

# A trust-region method for stochastic variational inference with applications to streaming data

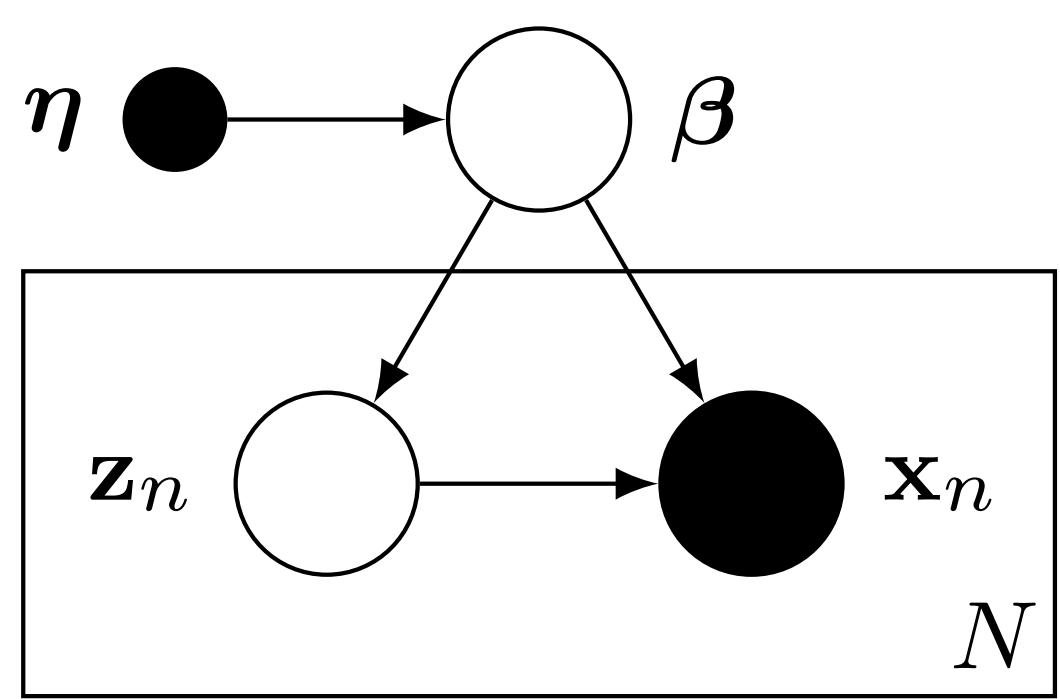
Lucas Theis<sup>1</sup>, Matthew D. Hoffman<sup>2</sup>

<sup>1</sup> Graduate School of Neural Information Processing, University of Tübingen  
<sup>2</sup> Adobe Research, San Francisco

## Introduction

**Stochastic variational inference** (SVI) has enabled fast posterior inference on massive datasets for a large class of complex Bayesian models. However, SVI can be sensitive to the choice of hyperparameters and is prone to local optima [1, 2, 3]. We introduce a new trust region based variant of SVI which is more robust. We also show how to apply SVI to streaming data and find that here our trust-region method is crucial.

## Models



We here assume an exponential family over **local parameters**  $z$  and observations  $x$  and a conjugate exponential family prior over **global parameters**  $\beta$ .

$$p(\beta) = h(\beta) \exp(\eta^\top t(\beta) - a(\eta))$$

$$p(x, z | \beta) = \prod_n h(x_n, z_n) \exp(t(\beta)^\top f(x_n, z_n))$$

Instances of this model include, for example, latent Dirichlet allocation, mixture models, HMMs, probabilistic matrix factorizations, and hierarchical linear and probit regression models.

