

# **Data Science in Focus**

Andrew Stewart

2025-04-12

# Table of contents

<b>Data Science In Focus</b>	<b>4</b>
<b>Introduction</b>	<b>5</b>
<b>Defining Data Science</b>	<b>6</b>
<b>What Is Data Science?</b>	<b>7</b>
A Discipline of Inquiry, Not Output . . . . .	7
The Objects of Study . . . . .	7
The Tools Are Not the Discipline . . . . .	8
Science Inside a Software Org . . . . .	8
In Summary . . . . .	8
<b>A Brief History of Data Science</b>	<b>10</b>
<b>Where Data Science Fits in Technology Organizations</b>	<b>11</b>
<b>Data Science as Applied Systems Science</b>	<b>12</b>
<b>The Data Scientists</b>	<b>13</b>
<b>A Refined Venn Diagram of Data Science</b>	<b>14</b>
<b>The T-Model of Skill Development</b>	<b>15</b>
<b>Data Science as a Team Sport</b>	<b>16</b>
<b>Research Management and Process</b>	<b>17</b>
<b>Managing Data Science Through Research Programs</b>	<b>18</b>
<b>The Scientific Method in Data Science</b>	<b>19</b>
<b>Using the Capability Maturity Model in Data Science</b>	<b>20</b>

<b>Tooling and Lab Practice</b>	<b>21</b>
<b>The Data Scientist's Lab Workbench</b>	<b>22</b>
<b>Documentation and Knowledge Repositories</b>	<b>23</b>
<b>Statistical Thinking</b>	<b>24</b>
<b>Probability and Statistical Inference</b>	<b>25</b>

# Data Science In Focus

# Introduction

**Data Science In Focus** is a reflective essay series on data science as a scientific discipline—what it is, how it’s practiced, and what it has become. As the field matures beyond its early hype cycles and into a coherent form of applied research, this series aims to sharpen our collective understanding of the work itself.

Where a bountiful collection of earlier works have laid the groundwork for a newly forming field, these essays here revisit core questions with the benefit of hindsight:

- What does it mean to practice data science as science?
- How should teams, tools, and systems support inquiry over output?
- What kind of knowledge does data science produce—and for whom?

Rooted in the scientific method, structured around the research lifecycle, and steeped in the evolving norms of modern tech orgs, this series puts the discipline itself into focus.

# **Defining Data Science**

# What Is Data Science?

In the era of tech monoculture, the term *data science* has been stretched to near incoherence—absorbing everything from analytics engineering to AI research under its inflated halo. But if we strip away the branding and job title inflation, what remains is something much older, much simpler, and much more principled: **data science is the application of the scientific method to the study of data-generating systems.**

It is not a subfield of software engineering. It is not a synonym for machine learning. It is not a placeholder for “person who works with data.” Data science, properly understood, is a **scientific discipline**—defined not by its tooling or domain, but by its epistemology. Its goal is to generate knowledge. Its process is experimental. Its currency is uncertainty. And its outputs are not products, but explanations.

## A Discipline of Inquiry, Not Output

What distinguishes data science from engineering is not the data—it’s the **orientation toward inquiry**.

- **Engineers** build systems that are designed to perform reliably and at scale.
- **Scientists** study systems to understand how and why they behave the way they do.

Data scientists may use engineering tools, work within engineering organizations, and produce artifacts that feed into engineering systems. But their foundational job is to ask—and answer—questions about system behavior. They design experiments, test hypotheses, analyze variation, and build explanatory models. In this sense, a data scientist is closer to a physicist studying turbulence than a developer deploying a feature.

This isn’t a hierarchy. It’s a **division of labor**—and misunderstanding it leads to broken workflows, misaligned expectations, and org charts that burn out good scientists by asking them to write production code full-time.

## The Objects of Study

Data science is concerned with **data-generating processes**, especially those that arise within technological systems. Some examples:

- How does user engagement change in response to a new design?
- What latent behaviors drive churn in a subscription model?
- Why did model performance degrade last week?
- What features of a marketplace system produce price instability?

These are not engineering problems. They're **systems questions**. Answering them requires conceptual models, uncertainty quantification, domain awareness, and often a blend of statistical inference and simulation. They also often involve dead ends, ambiguous results, and theoretical exploration—things that are normal in science but foreign to many software workflows.

## The Tools Are Not the Discipline

It's tempting to define data science by its stack: SQL, Python, pandas, Jupyter, etc. But that would be like defining chemistry by beakers and Bunsen burners. Tools enable the work—they aren't the work.

In fact, many of the tools used in data science are borrowed from engineering or software development. The difference is in **how they're used**. A data scientist doesn't write Python to deploy services; they use it to simulate a hypothesis, analyze system output, or validate statistical assumptions. SQL isn't a pipeline—it's a telescope.

## Science Inside a Software Org

One of the greatest challenges facing data scientists today is that they are often the only scientists inside engineering organizations. That creates cultural friction. Deadlines prioritize shipping over understanding. Metrics are flattened into KPIs. Curiosity becomes a liability. Documentation is seen as overhead rather than intellectual scaffolding.

But despite these tensions, data science has a crucial role to play: it helps organizations **understand themselves**. It maps the terrain, exposes the mechanisms, and builds the mental models that engineering and product teams rely on to make informed decisions.

