

Raport de Proiect

Modul spaCy pentru Limba Aromână (rup)

Decembrie 2025

Rezumat

Acest raport prezintă dezvoltarea unui modul complet de procesare a limbajului natural pentru limba aromână (macedo-română) în cadrul bibliotecii spaCy. Proiectul include implementarea componentelor esențiale: tokenizare, lematizare, lista de stop words, atrbute lexicale și utilitare pentru conversia între standardele ortografice. Limba aromână este o limbă romanică orientală vorbită în zona Balcanilor, având aproximativ 250.000 de vorbitori.

Cuprins

1 Introducere	3
1.1 Context și Motivație	3
1.2 Obiective	3
2 Arhitectura Proiectului	3
2.1 Structura Folderului <code>spacy_rup</code>	3
2.2 Componentele Principale	4
2.2.1 <code>__init__.py</code> - Clasele Language	4
3 Implementare Detaliată	4
3.1 Tokenizare	4
3.1.1 Reguli Specifice pentru Cliticele Aromâne	4
3.1.2 Excepții de Tokenizare	4
3.2 Lematizare	4
3.2.1 Lematizare Verbală	5
3.2.2 Lematizare Nominală	5
3.2.3 Reguli cu Sufixe	5
3.3 Stop Words	5
3.4 Atribute Lexicale	6
3.4.1 Detectare Numerală	6
3.5 Conversia Ortografică	6
3.5.1 Standardul DIARO (Caragiu-Marioțeanu)	6
3.5.2 Standardul Cunia (Tiberiu Cunia)	7
3.5.3 Funcțiile de Conversie	7

4 Integrare și Utilizare	7
4.1 Instalare	7
4.2 Exemple de Utilizare	8
4.2.1 Pipeline de Bază	8
4.2.2 Cu Lematizare	8
4.2.3 Cu Conversie Ortografică	8
5 Date și Corpus	8
6 Realizări și Metrici	9
6.1 Componentele Implementate	9
6.2 Statistici Lematizator	9
6.3 Acoperire Ortografică	10
7 Provocări și Soluții	10
7.1 Provocări Întâmpinate	10
7.2 Decizii de Design	10
8 Direcții Viitoare	11
8.1 Îmbunătățiri Pe Termen Scurt	11
8.2 Obiective Pe Termen Lung	11
9 Impact și Aplicații	12
9.1 Importanță pentru Limba Aromână	12
9.2 Aplicații Potențiale	12
10 Concluzii	12
10.1 Sinteza Lucrării	12
10.2 Contribuții Principale	13
10.3 Reflecții Finale	13
A Cod Sursă - Exemplu Lematizator	13
B Exemple de Text Procesate	14
B.1 Exemplu 1: Propoziție Simplă	14
B.2 Exemplu 2: Propoziție Complexă	14

1 Introducere

1.1 Context și Motivație

Limba aromână (cod ISO 639-3: `rup`), cunoscută și sub numele de macedo-română sau armăneashti, este o limbă romanică orientală vorbită în Grecia, Albania, Macedonia de Nord, România, Bulgaria și Serbia. În ciuda numărului semnificativ de vorbitori (aproxi-mativ 250.000), limba aromână este clasificată ca fiind în pericol și beneficiază de resurse NLP limitate.

Acest proiect a avut ca scop crearea unui modul complet pentru procesarea limbii aromâne în cadrul bibliotecii spaCy, una dintre cele mai populare biblioteci Python pentru procesarea limbajului natural.

1.2 Obiective

Obiectivele principale ale proiectului au fost:

- Implementarea unui tokenizer adaptat particularităților morfologice ale limbii aromâne
- Dezvoltarea unui lematizator funcțional pentru verbe și substantive
- Crearea unei liste complete de stop words (cuvinte funcționale)
- Implementarea atributelor lexicale (detectare numere)
- Dezvoltarea utilitarelor pentru conversia între standardele ortografice
- Integrarea completă cu ecosistemul spaCy 3.x

2 Arhitectura Proiectului

2.1 Structura Folderului `spacy_rup`

Modulul a fost organizat în mai multe fișiere specializate, fiecare având o responsabilitate bine definită:

Structura Proiectului

```
spacy_rup/
    __init__.py                  # Clasa principala Language
    lemmatizer.py                # Logica de lematizare
    lemma_component.py           # Componenta pipeline spaCy
    stop_words.py                # Lista de stop words (163+)
    tokenizer_exceptions.py      # Exceptii tokenizare
    punctuation.py               # Reguli punctuatie
    lex_attrs.py                 # Atribute lexicale
    orthography.py              # Conversie ortografica
    examples.py                  # Exemple de utilizare
    README.md                    # Documentatie
```

2.2 Componentele Principale

2.2.1 `__init__.py` - Clasele Language

Fișierul principal definește două clase fundamentale:

```
1 class AromanianDefaults(Language.Defaults):
2     tokenizer_exceptions = {...}
3     prefixes = TOKENIZER_PREFIXES
4     suffixes = TOKENIZER_SUFFIXES
5     infixes = TOKENIZER_INFIXES
6     lex_attr_getters = LEX_ATTRS
7     stop_words = STOP_WORDS
8
9 class Aromanian(Language):
10    lang = 'rup'
11    Defaults = AromanianDefaults
```

Acste clase implementează interfața standard spaCy pentru o nouă limbă, integrând toate componentele dezvoltate.

3 Implementare Detaliată

3.1 Tokenizare

3.1.1 Reguli Specifice pentru Cliticele Aromâne

Limba aromână are un sistem complex de clitice care se atașează verbelor și pronumelor. Tokenizatorul a fost configurat să trateze corect aceste structuri:

- **Clitice verbale:** s-, sh-, n-, lj-
- **Articole definite:** -lu, -a, -lji
- **Pronume enclitice:** -mi, -ti, -si

3.1.2 Excepții de Tokenizare

Am implementat excepții pentru:

- Numerale ordinale: 1-a, 2-lea, etc.
- Abrevieri specifice: d-lu (domnul), d-na (doamna), etc., nr.
- Expresii compuse frecvente

3.2 Lematizare

Lematizatorul este una dintre cele mai complexe componente ale proiectului, implementând:

Verb	Forme	Lema
a hi (a fi)	hiu, eshti, easti, eara, fu, fura	hiu
a avea	am, ai, are, avea, aveam, avura	am
a fac (a face)	fac, fatse,featse,featsira	fac
a dzac (a zice)	dzac, dzatse, dzatsea	dzac
a vrea	vrea, va, vor, vrura	vrea
a vedu (a vedea)	vedu, vedz, vidzu, vidzura	vedu
a yinu (a veni)	yin, yine, yinea, vinje	yinu
a ljau (a lua)	ljau, ljea, lo, loara	ljau

Tabela 1: Exemple de verbe neregulate lematizate

3.2.1 Lematizare Verbală

Am construit un dicționar complet de forme verbale pentru verbele neregulate cele mai frecvente:

Au fost implementate **peste 300 de forme verbale** pentru 30+ verbe neregulate.

3.2.2 Lematizare Nominală

Lematizatorul elimină articolele definite atașate:

Cu articol	Lema	Traducere
ficiarlu	ficiar	baiat
amiralu	amira	imparat
vulpea	vulpe	vulpe
calea	cale	drum
ocljilji	oclji	ochi (pl.)

Tabela 2: Exemple de lematizare nominală

3.2.3 Reguli cu Sufixe

Pentru formele regulate, am implementat reguli bazate pe sufixe:

- Eliminare articole: **-lu, -a, -lji, -le**
- Desinențe verbale: **-ă, -eashti, -escu**
- Plurale: **-uri, -i, -e**

3.3 Stop Words

Am compilat o listă de **163+ cuvinte funcționale** din limba aromână, incluzând:

- **Pronume:** a, aestu, aesta, aoa, eu, mini, tini, el, ea, noi, voi, elji
- **Prepoziții:** di, din, cu, cătră, după, prit, pi, pisti, tra

- **Conjuncții:** că, sh, ma, shi, ama, cama, ca
- **Articole:** un, una, unu, ună, lu, a
- **Adverbe:** ashi, multi, putsă, agărsita, iara, tora, ică
- **Verbe auxiliare:** am, hiu, eshti, easti, aveam, fu, era

