

# A Bayesian algorithm to identify zones of shared evolution in space

Peter Ranacher<sup>a</sup>

<sup>a</sup>*Department of Geography, University of Zurich*

---

## Abstract

TBD

*Keywords:* TBD

---

## 1. Introduction

## 2. Related Work

## 3. Methodology

Human traits, such as language, material culture, religion, mythology or genes diffuse in geographical space either through expansion or contact between societies. *Sbayes* analyses the spatial distribution of human traits and identifies zones of shared evolution. The algorithm proposes zones – geographical areas where traits might have diffused – (Section 3.1) and evaluates these against evidence from the data (Section 3.2). Thus, *Sbayes* computes a posterior distribution of zones where human traits are likely to have diffused (Section 3.3).

### 3.1. Proposing zones of shared evolution

*Sbayes* connects the spatial locations of all societies in the data by a Delaunay triangulation. The resulting auxiliary graph creates an adjacency between the societies making it easy to propose spatial zones in an efficient and fast way. *Sbayes* proposes an initial random zone - a connected subgraph in the auxiliary graph consisting of five societies (Figure 1a). Starting from the initial zone, the algorithm generates new random candidate zones by applying one of the following three steps:

- 20 – Grow step: A random neighbor is added to the zone (Figure 1b). Neigh-  
21 bors are societies adjacent to the current zone.
- 22 – Shrink step: A random node is removed from the zone (Figure 1c).  
23 The resulting candidate zone is not necessarily a connected subgraph  
24 of the auxiliary graph.
- 25 – Swap step: A shrink step followed by a grow step (Figure 1d).

### 26 3.2. The likelihood of a zone

27 *Sbayes* evaluates the likelihood of a zone in terms of evidence for shared  
28 evolution (feature likelihood) and spatial proximity (geo-likelihood). We im-  
29 plement two models :

- 30 – The *Generative model (GM)* assesses the likelihood of a zone to result  
31 from a common generative process.
- 32 – The *Particularity model (PM)* assesses the likelihood of a zone to be  
33 *distinctive*.

#### 34 Feature likelihood

35 *Sbayes* requires traits to be binary: a trait is either *present* and can be  
36 found in a society, or it is *absent* and cannot be found. We assume that only  
37 presence is indicative of a possible shared evolution.

38 In *GM* ... (Nico).

39 In *PM*, *Sbayes* evaluates how exceptional the shared evolution in the zone  
40 is, taking into account both global probabilities and the local presence of a  
41 trait. Let us assume a zone of size  $n$ , with the trait  $f$  present in  $k$  sites. The  
42 feature likelihood equals the inverse of the probability of finding a zone which  
43 is at least as *distinctive* in terms of evidence for shared evolution (Figure 2,  
44 red part of the histogram). Hence, in *PM* the feature likelihood is defined as

$$\mathcal{L}_f^{PM} = \frac{1}{Pr(K \geq k)}, \quad (1)$$

45 where  $Pr(K \geq k)$  is the probability that the trait is present in  $k$  societies  
46 in a random sample of size  $n$ , given the global probability of presence  $p$ . In a  
47 sense, *PM* emulates the implicit reasoning of a domain expert who attempts  
48 to determine the particularity of a sample. If  $f$  is hardly present there is

