# Sums of possibly associated multivariate indicator functions: The Conway–Maxwell-Multinomial distribution

**Joseph B. Kadane[a] and Zhi Wang[b]**

[a]*Carnegie Mellon University*
[b]*Columbia University*

**Abstract.** The Conway–Maxwell-Multinomial distribution is studied in this paper. Its properties are demonstrated, including sufficient statistics and conditions for the propriety of posterior distributions derived from it. An application is given using data from Mendel's ground-breaking genetic studies.

## 1 The Conway–Maxwell-Multinomial distribution

The Conway–Maxwell-Multinomial (COMM) Distribution has probability mass function (for fixed $m$)

$$P\{\boldsymbol{X} = \boldsymbol{k}|(\boldsymbol{p}, \nu)\} = \binom{m}{\boldsymbol{k}}^{\nu} \prod_{i=1}^{r} p_i^{k_i} \Big/ \sum_{\boldsymbol{j} \in D} \binom{m}{\boldsymbol{j}}^{\nu} \prod_{i=1}^{r} p_i^{j_i}, \qquad \boldsymbol{k} \in D,$$

where

$$\boldsymbol{p} = (p_1, \dots, p_r), \qquad p_i > 0, \qquad \sum_{i=1}^{r} p_i = 1,$$

$$\boldsymbol{k} = (k_1, \dots, k_r), \qquad k_i \geq 0, \qquad \sum_{i=1}^{r} k_i = m, \qquad k_i\text{'s integers,}$$

$$\binom{m}{\boldsymbol{k}} = \frac{m!}{k_1! k_2! \cdots k_r!}$$

and $D$ is the set of vectors of integers $\boldsymbol{j}$ satisfying $j_i \geq 0$ and $\sum_{i=1}^{r} j_i = m$.

This distribution is a generalization of other distributions as follows:

$$r = 2, \qquad \nu = 1 \qquad \text{binomial,}$$

$$r > 2, \qquad \nu = 1 \qquad \text{multinomial,}$$

$$r = 2, \qquad \nu \neq 1 \qquad \text{Conway–Maxwell Binomial.}$$

**Barycentric Display of the Conway−Maxwell Trinomial**

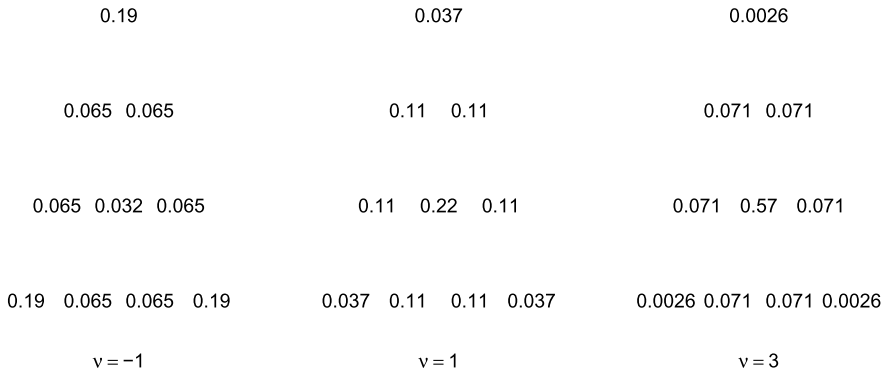| | | |
|---|---|---|
| 0.19 | 0.037 | 0.0026 |
| 0.065  0.065 | 0.11   0.11 | 0.071  0.071 |
| 0.065  0.032  0.065 | 0.11   0.22   0.11 | 0.071   0.57   0.071 |
| 0.19  0.065  0.065  0.19 | 0.037  0.11   0.11  0.037 | 0.0026 0.071  0.071 0.0026 |
| $\nu = -1$ | $\nu = 1$ | $\nu = 3$ |

**Figure 1**

Compared to the multinomial distribution, the novelty is in the parameter $\nu$. As illustrated in Figure 1 (for the special case $m = 3$, $p = (1/3, 1/3, 1/3)$, when $\nu > 1$, the center of the distribution is upweighted relative to the multinomial distribution, indicating negative association among the underlying multivariate indicator functions. Conversely, when $\nu < 1$, the tails of the distribution are upweighted relative to the multinomial distribution, indicating positive association.

The remainder of this paper is organized as follows: Section 2 gives an example of the use of the COMM distribution to a small part of the controversial Mendel data. Since the question there is whether Mendel's data are "too good to be true," *that is*, whether they fit the multinomial hypothesis too well, it seems ideal as an example. Section 3 gives several of the properties of the COMM distribution, results that are generalizations of the principal findings of Kadane (2016) for the Conway–Maxwell Binomial distributions. The proofs are in the Appendix.

## 2  An example: Mendel's data and Fisher's analysis

Gregor Mendel, an Augustinian monk published (1866) results of his experiments on garden peas, and proposed the model of genetics since known as Mendelian genetics. Fisher (1936), following Weldon (1902), did an analysis suggesting that Mendel's data were "too good to be true," that is, that they fit a binomial or multinomial distribution too well. This finding has led to intense speculation about how that might have occurred. However, since Mendel's papers were destroyed after his death, it is unlikely that additional evidence on this matter will now be discovered. Some recent work on the general matter can be found in Franklin et al. (2008) and Pires and Branco (2010).

