```
# personality trait prediction using Logistic Regression and
parsing resume
@app.route('/prediction' , methods = ['GET' , 'POST'])
def predict():
    # get form data
    if request.method == 'POST':
        fname = request.form['firstname'].capitalize()
        lname = request.form['lastname'].capitalize()
        age = int(request.form['age'])
        gender = request.form['gender']
        email = request.form['email']
        file = request.files['resume']
        path = './static/{}'.format(file.filename)
        file.save(path)
        val1 = int(request.form['openness'])
        val2 = int(request.form['neuroticism'])
        val3 = int(request.form['conscientiousness'])
        val4 = int(request.form['agreeableness'])
        val5 = int(request.form['extraversion'])
      # model prediction
        df = pd.read csv(r'static\trainDataset.csv')
        le = LabelEncoder()
        df['Gender'] = le.fit transform(df['Gender'])
        x_train = df.iloc[:, :-1].to_numpy()
        y_train = df.iloc[:, -1].to_numpy()
        lreg = LogisticRegression(multi class='multinomial',
solver='newton-cg',max iter =1000)
        lreg.fit(x_train, y_train)
        if gender == 'male':
            gender = 1
        elif gender == 'female':
            gender = 0
        input = [gender, age, val1, val2, val3, val4, val5]
        pred = str(lreg.predict([input])[0]).capitalize()
```

```
# get data from the resume
# data = ResumeParser(path).get_extracted_data()

result = {'Name':fname+' '+lname , 'Age':age ,
'Email':email, 'Predicted Personality':pred}

with open('./static/result.json' , 'w') as file:
    json.dump(result , file)

return render_template('questionPage.html')
```

Certainly! Let's break down the code line by line to understand each part:

```
python

@app.route('/prediction', methods=['GET', 'POST'])
def predict():
```

This line defines a route in the Flask web application. When a user accesses the URL
 '/prediction', the function 'predict()' will be executed. It accepts both GET and POST requests.

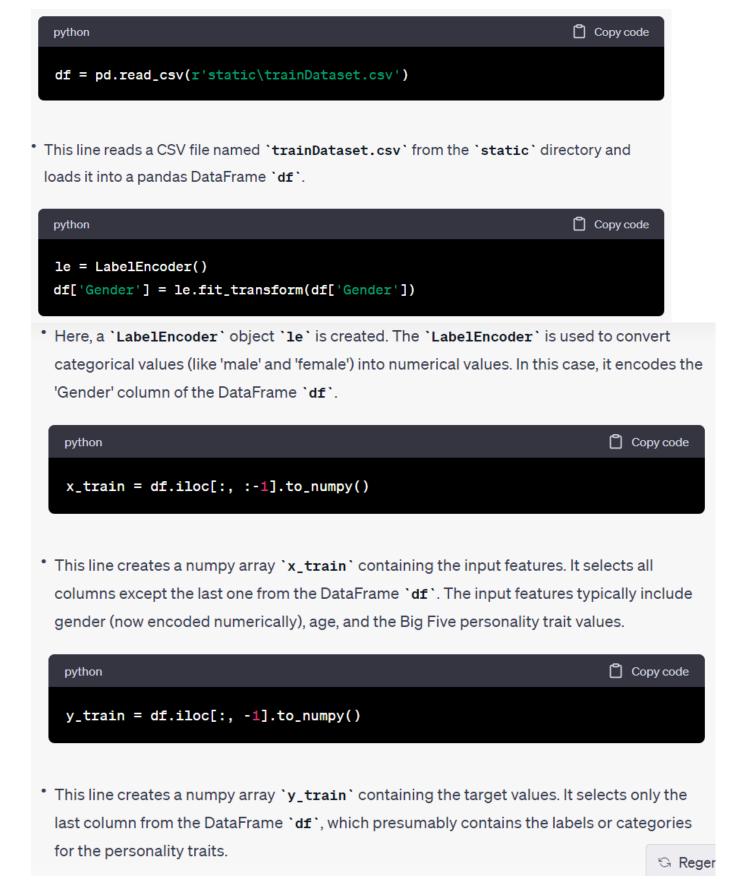
```
python

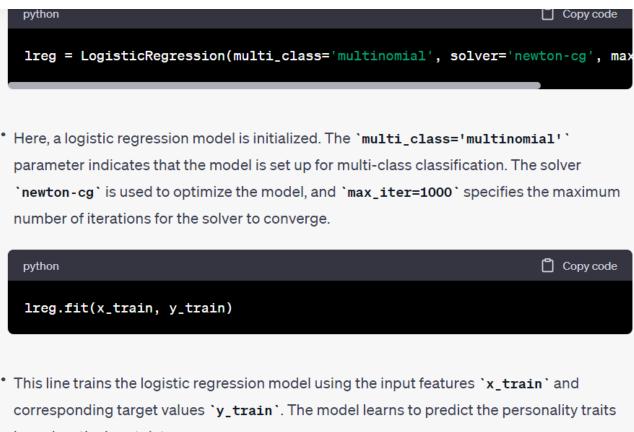
if request.method == 'POST':
```

