# Convert Dataturks NER JSON ouptut to Spacy training data.

A simple script to create dataset in Spacy NER format.

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, hence just use the below function to convert Dataturks JSON to Spacy training data.

## Dataturks to Spacy

```
1    ######################################  NOTE   #######################################
2    #
3    #            Creates NER training data in Spacy format from JSON downloaded from Dataturks.
4    #
5    #            Outputs the Spacy training data which can be used for Spacy training.
6    #
7    #####################################################################################
8    def convert_dataturks_to_spacy(dataturks_JSON_FilePath):
9        try:
10           training_data = []
11           lines=[]
12           with open(dataturks_JSON_FilePath, 'r') as f:
13               lines = f.readlines()
14
15           for line in lines:
16               data = json.loads(line)
17               text = data['content']
18               entities = []
19               for annotation in data['annotation']:
20                   #only a single point in text annotation.
21                   point = annotation['points'][0]
22                   labels = annotation['label']
23                   # handle both list of labels or a single label.
24                   if not isinstance(labels, list):
25                       labels = [labels]
26
27                   for label in labels:
28                       #dataturks indices are both inclusive [start, end] but spacy is not [start,
29                       entities.append((point['start'], point['end'] + 1 ,label))
30
31
32               training_data.append((text, {"entities" : entities}))
33
34           return training_data
35       except Exception as e:
36           logging.exception("Unable to process " + dataturks_JSON_FilePath + "\n" + "error = " +
37           return None
```

Here is a sample code on to train a Spacy model from the above data:

# Train Spacy

```python
1   import spacy
2   ################## Train Spacy NER.###########
3   def train_spacy():
4       TRAIN_DATA = convert_dataturks_to_spacy("dataturks_downloaded.json");
5       nlp = spacy.blank('en')  # create blank Language class
6       # create the built-in pipeline components and add them to the pipeline
7       # nlp.create_pipe works for built-ins that are registered with spaCy
8       if 'ner' not in nlp.pipe_names:
9           ner = nlp.create_pipe('ner')
10          nlp.add_pipe(ner, last=True)
11
12      # add labels
13      for _, annotations in TRAIN_DATA:
14          for ent in annotations.get('entities'):
15              ner.add_label(ent[2])
16
17      # get names of other pipes to disable them during training
18      other_pipes = [pipe for pipe in nlp.pipe_names if pipe != 'ner']
19      with nlp.disable_pipes(*other_pipes):  # only train NER
20          optimizer = nlp.begin_training()
21          for itn in range(1):
22              print("Statring iteration " + str(itn))
23              random.shuffle(TRAIN_DATA)
24              losses = {}
25              for text, annotations in TRAIN_DATA:
26                  nlp.update(
27                      [text],  # batch of texts
28                      [annotations],  # batch of annotations
29                      drop=0.2,  # dropout - make it harder to memorise data
30                      sgd=optimizer,  # callable to update weights
31                      losses=losses)
32              print(losses)
33
34      #do prediction
35      doc = nlp("Samsing mobiles below $100")
36      print ("Entities= " + str([""" + str(ent.text) + "_" + str(ent.label_) for ent in doc.ents])
```

You can also download a full python script to generate the Spacy training data and store it as a Pickle file. Download from: GitHub (https://gist.github.com/DataTurks/97ff613967e8139e57091f9299c3a104)

## About

> Privacy Policy (/privacy.php)
> Open Data Policy (/open-data-policy.php)
> Help (/help/help.php)
> FAQs (/faq.php)

## Blog

> All posts (/blog/blog.php)
> OCR APIs comparison (/blog/compare-image-text-recognition-apis.php)
> Face reco APIs comparison. (/blog/face-verification-api-comparison.php)
> Best image moderation APIs. (/blog/image-moderation-api-comparison.php)

## How to

> Image Bounding Box. (/help/image-rectangle-bounding-box-annotation.php)
> Document Annotation. (/help/document-annotation-POS-NER.php)
> Polygon Bounding Box. (/help/image-polygon-bounding-box-annotation.php)
> POS Tagging. (/help/pos-text-annotations.php)

## Features

> Image Annotations (/features/image-bounding-box.php)
> Text Annotations (/features/document-ner-annotation.php)
> API Docs (https://docs.dataturks.com/)

## Documentation

> Export to Pascal VOC (/help/ibbx_dataturks_to_pascal_voc_format.php)

› Export in TensorFlow Format (https://medium.com/dataturks/converting-dataturks-image-classifier-tools-output-to-tensorflow-format-6f569e085bf3)
› NER in Spacy Format (/help/dataturks-ner-json-to-spacy-train.php)
› Docs (/help/help.php)

## ML Tutorial

› Introduction to ML (/blog/intro-to-machine-learning-NER-deep-dive.php)
› ML based troll filter (https://medium.com/@dataturks/using-machine-learning-to-fight-cyber-trolls-9bf0fa1c5df9)
› ML and GDPR (https://medium.com/@dataturks/how-does-gdpr-impact-machine-learning-keystrokes-pascal-voc-and-much-more-1625b8a1147b)
› TensorFlow vs Keras? (https://hackernoon.com/tensorflow-vs-keras-comparison-by-building-a-model-for-image-classification-f007f336c519)

---

(+91) 080-331-72755, +91-99010-49915, +91-88614-08222

contact@dataturks.com

**Say Hi:**

---

Copyright © Trilldata Technologies Pvt Ltd 2018