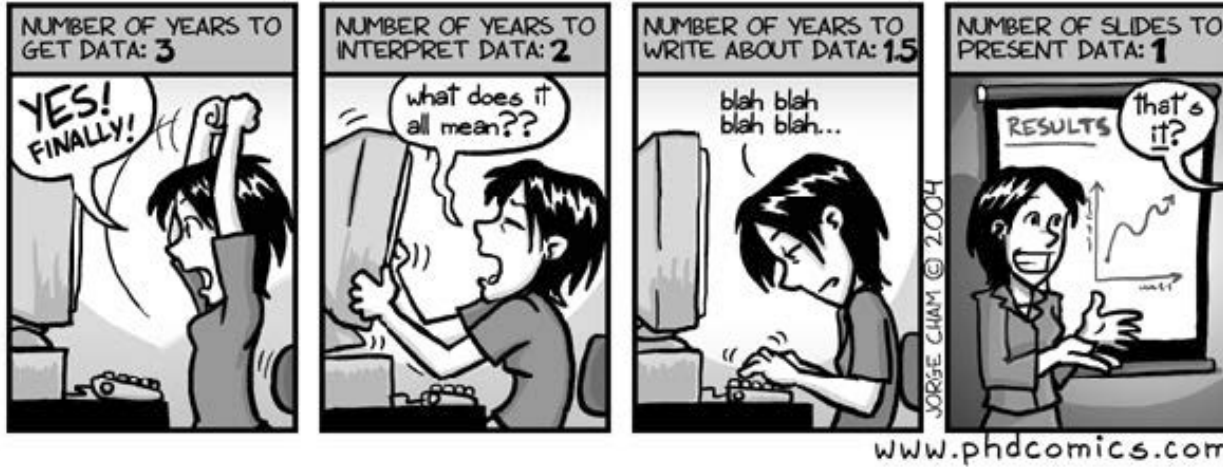


## DATA: BY THE NUMBERS



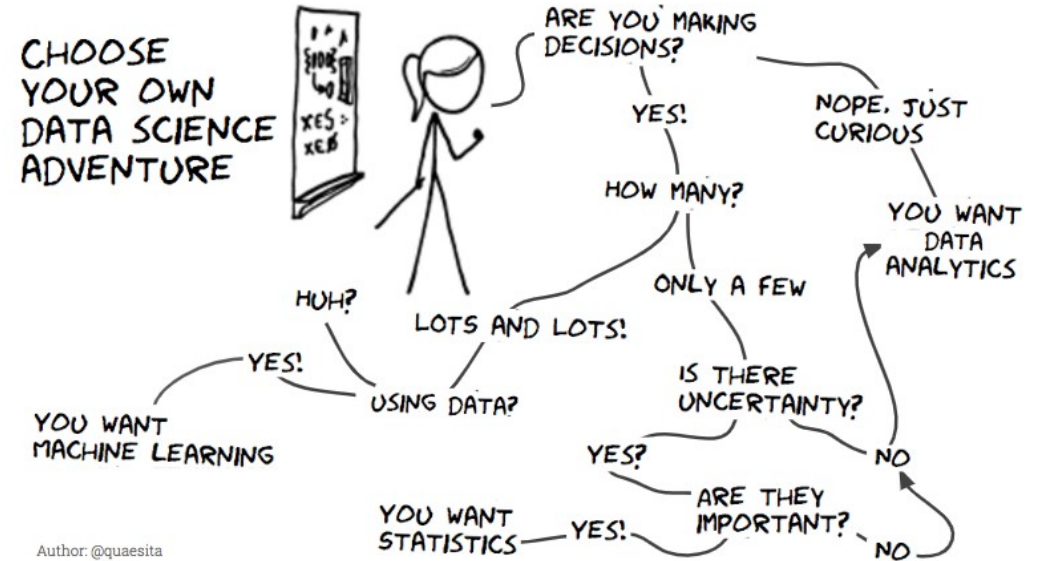
**Data Science:** intro, skills, works, explore, produce, interpret,...

...make sense for a particular problem, generating added value...

