

#	Title	Due Date	Grade Release Date
Assignment 02	Spell Correction using LM	March 07, 2022, AoE	March 14, 2022

The objectives of the assignments are to practice on topics covered in the lectures as well as improve the student's critical thinking and problem-solving skills in ad hoc topics that are closely related but not covered in the lectures. Lecture assignments also help students with research skills, including the ability to access, retrieve, and evaluate information (information literacy).

### Assignment

*Given a set of language models  $L$  that are trained on a corpus  $C$  and a spelling error corpus  $E$  for the English language, calculate the average success at  $k$  ( $s@k$ ) using each language model  $l \in L$  for  $E$ .*

- Train  $L = \{n\text{-Gram language models}\}$  on  $C = \{\text{news genre of Brown's corpus or any reasonable corpus}\}$  for  $n = \{1, 2, 3, 5, 10\}$ <sup>1</sup>.
- Use APPLING1DAT.643 in Birkbeck<sup>2</sup> spelling error corpus that includes the most common misspelled tokens, the correct spells, and the sentences that the misspelled token happened, in triples. For instance (stepped stepped when I first \*). *You can use any other well-know corpus.*
- Success at  $k$  ( $s@k$ ) measures whether the correct spell of the token happens to be in the top- $k$  (most probable) list of tokens that are retrieved by a language model. For instance, given 'when I first', the top-5 most probable tokens based on unigram language model would be ['went', 'saw', 'started', 'stepped', 'looked']. Then,  $s@1$  is 0 since the correct spell from Birkbeck is 'stepped' which is not happening at the first item. However,  $s@k$  for  $k \geq 4$  is 1. Report the average  $s@k$  for  $k = \{1, 5, 10\}$  using PyTrec\_Eval<sup>3</sup> for each  $l \in L$ .
- Hint: unplugging the MED (Assign 1) and plug the trained LM. There should be no change in evaluation.*

### Submission Guidelines

- Submission must be written as a report in English, in the current ACM two-column conference format in LaTeX. Overleaf templates<sup>4</sup> are available from the ACM Website<sup>5</sup> (use the sigconf proceedings template).
- The report must be 1 page in length, no more no less, including figures, tables, references, and single-authored by the student.
- The code should be available in an online repo (preferably Github) and the link should be mentioned as a footnote to the report's title. See the example below. The results reported in the report must be reproducible (multiple runs with the same setting should result in the same results.)
- Submission must be in one single zip file named COMP8730\_Assign02\_UWindId.zip, including:
  - the LaTeX files
  - the pdf file

A sample submission has been attached to this manual in Blackboard, also available online<sup>6</sup>.

<sup>1</sup> <https://www.kaggle.com/alvations/n-gram-language-model-with-nltk>

<sup>2</sup> <https://ota.bodleian.ox.ac.uk/repository/xmlui/handle/20.500.12024/0643>

<sup>3</sup> [https://github.com/cvangysel/pytrec\\_eval](https://github.com/cvangysel/pytrec_eval)

<sup>4</sup> <https://www.overleaf.com/gallery/tagged/acm-official>

<sup>5</sup> <https://www.acm.org/publications/proceedings-template>

<sup>6</sup> <https://www.overleaf.com/read/dbrhbpqghfjc>