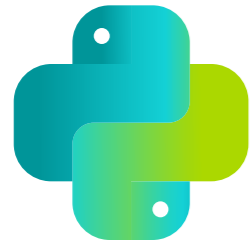2018

# Data Science Survey

In spring 2018, we polled over 1,600 people involved in Data Science and based in the US, Europe, Japan, and China, in order to gain insight into how this industry sector is evolving. Here's what we learned.

## Methodology

We distributed the survey via targeted ads on Facebook, Twitter, and LinkedIn. We screened respondents by excluding those who replied "I am not involved in data analysis." We collected 400 complete and valid responses from the US, Japan, and China. To represent Europe, we used quotas for select European countries to collect a set of responses which also totaled 400.

Some bias is likely present as JetBrains users may have been more willing on average to complete the survey.
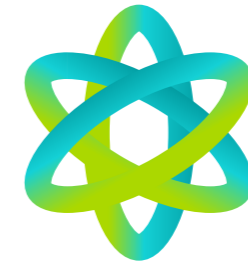
# Key Takeaways

## Most popular language

Python is currently the most popular language among data scientists.

## Primary language

Most people assume that Python will remain the primary programming language in the field for the next 5 years.

## R, Keras and Tableau

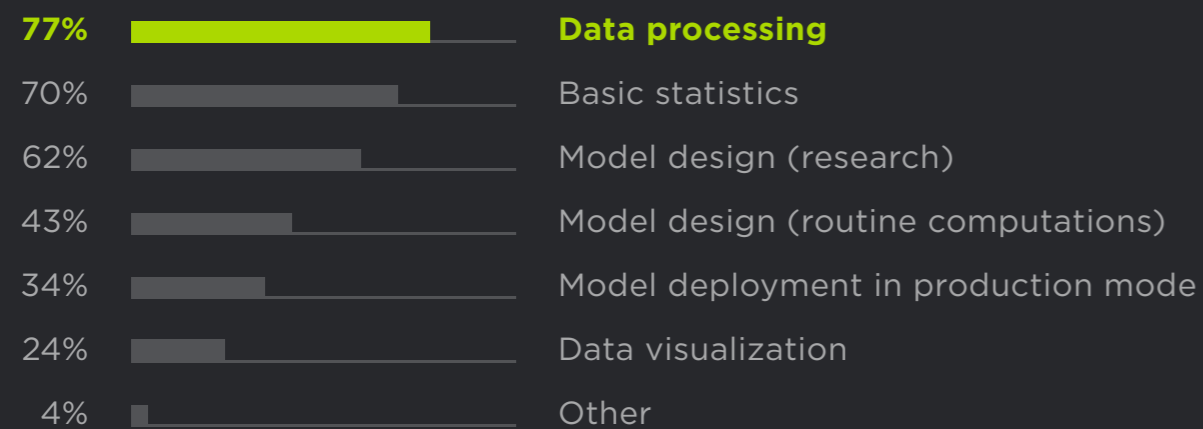Data Science professionals tend to use R, Keras, and Tableau, while amateur data scientists are more likely to prefer Microsoft Azure ML.

# 1

## Tasks

**Which of the following are you involved in?**

Number of answers: 1282

| | |
|---|---|
| **77%** | **Data processing** |
| 70% | Basic statistics |
| 62% | Model design (research) |
| 43% | Model design (routine computations) |
| 34% | Model deployment in production mode |
| 24% | Data visualization |
| 4% | Other |

This question was only answered by respondents who are professionally involved in data analysis.

**Natalia Vassilieva**
Head of Software and AI,
Hewlett Packard Labs

According to the very first table, a relatively small percentage of the respondents work on model deployment in production mode. This correlates with multiple market research studies which report that the majority of enterprises are just starting to explore machine learning and deep learning.

They have small teams working on PoCs*, and model deployment in production still needs to be addressed. But I expect this type of activity to become more and more visible within the next few years when more and more businesses will proceed from PoCs to production deployments.
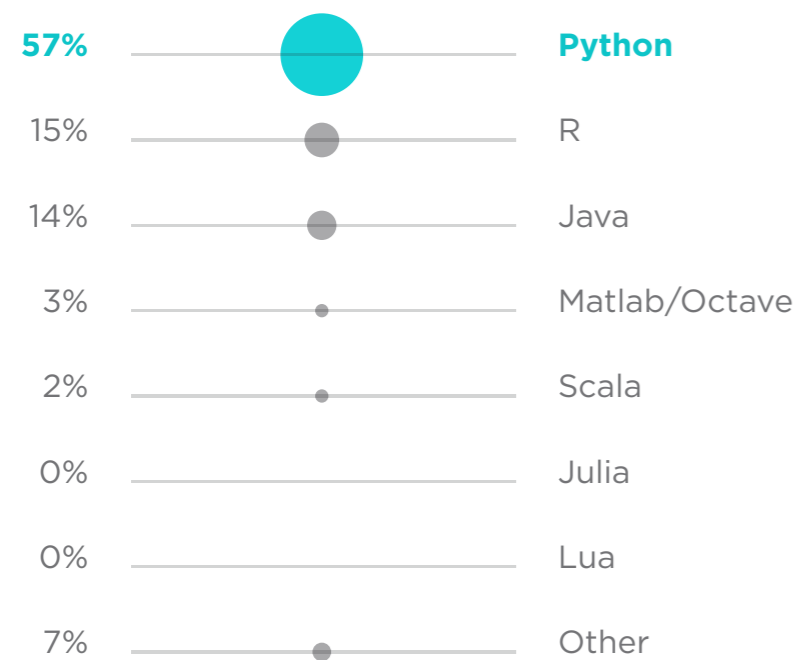
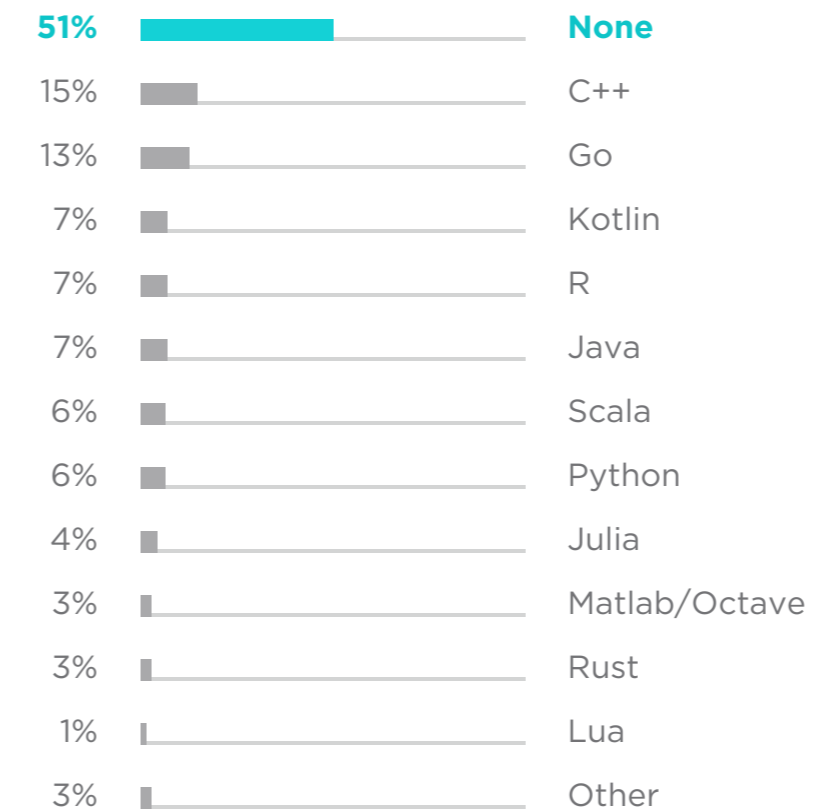*PoCs — Proof of concepts

