

DIGIT.B4 – Big Data PoC

DIGIT 01 - Social Media

D02.02 Technological Architecture

everis Spain S.L.U



Table of contents

1	Intro	oduction	5
2	Meth	nodological Approach	6
	2.1	Business understanding	7
	2.2	Data linguistic understanding	7
	2.3	Data preparation	8
	2.4	Modelling	8
	2.5	Evaluation 1	0
	2.6	Deployment1	0
3	Tech	nnical Architecture1	1
	3.1	Information gathering1	2
	3.2	Information Store 1	2
	3.3	Data Analytics 1	2
	3.4	Governance 1	2
	3.5	Decision Support/Utilization1	3



Table of tables

Table 1- Information gathering	12
Table 2 - Information Store	12
Table 4 - Data Analytics	12
Table 5 - Governance	13
Table 6 - Decision Support/Utilization	13



Table of figures

Figure 1 - Data analysis main flow	6
Figure 2 - Methodology phases flow	6
Figure 3 - Documents to corpus transformation	7
Figure 4 - Completeness and business standpoint analysis	7
Figure 5 - Data preparation process	8
Figure 7 - Clustering process	9
Figure 8 - Classification process	9
Figure 9 - Wordcloud	9
Figure 10 - Occurrence matrices	10
Figure 11 - Semantic network	
Figure 12 - System Architecture	11



1 INTRODUCTION

The construction of the proof of concept is based on the definition of a technical architecture. The best way to define it is first to ensure that the PoC is built by following a solid methodological approach, and then to define a technical architecture supporting it.

The elements above are analysed in this document, which includes the following sections:

- Methodological Approach.
- Definition of the Technical Architecture.



2 METHODOLOGICAL APPROACH

From a methodological point of view, the main tasks that will be carried out are:

- Text-mining treatment: transform the content of the posts into a format that serves as an input for modelling algorithms
- Clustering: make homogenous groups of similar posts
- Classification model: obtain the rules to classify a document in one of the categories defined
- Draw conclusions from the classification of documents



Figure 1 - Data analysis main flow

The activities are based on the CRISP-DM (Cross Industry Standard Process for Data Mining) methodology. This methodology is the most commonly used standard process model in both academic and industrial fields, and provides a solid framework for data mining projects.

The methodology displays a process in six phases which accounts for all the activities required to gather, classify, store and analyse the data. The CRISP-DM methodology has been chosen due to previous success cases in everis using it.







2.1 Business understanding

The objective of this phase is to understand the project objectives and requirements from a business perspective, and then convert this knowledge into a text mining problem definition.

The business orientation of the problem is to obtain a better understanding of the main topics that are being treated to support the orientation of research funds.

The available sources will be analyzed to establish the best suited solution to the problem.

2.2 Data linguistic understanding

Once de data is loaded, a corpus (large collections of texts representing a sample of a variety presented a machine -readable form) is created to make analysis in order to understand the information contained in the posts received and ensure the data quality and the project viability.

