

1. Problem

Most datasets used for cyber-security can be considered as mixtures of human and automated events. For example, it is estimated that the proportion of automated traffic in Network Flow data is approximately 95%. For statistical purposes, it is essential to correctly separate these two types of activity, in order to build sound models of normal behaviour of the network.

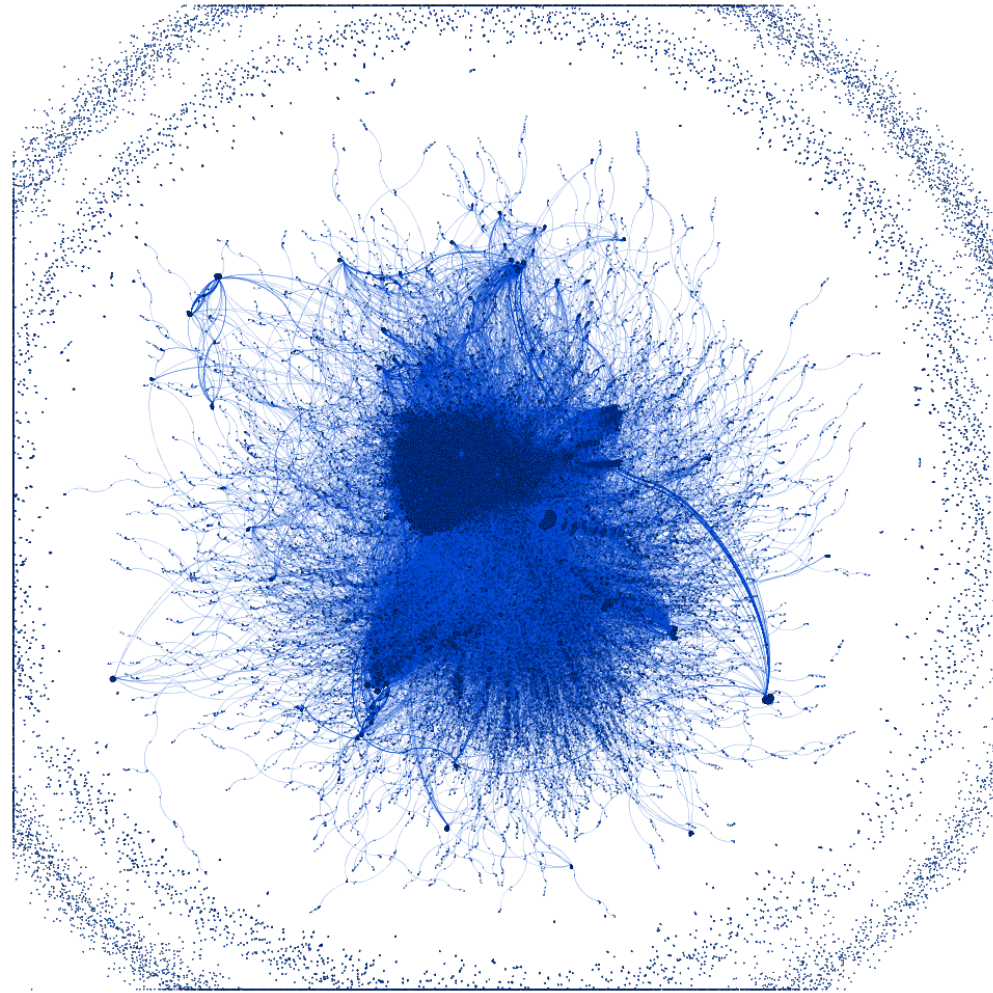


Figure 1: Imperial College network graph on June 7, 2017, 11:15 – 11:16am. Each node corresponds to an IP address, an edge is drawn if the two IPs have connected within the observation period.