Since Fisher is such a dominant figure in both statistics and genetics, it might be useful to give a fast gloss of his argument. Fisher uses chi-square measures of
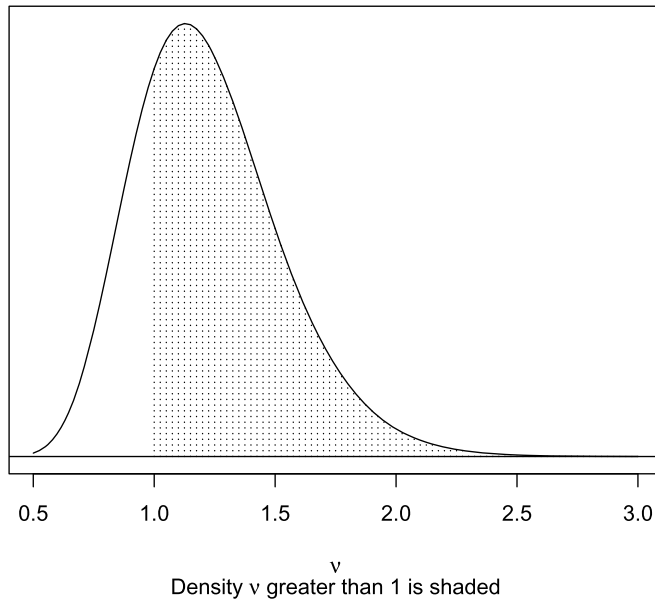
goodness of fit for Mendel's data, treating them as independent across different experiments. By adding the chi-squares and the degrees of freedom, he finds an extra-ordinary degree of coincidence between Mendel's theory and his data. The chi-square analysis of Fisher (started by E. Pearson about 1900) was (and still is) used as the basis of a test of significance. If the chi-square is too large, the null hypothesis is rejected. According to Fisher (1959), if the null hypothesis is rejected, "The force with which such a conclusion is supported is that of the simple disjunction: *Either* an exceptionally rare chance has occurred, *or* the theory of random distribution is not true" (p. 39). Fisher's theory does not permit one to say which of the two possibilities is the case, nor to give a probability for it. Furthermore, if significance is not achieved, nothing can be concluded. In order for the probability distribution that forms the basis of a chi-square test to be valid, the hypothesis to be tested must be declared before the data are examined.

Viewed in this light, there are several gaps between Fisher's calculations and his conclusion. Fisher is rejecting the multinomial null hypothesis if the chi-square is too small, which would be legitimate if the hypothesis test were declared before Weldon pointed the way, or if Fisher routinely used a two-tailed chi-square test. Neither is the case. And one still has Fisher's disjunction to contend with. Nonetheless, Fisher is a superb data-analyst, and we should not be interpreted as challenging his conclusion.

Our intent here is not to enter into these controversial waters, but simply to show the kind of contribution that the COMM distribution might make in this setting. As explained above, when $\nu = 1$ the usual multinomial distribution results. When $\nu < 1$, the tails of the distribution are upweighted relative to the center. Conversely, when $\nu > 1$ the center of the distribution is upweighted relative to the tails. Because $\nu$ is a continuous parameter ($-\infty \leq \nu \leq \infty$), the COMM distribution gracefully handles both positive and negative association.

In the case of the Mendel data, the notion that Mendel's data are "too good to be true" is translated into the hypothesis that $\nu > 1$. To illustrate how this works, we took a simple trinomial piece of Mendel's data, namely the first line of Fisher's Table 1 (1936). The data are 36, 60 and 28 for plants in a Bifactoral experiment, where the expected proportions are (1/4, 1/2, 1/4). Along with the COMM likelihood, we must declare a prior. In this case it makes sense to have an opinionated prior on ($p_1$, $p_2$ and $p_3$), namely that we are sure the correct values are, respectively, (1/4, 1/2, 1/4), as predicted by Mendelian genetics. For a prior on $\nu$, we'll take $\nu$ to have a unit normal distribution centered at 1. In this way, the prior on $\nu$ does not influence the analysis toward or away from the space $\nu > 1$.

The posterior distribution is displayed in Figure 2. Approximately 80% of the probability is above $\nu = 1$, confirming that, in this small piece of Mendel's data, the data are likely to be more negatively associated than they would be were the distribution trinomial. As this example shows, a Bayesian analysis with a COMM likelihood can deliver the probability that Mendel's data are more negatively associated than would be expected under the trinomial hypothesis. Again, this should

Posterior Density of ν for snippet of Mendel Data



ν
Density ν greater than 1 is shaded

**Figure 2**

not be taken as a full analysis of Mendel's data, and says nothing about how the data came to be the way they are. A fuller analysis of the Mendel data would address the complete dataset.

## 3 Properties of the Conway–Maxwell Multinomial distribution

To begin, we should explain how the names Conway and Maxwell came to be associated with this distribution. They wrote a short paper (Conway and Maxwell (1962)) in which they discuss a generalization of the Poisson distribution having pmf proportional to

$$\lambda^x / (x!)^\nu,$$

where $x = 0, 1, 2, \ldots$, and $\lambda$ and $\nu$ are parameters. When $\nu = 1$, the usual Poisson distribution results. This distribution turned out to be useful to model count data because it allows for heavier or lighter tails than the Poisson model (Boatwright, Borle and Kadane (2003), Borle et al. (2005)). Because Conway and Maxwell proposed it first, it became known as the Conway–Maxwell–Poisson distribution (Shmueli et al. (2004)).

It is a simple calculation to show that if $X_1$ has a Poisson distribution with parameter $\lambda_1$, and is independent of $X_2$, a Poisson random variable with parameter

$\lambda_2$, then the distribution of $X_1$ conditional on the event $X_1 + X_2 = m$ is Binomial with parameter $p = \lambda_1/(\lambda_1 + \lambda_2)$ and $n$. It is only a slightly more complicated calculation to show the multivariate generalization: If $X_1, X_2, \ldots, X_r$ are independent Poisson random variables with means $\lambda_1, \lambda_2, \ldots, \lambda_\nu$, then the distribution of $(X_1, X_2, \ldots, X_r)$ conditional on the event that $\sum_{i=1}^{r} X_i = m$ is multinomial with parameter vector $\boldsymbol{p} = (\lambda_1, \lambda_2, \ldots, \lambda_r)/\sum_{i=1}^{r} \lambda_i$ and $m$.

