

TSSB-DW: Two-level clustering analysis

Yinqiao Yan

Last updated on 22/05/2021

Contents

1	Introduction	1
2	Preliminaries	2
2.1	DP and stick breaking process	2
2.2	Tree-structured stick-breaking process	4
3	Our model	5
3.1	Posterior inference	5
4	Simulation study	7
4.1	Prepare data	7
4.2	Initialization and training	8
4.3	Results	9
5	MNIST dataset	12
5.1	Prepare the MNIST-mini data	12
5.2	Initialization and training	13
5.3	Results	14
6	Conclusion	15
7	Future work	15

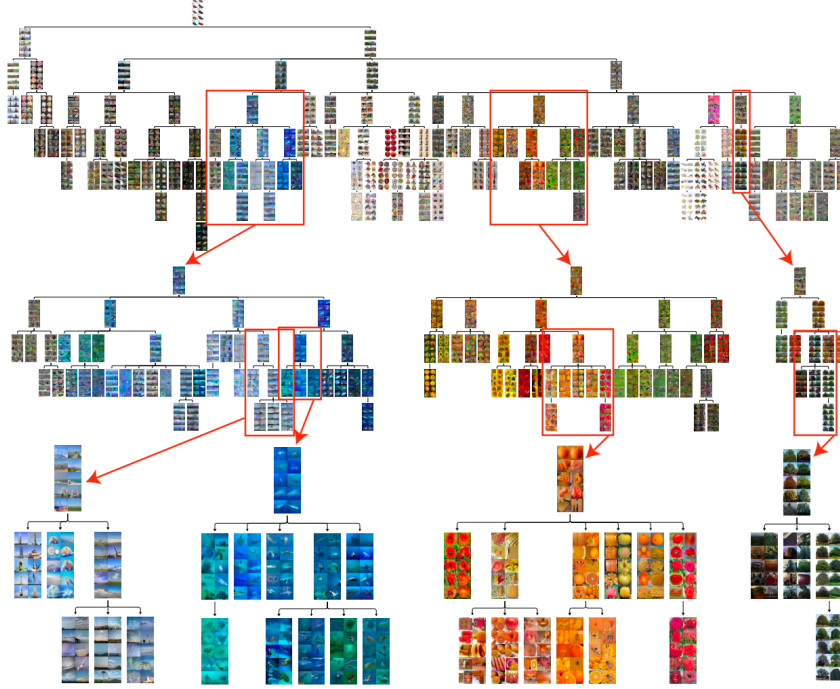
1 Introduction

Clustering is one of the most important unsupervised learning problem. For example, in single-cell analysis, cells are grouped into several cell types based on the cell profile data. This is really useful when we hope to analyze the potential various cell types existed in the original tissue, or explore some new types. In image clustering problems, images are clustered based on some features (extracted from PCA or deep neural networks) to find some underlying relationships or patterns behind the images in the same cluster.

Recently, tree-structured clustering methods are getting more attention in *relationship reconstruction problem among subclones*. We need to simultaneously cluster the data into several groups and reconstruct the underlying relationship among these groups. Adams et al. (2010) proposed a novel nonparametric Bayesian prior named tree-structured stick-breaking prior (TSSB). They extended the stick-breaking process of DP to a two-dimension case. This method has been implemented in tumor subclone phylogeny reconstruction problem (Deshwar et al., 2015; Yuan et al., 2015). TSSB has infinite width and depth, and each datum can locate in any internal node of the tree. In (Adams et al.2010), TSSB was applied to image hierarchical clustering problem (CIFAR-100). They firstly extracted the 256 binary features $x \in \{0, 1\}^{256}$ from a DNN, and then used the factored Bernoulli likelihood at each node. Note that the factored likelihood is based on

the assumption that the elements of the data are independent. This setting can avoid the high dimensionality curse.

$$f(x_n | \theta_\epsilon) = \prod_{d=1}^{256} \left(1 + \exp\{-\theta_\epsilon^{(d)}\}\right)^{-x_n^{(d)}} \left(1 + \exp\{\theta_\epsilon^{(d)}\}\right)^{1-x_n^{(d)}}$$



2 Preliminaries

Adams et al. (2010) proposed a nonparametric Bayesian prior to model the underlying tree structure behind the different clusters, which is referred to as tree-structured stick-breaking process (TSSB). Recently, TSSB has been applied in some hierarchical clustering problem, such as image clustering (Adams et al. 2010) and tumor phylogeny reconstruction (Yuan et al. 2015, Deshwar et al. 2015). This prior has infinite depth and width and each datum can live at any internal node of the tree, not only the leaf nodes (different from the hierarchical clustering method).

2.1 DP and stick breaking process

Brief definition of Dirichlet process (DP).

- A DP is a **random probability measure** G defined on the probability space S (prior on the space of probability measure, $\mathcal{M}(S)$). Denoted by $G \sim DP(\alpha, H)$.
- For measurable **finite** partition $\{B_1, \dots, B_k\}$ of S , the joint distribution of the vector $(G(B_1), \dots, G(B_k))$ is the **Dirichlet distribution** with parameters $(\alpha H(B_1), \dots, \alpha H(B_k))$.

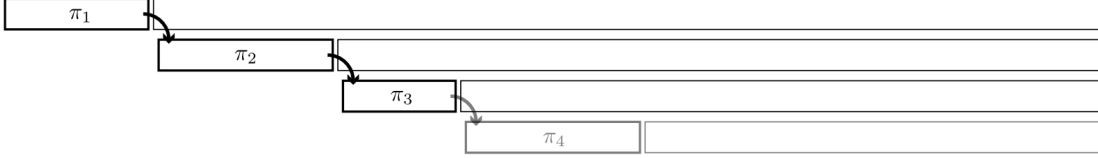
Stick breaking construction of DP (equivalent definition).

DP is almost surely discrete probability measure (natural property for clustering). It can be written as (Sethuraman 1994)

$$G(\cdot) = \sum_{k=1}^{\infty} \pi_k \delta_{Z_k}(\cdot)$$

- $\delta_{Z_k}(\cdot)$ denote a discrete measure concentrated at $Z_k \stackrel{\text{iid}}{\sim} H$.
- π_k are random weights, chosen to be independent of Z_k and satisfy $\sum_{k=1}^N \pi_k = 1$.
- $\pi_1 = V_1$ and $\pi_k = (1 - V_1)(1 - V_2) \cdots (1 - V_{k-1}) V_k$, $k \geq 2$.
- $V_k \stackrel{\text{iid}}{\sim} \text{Be}(1, \alpha)$.

