

PR 5 continuous-bag-of-words

It is a model that tries to predict words given the context of a few words before and a few words after the target word.

The Continuous Bag-of-Words model (CBOW) is frequently used in NLP deep learning.
It is a model that tries to predict words given the context of a few words before and a few words after the target word.
This is distinct from language modeling, since CBOW is not sequential and does not have to be probabilistic.
Typically, CBOW is used to quickly train word embeddings, and these embeddings are used to initialize the embeddings of some more complex models.
Usually, this is referred to as pretraining embeddings. It almost always helps performance a couple of percent.
This is the solution of the final exercise of this great tutorial on NLP in PyTorch.

```
import torch
import torch.nn as nn

def make_context_vector(context, word_to_ix):
    idxs = [word_to_ix[w] for w in context]
    return torch.tensor(idxs, dtype=torch.long)

CONTEXT_SIZE = 2  # 2 words to the left, 2 to the right
EMDEDDING_DIM = 100

raw_text = """We are about to study the idea of a computational process.
Computational processes are abstract beings that inhabit computers.
As they evolve, processes manipulate other abstract things called data.
The evolution of a process is directed by a pattern of rules
called a program. People create programs to direct processes. In effect,
we conjure the spirits of the computer with our spells.""".split()

# By deriving a set from `raw_text`, we deduplicate the array
vocab = set(raw_text)
vocab_size = len(vocab)

word_to_ix = {word:ix for ix, word in enumerate(vocab)}
ix_to_word = {ix:word for ix, word in enumerate(vocab)}

data = []
for i in range(2, len(raw_text) - 2):
    context = [raw_text[i - 2], raw_text[i - 1],
               raw_text[i + 1], raw_text[i + 2]]
    target = raw_text[i]
    data.append((context, target))

class CBOW(torch.nn.Module):
    def __init__(self, vocab_size, embedding_dim):
        super(CBOW, self).__init__()

        #out: 1 x emdedding_dim
        self.embeddings = nn.Embedding(vocab_size, embedding_dim)
        self.linear1 = nn.Linear(embedding_dim, 128)
        self.activation_function1 = nn.ReLU()

        #out: 1 x vocab_size
        self.linear2 = nn.Linear(128, vocab_size)
        self.activation_function2 = nn.LogSoftmax(dim = -1)

    def forward(self, inputs):
        embeds = sum(self.embeddings(inputs)).view(1,-1)
        out = self.linear1(embeds)
        out = self.activation_function1(out)
        out = self.linear2(out)
        out = self.activation_function2(out)
        return out

    def get_word_emdedding(self, word):
        word = torch.tensor([word_to_ix[word]])
        return self.embeddings(word).view(1,-1)

model = CBOW(vocab_size, EMDEDDING_DIM)

loss_function = nn.NLLLoss()
optimizer = torch.optim.SGD(model.parameters(), lr=0.001)

#TRAINING
for epoch in range(50):
    total_loss = 0

    for context, target in data:
```

```
context_vector = make_context_vector(context, word_to_ix)

log_probs = model(context_vector)

total_loss += loss_function(log_probs, torch.tensor([word_to_ix[target]]))

#optimize at the end of each epoch
optimizer.zero_grad()
total_loss.backward()
optimizer.step()

#TESTING
context = ['People','create','to', 'direct']
context_vector = make_context_vector(context, word_to_ix)
a = model(context_vector)

#Print result
print(f'Raw text: {" ".join(raw_text)}\n')
print(f'Context: {context}\n')
print(f'Prediction: {ix_to_word[torch.argmax(a[0]).item()]}')

print(vocab)

print(ix_to_word)
```

```
📄 Raw text: We are about to study the idea of a computational process. Computational processes are abstract beings that inhabit con

Context: ['People', 'create', 'to', 'direct']

Prediction: programs
{'abstract', 'computers.', 'direct', 'computer', 'are', 'that', 'pattern', 'computational', 'People', 'processes.', 'processes',
{0: 'abstract', 1: 'computers.', 2: 'direct', 3: 'computer', 4: 'are', 5: 'that', 6: 'pattern', 7: 'computational', 8: 'People',
```