

Treball de Fi de Grau

Grau en Enginyeria en Tecnologies Industrials

Desenvolupament d'una aplicació web per a NLP

ANNEXOS

Autor: Max Montoliu Torruella
Director: Alexandre Perera Lluna
Convocatòria: Juliol 2020



Escola Tècnica Superior
d'Enginyeria Industrial de Barcelona





ANNEX A

```

1. def clean_frase(line, stop=False, inp= False):
2.
3.
4. # cle a n the terms in di c by :
5. # - lower a l l l e t t e r s
6. # - conve r t numbers in l e t t e r
7. # - s e p a r a t e words divided by a -
8. # - elimi n a t e punc tu a tion and o the r non p ri n t a bl e c h a r a c t
   e r s
9. # - i f Stop = True , elimi n a t e s stopwords
10. # - i f inp = t ru e i t a l s o sep a r a t e d coupled words (w1-w2) t h a t a
   re no t in
11. # the l i s t coupled terms and i t j o i n the one t h a t a re
12.
13. import re
14. import string
15. import nltk
16.
17. from unicodedata import normalize
18. from num2words import num2words
19.
20. cleaned = []
21. # prepare regez for char filtering
22. re_print = re.compile('[^%s]' % re.escape(string.printable))
23.
24. remove = string.punctuation
25. remove = remove.replace("-", "") #don't remove hyphens
26. table = str.maketrans(remove, ' '*len(remove))
27. #normalize unicode characters (no vec que fa m'ho puc estalviar al haver aco
   nseguit lines ja oi??)
28. line = normalize('NFD', line).encode('ascii','ignore')
29. line = line.decode('UTF-8')
30.
31. #convert to lowercase
32. line = line.lower()
33.
34.
35. #translate ordinal number (e.g: 1st ---> first)
36. re_results = re.findall('(\d+(st|nd|rd|th))',line)
37. if re_results:
38.     for enitre_result, suffix in re_results:
39.         num = int(enitre_result[:-2])
40.         tmp = num2words(num, ordinal=True)
41.         tmp += ' '
42.         line = line.replace(enitre_result, tmp)
43. #convert numb-numb couples number to string (e.g.: 2-3 --> two three)
44. re_results = re.findall("\d[\d]*\d[\d]*", line)
45. if re_results:
46.     for enitre_result in re_results:
47.         re_results2 = re.findall('\d[\d]*',enitre_result)
48.         string_tmp = ' '
49.         for num in re_results2:
50.             tmp = num2words(float(num), ordinal = True)
51.             tmp += ' '
52.             string_tmp += tmp
53.         line = re.sub(enitre_result, string_tmp, line)
54.
55. #convert all other numbers to string
56. re_results = re.findall(" [-+]?[.]?[\d]+(?:,[\d\.\d]*)*[\.\d]?(?:[eE][[-
   +]?[d+])? ",line)
57. if re_results:
58.     re_results = sorted(re_results, key= len)
59.     for enitre_result in reversed(re_results):
60.         try:
61.             num = float(enitre_result)
62.             tmp = num2words(num)

