



# *EDISON Data Science Framework: Building the Data Science Profession*

## *Introduction to discussions*

Yuri Demchenko, EDISON  
University of Amsterdam



**EDISON**  
building the data  
science profession

EDISON Data Science Champions  
conference July 2016

13 July 2016, New Forest, Brockenhurst, UK

EDISON – Education for Data Intensive  
Science to Open New science frontiers

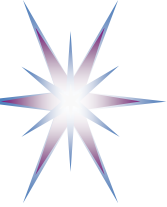
Grant 675419 (INFRASUPP-4-2015: CSA)



# Outline

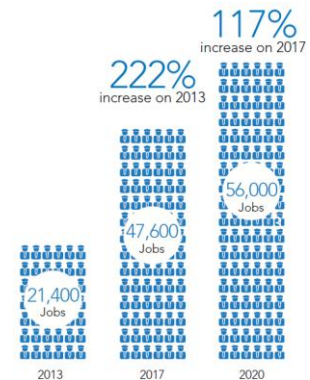
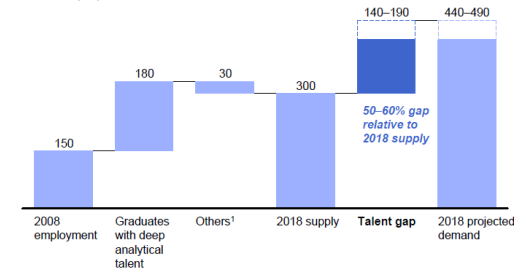
- Background and motivation
  - Demand for Data Science and data related professions
  - European initiatives related to Digital Single Market (DSM) and demand to data related competences and skills
- EDISON Data Science Framework
  - From Data Science Competences to Body of Knowledge and Model Curriculum
- Data Science Competence Framework (CF-DS)
  - Essential competences – Suggested use – Discussion questions
- Data Science Professions family and competence profiles (DSP)
  - Profiles definition and linking to CF-DS
- Data Science Body of Knowledge (DS-BoK)
  - Knowledge areas – Suggested use – Discussion questions
- Data Science Model Curriculum (MC-DS)
  - Learning Outcomes and Academic disciplines – Suggested use – Discussion questions





# Demand for Data Science and data related professions

- McKinsey Global Institute on Big Data Jobs (2011)  
[http://www.mckinsey.com/mgi/publications/big\\_data/index.asp](http://www.mckinsey.com/mgi/publications/big_data/index.asp)
  - Estimated gap of 140,000 - 190,000 data analytics skills by 2018
- UK Big Data skills report 2014
  - 6400 UK organisations with 100+ staff will have implemented Big Data Analytics by 2020
  - Increase of Big Data jobs from 21,400 (2013) to 56,000 (2017)
- IDC Report on European Data Market (2015)
  - Number of data workers 6.1 mln (2014)
  - increase 5.7% from 2013
  - Average number of data workers per company 9.5 - increase 4.4%
  - Gap between demand and supply 509,000 (2014) or 7.5%
- HLEG report on European Open Science Cloud (2016) identified need for data experts and data stewards
  - Recommendation: Allocate 5% from grant funding for Data management and preservation
  - Estimation: More than 80,000 data stewards (1 per every 20 scientists)
  - Core data experts need to be trained and their career perspective improved





# Recent European Commission Initiatives

Digitising European Industry: Reaping the full benefits of a **Digital Single Market**. COM(2016) 180 final, Brussels, 19.4.2016

- The need for new multidisciplinary and digital skills is exploding, including such as (Data Scientist) combining data analytics and business or engineering skills.
  - Expected rapidly growing demand will lead to more than 800 000 unfilled vacancies by 2020

New skills agenda for Europe. Working together to strengthen human capital, employability and competitiveness, COM(2016) 381 final, Brussels, 10.6.2016

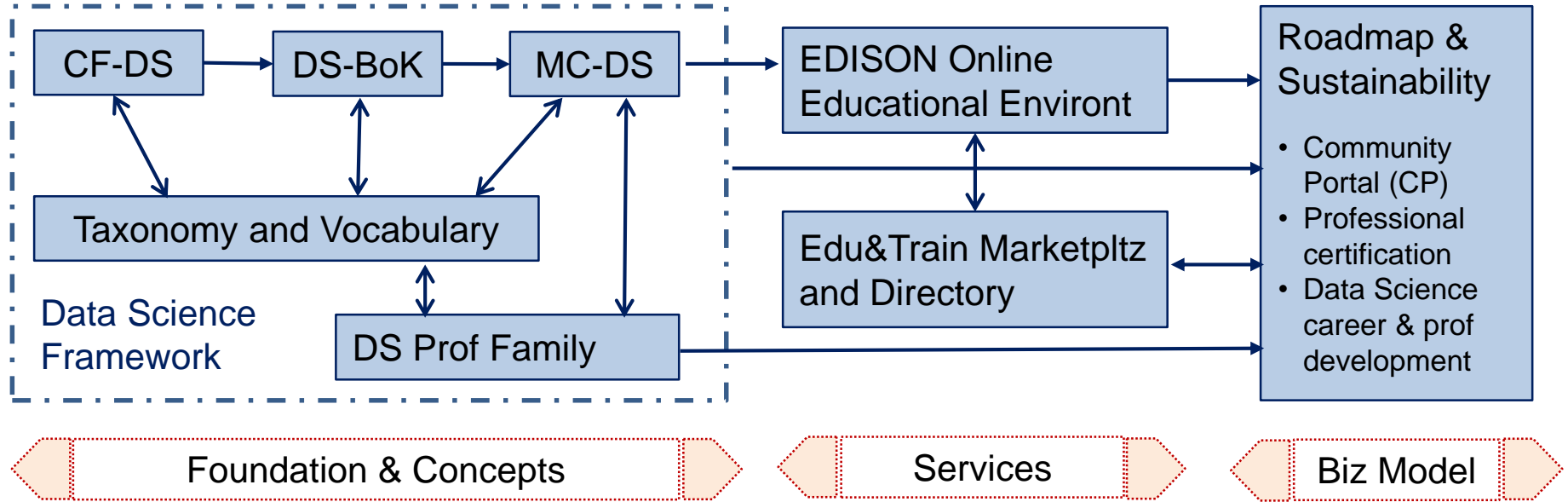
- Address the need for digital and complementary skills, ensure young talents flow into data driven research and industry, recognition and career development
- (Re-) Launch the **Digital Skills and Jobs Coalition (end of 2016)**
- **Develop comprehensive national digital skills strategies by mid-2017**

