

R vs Python for Data Science and Visualization

By Martin Anderson



First published November 23rd, 2020 at:

<https://www.iflexion.com/blog/r-vs-python-data-science>

[Web-archived version](#)

Python and R are frequently pitted against each other as if one or the other might eventually become the Betamax¹ of programming languages for data science analysis. If economy of scale and growth of market share are the only criteria, we can award Python as the winner without further consideration.

However, comparing R and Python in this way may be a false equivalency in the light of the GPU-accelerated machine learning revolution of the last ten years, and the evolving overlap between data science and machine learning.

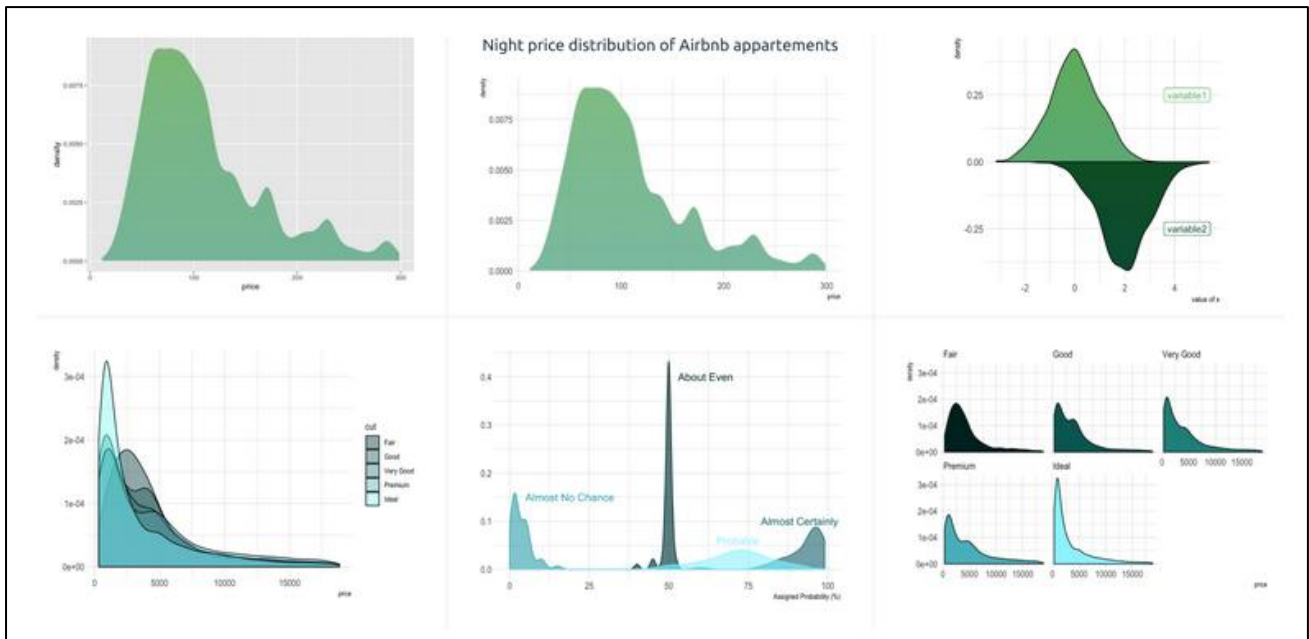
Neither does framing the languages' adherents as 'rivals' account for the interoperability that can be achieved between them, nor the arguable unique strengths of each.

In this article we'll consider what Python and R currently have to offer for the effective development and implementation of deep neural networks, and how emerging trends may affect their uptake.

R for Data Science

R evolved in the mid-1990s from S², a statistical programming language popular in the 1980s as statistical analysis became a driving feature of the new business computing revolution.

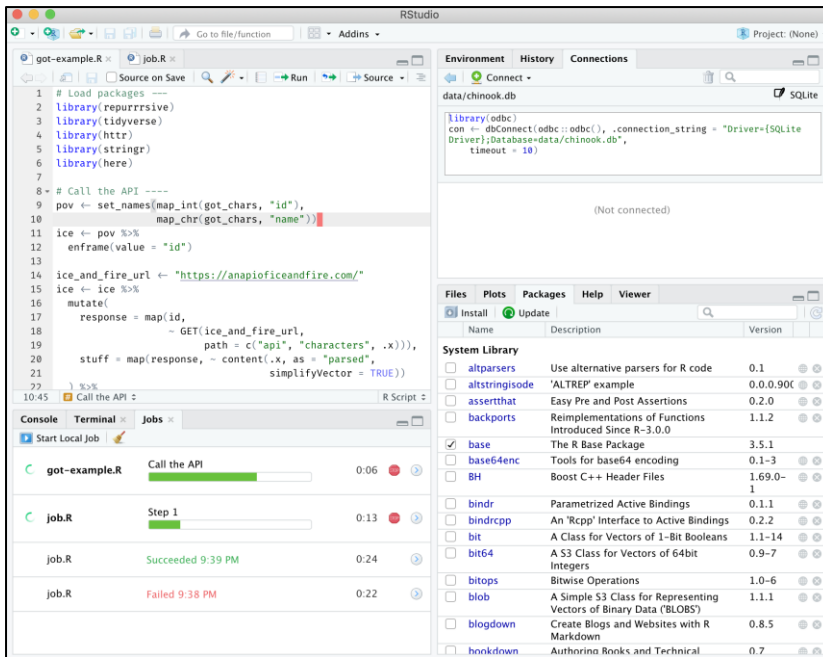
Among its core strengths, R provides very high quality and configurable graphing output even on a base install.



Density plot graphs using the ggplot2 data visualization package in R. Source: <https://www.r-graph-gallery.com/density-plot.html>

R is an extensible language, with more than 20,000³ available user-contributed extensions. Areas covered include finance, genetics, econometrics, medical imaging, machine learning, psychometrics and social sciences, among many others⁴. Packages are archived and distributed from the Comprehensive R Archive Network (CRAN⁵).

In 2010 the release of the RStudio⁶ integrated development environment (IDE) gave (CLI-based) R a more user-friendly but powerful interface, while the Tidyverse⁷ curated collection of R packages for data science became a standard environment for selecting and using the most powerful and popular adjuncts to R.



