

Deep Learning — Assignment 5

Fifth assignment for the 2023 Deep Learning course (NWI-IMC070) of the Radboud University.

Names:

Group:

Instructions:

- Fill in your names and the name of your group.
- Answer the questions and complete the code where necessary.
- Keep your answers brief, one or two sentences is usually enough.
- Re-run the whole notebook before you submit your work.
- Save the notebook as a PDF and submit that in Brightspace together with the `.ipynb` notebook file.
- The easiest way to make a PDF of your notebook is via File > Print Preview and then use your browser's print option to print to PDF.

Objectives

In this assignment you will

1. Construct a PyTorch `DataSet`
2. Train and modify a transformer network
3. Experiment with a translation dataset

Required software

If you haven't done so already, you will need to install the following additional libraries:

- `torch` for PyTorch,
- `d2l`, the library that comes with the [Dive into deep learning](#) book.
Note: if you get errors, make sure the right version of the d2l library is installed: `pip install d2l==1.0.0a1.post0`

All libraries can be installed with `pip install`.

```
In [ ]: %matplotlib inline
from d2l import torch as d2l
import math
from random import Random
from typing import List
import numpy as np
import torch
from torch import nn
from torch.utils.data import IterableDataset, DataLoader
import matplotlib.pyplot as plt

device = d2l.try_gpu()
```

5.1 Learning to calculate (5 points)

In this assignment we are going to train a neural network to do mathematics. When communicating between humans, mathematics is expressed with words and formulas. The simplest of these are formulas with a numeric answer. For example, we might ask what is `100+50`, to which the answer is `150`.

To teach a computer how to do this task, we are going to need a dataset.

Below is a function that generates a random formula. Study it, and see if you understand its parameters.

```
In [ ]: def random_integer(length: int, signed: bool = True, rng: Random = Random())
    max = math.pow(10, length)
    min = -max if signed else 0
    return rng.randint(min, max)

def random_formula(complexity: int, signed: bool = True, rng: Random = Random())
    """
    Generate a random formula of the form "a+b" or "a-b".
    complexity is the maximum number of digits in the numbers.
    """
    a = random_integer(complexity, signed, rng)
    b = random_integer(complexity, False, rng)
    is_addition = not signed or rng.choice([False, True])
    if is_addition:
        return (f"{a}+{b}", str(a + b))
    else:
        return (f"{a}-{b}", str(a - b))
```

```
In [ ]: seed = 123456
random_formula(3, rng=Random(seed))
```

Note that the `rng` argument allows us to reproduce the same random numbers, which you can verify by running the code below multiple times. But if you change the seed to `None` then the random generator is initialized differently each time.

```
In [ ]: def random_formulas(complexity, signed, count, seed):
        """
        Iterator that yields the given count of random formulas
        """
        rng = Random(seed)
        for i in range(count):
            yield random_formula(complexity, signed, rng=rng)

for q, a in random_formulas(3, True, 5, seed):
    print(f'{q} = {a}')
```

We are going to treat these expressions as sequences of tokens, where each character is a token. In addition we will need tokens to denote begin-of-sequence and end-of-sequence, as well as padding, for which we will use '<bos>', '<eos>', and '<pad>' respectively, as is done in the book.

d2l chapter 9.2 includes an example of tokenizing a string, and it also defines a `Vocab` class that handles converting the tokens to numbers.

For this dataset we know beforehand what the vocabulary will be.

Creating a vocabulary

(a) What are the tokens in this dataset? Complete the code below.

(1 point)

```
In [ ]: # TODO: fill in all possible tokens
vocab = d2l.Vocab([...], reserved_tokens=[...])
```

We can print the vocabulary to double check that it makes sense:

```
In [ ]: print('Vocabulary size:', len(vocab))
        print('Vocabulary:', vocab.idx_to_token)
```

Note that the d2l Vocab class includes a '<unk>' token, for handling unknown tokens in the input.

We are now ready to tokenize and encode formula.

(b) Complete the code below.

