Slide 1



What is "big"? And how does it apply to data collection?

Slide 2



**What is big data?**

Some other definition — 11%

Explosion of new data sources (social media, mobile device, and machine-generated devices) — 18%

Requirement to store and archive data for regulatory and compliance — 19%

Massive growth of transaction data, including data from customers and the supply chain — 28%

New technologies designed to address the volume, variety, and velocity challenges of Big Data — 24%

**Fig. 2.** Definitions of big data based on an online survey of 154 global executives in April 2012.

Source: Gandomi, A., and Haider, M., 2015. Beyond the hype: Big data concepts, methods, and analytics. *International Journal of Information Management.* 35(2), 137-144.

Big data obviously means different things to different people. But if you pay attention carefully to this diagram, which collected definitions from 154 executives, you will notice that they all address mostly data collection aspects, rather than, for example, data processing or modelling.

What is big data?

http://www.ibmbigdatahub.com/infographic/four-vs-big-data

Most practitioners today use the definition of the 4 Vs of big data: Volume, Velocity, Variety and Veracity.

"Volume" represents the sheer size of the dataset. The size is considered due to both the large number of variables that is collected, as well as the size of the observations' set for each variable.

"Velocity" reflects the speed at which these datasets are collected and analysed. It might be real-time collection from sensors, sales transactions or social media.

"Variety" represents the diverse structured and unstructured data sources and formats such as text, videos, networks, and images.

"Veracity" represents the inconsistency and unreliability of collected data. (George et al., 2016).

Simply put, big data is an amount and nature of data that we are uncomfortable (or less comfortable) to collect, visualize, store or analyse.

Slide 4



If only too few years ago, we were expecting our collected dataset to fit nicely into a spreadsheet, or in a typical relational database, the challenge of collecting data has completely changed, and changed with it our ability to make strong conclusions from it.

The study shown here, for example, is based on the collection of 90 million business articles, integrated with a dataset consisting of minute-by-minute stock prices from 2013 and 2014. This is far from just a sample.
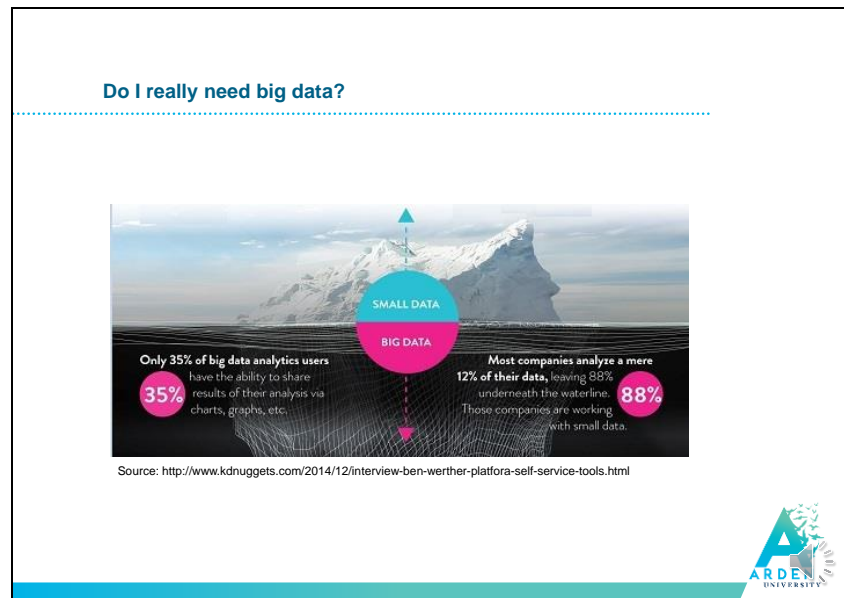
Big data is about collecting more data, that is being created more frequently, in a higher resolution (or granularity), and with much less ability to control its quality.

We were used to studies in which the unit of analysis is a human subject or a group of people, where data was collected using surveys and indirect observations.

In big data, however, the unit of analysis goes down to messages, clicks, or milliseconds of collected observations. This is yet another reason for which we must have automatic methods for data collection.

In this study, for example, the researchers collected more than one billion of clicks from Wikipedia during a single month.

In reality, still, most of the time, we do not have the big data to collect about our business questions.

And when we do, we do not necessarily need it.

And when we do, sometimes its lack of structure, or level of quality, does not allow us to get to the bottom of things.

In Lesson 7 of this module, we will look into analysis methods, which are most often related to "small data"; that is, data which is sampled.

Big data typically requires storage capacity which exceeds the capacity found in regular desktop computers or laptops.

In addition, the variety and veracity of big data is very large, which makes it hard to fit into a certain set of well-specified fields in a table.

