

ddtlcm: An R package for overcoming weak separation in Bayesian latent class analysis via tree-regularization

7 September 2023

Summary

Latent class models (LCMs) have been a popular tool used by social, behavioral, and medical researchers to cluster individuals based on categorical responses over a collection of items. Traditional applications of LCMs often focus on scenarios where clear class separation is evident, resulting in well-defined and easily distinguishable classes. However, in real-world applications, weak class separation is a common challenge, where the classes are less distinct and overlapping. For example, nutritional epidemiologists often encounter weak class separation when deriving dietary patterns from dietary intake assessment, where each class profile represents the probabilities of exposure to a set of diet components. LCM-derived dietary patterns can exhibit strong similarities, or weak separation, resulting in numerical and inferential instabilities that challenge scientific interpretation. This issue is exacerbated in small-sized subpopulations.

To address these issues, we develop an R package `ddtlcm` that empowers LCMs to account for weak class separation. This package implements a tree-regularized Bayesian LCM that leverages statistical strength between latent classes to make better estimates using limited data. With a tree-structured prior distribution over class profiles, classes that share proximity to one another in the tree are shrunk towards ancestral classes *a priori*, with the degree of shrinkage varying across pre-specified major item groups. This software package takes data on multivariate binary responses over items in pre-specified major groups, and generates statistics and visualizations based on the inferred tree structures and LCM parameters. Overall, `ddtlcm` provides tools specifically designed to enhance the robustness and interpretability of LCMs in the presence of weak class separation, particularly useful for small sample sizes.

Statement of Need

A number of packages are capable of fitting LCMs in R. For example, `poLCA` (Linzer & Lewis, 2011) is a fully featured package that fits LCMs and latent class regression on polytomous outcome variables. `BayesLCA` (White & Murphy, 2014) is designed for LCMs in a Bayesian setting, incorporating three algorithms: expectation-maximization, Gibbs sampling and variational Bayes approximation. `randomLCA` (Beath, 2017) provides tools to perform LCMs with individual-specific random effects. These packages focus on LCMs where the class profiles are well-separated. Directly applying these packages to data that suffer from weak separation may result in large standard deviations of the class profiles, tendency to merge similar classes, and inaccurate individual class membership assignments. These phenomena are exacerbated especially when the sample size is small.

The package `ddtlcm` implements the tree-regularized LCM proposed in Li et al. (2023), a general framework to facilitate the sharing of information between classes to make better estimates of parameters using limited data. The model addresses weak separation for small sample sizes by (1) sharing statistical strength between classes guided by an unknown tree, and (2) accounting for varying degrees of shrinkage across major item groups. The proposed model uses a Dirichlet diffusion tree (DDT) process (Neal, 2003) as a fully probabilistic device to specify a prior distribution for the class profiles on the leaves of an unknown tree (hence termed “DDT-LCM”). Classes positioned closer on the tree exhibit more profile similarities. The degrees of separation by major item groups are modeled by group-specific diffusion variances.

Usage

In the following code, we use an example list of model parameters, named “`data_hchs`”, that comes with the package to demonstrate the utility of `ddtlcm` in deriving weakly separated latent class profiles. We start with demonstrating how a simulated dataset is generated. We next apply the primary model fitting function to the simulated dataset. Finally we summarize the fitted model and visualize the result.

Data Loading The “`data_hchs`” contains model parameters obtained from applying DDT-LCM to dietary assessment data described in Li et al. (2023). Specifically, the dataset includes a tree named “`tree_phylo`” (class “`phylo`”), a list of $J = 78$ food item labels, and a list of $G = 7$ pre-defined major food groups to which the food items belong. The food groups are dairy, fat, fruit, grain, meat, sugar, and vegetables.

```
install.packages("ddtlcm")
library(ddtlcm)
# load the data
data(data_hchs)
```

```

# unlist the elements into variables in the global environment
list2env(setNames(data_hchs, names(data_hchs)), envir = globalenv())
# look at items in group 1
g <- 1
# indices of the items in group 1
item_membership_list[g]

[[1]]
[1] 1 2 3 4 5 6 7 8 9 10 11

# names of the items in group 1. The name of the list element is
# the major food group
item_name_list[g]

$Dairy
[1] "dairy_1" "dairy_2" "dairy_3" "dairy_4" "dairy_5"
    "dairy_6" "dairy_7" "dairy_8" "dairy_9" "dairy_10"
[11] "dairy_11"

```

Data Simulation Simulation of a dataset is handled by the `simulate_lcm_given_tree()` function. Following the dietary assessment example, we simulate a multivariate binary data matrix of $N = 496$ subjects over the $J = 78$ food items, from $K = 6$ latent classes along “tree_phylo.” The resulting class profiles are weakly separated. Note that the number of latent classes equals the number of leaves in “tree_phylo.”

```

# number of individuals
N <- 496
# random seed to generate node parameters given the tree
seed_parameter = 1
# random seed to generate multivariate binary observations from LCM
seed_response = 1
# simulate data given the parameters
sim_data <- simulate_lcm_given_tree(tree_phylo, N,
  class_probability, item_membership_list, Sigma_by_group,
  root_node_location = 0, seed_parameter = 1, seed_response = 1)