## Mean-field approximation

$$q(\beta, z) = q(\beta) \prod_{n,m} q(z_{nm}),$$

$$q(\beta) = h(\beta) \exp(\lambda^\top t(\beta) - a(\lambda)),$$

where  $q(z_{nm})$  is controlled by  $\phi_{nm}$ . Our goal is to optimize the lower bound on the log-likelihood:

$$\mathcal{L}(\lambda) = \mathbb{E}_q \left[ \log \frac{p(\beta)}{q(\beta)} \right] + \max_{\phi} \sum_{n=1}^N \mathbb{E}_q \left[ \log \frac{p(x_n, z_n | \beta)}{q(z_n)} \right]$$

## Stochastic variational inference

Stochastic variational inference optimizes the lower bound by following a stochastic approximation,

$$\mathcal{L}_n(\lambda) = \mathbb{E}_q \left[ \log \frac{p(\beta)}{q(\beta)} \right] + N \max_{\phi_n} \mathbb{E}_q \left[ \log \frac{p(x_n, z_n | \beta)}{q(z_n)} \right],$$

via natural gradient ascent,

$$\lambda_{t+1} = \lambda_t + \rho_t \mathcal{I}(\lambda_t)^{-1} \nabla \mathcal{L}_n(\lambda_t).$$

```

1. repeat
2.   select n
3.    $\phi_n^* \leftarrow \text{argmax}_{\phi_n} \mathcal{L}_n(\lambda_t, \phi_n)$ 
4.    $\lambda_{t+1} \leftarrow (1 - \rho_t)\lambda_t + \rho_t (\eta + N\mathbb{E}_{\phi^*}[f(x_n, z_n)])$ 

```

## Trust-region method

We propose to perform the following update instead:

$$\lambda_{t+1} = \text{argmax}_{\lambda} \{ \mathcal{L}_n(\lambda) - \xi_t D_{\text{KL}}(\lambda, \lambda_t) \}$$

