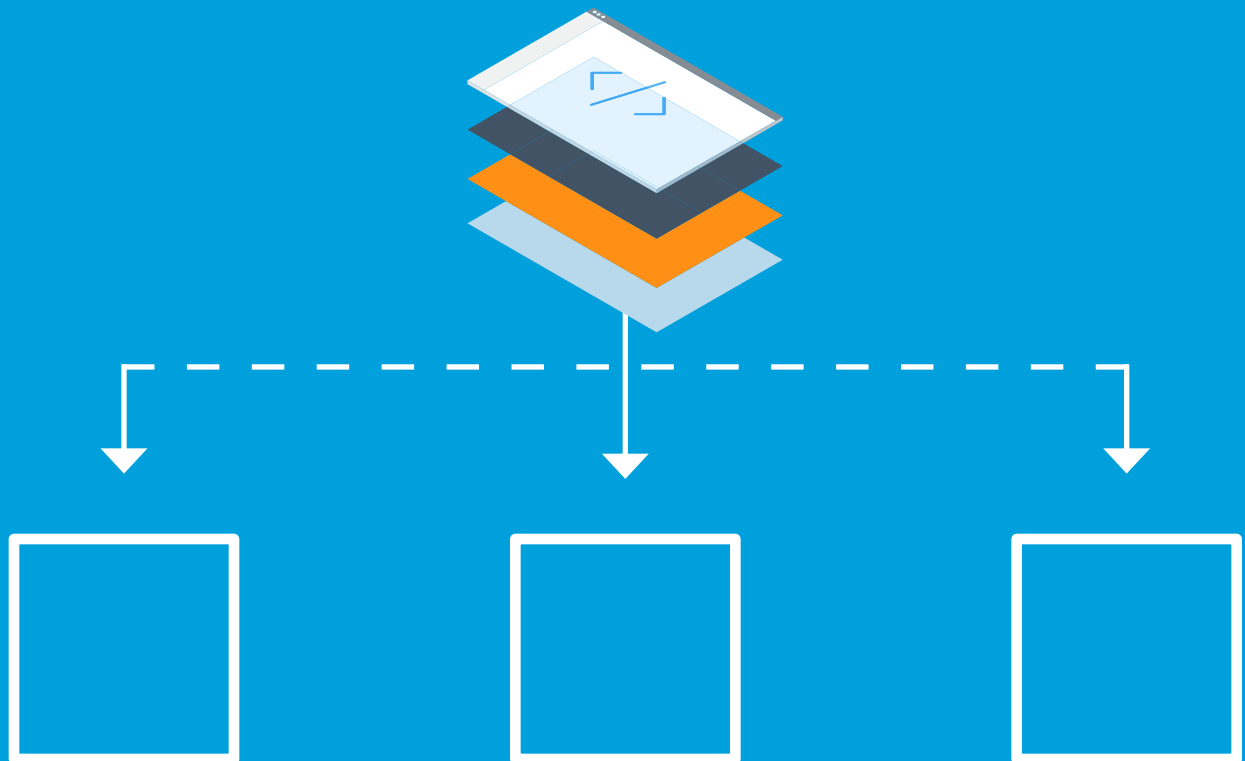


ORACLE® + DATASCIENCE.COM

# Data Science Platforms

Helping IT Drive Digital Transformation



# Introduction

Digital transformation is an increasingly prevalent buzzword, but what does it mean? George Westerman, principal research scientist with the MIT Sloan Initiative on the Digital Economy, provides an easily digestible definition: “Digital transformation is when companies use technology to radically change the performance or reach of an enterprise.”<sup>1</sup>

As technology advances at a remarkably quick pace, companies must shift and adapt their high-level strategies to keep up. Research attests to the importance of a digital transformation strategy. In a study conducted by marketing research firm Altimeter, 41% of the 500 strategists and executives surveyed reported increased market share thanks to their digital transformation efforts, 37% reported increased customer engagement in digital channels, and 30% reported increased customer revenue, among other benefits.<sup>2</sup>

Accurate data is the critical foundation for an effective digital transformation strategy, and hiring and expanding data science teams to leverage insights at the enterprise level has become a top priority. As enterprises increasingly scale their data science teams, it falls on IT to support them effectively.