We can also denote $(\pi_1, \pi_2, \dots) \sim \text{GEM}(\alpha)$, which is named after Griffiths, Engen, and McCloskey.



(a) Dirichlet process stick breaking

Truncation of DP.

Ishwaran and James (2001) designed a block Gibbs sampler based on truncation version of stick-breaking process, denoted by

$$\sum_{k=1}^N \pi_k \delta_{Z_k}(\cdot), \quad N < \infty$$

where

- $\pi_1 = V_1$, $\pi_k = V_k \prod_{j < k} (1 - V_j)$, $k = 2, \dots, N$.
- Set $V_N = 1$ to guarantee $\sum_{k=1}^N \pi_k = 1$ with probability 1.

Motivation of truncation.

1. Limitations of Polya urn Gibbs scheme (marginalizing) for infinite measure:
 - will lead to slowly mixing Markov chain (slowly converges to stationary distribution)
 - one-at-a-time updates
 - hard to solve non-conjugate cases
 - difficult to find a unified sampling framework (*Griffin (2016) pointed out the unavailability of a suitable Polya urn scheme for some priors.*)
2. Block Gibbs sampler:
 - simpler and more efficient
 - unified Gibbs framework
 - parallel computing (update a block of parameters each time)

Non-parametric mixture model. In (Ishwaran and James, 2001), the mixture model is given by

$$\begin{aligned} (X_i | Y_i) &\stackrel{\text{ind}}{\sim} \pi(X_i | Y_i), \quad i = 1, \dots, n, \\ (Y_i | P) &\stackrel{\text{iid}}{\sim} P \\ P &\sim \mathcal{P}_N. \end{aligned}$$

~~To implement a Polya urn scheme sampler, integrate out the nonparametric prior. (Not recommended)~~

To implement a block Gibbs sampler, introduce a index variable K . (*Data augmentation*)

$$\begin{aligned} (X_i | \mathbf{Z}, \mathbf{K}) &\stackrel{\text{iid}}{\sim} \pi(X_i | Z_{K_i}), \quad i = 1, \dots, n \\ (K_i | \mathbf{p}) &\stackrel{\text{iid}}{\sim} \sum_{k=1}^N p_k \delta_k(\cdot), \quad N < \infty \\ (\mathbf{p}, \mathbf{Z}) &\sim \pi(\mathbf{p}) \times H^N(\mathbf{Z}), \end{aligned}$$

where

- $Y_i = Z_{K_i}$, $Z_k \stackrel{\text{iid}}{\sim} H$, $k = 1, \dots, N$.
- $\mathbf{p} = (p_1, \dots, p_N) \sim \text{GEM}_N(\alpha)$ (belongs to generalized dirichlet distribution - *won't go into details here*)
- K_i are i.i.d and can be updated independently (parallel computing).

Block Gibbs sampler.

Update parameters **in blocks**.

$$\begin{aligned} &(\mathbf{Z} | \mathbf{K}, \mathbf{X}), \\ &(\mathbf{K} | \mathbf{Z}, \mathbf{p}, \mathbf{X}), \\ &(\mathbf{p} | \mathbf{K}), \end{aligned}$$

2.2 Tree-structured stick-breaking process

Extended to two breaking processes.

Adams et al. (2010) proposed this two-dimensional stick-breaking process, referred to as tree-structured stick-breaking process (TSSB).

- ν -breaks: decides the weight for staying at a node.

$$v_\epsilon \sim \text{Be}(1, \alpha(|\epsilon|))$$

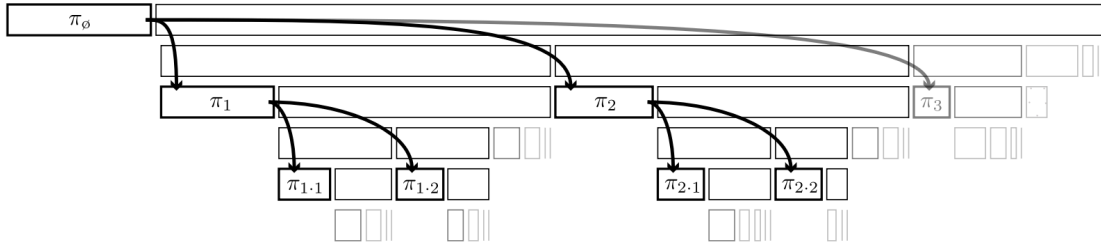
(Adams et al. suggested $\alpha(j) = \lambda^j \alpha_0$.)

- ψ -breaks: decides the weight for selecting a child node to descend.

$$\psi_\epsilon \sim \text{Be}(1, \gamma)$$

The random weights are then given by

$$\pi_\epsilon = \nu_\epsilon \varphi_\epsilon \prod_{\epsilon' \prec \epsilon} \varphi_{\epsilon'} (1 - \nu_{\epsilon'}) \quad \varphi_{\epsilon \epsilon_i} = \psi_{\epsilon \epsilon_i} \prod_{j=1}^{\epsilon_i-1} (1 - \psi_{\epsilon_j}) \quad \pi_\emptyset = \nu_\emptyset,$$



(b) Tree-structured stick breaking

Hierarchical priors for node parameters.

- In DP, Z_k are i.i.d. from the base measure H .

- In TSSB, the i.i.d. assumption on the node parameters is inappropriate, since the random weights have a hierarchical structure. Hence we need a **transition distribution** for the node parameters.

$$T(\theta_\epsilon \leftarrow \theta_{pa(\epsilon)})$$

where $\theta_{pa(\epsilon)}$ denote the parent node of θ_ϵ .

