

Model-based Rank Aggregation

Based on discussion with by Tom Dietterich and Andrew Emmott

Generative Process with Separate λ for each detector

1. Draw $\pi \sim \text{Beta}(\text{scale} = 0.05, \text{shape} = 200)$. This is the proportion of anomalous points
2. Draw $\lambda \triangleq \{\alpha, \beta\} \sim \text{Gamma}(\cdot, \cdot)$. This is the score “bonus” for being anomalous
3. For each point x ,
 - a. Draw $\theta_x \sim \text{Bern}(\pi)$
 - b. If $\theta_x = 0$ then x is “normal” and $\text{score}(x, i) \sim \text{Gamma}(\alpha_r, \beta_r)$
 - c. Else x is “anomalous”, and $\text{score}(x, i) \sim \text{Gamma}(\alpha_a, \beta_a)$
4. Sort the $\{x\}$ into descending order and assign ranks such that $\text{rank}(x, i)$ is the position of x in the sorted list.

Assume there are D detectors. $\mathbf{x}_i = \{x_{i1}, \dots, x_{iD}\}$ are the scores reported by detectors d_1, \dots, d_D . Let $\alpha_r = \{\alpha_{r1}, \dots, \alpha_{rD}\}$ be the mean scores for regular instances and let $\alpha_a = \{\alpha_{a1}, \dots, \alpha_{aD}\}$ be the means for anomaly scores. The likelihood of a score under the distribution for ‘regular’ scores is:

$$f_r(\mathbf{x}_i | \alpha_r, \beta_r) \sim \text{Gamma}(\alpha_r, \beta_r) = \prod_{d=1}^D \frac{\beta_{rd}^{\alpha_{rd}}}{\Gamma(\alpha_{rd})} x_{rd}^{\alpha_{rd}-1} e^{-\beta_{rd} x_{rd}} \quad (1)$$

And the likelihood of a score under the distribution for ‘anomalous’ scores is:

$$f_a(\mathbf{x}_i | \alpha_a, \beta_a) \sim \text{Gamma}(\alpha_a, \beta_a) = \prod_{d=1}^D \frac{\beta_{ad}^{\alpha_{ad}}}{\Gamma(\alpha_{ad})} x_{ad}^{\alpha_{ad}-1} e^{-\beta_{ad} x_{ad}} \quad (2)$$

The likelihood of each score assuming we know whether it is normal or anomalous is:

$$f(\mathbf{x}_i | \pi, \alpha_a, \beta_a, \alpha_r, \beta_r) = ((1 - \pi) f_r(\mathbf{x}_i | \alpha_r, \beta_r))^{1-\theta_i} (\pi f_a(\mathbf{x}_i | \alpha_a, \beta_a))^{\theta_i} \quad (3)$$

The priors on $\pi, \lambda_r = \{\alpha_r, \beta_r\}, \lambda_a = \{\alpha_a, \beta_a\}$ are:

$$f_\pi(\pi) = \frac{\pi^{\alpha-1} (1-\pi)^{\beta-1}}{B(\alpha, \beta)}, \quad f_{\lambda_a}(\alpha_{ad}, \beta_{ad} | p_{ad}, q_{ad}, r_{ad}, s_{ad}) \propto \frac{p_{ad}^{\alpha_{ad}-1}}{\Gamma(\alpha_{ad})^{r_{ad}} \beta_{ad}^{-\alpha_{ad} s_{ad}}} e^{-\beta_{ad} q_{ad}},$$

$$f_{\lambda_r}(\alpha_{rd}, \beta_{rd} | p_{rd}, q_{rd}, r_{rd}, s_{rd}) \propto \frac{p_{rd}^{\alpha_{rd}-1}}{\Gamma(\alpha_{rd})^{r_{rd}} \beta_{rd}^{-\alpha_{rd} s_{rd}}} e^{-\beta_{rd} q_{rd}}$$

Where we assume that $\alpha, \beta, \{p_{ad}, q_{ad}, r_{ad}, s_{ad}\}$ and $\{p_{rd}, q_{rd}, r_{rd}, s_{rd}\}$ are known constants; $B(\alpha, \beta)$ is the Beta function and $\Gamma(\alpha_{ad}), \Gamma(\alpha_{rd})$ are Gamma functions. We will denote $\lambda_{ad} = \{\alpha_{ad}, \beta_{ad}\}$ and $\lambda_{rd} = \{\alpha_{rd}, \beta_{rd}\}$.

