

DATA SOCI ETY:

Fundamentals of Data
Literacy

Participant Guide

Table of Contents

<i>Class activities</i>	3
Data governance assessment	4
Data-driven culture assessment	8
Project brainstorm.....	11
Activity: field trip	16
Activity: what would you viz?	17
<i>Additional resources</i>	21
Data science glossary.....	22
Comparison: popular analysis tools	26
Comparison: popular visualization tools	28
Additional reading & reference	29

Class activities



Data governance assessment

The Basic Maturity Assessment is a condensed version of the Stanford Maturity Measurement Tool that model focuses both on foundational and project aspects of data governance (DG). The foundational components (Awareness, Formalization, and Metadata) of the maturity model focus on measuring core DG competencies and development of critical program resources. The project components (Stewardship, Data Quality, and Master Data) measure how effectively DG concepts are applied during funded projects.

Additionally, it includes the three dimensions (People, Policies, and Capabilities) which further subdivide each of the six maturity components, focusing on specific aspects of component maturation.

Whether your organization uses the Stanford Maturity Measurement Tool or the Basic Maturity Assessment, it is imperative that the maturity model you choose is finalized and adopted early in the rollout of the DG program. Depending on where your organization is in the process of standing up the data governance program, it may be most appropriate to use the Basic Maturity Assessment to measure the baseline maturity of and resources available to the organization. Then, as the data governance program is fleshed out, perhaps you will find that a more robust maturity assessment is needed. In that case, because they are both based on the same component-dimensions, you can easily transition from using the Basic Maturity Assessment to using the full Stanford Maturity Measurement Tool.

Regardless of which tool you choose to use, or if you choose to use a combination of both, thoughtful input from across the organization will help assure the model's usefulness and long- term fitness.

Foundational components

Component: Awareness - The extent to which individuals within the organization have knowledge of the roles, rules and technologies associated with the DG program.		
	Objective	Rating
People	Are executives, employees, and stakeholders aware of the purpose or value of the DG program?	

Policies	Are existing data policies documented, consistently maintained, and available to stakeholders?	
Capabilities	Are stakeholders aware of the specific DG capabilities that are available at the organization?	
Component: Formalization - The extent to which roles are structured in an organization and the activities of the employees are governed by rules and procedures.		
	<i>Objective</i>	<i>Rating</i>
People	Have DG roles and responsibilities been defined and vetted with program sponsors?	
Policies	Are data policies around the governance of specific data defined as best practices?	
Capabilities	Are classes of DG capabilities defined and is there an available solution?	
Component: Metadata - Technical metadata describes data elements and other IT assets as well as their use, representation, context, and interrelations. Business metadata answers who, what, where, when, why, and how for users of the data and other IT assets.		
	<i>Objective</i>	<i>Rating</i>
People	Do executives, employees or stakeholders understand types and values of metadata?	
Policies	Are metadata best practices produced and made available?	
Capabilities	Is metadata consistently collected, consolidated, and available from a single portal?	

Project components

Component: Stewardship - The formalization of accountability for the definition, usage, and quality standards of specific data assets within a defined organizational scope.		
	<i>Objective</i>	<i>Rating</i>
People	Have DG or stewardship roles and responsibilities been defined within the organization?	
Policies	Have policies around data stewardship been defined within a functional area?	
Capabilities	Does a centralized location exist for consolidation of and/or access to stewardship related documentation?	
Component: Data Quality - The continuous process for defining the parameters for specifying acceptable levels of data quality to meet business needs, and for ensuring that data quality meets these levels.		
	<i>Objective</i>	<i>Rating</i>

People	Are people assigned to assess and ensure data quality within the scope of each project?	
Policies	Have data quality best practices been defined and adopted as official organizational data policies?	
Capabilities	Have basic data profiling tools been made available for use anywhere in the system development lifecycle?	
Component: Master Data - Business-critical data that is highly shared across the organization. Master data are often codified data, data describing the structure of the organization or key data entities.		
	<i>Objective</i>	<i>Rating</i>
People	Is there consistent understanding among stakeholders of the concepts and benefits of master data?	
Policies	Are there formal policies that define what data are considered institutional master data?	
Capabilities	Are master data identified, managed and provisioned?	