(1 point)

```
In [ ]: def tokenize_and_encode(string: str, vocab=vocab) -> List[int]:
        # TODO: Tokenize the string and encode using the vocabulary.
        #         Include an end-of-string token (but not a begin-of-string token)
```

Let's test it on a random formula:

```
In [ ]: q, a = random_formula(3, rng=Random(seed))
        print('The question', q, 'and answer', a)
        print('are encoded as', tokenize_and_encode(q), 'and', tokenize_and_encode(a))
```

```
# Check tokenize_and_encode
assert ''.join(vocab.to_tokens(tokenize_and_encode(q))) == q + '<eos>'
assert len(tokenize_and_encode(q)) == len(q) + 1
```

Padding and trimming

Next, to be able to work with a whole dataset of these encoded sequences, they all need to be the same length.

(c) Implement the function below that pads or trims the encoded token sequence as needed. (1 point)

Hint: see [d2l section 10.5.3](#) for a very similar function.

```
In [ ]: def pad_or_trim(tokens: List[int], target_length: int, vocab=vocab):
        # TODO return a padded or trimmed sequence
```

```
In [ ]: # pad or trim q to get a sequence of 10 tokens
        pad_or_trim(tokenize_and_encode(q), 10)
```

```
In [ ]: # Check pad_or_trim
        assert len(pad_or_trim([1,2,3,4,5],10)) == 10
        assert len(pad_or_trim(list(range(20)),10)) == 10
        assert vocab.to_tokens(pad_or_trim([1,2,3,4,5],10)[5:]) == ['<pad>', '<pad>',
        f"Incorrect padding tokens, found {vocab.to_tokens(pad_or_trim([1,2,3
```

Translating tokens

We can use `vocab.to_tokens` to convert the encoded token sequence back to something more readable:

```
In [ ]: vocab.to_tokens(pad_or_trim(tokenize_and_encode(q), 10))
```

For convenience, we define the `decode_tokens` function to convert entire lists or tensors:

```
In [ ]: def decode_tokens(t, vocab=vocab):
        # convert a list, tensor, or array of encoded tokens
        if isinstance(t, torch.Tensor):
            t = t.detach().cpu().numpy()
            t = np.asarray(t)
            return np.asarray(vocab.to_tokens(list(t.flatten()))).reshape(*t.shape)

        # convert all tokens at once
        print(decode_tokens([pad_or_trim(tokenize_and_encode('513+1323'), 10),
                                pad_or_trim(tokenize_and_encode('412+42'), 10)]))
```

Creating a dataset

The most convenient way to use a data generating function for training a neural network is to wrap it in a PyTorch `Dataset`. In this case, we will use an `IterableDataset`, which can be used as an iterator to walk over the samples in the dataset.

(d) Complete the code below.

(1 point)

```
In [ ]: class FormulaDataset(IterableDataset):
        def __init__(self, complexity, signed, count, seed=None, vocab=vocab):
            self.seed = seed
            self.complexity = complexity
            self.signed = signed
            self.count = count
            self.vocab = vocab
            self.max_question_length = 2 * complexity + 3
            self.max_answer_length = complexity + 2

        # TODO: Complete the class definition.
        #     See the documentation for IterableDataset for examples.
        #     Make sure that the values yielded by the iterator are pairs of t
        #     To create a repeatable dataset, always start with the same random
```

(e) Define a training set with 10000 formulas and a validation set with 5000 formulas, both with complexity 3.

(1 point)

Note: make sure that the training and validation set are different.

```
In [ ]: complexity = 3
        signed = True
        # TODO: Your code here.
        train_data = ...
        val_data = ...
```

