Distributed Inference for LDA D. Newman, A. Asuncion, P. Smyth, M. Welling

Async. Dist. Learning of Topic Models A. Asuncion, P. Smyth, M. Welling

Presented by: Golnar Sheikhshab

Outline

- Latent Dirichlet Allocation
- Distributed Inference for LDA
- Asynchronous Distributed Learning for Topic Models

Latent Dirichlet Allocation (LDA)

- State of the art topic modeling method
- Clusters the words into topics
 - Information Retrieval
 - Data visualization
 - . . .
- Also has applications in
 - Image Processing
 - Bioinformatics

LDA



- For each Document:
 - 1. Randomly choose a distribution over topics.
 - 2. For each word in the document
 - a) Randomly choose a topic from the distribution over topics in step #1.
 - b) Randomly choose a word from the corresponding distribution over the vocabulary.

 $\theta_{k,j} \sim D[\alpha] \qquad \phi_{w,k} \sim D[\eta] \qquad z_{ij} \sim \theta_{k,j} \qquad x_{ij} \sim \phi_{w,z_{ij}}$

LDA



LDA Parameter Estimation

- Gibbs Sampling
- Observed variables: wij
- Latent variables:

– Zij

- θkj and φwk (marginalized out)

$$P(z_{ij} = k | z^{\neg ij}, w) \propto \frac{N_{wk}^{\neg ij} + \eta}{\sum_{w} N_{wk}^{\neg ij} + W\eta} \left(N_{jk}^{\neg ij} + \alpha \right)$$

Distributed Inference Algorithms

- Approximate Distributed Inference (AD-LDA)
- Hierarchical Distributed Inference (HD-LDA)
 - Not to be confused with Hierarchical LDA

AD-LDA

 Run LDA on subsets of data on different processors + a synchronization step

 $p(z_{ijp} = k | \mathbf{z}_p^{\neg ijp}, \mathbf{x}, \alpha, \beta) \propto (\alpha + n_{k|jp}^{\neg ijp})(\beta + n_{x_{ij}|kp}^{\neg ijp})(W\beta + n_{kp}^{\neg ijp})^{-1}$

$$n_{w|k} \leftarrow n_{w|k} + \sum_{p} (n_{w|kp} - n_{w|k}), \qquad n_{w|kp} \leftarrow n_{w|k}$$

• Note: $n_{w|kp}$ is not the result of separate LDA models running on separate data.

$$\sum_{w,k} n_{w|kp} = N$$

HD-LDA



$$p(z_{ijp} = k | \mathbf{z}_p^{\neg ijp}, \mathbf{x}, \alpha, \beta, \Phi) \propto (\alpha_p + n_{k|jp}^{\neg ijp}) (\beta_k \Phi_{x_{ij}|k} + n_{x_{ij}|kp}^{\neg ijp}) (\beta_k + n_{kp}^{\neg ijp})^{-1}$$

$$\begin{split} \alpha_p \leftarrow & \frac{c - 1 + \alpha_p \sum_{jk} \left[\Psi(\alpha_p + n_{k|jp}) - \Psi(\alpha_p) \right]}{d + K \sum_j \left[\Psi(K\alpha_p + n_{jp}) - \Psi(K\alpha_p) \right]} \\ \beta_k \leftarrow & \frac{a - 1 + \beta_k \sum_{wp} \Phi_{w|k} \left[\Psi(\beta_k \Phi_{w|k} + n_{w|kp}) - \Psi(\beta_k \Phi_{w|k}) \right]}{b + \sum_p \left[\Psi(\beta_k + n_{kp}) - \Psi(\beta_k) \right]} \\ \Phi_{w|k} \leftarrow & \frac{\gamma - 1 + \sum_p \beta_k \Phi_{w|k} \left[\Psi(\beta_k \Phi_{w|k} + n_{w|kp}) - \Psi(\beta_k \Phi_{w|k}) \right]}{\gamma W - W + \sum_{wp} \beta_k \Phi_{w|k} \left[\Psi(\beta_k \Phi_{w|k} + n_{w|kp}) - \Psi(\beta_k \Phi_{w|k}) \right]} \end{split}$$

9

 α_p

c,d

Ρ

HD-LDA

- Hyper-parameter selection
 - Guided by prior experience with AD-LDA
 - Note: $\sum_{w,k} n_{w|kp} \approx \frac{N}{P}$ a and b were chosen such that $\beta_k = \frac{P-1}{P}N$ c and d chosen such that $\alpha_p = 0.1$ $\gamma = 1 + \frac{1}{T}$
 - Critical in getting good results

Experiments

- Why would AD-LDA work?
- Toy example 3 words, 2 topics



Figure 2: (Left) L_1 distance to the mode for LDA and for P=2 AD-LDA. (Center) Projection of topics onto simplex, showing convergence to mode. (Right) Same setup as center panel, but with P = 10 processors.

Prediction power



Figure 3: Test perplexity of models versus number of processors P for KOS (left) and NIPS (right). P=1 corresponds to LDA (circles), and AD-LDA (crosses), and HD-LDA (squares) are shown at P=10 and 100.

 Prediction power vs. iteration number and number of topics



Figure 4: (Left) Test perplexity versus iteration.

(Right) Test perplexity versus number of topics.

