

CS187 Lab 4-2: Semantic interpretation with SQL

November 11, 2024

```
[ ]: # 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/cs187-2024/lab4-2.git .tmp
    mv .tmp/tests ./
    mv .tmp/requirements.txt ./
    rm -rf .tmp
fi
pip install -q -r requirements.txt
""")
```

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

In the previous lab, you built a syntactic-semantic grammar for semantic interpretation of natural language into First-Order Logic (FOL) expressions.

In this lab, you'll use a similar approach for a standard NLP task, natural-language database query. You'll use the compositional semantics methods of the previous lab to convert natural-language queries into SQL queries, which can be actually executed against a SQL database. Familiarity with this task will be useful in the fourth project segment, where you'll be building systems for natural-language queries against the ATIS flight information database.

Preparation

```
[ ]: import os
import pprint
import sys
import wget

import nltk
import sqlite3
```

```
[ ]: # Prepare to download needed data
def download_if_needed(source, dest, filename, add_to_path=True):
    os.makedirs(dest, exist_ok=True) # ensure destination
    if add_to_path:
        sys.path.insert(1, dest) # add local to path
    if os.path.exists(f"./{dest}/{filename}"):
        print(f"Skipping {filename}")
    else:
        print(f"Downloading {filename} from {source}")
        wget.download(source + filename, out=dest)
        print("", flush=True)

remote_script_dir = "https://raw.githubusercontent.com/nlp-course/data/master/
↳scripts/trees/"
local_script_dir = "./scripts/"

# Download files
for filename in [
    "transform.py", # script for transforming grammars
]:
    download_if_needed(remote_script_dir, local_script_dir, filename)
```

```
[ ]: # Import functions for transforming augmented grammars
import transform as xform
```

Consider the phrase “flights from Boston to New York”. The representation of the property denoted by this phrase might be (as per last lab),

$$\lambda x. \text{Flight}(x) \wedge \text{Origin}(x, \text{Boston}) \wedge \text{Destination}(x, \text{NewYork})$$

If instead we had a SQL database with a `flight` relation with fields `flight_id`, `origin`, and `destination`, we might translate this phrase into the following query

```
SELECT flight_id from flight WHERE origin == "Boston" and destination == "NewYork"
```

which returns the flight IDs of all the “flights from Boston to New York”.

Then, we will be able to run the generated SQL query on a database of flights, to actually answer the query. This process can be described as:

NL query SQL query response

We will focus on the first transformation (NL question SQL query) using a syntactic-semantic grammar. The second transformation (SQL query response) will be executed automatically by the database.

1 Establishing the SQL database

First, we will initialize the SQL dataset. We will populate the dataset similarly to the flight world of the previous lab.

We initialize the same constants as in the previous lab:

```
[ ]: # Constants
Boston = "Boston"
NewYork = "New York"
TelAviv = "Tel Aviv"
DL10 = "DL10"
DL11 = "DL11"
DL13 = "DL13"
LY01 = "LY01"
LY12 = "LY12"
Morning = "Morning"
Evening = "Evening"
```

We populate the SQL dataset using a single table called **Flights**:

Flights:

flightid	origin	destination	departureTime	arrivalTime
DL10	Boston	NewYork	Morning	Evening
DL11	Boston	TelAviv	Evening	Morning
DL13	NewYork	Boston	Evening	Evening
LY01	TelAviv	NewYork	Evening	Morning
LY12	NewYork	TelAviv	Morning	Evening

To reset the database when working on the lab, you can re-run the cell below.

```
[ ]: def establish_database():
    conn = sqlite3.connect(":memory:")
    c = conn.cursor()
```

```

c.execute(
    "CREATE TABLE Flights (flightid TEXT, origin TEXT, destination TEXT,
    ↪departureTime TEXT, arrivalTime TEXT)"
)
c.executemany(
    "INSERT INTO Flights VALUES (?, ?, ?, ?, ?)",
    [
        (DL10, Boston, NewYork, Morning, Evening),
        (DL11, Boston, TelAviv, Evening, Morning),
        (DL13, NewYork, Boston, Evening, Evening),
        (LY01, TelAviv, NewYork, Evening, Morning),
    ],
)
return c

c = establish_database()

```

Let's query the table, to verify that it contains the proper rows:

```

[ ]: print('Flights:')
result = c.execute('SELECT * FROM Flights')
for row in result:
    print(row)

```

Complete the initialization of the database (using `c.execute('INSERT INTO...')`), for the last flight:

flightid	origin	destination	departureTime	arrivalTime
LY12	NewYork	TelAviv	Morning	Evening

Hint: Function `establish_database` provides an example of how to insert into the database.

```

[ ]: #TODO - Insert the final flight into the table
#
#     Note that you should either use the string "New York" (with space),
#     or the Python object `NewYork` (without space), but not the string
    ↪"NewYork".
#
...

