

Profiling Data Science Education and Training

based on EDISON Data Science Framework (EDSF)



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

Yuri Demchenko University of Amsterdam

AACSB Conference on Data Science at Business Schools Amsterdam, May 16 / 17, 2017

Grant 675419 (INFRASUPP-4-2015: CSA)

Project: Building Data Science Teams that work



EDISON Data Science Framework (EDSF) Release 1 (October 2016)



- 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 Professional profiles
 - Data Science Taxonomies and Scientific Disciplines Classification
 - EOEE EDISON Online Education Environment

EDSF: How CF-DS was constructed

- Background: Standards and Best Practices
- Jobs market analysis: Demanded Data Science Competences and Skills



- e-CFv3.0 European e-Competence Framework for IT
 - Structured by 4 Dimensions and organizational processes
 - Competence Areas: Plan Build Run Enable Manage
 - Competences: total defined 40 competences
 - Proficiency levels: identified 5 levels linked to professional education 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
 - Expected inclusion of the Data Science occupations family end 2017
- ACM Classification of Computer Science CCS (2012)
- ACM Computer Science Body of Knowledge (CS-BoK) and ACM and IEEE Computer Science Curricula 2013 (CS2013)



- e-CFv3.0 European e-Competence Framework
 - Structured by 4 Dimensions and
- Currently work on e-CF4 is moved to CEN TC 428 To be extended with Data Science competences Competer - ivianage ievels: identified 5 levels linked to professional education 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
 - Expected inclusion of the Data Science occupations family end 2017
- ACM Cla New Joint Initiative ACM, IEEE, ASA, AAAS, AIS, ACH CM and ACM Co To develop Data Science curriculum IEEE Computer Science Curricula 2010



Jobs market analysis: Demanded Data Science Competences and Skills

- Initial Analysis (period Aug Sept 2015) -> Continuous monitoring (in development)
 - 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
- Analysis methods
 - Data analytics methods: classification, clustering, feature extraction
 - Research methods: Data collection Hypothesis Artefact -Evaluation
 - Expert evaluation by EDISON Liaison Groups (ELG), multiple workshops

Data Science Professions Family



Icons used: Credit to [ref] https://www.datacamp.com/community/tutorials/data-science-industry-infographic

EDISON Data Science Framework

Data Science Competences Groups – Business



Data Science Competences include 5 groups

- Data Science Analytics
- Data Science Engineering
- Domain Knowledge and Expertise
- Data Management
- Scientific Methods or Business
 Process Management

Business Process Operations/Stages

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

Scientific Methods

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



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 21st century skills

- Critical thinking, personal, inter-personal communication, team work, professional network



Identified Data Science Competence Groups

	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	DSDA01 Use predictive analytics to analyse big data and discover new relations	DSDM01 Develop and implement data strategy, in particular, Data Management Plan (DMP)	DSENG01 Use engineering principles to design, prototype data analytics applications, or develop instruments, systems	DSRM01 Create new understandings and capabilities by using scientific/ research methods or similar domain related development methods	DSBPM01 Understand business and provide insight, translate unstructured business problems into an abstract mathematical framework
2	DSDA02 Use statistical techniq to deliver insights	DSDM02 Develop data models including metadata	DSENG02 Develop and apply computational solutions	DSRM02 Direct systematic study toward a fuller knowledge or understanding of the observable facts	DSBPM02 Participate strategically and tactically in financial decisions
3	DSDA03 Develop specialized	DSDM03 Collect integrate data	DSENG03 Develops specialized tools	DSRM03 Undertakes creative work	DSBPM03 Provides support services to other
4	DSDA04 Analyze complex data	DSDM04 Maintain repository	DSENG04 Design, build, operate	DSRM04 Translate strategies into actions	DSBPM04 Analyse data for marketing
5	DSDA05 Use different analytics	DSDM05 Visualise cmplx data	DSENG05 Secure and reliable data	DSRM05 Contribute to organizational goals	DSBPM05 Analyse optimise customer relatio

Amsterdam, 17 May 2017

EDISON Data Science Framework



Individual Competences Benchmarking



Individual Education/Training Path based on Competence benchmarking

- Red polygon indicates the chosen professional profile: Data Scientist (general)
- Green polygon indicates the candidate or practitioner competences/skills profile
- Insufficient competences (gaps) are highlighted in *red*
 - DSDA01 DSDA06 Data Science Analytics
 - DSRM01 DSRM05 Data Science Research Methods
- Can be use for team skills match marking and organisational skills management

[ref] For DSP Profiles definition and for enumerated competences refer to EDSF documents CF-DS and DSP Profiles.

