# THE COMPETITIVE VALUE OF DATA

**Steve Strongin**
steve.strongin@gs.com

**Amanda Hindlian**
amanda.hindlian@gs.com

**Sandra Lawson**
sandra.lawson@gs.com

**Sonya Banerjee**
sonya.banerjee@gs.com

The Goldman Sachs Group, Inc.

# Table of Contents

The Global Markets Institute is the research think tank within Goldman Sachs Global Investment Research. For other important disclosures, see the Disclosure Appendix.

# Executive summary

Data is now the lifeblood of many firms, particularly in the modern economy in which companies tend to focus on their narrow area of expertise while outsourcing the rest[1]. From organizing and optimizing complex multi-vendor production processes to customer acquisition, service and retention – these modern firms are almost entirely dependent on data. Naturally, trying to use data to establish a competitive edge has therefore become big business.

Anecdotes about data-driven successes abound, but experience suggests that it is actually quite difficult for businesses to use data to build a sustainable competitive advantage. In fact, pinpointing examples of companies that have successfully used data to maintain a competitive edge is a challenging task. This begs the following two questions: 1) why haven't more companies been able to build a sustainable competitive edge using data, and 2) when can data serve this purpose?

We address these two questions by building a conceptual framework that we refer to as the "learning curve." The learning curve helps us assess the factors that underpin when a company can use data to create a competitive edge – and perhaps more importantly, when it cannot.

Using the learning curve, we analyze four types of data-driven learning strategies:

- **Data-smart strategies** rely on a business's internally generated data as the foundation for producing data-based insights – or what can be thought of as learning. These insights can be used to optimize both a firm's operations as well as its output. An example of a business that uses a data-smart strategy is Amazon's logistics service.

- **Data-asset strategies** are dependent on a business's ability to build a proprietary dataset using secondary sources, for example by collecting (free or purchased) data from sensors, genetic labs or satellites. These proprietary datasets can be used to produce data services that are sold to others. An example of a business that uses a data-asset strategy is IBM Watson Health.

- **Data-feedback strategies** are applicable to businesses that collect user data. To employ this strategy, businesses collect the data that is generated by the users of their products or services, analyze it and leverage the resulting insights to enhance their products or services. Said another way, data-feedback strategies describe when a company leverages user data to create a feedback loop between its users and the goods or services it provides to those users. Examples of businesses that use data-feedback strategies include Spotify with its playlist suggestions, Amazon with its retail product recommendations or Google Maps.

- **Network strategies** are also applicable to businesses that collect user data. However, the purpose of a network strategy is to leverage user data to connect

---

[1]   Strongin et al, "The Everything-as-a-Service Economy" (Dec 2018)

users with one another. Examples of businesses that use network strategies include Uber, Lyft, Airbnb and Facebook.

While the economic models that underlie each type of learning strategy are unique, each one requires data accumulation to drive learning, which then serves as the primary source of potential competitive advantage. We also analyze the role of data decay and copy risk in determining the competitive value of a data-based advantage.
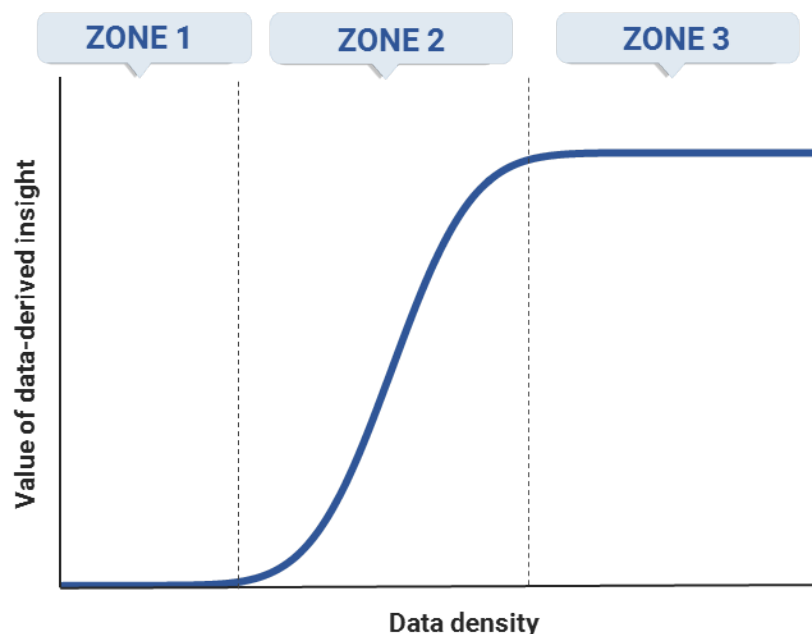
# The learning curve

