

Consistence Condition of Kernel Selection in Regular Linear Kernel Regression and Its Application in COVID-19 High-risk Areas Exploration

Lu xan, Ba lin

Department of Mathematic and Statistics, Central China Normal University, Luoyu, Wuhan, 430071, Hubei, China

Abstract

With the long-term outbreak of the COVID-19 around the world, identifying high-risk areas is becoming a new research boom. In this paper, we propose a novel regression method namely Regular Linear Kernel Regression (RLKR) for COVID-19 high-risk areas exploration. We explain in detail how the canonical linear kernel regression method is linked to the identification of high-risk areas for COVID-19. Furthermore, the consistence condition of Kernel Selection, which is closely related to the identification of high-risk areas, is given with two mild assumptions. Finally, the RLKR method was verified by simulation experiments and applied for COVID-19 high-risk area Exploration.

Keywords: Linear kernel • COVID-19 source exploration • Kernel selection • Superimposed fields

Introduction

COVID-19 was recognized in December 2019 [1]. Since then, the COVID-19 out-break has posed critical challenges for the public health, research, and medical communities. Currently, it was recognized that human to human transmission played a major role in the subsequent outbreak [2]. Based on this knowledge, in the data collection stage, we try to extract data related to the flow of people as much as possible. At present, the research on the COVID-19 mainly focuses on the treatment of the epidemic [3,4]; the transmission of the epidemic [5,6] and the prevention of the epidemic [7,8]. In F. Ndairou 2020, the concept of super-spreaders, who contribute disproportionately to a much larger number of cases in the ongoing COVID-19 pandemic, was proposed for modeling the COVID-19 transmission. We focus our attention on high-risk areas exploration from "sky" perspective. Which means that, we ignore different regional political factors and medical level factors? All the factors we can know come from aerial monitoring. We imagine ourselves as a space traveler. Before landing on the earth, determining which areas are high-risk areas and which areas are low-risk areas through detection in the sky is actually our task. In order to make a decision, we had to count the number of outbreaks in each region, as well as all the data related to the spread of the outbreak, e.g. population, population density per square kilometer, population transmission, people gathering frequency, distances between regions, etc. Given this information of different areas, we cannot decide the risk level of one area, just by the outbreaks. Because, for several adjacent areas, the surrounding areas are likely to be only the victim and only one area is the source of the outbreak. One the other hand, for areas of high volume of population transmission and people gathering frequency, its low-risk label may only be temporary. Along the continuous outbreak of COVID-19, those areas are much like to become high-risk areas.

The inspiration of Regular Linear Kernel Regression, which is used for high risk area detection of COVID-19 in this paper, comes from a variable selection method, called Lasso [9], which is firstly proposed for solving S-sparse linear regression model in high-dimensional statistics. Actually the

proof framework of Kernel Selection Consistence origins from the variable selection consistence of Lasso [10], but a slightly modification for suitable for Kernel Selection. To gain an intuition of the link between Kernel Selection and high-risk areas exploration, let us see inside of Kernel Regression. Denote $K_{\alpha}(x)=K(\alpha, x)$ be a measure of the distance between α and x , where the Gaussian RBF(Radial Basis Function) Kernel and Linear Inner Product Kernel are most widely used. In this paper, the kernel selection consistence condition of Regular Linear Kernel Regression was given. For a given group of kernels, $K_{\alpha}(\cdot)$, $\alpha \in \bar{\alpha}$, where $\bar{\alpha}$ stands for an index set, Kernel Regression model has the form,

$$y = \sum_{\alpha \in \bar{\alpha}} \theta_{\alpha} K_{\alpha}(x) + w, \quad (1)$$

