

Collapsed Variational Dirichlet Process Mixture Models*

Kenichi Kurihara

Dept. of Computer Science
Tokyo Institute of Technology, Japan
kurihara@mi.cs.titech.ac.jp

Max Welling

Dept. of Computer Science
UC Irvine, USA
welling@ics.uci.edu

Yee Whye Teh

Dept. of Computer Science
National University of Singapore
tehyw@comp.nus.edu.sg

Discussion led by
Youssef Barhomi

Motivation

- Gibbs sampling is not efficient
- Sampling requires careful monitoring of the convergence of the Markov chain



Variational Bayesian methods

Motivation

Variational methods:

- A good approximation of the DP
- Deterministic
- Handles “modern” datasets faster than Gibbs Sampling

TSB and FSD

TSB: Truncated Stick Breaking process with standard variational bayesian model

FSD: Finite Symmetric Dirichlet representation with standard variational bayesian model

I - TSB

TSB: Truncated Stick Breaking process with standard variational bayesian model

$$v_i \sim \mathcal{B}(v_i; 1, \alpha) \quad i = 1, \dots, T-1 \quad (1)$$

$$v_T = 1 \quad (2)$$

$$\pi_i = v_i \prod_{j < i} (1 - v_j) \quad i = 1, \dots, T \quad (3)$$

$$\pi_i = 0 \quad i > T \quad (4)$$

$$P(\mathbf{X}, \mathbf{z}, \mathbf{v}, \boldsymbol{\eta}) = \left[\prod_{n=1}^N p(\mathbf{x}_n | \eta_{z_n}) p(z_n | \boldsymbol{\pi}(\mathbf{v})) \right] \left[\prod_{i=1}^T p(\eta_i) \mathcal{B}(v_i; 1, \alpha) \right] \quad (5)$$

\mathbf{X} are data points, \mathbf{z} are the assignments, \mathbf{v} are the stick breaking weights, and $\boldsymbol{\eta}$ is cluster parameters

II - FSD

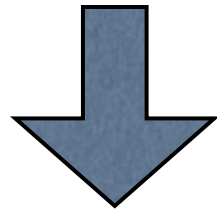
$$P(X, \mathbf{z}, \mathbf{v}, \boldsymbol{\eta}) = \left[\prod_{n=1}^N p(\mathbf{x}_n | \eta_{z_n}) p(z_n | \boldsymbol{\pi}(\mathbf{v})) \right] \left[\prod_{i=1}^T p(\eta_i) \mathcal{B}(v_i; 1, \alpha) \right] \quad (5)$$

+

$\boldsymbol{\pi} \sim \mathcal{D}(\boldsymbol{\pi}; \frac{\alpha}{K}, \dots, \frac{\alpha}{K})$ mixture weights following a symmetric dirichelet

+

assume a large number of clusters K



[Ishwaran and Zarepour, 2002]

$$P(X, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\eta}) = \left[\prod_{n=1}^N p(\mathbf{x}_n | \eta_{z_n}) p(z_n | \boldsymbol{\pi}) \right] \left[\prod_{i=1}^K p(\eta_i) \right] \mathcal{D}(\boldsymbol{\pi}; \frac{\alpha}{K}, \dots, \frac{\alpha}{K}) \quad (7)$$

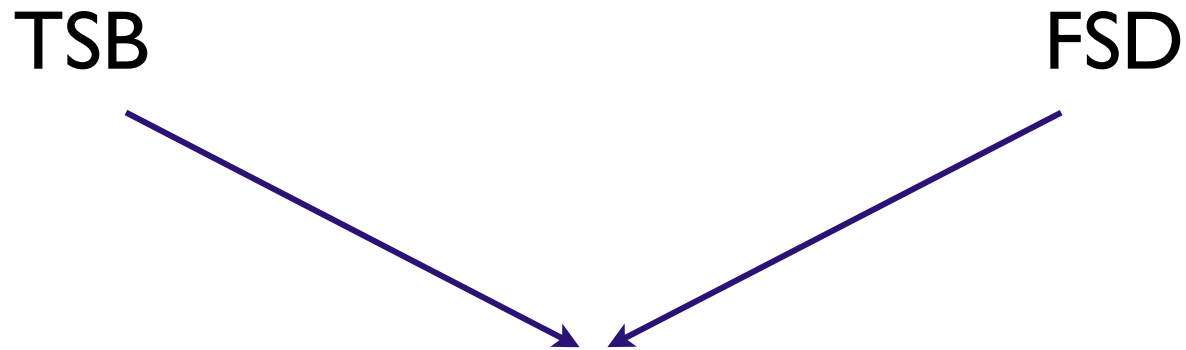
Marginalizing the mixture of weights

$$P(X, \mathbf{z}, \mathbf{v}, \boldsymbol{\eta}) = \left[\prod_{n=1}^N p(\mathbf{x}_n | \eta_{z_n}) p(z_n | \boldsymbol{\pi}(\mathbf{v})) \right] \left[\prod_{i=1}^T p(\eta_i) \mathcal{B}(v_i; 1, \alpha) \right] \quad (5)$$

