

# ENACTING DATA SCIENCE PIPELINES FOR EXPLORING GRAPHS

*From Libraries to Studios*

**Genoveva.Vargas-Solar**

*LIG-LAFMIA, CNRS, France*

[genoveva.vargas@imag.fr](mailto:genoveva.vargas@imag.fr)

*José-Luis Zechinelli-Martini, UDLAP, Mexico*

*Javier Espinosa, LAFMIA-U. Lyon 3, France*

BBIGAP Workshop, Lyon, 25<sup>th</sup> August, 2020

# DATIFICATION

2

Rendering into **data** aspects of the world that **have never been quantified**



Any individual can analyse huge amounts of data in short periods of time

- **Analytical knowledge:** most of the crucial algorithms are accessible
- Use rich data to make **evidence-based decisions** open to virtually **any person or company**

# CONNECTED DATA

3

**1** Observations structured as networks with their own interconnection rules determined by the variables characterising each observation

## Connected Enterprise

*Employees, Costumers  
& Partners*



## Internet Connected Things



## Graph



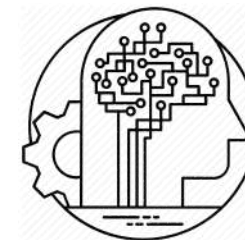
**2** *Mathematical Concept  
& operations  
Data structures  
& algorithms*

## Digital Mesh

*Hubs, Nodes  
& Links*



## Knowledge Network



# CONNECTED DATA

4

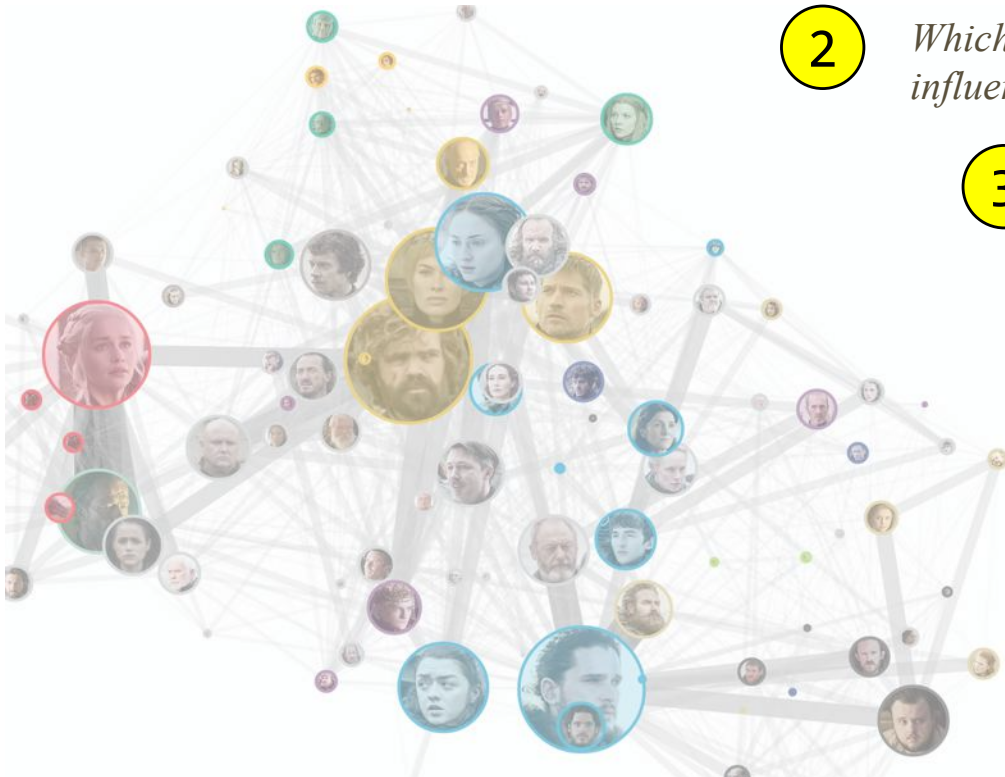
1 *How many characters are there in GOT?*

