# **Objectives and considerations of (Big) data science education in Life Sciences domain**

### 26 October 2018, Lukasz Grus

















100 years

### What is data science ?

**Data science** is the extraction of actionable knowledge directly from data through a process of discovery, or hypothesis formulation and hypothesis testing.

Data science [...] incorporates principles, techniques, and methods from many disciplines and domains including **(advanced) data management, data cleansing, analytics, visualization, engineering, etics, policy**, in the context of Big Data\*







### Why data science?





Devis Tuia, Benjamin Kellenberger (GRS)



100 years





Prediction generally > 90% using BN. Forecasting also possible







### Cell

#### ARTICLE | VOLUME 163, ISSUE 5, P1079-1094, NOVEMBER 19, 2015

#### Personalized Nutrition by Prediction of Glycemic

#### Responses

David Zeevi <sup>®</sup> • Tal Korem <sup>®</sup> • Niv Zmora <sup>®</sup> • David Israeli <sup>®</sup> • Daphna Rothschild • Adina Weinberger • ... Zamir Halpern • Eran Elinav <sup>A</sup> ° ⊡ • Eran Segal <sup>A</sup> ° ⊡ • Show all authors • Show footnotes

Open Archive • DOI: https://doi.org/10.1016/j.cell.2015.11.001 • (E) Check for updates

| Highlights   | View fulltext  |
|--------------|--|
| Summary      | Highlights   |
| Graphical    | ingingito  |
| Abstract     | High interpersonal variability in post-meal glucose observed in an 800-person cohort |
| References   | Using personal and microbiome features enables accurate glucose response prediction  |
| Article Info | Prediction is accurate and superior to common practice in an independent cohort      |
| Figures      | Short-term personalized dietary interventions successfully lower post-meal glucose   |







Cell. 2015 Nov 19;163(5):1079-1094. doi: 10.1016/j.cell.2015.11.001.

#### Personalized Nutrition by Prediction of Glycemic Responses.

Zeevi D<sup>1</sup>, Korem T<sup>1</sup>, Zmora N<sup>2</sup>, Israeli D<sup>3</sup>, Rothschild D<sup>1</sup>, Weinberger A<sup>1</sup>, Ben-Yacov O<sup>1</sup>, Lador D<sup>1</sup>, Aunt-Sapi T<sup>1</sup>, Lotan-Pomoan M<sup>1</sup>, Suez J<sup>4</sup>, Mahdi JA<sup>4</sup>, Matot E<sup>1</sup>, Malka G<sup>1</sup>, Kosower N<sup>1</sup>, Rein M<sup>1</sup>, Zilberman-Schapira G<sup>4</sup>, Dohnalová L<sup>4</sup>, Pevsner-Fischer M<sup>4</sup>, Bikovsky R<sup>1</sup>, Halpern Z<sup>5</sup>, Elinav E<sup>6</sup>, Segal E<sup>7</sup>. 8

### Wageningen Data Competence Centre







(in relation to education)





