



Neural Sparse Topical Coding

Min Peng¹, Qianqian Xie¹, Yanchun Zhang², Hua Wang², Xiuzheng Zhang³

¹School of Computer Science, Wuhan University

²Centre for Applied Informatics, Victoria University

³School of Science, RMIT University

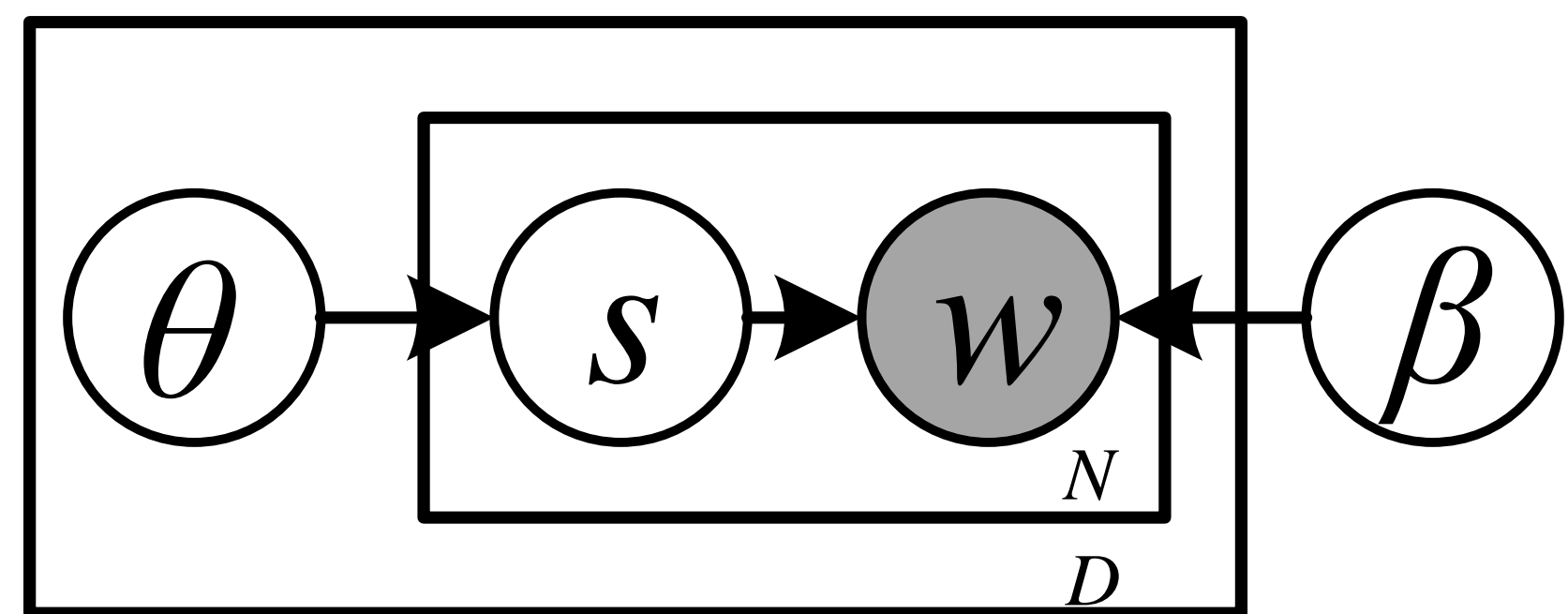


Abstract

- The poor scalability has become a bottleneck for the application and extension of the traditional topic models
- We propose Neural Sparse Topical Coding (*NSTC*) via adopting neural network to model the generation process of traditional sparsity-enhanced topic model *STC*, thus can improve its flexibility
- To illustrate the flexibility, we present three extensions base on *NSTC* without re-deduced inference algorithms