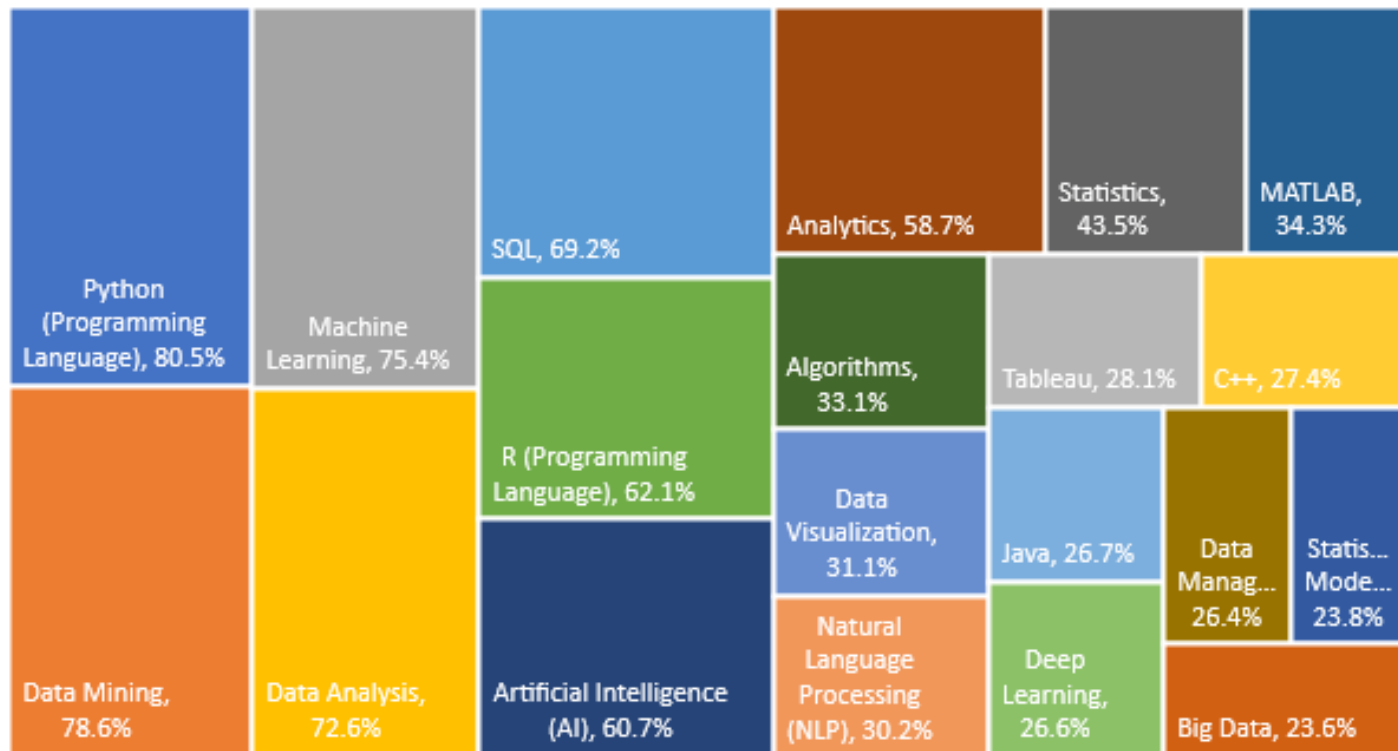
## Two pictures to understand our challenge!