Exhibit 1 is an illustrative depiction of a learning curve, which can be used to assess the scale-based economics of learning as well as the potential competitive impact. More specifically, and as the exhibit shows, the learning curve is a depiction of the potential value of data (PVD) – or the total value of what can be learned from data – as a function of the amount of usable data a business possesses.

Each unit on the y-axis represents the incremental value derived from analyzing data related to a specific question. Each unit on the x-axis represents the density (or volume) of usable data, which is dependent on the rate of data collection as well as the rate of data decay.

**Exhibit 1: The learning curve: the potential value of data (PVD) as a function of the amount of usable data a business possesses**
A conceptual framework for assessing the scale-based economics of learning



Source: Goldman Sachs Global Investment Research

From an economic standpoint, the central point of the learning curve is that data-derived knowledge is non-linear – meaning, it does not increase without bounds as the volume of data increases. This is for the simple reason that once there is sufficient data to answer the question or problem at hand, additional data only confirms what's already known – so the value of additional data and analysis is trivial. The total potential value of data is therefore constrained by the nature of the question (or questions) at hand.
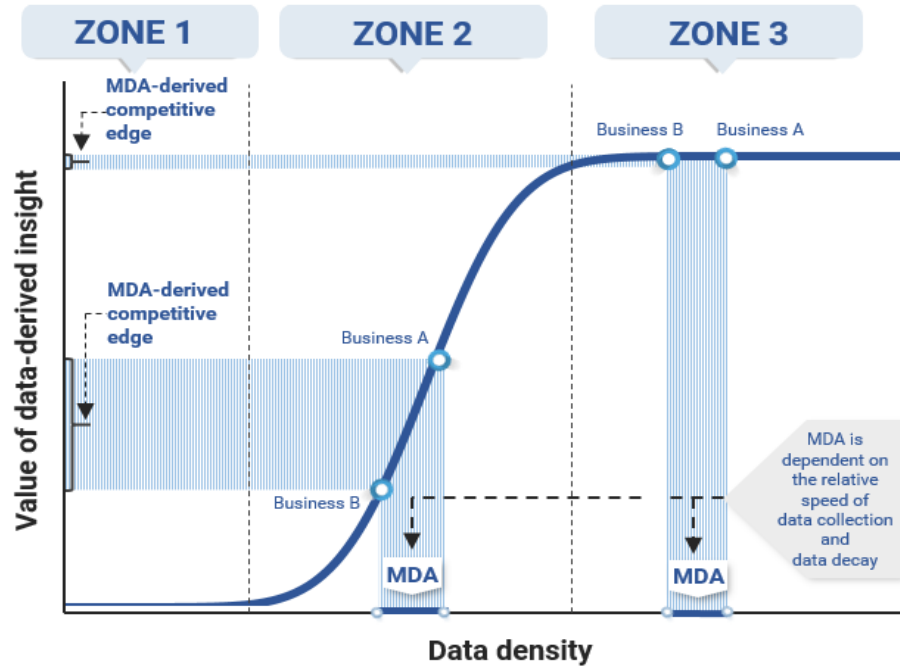
Thus, for each type of learning strategy, the uncertainty related to the PVD is a central question. The PVD must be large enough to justify the expense of building, buying or collecting the data. But, as is often the case, the actual value of data-based insights is largely unknown until after the underlying database is built and the analysis has been done.

Another central issue for all learning strategies is related to data scarcity. On the one hand, if there isn't enough data available to analyze, data-based strategies are likely to get trapped in zone 1, producing little value. On the other hand, if the data isn't scarce in some way, all market participants will likely reach zone 3, where data-based analysis does not provide meaningful competitive differentiation.

With this in mind, consider that the learning curve has a fairly specific shape that is common to all learning problems, and that it is comprised of three specific zones. Exhibit 2 illustrates these dynamics.