In the tables below, record the **average of each component** as calculated above. Review these scores and then take some time to note down your goals for the future.

Foundational components

	<i>Average Score</i>	<i>Goals</i>
Awareness		
Formalization		
Metadata		

Foundational components

	<i>Average Score</i>	<i>Goals</i>
Awareness		

Formalization		
Metadata		



Data-driven culture assessment

Choose the answers that best apply to you and your organization then scroll to the next page to see your results.

<i>I can easily access the data I need without asking others for help.</i> 0 - Not at all 1 - Only for some colleagues 2 - Only for some teams 3 - Organization-wide	
<i>I can easily access the data I need in a timely manner.</i> 0 - Not at all 1 - Only for some data 2 - Only for data in my team / related teams 3 - Organization-wide	
<i>Data is automatically collected and stored on a continuous basis.</i> 0 - Not at all 1 - Only at someone's request 2 - Regularly, a few times a year 3 - There is continuous data collection	
<i>The data we have is accurate and good quality (few missing entries, few duplicates, accurate measurements).</i> 0 - Not at all 1 - Only for some data 2 - Only for data in my team / related teams 3 - Organization-wide	
<i>Our data is stored securely either internally or offsite.</i> 0 - Not at all 1 - Only for some data 2 - Only for data in my team / related teams 3 - Organization-wide	

<p><i>My company routinely offers data trainings and other educational opportunities.</i></p> <p>0 - Not at all 1 - Occasionally 2 - Regularly, a few times a year 3 - There are continuous learning opportunities</p>	
<p><i>Most of my colleagues understand the importance of data.</i></p> <p>0 - Not at all 1 - Only for some colleagues 2 - Only for some teams 3 - Organization-wide</p>	
<p><i>Our organization has a set of data standards that reviews how data should be collected, stored, and analyzed.</i></p> <p>0 - Not at all 1 - Only for some colleagues 2 - Only for some teams 3 - Company-wide</p>	
<p><i>My organization emphasizes the importance of using data to track initiatives.</i></p> <p>0 - No one 1 - A few people across the company 2 - Some teams across the company 3 - Organization-wide</p>	
<p><i>I am expected to present data metrics when I explain conclusions and decisions.</i></p> <p>0 - Not at all 1 - Only for some colleagues 2 - Only for some teams 3 - Organization-wide</p>	

Results

Your data infrastructure score is _____ out of 15. This represents how well your organization does with data access, collection, and storage.

Your data literacy score is _____ out of 15. This represents how well your organization does with data knowledge, governance, and leadership.

Your overall data culture rating is based on these two scores. Check out the matrix below to see where you belong.

- Data infrastructure +	
- Data literacy +	Data literate Data infrastructure 0-7 Data literacy 8-15
	Data driven Data infrastructure 8-15 Data literacy 8-15
- Data literacy +	Data nascent Data infrastructure 0-7 Data literacy 0-7
	Data prepared Data infrastructure 8-15 Data literacy 0-7



