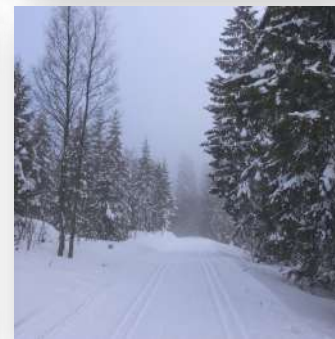
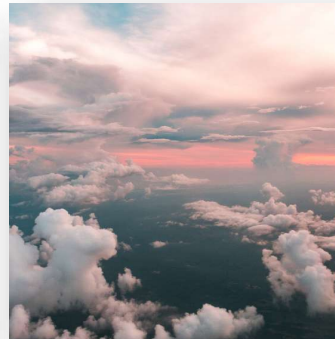


# Introduction to Data Science

Lecture as part of TERI-NORCE Research School //  
Funded by the Norwegian Ministry of Foreign Affairs



Dr. Michel d. S. Mesquita<sup>1,2,3</sup>



LinkedIn QR code

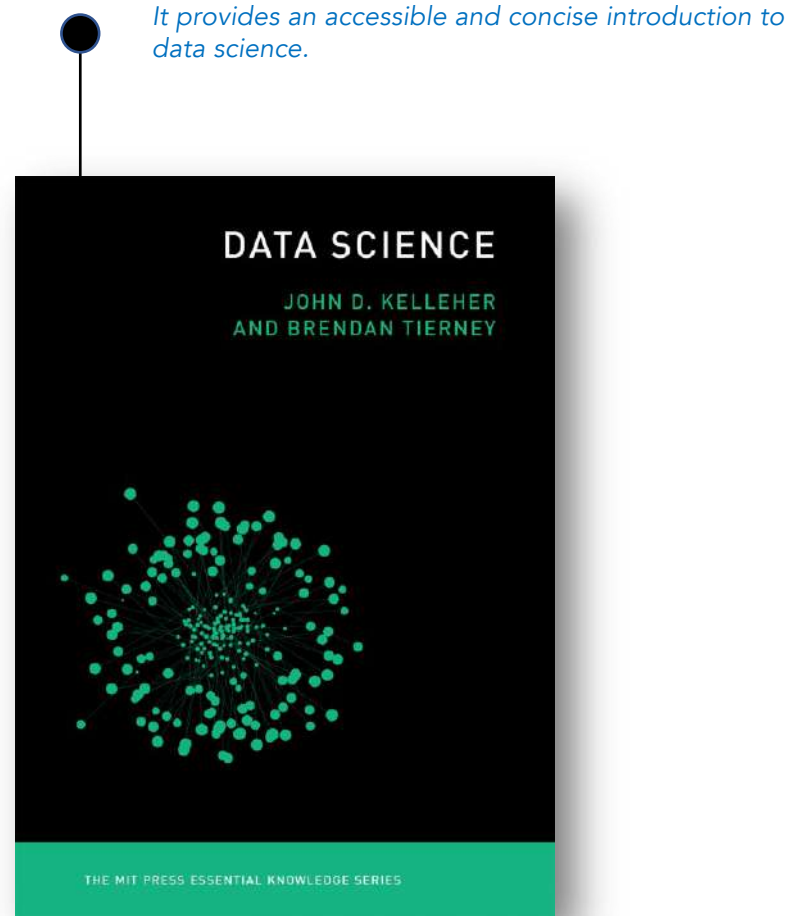
<sup>1</sup>NORCE Norwegian Research Centre, **Norway** / <sup>2</sup>Bjerknes Centre for Climate Research, **Norway** / <sup>3</sup>M2Lab.org, **Norway**



## Main reference //

This lecture is based on the book shown below  
'KT2018' hereafter

Kelleher and Tierney (2018) Data Science. Boston: The MIT Press Essential Knowledge Series. URL:  
<https://mitpress.mit.edu/books/data-science>



**Data science /**

«... a set of principles, problem definitions, algorithms, and processes for extracting nonobvious and useful patterns from large data sets»

**Machine learning (ML) /**

«... focuses on the design and evaluation of algorithms for extracting patterns from data»

**Data mining /**

«... generally deals with the analysis of structured data and often implies and emphasis on commercial applications»

KT2018



When to use data science? //

It has to be useful

When we deal with a large number of data examples or the patterns are too complex for humans to discover and extract manually



## Actionable insight /

Insight = the pattern should give us relevant information about the problem that is not obvious

Actionable = the insight we get should be something we can use in some way



Data science starts in the 1990s drawing knowledge from two historical fields: data collection and data analysis



### Data gathering /

Earlier marks of solstices and sunrise

Transactional data: advent of writing brought record keeping, e.g. Mesopotamia ~3200 BC

Nontransactional data: demographic data, e.g. earliest census in pharaonic Egypt ~3000 BC

Relational data model: Edgar F. Codd in 1970 publishes a paper on a model, which lays the foundation for modern databases, such as the *structure query language (SQL)*

Big data: today we have volume, variety, velocity; *NoSQL databases*

Data science starts in the 1990s drawing knowledge from two historical fields: data collection and data analysis



Thomas Bayes  
(Wikipedia)

### Data analysis /

Statistics: the branch of science that deals with the collection and analysis of data

Probability theory: 17th and 18th centuries, e.g. work of Blaise Pascal, Jakob Bernoulli, Thomas Bayes, and others. Probability distributions became the new tool to move beyond descriptive statistics

Statistical learning and modern data science: 19th century; new developments in probability theory facilitated statistical learning; e.g. Works by Pierre Laplace, Carl Friedrich Gauss, among others; method of least squares, which laid the foundation for linear regression and later artificial neural network

Data visualisation and exploratory data analysis: 1780-1820 William Playfair invented statistical graphics, such as the line chart, area chart, bar chart, pie chart; these led the foundation for modern methods, such as the t-distributed stochastic neighbor embedding (t-SNE) algorithm



Brief history of data science //  
The role of women

Diagram for the computation by the Engine of the Numbers of Bernoulli. See Note G. (page 729 of my.)

