

CS187 Project Segment 4: Semantic Interpretation – Question Answering

December 15, 2021

```
[1]: # Please do not change this cell because some hidden tests might depend on it.
import os

# Otter grader does not handle ! commands well, so we define and use our
# own function to execute shell commands.
def shell(commands, warn=True):
    """Executes the string `commands` as a sequence of shell commands.

    Prints the result to stdout and returns the exit status.
    Provides a printed warning on non-zero exit status unless `warn`
    flag is unset.
    """
    file = os.popen(commands)
    print (file.read().rstrip('\n'))
    exit_status = file.close()
    if warn and exit_status != None:
        print(f"Completed with errors. Exit status: {exit_status}\n")
    return exit_status

shell("""
ls requirements.txt >/dev/null 2>&1
if [ ! $? = 0 ]; then
    rm -rf .tmp
    git clone https://github.com/cs187-2021/project4.git .tmp
    mv .tmp/requirements.txt ./
    rm -rf .tmp
fi
pip install -q -r requirements.txt
""")
```

```
[2]: # Initialize Otter
import otter
grader = otter.Notebook()
```

```
[3]: from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

```
[4]: %cd /content/drive/MyDrive/project4-dkhimey-1
```

/content/drive/MyDrive/project4-dkhimey-1

1 CS187

1.1 Project 4: Semantic Interpretation – Question Answering

The goal of semantic parsing is to convert natural language utterances to a meaning representation such as a *logical form* expression or a *SQL query*. In the previous project segment, you built a parsing system to reconstruct parse trees from the natural-language queries in the ATIS dataset. However, that only solves an intermediary task, not the end-user task of obtaining answers to the queries.

In this final project segment, you will go further, building a semantic parsing system to convert English queries to SQL queries, so that by consulting a database you will be able to answer those questions. You will implement both a rule-based approach and an end-to-end sequence-to-sequence (seq2seq) approach. Both algorithms come with their pros and cons, and by the end of this segment you should have a basic understanding of the characteristics of the two approaches.

1.2 Goals

1. Build a semantic parsing algorithm to convert text to SQL queries based on the syntactic parse trees from the last project.
2. Build an attention-based end-to-end seq2seq system to convert text to SQL.
3. Improve the attention-based end-to-end seq2seq system with self-attention to convert text to SQL.
4. Discuss the pros and cons of the rule-based system and the end-to-end system.
5. (Optional) Use the state-of-the-art pretrained transformers for text-to-SQL conversion.

This will be an extremely challenging project, so we recommend that you start early.

2 Setup

```
[5]: import copy
import datetime
import math
import re
import sys
import warnings

import wget
import nltk
```

```

import sqlite3
import torch
import torch.nn as nn
import torchtext.legacy as tt

from cryptography.fernet import Fernet
from func_timeout import func_set_timeout
from torch.nn.utils.rnn import pack_padded_sequence as pack
from torch.nn.utils.rnn import pad_packed_sequence as unpack
from tqdm import tqdm
from transformers import BartTokenizer, BartForConditionalGeneration

```

```

[6]: # Set random seeds
seed = 1234
torch.manual_seed(seed)
# Set timeout for executing SQL
TIMEOUT = 3 # seconds

# GPU check: Set runtime type to use GPU where available
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print (device)

```

cpu

```

[7]: ## Download needed scripts and data
os.makedirs('data', exist_ok=True)
os.makedirs('scripts', exist_ok=True)
source_url = "https://raw.githubusercontent.com/nlp-course/data/master"

# Grammar to augment for this segment
if not os.path.isfile('data/grammar'):
    wget.download(f"{source_url}/ATIS/grammar_distrib4.crypt", out="data/")

# Decrypt the grammar file
key = b'bfksTY2BJ5VKKK9xZb1PDDLgGkdu7KCDFYfVePSEfGY='
fernet = Fernet(key)
with open('./data/grammar_distrib4.crypt', 'rb') as f:
    restored = Fernet(key).decrypt(f.read())
with open('./data/grammar', 'wb') as f:
    f.write(restored)

# Download scripts and ATIS database
wget.download(f"{source_url}/scripts/trees/transform.py", out="scripts/")
wget.download(f"{source_url}/ATIS/atis_sqlite.db", out="data/")

```

```

[7]: 'data//atis_sqlite (1).db'

```

```
[8]: # Import downloaded scripts for parsing augmented grammars
sys.path.insert(1, './scripts')
import transform as xform
```

3 Semantically augmented grammars

In the first part of this project segment, you'll be implementing a rule-based system for semantic interpretation of sentences. Before jumping into using such a system on the ATIS dataset – we'll get to that soon enough – let's first work with some trivial examples to get things going.

The fundamental idea of rule-based semantic interpretation is the rule of compositionality, that *the meaning of a constituent is a function of the meanings of its immediate subconstituents and the syntactic rule that combined them*. This leads to an infrastructure for specifying semantic interpretation in which each syntactic rule in a grammar (in our case, a context-free grammar) is associated with a semantic rule that applies to the meanings associated with the elements on the right-hand side of the rule.

3.1 Example: arithmetic expressions

As a first example, let's consider an augmented grammar for arithmetic expressions, familiar from lab 3-1. We again use the function `xform.parse_augmented_grammar` to parse the augmented grammar. You can read more about it in the file `scripts/transform.py`.

```
[9]: arithmetic_grammar, arithmetic_augmentations = xform.parse_augmented_grammar(
    """
    ## Sample grammar for arithmetic expressions

    S -> NUM                                : lambda Num: Num
      / S OP S                              : lambda S1, Op, S2: Op(S1, S2)

    OP -> ADD                               : lambda Op: Op
      / SUB
      / MULT
      / DIV

    NUM -> 'zero'                           : lambda: 0
      / 'one'                               : lambda: 1
      / 'two'                               : lambda: 2
      / 'three'                             : lambda: 3
      / 'four'                              : lambda: 4
      / 'five'                              : lambda: 5
      / 'six'                               : lambda: 6
      / 'seven'                             : lambda: 7
      / 'eight'                             : lambda: 8
      / 'nine'                              : lambda: 9
      / 'ten'                               : lambda: 10
```

```

ADD -> 'plus' / 'added' 'to'      : lambda: lambda x, y: x + y
SUB -> 'minus'                    : lambda: lambda x, y: x - y
MULT -> 'times' / 'multiplied' 'by' : lambda: lambda x, y: x * y
DIV -> 'divided' 'by'             : lambda: lambda x, y: x / y
"""
)

```

Recall that in this grammar specification format, rules that are not explicitly provided with an augmentation (like all the OP rules after the first `OP -> ADD`) are associated with the textually most recent one (`lambda Op: Op`).

The `parse_augmented_grammar` function returns both an NLTK grammar and a dictionary that maps from productions in the grammar to their associated augmentations. Let's examine the returned grammar.

```

[10]: for production in arithmetic_grammar.productions():
      print(f"{repr(production):25} {arithmetic_augmentations[production]}")

```

```

S -> NUM                <function <lambda> at 0x7fb3e1ed1950>
S -> S OP S             <function <lambda> at 0x7fb3e1ed19e0>
OP -> ADD                <function <lambda> at 0x7fb3e1ed1b00>
OP -> SUB                <function <lambda> at 0x7fb3e1ed1c20>
OP -> MULT               <function <lambda> at 0x7fb3e1ed1d40>
OP -> DIV                <function <lambda> at 0x7fb3e1ed1e60>
NUM -> 'zero'            <function <lambda> at 0x7fb3e1ed1f80>
NUM -> 'one'             <function <lambda> at 0x7fb3e1ef40e0>
NUM -> 'two'             <function <lambda> at 0x7fb3e1ef4200>
NUM -> 'three'           <function <lambda> at 0x7fb3e1ef4320>
NUM -> 'four'            <function <lambda> at 0x7fb3e1ef4440>
NUM -> 'five'            <function <lambda> at 0x7fb3e1ef4560>
NUM -> 'six'             <function <lambda> at 0x7fb3e1ef4680>
NUM -> 'seven'           <function <lambda> at 0x7fb3e1ef47a0>
NUM -> 'eight'           <function <lambda> at 0x7fb3e1ef48c0>
NUM -> 'nine'            <function <lambda> at 0x7fb3e1ef49e0>
NUM -> 'ten'             <function <lambda> at 0x7fb3e1ef4b00>
ADD -> 'plus'            <function <lambda> at 0x7fb3e1ef4cb0>
ADD -> 'added' 'to'      <function <lambda> at 0x7fb3e1ef4e60>
SUB -> 'minus'           <function <lambda> at 0x7fb3e1ef3050>
MULT -> 'times'          <function <lambda> at 0x7fb3e1ef3200>
MULT -> 'multiplied' 'by' <function <lambda> at 0x7fb3e1ef33b0>
DIV -> 'divided' 'by'    <function <lambda> at 0x7fb3e1ef3560>

```

We can parse with the grammar using one of the built-in NLTK parsers.

```

[11]: arithmetic_parser = nltk.parse.BottomUpChartParser(arithmetic_grammar)
      parses = [p for p in arithmetic_parser.parse('three plus one times four'.
      ↪split())]
      for parse in parses:

```

```
parse.pretty_print()
```

```

      S
    -----|-----
      S      |      |
    -----|-----
  S      OP  S      OP  S
  |      |   |      |   |
NUM  ADD NUM MULT NUM
  |      |   |      |   |
three plus one times four

```

```

      S
    -----|-----
  |      |      S
  |      |      -----|-----
  S      OP  S      OP  S
  |      |   |      |   |
NUM  ADD NUM MULT NUM
  |      |   |      |   |
three plus one times four

```

Now let's turn to the augmentations. They can be arbitrary Python functions applied to the semantic representations associated with the right-hand-side nonterminals, returning the semantic representation of the left-hand side. To interpret the semantic representation of the entire sentence (at the root of the parse tree), we can use the following pseudo-code:

to interpret a tree:

interpret each of the nonterminal-rooted subtrees

find the augmentation associated with the root production of the tree

(it should be a function of as many arguments as there are nonterminals on the right-hand side)

return the result of applying the augmentation to the subtree values

(The base case of this recursion occurs when the number of nonterminal-rooted subtrees is zero, that is, a rule all of whose right-hand side elements are terminals.)

Suppose we had such a function, call it `interpret`. How would it operate on, for instance, the tree (S (S (NUM three)) (OP (ADD plus)) (S (NUM one)))?

```

interpret (S (S (NUM three)) (OP (ADD plus)) (S (NUM one)))
  |->interpret (S (NUM three))
    |      |->interpret (NUM three)
    |      |      |->(no subconstituents to evaluate)
    |      |      |->apply the augmentation for the rule NUM -> three to the empty set of values
    |      |      |      (lambda: 3) () ==> 3
    |      |      \==> 3
    |      |->apply the augmentation for the rule S -> NUM to the value 3
    |      |      (lambda NUM: NUM)(3) ==> 3
    |      \==> 3

```

```

|->interpret (OP (ADD plus))
|      |...
|      \==> lambda x, y: x + y
|->interpret (S (NUM one))
|      |...
|      \==> 1
|->apply the augmentation for the rule S -> S OP S to the values 3, (lambda x, y: x + y), a
|      (lambda S1, Op, S2: Op(S1, S2))(3, (lambda x, y: x + y), 1) ==> 4
\==> 4

```

Thus, the string “three plus one” is semantically interpreted as the value 4.

We provide the `interpret` function to carry out this recursive process, copied over from lab 4-2:

```

[12]: def interpret(tree, augmentations):
      syntactic_rule = tree.productions()[0]
      semantic_rule = augmentations[syntactic_rule]
      child_meanings = [interpret(child, augmentations)
                        for child in tree
                        if isinstance(child, nltk.Tree)]
      return semantic_rule(*child_meanings)

```

Now we should be able to evaluate the arithmetic example from above.

```

[13]: interpret(parses[0], arithmetic_augmentations)

```

[13]: 16

And we can even write a function that parses and interprets a string. We’ll have it evaluate each of the possible parses and print the results.

```

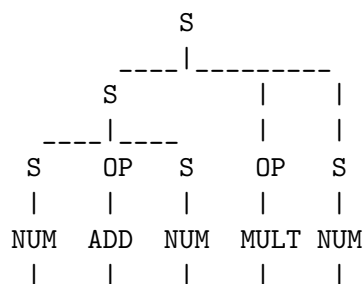
[14]: def parse_and_interpret(string, grammar, augmentations):
      parser = nltk.parse.BottomUpChartParser(grammar)
      parses = parser.parse(string.split())
      for parse in parses:
          parse.pretty_print()
          print(parse, "==>", interpret(parse, augmentations))

```

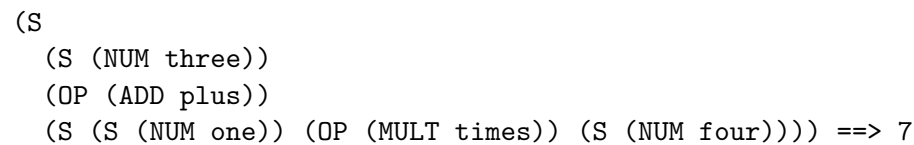
```

[15]: parse_and_interpret("three plus one times four", arithmetic_grammar,
      ↪arithmetic_augmentations)

```



```
(S
  (S (S (NUM three)) (OP (ADD plus)) (S (NUM one)))
  (OP (MULT times))
  (S (NUM four))) ==> 16
```



3.2 Some grammar specification conveniences

```
[16]: def constant(value):  
        """Return `value`, ignoring any arguments"""  
        return lambda *args: value
```

```
[17]: def first(*args):
      """Return the value of the first (and perhaps only) subconstituent,
      ignoring any others"""
      return args[0]
```

In the call to `parse_augmented_grammar` below, we pass in the global environment, extracted via a `globals()` function call, via the named argument `globals`. This allows the `parse_augmented_grammar` function to make use of the global bindings for `constant`, `first`, and the like when evaluating the augmentation expressions to their

values. You can check out the code in `transform.py` to see how the passed in `globals` bindings are used. To help understand what's going on, see what happens if you don't include the `globals=globals()`.

```
[18]: arithmetic_grammar_2, arithmetic_augmentations_2 = xform.  
      ↪ parse_augmented_grammar(  
          """  
          ## Sample grammar for arithmetic expressions  
  
          S -> NUM                                : first  
            / S OP S                              : lambda S1, Op, S2: Op(S1, S2)  
  
          OP -> ADD                                : first  
            / SUB  
            / MULT  
            / DIV  
  
          NUM -> 'zero'                            : constant(0)  
            / 'one'                                : constant(1)  
            / 'two'                                : constant(2)  
            / 'three'                              : constant(3)  
            / 'four'                               : constant(4)  
            / 'five'                               : constant(5)  
            / 'six'                                : constant(6)  
            / 'seven'                              : constant(7)  
            / 'eight'                              : constant(8)  
            / 'nine'                               : constant(9)  
            / 'ten'                                : constant(10)  
  
          ADD -> 'plus' / 'added' / 'to'           : constant(lambda x, y: x + y)  
          SUB -> 'minus'                           : constant(lambda x, y: x - y)  
          MULT -> 'times' / 'multiplied' / 'by'    : constant(lambda x, y: x * y)  
          DIV -> 'divided' / 'by'                  : constant(lambda x, y: x / y)  
          """,  
          globals=globals())
```

Finally, it might make our lives easier to write a template of augmentations whose instantiation depends on the right-hand side of the rule.

We use a reserved keyword `_RHS` to denote the right-hand side of the syntactic rule, which will be replaced by a **list** of the right-hand-side strings. For example, an augmentation `numeric_template(_RHS)` would be as if written as `numeric_template(['zero'])` when the rule is `NUM -> 'zero'`, and `numeric_template(['one'])` when the rule is `NUM -> 'one'`. The details of how this works can be found at [scripts/transform.py](#).

This would allow us to use a single template function, for example,

```
[19]: def numeric_template(rhs):
      """Ignore the subphrase meanings and lookup the first right-hand-side symbol
         as a number"""
      return constant({'zero':0, 'one':1, 'two':2, 'three':3, 'four':4, 'five':5,
                       'six':6, 'seven':7, 'eight':8, 'nine':9, 'ten':10}[rhs[0]])
```

and then further simplify the grammar specification:

```
[20]: arithmetic_grammar_3, arithmetic_augmentations_3 = xform.
      ↪ parse_augmented_grammar(
          """
          ## Sample grammar for arithmetic expressions

          S -> NUM : first
            / S OP S : lambda S1, Op, S2: Op(S1, S2)

          OP -> ADD : first
            / SUB
            / MULT
            / DIV

          NUM -> 'zero' / 'one' / 'two' : numeric_template(_RHS)
            / 'three' / 'four' / 'five'
            / 'six' / 'seven' / 'eight'
            / 'nine' / 'ten'

          ADD -> 'plus' / 'added' 'to' : constant(lambda x, y: x + y)
          SUB -> 'minus' : constant(lambda x, y: x - y)
          MULT -> 'times' / 'multiplied' 'by' : constant(lambda x, y: x * y)
          DIV -> 'divided' 'by' : constant(lambda x, y: x / y)
          """,
          globals=globals())
```

```
[21]: parse_and_interpret("six divided by three", arithmetic_grammar_3,
      ↪ arithmetic_augmentations_3)
```

```

      S
    -----|-----
    S      OP      S
    |      |      |
  NUM    DIV    NUM
    |      |      |
six divided      by three
```

```
(S (S (NUM six)) (OP (DIV divided by)) (S (NUM three))) ==> 2.0
```

3.3 Example: *Green Eggs and Ham* revisited

This stuff is tricky, so it's useful to see more examples before jumping in the deep end. In this simple GEaH fragment grammar, we use a larger set of auxiliary functions to build the augmentations.