2 *Which are the houses that challenge the throne and how influential are they?*

3 *How are the most popular characters?  
How does popularity evolve of main characters  
evolve along the story?*

4 *Determine a geographical organisation  
of the different countries given the trajectories  
performed by armies?*

5 *Predict who can be the next King/Queen?*



# QUERYING APPROACHES

5

|                    | Databases ; information retrieval  |  | Data Science   |
|--------------------|--|--|--|
| Query Types        | <ul style="list-style-type: none"> <li>- Relational, multi-dimensional, spatio-temporal, aggregation</li> <li>- Patterns, regular expressions</li> </ul> | <ul style="list-style-type: none"> <li>- (Dis)conjunctive</li> <li>- Navigational</li> </ul>   | <ul style="list-style-type: none"> <li>- Exploratory</li> <li>- Analytics: modelling &amp; predicting</li> </ul> |
| Execution model    | On demand/continuous   | On demand  | Step by step   |
| Results properties | <ul style="list-style-type: none"> <li>- Completeness (full/partial)</li> <li>- Fussiness</li> </ul>   | Approximation<br>Precision/recall   Probabilistic  | Approximation<br>with some error degree<br>Data, queries, samples  |
| Dataset content    | Intention model  | Extension model  | Extension (raw)  |
|                    | <i>Data structure</i><br>table, key-value, tuple, document, graph  | <i>Quantitative representation</i><br>Frequency matrix, Statistical profiling<br><br><i>Semantic representation</i><br>Ontology, Terms graph | <i>csv, XML, JSON, BLOB, ...</i>   |

# EXPLORATORY QUERYING

6

|                    |   |
|--------------------|---|
|                    | Data Science  |
| Query Types        | - Exploratory<br>- Analytics: modelling & predicting  |
| Execution model    | Step by step  |
| Results properties | Approximation<br>with some error degree<br>Data, queries, samples   |
| Dataset content    | Extension (raw)<br><i>csv, XML, JSON, BLOB, ...</i><br><i>Data structure</i><br><i>table, key-value, tuple, document, graph</i> |

- Methodologies weaving data management, greedy algorithms
- Programming models that must be tuned to be deployed in different target architectures

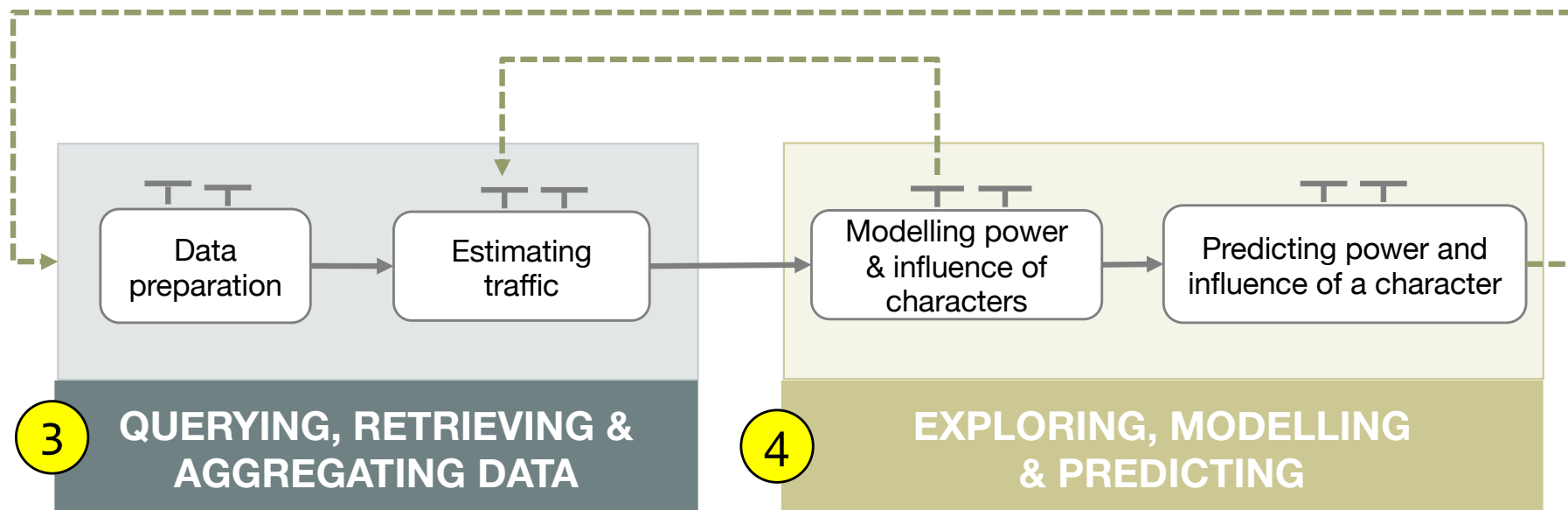
Data collections as backbone for conducting **experiments**, drive hypothesis and lead to “valid” conclusions, models, simulations, understanding