• This line checks if the incoming request method is POST. In the context of a web form this means the form has been submitted.

```
fname = request.form['firstname'].capitalize()
lname = request.form['lastname'].capitalize()
age = int(request.form['age'])
gender = request.form['gender']
email = request.form['email']
file = request.files['resume']
path = './static/{}'.format(file.filename)
file.save(path)
val1 = int(request.form['openness'])
val2 = int(request.form['neuroticism'])
val3 = int(request.form['conscientiousness'])
val4 = int(request.form['agreeableness'])
val5 = int(request.form['extraversion'])
```

* These lines extract the form data submitted via POST request. The data includes first name, last name, age, gender, email, resume file, and values for the Big Five personality traits. The resume file is saved on the server in the `./static/` directory.





based on the input data.

In summary, these lines of code read the training data from a CSV file, encode categorical variables, prepare the input features and target values, and then train a logistic regression model for predicting personality traits based on the provided data.

Model Initialization:

```
python

lreg = LogisticRegression(multi_class='multinomial', solver='newton-cg', max
```

- `multi_class='multinomial' `: This parameter indicates that the logistic regression model is set up for multi-class classification. In the context of personality traits prediction, this means that the model can predict multiple categories or classes of personality traits.
- * `solver='newton-cg' `: The solver is the algorithm used for optimization. In this case, `newton-cg` is one of the solvers supported by scikit-learn for logistic regression. It's an optimization algorithm that can handle multinomial loss (used in multi-class classification). Different solvers have different algorithms for finding the optimal weights for the logistic regression model.
- `max_iter=1000`: This parameter sets the maximum number of iterations for the solver to converge. Logistic regression is an iterative optimization algorithm. The `max_iter` parameter ensures that the solver doesn't run indefinitely, especially for complex datasets.

Model Training:



- 'x_train': This represents the input features of the training data. In the context of
 personality traits prediction, these features include gender, age, and the Big Five
 personality trait values.
- 'y_train': This represents the target values or labels corresponding to the input features.
 For each set of input features in 'x_train', 'y_train' contains the corresponding personality trait category or class that the model should learn to predict.

When the `fit` method is called with the training data (`x_train` and `y_train`), the logistic regression model adjusts its weights (coefficients) during the training process. It does this iteratively, using the specified solver (`newton-cg`) and stops after the maximum number of iterations (`max_iter=1000`) or when it converges to the optimal solution.

After this training process, the `lreg` model is ready to make predictions on new, unseen data based on the patterns it learned from the training set.

Introduction:

The personality of a human plays a major role in his personal and professional life. Many organizations have also started shortlisting the candidates based on their personality as this increase efficiency of work because the person is working in what he is good at than what he is forced to do.

The project is based on identifying the personality of an individual using machine learning algorithms and **Big Five Model** also known as Five-Factor Model which classify the personality of a person in some classes based on following behavourial data-

- Open to Experience: It involves various dimensions, like imagination, sensitivity, attentiveness, preference to variety, and curiosity.
- **Conscientiousness**: This trait is used to describe the carefulness and diligence of the person. It is the quality that describes how organized and efficient a person is.
- **Extraversion**: It is the trait that describes how the best candidates can interact with people that is how good are his/her social skills.
- **Agreeableness**: It is a quality that analyses the individual behavior based on the generosity, sympathy, cooperativeness and ability to adjust with people.
- **Neuroticism**: This trait usually describes a person to have mood swings and has extreme expressive power.