Number of Operations.	Number of Variables.	Variables used.	Variables receiving results.	Indication of change in the value of any Variable.	Statement of Results.	Data.										Working Variables.				Result Variables.			
						$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$	$x_{10}$	$y_1$	$y_2$	$y_3$	$y_4$	$z_1$	$z_2$	$z_3$	$z_4$
1.	1.	$x_1 \times x_2$	$x_3$	$x_3 = x_1 \times x_2$	$x_3 = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
2.	1.	$x_3 - x_1$	$x_4$	$x_4 = x_3 - x_1$	$x_4 = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
3.	1.	$x_4 + x_2$	$x_5$	$x_5 = x_4 + x_2$	$x_5 = 2 + 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
4.	1.	$x_5 - x_3$	$x_6$	$x_6 = x_5 - x_3$	$x_6 = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
5.	1.	$x_6 \times x_2$	$x_7$	$x_7 = x_6 \times x_2$	$x_7 = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
6.	1.	$x_7 - x_5$	$x_8$	$x_8 = x_7 - x_5$	$x_8 = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
7.	1.	$x_8 \times x_2$	$x_9$	$x_9 = x_8 \times x_2$	$x_9 = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
8.	1.	$x_9 - x_7$	$x_{10}$	$x_{10} = x_9 - x_7$	$x_{10} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
9.	1.	$x_{10} \times x_2$	$x_{11}$	$x_{11} = x_{10} \times x_2$	$x_{11} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
10.	1.	$x_{11} - x_9$	$x_{12}$	$x_{12} = x_{11} - x_9$	$x_{12} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
11.	1.	$x_{12} \times x_2$	$x_{13}$	$x_{13} = x_{12} \times x_2$	$x_{13} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
12.	1.	$x_{13} - x_{11}$	$x_{14}$	$x_{14} = x_{13} - x_{11}$	$x_{14} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
13.	1.	$x_{14} \times x_2$	$x_{15}$	$x_{15} = x_{14} \times x_2$	$x_{15} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
14.	1.	$x_{15} - x_{13}$	$x_{16}$	$x_{16} = x_{15} - x_{13}$	$x_{16} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
15.	1.	$x_{16} \times x_2$	$x_{17}$	$x_{17} = x_{16} \times x_2$	$x_{17} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
16.	1.	$x_{17} - x_{15}$	$x_{18}$	$x_{18} = x_{17} - x_{15}$	$x_{18} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
17.	1.	$x_{18} \times x_2$	$x_{19}$	$x_{19} = x_{18} \times x_2$	$x_{19} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
18.	1.	$x_{19} - x_{17}$	$x_{20}$	$x_{20} = x_{19} - x_{17}$	$x_{20} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
19.	1.	$x_{20} \times x_2$	$x_{21}$	$x_{21} = x_{20} \times x_2$	$x_{21} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
20.	1.	$x_{21} - x_{19}$	$x_{22}$	$x_{22} = x_{21} - x_{19}$	$x_{22} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
21.	1.	$x_{22} \times x_2$	$x_{23}$	$x_{23} = x_{22} \times x_2$	$x_{23} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
22.	1.	$x_{23} - x_{21}$	$x_{24}$	$x_{24} = x_{23} - x_{21}$	$x_{24} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
23.	1.	$x_{24} \times x_2$	$x_{25}$	$x_{25} = x_{24} \times x_2$	$x_{25} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
24.	1.	$x_{25} - x_{23}$	$x_{26}$	$x_{26} = x_{25} - x_{23}$	$x_{26} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
25.	1.	$x_{26} \times x_2$	$x_{27}$	$x_{27} = x_{26} \times x_2$	$x_{27} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
26.	1.	$x_{27} - x_{25}$	$x_{28}$	$x_{28} = x_{27} - x_{25}$	$x_{28} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
27.	1.	$x_{28} \times x_2$	$x_{29}$	$x_{29} = x_{28} \times x_2$	$x_{29} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
28.	1.	$x_{29} - x_{27}$	$x_{30}$	$x_{30} = x_{29} - x_{27}$	$x_{30} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
29.	1.	$x_{30} \times x_2$	$x_{31}$	$x_{31} = x_{30} \times x_2$	$x_{31} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
30.	1.	$x_{31} - x_{29}$	$x_{32}$	$x_{32} = x_{31} - x_{29}$	$x_{32} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
31.	1.	$x_{32} \times x_2$	$x_{33}$	$x_{33} = x_{32} \times x_2$	$x_{33} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
32.	1.	$x_{33} - x_{31}$	$x_{34}$	$x_{34} = x_{33} - x_{31}$	$x_{34} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
33.	1.	$x_{34} \times x_2$	$x_{35}$	$x_{35} = x_{34} \times x_2$	$x_{35} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
34.	1.	$x_{35} - x_{33}$	$x_{36}$	$x_{36} = x_{35} - x_{33}$	$x_{36} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
35.	1.	$x_{36} \times x_2$	$x_{37}$	$x_{37} = x_{36} \times x_2$	$x_{37} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
36.	1.	$x_{37} - x_{35}$	$x_{38}$	$x_{38} = x_{37} - x_{35}$	$x_{38} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
37.	1.	$x_{38} \times x_2$	$x_{39}$	$x_{39} = x_{38} \times x_2$	$x_{39} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
38.	1.	$x_{39} - x_{37}$	$x_{40}$	$x_{40} = x_{39} - x_{37}$	$x_{40} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
39.	1.	$x_{40} \times x_2$	$x_{41}$	$x_{41} = x_{40} \times x_2$	$x_{41} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
40.	1.	$x_{41} - x_{39}$	$x_{42}$	$x_{42} = x_{41} - x_{39}$	$x_{42} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
41.	1.	$x_{42} \times x_2$	$x_{43}$	$x_{43} = x_{42} \times x_2$	$x_{43} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
42.	1.	$x_{43} - x_{41}$	$x_{44}$	$x_{44} = x_{43} - x_{41}$	$x_{44} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
43.	1.	$x_{44} \times x_2$	$x_{45}$	$x_{45} = x_{44} \times x_2$	$x_{45} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
44.	1.	$x_{45} - x_{43}$	$x_{46}$	$x_{46} = x_{45} - x_{43}$	$x_{46} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
45.	1.	$x_{46} \times x_2$	$x_{47}$	$x_{47} = x_{46} \times x_2$	$x_{47} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
46.	1.	$x_{47} - x_{45}$	$x_{48}$	$x_{48} = x_{47} - x_{45}$	$x_{48} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
47.	1.	$x_{48} \times x_2$	$x_{49}$	$x_{49} = x_{48} \times x_2$	$x_{49} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
48.	1.	$x_{49} - x_{47}$	$x_{50}$	$x_{50} = x_{49} - x_{47}$	$x_{50} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
49.	1.	$x_{50} \times x_2$	$x_{51}$	$x_{51} = x_{50} \times x_2$	$x_{51} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
50.	1.	$x_{51} - x_{49}$	$x_{52}$	$x_{52} = x_{51} - x_{49}$	$x_{52} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
51.	1.	$x_{52} \times x_2$	$x_{53}$	$x_{53} = x_{52} \times x_2$	$x_{53} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
52.	1.	$x_{53} - x_{51}$	$x_{54}$	$x_{54} = x_{53} - x_{51}$	$x_{54} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
53.	1.	$x_{54} \times x_2$	$x_{55}$	$x_{55} = x_{54} \times x_2$	$x_{55} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
54.	1.	$x_{55} - x_{53}$	$x_{56}$	$x_{56} = x_{55} - x_{53}$	$x_{56} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
55.	1.	$x_{56} \times x_2$	$x_{57}$	$x_{57} = x_{56} \times x_2$	$x_{57} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
56.	1.	$x_{57} - x_{55}$	$x_{58}$	$x_{58} = x_{57} - x_{55}$	$x_{58} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
57.	1.	$x_{58} \times x_2$	$x_{59}$	$x_{59} = x_{58} \times x_2$	$x_{59} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
58.	1.	$x_{59} - x_{57}$	$x_{60}$	$x_{60} = x_{59} - x_{57}$	$x_{60} = 2 - 1$	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2	1	2
59.	1.	$x_{60} \times x_2$	$x_{61}$	$x_{61} = x_{60} \times x_2$	$x_{61} = 2 \times 1$	1	2	1	2	1	2	1	2	1	2	1							