As usual, we wrap each dataset in a `DataLoader` to create minibatches.

```
In [ ]: # Define data loaders
        batch_size = 125
        data_loaders = {
            'train': torch.utils.data.DataLoader(train_data, batch_size=batch_size),
            'val': torch.utils.data.DataLoader(val_data, batch_size=batch_size),
        }
```

```
In [ ]: # The code below checks that the datasets are defined correctly
        train_loader = data_loaders['train']
        val_loader = data_loaders['val']

        from typing import Tuple
        from typing_extensions import assert_type
        for (name, loader), expected_size in zip(data_loaders.items(), [10000, 5000]):
            first_batch = next(iter(loader))
            assert len(first_batch) == 2, \
                f"The {name} dataset should yield (question, answer) pairs when i
```

```

assert torch.is_tensor(first_batch[0]), \
    f"The questions in the {name} dataset should be torch.tensors"
assert tuple(first_batch[0].shape) == (batch_size, 2*complexity+3), \
    f"The questions in the {name} dataset should be of size (batch_si
assert first_batch[0].dtype in [torch.int32, torch.int64], \
    f"The questions in the {name} dataset should be encoded as intege
assert torch.equal(next(iter(loader))[0], next(iter(loader))[0]), \
    f"The {name} dataset should be deterministic, it should produce t
assert all([len(batch[0]) == batch_size for batch in iter(loader)]), \
    f"Batches should all have the right size. Perhaps the batch size
assert sum([len(batch[0]) for batch in iter(loader)]) == expected_size,
    f"{name} dataset does not have the right size, expected {expected
assert not torch.equal(next(iter(train_loader))[0], next(iter(val_loader))[0
    "The training data and validation data should not be the same"

```

5.2 Transformer inputs (10 points)

There is a detailed description of the transformer model in [chapter 11 of the d2l book](#). We will not use most the code from the book, and instead use [PyTorch's built-in Transformer layers](#).

However, some details we still need to implement ourselves.

Masks

Training a transformer uses masked self-attention, so we need some masks. Here are two functions that make these masks.

```

In [ ]: def generate_square_subsequent_mask(size, device=device):
    """
    Mask that indicates that tokens at a position are not allowed to attend
    tokens in subsequent positions.
    """
    mask = (torch.tril(torch.ones((size, size), device=device))) == 0
    return mask

def generate_padding_mask(tokens, padding_token):
    """
    Mask that indicates which tokens should be ignored because they are pad
    """
    return tokens == torch.tensor(padding_token)

```

(a) Generate a padding mask for a random encoded token string.

(1 point)

Hint: make sure that `tokens` is a torch.tensor.

```

In [ ]: q, a = random_formula(3, rng=Random(seed))
# TODO: your code here
tokens = ...
padding_mask = ...
print(tokens)

```

```
print(decode_tokens(tokens))
print(padding_mask)
```

```
In [ ]: # More tests
        assert list(generate_padding_mask(torch.tensor(pad_or_trim(tokenize_and_encode('The quick brown fox jumps over the lazy dog')))))
```

(b) How will this mask be used by a transformer? (1 point)

TODO: Your answer here.

The code below takes the first batch of data from the training set, and it generates a shifted version of the target values.

```
In [ ]: x, y = next(iter(train_loader))
bos = torch.tensor(vocab['<bos>']).expand(y.shape[0], 1)
y_prev = torch.cat((bos, y[:, :-1]), axis=1)

# print the first five samples
print(decode_tokens(y[:5]))
print(decode_tokens(y_prev[:5]))
```

(c) Look at the values for the example above. What is `y_prev` used for during training of a transformer model? (1 point)

TODO: Your answer here.

(d) Why do some rows of `y_prev` end in '`<eos>`', but not all? Is this a problem? (1 point)

TODO: Your answer here.

The code below illustrates what the output of `generate_square_subsequent_mask` looks like.

```
In [ ]: square_subsequent_mask = generate_square_subsequent_mask(y.shape[1])

print(square_subsequent_mask.shape)
print(square_subsequent_mask)
```

(e) How and why should this mask be used? State your answer in terms of `x`, `y` and/or `y prev`. (1 point)

TODO: Your answer here.

(f) Give an example where it could make sense to use a different mask in a transformer network, instead of the `square_subsequent_mask` ? (1 point)

TODO: Your answer here.

Embedding

Our discrete vocabulary is not suitable as the input for a transformer. We need an embedding function to map our input vocabulary to a continuous, high-dimensional space.

We will use the `torch.nn.Embedding` class to for this. As you can read in the [documentation](#), this class maps each token in our vocabulary to a specific point in embedding space, its embedding vector. We will use this embedding vector as the input features for the next layer of our model.

The parameters of the embedding are trainable: the embedding vector of each token is optimized along with the rest of the network.