# Programming Languages and Tools

**2**

## Main programming language for data analysis

Number of answers: 1522

| | |
|---|---|
| **57%** | **Python** |
| 15% | R |
| 14% | Java |
| 3% | Matlab/Octave |
| 2% | Scala |
| 0% | Julia |
| 0% | Lua |
| 7% | Other |

## Do you plan to adopt or migrate to other languages in the next 12 months?

Number of answers: 1522

| | |
|---|---|
| **51%** | **None** |
| 15% | C++ |
| 13% | Go |
| 7% | Kotlin |
| 7% | R |
| 7% | Java |
| 6% | Scala |
| 6% | Python |
| 4% | Julia |
| 3% | Matlab/Octave |
| 3% | Rust |
| 1% | Lua |
| 3% | Other |

15% of data scientists are going to adopt or migrate to C++ in the next 12 months. This is probably due to performance issues.

7% of respondents plan to start using the Kotlin language in the next 12 months.

# In your opinion, what programming language will be most used for data analysis in the next 5 years?

Number of answers: 1522

Most respondents believe that Python will remain on top for the next 5 years.

| | |
|---|---|
| **56%** | **Python** |
| 9% | R |
| 7% | Java |
| 6% | Go |
| 5% | C++ |
| 3% | Julia |
| 3% | Kotlin |
| 2% | Scala |
| 2% | Matlab/Octave |
| 1% | Rust |
| 0% | Lua |

Note that Go, C++, Julia, Kotlin, Rust, and Lua were unavailable as answer choices for regularly and most used languages.

Overall, people tend to choose the language they use. Of those who don't use a language they think will dominate, most want to start using it. Half of those who believe Kotlin will dominate are planning to adopt it in the nearest future.

**Natalia Vassilieva**
Head of Software and AI, Hewlett Packard Labs

No surprises with the programming languages. Traditional data scientists are some of the most likely to still use R, there are plenty of statistics libraries for R. The new generation of data scientists are choosing Python.
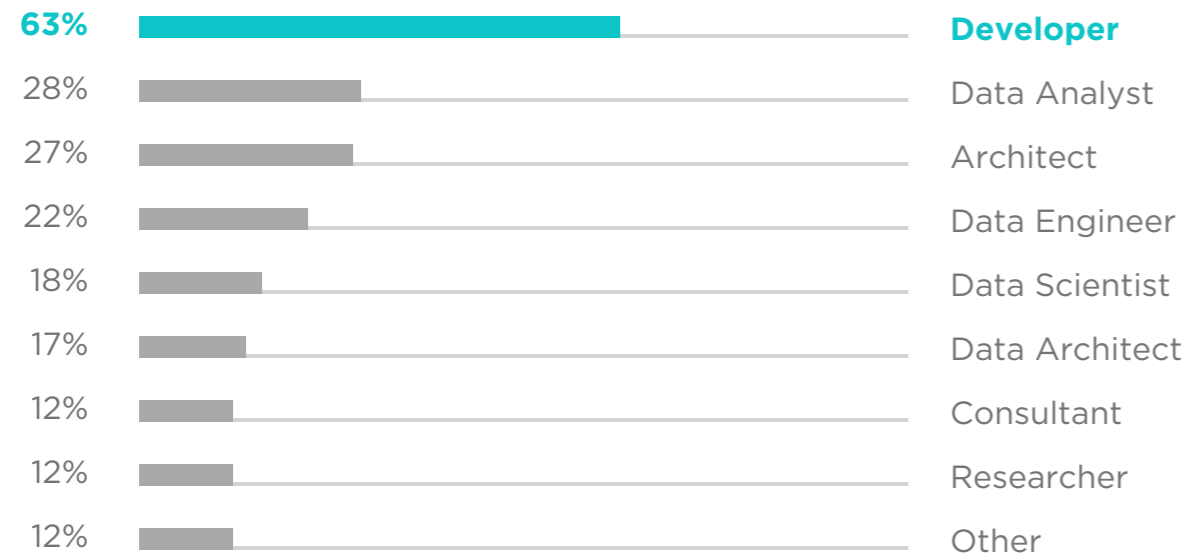
When it comes to high-performance data analytics, I'd expect to see C/C++ in the picture. Currently, we are observing that many HPC techniques and tools are being adopted and re-used for high-performance data analytics and deep learning.

# Kotlin adopters

7% of data analysts using a programming language want to adopt Kotlin for Data Science in the nearest future.

Note this was a question with checkboxes. Shares may total more than 100%.

## Which of the following best describes your job roles regardless of your position level?

Number of answers: 60

| | |
|---|---|
| **63%** | **Developer** |
| 28% | Data Analyst |
| 27% | Architect |
| 22% | Data Engineer |
| 18% | Data Scientist |
| 17% | Data Architect |
| 12% | Consultant |
| 12% | Researcher |
| 12% | Other |

## What programming languages do you regularly use for data analysis, if any?

Number of answers: 112

**72%**
**Python**

| 62% | 13% | 4% | 4% |
|---|---|---|---|
| Java | Scala | Julia | Other |

| 23% | 13% | 4% | 0% |
|---|---|---|---|
| R | Matlab/Octave | Lua | None |

# Kotlin Learning

**Vitaly Khudobakhshov**
Product Manager, JetBrains

Kotlin is a general-purpose language running on the Java virtual machine. It is concise and easily integrates with popular data processing frameworks such as Hadoop and Spark.

