Wiley Online Library

👤 Log in / Register

↵ Go to old article view

🔑 Get access

ORIGINAL ARTICLE

# A Bayesian mixture model for clustering and selection of feature occurrence rates under mean constraints

Qiwei Li,   Michele Guindani,   Brian J. Reich,   Howard D. Bondell,   Marina Vannucci ✉

(Am) score

## Abstract

In this paper, we consider the problem of modeling a matrix of count data, where multiple features are observed as counts over a number of samples. Due to the nature of the data generating mechanism, such data are often characterized by a high number of zeros and overdispersion. In order to take into account the skewness and heterogeneity of the data, some type of normalization and regularization is necessary for conducting inference on the occurrences of features across samples. We propose a zero-inflated Poisson mixture modeling framework that incorporates a model-based normalization through prior distributions with mean constraints, as well as a feature selection mechanism, which allows us to identify a parsimonious set of discriminatory features, and simultaneously cluster the samples into homogenous groups. We show how our approach improves on the accuracy of the clustering with respect to more standard approaches for the analysis of count data, by means of a simulation study and an application to a bag-of-words benchmark data set, where the features are represented by the frequencies of occurrence of each word.

🔑 **Get access to the full text of this article**

» **Article Information**

» **Data Accessibility**

⌄ **Related content**

## Articles related to the one you are viewing

The articles below have been selected for you based on the article you are currently viewing.

### Robustifying Bayesian nonparametric mixtures for count data

Antonio Canale, Igor Prünster

28 April 2016