# 科研費シンポジウム

# "Recent Progress in Spatial and/or Spatio-temporal Data Analysis"

### Keynote speaker: Prof. Daniel A. Griffith

### (Univ. of Texas at Dallas)

# Program and Abstracts

October 30, 2020

Tohoku University

# Program

Friday, Oct 30, 2020
10:00-11:00 Griffith, D. (Univ. Texas at Dallas)
Important considerations about space-time data: modeling, scrutiny and ratification

11:00-11:30 Murakami, D. (ISM)
Compositionally-warped additive mixed modeling for large non-Gaussian data: Application to COVID-19 analysis

11:30-12:00 Shimono, T. (ISM)
Space-time analysis of COVID-19 in Japan using mobile space statistics(R).

Lunch break

13:30-14:00 Yajima, Y. (Tohoku U.)
On estimation of intrinsically stationary random fields.

14:00-14:30 Sugasawa, S. (Univ. Tokyo)
Spatially clustered regression

14:30-15:00　Toda,H. (Nagoya Institute of Technology)
Post-selection Inference for spatio-temporal trajectory segmentation

break

15:15-15:45 Hirano, T. (Kanto Gakuin U,)
A multi-resolution approximation via linear projection for large spatial datasets

15:45-16:15 Matsui, M. (Nanzan U.)
Testing independence of continuous time stochastic processes
-- toward independence test for random fields --

16:15-16:45 Matsuda, Y. (Tohoku U.)
Space time ARMA model

ABRIDGED VERSION[1]—IMPORTANT CONSIDERATIONS ABOUT SPACE-TIME
DATA: MODELING, SCRUTINY, AND RATIFICATION

Daniel A. Griffith

Ashbel Smith Professor of Geospatial Information Sciences, U. of Texas at Dallas

1. Introduction

Space-time data analysis has a literature spanning many decades, including Cliff et al. (1975), Bennett (1979), Heuvelink and Griffith (2010), and Cressie and Wikle (2011), among others. This literature's most modern section includes matrix models and Markov chain analysis, dynamic geographic optimization, and numerous statistical entries by Gelfand and his colleagues, and Christakos and his colleagues, inter alia. This paper acknowledges, without reviewing, this literature; instead, it focuses on the evolution of the space-time autocorrelation concept, a vital property of space-time data.

Autocorrelation (i.e., correlation amongst a single variable's observations) characterizes some correlated data family members, including time series, space series, and space-time series (Griffith, 2020). Its all-inclusive literature covers nearly two centuries, beginning with temporal autocorrelation—correlation among a single variable's data values at consecutive points in time—followed by spatial autocorrelation (SA)—correlation among a single variable's data values at pairs of neighboring points in space—and culminating in space-time autocorrelation—correlation among a single variable's data values at two consecutive points in time as well as nearby points in space.

Although interest in space-time data dates back at least to the mid-1900s, almost exclusively with separate treatments of space and time, such socio-economic and demographic datasets remained scarce until the 2000s. Unfortunately, earlier widely embraced accessible datasets furnished by Andrews and Herzberg (1985), for example, are fraught with entry errors that almost certainly scramble their space and/or time orderings, highlighting that accuracy is a fundamental space-time data quality assessment requirement, one complicated by the relatively large size and complexity of many space-time datasets vis-à-vis solely their time series or space series components.

The purpose of this paper is to address two important issues. The first is space-time autocorrelation in terms of its individual temporal and SA constituents latent in, and accuracy assessment preliminaries for evaluating the veracity of, a space-time dataset. The second is furnishing insights, articulating autocorrelation relationships, and raising awareness about probable best practices when undertaking an analysis of space-time data. Exemplifications of discussions employ annual United States (US) county population data for 1969-2019, the overwhelming majority of which the policy setting US National Cancer Institute (NCI) Surveillance, Epidemiology, and End Results (SEER) Program uses. The data generating process here is a combination of the US decennial census survey and annual updates to these census figures extracted from official government records that add births, subtract deaths, and adjust for net migration as well as age cohort shifts attributable to the passing of time.

2. The data: a brief overview

The US NCI SEER dataset currently consists of annual county time series from 1969 to 2018. This is an appealing dataset because it satisfies expert opinion asserting that the simplest of analyses requires a time series with a minimum of about 50 observations (see Hanke and Wichern, 2013). The selected geographic sample of these data for evaluation here contains the following six states: Ohio (OH), Oregon (OR), Florida (FL), Maine (ME), South Dakota (SD), and Texas (TX).

3. Temporal, spatial, and space-time autocorrelation

On average, temporal autocorrelation ($\rho_T$) tends to be stronger than SA ($\rho_s$), frequently exceeding 0.95; the degree of SA for socio-economic and demographic data tends to be in the 0.4 to 0.6 range, with most remotely sensed data having a $\rho_s$ value in excess of 0.9. These two descriptions underscore the empirical tendency for both forms of autocorrelation to be positive, with negative SA being much rarer than negative temporal autocorrelation. An intermingling of these two kinds of autocorrelation in space-time data routinely results in temporal autocorrelation dominating SA.

Each of the sample states has n counties ($16 \leq n \leq 254$), with each of these counties having a time series of length 51 before differencing, and 50 after differencing. Autocorrelation in each of the differenced time series is well-described by a lag-one model specification; however, temporal autocorrela-

---

tion is not very concentrated for most of the individual states. Although the preponderance of time series display relatively high degrees of positive temporal autocorrelation, the considerable heterogeneity of parameter estimates here fails to support the parsimonious positing of a single temporal autocorrelation parameter in a space-time autocorrelation specification, even separately state-by-state.

Each of the states has 50 time-differenced log-population density (improving normality) geographic distributions. Simultaneous autoregressive (i.e., SAR) spatial regression model estimation results reveal that a rook adjacency definition coupled with a second-order model specification describes autocorrelation in each of the space series well, and that individual state SA also is not very concentrated. The considerable heterogeneity of parameter estimates here once more fails to support the parsimonious positing of a single SA parameter in a space-time autocorrelation specification.

The preceding discussion treats temporal and SA separately, when they actually coalesce in space-time situations. The conceptualization and description of this interacting combination tends to occur in three different ways. Cliff et al. (1975) furnish one of the first comprehensive treatments of this conceptualization, presenting a space-time autoregressive (STAR) model specification of the following two forms: space-time lag (i.e., SA arises in terms of the preceding time period), and contemporaneous (i.e., SA is instantaneous, arising from the current time period). Meanwhile, contemporary statistical mixed models theory (e.g., West et al., 2007) furnishes a third form by positing a random effects (RE) specification to describe space-time data. This RE conceptualization maintains that regression model residuals are the sum of a systematic component, arising from, say, missing variables, plus a stochastic component, arising from the independent and identically distributed (iid) regression errors assumption. A standard chorological regression analysis requires additional information to separate these two components. One ancillary information source is repeated measures, such as time series of annual population densities, and another is a spatial weights matrix (SWM) that allows the partitioning of the systematic part into two sub-parts, a spatially structured RE (SSRE), which relates to a SWM, and a spatially unstructured RE (SURE), which is geographically random in nature, and hence void of SA. This RE term is a time invariant map that repeats itself for each point in a time series; it is a common factor across time instilling temporal autocorrelation into a space-time series dataset.

