
Supplementary Materials for “An Alternative Prior Process for Nonparametric Bayesian Clustering”

Hanna M. Wallach

Department of Computer Science
University of Massachusetts Amherst

Shane T. Jensen

Department of Statistics
The Wharton School, University of Pennsylvania

Lee Dicker

Department of Biostatistics
Harvard School of Public Health

Katherine A. Heller

Engineering Department
University of Cambridge

1 Proof of Law for $\mathbb{E}(K_N | \text{UN})$

We start by defining $T_k = \inf \{m > T_{k-1}; X_m \notin \{X_1, \dots, X_{m-1}\}\}$. T_k is the “waiting time” (number of observations needed) until the k^{th} new cluster is generated by the uniform process. Under the uniform process, $T_k = \sum_{i=1}^k \tau_i$ where $\tau_i \sim \text{Geometric}(\theta / (\theta + i - 1))$ and the τ_i variables are independent, so

$$\mathbb{E}(T_k) = \sum_{i=1}^k \frac{\theta + i - 1}{\theta} = \frac{k^2}{2\theta} + k \left(1 - \frac{1}{2\theta}\right)$$

and

$$\begin{aligned} \text{Var}(T_k) &= \sum_{i=1}^k \frac{(\theta + i - 1)(i - 1)}{\theta^2} \\ &= \frac{k^3}{3\theta^2} + k^2 \frac{1}{2\theta} \left(1 - \frac{1}{\theta}\right) + k \frac{1}{2\theta} \left(\frac{1}{3\theta} - 1\right). \end{aligned} \quad (1)$$

In terms of T_k , $K_N = \max\{k; T_k \leq N\} = \sum_{k=1}^N \mathbb{I}(T_k \leq N)$. We first prove a strong law for the convergence of T_k . Let $\epsilon > 0$. From Chebychev’s inequality and (1), we have the following:

$$\mathbb{P}(|T_k - \mathbb{E}(T_k)| > \epsilon k^2) \leq \frac{\text{Var}(T_k)}{\epsilon^2 k^4} \leq \frac{C(\theta, \epsilon)}{k}. \quad (2)$$

From (2),

$$\mathbb{P}(|T_{k^2} - \mathbb{E}(T_{k^2})| > \epsilon k^4) \leq \frac{C(\theta, \epsilon)}{k^2},$$

and so by the Borel-Cantelli lemma, we have $\mathbb{P}(|T_{k^2} - \mathbb{E}(T_{k^2})| > \epsilon k^4) = 0$. Since $\epsilon > 0$ was chosen arbitrarily, it follows that $\frac{T_{k^2} - \mathbb{E}(T_{k^2})}{k^4} \rightarrow 0$ almost surely and hence $\frac{T_{k^2}}{k^4} \rightarrow \frac{1}{2\theta}$ almost surely. Now, let $m = \lfloor \sqrt{k} \rfloor$. Since T_k is increasing, we have:

$$\frac{T_{m^2}}{(m+1)^4} \leq \frac{T_k}{k^2} \leq \frac{T_{(m+1)^2}}{m^4}. \quad (3)$$

Since $\frac{m+1}{m} \rightarrow 1$, both sides of the inequality (3) converge to $(2\theta)^{-1}$ almost surely, and so

$$\frac{T_k}{k^2} \rightarrow \frac{1}{2\theta} \text{ almost surely.} \quad (4)$$

The strong law (4) implies a strong law for K_N as follows. $T_{K_N} \leq N < T_{K_N+1}$ and, consequently,

$$\frac{T_{K_N}}{K_N^2} \leq \frac{N}{K_N^2} < \frac{T_{K_N+1}}{K_N^2}.$$

Since $K_N \rightarrow \infty$ almost surely and $T_k / k^2 \rightarrow 1 / (2\theta)$ almost surely, it follows that the left and right hand side above both converge to $1 / (2\theta)$ almost surely. Thus, $K_N^2 / N \rightarrow 2\theta$ almost surely and so

$$\frac{K_N}{\sqrt{N}} \rightarrow \sqrt{2\theta} \text{ almost surely.} \quad (5)$$

From the strong law (5) and the dominated convergence theorem, we have the following:

$$\frac{\mathbb{E}(K_N)}{N} \rightarrow 0. \quad (6)$$

Combining (6) with following result from section 2,

$$\mathbb{E}(K_N^2) = \mathbb{E}(K_N) + 2\theta(N - \mathbb{E}(K_N)). \quad (7)$$

gives us

$$\frac{\mathbb{E}(K_N^2)}{N} \rightarrow 2\theta. \quad (8)$$

Finally, using (8) together with Fatou’s lemma and Jensen’s inequality, gives us the following:

$$\begin{aligned} \sqrt{2\theta} &\leq \liminf_{N \rightarrow \infty} \frac{\mathbb{E}(K_N)}{\sqrt{N}} \leq \limsup_{N \rightarrow \infty} \frac{\mathbb{E}(K_N)}{\sqrt{N}} \\ &\leq \limsup_{N \rightarrow \infty} \sqrt{\frac{\mathbb{E}(K_N^2)}{N}} = \sqrt{2\theta}. \end{aligned}$$

This then proves the result

$$\frac{\mathbb{E}(K_N)}{\sqrt{N}} \rightarrow \sqrt{2\theta}$$

under the uniform process.

