

# Project Warmup

## First Attempt on Slot Filling

**Due: Sunday 5pm 1/20**

Read the EC2 tutorial ([Link](#)) before working on this assignment. Don't forget to redeem the \$100 code (instructions in the tutorial).

In this assignment, you'll be asked to write a simple extractor based on regular expressions. It might not work very well but the process should hopefully get you familiar with the data and the task.

We prepared a prototype system for you written in Java. It implements a simple pipeline: reading queries, scanning sentences<sup>1</sup>, filling the slot for each query, and then evaluates the results against gold standards (*Assignment1.java*; hereafter the italicized font indicates a class name). Please read the data format document so that you will know how to process the data.

The code is located at <volume-mounted-location>/assignments/a1<sup>2</sup>. Change your current directory to there. To begin with, you need change the root path (*rootPath* at *src/tackbp/KbpConstants.java*) if you mount the volume not in "/kbp". Run "java -cp bin Assignment1"<sup>3</sup> and you will see some results for the "per:date\_of\_birth" slot. Your task is to build an extractor using any methods you can think of, and using as much of the preprocessed data as you want. (Documentation on the data formats is [HERE](#) and we'll cover it to you in class. But for this warmup part of the project, you don't need to use anything besides the token file if you don't want to or don't understand it). Specifically, you could mimic *RegexPerDateOfBirthFiller* by writing simple regular expressions and matching them on the sentences of the token file.

Specifically, create a class named "<SlotName>Filler" extending *sf.filler.Filler*. Remove the punctuations in the slot name and make it camel. e.g. for "per:date\_of\_birth", the class name is "PerDateOfBirthFiller".

Your assigned slot will be one of the following. Determine it by taking your student id mod 17.

- per:date\_of\_birth \* done in the example code
- 0. per:age
- 1. per:country\_of\_birth
- 2. per:stateorprovince\_of\_birth

---

<sup>1</sup> To shorten the running time, we prepared a special and much smaller data set where only documents that have answers are provided. In the real task, many more sentences will need to be processed.

<sup>2</sup> Individual jar files of the assignments will be available for download.

<sup>3</sup> Java 7 is required.

3. per:city\_of\_birth
4. per:date\_of\_death
5. per:country\_of\_death
6. per:stateorprovince\_of\_death
7. per:city\_of\_death
8. per:cause\_of\_death
9. per:religion
10. org:number\_of\_employees/members
11. org:founded
12. org:dissolved
13. org:country\_of\_headquarters
14. org:stateorprovince\_of\_headquarters
15. org:city\_of\_headquarters
16. org:website

Turn in the following to <https://catalyst.uw.edu/collectit/dropbox/xling/25619>:

1. The code you implemented for the slot filling.
2. A report that explains what relation you tackled and what you did in the code. For example, if you created a bunch of regular expressions, explain them in English. If you tried something more complicated, please also explain that. Report the performance. It's **absolutely** ok to have poor performance but some intuitions and analysis of the errors would be ideal (*sf.ErrorAnalysis* can help you identify the error for a given slot name.).
3. Good luck! And remember, we'll be discussing more advanced extraction methods in class, so our performance will improve.