

What is this about?

- Language technologies in the job of a translator are a must
- [...] rely on machine translation to support their regular endeavours
- Technology [...] based on AI, NLP, ML
- a translator that is not only a user, but could eventually take part in the development of the technology

Rough programme

- An introduction to language technologies
- A path to jump from translation into language technologies
- Understanding IBM M1, the core of statistical machine translation
- A quick overview of the python programming language

3

Into Language Technologies



1

Alberto Barrón-Cedeño

Alma Mater Studiorum - Università di Bologna

a.barron@unibo.it
@_albarron_

Session 1

19 December 2022

Into Language Technologies I

Department of Interpreting and Translation

- Born in 2012 (merging the SITLeC Dept. and the Scuola Superiore di Lingue Moderne per Interpreti e Traduttori)
- Emphasis in **applied research**, theoretical, practical, and didactic aspects of **translation** and interpreting

Degrees

- Bachelor in Intercultural and Linguistic Mediation
- Masters in Interpreting
- **Masters in Specialized Translation**
- PhD in Translation, Interpreting, Interculturality



Session 1

19 December 2022



2

Why Language Technologies?



Side-eye of the tiger (^_~ . -^)

made a lil meme about the job search

A collage of two images showing a person smiling over a tray of donuts. The left image shows a close-up of the person's face and the donuts. The right image is a slightly darker version of the same scene. Below the images, there is text on the left side and right side.

3:08 am · 18 May 2021 · Twitter Web App

Language data and project specialist

- knowledge, skills and competences -

DISCIPLINARY

- | |
|--|
| Knowledge of specific languages |
| Ability to conduct linguistic analysis |
| Translation, interpreting, |
| Proof-editing, localisation |

- (INTER)CULTURAL**

 - Awareness of specific cultural contexts and cultural differences
 - Cultural agility

TECHNICAL

- Understanding and use of **language technologies** and resources
 - Knowledge of a programming language
 - Understanding of **computational linguistics** / NLP

ORGANISATIONAL

- Entrepreneurship
 - Project management
 - Quality control
 - Planning
 - Teamwork

medium.com

1. Why Language Technologies?

2. How CS is being *plugged* into DIT
 3. Teaching initiatives
 4. Three student projects

7

Overview

1. Why Language Technologies?

2. How CS is being *plugged* into DIT
 3. Teaching initiatives
 4. Three student projects

4

Why Language Technologies?

medium.com

Adapted from Freepik Storyset,
<https://storyset.com/data>

Adapted from Freepik Storyset,
<https://storyset.com/data>

LT at DIT is (mostly) built on UPSKILLS



How CS has been (re)plugged into DIT

**UPgrading the
SKILLS of
Linguistics and
Language
Students**

An Erasmus+ strategic
partnership September 2020 -
August 2023

<https://upskillssproject.eu>

11

Hiring of computing scientists

Spring 2019
Professor with NLP background



قطرية للبحوث الحاسوبية
Qatar Computing Research Institute
Member of Qatar Foundation

Winter 2020
Research assistant with MT background



Microsoft
Research

* but not only!

12

Into Alma AI

DIT adhered to UniBO's Alma Mater Research Institute for Human Centered Artificial Intelligence



Foundations of AI	AI for health and well-being
AI and hard sciences	AI for law and governance
Humanistic AI	AI and education
AI for industry	AI and high performance computing

13

Initiatives with CS load

- Neural machine translation
- MT of academic websites
- Interaction with local companies
- Webinars and workshops on MT and related technologies
- Automatic identification of misconducts online (**antenna**)
- Automatic identification of MT suitability (**MTsweet**)
- Discussions on the curriculum in 5 years time
- Conveying that **MT is not an enemy**

15

- Official objective 1.** Giving a gentle introduction to programming in python to get students in the right position to go further on their own
- Official objective 2.** Paving the way to the **NLP and other courses**



say "**that would be so awesome! But stop... it needs programming**"

Teaching initiatives



DIT Python course 2020, 2021, 2022*

* 2020-2021: crash course, open to all; 2022: part of Profession-based research and PhD seminars

