

# What makes Data Science different? A discussion involving Statistics2.0 and Computational Sciences

Christophe Ley · Stéphane P.A. Bordas

Received: date / Accepted: date

**Abstract** Data Science is today one of the main buzzwords be it in business, industrial or academic settings. Machine learning, experimental design, data-driven modelling are all, undoubtedly, rising disciplines if one goes by the soaring number of research papers and patents appearing each year. The prospect of becoming a “Data Scientist” appeals to many. A discussion panel organised as part of the European Data Science Conference (European Association for Data Science (EuADS)<sup>1</sup>) asked the question: “What makes Data Science different?” In this paper we give our own, personal and multi-faceted view on this question, from an engineering and a statistics perspective. In particular, we compare Data Science to Statistics and discuss the connection between Data Science and Computational Science.

**Keywords** High-dimensional statistics · Interdisciplinary research · Data science · Computational science · Machine learning · Scientific computing · Data Analytics · Data-driven modelling · Modelling ·

S. P. A. Bordas

University of Luxembourg, Institute of Computational Engineering, 6 Avenue de la Fonte, 4362 Esch-sur-Alzette  
Cardiff University, Institute of Mechanics and Advanced Materials

University of Western Australia, School of Mechanical and Chemical Engineering, Intelligent Systems for Medicine Lab, The University of Western Australia (M050) 35 Stirling Highway CRAWLEY WA 6009, Australia

Tel.: +352-621-131048

E-mail: stephane.bordas@alum.northwestern.edu

C. Ley

Ghent University, Department of Applied Mathematics, Computer Science and Statistics

Krijgslaan 281, Campus Sterre, S9, B-9000 Ghent

Tel.L +32 92644908

Fax: +32 92644995

E-mail: christophe.ley@ugent.be

<sup>1</sup> <https://euads.org/edsc/>

Applied Mathematics · Simulation · Digital Twins · Training · Education · Research

## 1 Introduction

According to IBM, 90% of the data available today has been generated over the last two years [1]. We have been experiencing a data-flood, fuelled by a surge in (mobile) computing power which has enabled the creation of devices which can create, collect, store and transfer increasingly complex and large data sets. This accelerated data-gathering ability has been drastically changing the world of science and business. The “internet of things” and wearable technologies densely maculate our world with digital footprints. These massive amounts of data are continuously being gathered in geography, geophysics, medicine, genetics, social science (media), finance, climatology and engineering. Evidence suggests that the intensity of this surge will only increase with time. We are living in the “Big Data” era and this yet ill-defined concept is now ubiquitous, be it in science, business, healthcare, media, industry, business, politics or sports. The challenges posed by the Big Data phenomenon are numerous, and the discipline known as “Data Science” may well be a natural consequence of the data outpour we have been witnessing.

But what does Data Science actually stand for? What makes it different from other, well-established disciplines? Why has it become so popular over the past years? Is Data Science merely Statistics? Is it Computer Science, Machine Learning? Wikipedia provides the following answer to the first question:

*Data science, also known as data-driven science, is an interdisciplinary field about scientific methods, processes and systems to extract knowledge*

or insights from data in various forms, either structured or unstructured, similar to *Knowledge Discovery in Databases*. Wikipedia, accessed on February 23, 2017.

This definition, like many others, remains vague and is surely insufficient to differentiate this disciplines from its cousins. Many have attempted to define Data Science through articles [2,3], [4], [5] as well as numerous panel discussions at the highest level as well as conference/seminar presentations. There has also been significant discussions on teaching and education in data science [6], [7], approaches to building data science teams [8] and the use of data sciences for various applications ranging from social sciences [9] to the material genome initiative [10].

But, in spite of all these developments, what has really changed since the 1962 article entitled “*The future of data analysis*” by John W. Tukey<sup>2</sup> in *The Annals of Mathematical Statistics*?