Approach:

Five characteristics of different individuals commonly known as big five characteristics namely, openness, neuroticism, conscientiousness, agreeableness and extraversion are stored in a dataset along with gender and age of indivitual and used for training. Before training the model, data is preprocessed like handling missing values, data discretization, standardization etc. This preprocessed data is then used to train the model. User rates himself for different behavourial characteristics and based upon the information provided by user his/her personality is predicted using trained ML model.

Logistic Regression algorithm is used to train the model with highest accuracy of 75% because it is best to predict the categorical dependent variable/ class labels using a given set of independent variables. Since, there were more than two classes so here multinomial logistic regression is used with newton-cg which is best while working with big datasets.

Dataset Description:

Attributes description

- Gender [Male/Female]
- Age [17-28]
- Openness [1-8]
- Neuroticism [1-8]
- Conscientiousness [1-8]
- Agreeableness [1-8]

Extraversion [1-8]

Class Labels

- Extraverted
- Serious
- Responsible
- Lively
- Dependable

- **Multi-Faceted Analysis:**

- Utilizes the OCEAN model for comprehensive personality analysis: Openness, Conscientiousness, Extroversion, Agreeableness, and Neuroticism.
- Incorporates emotion analysis to gauge emotional states and responses.
- Integrates speech tone analysis for evaluating vocal cues and emotions.

- **Data Collection and Preprocessing:**

- Collects diverse and extensive datasets for training and testing purposes.
- Implements thorough data preprocessing techniques to ensure accuracy and reliability in analysis.

- **Advanced Machine Learning Models: **

- Employs state-of-the-art machine learning algorithms for personality prediction.
- Utilizes deep learning models for emotion and speech tone analysis, enhancing accuracy and granularity.

- **User-Friendly Interface:**

- Develops an intuitive and user-friendly interface for seamless user interaction.
- Enables users to input text or speech data for instant personality and emotion analysis.

- **Real-Time Analysis:**

- Supports real-time analysis of input data, allowing for instant feedback and insights.

- **Visualizations and Reports:**

- Generates visualizations and reports showcasing personality traits, emotional states, and speech tones.
- Provides detailed insights and interpretations for users to understand the analysis results.

to ensure the ethical deployment of the technology in various domains.

- **Future Enhancements:**

- Discusses potential future enhancements, such as incorporating additional psychological models or expanding the analysis to multi-modal inputs (text, speech, images).

- Considers integration with emerging technologies like natural language processing advancements or affective computing for more nuanced analysis.

```
import pandas as pd
import numpy as np
import keras
import tensorflow
from keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import
pad_sequences
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split
from keras.models import Sequential
from keras.layers import Embedding, Flatten, Dense

data = pd.read_csv("train.txt", sep=';')
data.columns = ["Text", "Emotions"]
data.head()
```

1. Imports:

- pandas is used for data manipulation and analysis.
- **numpy** provides support for arrays, matrices, and mathematical functions.
- **keras** and **tensorflow** are deep learning libraries used for building and training neural networks.
- **Tokenizer** from Keras is used for tokenizing text.
- pad sequences from TensorFlow is used for padding sequences.
- LabelEncoder from scikit-learn is used for encoding target labels.
- **train_test_split** from scikit-learn is used for splitting the data into training and testing sets.
- **Sequential**, **Embedding**, **Flatten**, and **Dense** are classes from Keras used to build neural network models.

