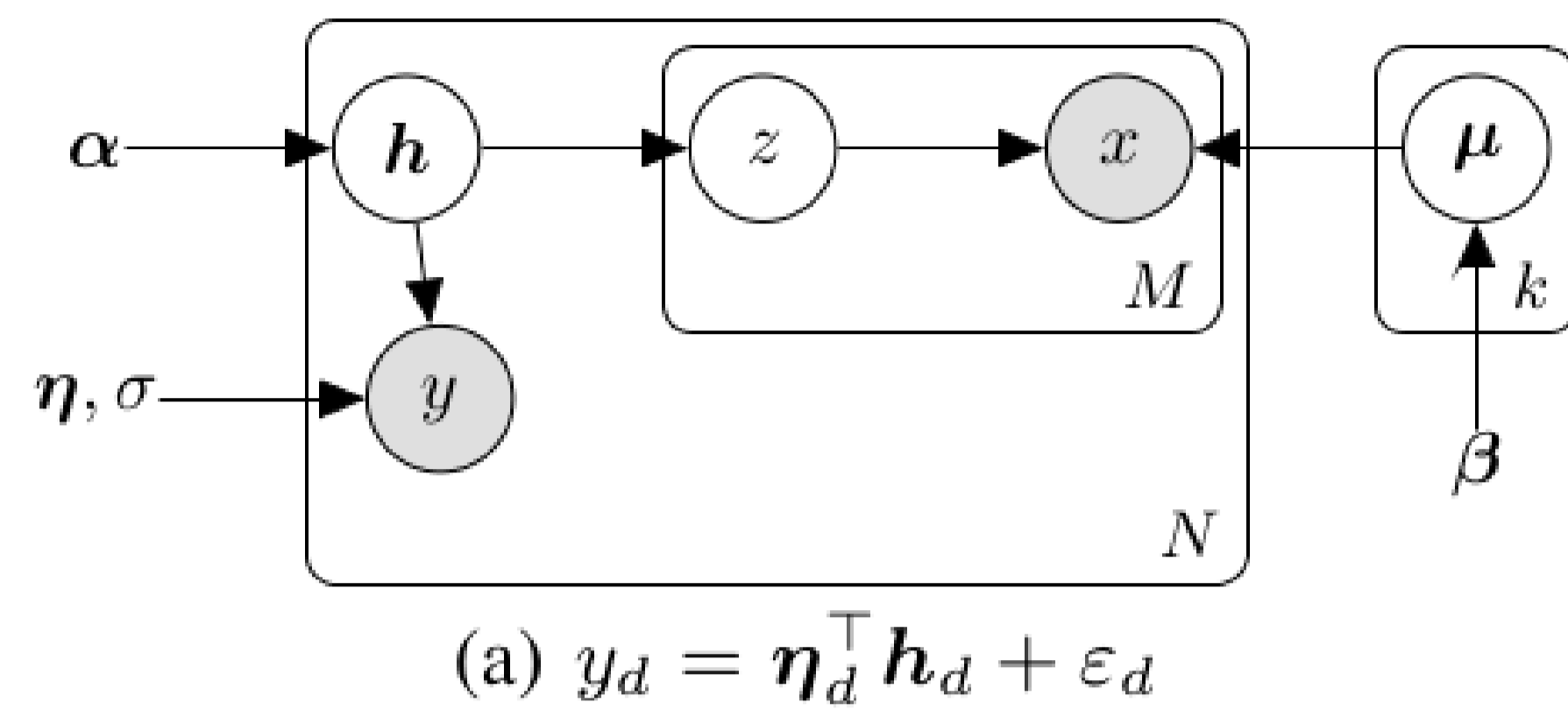


BACKGROUND OF sLDA AND TENSOR DECOMPOSITION

Supervised LDA (Blei and McAuliffe, 2007)



- Data: Document collection x_{dn} , response variables $y_d \in \mathbb{R}$.
- $\mathbf{O} = (\mu_1, \dots, \mu_k) \in \mathbb{R}^{V \times k}$: the topic dictionary.
- $h_d \in \Delta^{k-1}$: topic mixing vectors for each document.
- $z_{dn} \in \{1, 2, \dots, k\}$: topic assignments for each word.
- $\eta \in \mathbb{R}^k$: the linear regression model.
- $\alpha \in \mathbb{R}^k$: prior parameter for topic mixing vectors.

Orthogonal tensor decomposition

- For a p -th order tensor \mathbf{T} , find orthonormal basis $\{v_i\}_{i=1}^r$ and scalars $\{\lambda_i\}_{i=1}^r$ such that $\mathbf{T} = \sum_{i=1}^r \lambda_i v_i^{\otimes p}$.
- Tools: robust tensor power method (Anandkumar et al., 2012), ALS method.

