

# High to High Dimensional Multivariate Mixture Regression

Alex White

4/16/2021

# Idea

Goal: Correctly cluster observations & regress in high dimensional  $X$  &  $Y$ .

- ▶  $Y_{n \times q}$  Matrix of responses
- ▶  $X_{n \times p}$  Design matrix
- ▶  $A_{p \times q}$  Coefficient matrix (sparse in  $p$ )
- ▶  $k$  clusters

$$f(\mathbf{y}_i \mid \mathbf{x}_i; \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}_q(\mathbf{y}_i; \mathbf{x}_i A_k, \Sigma_k)$$

- ▶ Parameter space  $\boldsymbol{\theta} = \{\pi_k, A_k, \Sigma_k; k = 1 \dots K\}$  solved by general EM using SARRS to compute  $A_k$ .

# SARRS

- ▶ Excellent performance, even in high dimension ( $> 5000$ )
- ▶ Computationally very fast

---

**Algorithm 1:** Subspace Assisted Regression with Row Sparsity (SARRS)

---

**Input:** Observed response matrix  $Y$ , design matrix  $X$ , rank  $r$ , initial matrix  $V_{(0)}$  and penalty function  $\rho(\cdot; \lambda)$  with penalty level  $\lambda$ .

**Output:** Estimated coefficient matrix  $\hat{A}$ .

1 Group penalized regression

$$B_{(1)} = \arg \min_{B \in \mathbb{R}^{p \times r}} \left\{ \|YV_{(0)} - XB\|_F^2 / 2 + \rho(B; \lambda) \right\},$$

2 Compute the left singular vectors of  $XB_{(1)}$ , denoted by  $U_{(1)}$ .

3 Compute the right singular vectors of  $U_{(1)}U_{(1)}'Y$ , denoted by  $V_{(1)}$ .

4 Group penalized regression

$$B_{(2)} = \arg \min_{B \in \mathbb{R}^{p \times r}} \left\{ \|YV_{(1)} - XB\|_F^2 / 2 + \rho(B; \lambda) \right\},$$

5 Compute the estimated coefficient matrix by  $\hat{A} = B_{(2)}V_{(1)}'$ .

---

Figure 1: “SARRS Main Algorithm”

# HTH Mixture Algorithm

- ▶ Initialize:  $\pi_k^{(0)} = \frac{n_k^{(0)}}{n}$
- ▶ Randomly initialize observations into  $k$  clusters

While not converged ( $m = 1, \dots, M$ ) do:

- ▶ for  $k = 1, \dots, K$  apply SARRS on all observations in  $C_k^{(m-1)}$  to obtain  $A_k^{(m)}, \Sigma_k^{(m)}$
- ▶ compute  $\mu_{ik}^{(m)} = \mathcal{N}_p(\mathbf{y}_i; A_k^{(m)} \mathbf{x}_i, \Sigma_k^{(m)})$
- ▶  $C_k^{(m)} = \{i | \text{ML component } k\}$
- ▶  $\pi_k^{(m)} = \frac{n_k^{(m)}}{n}$

# HTH Mixture Algorithm

- ▶ Empirically, HTH Mixture reaches local maximum quickly
- ▶ Need to determine method for efficiently exploring the likelihood space with different random initialization.

# Likelihood Space Exploration Idea 1

- ▶ Mimic the idea of chains from MCMC
- ▶ In parallel, run large number of chains and select the one which reaches the largest likelihood.
- ▶ This method works well, although not the most efficient way to explore the space.

## Likelihood Space Exploration Idea 2

- ▶ To more efficiently explore the space, when a local maximum is reached, perturb the initialized state by a proportion  $p$  (e.g. 20%) and rerun.

# Data Simulation

- ▶  $X_k$  consists of iid random vectors sample from  $MVN(\mathbf{0}, \Sigma_k)$
- ▶  $\Sigma_k$  diagonal
- ▶ Noise matrix  $Z_k \in \mathbb{R}^{n \times q}$  has iid  $N(0, \sigma^2)$  entries
- ▶  $A_k = \begin{pmatrix} b_k B_{0_k} B_{1_k} \\ 0 \end{pmatrix}$ 
  - ▶ with  $b > 0$ ,  $B_0 \in \mathbb{R}^{s \times r}$ ,  $B_1 \in \mathbb{R}^{r \times q}$
- ▶  $Y_k = X_k A_k + Z_k$

Finally, combine  $X$  &  $Y$

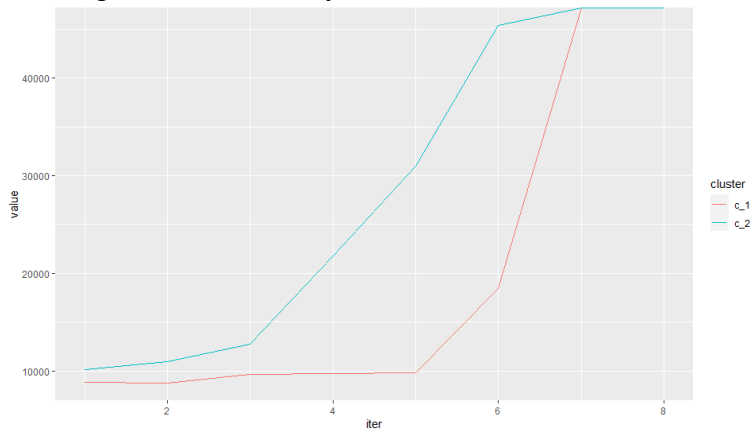


# Simulation

- ▶ N: 50 – 400
- ▶ s: 5 – 20
- ▶ rank: 1 – 2
- ▶ B: 1
- ▶  $p = m = 1000$
- ▶ Perturb: 20%
- ▶ Cluster with equal probability
- ▶ K: 2
- ▶  $\sigma$ : 1
- ▶ CHains: 500

