

Implement the Continuous Bag of Words (CBOW) Model. Stages can be:

a. Data preparation b. Generate training data c. Train model d. Output

-Tushar Bhagat Roll_No. 07

```
In [1]: import numpy as np
import pandas as pd
import tensorflow as tf
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Embedding, Dense, Flatten

# Sample data
corpus = [
    "I love machine learning",
    "Deep learning is a subset of machine learning",
    "Natural language processing is fascinating",
    "I enjoy learning new things",
    "Machine learning can be applied in various fields"
]

# Step 1: Data preparation
# Tokenization
tokenizer = Tokenizer()
tokenizer.fit_on_texts(corpus)
word_index = tokenizer.word_index
total_words = len(word_index) + 1 # Add 1 for padding

print("Word Index:", word_index)
print("Total Words:", total_words)
```

Word Index: {'learning': 1, 'machine': 2, 'i': 3, 'is': 4, 'love': 5, 'deep': 6, 'a': 7, 'subset': 8, 'of': 9, 'natural': 10, 'language': 11, 'processing': 12, 'fascinating': 13, 'enjoy': 14, 'new': 15, 'things': 16, 'can': 17, 'be': 18, 'applied': 19, 'in': 20, 'various': 21, 'fields': 22}
Total Words: 23

```
In [2]: # Step 2: Generate training data
def generate_cbow_data(corpus, window_size=2):
    input_data = []
    output_data = []

    for sentence in corpus:
        tokenized_sentence = tokenizer.texts_to_sequences([sentence])[0]
        for i in range(window_size, len(tokenized_sentence) - window_size):
            context = []
            for j in range(i - window_size, i + window_size + 1):
                if j != i: # Skip the target word
                    context.append(tokenized_sentence[j])
            input_data.append(context)
            output_data.append(tokenized_sentence[i])

    return np.array(input_data), np.array(output_data)

input_data, output_data = generate_cbow_data(corpus)
print("Input Data:", input_data)
print("Output Data:", output_data)
```

Input Data: [[6 1 7 8]
[1 4 8 9]
[4 7 9 2]
[7 8 2 1]
[10 11 4 13]
[3 14 15 16]
[2 1 18 19]
[1 17 19 20]
[17 18 20 21]
[18 19 21 22]]
Output Data: [4 7 8 9 12 1 17 18 19 20]

```
In [3]: # Step 3: Train model
# One-hot encoding of the output data
output_data = tf.keras.utils.to_categorical(output_data, num_classes=total_words)

# Define the model
model = Sequential()
model.add(Embedding(input_dim=total_words, output_dim=10, input_length=input_data.shape[1]))
model.add(Flatten())
model.add(Dense(total_words, activation='softmax'))
```

```
# Compile the model
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])

# Train the model
model.fit(input_data, output_data, epochs=100, verbose=1)
```

Epoch 1/100

C:\Users\tmbha\AppData\Local\Programs\Python\Python310\lib\site-packages\keras\src\layers\core\embedding.py:90: UserWarning: Argument `input_length` is deprecated. Just remove it.
warnings.warn(