© marketoonist.com



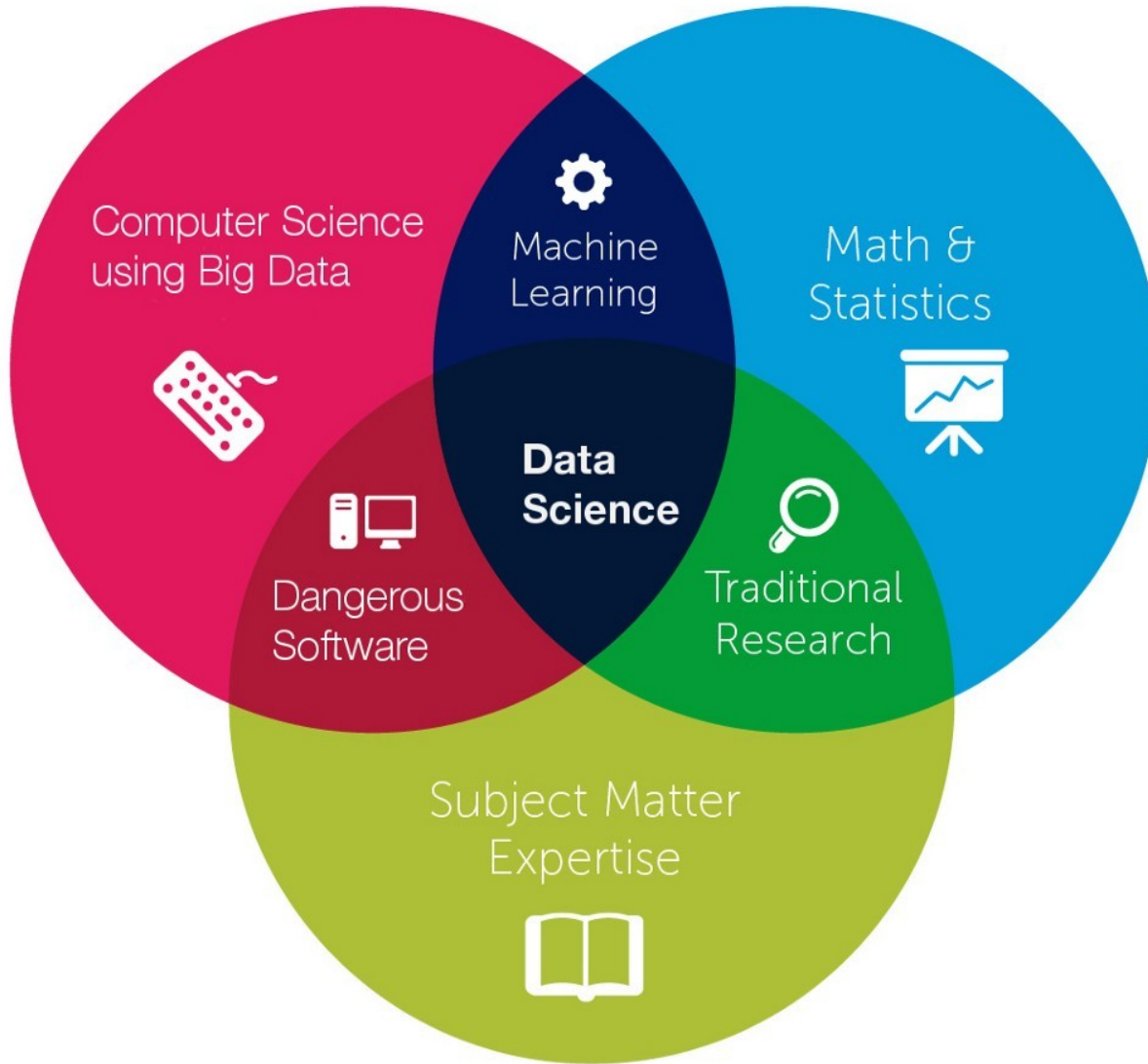
## Top 20 skills reported by Data Scientists



- Figure summarizes the **skills of 200 data scientists** from leading technology destinations such as the San Francisco Bay Area, New York City Metropolitan Area, Seattle, Dallas-Fort Worth Metroplex, Raleigh-Durham-Chapel Hill Area, Greater Chicago Area, Greater Boston, London Area, Bengaluru, and New Delhi.
- It is abundantly clear from Figure that the top skill that most data scientists have is Python, followed by data mining.

<https://www.red-gate.com/simple-talk/development/data-science-development/how-to-become-data-scientist-data-driven-approach-careers-data/>

## Data Science Venn Diagram



The **big** Picture



## where does big data come from?

The bulk of big data generated comes from three primary sources: social data, machine data and transactional data.

### Big Data comes from everywhere

Sensors used to gather climate information



Posts to social media sites



Software logs, Camera, Microphone



Digital pictures and videos



Emails, blogs and e-news



Traffic data and GPS Signals



The bulk of big data generated comes from three primary sources: **social data**, **machine data** and **transactional data**. In addition, companies need to make the distinction between data which is generated internally, that is to say it resides behind a company's firewall, and externally data generated which needs to be imported into a system.

Whether data is unstructured or structured is also an important factor. Unstructured data does not have a pre-defined data model and therefore requires more resources to make sense of it.





## KEY

SOME APIs  
NO APIs



INTERNAL



EXTERNAL



BOTH

## TERMINOLOGY

### SOME APIs

Data that has a standard Web service

### NO APIs

Data that has no standard Web service and requires alternative methods of integration

### INTERNAL

Data that resides behind an organization's firewall

### EXTERNAL

Data that resides outside of an organization's firewall

### UNSTRUCTURED

Data that does not have a pre-defined data model or is not organized in a pre-defined manner

### STRUCTURED

Data that resides in a fixed field within a record or file

### VELOCITY

The rate at which data is generated and changed

### VARIETY

The number of different data sources and types

### VOLUME

The average quantity of data units per category



### ARCHIVES

Archives of scanned documents, statements, insurance forms, medical record and customer correspondence, paper archives, and print stream files that contain original systems of record between organizations and their customers



### DOCS

XLS, PDF, CSV, email, Word, PPT, HTML, HTML 5, plain text, XML, JSON, etc.



### MEDIA

Images, videos, audio, Flash, live streams, podcasts, etc.



### DATA STORAGE

SQL, NoSQL, Hadoop, doc repository, file systems, etc.



### BUSINESS APPS

Project management, marketing automation, productivity, CRM, ERP content management systems, HR, storage, talent management, procurement, expense management, Google Docs, intranets, portals, etc.