```

```
1. def generar_vocab(corpus, freq_min = 1):
2.     '''
3.     Input:
4.     - corpus: llista de llistes per cada terme
5.     - freq_min
6.     Tokenitzo les els termes del corpus. Em carrego els menys freqüents
7.     '''
8.     print("Freqüència mínima del diccionari:", freq_min)
9.     freq_dict = dict()
10.    for frase in corpus:
11.        for paraula in frase:
12.            if paraula in freq_dict:
13.                freq_dict[paraula] += 1
14.            else:
15.                freq_dict[paraula] = 1
16.
17.    tokens_freq = [tokens for tokens in freq_dict if freq_dict[tokens] >= freq_min]
18.    idx = range(2, len(tokens_freq) + 2) # el 0 serà per paraules desconegudes
19.    vocab = dict(zip(tokens_freq, idx))
20.    vocab
21.    vocab['pad'] = 0
22.    vocab['desc'] = 1
23.    return freq_dict, vocab
```

```
1. def genera_tokens_processats(Terme, clean_frase = clean_frase):
2.     '''
3.     A partir d'un terme de l'ontologia genera un text no processat amb la informació necessària:
4.     * Nom
5.     * Definició
6.     * Sinònims
7.     * Pares
8.     Input: Terme
9.     Output: llista de les paraules preprocessades
10.    '''
11.
12.    it1 = "" if not Terme.synonyms else map(lambda Synonym: Synonym.description, Terme.synonyms)
13.    it2 = "" if not Terme.superclasses(1) else map(lambda Pare: Pare.name, Terme.superclasses(1))
14.    frase = " ".join([" " if not Terme.name else Terme.name, " " if not Terme.definition else Terme.definition, " ".join(it1), " ".join(it2)])
15.    frase = " ".join([" " if not Terme.name else Terme.name, " " if not Terme.definition else Terme.definition])
16.    frase_pre = clean_frase(frase)
17.    return frase_pre.split()
```

```

1. def genera_train_i_vocab():
2.     from pronto import Ontology
3.     import random
4.     import numpy as np
5.     from tensorflow.keras.preprocessing.sequence import pad_sequences
6.
7.     ls_ids = ['HP:0000707', 'HP:0000924', 'HP:0040064']
8.     X_y = []
9.     ms = Ontology("hp.owl")
10.    for i in range(len(ls_ids)):
11.        n = ms[ls_ids[i]]
12.        X_y += [(i,genera_tokens_processats(node)) for node in n.subclasses()]
13.
14.    y, corpus = list(zip(*X_y))## Aquí s'hauria d'aleatorietzar
15.    max_tokens = get_max(corpus)
16.
17.    print('Corpus:',len(corpus))
18.    print('Frase més llarga (num):',max_tokens)
19.
20.    freq,vocabulari = generar_vocab(corpus,freq_min = 1)
21.    desa_vocabulari(vocabulari,path='./')
22.
23.    y_train = np.array(y)
24.    print('Mida y_train:',y_train.shape)
25.    corpus_tokenitzat = list(map( lambda terme: [vocabulari[i] for i in terme
26.        if i in vocabulari.keys() ],corpus ))
27.    X_train = pad_sequences(corpus_tokenitzat,maxlen=max_tokens,padding='post',t
28.        runcating='post')
29.    print('Mida X_train:',X_train.shape, '-numtokens i max_tokens')
30.
31.    return X_train, y_train

```

```

1. def model2(num_tokens=8501,max_tokens=119):
2.     from tensorflow.keras.layers import (GRU, Dense, Dropout, Embedding, Flatten
3.     ,
4.         Input, Multiply, Permute, RepeatVector,
5.         Softmax,Embedding)
6.     from tensorflow.keras.models import Sequential
7.     from tensorflow.keras.utils import to_categorical
8.     from tensorflow.keras.optimizers import Adam
9.     .....
10.    definem model
11.    ...
12.    model = Sequential()
13.    model.add(Embedding(input_dim= num_tokens,
14.        output_dim = EMBEDDING_SIZE,
15.        input_length = max_tokens,
16.        name='layer_embedding'))
17.    model.add(GRU(units=16, name = 'gru_1',return_sequences=True))
18.    model.add(GRU(units=8, name='gru_2',return_sequences=True))
19.    model.add(GRU(units=4,name='gru_3'))
20.    model.add(Dense(units = 3,activation='softmax',name="dense_1"))
21.
22.    print(model.summary())
23.    return model

```

```
1. def model1(num_tokens=8501,max_tokens=119):
2.     from tensorflow.keras.layers import (LSTM, Dense, Dropout, Embedding, Flatten, Softmax,Embedding)
3.     from tensorflow.keras.models import Sequential
4.
5.     EMBEDDING_SIZE = 8
6.
7.
8.     #si s'executa més d'un cop podem els layers s'hauria de sistematitzar
9.     model = Sequential()
10.    model.add(Embedding(input_dim= num_tokens,
11.                        output_dim = EMBEDDING_SIZE,
12.                        input_length = max_tokens,
13.                        name='layer_embedding'))#l'embedding hauria de sortir fora és veu clar amb lo de les GRU
14.    model.add(LSTM(units=16, name = 'lstm_1',return_sequences=True))
15.    model.add(LSTM(units=8, name = 'lstm_2',return_sequences=True))
16.    model.add(LSTM(units=4, name = 'lstm_3'))
17.    model.add(Dense(units = 3,activation='softmax',name="dense_1"))
18.
19.    return model
```

```
1. def desa_vocabulari(vocabulari,path='./'):
2.     import os
3.     import csv
4.
5.     path_fitxer = os.path.join(path,'vocabulari.csv')
6.     f = open(path_fitxer, 'w')
7.
8.     with f:
9.         #fnames = ['first_name', 'last_name'], fieldnames=fname
10.        writer = csv.writer(f)
11.        for i in vocabulari.keys():
12.            writer.writerow([i, vocabulari[i] ])
13.        print('Vocabulari desat a:',path_fitxer)
```

```
1. def export_model(model,path='./model2'):
2.     import tensorflowjs as tfjs
3.     tfjs.converters.save_keras_model(
4.         model,
5.         path# './model2'
6.     )
```