Kotlin is statically typed and uses type inference that increases its reliability. These features all make Kotlin a handy instrument for data engineering and data science.

Thomas Nield has assembled a helpful collection of Kotlin resources for data science on his Github.

If you are new to Kotlin and are considering it for your next language, start from learning the basic syntax.

If you are already familiar with Java, you may want to have a play with Kotlin Koans.
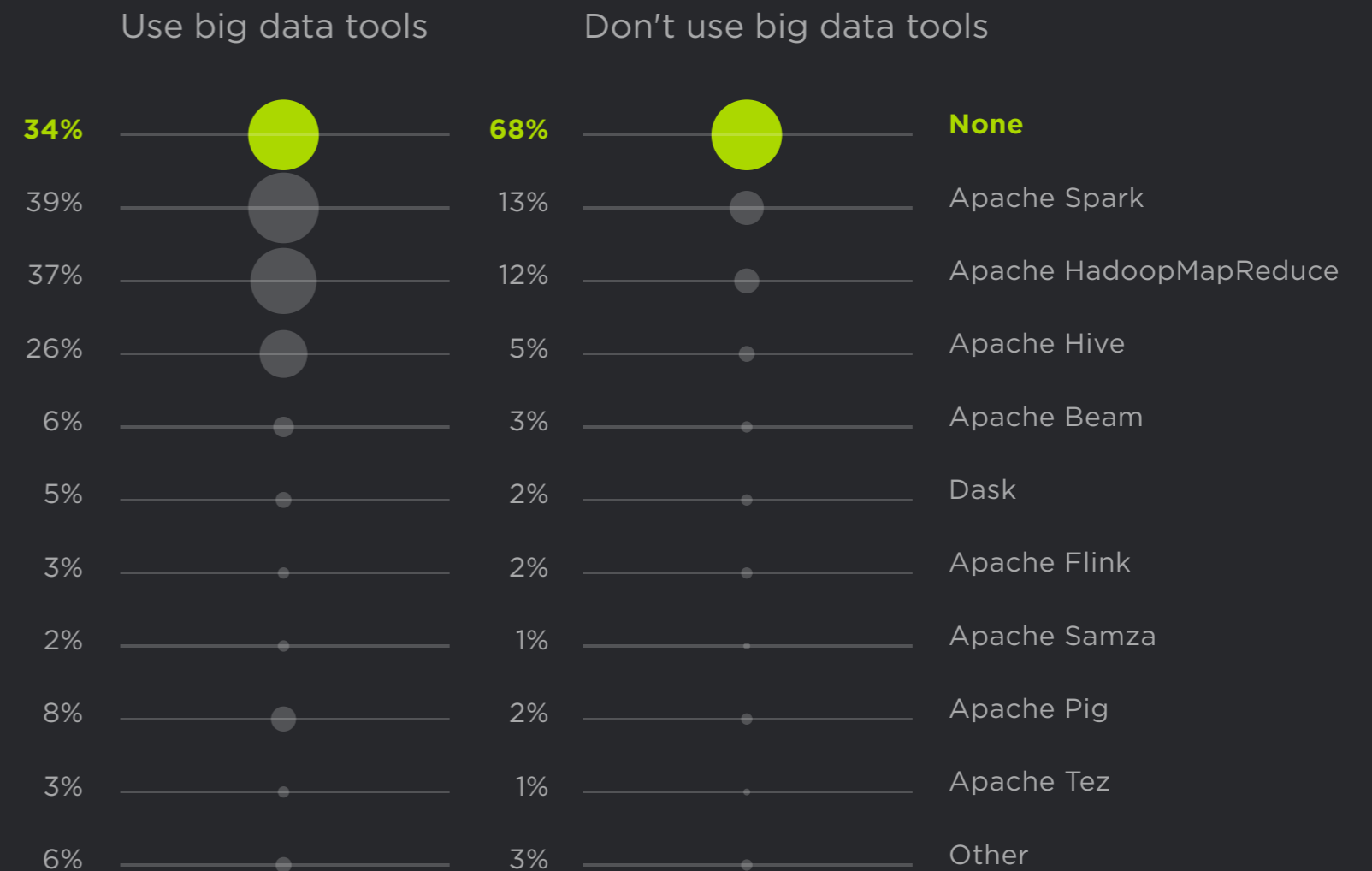
# Tools & Technologies & Editors

## Big Data tools

Number of answers: 1477

> Note this was a question with checkboxes.
> Shares may total more than 100%.

One third of those who say they work with big data
don't use any big data tools. Conversely, a third
of those who do NOT work with Big data DO use
some big data tools. Still, this self-identification
does correlate with formal factors.

| Use big data tools | Don't use big data tools | |
|---|---|---|
| 34% | 68% | None |
| 39% | 13% | Apache Spark |
| 37% | 12% | Apache HadoopMapReduce |
| 26% | 5% | Apache Hive |
| 6% | 3% | Apache Beam |
| 5% | 2% | Dask |
| 3% | 2% | Apache Flink |
| 2% | 1% | Apache Samza |
| 8% | 2% | Apache Pig |
| 3% | 1% | Apache Tez |
| 6% | 3% | Other |

