

Clustering with Multinomial Outcomes

Shiro Kuriwaki*

October 2019

1 Data Generating Process

Setup Index individuals by $i \in \{1, \dots, N\}$ and the universe of races excluding the top of the ticket as $j \in \{1, \dots, D\}$. The data we observe is a length- D vector of votes \mathbf{Y}_i . Y_{ij} is a discrete response value, $Y_{ij} \in \{0, \dots, L\}$. For now, let's use $L = 1$ so that each vote y_{ij} be a binary variable for splitting their ticket or not. $Y_{ij} = 1$ would mean voter i splitting their ticket in some office j , with reference to a top of the ticket office like the President or Governor.

Parameters Individuals are endowed with a cluster (or type) $k \in 1 : K$, which is drawn from

$$Z_i \sim \text{Cat}(\boldsymbol{\theta}),$$

Where the length- K simplex $\boldsymbol{\theta}$ is called the mixing proportion.

Let $\mu_{k,j} \in [0, 1]$ be the parameter that governs each observation through the cluster assignment. Therefore $\boldsymbol{\mu}$ is a $\{K \times D \times (L + 1)\}$ array, where

$$\Pr(Y_{ij} = \ell \mid Z_i = k) = \mu_{kj\ell}.$$

In other words, for each individual who belongs in type k , we can their observed vector $\mathbf{Y}_i \mid Z_i = k$ is governed by a length- D parameter vector $\boldsymbol{\mu}_k$. Therefore, we can express the joint density as follows.

$$\Pr(\mathbf{Y}_i \mid Z_i = k, \boldsymbol{\theta}) = \prod_{j=1}^D \text{Cat}(Y_{ij} \mid \boldsymbol{\mu}_k) = \prod_{j=1}^D \prod_{\ell=1}^L \mu_{kj\ell}^{\mathbf{1}(Y_{ij}=\ell)} \quad (1.1)$$

In our $L = 1$ case, the Categorical reduces to a Bernoulli and we can get rid of the ℓ index altogether to let $\mu_{jk} \equiv \mu_{jk,(\ell=1)} = 1 - \mu_{jk,(\ell=0)}$ and so

*Thanks to Shusei Eshima, Soohahn Shin, and Soichiro Yamauchi for their help.

$$\Pr(\mathbf{Y}_i \mid Z_i = k, \boldsymbol{\theta}) = \prod_{j=1}^D \mu_{jk}^{Y_{ij}} (1 - \mu_{jk})^{1-Y_{ij}} \quad (1.2)$$

The benefit of this modeling exercise over that from a naive sample of $N \times D$ Bernoullis is that we have captured the correlations between variables.¹

2 EM

This mixture model lends itself to clustering analysis such as K -means. Although we can estimate this model in a Bayesian fashion by setting a prior for $\boldsymbol{\theta}$ and $\boldsymbol{\mu}$, the **Stan** program cannot reliably estimate clustering models like this one by MCMC because of label-switching and multimodality.

Because the model is simple enough, we can derive an algorithm to obtain the global solution for the parameters.² An EM implementation makes it possible to handle extensions, such as systematic missing data, multinomial outcomes, and covariates.

Complete Likelihood If we knew the cluster assignment, we would be able to write the complete likelihood (\mathbf{L}_{comp}) and the complete log-likelihood ($\mathcal{L}_{\text{comp}}$),

$$\begin{aligned} \mathbf{L}_{\text{comp}}(\boldsymbol{\mu}, \boldsymbol{\theta} \mid \mathbf{Y}, \mathbf{Z}) &= \prod_{i=1}^N \prod_{j=1}^D \prod_{k=1}^K \Pr(Z_i = k) \prod_{\ell=1}^L \Pr(Y_{ij} = 1 \mid Z_i = k, \boldsymbol{\mu})^{\mathbf{1}(Y_{ij}=\ell \mid Z_i=k)} \\ \mathcal{L}_{\text{comp}}(\boldsymbol{\mu}, \boldsymbol{\theta} \mid \mathbf{Y}, \mathbf{Z}) &= \sum_{i=1}^N \sum_{j=1}^D \sum_{k=1}^K \sum_{\ell=0}^2 \mathbf{1}\{Y_{ij} = \ell, Z_i = k\} \left\{ \log \Pr(Y_{ij} = \ell \mid Z_i = k, \boldsymbol{\mu}) \right\} \\ &\quad + \sum_{i=1}^N \sum_{k=1}^K \mathbf{1}\{Z_i = k\} \log \Pr(Z_i = k \mid \boldsymbol{\theta}) \end{aligned} \quad (2.1)$$

We first take expectations over the latent variable Z_i ,

$$\begin{aligned} \mathbb{E}[\mathcal{L}_{\text{comp}}] &= \sum_{i=1}^N \sum_{j=1}^D \sum_{k=1}^K \sum_{\ell=0}^2 \mathbf{1}\{Y_{ij} = \ell\} \mathbb{E}[\mathbf{1}\{Z_i = k\}] \underbrace{\left\{ \log \Pr(Y_{ij} = \ell \mid Z_i = k, \boldsymbol{\mu}) \right\}}_{\equiv \log \boldsymbol{\mu}_{kj\ell}} \\ &\quad + \sum_{i=1}^N \sum_{k=1}^K \mathbb{E}[\mathbf{1}\{Z_i = k\}] \underbrace{\log \Pr(Z_i = k \mid \boldsymbol{\theta})}_{\equiv \log \boldsymbol{\theta}_k} \end{aligned} \quad (2.2)$$

