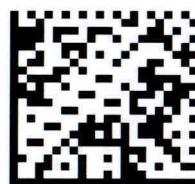


NOM BOUKHELOUA

Prénom Rémi

Promo L2S2

Date 16.04.



20150257: BOUKHELOUA Rémi

M1:

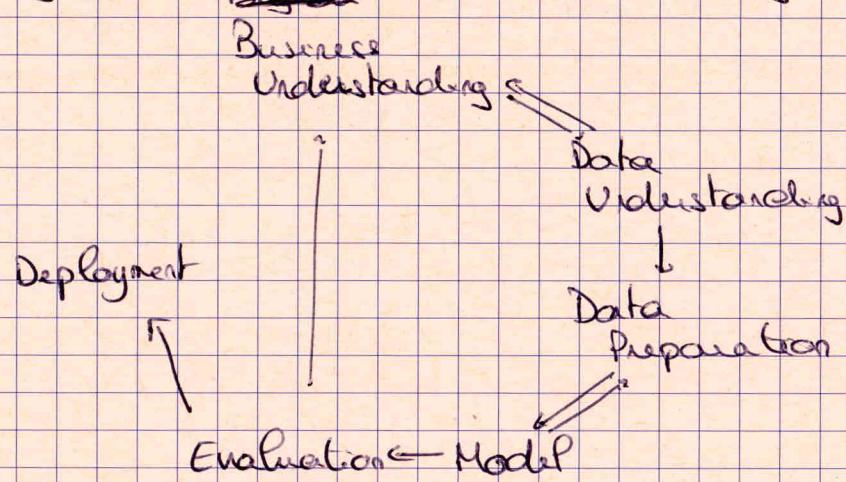
ST2DVA-DE (16/04/2019)

Amphi jaune

16

MATIÈRE Data Visualization (1/2).

In this we ~~we~~ followed CRISP DM method Cross Industrial Standard Process for Data Mining (CRISP DM) used to do ~~Big Data~~ ~~big data~~ ~~experimental~~ data mining. It follows this graph:



- ~~Business Understanding~~ is knowing what we need and why we need it.
- Data Understanding is knowing what data we acquired and what are the trends.
- Data Preparation is creating new features or arranging data according to the model we want to do.

- Model is creation of data mining model.
- Evaluation is testing our model to see if it fits business understanding.
- And if it fits, we deploy it to make it ~~useable~~ by everyone

Otherwise we have the SEMMA (Sample ~~Explore~~ ~~Model~~ ~~Modify~~ Model Assess) model which ~~consists~~ consists in:

- Sample: select usefull data regarding our needs.
- Explore: ~~explore~~ looking for trends, visualize data, basics analytics.
- Modify: transform the data to make it ~~useable~~ ~~for~~ our model
- Model: create model regarding our needs.
- Assess: Evaluate it.

Q3. Comparisons between Tableau, PowerBI, QlikView:

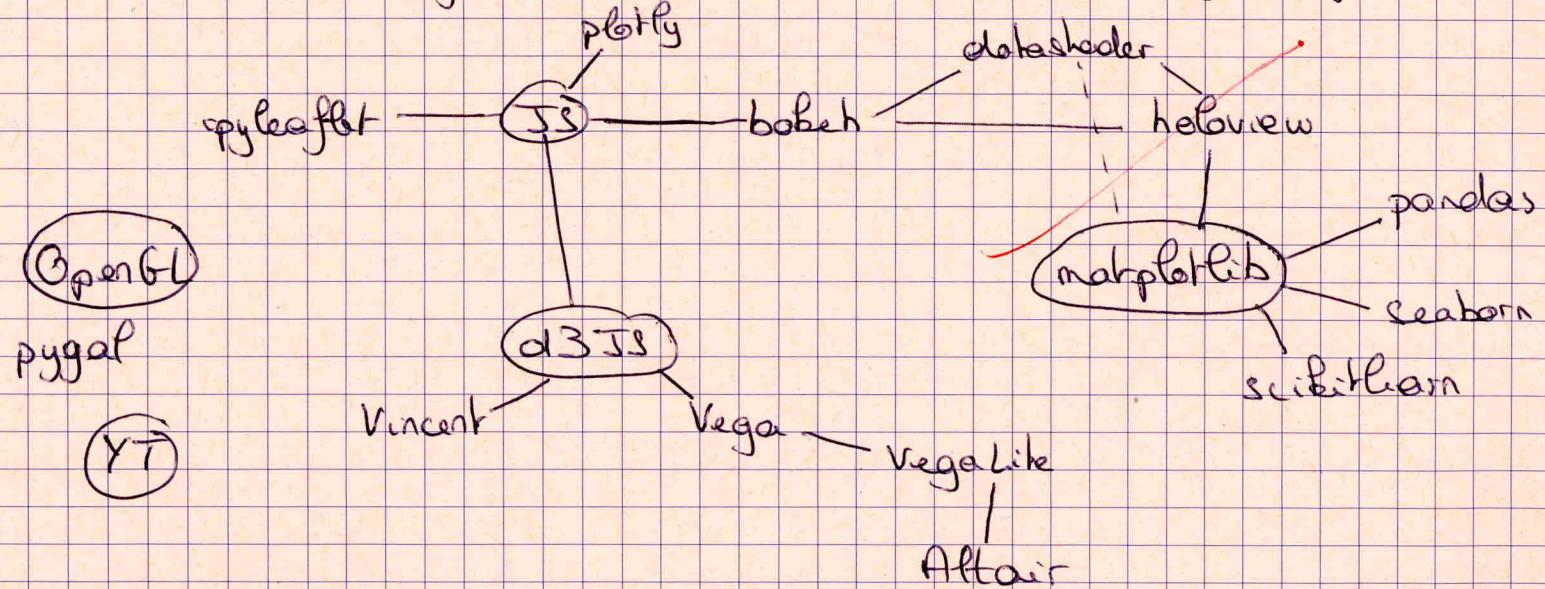
	Tableau	PowerBI	QlikView
Reputation	Mostly Used Tool.	2 nd most popular	
Usage	Good Visualization tool	Good Data manipulation tool.	Good to do in both.
Price	Expensive over \$8/month/user	Less expensive (around \$10/month)	Mid Range Price
Sources	Huge number of sources.	Most restricted number of sources.	
Limitations	Server's limit	10GB	Server's limit.
Speed	Efficient (RAM used)	Slower.	Efficient (RAM) used.
Others		<ul style="list-style-type: none"> • Produced By Microsoft • Lots of updates. • Good compatibility with Excel. 	

Q5. It exists 4 main point of views bias that ~~we have to~~ we have to ~~check avoid~~ if we want to present the data rightly.

- Conceptual POV: if you find a correlation between two features, it doesn't mean they are correlated
- Perceptual POV: even though the dataset is complex, we have to make sure our representations are clear, not ~~saturated~~ overloaded and straight to what we want to say.

- **Graphical Pov:** we have to make sure our representations fits the reality. To do so, we can use the Lie Factor $\frac{\text{size of object}}{\text{size of person}}$ to estimate proportions. Make sure of the size of objects or persons such as the doctor that get smaller and smaller (in %ages)
 - **Political Pov:** we have to be careful of any political orientation. For example, for one election, we could represent the winner ~~in~~ in every city over a map. But not every city has the same number of electors or senators. So it would be more accurate to represent ~~each~~ the winner per city when each city is sized regarding number of electors.

Q7. Some of most famous Python Landscape libraries ~~for~~ for Data Visualization



~~Q2~~ We used to follow we follow this cycle:

Business Process Definition: where we state our needs.

Research: where we observe what consumers do and what decide what we want to do.

Human Resource Assessment: can we ~~suggest~~ do this project?
or do we need to outsource it?

Data Acquisition: recuperating data. Thanks to a crawler for example.

Data Munging: transforming unstructured data to structured one.

Data Storage: storing almost raw data

Exploratory Data Analysis: exploring data, looking for trends.

Data preparation: preparation for modeling, selecting features,
creating new ones.

Model: creating the Model.

Implement: deploy the model.

Q4. Nowadays, Big Data can be described through 3V: Volume, Velocity
and ~~Versus~~ Variety. Data Visualization faces those three & as

challengers. We need to process more and new data from with different format but also need to represent them as we used to do. Having more and more ~~parameters~~ but the same number of axes. So we need ~~many~~ libraries that can support a lot of data and a lot of formats.

At the same time, we ~~need~~ want to have a view of how our product is ~~developing~~ used. So we need real time view such as the Singaporean heatmap representation. Tools like Jupyter Notebook allows us to do so but still not very comfortable to use for majority.

NOM BOUKHELOU A

Prénom Rimi

Promo 2010

Date 16.04.

MATIÈRE Data Visualization (2/2).

Q6 Big Data Visualization contains various styles such as
Data reduction where you try to aggregate many data
in a simple manner and a simple view. For example, ~~making~~
you could make a network between ~~main~~ characters of
Bible and mostly used words.

Moreover, there is one where you will try to not only
~~only~~ show facts but also explain ~~them~~ using heatmaps
and PCA. Be carefull of Conceptual POV!