4. Space-time data accuracy assessment

One goal traversing the exploratory diagnostic computation of n temporal autocorrelation estimates, $\hat{\rho}_T$, and T SA estimates, $\hat{\rho}_S$, is assessing whether or not their respective variations are within margins of design-based or model-based stochastic sampling error. Sufficiently narrow estimate ranges support a parsimonious space-time data description whose predictive power reinforcements sustains the fidelity of its implications. Another useful investigative task is to inspect the global mean and variation of a space-time dataset. Yet a third helpful examination ascertains the degree to which space-time dataset parts may be interpolation and extrapolation technique constructions. Inspecting data extremes as well as variance homogeneity are two additional interrelated standard considerations. A further appraisal concerns model overfitting and the quality of any recognized imputed values. These seven assessments constitute best practice procedures for initiating a data scrutiny and ratification plan.

5. Conclusions

This paper establishes seven beneficial best practices enabling a space-time analyst to become more familiar with a given dataset, to more easily address data debugging and remediation issues, and to better express the soundness of inferred generalizations coupled with more robust finding limitations.

6. References

Andrews, D., and A. Herzberg. 1985. *Data: A Collection of Problems from Many Fields for the Student and Research Worker*. New York: Springer-Verlag.

Cliff, A., P. Haggett, J. Ord, K. Bassett, and R. Davies. 1975. *Elements of Spatial Structure: A Quantitative Approach*. London: Cambridge U. Press.

Cressie, N., and C. Wikle. 2-11. *Statistics for Spatio-Temporal Data*. New York: Wiley.

Griffith, D. 2020. A family of correlated observations: From independent to strongly interrelated ones, *Stats*, 3: 166-184; doi:10.3390/stats3030014

Hanke, J., and D. Wichern. 2013. *Business Forecasting* (9th ed.). Upper Saddle River, NJ: Pearson.

Heuvelink, G., and D. Griffith. 2010. Space–time geostatistics for geography: A case study of radiation monitoring across parts of Germany, *Geographic Analysis*, 42: 161-170.

West, B., K. Welch, and A. Galecki. 2007. *Linear Mixed Models: A Practical Guide Using Statistical Software*. New York: Chapman & Hall/CRC.

# Compositionally-warped additive mixed modeling for large non-Gaussian data: Application to COVID-19 analysis

Daisuke Murakami

Institute of Statistical Mathematics, 10-3 Midoricho, Tachikawa, Tokyo, 190-8562, Japan
Email: dmuraka@ism.ac.jp

## 1. Outline

An increasing number of non-Gaussian geospatial data is becoming available. At the same time, the size of spatial data rapidly grows together with the development of sensing technology. Given these backgrounds, this study develops a flexible additive mixed modeling approach for large non-Gaussian data. The development is done by combining an additive mixed model (AMM), which accommodates spatial and other effects, with the compositionally-warped Gaussian process (CWGP; Rois and Tober, 2019) estimating the shape of data distribution that can be either Gaussian or non-Gaussian possibly have skewness, fat tail, and other properties. The proposed model, termed compositionally-warped additive mixed model (CAMM) is estimated through a restricted likelihood maximization balancing model accuracy and complexity. Monte Carlo experiments shows that the proposed approach accurately model a wide variety of non-Gaussian data accuracy without losing computational efficiency relative to the linear AMM. The developed CAMM is applied to a spatiotemporal analysis of COVID-19 in Japan. The developed approach will be implemented in an R package spmoran.

## 2. Compositionally-warped additive mixed model (CAMM)

The proposed model describes non-Gaussian explained variables $y_i | i \in \{1, \ldots, N\}$ as follows:

$$\varphi_{\boldsymbol{\Theta}}(y_i) = \sum_{k=1}^{K} x_{i,k}\beta_k + \sum_{l=1}^{L} f_l(z_{i,l}) + \varepsilon_i, \ldots \ldots \varepsilon_i \sim N(0, \sigma^2). \tag{1}$$

$f_k(z_{i,k})$ is a smooth function depending on $l$-th covariate $z_{i,l}$, accommodating a wide variety of effects. Just like the classical AMM, this term can capture linear/non-linear effects, spatial and/or temporal effects, and other effects.

$\varphi_{\boldsymbol{\Theta}}(\cdot)$ is a function transforming the non-Gaussian variable $y_i$ to a nearly Gaussian variable. Interestingly, a wide variety of non-Gaussian variables can be transformed to Gaussian variables without assuming data distribution if the $\varphi_{\boldsymbol{\Theta}}(\cdot)$ function is defined by concatenating $D$ transformation functions as below (see Rois and Tober, 2019):

$$\varphi_{\boldsymbol{\Theta}}(y_i) = \varphi_{\boldsymbol{\theta}_D}\big(\varphi_{\boldsymbol{\theta}_{D-1}}(\cdots \varphi_{\boldsymbol{\theta}_1}(y_i))\big), \tag{2}$$

where $\boldsymbol{\Theta} \in \{\boldsymbol{\theta}_D, \boldsymbol{\theta}_{D-1}, \ldots, \boldsymbol{\theta}_1\}$. The $d$-th transformation function $\varphi_{\boldsymbol{\theta}_d}(y_i)$ in Eq. (2) is specified as

$$\varphi_{\boldsymbol{\theta}_d}(y_i) = a_d + b_d \sinh(c_d \operatorname{arcsinh}(y_i) - d_d), \tag{3}$$

where $\boldsymbol{\theta}_d \in \{a_d, b_d, c_d, d_d\}$ are parameters characterizing the $d$-th transformation. Rios and Tober (2019) called Eq.(3) as SAL layer (Sinh-Arcsinh and Affine where the "L" comes from linear).As illustrated in Figure 1, I confirmed accuracy of this transformation. After preliminary analysis, we decided to specify the $\varphi_{\boldsymbol{\Theta}}(y_i)$ function as shown in Figure 2.

Although Rois and Tober (2019) proposed the SAL transformation, the transformation has never been applied to regression modelling. Our novelty is to combine the SAL transformation and AMM to enable us flexibly modelling a wide variety of non-Gaussian data using AMM without assuming data distribution. It considerably improves modelling accuracy. In addition, the transformation function can be estimated computationally quite efficiently.