```
[22]: def forward(F, A):
      """Forward application: Return the application of the first
          argument to the second"""
      return F(A)

      def backward(A, F):
          """Backward application: Return the application of the second
              argument to the first"""
          return F(A)

      def second(*args):
          """Return the value of the second subconstituent, ignoring any others"""
          return args[1]

      def ignore(*args):
          """Return `None`, ignoring everything about the constituent. (Good as a
              placeholder until a better augmentation can be devised.)"""
          return None
```

Using these, we can build and test the grammar.

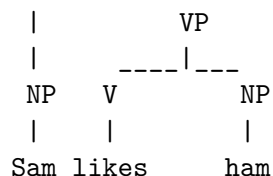
```
[23]: geah_grammar_spec = """
      ## Productions
      S -> NP VP          : backward
      VP -> V NP           : forward

      ## Lexicon
      V -> 'likes'         : constant(lambda Object: lambda Subject: λ
      ↪f"like({Subject}, {Object})")
      NP -> 'Sam' | 'sam'   : constant(_RHS[0])
      NP -> 'ham'
      NP -> 'eggs'
      """
```

```
[24]: geah_grammar, geah_augmentations = xform.
      ↪parse_augmented_grammar(geah_grammar_spec,
      ↪globals=globals())
```

```
[25]: parse_and_interpret("Sam likes ham", geah_grammar, geah_augmentations)
```

```
      S
    ____|____
```



(S (NP Sam) (VP (V likes) (NP ham))) ==> like(Sam, ham)

4 Semantics of ATIS queries

Now you're in a good position to understand and add augmentations to a more comprehensive grammar, say, one that parses ATIS queries and generates SQL queries.

In preparation for that, we need to load the ATIS data, both NL and SQL queries.

4.1 Loading and preprocessing the corpus

To simplify things a bit, we'll only consider ATIS queries whose question type (remember that from project segment 1?) is `flight_id`. We download training, development, and test splits for this subset of the ATIS corpus, including corresponding SQL queries.

```
[26]: # Acquire the datasets - training, development, and test splits of the
# ATIS queries and corresponding SQL queries
wget.download(f"{source_url}/ATIS/test_flightid.nl", out="data/")
wget.download(f"{source_url}/ATIS/test_flightid.sql", out="data/")
wget.download(f"{source_url}/ATIS/dev_flightid.nl", out="data/")
wget.download(f"{source_url}/ATIS/dev_flightid.sql", out="data/")
wget.download(f"{source_url}/ATIS/train_flightid.nl", out="data/")
wget.download(f"{source_url}/ATIS/train_flightid.sql", out="data/")
```

```
[26]: 'data//train_flightid (1).sql'
```

Let's take a look at the data: the NL queries are in `.nl` files, and the SQL queries are in `.sql` files.

```
[27]: shell("head -1 data/dev_flightid.nl")
shell("head -1 data/dev_flightid.sql")
```

```

what flights are available tomorrow from denver to philadelphia
SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service
airport_service_1 , city city_1 , airport_service airport_service_2 , city
city_2 , days days_1 , date_day date_day_1 WHERE flight_1.from_airport =
airport_service_1.airport_code AND airport_service_1.city_code =
city_1.city_code AND city_1.city_name = 'DENVER' AND ( flight_1.to_airport =
airport_service_2.airport_code AND airport_service_2.city_code =
city_2.city_code AND city_2.city_name = 'PHILADELPHIA' AND flight_1.flight_days
= days_1.days_code AND days_1.day_name = date_day_1.day_name AND date_day_1.year
= 1991 AND date_day_1.month_number = 1 AND date_day_1.day_number = 20 )

```

4.2 Corpus preprocessing

We'll use torchtext to process the data. We use two Fields: SRC for the questions, and TGT for the SQL queries. We'll use the tokenizer from project segment 3.

```
[28]: ## Tokenizer
tokenizer = nltk.tokenize.RegexpTokenizer('\d+|st\.|[\w-]+|\$[\d\.]+|\S+')
def tokenize(string):
    return tokenizer.tokenize(string.lower())

## Demonstrating the tokenizer
## Note especially the handling of "11pm" and hyphenated words.
print(tokenize("Are there any first-class flights from St. Louis at 11pm for
→less than $3.50?"))
```

```
['are', 'there', 'any', 'first-class', 'flights', 'from', 'st.', 'louis', 'at',
'11', 'pm', 'for', 'less', 'than', '$3.50', '?']
```

```
[29]: SRC = tt.data.Field(include_lengths=True,           # include lengths
                        batch_first=False,             # batches will be max_len x
                        →batch_size
                        tokenize=tokenize,             # use our tokenizer
                        )
TGT = tt.data.Field(include_lengths=False,
                    batch_first=False,                 # batches will be max_len x
                    →batch_size
                    tokenize=lambda x: x.split(),      # use split to tokenize
                    init_token="<bos>",                # prepend <bos>
                    eos_token="<eos>",                 # append <eos>
                    )
fields = [('src', SRC), ('tgt', TGT)]
```

Note that we specified `batch_first=False` (as in lab 4-4), so that the returned batched tensors would be of size `max_length x batch_size`, which facilitates seq2seq implementation.

Now, we load the data using torchtext. We use the TranslationDataset class here because our task is essentially a translation task: “translating” questions into the corresponding SQL queries. Therefore, we also refer to the questions as the *source* side (SRC) and the SQL queries as the *target* side (TGT).

```
[30]: # Make splits for data
train_data, val_data, test_data = tt.datasets.TranslationDataset.splits(
    ('_flightid.nl', '_flightid.sql'), fields, path='./data/',
    train='train', validation='dev', test='test')

MIN_FREQ = 3
SRC.build_vocab(train_data.src, min_freq=MIN_FREQ)
TGT.build_vocab(train_data.tgt, min_freq=MIN_FREQ)
```

```

print (f"Size of English vocab: {len(SRC.vocab)}")
print (f"Most common English words: {SRC.vocab.freqs.most_common(10)}\n")

print (f"Size of SQL vocab: {len(TGT.vocab)}")
print (f"Most common SQL words: {TGT.vocab.freqs.most_common(10)}\n")

print (f"Index for start of sequence token: {TGT.vocab.stoi[TGT.init_token]}")
print (f"Index for end of sequence token: {TGT.vocab.stoi[TGT.eos_token]}")

```

Size of English vocab: 421

Most common English words: [('to', 3478), ('from', 3019), ('flights', 2094), ('the', 1550), ('on', 1230), ('me', 973), ('flight', 972), ('show', 845), ('what', 833), ('boston', 813)]

Size of SQL vocab: 392

Most common SQL words: [('=', 38876), ('AND', 36564), ('.', 22772), ('airport_service', 8314), ('city', 8313), ('(', 6432), (')', 6432), ('flight_1.flight_id', 4536), ('flight', 4221), ('SELECT', 4178)]

Index for start of sequence token: 2

Index for end of sequence token: 3

Next, we batch our data to facilitate processing on a GPU. Batching is a bit tricky because the source and target will typically be of different lengths. Fortunately, `torchtext` allows us to pass in a `sort_key` function. By sorting on length, we can minimize the amount of padding on the source side, but since there is still some padding, we need to handle them with `pack` and `unpack` later on in the `seq2seq` part (as in lab 4-5).

```

[31]: BATCH_SIZE = 16 # batch size for training/validation
TEST_BATCH_SIZE = 1 # batch size for test, we use 1 to make beam search
    ↪ implementation easier

train_iter, val_iter = tt.data.BucketIterator.splits((train_data, val_data),
                                                    batch_size=BATCH_SIZE,
                                                    device=device,
                                                    repeat=False,
                                                    sort_key=lambda x: len(x.
    ↪src),
                                                    sort_within_batch=True)

test_iter = tt.data.BucketIterator(test_data,
                                   batch_size=TEST_BATCH_SIZE,
                                   device=device,
                                   repeat=False,
                                   sort=False,
                                   train=False)

```

Let's look at a single batch from one of these iterators.

```
[32]: batch = next(iter(train_iter))
train_batch_text, train_batch_text_lengths = batch.src
print (f"Size of text batch: {train_batch_text.shape}")
print (f"Third sentence in batch: {train_batch_text[:, 2]}")
print (f"Length of the third sentence in batch: {train_batch_text_lengths[2]}")
print (f"Converted back to string: {' '.join([SRC.vocab.itos[i] for i in
↳ train_batch_text[:, 2]])}")

train_batch_sql = batch.tgt
print (f"Size of sql batch: {train_batch_sql.shape}")
print (f"Third SQL in batch: {train_batch_sql[:, 2]}")
print (f"Converted back to string: {' '.join([TGT.vocab.itos[i] for i in
↳ train_batch_sql[:, 2]])}")
```

[illegible]

Alternatively, we can directly iterate over the raw examples:

```
[33]: for example in train_iter.dataset[:1]:
      train_text_1 = ' '.join(example.src) # detokenized question
      train_sql_1 = ' '.join(example.tgt) # detokenized sql
      print (f"Question: {train_text_1}\n")
      print (f"SQL: {train_sql_1}")
```

Question: list all the flights that arrive at general mitchell international from various cities

```
SQL: SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport airport_1
, airport_service airport_service_1 , city city_1 WHERE flight_1.to_airport =
airport_1.airport_code AND airport_1.airport_code = 'MKE' AND
flight_1.from_airport = airport_service_1.airport_code AND
airport_service_1.city_code = city_1.city_code AND 1 = 1
```

4.3 Establishing a SQL database for evaluating ATIS queries

The output of our systems will be SQL queries. How should we determine if the generated queries are correct? We can't merely compare against the gold SQL queries, since there are many ways to implement a SQL query that answers any given NL query.

Instead, we will execute the queries – both the predicted SQL query and the gold SQL query – on an actual database, and verify that the returned responses are the same. For that purpose, we need a SQL database server to use. We'll set one up here, using the [Python sqlite3 module](#).

```
[34]: @func_set_timeout(TIMEOUT)
def execute_sql(sql):
    conn = sqlite3.connect('data/atis_sqlite.db') # establish the DB based on
    ↳ the downloaded data
    c = conn.cursor() # build a "cursor"
    c.execute(sql)
    results = list(c.fetchall())
    c.close()
    conn.close()
    return results
```

To run a query, we use the cursor's `execute` function, and retrieve the results with `fetchall`. Let's get all the flights that arrive at General Mitchell International – the query `train_sql_1` above. There's a lot, so we'll just print out the first few.

```
[35]: predicted_ret = execute_sql(train_sql_1)

print(f"""
Executing: {train_sql_1}

Result: {len(predicted_ret)} entries starting with

{predicted_ret[:10]}
""")
```

```
Executing: SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport
airport_1 , airport_service airport_service_1 , city city_1 WHERE
flight_1.to_airport = airport_1.airport_code AND airport_1.airport_code = 'MKE'
AND flight_1.from_airport = airport_service_1.airport_code AND
```



```
airport_service_1.city_code = city_1.city_code AND 1 = 1
```

Result: 534 entries starting with

```
[(107929,), (107930,), (107931,), (107932,), (107933,), (107934,), (107935,),  
(107936,), (107937,), (107938,)]
```

For your reference, the SQL database we are using has a database schema described at <https://github.com/jkkummerfeld/text2sql-data/blob/master/data/atis-schema.csv>, and is consistent with the SQL queries provided in the various `.sql` files loaded above.

5 Rule-based parsing and interpretation of ATIS queries

First, you will implement a rule-based semantic parser using a grammar like the one you completed in the third project segment. We've placed an initial grammar in the file `data/grammar`. In addition to the helper functions defined above (`constant`, `first`, etc.), it makes use of some other simple functions. We've included those below, but you can (and almost certainly should) augment this set with others that you define as you build out the full set of augmentations.

```
[36]: def upper(term):  
    return ''' + term.upper() + '''  
  
def weekday(day):  
    return f"flight.flight_days IN (SELECT days.days_code FROM days WHERE days.  
    ↳day_name = '{day.upper()}')"  
  
def month_name(month):  
    return {'JANUARY' : 1,  
            'FEBRUARY' : 2,  
            'MARCH' : 3,  
            'APRIL' : 4,  
            'MAY' : 5,  
            'JUNE' : 6,  
            'JULY' : 7,  
            'AUGUST' : 8,  
            'SEPTEMBER' : 9,  
            'OCTOBER' : 10,  
            'NOVEMBER' : 11,  
            'DECEMBER' : 12}[month.upper()]  
  
def airports_from_airport_name(airport_name):  
    return f"(SELECT airport.airport_code FROM airport WHERE airport.airport_name_  
    ↳= {upper(airport_name)})"  
  
def airports_from_city(city):  
    return f"""
```

```

        (SELECT airport_service.airport_code FROM airport_service WHERE_
↪airport_service.city_code IN
        (SELECT city.city_code FROM city WHERE city.city_name = {upper(city)}))
        """

def null_condition(*args, **kwargs):
    return 1

def depart_around(time):
    return f"""
        flight.departure_time >= {add_delta(miltime(time), -15).strftime('%H%M')}
        AND flight.departure_time <= {add_delta(miltime(time), 15).strftime('%H%M')}
        """.strip()

def add_delta(tme, delta):
    # transform to a full datetime first
    return (datetime.datetime.combine(datetime.date.today(), tme) +
            datetime.timedelta(minutes=delta)).time()

def miltime(minutes):
    return datetime.time(hour=int(minutes/100), minute=(minutes % 100))

```

We can build a parser with the augmented grammar:

```

[37]: atis_grammar, atis_augmentations = xform.read_augmented_grammar('data/grammar',_
↪globals=globals())
atis_parser = nltk.parse.BottomUpChartParser(atis_grammar)

```

We'll define a function to return a parse tree for a string according to the ATIS grammar (if available).

```

[38]: def parse_tree(sentence):
        """Parse a sentence and return the parse tree, or None if failure."""
        try:
            parses = list(atis_parser.parse(tokenize(sentence)))
            if len(parses) == 0:
                return None
            else:
                return parses[0]
        except:
            return None

```

We can check the overall coverage of this grammar on the training set by using the `parse_tree` function to determine if a parse is available. The grammar that we provide should get about a 40% coverage of the training set.

