# Estimating a Dirichlet distribution
## Thomas P. Minka
2000 (revised 2003, 2009)

**Abstract**

The Dirichlet distribution and its compound variant, the Dirichlet-multinomial, are two of the most basic models for proportional data, such as the mix of vocabulary words in a text document. Yet the maximum-likelihood estimate of these distributions is not available in closed-form. This paper describes simple and efficient iterative schemes for obtaining parameter estimates in these models. In each case, a fixed-point iteration and a Newton-Raphson (or generalized Newton-Raphson) iteration is provided.

## 1 The Dirichlet distribution

The Dirichlet distribution is a model of how proportions vary. Let $\mathbf{p}$ denote a random vector whose elements sum to 1, so that $p_k$ represents the proportion of item $k$. Under the Dirichlet model with parameter vector $\boldsymbol{\alpha}$, the probability density at $\mathbf{p}$ is

$$p(\mathbf{p}) \quad \sim \quad \mathcal{D}(\alpha_1, ..., \alpha_K) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k p_k^{\alpha_k - 1} \tag{1}$$

$$\text{where } p_k \quad > \quad 0 \tag{2}$$

$$\sum_k p_k \quad = \quad 1 \tag{3}$$

The parameters $\boldsymbol{\alpha}$ can be estimated from a training set of proportions: $D = \{\mathbf{p}_1, ..., \mathbf{p}_N\}$. The maximum-likelihood estimate of $\boldsymbol{\alpha}$ maximizes $p(D|\boldsymbol{\alpha}) = \prod_i p(\mathbf{p}_i|\boldsymbol{\alpha})$. The log-likelihood can be written

$$\log p(D|\boldsymbol{\alpha}) \quad = \quad N \log \Gamma(\sum_k \alpha_k) - N \sum_k \log \Gamma(\alpha_k) + N \sum_k (\alpha_k - 1) \log \bar{p}_k \tag{4}$$

$$\text{where } \log \bar{p}_k \quad = \quad \frac{1}{N} \sum_i \log p_{ik} \tag{5}$$

This objective is convex in $\boldsymbol{\alpha}$ since Dirichlet is in the exponential family. This implies that the likelihood is unimodal and the maximum can be found by a simple search. A direct convexity proof has also been given by Ronning (1989). The gradient of the log-likelihood with respect to one $\alpha_k$ is

$$g_k \quad = \quad \frac{d \log p(D|\boldsymbol{\alpha})}{d\alpha_k} = N\Psi(\sum_k \alpha_k) - N\Psi(\alpha_k) + N \log \bar{p}_k \tag{6}$$

$$\Psi(x) \quad = \quad \frac{d \log \Gamma(x)}{dx} \tag{7}$$

$\Psi$ is known as the digamma function and is similar to the natural logarithm. As always with the exponential family, when the gradient is zero, the expected sufficient statistics are equal to the observed sufficient statistics. In this case, the expected sufficient statistics are

$$E[\log p_k] = \Psi(\alpha_k) - \Psi(\sum_k \alpha_k) \tag{8}$$

and the observed sufficient statistics are $\log \bar{p}_k$.

A fixed-point iteration for maximizing the likelihood can be derived as follows. Given an initial guess for $\boldsymbol{\alpha}$, we construct a simple lower bound on the likelihood which is tight at $\boldsymbol{\alpha}$. The maximum of this bound is computed in closed-form and it becomes the new guess. Such an iteration is guaranteed to converge to a stationary point of the likelihood—in fact it is the same principle behind the EM algorithm (Minka, 1998). For the Dirichlet, the maximum is the only stationary point.

As shown in appendix A, a bound on $\Gamma(\sum_k \alpha_k)$ leads to the following fixed-point iteration:

$$\Psi(\alpha_k^{new}) = \Psi(\sum_k \alpha_k^{old}) + \log \bar{p}_k \tag{9}$$

This algorithm requires inverting the $\Psi$ function—a procedure which is described in appendix C.

Another approach to finding a stationary point is Newton iteration. The second-derivatives, i.e. Hessian matrix, of the log-likelihood are given by

$$\frac{d \log p(D|\alpha)}{d\alpha_k^2} = N\Psi'(\sum_k \alpha_k) - N\Psi'(\alpha_k) \tag{10}$$

$$\frac{d \log p(D|\alpha)}{d\alpha_k d\alpha_j} = N\Psi'(\sum_k \alpha_k) \qquad (k \neq j) \tag{11}$$

