

Title: Integrated species distribution model improves spatial uncertainty of eDNA of Edomae-fish in Tokyo Bay

Yuki Kanamori^{1*}, Hiroshi Okamura², Shota Nishijima², Yuki Hongo², Yasuyuki Uto³, Hisatoku Mita⁴, Mitsuhiyo Ishii⁴, Kiyoharu Akimoto⁵, and Akane Kusano⁶

¹ Fisheries Resources Institute, Japan Fisheries Research and Education Agency, 25-259 Shimomekurakubo, Samemachi, Hachinohe, Aomori 031-0841, Japan

² Fisheries Resources Institute, Japan Fisheries Research and Education Agency, 2-12-4 Fukuura, Kanazawa, Yokohama, Kanagawa 236-8648, Japan

³

⁴

⁵

⁶

* Corresponding author

Email: kana.yuki@fra.affrc.go.jp

Abstract**Keywords**

Spatial distribution, Environmental DNA, data fusion, Spatial correlation, Ecology of eDNA

1 Introduction

Understanding of spatial distribution of species and underlying its mechanism is a major goal in ecology. Field surveys using environmental DNA (eDNA) are widely used for monitoring such as endangered species, invasive species (e.g., Dougherty et al. 2016; Larson et al. 2020), and biodiversity because the surveys of eDNA are easy to detect occurrence of target species, non-invasiveness, and high cost effectiveness rather than previous direct sampling method (Rees et al. 2014; Thomsen & Willerslev 2015). However, uncertainty of the occurrence of eDNA is a issue in eDNA studies because various factors, such as biological and environmental features, intricately affect the shedding (i.e., production of eDNA from organisms) and degradation (i.e., decay of eDNA from a system) of eDNA (Fig. 1). Therefore, it is necessary to consider the spatial uncertainty included in eDNA to infer the spatial distribution of species from eDNA occurrence data.

After the suggestion of the importance of understanding about origin, state, transport, and fate of eDNA ("ecology of eDNA") by Barnes & Turner (2016), the studies which used laboratory experiments and numerical hydrological models have been increased to overcome the uncertainty of eDNA. For example, the relationship between production rate of eDNA and biological factors (e.g., species and body mass) or environmental factors (e.g., temperature), and between decomposition of eDNA and the state of eDNA (e.g., length of eDNA) or environmental factors were measured. In numerical hydrological models, spatial distribution of eDNA was simulated in aquatic area and obtain the production, transport and decay of eDNA (e.g., Fukaya et al. 2020). However, keeps for all focal species in laboratory are not realistic and low cost-efficiency when multi species studies such as biodiversity

23 monitoring and fisheries management. Moreover, biologists monitoring species are not often
24 familiar with hydrological models and the simulation of eDNA seems to be hard for them.

25 Integrated species distribution models (IDMs) are now common spatial model to
26 predict spatial pattern of species (Issac et al. 2020). The models combine the different type
27 of data with strengths and weaknesses in a single model. For example, the model can use
28 both scientific survey data, which is high quality but less abundant due to restriction of cost,
29 and citizen data, which is widely collected and abundant but may be low quality due to not
30 using consistent field methods in a model. Combining both types of data can capitalize on
31 the strengths of each data and perform better prediction than models when using single data
32 (Pacifici et al. 2017; Miller et al. 2019). Hence, if we develop a IDM that uses not only the
33 occurrence data of eDNA but also the second data which contains the information that
34 species exists at a location, as well as environmental factors as covariance in the model, then
35 the inference from the model should represent the reduction of the spatial uncertainty of
36 eDNA and improvement (Fig. 1).

37 In this paper, to infer spatial distribution of species from the occurrence data of eDNA,
38 we first make a spatial distribution model which considers spatial uncertainty of eDNA by
39 using an integrated spatial distribution model. We then applied the model to both eDNA data
40 and catch weight data for four fish in Tokyo Bay, Japan, and predicted the spatial
41 distribution of these fish. We finally compared the differences in (i) the predicted spatial
42 distribution of eDNA among fish and (ii) the relationship between the occurrence of eDNA
43 and environmental factors among fish to discuss the effects of the biological factors on the
44 spatial uncertainty of eDNA.

45

46 **2 Materials and Methods**

47 **2.1 A general model to estimate species distribution from eDNA**

48 **Features of integrated spatial distribution model**

49 Integrated spatial distribution model that account for explicitly spatial autocorrelation in
50 occurrence were built by Pacifici et al. (2017), which shows three approaches to predict the
51 spatial distribution of species: the joint likelihood (shared), correlation, and covariate
52 methods. The joint likelihood method uses multiple data types to simultaneously estimate a
53 shared set of parameters with constraining that the likelihoods of shared set of parameters to
54 be equal across. The correlation method connects multiple data types indirectly through a
55 shared covariance matrix that captures similar patterns present in each data sources. The
56 covariate method incorporates information from a added dataset via a fixed effect.