$$P(X, \mathbf{z}, \boldsymbol{\pi}, \boldsymbol{\eta}) = \left[\prod_{n=1}^N p(\mathbf{x}_n | \eta_{z_n}) p(z_n | \boldsymbol{\pi}) \right] \left[\prod_{i=1}^K p(\eta_i) \right] \mathcal{D}(\boldsymbol{\pi}; \frac{\alpha}{K}, \dots, \frac{\alpha}{K}) \quad (7)$$

TSB

FSD


$$P(X, \mathbf{z}, \boldsymbol{\eta}) = \left[\prod_{n=1}^N p(\mathbf{x}_n | \eta_{z_n}) \right] p(\mathbf{z}) \left[\prod_{i=1}^{\infty} p(\eta_i) \right] \quad (9)$$

Collapsed model

Marginalizing the mixture of weights

$$p_{\text{TSB}}(\mathbf{z}) = \prod_{i < T} \frac{\Gamma(1 + N_i) \Gamma(\alpha + N_{>i})}{\Gamma(1 + \alpha + N_{\geq i})} \quad (10)$$

$$p_{\text{FSD}}(\mathbf{z}) = \frac{\Gamma(\alpha) \prod_{k=1}^K \Gamma(N_k + \frac{\alpha}{K})}{\Gamma(N + \alpha) \Gamma(\frac{\alpha}{K})^K} \quad (12)$$

with

$$N_i = \sum_{n=1}^N \mathbb{I}(z_n = i) \quad N_{>i} = \sum_{n=1}^N \mathbb{I}(z_n > i) \quad (11)$$

and $N_{\geq i} = N_i + N_{>i}$. For FSD we find instead,

Lower bound formulation

log marginal likelihood

$$\mathcal{L}(X) \geq \mathbf{B}(X) = \sum_{\mathbf{z}} \int_{d\theta} Q(\mathbf{z})Q(\theta) \log \frac{P(X, \mathbf{z}, \theta)}{Q(\mathbf{z})Q(\theta)} \quad (13)$$

lower bound

where θ is either $\{\eta, \mathbf{v}\}$, $\{\eta, \boldsymbol{\pi}\}$ or $\{\eta\}$

$$Q_{\text{TSB}}(\mathbf{z}, \boldsymbol{\eta}, \mathbf{v}) = \left[\prod_n^N q(z_n) \right] \left[\prod_{i=1}^T q(\eta_i)q(v_i) \right] \quad (14)$$

$$Q_{\text{FSD}}(\mathbf{z}, \boldsymbol{\eta}, \boldsymbol{\pi}) = \left[\prod_n^N q(z_n) \right] \left[\prod_{k=1}^K q(\eta_k) \right] q(\boldsymbol{\pi}) \quad (15)$$

Lower bound formulation

$$\mathbf{B}(X) = \sum_{n=1}^N \sum_{z_n} \int_{d\eta_{z_n}} q(z_n)q(\eta_{z_n}) \log p(\mathbf{x}_n|\eta_{z_n}) + \sum_i \int_{d\eta_i} q(\eta_i) \log \frac{p(\eta_i)}{q(\eta_i)} - \sum_{n=1}^N \sum_{z_n} q(z_n) \log q(z_n) + \text{Extra Term}$$

$$\text{Term}_{\text{TSB}} = \sum_{n=1}^N \sum_{z_n=1}^T q(z_n) \int_{d\mathbf{v}} \left[\prod_{i=1}^{z_n} q(v_i) \right] \log p(z_n|\mathbf{v}) + \sum_{i=1}^T \int_{dv_i} q(v_i) \log \frac{p(v_i)}{q(v_i)} \quad (17)$$

$$\text{Term}_{\text{FSD}} = \sum_n \sum_{z_n=1}^K \int_{d\pi} q(z_n)q(\pi) \log p(z_n|\pi) + \int_{d\pi} q(\pi) \log \frac{p(\pi)}{q(\pi)} \quad (18)$$

$$\text{Term}_{\text{CTSB/CFSD}} = \sum_{\mathbf{z}} \left[\prod_{n=1}^N q(z_n) \right] \log p(\mathbf{z}) \quad (19)$$