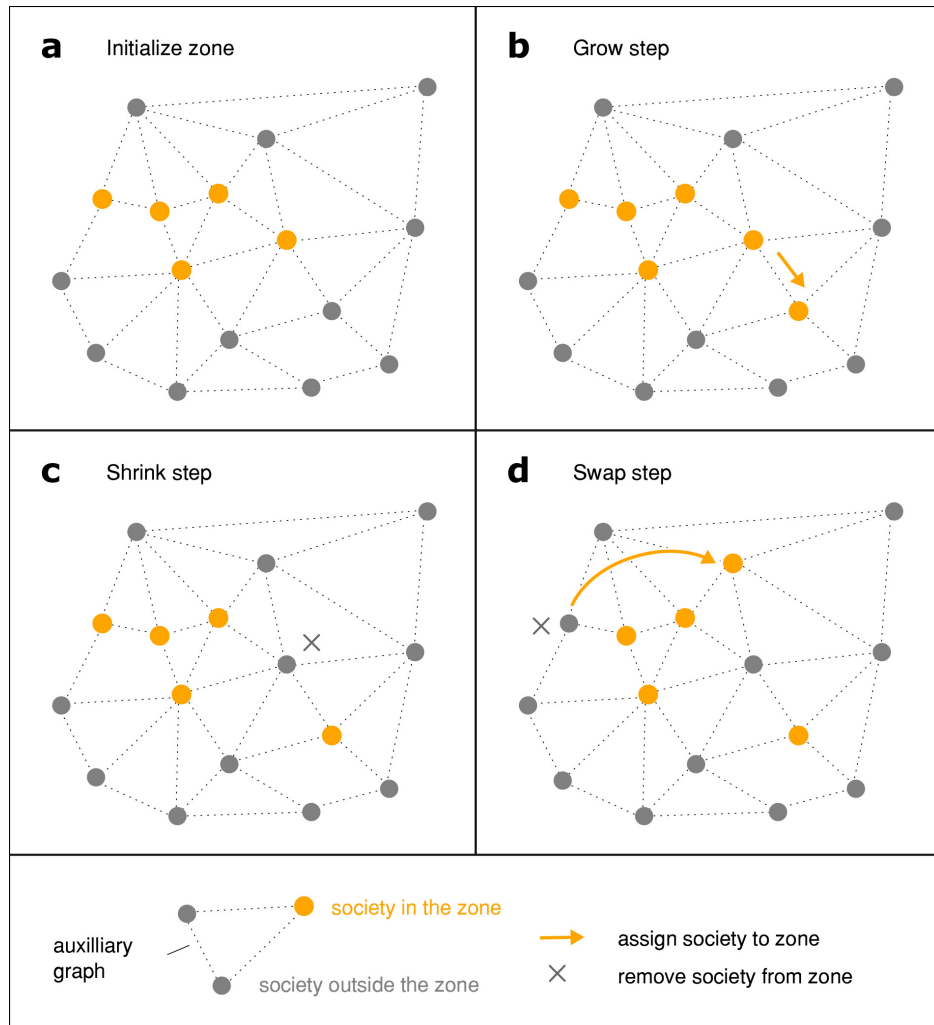


Figure 1: Proposing contact zones

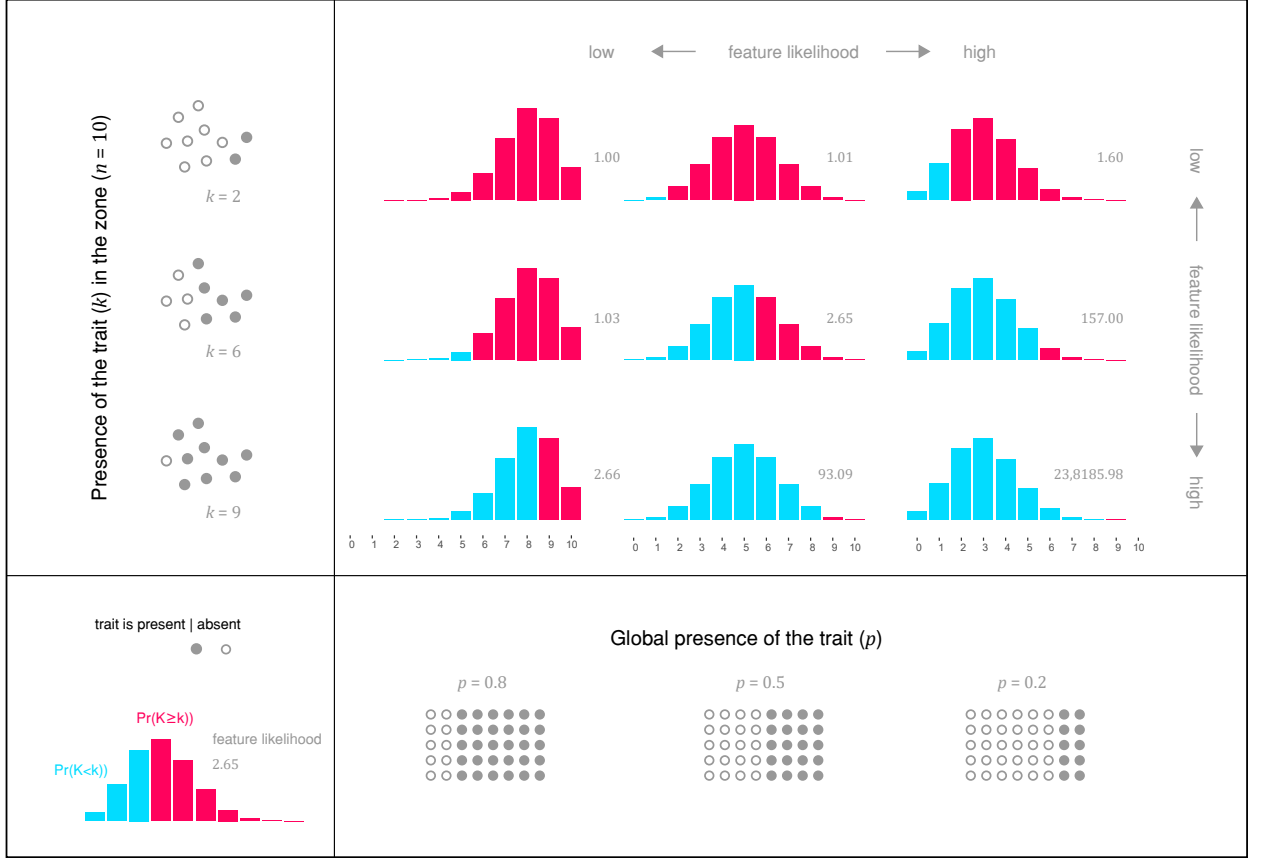


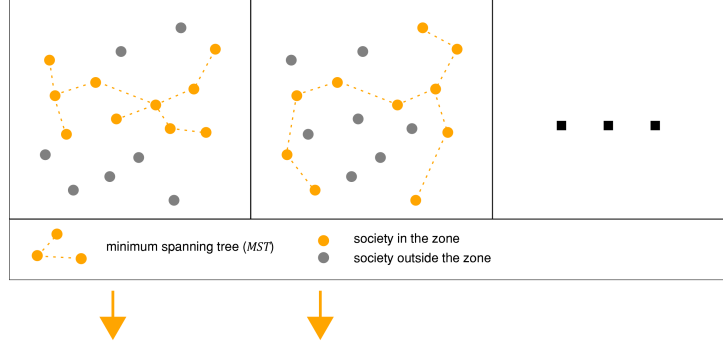
Figure 2: The feature likelihood in the *PM*

little to no evidence for shared evolution (Figure 2,  $k = 2$ ). If  $f$  is present the evidence is weaker for traits which are globally present (Figure 2,  $k = 9$  and  $p = 0.8$ ) and stronger for traits which are globally absent (Figure 2,  $k = 9$  and  $p = 0.2$ ). Simply put, it is not exceptional to find a common trait in a zone, whereas it is exceptional to find a rare one.

Usually, data will consist of several traits  $f_1, f_2, \dots, f_m$ . Each trait represents one piece of evidence for shared evolution. The overall feature likelihood is simply the product of the feature likelihood of all  $m$  independent traits:

$$\mathcal{L}_F = \mathcal{L}_{f_1} \cdot \dots \cdot \mathcal{L}_{f_m} \quad (2)$$