Data science starts in the 1990s drawing knowledge from two historical fields: data collection and data analysis



Claude Shannon  
(Wikipedia)

### Data analysis (cont.) /

20th century

Karl Pearson developed modern hypothesis testing; R. A. Fisher developed statistical methods for multivariate analysis and maximum likelihood estimate

Alan Turing in the Second World War led the invention of the electronic computer

Warren McCulloch and Walter Pitts in 1943 proposed the first mathematical model of a neural network

Claude Shannon in 1948 published «A Mathematical Theory of Communication» and founded *information theory*



ON COMPUTABLE NUMBERS, WITH AN APPLICATION TO  
THE ENTSCHEIDUNGSPROBLEM

By A. M. TURING.

[Received 28 May, 1936.—Read 12 November, 1936.]

The “computable” numbers may be described briefly as the real numbers whose expressions as a decimal are calculable by finite means. Although the subject of this paper is ostensibly the computable numbers, it is almost equally easy to define and investigate computable functions of an integral variable or a real or computable variable, computable predicates, and so forth. The fundamental problems involved are, however, the same in each case, and I have chosen the computable numbers for explicit treatment as involving the least cumbersome technique. I hope shortly to give an account of the relations of the computable numbers, functions, and so forth to one another. This will include a development of the theory of functions of a real variable expressed in terms of computable numbers. According to my definition, a number is computable if its decimal can be written down by a machine.

In §§9, 10 I give some arguments with the intention of showing that the computable numbers include all numbers which could naturally be regarded as computable. In particular, I show that certain large classes of numbers are computable. They include, for instance, the real parts of all algebraic numbers, the real parts of the zeros of the Bessel functions, the numbers  $\pi$ ,  $e$ , etc. The computable numbers do not, however, include all definable numbers, and an example is given of a definable number which is not computable.

Although the class of computable numbers is so great, and in many ways similar to the class of real numbers, it is nevertheless enumerable. In §8 I examine certain arguments which would seem to prove the contrary. By the correct application of one of these arguments, conclusions are reached which are superficially similar to those of Gödel†. These results

† Gödel, “Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme, I”, *Monatsh. Math. Phys.*, 38 (1931), 173–198.

MIND  
A QUARTERLY REVIEW  
OF  
PSYCHOLOGY AND PHILOSOPHY

1.—COMPUTING MACHINERY AND  
INTELLIGENCE

By A. M. TURING

1. *The Imitation Game.*

I propose to consider the question, ‘Can machines think?’ This should begin with definitions of the meaning of the terms ‘machine’ and ‘think’. The definitions might be framed so as to reflect so far as possible the normal use of the words, but this attitude is dangerous. If the meaning of the words ‘machine’ and ‘think’ are to be found by examining how they are commonly used it is difficult to escape the conclusion that the meaning and the answer to the question, ‘Can machines think?’ is to be sought in a statistical survey such as a Gallup poll. But this is absurd. Instead of attempting such a definition I shall replace the question by another, which is closely related to it and is expressed in relatively unambiguous words.

The new form of the problem can be described in terms of a game which we call the ‘imitation game’. It is played with three people, a man (A), a woman (B), and an interrogator (C) who may be of either sex. The interrogator stays in a room apart from the other two. The object of the game for the interrogator is to determine which of the other two is the man and which is the woman. He knows them by labels X and Y, and at the end of the game he says either ‘X is A and Y is B’ or ‘X is B and Y is A’. The interrogator is allowed to put questions to A and B thus:

C: Will X please tell me the length of his or her hair?  
Now suppose X is actually A, then A must answer. It is A’s



Alan Turing  
(Wikipedia)

Turing (1936, 1950)  
Papers by Alan Turing on the topics of computable numbers and artificial intelligence.