### Consideration 1 – Data science field is very broad





100 years

# Profile Data scientist - competences

### (adapted from EDISON project)



| Data Analytics   | Data Science Engineering   | Data Management  | Research Methods and<br>Project Management   | Domain related<br>Competences   |
|--|--|--|--|---|
| Use data analysis and statistical<br>techniques on data <i>to deliver insights<br/>into research problem</i>   | Use engineering principles (software<br>design and development) and modern<br>computer technologies (programming) to<br>research, design, implement new data<br>analytics applications.  | Develop and implement data<br>management strategy for data collection,<br>storage, preservation, and availability for<br>further processing          | Create new understandings and<br>capabilities by using the scientific<br>method  | Use domain knowledge to develop<br>relevant data analytics applications   |
| Use techniques such as Machine learning,<br>Data Mining, Prescriptive and Predictive<br>analytics, for complex data analysis through<br>the whole data cycle.  | Use engineering principles (general and<br>software) to research, <b>design</b> , <b>develop and</b><br><b>implement instruments and applications</b><br>for data collection, storage, analysis and<br>visualisation                                 | Develop and implement data strategy, in<br>particular, in a form of data management<br>policy and Data Management Plan (DMP)                         | Create new understandings by using the research methods  | Analyse information needs, assess<br>exisitng data and suggest/identify new<br>data required for specific context   |
| Apply statistics, time series analysis,<br>optimization, simulation, to deploy models<br>for analysis and prediction   | Develop and apply computational<br>solutions to domain related problems using<br>data analytics platforms, with the special<br>focus on Big Data technologies and cloud<br>based data analytics platforms  | Develop and implement <b>data models</b> , define<br><b>metadata</b> using common standards and<br>practices   | Direct systematic study toward understanding<br>of the observable facts, and discovers new<br>methods                            | Operationalise fuzzy concepts to enable key<br>performance indicators measurement to<br>validate the research results or business<br>analysis, identify and assess potential<br>challenges                          |
| Identify, extract, and pull together available<br>and pertinent heterogeneous data,<br>including modern data sources such as<br>social media data, open data,<br>governmental data                           | Develop and prototype specialised data<br>analysis applicaions, tools and supporting<br>infrastructures for data driven scientific<br>workflow.  | Integrate heterogeneous data from<br>multiple source and provide them for further<br>analysis and use  | Analyse domain related available data to<br>identify research questions and<br>formulate sound hypothesis                        | Deliver business focused analysis using<br>appropriate BA/BI methods and tools, identify<br>business impact from trends; make business<br>case as a result of organisational data<br>analysis and identified trends |
| Understand and use different performance<br>and accuracy metrics for model validation in<br>analytics projects, hypothesis testing, and<br>information retrieval   | Develop, deploy and operate large scale<br>data storage and processing solutions<br>using different distributed and cloud based<br>platforms for storing data (e.g. Data Lakes,<br>Hadoop, Hbase, Cassandra, MongoDB,<br>Accumulo, DynamoDB, others) | Maintain historical information on data<br>handling, including reference to published<br>data and corresponding data sources (data<br>provenance)    |  | Analyse opportunity and suggest use of<br>historical data available in the study field or<br>organization for creating new knowledge or<br>optimization   |
| Develop required data analytics for<br>organizational tasks, integrate data analytics<br>and processing applications into organization<br>workflow and business processes to enable<br>agile decision making | <b>Consistently apply data security</b><br>mechanisms and controls at each stage of the<br>data processing, including data<br>anonymisation, privacy and IPR protection.   | Ensure data quality, accessibility,<br>interoperability, compliance to<br>standards, and publication (data curation),<br>comply with FAIR principles | Design experiments which include data<br>collection (passive and active) for hypothesis<br>testing and problem solving           | Analyse customer relations data to<br>optimise/improve interacting with the specific<br>user groups or in the specific business sectors   |
| Visualise results of data analysis, design<br>dashboard and use storytelling methods   | Design, build, operate relational and non-<br>relational databases (SQL and NoSQL),<br>integrate them with the modern Data<br>Warehouse solutions, ensure effective ETL<br>(Extract, Transform, Load), OLTP, OLAP<br>processes for large datasets    | Develop and manage/supervise policies on<br>data protection, privacy, IPR and ethical<br>issues in data management                                   | Develop and guide <b>data driven</b> projects,<br>including project planning, experiment design,<br>data collection and handling | Analyse multiple data sources for marketing<br>purposes; identify effective marketing actions   |

All students share interest and skills in Life Sciences

Not all students share interest and skills in

(Big) Data Science

























Classic researcher: *T-shaped* Alteroped Supervised Supervised

S. Ceri / Statistics and Probability Letters 136 (2018) 68-72

Fig. 1. Representations of T vs. Pi-shaped education.

