

Kernel Approximate Bayesian Computation for Population Genetics

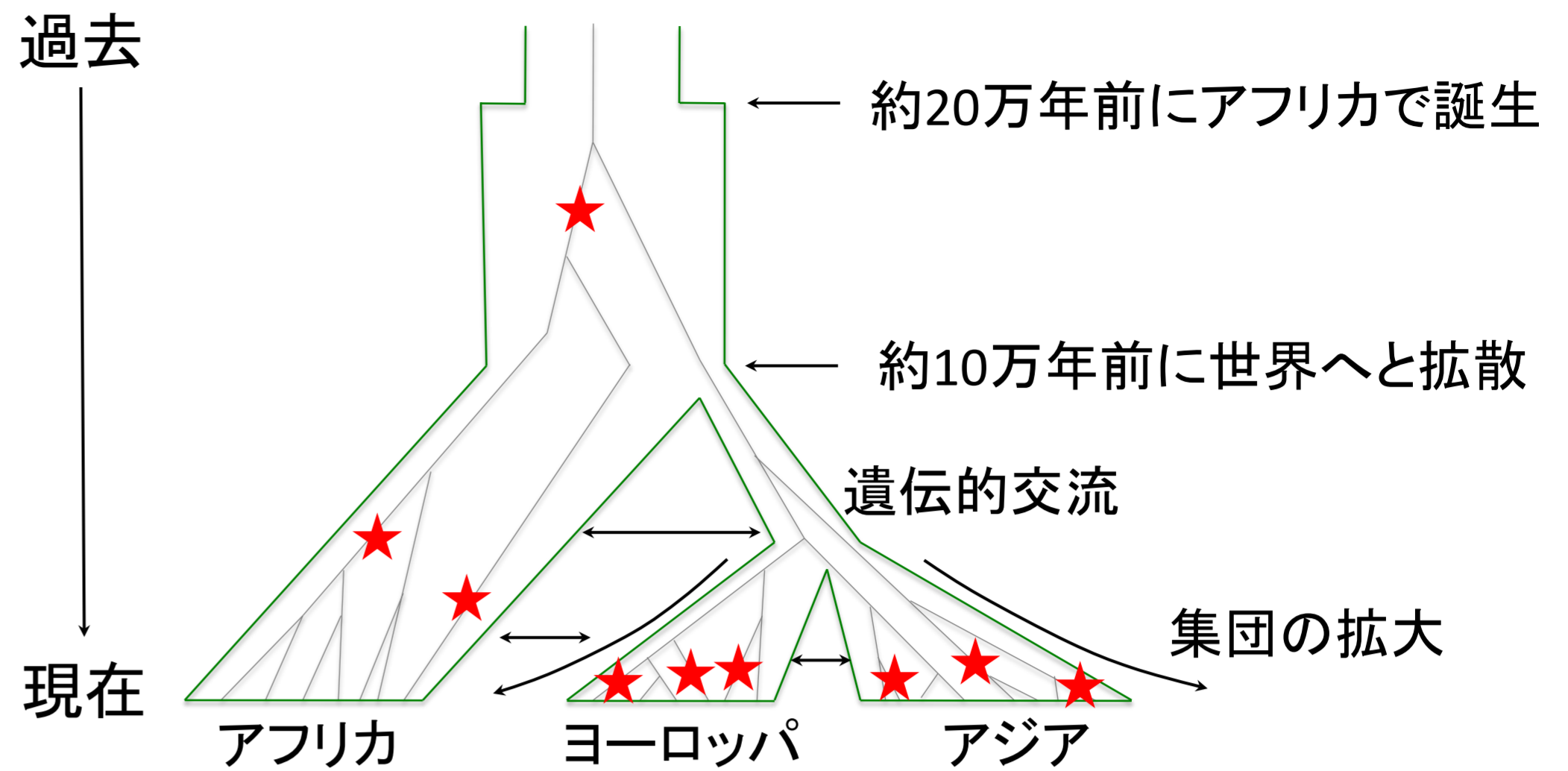
中込 滋樹 リスク解析戦略研究センター 外来研究員

【背景】



現在の遺伝的多様性から過去の歴史を推定する

人類集団の進化モデル



【方法】

ベイズ推論

$$\pi(\theta|D) \propto f(D|\theta)\pi(\theta)$$

D : data, θ : parameters

● 尤度関数 $f(D|\theta)$ が既知の場合

[A] Rejection-sampling method (Ripley 1987)

- A1. Generate θ_i from $\pi(\cdot)$.
- A2. Accept θ_i with probability $f(D|\theta_i)$, and go to A1.