We wish to contribute to this active discussion through the present paper, that is based on a panel discussion to which we contributed during the European Data Science Conference in November 2016 in Luxembourg. The originality of our approach is the combination of two apparently disjoint domains which Data Science may have already brought closer together, namely computational science and statistics.

Should the reader wish to delve more into the details of particular Big Data disciplines, the following review papers [1] and [11] are excellent sources of information.

## 2 A simple classification of data science approaches

Before looking at how data science approaches can be classified, let us first think of examples of typical data science problems. Data science answers sharp and quantitative questions such as:

Quantify: How many coffee bean futures should I order assuming the temperature in the tropics rises by 5 degrees? This is done using regression algorithms.

Detect anomaly: Has this credit card been stolen?

Classify and make predictions: Will this aircraft door fail within the next 2,000 flights? How likely is a returning customer to become a regular customer? Given images of a brain, what is the probability that the tumour is located within 10mm of an eloquent region of this brain?

<sup>2</sup> Note that the Cooley-Tukey algorithm is probably the most common Fast Fourier Transform algorithm, a common signal processing technique

Organize: How is the data organized? For example, clustering algorithms help organize data. This can be useful to predict behaviour and events.

Choose the next step: What innovation directions should this country follow next to maximise its GDP? These algorithms are known as “reinforcement learning” and can be used to control autonomous systems, for example self-driving cars or climate control systems. They learn by trial and error.

Clearly, attacking such problems in their full complexity requires a serious mathematical arsenal. The mathematical methods behind data science applications can seem mystical to the neophyte. We summarize here very briefly how we believe data science methods can be classified. We distinguish between bottom-up and top-down approaches.

In top-down approaches, a model is built which represents the information contained in the data. This is usually a statistical model, for example a regression approach, possibly Bayesian when information is scarce and sparse. With such a model in hand and statistical methods such as Monte-Carlo simulations, K-means clustering, Markov chains, decision trees, Bayesian inference, a data analyst can analyse data and make predictions. In practice, what makes a top-down data science algorithm successful is the craft with which the above statistical models are used in concert. We discuss some of this orchestration and how data science relates to statistics in the following section.

In bottom-up approaches, on the contrary, the starting point is the data and the model of this data is generated by a computer (and updated continuously as new data is acquired) to match observations. This approach is known as machine learning and its outcome is a statistical model which is able to account for complex relationships between entities. However powerful have machine learning methods become, (skilled) human intervention is still necessary to filter outliers, optimise the learning paradigm to ensure the accuracy of classifications, tune the parameters involved in the model, etc.

In fact, successful data science algorithms are usually a combination of top-down and bottom-up approaches. The top-down approach brings domain- (application-) knowledge which leads to significant savings in the computing power required by the bottom-up approaches, for example by accelerating classification.

### 3 How does Data Science relate to Statistics?

#### 3.1 Nomen est omen

The first author has recently asked the students of his Data Mining class what the word Data Science meant to them. After a long silence, the following answer came: “Data Science is the discipline that makes sense out of data”. For a statistician such an answer is surprising, as this is precisely what Statistics aims to do. What causes this difference in perception between professional statisticians and non(or not yet)-statisticians? The reason is simple: Data Science seems, just by its name, to be a more data-oriented area than Statistics. And more attractive. If you say to a random person on the street that you are a statistician, the typical reaction of that person is to think you are dealing with spreadsheets which can seem monotone. However, if you happen to say you are a data scientist, then that same person will have no clue about your job, yet he/she will have the feeling you are doing an exciting work. The core task is twice data analysis, but the marketing effect of the name Data Science is incontestable.

A similar effect happens to prevail also among people with an advanced understanding of data analysis. While Statistics appears to be a rigid field, filled with rules to follow and warnings of how to correctly quantify the uncertainty inherent to any data set, Data Science seems to invite theoreticians and practitioners to play around with data in an unrestricted way. This is, again, just a subjective impression.