-	Stephen Quest (Istephen_gon Plicting reverse mentoring by a @mjmorrow for the inspiration	ist - Jun 5 millenniais @EU_Commission. Thanks to # wDkGiTconf	0				
	Best digital innovation	n award for #digitconf, w DenisProst1 @webmutation		an Date	Teast .	7	Catana Catana C
	Bande inter werten	aconina (Beo_Commission		1222-1614-2023 Dave Streamer	Beit Eighte innerstan anweit ferstightent, in berefriedt Øvechruteten Øfte sefteten Anterseisenne.	Anges (Perdanter) or 10 martine to a Martin Robert France	14) 7 T
		Miteedance BDDC.ed		1221-18jun. 208 den Simeore	Best Signal innevation available to give the Densitive Signature Site of Site as Technikolamic contra- Site Announces	https://twitter.com/dein_comesnove/stanue/%108894871440	
	83 aft	Their Selections - House to 11, 2014 & 12 are		1231- Mpin. 307 Den Smann	Beit figter innereter anwerter registert, ur Derefrect, Dieberufeten Die auf bein Ansere annen	The Protectory of the second state and the second state of the sec	B 2 2
	and and a local state	A Letting by the end of the execution of a case and the only of the end of		1221-1834-2838 Den Smeane	Best dig the inner all on avaid for Heightend, wither all north give brouds on give as Techn Anterna comme	The first of the second second second second	
	-	with fairs (a) facilitation (Not) parameters and theme from whet see form of the server these on \$200 (DET a code to on 10 handlings for 2 converses) and they were teaching the a most exercise to find a code of the server form of the server.		1121-1614-1839 Dam Stream	beit Eigheil mersen bei anwerten Weigheit in Dereitfreit Beisberumen BTeise Tette Keterseigenes 600. Sommeren	The Original States and the supervision of the substates of the supervision of the superv	
		is the Mindes Sens Post and Septer Gent		1221-14jun. 200 Dem Simeone	Best digital innovation avant for Hdigitanit, w Denahoord, @vebmu/ation @Twee Tixtin Hinterial commo DEL. Sommassin	Max (Notes and the games to a Maria Statements)	
	6 13	Dem from		55.21- More. 2021 Dent Streame	Build gate in Novation associate Registerit, in Denahment Breathnyation Bhase Testin Rindersa comme Biol Commission	Man Sherber any is in a manifest that the set	
		Fair approach to Lothney faithand here a sine internation		12 21- 14 jun. 202 Dark Dimester	But ogste innersten avaraforstagsteint, widen afresti Bykabnutstein Bfilvaaffikten entarse eining.	THE OWNER CONTRACTOR OF THE NEW WORKS, WITHOUT AND	
		an alternity provide program thank on agent way were take worked provide a sector taken to be a sector to be a sector of the activity of the sector than a field from that the beaming the afficiency of the sector (coldy) we all sector to fine		1221-1814. 323 Devi Smerre	Seat Spiter innorman available Register II. v Denahmett Breasmusteren Bilv auf tenn Brearra genne	The Other and the provider of March 10000000000	
		sarker flows: - () - (*), 54% churd motion a sen is fanite spinduote with a flowing a server. (7)(in the the field that, working in the flow server mode.		12.21-14/vn.354 Den Smeane	Best digter monation available register if, in benefred i gradmustion gift as from entered comes grad dimension	https://twitter.com/dem_primeoneva/itatus/\$10009925440	14 I I
		-score of the Antonians are Astractly stored. This reasons Sizes for Concernation of the Astractic Store of the Astra-					

Figure 3 - Documents to corpus transformation

The main checks to make are:

- Completeness: the volume of documents and different words must be sufficient to allow reliable analysis.
- Quality: the words to be used for the analysis must be grammatically correct.
- Business point of view: the most frequent terms are analysed to validate quality from a business point of view.



Figure 4 - Completeness and business standpoint analysis



2.3 Data preparation

In this phase, the necessary tasks are performed to transform the original text in a final dataset (set of keywords) that can be use by the algorithms in the modelling phase.

		Corpus	Algorithm input dataset
Author	Date	Publication	Publication
Deni Simeonova	16/06/2015	Best digital innovation award for #digitconf, w DenisProst1 @webmutation @TweeTikitin #internalcomms @EU_Commission	 digital innovatiton award
Stephen Quest	05/06/2015	Piloting reverse mentoring by millennials @EU_Commission. Thanks to @mjmorrow for the inspiration! #DIGITconf	Piloting reverse mentoring millenials inspiration
Dace Kalnina	19/11/2014	Thanks for having so cool 3d printer stand at #DIGITconf everybody loved it!	3d printer

Figure 5 - Data preparation process

The steps to transform the corpus are:

- Delete stop words: words that don't give any information and have no effect on how the rest of the words are related with each other.
- Correction of spelling errors.
- Stemming: the word stem is kept to unify conjugations.
- Create dictionaries: used to unify different terms that are expressing the same concepts
- Synonym and hyperonymic: used to unify terms that often appear together (these relations should be consistent from the perspective of business).

2.4 Modelling

During this phase, the statistics models are built to achieve the main objectives of the project and to assign a cluster and a classification to each document.

The first step is building the term document matrix (TDM) that will be the input for the LSA algorithm (will be explained later). This matrix contains the frequency of each term in every document.

Using this TDM, the score of relevance will determine which words are significant for clustering documents.

With these inputs, latent semantic analysis (LSA) is used to analyze links between a set of documents and the terms they contain by producing a set of concepts related to the documents.

Once the main clusters are defined, they are associated to a concept classification (defined in the 2014 topics document).



Figure 6 - Clustering process

This clustering is used as an input for a classification models. The result of executing a classification model is a set of rules that let us classify a new document based on the terms that appear in it.



Figure 7 - Classification process

Before the execution of LSA and the classification model, other modelling tasks and analysis are made to know how to approach the main models:

Wordcloud: to understand most frequent terms.