2. Detection of periodicities

Methodology developed in **Heard, Rubin-Delanchy and Lawson (2014)**:

- $t_1, t_2, \dots, t_N \rightarrow$ timestamps of the NetFlow events involving a client X and a server Y ,
- $N(t), t \geq 0 \rightarrow$ counting process: number of NetFlow records involving the client X and the server Y at each time point t , starting from $t = 0$,
- Periodogram $\hat{S}(f)$ at frequency $f > 0$:

$$\hat{S}(f) = \frac{1}{T} \left| \sum_{t=1}^T \left(dN(t) - \frac{N(T)}{T} \right) e^{-2\pi i f t} \right|^2$$

where $dN(t) = N(t) - N(t-1)$.

- Fourier's g -test for the null H_0 of no periodicities:

$$g = \frac{\max_{1 \leq k \leq \lfloor T/2 \rfloor} \hat{S}(f_k)}{\sum_{1 \leq j \leq \lfloor T/2 \rfloor} \hat{S}(f_j)}, \quad f_k = \frac{k}{T\Delta t}$$

- Setting $\lambda = \min\{\lfloor 1/g \rfloor, \lfloor T/2 \rfloor\}$, the p -value is:

$$\mathbb{P}(g > g_*) = \sum_{j=1}^{\lambda} (-1)^{j-1} \binom{m}{j} (1 - jg_*)^{m-1}$$

3. Transforming the data

Suppose that an edge is periodic at significance level α with periodicity $p = T\Delta t / \arg\max_{1 \leq k \leq \lfloor T/2 \rfloor} \hat{S}(f_k)$. Let t_1, \dots, t_N be the sequence of **arrival times** on the edge. The quantity of interest for inference is a **latent assignment** z_i , defined as follows:

$$z_i = \begin{cases} 0 & \text{if } t_i \text{ is human} \\ 1 & \text{if } t_i \text{ is automated} \end{cases}$$

where $\mathbb{P}(z_i = 1) = \theta$ and $\mathbb{P}(z_i = 0) = 1 - \theta$.

Two quantities are used to model the arrival times:

- the **wrapped arrival time** x_i :

$$x_i = (t_i \bmod p) \times 2\pi/p$$

- the **daily arrival time** y_i :

$$y_i = (t_i \bmod 86400) \times 2\pi/86400$$

where 86400 is the number of seconds in one day.

4. The model

- For simplicity, assume $T \bmod 86400 = 0$ and $T \bmod p = 0$. Then the density of an arrival time can be decomposed as:

$$f(t_i | z_i) \propto f_A(x_i)^{z_i} f_H(y_i)^{1-z_i}$$

- Human events are modelled using the daily arrival time y_i , automated events using the wrapped arrival time x_i .
- *Fixed phase polling*: event times occur every p seconds plus a random zero-mean error.

$$x_i | (z_i = 1, \mu, \sigma^2) \stackrel{d}{\sim} \text{WN}_{[0, 2\pi)}(\mu, \sigma^2)$$

- Unknown density of the daily arrival times \rightarrow step function:

$$p(y_i | z_i = 0, \mathbf{h}, \boldsymbol{\tau}, \ell) = \sum_{j=1}^{\ell} \frac{h_j}{\tau_{(j+1)} - \tau_{(j)}} \mathbb{1}_{[\tau_{(j)}, \tau_{(j+1)})}(y_i)$$

where ℓ is the number of bins, $\boldsymbol{\tau} = (\tau_1, \dots, \tau_{\ell+1})$ are the bin locations, and $\mathbf{h} = (h_1, \dots, h_{\ell})$, $\sum_j h_j = 1, h_j \geq 0 \forall j$ are the bar heights.