#### 3.2 Statistics in the Big Data era

A core role of statistics is the quantification of the uncertainty accompanying any data analysis. Sir Ronald Fisher has laid a solid mathematical background for this endeavour in the beginning of the 20th century. Estimation, testing and regression procedures were devised on basis of this formalism. These methods, however, can no longer be blindly applied to 21st-century data which happen to be complex and occur in unprecedented quantities. We illustrate this statement through two classical statistical procedures:

- Linear regression: suppose we are interested in modeling the relationship between a one-dimensional outcome variable  $Y$  and  $p$  one-dimensional predictors  $X_1, \dots, X_p$ , and we have good reasons to believe the relationship to be of the form

$$Y = \beta_0 + \sum_{i=1}^p \beta_i X_i + \epsilon,$$

where  $\epsilon$  is an error term (typically assumed to follow a normal distribution) and  $\beta_0, \beta_1, \dots, \beta_p$  are the regression parameters we need to estimate. The standard solution to this estimation problem is least squares estimation. This approach works very well as long as the number  $n$  of observations  $((Y_1, X_{11}, \dots, X_{1p}), \dots, (Y_n, X_{n1}, \dots, X_{np}))$  is larger than the dimension  $p$ . However, in many datasets nowadays the situation is rather reversed, with  $p$  being larger than  $n$ . Think of genetics, where every single gene should in principle be taken into account to measure the impact of a new treatment. Least squares estimation breaks down in such a context because the empirical covariance matrix is no longer invertible.

As a reaction, variable selection methods have been proposed. This idea is rooted on the belief in sparsity: the majority of predictors, here genes, shall only have a very small, irrelevant impact on the outcome variable, hence should be discarded. Variable selection does precisely this: it focusses on a small number of predictors that really do matter in the linear regression. Linear regression combined with variable selection can deal with  $p > n$  situations. The perhaps most famous example is the so-called Lasso regression of [12].

- Hypothesis testing: suppose we have  $n$  data points  $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})'$ ,  $i = 1, \dots, n$  of dimension  $p$  and we wish to perform a typical hypothesis testing problem of the form  $\mathcal{H}_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$  versus the alternative  $\mathcal{H}_1 : \boldsymbol{\mu} \neq \boldsymbol{\mu}_0$  for  $\boldsymbol{\mu}_0$  some particular value of the parameter  $\boldsymbol{\mu}$  (which can be a parameter of location, scatter, skewness, etc.). Suppose that the classical (meaning  $n \rightarrow \infty$  while  $p$  remains small) asymptotic distribution of the associated test statistic  $T_p^n$  follows a chi-square distribution with  $p$  degrees of freedom, which we write  $\chi_p^2$ . In other words,  $\mathcal{H}_0$  is rejected whenever  $T_p^n > \chi_{p;1-\alpha}^2$ , the  $\alpha$ -upper quantile of the  $\chi_p^2$  distribution. Now, when the dimension  $p$  itself becomes very large, potentially larger than  $n$ , this test becomes worthless as the chi-square distribution will diverge (recall that its expectation is  $p$  and its variance is  $2p$ ). Consequently, the test statistic needs to be modified, for instance into

$$\tilde{T}_p^n := \frac{T_p^n - p}{\sqrt{2p}} \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \frac{X_p - p}{\sqrt{2p}} \xrightarrow[p \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, 1) \quad (1)$$

where  $X_p$  stands for a  $\chi_p^2$  random variable and  $\mathcal{D}$  means convergence in distribution. From (1) we see that comparing the modified test statistic  $\tilde{T}_p^n$  to quantiles of the standard normal distribution would allow us to have a new large- $p$  test for our hypothesis of interest, provided the so-called  $(n, p)$ -asymptotic

result from (1) holds true. Indeed, as both  $n$  and  $p$  grow large, there is no guarantee that the limit when both  $n$  and  $p$  go to infinity can be calculated by first letting  $n$  become large and then  $p$ . This must be formally proved. In certain cases it turns out to be a valid manipulation, but in other situations it does not and the initial test statistic must be changed more substantially. An example of such distinct situations is provided in the seminal paper by Ledoit and Wolf [13] who considered scatter matrices.