RStudio running background scripts to avoid modal lock. Source: <https://blog.rstudio.com/2019/03/14/rstudio-1-2-jobs/>

Python for Data Science

Python is an interpreted and dynamically-typed language developed from the ABC⁸ language in the late 1980s. Like R, Python is highly extensible, with its popularity driven by the extraordinary range of additional functionality from over 137,000⁹ available libraries.

Python's ability to create native and cohesive frameworks has driven it to the heart of the machine learning revolution over the last 15 years. Pivotal AI libraries such as TensorFlow, Pandas, Scikit-Learn, NumPy and Keras are either written in Python or support it as a priority over other languages, such as C.

In the same period, a confluence was forming between machine learning and data science, as the power of high-scale deep neural networks transformed a culture of 'local' and limited analysis into the prospect of developing an industrialized, AI-driven panopticon of statistical insight.

Thus, over the last decade, Python's strengths as an analysis platform have met up with its vanguard position in machine learning development, recently pushing the language to first position¹⁰ (57%) on GitHub's survey of machine learning languages, with rival R lagging in eighth place.

Besides those mentioned above, popular Python data science libraries include [PyTorch and Keras](#), as well as staples such as Matplotlib, Seaborn, Pydot (for the C-written GraphViz data visualization library), SciPy and the Beautiful Soup HTML parser, for web-scraping.

The Python Package Index (PyPi) is the central library repository and equivalent to R's CRAN. It's a crowded place, with a wider scope than data science, and therefore PyForest¹¹ offers a filter for the automatic import of PyPi data science libraries.

Interoperability

Running R From Python

The rpy2 repository¹² provides full-featured access to all R functionality from within Python, translating R's objects into Python functions, with transparent conversion to Pandas¹³ and NumPy¹⁴ data structures.

Running Python From R

Reticulate¹⁵ is the most popular method to access Python functionality from within R. Reticulate inserts a Python session directly within an R session, allowing calls to Python by various methods:

- Importing Python modules and gaining direct access to its functions.
- Using the R Markdown¹⁶ language engine as a bi-directional interchange between R and Python.
- Sourcing Python scripts with seamless object/function translation, allowing the same functionality as calling an R script.

Additionally, the RTorch repository¹⁷ provides a wrapper to PyTorch, effectively combining all the capabilities of each language.

Machine Learning and Statistical Libraries Across R and Python

Both R and Python offer a comprehensive array of importable or native libraries covering use cases for data science analysis, from the most popular functions to the most arcane. Since interoperability between the two languages is well-facilitated (see above), there is no appreciable gap in functionality. Some libraries, such as XGBoost, are equally available to either language.

In general Python relies on the machine learning library scikit-learn¹⁸ for the majority of data science-related functionality, while CRAN is the equivalent external source for R. Since R was created with statistical analysis in mind, it features a slightly higher range of related native functions.

	Python	R
Generalized Linear	SKLearn ¹⁹ Statsmodels ²⁰	Native (GLM ²¹)
Linear/Logistic/Poisson Regressions	Statsmodels	Native (base) with parameters
Bayesian	PyStan ²² PyMc ²³ Edward ²⁴	RStan ²⁵
Regularized Regression/ Partial Least Square	SKLearn Statsmodels	GLM Caret ²⁶

Ensembles	SKLearn: RandomForestClassifier ²⁷ XGBoost ²⁸	Native Package ²⁹ (gradient boosted trees) XGBoost ³⁰
Survival	LifeLines ³¹ Statsmodels ³²	Survreg ³³
Decomposition	SKLearn: Decomposition ³⁴	Native (PCA: 'prcom' function) CRAN: ICA ³⁵ (ICA) CRAN: NMF (Nonnegative Matrix Factorization) ³⁶
Kmeans and Hierarchical Clustering	SKLearn: KMeans ³⁷	Native (base)
DBScan (clustering)	SKLearn: DBScan ³⁸	CRAN: DBScan ³⁹
Affinity Propagation (clustering)	SKLearn: AffinityPropagation ⁴⁰	ApCluster ⁴¹
TimeSeries	Statsmodels (ARIMA, VAR, GARCH and exponential smoothing) Prophet ⁴²	CRAN: Tseries ⁴³ Prophet ⁴⁴
Decision Trees	SKLearn: ID3, C4.5, C5.0 and CART ⁴⁵ PyPi: CHAID ⁴⁶	Native (base) CRAN: RPart ⁴⁷ R: CHAID ⁴⁸ CRAN: DataTree ⁴⁹ (ID3) CRAN: Weka ⁵⁰ (J48 / C45)

Graphing

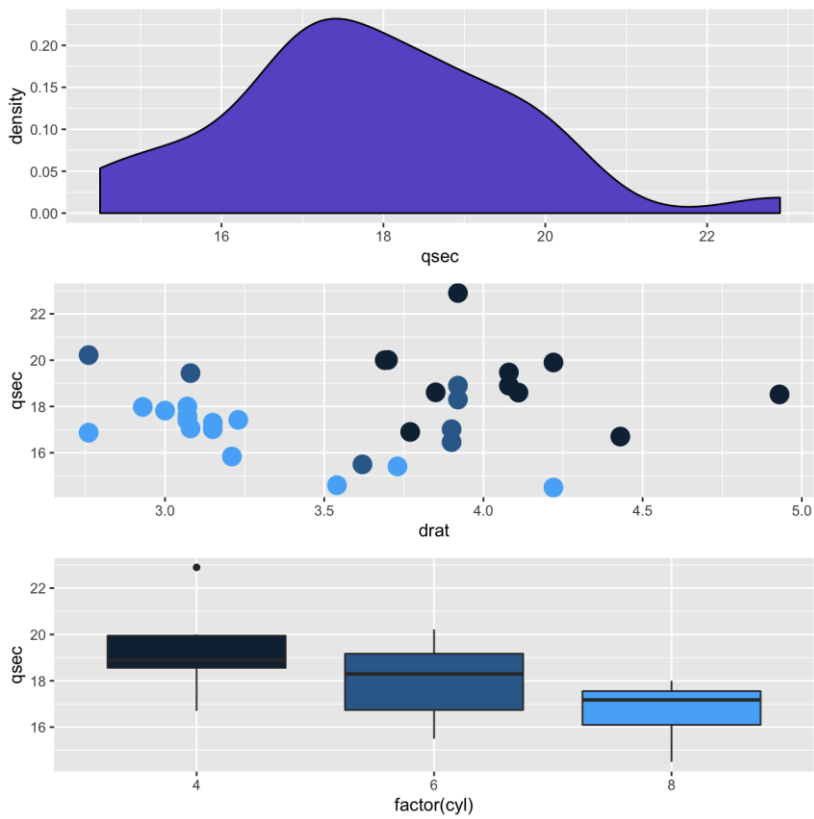
Graphing and Visual Data Exploration in R

Though R's native Graphics package provides 100 functions for generating histograms, scatterplots and boxplots, the far more sophisticated capabilities of the ggplot2^{51,52} package have arguably been the driver for uptake.

