FAST LISA E-Participation platform

15.09.2023

Chiara Catizone, Giulia Venditti, Monica Palmirani – CIRSFID monica.palmirani@unibo.it chiara.catizone2@unibo.it giulia.venditti2@unibo.it



Index

Introduction

CONSULTATION AD DATA GATHERING

- e-participation UI (consultation):
 - Piattaforma di aggregazione attività patrticipation, come arrivare alle consultazioni (chi le fornisce), descrizione UI pagina consultazione (main, metadati, form commenti)
- Moderation:
 - Dashboard wordpress moderatori flow:
 1. login (credenziali fornite), 2. moderazione commenti (azioni principali)

FASTILISA

- Cosa succede quando la consultazione viene chiusa?
 - 1.Metadati cambiano da status:Active a status:Closed 2. form comme nti viene chiuso, 3. dati scaricati ed elabortai da Al 4. Dati pubblicati in una dashboard interattiva.

DATA ELABORATION

- Flow elaborazione dati
- Preliminary studies
 - Hate Speech Classification
 - Taxonomy
 - Arguments analysis

Index

- Software developement
 - Flow
 - Classificators selection and fine-tuning
- Dissemination
 - JSON format
 - Upload on the server

RESULTS VISUALIZATION

 Dashboard Access and UI (header, help button, info button)

- Dashboard full visualization division in fore area
 - Left (»level of severity»)
 - central ("lexicon" and "hate speech classification"),
 - right ("purpose"),
 - bottom ("discussion monitoring")
- Report
 - Purpose
 - Structure

INTRODUCTION



Objective

The creation of an **online portal** at <u>http://fast-lisa.unibo.it</u> (based on WordPress), designed to ensure the **protection of personal data** while **fostering open public discussions** led by the **public administration** on topics closely linked to the everyday concerns of **citizens**.

The main objective is to assess the quality of these discussions, detect potential instances of hate speech, and provide decisionmakers with a concise summary of prevailing opinions.

This approach empowers public administrations to **evaluate policy effectiveness, monitor shifting citizen needs**, and **make more informed and impactful decisions**.

FASTLISA



- INTRODUCTION

System Architecture

- **1. Public Administration Operators** initiate topic-specific consultations for citizen engagement.
- 2. Citizens search and participate in specific consultations leaving comments stored in anonymous form.
- **3.** Debates are moderated by FAST LISA ambassadors for productive and respectful discussions.
- 4. Al algorithm analyse comments and messages for hate speech detection and cluster arguments into pros and cons.
- 5. Interactive Dashboards provide clear visualisations of aggregated hate speech data, promoting transparency and aid decision-making for public administration.



- INTRODUCTION

CONSULTATION & DATA GATTERING





FAST LISA e-Participation platform

- CONSULTATION & DATA GATTERING

FAST LISA e-Participation platform works as a hub of participation activities across different partner organisation.

Home page

Primary Entry Point. It works as central hub of information

Partner Organisation Customise in the respect of Public Organisation community cultural and

linguistic expectation



Public Consultation

- CONSULTATION & DATA GATTERING

Extensive explanation of the **discussion topic** and **participation rules** for informed and meaningful contributions.

Discussion forum for **anonymous commenting** allows for:

Filtering options to arrange comments chronologically.

Navigate to prioritize reacted ones, or focus on highlighted contributions.

Reply to others comments and leave a reaction (like/dislike)