DSP Profiles mapping to ESCO Taxonomy High Level Groups and CF-DS Competences



DSP Profiles mapping to corresponding CF-DS Competence Groups
 – Relevance level from 5 – maximum to 1 – minimum

Education and Training – Part of EDSF

- Foundation and methodological base
 - Data Science Body of Knowledge (DS-BoK)
 - Taxonomy and classification of Data Science related scientific subjects
 - Data Science Model Curriculum (MC-DS)
 - Set Learning Units mapped to CF-DS Learning and DS-BoK Knowledge Areas/Units
 - Instructional methodologies and teaching models
- Platforms and environment
 - Virtual labs, datasets, developments platforms
 - Online education environment and courses management
- Services
 - Individual benchmarking and profiling tools (competence assessment)
 - Knowledge evaluation tools
 - Certifications and training for self-made Data Scientists practitioners
 - Education and training marketplace: Courses catalog and repository

×

Outcome Based Educations and Training Model



From Competences and DSP Profiles to Learning Outcomes (LO) and to Knowledge Unites (KU) and Learning Units (LU)

 EDSF allow for customized educational courses and training modules design

Data Science Body of Knowledge (DS-BoK)

DS-BoK Knowledge Area Groups (KAG)

- KAG1-DSA: 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

Data Science Model Curriculum (MC-DS)

Data Science Model Curriculum includes

- Learning Outcomes (LO) definition based on CF-DS
 - LOs are defined for CF-DS competence groups and for all enumerated competences
- LOs mapping to Learning Units (LU)
 - LUs are based on CCS(2012) and universities best practices
 - Data Science university programmes and courses inventory (interactive) <u>http://edison-project.eu/university-programs-list</u>
- LU/course relevance: Mandatory Tier 1, Tier 2, Elective, Prerequisite
- Learning methods and learning models (in progress)

Example DS-BoK Knowledge Areas definition and mapping to existing BoKs and CCS (2012)

Knowledge Area Groups (KAG)	Knowledge (KA)	Areas	Suggested Kno	owledge Units (I	ки)	Mapping to CCS Science extension	2012 (includir ons) and existi	ng suggested Data ng BoKs				
KAG1 DSDA: Data Analytic	Theory of computatio	'n	Design and An	alysis of Algorith	nms CCS2012: Theory of computa Design and analysis			ion of algorithms tures design and				
group (including Machine	Knowledge Area Groups (KAG)	Knov (KA)	wledge Areas	Suggested Ki	nowledge Units (KU) Distributed Computer		Mapping to Science exte	CCS2012 (including suggested Data nsions) and existing BoKs				
statistical methods)	KAG2-DSENG Data Science Engineering	6: Com orga Big [nputer systems inisation for Data	Parallel and Architecture			CCS2012: Co Art	mputer systems organization chitectures Parallel architectures	-			
	group including Software an	Knowled Area Gro (KAG)	dge Knov oups (KA)	Knowledge Areas (KA) Data Management and Enterprise data infrastructure		sted Knowledge U	nits (KU)	Mapping to CCS2012 (including suggested Data Science extensions) and existing BoKs				
	infrastructu engineering		Data and I data			nanagement, inclu nce and Master D	ding ata	DM-BoK selected KAs (1) Data Governance, (2) Data Architecture,				
					Intelligence Data storage and operations Data archives/storage compliance and certification Metadata, linked data, provenance			 (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, and 				
Ι												
					Data ir and da	nfrastructure, data ta factories	registries	(11) Data Quality.				
 Mappir 	na suaaes	sted to	o CCS20	12	Data s	ecurity and protec	tion					
and ex	isting Bol	Ks			Data g Integra	overnance, data q ation and Interope	uality, data rability					

Example MC-DS Mapping Learning Units to DS-BoK and CCS (2012)