2. Loading and Preprocessing Data:

- The code reads data from a CSV file named "train.txt" with columns "Text" and "Emotions" separated by semicolons.
- The columns are renamed to "Text" and "Emotions" for clarity

The **head()** function is used to display the first few rows of the loaded dataset, allowing you to quickly inspect the data structure and values.

```
texts = data["Text"].tolist()
labels = data["Emotions"].tolist()

# Tokenize the text data
tokenizer = Tokenizer()
tokenizer.fit_on_texts(texts)

sequences = tokenizer.texts_to_sequences(texts)
max_length = max([len(seq) for seq in sequences])
padded_sequences = pad_sequences(sequences, maxlen=max_length)

# Encode the string labels to integers
label_encoder = LabelEncoder()
labels = label_encoder.fit_transform(labels)
# One-hot encode the labels
one_hot_labels = keras.utils.to_categorical(labels)
```

text corpus.

```
sequences = tokenizer.texts_to_sequences(texts)
max_length = max([len(seq) for seq in sequences])
padded_sequences = pad_sequences(sequences, maxlen=max_length)
```

- `tokenizer.texts_to_sequences(texts)` converts the `texts` into a list of sequences.
 Each sequence is a list of integers where each integer represents the index of a word in the tokenizer's word index.
- `max_length` is calculated by finding the length of the longest sequence in the
 `sequences` list.
- `pad_sequences(sequences, maxlen=max_length)` pads all sequences in the `sequences` list to have a length of `max_length`. Sequences shorter than `max_length` are padded with zeros at the beginning. This ensures that all sequences have the same length.

At this point, **padded_sequences** contains the preprocessed text data, ready to be used for training a neural network. The data is tokenized, converted to numerical sequences, and padded to ensure uniformity in sequence length, which is essential for training machine learning models, particularly neural networks, for natural language processing tasks.

```
python

label_encoder = LabelEncoder()
labels = label_encoder.fit_transform(labels)
```

- `LabelEncoder()` initializes a label encoder object. Label encoding is a technique to convert categorical labels (emotions in this case) into numerical values.
- `label_encoder.fit_transform(labels)` encodes the string labels into integers. Each
 unique label is mapped to a unique integer.

```
python

One_hot_labels = keras.utils.to_categorical(labels)
```

• `keras.utils.to_categorical(labels)` converts the integer-encoded labels to one-hot encoded vectors. One-hot encoding is a binary representation of labels where each integer is mapped to a binary vector. For example, if there are 3 classes and a label is 2, it will be represented as [0, 0, 1] in one-hot encoding.

After these steps, `padded_sequences` contains the preprocessed text data, and `one_hot_labels` contains the corresponding one-hot encoded emotion labels. These processed data can be used for training machine learning models, especially neural networks, for tasks like text classification, sentiment analysis, or any other NLP task that involves analyzing textual data with categorical labels.

Certainly! The line you provided is splitting the preprocessed data into training and testing sets using the `train_test_split` function from scikit-learn. Let me explain in detail:

- 'padded_sequences': This variable contains the preprocessed textual data after tokenization and padding. Each row represents a sequence of integers corresponding to words in the text.
- `one_hot_labels`: This variable contains the emotion labels encoded in one-hot format.
 Each row represents a one-hot encoded vector corresponding to the emotion label of the respective text.

- `train_test_split(padded_sequences, one_hot_labels, test_size=0.2)`: This function splits the data into random train and test subsets.
 - 'padded_sequences' and 'one_hot_labels' are the arrays to be split.
 - 'test_size=0.2' specifies that 20% of the data will be used for testing, and the remaining 80% will be used for training the model.

After executing this line, the data is split into the following variables:

- * `xtrain`: The training data (padded sequences) that the model will learn from.
- 'xtest': The testing data (padded sequences) that the model will be evaluated on to assess its performance.
- * 'ytrain': The corresponding one-hot encoded emotion labels for the training data.
- 'ytest': The corresponding one-hot encoded emotion labels for the testing data.

These variables are now ready to be used for training and evaluating your neural network model for text emotion classification.

why are we doing this??

Splitting the data into training and testing sets is a fundamental practice in machine learning and deep learning for several important reasons:

1. Evaluation of Model Performance:

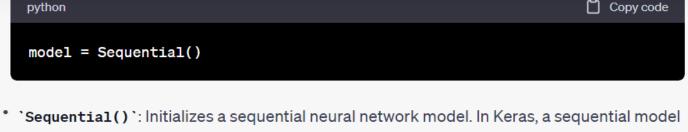
• By splitting the dataset into training and testing sets, you can train the machine learning model on one subset (training set) and evaluate its performance on another subset (testing set). This allows you to assess how well your model generalizes to new, unseen data. It helps you understand if your model has learned the underlying patterns in the data or if it has simply memorized the training data (overfitting).