1/1		1s	1s/step - accuracy: 0.0000e+00 - loss: 3.1428
Epoch 2/100			
1/1		0s	33ms/step - accuracy: 0.0000e+00 - loss: 3.1363
Epoch 3/100			
1/1		0s	35ms/step - accuracy: 0.0000e+00 - loss: 3.1298
Epoch 4/100			
1/1		0s	57ms/step - accuracy: 0.0000e+00 - loss: 3.1233
Epoch 5/100			
1/1		0s	43ms/step - accuracy: 0.0000e+00 - loss: 3.1168
Epoch 6/100			
1/1		0s	25ms/step - accuracy: 0.0000e+00 - loss: 3.1102
Epoch 7/100			
1/1		0s	45ms/step - accuracy: 0.1000 - loss: 3.1037
Epoch 8/100			
1/1		0s	15ms/step - accuracy: 0.2000 - loss: 3.0972
Epoch 9/100			
1/1		0s	39ms/step - accuracy: 0.2000 - loss: 3.0907
Epoch 10/100			
1/1		0s	42ms/step - accuracy: 0.3000 - loss: 3.0842
Epoch 11/100			
1/1		0s	42ms/step - accuracy: 0.3000 - loss: 3.0776
Epoch 12/100			
1/1		0s	39ms/step - accuracy: 0.3000 - loss: 3.0711
Epoch 13/100			
1/1		0s	49ms/step - accuracy: 0.4000 - loss: 3.0645
Epoch 14/100			
1/1		0s	45ms/step - accuracy: 0.4000 - loss: 3.0579
Epoch 15/100			
1/1		0s	44ms/step - accuracy: 0.6000 - loss: 3.0513
Epoch 16/100			
1/1		0s	51ms/step - accuracy: 0.7000 - loss: 3.0447
Epoch 17/100			
1/1		0s	53ms/step - accuracy: 0.8000 - loss: 3.0381
Epoch 18/100			
1/1		0s	49ms/step - accuracy: 0.8000 - loss: 3.0314
Epoch 19/100			
1/1		0s	41ms/step - accuracy: 0.8000 - loss: 3.0246
Epoch 20/100			
1/1		0s	93ms/step - accuracy: 0.8000 - loss: 3.0179
Epoch 21/100			
1/1		0s	47ms/step - accuracy: 0.8000 - loss: 3.0111
Epoch 22/100			
1/1		0s	41ms/step - accuracy: 0.8000 - loss: 3.0043
Epoch 23/100			
1/1		0s	48ms/step - accuracy: 0.8000 - loss: 2.9974
Epoch 24/100			
1/1		0s	48ms/step - accuracy: 0.8000 - loss: 2.9904
Epoch 25/100			
1/1		0s	43ms/step - accuracy: 0.8000 - loss: 2.9834
Epoch 26/100			
1/1		0s	32ms/step - accuracy: 1.0000 - loss: 2.9764
Epoch 27/100			
1/1		0s	63ms/step - accuracy: 1.0000 - loss: 2.9693
Epoch 28/100			
1/1		0s	70ms/step - accuracy: 1.0000 - loss: 2.9621
Epoch 29/100			
1/1		0s	39ms/step - accuracy: 1.0000 - loss: 2.9548
Epoch 30/100			
1/1		0s	29ms/step - accuracy: 1.0000 - loss: 2.9475
Epoch 31/100			
1/1		0s	48ms/step - accuracy: 1.0000 - loss: 2.9401
Epoch 32/100			
1/1		0s	34ms/step - accuracy: 1.0000 - loss: 2.9326
Epoch 33/100			
1/1		0s	48ms/step - accuracy: 1.0000 - loss: 2.9251
Epoch 34/100			
1/1		0s	38ms/step - accuracy: 1.0000 - loss: 2.9174
Epoch 35/100			
1/1		0s	31ms/step - accuracy: 1.0000 - loss: 2.9097
Epoch 36/100			
1/1		0s	27ms/step - accuracy: 1.0000 - loss: 2.9018
Epoch 37/100			
1/1		0s	40ms/step - accuracy: 1.0000 - loss: 2.8939
Epoch 38/100			
1/1		0s	25ms/step - accuracy: 1.0000 - loss: 2.8859
Epoch 39/100			
1/1		0s	55ms/step - accuracy: 1.0000 - loss: 2.8777
Epoch 40/100			
1/1		0s	39ms/step - accuracy: 1.0000 - loss: 2.8695
Epoch 41/100			
1/1		0s	59ms/step - accuracy: 1.0000 - loss: 2.8611
Epoch 42/100			
1/1		0s	49ms/step - accuracy: 1.0000 - loss: 2.8527
Epoch 43/100			
1/1		0s	34ms/step - accuracy: 1.0000 - loss: 2.8441
Epoch 44/100			
1/1		0s	45ms/step - accuracy: 1.0000 - loss: 2.8354
Epoch 45/100			
1/1		0s	47ms/step - accuracy: 1.0000 - loss: 2.8266
Epoch 46/100			