# IDEs and Editors

Number of answers: 1522

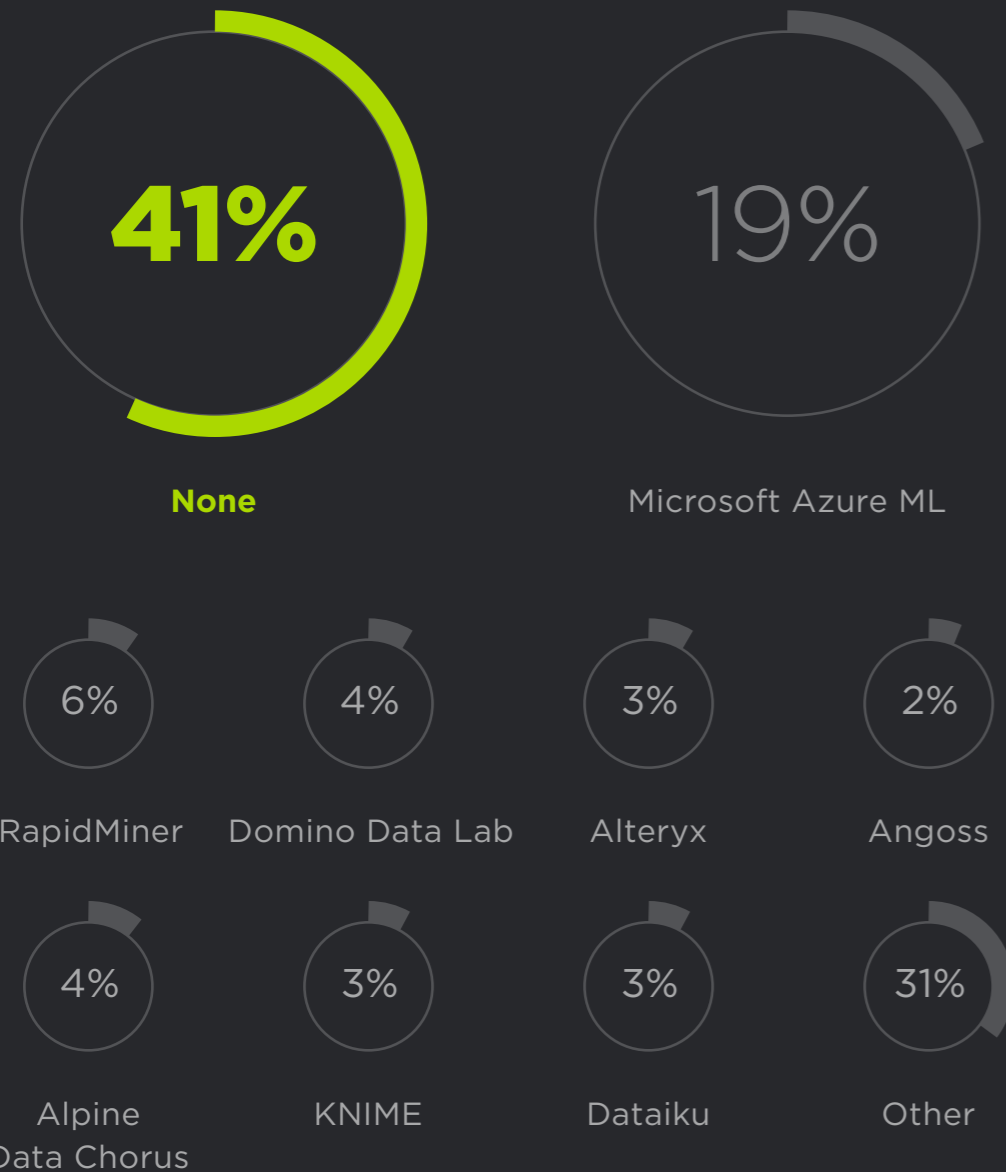| | |
|---|---|
| **43%** | **Jupyter/IPython notebook** |
| 38% | PyCharm |
| 26% | RStudio |
| 22% | IntelliJ IDEA |
| 18% | Atom |
| 15% | Visual Studio Code |
| 13% | Vim |
| 13% | Eclipse |
| 12% | Sublime |
| 11% | Visual Studio |
| 11% | JupyterLab |
| 10% | Google Cloud Datalab |
| 8% | Spyder IDE |
| 4% | Colaboratory |
| 3% | Zeppelin |
| 1% | Rodeo |
| 4% | Other |

Note this was a question with checkboxes. Shares may total more than 100%.

# Which tools do you use for data analysis, if any?

Number of answers: 1666

**41%**

**None**

19%

Microsoft Azure ML

6%

RapidMiner

4%

Domino Data Lab

3%

Alteryx

2%

Angoss

4%

Alpine
Data Chorus

3%

KNIME

3%

Dataiku

31%

Other

# What deep learning libraries do you use, if any?

Number of answers: 1666

| 41% | None |
| **47%** | **TensorFlow** |
| 26% | Keras |
| 11% | Torch |
| 6% | Theano |
| 7% | Other |

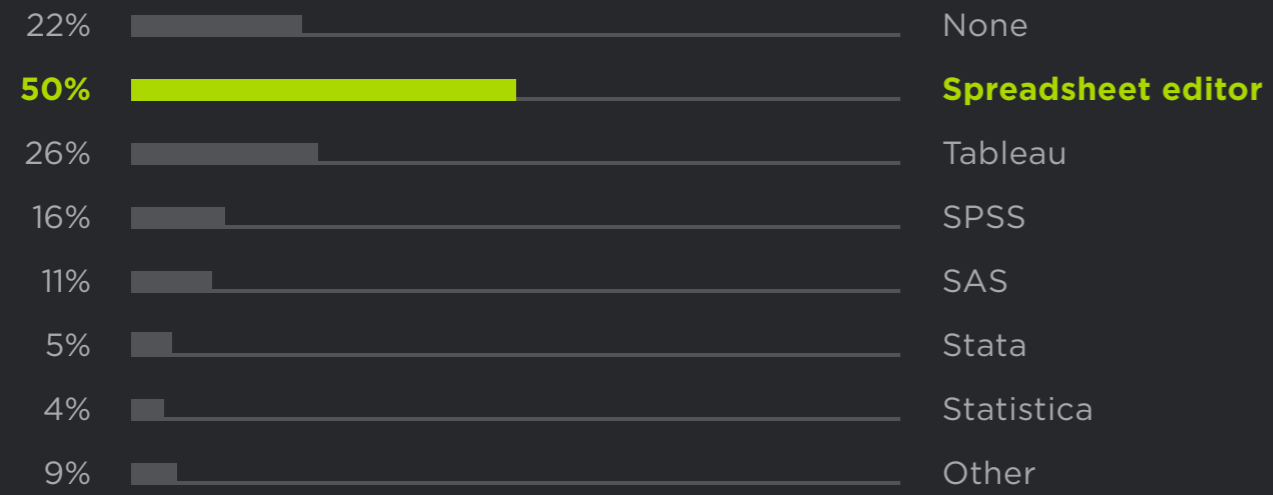Note this was a question with checkboxes. Shares may total more than 100%.