- In zone 1, depicted on the left-hand side of Exhibit 2, the learning curve is flat and the incremental value associated with data analysis is low. This means the gains associated with incremental data analysis and data density are limited. The fact that learning is slow in zone 1 is due to the fact that a certain amount of data must be collected before it can be effectively modeled.

- In zone 2, the learning curve begins to slope upward and becomes steeper, typically very steep. At this point, the nature of the data model has become clearer and is better defined, so the incremental value of data-derived information is high. As a result, in this zone, accumulating more data – particularly relative to competitors – can result in a maintainable data advantage (MDA) and can generate significant incremental value (as the middle portion of Exhibit 2 shows). The MDA refers to the pure advantage in the amount of data one business can collect relative to another; the learning curve can then be used to map that MDA to determine a business's relative competitive position given the value of its data-derived insights.

- In zone 3, the learning curve flattens because incremental data accumulation and analysis no longer result in significant value, which can be seen on the right-hand side of Exhibit 2. In this zone, the learning process is nearly complete since most of what can be learned from data to address a specific question or problem has already been learned; businesses in the same market segment that reach zone 3 are in essentially the same competitive position. Accordingly, as Exhibit 2 shows, the same MDA that resulted in a significant competitive advantage for Business A relative to Business B in zone 2 becomes a very small advantage if both businesses reach zone 3.

**Exhibit 2: The value of a maintainable data advantage (MDA) differs significantly by zone**

While not technically precise, it can be helpful to think of zone 1 as the model specification search, zone 2 as the model estimation and zone 3 as the model verification.

More broadly, the learning curve – with its characteristic S-curve shape – can be derived in a number of ways, but the derivation that is easiest to understand (both in terms of the underlying mathematics and the economics) comes from network theory. As a network is built by connecting nodes (at random in the simplest derivations), initially the connections only link two isolated nodes. This creates some incremental value per connection added, but not a tremendous amount (zone 1).

As more nodes are connected, a state is reached in which many of the nodes are already connected to other nodes. As a result, instead of linking two isolated nodes, new connections usually link clusters of already connected nodes. This means that incremental connections increase the average size of a network cluster much faster and thus the value of the network increases more rapidly as the number of connections increases (zone 2).

Once this stage of connecting clusters begins, the network quickly becomes one large cluster plus some isolated nodes and small clusters that remain unconnected to the big cluster. At that point, additional connections can no longer create much value, as they can only add on small clusters to the big cluster and often only connect nodes that are already connected through other paths (zone 3).
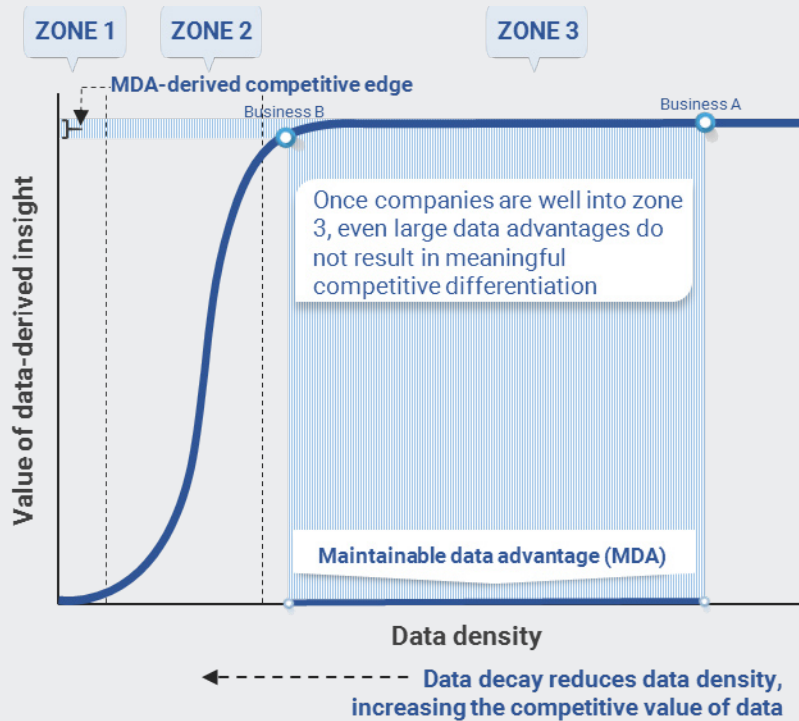
## Data density: the rate of collection vs. the rate of decay

The key to understanding whether user data can serve as a source of sustainable competitive advantage – and whether it cannot – is data density. Data density is driven by two separate processes: the rate of data collection and the rate of data decay.