3 Our model

Motivated by the work of Ishwaran and James (2001), in this report we propose a truncation version of TSSB, referred to as TSSB-DW (**TSSB** with finite **Depth** and **Width**). Factored normal likelihood is used to avoid the high-dimensionality problem, and the prior of node parameters θ_ϵ and σ_ϵ^2 are the most commonly used conjugate prior Normal-invGamma distribution. The hyper-parameter λ (sometimes called the drift parameter in $T(\theta_\epsilon \leftarrow \theta_{pa(\epsilon)})$) has inverse Gamma distribution. The hyper-parameter $\eta_{\mathcal{N}}$ and η_{Θ} are fixed.

$$\begin{aligned} (X_i \mid \theta, \Sigma, c_i = \epsilon) &\stackrel{\text{iid}}{\sim} \prod_{\ell=1}^L N(X_i^\ell \mid \theta_\epsilon^\ell, \eta_{\mathcal{N}}^{|\epsilon|} \sigma_\epsilon^{2\ell}), \quad i = 1, \dots, n \\ c_i \mid \pi &\stackrel{\text{iid}}{\sim} \sum_{\epsilon} \pi_\epsilon \delta_\epsilon \\ \pi &\sim \text{TSSB-DW}(\alpha_0, \rho, \gamma) \\ \theta_\emptyset^\ell &\sim N(\theta_\emptyset^\ell \mid \mu_0^\ell, \lambda^\ell), \quad \ell = 1, \dots, L \\ \theta_\epsilon^\ell \mid \theta_{pa(\epsilon)}^\ell, \eta_\Theta &\stackrel{\text{iid}}{\sim} N(\theta_\epsilon^\ell \mid \theta_{pa(\epsilon)}^\ell, \eta_\Theta^{|\epsilon|} \lambda^\ell) \\ \sigma_\epsilon^{2\ell} &\stackrel{\text{iid}}{\sim} \text{InvGamma}(v_{sig}, s_{sig}) \\ \lambda^\ell &\stackrel{\text{iid}}{\sim} \text{InvGamma}(v_{dft}, s_{dft}) \end{aligned}$$

Another option of the prior of drift is uniform distribution.

$$\lambda^\ell \stackrel{\text{iid}}{\sim} \text{Unif}(\min, \max)$$

3.1 Posterior inference

The posterior inference includes the following steps.

Conditional for θ_ϵ^ℓ , $\ell = 1, \dots, L$

$$\begin{aligned} p(\theta_\epsilon^\ell \mid -) &\propto N(\theta_\epsilon^\ell \mid \theta_{pa(\epsilon)}^\ell, \eta_\Theta^{|\epsilon|} \lambda^\ell) \prod_{ch(\epsilon)} N(\theta_{ch(\epsilon)}^\ell \mid \theta_\epsilon^\ell, \eta_\Theta^{|\epsilon|+1} \lambda^\ell) \prod_{\{i: c_i = \epsilon\}} N(X_i^\ell \mid \theta_\epsilon^\ell, \eta_{\mathcal{N}}^{|\epsilon|} \sigma_\epsilon^{2\ell}) \\ &\propto N(\mu_\epsilon^\ell, \tau_\epsilon^\ell) \end{aligned}$$

where

$$\begin{aligned} \mu_\epsilon^\ell &= \left(\frac{N_\epsilon \bar{X}_\epsilon^\ell}{\eta_{\mathcal{N}}^{|\epsilon|} \sigma_\epsilon^{2\ell}} + \frac{W \bar{\theta}_{ch(\epsilon)}^{\ell(t-1)} + \eta_\Theta \theta_{pa(\epsilon)}^{\ell(t)}}{\eta_\Theta^{|\epsilon|+1} \lambda^\ell} \right) \left(\frac{N_\epsilon}{\eta_{\mathcal{N}}^{|\epsilon|} \sigma_\epsilon^{2\ell}} + \frac{W + \eta_\Theta}{\eta_\Theta^{|\epsilon|+1} \lambda^\ell} \right)^{-1}, \\ \tau_\epsilon^\ell &= \left(\frac{N_\epsilon}{\eta_{\mathcal{N}}^{|\epsilon|} \sigma_\epsilon^{2\ell}} + \frac{W + \eta_\Theta}{\eta_\Theta^{|\epsilon|+1} \lambda^\ell} \right)^{-1}. \end{aligned}$$

and N_ϵ is the number of data stopping at node ϵ . If $N_\epsilon = 0$, then

$$\mu_\epsilon^\ell = \frac{W \bar{\theta}_{ch(\epsilon)}^{\ell(t-1)} + \eta_\Theta \theta_{pa(\epsilon)}^{\ell(t)}}{W + \eta_\Theta}, \quad \tau_\epsilon^\ell = \frac{\eta_\Theta^{|\epsilon|+1} \lambda^\ell}{W + \eta_\Theta}$$

Conditional for $\sigma_\varepsilon^{2\ell}$, $\ell = 1, \dots, L$

$$p(\sigma_\varepsilon^{2\ell} | -) \propto \prod_{\{i: c_i = \varepsilon\}} N(X_i^\ell | \theta_\varepsilon^\ell, \eta_{\mathcal{N}}^{|\varepsilon|} \sigma_\varepsilon^{2\ell}) \text{InvGamma}(v_{sig}, s_{sig}) \\ \propto \text{InvGamma}(\tilde{v}_{sig}^\ell, \tilde{s}_{sig}^\ell)$$

where

$$\tilde{v}_{sig}^\ell = v_{sig} + N_\varepsilon/2 \\ \tilde{s}_{sig}^\ell = s_{sig} + \frac{1}{2\eta_{\mathcal{N}}^{|\varepsilon|}} \sum_{\{i: c_i = \varepsilon\}} (X_i^\ell - \theta_\varepsilon^\ell)^2$$

Note that if there is no datum at the current node (i.e., $N_\varepsilon = 0$), then v, s do not change. The density of inverse Gamma distribution $\text{InvGamma}(v, s)$ is

$$p(x) = \frac{s^v}{\Gamma(v)} x^{-(v+1)} \exp\left(-\frac{s}{x}\right)$$

where $v > 0$ and $s > 0$ are called shape and scale parameter respectively.

Conditional for c

$$p(c_i | \theta, \Sigma, \pi, X) \stackrel{\text{ind}}{\sim} \sum_{\varepsilon} \beta_\varepsilon^{(i)} \delta_\varepsilon$$

where

$$\beta_\varepsilon^{(i)} \propto \pi_\varepsilon \prod_{\ell=1}^L N(X_i^\ell | \theta_\varepsilon^\ell, \eta_{\mathcal{N}}^{|\varepsilon|} \sigma_\varepsilon^{2\ell})$$

Conditional for π

The key is to update ν and ψ , and π_ε are given by

$$\pi_\varepsilon = \nu_\varepsilon^* \psi_\varepsilon^* \prod_{\varepsilon' \prec \varepsilon} \psi_{\varepsilon'}^* (1 - \nu_{\varepsilon'}^*)$$

where

$$\nu_\varepsilon^* \stackrel{\text{ind}}{\sim} \text{Be}(1 + N_\varepsilon, \alpha(|\varepsilon|) + N_{\varepsilon \prec \cdot}) \\ \psi_{\varepsilon e_i}^* \stackrel{\text{ind}}{\sim} \text{Be}(1 + N_{\varepsilon e_i \prec \cdot}, \gamma + \sum_{j > e_i} N_{\varepsilon j \prec \cdot})$$

$N_\varepsilon = \#\{i : c_i = \varepsilon\}$ is the number of data stopping at the node ε , $N_{\varepsilon \prec \cdot}$ is “the number of data that come down this path but does not stop at ε ”, and $N_{\varepsilon \preccurlyeq \cdot}$ is equal to $N_\varepsilon + N_{\varepsilon \prec \cdot}$.

