Correlated Topic Model with Transformer Embeddings

トランスフォーマーの埋め込みによる相関トピックモデル

by

Chun Wa Leung

梁俊華

A Master Thesis

修士論文

Submitted to

the Graduate School of the University of Tokyo

on December 9, 2021

in Partial Fulfillment of the Requirements

for the Degree of Master of Information Science and

Technology

in Computer Science

Thesis Supervisor: Akihiko Takano　高野明彦

Professor of Computer Science

## ABSTRACT

Topic modeling is one of the most common information retrieval task in natural language processing. In particular, Correlated Topic Model(CTM) is a topic model which captures the correlation between topics associated. However, such a classic statistical approach is not able to capture positional information from sequential input. At that point, traditional topic models may perform poorly in generating words from large number of topics. In this research, we introduce Correlated Topic Model with Transformer embeddings, a generative model where combine the advantage of using positional information of words and topic correlation. Specifically, transformer embeddings map topic words into latent space and further assign to its assigned topic. Our approach manages to add a covariance prior to the topic model, LKJ correlation prior to logistic-normal distribution, which aims to fit the correlation information from the data. In addition, we extend the model to handle time-series data integrated with Gaussian Process Latent Variable Model(GPLVM), which also captures temporal information from words occurrence of documents over time. The model which was optimized using Stochastic Variational Inference (SVI), allows handling massive data sets with mini-batching. As compared to empirical results from experiments, our approach performs a better fit of the data than existing generative topic model and exhibits a better capability in obtaining high quality topics.