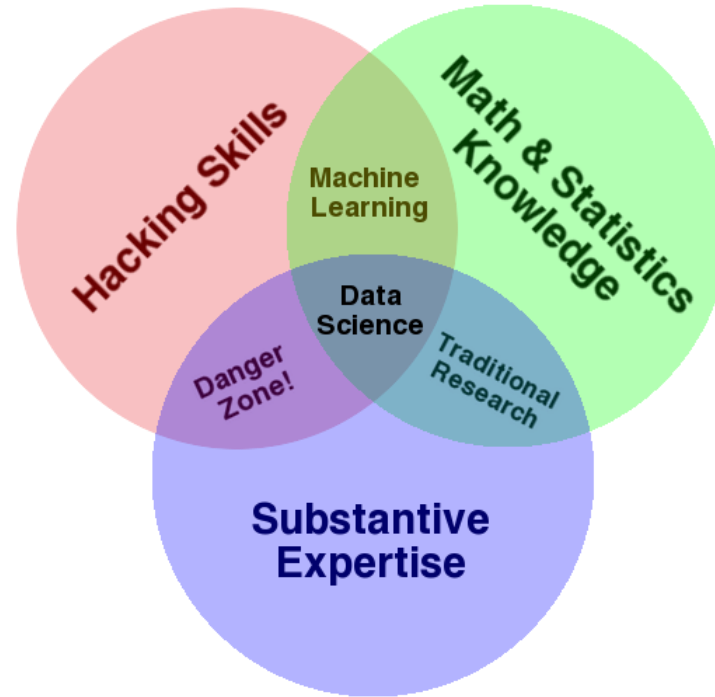


## Overview – Data Science at IMD

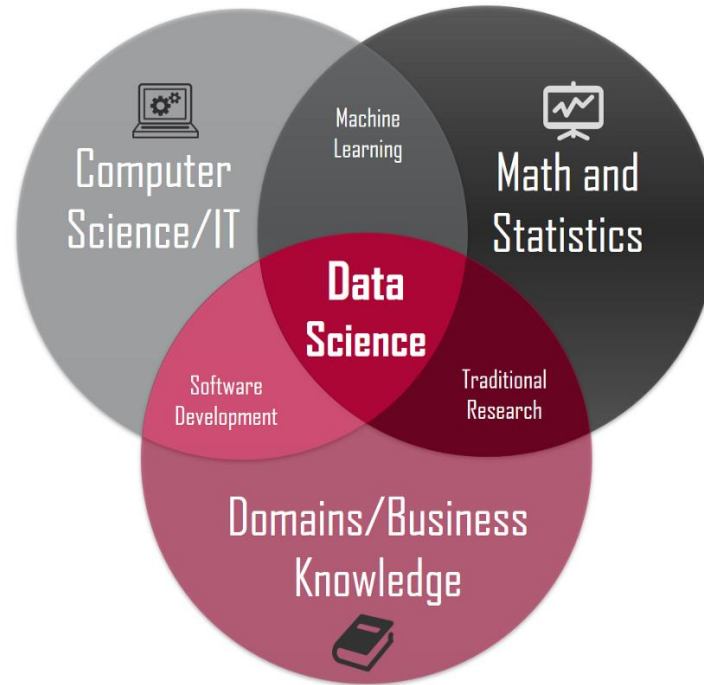
Jung PARK, PhD  
Research Fellow in Data Science

Version 2019 Aug 05

- Data Science is a multidisciplinary field that combines information technology (IT), statistics and management study.
- Due to the rapid advancement of IT, much more data and new analytic techniques became available.
- We need to balance it with sound statistical knowledge and business expertise to create useful insights.



Data science venn diagram  
(<http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>)