```
[39]: # Check coverage on training set
parsed = 0
with open("data/train_flightid.nl") as train:
    examples = train.readlines()[:]
    for sentence in tqdm(examples):
        if parse_tree(sentence):
            parsed += 1
        else:
            next

print(f"\nParsed {parsed} of {len(examples)} ({parsed*100/(len(examples)):.
    ↳2f}%)")
```

```
100%|          | 3651/3651 [00:21<00:00, 172.94it/s]
```

```
Parsed 1525 of 3651 (41.77%)
```

5.1 Goal 1: Construct SQL queries from a parse tree and evaluate the results

It's time to turn to the first major part of this project segment, implementing a rule-based semantic parsing system to answer flight-ID-type ATIS queries.

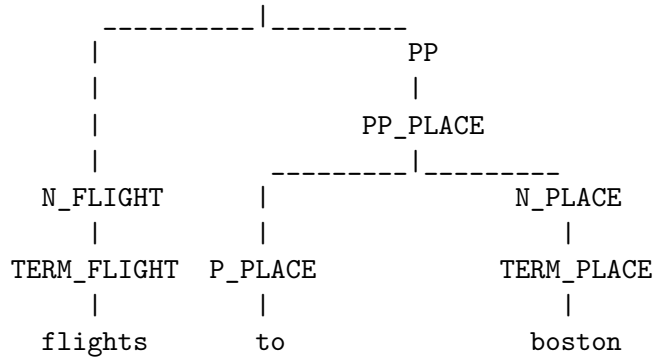
Recall that in rule-based semantic parsing, each syntactic rule is associated with a semantic composition rule. The grammar we've provided has semantic augmentations for some of the low-level phrases – cities, airports, times, airlines – but not the higher level syntactic types. You'll be adding those.

In the ATIS grammar that we provide, as with the earlier toy grammars, the augmentation for a rule with n nonterminals and m terminals on the right-hand side is assumed to be called with n positional arguments (the values for the corresponding children). The `interpret` function you've already defined should therefore work well with this grammar.

Let's run through one way that a semantic derivation might proceed, for the sample query "flights to boston":

```
[40]: sample_query = "flights to boston"
print(tokenize(sample_query))
sample_tree = parse_tree(sample_query)
sample_tree.pretty_print()
```

```
['flights', 'to', 'boston']
      S
      |
    NP_FLIGHT
      |
    NOM_FLIGHT
      |
    N_FLIGHT
```



```
(SELECT    airport_service.airport_code    FROM    airport_service    WHERE    air-
port_service.city_code IN (SELECT city.city_code FROM city WHERE city.city_name =
{upper(city)}))
```

Given a sentence, we first construct its parse tree using the syntactic rules, then compose the corresponding semantic rules bottom-up, until eventually we arrive at the root node with a finished SQL statement. For this query, we will go through what the possible meaning representations for the subconstituents of “flights to boston” might be. But this is just one way of doing things; other ways are possible, and you should feel free to experiment.

Working from bottom up:

1. The `TERM_PLACE` phrase “boston” uses the composition function template `constant(airports_from_city(' '.join(_RHS)))`, which will be instantiated as `constant(airports_from_city(' '.join(['boston'])))` (recall that `_RHS` is replaced by the right-hand side of the rule). The meaning of `TERM_PLACE` will be the SQL snippet

```
SELECT airport_service.airport_code
FROM airport_service
WHERE airport_service.city_code IN
  (SELECT city.city_code
   FROM city
   WHERE city.city_name = "BOSTON")
```

(This query generates a list of all of the airports in Boston.)

2. The `N_PLACE` phrase “boston” can have the same meaning as the `TERM_PLACE`.
3. The `P_PLACE` phrase “to” might be associated with a function that maps a SQL query for a list of airports to a SQL condition that holds of flights that go to one of those airports, i.e., `flight.to_airport IN (...)`.
4. The `PP_PLACE` phrase “to boston” might apply the `P_PLACE` meaning to the `TERM_PLACE` meaning, thus generating a SQL condition that holds of flights that go to one of the Boston airports:

```
flight.to_airport IN
  (SELECT airport_service.airport_code
   FROM airport_service
   WHERE airport_service.city_code IN
     (SELECT city.city_code
```

```
FROM city
WHERE city.city_name = "BOSTON")
```

5. The PP phrase “to Boston” can again get its meaning from the PP_PLACE.
6. The TERM_FLIGHT phrase “flights” might also return a condition on flights, this time the “null condition”, represented by the SQL truth value 1. Ditto for the N_FLIGHT phrase “flights”.
7. The N_FLIGHT phrase “flights to boston” can conjoin the two conditions, yielding the SQL condition

```
flight.to_airport IN
(SELECT airport_service.airport_code
FROM airport_service
WHERE airport_service.city_code IN
(SELECT city.city_code
FROM city
WHERE city.city_name = "BOSTON")
AND 1
```

which can be inherited by the NOM_FLIGHT and NP_FLIGHT phrases.

8. The S phrase “flights to boston” can use the condition provided by the NP_FLIGHT phrase to select all flights satisfying the condition with a SQL query like

```
SELECT DISTINCT flight.flight_id
FROM flight
WHERE flight.to_airport IN
(SELECT airport_service.airport_code
FROM airport_service
WHERE airport_service.city_code IN
(SELECT city.city_code
FROM city
WHERE city.city_name = "BOSTON")
AND 1
```

This SQL query is then taken to be a representation of the meaning for the NL query “flights to boston”, and can be executed against the ATIS database to retrieve the requested flights.

Now, it’s your turn to add augmentations to `data/grammar` to make this example work. The augmentations that we have provided for the grammar make use of a set of auxiliary functions that we defined above. You should feel free to add your own auxiliary functions that you make use of in the grammar.

```
[41]: #TODO: add augmentations to `data/grammar` to make this example work
atis_grammar, atis_augmentations = xform.read_augmented_grammar('data/grammar',
↪globals=globals())
atis_parser = nltk.parse.BottomUpChartParser(atis_grammar)
# print("interpreting...")
predicted_sql = interpret(sample_tree, atis_augmentations)
print("Predicted SQL:\n\n", predicted_sql, "\n")
```

Predicted SQL:

```
SELECT DISTINCT flight.flight_id FROM flight WHERE 1 AND flight.to_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE
airport_service.city_code IN
    (SELECT city.city_code FROM city WHERE city.city_name = "BOSTON"))
```

Verification on some examples With a rule-based semantic parsing system, we can generate SQL queries given questions, and then execute those queries on a SQL database to answer the given questions. To evaluate the performance of the system, we compare the returned results against the results of executing the ground truth queries.

We provide a function `verify` to compare the results from our generated SQL to the ground truth SQL. It should be useful for testing individual queries.

```
[42]: def verify(predicted_sql, gold_sql, silent=True):
    """
    Compare the correctness of the generated SQL by executing on the
    ATIS database and comparing the returned results.
    Arguments:
        predicted_sql: the predicted SQL query
        gold_sql: the reference SQL query to compare against
        silent: print outputs or not
    Returns: True if the returned results are the same, otherwise False
    """
    # Execute predicted SQL
    try:
        predicted_result = execute_sql(predicted_sql)
    except BaseException as e:
        if not silent:
            print(f"predicted sql exec failed: {e}")
        return False
    if not silent:
        print("Predicted DB result:\n\n", predicted_result[:10], "\n")

    # Execute gold SQL
    try:
        gold_result = execute_sql(gold_sql)
    except BaseException as e:
        if not silent:
            print(f"gold sql exec failed: {e}")
        return False
    if not silent:
        print("Gold DB result:\n\n", gold_result[:10], "\n")

    # Verify correctness
```

```

if gold_result == predicted_result:
    return True

```

Let's try this methodology on a simple example: "flights from phoenix to milwaukee". we provide it along with the gold SQL query.

```

[43]: def rule_based_trial(sentence, gold_sql):
    print("Sentence: ", sentence, "\n")
    tree = parse_tree(sentence)
    print("Parse:\n\n")
    tree.pretty_print()

    predicted_sql = interpret(tree, atis_augmentations)
    print("Predicted SQL:\n\n", predicted_sql, "\n")

    if verify(predicted_sql, gold_sql, silent=False):
        print('Correct!')
    else:
        print('Incorrect!')

[44]: # Run this cell to reload augmentations after you make changes to `data/grammar`
atis_grammar, atis_augmentations = xform.read_augmented_grammar('data/grammar',
↪globals=globals())
atis_parser = nltk.parse.BottomUpChartParser(atis_grammar)

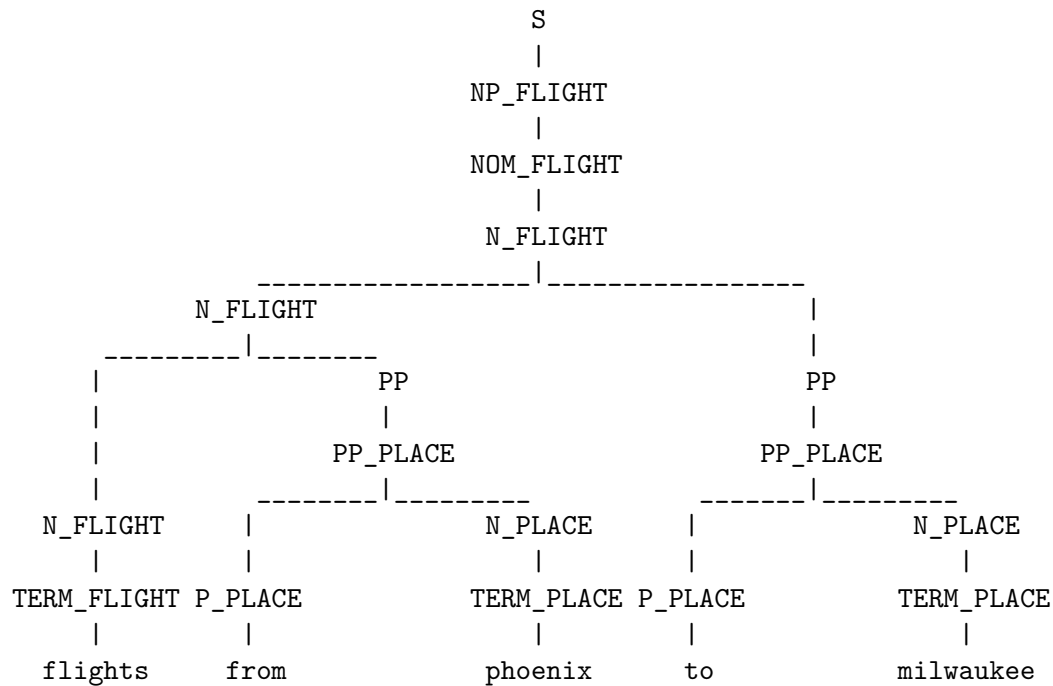
[45]: #TODO: add augmentations to `data/grammar` to make this example work
# Example 1
example_1 = 'flights from phoenix to milwaukee'
gold_sql_1 = """
SELECT DISTINCT flight_1.flight_id
FROM flight flight_1 ,
     airport_service airport_service_1 ,
     city city_1 ,
     airport_service airport_service_2 ,
     city city_2
WHERE flight_1.from_airport = airport_service_1.airport_code
     AND airport_service_1.city_code = city_1.city_code
     AND city_1.city_name = 'PHOENIX'
     AND flight_1.to_airport = airport_service_2.airport_code
     AND airport_service_2.city_code = city_2.city_code
     AND city_2.city_name = 'MILWAUKEE'
"""

rule_based_trial(example_1, gold_sql_1)

```

Sentence: flights from phoenix to milwaukee

Parse:



Predicted SQL:

```

SELECT DISTINCT flight.flight_id FROM flight WHERE 1 AND flight.from_airport IN
  (SELECT airport_service.airport_code FROM airport_service WHERE
airport_service.city_code IN
  (SELECT city.city_code FROM city WHERE city.city_name = "PHOENIX"))
AND flight.to_airport IN
  (SELECT airport_service.airport_code FROM airport_service WHERE
airport_service.city_code IN
  (SELECT city.city_code FROM city WHERE city.city_name = "MILWAUKEE"))

```

Predicted DB result:

```

[(108086,), (108087,), (301763,), (301764,), (301765,), (301766,), (302323,),
(304881,), (310619,), (310620,)]

```

Gold DB result:

```

[(108086,), (108087,), (301763,), (301764,), (301765,), (301766,), (302323,),
(304881,), (310619,), (310620,)]

```

Correct!

To make development faster, we recommend starting with a few examples before running the full evaluation script. We've taken some examples from the ATIS dataset including the gold SQL

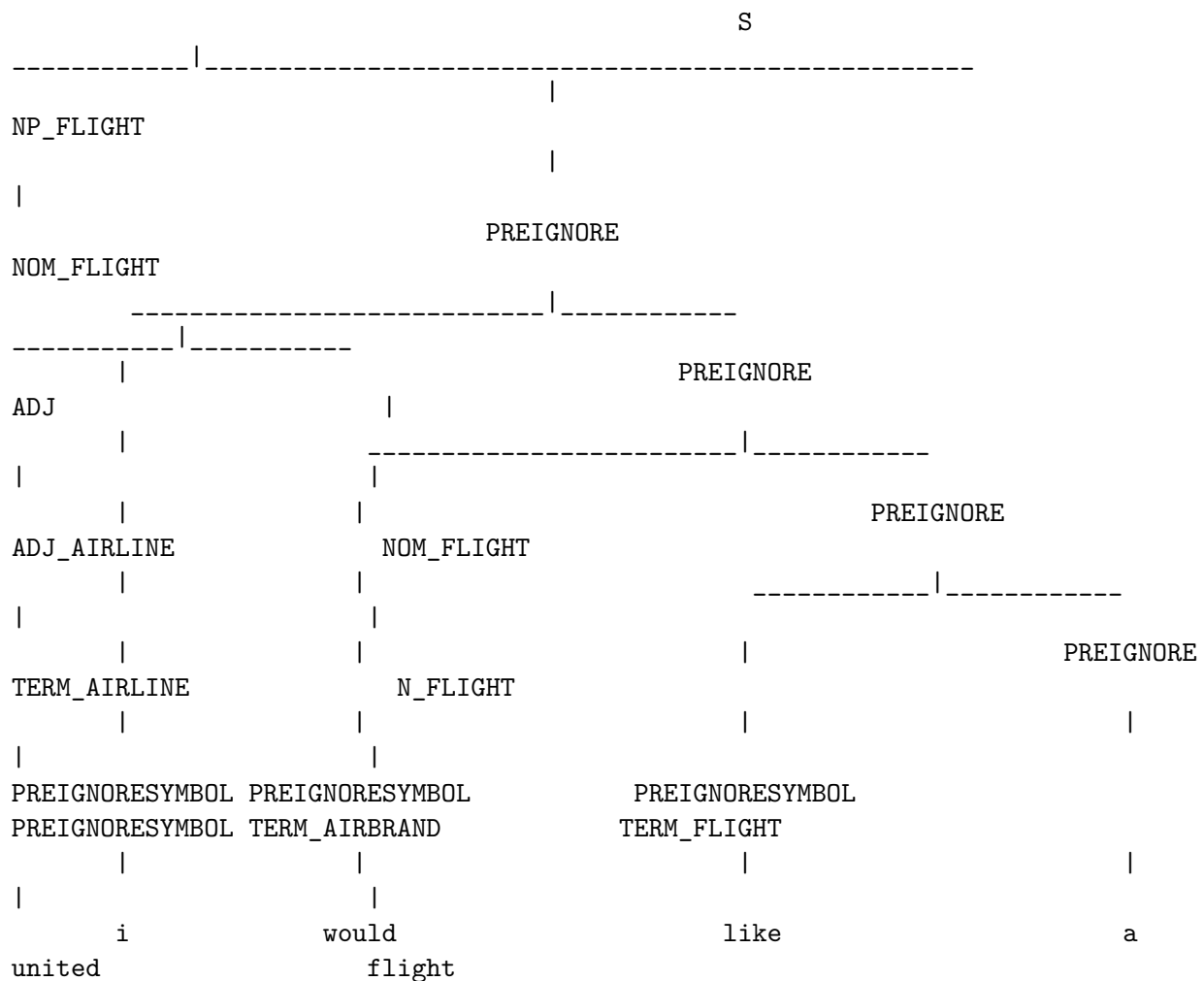
queries that they provided. Of course, yours (and those of the project segment solution set) may differ.

```
[46]: #TODO: add augmentations to `data/grammar` to make this example work
# Example 2
example_2 = 'i would like a united flight'
gold_sql_2 = """
SELECT DISTINCT flight_1.flight_id
FROM flight flight_1
WHERE flight_1.airline_code = 'UA'
"""

rule_based_trial(example_2, gold_sql_2)
```

Sentence: i would like a united flight

Parse:



Predicted SQL:

```
SELECT DISTINCT flight.flight_id FROM flight WHERE flight.airline_code = 'UA'
AND 1
```

Predicted DB result:

```
[(100094,), (100099,), (100145,), (100158,), (100164,), (100167,), (100169,),
(100203,), (100204,), (100296,)]
```

Gold DB result:

```
[(100094,), (100099,), (100145,), (100158,), (100164,), (100167,), (100169,),
(100203,), (100204,), (100296,)]
```

Correct!

```
[47]: #TODO: add augmentations to `data/grammar` to make this example work
# Example 3
example_3 = 'i would like a flight between boston and dallas'
gold_sql_3 = """
SELECT DISTINCT flight_1.flight_id
FROM flight flight_1 ,
     airport_service airport_service_1 ,
     city city_1 ,
     airport_service airport_service_2 ,
     city city_2
WHERE flight_1.from_airport = airport_service_1.airport_code
     AND airport_service_1.city_code = city_1.city_code
     AND city_1.city_name = 'BOSTON'
     AND flight_1.to_airport = airport_service_2.airport_code
     AND airport_service_2.city_code = city_2.city_code
     AND city_2.city_name = 'DALLAS'
"""

# Note that the parse tree might appear wrong: instead of
# `PP_PLACE -> 'between' N_PLACE 'and' N_PLACE`, the tree appears to be
# `PP_PLACE -> 'between' 'and' N_PLACE N_PLACE`. But it's only a visualization
# error of tree.pretty_print() and you should assume that the production is
# `PP_PLACE -> 'between' N_PLACE 'and' N_PLACE` (you can verify by printing out
# all productions).
rule_based_trial(example_3, gold_sql_3)
```

Sentence: i would like a flight between boston and dallas

Parse:

```

      S
-----|-----
NP_FLIGHT
|
|
NOM_FLIGHT
|
PREIGNORE
N_FLIGHT
-----|-----
|
|
PP
|
|
PP_PLACE
|
|
N_FLIGHT
|
N_PLACE
|
N_PLACE
|
PREIGNORE
PREIGNORESYPBOL PREIGNORESYPBOL
PREIGNORESYPBOL TERM_FLIGHT
|
|
TERM_PLACE TERM_PLACE
|
|
i would like a
flight between and boston dallas

```

Predicted SQL:

```

SELECT DISTINCT flight.flight_id FROM flight WHERE 1 AND flight.from_airport IN
  (SELECT airport_service.airport_code FROM airport_service WHERE
airport_service.city_code IN
  (SELECT city.city_code FROM city WHERE city.city_name = "BOSTON"))
AND flight.to_airport IN
  (SELECT airport_service.airport_code FROM airport_service WHERE
airport_service.city_code IN
  (SELECT city.city_code FROM city WHERE city.city_name = "DALLAS"))

```

Predicted DB result:

```

[(103171,), (103172,), (103173,), (103174,), (103175,), (103176,), (103177,),

```

(103178,), (103179,), (103180,)]

Gold DB result:

[(103171,), (103172,), (103173,), (103174,), (103175,), (103176,), (103177,),
(103178,), (103179,), (103180,)]

Correct!

```
[48]: #TODO: add augmentations to `data/grammar` to make this example work
# Example 4
example_4 = 'show me the united flights from denver to baltimore'
gold_sql_4 = """
SELECT DISTINCT flight_1.flight_id
FROM flight flight_1 ,
     airport_service airport_service_1 ,
     city city_1 ,
     airport_service airport_service_2 ,
     city city_2
WHERE flight_1.airline_code = 'UA'
      AND ( flight_1.from_airport = airport_service_1.airport_code
            AND airport_service_1.city_code = city_1.city_code
            AND city_1.city_name = 'DENVER'
            AND flight_1.to_airport = airport_service_2.airport_code
            AND airport_service_2.city_code = city_2.city_code
            AND city_2.city_name = 'BALTIMORE' )

"""

rule_based_trial(example_4, gold_sql_4)
```

Sentence: show me the united flights from denver to baltimore

Parse:

S

	----- -----	
NP_FLIGHT		
NOM_FLIGHT		
----- -----		
NOM_FLIGHT		

Correct!

Parse:




```

SELECT DISTINCT flight_1.flight_id
FROM flight flight_1 ,
     airport_service airport_service_1 ,
     city city_1 ,
     airport_service airport_service_2 ,
     city city_2 ,
     days days_1 ,
     date_day date_day_1
WHERE flight_1.from_airport = airport_service_1.airport_code
     AND airport_service_1.city_code = city_1.city_code
     AND city_1.city_name = 'TAMPA'
     AND ( flight_1.to_airport = airport_service_2.airport_code
         AND airport_service_2.city_code = city_2.city_code
         AND city_2.city_name = 'CHARLOTTE'
         AND flight_1.flight_days = days_1.days_code
         AND days_1.day_name = date_day_1.day_name
         AND date_day_1.year = 1991
         AND date_day_1.month_number = 8
         AND date_day_1.day_number = 27 )
"""

```

*# You might notice that the gold answer above used the exact date, which is
not easily implementable. A more implementable way (generated by the project
segment 4 solution code) is:*

```

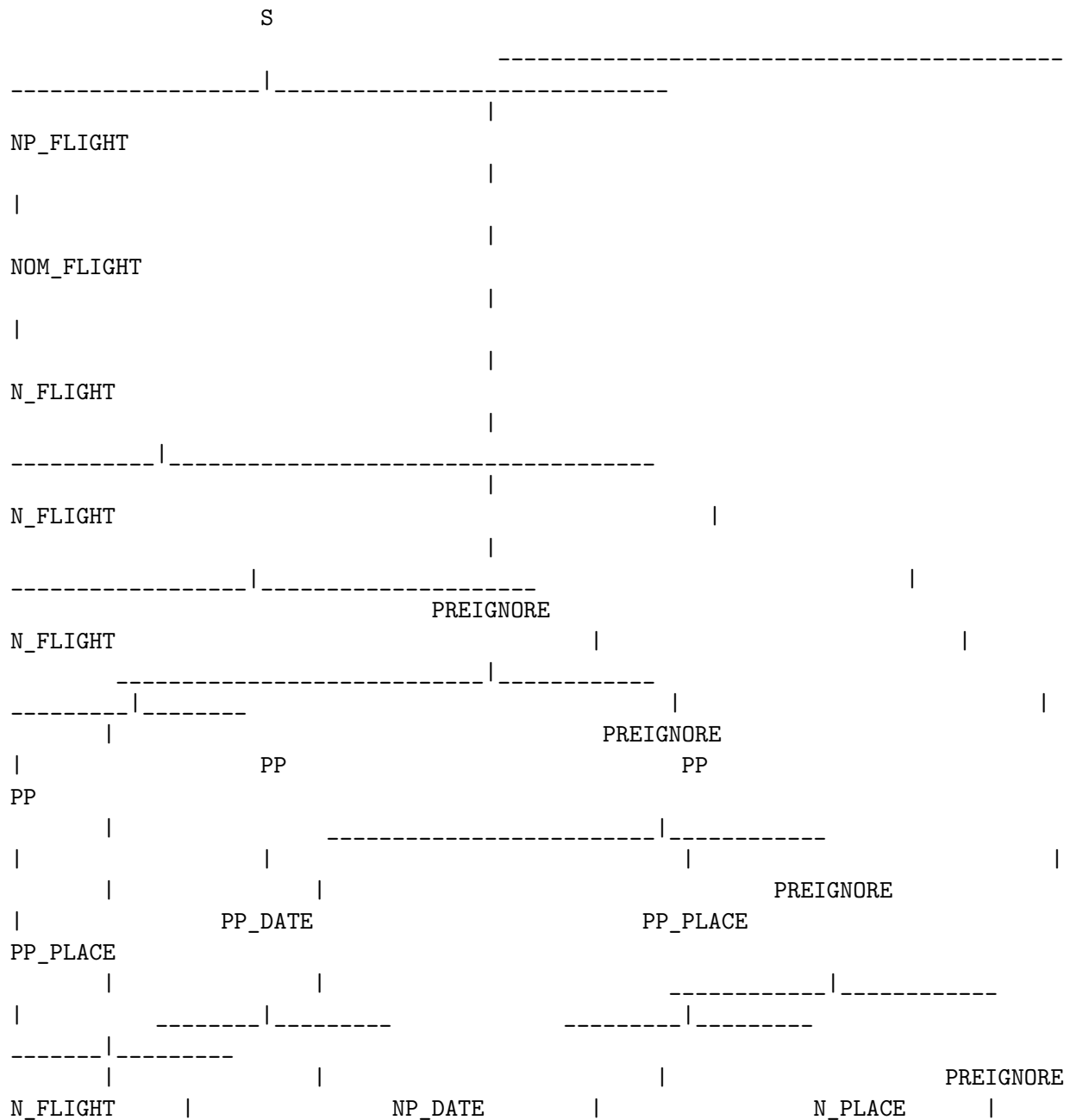
gold_sql_6b = """
SELECT DISTINCT flight.flight_id
FROM flight
WHERE (((1
        AND flight.flight_days IN (SELECT days.days_code
                                   FROM days
                                   WHERE days.day_name = 'SUNDAY')
        )
        AND flight.from_airport IN (SELECT airport_service.airport_code
                                       FROM airport_service
                                       WHERE airport_service.city_code IN
↳(SELECT city.city_code
                                       FROM
↳city
                                       )
↳WHERE city.city_name = "TAMPA"))))
        AND flight.to_airport IN (SELECT airport_service.airport_code
                                   FROM airport_service
                                   WHERE airport_service.city_code IN (SELECT
↳city.city_code
                                   FROM
↳city

```


WHERE_

Sentence: okay how about a flight on sunday from tampa to charlotte

Parse:



N_PLACE						
PREIGNORES	SYMBOL	PREIGNORES	SYMBOL	PREIGNORES	SYMBOL	
PREIGNORES	SYMBOL	TERM_FLIGHT	P_DATE	TERM_WEEKDAY	P_PLACE	
TERM_PLACE	P_PLACE	TERM_PLACE				
okay	on	how	sunday	about	tampa	a
flight	on		sunday	from	tampa	to
charlotte						

Predicted SQL:

```
SELECT DISTINCT flight.flight_id FROM flight WHERE 1 AND flight.flight_days IN
(SELECT days.days_code FROM days WHERE days.day_name = 'SUNDAY') AND
flight.from_airport IN
  (SELECT airport_service.airport_code FROM airport_service WHERE
airport_service.city_code IN
  (SELECT city.city_code FROM city WHERE city.city_name = "TAMPA"))
AND flight.to_airport IN
  (SELECT airport_service.airport_code FROM airport_service WHERE
airport_service.city_code IN
  (SELECT city.city_code FROM city WHERE city.city_name = "CHARLOTTE"))
```

Predicted DB result:

```
[(101860,), (101861,), (101862,), (101863,), (101864,), (101865,), (305231,)]
```

Gold DB result:

```
[(101860,), (101861,), (101862,), (101863,), (101864,), (101865,), (305231,)]
```

Correct!

```
[51]: #TODO: add augmentations to `data/grammar` to make this example work
# Example 7
example_7 = 'list all flights going from boston to atlanta that leaves before 7_
↳am on thursday'
gold_sql_7 = """
SELECT DISTINCT flight_1.flight_id
FROM flight flight_1 ,
      airport_service airport_service_1 ,
      city city_1 ,
```

```

        airport_service airport_service_2 ,
        city city_2 ,
        days days_1 ,
        date_day date_day_1
WHERE flight_1.from_airport = airport_service_1.airport_code
      AND airport_service_1.city_code = city_1.city_code
      AND city_1.city_name = 'BOSTON'
      AND ( flight_1.to_airport = airport_service_2.airport_code
            AND airport_service_2.city_code = city_2.city_code
            AND city_2.city_name = 'ATLANTA'
            AND ( flight_1.flight_days = days_1.days_code
                  AND days_1.day_name = date_day_1.day_name
                  AND date_day_1.year = 1991
                  AND date_day_1.month_number = 5
                  AND date_day_1.day_number = 24
                  AND flight_1.departure_time < 700 ) )

""""

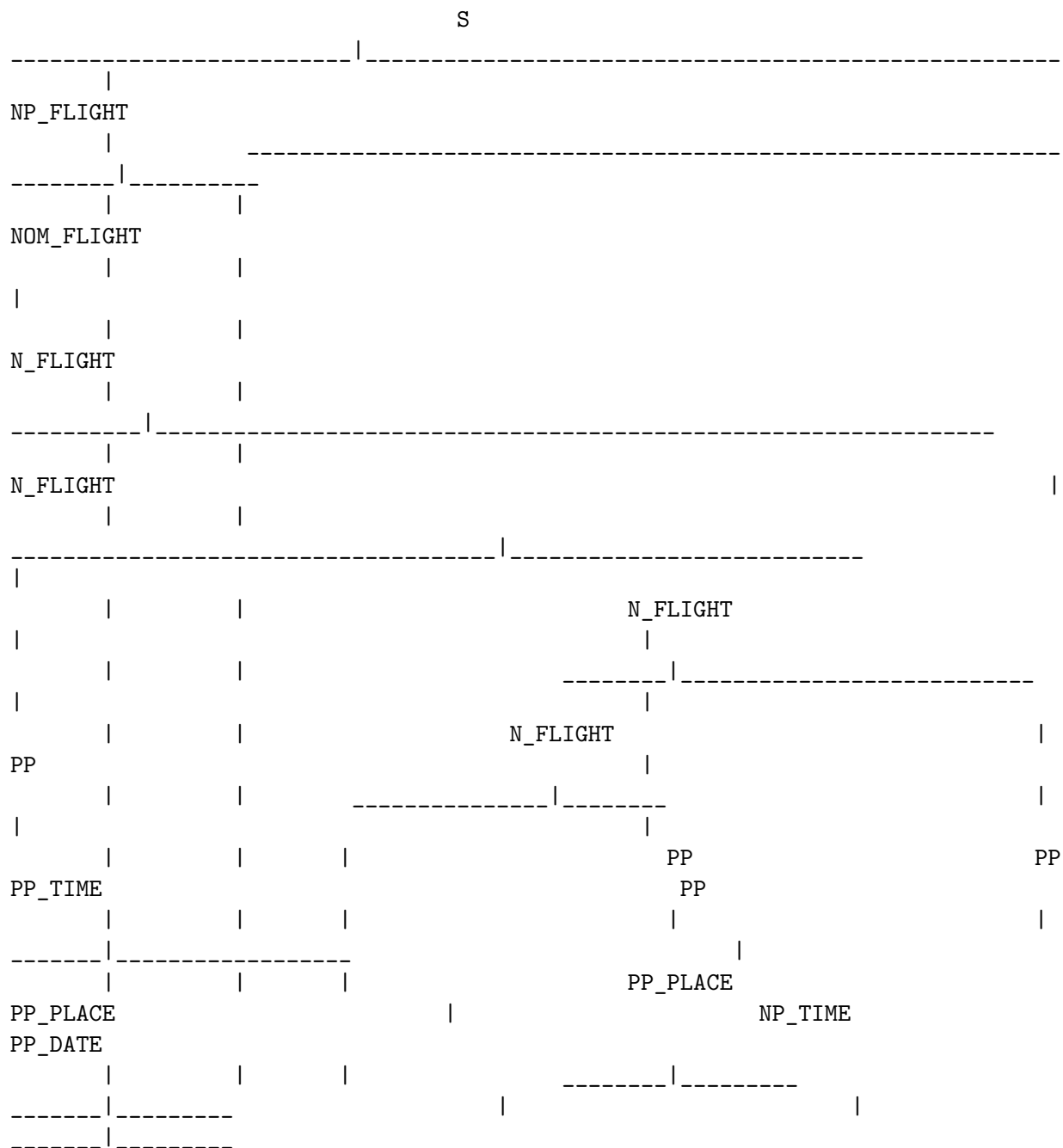
# Again, the gold answer above used the exact date, as opposed to the
# following approach:
gold_sql_7b = """
SELECT DISTINCT flight.flight_id
FROM flight
WHERE ((1
      AND (((1
            AND flight.from_airport IN (SELECT airport_service.
↪airport_code
                                     FROM airport_service
                                     WHERE airport_service.city_code_
↪IN (SELECT city.city_code
                                     FROM city
                                     WHERE city.city_name = "BOSTON"))))
            AND flight.to_airport IN (SELECT airport_service.airport_code
                                     FROM airport_service
                                     WHERE airport_service.city_code IN_
↪(SELECT city.city_code
                                     FROM city
                                     WHERE city.city_name = "ATLANTA"))))
      AND flight.departure_time <= 0700)
      AND flight.flight_days IN (SELECT days.days_code
                                FROM days
                                WHERE days.day_name = 'THURSDAY'))))

```

```
rule_based_trial(example_7, gold_sql_7b)
```

Sentence: list all flights going from boston to atlanta that leaves before 7 am on thursday

Parse:



precision, recall, and F1 metrics for the test set. It takes as argument a “predictor” function, which maps token sequences to predicted SQL queries. We’ve provided a predictor function for the rule-based model in the next cell (and a predictor for the seq2seq system below when we get to that system).

The rule-based system does not generate predictions for all queries; many queries won’t parse. The precision and recall metrics take this into account in measuring the efficacy of the method. The recall metric captures what proportion of *all of the test examples* for which the system generates a correct query. The precision metric captures what proportion of *all of the test examples for which a prediction is generated* for which the system generates a correct query. (Recall that F1 is just the geometric mean of precision and recall.)

Once you’ve made some progress on adding augmentations to the grammar, you can evaluate your progress by seeing if the precision and recall have improved. For reference, the solution code achieves precision of about 71% and recall of about 27% for an F1 of 40%.

```
[53]: def evaluate(predictor, dataset, num_examples=0, silent=True):
    """Evaluate accuracy of `predictor` by executing predictions on a
    SQL database and comparing returned results against those of gold queries.