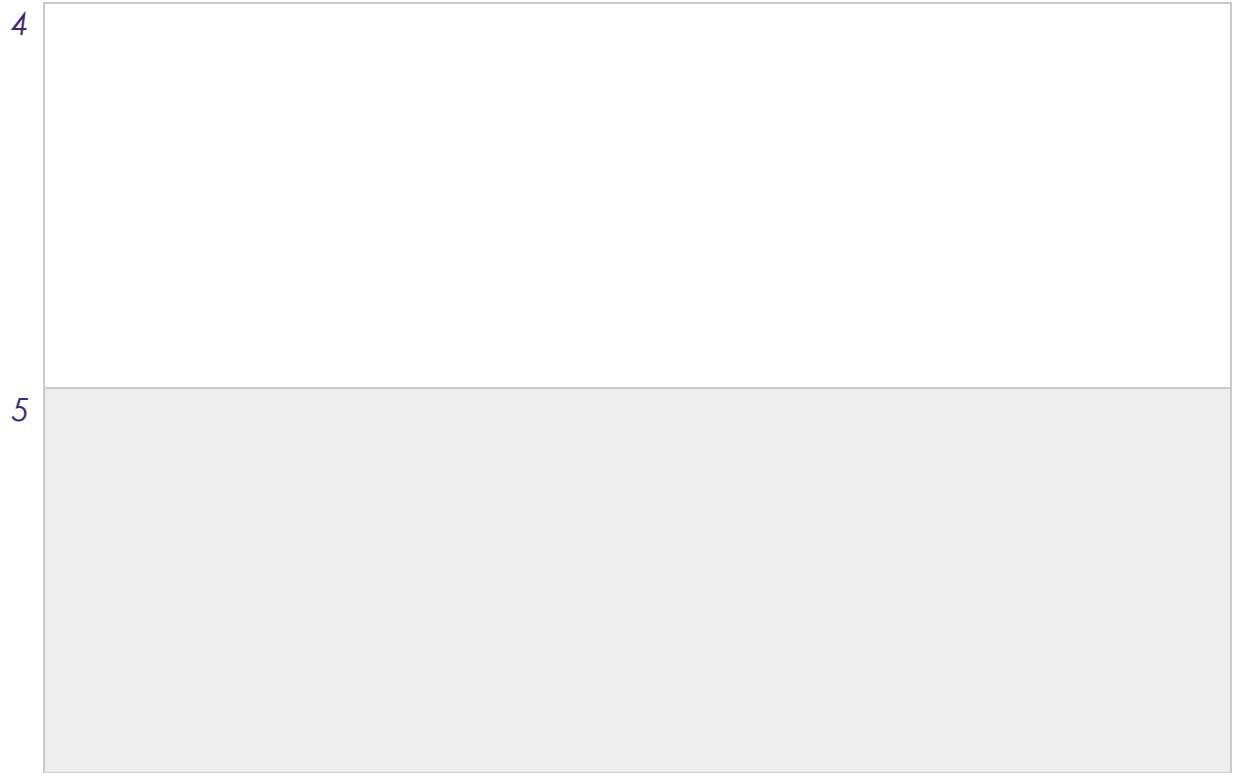
Project brainstorm

Identify 3-5 ideas for leveraging data in your workplace and write them below.