**a** Generate 10,000 random samples of size  $n = 10$



**b** Derive proximity distribution ( $D$ ) for zones of size  $n = 10$

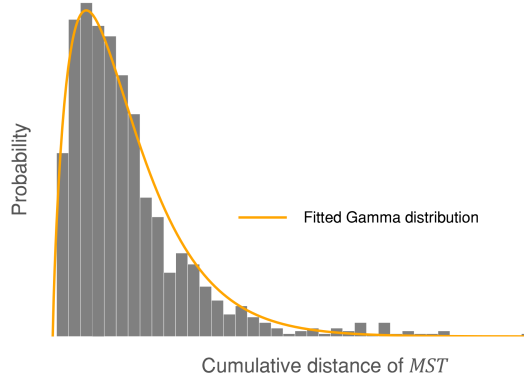


Figure 3: The proximity distribution for zones of size  $n=10$

### 58 *Geo-likelihood*

59 The geo-likelihood quantifies the spatial proximity of the societies in a  
60 zone.

61 In *GM*, ... (Nico)

62 In *PM*, *Sbayes* finds the minimum spanning tree (*MST*) of the zone and  
63 evaluates it against the proximity distribution ( $D$ ).  $D$  is derived empirically  
64 from the data for each sample size  $n$ : First, 10,000 random zones of size  $n$   
65 are generated. For each zone the *MST* is found and its cumulative distance  
66 is computed (Figure 3, a). Finally a gamma distribution is fitted to the  
67 resulting sample, which yields  $D$  (Figure 3, b). To compute the geo-likelihood  
68 *Sbayes* finds the *MST* of the societies in the zone, retrieves its cumulative  
69 distance  $d$ , and evaluates it against the proximity distribution for zones of

70 size  $n$ . Thus, in  $PM$  the geo-likelihood is defined as

$$\mathcal{L}_{geo}^{PM} = Pr(D \leq d). \quad (3)$$

71 The  $MST$  connects all societies in a zone minimizing the distance between  
 72 these. It reflects the minimum effort necessary for the societies in the zone to  
 73 diffuse or to have contact. The geo-likelihood shows how great this effort is,  
 74 taking into account the global proximity in the data. The geo-likelihood is  
 75 high for close societies (Figure 4, a and b) and low for distant ones (Figure 4,  
 76 c and d). Since the geo-likelihood derives from a cumulative distance measure  
 77 it is robust to single outliers.

### 78 3.3. Markov chain Monte Carlo Sampling

79 *Sbayes* computes the posterior probability for shared evolution in a zone,  
 80 according to Bayes' theorem:

$$Pr(\text{zone}|\text{data}) = \frac{Pr(\text{data}|\text{zone}) \cdot Pr(\text{zone})}{Pr(\text{data})}. \quad (4)$$

81 In equation 4,  $Pr(\text{zone})$  is the prior probability for shared evolution in a  
 82 zone. We assume that  $Pr(\text{zone})$  is uniform, i.e. a priori, shared evolution is  
 83 equally probable in each zone.  $Pr(\text{data}|\text{zone})$  is the combined feature and  
 84 geo-likelihood:

$$Pr(\text{data}|\text{zone}) = \mathcal{L} = \mathcal{L}_F^{w_F} \cdot \mathcal{L}_{geo}^{w_{geo}}, \quad (5)$$

85 where  $w_F$  and  $w_{geo}$  are weights applied to either likelihood.  $Pr(\text{data})$  is the  
 86 marginal likelihood of the model and acts as a normalizing constant of the  
 87 posterior.

88 There is no analytical solution for equation 4. *Sbayes* performs Markov-  
 89 chain-Monte Carlo (MCMC) sampling to calculate a numerical approxima-  
 90 tion. The MCMC first generates a random initial zone  $z$  and evaluates the  
 91 likelihood  $\mathcal{L}(z)$ . Then, a candidate zone  $z'$  is proposed, by either shrinking,  
 92 growing or swapping, and the candidate likelihood  $\mathcal{L}(z')$  is computed. The  
 93 MCMC evaluates the acceptance probability

$$\alpha = \frac{\mathcal{L}' \cdot Q(z|z')}{\mathcal{L} \cdot Q(z'|z)}, \quad (6)$$

94 where  $Q(z|z')$  is the proposal probability to move from the current zone  $z$   
 95 to the candidate zone  $z'$  and  $Q(z'|z)$  is the corresponding back probability