<https://towardsdatascience.com/introduction-to-statistics-e9d72d818745>



Source: <https://sourceable.net/degree-in-flying-cars-coming-soon/>

# Five processes in Data Science

Big Data Challenges and Solutions

Process	Challenges	Solutions	Key references
Data access and collection	<ul style="list-style-type: none"> <li>• Easy access to data offered in standardized formats. No practical limit to the size of these data offering unlimited scalability</li> <li>• Efficiently obtain detailed data for a large number of agents</li> <li>• Protocols on security, privacy, and data rights</li> </ul>	<ul style="list-style-type: none"> <li>• Sensors</li> <li>• Web scraping</li> <li>• Web traffic and communications monitoring</li> </ul>	<ul style="list-style-type: none"> <li>• Chaffin et al. (2015)</li> <li>• Sismeiro and Bucklin (2004)</li> </ul>
Data storage	<ul style="list-style-type: none"> <li>• Tools for data storage, matching and integration of different big datasets</li> <li>• Data reliability</li> <li>• Warehousing</li> </ul>	<ul style="list-style-type: none"> <li>• SQL, NoSQL, Apache Hadoop</li> <li>• Save essential information only and update in real time</li> </ul>	<ul style="list-style-type: none"> <li>• Varian (2014)</li> <li>• Prajapati (2013)</li> </ul>
Data processing	<ul style="list-style-type: none"> <li>• Use non-numeric data for quantitative analyses</li> </ul>	<ul style="list-style-type: none"> <li>• Text mining tools to transform text into numbers</li> <li>• Emotion recognition</li> </ul>	<ul style="list-style-type: none"> <li>• Manning, Raghavan, and Schütze (2009)</li> <li>• Teixeira, Wedel, and Pieters (2012)</li> </ul>
Data analysis	<ul style="list-style-type: none"> <li>• Large number of variables</li> <li>• Causality</li> <li>• Find latent topics and attach meaning</li> <li>• Data too large to process</li> </ul>	<ul style="list-style-type: none"> <li>• Ridge, lasso, principal components regression, partial least squares, regression trees</li> <li>• Topic modeling, latent Dirichlet allocation, entropy-based measures, and deep learning</li> <li>• Cross-validation and holdout samples</li> <li>• Field experiments</li> <li>• Parallelization, bags of little bootstrap, sequential analysis</li> </ul>	<ul style="list-style-type: none"> <li>• Hastie, Tibshirani, and Friedman (2009)</li> <li>• George and McCulloch (1993)</li> <li>• Archak, Ghose, and Ipeirotis (2011)</li> <li>• Tirunillai and Tellis (2012)</li> <li>• Blei, Ng, and Jordan (2003)</li> <li>• LeCun, Bengio, and Hinton (2015)</li> <li>• Lambrecht and Tucker (2013)</li> <li>• Wang, Chen, Schifano, Wu, and Yan (2015)</li> <li>• Wedel and Kannan (2016)</li> </ul>
Reporting and visualization	<ul style="list-style-type: none"> <li>• Facilitate interpretation, representation with external partners and knowledge users</li> <li>• Difficult to understand complex patterns</li> </ul>	<ul style="list-style-type: none"> <li>• Describe data sources</li> <li>• Describe methods and specifications</li> <li>• Bayesian analysis</li> <li>• Visualization and graphic interpretations</li> </ul>	<ul style="list-style-type: none"> <li>• Loughran and McDonald (2011)</li> <li>• Simonsohn, Simmons, and Nelson (2015)</li> </ul>

- Web scraping from public websites
  - How to collect information from websites
  - 10K reports from SEC.gov
  - EU Open Data
  - Open data Swiss
- Collecting data using API
  - Tweeter
  - Facebook
- Database in IMD library
  - Bloomberg
  - Thomson One
  - Datastream
  - Factiva: news media
- IoT (Internet of Things)
  - Smart building
  - Wearable devices
  - Web traffic

- Distributed computing
  - Do I need it?
  - What is Apache Hadoop
  - What is Apache Spark
  
- Other relevant information
  - SQL and NoSQL
  - Cloud technology



- Data wrangling
  - Merging the unstructured data
  - Text data
  - Images
  - Audios
  
- Contents analysis
  - NPL (Natural Language Processing)

- Big data specific (list from McKinsey Global Institute 2011)
  - A/B testing
  - Cluster analysis
  - Data fusion and integration
  - Data mining, genetic algorithms, machine learning
  - Natural language processing
  - Neural networks
  - Network analysis
  - Signal processing and spatial analysis, simulation
  - Time series analysis
  - Visualisation
- Econometrics
  - Time series analysis
  - Score matching techniques
  - Two-stage models for endogeneity problems
- Artificial Intelligence and Machine Learning
  - History of Artificial Intelligence
  - Supervised Machine Learning
  - Unsupervised Machine Learning/ Deep learning/ Neural networks

- Data Reporting
  - Robustness check
  - Information for reproducibility
  
- Related information
  - Open science
  - Ethics and privacy
  - GDPR (General Data Protection Regulation)

- George, G., Haas, M. R., & Pentland, A. (2014). FROM THE EDITORS BIG DATA AND MANAGEMENT. *Academy of Management Journal*, 57(2), 321-326.
- Drew Conway's data science venn diagram, accessed 5 Aug 2019 at <http://drewconway.com/zia/2013/3/26/the-data-science-venn-diagram>
- Michael Barber, accessed 5 Aug 2019 at <https://towardsdatascience.com/introduction-to-statistics-e9d72d818745>