Jupyter Notebook interface. On the left is a file browser with a sidebar containing icons for home, back, forward, and search. Below the sidebar is a list of local paths: anaconda, Applications, Desktop, Documents, Downloads, and metastore_db. To the right of the file browser is a context menu with the following options:

- Text File
- Folder
- Terminal
- Notebooks
- Python [conda root]
- Python [default] (this option is highlighted)
- Spark 2.1.0

Click on new to start new notebook. For every hands on exercise, start a new notebook.