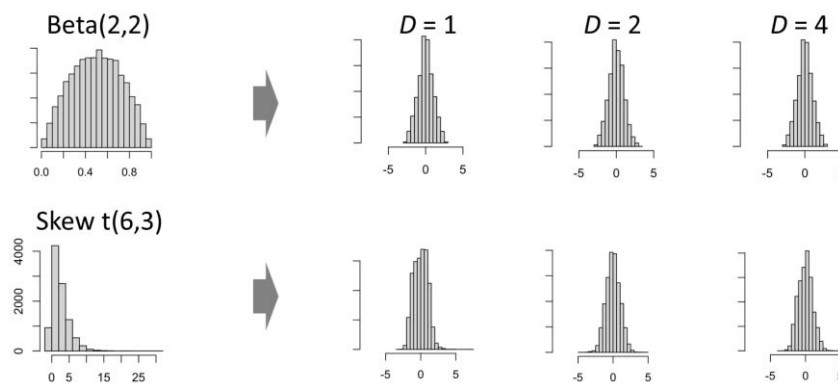In the presentation, I apply the developed CAMM to an COVID-19 analysis in Japan.



Figure 1: Fitting result of the for beta and skew t distributions. Left panels represent histograms of the simulated data and the right six panels show the histograms after the transformation.



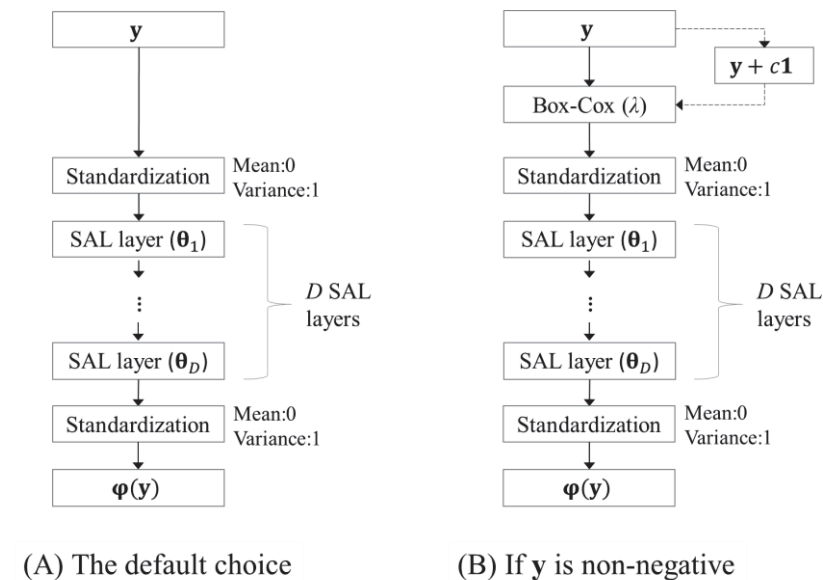(A) The default choice          (B) If **y** is non-negative

Figure 2: The $\varphi_{\boldsymbol{\Theta}}(\cdot)$ function for CAMM. (A) is the suggested default function while (B) is recommended if non-negative explained variables.

Reference

- Rios, G., & Tobar, F. (2019). Compositionally-warped Gaussian processes. *Neural Networks*, 118, 235-246.

# Space-time analysis of COVID-19 in Japan using Mobile Space Statistics®

Toshiyuki Shimono (ISM)

**Abstract:** Since around May 2020 COVID-19 is rampant in Japan nonetheless it is not so expanded as in the other severe countries. Two waves have been observed in the daily numbers of the reported patients in Japan of which peaks are on April and August. Looking into the patients report numbers classified by week, prefecture and age-group, the rampant prefectures so far is only a few overpopulated ones, and probably Tokyo is the only prefecture which failed to terminate the spread of the viruses during the state of emergency from April to May, and. We can look into the more detail geographically using the 500-meter grid population statistics all over Japan that is obtained by Mobile Space Statistics[1] provided by DOCOMO Insight Marketing, INC. It reveals that the infection largely has occurred probably only a few of very narrow geographic spaces in a few large cities because the mingling of people has happened effectively almost only there, concerning the people influx.

## 1. What the geographic time-space data reveal

(1) From the CSV formatted data obtained from https://gis.jag-japan.com/covid19jp/ (J.A.G Japan), one can see the age-demographic of the patients every day for each of 47 prefectures.

| 週の開始日 (日曜から土曜) | 1 北海道 | 2 青森 | 3 岩手 | 4 宮城 | 5 秋田 | 6 山形 | 7 福島 | 8 茨城 | 9 栃木 | 10 群馬 | 11 埼玉 | 12 千葉 | 13 東京 | 14 神奈川 | 15 新潟 | 16 富山 | 17 石川 | 18 福井 | 19 山梨 | 20 長野 | 21 岐阜 | 22 静岡 | 23 愛知 | 24 三重 | 25 滋賀 | 26 京都 | 27 大阪 | 28 兵庫 | 29 奈良 | 30 和歌山 | 31 鳥取 | 32 島根 | 33 岡山 | 34 広島 | 35 山口 | 36 徳島 | 37 香川 | 38 愛媛 | 39 高知 | 40 福岡 | 41 佐賀 | 42 長崎 | 43 熊本 | 44 大分 | 45 宮崎 | 46 鹿児島 | 47 沖縄 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2020-04-19 | 208 | 0 | 0 | 2 | 0 | 6 | 8 | 24 | 10 | 19 | 144 | 118 | 655 | 166 | 12 | 85 | 50 | 13 | 4 | 17 | 12 | 11 | 77 | 9 | 24 | 43 | 312 | 116 | 17 | 11 | 0 | 3 | 3 | 12 | 1 | 1 | 3 | 2 | 9 | 99 | 19 | 148 | 8 | 6 | 0 | 3 | 23 |
| 2020-04-26 | 227 | 4 | 0 | 3 | 0 | 2 | 7 | 5 | 1 | 6 | 78 | 44 | 663 | 140 | 10 | 34 | 34 | 1 | 2 | 4 | 1 | 10 | 17 | 0 | 2 | 40 | 181 | 35 | 9 | 5 | 0 | 6 | 1 | 13 | 3 | 0 | 1 | 0 | 1 | 46 | 5 | 0 | 2 | 1 | 0 | 0 | 8 |
| 2020-05-03 | 103 | 1 | 0 | 0 | 0 | 1 | 5 | 3 | 5 | 0 | 86 | 20 | 373 | 73 | 4 | 10 | 17 | 0 | 1 | 5 | 0 | 0 | 10 | 0 | 2 | 22 | 74 | 32 | 4 | 0 | 0 | 0 | 1 | 4 | 3 | 0 | 0 | 0 | 0 | 6 | 3 | 0 | 1 | 0 | 0 | 0 | 0 |
| 2020-05-10 | 44 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | 21 | 16 | 130 | 106 | 1 | 4 | 9 | 0 | 2 | 1 | 0 | 0 | 8 | 0 | 2 | 6 | 39 | 9 | 0 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 23 | 0 | 4 | 2 | 0 | 0 | 0 | 0 | 0 |
| 2020-05-17 | 38 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 4 | 2 | 13 | 6 | 41 | 64 | 0 | 2 | 9 | 0 | 3 | 0 | 0 | 2 | 1 | 0 | 1 | 0 | 13 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 9 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2020-05-24 | 37 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 5 | 5 | 96 | 39 | 0 | 0 | 3 | 0 | 4 | 1 | 1 | 0 | 12 | 0 | 0 | 0 | 0 | 2 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 84 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2020-05-31 | 29 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 2 | 8 | 7 | 140 | 27 | 0 | 1 | 0 | 3 | 1 | 1 | 4 | 3 | 1 | 0 | 1 | 2 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 52 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 2020-06-07 | 49 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 13 | 13 | 129 | 20 | 0 | 4 | 0 | 0 | 6 | 1 | 6 | 2 | 11 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 23 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

As shown above, only a few prefectures seemingly failed to stop COVID-19 during Emergency period from April to May in 2020, and based on the age-demographic of the patients the data also reveals that only Tokyo is probably the only prefecture if the people influx is concerned.