(g) Define an embedding that maps our vocabulary to a 5-dimensional space. (1 point)

```
In [ ]: # TODO: Your code here.
        embedding = ...
        print(embedding)
```

Let's apply the embedding to some sequences from our training set.

```
In [ ]: # take the first batch
        x, y = next(iter(train_loader))
        # take three samples
        x = x[:3]
        # print the shapes
        print(x)
        print(embedding(x))
        print(x.shape)
        print(embedding(x).shape)
```

(h) Explain the output shape. (1 point)

TODO: Your answer here.

The size of the embedding vectors, or the dimensionality of the embedding space, does not depend on the number of tokens in our vocabulary. We are free to choose an embedding size that fits our problem.

For example, let's try an embedding with 2 dimensions, and plot the initial embedding for the tokens in our vocabulary.

(i) Create an embedding with 2 dimensions and plot the embedding for all tokens. (no points)


```
In [ ]: # TODO: Your code here.
        embedding = ...

        # embed all tokens of our vocabulary
        x = torch.arange(len(vocab))
        emb = embedding(x).detach().cpu().numpy()

        plt.scatter(emb[:, 0], emb[:, 1]);
        for i, token in enumerate(vocab.idx_to_token):
            plt.annotate(token, (emb[i,0]+0.04, emb[i,1]))
```

As always, we need to balance the complexity of our networks: a larger embedding will increase the number of parameters in our model, but increase the risk of overfitting.

(j) Would this 2-dimensional embedding space be large enough for our problem? **(1 point)**

TODO: Your answer here.

Instead of using an embedding, we could also use a simple one-hot encoding to map the words in the vocabulary to feature vectors. However, practical applications of natural language processing never do this. Why not?

(k) Explain the practical advantage of embeddings over one-hot encoding. **(1 point)**

TODO: Your answer here.

5.3 torch.nn.Transformer (8 points)

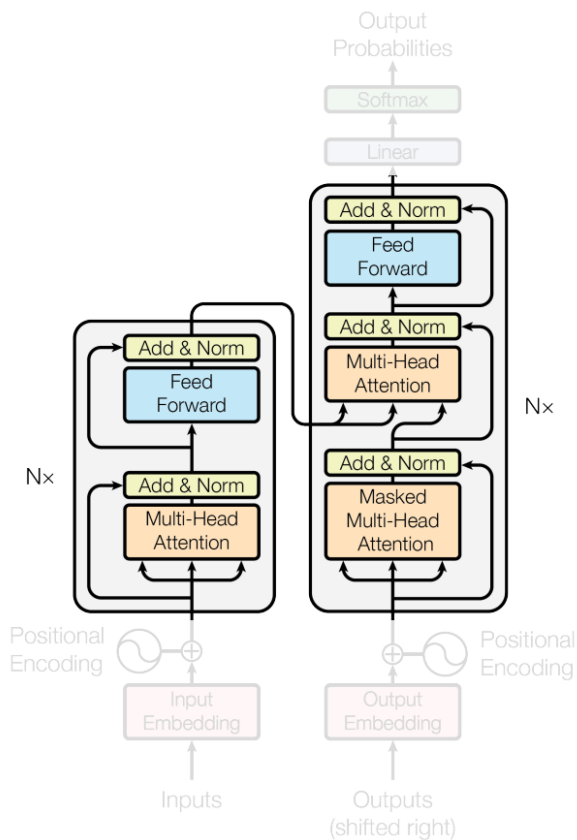
We now have all required inputs for our transformer.

Consult the documentation for the `torch.nn.Transformer` class of PyTorch. This class implements a full Transformer as described in "[Attention Is All You Need](#)", the paper that introduced this architecture.

The `Transformer` class implements the main part of the of the Transformer architecture, shown highlighted in the image on the left (see also Fig. 1 in "[Attention Is All You Need](#)").

For a given input sequence, it applies one or more encoder layers, followed by one or more decoder layers, to compute an output sequence that we can then process further.

Because the `Transformer` class takes care of most of the complicated parts of the model, we can concentrate providing the inputs and outputs: the grayed-out



areas in the image.

Check out the parameters for the `Transformer` class and the inputs and outputs of its `forward` function.