Update equations

$$q(\eta_i) \propto p(\eta_i) \exp \left(\sum_n q(z_n = i) \log p(\mathbf{x}_n | \eta_i) \right)$$

$$q(z_n) \propto \exp \left(\sum_{\mathbf{z}_{-n}} \prod_{m \neq n} q(z_m) \log p(z_n | \mathbf{z}_{-n}) \right) \times \exp \left(\int_{d\eta_{z_n}} q(\eta_{z_n}) \log p(\mathbf{x}_n | \eta_{z_n}) \right)$$

for TSB formulation:

$$p(z_n = i | \mathbf{z}_{-n}) = \frac{1 + N_i^{-n}}{1 + \alpha + N_{\geq i}^{-n}} \prod_{j < k} \frac{\alpha + N_{\geq i}^{-n}}{1 + \alpha + N_{\geq i}^{-n}}$$

for FSD formulation:

$$p(z_n = k | \mathbf{z}_{-n}) = \frac{N_k^{-n} + \frac{\alpha}{K}}{N^{-n} + \alpha}$$


Optimal labels re-ordering

$$v_i \sim \mathcal{B}(v_i; 1, \alpha) \quad i = 1, \dots, T - 1 \quad (1)$$

$$v_T = 1 \quad (2)$$

$$\pi_i = v_i \prod_{j < i} (1 - v_j) \quad i = 1, \dots, T \quad (3)$$

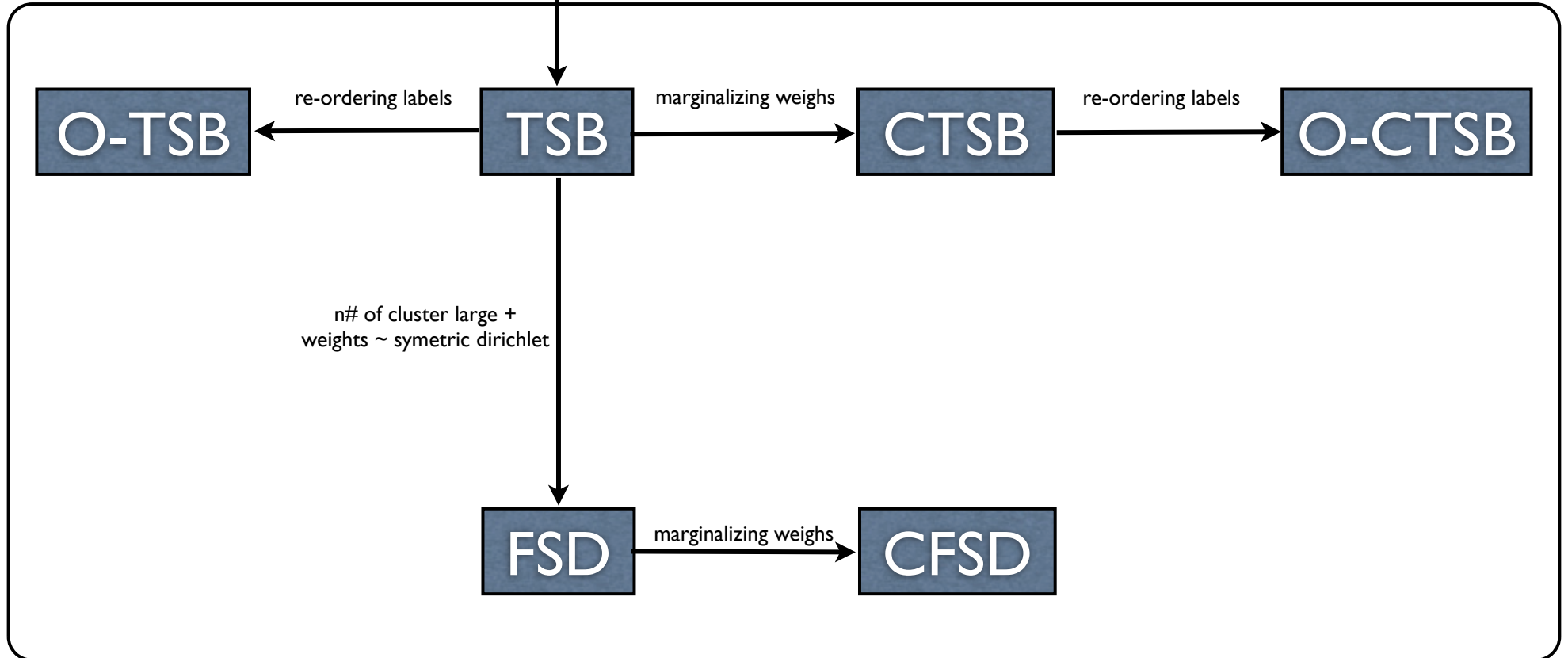
$$\pi_i = 0 \quad i > T \quad (4)$$

$$P(X, \mathbf{z}, \mathbf{v}, \boldsymbol{\eta}) = \left[\prod_{n=1}^N p(\mathbf{x}_n | \eta_{z_n}) p(z_n | \boldsymbol{\pi}(\mathbf{v})) \right] \left[\prod_{i=1}^T p(\eta_i) \mathcal{B}(v_i; 1, \alpha) \right] \quad (5)$$