### PUBLIC WEB

Government, weather, competitive, traffic, regulatory, compliance, health care services, economic, census, public finance, stock, OSINT, the World Bank, SEC/Edgar, Wikipedia, IMDb, and other Web services



### SOCIAL MEDIA

Twitter, LinkedIn, Facebook, Tumblr, Blog, SlideShare, YouTube, Google+, Instagram, Flickr, Pinterest, Vimeo, Wordpress, IM, RSS, Review, Chatter, Jive, Yammer, etc.



### MACHINE LOG DATA

Event logs, server data, application logs, business process logs, audit logs, call detail records (CDRs), mobile location, mobile app usage, clickstream data, etc.



### SENSOR DATA

Medical devices, smart electric meters, car sensors, road cameras, satellites, traffic recording devices, processors found within vehicles, video games, cable boxes or household appliances, assembly lines, office buildings, cell towers and jet engines, air conditioning units, refrigerators, trucks, farm machinery, etc.

USING BIG DATA, ORGANIZATIONS CAN GENERATE ACTIONABLE INSIGHTS THAT ENABLE THEM TO DRIVE THEIR BUSINESS FORWARD. RAPID INTEGRATION OF THE EVER-EXPANDING POOL OF DATA SOURCES AND TYPES IS OPENING A WHOLE NEW WORLD OF POSSIBILITIES.

Data compiled by the domain experts at Kapow Software, a Kofax company, and is based on almost a decade of experience helping hundreds of large global enterprises and innovative start-ups across industries leverage critical data from disparate internal and external sources to meet business objectives.

**kapow**  
SOFTWARE  
A Kofax Company



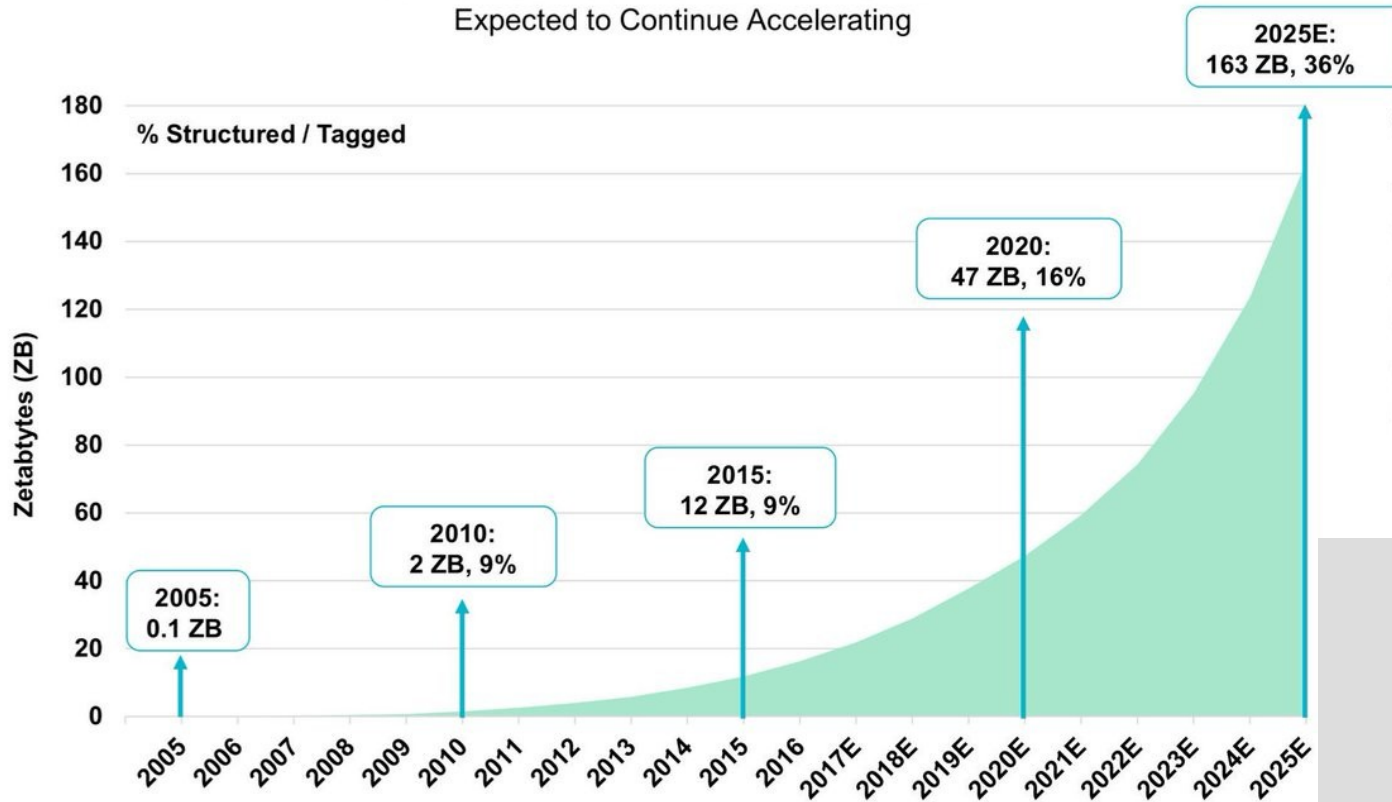
What amazing picture!