Therefore,

$$f(\mathbf{x}_i, \pi, \lambda_a, \lambda_r) \propto f(\mathbf{x}_i | \pi, \lambda_a, \lambda_r) f_{\lambda_a}(\lambda_a) f_{\lambda_r}(\lambda_r) f_\pi(\pi)$$

$$\propto \left(\pi \prod_{d=1}^D f_a(x_{id} | \lambda_{ad}) \right)^{\theta_i} \left((1 - \pi) \prod_{d=1}^D f_r(x_{id} | \lambda_{rd}) \right)^{1-\theta_i} \left(\prod_{d=1}^D \frac{p_{ad}^{\alpha_{ad}-1}}{\Gamma(\alpha_{ad})^{r_{ad}} \beta_{ad}^{-\alpha_{ad} s_{ad}}} e^{-\beta_{ad} q_{ad}} \frac{p_{rd}^{\alpha_{rd}-1}}{\Gamma(\alpha_{rd})^{r_{rd}} \beta_{rd}^{-\alpha_{rd} s_{rd}}} e^{-\beta_{rd} q_{rd}} \right) \frac{\pi^{\alpha-1} (1-\pi)^{\beta-1}}{B(\alpha, \beta)} \quad (4)$$

The complete-data likelihood of all scores across all detectors is (where n is the number of instances):

$$L(\mathbf{x}; \pi, \lambda_a, \lambda_r) \propto [\{\prod_{i=1}^n \prod_{d=1}^D f(x_{id} | \pi, \lambda_{ad})\} \{\prod_{d=1}^D f_{\lambda_a}(\lambda_{ad}) f_{\lambda_r}(\lambda_{rd})\}] f_\pi(\pi) \quad (5)$$

$$\begin{aligned}
&= \left[\prod_{i=1}^n \left\{ \left(\pi \prod_{d=1}^D f_a(x_{id}|\lambda_{ad}) \right)^{\theta_i} \left((1 - \pi) \prod_{d=1}^D f_r(x_{id}|\lambda_{rd}) \right)^{1-\theta_i} \right\} \right] \left(\prod_{d=1}^D \frac{p_{ad}^{\alpha_{ad}-1}}{\Gamma(\alpha_{ad})^{r_{ad}} \beta_{ad}^{-\alpha_{ad}s_{ad}}} e^{-\beta_{ad}q_{ad}} \frac{p_{rd}^{\alpha_{rd}-1}}{\Gamma(\alpha_{rd})^{r_{rd}} \beta_{rd}^{-\alpha_{rd}s_{rd}}} e^{-\beta_{rd}q_{rd}} \right) \frac{\pi^{\alpha-1}(1-\pi)^{\beta-1}}{B(\alpha, \beta)} \\
&\quad (6)
\end{aligned}$$

Instead of assigning the hard class labels $((1 - \theta_i), \theta_i)$ which will be hard to infer, we will use soft-assignments (i.e., responsibilities) denoted by z_{ai} and z_{ri} which refer to the probability of assigning i -th score of d -th detector to ‘anomaly’ and ‘normal’ classes respectively; $z_{ai} + z_{ri} = 1$. We rewrite $L(\mathbf{x}; \pi, \boldsymbol{\lambda}_a, \boldsymbol{\lambda}_r)$:

$$\begin{aligned}
&\propto \left[\prod_{i=1}^n \left\{ \left(\pi \prod_{d=1}^D f_a(x_{id}|\lambda_{ad}) \right)^{z_{ai}} \left((1 - \pi) \prod_{d=1}^D f_r(x_{id}|\lambda_{rd}) \right)^{z_{ri}} \right\} \right] \left[\prod_{d=1}^D \left\{ \prod_{k \in \{a, r\}} \frac{p_{kd}^{\alpha_{kd}-1}}{\Gamma(\alpha_{kd})^{r_{kd}} \beta_{kd}^{-\alpha_{kd}s_{kd}}} e^{-\beta_{kd}q_{kd}} \right\} \right] \frac{\pi^{\alpha-1}(1-\pi)^{\beta-1}}{B(\alpha, \beta)} \\
&\quad (7)
\end{aligned}$$