(a) Which parameter of the `Transformer` class should we base on our embedding? (1 point)

TODO: Your answer here.

(b) Given fixed input and output dimensions, which parameters of the `Transformer` can we use to change the complexity of our network? (1 point)

TODO: Your answer here.

(c) When using the `Transformer` class, where should we use the masks that we defined earlier? (1 point)

TODO: Your answer here.

Building a network

(d) Complete the code for the `TransformerNetwork`. (5 points)

Construct a network with the following architecture (see the image in the previous section for an overview):

1. An embedding layer that embeds the input tokens into a space of size `dim_hidden`.

2. A dropout layer (not shown in the image).
3. A **Transformer** with the specified parameters (`dim_hidden` , `num_heads` , `num_layers` , `dim_feedforward` , and `dropout`).
Note: you will need to pass `batch_first=True` , to indicate that the first dimension runs over the batch and not over the sequence.
4. A final linear prediction layer that takes the output of the transformer to `dim_vocab` possible classes.

Don't worry about positional encoding for now, we will add that later.

The `forward` function should generate the appropriate masks and combine the layers defined in `__init__` to compute the output.

```
In [ ]: class TransformerNetwork(torch.nn.Module):
    def __init__(self,
                  dim_vocab=len(vocab), padding_token=vocab['<pad>'],
                  num_layers=2, num_heads=4, dim_hidden=64, dim_feedforward=64,
                  dropout=0.01, positional_encoding=False):
        super().__init__()
        self.padding_token = padding_token
        # TODO: Your code here.
        self.embedding = ...
        self.dropout = ...
        self.transformer = ...
        self.predict = ...
        if positional_encoding:
            self.pos_encoding = ... # Fill this in later
        else:
            self.pos_encoding = torch.nn.Identity()

    def forward(self, src, tgt):
        # TODO: Your code here.
        # Combine self.embedding, self.dropout, self.transformer, self.predict
```

(e) Try the transformer with an example batch.

```
In [ ]: net = TransformerNetwork(dim_feedforward=72)
x, y = next(iter(train_loader))
bos = torch.tensor(vocab['<bos>']).expand(y.shape[0], 1)
y_prev = torch.cat((bos, y[:, :-1]), axis=1)

print('x.shape', x.shape)
print('y.shape', y.shape)
print('y_prev.shape', y_prev.shape)

y_pred = net(x, y_prev)
print('y_pred.shape', y_pred.shape)

# check the shape against what we expected
np.testing.assert_equal(list(y_pred.shape), [y.shape[0], y.shape[1], len(vocab)])
```

We can convert these predictions to tokens (but they're obviously random):

```
In [ ]: print(decode_tokens(torch.argmax(y_pred, dim=2))[:5])
```

```
In [ ]: # Check that the transformer is defined correctly
assert isinstance(net.embedding, torch.nn.Embedding)
assert isinstance(net.dropout, torch.nn.Dropout)
assert isinstance(net.transformer, torch.nn.Transformer)
assert isinstance(net.predict, torch.nn.Linear)
# Check parameters of transformer
assert net.transformer.d_model == 64
assert net.transformer.nhead == 4
assert net.transformer.batch_first == True
assert net.transformer.encoder.num_layers == 2
assert net.transformer.decoder.num_layers == 2
assert net.transformer.encoder.layers[0].linear1.out_features == 72
assert net.dropout.p == 0.01
assert net.transformer.encoder.layers[0].dropout.p == 0.01
# Check that the forward function behaves correctly
net.train(False)
assert torch.all(torch.isclose( \
    net(x, y_prev), \
    net(torch.cat((x, torch.tensor(vocab['<pad>'])).expand(x.shape[0],
    "Adding padding to x should not affect the output of the network. Che
assert torch.all(torch.isclose( \
    net(x, y_prev), \
    net(x, torch.cat((y_prev, torch.tensor(vocab['<pad>'])).expand(y.s
    "Adding padding to y should not affect the output of the network. Che
assert torch.all(torch.isclose( \
    net(x, y_prev)[:,:2], \
    net(x, y_prev[:,:2]), atol=1e-5)), \
    "The presence of later tokens in y should not affect the output for e
assert torch.all(torch.isclose( \
    net(x, y_prev), \
    net(torch.flip(x, [1]), y_prev), atol=1e-5)), \
    "Order of x should not matter for a transformer network. Check src_ma
assert not torch.all(torch.isclose( \
    net(x, torch.flip(y_prev, [1])), \
    torch.flip(net(x, y_prev), [1]), atol=1e-5)), \
    "Order of y should matter for a transformer network. Check tgt_mask."
```