- Permutation of cluster labels change the probability, therefore, an optimal reordering of the labels will maximize that probability

Big picture

Stick breaking



Experiments

Exp I:

- Synthetic data from a mixture of 10 Gaussians in 16 dimensions with a separation coefficient $c = 2$
- 30 independently sampled training/testing data, 1000 test datapoints

Exp II:

- MNIST dataset 28×28 images reduced to 50 dimensions with a PCA.
- 30 splits of the data, 5000 training and 10,000 testing.

Exp I

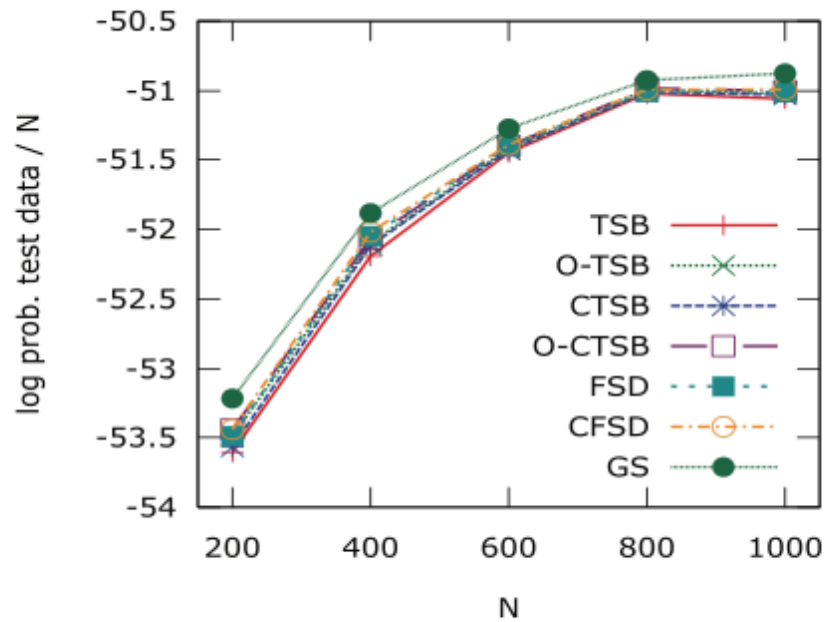


Figure 2: Average log probability per data-point for test data as a function of N .

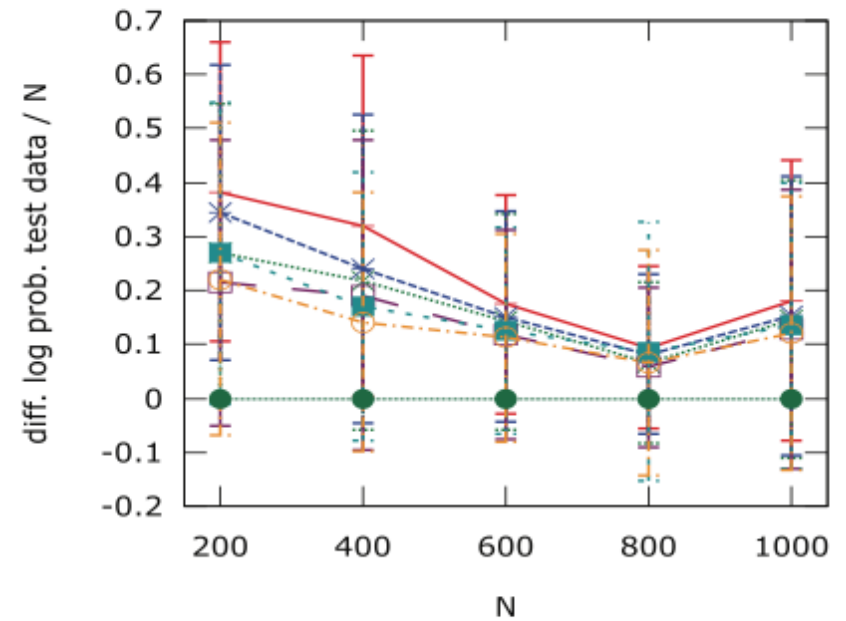


Figure 3: Relative average log probability per data-point for test data as a function of N .