The log-likelihood is:

$$\begin{aligned}
l(\mathbf{x}; \pi, \boldsymbol{\lambda}_a, \boldsymbol{\lambda}_r, \sigma^2) &= \log(L(\mathbf{x}; \pi, \boldsymbol{\lambda}_a, \boldsymbol{\lambda}_r, \sigma^2)) \\
&= \sum_{i=1}^n z_{ai} \left\{ \log(\pi) + \sum_{d=1}^D \log(f_a(x_{id}|\lambda_{ad})) \right\} + \sum_{i=1}^n z_{ri} \left\{ \log(1 - \pi) + \sum_{d=1}^D \log(f_r(x_{id}|\lambda_{rd})) \right\} \\
&\quad + \sum_{d=1}^D \sum_{k \in \{a, r\}} \{ (\alpha_{kd} - 1) \log(p_{kd}) + \alpha_{kd}s_{kd} \log(\beta_{kd}) - r_{kd} \log(\Gamma(\alpha_{kd})) - \beta_{kd}q_{kd} \} \\
&\quad + (\alpha - 1) \log(\pi) + (\beta - 1) \log(1 - \pi) - \log(B(\alpha, \beta)) \\
&\quad (8)
\end{aligned}$$

After substituting $f_a(x_{id}|\lambda_{ad})$ and $f_r(x_{id}|\lambda_{rd})$ in the above equation:

$$\begin{aligned}
l(\mathbf{x}; \pi, \lambda_{ad}, \lambda_{rd}) &= \sum_{i=1}^n z_{ai} \left\{ \log(\pi) + \sum_{d=1}^D \{ (\alpha_{ad} - 1) \log(x_{id}) + \alpha_{ad} \log(\beta_{ad}) - \log(\Gamma(\alpha_{ad})) - \beta_{ad}x_{id} \} \right\} \\
&\quad + \sum_{i=1}^n z_{ri} \left\{ \log(1 - \pi) + \sum_{d=1}^D \{ (\alpha_{rd} - 1) \log(x_{id}) + \alpha_{rd} \log(\beta_{rd}) - \log(\Gamma(\alpha_{rd})) - \beta_{rd}x_{id} \} \right\} \\
&\quad + \sum_{d=1}^D \sum_{k \in \{a, r\}} \{ (\alpha_{kd} - 1) \log(p_{kd}) + \alpha_{kd}s_{kd} \log(\beta_{kd}) - r_{kd} \log(\Gamma(\alpha_{kd})) - \beta_{kd}q_{kd} \} \\
&\quad + (\alpha - 1) \log(\pi) + (\beta - 1) \log(1 - \pi) - \log(B(\alpha, \beta)) + (\text{some constant}) \\
&\quad (9)
\end{aligned}$$

M-Step: Derive the MLE of parameters by differentiation.

MLE for π :

$$\frac{\partial l(\mathbf{x}; \pi, \lambda_{ad}, \lambda_{rd}, \sigma^2)}{\partial \pi} = \frac{1}{\pi} \sum_{i=1}^n z_{ai} - \frac{1}{1 - \pi} \sum_{i=1}^n z_{ri} + \frac{(\alpha - 1)}{\pi} - \frac{(\beta - 1)}{1 - \pi} = 0$$

$$\begin{aligned}
&\Rightarrow (1 - \pi) \sum_{i=1}^n z_{ai} - \pi \sum_{i=1}^n z_{ri} + (1 - \pi)(\alpha - 1) - \pi(\beta - 1) = 0 \\
&\Rightarrow \sum_{i=1}^n z_{ai} + (\alpha - 1) = \pi \left(\sum_{i=1}^n z_{ai} + \sum_{i=1}^n z_{ri} + (\alpha - 1) + (\beta - 1) \right) \\
&\Rightarrow \hat{\pi} = \frac{\sum_{i=1}^n z_{ai} + (\alpha - 1)}{\left(\sum_{i=1}^n z_{ai} + \sum_{i=1}^n z_{ri} + (\alpha - 1) + (\beta - 1) \right)} \\
&\Rightarrow \hat{\pi} = \frac{\sum_{i=1}^n z_{ai} + (\alpha - 1)}{(n + (\alpha - 1) + (\beta - 1))}
\end{aligned} \tag{10}$$

The above follows from the observation that: $z_{ai} + z_{ri} = 1$.

MLE for β_{kd} for $k \in \{a, r\}$:

$$\begin{aligned}
\frac{\partial l(\mathbf{x}; \pi, \lambda_{ad}, \lambda_{rd})}{\partial \beta_{kd}} &= \frac{\partial (\sum_{i=1}^n z_{ki} \{ \alpha_{kd} \log(\beta_{kd}) - \beta_{kd} x_{id} \} + \alpha_{kd} s_{kd} \log(\beta_{kd}) - \beta_{kd} q_{kd})}{\partial \beta_{kd}} = 0 \\
&\Rightarrow \frac{\partial (\alpha_{kd} (s_{kd} + \sum_{i=1}^n z_{ki}) \log(\beta_{kd}) - (q_{kd} + \sum_{i=1}^n z_{ki} x_{id}) \beta_{kd})}{\partial \beta_{kd}} = 0 \\
&\Rightarrow \frac{\alpha_{kd}}{\beta_{kd}} \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) = \left(q_{kd} + \sum_{i=1}^n z_{ki} x_{id} \right) \\
&\Rightarrow \hat{\beta}_{kd} = \alpha_{kd} \frac{s_{kd} + \sum_{i=1}^n z_{ki}}{q_{kd} + \sum_{i=1}^n z_{ki} x_{id}}
\end{aligned} \tag{11}$$

MLE for α_{kd} for $k \in \{a, r\}$:

$$\begin{aligned}
\frac{\partial l(\mathbf{x}; \pi, \lambda_{ad}, \lambda_{rd})}{\partial \alpha_{kd}} &= \frac{\partial}{\partial \alpha_{kd}} \left\{ \sum_{i=1}^n z_{ki} \log(\pi) + \sum_{i=1}^n z_{ki} \{ (\alpha_{kd} - 1) \log(x_{id}) + \alpha_{kd} \log(\beta_{kd}) - \log(\Gamma(\alpha_{kd})) - \beta_{ad} x_{id} \} \right. \\
&\quad \left. + (\alpha_{kd} - 1) \log(p_{kd}) + \alpha_{kd} s_{kd} \log(\beta_{kd}) - r_{kd} \log(\Gamma(\alpha_{kd})) - \beta_{kd} q_{kd} \right\} = 0 \\
&\Rightarrow \frac{\partial}{\partial \alpha_{kd}} \left\{ \sum_{i=1}^n z_{ki} \log(\pi) + \left(\log(p_{kd}) + \sum_{i=1}^n z_{ki} \log(x_{id}) \right) (\alpha_{kd} - 1) + \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) \alpha_{kd} \log(\beta_{kd}) \right. \\
&\quad \left. - \left(r_{kd} + \sum_{i=1}^n z_{ki} \right) \log(\Gamma(\alpha_{kd})) - \left(q_{kd} + \sum_{i=1}^n z_{ki} x_{id} \right) \beta_{kd} \right\} = 0 \\
&\Rightarrow \frac{\partial}{\partial \alpha_{kd}} \left\{ \sum_{i=1}^n z_{ki} \log(\pi) + \left(\log(p_{kd}) + \sum_{i=1}^n z_{ki} \log(x_{id}) \right) (\alpha_{kd} - 1) \right. \\
&\quad \left. + \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) \alpha_{kd} \log \left(\alpha_{kd} \frac{s_{kd} + \sum_{i=1}^n z_{ki}}{q_{kd} + \sum_{i=1}^n z_{ki} x_{id}} \right) - \left(r_{kd} + \sum_{i=1}^n z_{ki} \right) \log(\Gamma(\alpha_{kd})) \right. \\
&\quad \left. - \left(q_{kd} + \sum_{i=1}^n z_{ki} x_{id} \right) \alpha_{kd} \frac{s_{kd} + \sum_{i=1}^n z_{ki}}{q_{kd} + \sum_{i=1}^n z_{ki} x_{id}} \right\} = 0 \\
&\Rightarrow \left(\log(p_{kd}) + \sum_{i=1}^n z_{ki} \log(x_{id}) \right) + \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) \log(\alpha_{kd}) + \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) \\
&\quad + \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) \log \left(\frac{s_{kd} + \sum_{i=1}^n z_{ki}}{q_{kd} + \sum_{i=1}^n z_{ki} x_{id}} \right) - \left(r_{kd} + \sum_{i=1}^n z_{ki} \right) \Psi(\alpha_{kd}) - \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) = 0
\end{aligned}$$

$$\begin{aligned} \Rightarrow \sum_{i=1}^n z_{ki} \log(x_{id}) + \log(p_{kd}) + \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) \log(\alpha_{kd}) + \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) \log \left(\frac{s_{kd} + \sum_{i=1}^n z_{ki}}{q_{kd} + \sum_{i=1}^n z_{ki} x_{id}} \right) \\ - \left(\sum_{i=1}^n z_{ki} + r_{kd} \right) \Psi(\alpha_{kd}) = 0 \end{aligned} \quad (12)$$

Equation 12 can be solved for optimal α_{kd} using Newton Raphson. We should note that if this had been the case of a single Gamma distribution, then there would have been only one maxima. However, in the present context we have a mixture, and therefore the problem is no longer convex.

$$\begin{aligned} f'(\alpha_{kd}) &= \sum_{i=1}^n z_{ki} \log(x_{id}) + \log(p_{kd}) + \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) \log(\alpha_{kd}) - \left(s_{kd} + \sum_{i=1}^n z_{ki} \right) \log \left(\frac{q_{kd} + \sum_{i=1}^n z_{ki} x_{id}}{s_{kd} + \sum_{i=1}^n z_{ki}} \right) \\ &\quad - \left(\sum_{i=1}^n z_{ki} + r_{kd} \right) \Psi(\alpha_{kd}) \\ f''(\alpha_{kd}) &= \left(\sum_{i=1}^n z_{ki} + s_{kd} \right) \frac{1}{\alpha_{kd}} - \left(\sum_{i=1}^n z_{ki} + r_{kd} \right) \Psi'(\alpha_{kd}) \\ \alpha_{kd}^{new} &= \alpha_{kd} - \frac{f'(\alpha_{kd})}{f''(\alpha_{kd})} \end{aligned} \quad (13)$$

An alternative to Equation 13 was proposed by Minka using a local approximation [1]:

$$\frac{1}{\alpha_{kd}^{new}} = \frac{1}{\alpha_{kd}} + \frac{f'(\alpha_{kd})}{\alpha_{kd}^2 f''(\alpha_{kd})}$$

E-Step: Compute:

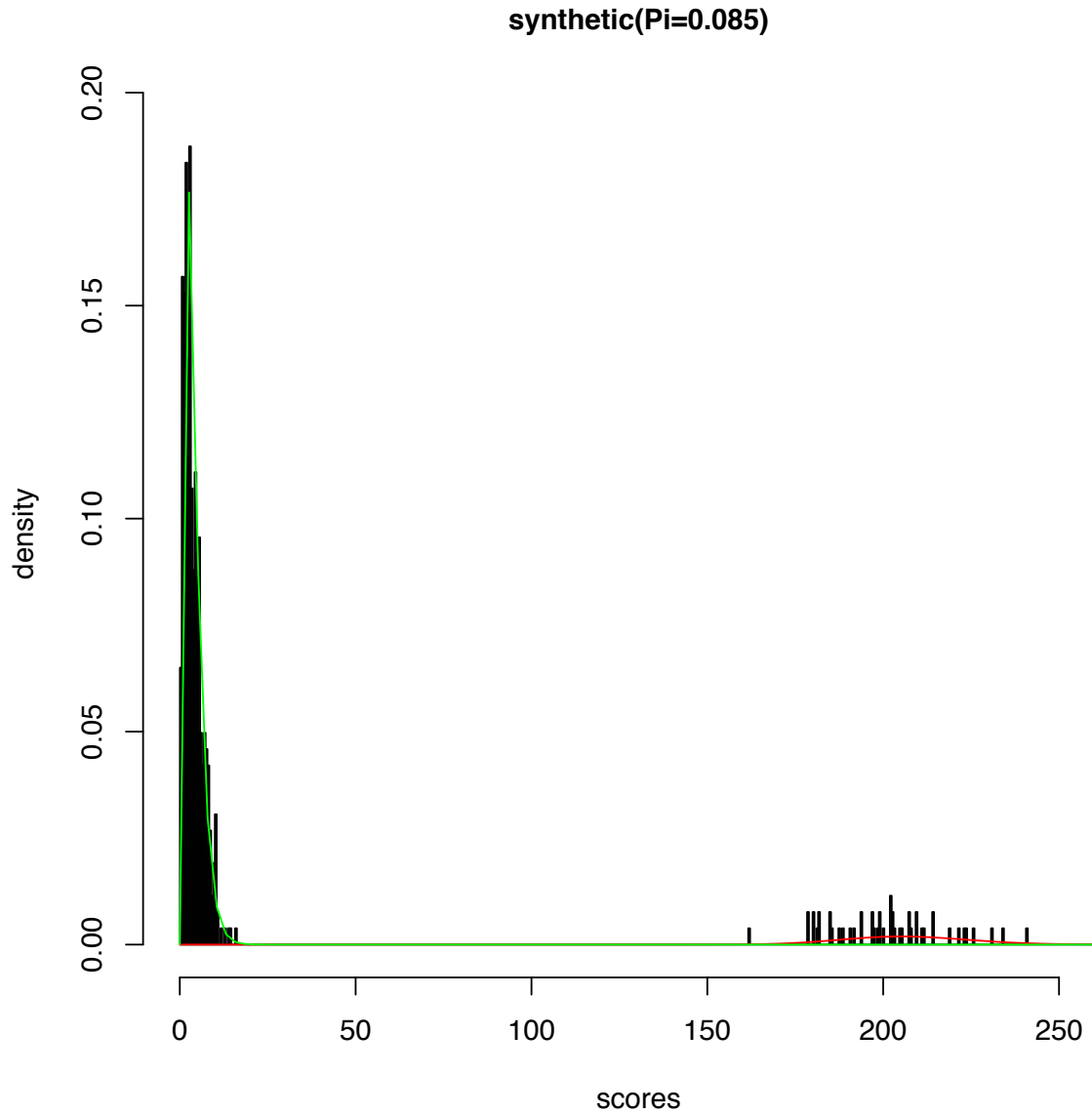
$$E[z_{ai} = 1 | \mathbf{x}] = P(z_{ai} = 1) = \frac{\pi \prod_{d=1}^D f_a(x_{id} | \lambda_{ad})}{\pi \prod_{d=1}^D f_a(x_{id} | \lambda_{ad}) + (1-\pi) \prod_{d=1}^D f_r(x_{id} | \lambda_{rd})} \quad (14)$$

and,

$$E[z_{ri} = 1 | \mathbf{x}] = P(z_{ri} = 1) = \frac{(1-\pi) \prod_{d=1}^D f_r(x_{id} | \lambda_{rd})}{\pi \prod_{d=1}^D f_a(x_{id} | \lambda_{ad}) + (1-\pi) \prod_{d=1}^D f_r(x_{id} | \lambda_{rd})} \quad (15)$$

Illustration of a fit on Synthetic Data

In this simple example, we assume that there is only one detector (in a real application we expect there to be more than one detector.) Hence, we are just fitting a mixture of two Gamma distributions on univariate data. The green distribution is the nominal distribution whereas red is the anomaly distribution. The estimated fraction of anomalies is 0.085.



References

[1] Thomas P. Minka. *Beyond newton's method*. research.microsoft.com/~minka/papers/newton.html, 2000.