# Data Science

## Introduction to Data Science

KEA – Data Science Elective PBA – Fall 2019 – Henrik Strøm

September 5, 2019

# Table of Contents

# About You and Me

**Henrik Strøm**

- Lecturer at KEA
- Ph.D. Fellow at Aalborg University
- M.Sc. in Engineering
- Master in Information and Communication Technologies
- Specialize in machine learning and data science
- Worked in Denmark, China, Thailand; projects in USA, Singapore
- Worked in major enterprise industry, financial sector, start-up
- Been at KEA since May 1st, 2018

**What about you guys – who are you, and what are your expectations?**

# Quick Poll

How many of you ...

- Know Python?
- Have taken class in artificial intelligence or machine learning?
- How many of you know about dimensionality reduction, cluster analysis, learning algorithms, how to handle bias and variance problems, and statistics?
- Have taken a high level math class?

## About the Course

**Objectives**
The objective of the module is to qualify the student to apply data science methods in a structured manner, extract inferential knowledge from a dataset, and make probabilistic forecasts. In addition, the student must be able to report on the findings and make use of visualization.

## About the Course

**Knowledge**
The objective is to give the student knowledge of:

- various definitions of the field of data science
- designing experiments for data collection
- using a framework for numerical computations
- using a framework for general data analysis
- using a framework for inferential statistical analysis
- using a framework for probabilistic forecasting
- a variety of data analysis algorithms and their applications

# About the Course

**Skills**

The objective is that the student will have acquired the ability to:

- collect data from a variety of sources
- organize data to prepare for analysis
- explore data to gain insights
- apply basic inferential statistics
- make forecasts using probabilistic machine learning tools
- use methods to facilitate reproducibility
- communicate findings in a written report, using visualizations

## About the Course

**Proficiency**
The objective is that the student will have acquired proficiency in participating and contributing in data science projects.

*Concepts and understanding is more important than syntax.*

# Course Workload

- This course is 10 ECTS points = 280 hours workload.
- It is important that you attend class – reading the material is not enough.
- Course routine and repetition.

# Mandatory Assignment and Exam

- Two written mandatory assignments: peer-review and hand-in.
- The peer-review is mandatory, you can't take the exam unless the mandatory assignments are accepted.
- Oral exam based on project.

# Defining Data Science

*Data science is gaining more and more and widespread attention, but no consensus viewpoint on what data science is has emerged. As a new science, its objects of study and scientific issues should not be covered by established sciences. In the present paper, data science is defined as the science of exploring datanature. We believe this is the most logical and accurate definition of data science (...) – (Zhu and Xiong, 2015)*
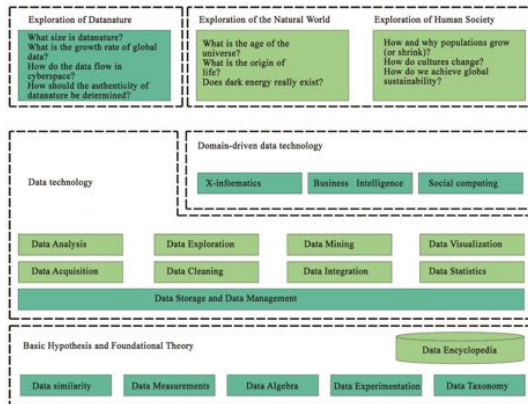
# Data Science Research Topics



Figure: Data Science research topics. (Zhu and Xiong, 2015)

# Elements of Data Science



Figure: Elements of Data Science (Goldstein, 2017)

# What's Most Important

- Transparency
- Reproducibility
- Correctness
- Predictive Power
- *Insights!*

Data Science is an application of scientific methods and principles to data processing.

## Modern Science

What is the difference between
the "science" of Artistole
and the science of Galileo Galilei?

# Computer Science – Levels of Knowledge

- Axiomatic knowledge
- Scientific Consensus
- Working application of best scientific standards (Development/Engineering)
- Software Construction
- Fooling Around

(Nunamaker Jr, Chen, and Purdin, 1990)

# Ethics

No, this is not going to be a social science class, however ...

*– ethics is important – both in how you conduct your research, and in how you report and communicate it. Data Science is not tabloid journalism or politics.*

With great power comes great responsibility: be honest and follow the evidence!

# Software Needed for the Course

- Git – use terminal or GUI
- Anaconda (Python 3 edition)
- Course materials:
  `https://bitbucket.org/henrikstroem/kea-public-datascience-2019-2`

# Python Primer

DEMO

## References I

📄 Goldstein, A. J. (2017). *Deconstructing Data Science: Breaking The Complex Craft Into It's Simplest Parts*. URL: https://medium.com/the-mission/deconstructing-data-science-breaking-the-complex-craft-into-its-simplest-parts-15b15420df21.

📄 Nunamaker Jr, Jay F, Minder Chen, and Titus D M Purdin (1990). "Systems development in information systems research". In: *Journal of management information systems* 7.3, pp. 89–106. URL: http://gkmc.utah.edu/7910F/papers/JMISsystemsdevelopmentinISresearch.pdf.

📄 Zhu, Yangyong and Yun Xiong (2015). "Defining data science". In: *arXiv preprint arXiv:1501.05039*. URL: https://arxiv.org/abs/1501.05039.