Resultats de l'entrenament:

Model amb LSTM

```

Epoch 1/15
213/213 [=====] - 26s 123ms/step - loss: 1.0009 - accuracy: 0.5354 - val_loss: 1.8799 - val_accuracy: 0.0000e+00
Epoch 2/15
213/213 [=====] - 26s 124ms/step - loss: 0.9866 - accuracy: 0.5354 - val_loss: 1.9221 - val_accuracy: 0.0000e+00
Epoch 3/15
213/213 [=====] - 26s 122ms/step - loss: 0.9868 - accuracy: 0.5354 - val_loss: 1.8668 - val_accuracy: 0.0000e+00
Epoch 4/15
213/213 [=====] - 27s 127ms/step - loss: 0.9869 - accuracy: 0.5354 - val_loss: 1.8599 - val_accuracy: 0.0000e+00
Epoch 5/15
213/213 [=====] - 28s 132ms/step - loss: 0.9867 - accuracy: 0.5354 - val_loss: 1.8772 - val_accuracy: 0.0000e+00
Epoch 6/15
213/213 [=====] - 27s 127ms/step - loss: 0.9864 - accuracy: 0.5354 - val_loss: 1.8201 - val_accuracy: 0.0000e+00
Epoch 7/15
213/213 [=====] - 27s 128ms/step - loss: 0.9865 - accuracy: 0.5354 - val_loss: 1.9159 - val_accuracy: 0.0000e+00
Epoch 8/15
213/213 [=====] - 28s 130ms/step - loss: 0.9866 - accuracy: 0.5354 - val_loss: 1.8776 - val_accuracy: 0.0000e+00
Epoch 9/15
213/213 [=====] - 27s 128ms/step - loss: 0.9862 - accuracy: 0.5354 - val_loss: 1.8337 - val_accuracy: 0.0000e+00
Epoch 10/15
213/213 [=====] - 28s 130ms/step - loss: 0.9866 - accuracy: 0.5354 - val_loss: 1.8510 - val_accuracy: 0.0000e+00
Epoch 11/15
213/213 [=====] - 28s 130ms/step - loss: 0.9866 - accuracy: 0.5354 - val_loss: 1.8738 - val_accuracy: 0.0000e+00
Epoch 12/15
213/213 [=====] - 28s 132ms/step - loss: 0.9863 - accuracy: 0.5354 - val_loss: 1.8509 - val_accuracy: 0.0000e+00
Epoch 13/15
213/213 [=====] - 27s 127ms/step - loss: 0.9864 - accuracy: 0.5354 - val_loss: 1.8618 - val_accuracy: 0.0000e+00
Epoch 14/15
213/213 [=====] - 27s 128ms/step - loss: 0.9862 - accuracy: 0.5354 - val_loss: 1.8700 - val_accuracy: 0.0000e+00
Epoch 15/15
213/213 [=====] - 28s 131ms/step - loss: 0.9865 - accuracy: 0.5354 - val_loss: 1.8633 - val_accuracy: 0.0000e+00