European Cloud Initiative - Building a competitive data and knowledge economy in Europe, COM(2016) 178 final, Brussels, 19.4.2016

- **European Open Science Cloud (EOSC)** and European digital research and data infrastructure
  - To offer 1.7 million European researchers and 70 million professionals in science and technology open and seamless services for storage, management, analysis and re-use of research data
  - Create incentives for academics, industry and public services to share their data
- Growing demand and shortage of data-related skills and lack of recognition of their value



# EDISON Data Science Framework (EDSF): Creating the Foundation for Data Science Profession



## EDISON Framework components

- CF-DS – Data Science Competence Framework
- DS-BoK – Data Science Body of Knowledge
- MC-DS – Data Science Model Curriculum
- DSP - Data Science Professions family and professional competence profiles
- EOEE - EDISON Online Education Environment

## Other components and services

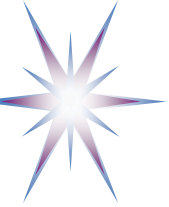
- EOEE - EDISON Online Education Environment
- Education and Training Marketplace and Resources Directory
- Data Science professional certification and training
- Community Portal (CP)

- Introduction to CF-DS
  - Background standards
  - How it was made
  - 5 main Data Science competences groups
  - Skills, tools and languages
- How it can be used
- Discussion questions



# Background Frameworks and Standards

- NIST SP1500 – 2015: Big Data Interoperability Framework (Volume 1-7)
  - Definitions of Data Science by NIST Big Data WG
- e-CFv3.0 - European e-Competence Framework for IT
  - Structured by 4 Dimensions and organizational processes
    - Competence Areas – Competences - Proficiency levels - Skills and Knowledge
- CWA 16458 (2012): European ICT Professional Profiles Family Tree
  - Defines 23 ICT profiles for common ICT jobs
- ESCO (European Skills, Competences, Qualifications and Occupations) framework
  - Standard for European job market since 2016
  - Intended inclusion of the Data Science occupations family – end of 2016
- ACM Classification of Computer Science – CCS (2012)
  - ACM Computer Science Body of Knowledge (CS-BoK) and ACM and IEEE Computer Science Curricula 2013 (CS2013)



# How it is made: Jobs market analysis and Community survey

## Demanded Data Science Competences and Skills

- Initial Analysis (period Aug – Sept 2015)
  - IEEE Data Science Jobs (World but majority US)
    - Collected > 120, selected for analysis > 30
  - LinkedIn Data Science Jobs (NL)
    - Collected > 140, selected for analysis > 30
  - Existing studies and reports + numerous blogs & forums
  - Automatic job market survey tool: to be operational in Fall 2016
- Analysis methods
  - Using data analytics methods: classification, clustering, expert evaluation
  - Research methods: Data collection - Hypothesis – Artefact - Evaluation
- Validation and community input
  - Survey on the general Data Science competences based on CF-DS
    - Domain related competences yet to be surveyed
  - Workshops and community feedback

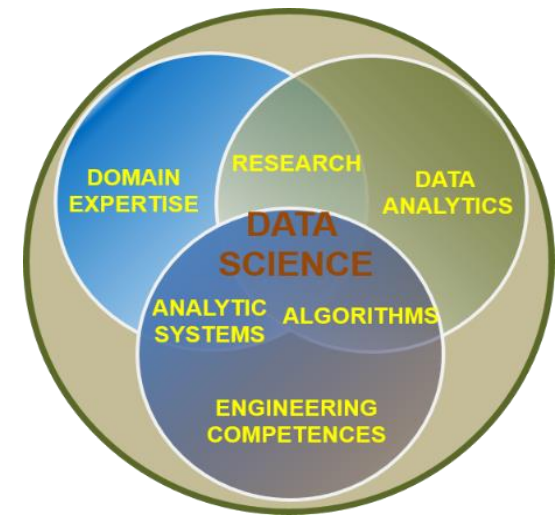




# Data Scientist definition by NIST

## Definitions by NIST Big Data WG (NIST SP1500 - 2015)

- *A **Data Scientist** is a practitioner who has sufficient knowledge in the overlapping regimes of expertise in business needs, domain knowledge, analytical skills, and programming and systems engineering expertise to manage the end-to-end scientific method process through each stage in the **big data lifecycle**.*
- ***Data Lifecycle in Big Data and Data Science***
- ***Data science** is the empirical synthesis of actionable knowledge and technologies required to handle data from raw data through the complete data lifecycle process.*