2. Preventing Overfitting:

Overfitting occurs when a model learns the training data too well, including its noise and outliers. Such a model performs well on the training data but fails to generalize to new, unseen data. By evaluating the model on a separate testing set, you can detect overfitting. If the model performs significantly worse on the testing set compared to the training set, it might be overfitting.

```
model.add(Flatten())
model.add(Dense(units=128, activation="relu"))
model.add(Dense(units=len(one_hot_labels[0]),
activation="softmax"))

model.compile(optimizer="adam",
loss="categorical_crossentropy", metrics=["accuracy"])
model.fit(xtrain, ytrain, epochs=15, batch_size=32,
validation_data=(xtest, ytest))
```



Sequential(): Initializes a sequential neural network model. In Keras, a sequential model
is a linear stack of layers, where you can simply add one layer at a time.

```
python

Copy code

1.add(Embedding(input_dim=len(tokenizer.word_index) + 1, output_dim=128, input
```

Embedding Layer:

- The embedding layer is the first layer of the model.
- It is used for word embedding, a technique where words are represented as dense vectors in a continuous vector space.
- `input_dim` is set to the size of the vocabulary (number of unique words + 1). Each word
 is represented by an integer index.
- `output_dim=128`: Specifies the dimensionality of the embedding vector. Each word will be represented as a vector of 128 dimensions.
- 'input_length=max_length': Sets the length of input sequences. All input sequences
 must have the same length ('max_length'), which is achieved through padding.

```
model.add(Flatten())
```

* Flatten Layer:

- * The flatten layer is used to flatten the 2D output from the embedding layer to a 1D array.
- This is necessary because many machine learning algorithms (like fully connected layers) work with flattened data.

```
python

Copy code

model.add(Dense(units=128, activation="relu"))
```

Dense Layer:

- Dense layers are fully connected layers where each neuron is connected to every neuron in the previous layer.
- `units=128`: Specifies the number of neurons (units) in the layer. More units allow the model to learn more complex patterns in the data.
- `activation="relu" `: Applies the Rectified Linear Unit (ReLU) activation function to introduce non-linearity. ReLU is often used in hidden layers as it helps the model learn complex patterns.

```
python

model.add(Dense(units=len(one_hot_labels[0]), activation="softmax"))
```

Output Layer:

- This is the final dense layer which produces the output.
- `units=len(one_hot_labels[0])`: The number of units in this layer is equal to the number of unique emotion labels. Each unit represents a class.
- `activation="softmax" `: The softmax activation function is used in the output layer for multi-class classification problems. It converts the raw predictions into probabilities.
 Each class gets a probability score, and the class with the highest probability is the predicted class.

model.compile(optimizer="adam", loss="categorical_crossentropy", metrics=["a

Compiling the Model:

- `optimizer="adam" `: The Adam optimizer is used. It's an adaptive learning rate
 optimization algorithm that's well-suited for deep learning tasks. It adjusts the learning
 rates of each parameter during training.
- `loss="categorical_crossentropy" `: Categorical cross-entropy loss function is used. It measures the difference between the predicted class probabilities and the true class probabilities. The aim is to minimize this value during training.
- 'metrics=["accuracy"]': Accuracy is chosen as the metric to evaluate the model's performance during training.

model.fit(xtrain, ytrain, epochs=15, batch_size=32, validation_data=(xtest,

Training the Model:

- 'xtrain' and 'ytrain' are the input features (padded sequences) and target labels (one-hot encoded) for training.
- 'epochs=15': Specifies the number of times the model will be trained on the entire training dataset.
- 'batch_size=32': Sets the number of samples used in each iteration for updating the model's weights. Smaller batch sizes often result in a more accurate model but require more computation.
- 'validation_data=(xtest, ytest)': Uses the testing data for validation during training. The model's performance on the validation data is evaluated after each epoch.

