

# Effective Split-Merge Monte Carlo Methods for Nonparametric Models of Sequential Data

Michael C. Hughes<sup>+</sup>, Emily B. Fox<sup>w</sup>, and Erik B. Sudderth<sup>+</sup>

## Model

Beta Process HMM [Fox et al. NIPS '09]  
Generates **collections** of sequential data

Global set of behaviors (features) # of behaviors learned from data

$\theta$  emission parameters

Each sequence uses a sparse subset of behaviors

$F$

$\theta$

$f_1, f_2, f_3, f_4$

$\theta_1, \theta_2, \theta_3, \theta_4, \theta_5$

$z_1, z_2, z_3, z_4$

$\eta_i$  seq.  $i$ 's HMM transition weights

$f_i$  seq.  $i$ 's binary feature assignment

$z_{it}$  active feature at time  $t$  in seq.  $i$

$x_{it}$  observed data at time  $t$  in seq.  $i$

$x_{it} \sim p(\cdot | \theta_{z_{it}})$

**BETA PROCESS**

$\theta_k$

time series  $i = 1 \dots N$

$k = 1, 2, \dots$

$X_{i1}, X_{i2}, X_{iT}$

## Data-Driven Birth/Death

Add or delete unique feature to **one** sequence

Reversible jump proposal for  $F, \theta$  ( $z$  marginalized out)

$F$   $\xrightarrow{\text{death}} F^*$   $\xleftarrow{\text{birth}} \theta^*$

**Propose from prior** [Fox et al. NIPS 2009]

$\theta_{k^*}^* \sim p(\theta)$

**Data-driven proposal**

- select random window  $W$  of sequence
- proposal: mixture of prior and posterior over  $W$

$\theta_{k^*}^* \sim \frac{1}{2}p(\theta) + \frac{1}{2}p(\theta|x_{it} : t \in W)$

Using mixture ensures good death move acceptance rate

Efficiently adds new behaviors informed by the data

## Split-Merge

Yields exact posterior samples with big changes to feature assignments via reversible Metropolis-Hastings proposal moves.

**Merge**

Select anchors  $i, j \sim \text{Unif}(\text{sequences})$

$f_n$  must own  $k_m$

$z_n$  block sampled given  $f_n$

**split**

Features available in  $f_n$  can appear anywhere in  $z_n$

**Must satisfy**

$f_{ik_i} = 1$

$f_{jk_j} = 1$

**NO**

**YES**

**SPLIT proposal**

$F^*, z^* \sim q_{\text{split}}(\cdot | k_i)$

HMM params  $\theta, \eta$  collapsed away

**MERGE proposal**

$F^*, z^* \sim q_{\text{merge}}(\cdot | k_i, k_j)$

**Joint probability**

$p(x, z^*, F^*) / p(x, z, F)$

**Proposal construction**

$q_{\text{merge}}(F, z | x, F^*, z^*, k_a, k_b) / q_{\text{split}}(F^*, z^* | x, F, z, k_m)$

$q_k(k_a, k_b | x, F^*, z^*, i, j) / q_k(k_m, k_m | x, F, z, i, j)$

**Feature selection**

**Both moves do not change sequences not in active set  $S$**

**3) Allocate anchors**

To ensure reversibility via a merge,  $i$  must own  $k_a$   $j$  must own  $k_b$

$f_i \sim \{ \text{black, blue} \}$   $f_j \sim \{ \text{black, red} \}$

This does not require that  $k_a$  appears in  $z_i$  or  $k_b$  appears in  $z_j$

## Problem

Existing MCMC inference slow to mix

- Requires long time to make big changes
  - Each update touches small subset of variables
- Rarely creates new features
  - Proposals from vague prior poorly-matched to data

## Contributions

- Split-Merge (SM) move**
  - Change many feature assignments at once
- Data-Driven (DD) birth/death move**
  - Propose new features consistent with observed data

Significantly improves mixing.  
Scales to 100+ sequences.

## Toy Experiments

Example toy sequences 5 dim. Gaussian obs. (1<sup>st</sup> & 2<sup>nd</sup> shown)

Starting at poor initialization of one feature,  $4 \times 10^4$  can MCMC recover all 8 true features?

log joint prob. vs. cpu time (sec)

Discovered emission parameters  $\theta$

best run prior birth

worst run DD birth

worst run SM

SM

DD

Prior

SM and DD moves find essential features  
Prior methods stuck in bad local optima

## Motion Capture

6 subjects performing 12 exercises (7 shared, 5 unique).

**Data:** 12 joint angle sensors. **Model:** 1<sup>st</sup>-order auto-regressive.

**Goal:** Compare BP-HMM segmentation to human annotation.

A: SM+DD, init=1 feature  
B: SM+DD, init=5 unique / seq.  
C: Prior RJ, init=5 unique / seq.

Subj. 1, Subj. 5, Subj. 6

Hamming dist. vs. cpu time (sec)

Segments of true jogging fragments. SM+DD finds dominant single behavior. Prior unable to merge due to local moves.

Previous work stuck with split behaviors even with clever initialization. New MCMC finds better segmentation starting with just one feature.

Example discovered behaviors on larger 124 sequence dataset

Ballet, Walk, Squat, Sword, Lambada, Dribble Basketball, Box, Climb, Bollywood Dance, Tai Chi

## Kitchen Video

126 subjects prepare 5 recipes in common kitchen. **Model:** Multinomial.

**Data:** Bag of visual words (motion + appearance), at 1 second intervals.

Usage Over Time for 8 Select Behaviors

Eggs, Salad, Pizza, Brownie

Flip Omelette, Slice/Chop, Grate Cheese, Pour Bowl, Stir Bowl 2, Stir Bowl 1, Open Fridge, Light Switch, 65 Others

Each row shows behaviors used throughout one video at final MCMC iteration.

Unsupervised. Recipe labels not given as input.

Not shown: Sandwich videos.

Example Key Frames:

Open Fridge, Grate Cheese, Light Switch, Stir Bowl 1