This introductory course in data science is built on three interrelated perspectives: inferential thinking, computational thinking, and real-world relevance. Given data arising from some real-world phenomenon, how does one analyze that data so as to understand that phenomenon? How does one collect data to answer questions that one is interested in? Inferential thinking refers to an ability to connect data to underlying phenomena and to the ability to think critically about the conclusions that are drawn from data analysis. Computational thinking refers to the ability to conceive of the abstractions and processes that allow inferential procedures to be embodied in computer programs, and to ensure that such programs are scalable, robust and understandable. In addition to teaching critical concepts and skills in computer programming and statistical inference, the course will involve the hands-on analysis of a variety of real-world datasets, including economic data, document collections, geographical data analysis, including issues of privacy and data ownership.



100 years

Fig. 2. Syllabus of the undergraduate course Foundations of Data Science, Berkeley University (2015 edition).

### Strong (Life Sciences) domain knowledge is an opportunity



| Data Analytics   | Data Science Engineering   | Data Management  | Research Methods and<br>Project Management   | Domain related<br>Competences   |
|--|--|--|--|---|
| Use data analysis and statistical<br>techniques on data to deliver insights<br>into research problem   | Use engineering principles (software<br>design and development) and modern<br>computer technologies (programming)<br>to research, design, implement new data<br>analytics applications.  | Develop and implement data<br>management strategy for data<br>collection, storage, preservation, and<br>availability for further processing          | Create new understandings and<br>capabilities by using the scientific<br>method  | Use domain knowledge to develop<br>relevant data analytics applications   |
| Use techniques such as Machine learning,<br>Data Mining, Prescriptive and Predictive<br>analytics, for complex data analysis through<br>the whole data cycle.  | Use engineering principles (general and<br>software) to research, design, develop and<br>implement new instruments and applications<br>for data collection, storage, analysis and<br>visualisation   | Develop and implement data strategy, in<br>particular, in a form of data management<br>policy and Data Management Plan (DMP)                         | Create new understandings by using the<br>research methods   | Analyse information needs, assess<br>exisiting data and suggest/identify new<br>data required for specific context  |
| Apply statistics, time series analysis,<br>optimization, simulation, to deploy models<br>for analysis and prediction   | Develop and apply computational solutions to<br>domain related problems using data analytics<br>platforms, with the special focus on Big Data<br>technologies and cloud based data analytics<br>platforms  | Develop and implement <b>data models</b> , define<br><b>metadata</b> using common standards and<br>practices   | Direct systematic study toward<br>understanding of the observable facts, and<br>discovers new methods                            | Operationalise fuzzy concepts to enable key<br>performance indicators measurement to<br>validate the research results or business<br>analysis, identify and assess potential<br>challenges                            |
| Identify, extract, and pull together available<br>and pertinent heterogeneous data,<br>including modern data sources such as<br>social media data, open data,<br>governmental data                           | Develop and prototype specialised data<br>analysis applicatons, tools and supporting<br>infrastructures for data driven scientific<br>workflow.  | Integrate heterogeneous data from<br>multiple source and provide them for further<br>analysis and use  | Analyse domain related available data to<br>identify research questions and<br>formulate sound hypothesis                        | Deliver business focused analysis using<br>appropriate BA/BI methods and tools,<br>identify business impact from trends; make<br>business case as a result of organicitational<br>data analysis and identified trends |
| Understand and use different performance<br>and accuracy metrics for model validation in<br>analytics projects, hypothesis testing, and<br>information retrieval   | Develop, deploy and operate large scale data<br>storage and processing solutions using<br>different distributed and cloud based<br>platforms for storing data (e.g. Data Lakes,<br>Hadoop, Hbase, Cassandra, MongoDB,<br>Accumulo, DynamoDB, others) | Maintain historical information on data<br>handling, including reference to published<br>data and corresponding data sources (data<br>provenance)    |  | Analyse opportunity and suggest use of<br>historical data available in the study field or<br>organization for creating new knowledge or<br>optimization   |
| Develop required data analytics for<br>organizational tasks, integrate data analytics<br>and processing applications into organization<br>workflow and business processes to enable<br>agile decision making | Consistently apply data security mechanisms<br>and controls at each stage of the data<br>processing, including data anonymisation,<br>privacy and IPR protection.  | Ensure data quality, accessibility,<br>interoperability, compliance to<br>standards, and publication (data<br>curation), comply with FAIR principles | Design experiments which include data<br>collection (passive and active) for hypothes<br>testing and problem solving             | Analyse customer relations data to<br>optimise/improve interacting with the specific<br>user groups or in the specific business<br>sectors  |
| Visualise results of data analysis, design<br>dashboard and use storytelling methods   | Design, build, operate relational and non-<br>relational databases (SQL and NoSQL),<br>integrate them with the modern Data<br>Warehouse solutions, ensure effective ETL<br>(Extract, Transform, Load), OLTP, OLAP<br>processes for large datasets    | Develop and manage/supervise policies on<br>data protection, privacy, IPR and ethical<br>issues in data management                                   | Develop and guide <b>data driven</b> projects,<br>including project planning, experiment<br>design, data collection and handling | Analyse multiple data sources for marketing<br>purposes; identify effective marketing actions   |



