

VAE Topic Detection

MATTEO POSENATO, TOMASO ERSEGHE*, Università degli studi di Padova, Italy

tbd

Additional Key Words and Phrases: tbd

ACM Reference Format:

Matteo Posenato, Tomaso Erseghe. 2023. VAE Topic Detection. *ACM Trans. Knowl. Discov. Data.* 1, 1 (May 2023), 3 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

[...]

2 NEURAL TOPIC MODEL BASED ON VAE

2.1 Generative model - Generation network

We review the topic detection model of [1]. We consider a corpus of D documents using a vocabulary of W words. Each document is represented by a (variable length) vector \mathbf{d}_i , $i = 1, \dots, D$ collecting the words occurrences in the document, so that $d_{i,n} \in \{1, \dots, W\}$. We let the corpus be organised in T topics, and denote with \mathbf{t}_i (vector of length T) the latent topic representation of document \mathbf{d}_i .

Our reference generative model starts from an hidden prior variable $\mathbf{z} \in \mathbb{R}^T$ normally distributed, i.e., with probability distribution function (PDF) $p(\mathbf{z}) = p_{\mathcal{N}}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ where

$$p_{\mathcal{N}}(\mathbf{x}; \mathbf{m}, \Sigma) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} e^{-\frac{1}{2}(\mathbf{x}-\mathbf{m})^T \Sigma^{-1}(\mathbf{x}-\mathbf{m})} \quad (1)$$

is the multivariate normal PDF. The latent topic representation $\mathbf{t} \in \mathbb{R}^T$ is approximated by a multilayer perceptron (MLP), to build a differentiable map of the form (if I understood correctly, but the text also mentions two MLPs!?)

$$\mathbf{t} = \mathcal{G}(\mathbf{z}) = \mathbf{W}_2 \tanh(\mathbf{W}_1 \mathbf{z} + \mathbf{b}_1) + \mathbf{b}_2 \quad (2)$$

The word-occurrence-pattern vector is then generated via softmax construction from the latent topic representation, that is (if I understood correctly, but the text also mentions a Gaussian softmax?!? also if this is correct then \mathbf{W}_2 and \mathbf{W}_3 are redundant so something is wrong)

$$\log(p_{d|t}(\mathbf{d}|\mathbf{t})) = \sum_n \log(s_{d_n}), \quad \mathbf{s} = \text{softmax}(\mathbf{W}_3 \mathbf{t} + \mathbf{b}_3)$$

* All authors contributed equally to this research.

Author's address: Matteo Posenato, Tomaso Erseghe, tomaso.erseghe@unipd.it, Università degli studi di Padova, Via VII Febbraio, 2, Padova, Italy, 35131.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Association for Computing Machinery.

Manuscript submitted to ACM

where s_{d_n} denotes the d_n th entry of \mathbf{s} . Hence, we have

$$p_\theta(\mathbf{d}|\mathbf{z}) = p_{d|t}(\mathbf{d}|\mathfrak{D}(\mathbf{z})) \quad (3)$$

with parameters $\theta = \{\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3\}$.

2.2 Approximation of posterior probability - Inference network

The posterior probability $q_\phi(\mathbf{z}|\mathbf{d})$ is approximated by the multivariate normal distribution

$$q_\phi(\mathbf{z}|\mathbf{d}) = p_{\mathcal{N}}(\mathbf{z}; \boldsymbol{\mu}_\phi(\mathbf{d}), \text{diag}(\boldsymbol{\sigma}_\phi^2(\mathbf{d}))) \quad (4)$$

where $\boldsymbol{\mu}_\phi$ and $\boldsymbol{\sigma}_\phi^2$ are differential maps generated through two MLPs. Specifically, we have (this is my guess from [2])

$$\begin{aligned} \boldsymbol{\mu}_\phi(\mathbf{d}) &= \mathbf{W}_5 \mathbf{h} + \mathbf{b}_5, & \mathbf{h} &= \tanh(\mathbf{W}_4 \mathbf{d} + \mathbf{b}_4) \\ \log(\boldsymbol{\sigma}_\phi^2(\mathbf{d})) &= \mathbf{W}_6 \mathbf{h} + \mathbf{b}_6 \end{aligned}$$

2.3 Target function

According to the variational auto encoder (VAE) approach of [2] we define a variational lower bound $f_{\theta,\phi}(\mathbf{d}) \leq \log p_\theta(\mathbf{d})$ as

$$\begin{aligned} f_{\theta,\phi}(\mathbf{d}) &= \log p_\theta(\mathbf{d}) - D_{\text{KL}}\left(q_\phi(\mathbf{z}|\mathbf{d}) \parallel p_\theta(\mathbf{z}|\mathbf{d})\right) \\ &= \int d\mathbf{z} q_\phi(\mathbf{z}|\mathbf{d}) \log \left(\frac{p_\theta(\mathbf{z}, \mathbf{d})}{q_\phi(\mathbf{z}|\mathbf{d})} \right) \\ &= \underbrace{\int d\mathbf{z} q_\phi(\mathbf{z}|\mathbf{d}) \log(p_\theta(\mathbf{d}|\mathbf{z}))}_{f_1} - \underbrace{\int d\mathbf{z} q_\phi(\mathbf{z}|\mathbf{d}) \log \left(\frac{p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{d})} \right)}_{f_2} \end{aligned} \quad (5)$$

with target function $f_{\theta,\phi}(\mathbf{d})$ to be maximized with respect to the parameters θ and ϕ . By exploiting (4), the target function can be rewritten in the form

$$\begin{aligned} f_1(\mathbf{d}) &= \int d\mathbf{u} p_{\mathcal{N}}(\mathbf{u}; \mathbf{0}, \mathbf{I}) \log \left(p_\theta(\mathbf{d} | \boldsymbol{\mu}_\phi(\mathbf{d}) + \boldsymbol{\sigma}_\phi(\mathbf{d}) \circ \mathbf{u}) \right) \\ &\simeq \frac{1}{L} \sum_{\ell=1}^L \log \left(p_\theta(\mathbf{d} | \boldsymbol{\mu}_\phi(\mathbf{d}) + \boldsymbol{\sigma}_\phi(\mathbf{d}) \circ \mathbf{u}_\ell) \right) \\ f_2(\mathbf{d}) &= -\frac{1}{2} \mathbf{1}^T \left(\mathbf{1} + \boldsymbol{\mu}_\phi^2(\mathbf{d}) + \boldsymbol{\sigma}_\phi^2(\mathbf{d}) + \log(\boldsymbol{\sigma}_\phi^2(\mathbf{d})) \right) \end{aligned} \quad (6)$$

where \circ stands for element-wise product, and where $\mathbf{u}_\ell \in \mathcal{N}(\mathbf{0}, \mathbf{I})$ are independent normal samples.

2.4 Topic detection

One important question. If I understood correctly \mathbf{W}_3 is the matrix that must be paired with the sentiment one, and this matrix identifies the topics. However, how can we assign documents to topics? After VAE we know $q_\phi(\mathbf{z}|\mathbf{d})$ so we can guess the statistics on \mathbf{t} via

$$p(\mathbf{t}|\mathbf{d}) = q_\phi(\mathfrak{D}^{-1}(\mathbf{t})|\mathbf{d})$$

but does it make any sense?

REFERENCES

- [1] Lin Gui, Jia Leng, Jiyun Zhou, Ruifeng Xu, and Yulan He. 2020. Multi task mutual learning for joint sentiment classification and topic detection. *IEEE Transactions on Knowledge and Data Engineering* 34, 4 (2020), 1915–1927.
- [2] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).

Received tbd