Conditional for the hyperparameter drift λ^ℓ , $\ell = 1, \dots, L$

$$p(\lambda^\ell | -) \propto \prod_{\varepsilon} N(\theta_\varepsilon^\ell | \theta_{pa(\varepsilon)}^\ell, \eta_{\Theta}^{|\varepsilon|} \lambda^\ell) \text{InvGamma}(v_{dft}, s_{dft}) \\ \propto \text{InvGamma}(\tilde{v}_{dft}^\ell, \tilde{s}_{dft}^\ell)$$

where

$$\tilde{v}_{dft}^\ell = v_{dft} + N_{nodes}/2, \quad N_{nodes} = 1 + \dots + W^D = \frac{W^{D+1} - 1}{W - 1} \\ \tilde{s}_{dft}^\ell = s_{dft} + \sum_{\varepsilon} \frac{1}{2\eta_{\Theta}^{|\varepsilon|}} (\theta_\varepsilon^\ell - \theta_{pa(\varepsilon)}^\ell)^2$$

Note that when $\varepsilon = \emptyset$, then $\theta_{pa(\varepsilon)}$ is equal to the initial mean μ_0 . (since the prior of θ_0 also includes the drift)

Conditional for the hyperparameters α_0, ρ, γ

In (Adams et al.2010), the authors used the slice sampler to update these hyperparameters. This is slice sampler-within-Gibbs framework.

$$p(\alpha_0, \lambda \mid \{\nu_\epsilon\}) \propto \mathbb{I}(\alpha_0^{\min} < \alpha_0 < \alpha_0^{\max}) \mathbb{I}(\lambda^{\min} < \lambda < \lambda^{\max}) \prod_{\epsilon} \text{Be}(\nu_\epsilon \mid 1, \lambda^{|\epsilon|} \alpha_0)$$

$$p(\gamma \mid \{\psi_\epsilon\}) \propto \mathbb{I}(\gamma^{\min} < \gamma < \gamma^{\max}) \prod_{\epsilon} \text{Be}(\psi_\epsilon \mid 1, \gamma)$$

Search for tree structure

In (Yuan et al. 2015), the authors added another swap-nodes step to propose a new tree structure. They switched the data Ids, node parameters (θ_ϵ and σ_ϵ^2) and ν -breaks ν_ϵ of the two chosen nodes. The proposal would be accepted if the new unnormalized posterior is large than the old one. Here we follow the idea in (Yuan et al. 2015).

4 Simulation study

4.1 Prepare data

Data: 7 classes normal data. Separable. Dims: 700×2 dims

Import the packages we will need.

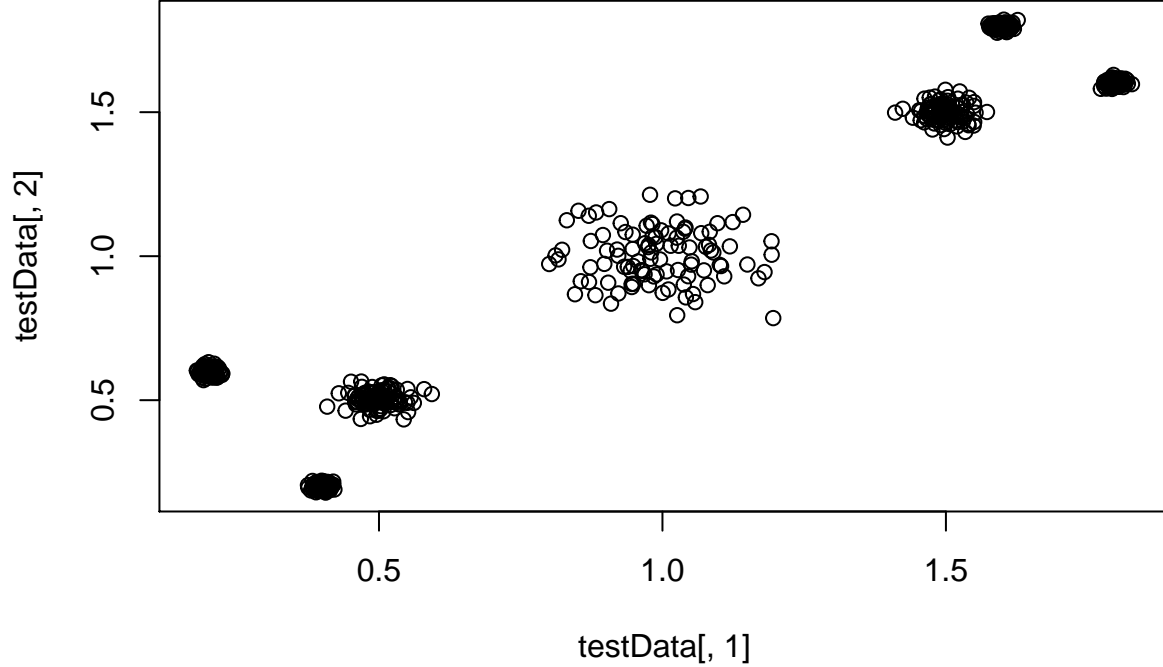
```
library(igraph)
library(ggplot2)
library(networkD3)
```

In this study, the data matrix are 700×2 . There are seven classes, each of which has 100 data.

```
set.seed(12)

m <- 700
dims <- 2
testData <- rbind(rmvnorm(m/7, mean = rep(1.0, dims), sigma = diag(0.10^2, dims, dims)),
                  rmvnorm(m/7, mean = rep(1.5, dims), sigma = diag(0.03^2, dims, dims)),
                  rmvnorm(m/7, mean = c(1.6, 1.8), sigma = diag(0.01^2, dims, dims)),
                  rmvnorm(m/7, mean = c(1.8, 1.6), sigma = diag(0.01^2, dims, dims)),

                  rmvnorm(m/7, mean = rep(0.5, dims), sigma = diag(0.03^2, dims, dims)),
                  rmvnorm(m/7, mean = c(0.4, 0.2), sigma = diag(0.01^2, dims, dims)),
                  rmvnorm(m/7, mean = c(0.2, 0.6), sigma = diag(0.01^2, dims, dims))
)
numOfData = nrow(testData)
plot(testData[,1], testData[,2])
```