Ggplot2 provides a layered and logical sequential approach to data graph development, with aesthetic components chosen and configured after essential elements such as axes and data positioning are established.

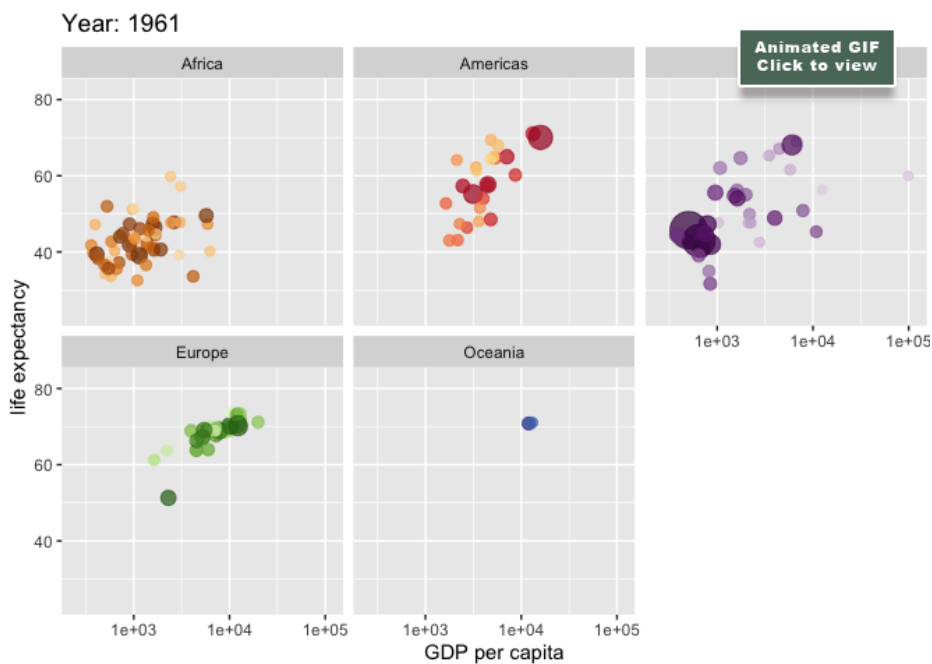
At the time of writing there are 80 registered extensions⁵³ for ggplot2, offering functionality as diverse as spellchecking, theme additions, HTML embedding, facility for survival curves and alluvial diagrams, time series visualizations, partial rasterizations and visualization of IP addresses and networks.

In addition to a versatile range of styles and configurations, graphs can be nested for single implementation into a document or interface:



The `arrangeGrob()` function in `ggplot2` allows for multiple nested charts within the same illustration. Source: <https://www.r-graph-gallery.com/261-multiple-graphs-on-same-page.html>

The `ganimate` library⁵⁴ can also compile sequential or time data into elegant animations, created as GIF files or passed to the FFMPEG renderer for video output.



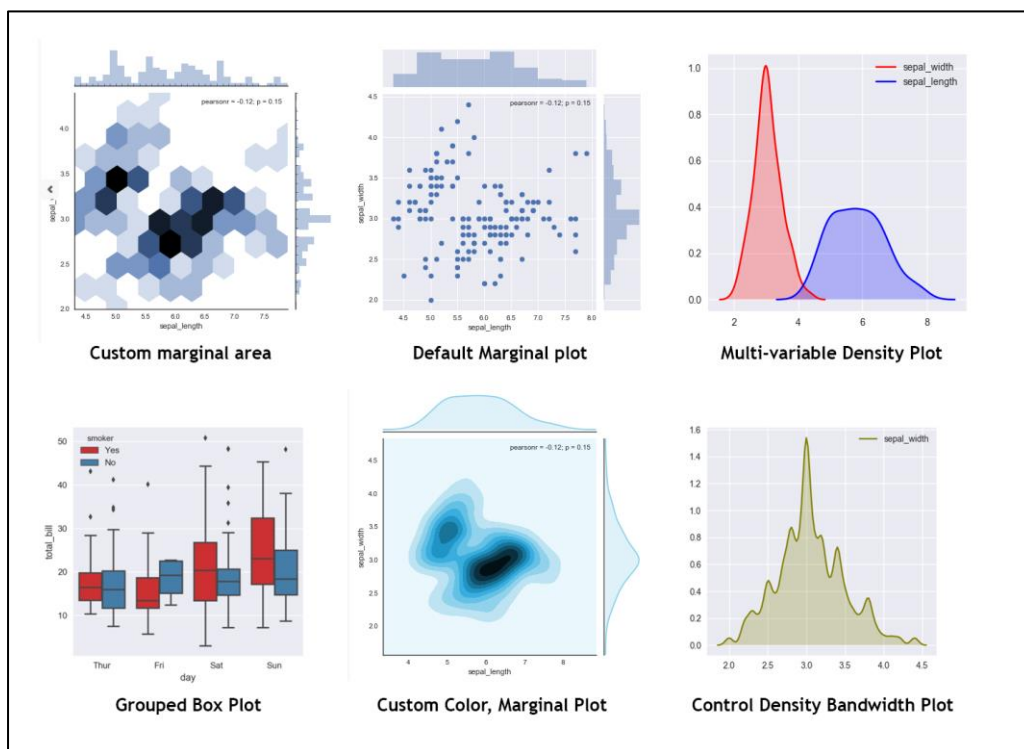
A `transition_time()` function rendered mobile in `ganimate`. Source: <https://github.com/thomasp85/ganimate/>

Graphing and Visual Data Exploration in Python

In keeping with its generalized scope as an all-purpose programming language, Python defers to external libraries such as Matplotlib⁵⁵ to provide equivalent graphing functionality to R.