Sparse Topical Coding

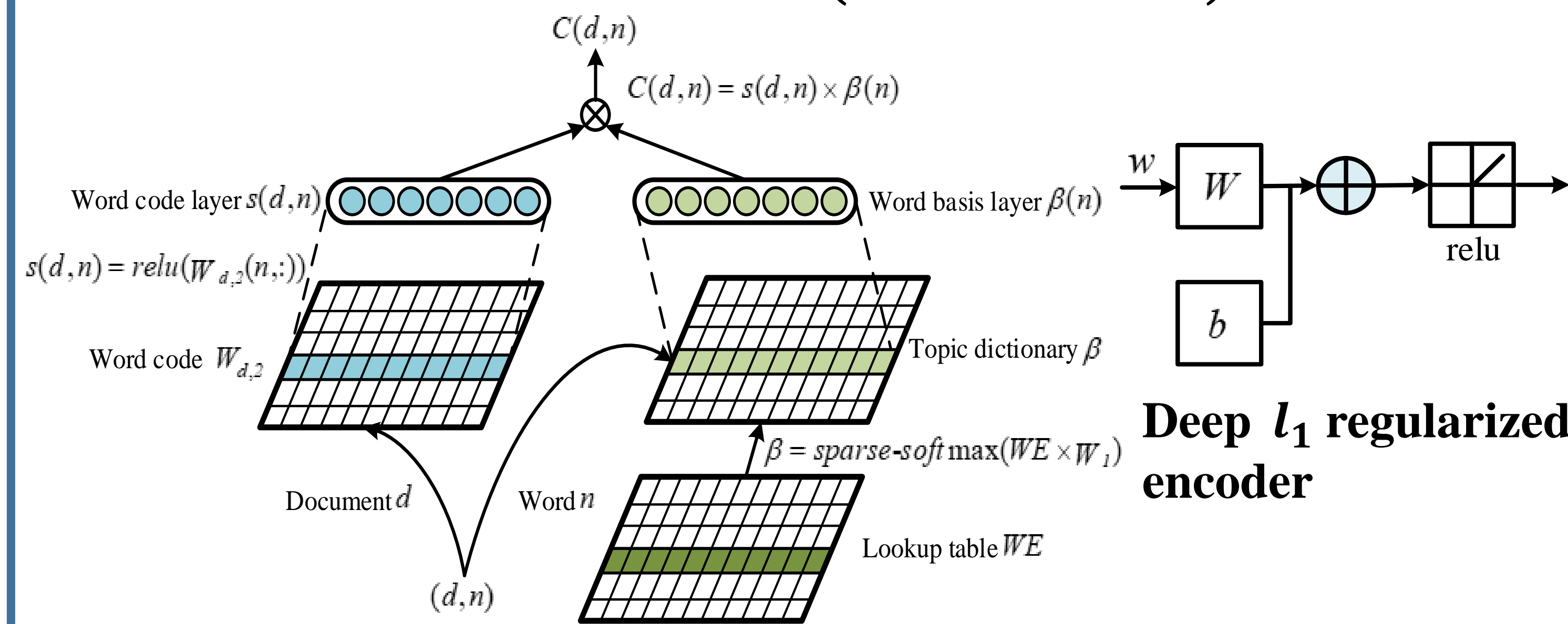
- Document code θ : the document representation of a document in topic space \mathbf{k}
- Word code \mathbf{s} : the word representation of a word in topic space \mathbf{k}
- Topic dictionary β : a global topic dictionary with \mathbf{K} bases



$$\min_{\theta, s, \beta} \sum_{d,n} l(w_{d,n}; s_{d,n}, \beta) + \lambda_1 \sum_d \|\theta_d\|_1 + \sum_{d,n} (\gamma \|s_{d,n} - \theta_d\|_2 + \rho \|s_{d,n}\|_1)$$

Our Method

- NSTC**: adopting neural network to model the generation process of *STC*
- The **cost function**: $L = l(w_{d,n}, c(d,n)) + \|s_{d,n}\|_1$