As Exhibit 3 shows, in zone 3 – regardless of the size – a maintainable advantage in data density (which we refer to as maintainable data advantage or MDA) has little value for a business since there is typically little competitive differentiation. In zone 2, however, even a small MDA can be very valuable and highly differentiating.

Depending on where a business is positioned along the learning curve, the value of data-derived insight changes and affects whether it is possible to turn a data advantage into a competitive one. To that end, the rate of data decay is critical. If the rate of data decay is low, then eventually all data collectors (even those with slow rates) will eventually end up in zone 3. If the rate of data decay is high, however, then the business is likely to be limited from progressing past zone 2 or zone 1. The key is that in zone 2 it becomes possible for the business to establish a competitive edge if it can maintain an advantage in data density.

**Exhibit 3: In zone 3, there is little competitive differentiation between businesses regardless of the maintainable data advantage (MDA)**



Source: Goldman Sachs Global Investment Research

As an example of how data decay works in practice, consider it in the context of navigational maps. Depending on the precise nature of a map's usage, the user's sensitivity to accuracy and to how recently the data underpinning it was collected changes – thus altering the effective rate of data decay.

Navigational maps that are used to locate places or roads generally have a slow rate of data decay since new places and new roads are relatively infrequent occurrences. For example, it's equally easy to locate

the Grand Canyon on a map of the United States today as it was 50 years ago. In past generations, it was common to find 10-year-old maps in cars that could be used during navigational emergencies.

Accordingly, in the case of simple navigation, the slow rate of data decay made it possible and relatively easy for all map providers to reach zone 3 on the learning curve (where little or no competitive advantage could be derived from differences in accuracy or how recently the data was collected). However, if we consider the use of maps in the context of a more demanding problem – for example, to find the fastest route home through a busy city during rush hour – the dynamic changes.

In the case of real-time traffic navigation applications, like Waze or Google Maps, the accumulated data is subject to very high rates of decay, such that reaching zone 3 in the learning curve is very difficult; this is particularly the case in terms of side routes, or when traffic patterns are changing rapidly. In this situation, the best vendor has a significant and self-reinforcing advantage.

This is because these services often collect and analyze users' location information to provide real-time navigation guidance. Thus the more users any one service can attract, the faster their rate of data collection and the more accurate their insights, which allows them to move up the curve in zone 2 and to stay there as users congregate around the best provider. What's more, the concentration of users on one vendor's service lowers competitors' rates of data collection, reducing the value of their data-derived insights (trapping competitors in zone 1). This further reinforces the lead vendor's edge, even on less used routes where data collection is more difficult.

A similar dynamic can be observed in web-based search. Early on, when web crawlers – a tool for indexing web pages to support search engines – were viewed as central to a vendor's competitiveness in the space, many vendors were willing to invest in developing the technology; the rate of change in web pages was sufficiently slow that reaching zone 3 was viewed as widely achievable. As the searches themselves – particularly recent searches with a short-lived relevancy – became more important to producing the best (most relevant) search results, a clear self-reinforcing dynamic took hold. This is especially true in the case of popular or trend-based searches.