```

Model amb GRU's

```

Epoch 1/15
213/213 [=====] - 32s 152ms/step - loss: 1.0090 - accuracy: 0.5354 - val_loss: 1.7434 - val_accuracy: 0.0000e+00
Epoch 2/15
213/213 [=====] - 32s 150ms/step - loss: 0.9875 - accuracy: 0.5354 - val_loss: 1.9277 - val_accuracy: 0.0000e+00
Epoch 3/15
213/213 [=====] - 32s 151ms/step - loss: 0.9874 - accuracy: 0.5354 - val_loss: 1.8799 - val_accuracy: 0.0000e+00
Epoch 4/15
213/213 [=====] - 32s 151ms/step - loss: 0.9864 - accuracy: 0.5354 - val_loss: 1.9175 - val_accuracy: 0.0000e+00
Epoch 5/15
213/213 [=====] - 32s 150ms/step - loss: 0.9869 - accuracy: 0.5354 - val_loss: 1.8833 - val_accuracy: 0.0000e+00
Epoch 6/15
213/213 [=====] - 32s 149ms/step - loss: 0.9867 - accuracy: 0.5354 - val_loss: 1.8592 - val_accuracy: 0.0000e+00
Epoch 7/15
213/213 [=====] - 32s 150ms/step - loss: 0.9866 - accuracy: 0.5354 - val_loss: 1.8271 - val_accuracy: 0.0000e+00
Epoch 8/15
213/213 [=====] - 31s 148ms/step - loss: 0.9871 - accuracy: 0.5354 - val_loss: 1.8683 - val_accuracy: 0.0000e+00
Epoch 9/15
213/213 [=====] - 32s 151ms/step - loss: 0.9867 - accuracy: 0.5354 - val_loss: 1.8992 - val_accuracy: 0.0000e+00
Epoch 10/15
213/213 [=====] - 33s 153ms/step - loss: 0.9870 - accuracy: 0.5356 - val_loss: 1.8675 - val_accuracy: 0.0000e+00
Epoch 11/15
213/213 [=====] - 32s 152ms/step - loss: 0.9864 - accuracy: 0.5356 - val_loss: 1.8513 - val_accuracy: 0.0000e+00
Epoch 12/15
213/213 [=====] - 32s 151ms/step - loss: 0.9860 - accuracy: 0.5356 - val_loss: 1.8368 - val_accuracy: 0.0000e+00
Epoch 13/15
213/213 [=====] - 32s 150ms/step - loss: 0.9473 - accuracy: 0.5768 - val_loss: 2.0229 - val_accuracy: 0.0000e+00
Epoch 14/15
213/213 [=====] - 32s 149ms/step - loss: 0.5461 - accuracy: 0.8118 - val_loss: 1.7481 - val_accuracy: 0.0000e+00
Epoch 15/15
213/213 [=====] - 32s 151ms/step - loss: 0.4523 - accuracy: 0.8304 - val_loss: 1.7023 - val_accuracy: 0.0000e+00