$\Psi'$ is known as the trigamma function. The Hessian can be written in matrix form as

$$\mathbf{H} = \mathbf{Q} + \mathbf{1}\mathbf{1}^{\mathrm{T}}z \tag{12}$$

$$q_{jk} = -N\Psi'(\alpha_k)\delta(j-k) \tag{13}$$

$$z = N\Psi'(\sum_k \alpha_k) \tag{14}$$

One Newton step is therefore

$$\alpha^{\mathrm{new}} = \alpha^{\mathrm{old}} - \mathbf{H}^{-1}\mathbf{g} \tag{15}$$

$$\mathbf{H}^{-1} = \mathbf{Q}^{-1} - \frac{\mathbf{Q}^{-1}\mathbf{1}\mathbf{1}^{\mathrm{T}}\mathbf{Q}^{-1}}{1/z + \mathbf{1}^{\mathrm{T}}\mathbf{Q}^{-1}\mathbf{1}} \tag{16}$$

$$(\mathbf{H}^{-1}\mathbf{g})_k = \frac{g_k - b}{q_{kk}} \tag{17}$$

$$\text{where } b = \frac{\mathbf{1}^{\mathrm{T}}\mathbf{Q}^{-1}\mathbf{g}}{1/z + \mathbf{1}^{\mathrm{T}}\mathbf{Q}^{-1}\mathbf{1}} = \frac{\sum_j g_j/q_{jj}}{1/z + \sum_j 1/q_{jj}} \tag{18}$$

Unlike some Newton algorithms, this one does not require storing or inverting the Hessian matrix explicitly. The same Newton algorithm was given by Ronning (1989) and by Naryanan (1991). Naryanan also derives a stopping rule for the iteration.

An approximate MLE, useful for initialization, is given by finding the density which matches the moments of the data. The first two moments of the density are

$$E[p_k] = \frac{\alpha_k}{\sum_k \alpha_k} \tag{19}$$

$$E[p_k^2] = E[p_k]\frac{1 + \alpha_k}{1 + \sum_k \alpha_k} \tag{20}$$

$$\sum_k \alpha_k = \frac{E[p_1] - E[p_1^2]}{E[p_1^2] - E[p_1]^2} \tag{21}$$

2

Multiplying (21) and (19) gives a formula for $\alpha_k$ in terms of moments. Equation (21) uses $p_1$, but any other $p_k$ could also be used to estimate $\sum_k \alpha_k$. Ronning (1989) suggests instead using all of the $p_k$'s via

$$\text{var}(p_k) = \frac{E[p_k](1 - E[p_k])}{1 + \sum_k \alpha_k} \tag{22}$$

$$\log \sum_k \alpha_k = \frac{1}{K-1} \sum_{k=1}^{K-1} \log \left( \frac{E[p_k](1 - E[p_k])}{\text{var}(p_k)} - 1 \right) \tag{23}$$

Another approximate MLE, specifically for the case $K = 2$, is given by Johnson & Kotz (1970):

$$\alpha_1 = \frac{1}{2} \frac{1 - \bar{p}_2}{1 - \bar{p}_1 - \bar{p}_2} \tag{24}$$

$$\alpha_2 = \frac{1}{2} \frac{1 - \bar{p}_1}{1 - \bar{p}_1 - \bar{p}_2} \tag{25}$$

## 2  Estimating Dirichlet mean and precision separately

The $\boldsymbol{\alpha}$ parameters of the Dirichlet can be understood by considering the following alternative representation:

$$s = \sum_k \alpha_k \tag{26}$$

$$\mathbf{m} = E[p_k] = \boldsymbol{\alpha}/s \tag{27}$$

Here $\mathbf{m}$ is the mean of the distribution for $\mathbf{p}$ and $s$ can be understood as the *precision*. When $s$ is large, $\mathbf{p}$ is likely to be near $\mathbf{m}$, and when $s$ is small, $\mathbf{p}$ is distributed more diffusely. Interpretation of $s$ and $\mathbf{m}$ suggests situations, such as hierarchical modeling, in which we may want to fix one parameter and only optimize the other. Additionally, $s$ and $\mathbf{m}$ are roughly decoupled in the maximum-likelihood objective, which means we can get simplifications and speedups by optimizing them alternately. Thus, in this section we will reparameterize the distribution with $(s, \mathbf{m})$ where

$$\alpha_k = sm_k \tag{28}$$

$$\sum_k m_k = 1 \tag{29}$$

### 2.1  Estimating Dirichlet precision

The likelihood for $s$ alone is

$$p(D|s) \propto \left( \frac{\Gamma(s) \exp(s \sum_k m_k \log \bar{p}_k)}{\prod_k \Gamma(sm_k)} \right)^N \tag{30}$$

whose derivatives are

$$\frac{d \log p(D|s)}{ds} = N\Psi(s) - N \sum_k m_k \Psi(sm_k) + N \sum_k m_k \log \bar{p}_k \tag{31}$$

$$\frac{d^2 \log p(D|s)}{ds^2} = N\Psi'(s) - N \sum_k m_k^2 \Psi'(sm_k) \tag{32}$$