As a result, Google – which pioneered the use of its repository of past searches to improve the applicability of its real-time search results – has been able to translate a data collection advantage into a competitive one in online search. Google's ability to anticipate users' keystrokes, highlight "hot" places to go or feature trending stories are examples of features that incorporate large volumes of data with high rates of decay.
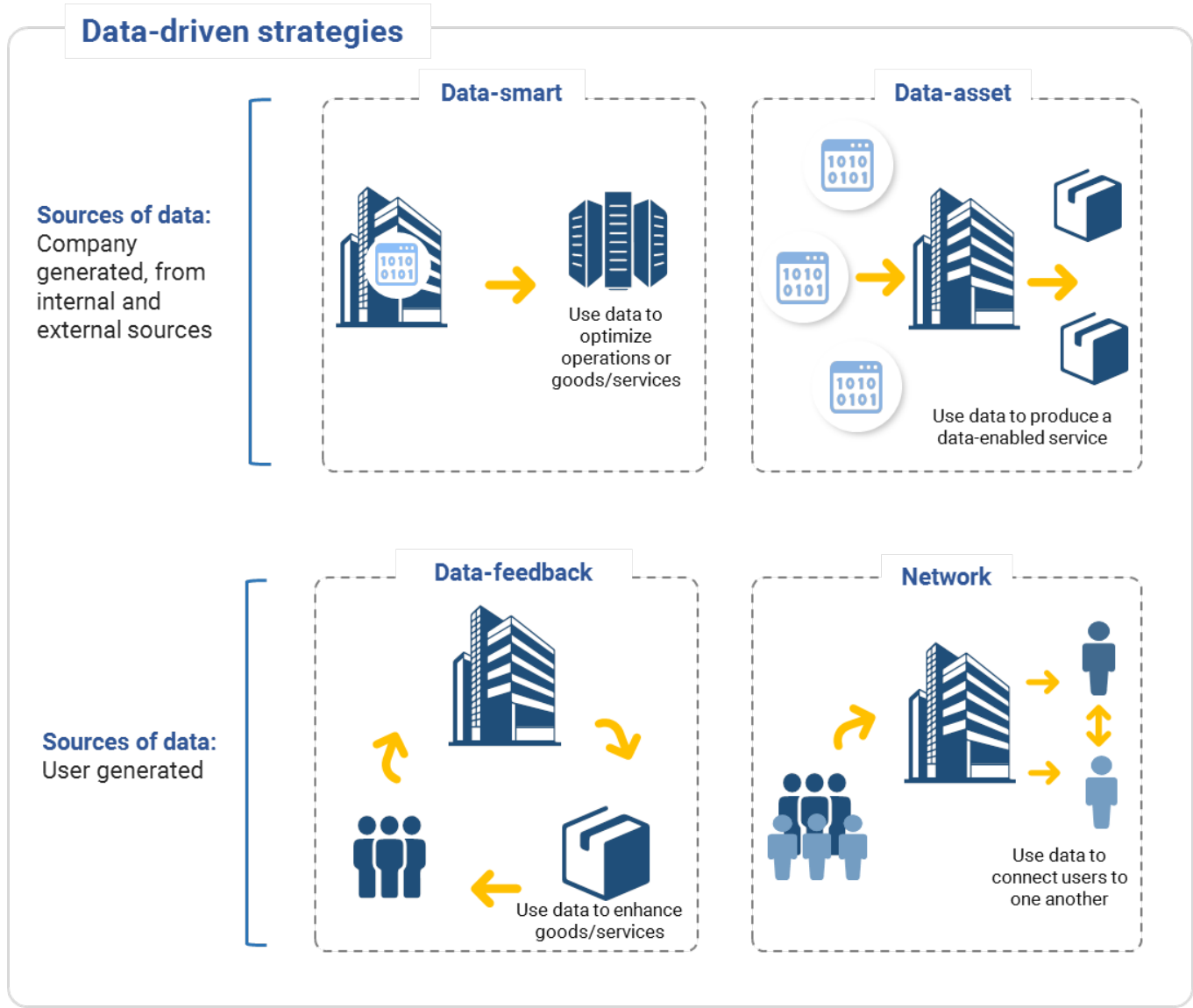
As we previously noted, unless the rate of data decay is high for a given problem, all businesses addressing that problem are likely to end up in zone 3 with little to no competitive differentiation. However, when the rate of data decay is high, a lead in data collection – an MDA – can become a self-sustaining and self-reinforcing advantage. Data-density advantages only translate into competitive advantage in zone 2 and can be maintained only if competitors (particularly the runner-up) don't make it to zone 3.

We believe this is why, despite the general perception of the importance of user data, there are more examples of successful uses of data-smart and data-asset strategies than of successful uses of data-feedback or network strategies. In the latter instances, it just isn't that easy to find examples where the runner-up doesn't eventually make it to zone 3.

# Data-driven learning strategies

As we noted at the outset of this report, we use the learning curve to analyze four types of data-driven learning strategies: data-smart, data-asset, data-feedback and network. Exhibit 4 illustrates how each strategy works, both in terms of how data is sourced and used in each case.

**Exhibit 4: The four types of data-driven learning strategies**

### The data-smart strategy

It's a popular refrain to believe that all companies should become data-smart – since collecting and analyzing one's own operational data may seem like low hanging fruit. In practice, however, it may not be relevant or possible to pursue this strategy.