Data science starts in the 1990s drawing knowledge from two historical fields: data collection and data analysis



Nils Nilsson  
(Wikipedia)

### Data analysis (cont.) /

Evelyn Fix and Joseph Hodges proposed a model for discriminatory analysis (classification or pattern-recognition problem), which became the basis for modern nearest-neighbor models

Establishment of the field of *artificial intelligence* at a workshop in Dartmouth College

The term *machine learning* was beginning to be used to describe programs that gave a computer the ability to learn from data

Nils Nilsson's book *Learning Machines* in 1965 showed how neural networks could be used to learn linear models for classification

Data science starts in the 1990s drawing knowledge from two historical fields: data collection and data analysis



Ensemble fatigue  
(Benestad et al., 2017,  
Nature Climate Change)

Data analysis (cont.) /

Earl Hunt, Janet Marin, and Philip Stone in 1966 developed the concept-learning system framework, which was the progenitor of an important family of ML algorithms that induced decision-tree models

A number of independent researchers developed and published earlier versions of the *k-means clustering* algorithm

Today: some of the most important developments include *ensemble models*, where predictions are made using a set of models, and *deep-learning neural networks*, which have multiple layers of neurons (with applications in machine vision, natural-language processing,...)

1990s: the term «data science» came to prominence in discussions of the need for statisticians to join with computer scientists to bring mathematical rigor to computational analysis of large data sets

1997: C.F. Jeff Wu's public lecture «Statistics = Data Science?»

2001: William S. Cleveland published an action plan for creating a university department in the field of data science

### Also in 2001 //

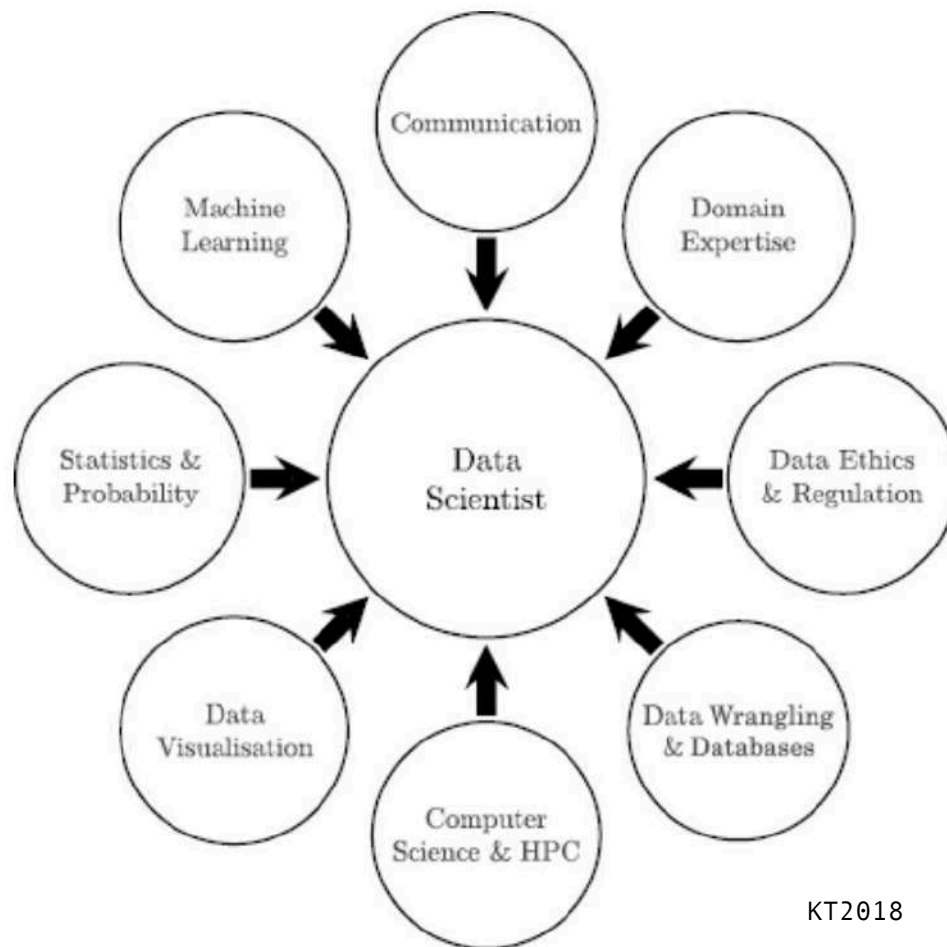
Leo Breiman published the paper «Statistical Modeling: the two cultures» /

His distinction between a statistical focus on models that explain the data versus an algorithmic focus on models that can actually predict the data highlights a core difference between statisticians and ML researchers

Today //

The role of a data scientist has become very broad

It is difficult for an individual to master all of the skill areas



KT2018



### **Datum /**

Or piece of information is an abstraction of a real-world entity (person, object, or event)

The terms *variable*, *feature*, and *attribute* are often used interchangeably to denote an individual abstraction

Each entity is typically described by a number of attributes. For instance, a book might have the following attributes: author, title, topic, genre, publisher, price, etc.