(2) Using the data at 12:00-13:00, July 15, 2020 from the visitor populations of 500-meter grids from Mobile Space Statistics, a clear tendency is that a prefecture is rampant or not with the threshold of 100 patients per week depends on if it has the most dense grid with 20,000 or 10,000 people, respectively, in April or in August, with a few exceptional prefectures.
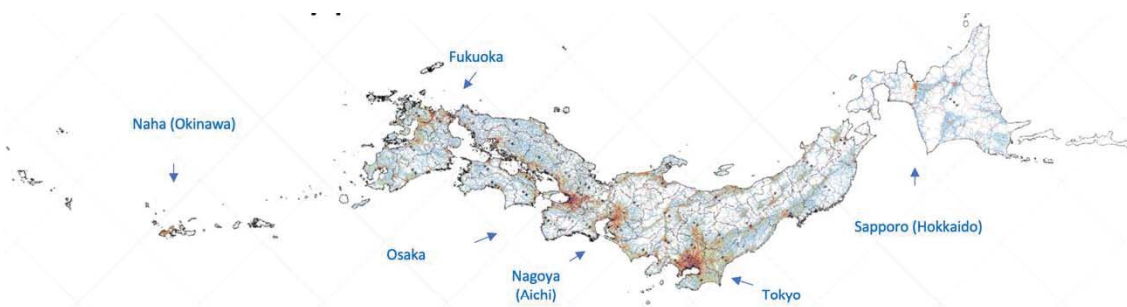


---

(3) The maps below of central areas of Tokyo represent (left) the populations whose residence is not in Tokyo or its neighbors, (middle) the demographics of all 47 residential prefectures, and (right) the demographics distinguishing about 2000 municipalities in all Japan. Each circle has the radius proportional to the square root of the corresponding population of the visitors, partitioned by corresponding to the ratios of different residential municipalities. In the middle map, blue means Saitama, cyan means Kanagawa, light green means Chiba, vast green means Tokyo residents, in each pie chart. As virtually only seen at and around large railway stations, the red partitions mean the visitors from the rest of 43 prefectures. In the right map, different colors just mean different municipalities (市区町村 in Japanese).



One can notice that the mingling of people largely occurs only at around Tokyo Station, Shinjuku, Shibuya, Ikebukuro, Ueno and Asakusa, each of six has the width of as small as about 1 to 2 kilometers. The similar mingling occurs at three spots in Osaka prefecture, two in Aichi, one in each prefecture with a large city such as Hokkaido, Miyagi, Okinawa.

## 2. Suggestions

Contrary to the intuitive perception that COVID-19 prevails everywhere so that the countermeasure resources should be invested extensively, the infection may largely occur very tiny spots in Japan even though sporadically small-scale clusters have been occurring already in all the 47 prefectures. This insight may play crucial role in modeling and analyzing the space time behavior of this disease to consider the spread and the flow of this disease, which may also require the age factor. Other infectious disease such as influenza, pneumonia, gastroenteritis also may have similar tendencies that could be surveyed from the data collected by the authority, which may give us the clues to control these infectious diseases.



(The colors of cyan to red here reflect $\mathrm{int}(\log_2(\text{population}))$ of the 500-meter grids with the data of MLIT, 国土数値情報)

# On Gaussian semiparametric estimation for
# two-dimensional intrinsic stationary random fields

Yoshihiro Yajima
Tohoku University
2020. Oct 30.

## Abstract

We propose two estimators of two-dimensional intrinsic stationary random fields (ISRFs) observed on a regular grid and derive their asymptotic properties. Originally they are proposed to estimate prameters of long memory models of stationary and nonstationary time series. One is the log-periodgram regression estimator (Robinson(1995a); Velasco(1995a)) and the other is the Gaussian semiparametric estimtor(the local Whittle estimator Robinson(1995b); Velasco(1999b)). We apply them to two dimensional ISRFs. These ISRFs include a fractional Brownian field, which is a Gaussian random field and is used to model many physical processes in space(Mandelbrot, B.B. and Van Ness, J.W. (1968)). The estimators are consistent and have the limiting normal distributions as the sample size goes to infinity. We also list some problems such as testing isotropy or applications to more general models that are to be solved in future.

**Keywords**:intrinsic stationary random fields; spatio-temporal models; local Whittle estimator; log periodogram regression; fractional Brownian field.

## References

Mandelbrot, B.B. and Van Ness, J.W. (1968). Fractional Brownian motion, fractal noises and applications. *SIAM. Rev.* **10** 422-437.

Robinson,P.M. (1995a). Log periodogram regression of time series with long range dependence. *Ann. Statist.* **23** 1048-1072.

Robinson,P.M. (1995b). Gaussian semiparametric estimation of long range dependence. *Ann. Statist.* **23** 1630-1661.

Velasco,C. (1999a). Non-stationary log-periodogram regression. *J. Econometrics* **91** 325-371.

Velasco,C. (1999). Gaussian semiparametric estimation of non-stationary time series. *J. Time Ser. Anal.* **20** 87-127.

# Spatially Clustered Regression

Shonosuke Sugasawa

Center for Spatial Information Science, The University of Tokyo

## 1 Introduction

Geographically weighted regression (GWR) is widely adopted for modeling possibly spatially varying regression coefficients. However, GWR is known to be numerically unstable and may produce extreme estimates of coefficients especially when covariates are spatially correlated. Recently, Li and Sang (2019) adopted fused lasso to shrink regression coefficients in neighboring locations, but it can be computationally intensive under large spatial data. In this work, we propose a new strategy for spatial regression that takes account of spatial heterogeneity in regression coefficients. The proposed method can be easily estimated via a simple iterative algorithm, and it can handle variable selection or semiparametric modeling.

## 2 Spatially Clustered Regression

Let $y_i$ be a response variable and $x_i$ is a vector of covariates in the $i$th location, for $i = 1, \ldots, n$, where $n$ is the number of samples. We suppose we are interested in the conditional distribution $f(y_i|x_i; \theta_i)$, where $\theta_i$ is a vector of unknown parameters. Here $\theta_i$ may change over different locations and represent spatial heterogeneity. For example, $f(y_i|x_i; \theta_i) = \phi(y_i; x_{i1}^t \theta_i + x_{i2}^t \gamma, \sigma_i^2)$. We assume that geographical information $s_i$ (e.g. longitude and latitude) is also available for the $i$th location. We further assume that $n$ locations are divided into $G$ groups and locations in the same group share the same parameter values of $\theta_i$. We introduce $g_i \in \{1, \ldots, G\}$, an unknown group membership variable for the $i$th location, and let $\theta_i = \theta_{g_i}$. Then, the distinct values of $\theta_i$'s reduce to $\theta_1, \ldots, \theta_G$, where $\theta = (\theta_1^t, \ldots, \theta_G^t)^t$ is the set of unknown parameters.