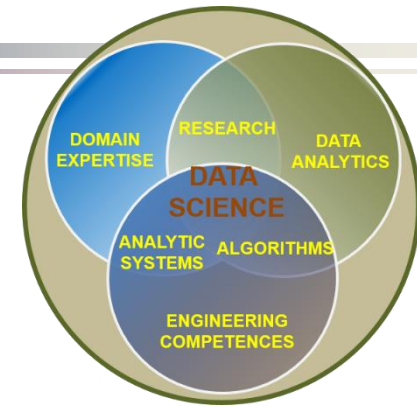


[ref] Legacy: NIST BDWG  
definition of Data Science



# Identified Data Science Competence Groups

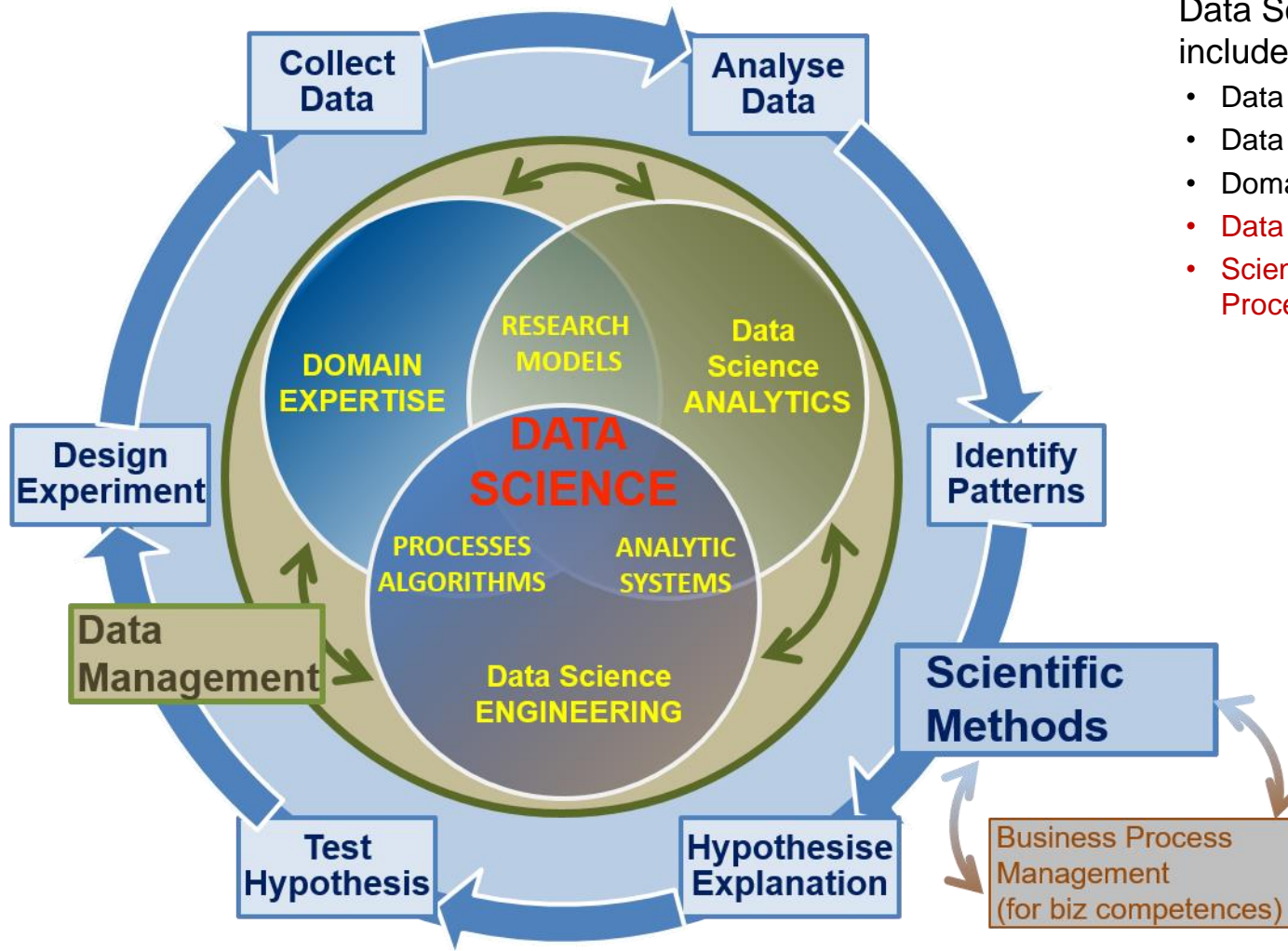
- Commonly accepted Data Science competences/skills groups include
  - Data Analytics or Business Analytics or Machine Learning
  - Engineering or Programming
  - Subject/Scientific Domain Knowledge
- EDISON identified 2 additional competence groups demanded by organisations
  - **Data Management, Curation, Preservation**
  - **Scientific or Research Methods and/vs Business Processes/Operations**
- Other skills commonly recognized aka “soft skills” or “social intelligence”
  - Inter-personal skills or team work, cooperativeness
- All groups need to be represented in Data Science curriculum and training programmes
  - Challenging task for Data Science education and training: multi-skilled vs team based
- Another aspect of integrating Data Scientist into organisation structure
  - General Data Science (or Big Data) **literacy** for all involved roles and management
  - Common agreed and understandable way of communication and **information/data presentation**
  - **Role of Data Scientist: Provide such literacy advice and guiding to organisation**



[ref] Legacy: NIST BDWG definition of Data Science



# Data Science Competence Groups - Research



Data Science Competence includes 5 areas/groups

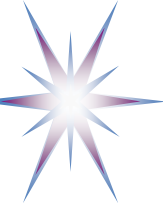
- Data Analytics
- Data Science Engineering
- Domain Expertise
- **Data Management**
- **Scientific Methods (or Business Process Management)**