Această listă a fost construită prin:

1. Analiza frecvențelor în corpus-ul de antrenare
2. Consultarea dicționarelor aromâne
3. Comparație cu stop words din alte limbi române

3.4 Atribute Lexicale

3.4.1 Detectare Numerală

Am implementat funcția `like_num()` care detectează:

- **Cifre arabe:** 1, 2, 100, 1.5, 3,14
- **Numerale cardinale:** un, doi, trei, patru, tsintsi, shase, shapte, optu, noauă, dzatsi
- **Numerale de la 11-19:** unsprădzatsi, doisprădzatse, etc.
- **Zecimale:** yinghits (20), treidzăts (30), patrudzăts (40)
- **Mii și milioane:** njilji, njiliună, miliună
- **Numerale ordinale:** protlu, doilea, treilea, ultimu
- **Fractii:** giumătati (jumătate), sfărtu (sfert)

Funcția suportă ambele ortografii (Cunia și DIARO).

3.5 Conversia Ortografică

Unul dintre aspectele cele mai complexe ale limbii aromâne este existența mai multor standarde ortografice. Am implementat utilitare complete pentru conversia între:

3.5.1 Standardul DIARO (Caragiu-Marioțeanu)

Folosit în contexte academice:

- Vocale centrale: ă, â, î
- Consoane speciale: dz, l (l palatal), n (n palatal), s, t
- Exemplu: *Buna dzua! Cum esti?*

3.5.2 Standardul Cunia (Tiberiu Cunia)

Practic pentru dactilografie digitală:

- Vocala centrală unificată: ā
- Digrame: dz, lj, nj, sh, ts
- Exemplu: *Buna dzua! Cum eshti?*

3.5.3 Funcțiile de Conversie

```
1 # Conversie Cunia -> DIARO
2 def cunia_to_diaro(text: str) -> str:
3     # Conversie vocale: ā -> ā
4     # Conversie digrame: sh->s, ts->t, lj->l, nj->n
5     ...
6
7 # Conversie DIARO -> Cunia
8 def diaro_to_cunia(text: str) -> str:
9     # Conversie inversă
10    ...
11
12 # Normalizare generală
13 def normalize_text(text: str, target="cunia") -> str:
14     # Curățarea și standardizarea
15     ...
```

Acste funcții gestionează:

- Conversii de caractere speciale
- Păstrarea majusculelor
- Gestionarea variantelor Unicode multiple
- Curățarea caracterelor problematice

4 Integrare și Utilizare

4.1 Instalare

Modulul poate fi instalat în două moduri:

```
1 # Mod dezvoltare
2 git clone https://github.com/username/spacy-rup.git
3 cd spacy-rup
4 pip install -e .
```

4.2 Exemple de Utilizare

4.2.1 Pipeline de Bază

```
1 import spacy
2
3 # Creare pipeline
4 nlp = spacy.blank('rup')
5
6 # Procesare text
7 doc = nlp("Buna dzua! Mini hiu arman.")
8
9 for token in doc:
10     print(token.text, token.is_stop)
```

4.2.2 Cu Lematizare

```
1 import spacy
2
3 nlp = spacy.blank('rup')
4 nlp.add_pipe('aromanian_lemmatizer')
5
6 doc = nlp("Ficiarlufeatseunlucrumultubun.")
7
8 for token in doc:
9     print(f"{token.text[:15]} -> {token.lemma_}")
10
11 # Output:
12 # Ficiarlu      -> ficiar
13 # featse       -> fac
14 # un           -> un
15 # lucru        -> lucru
16 # multu        -> multu
17 # bun          -> bun
```

4.2.3 Cu Conversie Ortografică

```
1 from spacy_rup.orthography import cunia_to_diaro
2
3 # Normalizare inainte de procesare
4 text = "Ficiarlufeatselucrubun."
5 text_diaro = cunia_to_diaro(text)
6 # Result: "Ficiarlufeatselucrubun."
7
8 doc = nlp(text_diaro)
```

5 Date și Corpus

Proiectul include un corpus bilingv română-aromână în folderul `data/`:

- **Corpus paralel:** Texte în română (ro) și aromână (rup)
- **Variante ortografice:** Cunia (rup_cun) și Standard (rup_std)
- **Split-uri:** Train, validation, test pentru antrenarea modelelor
- **Dimensiune:** Mii de perechi de propoziții aliniate

Structura Datelor

```
data/
  Tales.train.ro          # Antrenare romana
  Tales.train.rup          # Antrenare aromania
  Tales.train.rup_cun      # Antrenare Cunia
  Tales.train.rup_std      # Antrenare DIARO
  Tales.valid.ro / .rup    # Validare
  Tales.test.ro / .rup     # Testare
  unsplit/
    corpus.rup
    corpus.rup_cun
    corpus.rup_std
```

6 Realizări și Metrici

6.1 Componentele Implementate

Componenta	Status	Detalii
Tokenizer	✓	Reguli pentru clitice aromane
Stop Words	✓	163+ cuvinte funcționale
Lex Attrs	✓	Detectare numerale
Orthography	✓	Conversie Cunia ↔ DIARO
Lemmatizer	✓	300+ forme, reguli suffix
POS Tagger	✗	Necesita date de antrenare
NER	✗	Necesita date de antrenare

Tabela 3: Status implementare componente

6.2 Statistici Lematizator

- **Verbe neregulate:** 30+ verbe
- **Forme verbale:** 300+ forme mapate
- **Reguli de suffix:** 15+ pattern-uri
- **Eliminare articole:** 8 tipuri de articole definite

6.3 Acoperire Ortografică

- Suport complet pentru standardul Cunia
- Suport complet pentru standardul DIARO
- Conversie bidirectională fără pierdere
- Gestionarea variantelor Unicode multiple

7 Provocări și Soluții

7.1 Provocări Întâmpinate

1. Standardizare ortografică multiplă

- Problema: Existența a 2-3 standarde ortografice active
- Soluție: Implementare conversie automată, suport pentru ambele

2. Verbe neregulate complexe

- Problema: Verbele frecvente au forme foarte neregulate
- Soluție: Dicționare exhaustive cu toate formele

3. Clitice și agglutinare

- Problema: Limba aromână atașează multe elemente la verb
- Soluție: Reguli de tokenizare specializate pentru prefixe/sufixe

4. Lipsa resurselor existente

- Problema: Nu există corporuri mari adnotate
- Soluție: Construire manuală a resurselor de bază

7.2 Decizii de Design

- **Standardul implicit:** Am ales Cunia ca standard implicit pentru ușurință în digitizare
- **Lematizare conservatoare:** Pentru cuvinte necunoscute, returnăm forma originală
- **Modularitate:** Fiecare componentă poate fi folosită independent
- **Compatibilitate spaCy:** Respectarea strictă a interfețelor spaCy 3.x

8 Direcții Viitoare

8.1 Îmbunătățiri Pe Termen Scurt

1. Extinderea dicționarului de leme

- Adăugare mai multe verbe neregulate
- Îmbunătățirea lematizării adjetivelor

2. Reguli morfologice mai complexe

- Pattern-uri pentru derivare
- Suport pentru compunere

3. Testare extensivă

- Unit tests pentru toate componentele
- Validare pe corpus real

8.2 Obiective Pe Termen Lung

1. POS Tagger

- Adnotare corpus cu part-of-speech tags
- Antrenare model neural

2. Named Entity Recognition

- Identificare nume persoane, locații
- Entități specifice culturii aromâne

3. Dependency Parsing

- Adnotare dependențe sintactice
- Antrenare parser

4. Word Embeddings

- Antrenare word2vec/fastText pe corpus mare
- Integrare cu modelele spaCy

5. Traducere automată

- Folosirea corpus-ului paralel ro-rup
- Antrenare modele transformer

9 Impact și Aplicații

9.1 Importanță pentru Limba Aromână

Acest proiect reprezintă:

- **Prima implementare NLP completă** pentru limba aromână în ecosistemul Python/spaCy
- **O platformă pentru cercetare** în lingvistica computațională a limbilor romanice minore
- **Un instrument educațional** pentru predarea și învățarea limbii aromâne
- **O contribuție la preservarea unei limbi** în pericol