- NSTCE**: imposing deep l_1 regularized encoder on word code

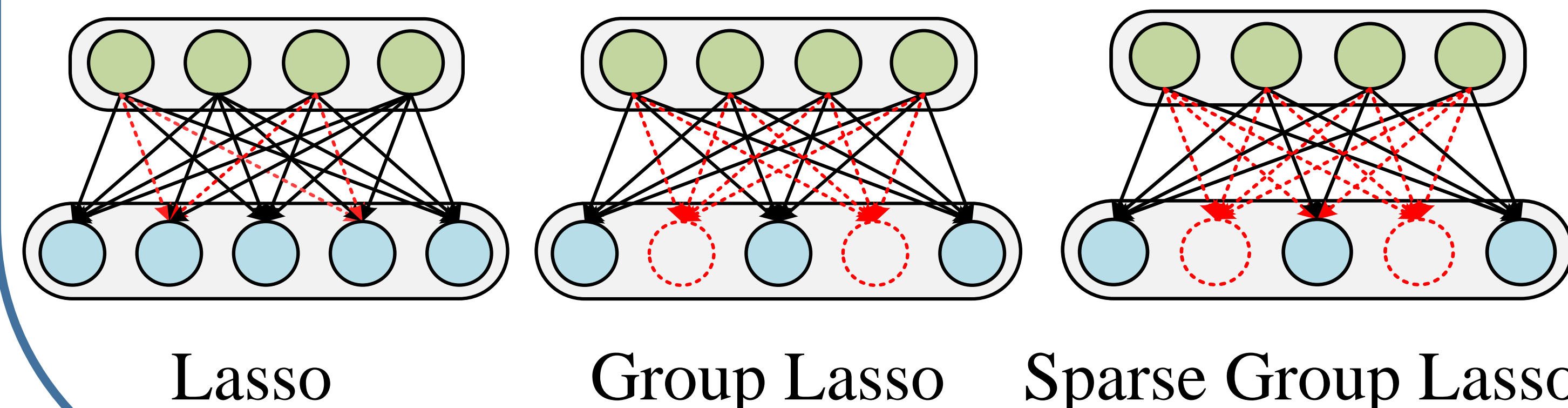
$$L = l(w_{d,n}, c(d,n)) + \|s_{d,n}\|_1 + \alpha \sum_{d,n} \|s_{d,n} - F(w_{d,n}; W, b)\|_2$$

- NGSTC**: imposing group sparse regularization on word code

$$L = l(w_{d,n}, c(d,n)) + \sum_{d,k} \lambda \|s_{d,k}\|_2$$

- NSTCSG**: imposing sparse group lasso regularization on word code

$$L = l(w_{d,n}, c(d,n)) + \sum_{d,n} \lambda_1 \|s_{d,n}\|_1 + \sum_{d,k} \lambda_2 \|s_{d,k}\|_2$$