Individual companies frequently have difficulty producing enough data on their own to be able to implement big-data types of analyses. Modern approaches to big data, AI and

the like require vast quantities of data to produce meaningful insights that can move a business from zone 1 to zone 2. Thus, in many cases, because the quantity of data is simply insufficient, being data-smart simply proves impossible.

But, when an individual company is able to generate enough operational data to successfully reach zone 2 or even zone 3, it is likely that the information that is collected is related to highly repetitive tasks (related to logistics, simple customer support or other basic operations, as examples). It is worth noting that the risk-to-reward associated with making significant investments in collecting and analyzing such data – based on the notion that doing so will reveal hidden or unknown insights – may be unfavorable; said another way, it is worth considering whether the PVD is sufficiently high relative to the investment involved.

Rather than trying to use this kind of data to optimize a product or service, instead what may be a better strategy – with a more favorable investment outcome – could be to focus on operational optimization. In this case, the first analysis can be run using the high initial volume of data, which can be complemented over time through high ongoing usage that allows even small efficiency improvements to accumulate with meaningful results.

Taken together, these factors suggest that feasible data-smart strategies are often likely to culminate in zone 3, which means they will generally be defensive in nature – they may be a cost of entry, for example. What's more, businesses that fail to realize the basic efficiencies associated with a data-smart strategy are likely to be at a significant competitive disadvantage relative to those businesses that have realized those efficiencies.

Another implication associated with this type of business strategy is that it may actually be better to be a second-mover rather than a first-mover from an investment perspective. Knowing another business has succeeded at uncovering meaningful efficiencies from a particular data-smart strategy significantly improves the related risk-to-reward ratio. It may actually be better to mimic the strategy that's already proven successful, rather than to pursue a novel data-smart strategy.

### The data-asset strategy
Unlike data-smart strategies, data-asset strategies are predicated on a business's ability to build robust proprietary databases – often using secondary sources of data – that allow it to produce data-driven services.

Constructing such a database typically requires a significant upfront investment associated with acquiring the necessary data, as does the related analysis. What's more, at the point when these investments are made, the business typically does not know how much data will be necessary to allow it to progress into zone 2 or zone 3, nor does it know the PVD of the data.

As we touched on earlier, second movers do not face the same level of uncertainty that first movers do, which means their investments are subject to a more favorable risk-to-reward tradeoff. Second movers not only know that valuable data-based insights

do exist, but they also have a general sense both for the volume of data necessary to extract these insights and for the magnitude of the related PVD.

At the same time, the second mover faces the risk of lower potential profitability. This is because when the second mover enters the market, the first mover is incentivized to cut prices well below the average cost for the simple reason that the marginal cost to deliver data-based services is lower than the fixed cost to develop the services in the first place.

Copy risk can be difficult to determine, particularly before a company knows how much value a particular data-asset strategy will generate to address a specific problem or question. This only reinforces the need for businesses pursuing data-asset strategies to diversify and to have sufficient capital to experiment again.

Broadly speaking, however, as businesses decide which strategies to invest in, there are two observations worth considering. First, if it is likely that the full investment (in both the data and the related analysis) will need to be replicated to produce the results, the investment is likely safer from a risk-to-reward perspective. Second, when it is likely that a second mover will be able to bypass the full investment and still arrive at the same results, the original strategy is more likely to be copied and the risk associated with the first mover's investment is high.

The nature of copy risk can be made clearer through examples. Consider a "safe" example first, meaning a case involving low copy risk. As in the case of IBM's Watson Health business, interpreting magnetic resonance imaging (MRI) data requires a large initial database of interpreted images and significant ongoing technology investments, both to receive and to interpret new MRI data. Thus an ongoing build of cross-checked interpretations would make replicating this data-asset strategy difficult.

Next, consider a less safe example involving a data-asset enabled maintenance service for elevators. This service is driven by data collected from sensor arrays or from past instances of elevator maintenance. On the one hand, if producing the maintenance service requires a complex assessment of the sensor input data, copying the strategy could be difficult – in which case the originating firm would be able to maintain a competitive advantage. On the other hand, if the maintenance service could be approximated through simpler rules, for example by counting hours of service rather than calendar time associated with the service, then the service could be copied at a lower cost – and the associated copy risk would be high.