# DATA SCIENCE PIPELINE

7

1 COMPLEX AND REPETITIVE PROCESSING & ANALYSIS TASKS

2 **Example:** *Predict who will be the last king/queen of GOT?*




# DATA SCIENCE PIPELINE

8


## ISSUES

- 1 **Artisanal design** depending on data scientist/engineers “*expertise*”
- 2 In-house programming using many different **libraries, stacks, tools** **difficult to integrate**

### DATA SCIENCE LABS

 [kaggle.com](https://www.kaggle.com)

 [Google Colab](https://colab.research.google.com/)

 [Azure Notebooks](https://azure.microsoft.com/en-us/services/notebooks/)

### DATA SCIENCE STACKS



### BD SERVICES PLATFORMS





# GRAPH STORES AND SYSTEMS

9

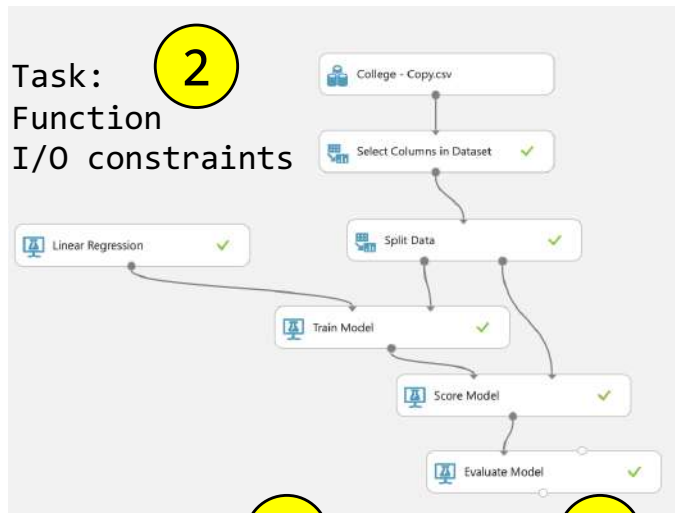
| Graph Operations  | Stores  | Analytics tools  |
|---|---|--|
| <ul style="list-style-type: none"><li>① - <b>Community detection:</b> <i>label propagation, Louvain, weakly/strongly connected components, triangle count, colouring, local clustering coefficient</i></li><li>② - <b>Centrality:</b> page rank, betweenness, closeness, degree, Eigenvector</li><li>③ - <b>Similarity:</b> node, nearest neighbours, cosine/Euclidean/Jaccard/ Overlap/Pearson</li><li>④ - <b>Pathfinding &amp; search:</b> adamic Adar, common/total neighbours, same community</li></ul> | <ul style="list-style-type: none"><li>① - GraphDB lite</li><li>- Neo4J</li><li>- OrientDB</li><li>- GraphEngine</li><li>- HyperGraphDB</li><li>- MapGraph</li><li>- ArangoDB</li><li>- Titan</li><li>- BrightStarDB</li><li>- CayLayGraph</li><li>- WhiteDB</li><li>- Orly</li><li>- CosmosDB</li></ul> | <ul style="list-style-type: none"><li>① <b>Python networkx:</b><br/>combined with numpy &amp; pandas<br/>shortest\_path(),<br/>shortest\_path\_length(),<br/>all\_pairs\_dijkstra\_path()</li><li>② <b>Spark Graphix</b></li><li>③ <b>Deep Learning Models:</b><br/>node classification, link prediction and graph clustering</li><li>④ <b>Graph processing systems (Pregel, Giraph):</b><br/>global pattern matching, shortest paths, max-flow or min-cut, minimum spanning trees, diameter, eccentricity, connected components, PageRank, traversals</li></ul> |
| <ul style="list-style-type: none"><li>⑤ - <b>Heuristic link prediction</b></li><li>- <b>Pattern discovery</b></li></ul>   | <ul style="list-style-type: none"><li>② - <b>Proprietary data structures</b> (RDF, matrices/tables, etc)</li><li>- <b>Query languages with built in functions</b></li></ul>   |  |