```


Carrega.js

```
1. // el carrego així per poder jugar-
   hi des del terminal d'inspecció
2.   async function carregaModel() {
3.     const MODEL_URL = 'http://127.0.0.1:8887/m
   odeljs/model.json';
4.     model = await tf.loadLayersModel(MODEL_URL
   );
5.     console.log(model.summary());
6.     return model
7.   };
8.   async function carregaVocabulari()
9.   {
10.    await $.ajax({
11.      url: 'http://127.0.0.1:8887/vocabulari.c
   sv',
12.      dataType: 'text',
13.      crossDomain : true}).done(success);
14.    };
15.    function success(data) {
16.      var wd_idx = new Object();
17.      lst = data.split(/\r?\n|\r/);
18.      for (var i = 0; i < lst.length; i++) {
19.        key = (lst[i]).split(',')[0]
20.        value = (lst[i]).split(',')[1]
21.
22.        if (key == "")
23.          continue
24.        wd_idx[key] = parseInt(value)
25.      }
26.      word_index = wd_idx
27.    }
28.
29.    function process(txt){
30.      txt= txt.replace(/^[a-zA-Z0-
   9\s]/, '')// s'ha de valorar
31.      txt= txt.trim().split(/\s+/)
32.      for (let i = 0; i < txt.length; i++)
33.        txt[i] = txt[i].toLowerCase();
34.      return txt// ja tenim l'array
35.    }
36.    const MAX_TOKENS = 119;
37.    function tokensaSeq(txt){
38.      txt= process(txt);
39.      let sequence= [];
40.      txt.forEach(function(word){
41.        let id = word_index[word];
42.        if (id == undefined){
43.          console.log(id);//per si incloem u
   nknown més endavant
44.          sequence.push(0)
45.        }else{
46.          sequence.push(id);
47.        }
48.        //afegim el padding
49.        if (sequence.length < MAX_TOKENS){
50.          let pad_array = Array(MAX_TOKENS -
   sequence.length)
51.          pad_array.fill(0);
52.          sequence = [...sequence,...pad_array];
53.        }
54.        return sequence
55.      }
56.
57.    async function pre_predict(){
58.      carregaVocabulari()
```

```

1. <!DOCTYPE html>
2. <html lang="en">
3. <head>
4.   <script src="https://cdn.jsdelivr.net/npm/@tensorflow/tfjs@latest"></scrip
t>
5.   <script src="https://cdn.jsdelivr.net/gh/nicolaspanel/numjs@0.15.1/dist/nu
mjs.min.js"></script>
6.   <script src="https://ajax.googleapis.com/ajax/libs/jquery/3.3.1/jquery.min
.js"></script>
7.   <script src="./carrega.js"></script>
8.   <meta charset="UTF-8">
9.   <title>Model LSTM</title>
10. </head>
11. <body>
12.
13.   <form id="form" onsubmit="return false;">
14.     <input style="position:absolute; top:20%; left:5%; width:70%; line-
height: 4em;" type="text" id="userInput" />
15.     <input style="position:absolute; top:35%; left:5%; width:70%;" type="s
ubmit" onclick="predict();" />
16.   </form>
17.   <div id="output1"></div>
18.   <div id="output2"></div>
19. </body>
20. </html>
21. <script >
22.
23. //prepararem model i diccionari
24. let word_index;
25. let model;
26. var t0 = performance.now()
27. pre_predict();// no és un problema de moment
28. var t1 = performance.now()
29. console.log("Tarda: " + (t1 - t0) + " mil·lisegons.")
30. alert("Ja s'ha carregat allò necessari per ajudar-
lo. Expliqui'ns que passa")
31.
32.
33. let pred;
34.
35. async function predict() {
36.   let tp1= performance.now()
37.   txt = document.getElementById("userInput").value;
38.   seq = tokensaSeq(txt)
39.   console.log("seqüència input:")
40.   console.log(seq);
41.   let iTensor = tf.tensor1d(seq, dtype = 'int32').expandDims(0);
42.   pred = model.predict(iTensor)
43.   document.getElementById("output1").innerText = hposMespropables(pred);
44.   let tp2 = performance.now()
45.   alert("Input tractat: " + txt+ " \nPredicció efectuada en: "+(tp2-
tp1)+" mil·lisegons")
46. }
47. let hpoArray = ['HP:0000707', 'HP:0000924', 'HP:0040064'];
48.
49. function hpoMesProbable(predictions){
50.   const predResult = predictions.dataSync();
51.   const idx = predResult.indexOf(Math.max(...predResult));//pot fer-
se amb tf segur
52.   return hpoArray[idx]
53. }

```

```
54. async function hposMespropables(predictions){
55.   const soglia = 0;
56.   const probabilitats = predictions.dataSync();
57.   const rati = tf.tensor1d([soglia])
58.   const mask = tf.squeeze(predictions.greater(rati));
59.   mask.print()
60.
61.   idxS = await tf.booleanMaskAsync(tf.range(0,hpoArray.length,1,dtype='int3
2'),mask);
62.   const outp = document.getElementById('output1');//apunto al div o2
63.   outp.innerHTML = "<br>Amb prob superior a :"+ soglia;
64.   for(const i of idxS.dataSync()){
65.     outp.innerHTML += "<br>"+hpoArray[i] + " : "+ probabilitats[i];
66.   }
67.
68.
69. }
70. </script>
```

Annex B

```

1. def genera_jsons_vega(ls_id_nom)
2.     """
3.     Generar els fitxers json necessaris per a la visualització del graf radial d
4.     e vega
5.     Input: Llista de tuples amb id i noms
6.     Output: Fitxers jsons amb nom de l'HPO que representent
7.     """
8.     import json
9.     ls_css_files = []
10.    for id,nom in ls_id_nom:
11.        nom = 'HP'+id[3:]+'.json'
12.        ls_css_files.append(nom)
13.        data_id = Graf_radial(ms,id)
14.        with open(nom, 'w') as fp:
15.            print("escrivint fitxer {}".format(nom))
16.            json.dump(data_id, fp,indent=4)