```

```

[ ]: grader.check("add_flight")

```

We can test that the row was properly added:

```
[ ]: result = c.execute('SELECT * FROM Flights')
    for row in result:
        print(row)
```

In the previous lab, you created a syntactic-semantic grammar that used the lambda calculus to build FOL expressions (as Python objects) to represent the meanings of queries. In this lab, you'll use the same syntactic productions, but instead, map constituents to functions that build SQL queries.

Complete the following grammar. The first rule is provided as an example.

Hints:

1. Recall that the semantic composition functions are functions from right-hand side meanings (there might be zero or more) to the meaning of the whole.
2. The general structure of SQL queries for this grammar will be:

```
SELECT DISTINCT flightid from Flights WHERE ...
```

For consistency, the queries will always have a `WHERE` clause. If there are no conditions required in the body of the `WHERE` clause, you can just use `TRUE` as the `WHERE` body. Again for consistency, you might want always to have a `TRUE` as the final condition, e.g.,

```
SELECT DISTINCT flightid from Flights
WHERE origin == "New York"
      AND arrivalTime == "Evening"
      AND TRUE
```

3. Use semantic types consistently. For instance, you might use the following typings for meanings:

Syntactic Type	Semantic Type	Example
Q	<code>str</code> (a full query)	'SELECT DISTINCT flightid FROM Flights WHERE TRUE'
NP, PP, PP_PLACE, PP_TIME	<code>str</code> (a <code>WHERE</code> body) -> <code>str</code> (a <code>WHERE</code> body)	<code>lambda P: f'origin == NewYork AND {P}'</code>
LOC	<code>str</code> (a place)	'NewYork'
TIME	<code>str</code> (a time)	'Morning'

By making NP and PP meanings functions from `WHERE` bodies to `WHERE` bodies, it's easy to have compound conditions like `origin == "New York" AND arrivalTime == "Evening"`.

Alternatively, a simpler but less general approach is to take NP and PP meanings to be typed as `WHERE` bodies.

Syntactic Type	Semantic Type	Example
Q	<code>str</code> (a full query)	'SELECT DISTINCT flightid FROM Flights WHERE ?
NP, PP, PP_PLACE, PP_TIME	<code>str</code> (a <code>WHERE</code> body)	<code>f'origin == NewYork AND {P}'</code>
LOC	<code>str</code> (a place)	'NewYork'

Syntactic Type	Semantic Type	Example
TIME	str (a time)	'Morning'

```
[ ]: # TODO - Add augmentations to the grammar to generate SQL queries.
      # We gave you a few to get started.
```

```
grammar_spec = """
    Q -> NP
    NP -> 'flights'
    NP -> NP PP

    PP -> PP_PLACE
    PP -> PP_TIME

    PP_PLACE -> 'from' LOC
                | 'leaving' LOC
                | 'to' LOC
                | 'arriving' 'at' LOC

    PP_TIME -> 'arriving' TIME
               | 'departing' TIME
               | 'leaving' TIME

    LOC -> 'Boston'           : lambda: Boston
    LOC -> 'New' 'York'       : lambda: NewYork
    LOC -> 'Tel' 'Aviv'       : lambda: TelAviv

    TIME -> 'in' 'the' 'morning' : lambda: Morning
    TIME -> 'in' 'the' 'evening'  : lambda: Evening
    """
```

```
[ ]: grammar, augmentations = xform.parse_augmented_grammar(grammar_spec,
    ↪globals=globals())
```

To test the grammar, we can parse a sample query:

```
[ ]: parser = nltk.parse.BottomUpChartParser(grammar)

      for parse in parser.parse('flights from Boston leaving in the morning'.split()):
          parse.pretty_print()
```

1.1 Semantically interpreting syntactic trees

With parse tree in hand, and the dictionary of semantic augmentations for each syntactic rule, we can recursively traverse the tree and compute its meaning.

Write a function `interpret`, which takes a parse tree and a dictionary of semantic augmentations

indexed by syntactic production (as returned by `xform.parse_augmented_grammar`) and returns the meaning for the tree.

The function will be naturally recursive, since to compute the meaning of the tree you'll need to apply the appropriate semantic composition function to the meanings of the subtrees (that's the recursive bit).

Hints: 1. To iterate over a tree's child subtrees, you can simply iterate as if it were a list: `[child for child in tree]`. Note that `child` can be either a string (for terminals), or an `nltk.Tree` object storing the subtree (for nonterminals).