Shmueli et al. (2004) generalizes the relationship between the Poisson and binomial distributions to the Conway–Maxwell–Poisson as follows: Suppose $X_1$ has a Conway–Maxwell–Poisson (CMP) distribution with parameters $(\lambda_1, \nu)$, and $X_2$ is independently distributed CMP $(\lambda_2, \nu)$. Then $X_1$ conditional on the event $X_1 + X_2 = m$ has a Conway–Maxwell Binomial distribution with parameters $p = \lambda_1/(\lambda_1 + \lambda_2)$, $\nu$ and $m$. It is reasonable to hope, then, that the conditional distribution of $r$ independent Conway–Maxwell–Poisson distributions with respective parameters $\lambda_1, \ldots, \lambda_r$ and $\nu$, would have a Conway–Maxwell Multinomial distribution. That this is the case is the content of Proposition 1:

**Proposition 1.** *Suppose $X_1, \ldots, X_r$ are independently distributed with probability mass function Conway–Maxwell Poisson $X_i \sim \mathrm{CMP}(\lambda_i, \nu)$:*

$$P\{X_i = s_i \mid \lambda_i, \nu\} = \frac{\lambda_i^{s_i}}{(s_i!)^\nu Z(\lambda_i, \nu)},$$

*where $Z(\lambda_i, \nu) = \sum_{j=0}^{\infty} \frac{\lambda_i^j}{(j!)^\nu}$.*

*Then $\boldsymbol{X} \mid \sum_{i=1}^{r} X_i = m$ has a COMM distribution with parameters $p_i = \lambda_i/\lambda$ and $\nu$, where $\lambda = \sum_{i=1}^{r} \lambda_i$.*

**Proof.** Let $S = \sum_{i=1}^{r} X_i$, $\boldsymbol{\lambda} = (\lambda, \ldots, \lambda_r)$ and $G(\boldsymbol{p}, \nu) = \sum_{\boldsymbol{j} \in D} \binom{m}{\boldsymbol{j}}^\nu \prod_{i=1}^{r} p_i^{j_i}$.
Then

$$P\{S = m\} = \sum_{\boldsymbol{j} \in D} \prod_{i=1}^{r} \frac{\lambda_i^{j_i}}{(j_i!)^\nu Z(\lambda_i, \nu)}$$

$$= \frac{\lambda^m}{(m!)^\nu} \cdot \frac{1}{\prod_{i=1}^{r} Z(\lambda_i, \nu)} \cdot \sum_{\boldsymbol{j} \in D} \prod_{i=1}^{r} (\lambda_i/\lambda)^{j_i} \binom{m}{\boldsymbol{j}}^\nu$$

$$= \frac{\lambda^m}{(m!)^\nu} \frac{1}{\prod_{i=1}^{r} Z(\lambda_i, \nu)} G(\boldsymbol{\lambda}/\lambda, \nu).$$

Hence,

$$P\{\boldsymbol{X} = \boldsymbol{k} \mid S = m\} = \prod_{i=1}^{r} \frac{\lambda_i^{k_i}}{(k_i!)^\nu Z(\lambda_i, \nu)} \bigg/ \frac{\lambda^m}{(m!)^\nu} G(\boldsymbol{\lambda}/\lambda, \nu) \prod_{i=1}^{r} \frac{1}{Z(\lambda_i, \nu)}$$

$$= \prod_{i=1}^{r} (\lambda_i/\lambda)^{k_i} \binom{m}{\boldsymbol{k}}^\nu \bigg/ G(\boldsymbol{\lambda}/\lambda, \nu) \qquad \text{for } \boldsymbol{k} \in D$$

which is the probability mass function of the COMM distribution with parameters $\boldsymbol{p} = \boldsymbol{\lambda}/\lambda$ and $\nu$.                                                                        □

Proposition 1 implies that there is a constructive way to generate Conway–Maxwell-Multinomial distributions for $\nu > 0$. However, the case for unrestricted $\nu$ is not so obvious. The next proposition shows a constructive way of generating (not-necessarily independent) multivariate discrete random variables that have an arbitrary sum.

**Proposition 2.** *Let $\boldsymbol{X}_1, \ldots, \boldsymbol{X}_n$ be $r$-dimensional indicator random variables, i.e., each is an $r$-vector of 0's and 1's. Let $P\{\boldsymbol{S} = \boldsymbol{k}\} = p_{\boldsymbol{k}} \geq 0$, where $\sum_{\boldsymbol{k} \in D} p_{\boldsymbol{k}} = 1$. Then there exists a unique distribution on $\boldsymbol{X}_1, \boldsymbol{X}_2, \ldots, \boldsymbol{X}_m$ such that $\boldsymbol{X}_1, \boldsymbol{X}_2, \ldots, \boldsymbol{X}_m$ are order $m$ exchangeable and $\sum_{i=1}^{m} \boldsymbol{X}_i$ has the same distribution as does $\boldsymbol{S}$.*

**Proof.** For each $\boldsymbol{k} \in D$, there are $\binom{m}{\boldsymbol{k}}$ different arrangements of 1's and 0's such that each vector component $i$ has $k_i$ 1's and $(m - k_i)$ 0's. Let each such arrangement have probability $p_{\boldsymbol{k}}/\binom{m}{\boldsymbol{k}}$. Then $P(\sum_{i=1}^{r} \boldsymbol{X}_i = \boldsymbol{k}) = p_{\boldsymbol{k}}$ and the $\boldsymbol{X}_i$'s are exchangeable of order $m$. To show uniqueness, if the sum of the probabilities of the sequences with $k_i$ 1's in the $i$th vector component for each $i$ were not $p_{\boldsymbol{k}}$, the sum constraint would not be met. If they did not have equal probability, order $m$ exchangeability would be violated.                                                                        □

It is also useful to display the COMM distribution as a member of the exponential family, and in particular to show its sufficient statistics. This is done in Proposition 3.