16

DIT Python course 2020

DIT Python course 2020

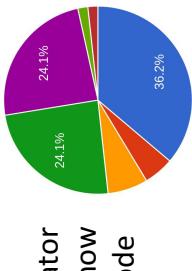
Session 1. The basics

- What is a programming language?
- What is an algorithm
- “Translating” from an algorithm into a program
- The characteristics of the python programming language
- Basic functions, variables, conditionals, loops

```
# my code
x = 0
while x < 50:
    for i in range(x):
        print('x', end='')
    print()
    x += 1
```

19

Pre-entry survey



Absolutely

Yes, if (s)he wants to do research (not for industry)

Yes, if (s)he wants to go for industry (n...
Yes, if (s)he is targetting the software/...

Maybe

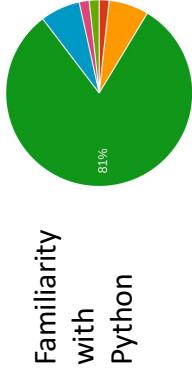
Nope. There is out-of-the-box software...
Yes, if (s)he is targeting the software/...

I hate it

Isn't it a snake?

last year / attended an online course 0...
So far I've only coded a few scripts for...

17



I use it on a regular basis

I can code a few routines, with a lot of effort

I have passive knowledge (I can read it, but I cannot produce it)

I've heard about it, but I don't know it...

I hate it

Isn't it a snake?

last year / attended an online course 0...
So far I've only coded a few scripts for...

17

DIT Python course 2020

Session 2. Python 4 Poets (1/2) (derived from K. Church's Unix for poets)

Course structure

Four+ sessions

- Opening text files
- Splitting into words (tokenisation)
- Obtaining vocabularies (types)
- Extracting n-grams

1. Presentation of concepts with the support of slides

2. Live on-screen coding of task-specific routines

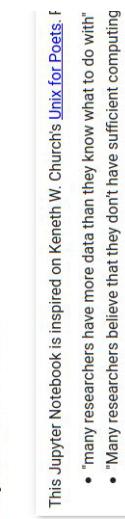
3. Take-home simple coding exercises

Coding platform

Jupyter notebooks on Google's **colab**

<http://colab.research.google.com>

Coding platform



- This Jupyter Notebook is inspired on Kenneth W. Church's [Unix for Poets](#) f
- "many researchers have more data than they know what to do with"
- "Many researchers believe that they don't have sufficient computing

20

18

NLP course

DIT Python course 2020

Learning outcomes. [...] basic theoretical aspects of computational linguistics [...] acquire practical skills [all the way to] **supervised models** (machine learning)

Specialised translation: optative

Translation and technology: core

Session 3. Python 4 Poets (2/2) (derived from K. Church's Unix for poets)

- Finding specific tokens/strings
- Finding palindromes
- String substitutions
- Functions
- Collocations

A collocation is a series of words or terms that co-occur more often than would be expected, following [Magedman and Marcus](#), in NLP it can be estimated as

$$MI(x, y) = \log_2 \frac{Pr(x,y)}{\frac{\sum_{j=1}^n Pr(x_j)}{\sum_{j=1}^n Pr(x_j)} \cdot \frac{\sum_{j=1}^n Pr(y_j)}{\sum_{j=1}^n Pr(y_j)}}$$

where \sum . is the sum over all instances of .

```
[ ] from math import log  
bigrams = ngrams(tokens, 2)  
unigrams = ngrams(tokens, 1)
```