			KAG/	Learning	g Unit (course name) ²		ie) ²	Type/relevance ³ Map to f						DS-BoK, CCS2012 and known BoKs					
			LU# *)				Tier 1	Tier 2	Elective Pre requi			CCS20	12 based academic sub	based academic subjects DS-Bok			oK and other BoKs		
				Software design	e requ	ireme	ents an	d						Extens	ions are suggested fror	n SWEBOK	SWEBOK Softv	selected ware requ	KAS Jirements
KA	G/ L	/ Larning Unit (course name) ² Type/relevance ³					ance ³			Map t	o DS-I	BoK, C(S2012 and	d known BoKs				1 ruction	
LU; *)	#					Tier Tier Elective		Pre requisite		CCS2012 based academic subjects				ojects	ts DS-BoK and other BoKs		g enance		
\leftarrow	In	nformation Aathematic	theory	ie								M	athema	atical analy	/sis				configuration
	E	xtensibility	point for	radding n	iew														engineering
	A	artificial Inte	elligence						\langle		Tomp	uting	metho	dologies	_	No specific BoK are d	lefined		eering process eering models and
	N	latural Lang	guage Pro	Cessing	Learn	ing H	nit (co	urso name	1^2	Tyr	oo/rolo	/relevance ³ Man to DS Bok (CCS2012 and known Boks							
	Ki Re	nowledge F easoning	edge Represer LU# *)			urse name		Tie	r Tie	r Ele	ective	Pre	CCS2012 based academic subjects D			DS-BoK a	DS-BoK and other BoKs		
	D di)ata mining liscovery	and knc	knc Data type registries, PID,			es, PID,			2			requirte	Extended with the general Data Management Knowledge Areas and related academic subjects. Data Data			General Data Management KA'		
	Te	ext analysis	s, Data n		metadata Rosoarch data infractructur			e.									Data Life Data arc	ifecycle Management rchives/storage complian	
	Text analytics includ linguistic, and struct techniques to analys				Open Science, Open Data, (Access, ORCID				Open								certification New KAs to support RDA recommendations and commu data management models (Op		
	ar M al	nd unstruct Aachine Lea Igorithms	tured da arning th	Extensibility point for addi courses					g new	/								Access, Open Data, etc) Data type registries, PIDs Data infrastructure and Data F	
	C	lassificatior	n metho															TBD – To commun	o follow RDA and ERA hity developments
		Research methodology, res cycle				earch						Extended with the g Methods subjects ar	eneral Scientific/Resea nd related academic su	arch Ibjects.	Suggester related o	ed KAs to develop DSRN competences:			
	Modelling an planning				and ex	periment										Research (e.g. 4 st	n methodology, researd ep model Hypothesis – 2 Methods – Artefact –		
					- ·	1.1	•	1. 10.										nesearci	- Methous - Artelatt -

Mapping suggested to ACM CCS2012, DS-BoK and other related BoKs

EDSF Recognition, Endorsement and Implementation

- **DARE (Data Analytics Rising Employment)** project by APEC (Asia Pacific Economic Cooperation)
 - DARE project Advisory Council meeting 4-5 May 2017, Singapore
 - Followed by Ministerial meeting on 14-15 May 2017 in Hanoi, Vietnam
- PcW and BHEF Report "Investing in America's data science and analytics talent" April 2017
 - Quotes EDSF and Amsterdam School of Data Science
- Dutch Ministry of Education recommended EDSF as a basis for university curricula on Data Science
 - Workshop "Be Prepared for Big Data in the Cloud: Dutch Initiatives for personalized medicine and health research & toward a national action programme for data science training", Amsterdam 28 June 2016
- European Champion Universities network
 - 1st Conference (13-14 July, UK), 2nd Conference (14-15 March, Madrid, Spain)
 - 3rd Conference 19-20 June 2017, Warsaw
- e-IRG workshop on Sustainability on 8-9 June 2017, Malta

Further developments and Next steps (1)