What are data? What is a data set? //

**Data set /**

It consists of the data relating to a collection of entitites, with each entity described in terms of a set of attributes

In its most basic form a data set is organized in an  $n*m$  data matrix called the *analytics record*

KT2018

-----						
ID	Title	Author	Year	Cover	Edition	Price
1	Emma	Austen	1815	Paperback	20th	\$5.75
2	Dracula	Stoker	1897	Hardback	15th	\$12.00
3	Ivanhoe	Scott	1820	Hardback	8th	\$25.00
4	Kidnapped	Stevenson	1886	Paperback	11th	\$5.00

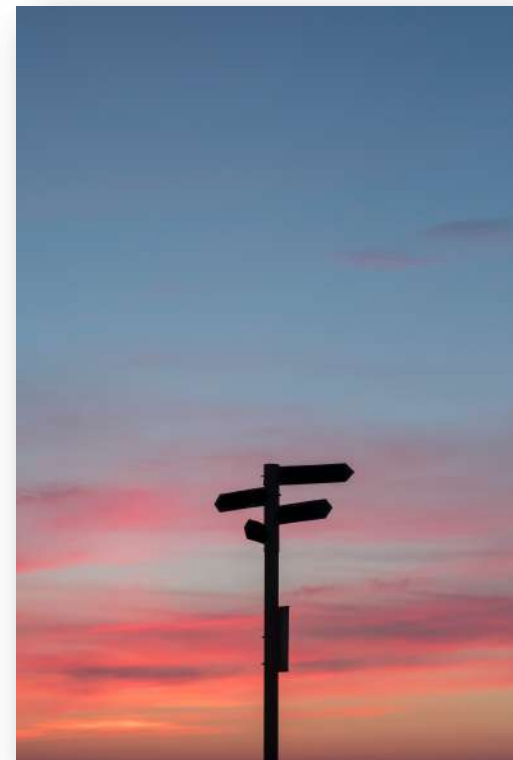
### Choice of attributes /

How do you choose the most appropriate attributes?


### Different attribute types /

Numeric, nominal and ordinal.

The data type of an attribute affects the methods we can use to analyze and understand the data, including both the basic statistics we can use to describe the distribution of values that an attribute takes and the more complex algorithms we use to identify the patterns of relationships between attributes