# MACHINE LEARNING STUDIOS

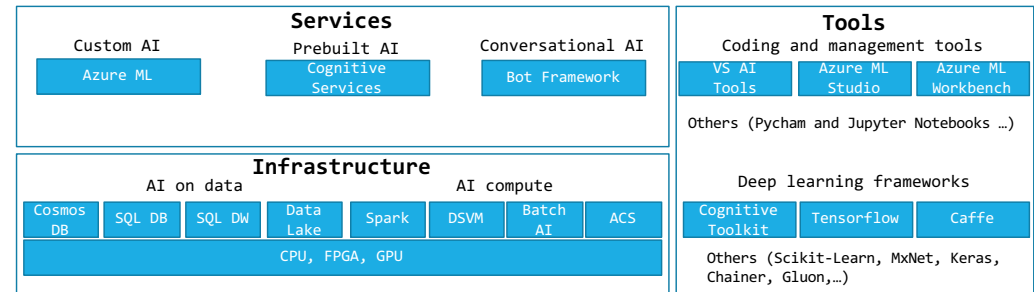
10

## Pipeline 1

Task:  
Function  
I/O constraints



## DS Tools



## Tracking 3

Record & query experiments:  
code, configs, results .. Etc

## Projects 4

Packaging format for reproducible runs on any platform

## Models 5

General model format that supports diverse deployment tools

Enactment (AI Platform, e.g. Microsoft, Databricks MLFlow, Google AI)

# CLOUD ML STUDIOS

|   | Amazon   | Microsoft   | Google                                   | IBM  |
|---|--|---|--|--|
| <b>Automated and semi-automated ML services</b> |  |   |  |  |
|   | AmazonML   | MS AzureML Studio                                       | Cloud AutoML                             | IBM Watson Model Builder   |
| Classification                                  | ✓  | ✓   | ✓  | ✓  |
| Regression                                      | ✓  | ✓   | ✓  | ✓  |
| Clustering                                      | ✓  | ✓   | ✗  | ✗  |
| Anomaly detection                               | ✗  | ✓   | ✗  | ✗  |
| Recommendation                                  | ✗  | ✓   | ✓  | ✗  |
| Ranking   | ✗  | ✓   | ✗  | ✗  |
| <b>Platforms for custom modelling</b>           |  |   |  |  |
|   | Amazon SageMaker   | Azure ML Services                                       | Google ML Engine                         | IBM WatsonM Studio   |
| Built-in algorithms                             | ✓  | ✗   | ✓  | ✓  |
| Supported Frameworks                            | Tensorflow, MXNet, Keras, Gluon, Pytorch, Caffe2, Chainer, Torch | Tensorflow, Scikit-Learn, MS Cognitive Toolkit, SparkML | Tensorflow, Scikit-Learn, XGBoost, Keras | Tensorfoow, SparkMLlib, Scikit-Learn, XGBoost, PyTorch, IBM SPSS, PMML |