# Data Volume Growth Continues @ Rapid Clip...

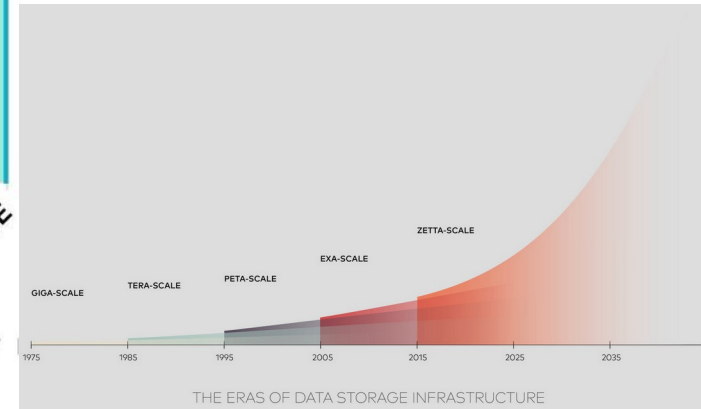
## % Structured / Tagged (~10%) Rising Fast...

**Information Created Worldwide =**  
Expected to Continue Accelerating



WHAT'S A ZETTABYTE?







1 kilobyte	1,000,000,000,000,000,000
1 megabyte	1,000,000,000,000,000,000,000
1 gigabyte	1,000,000,000,000,000,000,000,000
1 terabyte	1,000,000,000,000,000,000,000,000,000
1 petabyte	1,000,000,000,000,000,000,000,000,000,000
1 exabyte	1,000,000,000,000,000,000,000,000,000,000,000
1 zettabyte	1,000,000,000,000,000,000,000,000,000,000,000,000





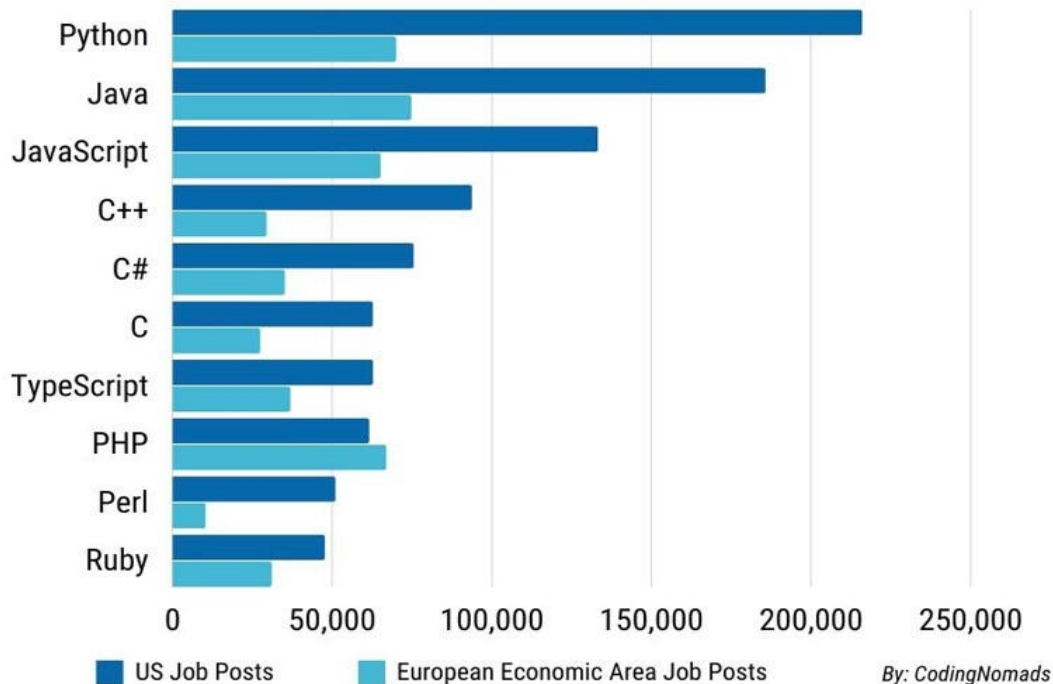
# The six Vs of big data

Big data is a collection of data from various sources, often characterized by what's become known as the 3Vs: *volume*, *variety* and *velocity*. Over time, other Vs have been added to descriptions of big data:

VOLUME	VARIETY	VELOCITY	VERACITY	VALUE	VARIABILITY
The amount of data from myriad sources.	The types of data: structured, semi-structured, unstructured.	The speed at which big data is generated.	The degree to which big data can be trusted.	The business value of the data collected.	The ways in which the big data can be used and formatted.
					