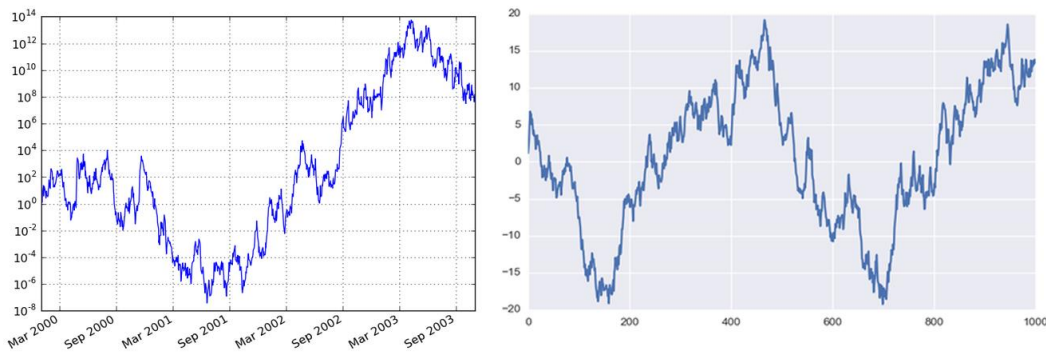
Matplotlib is an open source Python derivation of the proprietary MATLAB⁵⁶ programming environment. Similar to R's ggplot2 extensions repository (see above), it is greatly enhanced by a comprehensive range of third-party mapping toolkits⁵⁷.

Additional functionality available via these libraries includes the superimposition of data on map projections; facility for interactive data cursors; support for cross-platform visualization applications; enhanced table support (compared to matplotlib's native tables); support for annotated DNA and astronomy data maps; ridge maps; and non-contiguous graph axes.



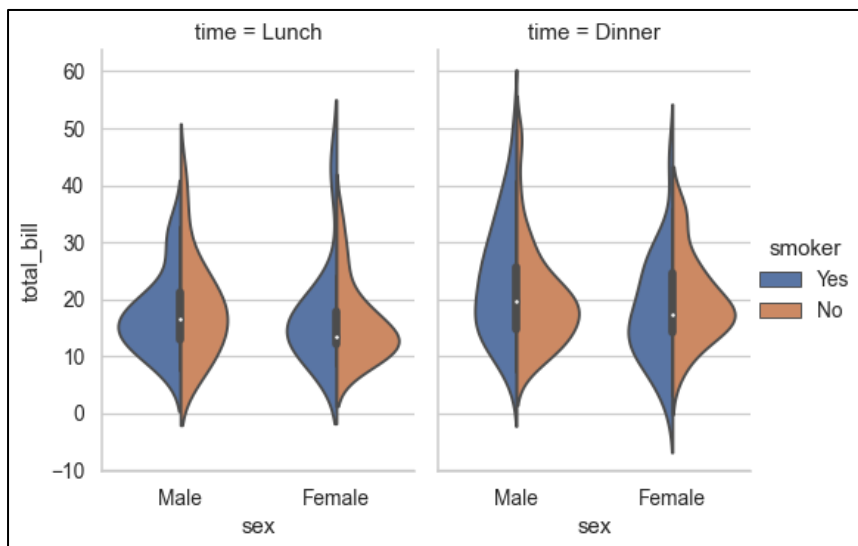
Graphs with varying functionality, obtained via Seaborn. Source: <https://python-graph-gallery.com/seaborn/#prettyPhoto>

However, arguably the most important additions are the ggplot library⁵⁸, which ports the entire graphics-driven philosophy and utility of R's ggplot2 to native Python; and Seaborn⁵⁹, which not only facilitates Trellis graphs (see below), but provides a super-layer of statistical graphics capability to Matplotlib via a high-level API that's easier to use and more up-to-date.



Data imported with no extra styling via Matplotlib (left) and Seaborn. Source: <https://www.quora.com/What-is-the-difference-between-Matplotlib-and-Seaborn-Which-one-should-I-learn-for-studying-data-science>

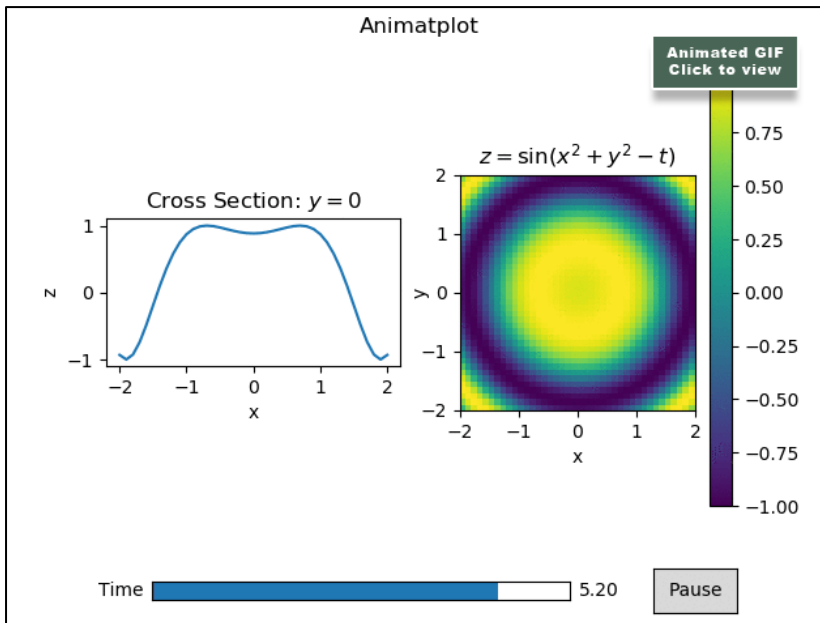
Seaborn offers greater control over color use in graphs, essential in differentiating data streams, as well as offering additional plot types, such as the Violin plot⁶⁰:



Source: <https://seaborn.pydata.org/generated/seaborn.violinplot.html>

Furthermore, Seaborn's pairplot() function⁶¹, among others, enables at-a-glance exploration of data prior to more sophisticated choices for graph output, and with less code than is necessary in R⁶²:

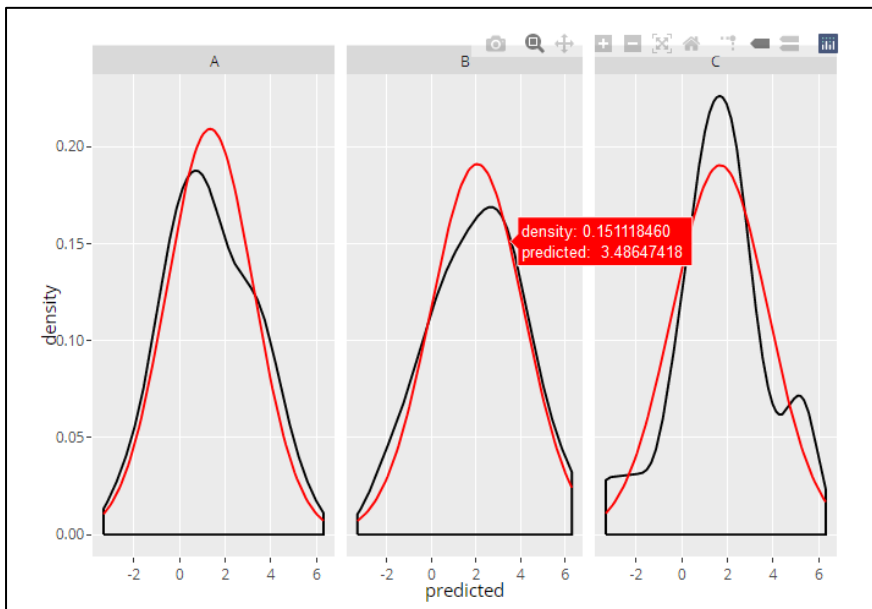
Animation is also catered for via the Animatplot^{63,64} library, an extension for Matplotlib, as well as being natively supported in Matplotlib itself^{65,66}.



Animated graphs via Animatplot. Source: <https://animatplot.readthedocs.io/en/stable/index.html>

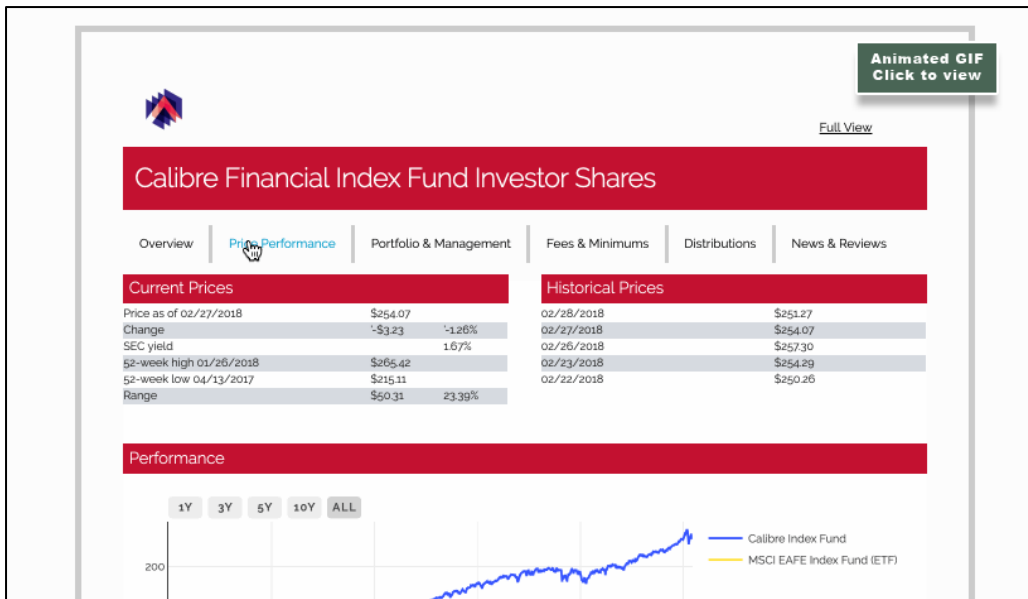
Plotly for R and Python

Canadian software company Plotly has developed various data visualization libraries for Python and R, including the Dash Python⁶⁷ and R (DashR⁶⁸) frameworks for the creation of interactive dashboards and graphics.



A multiple density plot created with the Plotly ggplot2 graphing library, featuring interactive on-mouseover tooltip. Source: https://plotly.com/ggplot2/geom_density/

Implementations of Plotly are derived from the plotly.js JavaScript repository⁶⁹, which can output still or animated raster graphics in addition to resolution-independent vector-based SVG files, allowing for more sophisticated and interactive applications, such as PDF-style dynamic documents:

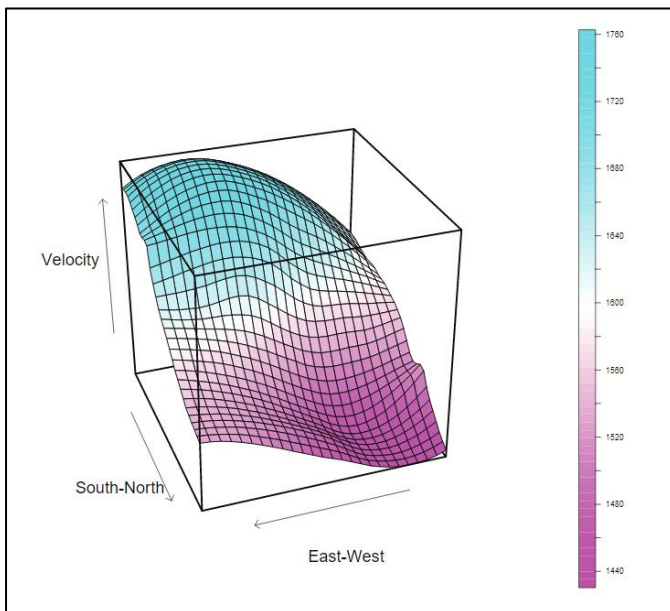


A vector-based, resolution-independent PDF-style interactive document developed from Plotly.js. Source: <https://github.com/plotly/plotly.js>

In addition to Python and R implementations, Plotly also supports the Julia programming language.