These two examples underline two novel challenges statisticians are facing when dealing with Big Data. The need to cope with such data has given rise to a popular new research direction, called high-dimensional statistics (see, e.g., [14]). Besides this new research line, the entire field of Statistics has undergone changes as a reaction to the new data paradigm<sup>3</sup>. Supervised and unsupervised learning, shrinkage techniques, graphical models, data mining, functional data analysis and methods to deal with intractable likelihood models are just a few of the new hot topics in statistical research. There is also an increasing trend towards Statistics occupying a central role in Science in general, as discussed in the next section.

### 3.3 Data Science = Statistics2.0

The idea of a statistician (or mathematician) working in an Ivory Tower is obsolete. Several fields are in need of statisticians to help them analyse their data; conversely, significant advances in statistics have been driven by such demands and the collaboration with experts having complementary knowledge. The Big Data era offers Statistics plenty of new possibilities and has brought this traditional field to the limelight of modern scientific research. The era of data may be that of the rebirth of Statistics. Hal Varian, chief economist of Google, said in 2009 “*I keep saying the sexy job in the next ten years will be statisticians. People think I’m joking, but who would’ve guessed that computer engineers would’ve been the sexy job of the 1990s?*”.

Where precisely lies the boundary between modern Statistics and machine learning? How much statistics is present in Computational Biology, in Bioinformatics? Health Sciences have benefitted enormously from tailor-made statistical research, see [16] for examples. The same holds true for systems biomedicine, finance and environmetrics, among many others. Diggle in [16] expresses his opinion that Statistics is actually the Data

Science of our modern times. We concur with him and like to say that Data Science is actually Statistics2.0, hereby underlining the new orientation Statistics has taken.

## 4 How does Data Science relate to Computational Science?

This soaring amount of data has brought a new life to Statistics, and by doing so has also opened new doors to the discipline known as “Computational Sciences” or “Scientific Computing.” We discuss briefly in this section how Data Science relates to Computational Sciences and how it may revolutionise the way we think about modelling, simulations and computations and enable a transformation of the engineering ecosystem.

First, let us agree that Science is defined as *the activity concerned with the systematic acquisition of knowledge and is an enterprise that builds and organises knowledge in the form of testable explanations and predictions about the universe*. Engineering we define as *the application of scientific and practical knowledge for the benefits of mankind*. For example, Theodore von Kármán, a leading mathematician, aerospace engineer and physicist, developed theories for aerodynamics, in particular supersonic and hypersonic airflow characterisation, which have been essential to the design and fabrication of modern jet engines and rockets. Computational Sciences have been an essential tool for such theories to bear upon modern design approaches.

To produce new knowledge and apply this knowledge to practical fields, scientists and engineers use the “scientific method” which tests statements that are logical consequences of scientific hypotheses (theories or computer models and simulations) through repeatable experiments and observations. This production of knowledge has been fuelled by a significant revolution which has taken place over the last 50 years, through which a new, inherently multi-disciplinary pillar of science has emerged to complement these theories and observations. Computational Sciences. Computational Sciences is the tri-disciplinary endeavour concerned with the use of computational methods and devices to enable scientific discovery and engineering applications in science.

In this new era, the wealth of Data has transformed the world of scientific discovery and engineering innovation. We believe that the fusion of computational science with data science will lie at the core of future scientific and engineering research. A new ability will play a central role, namely that of extracting knowledge from this wealth of information by storing, compressing, classifying, ordering and analyzing Data.

<sup>3</sup> Samworth in [15] provides a concise and very accessible overview on the new data-driven statistical research.

In particular, we will witness the emergence of smart systems, able to adapt to their environment through advanced data gathering and treatment approaches. These developments will be multi-disciplinary with mathematics, in particular statistics and numerical analysis, as well as computer science at its core.

In short, the fusion of Computational Sciences, a half-century old scientific field with Data Science, a modern embodiment of Statistics, will fuel the development of exciting new research, technological and business developments. The interested reader can refer to the two papers [17, 18].