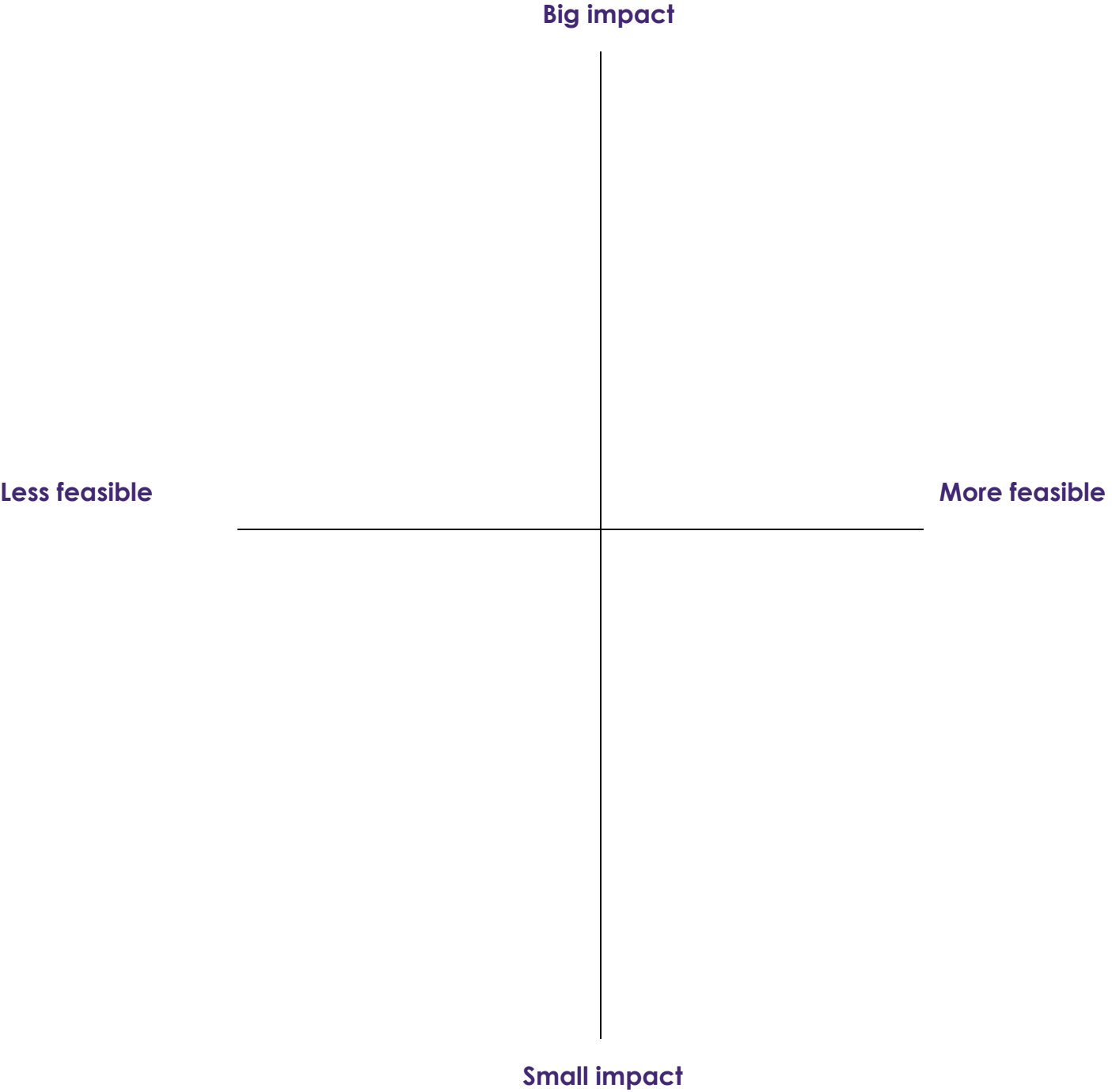
1

2

3



Then, on the next page, assess the feasibility and impact of your ideas and place them in the appropriate quadrant.



Worksheet: developing a plan

Now that you have ideas and know a little bit more about how data projects work, use the templates below and on the next page to further map out a plan.

What question do you want to answer?	
Data sources	
Data science methods	
Teams / staff	
Tools / budget	

Worksheet: developing a plan, cont'd

What challenges do you foresee? How can you solve them?

Challenges	Solutions



Activity: field trip

Visit <https://quickdraw.withgoogle.com/>

Click the “Let’s Draw!” button and play a round (6 drawings).

At the end of the round, visit the data to see why guesses were made.

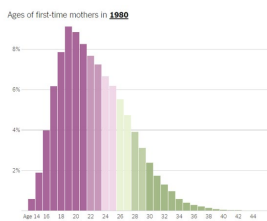
Also, make a note of how many of your drawings were guessed correctly.



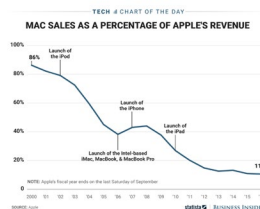
Activity: what would you viz?

For each of the dataset descriptions, choose which of the data visualizations you would select to represent it. Some have more than one possible answer. How would you visualize...

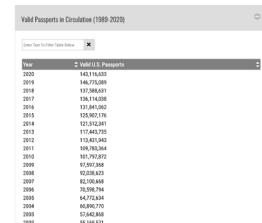
1) the number of applicants for a specific visa type by age in 2020?



a) histogram



b) line chart

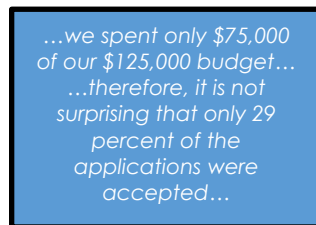


c) table

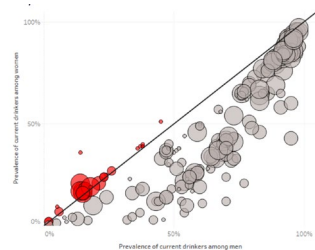
2) the relationship between the age of different embassies' electrical systems and the cost of annual repairs?