Trellis Graphs for R and Python

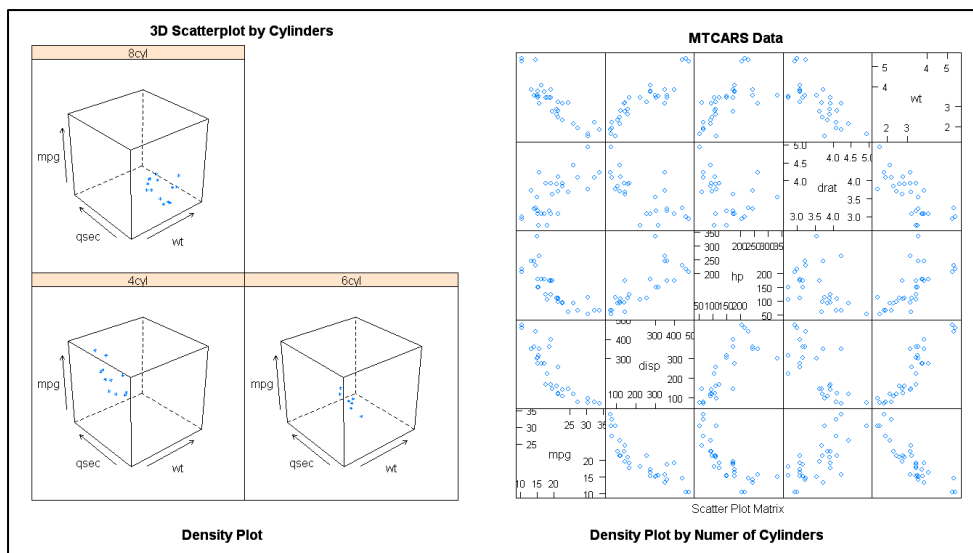
Trellis graphs⁷⁰ visualize the relationship between any number of variables, represented as a rectangular array of plots (histograms, box or scatter).



A Trellis graph manifested as a color wireframe plot.

Source: <http://ml.stat.purdue.edu/stat695t/writings/Trellis.User.pdf>

R's Lattice library^{71,72} is a high-level visualization resource specializing in multivariate data, and enabling the creation of highly stylized Trellis graphs, including impactful 3D wireframe graphs⁷³.



Two of the four methods available in Lattice. Source: <https://www.statmethods.net/advgraphs/trellis.html>

Available methods include bar chart, boxplot, 3D scatterplot, 3D contour plot, histogram, kernel density plot, 3D level plot, scatterplot matrix, parallel coordinates plot, strip plot, dotplot, scatterplot and 3D wireframe graph.

Besides interoperable builds where Python accesses R's Lattice library transparently (probably the easiest method), there are various ways to generate Trellis graphs in Python, including the LatticeDrawing repository⁷⁴ or PyPi's PythonLattice⁷⁵; accessing R's RPlot via Pandas⁷⁶; using Plotly's subplot capabilities⁷⁷; and using the FacetGrid class⁷⁸ in Seaborn.

Accessibility

R's tight focus on local data science analysis means that all tooling and all help resources will be in some way pertinent to the task, whereas the signal-to-noise ratio of the Python ecostructure can be an initial obstacle when developing dependencies in a new framework.

However, it's generally acknowledged that R has a steeper initial learning curve than Python, and is less widely applicable to other environments.

Python's object-oriented approach is less specious than R, and its learning curve shallower. Additionally, knowledge of other programming languages with a similar paradigm will speed the progress of a Python initiate substantially, but may be less useful when learning R.

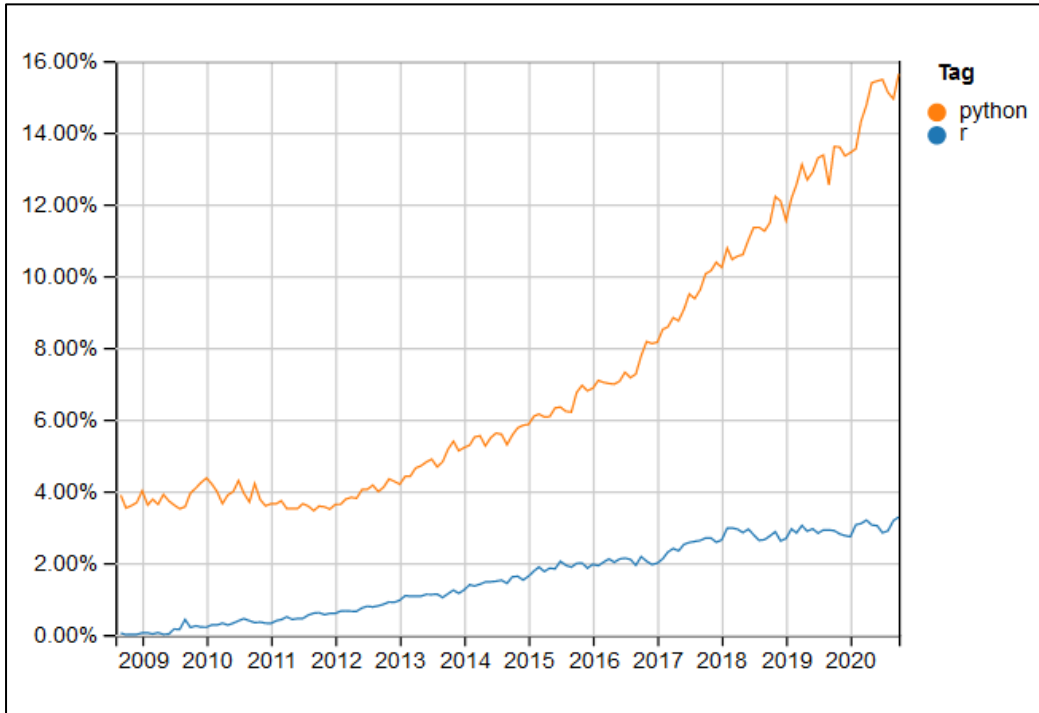
In terms of code complexity, there is very little difference⁷⁹ between R and Python.

Diffusion

We saw earlier that Python is the most popular machine learning language. According⁸⁰ to the Stack Overflow 2020 Developer Survey, Python maintains its 2019 position as the fourth most popular general programming language among professional developers at 41% share, while R drops one place to 17th position, with a 5.5% share.