Therefore, the unknown parameters in the model is the structural parameter $\theta$ and membership parameter $g = (g_1, \ldots, g_n)$. Regarding the membership parameter, it would be reasonable to consider that the membership in neighboring locations are likely to have the same memberships, which means that the fitted conditional distributions are likely to be the same in the neighboring locations. In order to encourage such structure, we propose the following penalized likelihood:

$$Q(\theta, g) \equiv \sum_{i=1}^{n} \log f(y_i | x_i; \theta_{g_i}) + \phi \sum_{i<j} w_{ij} I(g_i = g_j), \qquad (1)$$

where $w_{ij} = w(s_i, s_j) \in [0, 1]$, $w(\cdot, \cdot)$ is a weighting function, and $\phi$ controls strength of spatial similarity. The penalty function is motivated from the Potts model (Potts, 1952) and a similar penalty function is adopted in Sugasawa (2020). We define the estimator of $\theta$ and $g$ as the maximizer of the objective function $Q(\theta, g)$. The maximization can be easily carried out by a simple iterative algorithm similar to $k$-means algorithm. Owing to the simple formulation (1), the proposed strategy allows several important extensions. For example, variable selection can be done by introducing additional penalty function for $\theta$, and semiparametric form for the regression term can also be adopted.

We will report the numerical performance of the proposed method compared with existing methods such as GWR or method by Li and Sang (2019) through simulation studies and real data applications.

## References

Li, F. and H. Sang (2019). Spatial homogeneity pursuit of regression coefficients for large datasets,. *Journal of the American Statistical Association 114*, 1050–1062.

Potts, R. B. (1952). Some generalized order-disorder transformations. *Mathematical Proceedings of the Cambridge Philosophical Society 48*, 106–109.

Sugasawa, S. (2020). Grouped heterogeneous mixture modeling for clustered data. *Journal of the American Statistical Association*, to appear.

# Post-selection Inference for Spatio-temporal Trajectory Segmentation

Hiroki Toda[1], Yu Inatsu[2], and Ichiro Takeuchi[1, 2]

[1]Nagoya Institute of Technology, [2]Riken AIP

## 1 Introduction

Trajectory segmentation is a common task in spatio-temporal trajectory data analysis. It splits a sequence of locations with time stamps into a small number of sub-sequences or segments with respect to some criteria. Despite the development of a wide variety of methods [1], to date, little attention has been paid to quantifying the uncertainty of segment breakpoints identified by trajectory segmentation. In this study, we aim to develop inference tools that provide valid $p$-values for each breakpoint. The difficulty lies in the fact that the location of each breakpoint is selected by a segmentation algorithm, and this fact must be properly incorporated in the statistical inference. Unfortunately, if one uses classical statistical inference, the $p$-values or confidence intervals are not valid anymore in the sense that the false positive rate cannot be controlled at the desired significance level. To address this problem, we adopt the framework of Selective Inference (SI) [2] (also known as Post-selection Inference), a new statistical inference framework for data-driven hypotheses. This enables us to perform exact (non-asymptotic) inference conditioning on the selection procedure. Additionally, we introduce a parametric programing approach [3, 4] to solve the problem that SI has low statistical power due to over-conditioning, which was assumed to be one of the major drawbacks of SI. To the best of our knowledge, this study is the first application of SI to trajectory data analysis. In the talk, we will demonstrate the performance of the proposed methods when in animal trajectory data analysis.

## 2 Trajectory Segmentation and Statistical Tests

Let $T = [p_1, p_2, \ldots, p_n]$ denote a trajectory of length $n$, where each point $p_i = (x_i, y_i, t_i)$ consists of $(x, y)$ location (and possibly additional parameters) of a moving object at time $t_i$. A trajectory segmentation algorithm splits the trajectory $T$ into segments at breakpoints $\boldsymbol{\tau} = (\tau_1, \ldots, \tau_K)$ with known or unknown number of breakpoints $K$. Trajectory segmentation is mainly classified into two types, time series-based and topology-based algorithm. The former first converts the trajectory into univariate sequence $\boldsymbol{x}_{\mathrm{obs}} \in \mathbb{R}^N$ of a feature (e.g., speed, acceleration and direction), then change-point detection is applied to $\boldsymbol{x}_{\mathrm{obs}}$. The latter directly uses the locational data $\boldsymbol{x}_{\mathrm{obs}} = (x_1, \ldots, x_n, y_1, \ldots, y_n)^\top \in \mathbb{R}^N$ as an input. We assume that $\boldsymbol{x}_{\mathrm{obs}}$ is a single observation drawn from $X \sim \mathrm{N}(\boldsymbol{\mu}, \Sigma) \in \mathbb{R}^N$.

We develop inference tools for the two types of algorithms. For time series-based algorithm, we consider optimal change-point detection for a sequence with piecewise constant and piecewise linear mean. In the case of the piecewise constant, statistical test of interest might be $H_0^{(k)} : \mu_k = \mu_{k+1}$ v.s. $H_1^{(k)} : \mu_k \neq \mu_{k+1}$ for $k \in \{1, \ldots, K\}$, where $\mu_k$ is the population mean of $k$th segment. For topology-based algorithm, we consider Ramer-Douglas-Peucker (RDP) algorithm, which is traditional but popular because of its simplicity. Due to space limitations, we omit the details of the algorithms and statistical tests for each algorithm.

# 3 Selective Inference

The basic idea of SI [2] is to make inference conditional on the selection event, which allows us to derive the exact (non-asymptotic) sampling distribution of a test statistic. Given a statistical test $H_0$ and $H_1$, we assume that the test statistic can be written in the form of $\boldsymbol{\eta}^\top \boldsymbol{x}_{\mathrm{obs}}$ using some contrast vector $\boldsymbol{\eta} \in \mathbb{R}^N$. Then, we have the following (two-sided) selective $p$-value:

$$p := \mathrm{Pr}_{H_0}\left(|\boldsymbol{\eta}^\top \boldsymbol{X}| \geq |\boldsymbol{\eta}^\top \boldsymbol{x}_{\mathrm{obs}}| \mid \mathcal{M}(\boldsymbol{X}) = \mathcal{M}(\boldsymbol{x}_{\mathrm{obs}}), \mathcal{A}(\boldsymbol{X}) = \mathcal{A}(\boldsymbol{x}_{\mathrm{obs}}), \mathcal{P}_{\boldsymbol{\eta}}^\perp \boldsymbol{X} = \mathcal{P}_{\boldsymbol{\eta}}^\perp \boldsymbol{x}_{\mathrm{obs}}\right).$$