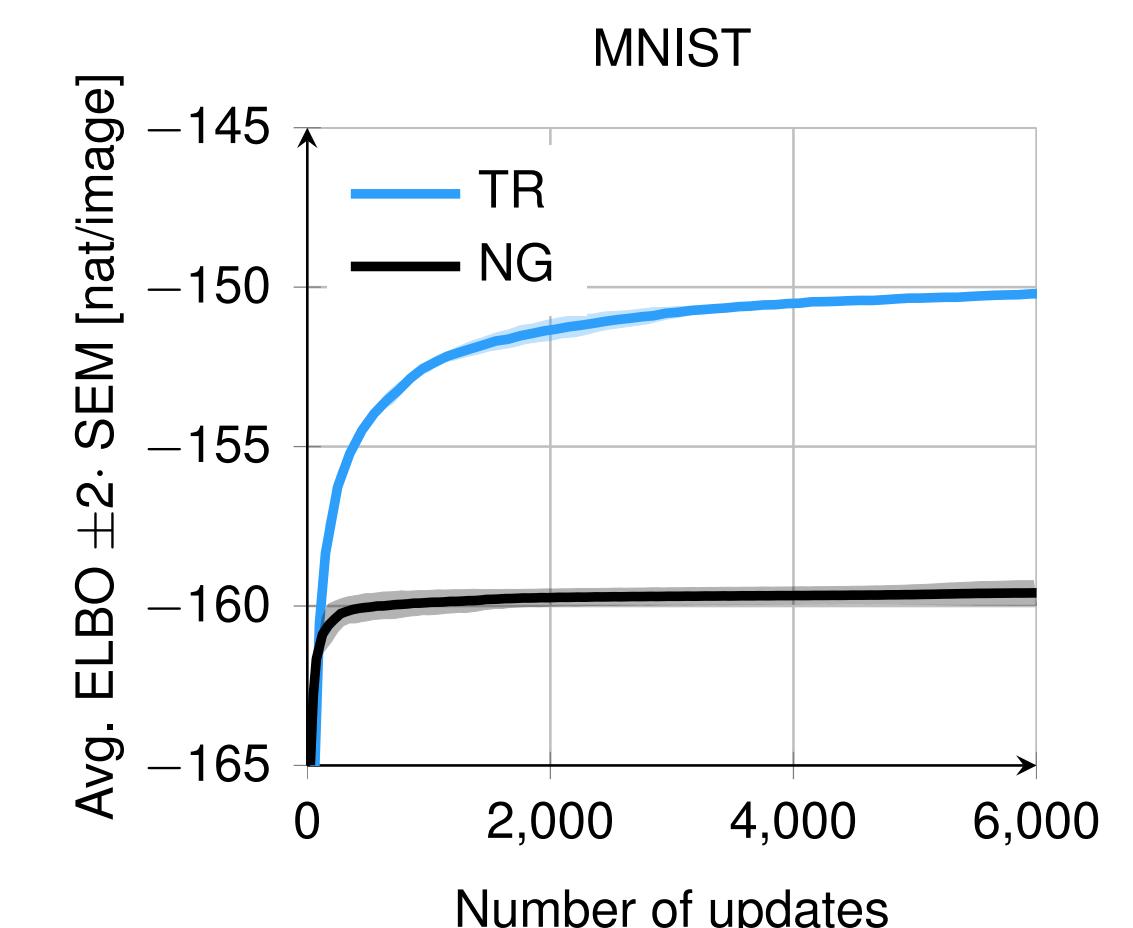
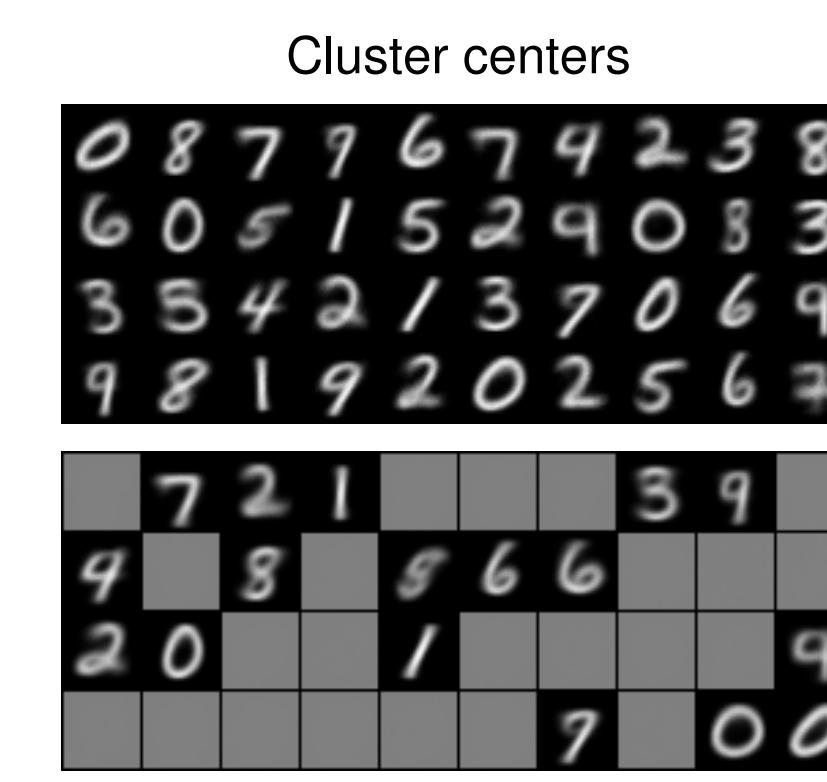
```

1. repeat
2.   select n
3.   initialize  $\phi_n^*$ 
4.   repeat
5.      $\lambda \leftarrow (1 - \rho_t)\lambda_t + \rho_t (\eta + N\mathbb{E}_{\phi^*}[f(x_n, z_n)])$ 
6.      $\phi_n^* \leftarrow \text{argmax}_{\phi_n} \mathcal{L}_n(\lambda, \phi_n)$ 
7.    $\lambda_{t+1} \leftarrow \lambda$ 