- The resulting model, assuming $T \bmod 86400 = 0$ and $\lfloor T/p \rfloor \gg T \bmod p$, is a mixture of the two components:

$$f(t_i | z_i) \propto \left(\frac{1}{\sqrt{2\pi\sigma^2}} \sum_{k=-\infty}^{\infty} \exp \left\{ -\frac{1}{2\sigma^2} (x_i + 2\pi k - \mu)^2 \right\} \right)^{z_i} \left(\sum_{j=1}^{\ell} \frac{h_j}{\tau_{(j+1)} - \tau_{(j)}} \mathbb{1}_{[\tau_{(j)}, \tau_{(j+1)})}(y_i) \right)^{1-z_i}$$

- Convenient choice of prior distributions for Bayesian inference:

- $(\mu, \sigma^2) \stackrel{d}{\sim} \text{NIG}(\mu_0, \lambda_0, \alpha_0, \beta_0)$,
- $\theta \stackrel{d}{\sim} \text{Beta}(\gamma_0, \delta_0)$,
- $\mathbf{h} | \ell, \boldsymbol{\tau} \stackrel{d}{\sim} \text{Dirichlet}[\eta(\tau_{(j+1)} - \tau_{(j)})]$,
- $\boldsymbol{\tau} | \ell \stackrel{d}{\sim} \text{Unif}[0, 2\pi]^\ell$,
- $\ell \stackrel{d}{\sim} \text{Geo}(\nu)$.

- A Collapsed Metropolis-within-Gibbs sampler with RJMCMC can be used to sample from the posterior distribution.

- The algorithm successfully separates human and automated activity in synthetic (labelled) datasets.

- Reasonable results on real edges, where the true labels are not available.