Economies of scope can easily play a significant role in driving data-asset efficiencies. The lessons learned and technologies developed in one data-asset project may result in new but related projects. Sensor-based data collection and interpretation, or image-based data processing and interpretation, as examples, could easily represent natural projects with scope efficiencies. In either instance, the related skills could be applied to many different databases and therefore could allow a business pursuing a data-asset strategy to become even better at both assessing the risks and lowering the cost associated with new ventures related to their particular area of expertise.

**Data-asset strategies: a comparison to deep-water drilling for oil or new drug development**

The risk-to-reward ratio of data-asset strategies is in many ways analogous to deep-water drilling for oil or to new drug development: there are high up-front costs and there is significant uncertainty associated with the potential discovery, but there is also a long tail of payments if the endeavor is successful. Another similarity is that data-asset strategies also require significant capital and diversification efforts to create a reasonable risk-to-reward tradeoff.

Accordingly, it is not surprising that well-established firms like IBM with its Watson Health business (and to a lesser extent, Google and Amazon) have led the way in the data-asset space, though there are start-up businesses that have made some inroads (as with Flatiron Health, for example, which was acquired by Roche Holdings).

There are also a number of important differences between data-asset strategies and oil platforms or the development of new drugs. Perhaps the most important difference is that data-asset firms, unlike oil platforms or pharmaceutical companies that develop new drugs, must assess copy risk (as discussed earlier), since potential competitors (for example, new entrants) face very different incentives and hurdles than the innovators that preceded them.

### The data-feedback strategy

The most complex – but perhaps most talked about – data strategy is one that relies on the collection of user data to refine the user experience, which we refer to as the data-feedback strategy. There are two key challenges related to pursuing a data-feedback strategy: the first is determining whether an advantage actually exists, and the second is determining whether it is possible to maintain that advantage.

Historical efforts suggest that finding a true advantage based on customer data isn't easy. Individuals' "discovered behavioral patterns" generally aren't complex or surprising. Amazon offers an illustrative example. The firm's use of its own data for logistics and inventory management (data-smart strategies) has been helpful. The firm also has one of the largest customer databases ever amassed, yet its product placement and sales strategies are often quite simple, to the point where third-party retailers can mimic Amazon's strategies and now outpace Amazon in terms of unit sales on Amazon's own retail platform.

To make this example more granular, consider that a business doesn't need to have Amazon's extensive customer database to realize that a consumer who is searching for ovens may want to purchase one. While an advertiser can use this information to display ads of ovens (showing ones that are better or cheaper, but similar to what the consumer has already viewed), for a merchant to serve this customer well, more often than not, it will simply need to stock the most popular oven models, which does not necessitate extensive customer-specific data or analysis. This is another kind of copy risk.

Accordingly, when an advantage can be found (when the PVD is high) the data-accumulation process must be sufficiently difficult that the business pursuing the data-feedback strategy is able to progress up the learning curve (and capture a

significant portion of that PVD). As well, competitors must be constrained from doing the same.

### The network strategy

Network strategies are similar to data-feedback strategies in that they also leverage user data in ways that reinforce the value of their products or services. The primary difference is that network strategies use this data to connect users to each other (while data-feedback strategies leverage user data to enhance the output that is provided to each user).

For businesses that use network strategies, data density is defined by the number of active users. The key driver of data decay typically has more to do with activity levels than a change in the data – as is the case for the other types of strategies we have identified.

The competitive issues associated with network strategies are similar to those associated with data-feedback strategies. For both types of strategies, progressing out of zone 1 involves significant hurdles related to amassing sufficient user data. After doing so, the business's ability to build a sustainable competitive advantage is dependent on whether competitors are also able to reach zone 3 – where there is little to no differentiation.

Businesses pursuing network strategies must consider: 1) how an active user in the space is defined, and 2) whether being an active user in one network precludes the user from being active in others. If businesses are competing for users' time (as with Netflix or Instagram), there is a natural constraint that forces the system toward dominant vendors. However, if the service is consumed based on specific needs (as with Uber and Lyft in terms of car service, or Airbnb and VRBO for temporary lodging), the market is more likely to have multiple vendors that are in ongoing competition, and the network alone is unlikely to result in a persistent advantage.