In summary, the code defines a sequential neural network model for text emotion classification. It starts with an embedding layer, followed by a flatten layer and two dense layers. The model is compiled with an Adam optimizer and categorical cross-entropy loss, and it's trained for 15 epochs on the provided training data, with validation performed on the testing data. The model's goal is to learn to accurately classify emotions based on the text data.

```
model_architecture = model.to_json()
with open("model_architecture.json", "w") as json_file:
        json_file.write(model_architecture)

model.save_weights("model_weights.h5")
```

In the given code snippet, the trained neural network model's architecture is saved to a JSON file, and the model's weights are saved to an H5 file. Let's break down each line:

```
python

model_architecture = model.to_json()
```

' `model.to_json()`: This function converts the model architecture (layers, connections, and configurations) into a JSON string. It captures the structure of the neural network.

```
python

with open("model_architecture.json", "w") as json_file:
    json_file.write(model_architecture)
```

- `with open("model_architecture.json", "w") as json_file: `: This line opens a file named "model_architecture.json" in write mode (`"w"`). The file will be created in the current working directory or overwritten if it already exists.
- 'json_file.write(model_architecture)': The model's architecture, in JSON format, is
 written to the opened file. This file will contain information about the layers and their
 configurations.

• `model.save_weights("model_weights.h5") `: This function saves the weights of the trained model to a file named "model_weights.h5". These weights contain the learned parameters (such as the connection weights and biases) of the model. Saving weights separately from the architecture is common practice; it allows you to rebuild the model's architecture and then load the trained weights to make predictions without retraining.

After executing these lines, you will have two files in your current directory:

- model_architecture.json: Contains the JSON representation of the neural network model's architecture. This file describes the layout of the layers and their configurations.
- model_weights.h5: Contains the learned weights of the trained model. These weights are
 crucial for making predictions because they represent the knowledge the model gained
 during training.

You can later use these saved files to reconstruct the model and make predictions on new data without having to retrain the entire model. This is especially useful when you want to deploy your trained model for real-world applications.

In neural networks, connection weights and biases are fundamental concepts that enable the network to learn and make predictions. Let's break down these concepts:

Connection Weights:

- Definition: Connection weights are numerical values associated with the connections between neurons in a neural network. Each connection between two neurons has a weight associated with it.
- **Purpose:** These weights signify the strength of the connections. A higher weight means a stronger connection, and a lower weight means a weaker connection.
- Role in Learning: During the training process, the neural network adjusts these weights to
 minimize the difference between predicted outputs and actual targets. Adjusting the
 weights allows the network to learn the patterns and relationships in the data.
- Mathematical Operation: In a neural network layer, the weighted sum of inputs is calculated
 for each neuron. This sum is then passed through an activation function to introduce nonlinearity. The weights determine how much influence each input has on the neuron's output.
- **Training:** During training, these weights are updated using algorithms like backpropagation and gradient descent. The goal is to find the optimal set of weights that minimizes the loss function, representing the disparity between predictions and actual values.

Biases:

- **Definition:** Biases are additional parameters in a neural network that are associated with each neuron. Every neuron (except the input neurons) has an associated bias term.
- Purpose: Biases allow neurons to adjust their output independently of the inputs. They
 provide flexibility and enable neurons to activate even when the weighted sum of inputs is
 not significantly different from zero.
- Role in Learning: Similar to weights, biases are adjusted during the training process. They allow the model to shift the activation function's output up or down, enabling the network to learn complex patterns in the data.
- Mathematical Operation: Biases are added to the weighted sum of inputs for each neuron before passing through the activation function. They provide an offset, allowing neurons to activate even when the sum of inputs is zero.

In summary, connection weights and biases are parameters that a neural network learns during the training process. Weights determine the strength of connections between neurons, while biases allow neurons to activate even when the inputs are not strong. The combination of weights and biases, along with activation functions, enables neural networks to approximate complex functions and learn intricate patterns in the data, making them powerful tools for various machine learning tasks.