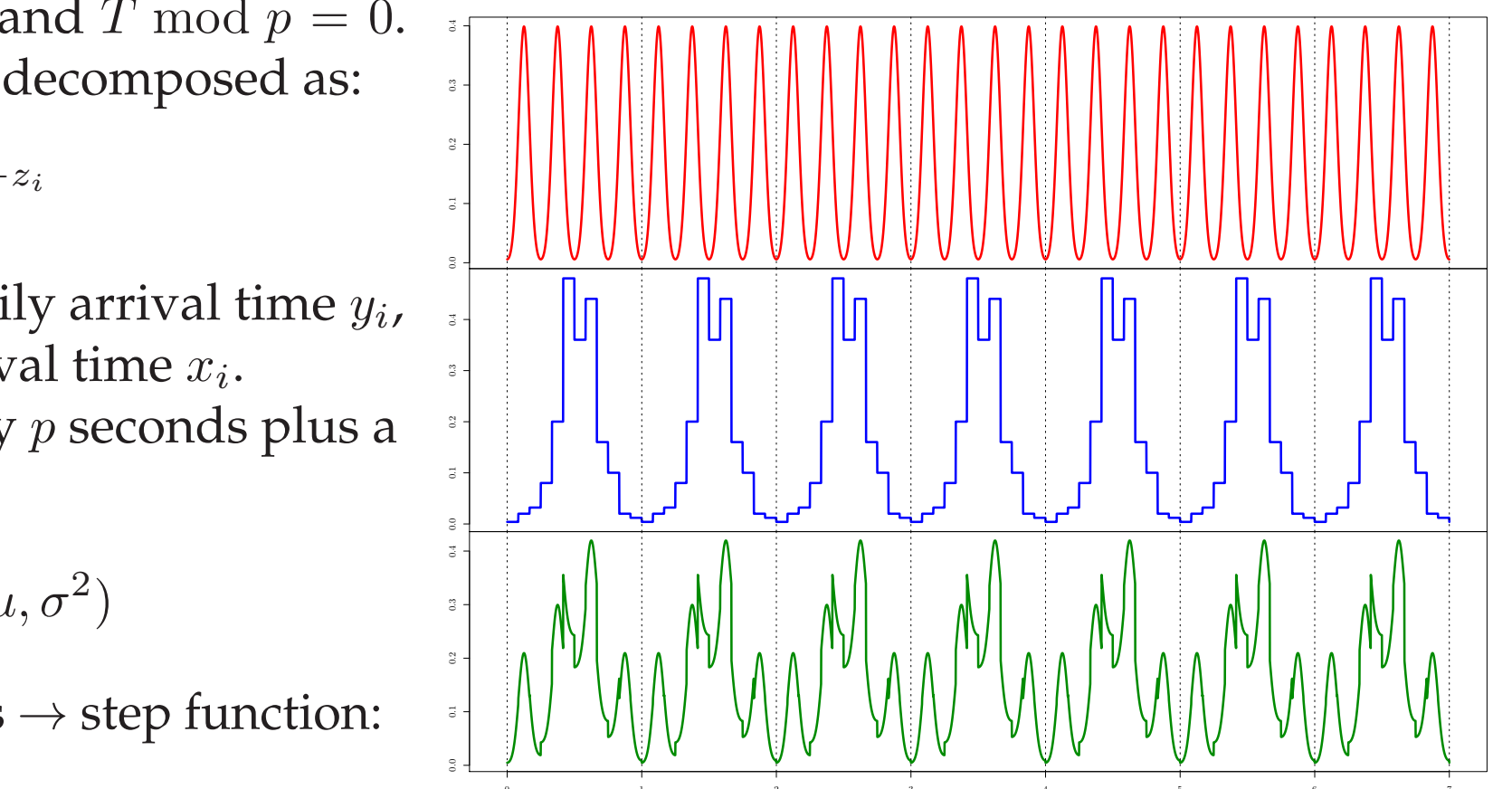


Figure 2: Example of the densities used in the model for $p = 6$ hours, $\mu = \pi$, $\sigma^2 = 1$, $\theta = 0.5$, $\ell = 12$, $\tau_j = \frac{2\pi j}{\ell}$, \mathbf{h} chosen to resemble a human-like distribution. Top plot (red): density of the automated events. Middle plot (blue): density of the human events. Bottom plot (green): density of the 50-50 mixture.