¹ That dependency can be expressed as $\mathbb{E}(\mathbf{y}) = \sum_{k=1}^K \theta_k \boldsymbol{\mu}_k$ and $\text{Cov}(\mathbf{y}) = \sum_k \theta_k (\boldsymbol{\Sigma}_k \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top) - \mathbb{E}(\mathbf{y}) \mathbb{E}(\mathbf{y})^\top$, where $\boldsymbol{\Sigma}_k = \text{diag}(\mu_{jk}(1 - \mu_{jk}))$.

² Thanks to Soichiro Yamauchi for deriving this algorithm in the original iteration.

Let's define this unknown quantity as

$$\zeta_{ik} \equiv \mathbb{E}[\mathbf{1}(Z_i = k)].$$

Then the E-step can be the normalized version of the posterior probability marginalized by the mixing proportion,

$$\hat{\zeta}_{ij} \propto \theta_k \prod_{j=1}^D \underbrace{\prod_{\ell=0}^2 (\mu_{kj,\ell})^{\mathbf{1}(Y_{ij}=\ell)}}_{\mu_{kj,Y_{ij}}} \quad (2.3)$$

The M-step is derived by taking the derivatives of $\mathbb{E}[\mathcal{L}_{\text{comp}}]$ with respect to the model parameters μ and θ . This leads to a MLE-like M-step

$$\hat{\theta}_k = \frac{1}{N} \sum_{i=1}^N \zeta_{ik} \quad (2.4)$$

$$\hat{\mu}_{jk\ell} = \frac{\sum_{i=1}^N \mathbf{1}(Y_{ij} = \ell) \zeta_{ik}}{\sum_{i=1}^N \zeta_{ik}} \quad (2.5)$$

EM Implementation We first need to set initial values for μ and θ . I do this by letting $\theta^{(0)} = (\frac{1}{K}, \dots, \frac{1}{K})$, randomly assigning an initial cluster assignment $Z'_i \sim \text{Cat}(\theta^{(0)})$, and setting the initial μ by the sample means of the data within those initial assignments, $\mu_{jk}^{(0)} = \frac{\sum_{i=1}^N \mathbf{1}(Y_{ij}=1) \mathbf{1}(Z'_i=k)}{\sum_{i=1}^N \mathbf{1}(Z'_i=k)}$.

Then we iterate as follows.

For each voter i , compute the probability that they belong in cluster k (E-step):

$$\zeta_{ik} \leftarrow \frac{\theta_k \prod_{j=1}^D \mu_{kj,Y_{ij}}}{\sum_{k'=1}^K \theta_{k'} \prod_{j=1}^D \mu_{k'j,Y_{ij}}} \quad (2.6)$$

Then given those cluster probabilities, update the parameters with the MLE (M-step)

$$\hat{\theta}_k \leftarrow \frac{1}{N} \sum_{i=1}^N \hat{\zeta}_{ik} \quad (2.7)$$

$$\hat{\mu}_{jk1} \leftarrow \frac{\sum_{i=1}^N \mathbf{1}(Y_{ij} = 1) \hat{\zeta}_{ik}}{\sum_{i=1}^N \hat{\zeta}_{ik}} \quad (2.8)$$

$$\hat{\mu}_{jk0} = 1 - \hat{\mu}_{jk1} \quad (2.9)$$

We repeat this until convergence.

Convergence We evaluate coverage by the observed log likelihood,

$$\begin{aligned}\mathbf{L}_{\text{obs}} &= \prod_{i=1}^N \sum_{k=1}^K \theta_k f_k(\mathbf{Y}_i | \boldsymbol{\mu}_k) \\ &= \prod_{i=1}^N \sum_{k=1}^K \theta_k \prod_{j=1}^D \mu_{k'j, Y_{ij}}\end{aligned}$$

So the log-likelihood is

$$\mathcal{L}_{\text{obs}} = \sum_{i=1}^N \left[\prod_{k=1}^K \left\{ \log \theta_k + \sum_{j=1}^D \sum_{\ell=1}^L \mathbf{1}(Y_{ij} = \ell) \log(\mu_{kj\ell}) \right\} \right] \quad (2.10)$$

Calculating (eq. 2.10) is computationally intensive, so a quick way to check convergence is to track the maximum of the change in parameters which are all on the probability scale, i.e. $\max \left\{ |\hat{\theta}_1^{(t+1)} - \hat{\theta}_1^{(t)}|, \dots, |\hat{\mu}_{K,D}^{(t+1)} - \hat{\mu}_{K,D}^{(t)}| \right\}$

Speed-Ups Because this EM algorithm deals with discrete data, the algorithm needs only sufficient statistics. In our setting the unique number of voting profiles is much smaller than the number of observations, because vote vectors follow a systematic pattern and most votes are straight-ticket votes. Therefore, we can re-format the dataset so that each row is a unique combination.