Where, y stands for outbreaks and w stands for observation noise. From the model, y is actually a linear combination of the characteristic Kernels at the point of x . We see a Kernel $K_{\alpha}(x)$ as a spread function of a high-risk area characterized by α . As stated, $K_{\alpha}(x)$ is a measure of the distance between α and x . For a long distance area x from high-risk area α , the influence from α , source of transmission, shall be small, which is coincident with the property of $K_{\alpha}(x)$. Upper to now, we explain the link between the spreading mode of COVID-19 and the Kernel Regression, that's high-risk areas were seen as emission sources characterized by $K_{\alpha}(\cdot)$ in the model and the outbreaks in one area is actually the combination of the influences of those emission sources. Still, there are problems that must be check carefully. First, the number of high-risk areas can be unlimited? For a given design matrix $X \in \mathbb{R}^{n \times d}$, and a candidate set $\text{Card}(\bar{\alpha})=m$, the number of high-risk areas s must be less than or equal to $\min(n, d, m)$, otherwise the true high-risk areas will not be check out entirely. Second, for a given estimation of (θ) with a support set of S^* , under which conditions the support set is consistent with the true support set S of true θ ? This question is crucial, as it means whether we have got the true high-risk areas or COVID-19 sources. As a non-zero θ_{α} means that the candidate area of $K_{\alpha}(\cdot)$ influences other areas.

Notice that, Interpolation learning has recently attracted growing attention in the machine learning community [11,12]. Linear Kernel

***Address for Correspondence:** Dr. Lu , Department of Mathematic and Statistics, Central China Normal University, Luoyu, Wuhan, 430071, Hubei, China; E-mail: LuChaoAnHui@126.com

Copyright: ©2023 Lu et al. This is an open-access article distributed under the terms of the creative commons attribution license which permits unrestricted use, distribution and reproduction in any medium, provided the original author and source are credited.

Received date: 18-March-2022, Manuscript No. VCRH-22-57653; **Editor assigned:** 21-March-2022, Pre QC No. VCRH-22-57653(PQ); **Reviewed:** 04-April-2022, QC No. VCRH-22-57653; **Accepted** 09-April-2022, Manuscript No. VCRH-22-57653(A); **Published date:** 16-April-2022, DOI: 10.37421/2736-657X.07.2023.003.

Regression model has been around for decades, whose continued popularity stems from its excellent properties, well generalization along with good interpolation, in high- dimensional statistics. Due to the derivability, Kernel Ridge Regression with l_2 -norm penalization is very popular in the High-dimensional area, including the converge rate [13], generalization risk [14] and variable selection [15]. While the Consistence condition of Kernel Regression with l_1 -norm penalization is the novelty of this paper, that's the topic of Kernel selection consistency under l_1 -norm regularization is first proposed in this paper. We compliment the consistence condition for Kernel Regression with l_1 -norm penalization, along with the generalization risk. Due to the interpolation of Kernel regression [16], the connection between Kernel selections with high-risk areas detection of COVID-19 is our innovation. Around those questions, we structured our paper as follows. In the introduction, we briefly introduce the background of COVID-19 transmission and give the inspiration of high-risk areas exploration. In the second section, some remarks were introduced and Linear Kernel Regression Model was proposed for the spreading mode of COVID-19 on high-risk areas influence view. As we have no hope for estimating an unique solution of $\hat{\theta}$ with the optimal least square method. Regular technique was suggested for estimating an unique solution. In the third section, the conditions of Kernel Selection Consistence were given. Finally the simulation experiment was conduct for supplements of consistence checking. Then the Kernel Selection Method was applied for COVID-19 high-risk areas detection with outbreaks of COVID-19 up to year 2021 along with some other features that influence the transmission of COVID-19.

Materials and Methods

Model construction

Suppose we are given relevant features in observation regions, $x_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$ and corresponding response values y_i , for $i=1, \dots, n$, which stands for the most concerned variable, e.g. the number of epidemic out breaks. A significant research is to find possible emission sources or high-risk areas from many candidate areas $\alpha_j = (\alpha_{j1}, \dots, \alpha_{jd})^T$, for $j=1, \dots, m$. The true high-risk areas α_S is supposed to be contained in those candidate areas $\alpha = (\alpha_1, \dots, \alpha_m)^T$ indexed by S and $S \subseteq \{1, \dots, m\}$ stands for the true support of high-risk areas. Here (m) is a simplified notation of $\{1, \dots, m\}$.

