

# Services to establish Data Science Profession: EDISON Data Science Framework and activities

Yuri Demchenko, EDISON  
University of Amsterdam

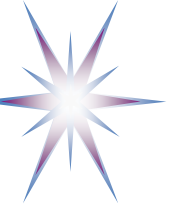


**Third Community Workshop on Open  
Science Cloud**

7 April 2016  
Science Park Amsterdam

EDISON – **E**ducation for **D**ata Intensive  
**S**cience to **O**pen **N**ew science frontiers

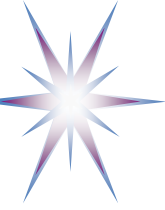
Grant 675419 (INFRASUPP-4-2015: CSA)



# Outline

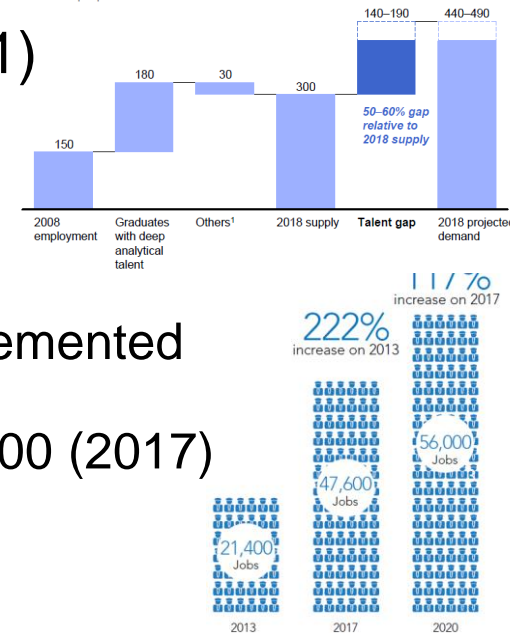
- Demand for Data Science and data related professions
- European initiatives related to Digital Single Market (DSM) and demand to data related competences and skills
- EDISON Data Science Framework
  - From Data Science Competences to Body of Knowledge and Model Curriculum
- Data Science Competence Framework: Essential competences and skills
- Taxonomy: Data Science occupations family
- Data Science Body of Knowledge (DS-BoK)
  - Knowledge areas and academic disciplines
- Further steps

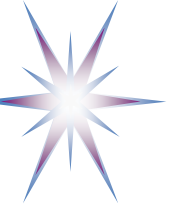




# Demand for Data Science and data related professions

- 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 EOSC (2016) identified need for data experts and data stewards
  - **Estimation: More than 500,000 data stewards (1 per every 20 scientists or 5% funding)**

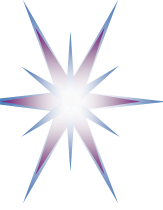




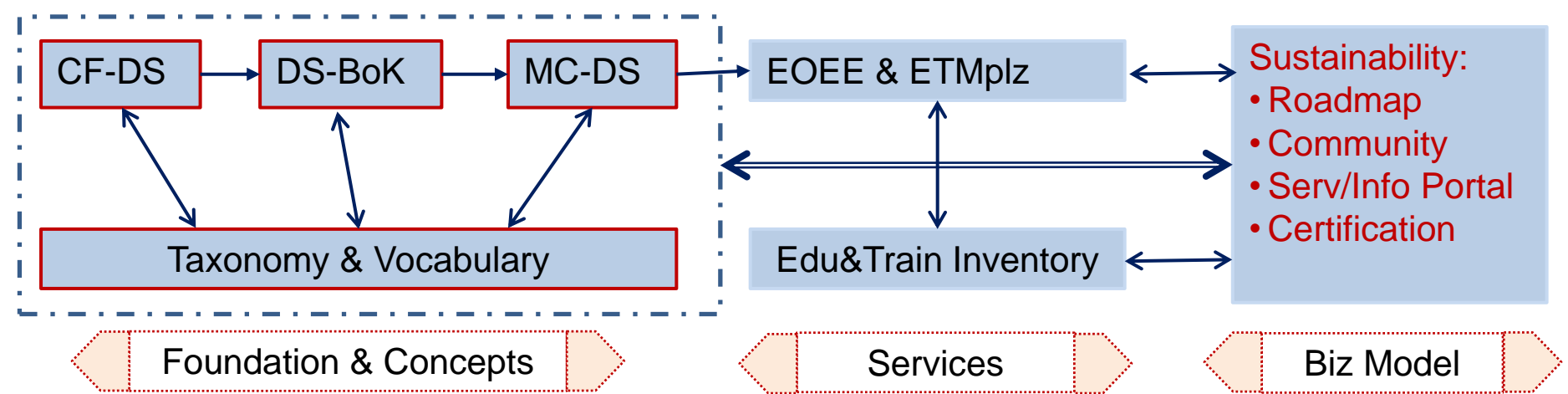
# Demand for Data Science and Data related professions