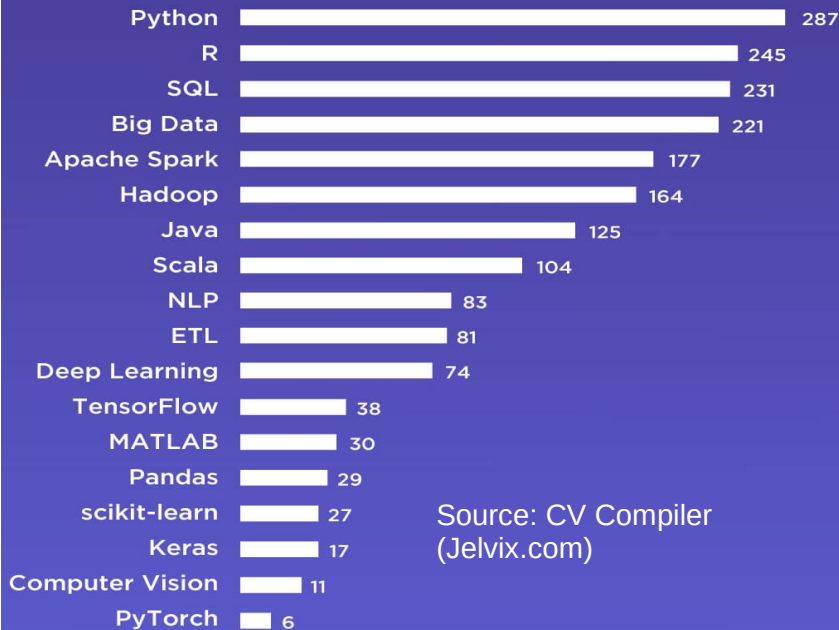
# Most in-demand programming languages of 2022

Based on LinkedIn job postings in the USA & Europe



## The skills Data Scientists need today

(based on 300 job listings from tech companies in June 2019)



### Colombia



### World



Source: [https://trends.google.com/trends/explore?date=today%205-y&geo=CO&q=%2Fm%2F012jm,%2Fm%2F05z1\\_](https://trends.google.com/trends/explore?date=today%205-y&geo=CO&q=%2Fm%2F012jm,%2Fm%2F05z1_)

# Data Scientist

Roadmap

- What do I have consolidated?
- What do I need to learn and what do I still need to learn more about?



## Mathematics

- Linear Algebra
- Analytics Geometry
- Matrix
- Vector Calculus
- Optimization
- Regression
- Dimensionality Reduction
- Density Estimation
- Classification

## Probability

- Discrete Distribution
  - Binomial
  - Bernoulli
  - Geometric etc
- Continuous Distribution
  - Uniform
  - Exponential
  - Gamma
- Normal Distribution
- Introduction to Probability
- 1D Random Variable
- Function of One Random Variable
- Joint Probability Distribution

## Statistics

- Introduction to Statistics
- Data Description
- Random Samples
- Sampling Distribution
- Parameter Estimation
- Hypotheses Testing
- ANOVA
- Reliability Engineering
- Stochastic Process
- Computer Simulation
- Design of Experiments
- Simple Linear Regression
- Correlation
- Multiple Regression
- Nonparametric Statistics
  - Sign Test
  - The Wilcoxon Signed-Rank Test
  - The Wilcoxon Rank Sum Test
  - The Kruskal-Wallis Test
- Statistical Quality Control
- Basic of Graphs

## Programming

- | Python   | R   |
|--|---|
| <b>Python Basics</b> <ul style="list-style-type: none"><li>• List</li><li>• Set</li><li>• Tuples</li><li>• Dictionary</li><li>• Function, etc.</li></ul> | <b>R Basic</b> <ul style="list-style-type: none"><li>• Vector</li><li>• List</li><li>• Data Frame</li><li>• Matrix</li><li>• Array, etc</li></ul> |
| <b>NumPy</b>   | <b>dplyr</b>  |
| <b>Pandas</b>  | <b>ggplot2</b>  |
| <b>Matplotlib/Seaborn, etc.</b>  | <b>Tidyr</b>  |
|  | <b>Shiny, etc.</b>  |
| <b>DataBase</b>  | <b>Other</b>  |
| <b>SQL</b>   | <b>Data Structure</b> <ul style="list-style-type: none"><li>• Array, etc</li></ul>  |
| <b>MongoDB</b>   | <b>Web Scraping</b>   |
|  | <b>Linux</b>  |
|  | <b>Git</b>  |