```

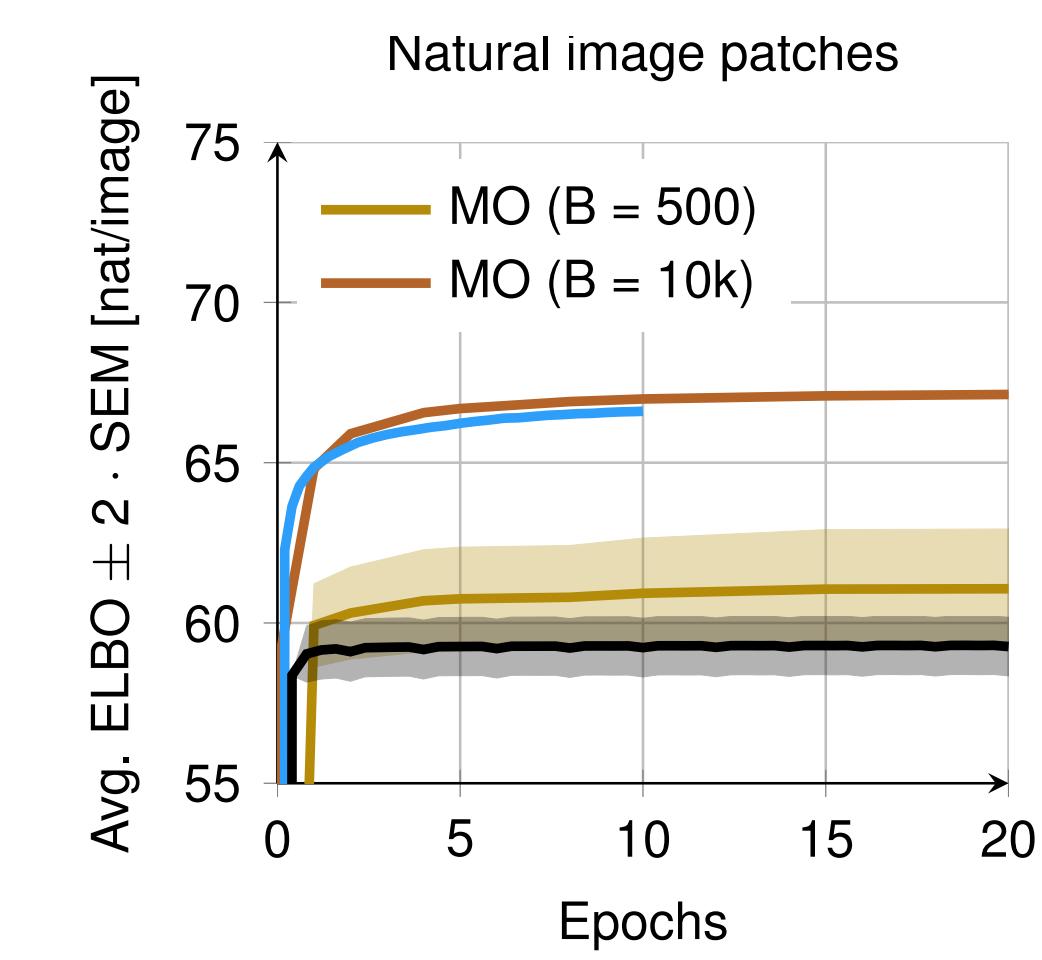
## MNIST

We trained a mixture of multivariate Bernoulli distributions on binarized MNIST digits. Left are the expected cluster means.