## 5 Interdisciplinarity aspects

Data Science is, by definition, an interdisciplinary field. It incorporates knowledge from Statistics, Computer Science and Mathematics which it brings to bear on challenging application domains which had remained out of reach because of a combined lack of data and computer power. In what follows we shall illustrate this interdisciplinary nature of data science by means of two case studies.

### 5.1 Case study 1: protein structure prediction

Predicting the correct three-dimensional structure of a protein given its one-dimensional protein sequence is a crucial issue in life sciences and bioinformatics. Massive databases of DNA and protein sequences have become available, and many research groups are actively pursuing their efforts to solve the protein folding problem.

A promising approach has been put forward by the research group of Prof. Thomas Hamelryck from the University of Copenhagen. It combines inputs from biology, statistics, machine learning, physics and computer science, and hence is a nice example of Data Science in action. One of their main ingredients are graphical models from Machine Learning such as dynamic Bayesian networks, which they analyze from a statistical physics standpoint. An essential part of every protein sequence are the dihedral angles between certain atoms. Predicting their most likely values is a key component in understanding the protein structure at a local level. These pairs of angles, however, are not typical quantities as 0 degrees and 360 degrees represent the same value, hence pairs of angles need to be represented as data points on a torus. Devising statistical models and methods for such data is part of a research stream called Directional Statistics (see the book [19] for a recent account) and requires, besides mathematics, also computer science skills. Finally, the Hamelryck

group uses probability kinematics to combine their findings on local and non-local structures in a meaningful way.

We refer the interested reader to the monograph [20] for details about this approach.

### 5.2 Case study 2: Digital Twins of Engineering devices

Our second case study is concerned with the problem of model selection in engineering and medical simulations. All systems devised today in engineering fall within the category of Complex Systems, i.e. *a system composed of many components which interact with each other*. Natural systems such as the human body or the environment are obvious examples of Complex Systems. It is not possible to study, design and optimise such systems using analytical methods, i.e. hand calculations and for this, recourse is always made to some type of mathematical model. This mathematical model, usually a set of partial differential equations (PDEs), is solved numerically using *discretisation methods* such as finite element methods [21–25], finite differences, meshfree methods [26] or isogeometric approaches [27, 28].

Although discretisation methods have been subject to a large amount of research, one of the most difficult tasks to perform reliable and predictive simulations of complex systems is not merely the necessary choice of discretisation methods, but that of a suitable mathematical model. In other words, computational engineers need to answer the question: “What is the best model of this system given computational constraints and the quantities I am interested in?”

Let us look at this problem of *model selection* through two connected examples. First, let us consider modern engineering materials, such as composite materials which have been developed to perform well in increasingly challenging environments<sup>4</sup>. The durability of gigantic composite structures such as the Airbus A380, over 70m in wingspan is influenced by physical phenomena occurring at the scale of carbon fibres which are around 5 microns in diameter. The brute-force approach consisting of including all carbon fibres in the simulation of one cubic millimetre of composite material would require solving a set of eight billion equations in eight billion unknowns, making the problem completely intractable over the size of the aircraft. The task of the computational engineer is therefore to select a model which is computationally cheaper whilst able to predict the behaviour of the structure at the carbon fibre level.

<sup>4</sup> in particular for space applications where not only mechanical but radiation and thermal effects become critical

Once a suitable model has been selected, the associated parameters must be identified in light of experimental observations, i.e. the model must be calibrated. In materials engineering, the traditional approach to this has been to perform experiments within laboratory conditions, which are most often far removed from those which the structure or system will undergo during its service life, in particular when harsh environmental effects are of interest. Statistical approaches can be used, but they only partially overcome the hurdle as they are reliant upon predefined statistical distributions, which do not account for “unknown unknowns” or in-service conditions which were not considered during the experimental campaigns.