## Scientific Methods

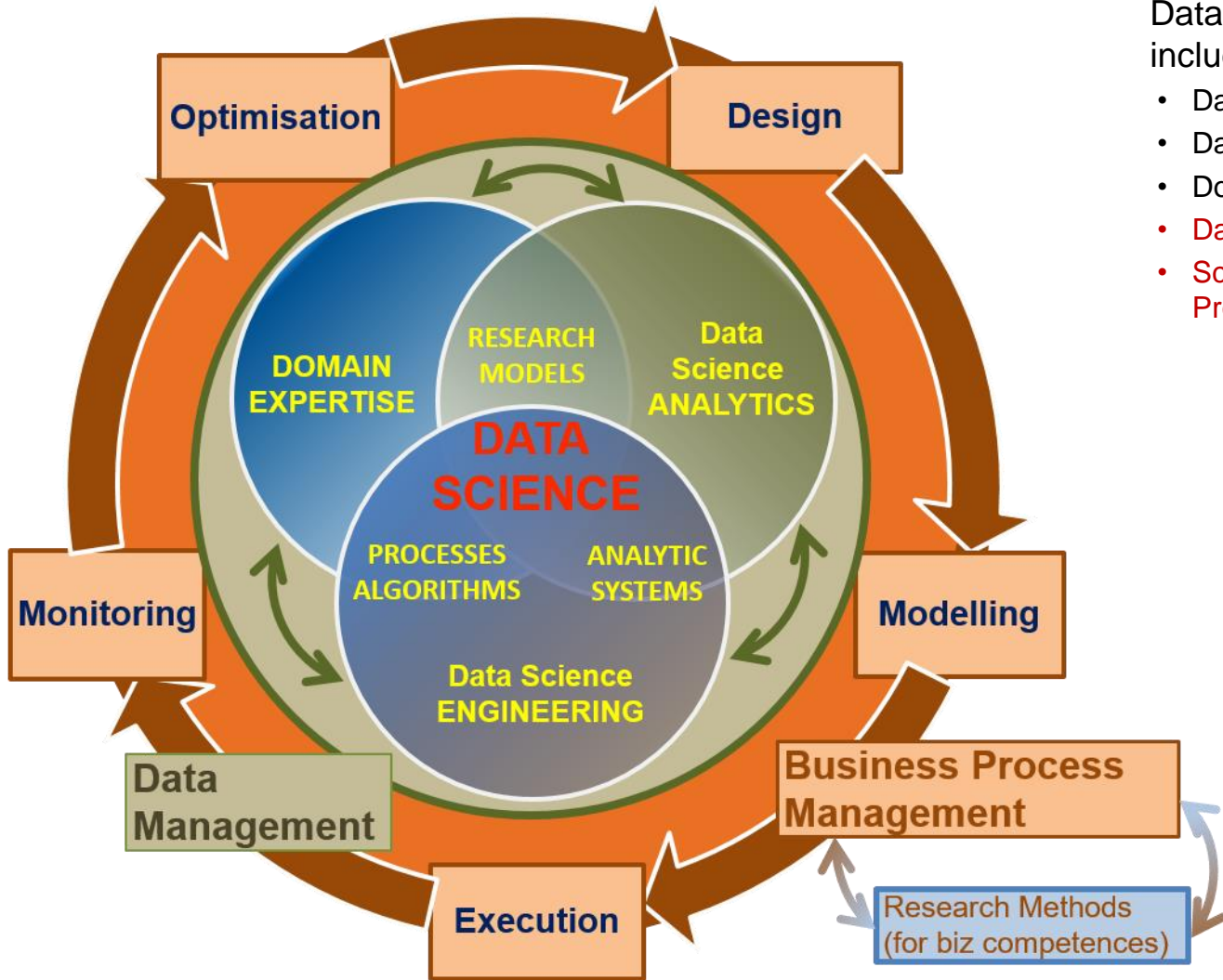
- Design Experiment
- Collect Data
- Analyse Data
- Identify Patterns
- Hypothesise Explanation
- Test Hypothesis

## Business Operations

- Operations Strategy
- Plan
- Design & Deploy
- Monitor & Control
- Improve & Re-design



# Data Science Competences Groups – Business



Data Science Competence includes 5 areas/groups

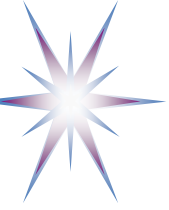
- Data Analytics
- Data Science Engineering
- Domain Expertise
- **Data Management**
- **Scientific Methods (or Business Process Management)**

## Scientific Methods

- Design Experiment
- Collect Data
- Analyse Data
- Identify Patterns
- Hypothesise Explanation
- Test Hypothesis

## Business Process Operations/Stages

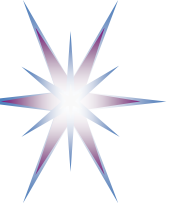
- Design
- Model/Plan
- Deploy & Execute
- Monitor & Control
- Optimise & Re-design



# Identified Data Science Competence Groups (Updated)

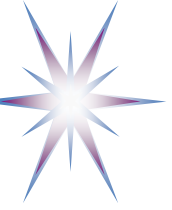
	Data Science Analytics (DSDA)	Data Management (DSDM)	Data Science Engineering (DSENG)	Research/Scientific Methods (DSRM)	Data Science Domain Knowledge, e.g. Business Processes (DSDK/DSBPM)
0	Use appropriate statistical techniques and predictive analytics on available data to deliver insights and discover new relations	Develop and implement data management strategy for data collection, storage, preservation, and availability for further processing.	Use engineering principles to research, design, develop and implement new instruments and applications for data collection, analysis and management	Create new understandings and capabilities by using the scientific method (hypothesis, test/artefact, evaluation) or similar engineering methods to discover new approaches to create new knowledge and achieve research or organisational goals	Use domain knowledge (scientific or business) to develop relevant data analytics applications, and adopt general Data Science methods to domain specific data types and presentations, data and process models, organisational roles and relations
1	<b>DSDA01</b> Use predictive analytics to analyse big data and discover new relations	<b>DSDM01</b> Develop and implement data strategy, in particular, Data Management Plan (DMP)	<b>DSENG01</b> Use engineering principles to design, prototype data analytics applications, or develop instruments, systems	<b>DSRM01</b> Create new understandings and capabilities by using scientific/ research methods or similar domain related development methods	<b>DSBPM01</b> Understand business and provide insight, translate unstructured business problems into an abstract mathematical framework
2	<b>DSDA02</b> Use statistical techniques to deliver insights	<b>DSDM02</b> Develop data models including metadata	<b>DSENG02</b> Develop and apply computational solutions	<b>DSRM02</b> Direct systematic study toward a fuller knowledge or understanding of the observable facts	<b>DSBPM02</b> Participate strategically and tactically in financial decisions
3	<b>DSDA03</b> Develop specialized ...	<b>DSDM03</b> Collect integrate data	<b>DSENG03</b> Develops specialized tools	<b>DSRM03</b> Undertakes creative work	<b>DSBPM03</b> Provides support services to other
4	<b>DSDA04</b> Analyze complex data	<b>DSDM04</b> Maintain repository	<b>DSENG04</b> Design, build, operate	<b>DSRM04</b> Translate strategies into actions	<b>DSBPM04</b> Analyse data for marketing
5	<b>DSDA05</b> Use different analytics	<b>DSDM05</b> Visualise complex data	<b>DSENG05</b> Secure and reliable data	<b>DSRM05</b> Contribute to organizational goals	<b>DSBPM05</b> Analyse optimise customer relations