## Machine Learning

- | Introduction  | Intermediate   |
|---|--|
| <ul style="list-style-type: none"><li>• How Model Works</li><li>• Basic Data Exploration</li><li>• First ML Model</li><li>• Model Validation</li><li>• Underfitting &amp; Overfitting</li><li>• Random Forests</li><li>• scikit-learn</li></ul> | <ul style="list-style-type: none"><li>• Handling Missing Values</li><li>• Handling Categorical Variables</li><li>• Pipelines</li><li>• Cross-Validation</li><li>• XGBoost</li><li>• Data Leakage</li></ul> |

## Deep Learning

- |  |  |
|--|--|
| <ul style="list-style-type: none"><li>• Artificial Neural Network</li><li>• Convolutional Neural Network</li><li>• Recurrent Neural Network</li><li>• Keras</li><li>• PyTorch</li><li>• TensorFlow</li></ul> | <ul style="list-style-type: none"><li>• A Single Neuron</li><li>• Deep Neural Network</li><li>• Stochastic Gradient Descent</li><li>• Overfitting and Underfitting</li><li>• Dropout Batch Normalization</li><li>• Binary Classification</li></ul> |
|--|--|

## Feature Engineering

- Baseline Model
- Categorical Encodings
- Feature Generation
- Feature Selection

## Natural language Processing

- Text Classification
- Word Vectors

## Data Visualization Tools

- Excel VBA
- BI (Business Intelligence)
  - Tableau
  - Power BI
  - Qlik View
  - Qlik Sense

## Deployment

- Microsoft Azure
- Heroku
- Google Cloud Platform
- Flask
- Django

## Other Points

- Domain Knowledge
- Communication Skill
- Reinforcement Learning
- Case Studies
  - Data Science at Netflix
  - Data Science at Flipkart
  - Project on Credit Card Fraud Detection
  - Project on Movie Recommendation , etc.

## Keep Practicing

# PYPL Popularity of Programming Language

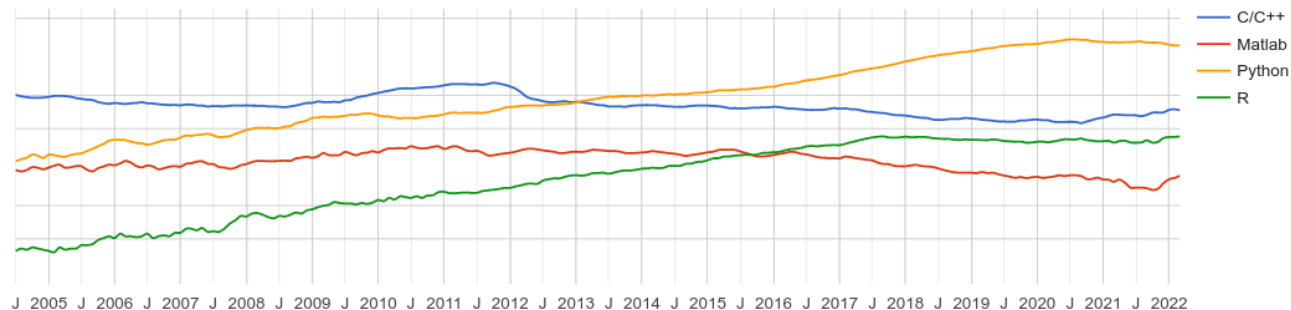
Worldwide, Mar 2022 compared to a year ago:

Rank	Change	Language	Share	Trend
1		Python	28.27 %	-2.0 %
2		Java	18.03 %	+0.8 %
3		JavaScript	8.86 %	+0.4 %
4		C#	7.51 %	+0.6 %
5		C/C++	7.32 %	+0.6 %
6		PHP	5.71 %	-0.4 %
7		R	4.23 %	+0.5 %
8		Objective-C	2.28 %	-1.2 %
9	↑	TypeScript	2.11 %	+0.3 %
10	↓	Swift	2.01 %	+0.2 %
11		Matlab	1.87 %	+0.2 %
12		Kotlin	1.57 %	-0.1 %

The PYPL Popularity of Programming Language Index is created by analyzing how often language tutorials are searched on Google.

The more a language tutorial is searched, the more popular the language is assumed to be. It is a leading indicator. The raw data comes from Google Trends.

If you believe in collective wisdom, the PYPL Popularity of Programming Language index can help you decide which language to study, or which one to use in a new software project.