## – EOSC and European RIs

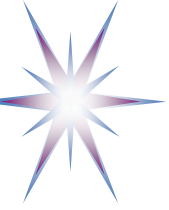
- HLEG report on EOSC identified need for data experts and data stewards
  - **Estimation: More than 500,000 data stewards**
    - 1 per every 20 scientists or 5% funding
  - Challenges and observation:
    - Shortage of data experts
    - Data literacy gap: 'Valley of death' between (e-)infrastructure providers and domain specialists
  - Requirements: Core data experts need to be trained and their career perspective significantly improved
  - Support: Professional data management and Long term data stewardship
  - **Implementation recommendations:**
    - I3: Fund a concentrated effort to locate and develop Data Expertise in Europe
    - I5: Make adequate data stewardship mandatory for all research proposals
    - II4: Train the data experts to bridge between e-Infra and ESFRI
    - II5: Assist data stewardship planning and exec tools for all researchers
  - **Open Science Presidency Conference 4-5 April 2016:**
    - Recognised data steward as a profession



# EDISON Framework: Concept, Services, Sustainability

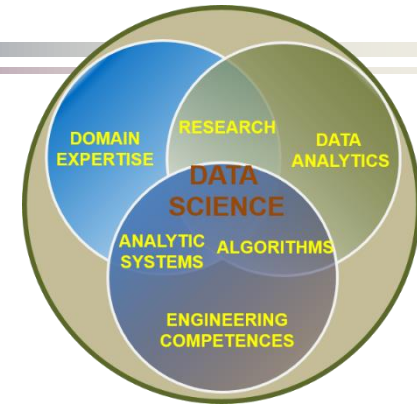


- EDISON Framework components
  - CF-DS – Data Science Competence Framework
  - DS-BoK – Data Science Body of Knowledge
  - MC-DS – Data Science Model Curriculum
  - Data Science Taxonomy and Scientific Disciplines Classification
  - EOEE - EDISON Online Education Environment
- Background: EU Competence Frameworks and Profiles
  - e-CFv3.0 - European e-Competence framework for IT
  - CWA 16458 (2012): European ICT Professional Profiles Family Tree
  - ESCO (European Skills, Competences, Qualifications and Occupations) framework