# Identified Data Science *Skills/Experience* Groups

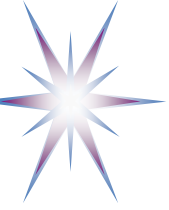
- **Group 1: Skills/experience related to competences**
  - Data Analytics and Machine Learning
  - Data Management/Curation (including both general data management and scientific data management)
  - Data Science Engineering (hardware and software) skills
  - Scientific/Research Methods or Business Process Management
  - Application/subject domain related (research or business)
  - Mathematics and Statistics
- **Group 2: Big Data (Data Science) tools and platforms**
  - Big Data Analytics platforms
  - Mathematics & Statistics applications & tools
  - Databases (SQL and NoSQL)
  - Data Management and Curation platform
  - Data and applications visualisation
  - *Cloud based platforms and tools*
- **Group 3: Programming and programming languages and IDE**
  - General and specialized development platforms for data analysis and statistics
- **Group 4: Soft skills or Social Intelligence**
  - Personal, inter-personal communication, team work, professional network



# Data Science Skill Groups related to Competences

	Data Analytics and Machine Learning	Data Management/ Curation	Data Science Engineering (hardware and software)	Scientific/ Research Methods	Personal/Inter-personal communication, team work	Application/subject domain (research or business), examples
1	Artificial intelligence, machine learning	Manipulating and analyzing complex, high-volume, high-dimensionality data from varying sources	Design efficient algorithms for accessing and analyzing large amounts of data	Analytical, independent, critical, curious and focused on results	Communication skills	Recommender or Ranking system
2	Machine Learning and Statistical Modelling	Data models and datatypes	Big Data solutions and advanced data mining tools	Confident with large data sets and ability to identify appropriate tools and algorithms	Inter-personal intra-team and external communication	Data Analytics for commercial purposes
3	Machine learning solutions and pattern recognition techniques	Experience of working with large data sets	Multi-core/distributed software, preferably in a Linux environment	Flexible analytic approach to achieve results at varying levels of precision	Network of contacts in Big Data community	Data sources and techniques for business insight and customer focus
4	Supervised and unsupervised learning	(non)relational and (un)-structured data	Databases, database systems, SQL and NoSQL	Interest in data science, exceptional analytical skills		Mechanism Design and/or Latent Dirichlet Allocation
5	Data mining	Cloud based data storage and data management	Statistical analysis languages and tooling			Game Theory
6	Markov Models, Conditional Random Fields	Data management planning	Cloud powered applications design			Copyright and IPR
7	Logistic Regression, Support Vector Machines	Metadata annotation and management				
8	Predictive analysis and statistics (including Kaggle platform)	Data citation, metadata, PID (*)				
9	(Artificial) Neural Networks					

**Mathematics foundation:**  
 Knowledge of mathematics, calculus, probability theory and statistics

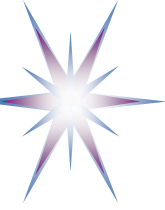


# Identified Big Data Tools and Programming Languages

	Big Data Analytics platforms	Math& Stats tools	Databases	Data/ applications visualization	Data Management and Curation platform
1	Big Data Analytics platforms	Advanced analytics tools (R, SPSS, Matlab, etc)	SQL and relational databases	Data visualization Libraries (D3.js, FusionCharts, Chart.js, other)	Data modelling and related technologies (ETL, OLAP, OLTP, etc)
2	Big Data tools (Hadoop, Spark, etc)	Data Mining tools: RapidMiner, others	NoSQL Databases	Visualisation software (D3, Processing, Tableau, <u>Gephi</u> , etc)	Data warehouses platform and related tools
3	Distributed computing tools a plus (Spark, MapReduce, Hadoop, Hive, etc.)	Mathlab	NoSQL, Mongo, Redis	<b>Online visualization tools (Datawrapper, Google Charts, Flare, etc)</b>	Data curation platform, metadata management (ETL, Curator's Workbench, DataUp, MIXED, etc)
4	Real time and streaming analytics systems (like Flume, Kafka, Storm)	Python	NoSQL, Teradata		Backup and storage management (iRODS, XArch, Nesstar, others)
5	Hadoop Ecosystem/platform	R, Tableau R	Excel		
6	Spotfire	SAS			
7	<b>Azure Data Analytics platforms (HDInsight, APS and PDW, etc)</b>	Scripting language, e.g. Octave			
8	<b>Amazon Data Analytics platform (Kinesis, EMR, etc)</b>	Statistical tools and data mining techniques			
9	<b>Other cloud based Data Analytics platforms, e.g. HortonWorks, Vertica LexisNexis HPC System</b>	Other Statistical computing and languages (WEKA, KNIME, IBM SPSS, etc)			

**Highlighted:  
Cloud based and online data analytics  
and data management platforms**





# Suggested Practical Application of the CF-DS

- Basis for the definition of the Data Science Body of Knowledge (DS-BoK) and Data Science Model Curriculum (MC-DS)
  - CF-DS => Learning Outcomes (MC-DS) => Knowledge Areas (DS-BoK)
  - CF-DS => Data Science taxonomy of scientific subjects and vocabulary
- Data Science professional profiles definition
  - Extend existing EU standards and occupations taxonomies: e-CFv3.0, ESCO, others
- Automatic job market monitor and survey tool
  - To be operational in Fall 2016
- Professional competence benchmarking (including CV and training programmes matching)
  - For customizable training and career development
- Professional certification
  - In combination with DS-BoK and professional competences benchmarking
- Vacancy construction tool for job advertisement (for HR)
  - Using controlled vocabulary and Data Science Taxonomy



# CF-DS - Discussion Questions

- DQ1: Collecting contribution from domain areas and experts
  - EDISON Survey for general Data Science competences and for target communities
- DQ2: Contributing to current standardisation activities
  - CEN e-CF workshop and CEN TC428 on standardisation of ICT competences and profiles
    - Extended mandate to define curriculum requirements/model
    - Part of New Skills Agenda for Europe
  - Extending CF-DS with dimensions on proficiency levels and skills and knowledge



# DSP – Data Science Professional Profiles

- Introduction to DSP
  - Data Science Professions family and ESCO taxonomy
    - ESCO – European
  - New 5 occupation/professional groups and 19 professional profiles
- How it can be used
- Discussion questions



# Data Science Professions Family



**Managers:** Chief Data Officer (CDO), Data Science (group/dept) manager, Data Science infrastructure manager, Research Infrastructure manager



**Professionals:** Data Scientist, Data Science Researcher, Data Science Architect, Data Science (applications) programmer/engineer, Data Analyst, Business Analyst, etc.