57 Although each methods estimate the spatial distribution of species using multiple data
58 sets, we need to select method depending on the data features for analysis because there are
59 strengths and weaknesses (Pacifici et al. 2017; Miller et al. 2018). The joint likelihood
60 method may be problematic when the second data is of poorly quality compared to
61 correlation and covariate methods because each data can directly inform the latent
62 occurrence state (probabilities?) and the weight given to estimate the parameters is naturally
63 determined by their relative size and quality. Thus, it is not the best method when our
64 second data is low quality while it is the best method when our second data is high quality
65 (vise versa). The correlation method is added robustness to the joint likelihood because the
66 second data indirectly inform the occurrence state. Thus, it is the best method when our
67 second data is low quality while it is inferior to the joint likelihood method when both data

are deemed reliable. The covariate method does not make full use of the information in the second data because the second data as a constructed covariate in the mean occurrence state. In addition, this method can reduce the computational cost because there are fewer parameters to estimate and the number of data locations can be reduced. Thus, it is the best method when the second data is low quality and/or there is computational limitation while it may not the best method when the information of the second data is needed.

The model for eDNA

When predicting the spatial distribution of species from eDNA using integrated species distribution model, the information that a species exists is needed as second data to consider spatial uncertainties of eDNA due to complex factors (Fig. 1). Hence, the second data is preferred to high quality as possible.

しかし、eDNA は直接的なモニタリングに比べて簡易的であるためより広い範囲で取得されている可能性が高く、eDNA のデータと同様の空間範囲で調査データのように質の高いデータを取得することは難しいかもしれない。その一方で、eDNA の空間的な不確実性を考慮するためには、種がいた証拠である 2 番目のデータの情報を eDNA のデータにしっかりと伝える必要がある。これらを考えると、integrated spatial distribution model を用いた eDNA からの空間分布の推定には、以下のような correlation method が適切である：

$$\begin{aligned} p_e(s_i) &= \alpha_e + \sum_k f_{e,k}(x_{e,k}(s_i)) + w\theta(s_i) + u_e(s_i) \\ p_a(s_i) &= \alpha_a + \sum_k f_{a,k}(x_{a,k}(s_i)) + \theta(s_i) + u_a(s_i) \end{aligned} \tag{1}$$

where α and $x_k(s_i)$ are the intercept and the covariates at sites i for occurrence probabilities at sites i of the added data (p_a) and eDNA data (p_e), respectively. $u(s_i)$ is spatial error that is

specific for each data following multivariate normal distributions $MVN(0, \mathbf{R})$, where the variance–covariance matrix \mathbf{R} is a Matérn correlation function. θ which is shared between two equations is the common spatial pattern between the two data, which cannot explain by each terms of the equations. That is, θ can be interpreted as "true" spatial distribution of species.

2.2 An application to a eDNA and catch data in Tokyo Bay

2.2.1 eDNA data

Field surveys

Field surveys were conducted by prefectural experimental station in Chiba, following the consistent sampling design at 14 sites in Tokyo Bay from April to December in 2018 (Fig. 1). In each sites, seawater and environmental data were simultaneously collected. For eDNA analysis, two litter of bottom seawater was collected using a Niskin water sampler, and then it was separated for two 1L samples for replicate. Each samples filtered glass fiber membrane GF/F ($0.7 \mu m$ pore size; Cytiva, Sheffield, UK) onboard and then the filters were frozen on a block of dry ice. These frozen filters were stored at -30° in the laboratory until eDNA extraction. To lower the levels of cross-contamination, equipments for eDNA sampling were changed new one or washed in each sites. During sampling the bottom seawater, seawater temperature, salinity, pH, and dissolved oxygen (DO) at the same depth of seawater sampling for eDNA were measured by CTD (メーカ一).

Laboratory experiments

In laboratory, eDNA extraction, eDNA amplification, and eDNA sequence were conducted.

109 Total eDNA was extracted from the frozen filters using a DNeasy Blood and Tissue Kit
110 (Qiagen, Hilden, Germany) following Yamamoto et al. 2019. Mitochondrial 12S rRNA
111 gene was amplified using MiFish universal primers referring to Miya et al. 2015 with slight
112 modification. The details was shown in Hongo et al. (受理されてないようだったら書くし
113 かない). eDNA sequence were

114 2.2.2 Catch statistics

115 A part of catch statistics of small-scale bottom trawl fisheries recorded by several
116 representative boats of Chiba Prefecture were provided by Chiba Prefecture. This data
117 included date, geographic location, efforts (number of tows), gear, and catch weight (kg) in
118 each fish. Almost of all gear was beam trawl although dredge net also used. The species
119 which also detected by eDNA was *Conger myriaster* (マアナゴ), *Kareius bicoloratus* (イシ
120 ガレイ), *Lateolabrax japonicus* (スズキ), and *Konosirus punctatus* (コノシロ). Thus, we
121 estimated the spatial distribution of these four species using the eDNA-IDM. マコガレイ,
122 カマス類, クロダイ, イシモチ類も解析できる??