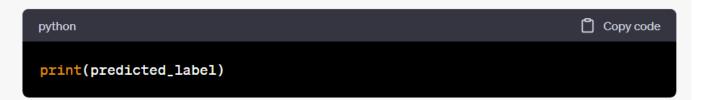
In the provided code, you're using the trained tokenizer, model, and label encoder to predict the emotion label for a new input text. Let's break down each line: Copy code python input_text = "I'm smart and talented and open to new experiences" `input_text`: This variable contains the new text input that you want to classify for emotion. Copy code python input_sequence = tokenizer.texts_to_sequences([input_text]) `tokenizer.texts_to_sequences([input_text])`: Tokenizes the input text using the same tokenizer that was used during the training phase. It converts the input text into a sequence of integers based on the vocabulary learned from the training data. The input is wrapped in a list because `texts_to_sequences` expects a list of texts. □ Regen Copy code python padded_input_sequence = pad_sequences(input_sequence, maxlen=max_length) 'pad_sequences(input_sequence, maxlen=max_length)': Pads the tokenized input sequence to match the required input length (`max_length`) of the neural network model. This ensures that the input has the same shape as the data the model was trained on. Copy code python prediction = model.predict(padded_input_sequence) `model.predict(padded_input_sequence) `: Uses the trained neural network model to predict the emotion label for the input text. The `padded_input_sequence` is passed to the

model, and it returns a prediction. The prediction is a probability distribution over the

different emotion classes due to the softmax activation function in the output layer.

predicted_label = label_encoder.inverse_transform([np.argmax(prediction[0])]

- `np.argmax(prediction[0])`: Finds the index of the class with the highest probability in the prediction array. This index represents the predicted emotion class.
- `label_encoder.inverse_transform([np.argmax(prediction[0])]) `: Uses the `inverse_transform` method of the label encoder to convert the predicted index back to the original label.



* `print(predicted_label) `: Prints the predicted emotion label for the input text.

In summary, this code takes a new input text, preprocesses it using the trained tokenizer, predicts the emotion label using the trained neural network model, and prints the predicted label. It's a typical process for making predictions with a trained NLP model after the training phase.

-----SUMMARY -----

In the provided code snippets and explanations, you are creating a **text classification model** using a **neural network** for **emotion classification**. Let me summarize the key components of this model:

Model Type:

• Neural Network Model: Specifically, a feedforward neural network, also known as a multilayer perceptron (MLP). It consists of an input layer, one or more hidden layers, and an output layer.

Problem Type:

• **Text Classification:** The model is designed to classify text inputs into specific categories or classes, which, in this case, represent different emotions.

Components of the Model:

1. Input Layer:

 The input layer processes tokenized and padded sequences of words. Each word is represented by an integer index obtained from tokenization.

2. Embedding Layer:

 Converts the integer indices of words into dense vectors of fixed size (embedding vectors). It captures semantic relationships between words based on their context within the dataset.

3. Flatten Layer:

• Flattens the 2D array output from the embedding layer into a 1D array. This flattening is necessary to connect to dense layers.

4. Dense Layers:

- One or more dense (fully connected) hidden layers, which learn complex patterns from the input data.
- The last dense layer has units equal to the number of unique emotion classes and uses the softmax activation function to output class probabilities.

Training Process:

- Loss Function: Categorical Cross-Entropy. This loss function is commonly used for multiclass classification problems.
- Optimizer: Adam Optimizer. It adapts the learning rates of each parameter, often leading to faster convergence during training.
- Metrics: Accuracy. It measures the proportion of correctly classified samples.

Data Preprocessing:

- Tokenization: The text data is tokenized, converting words into integer indices.
- Padding: Sequences of integers are padded to ensure uniform length, necessary for input into neural networks.
- Label Encoding: Emotion labels are encoded into integers.
- One-Hot Encoding: Integer-encoded labels are further converted into one-hot encoded vectors, suitable for multi-class classification.

Model Output:

 The model provides predictions in the form of class probabilities (thanks to the softmax activation function). The class with the highest probability is considered the predicted emotion for a given input text.

In summary, the model you're creating is a neural network designed for text classification, specifically for categorizing emotions based on input text. It's trained to learn the underlying patterns in the textual data and make predictions about the emotional content of new, unseen text samples.