Figure 8 - Wordcloud

 Term relation: to identify related terms to define broader concepts that increase the effectiveness of classification and clustering algorithms



Security		
Term	% of occurrence	
Authentication	35%	
Confidentiality	27%	
Integrity	13%	
Protect	7%	

Cloud and Hosting		
Term % of occurrence		
Hosting	50%	
Cloud	34%	
Migrating	10%	
Public Cloud	5%	

 Mobile

 Term
 % of occurrence

 Mobility
 73%

 Mobile device
 67%

 Mobile computing
 50%

 IoT
 43%

 ...
 ...

Figure 9 - Occurrence matrices

Semantic network: to understand relationships between all terms in the different posts.



Figure 10 - Semantic network

2.5 Evaluation

In this phase some analysis are performed to ensure the quality of the generated models. It runs in parallel with the modeling phase.

All analysis and results of the execution of algorithms are tested to ensure statistical validity and check that they make sense from a business standpoint.

2.6 Deployment

Transform the results obtained after modelling in the analysis and conclusions required for the project.

Support results for presentation and visualization.



3 TECHNICAL ARCHITECTURE

The system architecture of the DIGIT.B4 – Big Data PoC (Proof of Concept) for the Social Media analysis will be based on the NTT Data – everis reference architecture.

The reference architecture is a collection of building blocks which consist of a series of technologies and methodologies.

NTT DATA and everis reference architecture is composed of 2 levels:

- 1. **Layer:** is a group of building blocks which have the same technology/methodology. There are 8 layers that are categorized for all technologies/methodologies important for the architecture.
- Building Blocks: each one represents a technology/methodology. There are 38 building blocks in the architecture.



The blocks that will be used in the PoC for each layer have been highlighted.

Figure 11 - System Architecture

The technologies chosen for each building block are listed, and the descriptions for the use that the PoC will make of each highlighted block are listed below:



3.1 Information gathering

ETL (Extract, Transform and Load)		
Building Block description	Extract, Transformation and Load (ETL) is an industry standard term used to represent the data movement and transformation processes.	
Building Block PoC use	The files obtained from the data sources have to be transformed to other formats in order to be readable by the classification model.	
Building Block technologies	R, python	

Table 1- Information gathering

3.2 Information Store

Relational Data Base	
Building Block description	RDBMS (Relational Database Management System) is a software/middleware to manage relational data. It manages the data in the table that conforms to relational model.
Building Block PoC use	The gathered information will be stored in a structured manner.
Building Block technologies	Oracle Database, Oracle Exadata Machine, SAP IQ, SAP HANA, SQL Server, Parallel Data Warehouse, PureData, DB2, Teradata, Vertica, PIVOTAL, MySQL

Table 2 - Information Store

3.3 Data Analytics

Data Mining, Text Mining, Machine Learning			
	Data Mining : Data mining refers to the statistical techniques and rule-based models that aim at discovering new knowledge and trends from a large data set. In contrast to machine learning, data mining emphasizes on the process of analysing data rather than the accuracy of the used models.		
Building Block description	Text Mining : Text mining can be defined as an intensive process in which a user interacts with a document collection over time by using a suite of analysis tools to extract knowledge.		
	Machine Learning : Machine learning is a scientific discipline that explores the construction and study of algorithms that can learn from data. Such algorithms operate by building a model based on inputs that can be used to make predictions or decisions, rather than following only explicitly programmed instructions.		
Building Block PoC use	The classification model and algorithms used for the PoC will use techniques from all of these blocks.		
Building Block technologies	R, python		

Table 3 - Data Analytics

3.4 Governance

Quality of Data		
Building Block description	The goal of a data quality exercise is to establish high-quality data and maintain it. The main tools/techniques that help to be successful in data health are related to the Master Data Management (MDM).	
Building Block PoC use	The gathered data has to be treated to ensure its quality before sending it to the classification	



	model.
Building Block technologies	R, python

 Table 4 - Governance

3.5 Decision Support/Utilization

Visualization	
Building Block description	Data visualization means showing data in a graphical format. It involves the creation and study of the visual representation of data. The primary goal of data visualization is to communicate information clearly and efficiently to users via graphics, such as tables, charts and maps.
Building Block PoC use	The PoC results will be published in a web application where several charts and graphics will illustrate the outcome of the classification model.
Building Block technologies	R, JavaScript, HTML

 Table 5 - Decision Support/Utilization