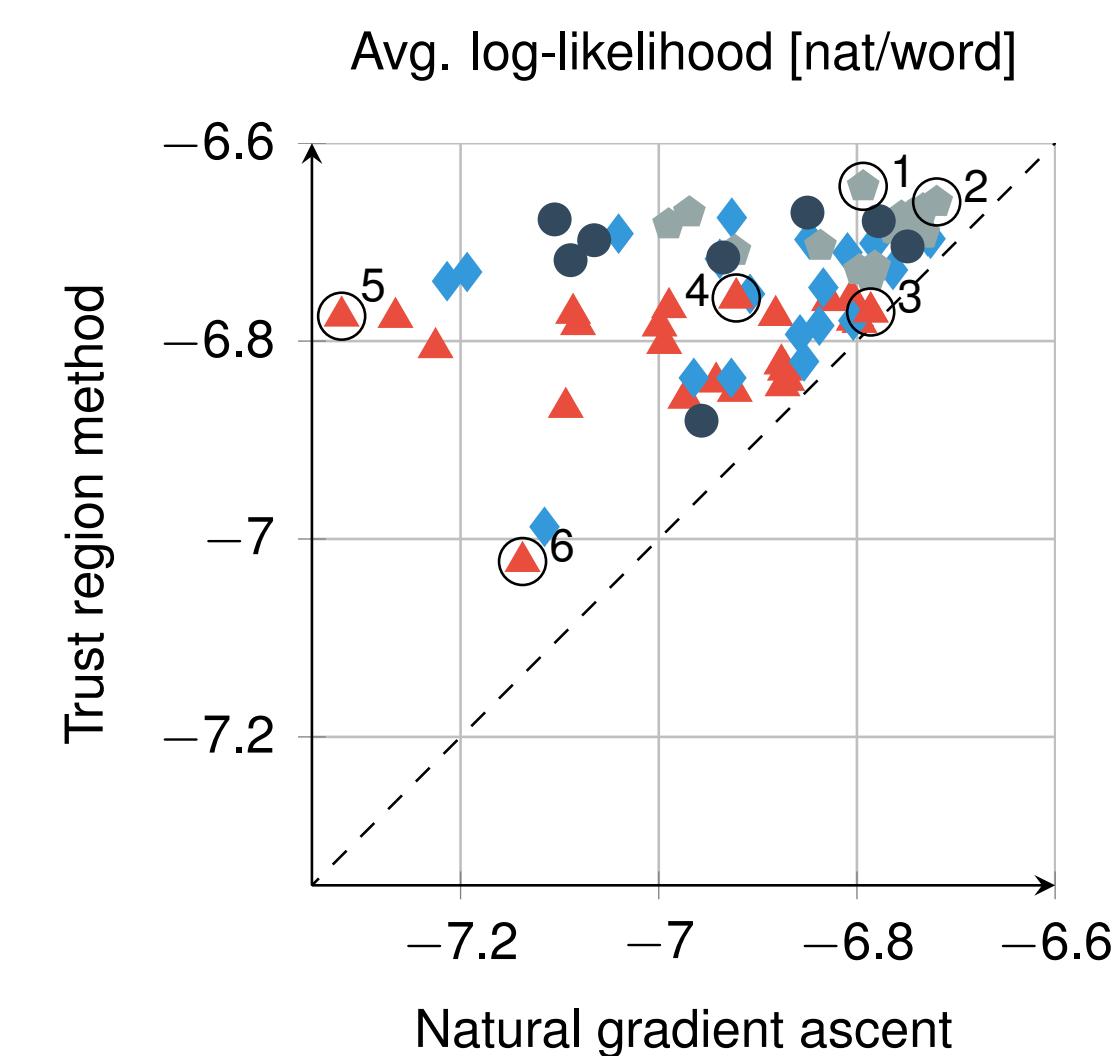
## Image patches

We trained Gaussian mixtures on 8x8 image patches and compared our method to **memoized variational inference** (MO) [2].