123 2.2.3 Estimation of spatial distribution

124 To estimate the spatial distribution of four focal species using eDNA and catch data by
125 considering with spatial uncertainties of eDNA, we fitted the model (equation 1) to the
126 presence/absence of eDNA and of catch collected in Tokyo Bay as follows:

$$\text{logit } p_e(s_i) = \alpha_e + \sum_k f_k(x_k(s_i)) + w\theta(s_i) + u_e(s_i)$$

127

$$\text{logit } p_c(s_i) = \alpha_c + \theta(s_i) + u_c(s_i)$$

128 where α is the intercept, and $x_k(s_i)$ is the covariates at sites i for occurrence probabilities of
129 eDNA at sites i . In the study, seawater temperature, salinity, pH, and DO were used as
130 covariates which effect on the occurrence of eDNA (i.e., $k = 4$). $u(s_i)$ is spatial error that is
131 specific for each data following multivariate normal distributions $MVN(0, \mathbf{R})$, where the
132 variance–covariance matrix \mathbf{R} is a Matérn correlation function. θ which is shared between
133 eDNA and catch is the common spatial pattern between the two data, which cannot explain
134 by each terms of the equations. That is, θ can be interpreted as the spatial distribution of
135 species we want to know. 共変量の非線形性について書く Parameters in this model was
136 estimated by Integrated Nested Laplace Approximation using using the R-INLA package
137 (Lindgren, 2012) in R 3.6.1 (R Development Core Team, 2019).

138

139 Acknowledgments

140 This research was financially supported by Grant-in-Aid for Fisheries Agency of Japan.

141 Authorship

142 YK conceived of the research idea. YH, YU, HM, MI, KA, and AK conducted field
143 sampling. YH performed the laboratory experiments. YK, HO, and SN designed statistical
144 analyses. YK wrote programs and performed the analyses. YK wrote the manuscript with
145 input from all co-authors' comments.

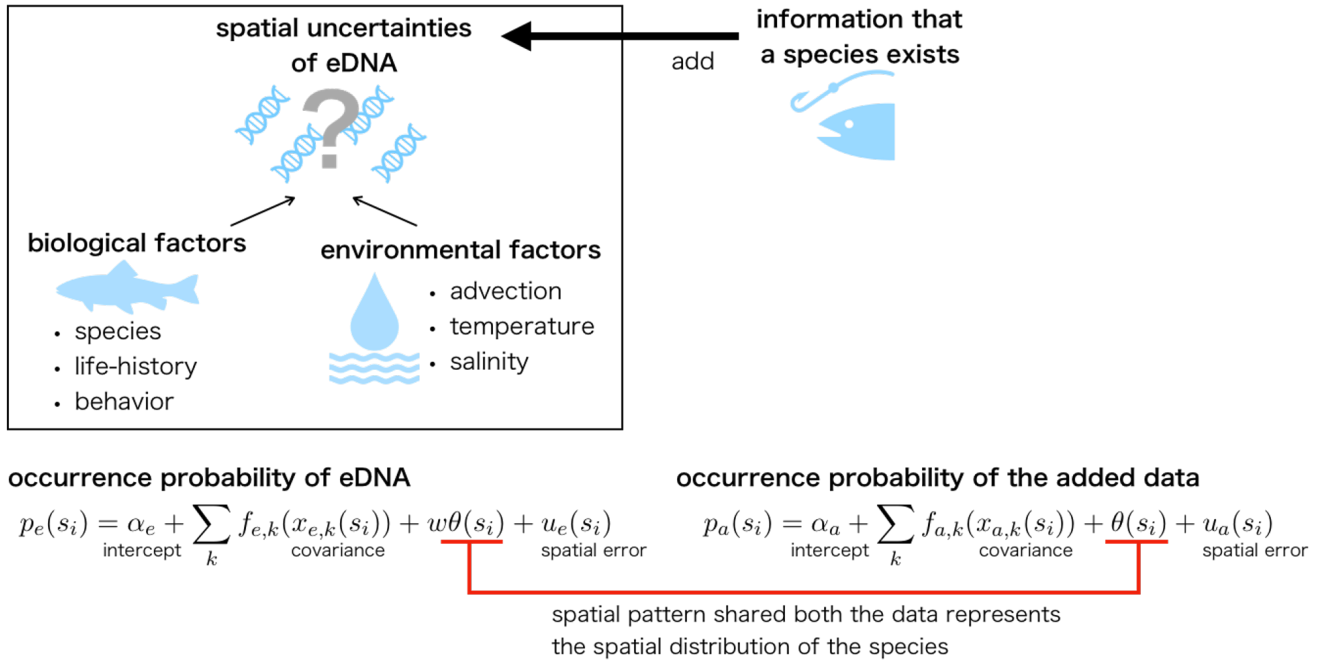


Fig. 1: Conceptual diagram of this study.