**Proposition 3.** *The COMM distribution has the following sufficient statistics:*

$$S_0 = \sum_{j=1}^{n} \log[k_{ij}! \cdots k_{rj}!],$$

$$S_i = \sum_{j=1}^{n} k_{ij}, \qquad i = 1, \ldots, r - 1,$$

*where $k_{ij}$ is the $i$th component of the $j$th sample. The COMM distribution is a member of the exponential family.*

**Proof.**

$$p(\boldsymbol{k}_1, \ldots, \boldsymbol{k}_n \mid \boldsymbol{p}, \nu) = \prod_{j=1}^{n} \left[ \binom{m}{\boldsymbol{k}_j}^{\nu} \prod_{i=1}^{r} p_i^{k_{ij}} / G(\boldsymbol{p}, \nu) \right]$$

$$\propto p_r^{nm} (m!)^{\nu n} \prod_{j=1}^{n} \left[ \prod_{i=1}^{r-1} (p_i/p_r)^{k_{ij}} \right] \left[ \prod_{i=1}^{r} k_{ij}! \right]^{-\nu}$$

$$\propto \exp\left(\sum_{i=1}^{r-1}\log(p_i/p_r)\sum_{j=1}^{n}k_{ij} - v\sum_{j=1}^{n}\log\left(\prod_{i=1}^{r}k_{ij}!\right)\right)$$

$$= e^{\sum_{i=1}^{r-1}\log(p_i/p_r)S_i - vS_0}.$$

This shows the sufficient statistics, and also shows that COMM is a member of the exponential family. □

Because the COMM distribution is in the exponential family, it has a conjugate distribution. However, when that distribution is proper is not so obvious, and is the content of Theorem 1 below.

Let $\boldsymbol{\psi} = (\log(p_1/p_r), \ldots, \log(p_{r-1}/p_r))$ and $t(\boldsymbol{k}) = -\log(\prod_{i=1}^{r}k_i!)$, then the COMM probability mass function can be re-expressed as

$$P(X = \boldsymbol{k}|\boldsymbol{\psi}, v) = e^{\boldsymbol{\psi}\cdot(k_1,k_2,\ldots,k_{r-1})+vt(\boldsymbol{k})-M(\boldsymbol{\theta})},$$

where, $\boldsymbol{\theta} = (\boldsymbol{\psi}, v)$, and $M(\boldsymbol{\theta}) = \log(\sum_{\boldsymbol{j}\in D} e^{\boldsymbol{\psi}\cdot(j_1,j_2,\ldots,j_{r-1})+vt(\boldsymbol{j})})$.

Consider a conjugate family of the form

$$\pi(\boldsymbol{\theta}|\boldsymbol{a}, b, c) \propto e^{\boldsymbol{\psi}\cdot\boldsymbol{a}+bv-cM(\boldsymbol{\theta})}, \tag{1}$$

where $\boldsymbol{a} = (a_1, a_2, \ldots, a_{r-1})$.

Then the updating of the hyper parameters using the sufficient statistics can be accomplished using

$$a_i' = a_i + S_i, \qquad i = 1, \ldots, r-1,$$
$$b' = b - S_0,$$
$$c' = c + n.$$

Let $a_r$ be a pseudo hyperparameter such that $\frac{a_r}{c} = m - \sum_{i=1}^{r-1}\frac{a_i}{c}$, then the conjugate prior $\pi$ can be written in a symmetric form with respect to $(\boldsymbol{p}, v)$.

The Jacobian matrix of transformation is $J = \frac{\partial\boldsymbol{\theta}}{\partial\boldsymbol{\rho}}$, $\boldsymbol{\rho} = (p_1, p_2, \ldots, p_{r-1}, v)$, with determinant $|J| = (\prod_{i=1}^{r}p_i)^{-1}$, using the matrix determinant lemma. Therefore, we have

$$h(\boldsymbol{p}, v|\boldsymbol{a}, a_r, b, c)$$

$$= \pi(\boldsymbol{\theta}|\boldsymbol{a}, b, c)\cdot\|J\|$$

$$\propto e^{\boldsymbol{\psi}\cdot\boldsymbol{a}+bv-cM(\boldsymbol{\theta})}\left(\prod_{i=1}^{r}p_i\right)^{-1}$$

$$= \left[\sum_{\boldsymbol{j}\in D}\frac{1}{(\prod_{i=1}^{r}j_i!)^v}p_1^{-\frac{a_1}{c}+j_1}\cdots p_{r-1}^{-\frac{a_{r-1}}{c}+j_{r-1}}p_r^{-\frac{a_r}{c}+j_r}e^{-\frac{b}{c}v}\right]^{-c}$$

$$\times\left(\prod_{i=1}^{r}p_i\right)^{-1}. \tag{2}$$

Note that when $\nu \neq 1$ the conjugate prior of this form depends on $m$, which can be regarded as the total number of trials. When $\nu = 1$, the COMM reduces to the multinomial distribution with conjugate prior the Dirichlet distribution, which does not depend on $m$.

**Theorem 1.** *The conjugate prior of COMM distribution in* (2) *is proper if and only if*

(i)   $c > 0$,

(ii)   $a_i > 0$,      $i = 1, 2, \ldots, r$,      $\displaystyle\sum_{i=1}^{r} \frac{a_i}{c} = m$,

(iii)   $\displaystyle -\log(m!) < \frac{b}{c} < -\sum_{i=1}^{r}\left[\left(\frac{a_i}{c} - \left\lfloor \frac{a_i}{c} \right\rfloor\right)\log\left\lceil \frac{a_i}{c} \right\rceil + \log\left\lfloor \frac{a_i}{c} \right\rfloor!\right].$

The proof of Theorem 1 is in the Appendix.

Sometimes conjugate prior distributions are not the most convenient. For example, the snippet of Mendel data is treated with a prior that is not conjugate. Theorem 2 shows that a finite moment generating function suffices.