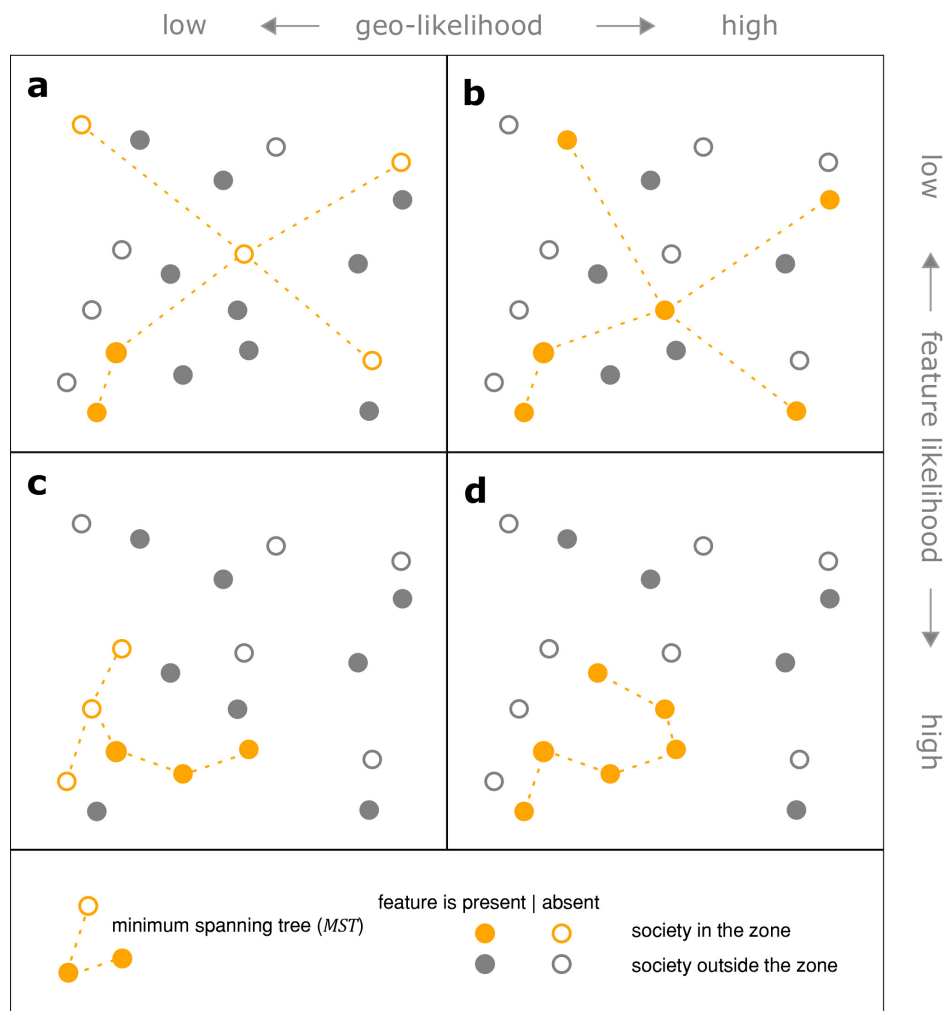


Figure 4: Feature and Geo-likelihood

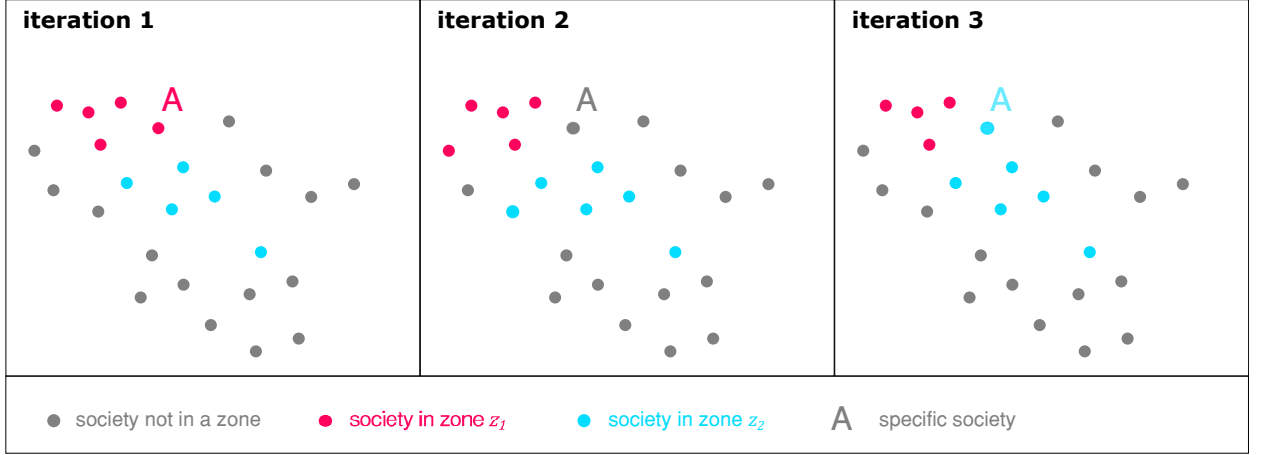


Figure 5: Parallel analysis

to move from  $z'$  to  $z$ . A move to the candidate zone is accepted with probability  $\min(1, \alpha)$  and rejected otherwise. The MCMC repeatedly proposes new candidate zones, evaluates the acceptance probability and accepts or rejects a move. The resulting Markov chain gives an approximation of the posterior probability for zones of shared evolution in space. For numerical convenience the MCMC uses the log-likelihood rather than the likelihood (*Missing: explain briefly why this does not matter*).

### 3.4. Parallel analysis

The diffusion of traits in space is a complex process in which expansion and contact might create multiple and possibly overlapping areas of shared evolution. Hence, it is not realistic to observe a posterior probability with a single mode resulting from one zone of shared evolution. *Sbayes* uses parallel zones to explore multiple modes in the data and to identify complex and possibly intertwined patterns of shared evolution.

In the parallel analysis, the MCMC initializes  $l$  individual Markov chains, to evaluate  $l$  parallel zones. Each chain explores the data, i.e. proposes and evaluates zones independently. However, a society can only belong to a single zone at a time. In Figure at iteration 1, the society  $A$  belongs to  $z_1$ . Only after  $A$  is removed from  $z_1$  (iteration 2) it can be assigned to  $z_2$  (iteration 3). The number of parallel zones is fixed.



116 *Missing: a method to estimate an appropriate number of parallel zones ( $l$ )*

- 117     • Compute the marginal likelihood of each chain
- 118     • Based on the marginal likelihood define a statistical procedure to esti-
- 119         mate  $l$  (AIC, BIC, scree plot, ...)

#### 120 **4. Test and evaluation**

121     In this section we test *Sbayes* on simulated data and evaluate its perfor-