3

A convergent fixed-point iteration for $s$ is

$$1/s^{new} = 1/s - \Psi(s) + \sum_k m_k \Psi(sm_k) - \sum_k m_k \log \bar{p}_k \tag{33}$$

**Proof**  Use the bound

$$\frac{\Gamma(s)}{\prod_k \Gamma(sm_k)} \geq \exp(sb + \log(s) + c) \tag{34}$$

$$b = \Psi(\hat{s}) - \sum_k m_k \Psi(\hat{s}m_k) - 1/\hat{s} \tag{35}$$

to get

$$\log p(D|s) \geq s \sum_k m_k \log \bar{p}_k + sb + \log(s) + (\text{const.}) \tag{36}$$

from which (33) follows.

This iteration is only first-order convergent because the bound only matches the first derivative of the likelihood. We can derive a second-order method using the technique of generalized Newton iteration Minka (2000). The idea is to approximate the likelihood by a simpler function, by matching the first two derivatives at the current guess:

$$\frac{\Gamma(s)}{\prod_k \Gamma(sm_k)} \approx \exp(sb + a\log(s) + c) \tag{37}$$

$$a = -\hat{s}^2(\Psi'(\hat{s}) - \sum_k m_k^2 \Psi'(\hat{s}m_k)) \tag{38}$$

$$b = \Psi(\hat{s}) - \sum_k m_k \Psi(\hat{s}m_k) - a/\hat{s} \tag{39}$$

Maximizing the approximation leads to the update

$$\frac{1}{s^{new}} = \frac{1}{s} + \frac{1}{s^2}\left(\frac{d^2 \log p(D|s)}{ds^2}\right)^{-1}\left(\frac{d \log p(D|s)}{ds}\right) \tag{40}$$

This update resembles Newton-Raphson, but converges faster.

For initialization of $s$, it is useful to derive a closed-form approximate MLE. Stirling's approximation to $\Gamma$ gives

$$\frac{\Gamma(s)\exp(s\sum_k m_k \log \bar{p}_k)}{\prod_k \Gamma(sm_k)} \approx \left(\frac{s}{2\pi}\right)^{(k-1)/2}\prod_k m_k^{1/2}\exp(s\sum_k m_k \log \frac{\bar{p}_k}{m_k}) \tag{41}$$

$$\hat{s} \approx \frac{(k-1)/2}{-\sum_k m_k \log \frac{\bar{p}_k}{m_k}} \tag{42}$$

## 2.2   Estimating Dirichlet mean

Now suppose we fix the precision $s$ and want to estimate the mean $\mathbf{m}$. The likelihood for $\mathbf{m}$ alone is

$$p(D|\mathbf{m}) \propto \left(\prod_k \frac{\exp(sm_k \log \bar{p}_k)}{\Gamma(sm_k)}\right)^N \tag{43}$$

4

Reparameterize with the unconstrained vector $\mathbf{z}$, to get the gradient:

$$m_k = \frac{z_k}{\sum_k z_k} \tag{44}$$

$$\frac{d\log p(D|\mathbf{m})}{dz_k} = \frac{Ns}{\sum_k z_k}\left(\log\bar{p}_k - \Psi(sm_k) - \sum_k m_k\left(\log\bar{p}_k - \Psi(sm_k)\right)\right) \tag{45}$$

The MLE can be computed by the fixed-point iteration

$$\Psi(\alpha_k) = \log\bar{p}_k - \sum_k m_k^{old}\left(\log\bar{p}_k - \Psi(sm_k^{old})\right) \tag{46}$$

$$m_k^{new} = \frac{\alpha_k}{\sum_k \alpha_k} \tag{47}$$

This update converges very quickly.

# 3 The Dirichlet-multinomial/Polya distribution

The Dirichlet-multinomial distribution is a compound distribution where $\mathbf{p}$ is drawn from a Dirichlet and then a sample of discrete outcomes $\mathbf{x}$ is drawn from a multinomial with probability vector $\mathbf{p}$. This compounding is essentially a Polya urn scheme, so the Dirichlet-multinomial is also called the *Polya distribution*. Let $n_k$ be the number of times the outcome was $k$, i.e.

