

In this notebook I will be using SPaCy to create model which will extract main points from a resume. We will train the model on almost 200 resumes.

After the model is ready, we will extract the text from a new resume and pass it to the model to get the summary.

### Importing the libraries

```
In [ ]: 1 import spacy
        2 import pickle
        3 import random
        4 import os
        5 for dirname, _, filenames in os.walk('/kaggle/input'):
        6     for filename in filenames:
        7         print(os.path.join(dirname, filename))
        8
```

### Loading the data

```
In [ ]: 1 train_data = pickle.load(open('/kaggle/input/resume-data/train_data.'
```

```
In [ ]: 1 # Let's see the data
        2 train_data[0]
```

### Training the data on the model

We will first load a blank SpaCy english model. Then we will write a function which will take the training data as the input. In the function, first we will add a ner i.e. Named Entity Recognition in the last position in the pipeline. Then we will add our custom labels in the pipeline.

```
In [ ]: 1 # Loading the blank SpaCy english model
        2 nlp = spacy.load('en')
        3
```

```

In [ ]: 1 nlp = spacy.blank('en')
        2
        3 def train_model(train_data):
        4     if 'ner' not in nlp.pipe_names:
        5         ner = nlp.create_pipe('ner')
        6         nlp.add_pipe(ner, last = True)
        7
        8     for _, annotation in train_data:
        9         for ent in annotation['entities']:
10             ner.add_label(ent[2])
11
12
13     other_pipes = [pipe for pipe in nlp.pipe_names if pipe != 'ner']
14     with nlp.disable_pipes(*other_pipes): # only train NER
15         optimizer = nlp.begin_training()
16         for itn in range(10):
17             print("Starting iteration " + str(itn))
18             random.shuffle(train_data)
19             losses = {}
20             index = 0
21             for text, annotations in train_data:
22                 try:
23                     nlp.update(
24                         [text], # batch of texts
25                         [annotations], # batch of annotations
26                         drop=0.2, # dropout - make it harder to memorize
27                         sgd=optimizer, # callable to update weights
28                         losses=losses)
29                 except Exception as e:
30                     pass
31
32             print(losses)

```

```

In [ ]: 1 # Let's train the model
        2 train_model(train_data)

```

```

In [ ]: 1 # Let's save the model for further use
        2 nlp.to_disk('nlp_model')

```

## Model testing

Let's check out how our model is performing. For this we will pass a new resume to this model.

```

In [ ]: 1 # Installing PyMuPDF for getting the text data from the resume pdf
        2 !pip install PyMuPDF

```

```
In [ ]: 1 import sys, fitz
        2
        3 fname = '/kaggle/input/resume-data/Alice Clark CV.pdf'
        4 doc = fitz.open(fname)
        5 text = ""
        6 for page in doc:
        7     text = text + str(page.getText())
        8
        9 tx = " ".join(text.split('\n')) # for removing the next line charac
       10 print(tx)
```

```
In [ ]: 1 # Now we will pass this extracted text to our model
        2 nlp_model = nlp.from_disk('/kaggle/input/resume-data/nlp_model/')
        3
        4 doc = nlp_model(tx)
        5 for ent in doc.ents:
        6     print(f'{ent.label_.upper():{30}}- {ent.text}')
```

***This is amazing. But we can make this model more accurate by training it on more data.***

And I will be constantly updating this notebook to make it better.