```

Model Fitting The primary model fitting function is `ddtlcm_fit()`, which function implements a hybrid Metropolis-Hastings-within-Gibbs algorithm to sample from the posterior distribution of model parameters. We assume that the number of latent classes $K = 6$ is known. To use `ddtlcm_fit()`, we need to specify the number of classes (K), a matrix of multivariate binary observations (`data`), a list of item group memberships (`item_membership_list`), and the number of posterior samples to collect (`total_iters`). For the purpose of quick illustration, here we specify a small number `total_iters = 100`.

```

set.seed(999)
# number of latent classes, or number of leaves on the tree

```

```
K <- 6
result_hchs <- ddtlcm_fit(K = K, data = sim_data$response_matrix,
  item_membership_list = item_membership_list, total_iters = 100)
print(result_hchs)
```

DDT-LCM with K = 6 latent classes run on 496 observations and 78 items in 7 major groups. 100 iterations of posterior samples drawn.

Model Summary We next summarize the posterior chain using the function `summary()`. We discard the first 50 iterations as burn-ins (`burnin = 30`). To deal with identifiability of finite mixture models, we perform post-hoc label switching using the Equivalence Classes Representatives (ECR) method by specifying `relabel = TRUE`. To save space in the document, we do not print the summary result here (`be_quiet = TRUE`).

```
burnin <- 50
summarized_result <- summary(result_hchs, burnin, relabel = TRUE,
  be_quiet = TRUE)
```

Visualization A simple `plot()` function is available in the package to visualize the summarized result. By specifying `plot_option = "all"`, we plot both the *maximum a posterior* (MAP) tree and the class profiles. Plotting only the tree or the class profiles is also available through `plot_option = "profile"` or `plot_option = "tree"`.

Figure 1 displays the result obtained from the above model summary. On the left shows the MAP tree structure over the latent classes. Numbers indicate the corresponding branch lengths. On the right shows class profiles for the $J = 78$ food items belonging to $G = 7$ pre-defined major food groups, distinguished by different colors. The description of individual items is provided in Table S6.1 in the supplement of (Li et al., 2023). The numbers after the class labels in the facets indicate class prevalences along with 95% credible intervals. Error bars show the 95% credible intervals of item response probabilities from the posterior distribution.

```
plot(x = summarized_result, item_name_list = item_name_list,
  plot_option = "all")
```

Interactive Illustration with RShiny

to be added

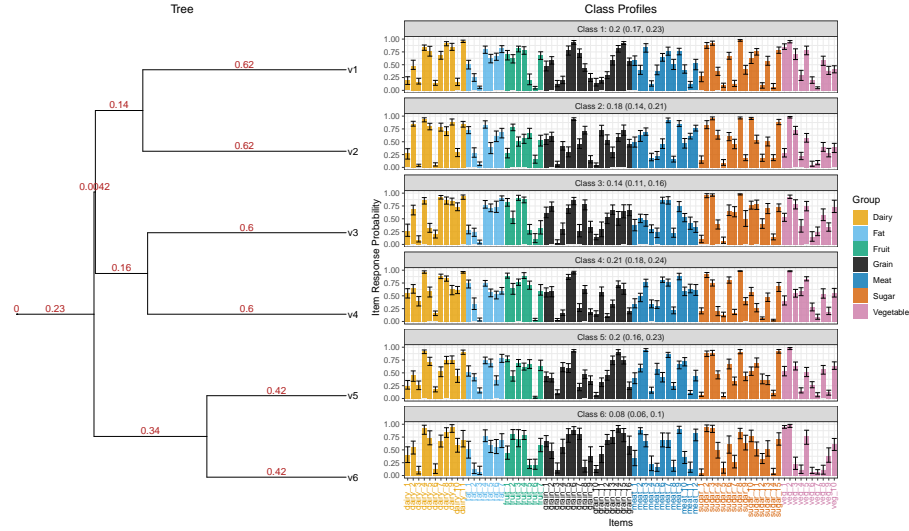


Figure 1: Posterior summary of DDT-LCM with $K = 6$ latent classes.

Conclusions

We developed an R package `ddtlcm` that solves weak class separation issue when deriving dietary patterns using LCMs. This paper offers a step-by-step case example that contextualizes the workflow with nutrition data. Details about usage and more elaborate examples can be found online at (<https://github.com/limengbinggz/ddtlcm>).

Acknowledgements

References

- Beath, K. J. (2017). randomLCA: An r package for latent class with random effects analysis. *Journal of Statistical Software*, 81, 1–25.
- Li, M., Stephenson, B., & Wu, Z. (2023). Tree-regularized bayesian latent class analysis for improving weakly separated dietary pattern subtyping in small-sized subpopulations. *arXiv Preprint arXiv:2306.04700*.
- Linzer, D. A., & Lewis, J. B. (2011). poLCA: An r package for polytomous variable latent class analysis. *Journal of Statistical Software*, 42, 1–29. <https://doi.org/10.18637/jss.v042.i10>
- Neal, R. M. (2003). Density modeling and clustering using Dirichlet diffusion trees. *Bayesian Statistics*, 7, 619–629.

White, A., & Murphy, T. B. (2014). BayesLCA: An r package for bayesian latent class analysis. *Journal of Statistical Software*, 61, 1–28. <https://doi.org/10.18637/jss.v061.i13>