$$n_k = \sum_j \delta(x_j - k) \tag{48}$$

Then the resulting distribution over $\mathbf{x}$, a vector of outcomes, is

$$p(\mathbf{x}|\boldsymbol{\alpha}) = \int_{\mathbf{p}} p(\mathbf{x}|\mathbf{p})p(\mathbf{p}|\boldsymbol{\alpha})d\mathbf{p} \tag{49}$$

$$= \frac{\Gamma(\sum_k \alpha_k)}{\Gamma(\sum_k n_k + \alpha_k)}\prod_k \frac{\Gamma(n_k + \alpha_k)}{\Gamma(\alpha_k)} \tag{50}$$

This distribution is also parameterized by $\boldsymbol{\alpha}$, which can be estimated from a training set of count vectors: $D = \{\mathbf{x}_1, ..., \mathbf{x}_N\}$. The likelihood is

$$n_i = \sum_k n_{ik} \tag{51}$$

$$p(D|\alpha) = \prod_i p(\mathbf{x}_i|\alpha) \tag{52}$$

$$= \prod_i \left(\frac{\Gamma(\sum_k \alpha_k)}{\Gamma(n_i + \sum_k \alpha_k)}\prod_k \frac{\Gamma(n_{ik} + \alpha_k)}{\Gamma(\alpha_k)}\right) \tag{53}$$

The gradient of the log-likelihood is

$$g_k = \frac{d\log p(D|\alpha)}{d\alpha_k} = \sum_i \Psi(\sum_k \alpha_k) - \Psi(n_i + \sum_k \alpha_k) + \Psi(n_{ik} + \alpha_k) - \Psi(\alpha_k) \tag{54}$$

The maximum can be computed via the fixed-point iteration

$$\alpha_k^{new} = \alpha_k \frac{\sum_i \Psi(n_{ik} + \alpha_k) - \Psi(\alpha_k)}{\sum_i \Psi(n_i + \sum_k \alpha_k) - \Psi(\sum_k \alpha_k)} \tag{55}$$

5

(see appendix B).

Alternatively, there is a simplified Newton iteration as in the Dirichlet case. The Hessian of the log-likelihood is

$$
\frac{d \log p(D|\alpha)}{d\alpha_k^2} = \sum_i \Psi'(\sum_k \alpha_k) - \Psi'(n_i + \sum_k \alpha_k) + \Psi'(n_{ik} + \alpha_k) - \Psi'(\alpha_k) \tag{56}
$$

$$
\frac{d \log p(D|\alpha)}{d\alpha_k d\alpha_j} = \sum_i \Psi'(\sum_k \alpha_k) - \Psi'(n_i + \sum_k \alpha_k) \qquad (k \neq j) \tag{57}
$$

The Hessian can be written in matrix form as

$$
\mathbf{H} = \mathbf{Q} + \mathbf{1}\mathbf{1}^\mathrm{T} z \tag{58}
$$

$$
q_{jk} = \delta(j-k) \sum_i \Psi'(n_{ik} + \alpha_k) - \Psi'(\alpha_k) \tag{59}
$$

$$
z = \sum_i \Psi'(\sum_k \alpha_k) - \Psi'(n_i + \sum_k \alpha_k) \tag{60}
$$

from which a Newton step can be computed as before. The search can be initialized with the moment matching estimate where $p_{ik}$ is approximated by $n_{ik}/n_i$.

Another approach is to reduce this problem to the previous one via EM; see appendix D.

A different method is to maximize the leave-one-out (LOO) likelihood instead of the true likelihood. The LOO likelihood is the product of the probability of each sample given the remaining data and the parameters. The LOO log-likelihood is

$$
f(\alpha) = \sum_{ik} n_{ik} \log\left( \frac{n_{ik} - 1 + \alpha_k}{n_i - 1 + \sum_k \alpha_k} \right) = \sum_{ik} n_{ik} \log(n_{ik} - 1 + \alpha_k) - \sum_i n_i \log(n_i - 1 + \sum_k \alpha_k) \tag{61}
$$

Note that it doesn't involve any special functions. The derivatives are

$$
\frac{df(\alpha)}{d\alpha_k} = \sum_i \frac{n_{ik}}{n_{ik} - 1 + \alpha_k} - \frac{n_i}{n_i - 1 + \sum_k \alpha_k} \tag{62}
$$

$$
\frac{df(\alpha)}{d\alpha_k^2} = \sum_i -\frac{n_{ik}}{(n_{ik} - 1 + \alpha_k)^2} + \frac{n_i}{(n_i - 1 + \sum_k \alpha_k)^2} \tag{63}
$$

$$
\frac{df(\alpha)}{d\alpha_k \alpha_j} = \sum_i \frac{n_i}{(n_i - 1 + \sum_k \alpha_k)^2} \tag{64}
$$

A convergent fixed-point iteration is

$$
\alpha_k^{new} = \alpha_k \frac{\sum_i \frac{n_{ik}}{n_{ik}-1+\alpha_k}}{\sum_i \frac{n_i}{n_i-1+\sum_k \alpha_k}} \tag{65}
$$