### Structured data

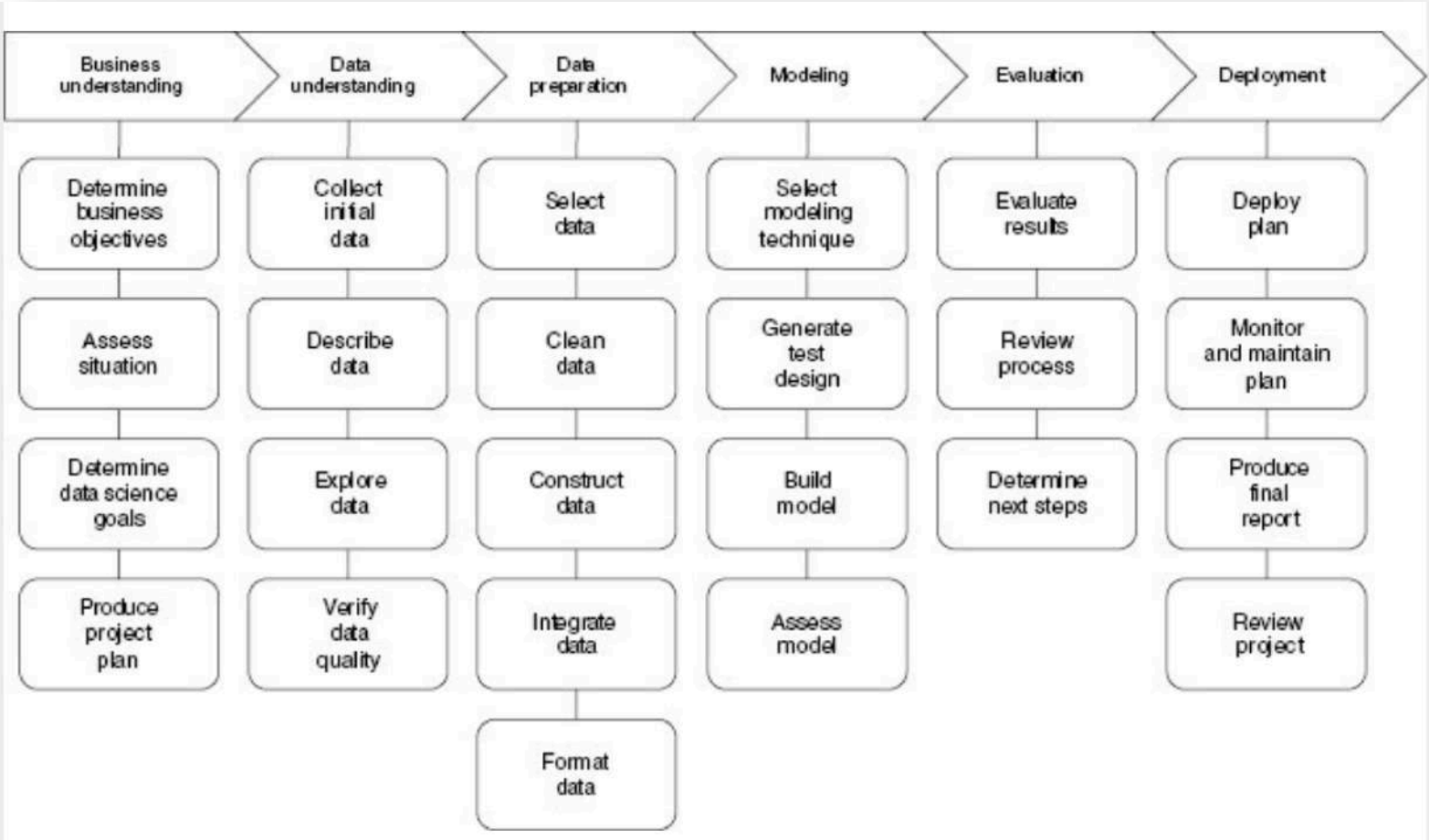


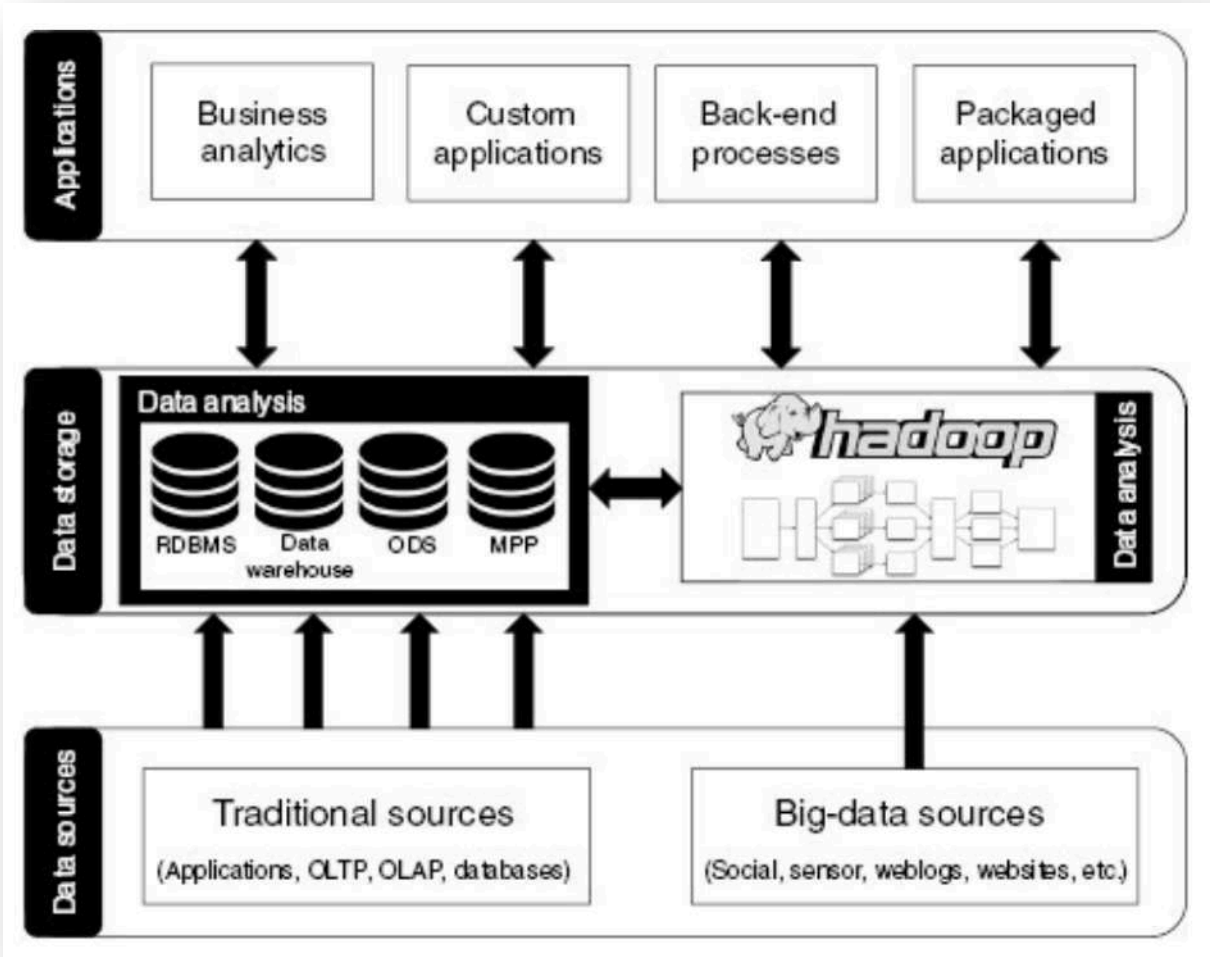
Can be stored in a table, and every instance in the table has the same structure (i.e., set of attributes)

### Unstructured data



Each instance in the data set may have its own internal structure, and this structure is not necessarily the same in every instance





**ML involves a two step process /**

Identify useful patterns given some data (modelling). Some methods include:

- Decision trees
- Regression models
- Neural networks

Once a model is created, it is used for analysis





### Supervised learning /

Learn a function that maps from the values of the attributes to the value of another (target) attribute.

«Supervised» because each of the instances in the data set lists both the input values and the output (target) value for each instance.

Each instance in the data set must be labeled with the value of the target attribute.

### Unsupervised learning /

There is no target attribute. The algorithm has the more general task of looking for regularities in the data.

- Cluster analysis
- Euclidean distance

«Supervised» because each of the instances in the data set lists both the input values and the output (target) value for each instance.

Each instance in the data set must be labeled with the value of the target attribute.

# Machine Learning Algorithms *(sample)*

	<u>Unsupervised</u>	<u>Supervised</u>
<u>Continuous</u>	<ul style="list-style-type: none"><li>• Clustering &amp; Dimensionality Reduction<ul style="list-style-type: none"><li>○ SVD</li><li>○ PCA</li><li>○ K-means</li></ul></li></ul>	<ul style="list-style-type: none"><li>• Regression<ul style="list-style-type: none"><li>○ Linear</li><li>○ Polynomial</li></ul></li><li>• Decision Trees</li><li>• Random Forests</li></ul>
<u>Categorical</u>	<ul style="list-style-type: none"><li>• Association Analysis<ul style="list-style-type: none"><li>○ Apriori</li><li>○ FP-Growth</li></ul></li><li>• Hidden Markov Model</li></ul>	<ul style="list-style-type: none"><li>• Classification<ul style="list-style-type: none"><li>○ KNN</li><li>○ Trees</li><li>○ Logistic Regression</li><li>○ Naive-Bayes</li><li>○ SVM</li></ul></li></ul>




Hypothetical dataset

Ten years of accumulated All India Rainfall in mm for the monsoon season and pre-monsoon DJF mean ENSO values for the corresponding year

Definitions /

Let  $x_i$  be the value of the explanatory variable for the  $i$ th datapoint, and  $Y_i$  the random variable representing the response for the same datapoint



All_India_Rainfall (mm)	El_Nino_Southern_Oscillation
1030	-1.50
950	0.50
980	0.80
1099	-1.80
1100	-1.80
1140	-1.90
999	0.90
950	0.85
988	0.88

