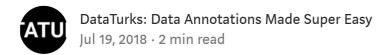
Convert Stanford CoreNLP training data to Dataturks NER JSON output



Generate a convenient JSON format for shifting from Stanford CoreNLP to SpaCy or other NER modules

Shameless plugin: We are a data annotation platform to make it super easy for you to build ML datasets. Just upload data, invite your team and build datasets super quick. Check us out!

Best online platform for your ML data annotation needs.

Just upload your data, invite your team members and start tagging. The best way to tag training/evaluation data for...

www.dataturks.com

. . .

A simple script to create dataset in Dataturks output format for NER, which can be very conveniently converted to any other training formats for NER API's like SpaCy, etc.

Dataturks NER output is very close to the format used by Spacy, just that Spacy used Python tuples which are not supported by JSON standard. So if you have your data in Stanford CoreNLP format and wish to check out the performance of SpaCy NER mode on your data, you can easily convert Stanford NER data to Dataturks output JSON format, which can be subsequently processed to SpaCy format.

We assume every word is assigned only one entity label in the script below.

Input File: Stanford CoreNLP NER Training Data.tsv File as shown below:

	0	
bought	0	
a	0	
hp	Brand	
spectre	ModelName	
x360	ModelName	
	0	

NOTE: Each training sample in Stanford NER format are separated by a newline.

Here's a script to convert Stanford CoreNLP format to Dataturks JSON format:

```
import json
    import logging
    import sys
    Creates NER training data in Dataturks format from Stanford Core NLP format.
              NOTE : This Function assumes that each token in the input fie is assigned only one ]
              Outputs the dataturks training data json format which can be used for conversion to
8
9
    10
    def Stanford_to_dataturks_format(path_to_stanford_ner_data,output_json_file_path,unknown_label)
11
12
       try:
              f=open(path to stanford ner data, 'r')#input file in Stanford NER format
13
              fp=open(output_json_file_path, 'w')#output format in Dtaturks JSON File
              data dict={}
              annotations=[]
              label dict={}
17
              s=''
18
19
              start=0
              for line in f:
                 if line!='\n':
21
22
                     word,entity=line.split('\t')
                     s+=word+" "
                     entity=entity[:len(entity)-1]
                     if entity!=unknown label:
25
                        d=\{\}
26
                        d['text']=word
                        #dataturks indices are both inclusive [start, end]
                        d['start']=start
30
                        d['end']=start+len(word)-1
                        try:
```

```
label dict[entity].append(d)
                             except:
                                  label dict[entity]=[]
                                  label_dict[entity].append(d)
                         #Increment start index and accomodate the position of an added space as wel
                         start+=len(word)+1
                     else:
                         data_dict['content']=s
                         S=' '
41
                         #create groups of points having same text
42
                         label_list=[]
                         for ents in list(label_dict.keys()):
43
                             for i in range(len(label_dict[ents])):
45
                                  if(label_dict[ents][i]['text']!=''):
                                      l=[ents,label_dict[ents][i]]
46
                                      for j in range(i+1,len(label_dict[ents])):
47
                                          if(label_dict[ents][i]['text']==label_dict[ents][j]['text']
49
                                            di={}
                                            di['start']=label_dict[ents][j]['start']
                                            di['end']=label_dict[ents][j]['end']
                                            di['text']=label_dict[ents][i]['text']
                                            1.append(di)
                                            label_dict[ents][j]['text']=''
                                      label_list.append(1)
                         for entities in label_list:
57
                             label={}
                             label['label']=[entities[0]]
                             label['points']=entities[1:]
                             annotations.append(label)
                         data dict['annotation']=annotations
63
                         annotations=[]
                         #one json object for a training sample written to output file
                         json.dump(data_dict, fp)
                         fp.write('\n')
                         data_dict={}
                         start=0
                         label_dict={}
70
         except Exception as e:
             logging.exception("Unable to process file" + "\n" + "error = " + str(e))
71
             return None
73
75
     if(len(sys.argv)<2):</pre>
         print("Please provide input and output data path and the label used for tagging as unknown
77
         exit(0)
     Stanford_to_dataturks_format(sys.argv[1],sys.argv[2],sys.argv[3])
```

The output Dataturks data would look something like he following:

```
1
     {
 2
        "content":"I bought a hp spectre x360 . ",
        "annotation":[
 3
            {
4
               "points":[
                  {
6
                      "start":11,
                      "end":12,
8
                      "text":"hp"
10
                  }
11
               ],
               "label":[
12
13
                  "Brand"
               ]
14
           },
            {
               "points":[
                  {
                      "start":14,
19
                      "end":20,
20
                      "text": "spectre"
21
                  }
22
23
               ],
               "label":[
25
                  "ModelName"
               ]
26
27
           },
28
           {
               "points":[
29
                  {
                      "start":22,
                      "end":25,
                      "text":"x360"
                  }
               "label":[
36
                  "ModelName"
37
               ]
38
           }
39
40
        ]
     }
```

have multi-word entities or consecutive words with same labels as shown below:

Have	0	
You	0	
repaired	0	
the	0	
home	Category	
theater	Category	
system	0	
?	0	

In such a case, it is wise to have 'home theater' as a span in the annotations.