**Proof** Use the bounds

$$
\log(n + x) \geq q \log x + (1-q) \log n - q \log q - (1-q) \log(1-q) \tag{66}
$$

$$
q = \frac{\hat{x}}{n + \hat{x}} \tag{67}
$$

$$
\log(x) \leq ax - 1 + \log \hat{x} \tag{68}
$$

$$
a = 1/\hat{x} \tag{69}
$$

to get

$$f(\alpha) \geq \sum_i n_{ik} q_{ik} \log \alpha_k - n_i a_i \sum_k \alpha_k + (\text{const.}) \tag{70}$$

leading to (65).

The LOO likelihood can be interpreted as the approximation

$$\frac{\Gamma(x+n)}{\Gamma(x)} \approx (x+n-1)^n \tag{71}$$

# 4 Estimating Polya mean and precision separately

The $\boldsymbol{\alpha}$ parameters of the Polya distribution can be decomposed into mean $\mathbf{m}$ and precision $s$, just as in the Dirichlet case, and optimization can be done separately on each part. The decomposition also leads to an interesting interpretation of the Polya, discussed in the next subsection.

## 4.1 A novel interpretation of the Polya distribution

The Polya distribution can be interpreted as a multinomial distribution over a modified set of counts, with a special normalizer. To see why, consider the log-probability of the outcomes $\mathbf{x}$ under the Polya versus multinomial:

$$\log p(\mathbf{x}|\boldsymbol{\alpha}) = \log \Gamma(s) - \log \Gamma(s+n) + \sum_k \log \Gamma(n_k + sm_k) - \log \Gamma(sm_k) \tag{72}$$

$$\log p(\mathbf{x}|\mathbf{p}) = \sum_k n_k \log p_k \tag{73}$$

The multinomial is an exponential family, and the Polya is not. But we can find an approximate exponential family representation of the Polya, by considering derivatives. In the multinomial case, the counts can be recovered from the expression

$$n_k = p_k \frac{d \log p(\mathbf{x}|\mathbf{p})}{dp_k} \tag{74}$$

In the Polya case, the analogous expression is

$$\tilde{n}_k = m_k \frac{d \log p(\mathbf{x}|\boldsymbol{\alpha})}{dm_k} \tag{75}$$

$$= \alpha_k \left( \Psi(n_k + \alpha_k) - \Psi(\alpha_k) \right) \equiv \nu(n_k, \alpha_k) \tag{76}$$

The log-probability of $\mathbf{x}$ under the Polya can thus be approximated by

$$\log p(\mathbf{x}|\boldsymbol{\alpha}) = \sum_k \tilde{n}_k \log m_k \tag{77}$$

When $s \to \infty$, the Dirichlet-multinomial becomes an ordinary multinomial with $\mathbf{p} = \mathbf{m}$, and therefore the 'effective' counts are the same as the ordinary counts:

$$\nu(n_k, \infty) = n_k \tag{78}$$

At the other extreme, as $s \to 0$, the Dirichlet-multinomial favors extreme proportions, and the effective counts are a binarized version of the original counts:

$$\nu(n_k, 0) \quad = \quad \begin{cases} 0 & \text{if } n_k = 0, \\ 1 & \text{if } n_k > 0. \end{cases} \tag{79}$$

For intermediate values of $\alpha$, the mapping $\nu$ behaves like a logarithm, reducing the influence of large counts on the likelihood (see figure 1). Thus the Polya can be understood as a multinomial with 'damped' counts.
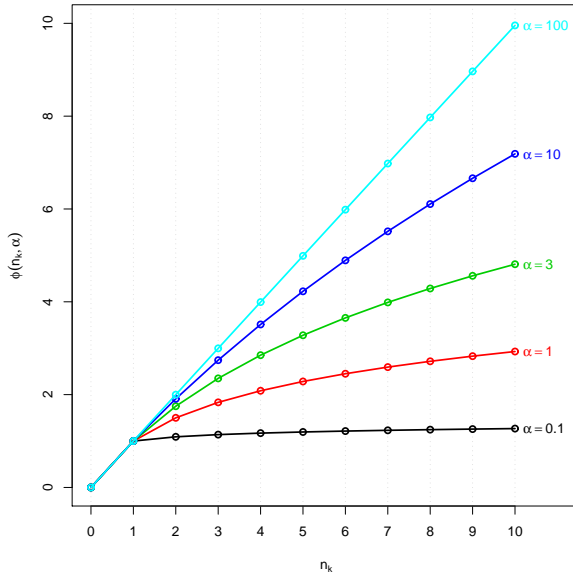


Figure 1: The 'effective' counts for the Polya distribution, as a function of the original count $n_k$ and the parameter $\alpha_k$.