4.2 Initialization and training

The settings of implementing the model are:

- $\eta_{\mathcal{N}} = 1$ and $\eta_{\Theta} = 0.5$
- Update order: i) Node parameters, ii) Data assignments, iii) Swap nodes
- burnIn = 100, Iter = 1000
- maxDepth = 3, maxWidth = 3
- “OnlyTree”

The drift λ is assigned the **inverse Gamma** prior. Then we **initialize** the Root Node and the corresponding Tssb.

```
set.seed(123)

empCov <- cov(t(testData))
priorSigmaScale = mean(diag(empCov))
priorDriftScale = priorSigmaScale

minDepth <- 0
maxDepth <- 3
maxWidth <- 3
etaNormal <- 1
etaTheta <- 0.3
Flag.onlyTree <- T

q0 <- Normal_DW_Factored_eta_ST$new(priorSigmaScale = priorSigmaScale,
  priorDriftScale = priorDriftScale,
  etaNormal = etaNormal, etaTheta = etaTheta,
  dataDims = ncol(testData))

tssbMCMC <- TssbMCMC_DW_Factored_eta_ST$new(q0, data = testData,
  dpAlpha = 1, dpGamma = 1, dpLambda = 1,
```



```

maxDepth = maxDepth, minDepth = minDepth,
maxWidth = maxWidth,
Flag.onlyTree = Flag.onlyTree)
#> Initialization: D = 3 and W = 3
#> No. 1 child of root is initializing...
#> No. 2 child of root is initializing...
#> No. 3 child of root is initializing...
#> Initialization Is Over!

```

After initialization, we update all the parameters using Gibbs sampler.

```

#> ==== iter: 0
#> Update node params: 0.015167
#> Update drift: 0.01378918
#> Update sticks: 0.02305794
#> Update 3 hypers: 0.08474994
#> Update assignments: 0.7528319
#> Swap-nodes move: 0.05140901
#> Update Keep tree: 0.01814198
#> ==== iter: 250
#> Update node params: 0.00420785
#> Update drift: 0.001023054
#> Update sticks: 0.000549078
#> Update 3 hypers: 0.001919985
#> Update assignments: 0.6108501
#> Swap-nodes move: 0.02652597
#> Update Keep tree: 0.0008840561
#> ==== iter: 500
#> Update node params: 0.003073931
#> Update drift: 0.0006411076
#> Update sticks: 0.0005190372
#> Update 3 hypers: 0.001432896
#> Update assignments: 0.5870321
#> Swap-nodes move: 0.02566195
#> Update Keep tree: 0.001096964

```

The total execution time is

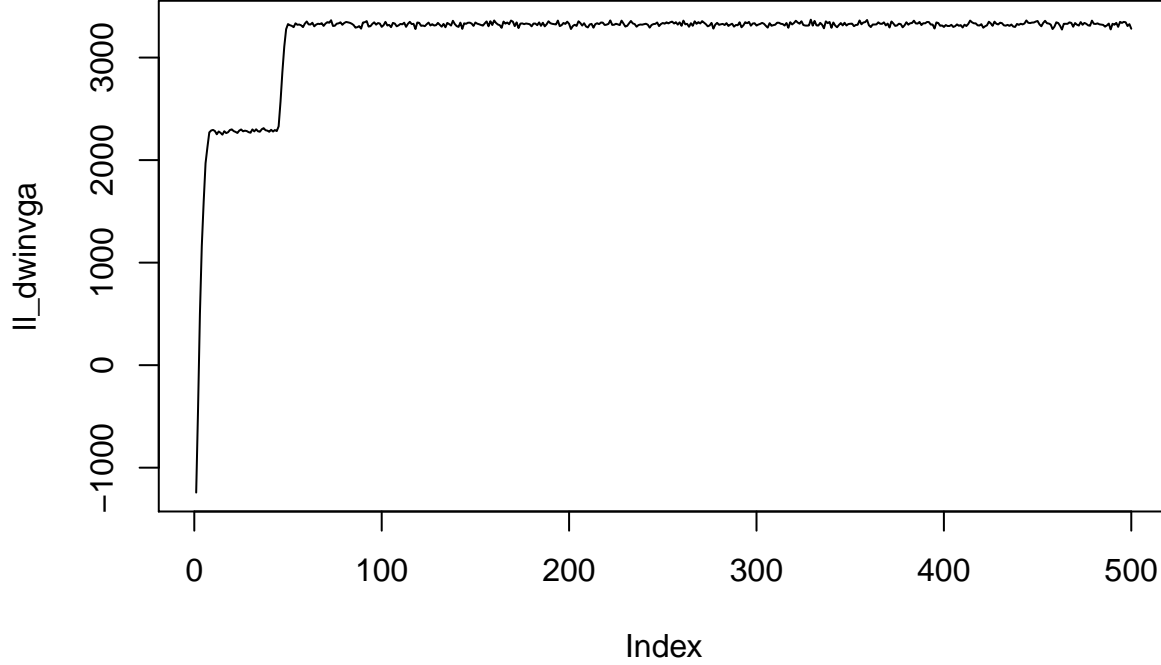
```

#> Time difference of 4.134636 mins

```

4.3 Results

we plot the log-likelihood curve to demonstrate the convergence of the Markov chain.