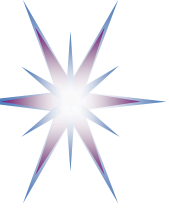


# Identified Data Science Competence Groups

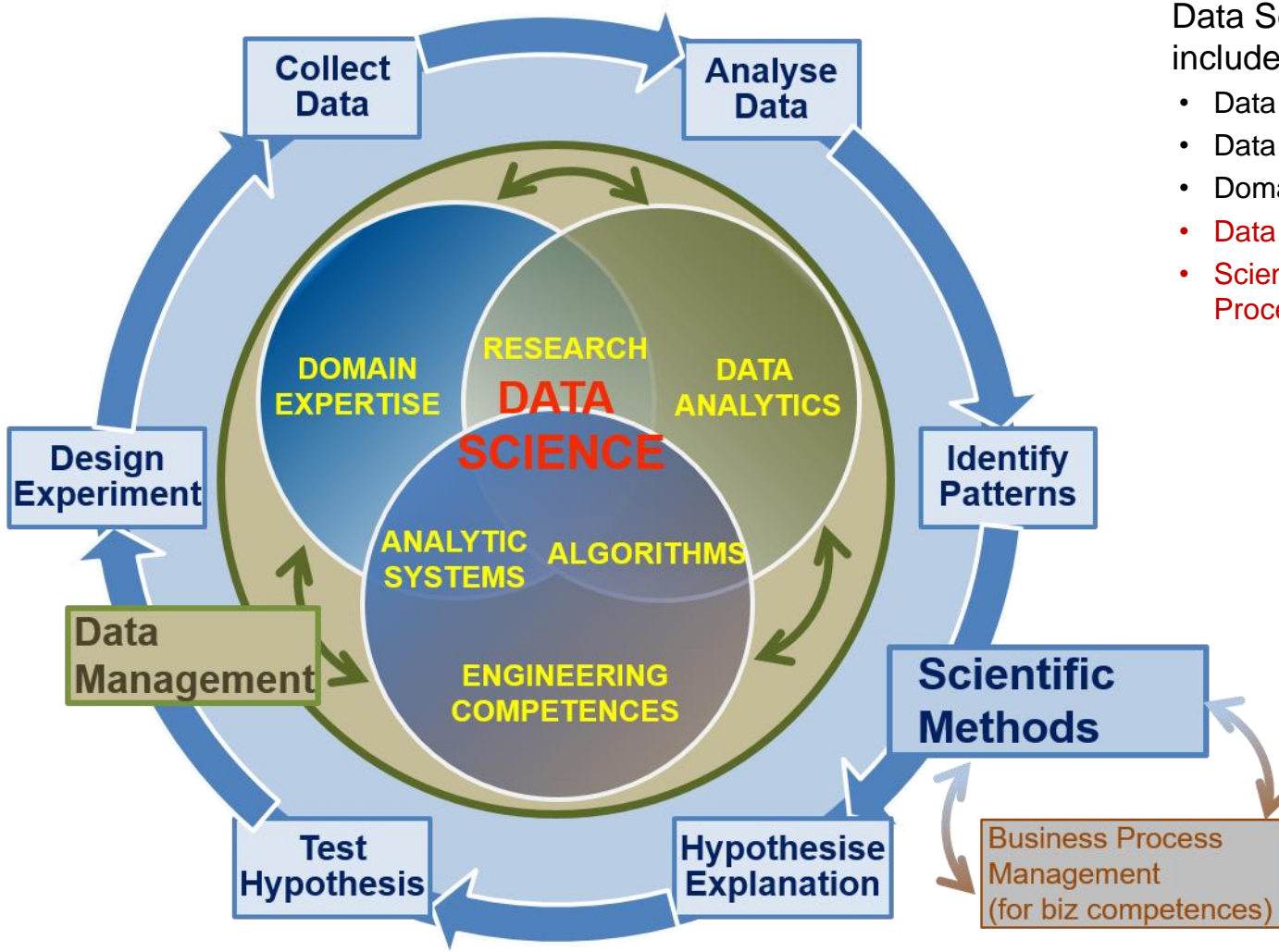
- Traditional/known Data Science competences/skills groups include
  - Data Analytics or Business Analytics or Machine Learning
  - Engineering or Programming
  - Subject/Scientific Domain Knowledge
- EDISON identified 2 additional competence groups demanded by organisations
  - **Data Management, Curation, 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
- 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

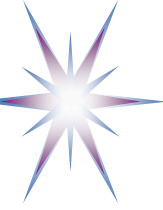
- Data Analytics
- Data Science Engineering
- Domain Expertise
- **Data Management**
- **Scientific Methods (or Business Process Management)**