1/1		0s	41ms/step	- accuracy: 1.0000	- loss: 2.8176
Epoch 47/100					
1/1		0s	36ms/step	- accuracy: 1.0000	- loss: 2.8085
Epoch 48/100					
1/1		0s	32ms/step	- accuracy: 1.0000	- loss: 2.7993
Epoch 49/100					
1/1		0s	46ms/step	- accuracy: 1.0000	- loss: 2.7900
Epoch 50/100					
1/1		0s	32ms/step	- accuracy: 1.0000	- loss: 2.7805
Epoch 51/100					
1/1		0s	67ms/step	- accuracy: 1.0000	- loss: 2.7709
Epoch 52/100					
1/1		0s	49ms/step	- accuracy: 1.0000	- loss: 2.7612
Epoch 53/100					
1/1		0s	36ms/step	- accuracy: 1.0000	- loss: 2.7513
Epoch 54/100					
1/1		0s	39ms/step	- accuracy: 1.0000	- loss: 2.7412
Epoch 55/100					
1/1		0s	48ms/step	- accuracy: 1.0000	- loss: 2.7310
Epoch 56/100					
1/1		0s	53ms/step	- accuracy: 1.0000	- loss: 2.7207
Epoch 57/100					
1/1		0s	31ms/step	- accuracy: 1.0000	- loss: 2.7102
Epoch 58/100					
1/1		0s	27ms/step	- accuracy: 1.0000	- loss: 2.6996
Epoch 59/100					
1/1		0s	32ms/step	- accuracy: 1.0000	- loss: 2.6888
Epoch 60/100					
1/1		0s	37ms/step	- accuracy: 1.0000	- loss: 2.6778
Epoch 61/100					
1/1		0s	42ms/step	- accuracy: 1.0000	- loss: 2.6667
Epoch 62/100					
1/1		0s	35ms/step	- accuracy: 1.0000	- loss: 2.6555
Epoch 63/100					
1/1		0s	47ms/step	- accuracy: 1.0000	- loss: 2.6441
Epoch 64/100					
1/1		0s	71ms/step	- accuracy: 1.0000	- loss: 2.6325
Epoch 65/100					
1/1		0s	48ms/step	- accuracy: 1.0000	- loss: 2.6207
Epoch 66/100					
1/1		0s	40ms/step	- accuracy: 1.0000	- loss: 2.6088
Epoch 67/100					
1/1		0s	46ms/step	- accuracy: 1.0000	- loss: 2.5967
Epoch 68/100					
1/1		0s	36ms/step	- accuracy: 1.0000	- loss: 2.5845
Epoch 69/100					
1/1		0s	30ms/step	- accuracy: 1.0000	- loss: 2.5721
Epoch 70/100					
1/1		0s	36ms/step	- accuracy: 1.0000	- loss: 2.5595
Epoch 71/100					
1/1		0s	40ms/step	- accuracy: 1.0000	- loss: 2.5467
Epoch 72/100					
1/1		0s	76ms/step	- accuracy: 1.0000	- loss: 2.5338
Epoch 73/100					
1/1		0s	77ms/step	- accuracy: 1.0000	- loss: 2.5207
Epoch 74/100					
1/1		0s	41ms/step	- accuracy: 1.0000	- loss: 2.5074
Epoch 75/100					
1/1		0s	40ms/step	- accuracy: 1.0000	- loss: 2.4940
Epoch 76/100					
1/1		0s	42ms/step	- accuracy: 1.0000	- loss: 2.4804
Epoch 77/100					
1/1		0s	33ms/step	- accuracy: 1.0000	- loss: 2.4666
Epoch 78/100					
1/1		0s	45ms/step	- accuracy: 1.0000	- loss: 2.4527
Epoch 79/100					
1/1		0s	40ms/step	- accuracy: 1.0000	- loss: 2.4385
Epoch 80/100					
1/1		0s	44ms/step	- accuracy: 1.0000	- loss: 2.4242
Epoch 81/100					
1/1		0s	45ms/step	- accuracy: 1.0000	- loss: 2.4098
Epoch 82/100					
1/1		0s	44ms/step	- accuracy: 1.0000	- loss: 2.3951
Epoch 83/100					
1/1		0s	42ms/step	- accuracy: 1.0000	- loss: 2.3803
Epoch 84/100					
1/1		0s	34ms/step	- accuracy: 1.0000	- loss: 2.3653
Epoch 85/100					
1/1		0s	31ms/step	- accuracy: 1.0000	- loss: 2.3502
Epoch 86/100					
1/1		0s	46ms/step	- accuracy: 1.0000	- loss: 2.3348
Epoch 87/100					
1/1		0s	82ms/step	- accuracy: 1.0000	- loss: 2.3193
Epoch 88/100					
1/1		0s	48ms/step	- accuracy: 1.0000	- loss: 2.3037
Epoch 89/100					
1/1		0s	41ms/step	- accuracy: 1.0000	- loss: 2.2878
Epoch 90/100					
1/1		0s	46ms/step	- accuracy: 1.0000	- loss: 2.2718
Epoch 91/100					

1/1 0s 50ms/step - accuracy: 1.0000 - loss: 2.2557

Epoch 92/100

1/1 0s 32ms/step - accuracy: 1.0000 - loss: 2.2393

Epoch 93/100

1/1 0s 46ms/step - accuracy: 1.0000 - loss: 2.2228

Epoch 94/100

1/1 0s 37ms/step - accuracy: 1.0000 - loss: 2.2062

Epoch 95/100

1/1 0s 45ms/step - accuracy: 1.0000 - loss: 2.1894

Epoch 96/100

1/1 0s 34ms/step - accuracy: 1.0000 - loss: 2.1724

Epoch 97/100

1/1 0s 40ms/step - accuracy: 1.0000 - loss: 2.1553

Epoch 98/100

1/1 0s 43ms/step - accuracy: 1.0000 - loss: 2.1380

Epoch 99/100

1/1 0s 38ms/step - accuracy: 1.0000 - loss: 2.1206

Epoch 100/100

1/1 0s 32ms/step - accuracy: 1.0000 - loss: 2.1030

Out[3]: <keras.src.callbacks.history.History at 0x1fbde20e950>

In [6]:

```
# Step 4: Output
def predict_word(context):
    context_seq = tokenizer.texts_to_sequences([context])[0]
    context_seq = pad_sequences([context_seq], maxlen=input_data.shape[1], padding='post')
    predicted = model.predict(context_seq)
    predicted_word_index = np.argmax(predicted, axis=-1)[0]
    for word, index in word_index.items():
        if index == predicted_word_index:
            return word

# Test prediction
context_words = ["I", "love", "language", "machine"]
predicted_word = predict_word(context_words)
print(f"Context: {context_words}, Predicted Word: {predicted_word}")
```

1/1 0s 38ms/step

Context: ['I', 'love', 'language', 'machine'], Predicted Word: learning