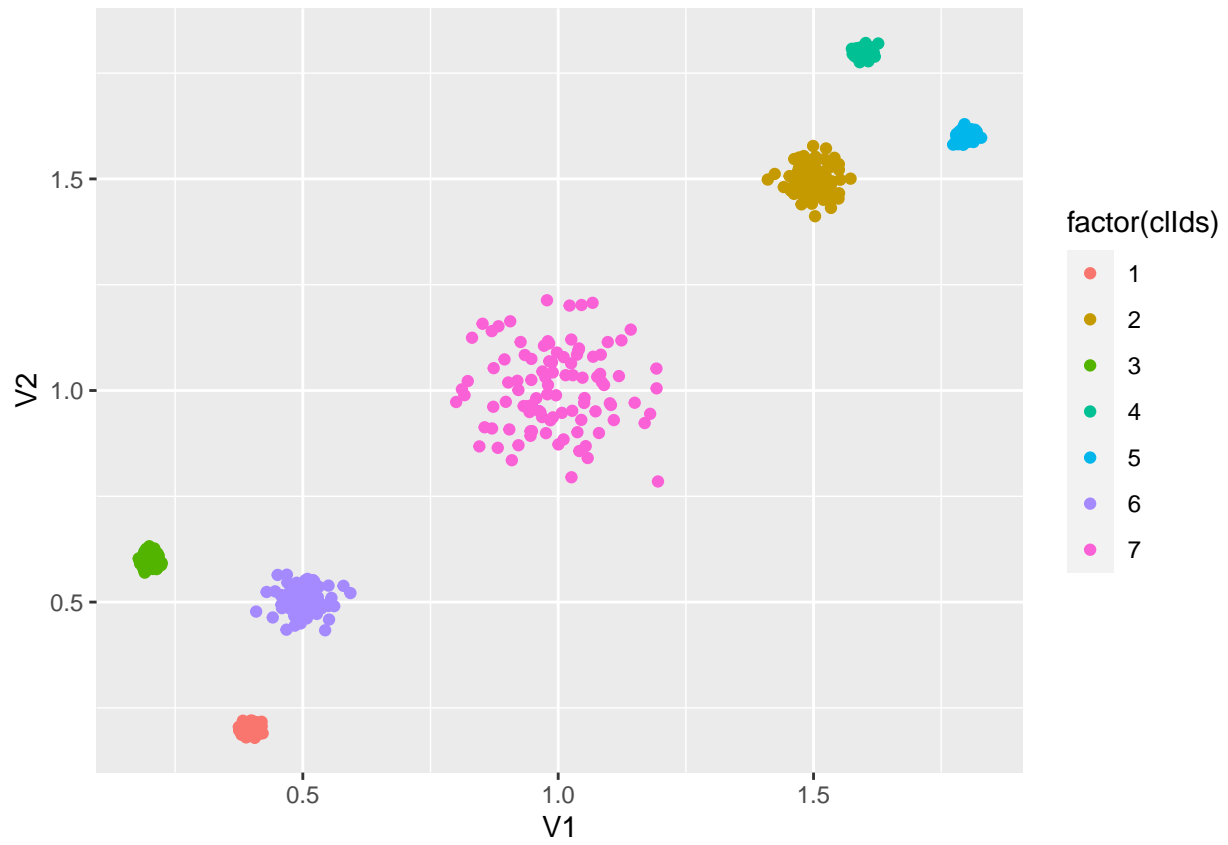


Finally, we plot the tree structure result using `forceNetwork` function in `networkD3` package.

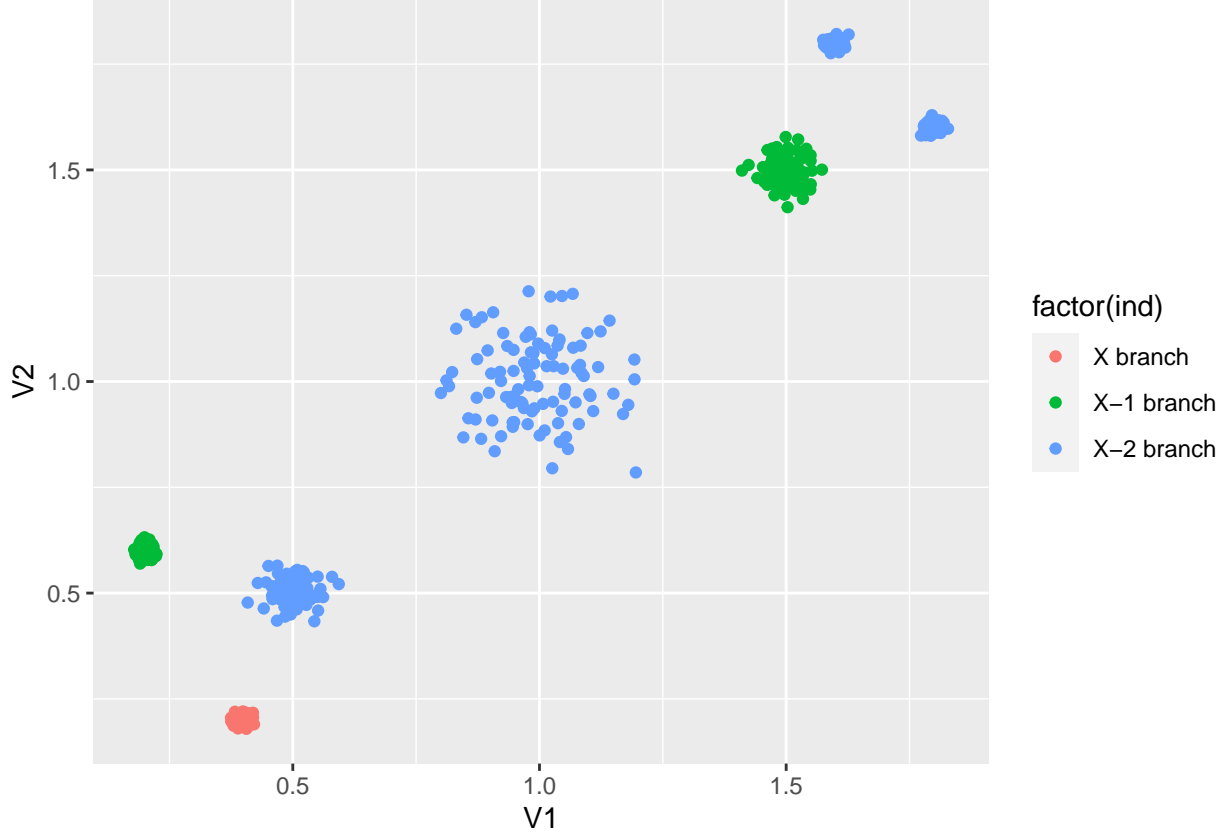
Select a single tree

To determine the final tree structure, (Yuan et al. 2015) provided an idea. They referred to the number of nodes whose weight is larger than 0.01 as the **big node number (BNN)**.

1. Group all MCMC samples into unique BNN categories.
2. Find out the most frequently occurring unique BNN group.
3. Choose the tree structure in this group which has the maximum marginal likelihood.



All the data are separated correctly based on their class.



We can observe that there exists a *trade-off* between “the clusters are very divided” and “the data of the same branch are more clustered together”. And this is controlled by the parameter η_{Θ} . By choosing an appropriate value of η_{Θ} , we can achieve these two tasks simultaneously. Unfortunately, in this study, the similarity within branches is not observed. We speculate that it may be related to the prior distribution of the drift λ . When η_{Θ} is smaller, the posterior of λ may tend to concentrate on larger values, which destroys the identifiability of the tree structure. In future research, we will try to use a bounded prior to solve this problem.

5 MNIST dataset

MNIST is a famous handwritten digits data set.