Note that $p$ satisfies $\mathrm{Pr}_{H_0}(p < \alpha) = \alpha$, $\forall \alpha \in [0, 1]$. The first condition $\mathcal{M}(\boldsymbol{X}) = \mathcal{M}(\boldsymbol{x}_{\mathrm{obs}}) = \hat{\boldsymbol{\tau}}$ indicates the event that breakpoints $\hat{\boldsymbol{\tau}}$ are selected. The second condition $\mathcal{A}(\boldsymbol{X}) = \mathcal{A}(\boldsymbol{x}_{\mathrm{obs}})$ indicates the algorithm-dependent nuisance selection event that we had no choice but to condition on for tractability in most cases. The condition $\mathcal{P}_{\boldsymbol{\eta}}^\perp \boldsymbol{X} = \mathcal{P}_{\boldsymbol{\eta}}^\perp \boldsymbol{x}_{\mathrm{obs}}$ is introduced for technical reasons, where $\mathcal{P}_{\boldsymbol{\eta}}^\perp = I_N - \boldsymbol{c}\boldsymbol{\eta}^\top$ is the orthogonal projection matrix with $\boldsymbol{c} = \Sigma\boldsymbol{\eta}(\boldsymbol{\eta}^\top\Sigma\boldsymbol{\eta})^{-1}$.

Conditioning not only on the selection of breakpoints but also on the algorithm procedure itself $\mathcal{A}(\boldsymbol{X}) = \mathcal{A}(\boldsymbol{x}_{\mathrm{obs}})$ makes the conditioning space smaller:

$$\{\boldsymbol{X} : \mathcal{M}(\boldsymbol{X}) = \mathcal{M}(\boldsymbol{x}_{\mathrm{obs}}), \mathcal{A}(\boldsymbol{X}) = \mathcal{A}(\boldsymbol{x}_{\mathrm{obs}})\} \subseteq \{\boldsymbol{X} : \mathcal{M}(\boldsymbol{X}) = \mathcal{M}(\boldsymbol{x}_{\mathrm{obs}})\}.$$

This leads low statistical power of SI. Existing exact SI framework has suffered from this over-conditioning problem.

Recently, we have developed a new SI framework [3, 4] that uses a parametric programming technique to circumvent the over-conditioning. We define the parameterized data $\boldsymbol{x}_{\mathrm{obs}}'(z) = \boldsymbol{c}z + \mathcal{P}_{\boldsymbol{\eta}}^\perp \boldsymbol{x}_{\mathrm{obs}}$ with a parameter $z \in \mathbb{R}$, then we have a valid selective $p$-value

$$p := \mathrm{Pr}_{H_0}\left(|\boldsymbol{\eta}^\top \boldsymbol{X}| \geq |\boldsymbol{\eta}^\top \boldsymbol{x}_{\mathrm{obs}}| \mid \mathcal{M}(\boldsymbol{X}) = \mathcal{M}(\boldsymbol{x}_{\mathrm{obs}}), \mathcal{P}_{\boldsymbol{\eta}}^\perp \boldsymbol{X} = \mathcal{P}_{\boldsymbol{\eta}}^\perp \boldsymbol{x}_{\mathrm{obs}}\right)$$
$$= \mathrm{Pr}_{H_0}\left(|z| \geq |\boldsymbol{\eta}^\top \boldsymbol{x}_{\mathrm{obs}}| \mid \mathcal{M}(\boldsymbol{x}_{\mathrm{obs}}'(z)) = \mathcal{M}(\boldsymbol{x}_{\mathrm{obs}})\right).$$

By using this, we are able to perform powerful testing. We also have developed the effective procedure to identify the all intervals of $z$ where the same breakpoints $\hat{\boldsymbol{\tau}}$ are obtained, by searching $z$ over $(-\infty, \infty)$.

## Acknowledgments

## References

[1] Hendrik Edelhoff, Johannes Signer, and Niko Balkenhol. Path segmentation for beginners: An overview of current methods for detecting changes in animal movement patterns. *Movement Ecology*, 4, 12 2016.

[2] Jason D. Lee, Dennis L. Sun, Yuekai Sun, and Jonathan E. Taylor. Exact post-selection inference, with application to the lasso. *The Annals of Statistics*, 44(3):907–927, Jun 2016.

[3] Vo Nguyen Le Duy, Hiroki Toda, Ryota Sugiyama, and Ichiro Takeuchi. Computing valid p-value for optimal changepoint by selective inference using dynamic programming. *arXiv:2002.09132*, 2020. (to appear in NeurIPS 2020).

[4] Vo Nguyen Le Duy and Ichiro Takeuchi. Parametric programming approach for more powerful and general lasso selective inference. *arXiv:2004.09749*, 2020.

# A multi-resolution approximation via linear projection for large spatial datasets

Toshihiro Hirano

Kanto Gakuin University

### Abstract

Recent technical advances in collecting spatial data have been increasing the demand for methods to analyze large spatial datasets. The statistical analysis for these types of datasets can provide useful knowledge in various fields. However, conventional spatial statistical methods, such as maximum likelihood estimation and kriging, are impractically time-consuming for large spatial datasets due to the necessary matrix inversions. To cope with this problem, we propose a multi-resolution approximation via linear projection ($M$-RA-lp). The $M$-RA-lp conducts a linear projection approach on each subregion whenever a spatial domain is subdivided, which leads to an approximated covariance function capturing both the large- and small-scale spatial variations. Moreover, we elicit the algorithms for fast computation of the log-likelihood function and predictive distribution with the approximated covariance function obtained by the $M$-RA-lp. Simulation studies and a real data analysis for air dose rates demonstrate that our proposed $M$-RA-lp works well relative to the related existing methods.

**Keywords:** Covariance tapering; Gaussian process; Geostatistics; Large spatial datasets; Multi-resolution approximation; Stochastic matrix approximation

## 1 Introduction

Advances in Global Navigation Satellite System (GNSS) and compact sensing devices have made it easy to collect a large volume of spatial data with coordinates in various fields such as environmental science, traffic, and urban engineering. The statistical analysis for these types of spatial datasets would assist in an evidence-based environmental policy and the efficient management of a smart city.

In spatial statistics, this type of statistical analysis, including model fitting and spatial prediction, has been conducted based on Gaussian processes. However, traditional spatial statistical methods, such as maximum likelihood estimation and kriging, are computationally infeasible for large spatial datasets, requiring $O(n^3)$ operations for a dataset

E-mail: 1hirano2@kanto-gakuin.ac.jp

of size $n$. This is because these methods involve the inversion of an $n \times n$ covariance matrix.