IT managers supporting a large team of data scientists in an enterprise setting are tasked with data governance, as well as providing the infrastructure and tools that data scientists need. The proliferation of data and data science tools and applications available provides opportunities as well as challenges. While data science teams want to use open source applications, and IT teams want to support open source applications, provisioning access, providing security audits, and resolving technical issues with open source applications can be a nightmare for IT to manage and secure.

Ultimately, if a company wants to stay ahead of the digital transformation curve, they need to have a product that supports reproducibility and standardization. A data science platform can help a company meet the demands of their team, but before we go into all of the problems that a platform solves, both from the IT manager perspective as well as across the larger organization, let's take a step back and discuss the state of data science today at the enterprise level.

---

<sup>1</sup> CIO, “What is digital transformation? A necessary disruption,” July 2017

<sup>2</sup> Altimeter, “The 2016 State of Digital Transformation,” 2016

# The State of Enterprise Data Science

By now, we know that data science, a discipline that uses data to inform business decisions, is more than just a trendy term. The discipline encompasses many job titles across different industries and organizations, from analytics officer, to actuary, to research scientist. But regardless of the exact title, all of these roles are united in their mission to unlock strategic insights from data, for which business demand is stronger than ever before. In fact, IBM predicts that demand for data scientists will soar 28% by 2020,<sup>3</sup> and IDC predicts that revenue from the sales of big data and analytics tools, applications, and services will increase to more than \$187 billion in 2019, up from \$122 billion in 2015.<sup>4</sup>

We've already seen many companies implement digital transformation strategies with a strong emphasis on data science, such as JetBlue. The company launched a subsidiary, JetBlue Technology Ventures, to invest in startups that are bringing machine learning and analytics to the travel market.<sup>5</sup> In March 2016, the company made its first investment in Flyr, a startup that provides predictive analytics and machine learning software to let travelers know when to purchase an airline ticket.<sup>6</sup> Discussing the reasoning behind the investment, President of JetBlue Technology Ventures Bonny Shimi said, "They share our belief that predictive analytics can provide value to travellers and will change the travel experience in ways we have yet to imagine."

Another industry harnessing the power of big data? Retail. "The battle against e-commerce is putting new pressure on brick-and-mortar retailers to fix up stores and deliver a more pleasant experience for shoppers. This means studying data – lots and lots of data,"<sup>7</sup> writes Taylor Cromwell, from Bloomberg. Indeed, a range of retail companies, from Walmart to Nordstrom, are utilizing data-driven insights to implement real-time dynamic pricing and recommendation engines, among other applications.

The evidence is overwhelming: The marriage of big data and data science is prevalent in every field and sector, and every company wants to leverage data to transform how it does business. However, most aren't. According to Forrester, only 22% of companies are actually leveraging big data well enough to get ahead of their competition.<sup>8</sup> What's holding businesses back? In the United States, over a third of companies consider the complexities of IT to be the biggest hurdle to digital transformation.<sup>9</sup> In the next section, we will explore some of the pain points IT teams face as they support increasing numbers of data scientists, and how a data science platform can help.

---

<sup>3</sup> Forbes, "IBM Predicts Demand for Data Scientists Will Soar 28% by 2020," May 2017

<sup>4</sup> InformationWeek, "Big Data, Analytics Sales Will Reach \$187 Billion by 2019," June 2016

<sup>5</sup> CIO, "An Inside Look at 10 real-world digital transformation success stories," June 2017

<sup>6</sup> CIO, "JetBlue CIO pilots VC arm in search of revenue growth," May 2016

<sup>7</sup> Bloomberg, "Here's a Retail Job That's Still in High Demand: Data Scientist," August 2017

<sup>8</sup> Forrester, "Data Science Platforms Help Companies Turn Data Into Business Value," December 2016

<sup>9</sup> Dynatrace, "The Global Digital Performance & Transformation Audit," 2017

## Problem: Data Silos

Data silos occur at the enterprise level when multiple teams set up their own data stores based on a use case or for the purposes of isolating data access. IT managers will only grant access to those who truly need it, and also keep databases separated by project, to ensure there won't be resource contention. However, even though data silos were built with the best intentions in mind, they often yield more problems than solutions, as they contradict the analytical needs of many organizations where querying across data is required. Harvard Business Review refers to them as "a big costly demon" that makes it "prohibitively costly to extract data and put it to other uses."<sup>10</sup> Implementing a data lake, a storage repository that holds a vast amount of data in its native format until it is needed,<sup>11</sup> is the first step towards dealing with the issue of a data silo.