FAST-LISA Fighting hAte Speech Through a Legal, ICT and Sociolinguistic approach	
FAST-LISA Comune di Ravenna Santa Coloma de Gramenet Pro Arbeit	
« BACK	
Ravenna / Consultazione pubblica / A casa mia?	STATUS: Active
A casa mia?	WHO: Famiglie
	HOW: Accoglienza
31 JULY 2023	WHAT: Integrazione rifugiati
Il suo compito è favorire l'autonomia dei cittadini immigrati in un'ottica di coesione sociale.	
Attraverso servizi, attività e una costante collaborazione con cittadini, istituzioni e associazioni si vuole promuovere una maggiore consapevolezza dei processi migratori ed una cultura	TIME: 01/08/2023 9:00 am
dell'uguaglianza come fondamenta di una convivenza civile in una società sempre più interculturale.	START: 01/08/2023
Fante famiglie hanno sperimentato l'accoglienza di rifugiati in casa propria. Jivere con persone del luogo, infatti, può ajutarli a sentirsi parte di una comunità a creare.	END: 31/08/2023
una rete di relazioni, ad accrescere le competenze linguistiche e a potenziare l'accesso a migliori opportunità lavorative. Chi apre le porte di casa ha l'opportunità di conoscere una	CONTACT: accoglienzastranieri@comune.ra.it ^[2]
nuova cultura, di diventare un cittadino più consapevole e attivo, di vivere un'esperienza umana che arricchisce.	PROPONENT: U.O. Politiche per l'immigrazione e la
Puoi visionare il progetto promosso dal Comune di Ravenna qui. 🗹	

Consultation Metadata

Activation status (active or closed)

- CONSULTATION & DATA GATTERING



Debate Moderation

- CONSULTATION & DATA GATTERING

FAST LISA Ambassadors, thanks to provided credential, can access the WordPress Dashboard where they can perform a series of action on users comments:

- Approve/Unapprove
- **Reply** to comment to foster and guide dialog
- **Stick** at the top of the conversation
- Edit comment for example when personal data are shared
- **Close** reply to a particular comment
- Mark as Spam
- Trash the comment





Consultation Closure

- CONSULTATION & DATA GATTERING

- 1. Consultation Status is changed to Closed.
- 2. Comment form is deactivated but posted comments are steal visible and navigable.
- 3. Data gattered are processed by the AI algorithm
- **4. Results** are **published** and accessible through the "Dashboard" and "Report" buttons.

	omune di Ravenna	Santa Coloma de Gramenet Pro	Arbeit	
« BACK				
Ravenna / Consultaz	zione pubblica / [TEM	/A DI PROVA] Ricon		
		<i>ا</i> م		Dashboard 🕒 Report 🗎
		A]		
Riconoscimento legale dei				status: Closed
matrimoni tra persone dello			WHO: cittadini	
stesso sesso			HOW: unioni civili	
30 MAY 2023				WHAT: unioni civili tra persone dello stesso sesso
Proponente: Uffici	o di Stato Civile			
L'Ufficio di Stato Civili	e è un ente all'interno	li un comune o di un municipio che si occupa vita di una parsona, como passita, matrimon	a di i e decessi	TIME: 22/05/2023 12:00 am
registrare gli eventi ci	ivili fondamentali della	VITA OF DUA DELSONA COME DASCHE THAT HUDDI		

DATA ELABORATION: attori e workfow





Workflow

- DATA ELABORATION

Premises

Focus on detecting and describing **hate speech** and the **flow** of the discussion

- Legal, social and linguistic basis for the taxonomy
- **Pragmatic aspect** of argumentation must be highlighted

Methodological pillars:

- Lexicon oriented to the legal domain Classes extracted from legal documents to support the creation of a trans-border taxonomy for classifying online hate speech.
- Situations of the hate speech (pragmatic aspect)
 Designing a methodology for extracting arguments from raw text to understand the users' inclinations when expressing their opinion in institutional context.



- DATA ELABORATION

Hate Speech Classification – I

Online Hate Speech can be automatcally descripted by AI tools trained in classification tasks.

To automate the descriptive process we first needed to set a transborder taxonomy to classify data in an harmonyzed – yet legallly valid – way.

"A taxonomy (or taxonomical classification) is a scheme of classification, especially a hierarchical classification, in which things are organized into groups or types."

- Wikipedia

In literature research step we gathered existingtaxonomies used for describing hate speech.

Unfortunately they have no legal value, necessary for our project as specified in the premises.

The <u>Future of free speech's taxonomy</u>, unique with legal value, because it wasvbuilt starting from European Court of Human Rights (ECHR) documents.