The world of Big Data has changed the way engineers look at model selection. The Computational Engineering community has been working on alternative paradigms to the traditional experimental-based model calibration approach about a decade now, by borrowing ideas from the Statistics community (namely Bayesian inference) and by working with Computer Science teams on machine learning methods [29–34]. The Bayesian paradigm enables the enrichment of prior knowledge with new data, as it is being acquired. The need for and possibility of such model selection and parameter identification approaches can be considered as one of the fruits of the advent of the Data Science era. A discussion of pros and cons of the Bayesian approach for model calibration is provided here [35].

Our second example is personalised medicine. Whilst important in engineering, the need to update models on the fly as new data becomes available is necessary in personalised medicine where all patients are different and experiments are simply not possible. In this field, it is necessary to infer the best possible model for a patient from a priori knowledge obtained from other patients. Successful approaches have been recently published [18] [34] which enable predictive science in medicine, for example for laser-treatment of tumours [29]. The reader is referred to [36] for a recent discussion of the emerging field known as “Computer-Guided Predictive Medicine” [36].

Such an aptitude for on-the-fly data fusion has been fuelling the development of a revolutionary paradigm known as the “digital-twin concept” enabling predictive, high-fidelity models to learn from real-time data acquired during the life of the system, accounting for “real” environmental conditions during predictions. In turn, digital twins may allow us to move beyond the “factors of safety” era in engineering, where uncertainties are lumped into global correction factors leading to over-engineering. The structures and systems we will develop will be able to adapt to their environment.

Yet, these digital twins will remain a dream unless data science approaches are harnessed by computational engineers. This will require significant efforts in educating the next generation engineers and data scientists. We expect this to happen at the interface between mathematics, statistics, engineering and application areas. Exciting futures are in sight if we harness such complexities and build the required multi-disciplinary teams.

## 6 Conclusions and discussion

We discussed in this paper what we believe makes data science different. We offered various interpretations of data science and differentiated between bottom-up and top-down data science approaches. We also defined science, engineering and computational sciences/scientific computing and attempted to relate data science to these more established disciplines. Through personal examples and two case studies we provided possible explanations for the singularity of data science.

In short, we conclude that data science enhances the traditional and more conservative world of statistics with advanced machine learning algorithms to enable us to make sense out of soaring amounts of data. Here are our conclusions:

- Data science fuses statistics, machine learning, mathematics and computer science. Computers are of key importance in data science, in particular for bottom-up approaches, but the creation of suitable models, mandatory to make machine learning approaches computationally tractable, requires expert knowledge which we believe will be brought forward by statisticians.
- Data science has the potential to have strong impact in application domains, in particular on engineering and medicine. Some of the exciting applications of data science include the delivery of the next generation smart and autonomous devices able to learn from and adapt to their environment.
- Through a craft coupling with Computational Sciences, Data Science Can help create “digital twins” of complex systems. Those are replicas of the actual system which live a parallel, virtual/digital life and can be interrogated in order to make decisions on the (cyber-)physical system itself.
- Data science is an attractive name which Data Science sounds young, exciting, innovative, and partially mysterious. This may endow those entering this field with a particularly creative and less conservative mindset than in other, more established disciplines.

– Data science is the right discipline at the right time: the data deluge creates urgent needs and challenging problems, both in academia, industry and business. Spurred by a rapid increase in compute power and the ability of mobile devices to generate large amounts of data everywhere we leave our digital footprints, Data Science appears to be the tailor-made discipline to help make sense out of large amounts of data.

Having made the above reflections, there are a number of points which seem important to us going forward in the world of data science:

- Ensuring that we do not fall for the “hype of data science” and ignore theories to the benefit of machine learning algorithms. There is need for a “scientist in the loop” even when bottom-up approaches are advocated.
- Devise suitable training programmes at all levels, in particular through continuing education, in order to help create sound careers for data scientists, at the interface between statistics and computer science, with robust mathematical foundations.
- Nurture an intellectually coherent core relying on Mathematics, Statistics and Computer Science to provide rigorous abstractions to application domains and receiving in return stimulating problems and challenges to address.
- Develop research and teaching programmes at the interface between Computational Science and Data Science.
- Foster communication between the disciplines at play by encouraging jargon-free discussions and joint conferences.