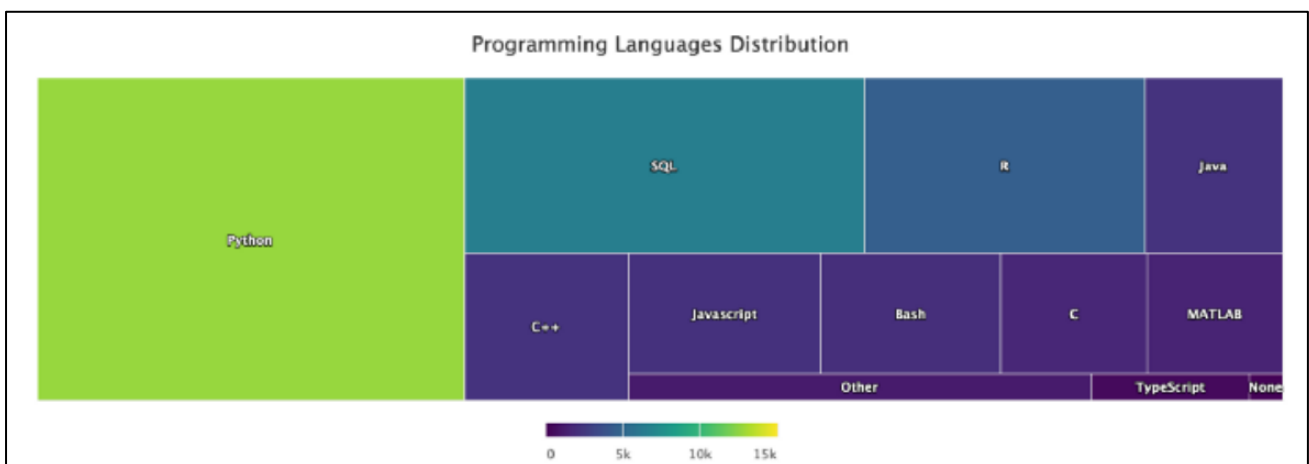
However, it should be considered that R is a dedicated data science analysis language and Python a heterogeneous programming language with wider scope, particularly in machine learning — a later consideration in the development of R.

Nonetheless the ascending trajectories of the two languages over ten years is a broad indicator of talent availability and market share in favor of Python:



Source: <https://insights.stackoverflow.com/trends?tags=r%2Cpython>

A recent Kaggle survey indicates⁸¹ that Python comfortably dominates R in terms of adoption for deep learning, covering the majority of work in Computer Vision and Natural Language Processing (NLP), with wider diffusion and uptake in the data science community.



Source: <https://towardsdatascience.com/python-vs-r-for-data-science-6a83e4541000>

R vs Python: Key Points for Comparison

Since, as we have seen, there is no notable shortfall in functionality for data science analysis between Python and R, there are only a limited number of convincing reasons to choose R:

- It has a shallower initial learning curve for pure data science projects.
- The Shiny R package⁸² facilitates the easy creation of interactive web applications and dashboards, and is a superior and more mature product than Python's Dash⁸³ equivalent.
- R uses more native functions, as opposed to Python's classes and secondary libraries.
- R offers a large ecostructure focused on data science, and a larger number of libraries than Python for more marginal use cases.
- RStudio is arguably the most mature and complete IDE for data science analysis.
- R has a clearer and more usable versioning system than Python, with much less risk of technical debt⁸⁴ or the need to support multiple installations or instances.
- Any pre-existing data science talent on your team may already be familiar with it.
- Suitable in cases where the ultimate scope and breadth of the project are each a known quantity, and affordable turnkey Python-based solutions are not going to be needed later when the scope develops or the data scales up significantly.
- R has good market share and loyal entrenchment in specific sectors of the scientific community, with comprehensive and focused help resources.

Conversely, there are stronger indicators for choosing Python over R:

- Clear market saturation makes Python software development talent more affordable⁸⁵ and available than for R.
- Python is a general modular programming language, meaning that infrastructure and deployments can take place in the same popular development environment as the core code.
- Superior package management, including dependency management, with Pip, as opposed to Packrat and other R package managers⁸⁶.
- Easier development⁸⁷ of high-level C/C++ interfaces.
- Python has demonstrated that it has a sufficient community will to maintain or exceed feature parity with R, while the converse is likely to remain logistically difficult for R in the future.
- A high likelihood that a company's technical teams already have some grounding in Python, even for purposes unrelated to data science.
- Faster performance in typical machine learning and general analytical tasks^{88,89}.
- Core Python-facing libraries such as Matplotlib, Scikit-learn, Keras, SciPy, NumPy and Pandas have grown massively in well-funded industry support and community impetus in the age of GPU-driven machine learning, assuring future support and updates, and fewer brittle or precarious interdependencies in future projects.
- Wider variety of mature and well-supported IDEs^{90,91}, including RStudio itself, since 2018⁹².

Conclusion

We began this article by considering whether Python or R might become an outmoded language in data science, or whether the two languages will maintain their current equilibrium. With feature parity between Python and R now a moot point, it can be argued that R will remain on the back foot, and is indeed becoming a 'legacy' approach when considered against the ascendancy of Python in data science analysis, and that R's

specialization in this field is no longer a compelling reason to prefer it.