5.4 Training (10 points)

Training loop

We will base the training code on last week's code. A complication in computing the loss and accuracy are the padding tokens. So, before we work on the training loop itself, we need to update the `accuracy` function so it ignores these `<pad>` tokens. Let's do this in a generic way

(a) Copy the `accuracy` function from last week, and add a parameter `ignore_index`. The tokens with `y == ignore_index` should be ignored.

(1 point)

Hint: you can select elements from a tensor with

`some_tensor[include]` where `include` is a tensor of booleans.

```
In [ ]: def accuracy(y_hat, y, ignore_index=None):
        # TODO: Your code here.
        # Hint: See assignment 4.
```

```
In [ ]: # Test the accuracy function.
assert accuracy(torch.tensor([[1,0,0],[0.4,0.5,0.1],[0,1,0],[0.4,0.1,0.5]]),
                torch.tensor([[1,0,0],[0.4,0.5,0.1],[0,1,0],[0.4,0.1,0.5]])),
assert accuracy(torch.tensor([[1,0,0],[0.4,0.5,0.1],[0,1,0],[0.4,0.1,0.5]]),
                torch.tensor([[1,0,0],[0.4,0.5,0.1],[0,1,0],[0.4,0.1,0.5]])),
assert accuracy(torch.tensor([[1,0,0],[0.4,0.5,0.1],[0,1,0],[0.4,0.1,0.5]]),
                torch.tensor([[1,0,0],[0.4,0.5,0.1],[0,1,0],[0.4,0.1,0.5]])),
assert accuracy(torch.tensor([[1,0,0],[0.4,0.5,0.1],[0,1,0],[0.4,0.1,0.5]]),
                torch.tensor([[1,0,0],[0.4,0.5,0.1],[0,1,0],[0.4,0.1,0.5]])),
```

(b) Write a training loop for the transformer model.

(4 points)

See last week's assignment for inspiration. The code is mostly the same with the following changes:

- The cross-entropy loss function and accuracy should ignore all `<pad>` tokens. (Use `ignore_index`, see the [documentation of CrossEntropyLoss](#).)
- The network expects `y_prev` as an extra input.
- The output of the network contains a batch of N samples, with maximum length L, and gives logits over C classes, so it has size (N,L,C). But `CrossEntropyLoss` and `accuracy` expect a tensor of size (N,C,L). You can use `torch.Tensor.transpose` to change the output to the right shape.

```
In [ ]: def train(net, data_loaders, epochs=100, lr=0.001, device=device):
        """
        Trains the model net with data from the data_loaders['train'] and data_l
        """
        # TODO: Your code here.
        # Hint: See assignment 4.
```

Experiment

(c) Train a transformer network. Use 100 epochs with a learning of 0.001

(no points)

```
In [ ]: # TODO: your answer here
```

(d) Briefly discuss the results. Has the training converged? Is this a good calculator?

(1 point)

TODO: Your answer here.

(e) Run the trained network with input "123+123" and "321+321" .

(1 point)

```
In [ ]: def predict(net, q, a):
        # Run net to predict the output given the input `q` and y_prev based on
        # Return predicted y
        with torch.no_grad():
            # TODO: Your code here.
            pass

        for src, tgt in [('123+123', '246'), ('321+321', '642')]:
            print(f'For {src}={tgt}')
            y_pred = predict(net, src, tgt)
            print(' y_pred[0]', y_pred[0])
            print(' encoded', torch.argmax(y_pred, dim=-1))
            print(' tokens', decode_tokens(torch.argmax(y_pred, dim=-1)))
            print()
```

(f) Compare the predictions for the first element of y with the two different inputs. Can you explain what happens?

(1 point)

TODO: Your answer here.