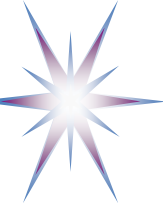
**Professional (database):** Large scale (cloud) database designers and administrators, scientific database designers and administrators



**Professional and clerical (data handling/management):** Data Stewards, Digital Data Curator, Digital Librarians, Data Archivists



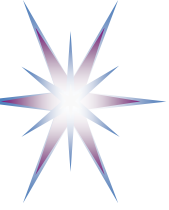
**Technicians and associate professionals:** Big Data facilities operators, scientific database/infrastructure operators



# Data Science Occupations: Extension for the ESCO (2016) taxonomy (1)

Professionals				
	Science and engineering professionals	<b>Data Science Professionals</b>	Data Science professionals not elsewhere classified	DSP04 Data Scientist
				DSP05 Data Science Researcher
				DSP08 (Big) Data Analyst
				DSP07 Data Science (Application) Programmer
				DSP09 Business Analyst
		Database and network professionals	Large scale (cloud) data storage designers and administrators	DSP14 Large scale (cloud) database designer*)
			Database designers and administrators	DSP15 Large scale (cloud) database administrator*)
			Database and network professionals not elsewhere classified	DSP16 Scientific database administrator*)
	Information and communications technology professionals	<b>Data Science technology professionals</b>	Data handling professionals not elsewhere classified	DSP12 Digital Librarian
				DSP13 Data Archivist
				DSP10 Data Steward
				DSP11 Data curator

19 DSP# Enumerated Data Science profiles defined by EDISON Framework

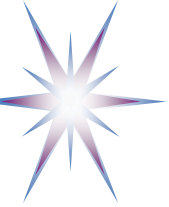


# Data Science Occupations: Extension for the ESCO taxonomy (2)

Technicians and associate professionals				
	Science and engineering associate professionals	<b>Data Science Technology Professionals</b>	Data Infrastructure engineers and technicians	DSP17 Big Data facilities Operators
				DSP18 Large scale (cloud) data storage operators
			Database and network professionals not elsewhere classified	DSP19 Scientific database operator*)
Managers				
	Production and specialised services managers	<b>Data Science/Big Data Infrastructure Managers</b>		DSP01/DSP02 Data Science/Big Data Infrastructure Manager
			Research Infrastructure Managers	DSP03 RI Manager
				DSP03 RI Data storage facilities manager
Clerical support workers				
	General and keyboard clerks			
	<b>Data handling support workers (alternative)</b>	<b>Data and information entry and access</b>	Digital Archivists and Librarians	DSP12* Digital Librarian
				DSP13* Data Archivist
				DSP10* Data Steward
				DSP11* Data curator

Profile ID	Data Science Profile title	Data Science Competences Groups (relevance 1 - low, 5 – high)				
		Data Analytics	Data Management	Data Science Engineering	Research Methods, Business methods	DS Subject Domain
<b>Managers</b>						
DSP01	Data Science (group) Manager	3	4	3	3	2
DSP02	Data Science Infrastructure Manager	2	4	4	2	2
DSP03	Research Infrastructure Manager	2	4	4	3	2
<b>Professionals</b>						
DSP04	Data Scientist	5	3	4	5	3
DSP05	Data Science Researcher	4	3	2	5	4
DSP06	Data Science Architect	4	3	5	3	3
DSP07	Data Science (Application) Programmer/Engineer	4	2	5	3	4
DSP08	Data Analyst	5	3	3	3	4
DSP09	Business Analyst	5	3	3	4	5
<b>Professional (data handling/ management)</b>						
DSP10	Data Stewards	3	5	3	3	3
DSP11	Digital data curator	1	5	2	2	3
DSP12	Digital Librarians	2	5	2	2	3
DSP13	Data Archivists	1	5	1	1	3
<b>Professional (database)</b>						
DSP14	Large scale (cloud) database designer	2	4	4	3	3
DSP15	Large scale (cloud) database administrator	2	4	3	2	3
DSP16	Scientific database administrator	2	4	3	2	3
<b>Technicians and associate professionals</b>						
DSP17	Big Data facilities Operator	1	4	4	2	3
DSP18	Large scale (cloud) data storage operator	1	4	3	1	1
DSP19	Scientific database operator	1	4	3	2	3

- Example mapping DSP profiles to competences
  - To be revised by experts and practitioners



# Data Science or Data Management Group/Department: Organisational structure and staffing - EXAMPLE

## Data Science or Data Management Group/Department

- Group Manager
  - Reporting to CDO/CTO/CEO
  - Providing cross-organizational services
  - Maintaining Data Value Chain
- Data Science Architect
- Data Analyst
- Data Science Application programmer
- Data Infrastructure/facilities administrator/operator: storage, cloud, computation
- Data stewards





# Data Science or Data Management Group/Department: Organisational structure and staffing - EXAMPLE

## Data Science or Data Management Group/Department

- >> Reporting to CDO/CTO/CEO
  - Providing cross-organizational services
  - Maintaining Data Value Chain

- (Managing) Data Science Architect (1)
- Data Analyst (1)
- Data Science Application programmer (2)
- Data Infrastructure/facilities administrator/operator: storage, cloud, computing (1)
- Data stewards, curators, archivists (3-5)

Estimated: Group of 10-12 specialists for research institution of 200-300 research staff.