**Theorem 2.** *Let* $g(\boldsymbol{\psi}, \nu)$ *be a probability distribution that has a finite moment generating function. Then the prior proportional to*

$$\exp(\boldsymbol{\psi} \cdot a + \nu b - cM(\boldsymbol{\psi}, \nu))g(\boldsymbol{\psi}, \nu)$$

*is proper for all* $a$, $b$, *and* $c$.

The proof of Theorem 2 is in the Appendix.

The generating functions are straightforward, and worth recording as well:

Let $\boldsymbol{t} = (t_1, \ldots, t_r)$. Then the probability generating function is

$$E(\boldsymbol{t}^{\boldsymbol{x}}) = \sum_{\boldsymbol{j} \in D} \boldsymbol{t}^{\boldsymbol{j}} \boldsymbol{p}^{\boldsymbol{j}} \binom{m}{\boldsymbol{j}}^{\nu} / G(\boldsymbol{p}, \nu)$$

$$= G(\boldsymbol{t}\boldsymbol{p}, \nu) / G(\boldsymbol{p}, \nu).$$

Similarly, the moment generating and characteristic functions are, respectively,

$$G(e^{\boldsymbol{t}} \boldsymbol{p}, \nu) / G(\boldsymbol{p}, \nu)$$

and

$$G(e^{i\boldsymbol{t}} \boldsymbol{p}, \nu) / G(\boldsymbol{p}, \nu).$$

## 4 Conclusion

The results of Kadane (2016) extend to the multivariate case (except his Theorem 2). The hard work is in extending Theorem 1 to the multivariate case. The Conway–Maxwell Multinomial distribution can be part of a statistician's toolkit.

## Appendix

Diaconis and Ylvisaker (1979) showed that for exponential family $\{P(\boldsymbol{\theta}) : dP(\boldsymbol{\theta}) = e^{\boldsymbol{x}\boldsymbol{\theta}-M(\boldsymbol{\theta})} d\mu(\boldsymbol{x}), \theta \in \Theta, \boldsymbol{x} \in \mathbb{R}^d\}$ where $\mu$ a is $\sigma$-finite measure on $\mathcal{B}(\mathbb{R}^d)$ and $\Theta = \{\theta : M(\theta) < \infty\}$, its conjugate prior family $\{\pi_{n_0,\boldsymbol{x}_0}(\boldsymbol{\theta}) : d\pi_{n_0,\boldsymbol{x}_0}(\boldsymbol{\theta}) = e^{n_0\boldsymbol{x}_0-n_0 M(\boldsymbol{\theta})} d\boldsymbol{\theta}\}$ has the following property: if $\Theta = \mathbb{R}^d$

$$\pi_{n_0,\boldsymbol{x}_0}(\Theta) < \infty \quad \Leftrightarrow \quad \boldsymbol{x}_0 \in \mathcal{X}, \qquad n_0 > 0,$$

where $\mathcal{X}$ is the interior of the convex hull of $\mu$'s support. For our problem, $\mathcal{X}$ is the interior of the convex hull of $\tilde{A} = \{(k_1, k_2, \ldots k_{r-1}, t(\boldsymbol{k}) : \boldsymbol{k} \in D\}$. In order to take advantage of the simplicity of symmetric forms, we first add the dimension $k_r$ and consider $A = \{(\boldsymbol{k}, t(\boldsymbol{k})) : \boldsymbol{k} \in D\}$ in the lemmas to follow, and then relate $A$ to $\tilde{A}$ in the proof of Theorem 1.

Here are some notations: $L = \{(\boldsymbol{z}, h) : \sum_{i=1}^r z_i = m\}$ is a subspace of $\mathbb{R}^{r+1}$ and thus we have $A \subseteq \text{Conv}(A) \subseteq L$. Let $H$ denote the interior of the convex hull of $A$ under the subspace topology, then its closure $\bar{H} = \text{Conv}(A)$. Let $\Omega = \{\boldsymbol{x} \in \mathbb{N}^r : m - r + 1 \le \sum_{i=1}^r x_i \le m - 1\}$. If $\boldsymbol{x} \in \Omega$, then the following set $C_{\boldsymbol{x}}$ generated by $\boldsymbol{x}$ is not empty,

$$C_{\boldsymbol{x}} = \left\{ \boldsymbol{k} : k_j = x_j \text{ or } x_j + 1, \sum_{j=1}^r k_j = m \right\}.$$

It is clear that $|C_{\boldsymbol{x}}| \ge \binom{r}{1} = r$ and $C_{\boldsymbol{x}} \subseteq D$.

The geometric intuition behind our proofs is that we first divide the "surface" of the convex hull into pieces of facets, then study the hyperplanes that contain these facets. The general steps of proof are:

1. Show that $\{\text{Conv}(C_{\boldsymbol{x}}) : \boldsymbol{x} \in \Omega\}$ exhausts $\text{Conv}(D)$

2. Show that for every $\boldsymbol{z} \in \text{Conv}(C_{\boldsymbol{x}})$, $(\boldsymbol{z}, h) \in \bar{H}$ if and only if $a(\boldsymbol{z}) \le h \le b_{\boldsymbol{x}}(\boldsymbol{z})$. The graphs of $a(\boldsymbol{z})$ and $b_{\boldsymbol{x}}(\boldsymbol{z})$ are, respectively, the bottom facet and a dominant facet indexed by $\boldsymbol{x}$.

3. Relate that above characterization of $\bar{H}$ to the necessary and sufficient conditions for the priors to be proper.

**Lemma 1.** $\{\text{Conv}(C_{\boldsymbol{x}}) : \boldsymbol{x} \in \Omega\}$ *exhausts* $\text{Conv}(D)$, *that is,* $\bigcup_{\boldsymbol{x} \in \Omega} \text{Conv}(C_{\boldsymbol{x}}) = \text{Conv}(D)$.

**Proof.** Suffices to show that $\text{Conv}(D) \subseteq \bigcup_{x \in \Omega} \text{Conv}(C_x)$. If $z \in \text{Conv}(D)$, then $\sum_{i=1}^r z_i = m$. We need to find $x \in \Omega$ such that $z \in \text{Conv}(C_x)$.