● 尤度関数 $f(D|\theta)$ が未知の場合

[B] Approximate Bayesian Computation (ABC) method (Beaumont et al. 2002)

- B1. Generate θ_i from $\pi(\cdot)$.
- B2. Simulate data \mathcal{D}_i by the model using θ_i .
- B3. Accept θ_i if $d(\mathcal{D}, \mathcal{D}_i) < \delta$, and go to B1.

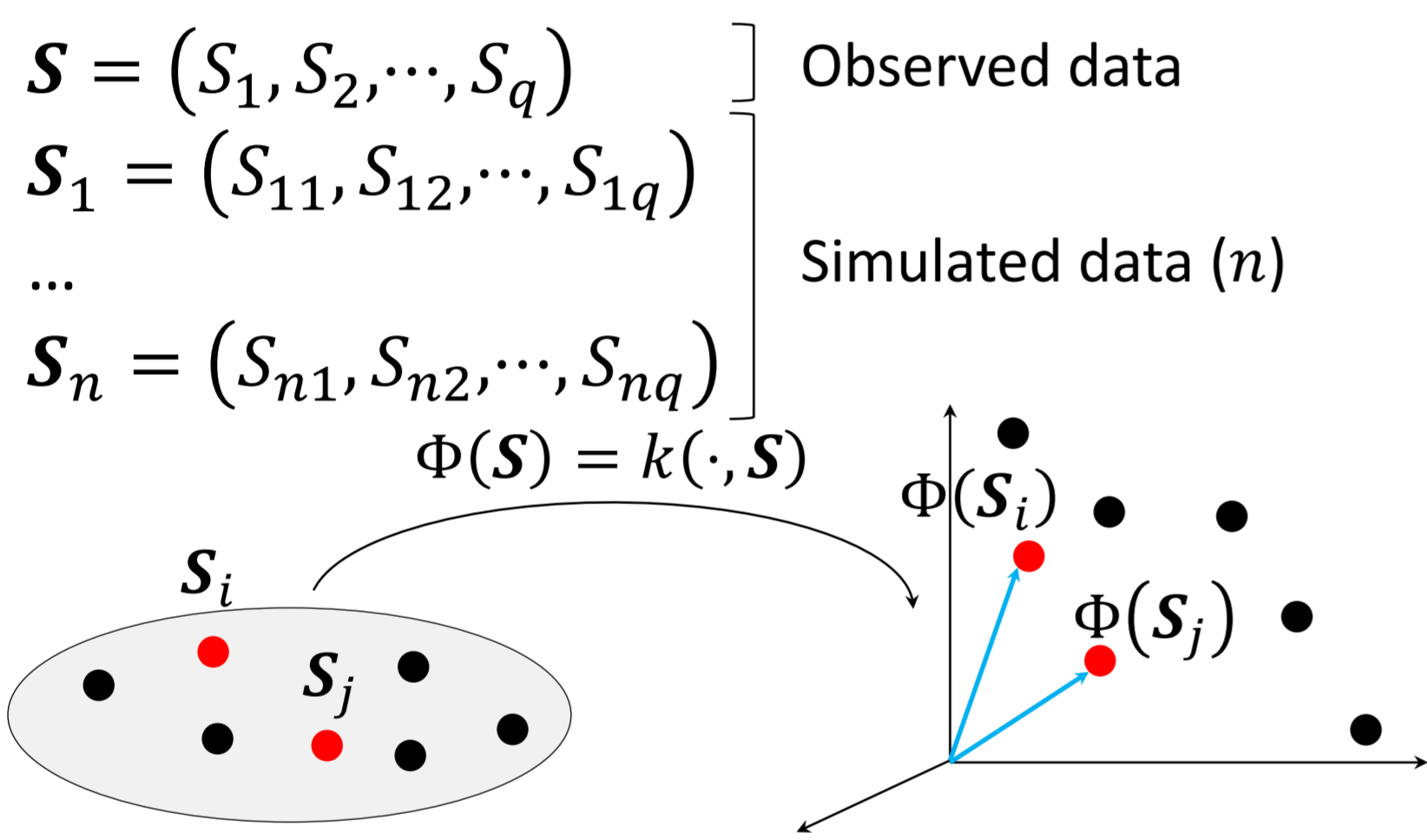
→ B3'. Accept θ_i if $d(s, s_i) < \delta$.