Hirano (2020) proposed a multi-resolution approximation via linear projection ($M$-RA-lp) of Gaussian processes observed at irregularly spaced locations. The $M$-RA-lp implements the linear projection on each subregion obtained by partitioning the spatial domain recursively, resulting in an approximated covariance function that captures both the large- and small-scale spatial variations unlike the covariance tapering and some low rank approaches. Additionally, we derive algorithms for fast computation of the log-likelihood function and predictive distribution with the approximated covariance function obtained by the $M$-RA-lp. Also, these algorithms can be parallelized. Our proposed $M$-RA-lp is regarded as a combination of the two recent low rank approaches: a modified linear projection (MLP) (Hirano, 2017) and a multi-resolution approximation ($M$-RA) (Katzfuss, 2017). The $M$-RA-lp extends the MLP by introducing multiple resolutions based on the idea of Katzfuss (2017), leading to better approximation accuracy of the covariance function than that by the MLP. Particularly, when the variation of the spatial correlation around the origin is smooth like the Gaussian covariance function, the approximation accuracy of the covariance function by the MLP often degrades. In contrast, the $M$-RA-lp avoids this problem. Additionally, the $M$-RA-lp is regarded as an extension of the $M$-RA and enables not only to alleviate the knot selection problem but also to increase empirically numerical stability in specific steps of fast computation algorithms of the $M$-RA. Simulation studies and a real data analysis for air dose rates generally support the effectiveness of our proposed $M$-RA-lp in terms of computational time, estimation of model parameters, and prediction at unobserved locations when compared with the MLP and $M$-RA.

# References

Hirano, T. (2017). Modified linear projection for large spatial datasets. *Communications in Statistics - Simulation and Computation*, 46:870–889.

Hirano, T. (2020). A multi-resolution approximation via linear projection for large spatial datasets. To appear in *Japanese Journal of Statistics and Data Science*.

Katzfuss, M. (2017). A multi-resolution approximation for massive spatial datasets. *Journal of the American Statistical Association*, 112:201–214.

# Testing independence of continuous time stochastic processes − toward independence test for random fields −

Nanzan University    Muneya Matsui

## Abstract

Firstly we give a talk about independence test of a pair of stochastic processes. As a measure of independence, we construct distance covariance (DC) and distance correlation (DCR) based on approximations of the component processes at finitely many discretization points. Assuming that the mesh of the discretization converges to zero as a suitable function of the sample size, we show that the sample distance covariance and correlation converge to limits which are zero if and only the component processes are independent. In the talk, we moderately explain theoretical results and spare more time for numerical studies.

Secondly several ideas toward independence test for random fields are given. Especially, we state differences in sampling scheme between stochastic processes and random fields.

## Definitions

For two processes $X$ and $Y$ on $[0,1]$ with some mild conditions, we define DC for processes

$$
\begin{aligned}
T_\beta(X,Y) &= \mathbb{E}\big[\|X_1 - X_2\|_2^\beta \|Y_1 - Y_2\|_2^\beta\big] + \mathbb{E}\big[\|X_1 - X_2\|_2^\beta\big]\mathbb{E}\big[\|Y_1 - Y_2\|_2^\beta\big] \\
&\quad -2\,\mathbb{E}\big[\|X_1 - X_2\|_2^\beta \|Y_1 - Y_3\|_2^\beta\big], \qquad \beta \in (0,2],
\end{aligned}
$$

where $\|\xi\|_2$ denotes the $L^2$-norm of a process $\xi$ on $[0,1]$. Of course, $T_\beta(X,Y) = 0$ for independent $X,Y$. The converse is not obvious; we prove it in Theorem 0.2. The corresponding distance correlation is given by $R_\beta(X,Y) = T_\beta(X,Y)/\sqrt{T_\beta(X,X)\cdot T_\beta(Y,Y)}$. Since the whole path of a process $Z$ on $[0,1]$ is unavailable in reality, we consider discretizations of the process at a partition $0 = t_0 < t_1 < \cdots < t_p = 1$ of $[0,1]$. Assuming that $p = p_n \to \infty$ as $n \to \infty$ and the mesh satisfies $\delta_n = \max_{i=1,\dots,p}(t_i - t_{i-1}) \to 0$, $n \to \infty$, we normalize the points $Z(t_i)$ by $\sqrt{t_i - t_{i-1}}$. Writing for any partition $(t_i)$, $\Delta_i = (t_{i-1}, t_i]$, $|\Delta_i| = t_i - t_{i-1}$, $i = 1,\dots,p$, we consider a vector of weighted discretizations $\mathbf{Z}_p = \big(|\Delta_1|^{1/2}Z(t_1),\dots,|\Delta_p|^{1/2}Z(t_p)\big)$ and define the discretization of the process $Z^{(p)}(t) = \sum_{i=1}^p Z(t_i)\mathbf{1}(t \in \Delta_i)$. For stochastically continuous, measurable and bounded processes $Z$ and $Z'$ we have as $p \to \infty$,

$$
|\mathbf{Z}_p - \mathbf{Z}'_p|^2 = \sum_{i=1}^p (Z(t_i) - Z'(t_i))^2 |\Delta_i| = \|Z^{(p)} - (Z')^{(p)}\|_2^2 \to \int_0^1 (Z(t) - Z'(t))^2\,dt = \|Z - Z'\|_2^2,
$$

in probability. Therefore, we could approximate $T_\beta(X,Y)$ by $T_\beta(X^{(p)}, Y^{(p)})$ properly.

The sample analog of $T_\beta(X,Y)$ and $R_\beta(X,Y)$ are respectively given by

$$
\begin{aligned}
T_{n,\beta}(X,Y) &= \frac{1}{n^2}\sum_{k,l=1}^n \|X_k - X_l\|_2^\beta \|Y_k - Y_l\|_2^\beta + \frac{1}{n^2}\sum_{k,l=1}^n \|X_k - X_l\|_2^\beta \frac{1}{n^2}\sum_{k,l=1}^n \|Y_k - Y_l\|_2^\beta \\
&\quad -2\frac{1}{n^3}\sum_{k,l,m=1}^n \|X_k - X_l\|_2^\beta \|Y_k - Y_m\|_2^\beta,
\end{aligned}
$$

and $R_{n,\beta}(X,Y) = T_{n,\beta}(X,Y)/\sqrt{T_{n,\beta}(X,X)\cdot T_{n,\beta}(Y,Y)}$.

## Main results

**1**. Asymptotics of test statistic.

**Theorem 0.1.** *Assume some moment and smoothness conditions of $(X, Y)$ and the growth condition on $p = p_n \to \infty$. Then under the null hypothesis ($X$ and $Y$ are independent),*

$$R_{n,\beta}(X^{(p)}, Y^{(p)}) \xrightarrow{p} 0, \quad and \quad n\, R_{n,\beta}(X^{(p)}, Y^{(p)}) \xrightarrow{d} \sum_{i=1}^{\infty} \lambda_i(N_i^2 - 1) + c$$

*for an iid sequence of standard normal random variables $(N_i)$, a constant c, and a square summable sequence $(\lambda_i)$.*

For proving the second quantity, we notice that $T_{n,\beta}(X, Y)$ has representation as a $V$-statistics of order 4 with a 1-degenerate symmetric kernel $h_4 = h(x_1, x_2, x_3, x_4)$.