This representation of the Polya also arises in the estimation of $\mathbf{m}$ when $s$ is fixed (section 5).

## 4.2   Estimating Polya precision

The likelihood for $s$ alone is

$$p(D|s) \propto \prod_i \left( \frac{\Gamma(s)}{\Gamma(n_i + s)} \prod_k \frac{\Gamma(n_{ik} + sm_k)}{\Gamma(sm_k)} \right) \tag{80}$$

The derivatives are

$$\frac{d \log p(D|s)}{ds} \quad = \quad \sum_i \Psi(s) - \Psi(n_i + s) + \sum_k m_k \Psi(n_{ik} + sm_k) - m_k \Psi(sm_k) \tag{81}$$

$$\frac{d^2 \log p(D|s)}{ds^2} \quad = \quad \sum_i \Psi'(s) - \Psi'(n_i + s) + \sum_k m_k^2 \Psi'(n_{ik} + sm_k) - m_k^2 \Psi'(sm_k) \tag{82}$$

A convergent fixed-point iteration is

$$s^{new} = s \frac{\sum_{ik} m_k \Psi(n_{ik} + sm_k) - m_k \Psi(sm_k)}{\sum_i \Psi(n_i + s) - \Psi(s)} \tag{83}$$

(the proof is similar to (55)). However, it is very slow. We can get a fast second-order method as follows. When $s$ is small, i.e. the gradient is positive, use the approximation

$$\log p(D|s) \approx a\log(s) + cs + k \tag{84}$$
$$a = -s_0^2 f''(s_0) \tag{85}$$
$$c = f'(s_0) - a/s_0 \tag{86}$$

to get the update

$$s^{new} = -a/c = s/(1 + \frac{f'(s)}{sf''(s)}) \tag{87}$$

except when $c \geq 0$, in which case the solution is $s = \infty$. When $s$ is large, i.e. the gradient is negative, use the approximation

$$\log p(D|s) \approx \frac{a}{2s^2} + \frac{c}{s} \tag{88}$$
$$a = s_0^3(s_0 f''(s_0) + 2f'(s_0)) \tag{89}$$
$$c = -(s_0^2 f'(s_0) + a/s_0) \tag{90}$$

to get the update

$$s^{new} = -a/c = s - \frac{f'(s)}{f''(s) + 3f'(s)/s} \tag{91}$$

For large $s$, the value of $a$ tends to be numerically unstable. If $sf''(s) + 2f'(s)$ is within machine epsilon, then it is better to substitute the limiting value:

$$a \to \sum_i \frac{n_i(n_i-1)(2n_i-1)}{6} - \sum_{ik} \frac{n_{ik}(n_{ik}-1)(2n_{ik}-1)}{6m_k^2} \tag{92}$$

A even faster update for large $s$ is possible by using a richer approximation:

$$\log p(D|s) \approx c\log\left(\frac{s}{s+b}\right) + \frac{e}{s+b} \tag{93}$$
$$c = \sum_{ik} \delta(n_{ik} > 0) - \sum_i \delta(n_i > 0) \tag{94}$$
$$e = -\frac{s_0 + b}{s_0}(s_0(s_0 + b)f'(s_0) - cb) \tag{95}$$
$$b = \text{RootOf}(a_2 b^2 + a_1 b + a_0) \tag{96}$$
$$a_2 = s_0^3(s_0 f''(s_0) + 2f'(s_0)) \tag{97}$$
$$a_1 = 2s_0^2(s_0 f''(s_0) + f'(s_0)) \tag{98}$$
$$a_0 = s_0^2 f''(s_0) + c \tag{99}$$

The approximation comes from setting $c\log s$ to match $\log p(D|s)$ as $s \to 0$ and then choosing $(b, e)$ to match the first two derivatives of $f$ at the current $s$. The resulting update is

$$s^{new} = \frac{cb^2}{e - cb} = \left(\frac{1}{s} - \frac{f'(s)(s+b)^2}{cb^2}\right)^{-1} \tag{100}$$

Note that $a_2$ is equivalent to $a$ above and should be corrected for stability via the same method.

**The case of large dimension**

An interesting special case arises when $K$ is very large. The precision can be estimated simply by counting the number of *singleton* elements in each $\mathbf{x}$. Because the precision acts like a smoothing

parameter on the estimate of $\mathbf{m}$, this result is reminiscent of smoothing methods in document modeling which are based on counting singletons.