Question: Can we use tensor decomposition based methods to obtain consistent estimates of sLDA parameters?

KEY TECHNIQUES

Observable moments

- $\mathbf{M}_1 = \frac{1}{\alpha_0} \sum_{i=1}^k \alpha_i \mu_i$; $\mathbf{M}_2 = \frac{1}{\alpha_0(\alpha_0+1)} \sum_{i=1}^k \alpha_i \mu_i \mu_i^\top$;
- $\mathbf{M}_3 = \frac{2}{\alpha_0(\alpha_0+1)(\alpha_0+2)} \sum_{i=1}^k \alpha_i \mu_i^{\otimes 3}$;
- $\mathbf{M}_y = \frac{2}{\alpha_0(\alpha_0+1)(\alpha_0+2)} \sum_{i=1}^k \alpha_i y_i \mu_i \mu_i^\top$.

Simultaneous diagonalization:

- First, find \mathbf{W} that whitens \mathbf{M}_2 .
- Obtain $v_i = \mathbf{W} \mu_i$ via orthogonal tensor decomposition on whitened tensor $\mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})$.
- Recover μ_i by multiplying \mathbf{W}^\dagger .

Power update for η : $\eta_i \approx v_i^\top \mathbf{M}_y(\mathbf{W}, \mathbf{W}) v_i$.

THE LEARNING ALGORITHM AND SAMPLING COMPLEXITY ANALYSIS

- 1: **Input:** Document collection x_{dn} , response variables y_d and $\alpha_0 = \|\alpha\|_1$.
- 2: Compute empirical moments and obtain $\widehat{\mathbf{M}}_2, \widehat{\mathbf{M}}_3$ and $\widehat{\mathbf{M}}_y$.
- 3: Find $\widehat{\mathbf{W}} \in \mathbb{R}^{n \times k}$ such that $\widehat{\mathbf{M}}_2(\widehat{\mathbf{W}}, \widehat{\mathbf{W}}) = \mathbf{I}_k$.
- 4: Find robust eigenvalues and eigenvectors $(\widehat{\lambda}_i, \widehat{v}_i)$ of $\widehat{\mathbf{M}}_3(\widehat{\mathbf{W}}, \widehat{\mathbf{W}}, \widehat{\mathbf{W}})$ using the robust tensor power method (Anandkumar et al., 2012).
- 5: Recover parameters: $\widehat{\alpha}_i \leftarrow \frac{4\alpha_0(\alpha_0+1)}{(\alpha_0+2)^2 \widehat{\lambda}_i^2}$, $\widehat{\mu}_i \leftarrow \frac{\alpha_0+2}{2} \widehat{\lambda}_i (\widehat{\mathbf{W}}^\dagger)^\top \widehat{v}_i$, $\widehat{\eta}_i \leftarrow \frac{\alpha_0+2}{2} \widehat{v}_i^\top \widehat{\mathbf{M}}_y(\widehat{\mathbf{W}}, \widehat{\mathbf{W}}) \widehat{v}_i$.
- 6: **Output:** $\widehat{\eta}, \widehat{\alpha}$ and $\{\widehat{\mu}_i\}_{i=1}^k$.

THEOREM (SAMPLE COMPLEXITY ANALYSIS)

Suppose we have an sLDA model $\mathcal{M} = (\{\mu_i\}_{i=1}^k, \alpha, \eta)$ and n documents i.i.d. sampled from \mathcal{M} . Let $\mathbf{O} = (\mu_1, \dots, \mu_k)$. Then the estimates $\widehat{\mathcal{M}} = (\{\widehat{\mu}_i\}_{i=1}^k, \widehat{\alpha}, \widehat{\eta})$ is ϵ -close to \mathcal{M} with high probability given that $n = \Omega(\text{poly}(\epsilon, k, \sigma_k(\mathbf{O})^{-1}, \alpha_{\max}, \alpha_{\min}^{-1}, \|\eta\|))$.