Hate Speech Case Database





Hate Speech Classification - II

All the legal documents have been gathered by domain experts.

Analysis tool:

KNIME, an Open Source software for data analysis

Following (Salminen et al. 2018) this methodology is based on the **open coding technique**, in which classes emerge from the material.

- DATA ELABORATION

Features:

- Information retrieval techniques like TF-IDF, which has been useful to retrieve important terms (Rajaraman and Ullman 2011; Luhn 1957; Robertson 2004; Vrysis et al. 2021)
- Latent Dirichlet Allocation was used to model topics and extract them atics proposed by (Blei 2003; Salminen et al. 2018)
- We used embeddings (Doc2Vec and Word2Vec) and word distances to support the linguistic categorization of terms related to hate speech (Lewandowska-Tomaszczyk et al. 2021).

- DATA ELABORATION

Hate Speech Classification – III

In brief, for each investigated country we were able to retrieve the protected characteristics constituting criminal grounds –i.e., bias motivation– in the national legislation and a lexicon of abusive behaviors characterizing the crimes targeting minorities.

Racism and ethnic hatred are the most common bias motivation among the investigated countries



FASTLISA

- DATA ELABORATION



- DATA ELABORATION

Argument Mining -I

Objective:

How to effectively structure argumentations from an online institutional environment?

Following the principle by (C. Joshi 2019) "bigger data doesn't always translate into better decisions", we employed text mining techniques to explore and compare the structure of argumentations and the behaviour of users to understand how to develop the argumentation mining process.



Argument Mining - I

Source:

European Commission's consultations collection (European Commission 2023b)

Data Preprocessing:

- 1. Reduced the number of total topics from 91 to 40
- 2. Anonymized the data
- 3. Remove unnecessary words
- 4. Random sapling (500 comments each lang)



- DATA ELABORATION

Argument Mining - III

Considerations on Data Exploration

- Non expert users are mainly Italian and German
- Regulation and Transport have the highest number of paticipants
- Spain is the country with the lowest participation rate



COnsiderations on TF and TF-**IDF** results

•Automatically detect and extrac t frequently occurring phrases from the responses Discourse indicators like "because" and "however" + mod al verbs emerge from all corpora •European citizens from differen t countries share similar worries and preferences towards the activities





- DATA ELABORATION

Final Remarks

Arguments analysis was a valuable tool for understanding how to classify and manage the commentsts, independently from their language to gather important consideratios:

- We base our work on a theory that studies argumentation schemes, where only premises and conclusions are classified (Mochales and Moens 2011)
- Explicit if a unit is in support or against the main thematic
- Explicit their relations using similarities

Software developement

- DATA ELABORATION

General workflow

Do not start building models froms scratch, reuse (models and training data) when possible.

Five different classification tasks :

- 1. Hate speech presence
- 2. Target identification
- 3. Bias motivation identification
- 4. Claims and premise s identification
- 5. Support/contro relationships identifification





Software development

- DATA ELABORATION

Classificators Selection

Testing of available hate speech detectors in each language (it, de, es).

Selected models for HS classification:

The model we relieve on for argument mining task:

<u>mDeBERTa-v3-base-xnli-multilingual-nli-</u> <u>2mil7</u>

Further training is needed for...

- Target
- Bias motivation
- Claims and premises

Models used for finetuning

Software development

- DATA ELABORATION

Steps:

- Detect hate speech presence (pie-chart)
- Detect hate speech targets (bar-chart, tree-map)
- Detect hate speech bias motivations (heat-map, treemap, sankey diagram)
- Detect argument units (claims and premises) (arguments map)

- Predict support or cotraddiction relationships between argument units (sankey diagram, arguments map)
- Compute duration time of the discussion, log data of the interactions. (line-chart and counters)