a) boxplot

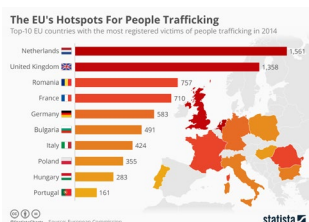


b) simple text

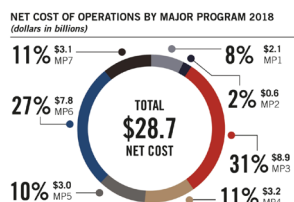


c) scatterplot

3) total annual energy consumption for a given country by type (i.e. natural gas, crude oil, solar, hydroelectric, etc.)?



a) bar chart



b) donut chart

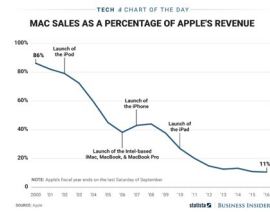
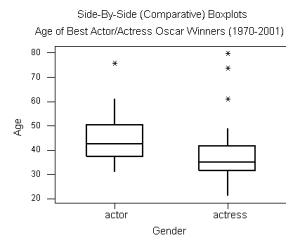


c) map

4) the shift in sales of a certain type of vehicle over a span of 30 years?

Valid Prospects in Circulation (1989-2020)

Year	Valid Prospects in Circulation
2020	142,174,025
2019	140,174,025
2018	137,588,031
2017	134,714,025
2016	129,497,062
2015	125,802,176
2014	121,313,241
2013	117,440,795
2012	113,440,795
2011	109,761,344
2010	107,761,344
2009	103,801,344
2008	102,801,344
2007	102,801,344
2006	102,801,344
2005	102,801,344
2004	102,801,344
2003	102,801,344
2002	102,801,344

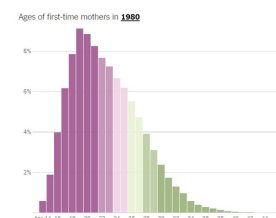
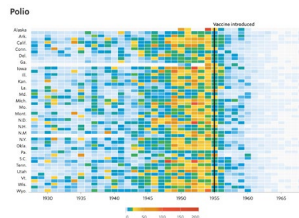
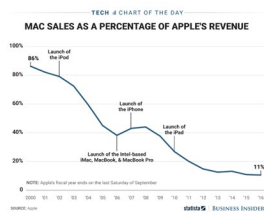


a) table

b) boxplot

c) line chart

5) the intensity in web traffic by hour and day per week?



a) line chart

b) heatmap

c) histogram

For answers and explanations, see the following page.

Answers:

- 1) To visualize **the number of applicants for a specific visa type by age in 2020**, it makes the most sense to use a **histogram**. Histograms are used to show the **distribution** of data. In this visualization, the x-axis would be different ages of applicant (or bins of ages). The y-axis would be the number of applicants.

You could also use a simple **table** with two columns. One column would show the applicant ages (or bins of ages). The other would show the number of applicants.

Line charts are useful for showing change over time, but this dataset is only for a single year.

a and **c** are acceptable answers.

- 2) To visualize **the relationship between the age of different embassies' electrical systems and the cost of annual repairs**, it makes the most sense to use a **scatterplot**. Scatterplots are used to show the **association** or relationship between two different variables. In this case, one axis could represent age in years, while the other could represent cost in dollars. Each embassy would be plotted as a single point on the chart.

Boxplots are used to show the distribution of data by quartile, as well as the median. But because a boxplot can only represent a single variable, it will not convey a relationship.

Finally, using text to convey an association like this for many data points is very inefficient.

c is the acceptable answer.

- 3) To visualize **total annual energy consumption for a given country by type**, it makes the most sense to use a **bar chart**. Bar charts are useful for comparing different **categories** of data, like different types of energy. In a horizontally oriented bar chart, the y-axis would show different categories of energy, and the x-axis would show a range of quantities.