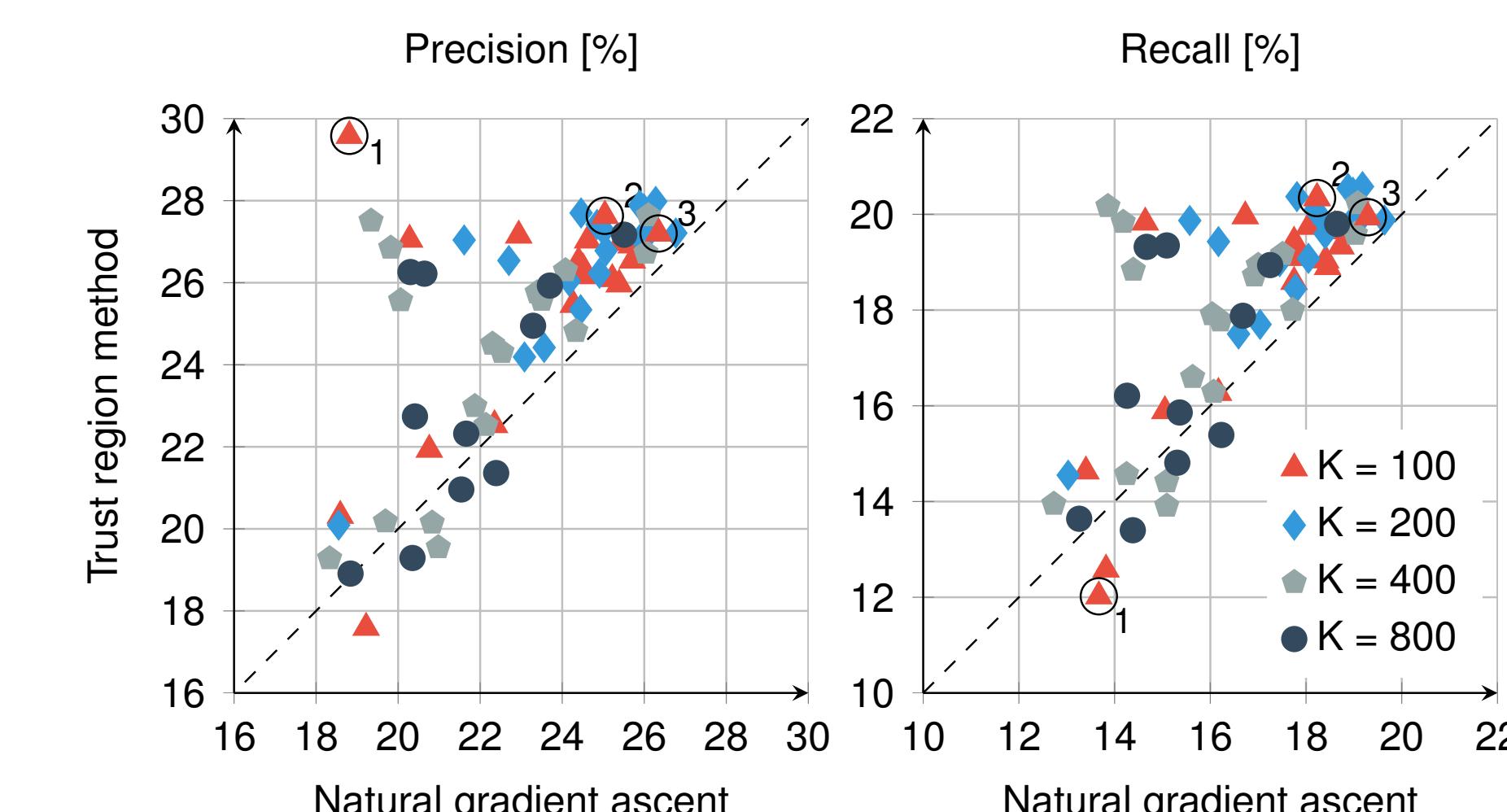
## Topic modeling

We trained **latent Dirichlet allocation** (LDA) on Wikipedia articles using the same random hyperparameter settings for natural gradient and trust-region updates.