    Arguments:
        predictor:    a function that maps a token sequence (provided by ↵
        ↵ torchtext)
                     to a predicted SQL query string
        dataset:      the dataset of token sequences and gold SQL queries
        num_examples: number of examples from `dataset` to use; all of
                     them if 0
        silent: if set to False, will print out logs
    Returns: precision, recall, and F1 score
    """
    # Prepare to count results
    if num_examples <= 0:
        num_examples = len(dataset)
    example_count = 0
    predicted_count = 0
    correct = 0
    incorrect = 0

    # Process the examples from the dataset
    for example in tqdm(dataset[:num_examples]):
        example_count += 1
        # obtain query SQL
        predicted_sql = predictor(example.src)
        if predicted_sql == None:
            continue
        predicted_count += 1
        # obtain gold SQL
        gold_sql = ' '.join(example.tgt)
```



```

    # check that they're compatible
    if verify(predicted_sql, gold_sql):
        correct += 1
    else:
        incorrect += 1

    # Compute and return precision, recall, F1
    precision = correct / predicted_count if predicted_count > 0 else 0
    recall = correct / example_count
    f1 = (2 * precision * recall) / (precision + recall) if precision + recall > 0
    else 0
    return precision, recall, f1

```

```

[54]: def rule_based_predictor(tokens):
    query = ' '.join(tokens)    # detokenized query
    tree = parse_tree(query)
    if tree is None:
        return None
    try:
        predicted_sql = interpret(tree, atis_augmentations)
    except Exception as err:
        return None
    return predicted_sql

```

```

[55]: precision, recall, f1 = evaluate(rule_based_predictor, test_iter.dataset,
    num_examples=0)
print(f"precision: {precision:3.2f}")
print(f"recall:    {recall:3.2f}")
print(f"F1:       {f1:3.2f}")

```

```
100%|          | 332/332 [00:04<00:00, 81.09it/s]
```

```

precision: 0.66
recall:    0.23
F1:        0.35

```

6 End-to-End Seq2Seq Model

In this part, you will implement a seq2seq model **with attention mechanism** to directly learn the translation from NL query to SQL. You might find labs 4-4 and 4-5 particularly helpful, as the primary difference here is that we are using a different dataset.

Note: We recommend using GPUs to train the model in this part (one way to get GPUs is to use [Google Colab](#) and clicking Menu -> Runtime -> Change runtime type -> GPU), as we need to use a very large model to solve the task well. For development we recommend starting with a smaller

model and training for only 1 epoch.

6.1 Goal 2: Implement a seq2seq model (with attention)

In lab 4-5, you implemented a neural encoder-decoder model with attention. That model was used to convert English number phrases to numbers, but one of the biggest advantages of neural models is that we can easily apply them to different tasks (such as machine translation and document summarization) by using different training datasets.

Implement the class `AttnEncoderDecoder` to convert natural language queries into SQL statements. You may find that you can reuse most of the code you wrote for lab 4-5. A reasonable way to proceed is to implement the following methods:

- **Model**

1. `__init__`: an initializer where you create network modules.
2. `forward`: given source word ids of size `(max_src_len, batch_size)`, source lengths of size `(batch_size)` and decoder input target word ids `(max_tgt_len, batch_size)`, returns logits `(max_tgt_len, batch_size, V_tgt)`. For better modularity you might want to implement it by implementing two functions `forward_encoder` and `forward_decoder`.

- **Optimization**

3. `train_all`: compute loss on training data, compute gradients, and update model parameters to minimize the loss.
4. `evaluate_ppl`: evaluate the current model's perplexity on a given dataset iterator, we use the perplexity value on the validation set to select the best model.

- **Decoding**

5. `predict`: Generates the target sequence given a list of source tokens using beam search decoding. Note that here you can assume the batch size to be 1 for simplicity.

```
[56]: def attention(batched_Q, batched_K, batched_V, mask=None):
    """
    Performs the attention operation and returns the attention matrix
    `batched_A` and the context matrix `batched_C` using queries
    `batched_Q`, keys `batched_K`, and values `batched_V`.

    Arguments:
        batched_Q: (q_len, bsz, D)
        batched_K: (k_len, bsz, D)
        batched_V: (k_len, bsz, D)
        mask: (bsz, q_len, k_len). An optional boolean mask *disallowing*
            attentions where the mask value is *False*.

    Returns:
        batched_A: the normalized attention scores (bsz, q_len, k_len)
        batched_C: a tensor of size (q_len, bsz, D).
    """
```

```

# Check sizes
D = batched_Q.size(-1)
bsz = batched_Q.size(1)
q_len = batched_Q.size(0)
k_len = batched_K.size(0)
assert batched_K.size(-1) == D and batched_V.size(-1) == D
assert batched_K.size(1) == bsz and batched_V.size(1) == bsz
assert batched_V.size(0) == k_len

if mask is not None:
    assert mask.size() == torch.Size([bsz, q_len, k_len])

q = batched_Q.transpose(0, 1) # (bsz, q_len, D)
k = batched_K.transpose(0, 1).transpose(1, 2) # (bsz, D, k_len)
scores = torch.bmm(q, k)

if mask is not None:
    scores = scores.masked_fill_(mask == False, -float('inf'))

batched_A = torch.softmax(scores, dim=-1) # (bsz, q_len, k_len)
batched_C = torch.bmm(batched_A, batched_V.transpose(0, 1)).transpose(0, 1)
→ # (q_len, bsz, D)
# Verify that things sum up to one properly.
assert torch.all(torch.isclose(batched_A.sum(-1),
                                torch.ones(bsz, q_len).to(device)))

return batched_A, batched_C

```

```

[57]: class Beam():
    """
    Helper class for storing a hypothesis, its score and its decoder hidden state.
    """
    def __init__(self, decoder_state, tokens, score):
        self.decoder_state = decoder_state
        self.tokens = tokens
        self.score = score

class BeamSearcher():
    """
    Main class for beam search.
    """
    def __init__(self, model):
        self.model = model
        self.bos_id = model.bos_id
        self.eos_id = model.eos_id
        self.padding_id_src = model.padding_id_src
        self.V = model.V_tgt

```

```

def beam_search(self, src, src_lengths, K, max_T):
    """
    Performs beam search decoding.
    Arguments:
        src: src batch of size (max_src_len, 1)
        src_lengths: src lengths of size (1)
        K: beam size
        max_T: max possible target length considered
    Returns:
        a list of token ids and a list of attentions
    """
    finished = []
    all_attns = []
    # Initialize the beam
    self.model.eval()
    #TODO - fill in `memory_bank`, `encoder_final_state`, and `init_beam` below
    memory_bank, encoder_final_state = self.model.forward_encoder(src,
↪src_lengths)
    init_beam = Beam(encoder_final_state, [torch.LongTensor(1).fill_(self.
↪bos_id).to(device)], 0)
    beams = [init_beam]

    with torch.no_grad():
        for t in range(max_T): # main body of search over time steps

            # Expand each beam by all possible tokens  $y_{t+1}$ 
            all_total_scores = []
            for beam in beams:
                y_1_to_t, score, decoder_state = beam.tokens, beam.score, beam.
↪decoder_state
                y_t = torch.LongTensor([y_1_to_t[-1]]).to(device)
                #TODO - finish the code below
                # Hint: you might want to use `model.forward_decoder_incrementally`
↪with `normalize=True`
                src_mask = src.ne(self.padding_id_src)
                logits, decoder_state, attn = self.model.
↪forward_decoder_incrementally(decoder_state, y_t, memory_bank, src_mask,
↪normalize=True)
                total_scores = logits + score
                all_total_scores.append(total_scores)
                all_attns.append(attn) # keep attentions for visualization
                beam.decoder_state = decoder_state # update decoder state in the beam
            all_total_scores = torch.stack(all_total_scores) # (K, V) when  $t > 0$ , (1,
↪V) when  $t = 0$ 

```

```

        # Find K best next beams
        # The code below has the same functionality as line 6-12, but is more
        ↪ efficient
        all_scores_flattened = all_total_scores.view(-1) # K*V when t>0, 1*V
        ↪ when t=0
        topk_scores, topk_ids = all_scores_flattened.topk(K, 0)
        beam_ids = topk_ids.div(self.V, rounding_mode='floor')
        next_tokens = topk_ids - beam_ids * self.V
        new_beams = []
        for k in range(K):
            beam_id = beam_ids[k] # which beam it comes from
            y_t_plus_1 = next_tokens[k] # which y_{t+1}
            score = topk_scores[k]
            beam = beams[beam_id]
            decoder_state = beam.decoder_state
            y_1_to_t = beam.tokens
            #TODO
            new_beam = Beam(decoder_state, y_1_to_t + [y_t_plus_1], score)
            new_beams.append(new_beam)
        beams = new_beams

        # Set aside completed beams
        # TODO - move completed beams to `finished` (and remove them from
        ↪ `beams`)
        new_beams = []
        for beam in beams:
            if beam.tokens[-1] == self.eos_id:
                finished.append(beam)
            else:
                new_beams.append(beam)
        beams = new_beams

        # Break the loop if everything is completed
        if len(beams) == 0:
            break

        # Return the best hypothesis
        if len(finished) > 0:
            finished = sorted(finished, key=lambda beam: -beam.score)
            return finished[0].tokens, all_attns
        else: # when nothing is finished, return an unfinished hypothesis
            return beams[0].tokens, all_attns

```

```

[58]: #TODO - implement the `AttnEncoderDecoder` class.
class AttnEncoderDecoder(nn.Module):
    def __init__(self, src_field, tgt_field, hidden_size=64, layers=3):
        """

```

Initializer. Creates network modules and loss function.

Arguments:

src_field: src field

tgt_field: tgt field

hidden_size: hidden layer size of both encoder and decoder

layers: number of layers of both encoder and decoder

"""

`super().__init__()`

`self.src_field = src_field`

`self.tgt_field = tgt_field`

Keep the vocabulary sizes available

`self.V_src = len(src_field.vocab.itos)`

`self.V_tgt = len(tgt_field.vocab.itos)`

Get special word ids

`self.padding_id_src = src_field.vocab.stoi[src_field.pad_token]`

`self.padding_id_tgt = tgt_field.vocab.stoi[tgt_field.pad_token]`

`self.bos_id = tgt_field.vocab.stoi[tgt_field.init_token]`

`self.eos_id = tgt_field.vocab.stoi[tgt_field.eos_token]`

Keep hyper-parameters available

`self.embedding_size = hidden_size`

`self.hidden_size = hidden_size`

`self.layers = layers`

Create essential modules

`self.word_embeddings_src = nn.Embedding(self.V_src, self.embedding_size)`

`self.word_embeddings_tgt = nn.Embedding(self.V_tgt, self.embedding_size)`

RNN cells

`self.encoder_rnn = nn.LSTM(`

`input_size = self.embedding_size,`

`hidden_size = hidden_size // 2, # to match decoder hidden size`

`num_layers = layers,`

`bidirectional = True # bidirectional encoder`

`)`

`self.decoder_rnn = nn.LSTM(`

`input_size = self.embedding_size,`

`hidden_size = hidden_size,`

`num_layers = layers,`

`bidirectional = False # unidirectional decoder`

`)`

Final projection layer

`self.hidden2output = nn.Linear(2*hidden_size, self.V_tgt) # project the`

`→ concatenation to logits`

```

    # Create loss function
    self.loss_function = nn.CrossEntropyLoss(reduction='sum',
                                              ignore_index=self.
→padding_id_tgt)

    def forward_encoder(self, src, src_lengths):
        """
        Encodes source words `src`.
        Arguments:
            src: src batch of size (max_src_len, bsz)
            src_lengths: src lengths of size (bsz)
        Returns:
            memory_bank: a tensor of size (src_len, bsz, hidden_size)
            (final_state, context): `final_state` is a tuple (h, c) where h/c
→is of size
                                (layers, bsz, hidden_size), and `context`
→is `None`.
        """
        context = None

        embeddings = self.word_embeddings_src(src)
        packed_embeddings = pack(embeddings, src_lengths.to('cpu'))
        out, (h,c) = self.encoder_rnn(packed_embeddings)
        out, _ = unpack(out)

        h_split = h.reshape(h.shape[0]//2, 2, h.shape[1], h.shape[2])
        c_split = c.reshape(c.shape[0]//2, 2, h.shape[1], c.shape[2])
        h_new = torch.cat([h_split[:, 0], h_split[:, 1]], dim=-1)
        c_new = torch.cat([c_split[:, 0], c_split[:, 1]], dim=-1)

        memory_bank = out
        final_state = (h_new, c_new)

        return memory_bank, (final_state, context)

    def forward_decoder_incrementally(self, prev_decoder_states, tgt_in_onestep,
                                     memory_bank, src_mask,
                                     normalize=True):
        """
        Forward the decoder for a single step with token `tgt_in_onestep`.
        This function will be used both in `forward_decoder` and in beam search.
        Note that bsz can be greater than 1.
        Arguments:
            prev_decoder_states: a tuple (prev_decoder_state, prev_context).
→`prev_context`
                                is `None` for the first step

```

```

        tgt_in_onestep: a tensor of size (bsz), tokens at one step
        memory_bank: a tensor of size (src_len, bsz, hidden_size), encoder_
→ outputs
                        at every position
        src_mask: a tensor of size (src_len, bsz): a boolean tensor,
→ `False` where
                        src is padding (we disallow decoder to attend to those_
→ places).
        normalize: use log_softmax to normalize or not. Beam search needs_
→ to normalize,
                        while `forward_decoder` does not

Returns:
    logits: log probabilities for `tgt_in_token` of size (bsz, V_tgt)
    decoder_states: (`decoder_state`, `context`) which will be used for_
→ the
                        next incremental update
        attn: normalized attention scores at this step (bsz, src_len)
    """
    prev_decoder_state, prev_context = prev_decoder_states
    #TODO
    embeddings = self.word_embeddings_tgt(tgt_in_onestep)
    if prev_context is not None:
        embeddings += prev_context
    embeddings = embeddings.unsqueeze(0)

    decoder_output, decoder_state = self.decoder_rnn(embeddings,
→ prev_decoder_state)

    src_mask_new = torch.transpose(src_mask, 0, 1)
    src_mask_new = src_mask_new.unsqueeze(1)

    attn, context = attention(decoder_output, memory_bank, memory_bank,
→ mask=src_mask_new)

    attn = attn.squeeze(1)
    context = context.squeeze(0)
    decoder_output = decoder_output.squeeze(0)

    logits = self.hidden2output(torch.cat([decoder_output, context], 1))
    decoder_states = (decoder_state, context)
    if normalize:
        logits = torch.log_softmax(logits, dim=-1)

    return logits, decoder_states, attn

```



```

def forward_decoder(self, encoder_final_state, tgt_in, memory_bank,
src_mask):
    """
    Decodes based on encoder final state, memory bank, src_mask, and ground
truth
target words.
Arguments:
encoder_final_state: (final_state, None) where final_state is the
encoder
final state used to initialize decoder. None
is the
initial context (there's no previous context
at the
first step).
tgt_in: a tensor of size (tgt_len, bsz)
memory_bank: a tensor of size (src_len, bsz, hidden_size), encoder
outputs
at every position
src_mask: a tensor of size (src_len, bsz): a boolean tensor,
False where
src is padding (we disallow decoder to attend to those
places).
Returns:
Logits of size (tgt_len, bsz, V_tgt) (before the softmax operation)
"""
max_tgt_length = tgt_in.size(0)

# Initialize decoder state, note that it's a tuple (state, context) here
decoder_states = encoder_final_state

all_logits = []
for i in range(max_tgt_length):
    logits, decoder_states, attn = \
        self.forward_decoder_incrementally(decoder_states,
tgt_in[i],
memory_bank,
src_mask,
normalize=False)

    all_logits.append(logits) # list of bsz, vocab_tgt
all_logits = torch.stack(all_logits, 0) # tgt_len, bsz, vocab_tgt
return all_logits

def forward(self, src, src_lengths, tgt_in):
    """
    Performs forward computation, returns logits.
Arguments:

```

```

        src: src batch of size (max_src_len, bsz)
        src_lengths: src lengths of size (bsz)
        tgt_in: a tensor of size (tgt_len, bsz)
    """
    src_mask = src.ne(self.padding_id_src) # max_src_len, bsz
    # Forward encoder
    memory_bank, encoder_final_state = self.forward_encoder(src,
→src_lengths)
    # Forward decoder
    logits = self.forward_decoder(encoder_final_state, tgt_in, memory_bank,
→src_mask)
    return logits

def evaluate_ppl(self, iterator):
    """Returns the model's perplexity on a given dataset `iterator`."""
    # Switch to eval mode
    self.eval()
    total_loss = 0
    total_words = 0
    for batch in iterator:
        # Input and target
        src, src_lengths = batch.src
        tgt = batch.tgt # max_length_sql, bsz
        tgt_in = tgt[:-1] # remove <eos> for decode input (y_0=<bos>, y_1,
→y_2)
        tgt_out = tgt[1:] # remove <bos> as target (y_1, y_2,
→y_3=<eos>)
        # Forward to get logits
        logits = self.forward(src, src_lengths, tgt_in)
        # Compute cross entropy loss
        loss = self.loss_function(logits.view(-1, self.V_tgt), tgt_out.
→view(-1))
        total_loss += loss.item()
        total_words += tgt_out.ne(self.padding_id_tgt).float().sum().item()
    return math.exp(total_loss/total_words)

def predict(self, tokens, K, max_T):
    # print("predicting...")
    beam_searcher = BeamSearcher(self)
    src = torch.tensor([[self.src_field.vocab.stoi[t]] for t in tokens],
→device=device) # max_src_len x 1
    src_length = len(src)
    # Predict
    # print("getting predictions")
    prediction, _ = beam_searcher.beam_search(src, torch.
→tensor([src_length]), K, max_T)

```

```

        # Convert to string
        # print("converting predictions")
        prediction = ' '.join([TGT.vocab.itos[token] for token in prediction])
        prediction = prediction.lstrip('<bos>').rstrip('<eos>').strip()
        return prediction

def train_all(self, train_iter, val_iter, epochs=10, learning_rate=0.001):
    """Train the model."""
    # Switch the module to training mode
    self.train()
    # Use Adam to optimize the parameters
    optim = torch.optim.Adam(self.parameters(), lr=learning_rate)
    best_validation_ppl = float('inf')
    best_model = None
    # Run the optimization for multiple epochs
    for epoch in range(epochs):
        total_words = 0
        total_loss = 0.0
        for batch in tqdm(train_iter):
            # Zero the parameter gradients
            self.zero_grad()
            # Input and target
            src, src_lengths = batch.src # text: max_src_length, bsz
            tgt = batch.tgt # max_tgt_length, bsz
            tgt_in = tgt[:-1] # Remove <eos> for decode input (y_0=<bos>,
→ y_1, y_2)
            tgt_out = tgt[1:] # Remove <bos> as target (y_1, y_2,
→ y_3=<eos>)
            bsz = tgt.size(1)
            # Run forward pass and compute loss along the way.
            logits = self.forward(src, src_lengths, tgt_in)
            loss = self.loss_function(logits.view(-1, self.V_tgt), tgt_out.
→ view(-1))
            # Training stats
            num_tgt_words = tgt_out.ne(self.padding_id_tgt).float().sum().
→ item()
            total_words += num_tgt_words
            total_loss += loss.item()
            # Perform backpropagation
            loss.div(bsz).backward()
            optim.step()

        # Evaluate and track improvements on the validation dataset
        validation_ppl = self.evaluate_ppl(val_iter)
        self.train()
        if validation_ppl < best_validation_ppl:

```

```

        best_validation_ppl = validation_ppl
        self.best_model = copy.deepcopy(self.state_dict())
    epoch_loss = total_loss / total_words
    print (f'Epoch: {epoch} Training Perplexity: {math.exp(epoch_loss):.
→4f} '
        f'Validation Perplexity: {validation_ppl:.4f}')

```

We provide the recommended hyperparameters for the final model in the script below, but you are free to tune the hyperparameters or change any part of the provided code.

For quick debugging, we recommend starting with smaller models (by using a very small `hidden_size`), and only a single epoch. If the model runs smoothly, then you can train the full model on GPUs.