- Lots of bottom-up (research) and very specific DS education initiatives.
- Need for generic Data Science education offerings

#### SSB-53306 Life Sciences Information - Integrating, combining organizing and analyzing heterogeneous Biological information

#### Learning outcomes:

- After successful completion of this course students are expected to be able to:
- create a FAIR Data Management plan for complex biological experiments;
- design and perform advanced queries of (online) biological databases using SPARQL;
- design a custom-made ontology for different kinds of biological data;
- design a statistical experimental plan to address a complex biological question;
- extract relevant biological information using inferential approaches using R;
- apply data reduction techniques to analyse and investigate large biological data set.

### INF-22803 Intro to Data Structures and Algorithms for Health

#### Learning outcomes:

After successful completion of this course students are expected to be able to:

- show insight in basic computer science terminology and concepts;
- identify several approaches to problem solving;
- explain key concepts of data structures and algorithms;
- compare and select suitable basic data structures for health related data;
- implement a proper data structure for health related data;
- analyse health related data problems using selected algorithms;
- solve health related data problems computationally;

#### INF-33306 Linked Data

#### Learning outcomes:

After successful completion of this course students are expected to be able to:

- explain the research data life-cycle;
- explain linked data technologies (Internet, Semantic Web, ontologies, graph databases) and standards (as URI, XML, RDF, OWL, SPARQL) and their use in forming a Web of Data:
- use Linked Data technologies for retrieving information in the Semantic Web (i.e use of SPARQL);
- use of existing vocabularies for annotating research data and endpoints for finding data from the Life Sciences domain ;
- develop own ontologies and demonstrate the use of reasoners;







It is essential that the data science skills and knowledge will be combined and applied in specific domains such as plant, animal, environmental, social and agrotechnology & food sciences. The WUR education offering of data science courses focusses on this domain application ambition.





Data Science Courses Data Science Educational Materials



How to become a Data Scientist?

#### Events

- > 24 October 2018 -Webinar The State of Open Data
- > 29 November 2018 -Course
   Masterclass: datadriven agriculture and food production









100 years

| General education                          |                          |                                 |                         |                          |  |  |
|--|--------------------------|---------------------------------|-------------------------|--------------------------|--|--|
| Agrotechnology<br>& Food Sciences<br>Group | Animal Sciences<br>Group | Environmental<br>Sciences Group | Plant Sciences<br>Group | Social Sciences<br>Group |  |  |
| MSc′s<br>BSc′s                             | MSc's<br>BSc's           | MSc's<br>BSc's                  | MSc's<br>BSc's          | MSc's<br>BSc's           |  |  |
| Specific                                   | Specific                 | Specific                        | Specific                | Specific                 |  |  |
| education                                  | education                | education                       | education               | education                |  |  |
|  |                          |                                 |                         |                          |  |  |