**2**. The condition $T_\beta(X, Y) = 0$ and independence of $X$ and $Y$.
Let $B_1$, $B_2$ be independent Brownian motions (BMs) on $[0, 1]$, independent of $(X, Y)$. The stochastic integrals $Z_1 = \int_0^1 X dB_1$ and $Z_2 = \int_0^1 Y dB_2$ are well defined (and are, given $(X, Y)$, independent normal random variables). Let $\mathcal{F}_B$ denote the $\sigma$-field generated by $B = (B_1, B_2)$. The quantity $T_\beta(X, Y)$ is shown to be contracted from the stochastic integrals $Z_1$, $Z_2$.

**Theorem 0.2.** *If the stochastic integrals $Z_1$ and $Z_2$ are a.s. conditionally independent given $\mathcal{F}_B$ then $X, Y$ are independent. In particular, if $\beta \in (0, 2)$ and $\mathbb{E}[\|X\|_2^\beta + \|Y\|_2^\beta + \|X\|_2^\beta \|Y\|_2^\beta] < \infty$, then $T_\beta(X, Y) = 0$ if and only if $X, Y$ are independent. Then we have*

**3**. The bootstrap for the sample distance covariance.
The bootstrap can be made to work for the degenerate $V$-statistic $T_{n,\beta}(X, Y)$. We validate that the bootstrap version of $n\, T_{n,\beta}(X^{(p)}, Y^{(p)})$ approximates the bootstrap distribution of $T_{n,\beta}(X, Y)$.

**4**. Simulations.
We illustrate the theoretical results in a small simulation study using typical processes such as fractional Brownian motions, $\alpha$-stable Lévy motions, etc. With various boxplots, we see the convergences of $T_{n,\beta}(X^{(p)}, Y^{(p)})$ to theoretical limits assuming $X, Y$ are independent/(weak/strong)dependent. We have also conducted a simulation study to illustrate the performance of the bootstrap procedure for the distance correlation based test for independence. Specifically, we have tested for independence of two BMs and two $\alpha$-stable Lévy motions $X$, $Y$.

All details of results are given in [2], see also DCR for time sires [1], and another version of DC for stochastic processes [3].
In the talk we present several sampling schemes for approximating DCR for random fields.

## References

[1] Davis, R.A., Matsui, M., Mikosch, T. and Wan, P. (2018) Applications of distance correlation to time series. *Bernoulli* **24**, 3087–3116.

[2] Dehling, H., Matsui, M., Mikosch, T., Samorodnitsky, G. and Tafakori, L. (2018) Distance covariance for discretized stochastic processes. *Bernoulli* (to appear).

[3] Matsui, M., Mikosch, T. and Samorodnitsky, G. (2017) Distance covariance for stochastic processes. *Probab. Math. Statist.* **37**, 355–372.

# SPACE-TIME
# AUTOREGRESSIVE MOVING AVERAGE MODELS

YASUMASA MATSUDA

## 1. Abstract

In this talk, we propose a space-time autoregressive and moving average (ST-ARMA) model for spatio-temporal data, a discrete time series observation of irregularly spaced data, denoted as $X_t(s), s \in \mathbb{R}^2, t = 1.2.\ldots.$ Figure 1 shows observation points in US to record monthly precipitation, providing a typical example of spatio-temporal data. Regarding $X_t(s)$ as a $L^2(\mathbb{R}^2)$-valued time series, we construct a space-time ARMA($p,q$) model, given by

$$(1) \qquad X_t(s) = \sum_{j=1}^p \int_{\mathbb{R}^2} \phi_j(s-u) X_{t-j}(u) du + \sum_{j=0}^q \int_{\mathbb{R}^2} \theta_j(s-u) L_{t-j}(du),$$

$$s \in \mathbb{R}^2, t = 0, \pm 1, \pm, 2, \ldots,$$

where $L_t(u)$ is a Lévy sheet on $\mathbb{R}^2$ independent across $t$, and $\phi(s)$ and $\theta(s)$ are CARMA kernels in $L^2(\mathbb{R}^2)$. It is an temporal extension of continuous ARMA random fields of Brockwell and Matsuda[2] by a convolutional operator of $\phi$ and $\theta$ on $L^2(\mathbb{R}^2)$.

A space-time ARMA model is a kind of model for functional time series, a $H^2$-valued time series for a Hilbert space $H^2$. See Ramsay and Silverman [5] for independent cases, Bosq[1] for stationary time series cases, Liu et al. [4] for a pure AR model for $L^2[0,1]$ valued time series, and Li et al. [3] for a semiparametric method to detect a long memory property in $L^2[0,1]$ valued time series. One feature of ST-ARMA model in (1) is that it is a $L^2(\mathbb{R}^2)$-valued time series model. It scauses several difficulty in establishing ST-ARMA model properties that infinite region $\mathbb{R}^2$ rather than the fixed interval $[0,1]^2$ over which square integrable functions of Hilbert space is defined.

We shall introduce the basic properties of space-time ARMA models given as

- stationary conditions, more specifically, causal and invertible conditions,
- explicit form of spectral density functions,

1

Figure 1. Weather stations in United States

- Whittle estimation for parameters that specify CARMA kernels of $\phi(\cdot)$ and $\theta(\cdot)$,
- forecasting,
- empirical applications to US precipitation data.

The striking features are the explicit derivation of spectral density functions which makes it possible to conduct a parametric estimation by Whittle likelihood function and forecasting of future values by the estimated CRMA kernels.

## References

[1] Bosq, D. (2000) *Linear Processes in Function Spaces.* Springer, New York.
[2] Brockwell, P. and Matsuda, Y. (2017) Continuous auto-regressive moving average random fields on $\mathbb{R}^n$. *Journal of the Royal Statistical Society: Series B*, 79, 833-857.
[3] Li, D., Robinson, P. and Shang, H. L. (2020) Long-Range Dependent Curve Time Series. *Journal of the American Statistical Association* , 115, 957-971.
[4] Liu, X., Xiao, H and Chen, R. (2016) Convolutional autoregressive models for functional time series. *Journal of Econometrics*, 194, 263-282.
[5] Ramsay, J. O. and Silverman, B. W. (1997) *Functional Data Analysis.* Springer, New York.

Graduate School of Economics and Management, Tohoku University, 27-1 Kawauchi, Aoba ward, Sendai 980-8576, Japan
    *E-mail address*: yasumasa.matsuda.a4@tohoku.ac.jp

# List of Participants

Daniel A. Griffith (Univ. of Texas at Dallas)

Daisuke Murakami (ISM)

Toshiyuki Shimono (ISM)

Yoshihiro Yajima (Tohoku University)

Shonosuke Sugasawa (University of Tokyo)

Hiromi Toda (Nagoya Institute of Technology)

Toshihiro Hirano (Kanto Gakuin University)

Muneya Matsui (Nanzan University)

Yasumasa Matsuda (Tohoku University)