

# **Dealing with overdispersion**

# in multivariate count data

Noemi Corsini & Cinzia Viroli

Statistical methods and models for complex data — Poster Session: 21 September 2022

# **Framework and motivation**

The overdispersion or extra-variation is a recurring phenomenon when dealing with **counts and categorical data**.

 $\hookrightarrow$  When fitting a binomial, a multinomial or a Poisson model if the data exhibit a larger variability than that the model is able to explain, the sampling variation will be greater than the estimated variation accounted by the model.

Overdispersion has specific causes and consequences



In both scenarios it is evident that the Deep Dirichlet-Multinomial distribution better and better approximates the empirical true variance as K increases.

#### The price of flexibility is a greater computational burden

	First simulation		Second simulation	
		S _		
2 -		— 52		

- **Causes** may be the result of data aggregation such as clumped sampling, the correlation between individual responses or the additional experimental variability.
- Inferential consequences are imprecise estimates and biased standard errors that make model selection, interpretability and prediction unreliable.

Our focus is on **multivariate count data**, whose natural probabilistic model is the multinomial distribution. However, in presence of overdispersion, this model may lead to a nominal variance well below the empirical one.

### How to cope with overdispersion?

- Quasi-likelihood approach where second order hypotheses allow for relaxing constraints on the variance structure;
- Likelihood-based approach exploiting distributions that differ in the variance and correlation structure, that can be
  - **Negative** Multinomial (MN), Dirichlet-Multinomial (DM), Random Clumped Multinomial (RCM);

Figure 1. Structure of DDM model. A hidden layer of nodes is introduced to better capture the overdispersion.

Let  $\boldsymbol{\theta}_k = \boldsymbol{\beta}(1 + \boldsymbol{\alpha}_k), \ \theta_{Ok} = \sum_{j=1}^p \beta_j (1 + \alpha_{jk}), \ \boldsymbol{\pi}_k = \frac{\boldsymbol{\theta}_k}{\theta_{Ok}} \text{ and } \rho_k^2 = \boldsymbol{\theta}_k$  $1/(1 + \theta_{k0})$ , then we can define expectation and variance as:  $\mathbb{E}\left[\mathbf{Y}\right] = \sum \omega_k m \boldsymbol{\pi}_k,$  $\mathbb{V}[\mathbf{Y}] = \sum \omega_k m\{diag(\boldsymbol{\pi}_k) - \boldsymbol{\pi}_k \boldsymbol{\pi}'_k\}(1 + \rho_k^2(m-1))$  $+\sum_{k=1}^{K}\omega_{k}m^{2}\boldsymbol{\pi}_{k}\boldsymbol{\pi}_{k}^{\prime}-m^{2}\left(\sum_{k=1}^{K}\omega_{k}\boldsymbol{\pi}_{k}\right)\left(\sum_{k=1}^{K}\omega_{k}\boldsymbol{\pi}_{k}\right)^{\top}.$ 

The variance can be split into two components:

- A weighted sum of **within variances**;
- between variance part that captures both over and underdispersion.



**Figure 3.** Evolution of the empirical true variance - solid black line - with respect to the estimated variances of the different models as the number of zeros in the dataset increases.

# **DDM** asymptotic behavior

Exploiting the same data generating process, we analyze the dynamic behavior of the DDM variance as function of K in the case of maximum overdispersion.



- **Positive** Negative Multinomial (NM);
- **General** Generalized Dirichlet Multinomial (GDM).

#### Main objectives

- **Comprehensive comparison** of probabilistic models that capture extra-variation in multivariate count data.
- Introduction of a **new model** that extends the Dirichlet-Multinomial in a deep fashion.
- **3.** Analyses of **empirical performances and properties** through a broad simulation study.

# **Deep Dirichlet Multinomial**

In order to deal with overdispersion, we propose a new model called **Deep Dirichlet-Multinomial (DDM)**. Derived from [2], it is a mixture of Dirichlet-Multinomial distributions with restrictions on the parameters.

Let  $DM(\boldsymbol{\theta}, m)$  be the probability mass function of a Dirichlet Multinomial distribution with parameters  $\theta$  and size m, then This model has a flexible correlation structure among vari**ables** that comes at the price of an high number of parameters to be estimated which largely increase with K.

The model parameters are estimated through a generalized EM algorithm with a quasi-Newton optimization step.

#### Simulation study

The capabilities of the proposed model are compared to the likelihood-based models in **two simulation studies** that differ in the way the overdispersion is introduced.

- Randomly add zeros into the data;
- 2. Gradually add zeros by replacing the smallest counts starting from cells with frequency one.

Different scenarios are simulated from a  $MN(\pi, m = 100)$ , starting from the case of lack of extra-variation (0% of zeros added), to the case of 90% of zeros added.



**Figure 4.** The red line describes the evolution of the fitted DDM variance when K increases to  $+\infty$ . The dashed black line is the sample variance.

The empirical analysis suggests that the **estimated variance** tends to the computed variance when the number of elements K of the mixture goes to  $+\infty$ .

#### Conclusion

- A new approach, the Deep Dirichlet Multinomial, is proposed. Its variance shows desirable properties  $\rightarrow$  It can ideally be split in within and between variance;
  - $\longrightarrow$  It tends to the computed one when  $K \rightarrow \infty$ .
- Simulations show that the new model deals better with overdispersed data, compare to other likelihood based solutions.

# **Further readings**

the probability distribution of the DDM model is defined as

$$\mathbb{P}(\mathbf{Y} = \mathbf{y}) = \sum_{k=1}^{K} \omega_k \operatorname{DM}(\boldsymbol{\beta}(1 + \boldsymbol{\alpha}_k), m)$$

where  $\omega_k$ , for  $k = 1, \ldots, K$ , are the weights of the mixture defined such that they satisfy  $0 < \omega_k < 1$  and  $\sum_{k=1}^{K} \omega_k = 1$ while  $\beta > 0$  is a vector of length *p*.

Each  $\alpha_k \in (-1, 1)$  is a *p*-dimensional vector that can be interpreted as **perturbation parameter** used to adjust  $\beta$  and **get** a more flexible model that behaves better in case of overdispersion.

- Positive values lead to a larger effect of  $\theta_i = \beta_i (1 + \alpha_i)$ ;
- Negative values show which categories have more zeros.



**Figure 2.** Average euclidean distances between the empirical and estimated variances in the two simulations.

As the mixture components increase from K = 2 to K = 20, the **DDM model** has a smaller and smaller euclidean distance.

Noemi Corsini and Cinzia Viroli. "Dealing with overdispersion in multivariate count data". In: Computational Statistics & Data Analysis 170 (2022), p. 107447.

Cinzia Viroli and Laura Anderlucci. "Deep mixtures of un-[2] igrams for uncovering topics in textual data". In: Statistics and Computing 31.3 (2021), pp. 1–10.

# **Contact information**

- Noemi Corsini, Ph.D. student
- University of Padua
- noemi.corsini@phd.unipd.it
- https://noecors.github.io 0





2022 - Statistical methods and models for complex data, Padova