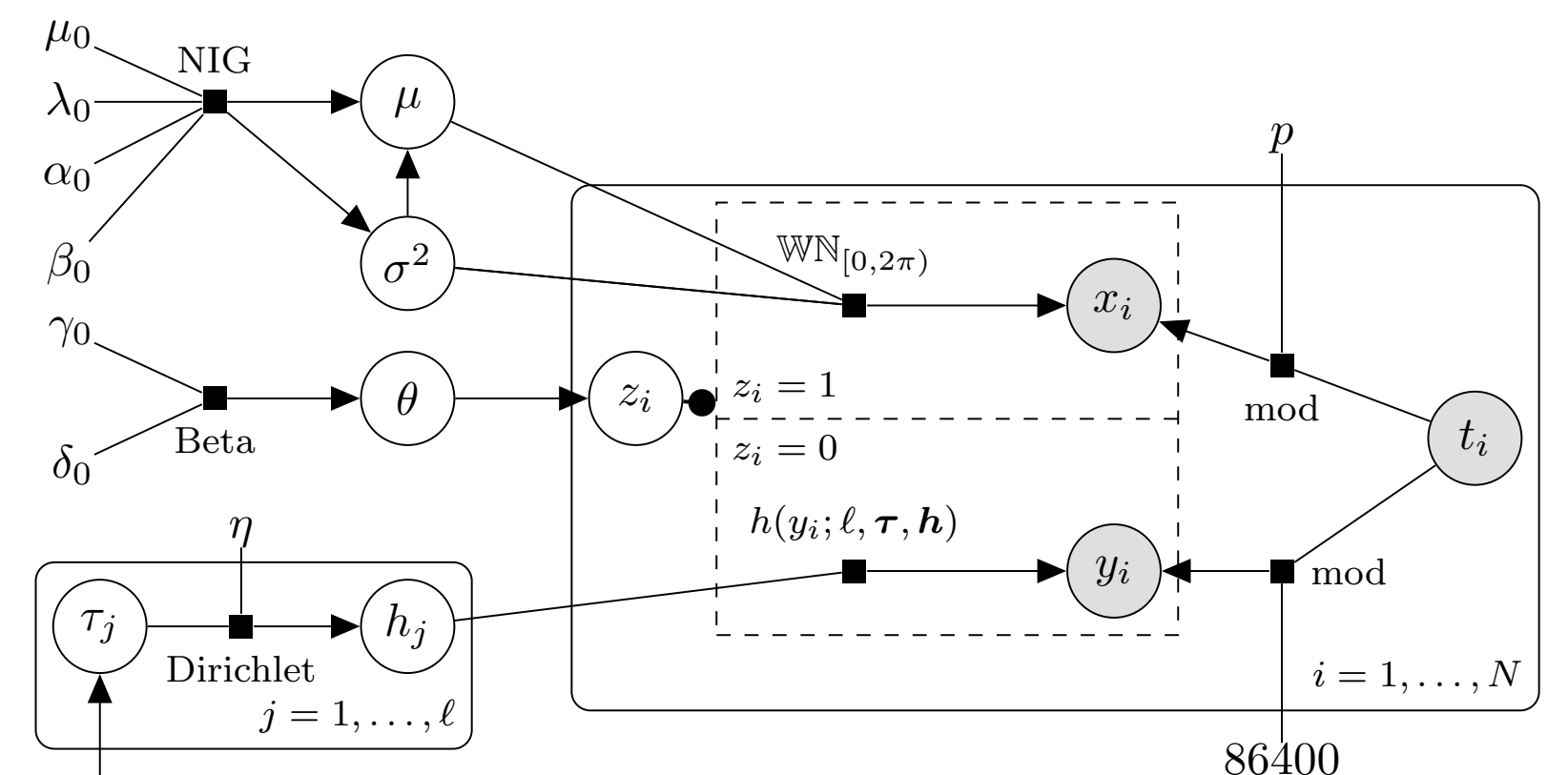


Figure 3: Graphical representation of the joint Bayesian model.

5. Results on a real edge

- 8 days of connections between an IP X and the outlook.com IP 13.107.42.11.
- 7375 events, 1329 filtered human connections.
- The activity slightly increases during the day, suggesting a mixture of human and automated events.
- The distribution of human events obtained from the model shows a clear diurnal pattern, with reduced activity during the night.
- Events are not labeled in this example, but encouraging results have been obtained on synthetic labeled data.

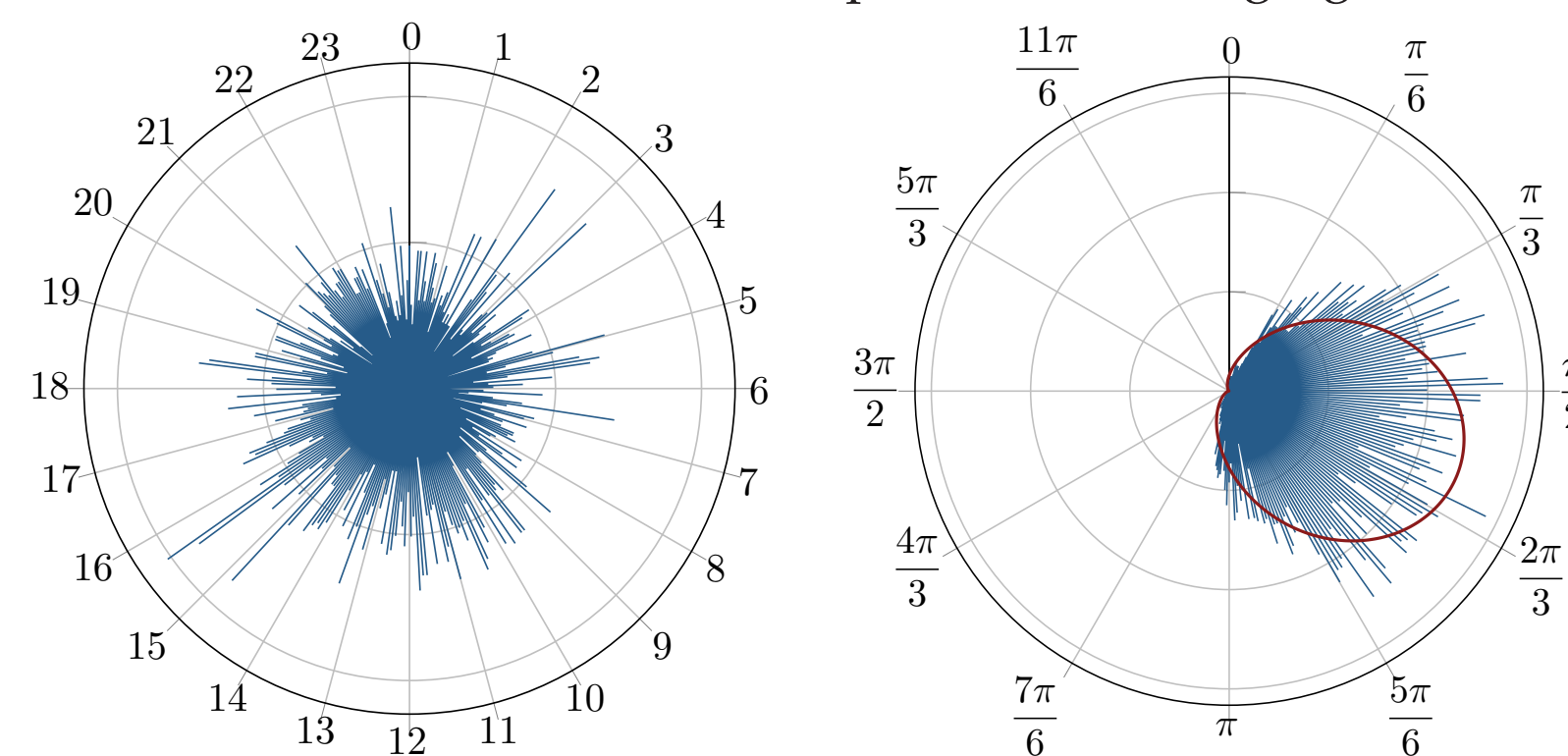


Figure 4: Daily distribution of the data, slight evidence of increased activity during working hours.

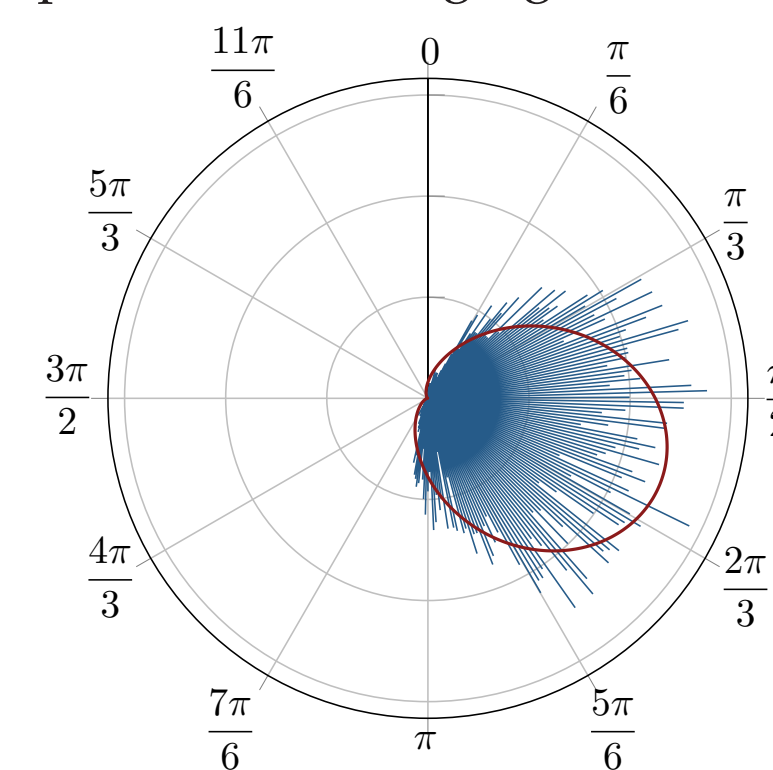


Figure 5: Distribution of the wrapped data, $p = 8s$ and model fit (MAP estimates of μ and σ^2).