```

[ ]: EPOCHS = 15 # epochs; we recommend starting with a smaller number like 1
LEARNING_RATE = 1e-4 # learning rate

# Instantiate and train classifier
model = AttnEncoderDecoder(SRC, TGT,
    hidden_size    = 1024,
    layers         = 1,
).to(device)

model.train_all(train_iter, val_iter, epochs=EPOCHS,
→learning_rate=LEARNING_RATE)
model.load_state_dict(model.best_model)

# Evaluate model performance, the expected value should be < 1.2
print (f'Validation perplexity: {model.evaluate_ppl(val_iter):.3f}')

```

```

100%|      | 229/229 [03:06<00:00,  1.23it/s]
Epoch: 0 Training Perplexity: 4.6154 Validation Perplexity: 1.8542
100%|      | 229/229 [03:05<00:00,  1.23it/s]
Epoch: 1 Training Perplexity: 1.5739 Validation Perplexity: 1.4608
100%|      | 229/229 [03:07<00:00,  1.22it/s]
Epoch: 2 Training Perplexity: 1.3403 Validation Perplexity: 1.3263
100%|      | 229/229 [03:05<00:00,  1.23it/s]
Epoch: 3 Training Perplexity: 1.2451 Validation Perplexity: 1.2480
100%|      | 229/229 [03:04<00:00,  1.24it/s]
Epoch: 4 Training Perplexity: 1.1867 Validation Perplexity: 1.2043
100%|      | 229/229 [03:06<00:00,  1.23it/s]
Epoch: 5 Training Perplexity: 1.1491 Validation Perplexity: 1.1849

```

```

100%|      | 229/229 [03:06<00:00, 1.22it/s]
Epoch: 6 Training Perplexity: 1.1230 Validation Perplexity: 1.1564
100%|      | 229/229 [03:05<00:00, 1.24it/s]
Epoch: 7 Training Perplexity: 1.1009 Validation Perplexity: 1.1418
100%|      | 229/229 [03:05<00:00, 1.24it/s]
Epoch: 8 Training Perplexity: 1.0830 Validation Perplexity: 1.1235
100%|      | 229/229 [03:05<00:00, 1.24it/s]
Epoch: 9 Training Perplexity: 1.0697 Validation Perplexity: 1.1293
100%|      | 229/229 [03:01<00:00, 1.26it/s]
Epoch: 10 Training Perplexity: 1.0606 Validation Perplexity: 1.1200
100%|      | 229/229 [03:05<00:00, 1.23it/s]
Epoch: 11 Training Perplexity: 1.0515 Validation Perplexity: 1.1223
100%|      | 229/229 [03:05<00:00, 1.24it/s]
Epoch: 12 Training Perplexity: 1.0463 Validation Perplexity: 1.1126
100%|      | 229/229 [03:07<00:00, 1.22it/s]
Epoch: 13 Training Perplexity: 1.0396 Validation Perplexity: 1.1064
100%|      | 229/229 [03:03<00:00, 1.25it/s]
Epoch: 14 Training Perplexity: 1.0341 Validation Perplexity: 1.1008
Validation perplexity: 1.101

```

```
[ ]: torch.save(model.state_dict(), "AttnEncoderDecoderWeights")
```

```

[61]: model = AttnEncoderDecoder(SRC, TGT,
    hidden_size    = 1024,
    layers         = 1,
    ).to(device)

model.load_state_dict(torch.load("AttnEncoderDecoderWeights",
    ↪map_location=torch.device('cpu'))))

```

[61]: <All keys matched successfully>

With a trained model, we can convert questions to SQL statements. We recommend making sure that the model can generate at least reasonable results on the examples from before, before evaluating on the full test set.

```

[62]: def seq2seq_trial(sentence, gold_sql):
    print("Sentence: ", sentence, "\n")
    tokens = tokenize(sentence)

```

```

predicted_sql = model.predict(tokens, K=1, max_T=400)
print("Predicted SQL:\n\n", predicted_sql, "\n")

if verify(predicted_sql, gold_sql, silent=False):
    print ('Correct!')
else:
    print ('Incorrect!')

```

```
[63]: seq2seq_trial(example_1, gold_sql_1)
```

Sentence: flights from phoenix to milwaukee

Predicted SQL:

```

SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service
airport_service_1 , city city_1 , airport_service airport_service_2 , city
city_2 WHERE flight_1.from_airport = airport_service_1.airport_code AND
airport_service_1.city_code = city_1.city_code AND city_1.city_name = 'PHOENIX'
AND flight_1.to_airport = airport_service_2.airport_code AND
airport_service_2.city_code = city_2.city_code AND city_2.city_name =
'MILWAUKEE'

```

Predicted DB result:

```

[(108086,), (108087,), (301763,), (301764,), (301765,), (301766,), (302323,),
(304881,), (310619,), (310620,)]

```

Gold DB result:

```

[(108086,), (108087,), (301763,), (301764,), (301765,), (301766,), (302323,),
(304881,), (310619,), (310620,)]

```

Correct!

```
[64]: seq2seq_trial(example_2, gold_sql_2)
```

Sentence: i would like a united flight

Predicted SQL:

```

SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service
airport_service_1 , city city_1 , days days_1 , date_day date_day_1 WHERE
flight_1.airline_code = 'UA' AND ( flight_1.from_airport =
airport_service_1.airport_code AND airport_service_1.city_code =
city_1.city_code AND city_1.city_name = 'DENVER' AND ( flight_1.to_airport =
airport_service_2.airport_code AND airport_service_2.city_code =

```

```
city_2.city_code AND city_2.city_name = 'DENVER' AND flight_1.flight_days =
days_1.days_code AND days_1.day_name = date_day_1.day_name AND date_day_1.year =
1991 AND date_day_1.month_number = 9 AND date_day_1.day_number = 1 )
```

predicted sql exec failed: incomplete input
Incorrect!

[65]: seq2seq_trial(example_3, gold_sql_3)

Sentence: i would like a flight between boston and dallas

Predicted SQL:

```
SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service
airport_service_1 , city city_1 , airport_service airport_service_2 , city
city_2 WHERE flight_1.from_airport = airport_service_1.airport_code AND
airport_service_1.city_code = city_1.city_code AND city_1.city_name = 'BOSTON'
AND flight_1.to_airport = airport_service_2.airport_code AND
airport_service_2.city_code = city_2.city_code AND city_2.city_name = 'DALLAS'
```

Predicted DB result:

```
[(103171,), (103172,), (103173,), (103174,), (103175,), (103176,), (103177,),
(103178,), (103179,), (103180,)]
```

Gold DB result:

```
[(103171,), (103172,), (103173,), (103174,), (103175,), (103176,), (103177,),
(103178,), (103179,), (103180,)]
```

Correct!

[66]: seq2seq_trial(example_4, gold_sql_4)

Sentence: show me the united flights from denver to baltimore

Predicted SQL:

```
SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service
airport_service_1 , city city_1 , airport_service airport_service_2 , city
city_2 WHERE flight_1.airline_code = 'UA' AND ( flight_1.from_airport =
airport_service_1.airport_code AND airport_service_1.city_code =
city_1.city_code AND city_1.city_name = 'DENVER' AND flight_1.to_airport =
airport_service_2.airport_code AND airport_service_2.city_code =
city_2.city_code AND city_2.city_name = 'BALTIMORE' )
```

Predicted DB result:

```
[(101231,), (101233,), (305983,)]
```

Gold DB result:

```
[(101231,), (101233,), (305983,)]
```

Correct!

[67]: seq2seq_trial(example_5, gold_sql_5)

Sentence: show flights from cleveland to miami that arrive before 4pm

Predicted SQL:

```
SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service
airport_service_1 , city city_1 , airport_service airport_service_2 , city
city_2 WHERE flight_1.from_airport = airport_service_1.airport_code AND
airport_service_1.city_code = city_1.city_code AND city_1.city_name =
'CLEVELAND' AND ( flight_1.to_airport = airport_service_2.airport_code AND
airport_service_2.city_code = city_2.city_code AND city_2.city_name = 'MIAMI'
AND flight_1.arrival_time < 1600 )
```

Predicted DB result:

```
[(107698,), (301117,)]
```

Gold DB result:

```
[(107698,), (301117,)]
```

Correct!

[68]: seq2seq_trial(example_6, gold_sql_6b)

Sentence: okay how about a flight on sunday from tampa to charlotte

Predicted SQL:

```
SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service
airport_service_1 , city city_1 , airport_service airport_service_2 , city
city_2 , days days_1 , date_day date_day_1 WHERE flight_1.from_airport =
airport_service_1.airport_code AND airport_service_1.city_code =
city_1.city_code AND city_1.city_name = 'TAMPA' AND ( flight_1.to_airport =
airport_service_2.airport_code AND airport_service_2.city_code =
city_2.city_code AND city_2.city_name = 'CHARLOTTE' AND flight_1.flight_days =
days_1.days_code AND days_1.day_name = date_day_1.day_name AND date_day_1.year =
1991 AND date_day_1.month_number = 8 AND date_day_1.day_number = 27 )
```


Predicted DB result:

```
[(101860,), (101861,), (101862,), (101863,), (101864,), (101865,), (305231,)]
```

Gold DB result:

```
[(101860,), (101861,), (101862,), (101863,), (101864,), (101865,), (305231,)]
```

Correct!

[69]: seq2seq_trial(example_7, gold_sql_7b)

Sentence: list all flights going from boston to atlanta that leaves before 7 am on thursday

Predicted SQL:

```
SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service
airport_service_1 , city city_1 , airport_service airport_service_2 , city
city_2 , days days_1 , date_day date_day_1 WHERE flight_1.from_airport =
airport_service_1.airport_code AND airport_service_1.city_code =
city_1.city_code AND city_1.city_name = 'BOSTON' AND ( flight_1.to_airport =
airport_service_2.airport_code AND airport_service_2.city_code =
city_2.city_code AND city_2.city_name = 'ATLANTA' AND ( flight_1.flight_days =
days_1.days_code AND days_1.day_name = date_day_1.day_name AND date_day_1.year =
1991 AND date_day_1.month_number = 5 AND date_day_1.day_number = 24 AND
flight_1.departure_time < 700 ) )
```

Predicted DB result:

```
[(100014,)]
```

Gold DB result:

```
[(100014,)]
```

Correct!

[70]: seq2seq_trial(example_8, gold_sql_8)

Sentence: list the flights from dallas to san francisco on american airlines

Predicted SQL:

```
SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service
airport_service_1 , city city_1 , airport_service airport_service_2 , city
city_2 WHERE flight_1.airline_code = 'AA' AND ( flight_1.from_airport =
airport_service_1.airport_code AND airport_service_1.city_code =
```

```
city_1.city_code AND city_1.city_name = 'DALLAS' AND flight_1.to_airport =
airport_service_2.airport_code AND airport_service_2.city_code =
city_2.city_code AND city_2.city_name = 'SAN FRANCISCO' )
```

Predicted DB result:

```
[(108452,), (108454,), (108456,), (111083,), (111085,), (111086,), (111090,),
(111091,), (111092,), (111094,)]
```

Gold DB result:

```
[(108452,), (108454,), (108456,), (111083,), (111085,), (111086,), (111090,),
(111091,), (111092,), (111094,)]
```

Correct!

6.1.1 Evaluation

Now we are ready to run the full evaluation. A proper implementation should reach more than 35% precision/recall/F1.

```
[71]: def seq2seq_predictor(tokens):
        prediction = model.predict(tokens, K=1, max_T=400)
        return prediction

[72]: precision, recall, f1 = evaluate(seq2seq_predictor, test_iter.dataset,
    ↪ num_examples=0)
print(f"precision: {precision:3.2f}")
print(f"recall:    {recall:3.2f}")
print(f"F1:       {f1:3.2f}")
```

```
100%|          | 332/332 [02:54<00:00, 1.91it/s]
```

```
precision: 0.37
```

```
recall:    0.37
```

```
F1:        0.37
```

6.2 Goal 3: Implement a seq2seq model (with cross attention and self attention)

In the previous section, you have implemented a seq2seq model with attention. The attention mechanism used in that section is usually referred to as “cross-attention”, as at each decoding step, the decoder attends to encoder outputs, enabling a dynamic view on the encoder side as decoding proceeds.

Similarly, we can have a dynamic view on the decoder side as well as decoding proceeds, i.e., the decoder attends to decoder outputs at previous steps. This is called “self attention”, and has been found very useful in modern neural architectures such as transformers.

Augment the seq2seq model you implemented before with a decoder self-attention mechanism as class `AttnEncoderDecoder2`. A model diagram can be found below:

At each decoding step, the decoder LSTM first produces an output state o_t , then it attends to all previous output states o_1, \dots, o_{t-1} (decoder self-attention). You need to special case the first decoding step to not perform self-attention, as there are no previous decoder states. The attention result is added to o_t itself and the sum is used as q_t to attend to the encoder side (encoder-decoder cross-attention). The rest of the model is the same as encoder-decoder with attention.