Apparently, if you do deep learning, you do it with TensorFlow. Nearly 80% of those using deep learning libraries use TensorFlow, and almost all Keras users use it alongside TensorFlow.

## Which statistics packages do you use to analyze and visualize data, if any?

Number of answers: 1666

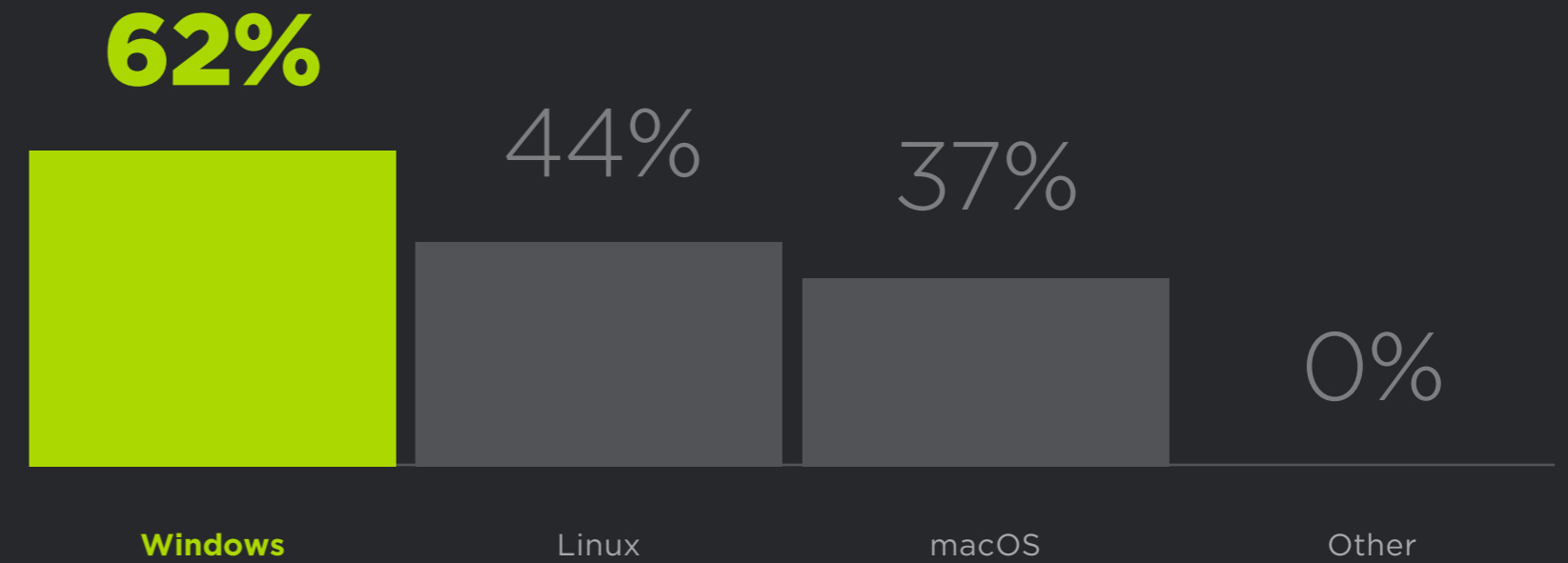| | |
|---|---|
| 22% | None |
| 50% | Spreadsheet editor |
| 26% | Tableau |
| 16% | SPSS |
| 11% | SAS |
| 5% | Stata |
| 4% | Statistica |
| 9% | Other |

Note this was a question with checkboxes.
Shares may total more than 100%.

## What operating systems do you use as your work environment for data analysis?
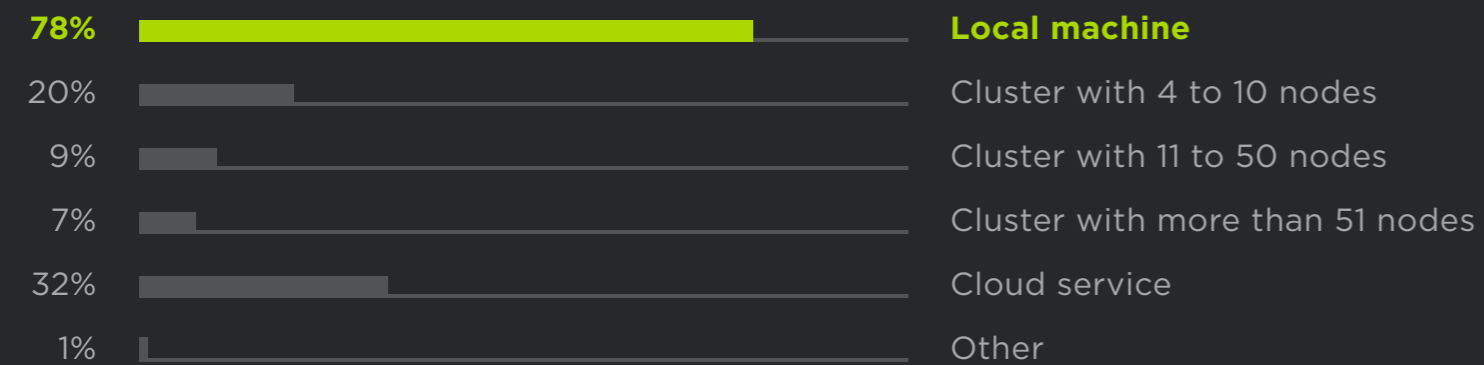
Number of answers: 1666



| 62% | 44% | 37% | 0% |
|---|---|---|---|
| Windows | Linux | macOS | Other |

Note this was a question with checkboxes.
Shares may total more than 100%.

# What do you use to perform computations?

Number of answers: 1666

| | |
|---|---|
| **78%** | **Local machine** |
| 20% | Cluster with 4 to 10 nodes |
| 9% | Cluster with 11 to 50 nodes |
| 7% | Cluster with more than 51 nodes |
| 32% | Cloud service |
| 1% | Other |

Note this was a question with checkboxes. Shares may total more than 100%.

78% data science specialists perform computations on local machines.