## Scientific Methods

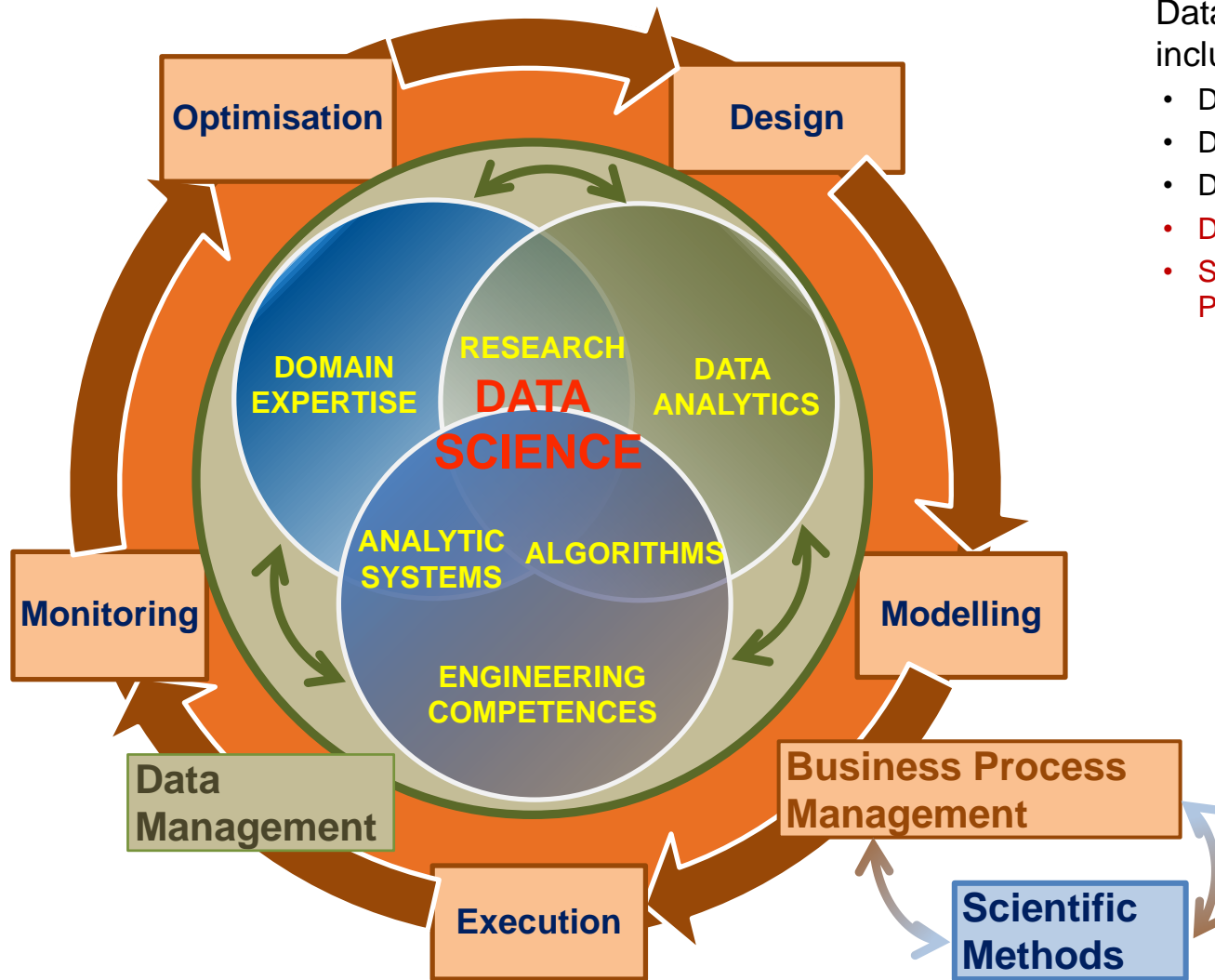
- Design Experiment
- Collect Data
- Analyse Data
- Identify Patterns
- Hypothesise Explanation
- Test Hypothesis

## Business Operations

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



# Data Science Competences Groups – Business



Data Science Competence includes 5 areas/groups

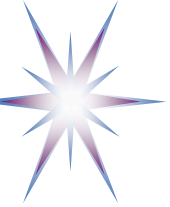
- Data Analytics
- Data Science Engineering
- Domain Expertise
- **Data Management**
- **Scientific Methods (or Business Process Management)**