EXPERIMENTS

Synthetic datasets

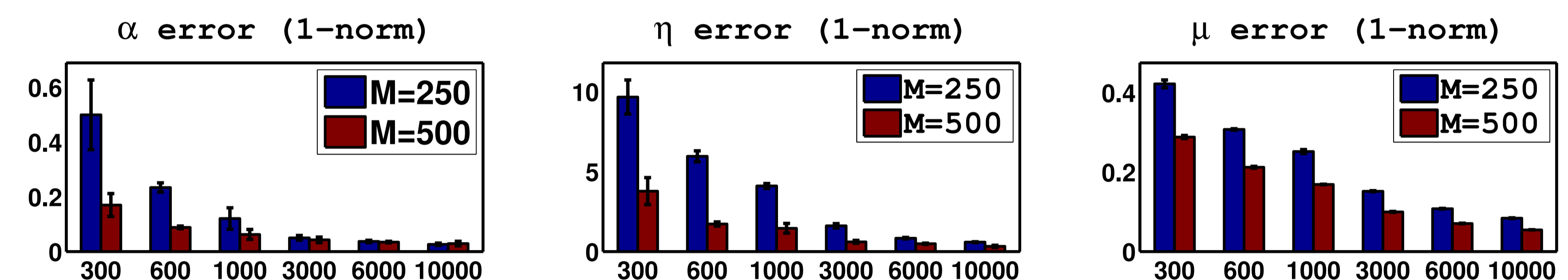


Figure: Parameter estimation error (1-norm) on synthetic datasets.

Amazon movie review dataset (McAuley and Leskovec, 2013)

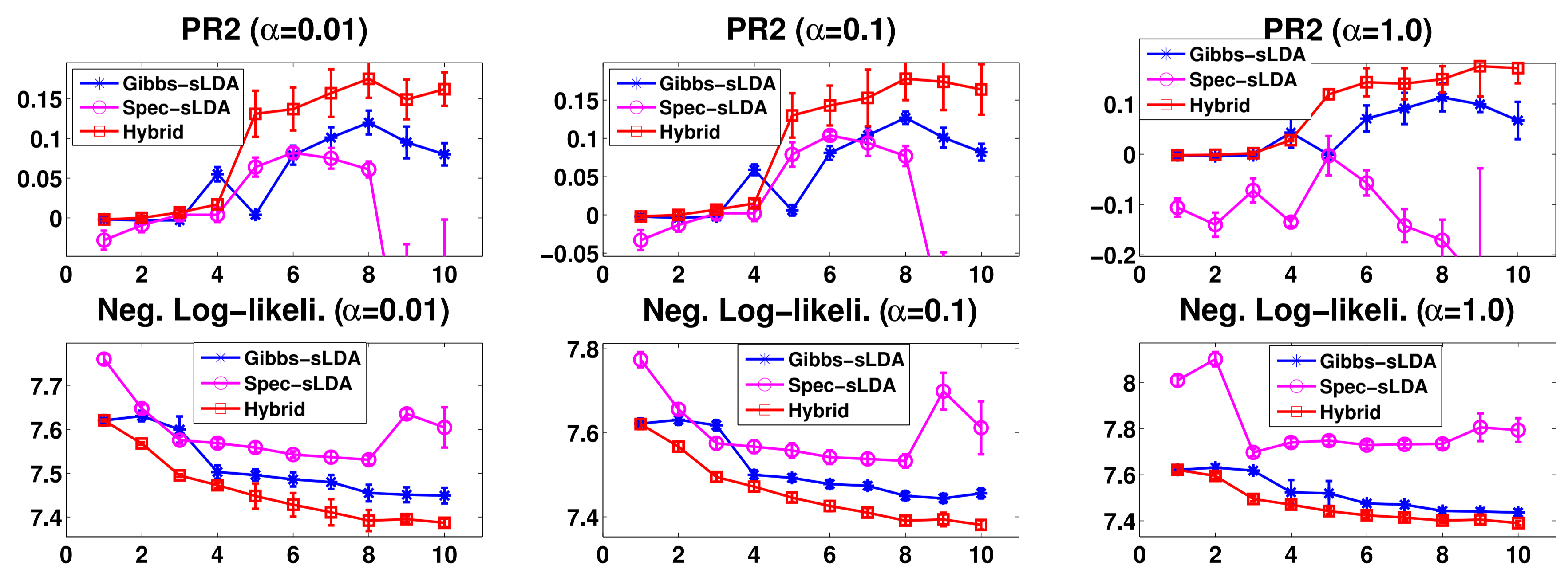


Figure: pR^2 scores (top) and negative per-word log-likelihood (bottom) on Amazon movie review dataset.

Table: Running time for Gibbs sampling and spectral method.

	$k = 10$					$k = 50$				
$n (\times 10^4)$	1	5	10	50	100	1	5	10	50	100
Gibbs-sLDA	0.6	3.0	6.0	30.5	61.1	2.9	14.3	28.2	145.4	281.8
Spec-sLDA	1.5	1.6	1.7	2.9	4.3	3.1	3.6	4.3	9.5	16.2