23

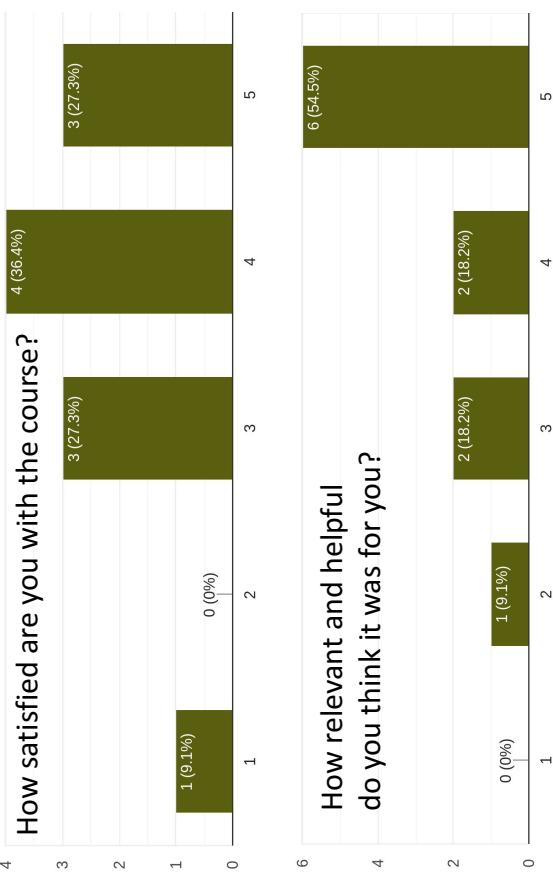
21

NLP course

Course structure

1. Presentation of concepts with the support of slides
2. Live on-screen coding with simple running routines + voluntary homework
3. Evaluation based on one final project
(with potential to become a publication)

Closing perception



24

22

AriEmozione

NLP course

Rough contents: **coding, statistics, and machine learning applied to text**

Objective. Identifying the **emotion** transmitted in 17th/18th-century Italian opera arias at the verse level



Francesco Fernicola

Specialized Translation Masters

Developed in the context of Unibo's Centro per l'Interazione con le Industrie Culturali e Creative (<https://site.unibo.it/cricc/it>)



Shifeng Zhang

Building blocks...

1. Introduction to natural language processing (algebra)
2. Tokens and the vector space model (statistics)
3. The Naïve Bayes classifier (machine learning)
4. The training and evaluation process (algebra)
5. Word vectors (algebra)
6. Latent semantic analysis (semantics)
7. Neural networks (machine learning)
8. Word embeddings (semantics)
9. Convolutional neural networks (machine learning)
10. Recurrent neural networks (machine learning)
11. Text generation

27

25

AriEmozione

AriEmozione 1.0 corpus

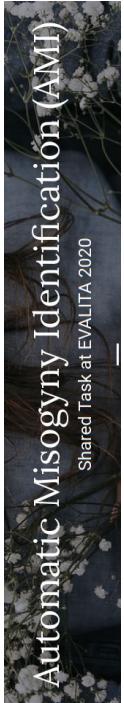
- 678 operas composed between 1655 and 1765
- All texts are written in Italian of the period and articulated in verses
- 2,473 verses manually annotated in six classes
 - Amore (Love)
 - Gioia (Joy)
 - Ammirazione (Admiration)
 - Rabbia (Anger)
 - Tristezza (Sadness)
 - Paura (Fear)



Three student projects

28

AriEmozione



Objective. Identifying whether a tweet is misogynous and, if it is, whether it is aggressive

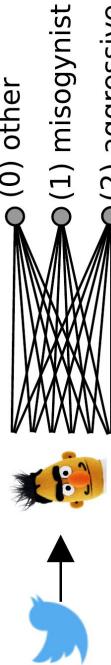
**Language, Society,
and Communication**
Master's

UniBO @ AMI 2020

Two classification tasks

- | | |
|---|--------|
| 1. Is this tweet misogynous? | YES/NO |
| 2. Is this misogynous tweet aggressive? | |

YES/NO



Our solution: one multi-class network built on top of ALBERTo

Our solution: one multi-class
(an Italian version of BERT)