It is possible to use a **donut chart** to convey this information, too, with a couple of caveats. First, donut charts emphasize the relationship of **part to whole**, so using one would imply that the *total* energy consumed is a significant statistic for the reader. Second, donut charts lose their efficacy

when representing more than about 3 categories, so if the dataset reflected 12 different types of energy, the chart could be difficult to read.

A map does not make sense for this dataset, since there is no obvious geographical underpinning to different types of energy consumed.

a and **b** are acceptable answers.

- 4)** To visualize **the shift in sales of a certain type of vehicle over a span of 30 years**, it makes the most sense to use a **line chart**. Line charts are used to show **change over time**, as in the ups and downs of the sales of a type of car or plane. The x-axis would show the year, and the y-axis would show the quantity.

You might be able to get away with using a table, but you would need to bin the time units to a much smaller overall number than 1 row per year in order to increase readability at a glance.

A boxplot could not represent this information because it is meant to show the distribution of a dataset, not a trend.

c is the acceptable answer.

- 5)** To visualize **the intensity in web traffic by hour and day per week**, it makes the most sense to use a **heatmap**. Heatmaps are used to show **density**, using color coding to indicate when the most intense activity has occurred. In a heatmap of web traffic, the x-axis could indicate time of day, and the y-axis could indicate day of the week. Each hour per day could be color coded to reflect a certain “volume” of users via a gradient.

A single line chart, with the x-axis representing time of day over the course of a whole week, would likely contain too much information to be easily interpreted. In theory, you could use it to show peaks, but it isn't advisable.

A histogram could not represent this information because it is meant to show the distribution of a dataset, not a trend.

b is the acceptable answer.

Additional resources

Data science glossary

Algorithm – An algorithm is a series of steps to accomplish a task (a set of directions that gets you from point A to point B, if you will). It can be a thought exercise, a series of mathematical expressions, or a piece of code (or pseudo-code). Algorithms are used in everyday life, and many industries. They are an essential building block of data science, analytics, and any other quantitative field.

API – An application programming interface (API) is an “entryway” to a computer system (such as a database) that allows you to access, retrieve, and edit its components. Most often APIs are used as data-communication mediums between applications. They are an essential part of scalable data-centric applications, research projects that are built around data, and anything else that requires automated data access.

Artificial Intelligence – Artificial intelligence (AI) is the apparent ability of machines to act “intelligently.”

Bayes Theorem – Bayes' theorem is used to calculate conditional probability of an event given the knowledge of conditions that might be related to the event. Conditional probability is the probability of an event occurring given another related event has already occurred.

Example: We want to know the probability of A(ge), given the diagnosis of C(ancer). This quantity is unknown to us and hard to estimate. We could use Bayes Theorem to do that, because we have the probability of C(ancer) given the A(ge), the probability of A(ge), and the probability of C(ancer) like so:

$$P(\text{Age given Cancer}) = P(\text{Cancer given Age}) * P(\text{Age})P(\text{Cancer})$$

Big Data – Big Data is a term that describes a large volume of data—both structured and unstructured. It is produced in such quantities that it cannot be ingested or processed by a single machine at once (even a large machine like a supercomputer), so special tools and techniques are needed to process and store it.

Note: This term is often misused; most data that is called Big Data is not really that big. True Big Data is produced by things like Google Searches, satellite imagery of the Earth, climate simulations used by weather services, sensor data generated by cell towers, etc.

Classification – Classification is a supervised machine learning method where the output variable is a category, such as “yes” or “no.”

Clustering – Clustering is an unsupervised machine learning method used to discover groupings that are inherent in the data.

Clickstream – A clickstream is a record of a user’s activity on the Internet, including every Web site that the user visits, how long the user was on a page or site, in what order the pages were visited, and newsgroups that the user participates in and even the e-mail addresses of mail that the user sends and receives.

Data governance – The management of the overall quality, integrity, relevance, and security of available data.

Data Lake – A data lake is a system or repository of data stored in its natural/raw format, usually object blobs or files.