(1) $\sum_{j=1}^r \lfloor z_j \rfloor \leq m - 1$.
Let $x = (\lfloor z_1 \rfloor, \ldots, \lfloor z_r \rfloor)$, then $x \in \Omega$ as $m - 1 \geq \sum_{j=1}^r \lfloor z_j \rfloor > \sum_{j=1}^r (z_j - 1) = m - r$. We will show that $\text{Conv}(C_x) = \{z | x_j \leq z_j \leq x_j + 1, \sum_{j=1}^r z_j = m\} := R_x$.

For all $z_1, z_2 \in R_x$ and $\lambda \in [0, 1]$, we have $\lambda z_1 + (1 - \lambda)z_2 \in R_x$, therefore, $R_x$ is a convex set. $C_x \subseteq R_x \Rightarrow \text{Conv}(C_x) \subseteq R_x$.

Next, we prove that elements in $R_x \setminus C_x$ can not be vertices of $R_x$. Assume $\tilde{z} \in R_x \setminus C_x$ and is a vertex of $R_x$. Without loss of generality, let $\tilde{z}_1, \tilde{z}_2$ be such that $\tilde{z}_1 \in (x_1, x_1 + 1), \tilde{z}_2 \in (x_2, x_2 + 1)$. For $\epsilon > 0$ sufficiently small, both $z_\epsilon^+ = (\tilde{z}_1 + \epsilon, \tilde{z}_2 - \epsilon, \ldots, \tilde{z}_r)$ and $z_\epsilon^- = (\tilde{z}_1 - \epsilon, \tilde{z}_2 + \epsilon, \ldots, \tilde{z}_r)$ are elements of $R_x$. Consequently, $\tilde{z} = \frac{1}{2}(z_\epsilon^+ + z_\epsilon^-)$, which contradicts with the assumption that $\tilde{z}$ is a vertex of $R_x$. Hence, $R_x \subseteq \text{Conv}(C_x)$.

As a result, $z \in R_x = \text{Conv}(C_x)$.

(2) $\sum_{j=1}^r \lfloor z_j \rfloor = m$.
We know $\sum_{j=1}^r z_j = m$, so $z_1, z_2, \ldots, z_r$ must be integers. As $m \geq 1$, there exists $z_i \geq 1$ for some $i$. Let $x = (z_1, \ldots, z_{i-1}, z_i - 1, z_{i+1}, \ldots, z_r) \in \Omega$, then $z \in C_x \subseteq \text{Conv}(C_x)$. $\qquad\square$

The following proposition is a step towards the proof of Lemma 1.

**Proposition 4.** *For every $x \in \Omega$ and $z \in D$,*

$$t(z) \leq b_x(z) = t(x) - \sum_{i=1}^r (z_i - x_i) \log(x_i + 1).$$

*Equality holds if and only if $z \in C_x$.*

**Proof.** Let $g(x) = -\log(\Gamma(x + 1))$, then $-\log(x_i + 1) = g(x_i + 1) - g(x_i)$, $1 \leq i \leq r$. Also, $t(x) = -\log(\prod_{i=1}^r \Gamma(x_i + 1)) = -\sum_{i=1}^r \log(\Gamma(x_i + 1)) = \sum_{i=1}^r g(x_i)$. Note that $g$ is strictly concave, so $\forall x \in \Omega, z \in D$ and $i \in 1, \ldots, r$

$$(z_i - x_i)\big(g(x_i + 1) - g(x_i)\big) \geq g(z_i) - g(x_i).$$

Equality holds if and only if $z_i = x_i$ or $x_i + 1$. Therefore,

$$\sum_{i=1}^r \big(g(x_i) + (z_i - x_i)(g(x_i + 1) - g(x_i))\big) \geq \sum_{i=1}^r g(z_i).$$

Consequently,

$$t(x) + \sum_{i=1}^r (z_i - x_i)\big(g(x_i + 1) - g(x_i)\big) \geq t(z).$$

Equality holds if and only if

$$
\begin{cases}
z_i = x_i \quad \text{or} \quad z_i + 1, \qquad i = 1, 2, \ldots, r, \\
\sum_{i=1}^{r} z_i = m
\end{cases}
\quad \Leftrightarrow \quad z \in C_x.
$$

$\square$

**Lemma 2.** *For every* $z \in \text{Conv}(C_x)$, $(z, h) \in \bar{H}$ *if and only if* $-\log(m!) \leq h \leq b_x(z)$.

**Proof.** Assume $(q_1, t(q_1)), \ldots, (q_{|A|}, t(q_{|A|}))$ are the elements in $A$.

(i) $\forall k \in D$, we have $\binom{m}{k} \geq 1$. Hence, $t(k) = \log \binom{m}{k} - \log(m!) \geq -\log(m!)$. If $(z, h) \in \bar{H} = \text{Conv}(A)$, then there exists $\alpha_j \geq 0$, $\sum_{j=1}^{|A|} \alpha_j = 1$ such that $(z, h) = \sum_{j=1}^{|A|} \alpha_j(q_j, t(q_j))$. We have $h = \sum_{j=1}^{|A|} \alpha_j t(q_j) \geq -\log(m!) \sum_{j=1}^{|A|} \alpha_j = -\log(m!)$. Again, by Proposition 4 and linearity of $b_x(z)$, we get $b_x(z) = b_x(\sum_{j=1}^{|A|} \alpha_j q_j) = \sum_{j=1}^{|A|} \alpha_j b_x(q_j) \geq \sum_{j=1}^{|A|} \alpha_j t(q_j) = h$.