In that case, it is possible to combine the consecutive words with same label with an added script in the previous code as follows:

```
import json
1
    import logging
    import sys
    Creates NER training data in Dataturks format from Stanford Core NLP format.
6
              NOTE :This Function assumes that each token in the input fie is assigned only one ]
              Outputs the dataturks training data json format which can be used for conversion to
9
    10
11
    def Stanford_to_dataturks_format(path_to_stanford_ner_data,output_json_file_path,unknown_label)
12
       try:
              f=open(path_to_stanford_ner_data,'r')#input file in Stanford NER format
13
14
              fp=open(output json file path, 'w')#output format in Dtaturks JSON File
              data_dict={}
15
              annotations=[]
17
              label_dict={}
              s=''
18
              start=0
              for line in f:
                 if line!='\n':
21
                    word,entity=line.split('\t')
                    s+=word+" "
23
                    entity=entity[:len(entity)-1]
                    if entity!=unknown_label:
26
                        d=\{\}
                        d['text']=word
27
                        #dataturks indices are both inclusive [start, end]
                        d['start']=start
```

```
d['end']=start+len(word)-1
31
                              try:
32
                                  label_dict[entity].append(d)
33
                              except:
                                  label_dict[entity]=[]
                                  label_dict[entity].append(d)
                         #Increment start index and accomodate the position of an added space as wel
                         start+=len(word)+1
                     else:
                         data_dict['content']=s
                         s=''
41
42
                         #combine consecutive words with same labels
43
                         for entities in list(label dict.keys()):
                              for i in range(len(label_dict[entities])-1):
                                  if(len(list(label dict[entities][i].keys()))!=0):
46
                                      c=i+1
                                      while(c<len(label_dict[entities]) and label_dict[entities][i][</pre>
47
                                          label_dict[entities][i]['end']=label_dict[entities][c]['end
48
                                          label_dict[entities][i]['text']+=" "+label_dict[entities][
49
                                          label_dict[entities][c]={}
51
                                          c+=1
                              label_dict[entities]=list(filter(lambda a: len(list(a.keys()))!=0, label_dict[entities]
53
                         #create groups of points having same text
                         label_list=[]
                         for ents in list(label dict.keys()):
                              for i in range(len(label_dict[ents])):
                                  if(label_dict[ents][i]['text']!=''):
57
                                      l=[ents,label_dict[ents][i]]
59
                                      for j in range(i+1,len(label_dict[ents])):
                                          if(label_dict[ents][i]['text']==label_dict[ents][j]['text']
61
                                            di={}
                                            di['start']=label dict[ents][j]['start']
63
                                            di['end']=label dict[ents][j]['end']
                                            di['text']=label_dict[ents][i]['text']
                                            1.append(di)
65
                                            label_dict[ents][j]['text']=''
                                      label list.append(1)
67
68
                         for entities in label list:
                             label={}
71
                             label['label']=[entities[0]]
                             label['points']=entities[1:]
72
                             annotations.append(label)
                         data_dict['annotation']=annotations
                         annotations=[]
76
                         #one json object for a training sample written to output file
```

```
77
                          json.dump(data_dict, fp)
                          fp.write('\n')
                          data_dict={}
79
                          start=0
                          label_dict={}
         except Exception as e:
             logging.exception("Unable to process file" + "\n" + "error = " + str(e))
             return None
84
85
86
     if(len(sys.argv)<2):</pre>
87
         print("Please provide input and output data path and the label used for tagging as unknown
88
         exit(0)
89
```

TOHOW THE TOHOWING CHICHAI.

Creates NER training data in Spacy format from JSON downloaded from Dataturks.

Dataturks NER output is very close to the format used by Spacy, just that Spacy used Python tuples which are not...

dataturks.com

If you have any queries or suggestions, I would love to hear about it. Please write to me at abhishek.narayanan@dataturks.com.

Shameless plugin: We are a data annotation platform to make it super easy for you to build ML datasets. Just upload data, invite your team and build datasets super quick. **Check us out!**

Best online platform for your ML data annotation needs.

Just upload your data, invite your team members and start tagging. The best way to tag training/evaluation data for...

www.dataturks.com

Get the Medium app