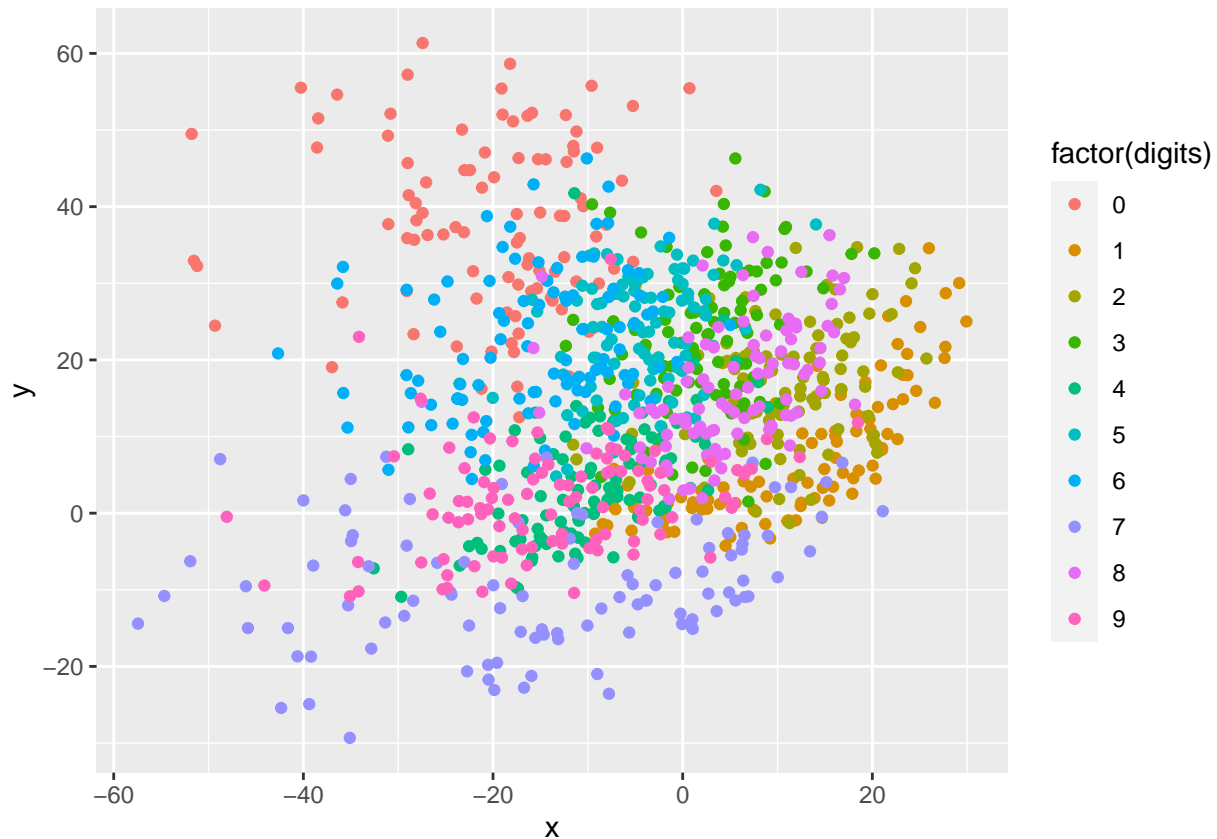
- Training: 60,000 images, Test: 10,000 images.
- Each image has 28x28 pixels.
- Each pixel of the image is in $[0, 255]$.

5.1 Prepare the MNIST-mini data

We extracted **5-dim features** for all training data from an *AutoEncoder* pre-trained using the whole training set. Then we sampled a subset of the training set where **each digit has 100 samples**. Therefore, the dimensions of data matrix we finally used is 1000×5 . (*The features extracted from AutoEncoder are partly separated, which is helpful for subsequent clustering.*)

The settings are the same as in the simulation study, except that

- burnIn = 500, Iter = 3,000.



5.2 Initialization and training

```
#> Initialization: D = 3 and W = 3
#> No. 1 child of root is initializing...
#> No. 2 child of root is initializing...
#> No. 3 child of root is initializing...
#> Initialization Is Over!
```

After initialization, we update all the parameters using Gibbs sampler.

```
#> ==== iter: 0
#> Update node params: 0.01248503
#> Update drift: 0.03904581
#> Update sticks: 0.004839897
#> Update 3 hypers: 0.004785061
#> Update assignments: 1.027414
#> Swap-nodes move: 0.02969909
#> Update Keep tree: 0.001091003
#> ==== iter: 1000
#> Update node params: 0.009140968
#> Update drift: 0.0006859303
#> Update sticks: 0.0007860661
#> Update 3 hypers: 0.001587868
#> Update assignments: 0.8346329
#> Swap-nodes move: 0.02141404
#> Update Keep tree: 0.0008919239
#> ==== iter: 2000
#> Update node params: 0.005545139
```

```

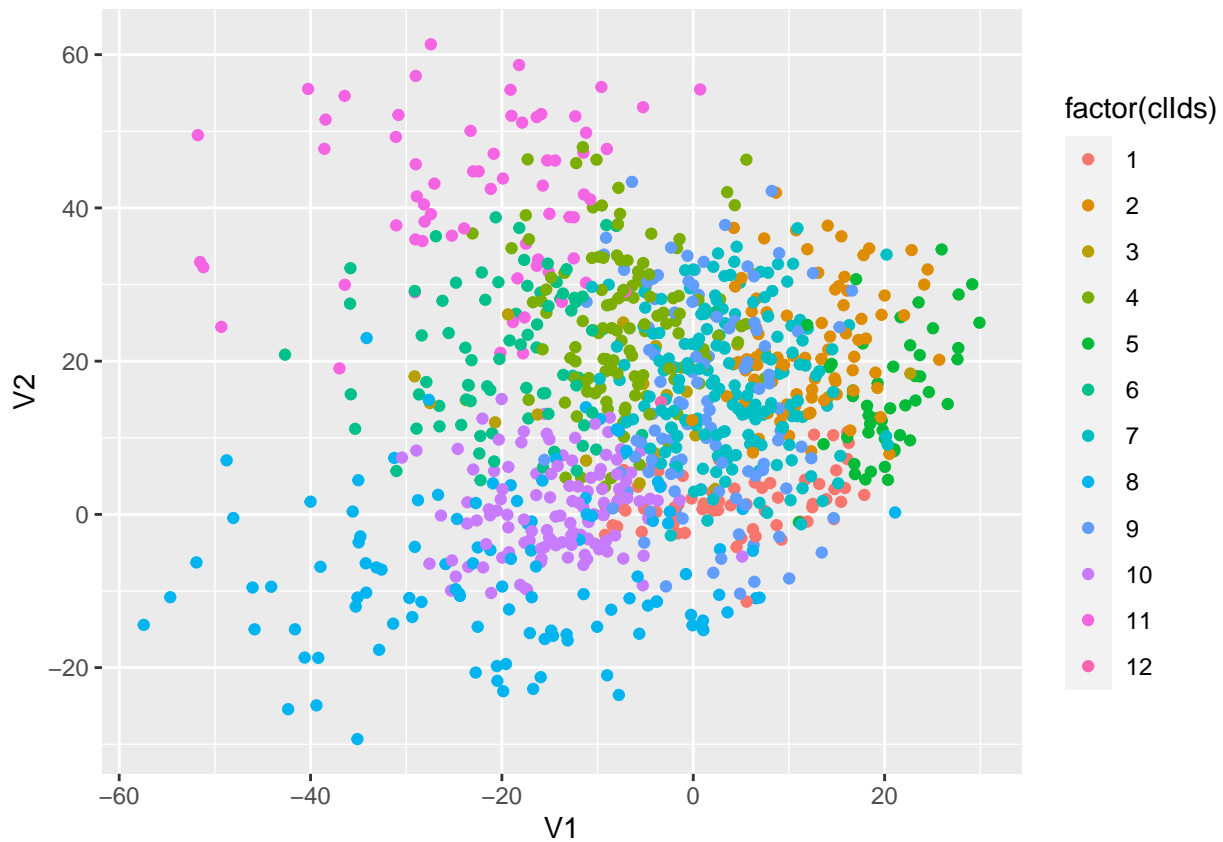
#> Update drift: 0.0006380081
#> Update sticks: 0.0007190704
#> Update 3 hypers: 0.001183987
#> Update assignments: 0.816386
#> Swap-nodes move: 0.03262591
#> Update Keep tree: 0.000921011
#> ==== iter: 3000
#> Update node params: 0.007287025
#> Update drift: 0.0006690025
#> Update sticks: 0.000783205
#> Update 3 hypers: 0.001724005
#> Update assignments: 0.8229921
#> Swap-nodes move: 0.027915
#> Update Keep tree: 0.000934124

#> Total execution time:
#> Time difference of 36.48625 mins