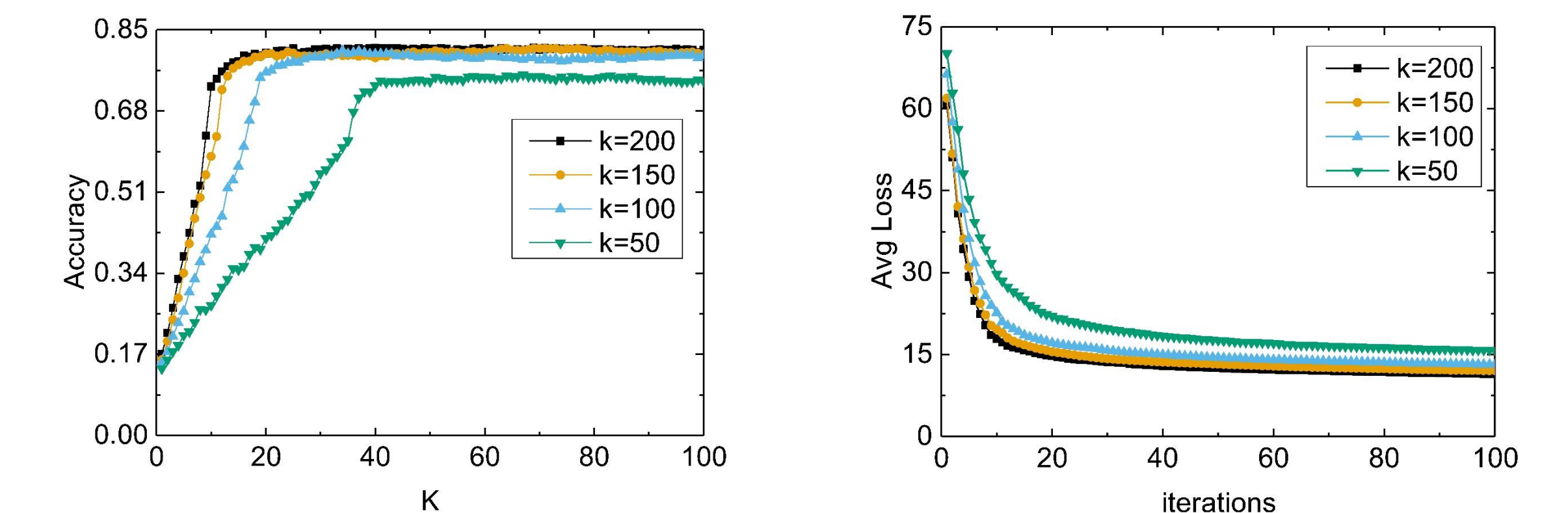


Experiments

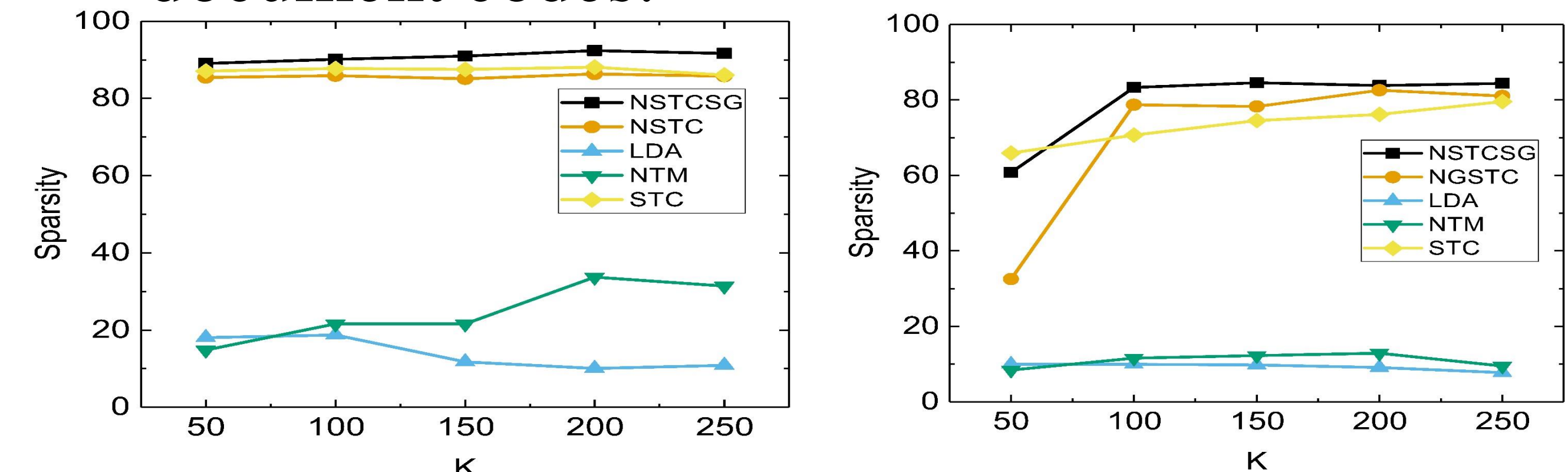
- Datasets**: 20Newsgroups, web snippet

Data	label	docs	avg words	vocab
20NG	20	18775	135	60698
snippet	8	12265	10.72	5581

- The loss and accuracy curves:



- The average sparsity ratio of word and document codes:



- The perplexity on 20NG test data:

Model	LDA	STC	DNADE	Topicvec	NSTC
ppl	1091	611	896	650	517

- t-SNE projection of the estimated document codes.:

