

# MDI Conditional probability for the context similarity parameter

Stephen Coleman<sup>1,\*</sup>

<sup>1</sup>MRC Biostatistics Unit, Cambridge, UK

\*stephen.coleman@mrc-bsu.cam.ac.uk

## ABSTRACT

The derivation of the conditional probability for the context similarity parameter  $\phi_{12}$  between 2 different contexts. The final form is a mixture of Gamma distributions.

### 1 Context similarity parameter ( $\phi_{12}$ )

For multiple dataset integration (**MDI**) in the case of  $n$  observations in 2 datasets (also referred to as *contexts*):

$$p(\{c_{i1}, c_{i2}\}_{i=1}^n, v) \propto (1 + \phi_{12}\mathbb{I}(c_{i1} = c_{i2})) \prod_{k=1}^2 \gamma_{c_{ik}k} \quad (1)$$

We assume priors of  $\gamma_{1k}, \dots, \gamma_{nk} \stackrel{i.i.d}{\sim} \text{Gamma}(\alpha_k/N, 1) \forall k \in \{1, 2\}$  where  $N$  is the smaller of the number of clusters in the two contexts. Similarly  $\phi_{12} \sim \text{Gamma}(a, b)$ .

From (1) we calculate the normalising constant  $Z$ , and find:

$$Z = \sum_{j_1=1}^N \sum_{j_2=1}^N \left( (1 + \phi_{12}\mathbb{I}(j_1 = j_2)) \prod_{k=1}^2 \gamma_{j_k k} \right) \quad (2)$$

The joint density is hence:

$$p(\{c_{i1}, c_{i2}\}_{i=1}^n, v) = \frac{1}{Z} \prod_{i=1}^n \left( (1 + \phi_{12}\mathbb{I}(c_{i1} = c_{i2})) \prod_{k=1}^2 \gamma_{c_{ik}k} \right) \quad (3)$$

Introduce a strategic latent variable  $v$  such that the form is:

$$p(\{c_{i1}, c_{i2}\}_{i=1}^n, v) = \frac{v^{n-1} \exp(-vZ)}{(n-1)!} \prod_{i=1}^n \left( (1 + \phi_{12}\mathbb{I}(c_{i1} = c_{i2})) \prod_{k=1}^2 \gamma_{c_{ik}k} \right) \quad (4)$$

Where  $Z$  remains as in (2).

#### 1.1 Conditional likelihood

Consider the conditional probability of  $\phi$ , then from (4) and expanding  $Z$ :

$$p(\phi_{12} | \{c_{i1}, c_{i2}\}_{i=1}^n, v) \propto \exp \left( -v \sum_{j_1=1}^N \sum_{j_2=1}^N \left( (1 + \phi_{12}\mathbb{I}(j_1 = j_2)) \prod_{k=1}^2 \gamma_{j_k k} \right) \right) \prod_{i=1}^n \left( (1 + \phi_{12}\mathbb{I}(c_{i1} = c_{i2})) \prod_{k=1}^2 \gamma_{c_{ik}k} \right) \quad (5)$$

Now, consider the coefficients of the two occurrences of  $\phi_{12}$ :

$$a = \prod_{i=1}^n \left( (1 + \phi_{12} \mathbb{I}(c_{i1} = c_{i2})) \prod_{k=1}^2 \gamma_{c_{ik}k} \right) \quad (6)$$

$$b = \sum_{j_1=1}^N \sum_{j_2=1}^N \left( (1 + \phi_{12} \mathbb{I}(j_1 = j_2)) \prod_{k=1}^2 \gamma_{j_kk} \right) \quad (7)$$

Beginning with  $a$  from above:

$$a = \prod_{i=1}^n \gamma_{c_{i1}1} \gamma_{c_{i2}2} (1 + \phi_{12} \mathbb{I}(c_{i1} = c_{i2})) \quad (8)$$

$$= (1 + \phi_{12} \mathbb{I}(c_{i1} = c_{i2}))^n \prod_{i=1}^n \gamma_{c_{i1}1} \gamma_{c_{i2}2} \quad (9)$$

$$\propto (1 + \phi_{12} \mathbb{I}(c_{i1} = c_{i2}))^n \quad (10)$$

$$= (1 + \phi_{12})^{\sum_{i=1}^n \mathbb{I}(c_{i1} = c_{i2})} \quad (11)$$

$$= \sum_{r=0}^{\sum_{i=1}^n \mathbb{I}(c_{i1} = c_{i2})} \binom{\sum_{i=1}^n \mathbb{I}(c_{i1} = c_{i2})}{r} \phi_{12}^r \quad (\text{from the binomial theorem}) \quad (12)$$

Here  $\sum_{i=1}^n \mathbb{I}(c_{i1} = c_{i2})$  is the count of observations assigned to the same cluster in both contexts and will be called  $c$ .  
Now consider  $b$ :

$$b = \exp \left( -v \sum_{j_1=1}^N \sum_{j_2=1}^N \left( (1 + \phi_{12} \mathbb{I}(j_1 = j_2)) \prod_{k=1}^2 \gamma_{j_kk} \right) \right) \quad (13)$$

We see that for our conditional we can ignore all cases when  $j_1 \neq j_2$  as  $\phi_{12}$  is not present in these. This simplifies  $b$  to:

$$b \propto \sum_{j=1}^N \gamma_{j1} \gamma_{j2} (1 + \phi_{12}) \quad (14)$$

$$\propto \sum_{j=1}^N \gamma_{j1} \gamma_{j2} \phi_{12} \quad (15)$$

Thus updating (5) accordingly gives us:

$$p(\phi_{12} | \{c_{i1}, c_{i2}\}_{i=1}^n, v) \propto \exp \left( -v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} \phi_{12} \right) \sum_{r=0}^c \binom{c}{r} \phi_{12}^r \quad (16)$$

We notice this has the structure similar to a mixture of Gamma distributions. We thus have:

$$p(\{c_{i1}, c_{i2}\}_{i=1}^n, v | \phi_{12}) \propto \sum_{r=0}^c \binom{c}{r} \frac{r!}{(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2})^{r+1}} \frac{(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2})^{r+1}}{r!} \phi_{12}^r \exp \left( -v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} \phi_{12} \right) \quad (17)$$

$$= \sum_{r=0}^c \binom{c}{r} \frac{r!}{(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2})^{r+1}} \text{Gamma} \left( r+1, v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} \right) \quad (18)$$

As we know that  $p(\phi_{12} | \{c_{i1}, c_{i2}\}_{i=1}^n, v)$  must integrate over  $\phi_{12}$  to 1, we know the normalising constant must be the sum of the integrals of the Gamma distributions, i.e.:

$$Z_{\phi_{12}} = \sum_{r=0}^c \binom{c}{r} \frac{r!}{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2}\right)^{r+1}} \int \frac{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2}\right)^{r+1}}{r!} \phi_{12}^r \exp\left(-v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} \phi_{12}\right) d\phi_{12} \quad (19)$$

$$= \sum_{r=0}^c \binom{c}{r} \frac{r!}{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2}\right)^{r+1}} \quad (20)$$

Combining these gives:

$$p(\phi_{12} | \{c_{i1}, c_{i2}\}_{i=1}^n, v) = \frac{1}{Z_{\phi_{12}}} \sum_{r=0}^c \binom{c}{r} \frac{r!}{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2}\right)^{r+1}} \text{Gamma}\left(r+1, v \sum_{j=1}^N \gamma_{j1} \gamma_{j2}\right) \quad (21)$$

## 1.2 Posterior distribution

Now if we consider a prior of  $\text{Gamma}(a_0, b_0)$  on the  $\phi_{12}$ , we have a prior probability of:

$$p(\phi_{12}) = \frac{b_0^{a_0}}{(a_0 - 1)!} \phi_{12}^{a_0-1} \exp(-b_0 \phi_{12}) \quad (22)$$

Thus our posterior conditional is:

$$p(\phi_{12} | \cdot) \propto p(\phi_{12}) p(\{c_{i1}, c_{i2}\}_{i=1}^n, v | \phi_{12}) \quad (23)$$

$$\propto \frac{b_0^{a_0}}{(a_0 - 1)!} \phi_{12}^{a_0-1} \exp(-b_0 \phi_{12}) \sum_{r=0}^c \binom{c}{r} \frac{r!}{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2}\right)^{r+1}} \frac{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2}\right)^{r+1}}{r!} \phi_{12}^r \exp\left(-v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} \phi_{12}\right) \quad (24)$$

$$\propto \sum_{r=0}^c \binom{c}{r} \phi_{12}^{r+a_0-1} \exp\left(\left(-v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} - b_0\right) \phi_{12}\right) \quad (25)$$

$$\propto \sum_{r=0}^c \binom{c}{r} \frac{(r+a_0-1)!}{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right)^{r+a_0}} \frac{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right)^{r+a_0}}{(r+a_0-1)!} \phi_{12}^{r+a_0-1} \exp\left(-\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right) \phi_{12}\right) \quad (26)$$

$$= \sum_{r=0}^c \binom{c}{r} \frac{(r+a_0-1)!}{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right)^{r+a_0}} \text{Gamma}\left(r+a_0, v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right) \quad (27)$$

For the normalising constant, we have, similarly to (20):

$$Z'_{\phi_{12}} = \sum_{r=0}^c \binom{c}{r} \frac{(r+a_0-1)!}{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right)^{r+a_0}} \int \frac{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right)^{r+a_0}}{(r+a_0-1)!} \phi_{12}^{r+a_0-1} \exp\left(-\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right) \phi_{12}\right) d\phi_{12} \quad (28)$$

$$= \sum_{r=0}^c \binom{c}{r} \frac{(r+a_0-1)!}{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right)^{r+a_0}} \quad (29)$$

Thus our final posterior on the context similarity parameter  $\phi_{12}$  is:

$$p(\phi_{12} | \{c_{i1}, c_{i2}\}_{i=1}^n, v) = \frac{1}{Z'_{\phi_{12}}} \sum_{r=0}^c \binom{c}{r} \frac{(r+a_0-1)!}{\left(v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right)^{r+a_0}} \text{Gamma}\left(r+a_0, v \sum_{j=1}^N \gamma_{j1} \gamma_{j2} + b_0\right) \quad (30)$$