Note: Data lakes and data warehouses are both widely used for storing big data, but they are not interchangeable terms. A data lake is a vast pool of raw data, the purpose for which is not yet defined. A data warehouse is a repository for structured, filtered data that has already been processed for a specific purpose. (Source: <https://www.talend.com/resources/data-lake-vs-data-warehouse/>)

Data Mining – Data mining is a study of extracting information from structured/unstructured data taken from various sources.

Data Science – Data science is an interdisciplinary field that combines elements of mathematics, statistics, logical thinking, programming, and a range of domain knowledge from various fields, which is used to solve day-to-day problems using a combination of methods and tools from the above disciplines.

Data Warehouse – A data warehouse is electronic storage of a large amount of data, which is designed for query and analysis. It is typically used to connect and analyze data from heterogeneous sources.

Note: Data lakes and data warehouses are both widely used for storing big data, but they are not interchangeable terms. A data lake is a vast pool of raw data, the purpose for which is not yet defined. A data warehouse is a repository for structured, filtered data that has already been processed for a specific purpose. (Source: <https://www.talend.com/resources/data-lake-vs-data-warehouse/>)

Data Visualization – Any attempt to make data more easily digestible by rendering it in a visual context (e.g., charting, graphing, etc.).

Database – A database is a structured collection of data. The collected information is organized in a way such that it is easily accessible by the computer. Databases are built and managed by using database programming languages.

Dataset – A dataset is a collection of data, which is organized into some type of data structure. Several characteristics define a dataset's structure and properties. These include the number and types of the attributes or variables, and various statistical measures applicable to them, such as standard deviation.

Deep Learning – Deep learning is a branch of machine learning that uses deep artificial neural networks. Artificial neural networks are types of machine learning algorithms that are built based on the idea of how a brain works with its neurons connected to each other and transmitting signals. Artificial neural networks are usually built in layers, those that have 3+ layers are considered “deep” and fall into the deep learning bucket.

Graph Analysis – Also known as **network analysis**. Graph analysis is an analysis of structures that model pairwise relationships between objects. The objects in graph analysis are called nodes. Their relationships (a.k.a. connections) are called edges. Graph analysis can be used across many industries, but the most common use case is analysis of social networks.

Machine Learning – Machine learning is the computational process wherein a machine “learns” and adjusts its behaviors based on feedback from data. Usually manifesting as an adaptable algorithm, machine learning helps computers predict outcomes without explicit human input.

Model – A model is a simplified replica of any real-life phenomenon / object at smaller scale. In quantitative disciplines, a model is usually a mathematical description of such a phenomenon / object.

Example: An example is a model to predict the level of education based on age. In a simplified way it looks like this:

Level of education = some quantity + some quantity * age (a linear model)

Although in real life there are a many more factors and variables, this model could potentially display a general trend.

Network Analysis – See graph analysis.

NoSQL Database – A NoSQL database provides storage and the ability to retrieve data that is modeled in means other than the tabular relations used in relational databases.

Open Data – Open data is data that can be freely used, shared and built-on by anyone, anywhere, for any purpose.

Outlier – An outlier is an observation that appears far away and diverges from an overall pattern in a sample.

Regression – Regression is a supervised learning method where the output variable is a real value, such as “amount” or “weight” and the input variable(s) is a real value, such as “size” or “height.” The output variable depends on the input variable(s). The relationship between the input and the output are usually recorded in a mathematical model.

Structured Data – Structured data is data that has been organized into a formatted repository, typically a database, so that its elements can be made addressable for more effective processing and analysis.

Supervised Learning – Supervised learning is a type of learning in which we teach or train the machine using data that is well labeled, meaning some data is already tagged with the correct answer class (or group to which it belongs).

Text Mining – Text mining is the process of converting unstructured text data into meaningful and actionable information.

Unstructured Data – Unstructured data is information that either does not have a pre-defined data model or is not organized in a pre-defined manner.

Unsupervised Learning – Unsupervised learning is the training of machines using information that is neither classified nor labeled and allowing the algorithm to act on that information without guidance.