s : summary statistics for \mathcal{D}

d : metric to measure dissimilarity between s and s_i

カーネル法を用いたABC methodの改良

Nakagome et al. (2012)



Ω : space of summary statistics \mathcal{H}_S : reproducing kernel Hilbert space

The inner product between mappings is given by

$$\langle \Phi(s_i), \Phi(s_j) \rangle = k(s_i, s_j)$$

[C] Kernel ABC method

- C1. Generate θ_i from $\pi(\cdot)$.
 - C2. Simulate data \mathcal{D}_i by the model using θ_i .
 - C3. Compute s_i for \mathcal{D}_i , and return to C1.
- $\{(\theta_i, s_i)\}_{i=1}^n$

Kernel Bayes Rule (Fukumizu et al. 2011)

The empirical estimator of the kernel posterior mean is given by

$$\hat{m}_{\theta|s} = \sum_{i=1}^n w_i k(\cdot, \theta_i),$$

The weighted coefficient is given by

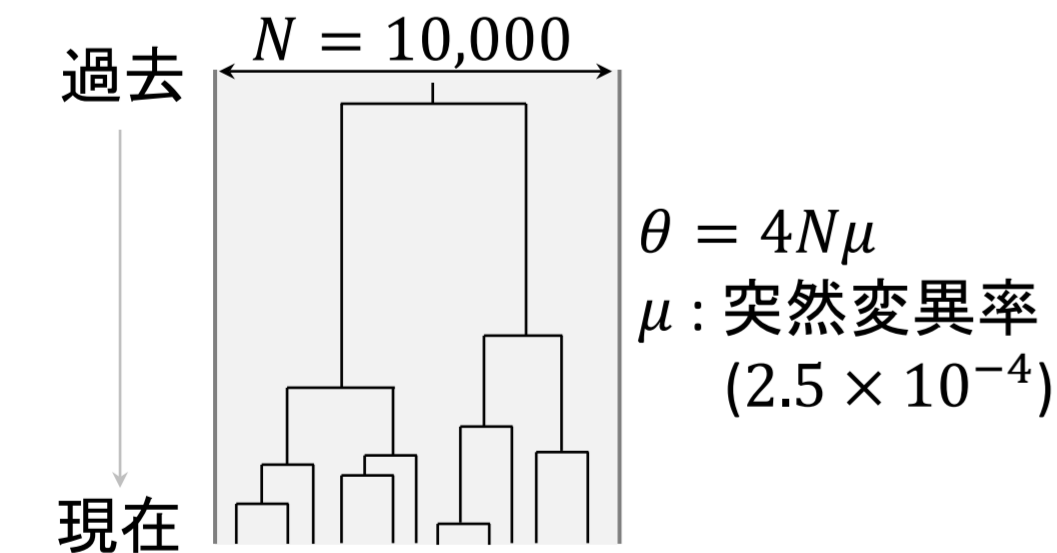
$$w_i(s) = \sum_{i=1}^n (G_S + n\epsilon_n I_n)^{-1}_{ij} k(s_j, s),$$

