

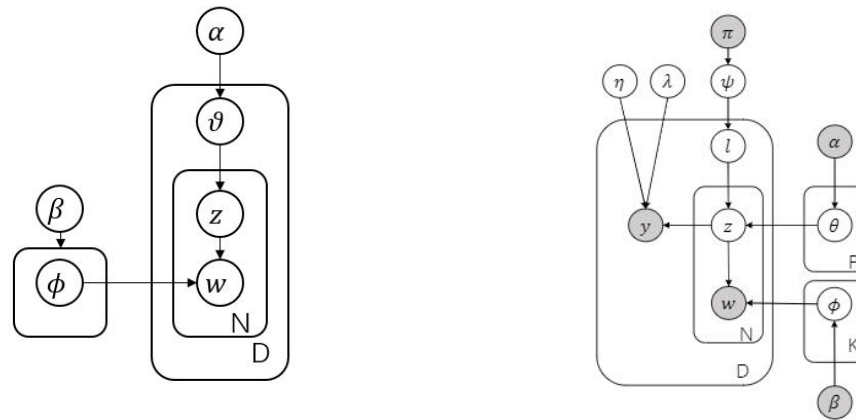
PSLDA: A novel supervised pseudo document-based topic model for short texts

**Mingtao SUN, Xiaowei ZHAO, Jingjing LIN, Jian JING,
Deqing WANG, Guozhu JIA**

Frontiers of Computer Science, DOI: [10.1007/s11704-021-0606-3](https://doi.org/10.1007/s11704-021-0606-3)

Problems & Ideas

- Problems of conventional short text topic model:
 - Sparseness of word-occurrence and the diversity of topics.
 - Previous methods suffer from high time complexity because of the pre-trained processing on the external large-scale dataset.
- Ideas: Assume that short texts are generated from the normal size latent pseudo documents, and the topic distributions are sampled from the pseudo documents.



The graphical representation of traditional topic model and PSLDA.

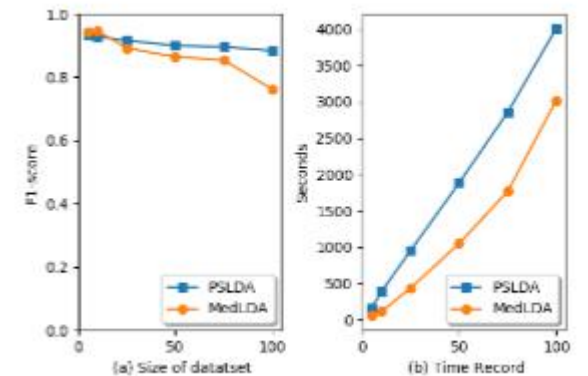
Main Contributions

- Contributions:

- A Pseudo-document-based Topic Model (PTM) for observed short texts
- An expected classifier to the model, which is an alternative formulation of max-margin supervised topic model

Method	DBLP			NEWS			TWEET		
	Precision	Recall	F-measure	Precision	Recall	F-measure	Precision	Recall	F-measure
PSLDA	0.6356 \pm 2.2e-6	0.6376 \pm 2.2e-6	0.6344 \pm 6.4e-7	0.8884 ^{**} \pm 2.1e-5	0.8824 ^{**} \pm 1.8e-5	0.8836 ^{**} \pm 2.3e-5	0.7252 ^{**} \pm 7.6e-7	0.7092 ^{**} \pm 7.4e-6	0.7136 ^{**} \pm 2.2e-4
MedLDA	0.6664 \pm 1.1e-2	0.658 \pm 9.9e-3	0.6572 \pm 1.1e-2	0.7344 \pm 5.3e-4	0.8008 \pm 1.6e-4	0.7616 \pm 3.0e-4	0.544 \pm 6.9e-3	0.5548 \pm 5.7e-3	0.5456 \pm 6.4e-3
SLDA	0.652 \pm 6.9e-5	0.6567 \pm 4.9e-5	0.6567 \pm 7.3e-5	0.853 \pm 8.3e-5	0.854 \pm 8.6e-5	0.854 \pm 8.6e-5	0.63 \pm 1.2e-4	0.643 \pm 8.9e-5	0.6313 \pm 1.3e-4
LDA+SVM	0.604 \pm 1.2e-4	0.59 \pm 2.2e-4	0.594 \pm 3.5e-4	0.7384 \pm 7.9e-5	0.7348 \pm 3.9e-4	0.7356 \pm 5.7e-5	0.5536 \pm 6.8e-3	0.5168 \pm 7.5e-4	0.5292 \pm 7.3e-4
SPTM	0.661 \pm 2.4e-4	0.667 \pm 7.6e-4	0.663 \pm 7.2e-3	0.760 \pm 1.8e-3	0.761 \pm 7.9e-5	0.759 \pm 3.1e-3	0.551 \pm 6.2e-4	0.558 \pm 4.1e-3	0.550 \pm 1.1e-2
SATM	0.657 \pm 6.3e-5	0.662 \pm 5.1e-6	0.654 \pm 4.3e-3	0.697 \pm 7.6e-4	0.702 \pm 8.1e-3	0.686 \pm 8.4e-3	0.599 \pm 6.0e-4	0.605 \pm 7.1e-5	0.594 \pm 2.5e-4
PTM	0.667 ^{**} \pm 4.1e-3	0.672 ^{**} \pm 2.9e-2	0.668 ^{**} \pm 6.1e-4	0.755 \pm 5.7e-3	0.757 \pm 5.4e-4	0.754 \pm 1.5e-2	0.561 \pm 8.1e-4	0.568 \pm 7.6e-3	0.559 \pm 4.7e-3

The Precision, Recall and F-measure on three datasets.



Effect of dataset size on (a) F1-score and (b) training time