## Scientific Methods

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

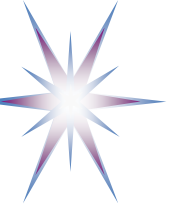
## Business Process Operations/Stages

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



# Identified Data Science Competence Groups

	Data Analytics (DA)	Data Management/ Curation (DM)	DS Engineering (DSE)	Scientific/Research Methods (DSRM)	DS Domain Knowledge (including Business Apps)
1	Use appropriate statistical techniques on available data to deliver insights	<b>Develop and implement data strategy</b>	Use engineering principles to research, design, or develop structures, instruments, machines, experiments, processes, systems, theories, or technologies	<b>Create new understandings and capabilities by using the scientific method's hypothesis, test, and evaluation techniques; critical review; or similar engineering research and development methods</b>	Understand business and provide insight, translate unstructured business problems into an abstract mathematical framework
2	Use predictive analytics to analyse big data and discover new relations	<b>Develop data models including metadata</b>	Develops specialized data analysis tools to support executive decision making	<b>Direct systematic study toward a fuller knowledge or understanding of the fundamental aspects of phenomena and of observable facts, and discovers new approaches to achieve goals</b>	Use data to improve existing services or develop new services
3	Research and analyze complex data sets, combine different sources and types of data to improve analysis.	<b>Integrate different data source and provide for further analysis</b>	Design, build, operate relational non-relational databases	<b>Undertake creative work, making systematic use of investigation or experimentation, to discover or revise knowledge of reality, and uses this knowledge to devise new applications</b>	Participate strategically and tactically in financial decisions that impact management and organizations
4	Develop specialized analytics to enable agile decision making	<b>Develop and maintain a historical data repository of analysis</b>	Develop and apply computational solutions to domain related problems using wide range of data analytics platforms	<b>Apply ingenuity to complex problems, develop innovative ideas</b>	Recommends business related strategic objectives and alternatives and implements them
5		<b>Collect and manage different source of data</b>	Develop solutions for secure and reliable data access	<b>Ability to translate strategies into action plans and follow through to completion.</b>	Provides scientific, technical, and analytic support services to other organisational roles
6		<b>Visualise complex and variable data.</b>	Develop algorithms to analyse multiple source of data	<b>Influence the development of organizational objectives</b>	Analyse multiple data sources for marketing purposes
7			Prototype new data analytics applications		Analyse customer data to identify/optimize customer relations actions

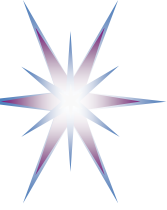


# Identified Data Science *Skills/Experience* Groups

- **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
  - Application/subject domain related (research or business)
  - Mathematics and Statistics
- **Group 2: Big Data (Data Science) tools and platforms**
  - Big Data Analytics platforms
  - Math & Stats apps & 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 (also called social intelligence or soft skills)

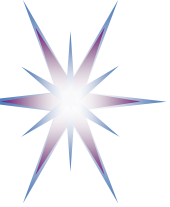
## Big Data Tools and Programming Languages

- Big Data Analytics platforms
- Math& Stats tools
- Databases
- Data/applications visualization
- Data Management and Curation



# Identified Data Science Skill Groups

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



# Suggested e-CF Competences for Data Science:

## Next eCF Workshop meeting – 14 April 2016

### A. PLAN and Design

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

### B. BUILD: Develop and Deploy/Implement

- 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

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

### D. ENABLE: Use/Utilise

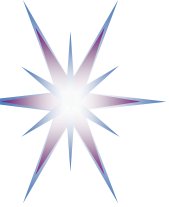
- D10. Information and Knowledge Management (powered by DS)
- D.13\* Data presentation/visualisation, actionable data extraction
- D.14\* Support business processes/roles with data and insight (support to D.5, D.6, D.7, D.12)
- D.15\* Data management/preservation/curation with data and insight

### E. MANAGE

- 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)

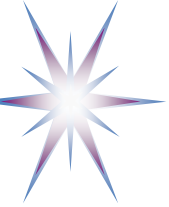
15 Data Science Competences proposed covering different organizational roles and workflow stages

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



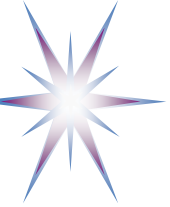
# Possible Data Scientist profiles/roles as extension to CWA16458 (2012)