# Cloud services

Number of answers: 527

**56%**
**Amazon Web Services (AWS)**

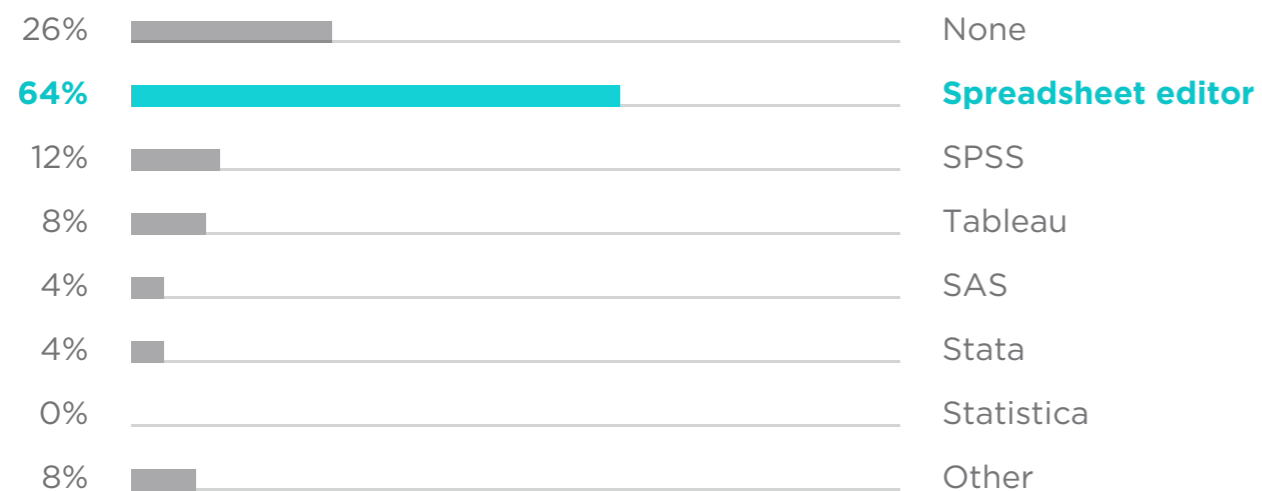41%
Google Cloud Platform

28%
Microsoft Azure

14%
Other

# Non-programmers

We received 77 responses from people who don't use any programming languages and aren't about to adopt any (5% of all data analysts who responded).

These respondents use spreadsheet editors more often than average, and most of them work in non-IT industries. They also tend to use data analysis tools less often.

## Which statistics packages do you use to analyze and visualize data, if any?

Number of answers: 77

| Percent | Package |
|---|---|
| 26% | None |
| **64%** | **Spreadsheet editor** |
| 12% | SPSS |
| 8% | Tableau |
| 4% | SAS |
| 4% | Stata |
| 0% | Statistica |
| 8% | Other |

## What is the industry you primarily analyze data for?

Number of answers: 43

**81%**
**A non-IT industry**

19%
IT

**5**

# Manager's expertise

## What is your manager's level of expertise in data analysis?

Number of answers: 924

| | |
|---|---|
| 19% | Has no qualifications in data analysis |
| 25% | Has basic qualifications in data analysis |
| **27%** | **Has average qualifications in data analysis** |
| 22% | Highly qualified data analysis specialist |
| 7% | Top expert in data analysis |

Almost half of all respondents report to managers who have little or no qualifications in data analysis.
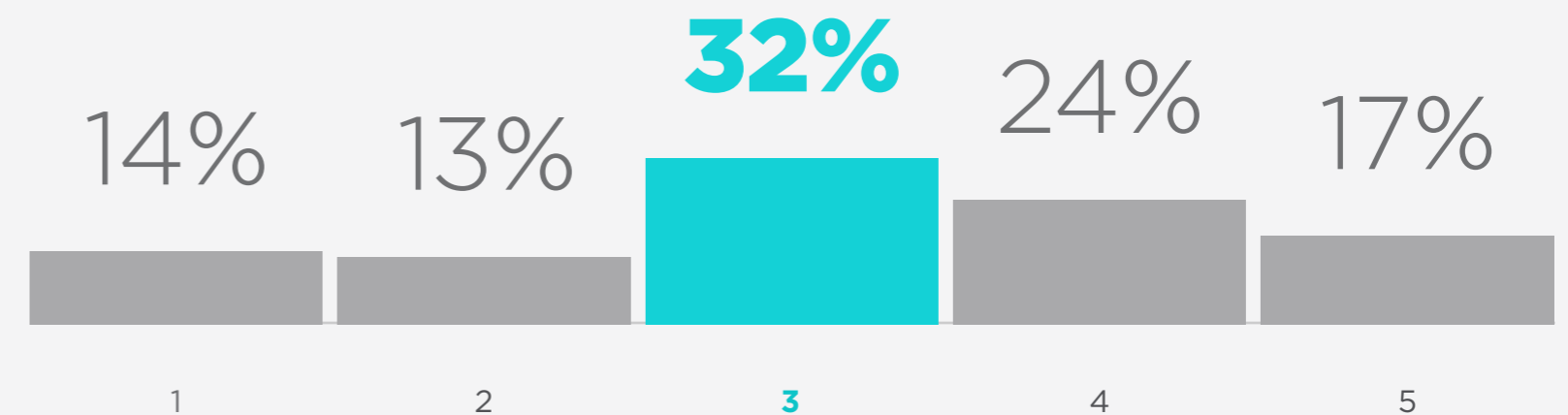
## To what extent do you associate the following phrase with your manager?: "My manager gives me realistic assignments that are relevant to my skills and responsibilities, with a clear and specific description of the requirements."

Number of answers: 918

**Answers:**
**1** = not at all
**5** = a great deal

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| 14% | 13% | **32%** | 24% | 17% |

# Correlation of Manager's expertise and assignment score

Number of answers: 924

| | Score no qualifications<br>N: 177 | Basic qualifications<br>N: 230 | Average qualifications<br>N: 246 | Highly qualified<br>N: 199 | Top expert<br>N: 66 |
|---|---|---|---|---|---|
| **1** | **31%** | 16% | 7% | 8% | 11% |
| **2** | 15% | 20% | 10% | 9% | 3% |
| **3** | 29% | **34%** | **37%** | 30% | 17% |
| **4** | 16% | 18% | 27% | **32%** | 29% |
| **5** | 10% | 11% | 18% | 22% | **41%** |