122     mance.

123 *Missing: a structured approach to test the algorithm on simulated data*

- 124     • How do we simulate zones?
  - 125         Draw zones by hand? Make an algorithm propose zones? Curdin?
  - 126         Nico?
- 127     • Single mode
  - 128         We simulate shared evolution in space and hope that the algorithm
  - 129         finds it. We can vary:
    - 130             – the approach (GM or PM)
    - 131             – the size of the simulated zones (i.e. number of societies)
    - 132             – their shape (i.e. elongated, compact, tree-like, ...)
    - 133             – the intensity in terms of how many features are indicative of a
    - 134                 shared evolution
    - 135             – the intensity in terms of how many societies in the simulated zones
    - 136                 are indicative of a shared evolution
- 137     • Multi-mode (parallel analysis)
  - 138         – all of the above
  - 139         – the number of simulated zones
  - 140         – the degree of how much the zones are intertwined
- 141     • Define a test statistics to evaluate the performance of the algorithm
- 142         (Curdin?), either quantitative or qualitative

## 143 5. Empirical Analysis

144 *Missing: Rik and Manuel, should we test the algorithm on the Amazon data,*  
145 *on the Asian data or on both? - It is strongest if we do both - Rik*

146 This depends a bit on the consecutive papers that we plan to publish and  
147 the story that we might or might not want to anticipate.

## 148 6. Discussion

- 149 • Explain how the marginal likelihood of a model can be used for spatial  
150 inference.
- 151 • Compare the algorithm to DBSCAN, Structure, TESS, the geo-model  
152 in BayesTraits ...
- 153 • Explain the difference between GM and PM.
- 154 • Explain how the results of the algorithm depend on the given set of  
155 societies and the given spatial extent.