Specialized Translation Masters



6

model	representation	10-fold CV			test
		F ₁	Acc	F ₁	Acc
kNN	char 3-grams	0.38	38.51	0.35	35.15
	words	0.36	36.08	0.35	34.73
	LDA char	0.30	29.97	0.31	30.54
SVM-RBF	char 3-grams	0.44	43.70	0.43	43.00
	words	0.42	42.00	0.44	44.00
	LDA char	0.28	28.00	0.30	30.00
Log reg	char 3-grams	0.44	45.57	0.42	43.10
	words	0.41	43.20	0.41	43.10
	LDA char	0.28	30.63	0.29	30.96
2-layers NN	char 3-grams	0.42	43.61	0.47	46.86
	words	0.42	42.91	0.43	43.10
	LDA char	0.27	29.56	0.27	31.80
3-layers NN	char 3-grams	0.49	41.86	0.40	41.84
	words	0.47	42.60	0.40	41.84
	LDA char	0.26	31.41	0.30	31.80
FastText	char 3-grams	0.43	45.00	0.41	42.37
	pre-trained chars	0.43	47.00	0.41	41.00
	words	0.42	42.56	0.39	44.07
	pre-trained words	0.38	41.00	0.40	42.00

29

Papers published in CLLC-it 2020 and IJCOL

AriEmozione: Identifying Emotions in Opera Verses
Francesco Fenicola¹, Shihengfeng Zhang¹, Federico Garcea¹
Paolo Bonora², and Alberto Barron-Cedeno¹
¹Department of Interpreting and Translation
Università di Bologna, Forlì, Italy
²Department of Classical Philology and Italian Studies
Università di Bologna, Bologna, Italy

Francesco Fernicola
Università di Bologna

Francesco Fornicola
Università di Bologna

Paolo Bonora[†]
Università di Bologna

Code available on github

Are fictional voices different?

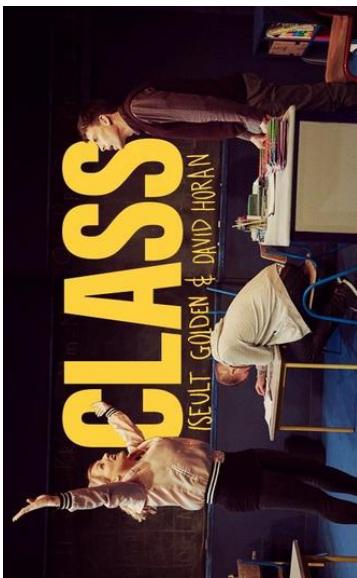
UniBO @ AMI 2020

team	run	constrained	score
UniBO ^a	2	yes	0.7438
jigsaw	2	no	0.7406
fabsam	1	no	0.7380
YNU_OXZ	1	yes	0.7343
fabsam	1	no	0.7314
NoPlaceForHateSpeech	2	yes	0.7309
YNU_OXZ	2	yes	0.7167
fabsam	2	no	0.7015
NoPlaceForHateSpeech	3	yes	0.6948
AMI.the_winner	1	yes	0.6934
MDD	3	no	0.6869
PoliTeam	3	yes	0.6844
MDD	1	yes	0.6820
PoliTeam	1	yes	0.6810
MDD	2	no	0.6679
AMI.the_winner	1	yes	0.6653
PoliTeam	2	yes	0.6473
UniBO ^b	1	yes	0.6343
AMI.the_winner	3	yes	0.6259
NoPlaceForHateSpeech	3	yes	0.4902

Objective. Finding out if the authors of a specific play managed to create recognisable fixional voices



Ettore Galletti



Specialized Translation Masters

Misogyny: 0.8102
Aggressiveness: 0.6774

35

33

Are fictional voices different?

Identifying characters

- McCafferty (teacher)
- Brian (father)
- Donna (mother)
- Jayden (son)
- Kaylie (Jayden's schoolmate)
- None (stage directions)

vs Identifying groups of characters

- Male
- Female
- None

Main drawbacks: still weak against aggressiveness

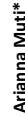
Further efforts: multimodal (memes) and multilingual models; implicit hate speech

UniBO @ AMI 2020 Task A

Papers at Evalita 2020, SemEval 2022, and more

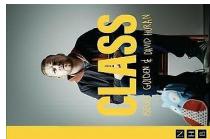
- UniBO@AMI: A Multi-Class Approach to Misogyny and Aggressiveness Identification on Twitter Posts Using ALBERTo
Arianna Mutti
Alberto Barron-Cedeno
Department of Modern Languages,
Literatures and Cultures - LILFC
University of Bologna
Bologna, Italy
arianna.mutti@scienze.unibo.it