Exp I

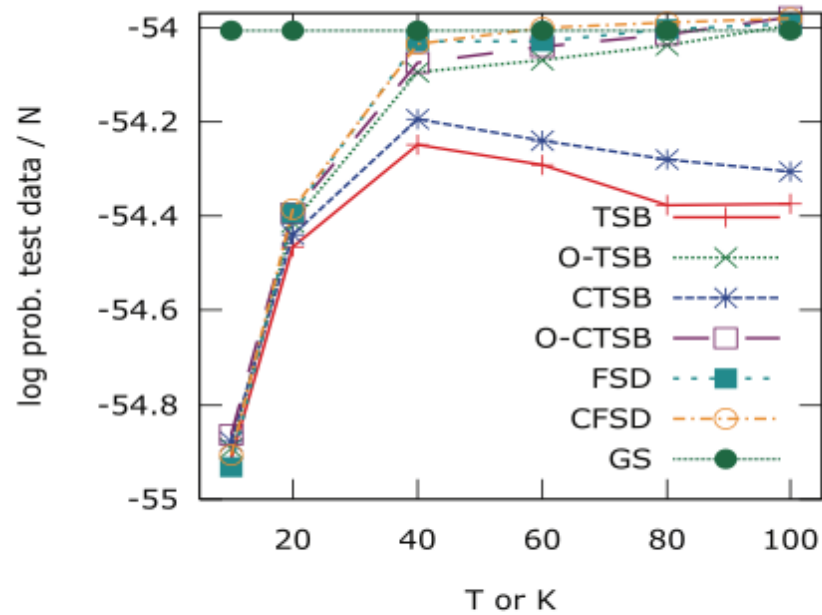


Figure 4: Average log probability per data-point for test data as a function of T (for TSB methods) or K (for FSD methods).

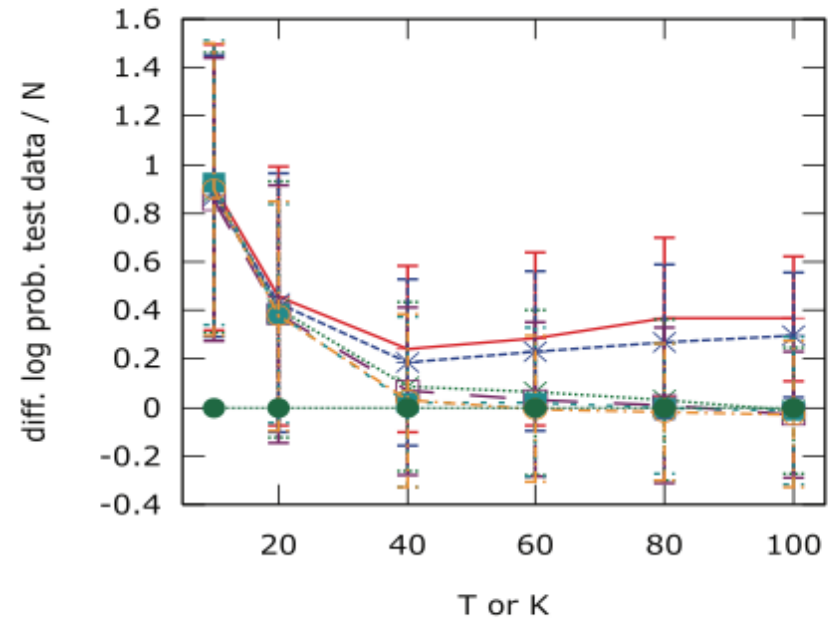


Figure 5: Relative average log probability per data-point for test data as a function of T (for TSB methods) or K (for FSD methods).

Exp II

TSB	6	3	9	1	9	5	1	8	1	0	0	0	2	2	4
O-TSB	6	3	9	1	9	5	1	8	1	0	0	0	2	2	4
CTSB	6	9	3	1	9	1	5	8	1	0	0	2	0	2	4
O-CTSB	6	3	9	1	9	1	5	8	1	0	0	0	2	2	4
FSD	6	3	9	1	9	1	5	8	1	0	0	0	2	2	4
CFSD	6	3	9	1	9	1	8	5	0	1	0	0	2	2	4

Conclusion

- There is little difference between TSB and FSD.
- Label re-ordering is important for the stick breaking representation (especially when we have no clue about how many clusters we may have).
- Variational bayesian algorithms are much more efficient computationally than Gibbs sampling, with almost no loss in accuracy.