**Answers:**

**1** = not at all
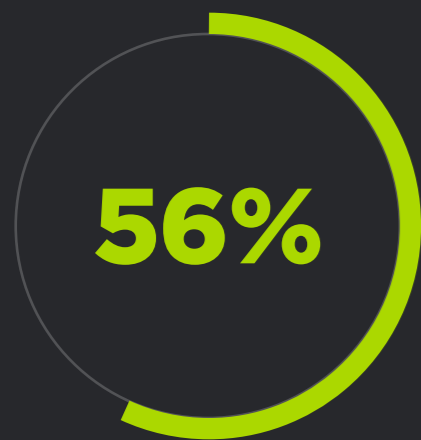
**5** = a great deal

**6**

# Industry & Demographic

## What is the industry you primarily analyze data for?

Number of answers: 1666

**56%**
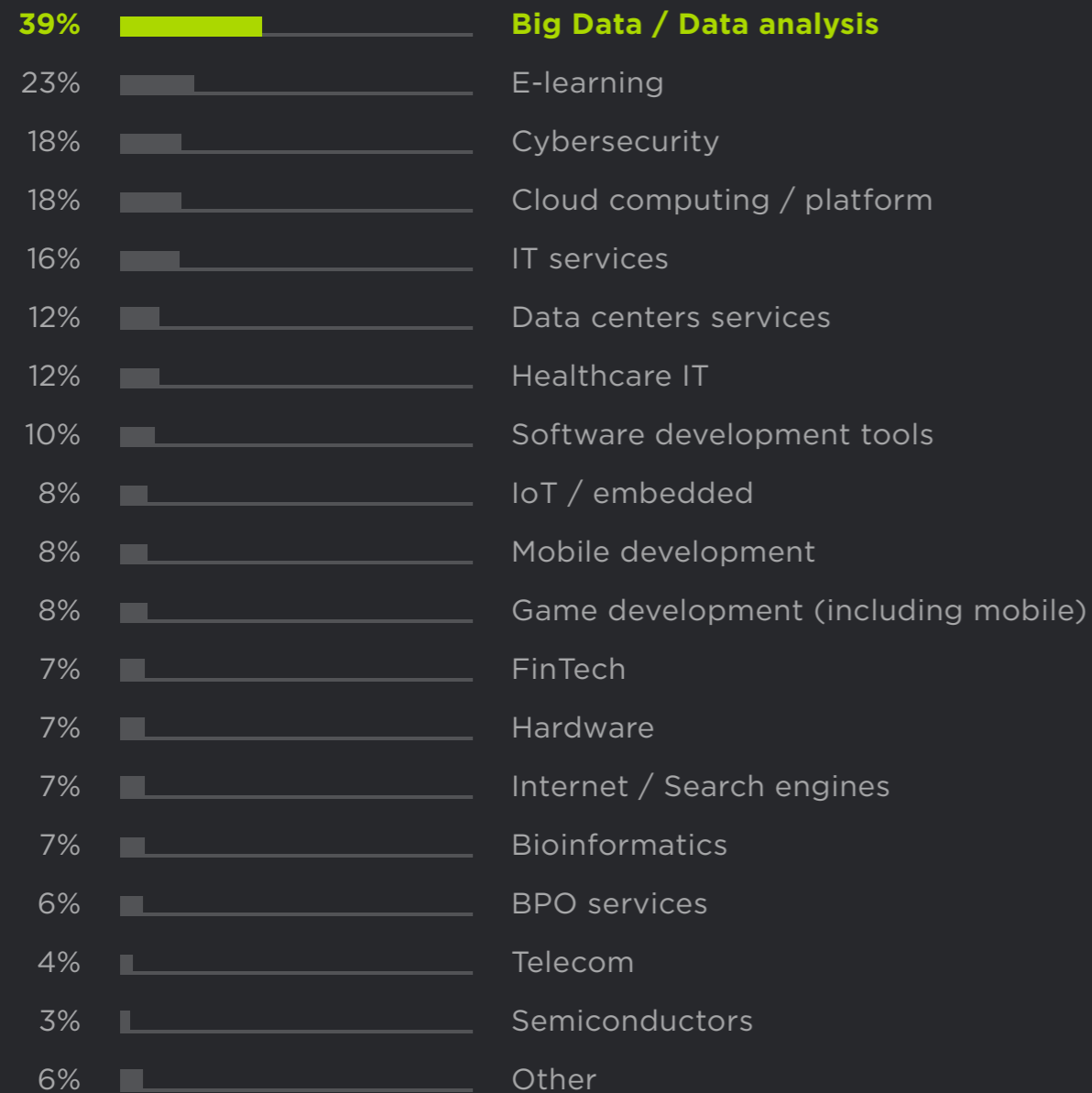
**44%**

**A non-IT industry**

IT

This was an optional question.

## Which fields of IT do you primarily analyze data for?

Number of answers: 703

| | |
|---|---|
| **39%** | **Big Data / Data analysis** |
| 23% | E-learning |
| 18% | Cybersecurity |
| 18% | Cloud computing / platform |
| 16% | IT services |
| 12% | Data centers services |
| 12% | Healthcare IT |
| 10% | Software development tools |
| 8% | IoT / embedded |
| 8% | Mobile development |
| 8% | Game development (including mobile) |
| 7% | FinTech |
| 7% | Hardware |
| 7% | Internet / Search engines |
| 7% | Bioinformatics |
| 6% | BPO services |
| 4% | Telecom |
| 3% | Semiconductors |
| 6% | Other |

Note this was a question with checkboxes. Shares may total more than 100%.

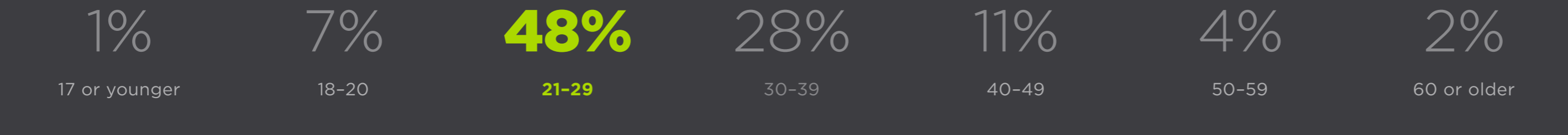## Which industry or industries do you primarily analyze data for?