```

```

1. def Graf_radial(ms,id):
2.     import json
3.     import urllib.request
4.     from altair import VegaLite
5.
6.     instruments = set(ms[id].subclasses())
7.     data = []
8.     for term in instruments:
9.         value = {"id": int(term.id[3:]), "name":term.id,"desc":term.name}
10.        parents = instruments.intersection(term.relationships.get(ms['is_a'],set()
11.        ))
12.        if parents:
13.            value['parent']=int(parents.pop().id[3:])
14.            data.append(value)
15.
16.        # Let's use the Vega radial tree example as a basis of the visualization
17.        view = json.load(urllib.request.urlopen("https://vega.github.io/vega/example
18.        s/radial-tree-layout.vg.json"))
19.        # First replace the default data with our own
20.        print(view['data'][0].pop('url'))
21.        view['data'][0]['values'] = data
22.        view['marks'][1]['encode']['enter']['tooltip'] = {"signal": "datum.desc"}
23.
24.        view['signals'][4]['value'] = 'cluster'
25.        # Render the clustered tree
26.        return view

```

```
1. def terme2json(id):
2.     """
3.     Funció per a generar l'estructura json necessària per al dendograma
4.     """
5.     from pronto import Ontology
6.     ms = Ontology("hp.owl")#carrego l'ontologia això només cal fer-ho un cop
7.     instruments = set(ms[id].subclasses())
8.     data=[]
9.     for term in instruments:
10.         value = [ int(term.id[3:]),term.name]
11.         parents = instruments.intersection(term.relationships.get(ms['is_a'],set()
12.         ))
13.         if parents:#no és el node arrel
14.             value.append(int(parents.pop().id[3:]))
15.         else:
16.             value.append('')
17.         data.append(value)
18.     d = {"data":data}
19.     return data
```

ANNEX C

```

1.  {
2.    "$id": "https://example.com/input.schema.json",
3.    "$schema": "https://json-schema.org/draft/2019-09/schema",
4.    "title": "Predictor Output",
5.    "type": "object",
6.    "properties": {
7.      "data": {
8.        "type": "array",
9.        "items": [
10.         {
11.           "type": "object",
12.           "properties": {
13.             "id": {
14.               "type": "integer",
15.               "description": "Numeric index of the prediction."
16.             },
17.             "data": {
18.               "type": "array",
19.               "description": "The data associated with the prediction subgraph
20. .",
21.               "items": [
22.                 {
23.                   "type": "object",
24.                   "properties": {
25.                     "name": {
26.                       "type": "string",
27.                       "description": "Name of the parent (strictly speaking ch
28. ildren) node of the predicted node."
29.                     },
30.                     "children": {
31.                       "type": "array",
32.                       "description": "Name of the children (strictly speaking
33. parent) nodes of the top-level node.",
34.                       "items": [
35.                         {
36.                           "type": "object",
37.                           "properties": {
38.                             "is_predicted": {
39.                               "type": "boolean",
40.                               "description": "if True this is the predicted no
41. de."
42.                             },
43.                             "name": {
44.                               "type": "string",
45.                               "description": "Name of the node(s) at the secon
46. d level (including predicted)."
47.                             },
48.                             "children": {
49.                               "type": "array",
50.                               "description": "Name of the children (strictly s
51. peaking parent) nodes of the second-level nodes (including predicted)."
52.                             }
53.                           },
54.                           "required": [ "is_predicted", "name" ]
55.                         }
56.                       ],
57.                       "minItems": 1,
58.                       "uniqueItems": true
59.                     }
60.                   },
61.                   "required": [ "name", "children" ]
62.                 }
63.               ],
64.               "minItems": 1,
65.               "uniqueItems": true
66.             }
67.           }
68.         }
69.       ]
70.     }

```

```
1. {
2.   "$id": "https://example.com/input.schema.json",
3.   "$schema": "https://json-schema.org/draft/2019-09/schema",
4.   "title": "Predictor Input",
5.   "type": "object",
6.   "properties": {
7.     "content": {
8.       "type": "string",
9.       "description": "The content introduced in the text box that needs to be translated.",
10.      "minLength": 2
11.    },
12.    "limit": {
13.      "type": "integer",
14.      "description": "Length of the returned list.",
15.      "minimum": 1
16.    }
17.  },
18.  "required": [ "content", "limit" ]
19. }
```