OpenRiskNet

RISK ASSESSMENT E-INFRASTRUCTURE

Demonstration on data curation and creation of pre-reasoned datasets in the OpenRiskNet framework

Noffisat Oki (Edelweiss Connect GmbH, Switzerland), **Danyel Jennen** (Maastricht University, The Netherlands), **Marc Jacobs** (Fraunhofer Institute, Germany), **Tim Dudgeon** (Informatics Matters Ltd, UK)

Webinar - 18 March 2019

OpenRiskNet: Open e-Infrastructure to Support Data Sharing, Knowledge Integration and *in silico* Analysis and Modelling in Risk Assessment Project Number 731075



About the project

OpenRiskNet is a 3-year EU Horizon 2020 project with the main objective to develop an open e-infrastructure providing resources and services to a variety of communities requiring risk assessment, including chemicals, cosmetic ingredients, therapeutic agents and nanomaterials.



Main components:

- → Case-study-driven development examples of tools to be integrated are selected based on the case study needs. More information: <u>https://openrisknet.org/e-infrastructure/development/case-studies/</u>
- → Solutions for all areas by **integrating existing tools** from consortium and associated partners (via the implementation challenge)
- → Integrated approach combining experimental data (*in vivo, in vitro, in chemico*) with analysis, modelling and simulation tools into risk assessment workflows



Webinars series

Live demonstrations on the e-infrastructure deployment and the risk assessment case studies

	Торіс	Date & Time
Past events	Introduction sessions to the OpenRiskNet e-infrastructure	See Webinar recordings: • Session 1 (24 Sep 2018) • Session 2 (27 Sept 2018) • Session 3 (4 Oct 2018) • Session 4 (30 Oct 2018)
	Learn how to deploy the OpenRiskNet virtual research environment	See <u>Webinar recordings</u> (25 Feb 2019)
Current event	Demonstration on data curation and creation of pre-reasoned datasets in the OpenRiskNet framework	Monday, 18 March 2019 16:00 CET
	Identification and linking of data related to AOPWiki (an OpenRiskNet case study)	Tuesday, 26 March 2019 17:00 CET
Future events	Semantic annotation	Monday, 1 April 2019 16:00 CET (to be confirmed)
	The Adverse Outcome Pathway Database (AOP-DB)	Monday, 8 April 2019 16:00 CET
	Additional demo on case studies	May - June 2019 (to be announced)



https://openrisknet.org/events/



DataCure Webinar Content

- 1.) What is DataCure and how is data curation on OpenRiskNet implemented?
- 2.) Demonstrations on various aspects of DataCure showing practical examples from the DataCure case study
 - a.) Transcriptomics data extraction and metadata annotation
 - b.) Text Mining workflow for metadata extraction for carcinogenicity prediction
 - c.) Data extraction and curations for liver toxicity modeling
- 3.) Q&A/Discussion



DataCure case study

• DataCure establishes a process for data curation and annotation that makes use of APIs (eliminating the need for manual file sharing) and semantic annotations for a more systematic and reproducible data curation workflow.

• The aim is to deliver curated and annotated datasets for OpenRiskNet service users as well as preparation and development of tools and workflows that provide examples of useful toxicogenomic data analysis methods will also allow users perform their own data curation and analysis.



Implementation

- 1.) **Data sources** including Physchem, toxicological and omics databases (e.g. EdelweissData Explorer, PubChem, Liver Toxicology Knowledge Base (LTKB), etc), ontology/terminology/annotation tools (e.g. SCAIView), and literature databases (e.g. Pubmed and Toxplanet).
- 2.) **Data Extraction** from the EdelweissData Explorer and all other resources using API calls and text mining workflows.
- 3.) **Data Searching** using workflows that employ text mining capabilities for searching, refinement and curation of the extracted data.
- 4.) **Data curation and reasoning** using workflows provided in Jupyter or Squonk notebooks and which cover exercises such as extraction of specific data, merging of datasets, and annotation of data for downstream analysis purposes.
- 5.) Resubmission to data source



Use case: finding toxicological information

Scientist: Is this chemical carcinogenic? Or maybe hepatotoxic? Or ...? Is there an easy way to retrieve toxicological information?



Click-click-read-click-read-click-read → answer?

OpenRiskNet RISK ASSESSMENT E-INFRASTRUCTURE



www.openrisknet.org

Manual annotation >300 chemicals, >200 references

				• • •										Y				Y		•===				••••	•••••			E -	•	-
		•		<i>M</i>			÷.			 •							• • • • • •						782		1926 -		er tan i			
	iii		8						. *		9			ŀ	240°	,														Û
1 mm				ar e	-m-									ŀ	veor					<u>a. 19</u>				W.				handan da	£	
	····	•							200			,		1			4:0	w.				:		l						
	··· ···	•		201 66								127 14	•	:	6 ,			6					4.55. West	ter etile						
1 2 m 🟥 🚥 🛊 🐄 🚥 🕮 🐜 🛲 💷 🖷 🕮	···	•		10. 191	10.000	-1:- 1- 1 1			· ••••	 							3 - 3 -		2.,	••••••			.uq.n. 1722	k:	1999 -	-				
	III	•				Ŀ.				 •			A			37	345	a l							1	-		Annedatory		
	II	•	10			S			à			, Č				X)		•• has of			¹	11. 11.						
		•				7						X				3	<u>.</u>				*;;;#			·		-		9		1
	×	• • •	10.	100 W		\$ 'ar	4. die		9 . tait		X	7			11		5	-	¢ •		·			e- 425.	1 II		w 10-1		ž	
	iii				S		- 490 m	- 3	a		<u>s</u>	1998 -	(w)	•		E D	ð 👘	w		, traggerer,	Y			e see			u da			
		•		ð,	0	•				 . 0	0					2		2		5						-				1
			10	0	•••				· 2.	 0					Ċ													"		
	#												-		6	- <u>t</u> :	3-3-					_	nd.			*				
		• •				··· · <										15.						2005	j::##:	tii,	THE C	5				
			Y										2008 2008						1		i i		ĺ							ļ.
	. ·	•			1											<u></u>	•					_	ndine ann			- I ·		trtoggadae	5.	
11 mu		•		1.jn ii		w	·w		· 2					ł				2					ri <i>ș</i> te _e teș			-			Ë.	
	III	· .			•	•												8							<u>.</u>	-				
					æ -										11						'		tatu							
				dat 11	1.11 mar						i si			八				人							·; •					