Code available on github



Arianna Mutti*

Language, Society,
and Communication
Masters



UniBO at Evita 2020 AMI Task A
The paper is being presented at UniBO at Evita 2020

National media coverage



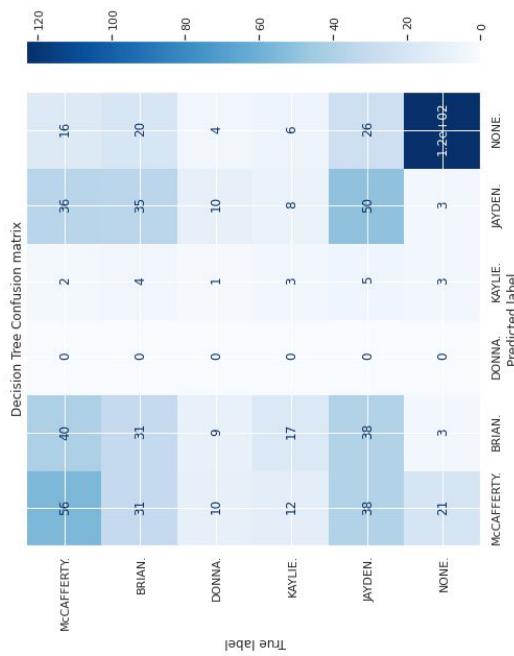
L'algo-ritmo di Arianna: "Cosa su Twitter dà la faccia ai post contro le donne"
di Giacomo Saccoccia

* PhD, UniBO

Are fictional voices different?

Are fictional voices different?

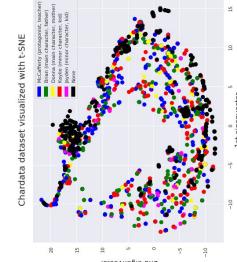
Supervised approach. Build a multi-class character classifier and study whether it manages to label the interventions accurately



39

Are fictional voices different?

Results: non-conclusive



Main drawbacks: we observe some hints, but we need to study the problem further

Italian translation: the characters are (a bit) more discriminative

Core idea

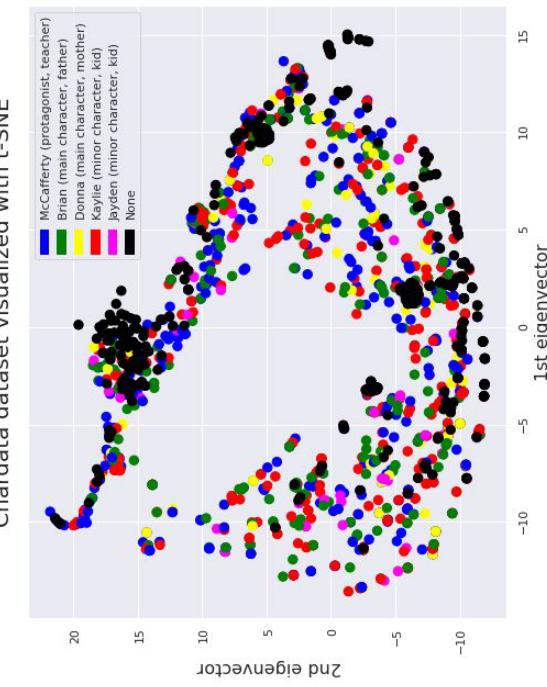
1. Build a topic-independent representation of every character intervention
2. Observe if the representations of all interventions make the characters (clearly) differentiable

Core idea

37

Are fictional voices different?

Unsupervised approach. Cluster all the instances and analyse at what extent the clusters correlate with the characters



40

38

Next...

Are fictional voices *different*?

上海(shang4 ha3) 浦东(pu3 dong1) 开发(kai1 fa1)
the development of shanghai's pudong



Ettore Galletti*



Python

IBM M1

Specialized
Translation
Masters

43

41

* Ufficio Rettorato,
UniBO

Interested, questions?



a.barron@unibo.it

@_allbarron_

Thanks!



Automatic identification of propaganda in text



CheckThat!
Automatic prioritisation and verification of claims

44

42

Other (non-student-led) efforts