Results elaboration

- DATA ELABORATION

		{
2	>	"pie_chart" : [[…
	>	"bar_chart" : […
14],
15	>	"heat_map" : […
18],
19	>	"line_chart" : […
34],
35		"sankey_diagram" : [
36		{"data": [
37		["Pro", "HS", 10],
38		["HS", "Class 1", 1],
39		["HS", "Class 2", 3],
40		["HS", "Class 3", 5],
41		["Pro", "NO HS", 90],
42		["Contro", "HS", 180],
43		["HS", "Class 4", 80],
44		["HS", "Class 6", 30],
45		["HS", "Class 8", 70],
46		["Contro", "NO HS", 9]
47		1}
48],

Dissemination

- Save the results in JSON format on local machine
- Upload the JSON file to FAST LISA's server

Limitations:

A future fine-tuning iteration might be needed for improving of the system. W deal with low reasource language, i.e., lack of ready-to-consume NLP tool.

RESULTS VISUALIZATION



Interactive Dashboard

- **RESULTS VISUALIZATIONS**

Context: Displaying relevant information about the discussion, including the title, proponent, date, and a brief summary of the central argument.

Transparency: Info buttons providing detailed information about the purpose and data represented in each visualization.

Help Section: Offering support and information about the dashboard's features and functions.

Privacy: Presenting data in an aggregated and anonymized manner.

Solution to the performance of the performance o



Left Column

- **RESULTS VISUALIZATIONS**

"level of severity"

A **pie chart** provides a clear ratio between Hate Speech and neutral comments.

Quickly grasping the **conversation's tone** and **hateful language presence**.



A **sankey diagram** displays the composition of hate speech within pro and counter comments.

Effectively illustrating the flow and **distribution of comments**. Highlights **hate speech presence** and **different classifications**.

Simplifies hate speech analysis, aiding in comprehension and discussion argumentations' dynamics.



Central Area

"lexicon" & "hate speech classification"

A tree map displays hierarchical relationships between Hate Speech classes and frequently used words in classified Hate Speech comments. Helps users understand the content associated with each class (e.g., racial, religious).



FASTLISA

An argument map highlights primary arguments and their relationships within the discussion. Offers a critical overview of the discussion's progression. Enables users to form well-informed opinions by visualizing key debating points.



- RESULTS VISUALIZATIONS

FASTLISA

Motivazioni linguaggio abusivo

1

0

0

3

0

1

4

5

3

2

0

3

7

20

15

A heat map illustrates the relationships between bias motivations within the discussion.

"purpose"

Right Column

Reveals correlations among various biases, aiding in understanding their presence and overlap.

Utilizes a classification system categorizing biases, including ethnic origin, race, religion, nationality, gender, and migrant status.

 $\equiv \mathbf{0}$

Target linguaggio abusivo

0

10 20 30 40 50 60 70

Highcharts.com

 $\equiv \mathbf{0}$

61

52

The target group victim of the offense is examined through a bar chart. It showcases the distribution of Hate Speech comments based on their respective targets.

Distinguishes between single and group targets, including specific characteristics like gender, race, and age.

- RESULTS VISUALIZATIONS

- **RESULTS VISUALIZATIONS**

Bottom Area

"discussion monitoring"

A **line chart** tracks trends in both neutral and Hate Speech comments throughout the entire discussion. Offers a comprehensive view of the conversation's progression, identifying points of increased Hate Speech prevalence or vice versa. Enhances understanding of discussion dynamics and patterns, aiding in analysis and critical evaluation.

FASTLISA

A representation of **comments engagement** showcases the number of likes and replies received by classified Hate Speech comments. Enables users to gauge the influence of Hate Speech on the discussion and the attention it receives from other users. Valuable for understanding Hate Speech's potential impact on the community and promoting awareness and critical evaluation.



Report

Purpose

Support the understanding of data through brief descriptions of the visualized data

- RESULTS VISUALIZATION

Structure

- 1. Introduction
 - Discussion name
 - Consultation period (start and end-date)
 - Proposing office/comunality
- 2. Dashboard purposes
- 3. Introduction to classes
- 4. Description (with picture) of each viz
- 5. Conclusion/duiscussion

Thank you.