2. To get the syntactic rule at the root of the tree, you can use `tree.productions()[0]`

3. You'll want to know about Python's `*` operator. If you want to apply a function that takes multiple arguments and you have a list of its arguments, you can "unpack" the list using the `*` operator. For example, the following code is valid and works as expected:

```
def f(a, b, c):
    return a + b + c
my_arguments = [1,2,3]
f(*my_arguments)
```

This might be useful when calling a semantic composition function, since how many arguments it takes is only known at runtime.

4. If you want to check whether the root of a subtree `t` is a nonterminal or a terminal, you can use `isinstance(t, nltk.Tree)`. This will return `True` for trees rooted in a nonterminal, and `False` for terminals (because terminals are just strings in NLTK).

5. The solution is only a few lines of code.

```
[ ]: # TODO - write the `interpret` function
def interpret(tree, augmentations):
    """Returns a string containing an SQL query for the parse `tree`
    as interpreted by the grammar `augmentations`.
    """
    result = ...
    return result
```

```
[ ]: grader.check("interpreting")
```

2 Putting it all together

Now we can put everything together to translate an NL query to SQL and execute the query against the database.

```
[ ]: def query_to_sql(nl_query, parser, augmentations):
    """Parses a natural language query `nl_query`, interprets it as SQL, and
    executes and returns the result of the query.
    """
    sentence = nl_query.split()
```

```

parses = [p for p in parser.parse(sentence)]
for tree in parses:
    tree.pretty_print()
return [interpret(parse, augmentations) for parse in parses]

def query_to_answer(nl_query, parser, augmentations):
    """Parses and interprets a natural language query `nl_query` to a SQL query
    as per the provided `parser` and `augmentations` and executes the SQL
    query on the database, printing some useful information and returning
    the query results.
    """
    sql_queries = query_to_sql(nl_query, parser, augmentations)
    for query in sql_queries:
        print(f"SQL query: {query}")
        print("Result:")
        res = list(c.execute(query))
        for row in res:
            print(row)
    return res

```

You can now test your parser by running some queries all of the way through. The expected SQL query for flights from Boston is: `SELECT DISTINCT flightid FROM Flights WHERE origin == "Boston" AND TRUE` (or some such).

Note the quotation marks around `Boston`. These are required so that SQL interprets it as a field value rather than a field name.

```

[ ]: res1 = query_to_answer('flights from Boston', parser, augmentations)

[ ]: res2 = query_to_answer('flights from Boston to New York', parser, augmentations)

[ ]: res3 = query_to_answer('flights from New York arriving in the evening', parser,
    ↪augmentations)

[ ]: res4 = query_to_answer('flights from New York to Tel Aviv departing in the
    ↪morning arriving in the evening', parser, augmentations)

[ ]: res5 = query_to_answer('flights from Tel Aviv arriving at New York leaving in
    ↪the evening arriving in the morning', parser, augmentations)

```

Write your own sentence (that our syntactic grammar can parse) and check its results:

```

[ ]: #TODO - Write your own sentence that is parsable
    your_query = ...
    query_to_answer(your_query, parser, augmentations)

[ ]: grader.check("your_own")

```


Our semantic parser requires that the natural language input sentence be syntactically well-formed according to the grammar. But what happens if the sentence fails to parse syntactically?

For example, our parser can parse the sentence `flights from Boston`, but it cannot parse the sentence `show me flights from Boston`:

```
[ ]: try:
      query_to_answer("show me flights from Boston", parser, augmentations)
except ValueError as e:
    print(e)
```

If the sentence does not parse syntactically, our semantic parser won't be able to interpret it semantically.

One possible way to address this problem is to use a *partial* syntactic parser – a parser that provides a parse tree for the longest subsequence that it can parse. In the case of `show me flights from Boston`, such a partial parser will drop the words `show me` and provide the parse tree only for `flights from boston`.

Another possible solution is to use neural sequence-to-sequence methods, that will naturally address this problem by having an embedding for “unknown” words (usually referred to as `<UNK>`). We'll turn to those in later labs.

3 Lab debrief

Question: We're interested in any thoughts you have about this lab so that we can improve this lab for later years, and to inform later labs for this year. Please list any issues that arose or comments you have to improve the lab. Useful things to comment on might include the following, but you're not restricted to these:

- Was the lab too long or too short?
- Were the readings appropriate for the lab?
- Was it clear (at least after you completed the lab) what the points of the exercises were?
- Are there additions or changes you think would make the lab better?.

Type your answer here, replacing this text.

Question: What *specific single* change to the lab would have made your learning more efficient? This might be an addition of a concept that was not explained, or an example that would clarify a concept, or a problem that would have captured a concept in a better way, or anything else you can think of that would have made this a better lab.

Type your answer here, replacing this text.

End of Lab 4-2

To double-check your work, the cell below will rerun all of the autograder tests.

```
[ ]: grader.check_all()
```