In our opinion, an exciting research direction lies at the interface between bottom-up (machine learning) and top-down (statistical modelling) approaches. In many systems, pure computing power and machine learning algorithms are insufficient to obtain results within a reasonable time frame. At the same time, full mathematical models involving the full complexity of the system at hand are also computationally intractable, for example in quantum physics [37–41]. Building such hybrid strategies, we expect, will continue to be exciting research directions, at the interface between statistics, machine learning, and application domains, e.g. [42] [43] [44]. These hybrid approaches will provide users with a new way to design experiments, based on data acquired on the fly [45, 46].

We presented what are but our personal opinions. The reader is free to disagree with us. We hope nonetheless to have contributed a fresh and multi-disciplinary

view to the understanding of what makes Data Science different and hence so popular as discipline.

## References

1. S. Sagioglu, D. Sinanc, in *Collaboration Technologies and Systems (CTS), 2013 International Conference on* (IEEE, 2013), pp. 42–47
2. C. Hayashi, in *Data Science, Classification, and Related Methods* (Springer, 1998), pp. 40–51
3. M. Loukides, *What is data science?* (" O'Reilly Media, Inc.", 2011)
4. R. Akerkar, P.S. Sajja, in *Intelligent Techniques for Data Science* (Springer, 2016), pp. 1–30
5. R. Kitchin, *The SAGE Handbook of Social Media Research Methods* p. 27 (2017)
6. S.C. Hicks, R.A. Irizarry, arXiv preprint arXiv:1612.07140 (2016)
7. R. Tang, W. Sae-Lim, *Education for Information* **32**(3), 269 (2016)
8. D. Patil, *Building data science teams* (" O'Reilly Media, Inc.", 2011)
9. C. Cioffi-Revilla, *Browser Download This Paper* (2016)
10. D.L. McDowell, S.R. Kalidindi, *MRS Bulletin* **41**(04), 326 (2016)
11. L. Cao, Submitted to *ACM Computing Survey* pp. 1–37 (2016)
12. R. Tibshirani, *Journal of the Royal Statistical Society. Series B (Methodological)* pp. 267–288 (1996)
13. O. Ledoit, M. Wolf, *Annals of statistics* pp. 1081–1102 (2002)
14. P. Bühlmann, S. Van De Geer, *Statistics for high-dimensional data: methods, theory and applications* (Springer Science & Business Media, 2011)
15. R.J. Samworth,
16. P.J. Diggle, *Journal of the Royal Statistical Society: Series A (Statistics in Society)* **178**(4), 793 (2015)
17. K. Willcox, G. Bounova, in *Proceedings of the 2004 ASEE Annual Conference & Exposition, Session* (Citeseer, 2004), 2465
18. U. Råde, K. Willcox, L.C. McInnes, H. De Sterck, G. Biros, H. Bungartz, J. Coronos, E. Cramer, J. Crowley, O. Ghattas, et al., arXiv preprint arXiv:1610.02608 (2016)
19. C. Ley, T. Verdebout, *Modern Directional Statistics* (Chapman and Hall/CRC, 2017)
20. T. Hamelryck, K. Mardia, J. Ferkinghoff-Borg, *Bayesian methods in structural bioinformatics* (Springer Science & Business Media, 2012)
21. G. Strang, G.J. Fix, *An analysis of the finite element method*, vol. 212 (Prentice-hall Englewood Cliffs, NJ, 1973)
22. O.C. Zienkiewicz, R.L. Taylor, R.L. Taylor, *The finite element method*, vol. 3 (McGraw-hill London, 1977)
23. K.J. Bathe, *Finite element method* (Wiley Online Library, 2008)
24. G. Dhatt, E. Lefrançois, G. Touzot, *Finite element method* (John Wiley & Sons, 2012)
25. T.J. Hughes, *The finite element method: linear static and dynamic finite element analysis* (Courier Corporation, 2012)
26. V.P. Nguyen, T. Rabczuk, S. Bordas, M. Dufloot, *Mathematics and computers in simulation* **79**(3), 763 (2008)
27. T.J. Hughes, J.A. Cottrell, Y. Bazilevs, *Computer methods in applied mechanics and engineering* **194**(39), 4135 (2005)