(ii) Here we show the "only if" part of the proof. Suffices to show that $B = \{(z, h) : z \in \text{Conv}(D), h = -\log(m!)\} \subseteq \bar{H}$ and $S_x = \{(z, h) : z \in \text{Conv}(C_x), h = b_x(z)\} \subseteq \bar{H}$. Let $e_i$ denote the unit vector with a 1 in the $i$th coordinate and 0 elsewhere. Now consider $me_i \in D, 1 \leq i \leq r$. Note that $t(me_i) = -\log(m!)$ and $\text{Conv}(me_i : 1 \leq i \leq r) = \text{Conv}(D)$, thus we have $B \subseteq \bar{H}$. From Proposition 4, we know for every $z \in C_x, t(z) = b_x(z)$. Again, $z$ can be written as $\sum_j \lambda_j k_j$, where $k_j \in C_x$, $\lambda_j \geq 0$ and $\sum_j \lambda_j = 1$. Hence, $b_x(z) = b_x(\sum_j \lambda_j k_j) = \sum_j \lambda_j b_x(k_j) = \sum_j \lambda_j t(k_j)$, which implies $S_x \subseteq \text{Conv}((k_j, t(k_j)) : k_j \in C_x) \subseteq \bar{H}$. $\square$

Using these results, we now resume the proof of Theorem 1.

**Proof of Theorem 1.** Because there are only finite number of elements in $D$, the parameter space $\Theta = \{\theta | M(\theta) < \infty\} = \mathbb{R}^r$. Let $\tilde{H}$ denote the interior of the convex hull of $\tilde{A}$. The conjugate prior is proper, according to Diaconis and Ylivisacker's theorem, if and only if

$$
\text{(i)} \quad c > 0, \qquad \text{(ii)} \quad \left( \frac{a_1}{c}, \frac{a_2}{c}, \ldots, \frac{a_{r-1}}{c}, \frac{b}{c} \right) \in \tilde{H}. \tag{A.1}
$$

The following mapping $T$ between $\tilde{A}$ and $A$ is bijective,

$$
T : \tilde{A} \to A,
$$

$$
\tilde{q}_i = \left( k_1^{(i)}, \ldots, k_{r-1}^{(i)}, t(k^{(i)}) \right) \mapsto q_i = \left( k^{(i)}, t(k^{(i)}) \right).
$$

We can therefore convert the problem to a symmetric one by adding one more dimension. As $T$ is bijective, we have $|\tilde{A}| = |A| := K$.

$$(\tau_1, \tau_2, \ldots, \tau_r) \in \text{Conv}(\tilde{A})$$

$$\Leftrightarrow \quad \exists \alpha_1, \alpha_2, \ldots, \alpha_K \geq 0, \qquad \sum_{j=1}^{K} \alpha_j = 1$$

$$\text{s.t. } (\tau_1, \tau_2, \ldots, \tau_r) = \sum_{j=1}^{K} \alpha_j \tilde{\boldsymbol{q}}_j$$

$$\Leftrightarrow \quad \exists \alpha_1, \alpha_2, \ldots, \alpha_K \geq 0, \qquad \sum_{j=1}^{K} \alpha_j = 1$$

$$\text{s.t. } \left( \tau_1, \tau_2, \ldots, \tau_{r-1}, m - \sum_{j=1}^{r-1} \tau_j, \tau_r \right) = \sum_{j=1}^{K} \alpha_j \boldsymbol{q}_j$$

$$\Leftrightarrow \quad \left( \tau_1, \tau_2, \ldots, \tau_{r-1}, m - \sum_{j=1}^{r-1} \tau_j, \tau_r \right) \in \text{Conv}(A).$$

Hence,

$$(\tau_1, \tau_2, \ldots, \tau_r) \in \tilde{H}$$

$$\Leftrightarrow \quad \exists \tilde{\varepsilon} > 0, \qquad \text{s.t. } B_{\tilde{\varepsilon}}(\tau_1, \tau_2, \ldots, \tau_r) \subseteq \text{Conv}(\tilde{A})$$

$$\Leftrightarrow \quad \exists \varepsilon > 0, \qquad \text{s.t. } L \cap B_{\varepsilon}\left( \tau_1, \ldots, \tau_{r-1}, m - \sum_{j=1}^{r-1} \tau_j, \tau_r \right) \subseteq \text{Conv}(A)$$

$$\Leftrightarrow \quad \left( \tau_1, \ldots, \tau_{r-1}, m - \sum_{j=1}^{r-1} \tau_j, \tau_r \right) \in H,$$

where $B_{\tilde{\varepsilon}}(\tau_1, \tau_2, \ldots, \tau_r)$ is an open $r$ dimensional sphere centered at $(\tau_1, \tau_2, \ldots, \tau_r)$ with radius $\tilde{\varepsilon}$ and $B_{\varepsilon}(\tau_1, \ldots, \tau_{r-1}, m - \sum_{j=1}^{r-1} \tau_j, \tau_r)$ is an open $r + 1$ dimensional sphere centered at $(\tau_1, \ldots, \tau_{r-1}, m - \sum_{j=1}^{r-1} \tau_j, \tau_r)$ with radius $\varepsilon$. As a result,

$$\left( \frac{a_1}{c}, \frac{a_2}{c}, \ldots, \frac{a_{r-1}}{c}, \frac{b}{c} \right) \in \tilde{H} \quad \Leftrightarrow \quad \left( \frac{a_1}{c}, \frac{a_2}{c}, \ldots, \frac{a_r}{c}, \frac{b}{c} \right) \in H. \qquad (A.2)$$

Let $F(z) = \sum_{i=1}^{r} [(z_i - \lfloor z_i \rfloor) \log \lceil z_i \rceil + \log \lfloor z_i \rfloor!]$. We claim that

$$\bar{H} = \left\{ (z, h) : -\log(m!) \leq h \leq F(z), z_j \geq 0, \sum_{j=1}^{r} z_j = m \right\}.$$

To show $\bar{H} \subseteq RHS$, assume $(z, h) \in \bar{H} = \text{Conv}(A)$, then $z \in \text{Conv}(D) = \{z : \sum_{j=1}^{r} z_j = m, z_j \geq 0\}$. According to Lemma 1, there exists $x \in \Omega$ such that $z \in \text{Conv}(C_x)$. Apply Lemma 2, we get $-\log(m!) \leq h \leq b_x(z)$. Now, choose $x$ as suggested by the proof of Lemma 1. If $\sum_{j=1}^{r} \lfloor z_j \rfloor \leq m - 1$, let $x = (\lfloor z_1 \rfloor, \ldots, \lfloor z_r \rfloor)$, then we have $b_x(z) = F(z)$. Otherwise, $\sum_{j=1}^{r} \lfloor z_j \rfloor = m$, we have $z \in C_x$ for every possible choice of $x$, so $b_x(z) = t(z) = F(z)$ holds.