(g) Does the validation accuracy estimate how often the model is able to answers formulas correctly? Explain your answer.

(1 point)

TODO: Your answer here.

(h) If the forward function takes the shifted output y_prev as input, how can we use it if we don't know the output yet?

(1 point)

TODO: Your answer here.

5.5 Positional encoding (5 points)

We did not yet include positional encoding in the network. PyTorch does not include such an encoder, so here we copied the code from the book (slightly modified):

```
In [ ]: class PositionalEncoding(nn.Module):
        """Positional encoding."""
        def __init__(self, num_hiddens, max_len=1000):
            super().__init__()
            # Create a long enough P
            self.P = torch.zeros((1, max_len, num_hiddens))
            X = torch.arange(max_len, dtype=torch.float32).reshape(
                -1, 1) / torch.pow(10000, torch.arange(
                    0, num_hiddens, 2, dtype=torch.float32) / num_hiddens)
            self.P[:, :, 0::2] = torch.sin(X)
            self.P[:, :, 1::2] = torch.cos(X)
```

```
def forward(self, X):  
    return X + self.P[:, :X.shape[1], :].to(X.device)
```

(a) Add positional encoding to the TransformerModel.

(point given in earlier question)

```
In [ ]: # TODO: See over there.
```

(b) Construct and train a network with positional encoding

(1 point)

```
In [ ]: # TODO: your answer here  
net_pos = ...
```

(c) How does the performance of a model with positional encoding compare to a model without?

(1 point)

TODO: Your answer here.

(d) Run the trained network with input "123+123" and "321+321" .

(no points)

```
In [ ]: # TODO: Your code here.
```

(e) Compare the predictions for the first element of y with what you found earlier. Can you explain what happens?

(1 point)

TODO: Your answer here.

(f) Explain in your own words why positional encoding is used in transformer networks.

(1 point)

TODO: Your answer here.

(g) Look at the learning curve. Can you suggest a way to improve the model?

(1 point)

TODO: Your answer here.

(h) Optional: if time permits, try to train an even better model

5.6 Predicting for new samples (5 points)

Predicting an output given a new sample requires an appropriate search algorithm (see [d2l chapter 10.8](#)). Here, we will implement the simplest form: a greedy search algorithm that selects the token with the highest probability at each time step.

(a) Describe this search strategy in pseudo-code. (1 point)

TODO: Your answer here.

(b) Implement a greedy search function to predict a sequence using net_pos . (2 points)

```
In [ ]: def predict_greedy(net, src, length):
        # predict an output sequence of the given (maximum) length given input s
        with torch.no_grad():
            # TODO: Your code here.
            pass

        predicted_sequence = predict_greedy(net_pos, '123+123', 6)
        print(decode_tokens(predicted_sequence))
```

(c) Does this search strategy give a high-quality prediction? Why, or why not? (1 point)

TODO: Your answer here.

(d) What alternative search strategy could we use to improve the predictions? Why would this help? (1 point)

TODO: Your answer here.

5.7 Discussion (4 points)

Last week, we looked at recurrent neural networks such as the LSTM. Both recurrent neural networks and transformers work with sequences, but in recent years the transformer has become more popular than the recurrent models.

(a) An advantage of transformers over recurrent neural is that they can be faster to train. Why is that? (1 point)

TODO: Your answer here.

(b) Does this advantage also hold when predicting outputs for new sequences? Why, or why not? (1 point)

TODO: Your answer here.

(c) Why is positional encoding often used in transformers, but not in convolutional or recurrent neural networks? (1 point)

TODO: Your answer here.

The structure of a recurrent neural network makes it very suitable for online predictions, such as real-time translation, because it only depends on prior inputs. You can design an architecture where the RNN produces an output token for every input token given to it, and it can produce that output without having to wait for the rest of the input.

Note: 'online' means producing outputs continuously as new input comes in, as opposed to collecting a full dataset and analyzing it afterwards, it has nothing to do with the internet.

(d) How would a transformer work in an online application? Do you need to change the architecture? (1 point)

TODO: Your answer here.

The end

Well done! Please double check the instructions at the top before you submit your results.

This assignment has 47 points.

Version 3c66915 / 2023-10-04