-
- ¹ <https://theconversation.com/how-betamax-bit-the-dust-and-other-tales-of-forgotten-tech-50619>
 - ² https://archive.org/details/sinteractiveenvi00beck_0
 - ³ <https://www.rdocumentation.org/>
 - ⁴ <https://cran.r-project.org/web/views/>
 - ⁵ <https://cran.r-project.org/>
 - ⁶ <https://rstudio.com/>
 - ⁷ <https://www.tidyverse.org/>
 - ⁸ <https://homepages.cwi.nl/~steven/abc/>
 - ⁹ <https://www.mygreatlearning.com/blog/open-source-python-libraries/>
 - ¹⁰ <https://github.blog/2019-01-24-the-state-of-the-octoverse-machine-learning/>
 - ¹¹ <https://pypi.org/project/pyforest/>
 - ¹² <https://pypi.org/project/rpy2/>
 - ¹³ <https://rpy2.github.io/doc/v3.3.x/html/pandas.html>
 - ¹⁴ <https://rpy2.github.io/doc/v3.3.x/html/numpy.html>
 - ¹⁵ <https://rstudio.github.io/reticulate/>
 - ¹⁶ <https://rmarkdown.rstudio.com/>
 - ¹⁷ <https://f0nzie.github.io/rTorch/>
 - ¹⁸ <https://scikit-learn.org/stable/>
 - ¹⁹ https://scikit-learn.org/stable/modules/linear_model.html
 - ²⁰ <https://www.statsmodels.org/stable/index.html>
 - ²¹ <https://www.rdocumentation.org/packages/stats/versions/3.6.1/topics/glm>
 - ²² <https://pystan.readthedocs.io/en/latest/>
 - ²³ <https://docs.pymc.io/>
 - ²⁴ <http://edwardlib.org/>
 - ²⁵ <https://mc-stan.org/users/interfaces/rstan>
 - ²⁶ <https://topepo.github.io/caret/>
 - ²⁷ <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>
 - ²⁸ https://xgboost.readthedocs.io/en/latest/python/python_intro.html
 - ²⁹ <https://www.rdocumentation.org/packages/randomForest/versions/4.6-14/topics/randomForest>
 - ³⁰ <https://xgboost.readthedocs.io/en/latest/R-package/xgboostPresentation.html>
 - ³¹ <https://lifelines.readthedocs.io/en/latest/>
 - ³² <https://www.statsmodels.org/stable/duration.html#regression-methods>
 - ³³ <https://www.rdocumentation.org/packages/survival/versions/2.44-1.1/topics/Survreg>
 - ³⁴ <https://scikit-learn.org/stable/modules/classes.html#module-sklearn.decomposition>
 - ³⁵ <https://cran.r-project.org/web/packages/ica/index.html>
 - ³⁶ <https://cran.r-project.org/web/packages/NMF/index.html>
 - ³⁷ <https://scikit-learn.org/stable/modules/clustering.html#k-means>
 - ³⁸ <https://scikit-learn.org/stable/modules/generated/sklearn.cluster.DBSCAN.html#sklearn.cluster.DBSCAN>
 - ³⁹ <https://cran.r-project.org/web/packages/dbscan/index.html>
 - ⁴⁰ <https://scikit-learn.org/stable/modules/clustering.html#affinity-propagation>
 - ⁴¹ <https://cran.r-project.org/web/packages/apcluster/index.html>
 - ⁴² <https://facebook.github.io/prophet/>
 - ⁴³ <https://cran.r-project.org/web/packages/tseries/index.html>
 - ⁴⁴ <https://facebook.github.io/prophet/>
 - ⁴⁵ <https://scikit-learn.org/stable/modules/tree.html#tree-algorithms-id3-c4-5-c5-0-and-cart>
 - ⁴⁶ <https://pypi.org/project/CHAID/>
 - ⁴⁷ <https://cran.r-project.org/web/packages/rpart/index.html>
 - ⁴⁸ https://r-forge.r-project.org/R/?group_id=343
 - ⁴⁹ <https://cran.r-project.org/web/packages/data.tree/index.html>
 - ⁵⁰ <https://cran.r-project.org/web/packages/RWeka/index.html>
 - ⁵¹ <https://ggplot2.tidyverse.org/>
 - ⁵² <https://cran.r-project.org/web/packages/ggplot2/index.html>
 - ⁵³ <https://exts.ggplot2.tidyverse.org/gallery/>
 - ⁵⁴ <https://cloud.r-project.org/web/packages/gganimate/index.html>
 - ⁵⁵ <https://matplotlib.org/>
 - ⁵⁶ <https://www.mathworks.com/products/matlab.html>
 - ⁵⁷ <https://matplotlib.org/thirdpartypackages/index.html>
 - ⁵⁸ <https://github.com/yhat/ggpy>
 - ⁵⁹ <https://seaborn.pydata.org/>
 - ⁶⁰ <https://seaborn.pydata.org/generated/seaborn.violinplot.html>
 - ⁶¹ <https://seaborn.pydata.org/generated/seaborn.pairplot.html>

- ⁶² <https://www.inwt-statistics.com/read-blog/data-visualization-R-versus-python.html>
- ⁶³ <https://github.com/t-makaro/animatplot>
- ⁶⁴ <https://animatplot.readthedocs.io/en/stable/index.html>
- ⁶⁵ <https://towardsdatascience.com/how-to-create-animated-graphs-in-python-bb619cc2dec1>
- ⁶⁶ <https://jingwen-z.github.io/draw-animated-graphs-with-matplotlib/>
- ⁶⁷ <https://github.com/plotly/dash>
- ⁶⁸ <https://github.com/plotly/dashR>
- ⁶⁹ <https://github.com/plotly/plotly.js>
- ⁷⁰ https://web.njit.edu/~dhar/mth664/trellis_gph_cmmd.pdf
- ⁷¹ <https://cran.r-project.org/web/packages/lattice/index.html>
- ⁷² <https://cran.r-project.org/web/packages/lattice/lattice.pdf>
- ⁷³ <https://www.statmethods.net/advgraphs/trellis.html>
- ⁷⁴ <https://github.com/coldlaugh/LatticeDrawing>
- ⁷⁵ <https://pypi.org/project/python-lattice/>
- ⁷⁶ <https://pandas.pydata.org/pandas-docs/version/0.14.1/rplot.html>
- ⁷⁷ <https://plotly.com/python/facet-plots/>
- ⁷⁸ <https://seaborn.pydata.org/generated/seaborn.FacetGrid.html>
- ⁷⁹ <https://www.dataquest.io/blog/python-vs-r/>
- ⁸⁰ <https://insights.stackoverflow.com/survey/2020#technology-programming-scripting-and-markup-languages-professional-developers>
- ⁸¹ <https://towardsdatascience.com/kaggle-user-survey-2019-326e187ff207>
- ⁸² <https://shiny.rstudio.com/>
- ⁸³ <https://plotly.com/dash/>
- ⁸⁴ <https://blog.ephorie.de/why-r-for-data-science-and-not-python>
- ⁸⁵ <https://michaeltOTH.me/r-programmers-earn-more-than-python-programmers.html>
- ⁸⁶ <https://ellisvalentiner.com/post/dpendency-hell/>
- ⁸⁷ <https://docs.python.org/3.1/extending/extending.html>
- ⁸⁸ <https://julialang.org/benchmarks/>
- ⁸⁹ <https://datascienceplus.com/loops-in-r-and-python-who-is-faster/>
- ⁹⁰ <https://realpython.com/python-ides-code-editors-guide/>
- ⁹¹ <https://datascience.stackexchange.com/questions/5345/ide-alternatives-for-r-programming-rstudio-intellij-idea-eclipse-visual-stud>
- ⁹² <https://blog.rstudio.com/2018/10/09/rstudio-1-2-preview-reticulated-python/>