# Results

- HTH algorithm consistently reaches local maximum



# Results

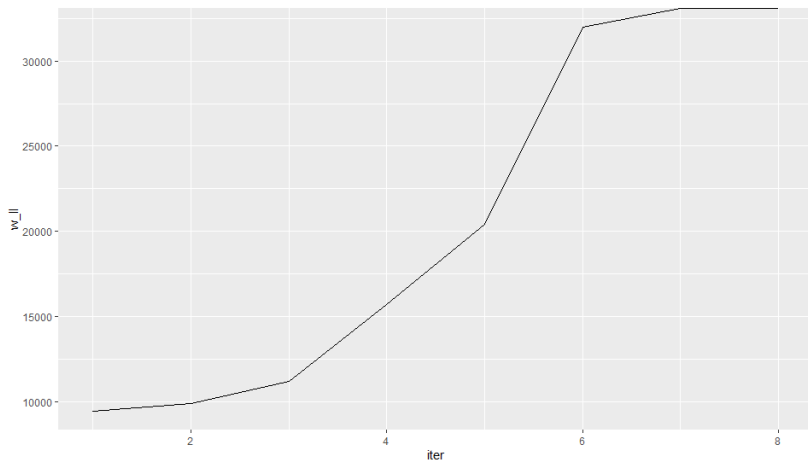
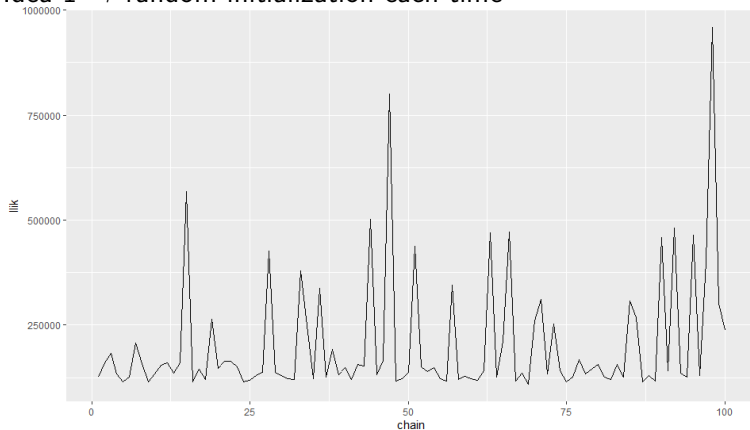


Figure 2: Weighted Likelihood

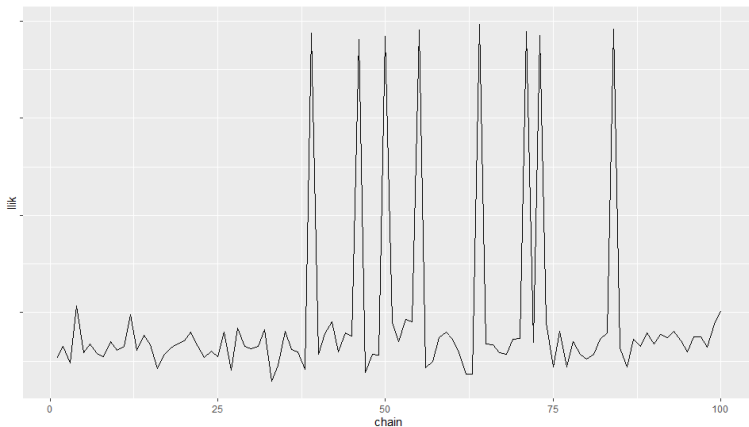
# Results

- Idea 1 → random initialization each time



# Results

- ▶ Idea 2  $\rightarrow$  perturb local maximum
- ▶ Appears to be more efficient



## Results

N	s	R	hthmix	kmeans
50	5	1	0.90	0.70
100	5	1	1.00	0.66
200	5	1	1.00	0.51
400	5	1	1.00	0.54
50	10	1	0.84	0.62
100	10	1	0.92	0.50
200	10	1	1.00	0.66
400	10	1	1.00	0.52
50	20	1	0.80	0.72
100	20	1	0.80	0.53
200	20	1	0.98	0.68
400	20	1	1.00	0.51

## Results

N	s	R	hthmix	kmeans
50	5	2	1.00	0.64
100	5	2	1.00	0.72
200	5	2	1.00	0.52
400	5	2	1.00	0.64
50	10	2	0.86	0.58
100	10	2	1.00	0.55
200	10	2	1.00	0.53
400	10	2	1.00	0.68
50	20	2	0.84	0.64
100	20	2	0.78	0.53
200	20	2	1.00	0.71
400	20	2	1.00	0.65

# Performance

- ▶ In simulated data, current algorithm clusters well (perfectly in most cases given enough chains)
- ▶ Challenges:
  - ▶ Likelihood space is explored inefficiently, sometimes requires large number of chains
  - ▶ Computationally cumbersome/inefficient