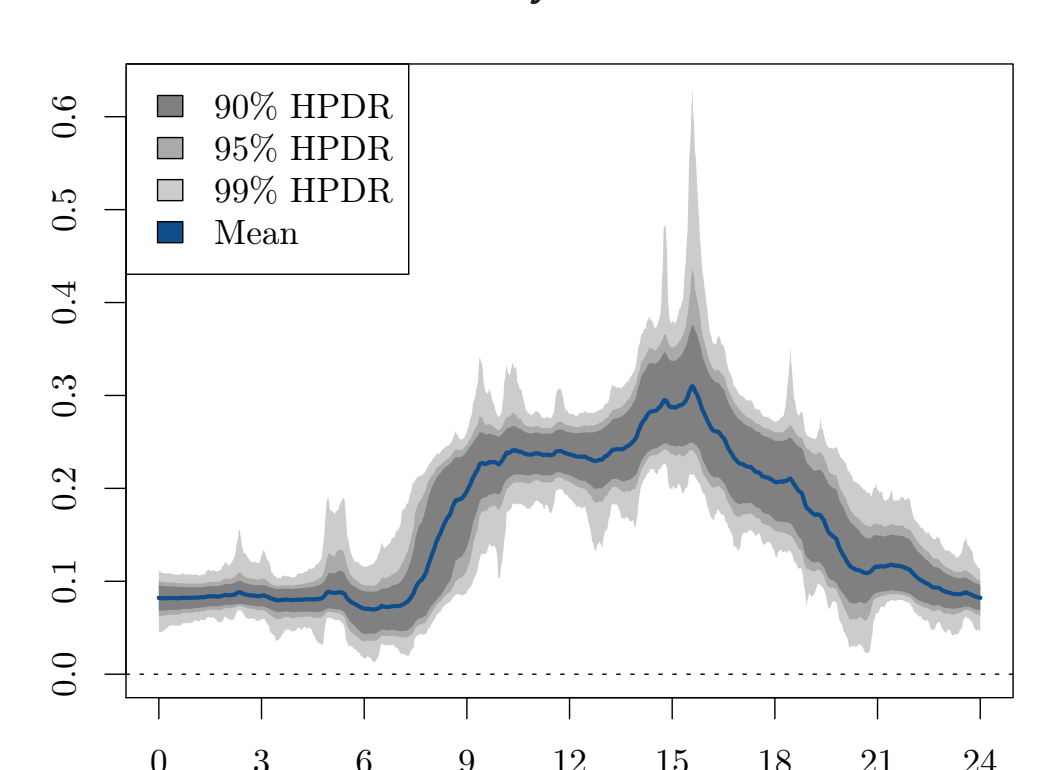


Figure 6: Estimated density of human events. Clear diurnal pattern, activity mostly concentrated in working hours.

6. Comments

- Simple algorithm to separate human and automated activity on a single edge in a computer network.
- Gibbs sampler with conjugate priors \rightarrow scalable to multiple edges and nodes across the entire network.
- Results on multiple real and simulated dataset show good performance of the model.

References and github



References:

- Heard, N.A, P.T.G. Rubin-Delanchy, and D.J. Lawson (2014), "Filtering automated polling traffic in computer network flow data". In: *Proceedings of the IEEE JISIC 2014*, pp. 268-271 (2014).

Link to the github repository in the QR code.