This lack of common structure of a large integrated dataset can be answered by NoSQL storage platforms (such as Hadoop's or MongoDB).

NoSQL databases can accommodate data that does not have a pre-defined structure.

Instead of a tabular (or relational) format, NoSQL formats are pairs of key-value (see in the image), graph-based, or document-based, all of which enable a dynamic structure of data.

This approach makes it more flexible in terms of different pre-defined schema, or format, of the collected data.

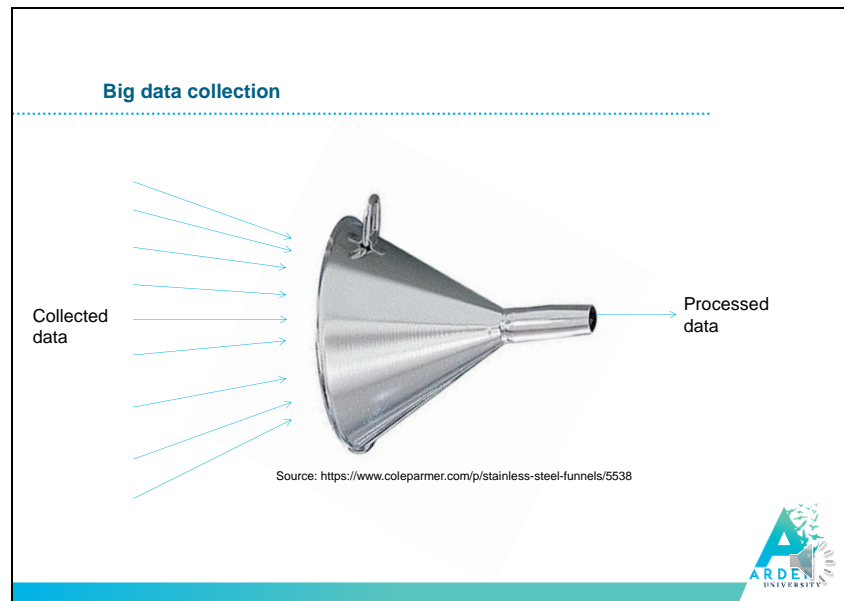Data collection is of course not the only aspect that needs consideration when it to comes to big data.

Big data analysis, for example, requires platforms that can handle large amounts of data. For example, social media analytics, video analytics and audio analytics.

Variable selection is another very typical issue. Big data may contain a large number of variables. If we are not sure in advance which exact variables we are going to need, we are going to need to use statistical methods for variable selection.

Other aspects exist, such as reporting, visualization and ethics.

However, the reality today is that much more emphasis is put on big data collection than on other aspects.
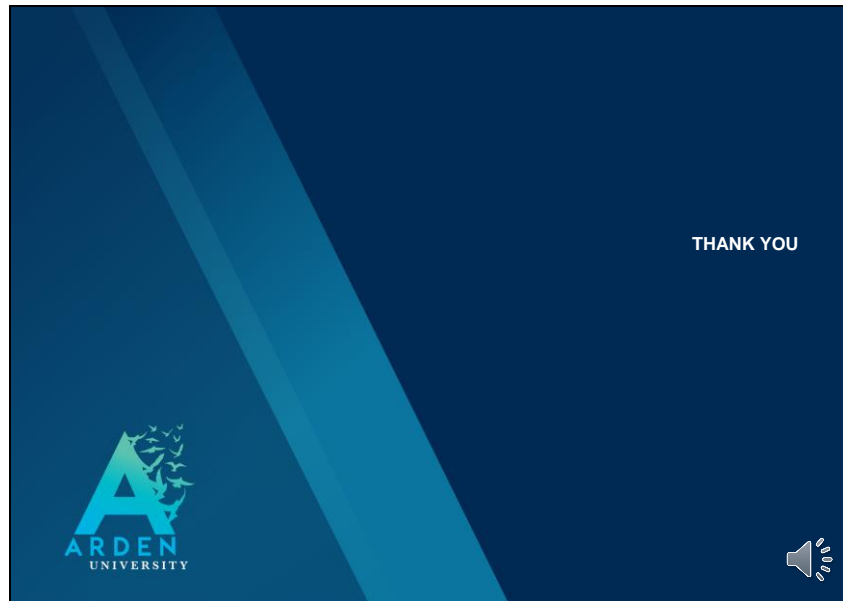
Slide 9



The amount and variety of data that is collected significantly exceeds the amount of data that is asked about and eventually processed.

Unstructured data, for example (such as video or image data), is still very challenging to analyze, and requires large amounts of storage.

Thus, automatic data collection methods are of great significance, specifically when it comes to big data.

Slide 10



Thank you.