

To run the provided code, we need the following Python libraries installed:

1. PyPDF2: Used for reading PDF files.

pip install PyPDF2

2. nltk: Used for natural language processing tasks.

pip install nltk

3. docx2txt: Used for extracting text from DOCX files.

pip install docx2txt

4. python-docx: Used for creating and manipulating Word documents.

pip install docx2txt

pip install python-docx

Additionally, you need to have a compatible version of Python installed. The code is compatible with Python 3.x versions.

The provided code is a Python script that processes a resume file, extracts relevant information such as name, email, mobile number, and parsed resume content, and saves the parsed information to a Word document.

Here is a step-by-step explanation of the code:

1. The script starts by importing the required libraries, including PyPDF2, re, string, nltk, word_tokenize, stopwords, docx2txt, and Document from the python-docx library.

2. It downloads necessary resources from the NLTK library, including stopwords and punkt tokenizer.

3. The code defines a function named ``get_pdf_content`` that takes a file path as input and extracts the text content from a PDF file using the PyPDF2 library. The content is then cleaned and normalized by removing extra spaces and special characters.

4. The code defines another function named ``get_text_from_docx`` that takes a filename as input and uses the docx2txt library to extract the text content from a DOCX file.

5. The main function ``process_resume`` is defined, which takes a filename as input. It checks the file format based on the file extension and calls the appropriate function (``get_pdf_content`` or ``get_text_from_docx``) to extract the text content from the resume file.

6. The script tokenizes the resume content into individual words using the ``word_tokenize`` function from the nltk library.

7. It defines a list of punctuation marks and a set of stopwords using the `nltk.corpus.stopwords` module.

8. The code filters out stopwords and punctuation marks from the tokens using list comprehensions and creates a new list called ``filtered`` containing the cleaned tokens.

9. The script extracts the name from the filtered tokens assuming that the first name and last name are the first two tokens.

10. It extracts the email address from the resume using regular expressions (``re`` module) by matching patterns that represent an email address.

11. Similarly, it extracts the mobile number from the resume using regular expressions by matching patterns that represent a mobile number.

12. The filtered tokens are joined back into a single string representing the parsed resume content.

13. The script generates shingles, which are sequences of consecutive words, from the filtered tokens using the ``ngrams`` function from the nltk library.

14. It creates a new Word document using the ``Document`` class from the python-docx library.

15. The parsed information, including name, email, mobile number, and parsed resume content, is added to the Word document as paragraphs.

16. The Word document is saved with the filename "output.docx".

17. Finally, the script prompts the user to enter the filename or path of the resume file and calls the ``process_resume`` function to process the resume and save the parsed information to the output.docx file.

The code utilizes various libraries and modules to extract information from different file formats, clean and normalize the text content, perform tokenization and filtering, and create a Word document to store the parsed information. It demonstrates how to use Python libraries for text processing and document manipulation tasks.