- Data Analyst, Business Analyst
  - Data Mining
  - Machine Learning
- Digital Librarian, Data Archivist, Data Curator, Data Steward
  - Data Management related competences
- Data Science Engineer/Administrator/Programmer
  - Data analytics applications development
  - Scientific programming
  - Data Science/Big Data Infrastructure development/operation
- Data Science Researcher
  - Data Science research methods
  - Data models and structures
- Data Scientist in subject/research domain
- Research e-Infrastructure brings its own specifics to required competences and skills definition



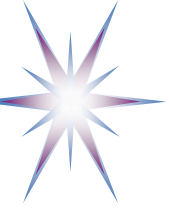
# Data Science occupations in ESCO taxonomy (1)

Professionals				
	Science and engineering professionals	<b>Data Science Professionals</b>	Data Science professionals not elsewhere classified	Data Scientist
				Data Science Researcher
				(Big) Data Analyst
				Data Science (Application) Programmer
				Business Analyst
		Database and network professionals	Large scale (cloud) data storage designers and administrators	Large scale (cloud) database designer*)
			Database designers and administrators	Large scale (cloud) database administrator*)
			Database and network professionals not elsewhere classified	Scientific database administrator*)
	Information and communications technology professionals	<b>Data Science technology professionals</b>	Data handling professionals not elsewhere classified	Digital Librarian
				Data Archivist
				Data Steward
				Data curator



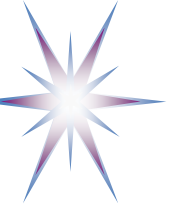
# Data Science occupations in ESCO taxonomy (2)

Technicians and associate professionals				
	Science and engineering associate professionals	<b>Data Science Technology Professionals</b>	Data Infrastructure engineers and technicians	Big Data facilities Operators
				Large scale (cloud) data storage operators
			Database and network professionals not elsewhere classified	Scientific database operator*)
Managers				
	Production and specialised services managers	<b>Data Science/Big Data Infrastructure Managers</b>		Data Science/Big Data Infrastructure Manager
			Research Infrastructure Managers	RI Manager
				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	Digital Librarian
				Data Archivist
				Data Steward
				Data curator
EOSC-3 Wsh - 7 April 2016		Building Data Science Profession		15



# Data Science or Data Management Support Group: Organisational structure and staffing

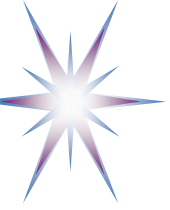
- Manager
- Data Science Architect
- Data Analyst
- Data Science Application programmer
- Data Infrastructure/facilities administrator/operator: storage, cloud, computation
- Data stewards



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

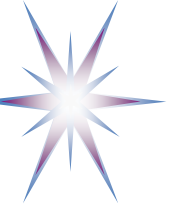
- Manager (1)
- Data Science Architect (1)
- Data Analyst (2)
- Data Science Application programmer (3)
- Data Infrastructure/facilities administrator/operator: storage, cloud, computation (1-2)
- Data stewards (5-10)

Group of 14-20 specialists



# Education and Training

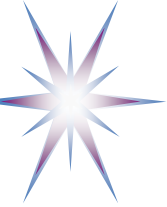
- Foundation and methodological base
  - Data Science Body of Knowledge (DS-BoK)
    - Taxonomy and classification of Data Science related scientific subjects
  - Instructional methodologies and teaching models
- Platforms and environment
  - Virtual labs, datasets, developments platforms
  - Access control and accounts/identity management
  - Online education environment and courses management
- Services
  - Individual profiling tools (competence assessment)
  - Knowledge evaluation tools
  - Certifications and training
  - Courses catalog and repository: Education and training marketplace



# 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*

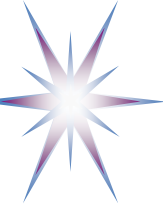
DM-BoK version 2 “Guide for performing data management”

– 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



# Topics considered for the Data Management (Literacy) Training – Working draft

## **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

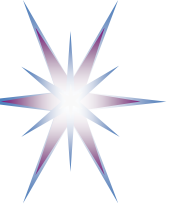
## **E. Hands on:**

### a) Data Management Plan design

- Why Create a Data Management Plan?
- Elements of a Data Management Plan
- Organization and Standards
- Data management plan checklist