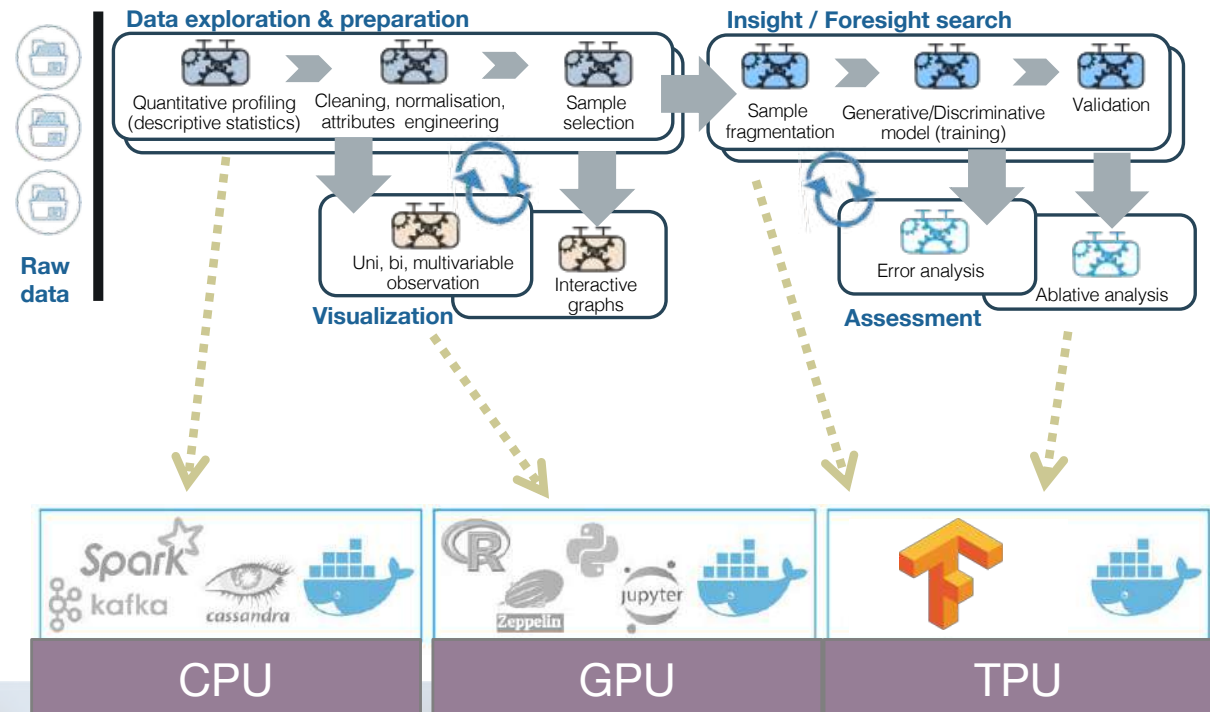
# CONCLUSION & OUTLOOK

12

Next data science querying environments providing services at scale addressing data and processing divide are yet to come

1 Data science queries as **services coordinations** for processing graphs

2 **Deployment** on target architectures (cloud, GPU) for large-scale processing





Genoveva **VARGAS-SOLAR.COM**

Senior Scientist, CNRS, France



**José-Luis Zechinelli - Martini**

Universidad de las Américas Puebla, Mexico

*joseluis.zechinelli@udlap.mx*



**Javier Espinosa**

University Jean Moulin Lyon 3, France

*javier.espinosa@imag.fr*