Communities of users – and the relevant boundaries – play an important role in driving the economics of network strategies. In some circumstances, networks naturally divide into communities in which there is an advantage in specializing in providing network services to a specific community rather than to the general population. Modern dating applications – like Bumble and e-Harmony – are examples of businesses leveraging network strategies and where success is determined by active users.

For businesses using network strategies, the ability to monitor and regulate membership can become a sustainable advantage. As examples, the ability to offer high-quality drivers, dynamic rental spaces, verified vendors, or other specific community affiliations may represent key competitive strengths. In such cases, the business may be mixing two types of learning strategies: a network strategy (based on the directory of users) with a data-asset strategy (using reviews, background checks) to police the service. The mix can result in a hard-to-replicate business model.

# The four-part test

In summary, analysis of the learning curve leads to a four-part test businesses can take to determine whether data-based strategies can produce a sustainable competitive advantage for them:

1. Is there sufficient data to analyze?

2. Are the insights gained from this data analysis novel enough to have significant value?

3. Is implementing those insights complex enough to do that competitors cannot simply copy the same approach?

4. Is the data scarce enough that a competitor cannot repeat the same analysis?

If each of these questions elicits an affirmative response, building a sustainable competitive edge through data is possible. More often than not, however, this is unlikely to be the case. As a result, robust second mover strategies may be more cost effective than first mover ones.

# Disclosure Appendix

## Disclosures

This report has been prepared by the Global Markets Institute, the research think tank within the Global Investment Research Division of The Goldman Sachs Group, Inc. ("Goldman Sachs").

Prior to publication, this report may have been discussed with or reviewed by persons outside of the Global Investment Research Division. While this report may discuss implications of legislative, regulatory and economic policy developments for industry sectors and the broader economy, may include strategic corporate advice and may have broad social implications, it does not recommend any individual security or an investment in any individual company and should not be relied upon in making investment decisions with respect to individual companies or securities.

### Distributing entities

The Global Investment Research Division of Goldman Sachs produces and distributes research products for clients of Goldman Sachs on a global basis. Analysts based in Goldman Sachs offices around the world produce equity research on industries and companies, and research on macroeconomics, currencies, commodities and portfolio strategy. This research is disseminated in Australia by Goldman Sachs Australia Pty Ltd (ABN 21 006 797 897); in Brazil by Goldman Sachs do Brasil Corretora de Títulos e Valores Mobiliários S.A.; Ombudsman Goldman Sachs Brazil: 0800 727 5764 and / or ouvidoriagoldmansachs@gs.com. Available Weekdays (except holidays), from 9am to 6pm. Ouvidoria Goldman Sachs Brasil: 0800 727 5764 e/ou ouvidoriagoldmansachs@gs.com. Horário de funcionamento: segunda-feira à sexta-feira (exceto feriados), das 9h às 18h; in Canada by either Goldman Sachs Canada Inc. or Goldman Sachs & Co. LLC; in Hong Kong by Goldman Sachs (Asia) L.L.C.; in India by Goldman Sachs (India) Securities Private Ltd.; in Japan by Goldman Sachs Japan Co., Ltd.; in the Republic of Korea by Goldman Sachs (Asia) L.L.C., Seoul Branch; in New Zealand by Goldman Sachs New Zealand Limited; in Russia by OOO Goldman Sachs; in Singapore by Goldman Sachs (Singapore) Pte. (Company Number: 198602165W); and in the United States of America by Goldman Sachs & Co. LLC. Goldman Sachs International has approved this research in connection with its distribution in the United Kingdom and European Union.

**European Union:** Goldman Sachs International authorised by the Prudential Regulation Authority and regulated by the Financial Conduct Authority and the Prudential Regulation Authority, has approved this research in connection with its distribution in the European Union and United Kingdom; Goldman Sachs AG and Goldman Sachs International Zweigniederlassung Frankfurt, regulated by the Bundesanstalt für Finanzdienstleistungsaufsicht, may also distribute research in Germany.

Goldman Sachs conducts a global full-service, integrated investment banking, investment management and brokerage business. It has investment banking and other business relationships with governments and companies around the world, and publishes equity, fixed income, commodities and economic research about, and with implications for, those governments and companies that may be inconsistent with the views expressed in this report. In addition, its trading and investment businesses and asset management operations may take positions and make decisions without regard to the views expressed in this report.