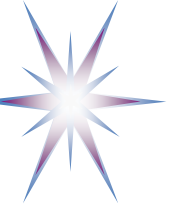
### b) Metadata and tools

### c) Selection of licenses for open data and contents (Creative Common and Open Database License)




# Further Steps

- Define a taxonomy and classification for DS competences and skills as a basis for more formal CF-DS definition
  - Closer look at skills, tools and platforms
- **Run surveys for target communities**  
[https://www.surveymonkey.com/r/EDISON\\_project - Defining Data science profession](https://www.surveymonkey.com/r/EDISON_project_-_Defining_Data_science_profession)
  - **Plan a number of key interviews, primarily experts and top executives at universities and companies**
- Proceed with suggested e-CF3.0 extensions and participate in the next e-CF meetings
  - Talk to national e-CF bodies or adopters if available
- Provide feedback and contribution to ESCO with the definition of the Data Science professions family
- Suggest ACM2012 Classification extensions and contact ACM people
- Provide input to DS-BoK definition following from CF-DS
  - Link/Map to taxonomy of academic and educational and training courses
- Create open community forum to collect contribution
  - **CF-DS and DS-BoK documents are on public comments available from EDISON website**  
<http://www.edison-project.eu/data-science-competence-framework-cf-ds>  
<http://www.edison-project.eu/data-science-body-knowledge-ds-bok>
  - Start related Social Network groups to promote already obtained results and obtain feedback and community contribution



# Discussion

- Questions
- Observations
- Survey: Invitation to participate

 **EDISON**  
building the data  
science profession

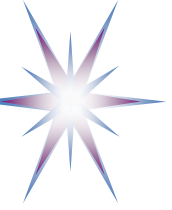
**EDISON project: Defining Data science profession**

**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"/>

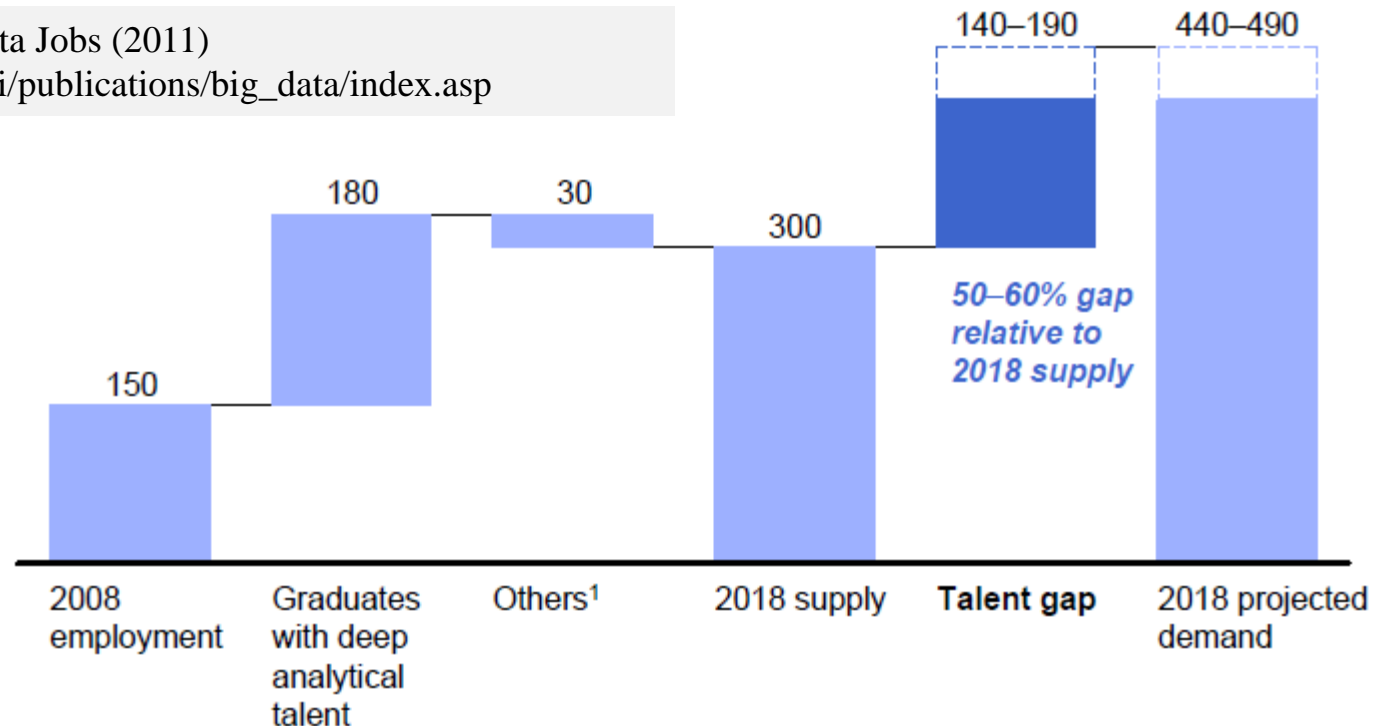
[https://www.surveymonkey.com/r/EDISON\\_project - Defining Data science profession](https://www.surveymonkey.com/r/EDISON_project_-_Defining_Data_science_profession)