## Assumptions /

$$E(Y_i) = \alpha + \beta x_i$$

The mean of the response variable  $Y_i$  depends on the value of  $x_i$  of the explanatory variable in a linear fashion

$$Y_i = \alpha + \beta x_i + \epsilon_i$$

The variation of the response variable  $Y_i$  about the mean is represented by a random variable  $\epsilon_i$

$$E(\epsilon_i) = 0$$

Note that for the mean of  $Y_i$  to be as equation 1 above,  $\epsilon_i$  has to be zero

$$E(\epsilon_i) \sim N(0, \sigma^2)$$

The variance of  $\epsilon_i$  is the same for all values of the explanatory variable (error variance). Also the different  $\epsilon_i$  ( $i=1,2,\dots,n$ ), are independent of each other

$$aX + b \sim N(a\mu + b, a^2\sigma^2)$$


Property of distributions. Also, assuming that  $Y_i$  is a fixed quantity  $\alpha + \beta x_i$  plus a random variable  $N(0, \sigma^2)$ , then:

$$Y_i = N(\alpha + \beta x_i, \sigma^2)$$

This is the simple linear regression model



Sometimes we need to transform the data. Provided that  $x > 1$ , we can use the following transformation alternatives:

$$\dots, x^{-2}, x^{-1}, x^{-1/2}, \log x, x^{1/2}, x^1, x^2, \dots$$


Reduce the high values in a dataset relative to low values

Stretch out high values relative to low ones



deeplearning.ai

# Introduction to Deep Learning

---

## What is a Neural Network?

*Deeplearning.ai*  
*Video by Andrew Ng*  
*URL: <https://youtu.be/n1l-9lIMW7E>*

# MATLAB® Tech Talks

*MATLAB Tech Talks*  
*Video by Shyamal Patel*

URL: <https://www.youtube.com/watch?v=3cSjsTKtN9M>

TERI-NORCE CRS, October 2019 - New Delhi

**Example //**  
Recognising objects with Convolutional Neural Networks  
Michel's solution, based on video from Deeplearning.ai



contributed articles

DOI:10.1145/2500499

**Big data promises automated actionable knowledge creation and predictive models for use by both humans and computers.**

BY VASANT DHAR

# Data Science and Prediction

USE OF THE term “data science” is increasingly common, as is “big data.” But what does it mean? Is there something unique about it? What skills do “data scientists” need to be productive in a world deluged by data? What are the implications for scientific inquiry? Here, I address these questions from the perspective of predictive modeling.

The term “science” implies knowledge gained through systematic study. In one definition, it is a systematic enterprise that builds and organizes knowledge in the form of testable explanations and predictions.<sup>1)</sup> Data science might therefore imply a focus involving data and, by extension, statistics, or the systematic study of the organization, properties, and analysis of data and its role in inference, including our confidence in the inference. Why then do we need a new term like data science when we have had statistics for centuries? The fact that we now have huge amounts of data should not in and of itself justify the need for a new term.

The short answer is data science is different from statistics and other existing disciplines in several important ways. To start, the raw material, the “data”

part of data science, is increasingly heterogeneous and unstructured—text, images, video—often emanating from networks with complex relationships between their entities. Figure 1 outlines the relative expected volumes of unstructured and structured data from 2008 to 2015 worldwide, projecting a difference of almost 200 petabytes (PB) in 2015 compared to a difference of 50 PB in 2012. Analysis, including the combination of the two types of data, requires integration, interpretation, and sense making that is increasingly derived through tools from computer science, linguistics, econometrics, sociology, and other disciplines. The proliferation of markup languages and tags is designed to let computers interpret data automatically, making them active agents in the process of decision making. Unlike early markup languages (such as HTML) that emphasized the display of information for human consumption, most data generated by humans and computers today is for consumption by computers; that is, computers increasingly do background work for each other and make decisions automatically. This scalability in decision making has become possible because of big data that serves as the raw material for the creation of new knowledge. Watson, IBM’s “Jeopardy!” champion, is a prime illustration of an emerging machine intelligence fueled by data and state-of-the-art analytics.

>> key insights

- Data science is the study of the generalizable extraction of knowledge from data.
- A common epistemic requirement in assessing whether new knowledge is actionable for decision making is its predictive power, not just its ability to explain the past.
- A data scientist requires an integrated skill set spanning mathematics, machine learning, artificial intelligence, statistics, databases, and optimization, along with a deep understanding of the craft of problem formulation to engineer effective solutions.

64

COMMUNICATIONS OF THE ACM

DECEMBER 2013 • VOL. 54 • NO. 12

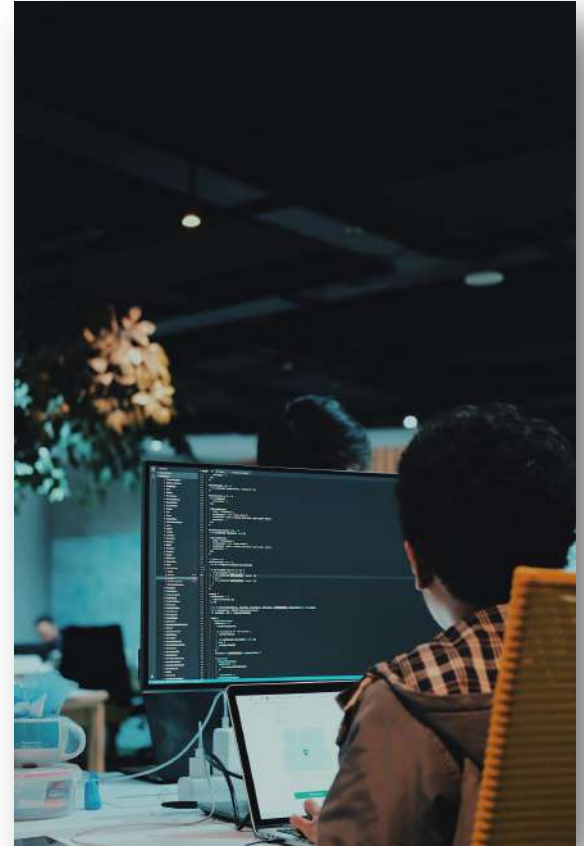
A photograph showing a person in a plaid shirt gesturing with their hands while talking to another person. In the foreground, a laptop and a smartphone are on a wooden table. The background is slightly blurred, showing another person.