Number of answers: 933

| | |
|---|---|
| **16%** | **Accounting / Finance / Insurance** |
| **16%** | **Science** |
| 15% | Education / Training |
| 12% | Bioinformatics |
| 11% | Logistics/ Transportation |
| 11% | Administration / Management / Business Development |
| 11% | Medicine / Health |
| 10% | Machinery |
| 10% | Business / Strategic Management |
| 10% | Sales / Distribution / Retail |
| 9% | Banking / Real Estate / Mortgage Financing |
| 8% | Manufacturing |
| 8% | Marketing |
| 6% | Design |
| 6% | Restaurants / Hospitality / Tourism |
| 6% | Construction / Architecture |
| 5% | Entertainment / Mass media and information / Publishing |
| 5% | Human Resources |
| 5% | Non-profit |
| 5% | Customer Support |
| 4% | Security |
| 4% | Service / Maintenance |
| 3% | Law |
| 8% | Other |

# Demographics

## Age range

Number of answers: 1666

| 1% | 7% | **48%** | 28% | 11% | 4% | 2% |
|---|---|---|---|---|---|---|
| 17 or younger | 18–20 | **21–29** | 30–39 | 40–49 | 50–59 | 60 or older |

## What is your main employment status?

Number of answers: 1666

| | |
|---|---|
| **61%** | **Fully employed by a company / organization** |
| 25% | Student |
| 5% | Freelancer |
| 4% | Partially employed by a company / organization |
| 3% | Looking for a job |
| 1% | Retired |
| 1% | Other |

# What is your main employment status?

Number of answers: 1666

**44%**

**Bachelor's degree (BA, BSc, B.Eng., etc.)**

3%

Professional degree (JD, MD, etc.)

36%

Master's degree (MA, MSc, M.Eng., MBA, etc.)

7%

Doctoral degree (Ph.D, Ed.D., etc.)

10%

Other

# Work experience

Number of answers: 924

17%

**42%**

23%

9%

10%

< 1 year

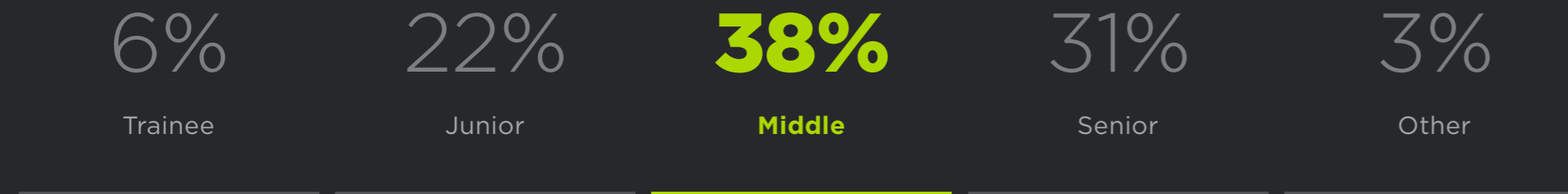**1 to 3 years**

3 to 6 years

6 to 10 years

> 10 years

This question was directed to professionals, that is, people professionally involved in data science or data analysis and working full-time or part-time.

# Work environment and employment
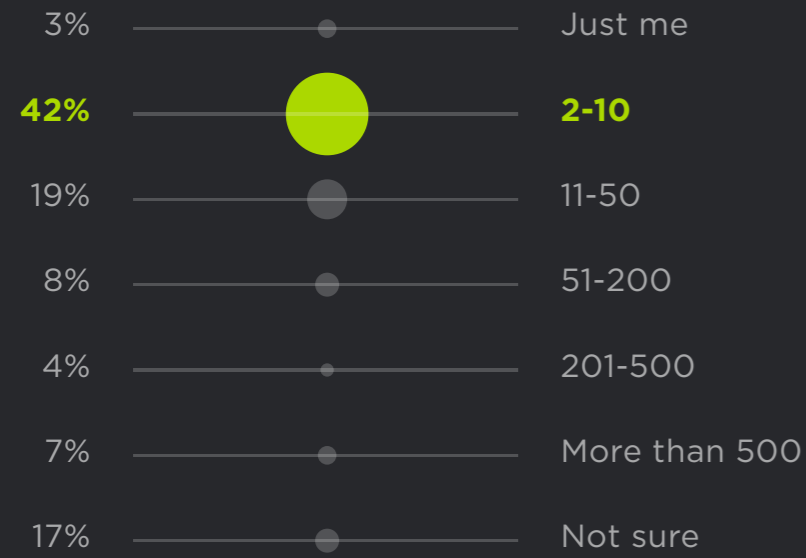
## Position level

Number of answers: 1014

| 6% | 22% | **38%** | 31% | 3% |
|---|---|---|---|---|
| Trainee | Junior | **Middle** | Senior | Other |

## Company size

Number of answers: 1086

| 1% | Just me |
|---|---|
| 9% | 2–10 |
| 19% | 11–50 |
| **28%** | **51–500** |
| 8% | 501–1,000 |
| 13% | 1,001–5,000 |
| 22% | More than 5,000 |
| 1% | Not sure |

## Number of data analysts in company

Number of answers: 1666

| | |
|---|---|
| 3% | Just me |
| **42%** | **2-10** |
| 19% | 11-50 |
| 8% | 51-200 |
| 4% | 201-500 |
| 7% | More than 500 |
| 17% | Not sure |

## Job Role

Number of answers: 917

| | |
|---|---|
| **37%** | **Developer** |
| 33% | Data Analyst |
| 32% | Data Scientist |
| 21% | Data Engineer |
| 19% | Researcher |
| 17% | Business Analyst |
| 14% | Consultant |
| 14% | Architect |
| 10% | Data and Analytics Manager |
| 9% | Data Architect |
| 8% | Statistician |
| 8% | General / Product Manager |
| 5% | Other |

Note this was a question with checkboxes. Shares may total more than 100%.

# JetBrains products for data science and big data



PyCharm Professional Edition is a Python IDE that enables Data Scientists and Web developers to become far more productive.

It offers in-depth Python code analysis and integrates with various libraries, frameworks, and tools. PyCharm's scientific tools are designed specifically with professional data analysts in mind and include a scientific development mode, integration with conda, code cells, Jupyter Notebook support, and much more. There is first-class support available for SQL databases as well.

jetbrains.com/pycharm

# JetBrains products for data science and big data

Datalore is an intelligent web application for data analysis and visualization for Python, with built-in tools and libraries for machine learning all included.

datalore.io

The smart Python code editor helps users write better code with suggestions, autocompletion, and syntax highlighting. Incremental recalculation enables dependencies between multiple computations to be followed, so users don't have to track which parts of the code were affected by recent edits. And there is access to the extended data storage and high-performance computational resources (including GPU instances) for an enhanced exploration experience.

# Thank you for your time!
# We hope you found our report useful.

If you have any questions or suggestions,
please contact us at surveys@jetbrains.com.