CSCI 544: Applied Natural Language Processing

Assignment 4: Parsing

Jonathan May (adapted from David Chiang)

Out: 20. September 2017
Due: 4. October 2017

Edit of 2017-09-06 16:24:16-07:00

In this assignment you will build a simple constituency parser trained from the ATIS portion of the Penn Treebank. Before you begin, please download hw4.tgz from Vocareum. You are free to use any of the Python code in this archive when you program your own solutions to this assignment, however as long as you produce proper answers to the 5 written questions and a runnable shell script to Vocareum in Q6 you may use other languages or your own code. As usual, all work must be your own and may not be copied from or modified from others, including those you know or don't know.

Preparing the data

The file train.trees contains a sequence of trees, one per line, each in the following format:

```
(TOP (S (VP (VB Book) (NP (DT that) (NN flight)))) (PUNC .))
```

Run train.trees through preprocess.py and save the output to train.trees.pre. This script makes the trees strictly binary branching. When it binarizes, it inserts nodes with labels of the form X*, and when it removes unary nodes, it fuses labels so they look like X_Y . Run train.post through postprocess.py and verify that the output is identical to the original train.trees. This script reverses all the modifications made by preprocess.py.

Run train.trees.pre through unknown.py and save the output to train.trees.pre.unk. This script replaces all words that occurred only once with the special symbol <unk>.

Learning a grammar

Write code to learn a probabilistic CFG from trees, and store it in the following format:

```
NP -> DT NN # 0.5
NP -> DT NNS # 0.5
DT -> the # 1.0
NN -> boy # 0.5
NN -> girl # 0.5
NNS -> boys # 0.5
NNS -> girls # 0.5
```

Run your code on train.trees.pre.unk.

Q1. (15 pts) How many rules are there in your grammar? What is the most frequent rule, and how many times did it occur?

Parsing sentences

Now write a parser that takes your grammar and a sentence as input, and outputs the highest-probability parse. Don't forget to replace unknown words with <unk>. Don't forget to use log-probabilities to avoid underflow (or use bigfloat.py). Run your parser on dev.strings and save the output to dev.parses.

- Q2. (10 pts) Show the output of your parser on the first line of dev.strings, along with its log-probability (base 10).
- Q3. (10 pts) Show a plot of parsing time (y axis) versus sentence length (x axis). Use a log-log scale. Estimate the value of k for which $y \approx x^k$ (you can do a least-squares fit or just guess). Is it close to 3, and why or why not?
- Q4. (10 pts) Run your parser output through postprocess.py and save the output to dev.parses.post. Evaluate your parser output against the correct trees in dev.trees using the command:

python evalb.py dev.parses.post dev.trees

Show the output of this script, including your F1 score.

Improving your parser

Make at least three modifications to try to improve the accuracy of your parser. You can try the techniques described in class, or any other techniques you can think of or find in the literature. Remember that if you modify preprocess.py, you should also modify postprocess.py accordingly.

- Q5. (45 pts) Describe your modifications and report your new F1 score on dev.strings. What helped, or what didn't? Why?
- Q6. (10 pts) Modify the runme batch script in vocareum such that all the pre- and post-processing needed to train your grammar on an arbitrary training corpus and evaluate it on an arbitrary test corpus are done (the included batch script will run preprocess.py and postprocess.py, will generate an artificial grammar, and will run a trivial parser; replace the dummy code with your own code!). When you submit your batch script will report a run that trains on train.trees and decodes dev.strings, evaluating against dev.trees, so you can verify proper behavior. After the submission window closes, your parser will be trained on a corpus similar to train.trees and then evaluated on a hidden test.strings. Your score on this question (maximum 10 points out of 120) will depend on your F1 score relative to the rest of the class.

What to turn in

- On CrowdMark: Your answers to questions Q1–Q5.
- On Vocareum: Code, including the runme batch script, in order to evaluate Q6 on a hidden test corpus.