9.2 Aplicații Potențiale

1. Educație

- Platforme de învățare asistată de calculator
- Verificare ortografică și gramaticală
- Generare exerciții automate

2. Traducere

- Sisteme de traducere aromână-română
- Dicționare digitale inteligente
- Alinierea textelor paralele

3. Cercetare

- Analiză dialectală
- Studii de variație ortografică
- Comparații cu alte limbi române

4. Digitalizare patrimoniu

- OCR pentru texte vechi aromâne
- Indexare corpus-uri literare
- Procesare folklore și texte tradiționale

10 Concluzii

10.1 Sinteza Lucrării

Acest proiect a dezvoltat cu succes un modul complet spaCy pentru limba aromână, inclusiv:

- Tokenizare specializată pentru morfologia aromână

- Lematizator funcțional cu 300+ forme verbale
- Listă de 163+ stop words
- Detectare numerale în ambele ortografi
- Conversie între standardele Cunia și DIARO
- Integrare completă cu spaCy 3.x

10.2 Contribuții Principale

1. **Tehnică:** Primul modul NLP complet pentru aromână în Python
2. **Lingvistică:** Compilarea resurselor lexicale și morfologice
3. **Practică:** Instrument funcțional pentru procesarea textelor aromâne
4. **Comunitară:** Platformă open-source pentru dezvoltări viitoare

10.3 Reflecții Finale

Dezvoltarea acestui modul a evidențiat complexitatea limbii aromâne și necesitatea instrumentelor NLP pentru limbile minore. Proiectul oferă o bază solidă pentru cercetări și aplicații viitoare, contribuind la eforturile de preservare și revitalizare a unei limbi în pericol.

Codul este disponibil open-source, permitând comunității aromâne și cercetătorilor să continue dezvoltarea și îmbunătățirea acestor instrumente.

Referințe

1. spaCy - Industrial-Strength Natural Language Processing. <https://spacy.io/>
2. Ethnologue: Languages of the World - Aromanian. <https://www.ethnologue.com/language/rup>
3. Tiberiu Cunia (1997). *Dicțiunari a Limbălor Armănești*
4. Matilda Caragiu Marioțeanu (1997). *Dictionar aromân (macedo-vlah)*
5. AroTranslate Project. <https://github.com/arotranslate/AroTranslate>
6. Aromanian Language Resources. <https://github.com/senisioi/aromanian>

A Cod Sursă - Exemplu Lematizator

```

1 from typing import Optional
2
3 VERB_LEMMAS = {
4     "hiu": "hiu", "eshti": "hiu", "easti": "hiu",
5     "am": "am", "ai": "am", "are": "am",

```

```

6     "fac": "fac", "fatse": "fac", "featse": "fac",
7     # ... 300+ forme
8 }
9
10 def lemmatize_verb(word: str) -> Optional[str]:
11     """Lematizeaza verbe neregulate."""
12     word_lower = word.lower()
13     if word_lower in VERB_LEMMAS:
14         return VERB_LEMMAS[word_lower]
15     return None
16
17 def lemmatize_noun(word: str) -> str:
18     """Elimina articolele definite."""
19     if word.endswith('lu'):
20         return word[:-2]
21     elif word.endswith('a') or word.endswith('lji'):
22         return word[:-1]
23     return word

```

B Exemplu de Text Procesate

B.1 Exemplu 1: Propoziție Simplă

Input (Cunia): *Ficiorlu featse un lucru multu bun.*

Input (DIARO): *Ficiorlu featse un lucru multu bun.*

Tokenizare: [Ficiorlu] [featse] [un] [lucru] [multu] [bun] [.]

Lemmatizare: ficiar, fac, un, lucru, multu, bun, .

Stop words: un

B.2 Exemplu 2: Propoziție Complexă

Input: *Eara oară un amira cari avea doi ficiori.*

Traducere: *Era odată un împărat care avea doi fi.*

Analiză:

- Eara → hiu (verb ”a fi”, imperfect)
- oară → oară (substantiv)
- un → un (articol) [stop word]
- amira → amira (substantiv)
- cari → cari (pronume relativ) [stop word]
- avea → am (verb ”a avea”, imperfect)
- doi → doi (numeral) [like_num=True]
- ficiori → ficiar (substantiv, plural)