## Solution: Data Lakes, Data Warehouses, and Hadoop-based Systems

To combat the issues surrounding data silos, many enterprises combine, at regular intervals, the data from each separate store into a data lake, which could include a Hadoop Distributed File System (HDFS), S3, or another Hadoop-compatible file system. Another option is to combine the data from each store into a data warehouse (Redshift, Vertica, etc.). Data lakes and data warehouses are not meant to replace isolated data stores, but are rather a solution to combine all data so it can be queried holistically.

Data lakes, which have "been well received by enterprises to help capture and store raw data of many types at scale and low cost to perform data management transactions, processing and analytics based on special use cases,"<sup>12</sup> are largely the preferred method of data storage over data warehouses.<sup>13</sup> When an enterprise implements a data lake, they are referring to a place to dump all data into a common location. Hadoop, an open source Java-based programming framework that supports the processing and storage of extremely large data sets, can store the data from a data lake in a HDFS. Once it has been collected into a Hadoop-compatible file system, there are a standard set of tools available to query and combine that data. In addition, Hadoop distributed file systems also address the security issues that data silos aim to solve, as there are a number of technologies that are able to secure a Hadoop environment. Most commonly, Kerberos is used for user-level authentication in Hadoop based environments. Additionally, tools like Apache Sentry, Knox, and Ranger are used for more fine grain authorization to access Hadoop data.

By authenticating through Kerberos, users can deploy Oracle's DataScience.com Platform on top of data lakes that are stored in Hadoop distributed file systems, ensuring that data is secure and can be queried holistically across the platform.

---

<sup>10</sup> Harvard Business Review, "Breaking Down Data Silos," December 2016

<sup>11</sup> TechTarget, "Definition: Data Lake"

<sup>12</sup> INFORMS Analytics, "DATA LAKES: The Biggest Big Data Challenges," November/December 2016

<sup>13</sup> KDNuggets, "Data Lake vs. Data Warehouse: Key Differences," September 2015

## Problem: Lack of Standardization

One way of standardizing data analysis and data science is to bring all of the tools, programming languages, package markers, and software dependencies into one centralized platform. However, due to the size of many enterprises, which can be comprised of hundreds of data scientists, as well as a wide breadth of projects, standardization can bring significant challenges and less-than-ideal scenarios, outlined below.

### **Scenario A: Large, Shared Remote Machines**

Data scientists work on remote machines provisioned by IT. IT installs all of the packages needed by data scientists throughout the enterprise. This results in management challenges and difficulty adding new tools, packages and dependencies as needs diverge across teams.

### **Scenario B: Low Standardization, Lack of Reproducibility**

Alternatively, data scientists are working on individual machines per team. In these cases, data scientists have additional flexibility, and may be able to configure the tools that they need for their specific task. However, these environments are rarely available across teams or lack the oversight from IT.

## Solution: Container Technologies

The use of containerization technologies, such as Docker, capture system dependencies in a lightweight, reproducible way that can be shared across teams. IT can implement system configuration in a Docker file and store Docker images in a central repository, allowing individuals to quickly launch the environment they need depending on the project they are working on. Using the DataScience.com Platform, IT can set up base environments with the packages, languages, and tools that data scientists need. This gives IT the governance and management over tools and applications, while also empowering data scientists to run self-serve analyses.

As the explosion of open source applications infiltrate the marketplace, it's more important than ever that IT provision settings so that data scientists working across an enterprise have access to the same versions of tools.

## Problem: Lack of Resource Management

All data science projects require compute resources, whether you're using a local laptop or a powerful cluster of cloud-based servers. Often in an enterprise, many large servers are shared across a team of data scientists. From an IT administration standpoint, provisioning and managing these servers can be a time-consuming endeavor. Because data science projects vary so much in resource requirements, a single user could easily be consuming all of a machine's compute power, leaving it unavailable for other users.

