

# TP2 de SPL

Filipe and J.João

## Contents

1	Emoji	1
2	Sentiment analysis	1
3	Word embeddings multilingua	2
4	TMX editor	2
5	Melhorar LF-aligner	2
6	Fluxo para préprocessamento e tratamento de corpus - multidocumento	2
7	Word embeddings shell	2
8	Multi-word terms processor (NL-flex)	3

## 1 Emoji

- taxonomia
  - bandeiras, animais, emoção, objectos, paisagem
- features
  - emoção (feliz, raiva, preocupado,...) / polaridade
  - masc, feminino
- emojis compostos, apresentação preferenciais
- tokenisadores de textos com emojis
- normalização
- corpus (Ex: twitter da aula?)
- embeddings envolvendo emoji
- emoji textual " :) :-) :D "
- pares e conversores
  - emoji → português “:smile:”
  - emoji textual → emoji

## 2 Sentiment analysis

- baseada em regras
- listas
  - pal → polaridade (-101) ou prob {fea→ ...}
- texto → sent
- frases → sent
- lista de multiword elements (“deix.\* muito a desejar”, “ponto forte”)
- negadores (muito X, bem X, muito pouco X, não X, nada X, não tem X)
- Módulo parametrizado por algumas listas;
  - pal → polaridade, negadores, ...

- sent: frase|texto  $\rightarrow$  real
- diagrama de sentimento: livro(=cap\*)  $\rightarrow$  (cap  $\rightarrow$  sent)

### 3 Word embeddings multilingua

- PT, EN
- transvec module (experimentalar...)
- a partir de TMX (ver pasta TMX no git): Como?
  - tmx  $\rightarrow$  frases embricadas  $\rightarrow$  word-emb
- alguns testes a medir / validar os resultados?
  - similar(cão , dog)
  - casamentos( gato chair cat cão dog cadeira table mesa )
  - analogias?

### 4 TMX editor

- dada uma tmx (ver pasta TMX no git)
  - editor favorito
  - split dentro de unidades alinhadas # (dog, cat and bird = cão, gato e pássaro) # virar... (dog = cão), (cat = gato), (bird = pássaro)
  - corrigir alinhamentos interactivamente (Teclas programadas?, interfaces)

### 5 Melhorar LF-aligner

- baseado no Hun-aligner
- ferramentas auxiliares (html,pdf,docx  $\rightarrow$  txt)
- recursos
- toolkit
- versão instalável! (pip install lf-aligner?)

### 6 Fluxo para préprocessamento e tratamento de corpus - multidoc- umento

- toolkit para préprocessar texto antes de
- spacy
  - (se pretendido) (V  $\rightarrow$  Lema)
  - (se pretendido) NER
  - (se pretendido) remoção de pontuação, stopword, numeros
- aglutinação/marcação de termos multipalavra
- ... domain specific language para definição das reescritas de transformanção (?)
- restauro de acentuações, capitalização (?)
- spell-checking ( módulo spacy, hunspell?)
- outras...
- ... Wemb

### 7 Word embeddings shell

gato : mamífero :: galinhas : ?

- importar algebra de gensim
  - similar
  - most similar
  - doesnt match
  - n\_similar
- funções “nossas”
- calculadora

## 8 Multi-word terms processor (NL-flex)

DSL sistema de reescrita (padrão → acção python ou texto de substituição)

```
pequeno-almoço | café da manhã → refeicao_PA  
(\w+), (\w+) e (\w+) ==> insere_irmão( \1,\2,\3)
```