A composite image. On the left, a close-up of a person's face looking upwards with a glowing blue light effect. On the right, a snippet of Python code for a class named 'RequestHandler'.

TERI-NORCE CRS, October 2019 - New Delhi

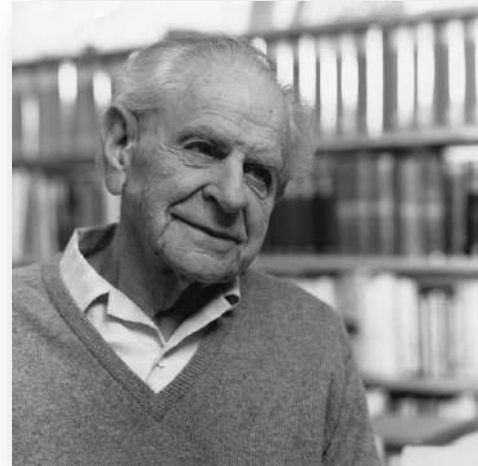
31

In which ways is data  
science different from  
statistics?





How much do you  
agree/disagree with Karl  
Popper's argument for  
evaluating a theory and  
scientific progress?



● *Karl Popper*  
(Wikipedia)

What powerful capability  
is machine learning able  
to give us?



*«...simpler models are more likely to hold up on future observations than more complex ones, all else being equal» (p.68)*

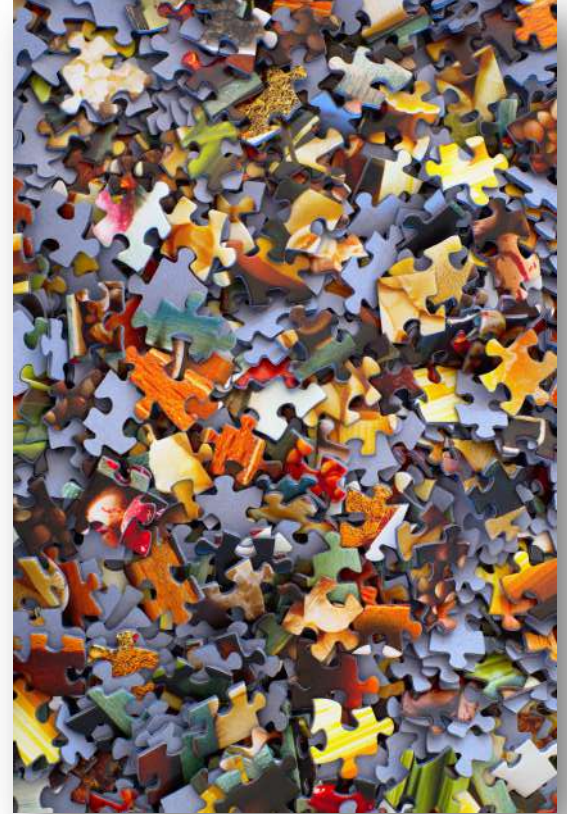
How much do you  
agree/disagree with this  
statement?



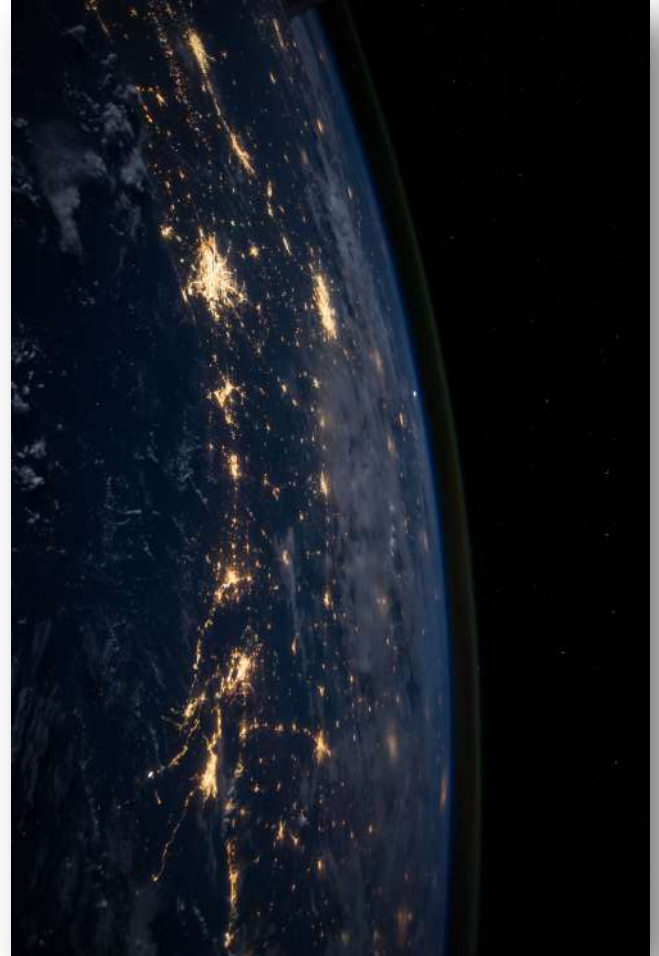
Why are problem  
formulation skills so  
important?



Where do errors come from?



How does the Internet  
contribute to inexpensive  
large-scale randomized  
experiments?



Why is earth science  
considered to be one of  
the «data-starved areas of  
inquiry»?



What should universities  
do to improve their  
students' data science  
skills?





How would you like to  
integrate data science in  
your own research work?



Thank you for your attention!  
Feel free to get in touch with us //

mido@norceresearch.no  
+47 4760 2340

In conclusion //

Data science has become a growing field, thanks to many past researchers in the field. Climate science will benefit even more in the years to come, as more and more researchers and students learn how to apply different algorithms to solve climate problems.

We hope this Climate Research School can inspire you as well! In the next days, you will be learning several tools that can become part of your professional toolbox!