Comparison: popular analysis tools

	R	Python	Excel
Learning curve	Steeper learning curve for people without a programming background	Can be easy to learn for people without a programming background	Easy to learn for any analyst
Automation	Once you write commands you won't have to re-do the work, just upload a new data set	Once you write commands you won't have to re-do the work, just upload a new data set	New data sets are not always plug and play with your analysis
Analyzing data	Lots of libraries contributed by a broad user community	Over 6,500 packages contributed by the community including top academics	Limited to any particular version
Speed	In-memory, only limited by your computer's RAM	In-memory, only limited by your computer's RAM	Can be slower unless an enterprise configuration is used
Type of data	Reads data of almost any type	Reads data of almost any type	Limited to xlsx and csv files unless macros are used
Compatibility	Compatible with almost any data output, storage or processing platform	Not as compatible as Python with some data architecture systems, may need custom-built interfaces	While macros enhance compatibility, Excel is comparatively limited

	R	Python	Excel
Data manipulation	Very flexible data manipulation	Very flexible data manipulation, augmented by numerous data processing and manipulation packages	Color-coded formulas can be easier to use but have a more limited functionality
Seeing data	Command line presentation unless visualization libraries are used	Spreadsheet-like view function that can be less intuitive to navigate	Easy to navigate spreadsheet
Graphics	Cutting edge graphics, however advanced coding and JavaScript knowledge may be necessary	Cutting edge graphics including dynamic visualizations, maps, network graphs, etc.	More limited options based on pre-set drop-down menus (unless macros are used)
Cost & platform	Free, any platform	Free, any platform	Hundreds of dollars, functionality on a Mac does not always mimic a PC

Comparison: popular visualization tools

Tool	Description
Microsoft Excel	<ul style="list-style-type: none">• Create basic chart types such as pie, line, bar, scatter and more• Charts created in Excel can easily be ported to PowerPoint and Word
Google Charts	<ul style="list-style-type: none">• Free and open source• Has more options than Excel (e.g., interactive, animated, and geospatial graphics)• Includes a rich gallery, fully customizable controls and dashboards, and HTML5
Tableau	<ul style="list-style-type: none">• Tool for creating powerful and insightful visuals• No programming required; drag and drop• Share and collaborate on premise or in the cloud• Platform can be used department or organization wide
R and R Studio	<ul style="list-style-type: none">• Free and open source• Programming tool• Mainly used for statistical analysis, but has sophisticated packages (code contributed by users) to create interactive dashboards as well
Python	<ul style="list-style-type: none">• Free and open source• Programming tool• You'll find libraries for practically every data visualization need
Power BI	<ul style="list-style-type: none">• Includes interactive visualizations and business intelligence capabilities• Simple interface• Create visualizations and dashboards

Additional reading & reference

Doing Data Science by Cathy O'Neil & Rachel Schutt

Data Science for Business by Foster Provost & Tom Fawcett

Data Smart by John W. Foreman

Mining the Social Web by Matthew A. Russell

Predictive Analytics by Eric Siegel

Analyzing the Analyzers by Harlan Harris, Sean Murphy and Marck Vaisman

Use cases

[BNY Mellon advances artificial intelligence tech across operations](#)

[Combining Satellite Imagery and Machine Learning to Predict Poverty](#)

[New York City uses “nudges” to reduce missed court dates](#)

[Proof of concept: Using predictive modeling to prioritize building inspections](#)

[Recruiting Chatbots in 2021: In-Depth Guide](#)

[Tactical Institute: Protecting people and communities with pre-emptive social media threat analytics](#)

[What Wal-Mart Knows About Customers' Habits](#)

Federal resources

[Data.gov](#) – “The home of the U.S. Government’s open data”

[Federal Data Strategy](#) - Provides a whole-of-government vision and offers guidance on how agencies should manage and use Federal data

[Federal Data Strategy: Data Ethics Framework](#) - Guides the data activities of agencies, providing the foundation for the ethical acquisition, management, and use of data for any federal purpose.

[Federal Data Strategy: Data Governance Playbook](#) – Helps agencies implement the Federal Data Strategy and fulfill the requirements of the Evidence Act by improving their organizational leadership for leveraging data as a strategic asset.

analytics.usa.gov – “a window into how people are interacting with the government online”