# Data Science Notes

Barclays Research's Data & Investment Sciences blog, where we dive into some of the frequently asked questions about data science and how it is reshaping investing.

**Data Science**
Adam Kelleher

**BARCLAYS**

26 June 2020

## Data Science 101

We frequently field the question of how investors should build a data science capability. In this entry, we share the framework we have used for making business decisions on data science.

### Below are the key dimensions to consider:

- **People.** Hiring data scientists is different from hiring investment professionals, and there is a wide range of levels of domain expertise, engineering skill, and quantitative acumen. Still, most investment firms employ plenty of people with the skill and aptitude needed to develop into data scientists. Given the freedom to pick up computer science skills, these people can grow to become data scientists.

- **Infrastructure.** The systems and facilities needed for data science start with advanced versions of what most investors are already using in their IT (with some entirely new capabilities at the very top end), but generally some spending will be required for automated data ingestion and analysis.

- **Data.** Most investors could add some data science capabilities by just doing new things with data they already own. At the more sophisticated (and expensive) end, the data requirements can be an entirely new discipline that requires a dramatically different set of capabilities across the organization, including data scientists, IT engineers, and legal and compliance expertise.

## How companies choose to address each one of these dimensions depends on answering the following questions:

- What are you trying to get out of your data science capability, and what is the value to your business?

- What capabilities do you already have?

- What can you reasonably accomplish; put another way, how much are you willing to spend?

This is a high-dimensionality problem, so there are innumerable ways to assemble staff, infrastructure, and data to address the three questions. That said, we think there are a few focal point answers: Figure 1 presents three reasonable options that cover the spectrum of how investors can approach the problem from a resourcing perspective.

It is more efficient to use some data science terms of art in these descriptions, so we define some key ones at the outset. These definitions reflect how we use them in our writing and may differ from how they are used elsewhere. **Data science** is statistics implemented using a heavy dose of computer science (or sometimes statistics as rediscovered by computer scientists). **Data scientists** are its researchers and practitioners (and usually have come to it from a spectrum between pure computer science and pure statistics). **Python** is a scripting language that is popular for data scientists, but one can substitute it for any similar language – R, matlab, etc. Jupyter notebooks are a popular software package that allows one to write Python code and see the results it produces as one goes. **Spark** is a framework for large-scale computing (allocating work across many computers) and can be used within Python. **Production-grade** processes are when computers are able to complete a task beginning to end with no human intervention (unless something goes wrong). **Development** is the process of building production processes, and **ad-hoc analyses** are when humans use data science tools to complete one-off tasks such as answering a question.

FIGURE 1

**Our view of some reasonable options for starting a Data Science team at different scales of investment**

| Scale | People | Infrastructure | Data | Direct Costs | Positives and Considerations |
|---|---|---|---|---|---|
| **Efficient**<br><br>Can do a lot with limited resources | • Assign existing analysts to integrate some alternative data into the research process | • Bring pre-packaged alternative data signals into Excel<br><br>• Do some work with Python in Jupyter notebooks with larger datasets<br><br>• Engineering requirements are limited to building and maintaining data stores, with limited "productionization" of data products. | • Machine-readable foundational data (prices, corporate fundamentals, etc)<br><br>• Use pre-packaged signals and APIs | • $0 for people, but some lost productivity as they ramp up<br><br>• $10s of k for database servers, plus some IT support<br><br>• <$100k (often much less) per dataset | + Cheap and quick to get running<br><br>+ Making "table stakes" investments in alt data<br><br>+ Integrating organically into culture<br><br>+ Fast delivery<br><br>- Limited scope for competitive advantage or marketing potential<br><br>- Reliance on vendors for data processing |

| Scale | People | Infrastructure | Data | Direct Costs | Positives and Considerations |
|---|---|---|---|---|---|
| **Innovative**<br><br>Some real competitive advantage with a few more resources | • Hire a dedicated data scientist to support investment and research teams<br><br>• Allocate dedicated IT resources (people and computers) to production tasks | • Work primarily with Python in Jupyter notebooks<br><br>• Build and maintain data stores, with some "production-ization" of data products.<br><br>• Add some specialty IT support or cloud services for production engineering<br><br>• Can use vendors for many intermediate tools | • Machine-readable foundational data (prices, corporate fundamentals, etc)<br><br>• Make better use of pre-packaged signals and APIs<br><br>• Purchase some smaller-scale raw sources of data | • A budget for data scientist and IT salaries comparable to your data budget~$[100k] for servers and cloud tools<br><br>• Typically want 1-3 people<br><br>• $50-100k per dataset for sector specific data, or $100k->$1mn for each general purpose datasets (transactions, geolocation, text) | + You can build some of your own custom insights and signals that nobody else might have<br><br>+ You have a data science team, so there is an opportunity for differentiation and marketing<br><br>+ Still lean and pretty fast<br><br>- Team will have limited capacity<br><br>- Capabilities depend a lot on who is hired as a data scientist; the best ones are more expensive<br><br>- Success will depend on integrating a new capability into the team culture |
| **Resourced**<br><br>Hire a full team, and buy new data to do truly novel work | • Hire a data scientist leader, who builds out a full team<br><br>• Multiple dedicated IT resources to support production tasks | • Python in Jupyter notebooks<br><br>• Spark on clusters<br><br>• Full IT build of servers, computing clusters, etc; or full IT support for cloud services, possibly including new service design patterns and large-scale compute models | • Machine readable foundational data (prices, corporate fundamentals, etc)<br><br>• Consume significant raw data, including the development of custom production algorithms needed to transform data to insights | • A budget for data scientist and IT salaries comparable to your data budget At least $100s of k for servers and cloud tools<br><br>• 3 people likely the practical minimum<br><br>• At least $1mn for multiple general use + specific use datasets<br><br>• Additional need for operational support to build and manage the business cases for different types of investments | + Potential for novel products that are truly differentiating<br><br>+ Can build entirely novel datasets<br><br>+ This is a real investment area now, so real opportunity for differentiation and marketing<br><br>- Spending a lot on it, which is hard to justify if results are not produced<br><br>- Scale, complexity, and difficulty of hiring can result in a slow ramp<br><br>- A substantial amount of what is produced will be invisible to customers<br><br>- More potential for culture friction<br><br>- When things break, you have to fix them yourself. - The more systems you have, the higher the cost – and it grows over time! |

Source: Barclays Research