```
[73]: #TODO - implement the `AttnEncoderDecoder2` class.
class AttnEncoderDecoder2(nn.Module):
    def __init__(self, src_field, tgt_field, hidden_size=64, layers=3):
        """
        Initializer. Creates network modules and loss function.
        Arguments:
            src_field: src field
            tgt_field: tgt field
            hidden_size: hidden layer size of both encoder and decoder
            layers: number of layers of both encoder and decoder
        """
        super().__init__()
        self.src_field = src_field
        self.tgt_field = tgt_field

        # Keep the vocabulary sizes available
        self.V_src = len(src_field.vocab.itos)
        self.V_tgt = len(tgt_field.vocab.itos)

        # Get special word ids
        self.padding_id_src = src_field.vocab.stoi[src_field.pad_token]
        self.padding_id_tgt = tgt_field.vocab.stoi[tgt_field.pad_token]
        self.bos_id = tgt_field.vocab.stoi[tgt_field.init_token]
        self.eos_id = tgt_field.vocab.stoi[tgt_field.eos_token]

        # Keep hyper-parameters available
        self.embedding_size = hidden_size
        self.hidden_size = hidden_size
        self.layers = layers

        # Create essential modules
        self.word_embeddings_src = nn.Embedding(self.V_src, self.embedding_size)
        self.word_embeddings_tgt = nn.Embedding(self.V_tgt, self.embedding_size)

        # RNN cells
        self.encoder_rnn = nn.LSTM(
            input_size = self.embedding_size,
            hidden_size = hidden_size // 2, # to match decoder hidden size
            num_layers = layers,
```

```

        bidirectional = True                # bidirectional encoder
    )
    self.decoder_rnn = nn.LSTM(
        input_size    = self.embedding_size,
        hidden_size   = hidden_size,
        num_layers    = layers,
        bidirectional = False                # unidirectional decoder
    )

    # Final projection layer
    self.hidden2output = nn.Linear(2*hidden_size, self.V_tgt) # project the
    ↪ concatenation to logits

    # Create loss function
    self.loss_function = nn.CrossEntropyLoss(reduction='sum',
                                              ignore_index=self.
    ↪ padding_id_tgt)

    def forward_encoder(self, src, src_lengths):
        """
        Encodes source words `src`.
        Arguments:
            src: src batch of size (max_src_len, bsz)
            src_lengths: src lengths of size (bsz)
        Returns:
            memory_bank: a tensor of size (src_len, bsz, hidden_size)
            (final_state, context): `final_state` is a tuple (h, c) where h/c
    ↪ is of size
                                   (layers, bsz, hidden_size), and `context`
    ↪ is `None`.
        """
        context = None

        embeddings = self.word_embeddings_src(src)
        packed_embeddings = pack(embeddings, src_lengths.to('cpu'))
        out, (h,c) = self.encoder_rnn(packed_embeddings)
        out, _ = unpack(out)

        h_split = h.reshape(h.shape[0]//2, 2, h.shape[1], h.shape[2])
        c_split = c.reshape(c.shape[0]//2, 2, h.shape[1], c.shape[2])
        h_new = torch.cat([h_split[:, 0], h_split[:, 1]], dim=-1)
        c_new = torch.cat([c_split[:, 0], c_split[:, 1]], dim=-1)

        memory_bank = out
        final_state = (h_new, c_new)

        return memory_bank, (final_state, context, None)

```

```

def forward_decoder_incrementally(self, prev_decoder_states, tgt_in_onestep,
                                memory_bank, src_mask,
                                normalize=True):
    """
    Forward the decoder for a single step with token `tgt_in_onestep`.
    This function will be used both in `forward_decoder` and in beam search.
    Note that bsz can be greater than 1.
    Arguments:
        prev_decoder_states: a tuple (prev_decoder_state, prev_context).
        ↪ `prev_context`
            is `None` for the first step
        tgt_in_onestep: a tensor of size (bsz), tokens at one step
        memory_bank: a tensor of size (src_len, bsz, hidden_size), encoder
        ↪ outputs
            at every position
        src_mask: a tensor of size (src_len, bsz): a boolean tensor,
        ↪ `False` where
            src is padding (we disallow decoder to attend to those
        ↪ places).
        normalize: use log_softmax to normalize or not. Beam search needs
        ↪ to normalize,
            while `forward_decoder` does not
    Returns:
        logits: log probabilities for `tgt_in_token` of size (bsz, V_tgt)
        decoder_states: (`decoder_state`, `context`) which will be used for
        ↪ the
            next incremental update
        attn: normalized attention scores at this step (bsz, src_len)
    """

    prev_decoder_state, prev_context, prev_decoder_output =
    ↪ prev_decoder_states
    # TODO
    embeddings = self.word_embeddings_tgt(tgt_in_onestep).unsqueeze(0)
    if prev_context is not None:
        embeddings += prev_context

    decoder_output, decoder_state = self.decoder_rnn(embeddings,
    ↪ prev_decoder_state)

    # self attention
    if prev_decoder_output is not None:
        attn1, new_context1 = attention(decoder_output, prev_decoder_output,
    ↪ prev_decoder_output)

```

```

        next_decoder_output = torch.cat((prev_decoder_output,
→decoder_output), dim=0)
        new_decoder_output = decoder_output + new_context1
    else:
        next_decoder_output = decoder_output
        new_decoder_output = decoder_output

    # cross attention
    src_mask = torch.transpose(src_mask, 0, 1).unsqueeze(1)
    attn2, new_context2 = attention(new_decoder_output, memory_bank,
→memory_bank, mask=src_mask)

    decoder_states = (decoder_state, new_context2, new_decoder_output)

    concat_prev_states = torch.cat((new_decoder_output, new_context2),
→dim=2)
    logits = self.hidden2output(concat_prev_states).squeeze(0)

    if normalize:
        logits = torch.log_softmax(logits, dim=-1)

    attn = attn2.squeeze(1)

    return logits, decoder_states, attn

def forward_decoder(self, encoder_final_state, tgt_in, memory_bank,
→src_mask):
    """
    Decodes based on encoder final state, memory bank, src_mask, and ground
→truth
    target words.
    Arguments:
        encoder_final_state: (final_state, None) where final_state is the
→encoder
                                final state used to initialize decoder. None
→is the
                                initial context (there's no previous context
→at the
                                first step).
        tgt_in: a tensor of size (tgt_len, bsz)
        memory_bank: a tensor of size (src_len, bsz, hidden_size), encoder
→outputs
                                at every position
        src_mask: a tensor of size (src_len, bsz): a boolean tensor,
→`False` where

```

```

        src is padding (we disallow decoder to attend to those
→places).
    Returns:
        Logits of size (tgt_len, bsz, V_tgt) (before the softmax operation)
    """
    max_tgt_length = tgt_in.size(0)

    # Initialize decoder state, note that it's a tuple (state, context) here
    decoder_states = encoder_final_state

    all_logits = []
    for i in range(max_tgt_length):
        logits, decoder_states, attn = \
            self.forward_decoder_incrementally(decoder_states,
                                                tgt_in[i],
                                                memory_bank,
                                                src_mask,
                                                normalize=False)
        all_logits.append(logits) # list of bsz, vocab_tgt
    all_logits = torch.stack(all_logits, 0) # tgt_len, bsz, vocab_tgt
    return all_logits

def forward(self, src, src_lengths, tgt_in):
    """
    Performs forward computation, returns logits.
    Arguments:
        src: src batch of size (max_src_len, bsz)
        src_lengths: src lengths of size (bsz)
        tgt_in: a tensor of size (tgt_len, bsz)
    """
    src_mask = src.ne(self.padding_id_src) # max_src_len, bsz
    # Forward encoder
    memory_bank, encoder_final_state = self.forward_encoder(src,
→src_lengths)
    # Forward decoder
    logits = self.forward_decoder(encoder_final_state, tgt_in, memory_bank,
→src_mask)
    return logits

def evaluate_ppl(self, iterator):
    """Returns the model's perplexity on a given dataset `iterator`."""
    # Switch to eval mode
    self.eval()
    total_loss = 0
    total_words = 0
    for batch in iterator:
        # Input and target

```

```

src, src_lengths = batch.src
tgt = batch.tgt # max_length_sql, bsz
tgt_in = tgt[:-1] # remove <eos> for decode input (y_0=<bos>, y_1,
→y_2)
tgt_out = tgt[1:] # remove <bos> as target (y_1, y_2,
→y_3=<eos>)
# Forward to get logits
logits = self.forward(src, src_lengths, tgt_in)
# Compute cross entropy loss
loss = self.loss_function(logits.view(-1, self.V_tgt), tgt_out.
→view(-1))
total_loss += loss.item()
total_words += tgt_out.ne(self.padding_id_tgt).float().sum().item()
return math.exp(total_loss/total_words)

def predict(self, tokens, K, max_T):
    # print("predicting...")
    beam_searcher = BeamSearcher(self)
    src = torch.tensor([[self.src_field.vocab.stoi[t]] for t in tokens],
→device=device) # max_src_len x 1
    src_length = len(src)
    # Predict
    # print("getting predictions")
    prediction, _ = beam_searcher.beam_search(src, torch.
→tensor([src_length]), K, max_T)
    # Convert to string
    # print("converting predictions")
    prediction = ' '.join([TGT.vocab.itos[token] for token in prediction])
    prediction = prediction.lstrip('<bos>').rstrip('<eos>').strip()
    return prediction

def train_all(self, train_iter, val_iter, epochs=10, learning_rate=0.001):
    """Train the model."""
    # Switch the module to training mode
    self.train()
    # Use Adam to optimize the parameters
    optim = torch.optim.Adam(self.parameters(), lr=learning_rate)
    best_validation_ppl = float('inf')
    best_model = None
    # Run the optimization for multiple epochs
    for epoch in range(epochs):
        total_words = 0
        total_loss = 0.0
        for batch in tqdm(train_iter):
            # Zero the parameter gradients
            self.zero_grad()
            # Input and target

```



```

src, src_lengths = batch.src # text: max_src_length, bsz
tgt = batch.tgt # max_tgt_length, bsz
tgt_in = tgt[:-1] # Remove <eos> for decode input (y_0=<bos>,  

→y_1, y_2)
tgt_out = tgt[1:] # Remove <bos> as target (y_1, y_2,  

→y_3=<eos>)
bsz = tgt.size(1)
# Run forward pass and compute loss along the way.
logits = self.forward(src, src_lengths, tgt_in)
loss = self.loss_function(logits.view(-1, self.V_tgt), tgt_out.  

→view(-1))
# Training stats
num_tgt_words = tgt_out.ne(self.padding_id_tgt).float().sum().  

→item()
total_words += num_tgt_words
total_loss += loss.item()
# Perform backpropagation
loss.div(bsz).backward()
optim.step()

# Evaluate and track improvements on the validation dataset
validation_ppl = self.evaluate_ppl(val_iter)
self.train()
if validation_ppl < best_validation_ppl:
    best_validation_ppl = validation_ppl
    self.best_model = copy.deepcopy(self.state_dict())
epoch_loss = total_loss / total_words
print (f'Epoch: {epoch} Training Perplexity: {math.exp(epoch_loss):.  

→4f} '
      f'Validation Perplexity: {validation_ppl:.4f}')

```

```

[ ]: EPOCHS = 15 # epochs, we recommend starting with a smaller number like 1
LEARNING_RATE = 1e-4 # learning rate

# Instantiate and train classifier
model2 = AttnEncoderDecoder2(SRC, TGT,
    hidden_size    = 1024,
    layers         = 1,
).to(device)

model2.train_all(train_iter, val_iter, epochs=EPOCHS,  

→learning_rate=LEARNING_RATE)
model2.load_state_dict(model2.best_model)

# Evaluate model performance, the expected value should be < 1.2
print (f'Validation perplexity: {model2.evaluate_ppl(val_iter):.3f}')

```

100%| | 229/229 [04:03<00:00, 1.07s/it]
Epoch: 0 Training Perplexity: 31.0557 Validation Perplexity: 13.5280
100%| | 229/229 [03:49<00:00, 1.00s/it]
Epoch: 1 Training Perplexity: 10.4912 Validation Perplexity: 7.7536
100%| | 229/229 [03:47<00:00, 1.01it/s]
Epoch: 2 Training Perplexity: 6.4168 Validation Perplexity: 6.4778
100%| | 229/229 [03:47<00:00, 1.01it/s]
Epoch: 3 Training Perplexity: 5.4294 Validation Perplexity: 4.9064
100%| | 229/229 [03:48<00:00, 1.00it/s]
Epoch: 4 Training Perplexity: 4.2455 Validation Perplexity: 4.6522
100%| | 229/229 [03:49<00:00, 1.00s/it]
Epoch: 5 Training Perplexity: 3.4722 Validation Perplexity: 3.1353
100%| | 229/229 [03:46<00:00, 1.01it/s]
Epoch: 6 Training Perplexity: 2.9147 Validation Perplexity: 2.7828
100%| | 229/229 [03:44<00:00, 1.02it/s]
Epoch: 7 Training Perplexity: 2.5609 Validation Perplexity: 3.1782
100%| | 229/229 [03:46<00:00, 1.01it/s]
Epoch: 8 Training Perplexity: 2.5760 Validation Perplexity: 3.0606
100%| | 229/229 [03:46<00:00, 1.01it/s]
Epoch: 9 Training Perplexity: 2.4758 Validation Perplexity: 2.1693
100%| | 229/229 [03:45<00:00, 1.02it/s]
Epoch: 10 Training Perplexity: 2.1312 Validation Perplexity: 2.0781
100%| | 229/229 [03:47<00:00, 1.01it/s]
Epoch: 11 Training Perplexity: 1.9050 Validation Perplexity: 2.1166
100%| | 229/229 [03:48<00:00, 1.00it/s]
Epoch: 12 Training Perplexity: 1.8267 Validation Perplexity: 1.8021
100%| | 229/229 [03:49<00:00, 1.00s/it]
Epoch: 13 Training Perplexity: 1.7214 Validation Perplexity: 2.3508
100%| | 229/229 [03:48<00:00, 1.00it/s]
Epoch: 14 Training Perplexity: 1.7965 Validation Perplexity: 1.6273
Validation perplexity: 1.627

```
[ ]: torch.save(model2.state_dict(), "AttnEncoderDecoder2Weights")
model2.load_state_dict(torch.load("AttnEncoderDecoder2Weights"))
```

```
[ ]: model2.train_all(train_iter, val_iter, epochs=EPOCHS,
    ↪learning_rate=LEARNING_RATE)
model2.load_state_dict(model2.best_model)
torch.save(model2.state_dict(), "AttnEncoderDecoder2Weights")
```

```
100%|      | 229/229 [03:48<00:00, 1.00it/s]
Epoch: 0 Training Perplexity: 1.6624 Validation Perplexity: 1.7390
100%|      | 229/229 [03:49<00:00, 1.00s/it]
Epoch: 1 Training Perplexity: 1.5744 Validation Perplexity: 1.5839
100%|      | 229/229 [03:49<00:00, 1.00s/it]
Epoch: 2 Training Perplexity: 1.5339 Validation Perplexity: 1.4643
100%|      | 229/229 [03:46<00:00, 1.01it/s]
Epoch: 3 Training Perplexity: 1.4520 Validation Perplexity: 1.5803
100%|      | 229/229 [03:44<00:00, 1.02it/s]
Epoch: 4 Training Perplexity: 1.4305 Validation Perplexity: 1.4715
100%|      | 229/229 [03:43<00:00, 1.02it/s]
Epoch: 5 Training Perplexity: 1.3760 Validation Perplexity: 1.3989
100%|      | 229/229 [03:50<00:00, 1.01s/it]
Epoch: 6 Training Perplexity: 1.3749 Validation Perplexity: 1.4703
100%|      | 229/229 [03:50<00:00, 1.01s/it]
Epoch: 7 Training Perplexity: 1.3906 Validation Perplexity: 1.4017
100%|      | 229/229 [03:47<00:00, 1.01it/s]
Epoch: 8 Training Perplexity: 1.3268 Validation Perplexity: 4.9753
100%|      | 229/229 [03:48<00:00, 1.00it/s]
Epoch: 9 Training Perplexity: 1.3749 Validation Perplexity: 2.5057
100%|      | 229/229 [03:51<00:00, 1.01s/it]
Epoch: 10 Training Perplexity: 1.3410 Validation Perplexity: 1.2989
100%|      | 229/229 [03:52<00:00, 1.01s/it]
Epoch: 11 Training Perplexity: 1.2239 Validation Perplexity: 1.3987
100%|      | 229/229 [03:48<00:00, 1.00it/s]
Epoch: 12 Training Perplexity: 1.2237 Validation Perplexity: 1.3383
```

```
100%|      | 229/229 [03:47<00:00, 1.01it/s]
Epoch: 13 Training Perplexity: 1.2433 Validation Perplexity: 1.3537
100%|      | 229/229 [03:44<00:00, 1.02it/s]
Epoch: 14 Training Perplexity: 1.2361 Validation Perplexity: 1.3045
```

```
[ ]: <All keys matched successfully>
```

```
[ ]: model2.load_state_dict(torch.load("AttnEncoderDecoder2Weights"))
      EPOCHS=5
      model2.train_all(train_iter, val_iter, epochs=EPOCHS,
      ↪learning_rate=LEARNING_RATE)
      model2.load_state_dict(model2.best_model)
      torch.save(model2.state_dict(), "AttnEncoderDecoder2Weights")
```

```
100%|      | 229/229 [03:49<00:00, 1.00s/it]
Epoch: 0 Training Perplexity: 1.2675 Validation Perplexity: 1.4161
100%|      | 229/229 [03:49<00:00, 1.00s/it]
Epoch: 1 Training Perplexity: 1.2448 Validation Perplexity: 1.3015
100%|      | 229/229 [03:45<00:00, 1.02it/s]
Epoch: 2 Training Perplexity: 1.2227 Validation Perplexity: 1.3076
100%|      | 229/229 [03:43<00:00, 1.02it/s]
Epoch: 3 Training Perplexity: 1.2247 Validation Perplexity: 1.2745
100%|      | 229/229 [03:46<00:00, 1.01it/s]
Epoch: 4 Training Perplexity: 1.2484 Validation Perplexity: 1.3126
```

```
[ ]: <All keys matched successfully>
```

```
[ ]: model2.load_state_dict(torch.load("AttnEncoderDecoder2Weights"))
      EPOCHS=5
      model2.train_all(train_iter, val_iter, epochs=EPOCHS,
      ↪learning_rate=LEARNING_RATE)
      model2.load_state_dict(model2.best_model)
      torch.save(model2.state_dict(), "AttnEncoderDecoder2Weights")
```

```
100%|      | 229/229 [04:58<00:00, 1.30s/it]
Epoch: 0 Training Perplexity: 1.2207 Validation Perplexity: 1.2621
100%|      | 229/229 [04:53<00:00, 1.28s/it]
Epoch: 1 Training Perplexity: 1.1940 Validation Perplexity: 1.2667
100%|      | 229/229 [05:01<00:00, 1.32s/it]
```

Epoch: 2 Training Perplexity: 1.1769 Validation Perplexity: 1.2470
 100%| | 229/229 [05:00<00:00, 1.31s/it]
 Epoch: 3 Training Perplexity: 1.1942 Validation Perplexity: 1.2224
 100%| | 229/229 [04:56<00:00, 1.29s/it]
 Epoch: 4 Training Perplexity: 1.1674 Validation Perplexity: 1.2448

[]: <All keys matched successfully>