Conversely, we show that $RHS \subseteq \bar{H}$. If $\sum_{j=1}^{r} z_j = m$ and $z_j \geq 0$ for $1 \leq j \leq r$, then $z \in \text{Conv}(D)$ which implies $z \in \text{Conv}(C_x)$ for some $x$ by Lemma 1. We have seen that $F(z) = b_x(z)$. Again, use Lemma 2, we get $(z, h) \in \bar{H}$.

Consequently,

$$H = \left\{ (z, h) \,\middle|\, -\log(m!) < h < F(z), z_j > 0, \sum_{j=1}^{r} z_j = m \right\}. \qquad \text{(A.3)}$$

Now, let $z = (\frac{a_1}{c}, \ldots, \frac{a_r}{c})$, we have

$$F\left(\frac{a_1}{c}, \frac{a_2}{c}, \ldots, \frac{a_r}{c}\right) = -\sum_{i=1}^{r} \left[ \left(\frac{a_i}{c} - \left\lfloor \frac{a_i}{c} \right\rfloor\right) \log \left\lceil \frac{a_i}{c} \right\rceil + \log \left\lfloor \frac{a_i}{c} \right\rfloor! \right]. \qquad \text{(A.4)}$$

By combining equation (A.1), (A.2), (A.3) and (A.4), we get Theorem 1. $\qquad \square$

**Proof of Theorem 2.** Let $h_j(\theta) = \psi \cdot (j_1, j_2, \ldots, j_{r-1}) + vt(j)$, which is linear in $\theta$. Let $R_j = \{\theta \mid \text{sgn}(c)h_j(\theta) \text{ is minimized over } j \in D\}$.

The $R_j$'s are exhaustive and measurable. Let $D^* = \{j \mid R_j(\theta) \neq \varnothing\}$. $D^*$ is not empty.

Choose $j^* \in D^*$. For $\theta \in R_{j^*}$,

$$\text{sgn}(c) \sum_{j \in D} \exp(h_j(\theta)) \geq \text{sgn}(c) |D| \exp(h_{j^*}(\theta)),$$

so

$$\text{sgn}(c) M(\theta) \geq \text{sgn}(c) K h_{j^*}(\theta),$$

where $K = \log(|D|) > 0$.

Hence, $\int_{R_{j^*}} \exp(\psi \cdot a + vb - cM(\theta)) g(\theta) \, d\theta \leq \int \exp(\psi \cdot a + vb - cK h_{j^*}(\theta)) \times g(\theta) \, d\theta < \infty$ because $\psi \cdot a + vb - cK h_{j^*}(\theta)$ is linear in $\theta$ and $g$ has a finite moment generating function. It follows that $\int \exp(\psi \cdot a + vb - cM(\theta)) g(\theta) \, d\theta \leq \sum_{j^* \in D^*} \int_{R_{j^*}} \exp(\psi \cdot a + vb - cM(\theta)) g(\theta) \, d\theta < \infty$. $\qquad \square$

## Acknowledgment

# References

Boatwright, P., Borle, S. and Kadane, J. B. (2003). A model of the joint distribution of purchase quantity and timing. *J. Amer. Statist. Assoc.* **98**, 564–572. MR2011672

Borle, S., Boatwright, P., Kadane, J. B., Nunes, J. C. and Shmueli, G. (2005). The effect of product assortment changes on customer retention. *Mark. Sci.* **4**, 616–622.

Conway, R. and Maxwell, W. (1962). A queing model with state dependent service rates. *J. Ind. Eng.* **12**, 132–136.

Diaconis, P. and Ylvisaker, D. (1979). Conjugate priors for exponential families. *Ann. Statist.* **7**, 269–281. MR0520238

Fisher, R. A. (1936). Has Mendel's work been rediscovered? *Ann. of Sci.* **1**, 115–137.

Fisher, R. A. (1959). *Statistical Methods and Scientific Inference*, 2nd ed. Edinburgh: Oliver and Boyd.

Franklin, A., Edwards, A. W. F., Fairbanks, D. J., Hartl, D. L. and Seidenfeld, T. (2008). *Ending the Mendel–Fisher Controversy*. Pittsburgh: University of Pittsburgh Press.

Kadane, J. B. (2016). Sums of possibly associated Bernoulli variables: The Conway–Maxwell-Binomial distribution. *Bayesian Anal.* **11**, 403–420.

Pires, A. M. and Branco, J. A. (2010). A statistical model to explain the Mendel–Fisher controversy. *Statist. Sci.* **25**, 545–565. MR2807770

Shmueli, G., Minka, T. P., Kadane, J. B., Borle, S. and Boatwright, P. (2004). A useful distribution for fitting discrete data: Revival of the COM-Poisson. *J. R. Stat. Soc., Ser. C, Appl. Stat.* **54**, 127–142. MR2134602

Weldon, W. R. F. (1902). Mendel's law of alternative inference in peas. *Biometrika* **1**, 228–254.

Department of Statistics
Carnegie Mellon University
Pittsburgh, Pennsylvania 15213
USA
E-mail: kadane@stat.cmu.edu

Department of Statistics
Columbia University
New York, New York 10027
USA
E-mail: zw2393@columbia.edu