# Data Scientist: New Profession and Opportunities

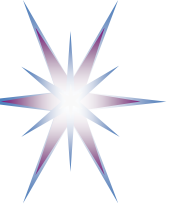
McKinsey Institute on Big Data Jobs (2011)

[http://www.mckinsey.com/mgi/publications/big\\_data/index.asp](http://www.mckinsey.com/mgi/publications/big_data/index.asp)



- There will be a shortage of talent necessary for organizations to take advantage of Big Data.
  - By 2018, the United States alone could face a shortage of 140,000 to 190,000 people with deep analytical skills as well as
  - 1.5 million managers and analysts with the know-how to use the analysis of big data to make effective decisions

SOURCE:US Bureau of Labor Statistics; US Census; Dun & Bradstreet; company interviews; McKinsey analysis



# EXAMPLE: Use of e-CF3.0 for Defining Profile of RI Technical (part of RDA IG-ETRD work)

## A. PLAN and DESIGN

- A.2. Service Level Management
- A.3. Product / Service Planning
- A.5. Application Design
- A.4. Architecture Design

Additional

- A.6. Sustainable Development
- A.7. Innovating and Technology Trend Monitoring
- A.8. Business/Research Plan Development and Grant application
- A.1. RI and Research Strategy Alignment

## B. BUILD: DEVELOP and DEPLOY/IMPLEMENT

- B.1. Application Development (Reqs Engineering, Function Specs, API, HCI)
- B.2. Component Integration
- B.3. Testing (RI services and Scientific Apps)
- B.4. Solution/Apps Deployment

Additional

- B.5. Documentation Production
- B.6. Systems Engineering (DevOps)

## C. OPERATE (RUN)

- C.1. User Support
- C.2. Service Delivery
- C.3. Problem Management

Additional

- C.4. Change Support (Upgrade/Migration)

## D. USE: UTILISE (ENABLE)

- D.1. Scientific Applications Integration (on running RI)
- D.5. Data collection and preservation
- D.4. New requirements and change Identification
- D.6. Education and Training Provision

Additional

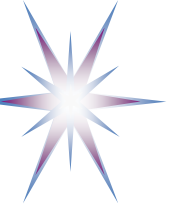
- D.2. Information Security Strategy Development
- D.3. RI/ICT Quality Strategy Development
- D.7. Purchasing/Procurement
- D.8. Contract Management
- D.9. Personnel Development
- D.10. Dissemination and outreach

## E. MANAGE

- E.1. Overall RI management (by systems and components)
- E.5. Information/Data Security Management

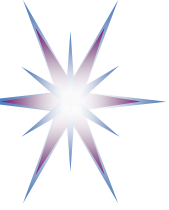
Additional

- E.6. Data Management (including planning and lifecycle management, curation)
- E.4. RI Security and Risk/Dependability Management
- E.2. Project and Portfolio Management
- E.3. ICT Quality Management and Compliance
- E.7. RI/IS Governance



# 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



# Data Science and Subject Domains