# Discussion Questions

- DQ1: Mapping DSP profiles to competences
  - Collecting input from different professional groups to define specific competences
  - Using CF-DS competences and running targeted surveys
- DQ2: Including DSP family into the next version of ESCO
- DQ3: Mapping DSP profiles to MC-DS and academic disciplines
  - Link to offered education and training and required certification



# DS-BoK – Data Science Body of Knowledge

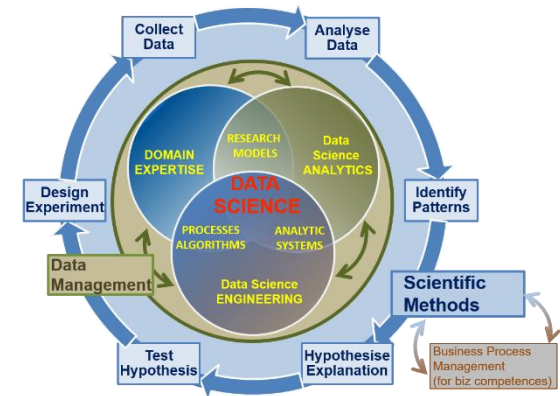
---

- Introduction to DS-BoK
  - Knowledge Area Groups
- How it can be used
- Discussion questions

# Data Science Body of Knowledge (DS-BoK)

## DS-BoK Knowledge Area Groups (KAG)

- KAG1-DNA: Data Analytics group including Machine Learning, statistical methods, and Business Analytics
- KAG2-DSE: Data Science Engineering group including Software and infrastructure engineering
- KAG3-DSDM: *Data Management group including data curation, preservation and data infrastructure*
- KAG4-DSRM: *Scientific/Research Methods group*
- KAG5-DSBP: Business process management group
- Data Science domain knowledge to be defined by related expert groups





# KAG3-DSDM: *Data Management group: data curation, preservation and data infrastructure*

DM-BoK version 2 “Guide for performing data management”

– 11 Knowledge Areas

(1) Data Governance

(2) Data Architecture

(3) Data Modelling and Design

(4) Data Storage and Operations

**(5) Data Security**

(6) Data Integration and Interoperability

**(7) Documents and Content**

(8) Reference and Master Data

(9) Data Warehousing and Business Intelligence

**(10) Metadata**

(11) Data Quality

Other Knowledge Areas motivated by RDA, European Open Data initiatives, European Open Data Cloud

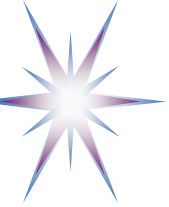
**(12) PID, metadata, data registries**

**(13) Data Management Plan**

**(14) Open Science, Open Data, Open Access, ORCID**

**(15) Responsible data use**

- Highlighted in red: Considered Research Data Management literacy (minimum required knowledge)



# Research Data Management Model Curriculum – Part of the EDISON Data Literacy Training

## A. Use cases for data management and stewardship

- Preserving the Scientific Record

## B. Data Management elements (organisational and individual)

- Goals and motivation for managing your data
- Data formats
- Creating documentation and metadata, metadata for discovery
- Using data portals and metadata registries
- Tracking Data Usage
- Handling sensitive data
- Backing up your data
- Data Management Plan (DMP) - to be a part of hands on session

## C. Responsible Data Use Section (Citation, Copyright, Data Restrictions)

## D. Open Science and Open Data (Definition, Standards, Open Data use and reuse, open government data)

- Research data and open access
- Repository and self- archiving services
- ORCID identifier for data
- Stakeholders and roles: engineer, librarian, researcher
- Open Data services: ORCID.org, Altmetric Doughnut, Zenodo

## E. Hands on:

- a) Data Management Plan design
- b) Metadata and tools
- c) Selection of licenses for open data and contents (e.g. Creative Common and Open Database)

To be supported by RDA WG on RDM Literacy

- BoF at RDA8 16 Sept 2016, Denver
- Contribution: Europe, US, AP
- Modular, customizable
- Localised: resources and languages
- Open Source under Creative Common Attribution



# MC-DS – Data Science Model Curriculum

---

- Introduction to MC-DS
  - Learning Outcomes based on CF-DS
  - Academic disciplines and courses based on CCS2013
- How it can be used to design a curriculum
- Discussion questions



# MC-DS – Data Science Model Curriculum

---

- Switch to website or document view





# MC-DS: Discussion Questions

- DQ1: Background knowledge/prerequisite
  - How to impose necessary mathematics and computer knowledge and skills in all Data Science programmes?
  - How to teach Scientific and Research Methods in the most effective way?
- DQ2: MC-DS and curriculum design for dynamically changing technology landscape
  - Creating mentality of self-reskilling workforce
- DQ 3: Top down vs bottom up approach in developing Data Science curricula
  - Universities practices
- DQ6. Benefits and challenges in adopting ACM Classification Computer Science (2012) and ACM/IEEE CS curriculum guidelines
  - How to extend current ACM classification to reflect required competences for Data Science? Including Domain related?



# Certification and Accreditation

- Certification scheme to be delivered by Sept 2016
- To be based on CF-DS and DS-BoK
- Using experience of the FitSM by EGI
- Talk to Małgorzata Krakowian (EGI)



# Data Science Professional Portal and Sustainability

---

- Switch to separate presentation by Andrea Manieri (Engineering)



# Discussion

- Questions
- Observations
- Suggestions

- Survey Data Science Competences [1]: Invitation to participate [https://www.surveymonkey.com/r/EDISON\\_project - Defining Data science profession](https://www.surveymonkey.com/r/EDISON_project_-_Defining_Data_science_profession)

- Community discussion documents: Request for comments
  - Data Science Competence Framework <http://edison-project.eu/data-science-competence-framework-cf-ds>
  - Data Science Body of Knowledge <http://edison-project.eu/data-science-body-knowledge-ds-bok>
  - Data Science Model Curriculum <http://edison-project.eu/data-science-model-curriculum-mc-ds>
  - Data Science Professional Profiles <http://edison-project.eu/data-science-professional-profiles>

EDISON building the data science profession

EDISON project: Defining Data science profession

Data Analytics skills and competencies for data science profession

\* 19. What are the competences and skills a data scientist should have on data analytics:

	Not relevant	Factual and theoretical knowledge	Comprehensive factual and theoretical knowledge	Advanced knowledge of a field, critical understanding of theories and principles	Highly specialized knowledge, Critical awareness, interface between different fields	Knowledge at the most advanced frontier of a field
Use appropriate statistics to provide insight on data	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Use appropriate techniques for analysing data (A/B Testing, Association rule Learning, Crowdsourcing, Data fusion and integration, Data Mining, Ensemble learning, Machine learning)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Use Predictive analytics to analyse big data and discover new relation	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Research and analyse complex data sets, combine different sources of data to improve analysis	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Develop specialised analytics to enable agile decision making	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>