#### Master's programmes

- > Agroecology (European)
- > Animal Sciences
- Aquaculture and Marine Resource Management
- > Biobased Sciences

> Bioinformatics

> Biology

> Biosystems Engineering

> Biotechnology

> Climate Studies

> Communication, Health and Life Sciences

- > Development and Rural Innovation
- > Earth and Environment
- > Environmental Sciences
- > Food Quality Management
- > Food Safety
- > Food Studies (European)
- > Food Technology
- > Food Technology (online)
- > Forest and Nature Conservation





- > Geographical Information Management and Applications
- Geo-information Science
- > International Development Studies
- > International Land and Water Management
- > Landscape Architecture and Planning
- > Leisure, Tourism and Environment
- > Management, Economics and Consumer Studies
- > Metropolitan Analysis, Design, and Engineering
- > Molecular Life Sciences
- > Nutritional Epidemiology and Public Health (online)
- > Nutrition and Health
- > Organic Agriculture
- > Plant Biotechnology
- > Plant Breeding (online)
- > Plant Sciences
- > Statistical Sciences for the Life and Behavioural Sciences
- > Urban Environmental Management
- > Water Technology

### Data science related courses

#### Information Technology (INF)

- > Data Management
- > Programming in Python
- > Big Data
- > Linked Data
- > Applied Information Technology
- > Software Engineering
- Computer Literacy
- > Agent-Based Modelling of Complex Adaptive Systems
- Modelling and Simulation of Complex Socio-Technical Systems

Amsterdam Institute for Advanced Metropolitan Solutions (YMS)

- > Metropolitan Data I
- > Metropolitan Data II

Human nutrition (HNE)

> Applied Data Analysis

Systems and Synthetic Biology (SSB) > Bioinformation Technology

oovears

#### Mathematical and Statistical Methods (MAT)

- > Mathematics
- > (Advanced) Statistics
- > Advanced Statistics for Nutritionists
- > R for Statistics

#### Library (LIB)

- > Information Literacy
- > Research Data Management (PhD / postdoc)

#### Farm Technology (FTE)

- > Machine Learning
- > Data Analysis Biosystems Engineering

Crop Systems Analysis (CSA)

Ecological Modelling and Data Analysis in R

#### Bioinformatics (BIF)

- > Advanced Bioinformatics
- > Algorithms in Bioinformatics
- Biological Data Analysis and Visualisation
- > Practical Computing for Biologists

### Geo-information Science and Remote Sensing (GRS)

- > Geo Scripting
- > Geo-Information Science for Society
- > Spatial Modelling and Statistics
- > Geo-Information Science in Context
- Machine Learning for Spatial Data (PhD / postdoc)

#### Research Methodology (YRM)

- > Data Analysis for Health and Society
- > Quantitative Data Analysis: Multivariate Techniques



### Data Science minor

- Employ a **data lifecycle approach** to organize their data driven research;

- Identify appropriate **analysis techniques** when confronted with new datasets and questions;

- Design and **implement databases** for applications in the WUR domain;
- Write efficient and well-documented computer programs;

- Apply and evaluate relevant visualization and quantitative data analysis methods;

- Communicate research findings by data storytelling;
- Consider broader issues related to **data collection**, processing and dissemination, including licensing and privacy issues.



INF-

22306

GRS

51306

Programming in Python

Geo-information Science

for Society

1AF

1AF

RO

RO

- > Bioinformatics
- > Plant Biotechnology
- > Systems Biology

### Data science **tracks** embedded in an application domains







### Summary

- Data science is a broad domain;
- Strong domain knowledge as an opportunity for DS
- Not everyone must like it but most will probably need it;
- Dillema which knowledge level;
- T  $\rightarrow$   $\Pi$  education;
- Research groups recognize already the need for DS education
  € € €