<https://pypl.github.io/PYPL.html>



***Getting started!***



*Next...*

# Creating environments and packages.....

- First: Install Acaconda (Windows)  
<https://docs.anaconda.com/anaconda/install/windows/>
- Next: Install dependences (Python)  
<https://medium.com/@GalarnykMichael/install-python-anaconda-on-windows-2020-f8e188f9a63d>
- Next: Install Jupyter Notebook  
<https://www.geeksforgeeks.org/how-to-install-jupyter-notebook-in-windows/>
- Next: Install another kernel  
<https://datatofish.com/r-jupyter-notebook/>

Anaconda <img alt="Anaconda logo" data-bbox="855 105 870 135"/>  
Computer program



Anaconda is a distribution of the Python and R programming languages for scientific computing, that aims to simplify package management and deployment. The distribution includes data-science packages suitable for Windows, Linux, and macOS.

[Wikipedia](#)

**License:** [Freemium](#) (Miniconda and the Individual Edition are free software, but the other editions are software as a service)

**Operating system:** [Windows](#), [macOS](#), [Linux](#)

**Stable release:** 2021.11 / 17 November 2021; 3 months ago

**Developer(s):** Anaconda, Inc. (previously [Continuum Analytics](#))

**Initial release:** 0.8.0/17 July 2012; 9 years ago

**Programming language:** [Python](#)

People also search for

[View 10+ more](#)

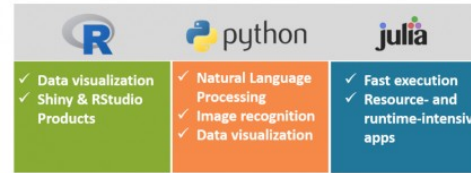


PyCharm

Visual  
Studio C...

Conda

Spyder



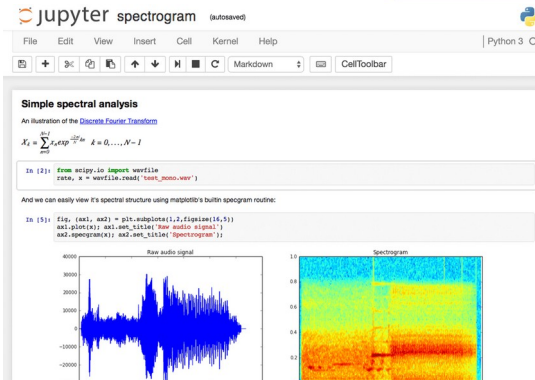
✓ Data visualization  
✓ Shiny & RStudio  
Products



✓ Natural Language  
Processing  
✓ Image recognition  
✓ Data visualization



✓ Fast execution  
✓ Resource- and  
runtime-intensive  
apps



# We need GitHub too...

## What is Git and what are its origins?

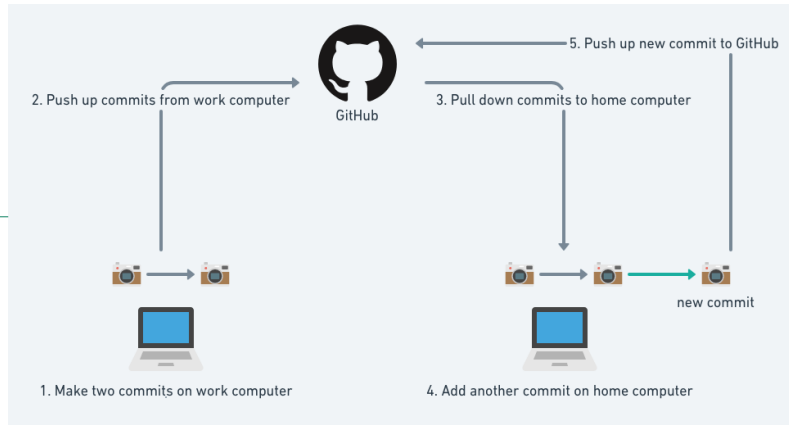
Have you ever had a document named something like report\_final\_draft\_final(3).doc? If so, you've felt the pain of managing and sharing files. Keeping track of the changes to a file over time is difficult but important. Git is a Version Control System (VCS) – a tool that helps us to keep track of differences in a file or collection of files over time.

## What is Git used for?

Git is used for managing the changes to a project over time. A project might be just a single file, a handful of files, or thousands of files. Those files can be anything from plain text to images or videos.

## What is Github and what are its origins? How did its creation change the way people collaborate?

GitHub, developed in 2008, is a web application that hosts Git repositories. The team that started GitHub saw that Git could solve important problems for many teams – but Git itself is often difficult to use. GitHub adds a bunch of collaboration and exploration tools on top of Git to help you (and your team) be more productive.



## WHAT ARE GIT & GITHUB?



**Git** is a version control system to keep track of changes to files and projects over time.

**GitHub** is a website that hosts Git repositories online, making it easier for developers to share code.

**Repositories** (or "repos") are folders which contain intentional snapshots of progress called commits.

## COMPANIES WHICH USE GITHUB



COURSE  
REPORT