# Suggested e-CF Competences for Data Science:

Presented at eCF Workshop meeting on 14 April 2016

## A. PLAN and Design (9 basic competences)

- A.10\* Organisational workflow/processes model definition/formalisation
- A.11\* Data models and data structures

## B. BUILD: Develop and Deploy/Implement (6 basic competences)

- B.7\* Apply data analytics methods (to organizational processes/data)
- B.8\* Data analytics application development
- B.9\* Data management applications and tools
- B.10\* Data Science infrastructure deployment

## C. RUN: Operate (4 basic competences)

- C.5\* User/Usage data/statistics analysis
- C.6\* Service delivery/quality data monitoring

## D. ENABLE: Use/Utilise (12 basic competences)

- D10. Information and Knowledge Management (powered by DS)
- D.13\* Data presentation/visualisation, actionable data extraction
- D.14\* Support business processes/roles with data and insight (support to D.5, D.6, D.7, D.12)
- D.15\* Data management/preservation/curation with data and insight

## E. MANAGE (9 basic competences)

- E.10\* Support Management and Business Improvement with data and insight (support to E.5, E.6)
- E.11\* Data analytics for (business) Risk Analysis/Management (support to E.3)
- E.12\* ICT and Information security monitoring and analysis (support to E.8)

**15 Data Science Competences proposed covering different organizational roles and workflow stages**

- Data Scientist roles are crossing multiple org roles and workflow stages

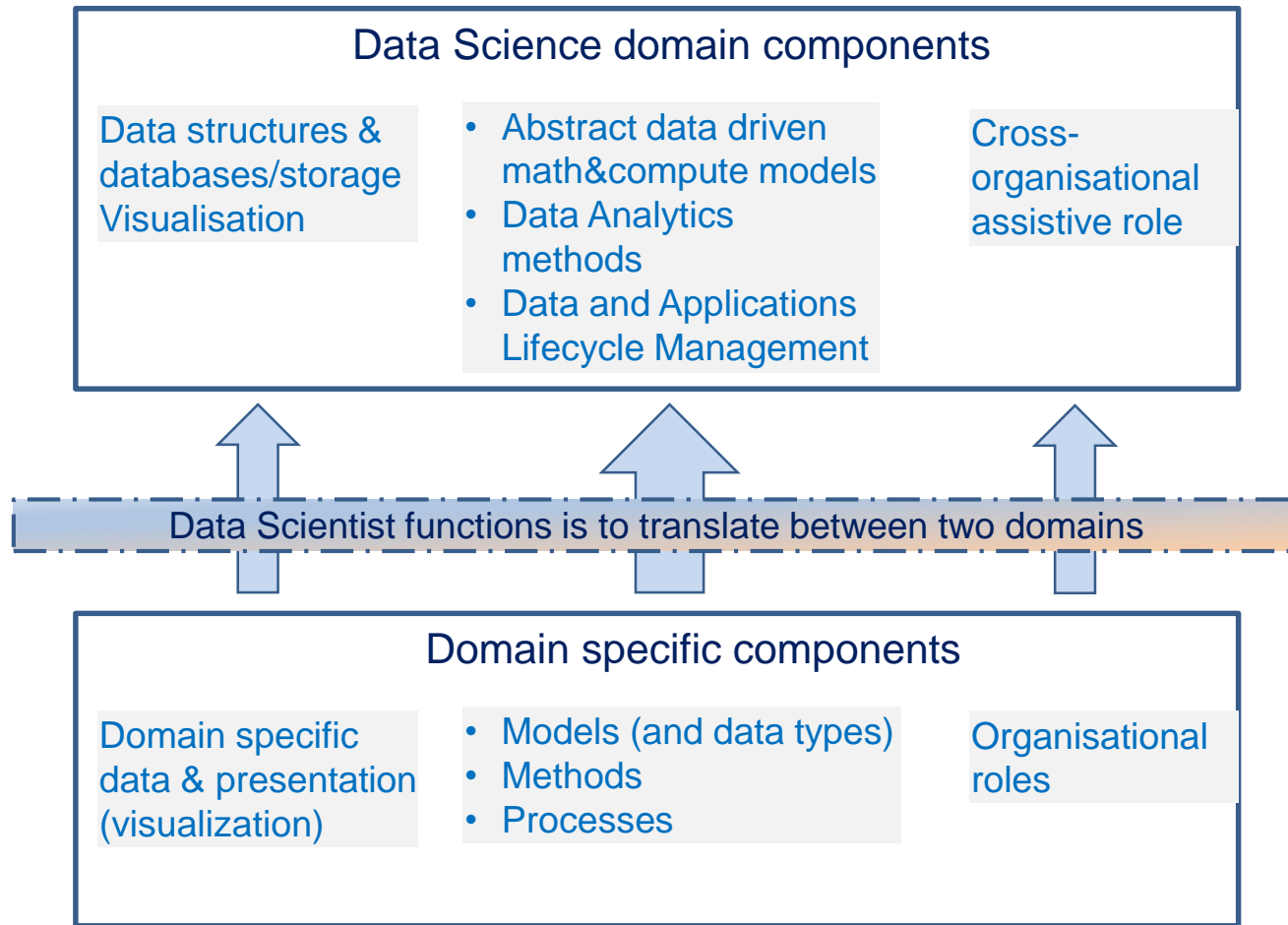


# Data Scientist and Subject Domain Specialist

- **Subject domain components**
  - Model (and data types)
  - Methods
  - Processes
  - Domain specific data and presentation/visualization methods
  - Organisational roles and relations
- **Data Scientist is an assistant to Subject Domain Specialists**
  - Translate subject domain Model, Methods, Processes into abstract data driven form
  - Implement computational models in software, build required infrastructure and tools
  - Do (computational) analytic work and present it in a form understandable to subject domain
  - Discover new relations originated from data analysis and advice subject domain specialist
  - Present/visualise information in domain related actionable way
  - Interact and cooperate with different organizational roles to obtain data and deliver results and/or actionable data
- **Overall goal: Maintain the Data Value Chain:**
  - Data Integration => Organisation/Process/Business Optimisation => **Innovation**



# Data Science and Subject Domains



**Data Scientist role is to maintain the Data Value Chain (domain specific):**

- Data Integration => Organisation/Process/Business Optimisation => **Innovation**