- Next EDSF release 2 (planned for June 2017) will link competences to skills and knowledge
- Final EDSF project deliverables (due August 2017) will include:
 - Data Science Education Sustainability Roadmap
 - Will involve wide consultation with experts community and also with EU policy makers
 - Will be reviewed by the EDISON Liaisons Groups (ELG)
 - Certification Framework for at least two levels of Data Science competences proficiency
 - Consultation with few certification providers is in the progress
- Toward EDSF and Data Science profession standardisation
 - ESCO (European Skills, Competences and Occupations) taxonomy extending with the Data Science related occupations, competences and skills
 - CEN TC428 (European std body) Extending current eCFv3.0 and ICT profiles towards e-CF4 with Data Science related competences
 - Work with the IEEE and ACM curriculum workshop to define Data Science Curriculum and extend current CCS2012 (Classification Computer Science 2012)
- Number of Case studies is planned in cooperation with active EU projects EDSA, EOSCpilot, BDVe, etc. (not limited to the project lifetime)

Further developments and Next steps (2)

- The EDISON project legacy will include (linked to the current project website and migrated to CP in the future)
 - EDSF EDISON Data Science Framework
 - Data Science Community Portal (CP) http://datasciencepro.eu/
 - EDISON project network including
 - EDISON Liaison Groups
 - Data Science Champions conference
 - Cooperative networks with European Research Infrastructures (e.g. HEP, Bioinformatics, Environment and Biodiversity, Maritime, etc),
 - International cooperative links BHEF, APEC, IEEE, ACM
- Applications and tools development
 - Prototypes will be produced in the timeline of the project but further development is a subject to additional funding
- Sustainability of the project legacy/products will be ensured by the project partners voluntarily for the period at least 3 yrs
 - EDSF will be maintained by UvA
 - CP by Engineering (Italy)



Further developments and Next steps (3)

- Further dissemination, engagement and outreach activity
 - Publishing final deliverables as BCP and books
 - Data Science Manifesto Primarily focused on professional and ethical issues in Data Science, new type of professional
 - Inter-universities initiative "Data Science for UN's Sustainable Development Goals" to focus in-curricula research (projects) on UN priority goals



Summary: Services and References

- EDISON Website <u>http://edison-project.eu/</u>
- EDISON Data Science Framework (EDSF) <u>http://edison-project.eu/edison/edison-data-science-framework-edsf</u>
- Directory of University programs
 <u>http://edison-project.eu/university-programs-list</u>
- Community Portal <u>http://datasciencepro.eu/</u>
- Competences benchmarking and tailored training for practitioners
- Data Science Curriculum advice and design for universities
- Data Science team building and organizational roles profiling



DATASCIENCEPRO



Links to EDISON Resources

- EDISON project website <u>http://edison-project.eu/</u>
- EDISON Data Science Framework Release 1 (EDSF)
 <u>http://edison-project.eu/edison-data-science-framework-edsf</u>
 - Data Science Competence Framework
 <u>http://edison-project.eu/data-science-competence-framework-cf-ds</u>
 - Data Science Body of Knowledge
 <u>http://edison-project.eu/data-science-body-knowledge-ds-bok</u>
 - Data Science Model Curriculum
 <u>http://edison-project.eu/data-science-model-curriculum-mc-ds</u>
 - Data Science Professional Profiles
 <u>http://edison-project.eu/data-science-professional-profiles-definition-dsp</u>
- Survey Data Science Competences: Invitation to participate
 <u>https://www.surveymonkey.com/r/EDISON_project_-_Defining_Data_science_profession</u>



Other related links

- Amsterdam School of Data Science
 - <u>https://www.schoolofdatascience.amsterdam/</u>
 - <u>https://www.schoolofdatascience.amsterdam/education/</u>
- Research Data Alliance interest Group on Education and Training on Handling of Research Data (IG-ETHRD)
 - https://www.rd-alliance.org/groups/education-and-training-handling-research-data.html
- PwC and BHEF report "Investing in America's data science and analytics talent: The case for action" (April 2017)
 - <u>http://www.bhef.com/publications/investing-americas-data-science-and-analytics-talent</u>
- Burning Glass Technology, IBM, and BHEF report "The Quant Crunch: How the demand for Data Science Skills is disrupting the job Market" (April 2017)
 - <u>http://www.bhef.com/publications/quant-crunch-how-demand-data-science-skills-disrupting-job-market</u>
 - https://public.dhe.ibm.com/common/ssi/ecm/im/en/iml14576usen/IML14576USEN.PDF