```
[ ]: model2.load_state_dict(torch.load("AttnEncoderDecoder2Weights"))
      EPOCHS=5
      model2.train_all(train_iter, val_iter, epochs=EPOCHS,
      ↪learning_rate=LEARNING_RATE)
      model2.load_state_dict(model2.best_model)
      torch.save(model2.state_dict(), "AttnEncoderDecoder2Weights")
```

100%| | 229/229 [05:43<00:00, 1.50s/it]
 Epoch: 0 Training Perplexity: 1.1915 Validation Perplexity: 1.3454
 100%| | 229/229 [05:41<00:00, 1.49s/it]
 Epoch: 1 Training Perplexity: 1.1817 Validation Perplexity: 1.2397
 100%| | 229/229 [05:42<00:00, 1.50s/it]
 Epoch: 2 Training Perplexity: 1.1408 Validation Perplexity: 1.2619
 100%| | 229/229 [05:35<00:00, 1.47s/it]
 Epoch: 3 Training Perplexity: 1.1578 Validation Perplexity: 2.3599
 100%| | 229/229 [04:59<00:00, 1.31s/it]
 Epoch: 4 Training Perplexity: 1.2019 Validation Perplexity: 1.3723

```
[ ]: model2.load_state_dict(torch.load("AttnEncoderDecoder2Weights"))
      EPOCHS=5
      model2.train_all(train_iter, val_iter, epochs=EPOCHS,
      ↪learning_rate=LEARNING_RATE)
      model2.load_state_dict(model2.best_model)
      torch.save(model2.state_dict(), "AttnEncoderDecoder2Weights")
```

100%| | 229/229 [04:58<00:00, 1.30s/it]
 Epoch: 0 Training Perplexity: 1.1384 Validation Perplexity: 1.2494
 100%| | 229/229 [04:56<00:00, 1.30s/it]
 Epoch: 1 Training Perplexity: 1.1345 Validation Perplexity: 1.2397
 100%| | 229/229 [05:02<00:00, 1.32s/it]
 Epoch: 2 Training Perplexity: 1.1413 Validation Perplexity: 1.5764

```
100%|      | 229/229 [05:05<00:00, 1.34s/it]
```

```
Epoch: 3 Training Perplexity: 1.1496 Validation Perplexity: 1.2146
```

```
100%|      | 229/229 [05:00<00:00, 1.31s/it]
```

```
Epoch: 4 Training Perplexity: 1.1052 Validation Perplexity: 1.1999
```

```
[74]: model2 = AttnEncoderDecoder2(SRC, TGT,
    hidden_size    = 1024,
    layers         = 1,
).to(device)
model2.load_state_dict(torch.load("AttnEncoderDecoder2Weights",
    ↪map_location=torch.device('cpu')))
```

```
[74]: <All keys matched successfully>
```

6.2.1 Evaluation

Now we are ready to run the full evaluation. A proper implementation should reach more than 35% precision/recall/F1.

```
[75]: def seq2seq_predictor2(tokens):
    prediction = model2.predict(tokens, K=1, max_T=400)
    return prediction
```

```
[76]: precision, recall, f1 = evaluate(seq2seq_predictor2, test_iter.dataset,
    ↪num_examples=0)
print(f"precision: {precision:3.2f}")
print(f"recall:    {recall:3.2f}")
print(f"F1:        {f1:3.2f}")
```

```
100%|      | 332/332 [02:22<00:00, 2.33it/s]
```

```
precision: 0.39
```

```
recall:    0.39
```

```
F1:        0.39
```

7 Discussion

7.1 Goal 4: Compare the pros and cons of rule-based and neural approaches.

Compare the pros and cons of the rule-based approach and the neural approaches with relevant examples from your experiments above. Concerning the accuracy, which approach would you choose to be used in a product? Explain.

Since the rule-based approach depends on the grammar and its augmentations, it is unable to handle phrases that cannot be parsed by the grammar. The neural approaches, however, are more likely to generate correct predictions for phrases that cannot be parsed by the grammar. That said, the

grammar implemented above has a large coverage of the dataset and covers a wide variety of phrases. It does not take any time to train the rule-based model, and precision is higher for the rule-based model compared to the precision of either neural model. Recall for the rule-based model is low which leads all three models to have comparable F1 scores. Between the two neural models, the Sequence-to-Sequence model took considerably less epochs (15 epochs) to reach a similar performance to the Sequence-to-Sequence model with self-attention (50 epochs). This means that training the model with self-attention takes a much longer time if there is no access to parallelizable computing. Having access to multiple processors would allow the model with self-attention to run faster, despite the larger number of epochs necessary to achieve a good performance. The model I would choose to be used in a product would depend on the machine on which the model will be trained. If running on a personal laptop with only CPU capabilities, running the neural models would not be practical and a rule-based approach would be best. If there is access to GPU's, the Sequence-to-Sequence model would work best. If the model is being trained on a computer cluster, I would choose the Sequence-to-Sequence model with self-attention.

7.2 (Optional) Goal 5: Use state-of-the-art pretrained transformers

The most recent breakthrough in natural-language processing stems from the use of pretrained transformer models. For example, you might have heard of pretrained transformers such as [GPT-3](#) and [BERT](#). (BERT is already used in [Google search](#).) These models are usually trained on vast amounts of text data using variants of language modeling objectives, and researchers have found that finetuning them on downstream tasks usually results in better performance as compared to training a model from scratch.

In the previous part, you implemented an LSTM-based sequence-to-sequence approach. To “upgrade” the model to be a state-of-the-art pretrained transformer only requires minor modifications.

The pretrained model that we will use is [BART](#), which uses a bidirectional transformer encoder and a unidirectional transformer decoder, as illustrated in the below diagram (image courtesy <https://arxiv.org/pdf/1910.13461>):

We can see that this model is strikingly similar to the LSTM-based encoder-decoder model we’ve been using. The only difference is that they use transformers instead of LSTMs. Therefore, we only need to change the modeling parts of the code, as we will see later.

First, we download and load the pretrained BART model from the [transformers](#) package by Huggingface. Note that we also need to use the “tokenizer” of BART, which is actually a combination of a tokenizer and a mapping from strings to word ids.

```
[ ]: pretrained_bart = BartForConditionalGeneration.from_pretrained('facebook/
    ↪bart-base')
    bart_tokenizer = BartTokenizer.from_pretrained('facebook/bart-base')
```

Below we demonstrate how to use BART’s tokenizer to convert a sentence to a list of word ids, and vice versa.

```
[ ]: # BART uses a predefined "tokenizer", which directly maps a sentence
    # to a list of ids
    def bart_tokenize(string):
```

```

    return bart_tokenizer(string)['input_ids'][:1024] # BART model can process at
    ↳most 1024 tokens

def bart_detokenize(token_ids):
    return bart_tokenizer.decode(token_ids, skip_special_tokens=True)

## Demonstrating the tokenizer
question = 'Are there any first-class flights from St. Louis at 11pm for less
    ↳than $3.50?'

tokenized_question = bart_tokenize(question)
print('tokenized:', tokenized_question)

detokenized_question = bart_detokenize(tokenized_question)
print('detokenized:', detokenized_question)

```

We need to reprocess the data using our new tokenizer. Note that here we set `batch_first` to `True`, since that's the expected input shape of the transformers package.

```

[ ]: SRC_BART = tt.data.Field(include_lengths=True,      # include lengths
                             batch_first=True,        # batches will be batch_size x
    ↳max_len
                             tokenize=bart_tokenize,  # use bart tokenizer
                             use_vocab=False,         # bart tokenizer already
    ↳converts to int ids
                             pad_token=bart_tokenizer.pad_token_id
                             )
TGT_BART = tt.data.Field(include_lengths=False,
                          batch_first=True,          # batches will be batch_size x
    ↳max_len
                          tokenize=bart_tokenize,    # use bart tokenizer
                          use_vocab=False,           # bart tokenizer already
    ↳converts to int ids
                          pad_token=bart_tokenizer.pad_token_id
                          )
fields_bart = [('src', SRC_BART), ('tgt', TGT_BART)]

# Make splits for data
train_data_bart, val_data_bart, test_data_bart = tt.datasets.TranslationDataset.
    ↳splits(
    ('_flightid.nl', '_flightid.sql'), fields_bart, path='./data/',
    train='train', validation='dev', test='test')

BATCH_SIZE = 1 # batch size for training/validation
TEST_BATCH_SIZE = 1 # batch size for test, we use 1 to make beam search
    ↳implementation easier

```



```

train_iter_bart, val_iter_bart = tt.data.BucketIterator.
    ↪ splits((train_data_bart, val_data_bart),

                                                    batch_size=BATCH_SIZE,
                                                    device=device,
                                                    repeat=False,
                                                    sort_key=lambda x: len(x.
    ↪ src),

                                                    sort_within_batch=True)
test_iter_bart = tt.data.BucketIterator(test_data_bart,
                                        batch_size=1,
                                        device=device,
                                        repeat=False,
                                        sort=False,
                                        train=False)

```

Let's take a look at the batch. Note that the shape of the batch is `batch_size x max_len`, instead of `max_len x batch_size` as in the previous part.

```

[ ]: batch = next(iter(train_iter_bart))
train_batch_text, train_batch_text_lengths = batch.src
print (f"Size of text batch: {train_batch_text.shape}")
print (f"First sentence in batch: {train_batch_text[0]}")
print (f"Length of the third sentence in batch: {train_batch_text_lengths[0]}")
print (f"Converted back to string: {bart_detokenize(train_batch_text[0])}")

train_batch_sql = batch.tgt
print (f"Size of sql batch: {train_batch_sql.shape}")
print (f"First sql in batch: {train_batch_sql[0]}")
print (f"Converted back to string: {bart_detokenize(train_batch_sql[0])}")

```

Now we are ready to implement the BART-based approach for the text-to-SQL conversion problem. In the below BART class, we have provided the constructor `__init__`, the `forward` function, and the `predict` function. Your job is to implement the main optimization `train_all`, and `evaluate_ppl` for evaluating validation perplexity for model selection.

Hint: you can use almost the same `train_all` and `evaluate_ppl` function you implemented before, but here a major difference is that due to setting `batch_first=True`, the batched source/target tensors are of size `batch_size x max_len`, as opposed to `max_len x batch_size` in the LSTM-based approach, and you need to make changes in `train_all` and `evaluate_ppl` accordingly.

```

[ ]: #TODO - finish implementing the `BART` class.
class BART(nn.Module):
    def __init__(self, tokenizer, pretrained_bart):
        """
        Initializer. Creates network modules and loss function.
        Arguments:
            tokenizer: BART tokenizer

```

```

        pretrained_bart: pretrained BART
    """
    super(BART, self).__init__()

    self.V_tgt = len(tokenizer)

    # Get special word ids
    self.padding_id_tgt = tokenizer.pad_token_id

    # Create essential modules
    self.bart = pretrained_bart

    # Create loss function
    self.loss_function = nn.CrossEntropyLoss(reduction="sum",
                                              ignore_index=self.padding_id_tgt)

    def forward(self, src, src_lengths, tgt_in):
        """
        Performs forward computation, returns logits.
        Arguments:
            src: src batch of size (batch_size, max_src_len)
            src_lengths: src lengths of size (batch_size)
            tgt_in: a tensor of size (tgt_len, bsz)
        """
        # BART assumes inputs to be batch-first
        # This single function is forwarding both encoder and decoder (w/ cross_
        ↪attn),
        # using `input_ids` as encoder inputs, and `decoder_input_ids`
        # as decoder inputs.
        logits = self.bart(input_ids=src,
                           decoder_input_ids=tgt_in,
                           use_cache=False
                           ).logits

        return logits

    def evaluate_ppl(self, iterator):
        """Returns the model's perplexity on a given dataset `iterator`."""
        #TODO - implement this function
        ...
        ppl = ...
        return ppl

    def train_all(self, train_iter, val_iter, epochs=10, learning_rate=0.001):
        """Train the model."""
        #TODO - implement this function
        ...

```

```

def predict(self, tokens, K=1, max_T=400):
    """
    Generates the target sequence given the source sequence using beam search,
    ↳ decoding.
    Note that for simplicity, we only use batch size 1.
    Arguments:
        tokens: a list of strings, the source sentence.
        max_T: at most proceed this many steps of decoding
    Returns:
        a string of the generated target sentence.
    """
    string = ' '.join(tokens) # first convert to a string
    # Tokenize and map to a list of word ids
    inputs = torch.LongTensor(bart_tokenize(string)).to(device).view(1, -1)
    # The `transformers` package provides built-in beam search support
    prediction = self.bart.generate(inputs,
                                    num_beams=K,
                                    max_length=max_T,
                                    early_stopping=True,
                                    no_repeat_ngram_size=0,
                                    decoder_start_token_id=0,
                                    use_cache=True)[0]

    return bart_detokenize(prediction)

```

The code below will kick off training, and evaluate the validation perplexity. You should expect to see a value very close to 1.

```

[ ]: EPOCHS = 5 # epochs, we recommend starting with a smaller number like 1
LEARNING_RATE = 1e-5 # learning rate

# Instantiate and train classifier
bart_model = BART(bart_tokenizer,
                  pretrained_bart
                  ).to(device)

bart_model.train_all(train_iter_bart, val_iter_bart, epochs=EPOCHS,
                    ↳ learning_rate=LEARNING_RATE)
bart_model.load_state_dict(bart_model.best_model)

# Evaluate model performance, the expected value should be < 1.2
print (f'Validation perplexity: {bart_model.evaluate_ppl(val_iter_bart):.3f}')

```

As before, make sure that your model is making reasonable predictions on a few examples before evaluating on the entire test set.

```

[ ]: def bart_trial(sentence, gold_sql):
    print("Sentence: ", sentence, "\n")

```

```

tokens = tokenize(sentence)

predicted_sql = bart_model.predict(tokens, K=1, max_T=300)
print("Predicted SQL:\n\n", predicted_sql, "\n")

if verify(predicted_sql, gold_sql, silent=False):
    print('Correct!')
else:
    print('Incorrect!')

```

```
[ ]: bart_trial(example_1, gold_sql_1)
```

```
[ ]: bart_trial(example_2, gold_sql_2)
```

```
[ ]: bart_trial(example_3, gold_sql_3)
```

```
[ ]: bart_trial(example_4, gold_sql_4)
```

```
[ ]: bart_trial(example_5, gold_sql_5)
```

```
[ ]: bart_trial(example_6, gold_sql_6b)
```

```
[ ]: bart_trial(example_7, gold_sql_7b)
```

```
[ ]: bart_trial(example_8, gold_sql_8)
```

7.2.1 Evaluation

The code below will evaluate on the entire test set. You should expect to see precision/recall/F1 greater than 40%.

```

[ ]: def seq2seq_predictor_bart(tokens):
    prediction = bart_model.predict(tokens, K=4, max_T=400)
    return prediction

[ ]: precision, recall, f1 = evaluate(seq2seq_predictor_bart, test_iter.dataset,
    ↪ num_examples=0)
print(f"precision: {precision:3.2f}")
print(f"recall:    {recall:3.2f}")
print(f"F1:       {f1:3.2f}")

```

8 Debrief

Question: We're interested in any thoughts you have about this project segment so that we can improve it for later years, and to inform later segments for this year. Please list any issues that arose or comments you have to improve the project segment. Useful things to comment on might include the following:

- Was the project segment clear or unclear? Which portions?
- Were the readings appropriate background for the project segment?
- Are there additions or changes you think would make the project segment better?

but you should comment on whatever aspects you found especially positive or negative.

The project segment was clear, but I wish more time was spent in class learning about transformer models. I also did not understand how the transformer model with self attention in lab 4-5 connected with the models we implement here. This project was a good length, I would not make any changes.

9 Instructions for submission of the project segment

This project segment should be submitted to Gradescope at <http://go.cs187.info/project4-submit-code> and <http://go.cs187.info/project4-submit-pdf>, which will be made available some time before the due date.

Project segment notebooks are manually graded, not autograded using otter as labs are. (Otter is used within project segment notebooks to synchronize distribution and solution code however.)

We will not run your notebook before grading it. Instead, we ask that you submit the already freshly run notebook. The best method is to “restart kernel and run all cells”, allowing time for all cells to be run to completion. You should submit your code to Gradescope at the code submission assignment at <http://go.cs187.info/project4-submit-code>. Make sure that you are also submitting your `data/grammar` file as part of your solution code as well.

We also request that you **submit a PDF of the freshly run notebook**. The simplest method is to use “Export notebook to PDF”, which will render the notebook to PDF via LaTeX. If that doesn’t work, the method that seems to be most reliable is to export the notebook as HTML (if you are using Jupyter Notebook, you can do so using `File -> Print Preview`), open the HTML in a browser, and print it to a file. Then make sure to add the file to your git commit. Please name the file the same name as this notebook, but with a `.pdf` extension. (Conveniently, the methods just described will use that name by default.) You can then perform a git commit and push and submit the commit to Gradescope at <http://go.cs187.info/project4-submit-pdf>.

End of project segment 4