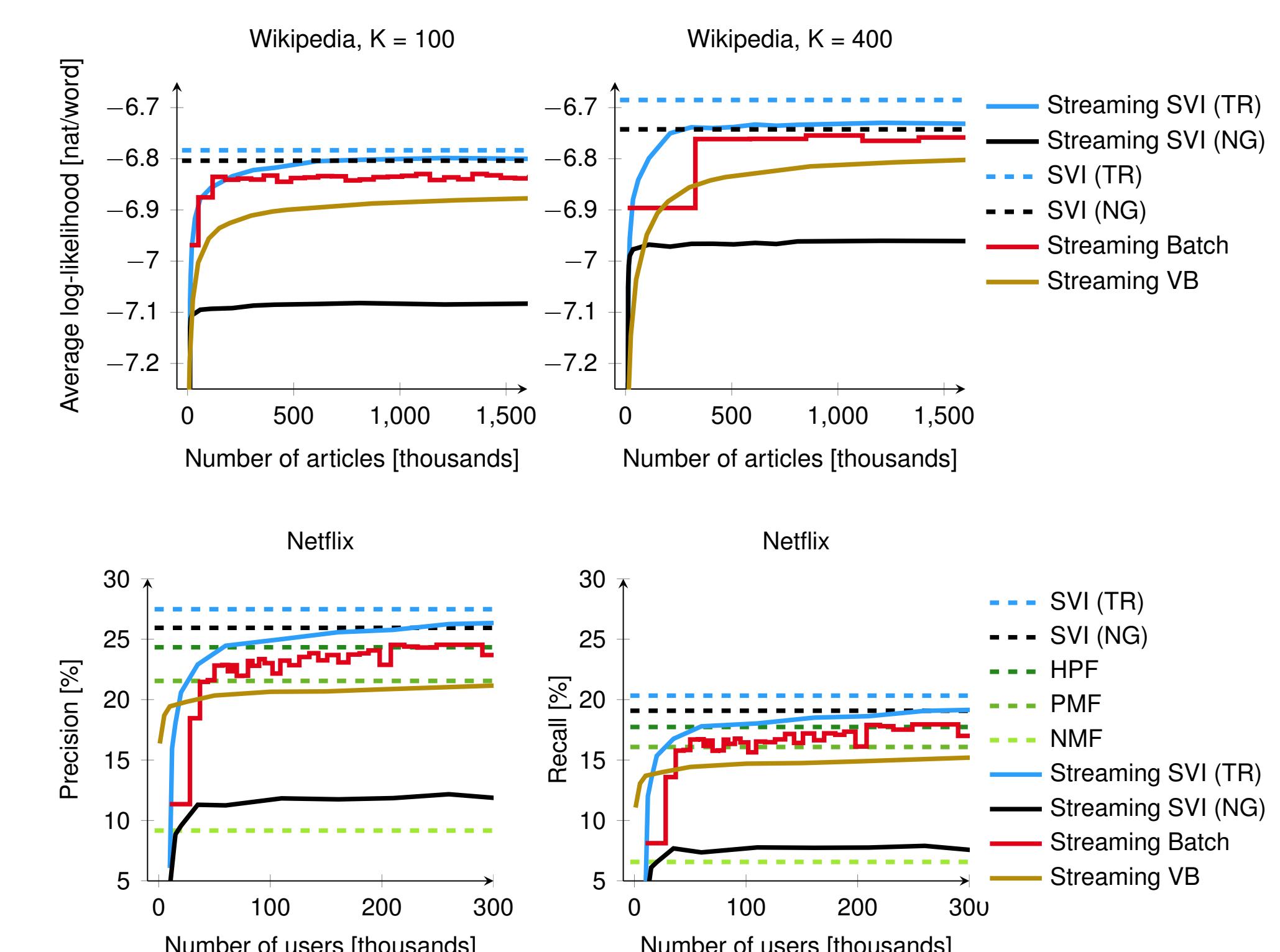
## Collaborative filtering

As above, but this time applying LDA to Netflix.



## Streaming SVI

Instead of applying SVI to a fixed dataset, we apply SVI to a continuously growing dataset (word by word, or movie by movie). We compare with **streaming variational Bayes** [4]. **Empirical Bayes** was used to automatically tune some of the hyperparameters.



## Conclusions and remarks

- SVI with trust-region updates is more **robust** to the choice of hyperparameters and generally performs a little better.
- SVI works well for **streaming** data with trust-region updates but not with natural gradient steps.
- What's missing: **automatic tuning of learning rates** for streaming data.

## Resources

Code for training latent Dirichlet allocation:

<http://github.com/lucastheis/trlida/>

## References

- [1] R. Ranganath et al., ICML, 2013
- [2] M.C. Hughes and E. B. Sudderth, NIPS, 2013
- [3] M. D. Hoffman and D. M. Blei, JMLR, 2014
- [4] T. Broderick et al., NIPS, 2013