The answer

OpenRiskNet



www.openrisknet.org

Data Curation Workflow

Presented by Noffisat Oki (Edelweiss Connect Inc, Durham, NC, USA)



Data curation workflows

Data is messy.

It nearly always needs to be cleaned up and manipulated before being used.

We will illustrate two common needs being performed on the ORN infrastructure.





Text Mining Workflow Presented by Marc Jacobs (Fraunhofer Institute, Germany)



Text Mining Workflow



https://monographs.iarc.fr/agents-classified-by-the-iarc/

1) https://openrisknet.org/e-infrastructure/development/case-studies/case-study-datacure/



Text Mining workflow

Task:

- Identify the concept of acetaminophen (definition, identifiers, synonyms)
- Find all relevant documents in the context of acetaminophen and carcinogenity
- What are the most relevant statements?

Technology:

- Semantic index of PubMed/PMC (> 20 terminologies)
- Solr index + OLS index + UIMA pipeline

OpenRiskNet





Data extraction and curations for liver toxicity modeling

Presented by Tim Dudgeon (Informatics Matters, UK)

OpenRiskNet RISK ASSESSMENT E-INFRASTRUCTURE



Data curation workflows

Data is messy.

It nearly always needs to be cleaned up and manipulated before being used.

We will illustrate two common needs being performed on the ORN infrastructure.





CPSign - LTKB modeling using Conformal Prediction and SVM Example: nortriptyline

Model building

SVM with RBF kernel; exhaustive parameter sweep for C and γ Evaluation using 10-fold cross validation Finding the most **efficient** model Final model packaged as JAR archive

Web service

OpenShift templating used for deploying the JAR as a web service Automatically generates an OpenAPI definition for each model e.g. <u>http://ltkb-no-vs-most-cpsign.prod.openrisknet.org/</u>





Jaqpot – Example modelling the LTKB data using a Multi-layer Perceptron classifier

Import JaqpotPy package

import pandas as pd from jaqpotpy import Jaqpot from sklearn import linear_model from sklearn import neural_network

Train NN using your preferred algorithm

```
[ ] nn=neural_network.MLPClassifier()
```

```
[ ] neural_model=nn.fit(X,Y)
```

```
[ ] nn.predict(X)
```

Turn the model into web service

- **any** *scikit-learn* algorithm under Python
- 1 command gets model uploaded
- Instant
- Model becomes a web service

j jaqpot.deploy_neural_network(neural_model,X,Y,"Neural network)



Jaqpot - predict using web service

= Jaqpot 🖄					Q	▲ ● ⊖
	Overview	Data	Predict / Validate	Discussion		
MODEL Title: Neural network model predicting DILI Owner: filipposd	Choose method Predict Upload dataset with t	he required indep	User pro for p using provide	ovides dataset oredictions g template ed by Jaqpot		
Predicting DILI	Input values for the in SeverityClass	dependent feature	ALo	gp2	fragC	



www.openrisknet.org





Acknowledgements

OpenRiskNet (Grant Agreement 731075) is a project funded by the European Commission within Horizon 2020 Programme

Project partners:



- P1 Douglas Connect GmbH, Switzerland (DC)
- P2 Johannes Gutenberg-Universität Mainz, Germany (JGU)
- P3 Fundacio Centre De Regulacio Genomica, Spain (CRG)
- P4 Universiteit Maastricht, Netherlands (UM)
- P5 The University Of Birmingham, United Kingdom (UoB)
- P6 National Technical University Of Athens, Greece (NTUA)
- P7 Fraunhofer Gesellschaft Zur Foerderung Der Angewandten Forschung E.V., Germany (Fraunhofer)
- P8 Uppsala Universitet, Sweden (UU)
- P9 Medizinische Universität Innsbruck, Austria (MUI)
- P10 Informatics Matters Limited, United Kingdom (IM)
- P11 Institut National De L'environnement Et Des Risques INERIS, France (INERIS)
- P12 Vrije Universiteit Amsterdam, Netherlands (VU)



Webinars series

Live demonstrations on the e-infrastructure deployment and the risk assessment case studies

	Торіс	Date & Time							
Past events	Introduction sessions to the OpenRiskNet e-infrastructure	See Webinar recordings: • Session 1 (24 Sep 2018) • Session 2 (27 Sept 2018) • Session 3 (4 Oct 2018) • Session 4 (30 Oct 2018)							
	Learn how to deploy the OpenRiskNet virtual research environment	See <u>Webinar recordings</u> (25 Feb 2019)							
Current event	Demonstration on data curation and creation of pre-reasoned datasets in the OpenRiskNet framework	Monday, 18 March 2019 16:00 CET							
	Identification and linking of data related to AOPWiki (an OpenRiskNet case study)	Tuesday, 26 March 2019 17:00 CET							
Future events	Semantic annotation	Monday, 1 April 2019 16:00 CET (to be confirmed)							
	The Adverse Outcome Pathway Database (AOP-DB)	Monday, 8 April 2019 16:00 CET							
	Additional demo on case studies	May - June 2019 (to be announced)							



https://openrisknet.org/events/

