

Controllable Generation Using kNN Language Models

Master/Bachelor Thesis

Motivation

Language modeling forms the foundation of many language processing problems and involves predicting which words come next in a sequence. Recent work on kNN language models has shown improved perplexity by storing encodings of sentence contexts and retrieving similar contexts to alter the probability distribution when predicting the next token (Khandelwal et al. 2020). If we store not only the encoding, but also labeled stylistic attributes of a sentence (e.g. politeness, formality, toxicity), can we use this to control generated language to contain these attributes to higher or lower degrees while preserving fluency? How can these stored attributes and encodings be leveraged most effectively?

Difficulty

Analysis

Programming

Task Description

on formality, We datasets politeness, offensive lancan use toxicity prompts, or others we find interesting to guage, exto show improved generation control while plore. We hope maintaining with automatic and fluency human evaluations. One potential foreseeable difficulty or Representations $\begin{array}{l} \textbf{Aggregation} \\ p_{\text{kNN}}(y) = \sum 1_{y=v_i} p(k_i) \end{array}$ Training Contexts Distances Nearest k Targets Normalization $d_i = d(q, k_i)$ $k_i = f(c_i)$ $p(k_i) \propto \exp(-d_i)$ v_i subproblem is mem-Hawaii 0.8 Illinois 0.2 Hawaii 3 -----Hawaii 0.7 Obama was senator for Illinois Illinois 0.2 Barack is married to Michelle 100 Illinois 4 ory usage for stored Obama was born in Hawaii Hawaii 5 encodings. This Obama is a native of Hawaii Interpolation Classification $y = \lambda p_{kNN}(y) + (1-\lambda)p_{LM}(y)$ $p_{LM}(y)$ can be mitigated Representation Test Context Target Hawaii 0.6 Hawaii 0.2 q = f(x)Illinois 0.2 Illinois 0.2 by only storing en-Obama's birthplace is codings of sentences



with relevant attributes and could be an interesting direction for further exploration. Work can be submitted for publication upon completion.

References

Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generaliza-1 tion through Memorization: Nearest Neighbor Language Models. In International Conference on Learning Representations (ICLR), 2020.

[2] Sudha Rao and Joel Tetreault. Dear sir or madam, may I introduce the GYAFC dataset: Corpus, benchmarks and metrics for formality style transfer. In NAACL-HLT, 2018.



[3] Rishav Hada, Sohi Sudhir, Pushkar Mishra, Helen Yannakoudakis, Saif M. Mohammad, and Ekaterina Shutova. Ruddit: Norms of offensiveness for English Reddit comments. In ACL, 2021.

[4] Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In EMNLP, 2020.

Cristian Danescu-Niculescu-Mizil, Moritz Sudhof, Dan Jurafsky, Jure Leskovec, and Christo-|5| pher Potts. A computational approach to politeness with application to social factors. In ACL, August 2013.