```

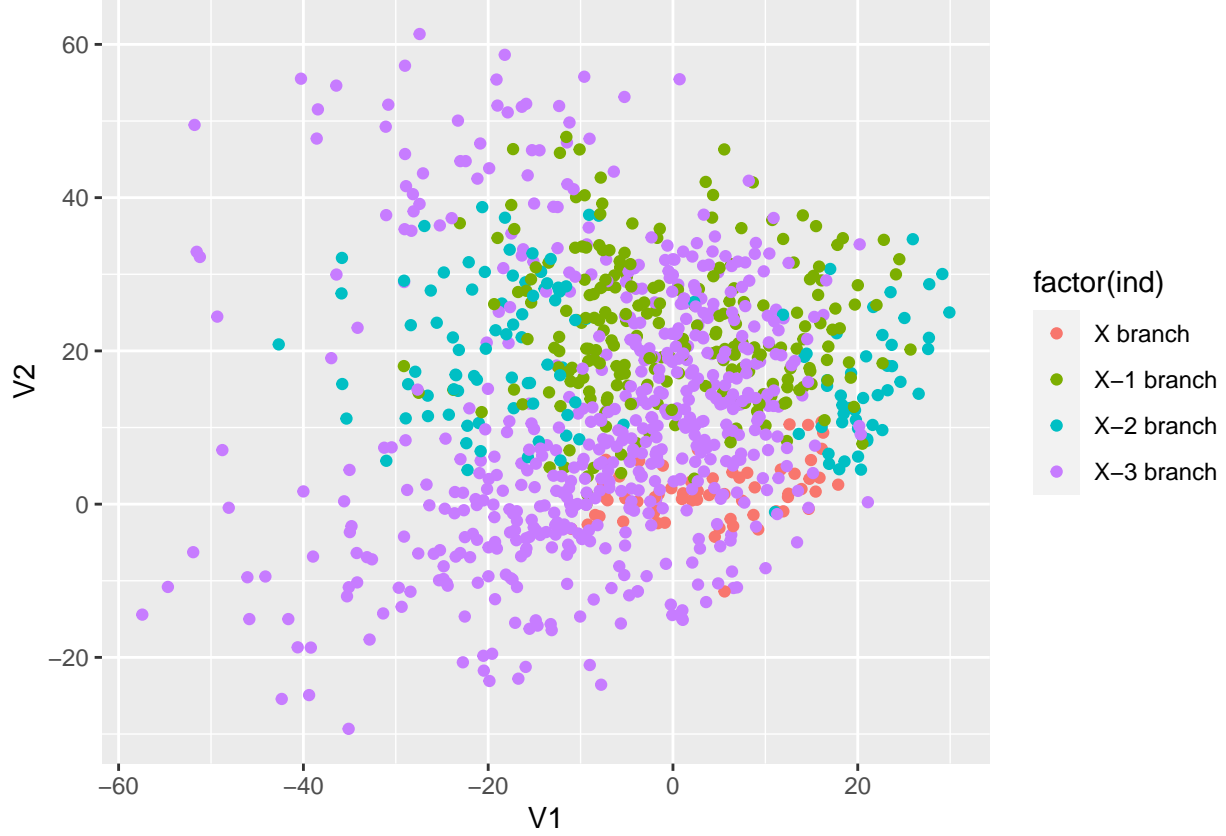
5.3 Results

The clustering result is shown in the following figure.



Result 1. Data are clustered almost based on their true classes, although there is a large overlap area. This demonstrates the ability of our model to cluster data with overlapping.

Next, we aggregate clusters on the same branch together to show the “similarity” (closely-located) within each branch.



Result 2. Data in the same branches are closer than others. This shows that our model is able to discover some unknown relationship behind clusters.

6 Conclusion

In this report, we proposed a nonparametric tree-structured prior, TSSB-DW, which was a truncation version of TSSB. Based on this prior, our model was able to achieve the two important tasks in two-level clustering analysis:

- clustering data into different groups
- discovering the tree-structured relationship behind groups.

The clustering abilities of our model were demonstrated by a simulation study and the application to the famous MNIST data set.

7 Future work

1. Compare the performances of clustering among our model and the baseline methods (hierarchical clustering and K-centroid), using some statistics such as **v-measure**.
2. Apply our model to some general real data.