When practiced as science, data science becomes the **epistemic engine** of a tech company. It gives us confidence not just in what we're building, but in what we believe.

## In Summary

- Data science is a **scientific discipline** rooted in the study of complex systems through data.



- Its central purpose is **explanation**, not output.
- Its methods are driven by **hypothesis, experimentation, and inference**.
- Its work supports and complements engineering by providing **clarity, context, and insight**.

By treating data science as science, we restore its rightful posture—an experimental partner to engineering, a conceptual partner to product, and a critical lens for understanding the systems we build and inhabit.

# A Brief History of Data Science

Data science didn't emerge from a single lineage—it was born from a convergence. Statistics, scientific computing, database management, and machine learning all played formative roles in shaping the field. During the early 2000s, as computational infrastructure and open-source tools proliferated, organizations began to realize that “data” wasn't just a byproduct of digital systems—it was a source of insight.

The rise of the “data scientist” as a role came partly from pragmatism: organizations needed generalists who could both analyze and code. But this era was also defined by confusion. Data science became a catch-all title encompassing analysts, engineers, ML researchers, and more. The result was a temporary inflation of the field—a golden age of generalists with ambiguous scope.

Now, a decade later, the landscape is maturing. Specialized engineering roles have splintered off. Machine learning engineering, analytics engineering, and decision science have defined clearer scopes. What remains in the heart of “data science” is its original essence: a scientific discipline grounded in inquiry, system behavior, and the pursuit of explanatory knowledge.

# Where Data Science Fits in Technology Organizations

Data science occupies a unique niche within technology orgs. It bridges engineering, product, and strategy—but it is fundamentally distinct from all three. Its core deliverable is understanding, not execution.

The organizational structure of data science varies:

- **Centralized teams** emphasize consistency and shared standards.
- **Embedded models** prioritize domain intimacy and responsiveness.
- **Hybrid approaches** attempt to balance autonomy and alignment.

Tensions often arise when data science is treated as a service function, expected to deliver dashboards or one-off analyses on demand. But a scientific function thrives on longer-term research questions, context-rich collaboration, and space to explore uncertainty. Successful teams recognize this and give data scientists both embedded partnerships and protected research capacity.

The most effective data science orgs invest in career ladders, program-based workstreams, and a clear cultural distinction between scientists and engineers—while encouraging tight collaboration across them.

# Data Science as Applied Systems Science

At its core, data science is the study of complex, adaptive systems. These systems—markets, platforms, recommendation engines, networks—are not static. They evolve, respond to feedback, and often exhibit emergent behavior.

To understand such systems, data scientists borrow heavily from adjacent fields:

- **Control theory** to manage dynamic processes
- **Cybernetics** to study feedback loops
- **Complexity science** to model emergence
- **Information theory** to reason about signal and noise

This systems-thinking mindset distinguishes data science from business analytics or product instrumentation. It frames metrics as proxies, not truths. It treats models as lenses, not deliverables. And it embraces the recursive nature of systems—where measurement affects behavior, and knowledge must continually update in response.

# **The Data Scientists**

# A Refined Venn Diagram of Data Science

Earlier attempts to define data science visually—such as Conway’s Venn diagram—centered on hacking skills, statistics, and domain knowledge. While influential, that framing is now outdated.

A more precise picture situates data science at the intersection of:

- **Statistical modeling** — tools for inference, uncertainty, and causality
- **Scientific computing** — numerical methods, simulations, and computation
- **Systems research** — studying complex, interactive software environments

This triad places data science alongside fields like econometrics, computational physics, and quantitative social science. It emphasizes rigor and replicability. And it highlights that data science is not just about *using* data, but about *understanding the systems* that produce it.

# The T-Model of Skill Development

Data scientists develop along a T-shaped trajectory:

- The **horizontal bar** represents breadth—exposure to tools, methods, and adjacent domains.
- The **vertical bar** represents depth—specialization in one or more areas like causal inference, simulation, or optimization.

This model supports differentiated roles:

- A generalist might help design broad research programs.
- A specialist might focus on methodological innovation or platform-level modeling.

Title tracks should reflect this diversity. Rather than forcing all scientists into a “full-stack” mold, orgs can recognize distinct paths:

- Researcher
- Methodologist
- Domain expert
- Tool-builder

This structure also helps build effective teams—ones that combine varied skill sets across the disciplinary landscape.

# Data Science as a Team Sport

(Organizing scientists into labs with complementary strengths across the venn.)



# **Research Management and Process**

# **Managing Data Science Through Research Programs**

(Differentiating research from projects; long-term, iterative, institutional memory.)

# **The Scientific Method in Data Science**

(Hypothesis-driven study design, using Jira and other tools to support structured inquiry.)

# Using the Capability Maturity Model in Data Science

(Applying the CMM to assess and scaffold experimentation and system design.)

# **Tooling and Lab Practice**

# The Data Scientist's Lab Workbench

(SQL as data collection, Python as instrumentation, Jupyter/Markdown/Quarto as lab notebooks.)

# Documentation and Knowledge Repositories

(Writing for different audiences, synthesis vs analysis, memory and communication.)

# **Statistical Thinking**



# Probability and Statistical Inference

(Quantifying uncertainty, validating observations, and the logic of belief.)