In order to avoid the misallocation of resources, the IT team should have complete control over the cluster, and data scientists should only be able to choose from a pre-determined bucket of compute sizes. Additionally, IT should be able to scale out the available resources as needed without downtime, a problem that is currently costing businesses \$700 billion a year.<sup>14</sup>

## Solution: Using Docker Swarm in the DataScience.com Platform

The DataScience.com Platform leverages Docker Swarm for container orchestration. This allows IT to provision a pool of servers and configure the sizing options from which a data scientist can choose.

The DataScience.com Platform allows for on-demand servers for our customers running on public clouds, like Amazon Web Services (AWS). The platform provides resource management tools for IT to show which users are using resources within particular projects, as well as giving IT the governance to shut down a user's container if it is left running too long or taking up too many resources.

---

<sup>14</sup> IHS Markit, "Businesses Losing \$700 Billion a Year Due to Downtime, Says IHS," January 2016

## Problem: Low Reproducibility

Imagine the scenario where a data scientist leaves a company unexpectedly, and there is no documentation made available regarding the best practices surrounding all of the work he or she did. When a new data scientist joins the team and is assigned that same project, he or she will have no idea what to do, and will likely lean heavily on IT to get up to speed. For many IT teams supporting a large team of data scientists, this scenario is all too realistic. The underlying pain point? A lack of reproducibility. As DataScience.com Chief Strategy Officer William Merchan explains: "A critical step to gaining reproducibility in data science is defining all of the different steps involved, from ingesting the first piece of data to deploying a model in production. A data science platform will help to create and manage the workflows involving these steps."

Without a standardized infrastructure and workflow, as well as a means of configuring, centralizing, discovering, and deploying these components, enterprises can't achieve reproducibility.

## Solution: An Enterprise Data Science Platform

Enterprise data science platforms, such as the DataScience.com Platform, effectively tackle the issue of reproducibility by:

- Centralizing all of the assets (tools, code, data) needed to do data science across the enterprise.
- Providing version control functionality so that data scientists can track versions of code, as well as deploy multiple versions and test them against each other.
- Showcasing work around the enterprise and making it accessible with features like projects, centralized outputs, and search. IT has governance over this functionality by implementing role-based access control.
- Providing standardized mechanisms for promoting work to production. This includes running scripts as scheduled jobs, or providing microservices, which is a way of breaking up a huge job into smaller parts that can be easily maintained. In the case of the DataScience.com Platform, this means deploying a predictive model or code as an API, which eliminates the burden on IT and engineering on a per-project basis. Otherwise, there would otherwise be a number of initiatives involved before deploying to production, including refactoring the code, rewriting the production stack language, and testing the performance, among other steps.<sup>15</sup>

---

<sup>15</sup> DataScience.com, Navigating the Pitfalls of Model Deployment, December 2016

# Scaling Successfully

As IT managers continue to support larger teams of data scientists within organizations, data science platforms will continue to rise in prominence as an effective and necessary means of scaling a digital transformation strategy. “An enterprise data platform helps data scientists get the most value out of their data by managing the data analytics lifecycle and standardizing routine processes while enforcing security and governance,” said Oracle's DataScience.com CEO Ian Swanson.<sup>16</sup> And, from the IT manager perspective, a data science platform also resolves crucial issues surrounding data silos, standardization, resource management, and reproducibility that currently prevent many companies from realizing the full extent of their revenue potential.

---

<sup>16</sup> Datanami, “Data Science Platforms Seen as Decision-Makers,” January 2017



ORACLE® + DATASCIENCE.COM

Connect with us on social media:

 @DataScienceInc

 <https://www.linkedin.com/company/datascienceinc>

 @DataScienceInc

 <https://www.facebook.com/datascience>

Oracle Corporation, World Headquarters 500 Oracle Parkway Redwood Shores, CA 94065 USA

Copyright © 2018, Oracle and/or its affiliates. All rights reserved. This document is provided for information purposes only, and the contents hereof are subject to change without notice. This document is not warranted to be error-free, nor subject to any other warranties or conditions, whether expressed orally or implied in law, including implied warranties and conditions of merchantability or fitness for a particular purpose. We specifically disclaim any liability with respect to this document, and no contractual obligations are formed either directly or indirectly by this document. This document may not be reproduced or transmitted in any form or by any means, electronic or mechanical, for any purpose, without our prior written permission. Oracle and Java are registered trademarks of Oracle and/or its affiliates. Other names may be trademarks of their respective owners.