2 Result relating $\mathbb{E}(K_N)$ to $\mathbb{E}(K_N^2)$

Recall the definition of T_k from above and now define $M_N = K_N + 1$. Consider the “waiting time” T_{M_N} until the observation that creates the $(K_N + 1)^{\text{th}}$ unique cluster. We relate $\mathbb{E}(K_N)$ to $\mathbb{E}(K_N^2)$ by calculating $\mathbb{E}(T_{M_N})$ in two different ways. First, observe that

$$\begin{aligned} \mathbb{E}(T_{M_N}) &= \mathbb{E}\left(\sum_{k=1}^{\infty} \tau_k \cdot \mathbf{I}(k \leq M_N)\right) \\ &= \frac{\theta - 1}{\theta} \sum_{k=1}^{\infty} \mathbb{P}(k \leq M_N) \\ &\quad + \frac{1}{\theta} \sum_{k=1}^{\infty} k \cdot \mathbb{P}(k \leq M_N) \\ &= \frac{\theta - 1}{\theta} \mathbb{E}(M_N) + \frac{1}{2\theta} \mathbb{E}(M_N(M_N + 1)), \end{aligned}$$

which, since $M_N = K_N + 1$, simplifies to

$$\mathbb{E}(T_{M_N}) = 1 + \mathbb{E}(K_N) \left(1 + \frac{1}{2\theta}\right) + \mathbb{E}(K_N^2) \frac{1}{2\theta}. \quad (9)$$

Now $T_{M_N} = N + \sum_j \mathbf{I}(M_{N+j} = M_N)$ and so $\mathbb{E}(T_{M_N}) = N + \sum_j \mathbb{P}(M_{N+j} = M_N)$ where

$$\begin{aligned} \mathbb{P}(M_{N+j} = M_N) &= \sum_k \mathbb{P}(T_k \leq N, N + j < T_{k+1}) \\ &= \sum_k \mathbb{P}(M_N = k + 1) \mathbb{P}(j < \tau_{k+1}). \end{aligned}$$

It follows that

$$\begin{aligned} \mathbb{E}(T_{M_N}) &= n + \sum_j \sum_k \mathbb{P}(M_N = k + 1) \mathbb{P}(j < \tau_{k+1}) \\ &= N + \sum_k \mathbb{P}(M_N = k + 1) \mathbb{E}(\tau_{k+1}) \\ &= N + \sum_k \mathbb{P}(K_N = k) \frac{k + \theta}{\theta}, \end{aligned}$$

which can be simplified to

$$\mathbb{E}(T_{M_N}) = N + 1 + \mathbb{E}(K_N) \frac{1}{\theta}. \quad (10)$$

Combining (9) and (10) gives (7):

$$\mathbb{E}(K_N^2) = \mathbb{E}(K_N) + 2\theta(N - \mathbb{E}(K_N)).$$

3 Evaluation Algorithm

The evaluation algorithm used to approximate $\log P(\mathcal{W}^{\text{test}} | \mathcal{W}^{\text{train}}, \mathbf{c}^{\text{train}}, \theta, \beta)$ is based on the “left-to-right” evaluation algorithm introduced by Wallach *et al.* (2009), adapted to marginalize out test cluster assignments. Pseudocode is given in algorithm 1.

Acknowledgements

This work was supported in part by the Center for Intelligent Information Retrieval, in part by CIA, NSA and NSF under NSF grant #IIS-0326249, and in part by subcontract #B582467 from Lawrence Livermore National Security, LLC, prime contractor to DOE/NNSA contract #DE-AC52-07NA27344. Any opinions, findings and conclusions or recommendations expressed in this material are the authors’ and do not necessarily reflect those of the sponsor.

References

Wallach, H., Murray, I., Salakhutdinov, R., and Mimno, D. (2009). Evaluation methods for topic models. In *26th International Conference on Machine Learning*.

```

initialize  $l := 0$ 
for each document  $d$  in  $\mathcal{W}^{\text{test}}$  do
  initialize  $p_d := 0$ 
  for each particle  $r = 1$  to  $R$  do
    for  $d' < d$  do
       $c_{d'}^{(r)} \sim P(c_{d'}^{(r)} \mid \mathcal{W}_{<d}^{\text{test}}, \{c_{<d}^{(r)}\}_{\setminus d'}, \mathcal{W}^{\text{train}}, \mathbf{c}^{\text{train}}, \theta, \beta)$ 
    end for
     $p_d := p_d + \sum_c P(\mathbf{w}_d^{\text{test}}, c_d^{(r)} = c \mid \mathcal{W}_{<d}^{\text{test}}, \mathbf{c}_{<d}^{(r)}, \mathcal{W}^{\text{train}}, \mathbf{c}^{\text{train}}, \theta, \beta)$ 
     $c_d^{(r)} \sim P(c_d^{(r)} \mid \mathbf{w}_d^{\text{test}}, \mathcal{W}_{<d}^{\text{test}}, \mathbf{c}_{<d}^{(r)}, \mathcal{W}^{\text{train}}, \mathbf{c}^{\text{train}}, \theta, \beta)$ 
  end for
   $p_n := p_n / R$ 
   $l := l + \log p_n$ 
end for
 $\log P(\mathcal{W}^{\text{test}} \mid \mathcal{W}^{\text{train}}, \mathbf{c}^{\text{train}}, \theta, \beta) \simeq l$ 

```

Algorithm 1: “Left-to-right” evaluation algorithm for computing $\log P(\mathcal{W}^{\text{test}} \mid \mathcal{W}^{\text{train}}, \mathbf{c}^{\text{train}}, \theta, \beta)$.