If $m_k$ is roughly uniform and $K >> 1$, then $\alpha_k << 1$ and we can use the approximations

$$\Gamma(n_k + \alpha_k) \approx \Gamma(n_k) \tag{101}$$

$$\Gamma(\alpha_k) \approx 1/\alpha_k \tag{102}$$

$$p(\mathbf{x}|s) \approx \frac{\Gamma(s)}{\Gamma(n+s)} \prod_{n_k>0} s m_k \Gamma(n_k) \tag{103}$$

$$\propto \frac{\Gamma(s) s^{\hat{K}}}{\Gamma(n+s)} \tag{104}$$

where $\hat{K}$ is the number of unique observations in $\mathbf{x}$. The approximation does not hold if $s$ is large, which can happen when $\mathbf{m}$ is a good match to the data. But if the dimensionality is large enough, the data will be too sparse for this to happen. The derivatives become

$$\frac{d \log p(D|s)}{ds} \approx \sum_i \Psi(s) - \Psi(n_i + s) + \hat{K}_i/s \tag{105}$$

$$\frac{d^2 \log p(D|s)}{ds^2} \approx \sum_i \Psi'(s) - \Psi'(n_i + s) - \hat{K}_i/s^2 \tag{106}$$

Newton iteration can be used as long as the maximum for $s$ is not on the boundary of $(0, \infty)$. These cases occur when $\hat{K} = 1$ and $\hat{K} = n$.

When the gradient is zero, we have

$$\hat{K} = s(\Psi(n+s) - \Psi(s)) = E[K|s, n] \tag{107}$$

A convergent fixed-point iteration is

$$s^{new} = \frac{\sum_i \hat{K}_i}{\sum_i \Psi(n_i + s^{old}) - \Psi(s^{old})} \tag{108}$$

**Proof** Use the bound

$$\frac{\Gamma(s)}{\Gamma(n+s)} \geq \frac{\Gamma(\hat{s}) \exp((\hat{s} - s)b)}{\Gamma(n+\hat{s})} \tag{109}$$

$$b = \Psi(n+\hat{s}) - \Psi(\hat{s}) \tag{110}$$

to get

$$p(D|s) \geq -s \sum_i b_i + \sum_i \hat{K}_i \log s + (\text{const.}) \tag{111}$$

leading to (108).

Applying the large $K$ approximation to the LOO likelihood gives

$$t = \sum_{ik} \delta(n_{ik} - 1) \quad \text{(number of singletons)} \tag{112}$$

$$f(s) = t \log s - \sum_i n_i \log(n_i - 1 + s) \tag{113}$$

$$\frac{df(s)}{ds} = \frac{t}{s} - \sum_i \frac{n_i}{n_i - 1 + s} \tag{114}$$

For $N = 1$:

$$s \;=\; \frac{t(n-1)}{n-t} \tag{115}$$

$$\frac{s}{s+n} \;=\; \frac{t(n-1)}{n^2-t} \approx \frac{t}{n} \tag{116}$$

which is the result we wanted.

# 5  Estimating Polya mean

The likelihood for **m** only is

$$p(D|\mathbf{m}) \propto \prod_{ik} \frac{\Gamma(n_{ik} + sm_k)}{\Gamma(sm_k)} \tag{117}$$

The maximum can be computed by the fixed-point iteration

$$m_k^{new} \propto \sum_i \nu(n_{ik}, sm_k) \tag{118}$$

where $\nu$ is defined in (76). This update can be understood intuitively as the maximum-likelihood estimate of a multinomial distribution from effective counts $\tilde{n}_{ik} = \nu(n_{ik}, sm_k)$. The proof of this iteration is similar to (55).

For a Newton-Raphson iteration, reparameterize to get

$$m_K \;=\; 1 - \sum_{k=1}^{K-1} m_k \tag{119}$$

$$g_k = \frac{d \log p(D|\mathbf{m})}{dm_k} \;=\; s \sum_i \Psi(n_{ik} + sm_k) - \Psi(sm_k) - \Psi(n_{iK} + sm_K) + \Psi(sm_K) \tag{120}$$

$$\frac{d^2 \log p(D|\mathbf{m})}{dm_k^2} \;=\; s^2 \sum_i \Psi'(n_{ik} + sm_k) - \Psi'(sm_k) + \Psi'(n_{iK} + sm_K) - \Psi'(sm_K) \tag{121}$$

$$\frac{d^2 \log p(D|\mathbf{m})}{dm_k m_j} \;=\; s^2 \sum_i \Psi'(n_{iK} + sm_K) - \Psi'(sm_K) \tag{122}$$

The search should be initialized at $m_k \propto \sum_i n_{ik}$, since for large $s$ this is the exact optimum.

# References

Johnson, N. L., & Kotz, S. (1970). *Distributions in statistics: Continuous univariate distributions.* New York: Hougton Mifflin.

