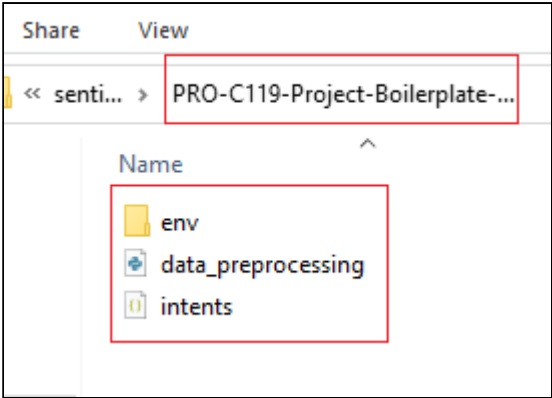


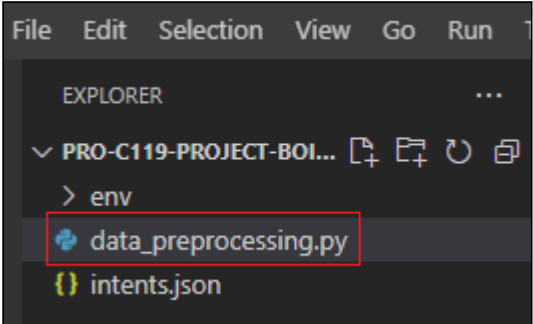




1. Open the Boilerplate [link](#) and download all the files within a new **folder** on your system.
2. Open the **command prompt**, traverse to that folder and create a python virtual environment inside it in such a way, so that the **virtual environment**, **intents.json** and **data\_preprocessing.py** files are within the same folder

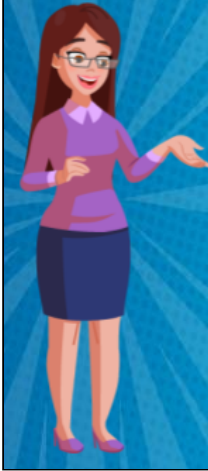


3. Activate the virtual environment and install the **nlTK** and **Tensorflow** library in it, using **pip install nlTK** and **pip install tensorflow==2.5.0**.
4. Open the **folder** in Visual Studio code, and click on the **data\_preprocessing.py** file.



Specific Tasks to complete the Project:


Step 1



Create an instance of class **PorterStemmer**, so that you can stem words.

```
# create an instance of class PorterStemmer
```

Step 2




Create a method and name it as **get\_stem\_words()**, which will accept 'list of **words** to be stemmed' and 'list of words that ought to be **ignored** as arguments and return a list of stemmed words.



```
# creating function to stem words
def get_stem_words(words, ignore_words):
    stem_words = []
    for word in words:
        # write stemming algorithm:
        ...
        Check if word is not a part of stop word:
        1) lowercase it
        2) stem it
        3) append it to stem_words list
        4) return the list
        ...
        # Add code here #
    return stem_words
```

Step 3



Write the appropriate code to preprocess the data.

```
# add the tokenized words to the words list

# add the 'tokenized word list' along with the 'tag' to pattern_word_tags_list
```

Step 4



After creating a list of stemmed words, remove duplicate words from it and sort it.

```
# Remove duplicate words from stem_words

# sort the stem_words list and classes list
```

Step 5



Write an algorithm to create Bag of Words for all patterns.

```
# Input data encoding
...
Write BOW algo :
1) take a word from stem_words list
2) check if that word is in stemmed_pattern_word
3) append 1 in BOW, otherwise append 0
...
bag.append(bag_of_words)
```



Step 6

Preserve the stem\_word and classes list using the pickle module.

```
# Convert Stem words and Classes to Python pickel file format
```

Step 7

Save, Compile and Run your code to see the list of stemmed words and training data.

Submitting the Project:

- 1. **SAVE** all the changes made to the project.
- 2. Click on "**Run**" once to check if it is working.
- 3. Open GitHub and create a repository named **Project119**.
- 4. Upload files and click **Commit Changes**.
- 5. Copy the link and submit it in the Student Dashboard Projects panel against the correct class number.

Hints:

- 1. To run the program, activate the virtual environment, traverse to the folder create and use the command **python data\_preprocessing.py [Step 7]**
- 2. In step 3, after tokenizing the patterns, add the tokenized words to the **words list** and the tokenized word to the pattern\_word\_tags \_list along with their tags, using suitable list methods.
- 3. In step 4, sort both the list of **stem\_words** and **classes**, using the appropriate list methods.