
Supplementary Material for “Distributed Online Learning for Latent Dirichlet Allocation”

JinYeong Bak, Dongwoo Kim and Alice Oh
Department of Computer Science
Korea Advanced Institute of Science and Technology
Daejeon, South Korea
{jy.bak, dw.kim}@kaist.ac.kr, alice.oh@kaist.edu

This supplementary material is for “Distributed Online Learning for Latent Dirichlet Allocation” which is submitted to NIPS 2012 workshop on parallel and large-scale machine learning.

1 Distributed Online Learning for LDA Algorithm

Here are the algorithms for the three parts of Distributed Online Learning for LDA.

Algorithm 1 Distributed Online Learning for LDA - Driver

- 1: **repeat**
 - 2: Set $\alpha, \eta, \lambda, \rho$ to distributed cache
 - 3: Set mini-batch size S , number of topic K and number of words W to distributed cache
 - 4: Create new MapReduce job for new documents
 - 5: Run MapReduce job
 - 6: Get new document's γ_d and λ_{new} from MapReduce job output
 - 7: Update α, η and change λ_{new} to λ
 - 8: **until** New documents are not observed
-

Algorithm 2 Distributed Online Learning for LDA - Mapper

Input

- 1: KEY - document ID $d \in \{1, \dots, S\}$
- 2: VALUE - document content

Configure

- 1: Load α, λ , number of topic K and words W from distributed cache
- 2: Calculate $\forall k, w \exp E_q[\log \beta_{kw}]$

Map

- 1: Initialize $\gamma_{dk} = 1$
- 2: Read document d content
- 3: **repeat**
- 4: Set $\phi_{dwk} \propto \exp\{E_q[\log \theta_{dk}] + E_q[\log \beta_{kw}]\}$
- 5: Set $\gamma_{dk} = \alpha + \sum_w n_{dw} \phi_{dwk}$
- 6: **until** convergence with γ

Output

- 1: variational parameter for document d 's $\theta_d : \gamma_d$
 - 2: sufficient statistics for $\lambda : n_{dw} \phi_{dwk}$
-

Algorithm 3 Distributed Online Learning for LDA - Reducer

Input

- 1: KEY - sufficient statistics
- 2: VALUE - $\forall d, w, k n_{dw}\phi_{dwk}$

Configure

- 1: Load η , λ , ρ , mini-batch size S , number of topic K and document D from distributed cache

Reduce

- 1: Set $\tilde{\lambda}_{kw} = \eta + \frac{D}{S} \sum_s n_{sw}\phi_{swk}$
- 2: Set $\lambda_{new} = (1 - \rho)\lambda + \rho\tilde{\lambda}$

Output

- 1: variational parameter for topic : λ_{new}
-

2 Twitter Conversation Topics

In this section, we show some topics from the Twitter conversation corpus discovered by DoLDA. These topics are from the settings of $k = 100$.

Table 1: Several topics from Twitter conversation corpus using distributed online LDA

Topic 15	Topic 21	Topic 36	Topic 38	Topic 45
tea	dm	breakfast	gym	black
party	text	toast	body	wear
vote	number	butter	weight	white
country	send	cereal	running	dress
tax	email	pancakes	workout	shirt
labour	sent	branch	chest	clothes
soup	phone	pjs	knee	pants
government	ok	jar	diet	light
ed	check	nutella	leg	look
political	skype	cinnamon	exercise	suit
Topic 54	Topic 60	Topic 74	Topic 79	Topic 95
game	xx	drink	movie	room
play	xxx	drinking	film	clean
team	aw	apple	singing	mood
football	babe	water	harry	dirty
mate	hun	bottle	read	shower
player	aww	orange	michael	kitchen
match	lovely	juice	book	house
think	thankyou	tall	films	wash
beat	xo	pie	potter	laundry
good	ahh	blonde	cinema	basement

References

- [1] M. Hoffman, D. Blei, and F. Bach. Online learning for latent dirichlet allocation. *Advances in Neural Information Processing Systems*, 23:856–864, 2010.