Minka, T. P. (1998). Expectation-Maximization as lower bound maximization. http://research.microsoft.com/~minka/papers/em.html.

Minka, T. P. (2000). Beyond newton's method. http://research.microsoft.com/~minka/papers/newton.html.

Naryanan, A. (1991). Algorithm as 266: Maximum likelihood estimation of the parameters of the dirichlet distribution. *Applied Statistics*, *40*, 365–374. `http://www.psc.edu/~burkardt/dirichlet.html`.

Ronning, G. (1989). Maximum-likelihood estimation of dirichlet distributions. *Journal of Statistical Computation and Simulation*, *32*, 215–221.

# A  Proof of (9)

Use the bound

$$\Gamma(x) \geq \Gamma(\hat{x})\exp((x-\hat{x})\Psi(\hat{x})) \tag{123}$$

to get

$$\frac{1}{N}\log p(D|\alpha) \geq (\sum_k \alpha_k)\Psi(\sum_k \alpha_k^{old}) - \sum_k \log\Gamma(\alpha_k) + \sum_k (\alpha_k - 1)\log\bar{p}_k + (\text{const.}) \tag{124}$$

leading to (9).

# B  Proof of (55)

Use the bound

$$\frac{\Gamma(x)}{\Gamma(n+x)} \geq \frac{\Gamma(\hat{x})\exp((\hat{x}-x)b)}{\Gamma(n+\hat{x})} \tag{125}$$

$$b = \Psi(n+\hat{x}) - \Psi(\hat{x}) \tag{126}$$

and the bound

$$\frac{\Gamma(n+x)}{\Gamma(x)} \geq cx^a \quad \text{if } n \geq 1 \tag{127}$$

$$a = (\Psi(n+\hat{x}) - \Psi(\hat{x}))\hat{x} \tag{128}$$

$$c = \frac{\Gamma(n+\hat{x})}{\Gamma(\hat{x})}\hat{x}^{-a} \tag{129}$$

to get

$$\log p(D|\alpha) \geq -(\sum_k \alpha_k - 1)\sum_i b_i + \sum_k a_{ik}\log\alpha_k + (\text{const.}) \tag{130}$$

leading to (55).

# C  Inverting the $\Psi$ function

This section describes how to compute a high-accuracy solution to

$$\Psi(x) = y \tag{131}$$

for $x$ given $y$. Given a starting guess for $x$, Newton's method can be used to find the root of $\Psi(x) - y = 0$. The Newton update is

$$x^{new} = x^{old} - \frac{\Psi(x) - y}{\Psi'(x)} \tag{132}$$

To start the iteration, use the following asymptotic formulas for $\Psi(x)$:

$$\Psi(x) \approx \begin{cases} \log(x - 1/2) & \text{if } x \geq 0.6 \\ -\frac{1}{x} - \gamma & \text{if } x < 0.6 \end{cases} \tag{133}$$

$$\gamma = -\Psi(1) \tag{134}$$

to get

$$\Psi^{-1}(y) \approx \begin{cases} \exp(y) + 1/2 & \text{if } y \geq -2.22 \\ -\frac{1}{y+\gamma} & \text{if } y < -2.22 \end{cases} \tag{135}$$

With this initialization, five Newton iterations are sufficient to reach fourteen digits of precision.

# D    EM for estimation from counts

Any algorithm for estimation from probability vectors can be turned into an algorithm for estimation from counts, by treating the $\mathbf{p}_i$ as hidden variables in EM. The E-step computes a posterior distribution over $\mathbf{p}_i$:

$$q(\mathbf{p}_i) \sim \mathcal{D}(n_{ik} + \alpha_k) \tag{136}$$

and the M-step maximizes

$$E[\sum_i \log p(\mathbf{p}_i | \alpha)] = N \log \Gamma(\sum_k \alpha_k) - N \sum_k \log \Gamma(\alpha_k) + N \sum_k (\alpha_k - 1) \log \bar{p}_k \tag{137}$$

$$\text{where } \log \bar{p}_k = \frac{1}{N} \sum_i E[\log p_{ik}] \tag{138}$$

$$= \frac{1}{N} \sum_i \Psi(n_{ik} + \alpha_k^{old}) - \Psi(n_i + \sum_k \alpha_k^{old}) \tag{139}$$

This is the same optimization problem as in section 1, with a new definition for $\bar{p}$. It is not necessary or desirable to reach the exact maximum in the M-step; a single Newton step will do. The Newton step will end up using the old Hessian (10) but the new gradient (54). Compared to the exact Newton algorithm, this uses half as much computation per iteration, but usually requires more than twice the iterations.