28. V.P. Nguyen, C. Anitescu, S.P. Bordas, T. Rabczuk, *Mathematics and Computers in Simulation* **117**, 89 (2015)
29. D. Fuentes, J. Oden, K. Diller, J. Hazle, A. Elliott, A. Shetty, R. Stafford, *Annals of biomedical engineering* **37**(4), 763 (2009)
30. T. Oden, R. Moser, O. Ghattas, *SIAM News* **43**(9), 1 (2010)
31. J.T. Oden, S. Prudhomme, *International Journal for Numerical Methods in Engineering* **87**(1-5), 262 (2011)
32. A. Hawkins-Daarud, S. Prudhomme, K.G. van der Zee, J.T. Oden, *Journal of mathematical biology* **67**(6-7), 1457 (2013)
33. E. Prudencio, P. Bauman, D. Faghihi, K. Ravi-Chandar, J. Oden, *International Journal for Numerical Methods in Engineering* **102**(3-4), 379 (2015)
34. J.T. Oden, E.A. Lima, R.C. Almeida, Y. Feng, M.N. Rylander, D. Fuentes, D. Faghihi, M.M. Rahman, M. DeWitt, M. Gadde, et al., *Archives of Computational Methods in Engineering* **23**(4), 735 (2016)
35. H. Rappel, L.A. Beex, J.S. Hale, S. Bordas, arXiv preprint arXiv:1606.02422 (2016)
36. J.T. Oden, E.A. Lima, R.C. Almeida, Y. Feng, M.N. Rylander, D. Fuentes, D. Faghihi, M.M. Rahman, M. DeWitt, M. Gadde, et al., *Archives of Computational Methods in Engineering* pp. 1–45 (2015)
37. A. Tkatchenko, M. Scheffler, *Physical review letters* **102**(7), 073005 (2009)
38. G. Montavon, K. Hansen, S. Fazli, M. Rupp, F. Biegler, A. Ziehe, A. Tkatchenko, A.V. Lilienfeld, K.R. Müller, in *Advances in Neural Information Processing Systems* (2012), pp. 440–448
39. M. Rupp, A. Tkatchenko, K.R. Müller, O.A. Von Lilienfeld, *Physical review letters* **108**(5), 058301 (2012)
40. G. Montavon, M. Rupp, V. Gobre, A. Vazquez-Mayagoitia, K. Hansen, A. Tkatchenko, K.R. Müller, O.A. Von Lilienfeld, *New Journal of Physics* **15**(9), 095003 (2013)
41. K. Hansen, F. Biegler, R. Ramakrishnan, W. Pronobis, O.A. Von Lilienfeld, K.R. Müller, A. Tkatchenko, *The journal of physical chemistry letters* **6**(12), 2326 (2015)
42. P. Kerfriden, K.M. Schmidt, T. Rabczuk, S.P.A. Bordas, *International Journal for Multiscale Computational Engineering* **11**(3) (2013)
43. T. Mueller, A.G. Kusne, R. Ramprasad, *Reviews in Computational Chemistry* **29**, 186 (2016)
44. P. Antony, P. Manujesh, N. Jnanesh, in *Recent Trends in Electronics, Information & Communication Technology (RTEICT), IEEE International Conference on* (IEEE, 2016), pp. 69–73
45. F. Viti, W. Verbeke, C. Tampère, *Transportation Research Record: Journal of the Transportation Research Board* (2049), 103 (2008)
46. A. Fonzone, J.D. Schmöcker, F. Viti. New services, new travelers, old models? directions to pioneer public transport models in the era of big data (2016)