G_S : Gram matrix $(k(s_i, s_j))_{i,j=1}^n$

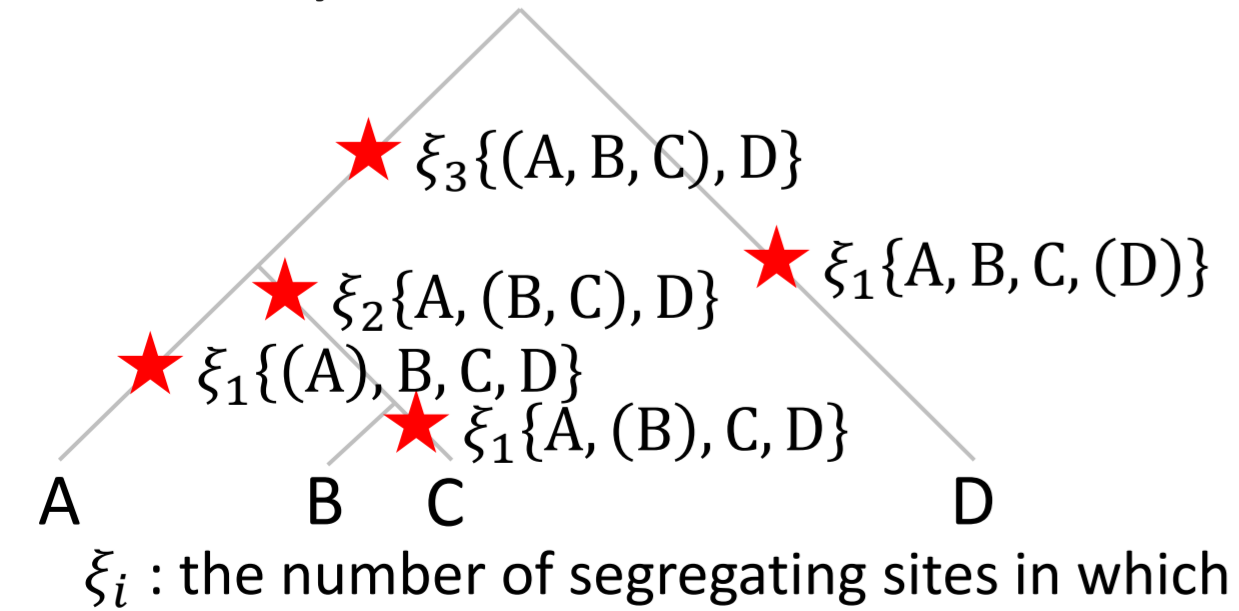
An estimator of posterior expectation of a function $f(\theta)$ is given as

$$E[f(\theta)|s] = \langle f(\cdot), \hat{m}_{\theta|s} \rangle_{\mathcal{H}_S} = \sum_{i=1}^n w_i f(\theta_i).$$

Evolutionary model



Summary statistics



1. The number of segregating sites (S_{Seg})
→ $\sum_{i=1}^3 \xi_i = 5$
2. Site frequency spectrum (S_{SFS})
→ $(\xi_1, \xi_2, \xi_3) = (3, 2, 1)$

【結果】

- $f(D|\theta)$ is computed by importance sampling (Griffiths, 2007),
- Gaussian RBF kernel
 $k(x, y) = \exp\left(-\frac{1}{2\sigma^2} \|x - y\|^2\right)$,
(σ : median of pairwise Euclidean distances)

Table 1. Comparison of posterior estimates of θ given \mathcal{D} , S_{SFS} , and S_{Seg} .

	$m_{\theta \mathcal{D}}^*$	$\hat{m}_{\theta S_{SFS}}^{**}$	$\hat{m}_{\theta S_{Seg}}^{**}$
Mean	10.498	10.510	9.677
S.D.	0.067	0.044	0.041

*The posterior mean given \mathcal{D} is generated by the rejection-sampling method.

**The kernel posterior means are obtained from 16,000 simulated samples.

【引用文献】

- Ripley, B. D. (1987). Stochastic simulation. John Wiley & Sons, New York.
- Beaumont, M. A. et al. (2002). Approximate Bayesian computation in population genetics. *Genetics* 162, 2025-2035.
- Nakagome, S. et al. (2012) Kernel Approximate Bayesian Computation for Population Genetics. arXiv: 1205.3246.
- Fukumizu, K. et al. (2011). Kernel Bayes' rule. *Advances in Neural Information Processing Systems* 24, 1549-1557.
- Griffiths, R. C. (2007). GENETREE version 9.0 <http://www.stats.ox.ac.uk/~griff/software.html>.