IR application TREC'S AP and FR collections



Figure 5: (Left) Precision/recall results.

- Speedup experiment:
 - On one 16-processor machine: not reflecting the communication time



(Right). Parallel speedup results.

Asynchronous Methods

- Advantages
 - No global synchronization step is required
 - Extremely fault-tolerant
 - Heterogeneous machines can be used
 - New processors and data can arrive any time
- Empirically
 - Robust
 - Converging fast

Asynchronous methods (Cont.)

- Still 2 steps:
 - Gibbs sampling
 - Synchronization
 - Gossip-based

$$P(z_{pij} = k | z_p^{\neg ij}, w_p) \propto \frac{(N^{\neg p} + N^p)_{wk}^{\neg ij} + \eta}{\sum_w (N^{\neg p} + N^p)_{wk}^{\neg ij} + W\eta} \left(N_{pjk}^{\neg ij} + \alpha \right)$$

Algorithm 1 Async-LDAfor each processor
$$p$$
 in parallel do
repeat
Sample z^p locally (Equation 2)
Receive N_{wk}^g from random proc g
Send N_{wk}^p to proc g
if p has met g before then
 $N_{wk}^{\neg p} \leftarrow N_{wk}^{\neg p} - \tilde{N}_{wk}^g + N_{wk}^g$
else
 $N_{wk}^{\neg p} \leftarrow N_{wk}^{\neg p} + N_{wk}^g$
end if
until convergence
end for

Asynchronous methods (Cont.)

 One over Hierarchical Dirichlet Process (HDP)

• HDP

$$\begin{aligned} \beta_k &\sim D[\gamma/L] & \theta_{k,j} \sim D[\alpha\beta_k] \\ z_{ij} &\sim \theta_{k,j} & x_{ij} \sim \phi_{w,z_{ij}} \end{aligned}$$

When L goes to infinity



Parallel-HDP & Async-HDP

Algorithm 2 Parallel-HDP	Algorithm 3 Async-HDP
repeat for each processor p in parallel do Sample z^p locally Send N_{wk}^p , β_k^p to master node end for $N_{wk} \leftarrow \sum_p N_{wk}^p$ $\beta_k \leftarrow (\sum_p \beta_k^p) / P$ Resample α , β_k , γ globally Distribute N_{wk} , α , β_k , γ to all processors until convergence	for each processor p in parallel do repeat Sample z^p and then α^p , β_k^p , γ^p locally Receive N_{wk}^g , α^g , β_k^g from random proc g Send N_{wk}^p , α^p , β_k^p to proc g if p has met g before then $N_{wk}^{\neg p} \leftarrow N_{wk}^{\neg p} - \tilde{N}_{wk}^g + N_{wk}^g$ else $N_{wk}^{\neg p} \leftarrow N_{wk}^{\neg p} + N_{wk}^g$ end if $\alpha^p \leftarrow (\alpha^p + \alpha^g)/2$ and $\beta_k^p \leftarrow (\beta_k^p + \beta_k^g)/2$ until convergence
	end for
$P(z_{pij} = k z_p^{\neg ij}, w_p) \propto \begin{cases} rac{(N^{\neg p} + N)}{\sum_w (N^{\neg p} + N)} & rac{lpha^p \beta_{ m new}^p}{W}, \end{cases}$	$\frac{N^p)_{wk}^{\neg ij} + \eta}{(N^p)_{wk}^{\neg ij} + W\eta} \left(N^{\neg ij}_{pjk} + \alpha^p \beta_k^p \right), \text{if } k \le K_p$
	if k is new 19

Parallel-HDP & Async-HDP (Cont.)

- New topics might be introduced while infering
- Bipartite matching across two topic sets might be required
- They don't do it but it works!
 - They just combine the sets by topic ids

Experiments

- Async-LDA
- Prediction power



Figure 2: (a) Left: Async-LDA perplexities on KOS. (b) Middle: Async-LDA perplexities on NIPS. (c) Right: Async-LDA perplexities on KOS with many procs. Cache=5 when $P \ge 100$. 3000 iterations run when $P \ge 500$.

- Async-LDA
- Speed analysis



Figure 3: (a) Left: Convergence plot for Async-LDA on KOS, K=16. (b) Middle: Same plot with x-axis as relative time. (c) Right: Speedup results for NYT and PUBMED on a cluster, using Message Passing Interface.

Parallel-HDP & Async-HDP



Figure 4: (a) Left: Perplexities for Parallel-HDP and Async-HDP. Cache=5 used for Async-HDP P=100. (b) Middle: Convergence plot for Parallel-HDP on KOS. (c) Right: Convergence plot for Async-HDP on KOS.

- Parallel-HDP & Async-HDP
 - The number of topics does stabilize after thousands of iterations
 - Topics are generated at a slightly faster rate for Async-HDP than for Parallel-HDP
 - Because Async-HDP takes a less aggressive approach on pruning small topics since processors need to be careful when pruning topics locally.

• 3 Different experiments here



Figure 5: (a) Left: Online learning for Async-LDA (K=16) on KOS. (b) Middle: Comparing random vs. non-random distribution of documents for Async-LDA on NIPS, K=20. (c) Right: Async-LDA on KOS, K=16, where processors have varying amounts of data. In all 3 cases, Async-LDA converges to a good solution.

Conclusions

• Go Async-LDA!