Intuitively, the response variable y at observation point x is the linear superimposed effects of each real high-risk areas, that's

$$y = \sum_{j \in S} \theta_j \langle x, \alpha_j \rangle + w \tag{2}$$

where $S \subseteq \{1, \dots, m\}$ stands for the supports of truth emission sources, and w stands for observation noise. $\langle \cdot, \cdot \rangle$ stands for inner product. Since S is unknown in practice, the estimation problem of S will be more meaningful, except for the estimation of θ_j , $j \in S$. Actually if the support S of the truth emission sources was given, we can estimate $\hat{\theta}$ very well. With the true S unknown, there is no hope for obtaining an unique solution of θ_j , when $m \geq \min(n, d)$, from Lemma 1 in Appendix. Without uniqueness, there is no way to talk about the Kernel Selection Consistence. Inspired by Lasso of l_1 norm regularization, consider the following optimization,

$$\hat{\theta} \leftarrow \arg \min \sum_{i=1}^n \frac{1}{2n} (y_i - \sum_{j \in [m]} \theta_j \langle x_i, \alpha_j \rangle)^2 + \lambda_n \|\theta\|_1 \tag{3}$$

Here, $[m] = \{1, \dots, m\}$ is a simple notation and λ_n is a given parameter. Denote $X = [x_1, \dots, x_n]^T \in \mathbb{R}^{n \times d}$, $\alpha = [\alpha_1, \dots, \alpha_m]^T \in \mathbb{R}^{m \times d}$, $\theta = [\theta_1, \dots, \theta_m]^T \in \mathbb{R}^m$ and $Y = [y_1, \dots, y_n]^T \in \mathbb{R}^n$. Therefore a matrix form alternative of equation (3), namely Regular Linear Kernel Regression, is given by,

$$\hat{\theta} \leftarrow \arg \min \frac{1}{2n} \|Y - X\alpha^T \theta\|_2^2 + \lambda_n \|\theta\|_1 \tag{4}$$

For the not differentiable of l_1 - norm, due to its sharp point at the origin, there's no hope for a close form solution. Since this is a convex optimization problem, many optimization methods can be used to solve the problem, e.g. Gradient descent or Least Angel Regression.

So the method of estimating $\hat{\theta}$ is not the focus of this paper. The explanation of the estimation result of $\hat{\theta}$ and the condition of consistence between the estimation and the truth θ is our focus. Here we give the explanation of $\hat{\theta}$. From the model(2), for a zero $\theta_j=0$, which means y is irrelevant with α_j , the location of α_j shall not be the emission source. For a non-zero θ_j , which means y is relevant with α_j , the location of α_j shall be the emission source. For a positive θ_j , the corresponding location of α_j shall be seen as pollution source in the pathogen exploration case. And a negative θ_j , the corresponding location of α_j shall be seen as purification source.

Results

Condition of consistence

Consider an S -sparse linear kernel regression model and the true support of θ^* is S . Given an estimation $\hat{\theta}$ with support S^* . The question is by which condition the estimation S^* is consistent with S . We begin by assuming X be deterministic design matrices. Let S stands for the supports of truth θ^* and α_S stands for the subset of $\alpha = [\alpha_1, \dots, \alpha_m]^T$ with indices of S . First the truth emission source α_S cannot be linearly correlated, in order to ensure that the model is identifiable, even if the support set S were known a priori. Second the fake source α_{S^c} shall be irrelevant with the truth emission source α_S , which makes it possible to exclude irrelevant areas. In general, we cannot expect this orthogonality to hold, but a type of approximate orthogonality. In particular, consider the following conditions:

Lower eigenvalue: The smallest eigenvalue of the sample covariance submatrix indexed by S is bounded below:

$$\gamma_{\min} \left(\frac{\alpha_S X^T X \alpha_S^T}{n} \right) \geq c_{\min} > 0 \tag{5}$$

Mutual incoherence: There exists some $\rho \in [0, 1)$ such that

$$\max_{j \in S^c} \left\| \left(\alpha_S X^T X \alpha_S^T \right)^{-1} \alpha_S X^T X \alpha_j \right\|_1 < \rho \tag{6}$$

With this set-up, the following results follow applied to the lagrange form Linear Kernel Lasso (4) when applied to an instance of the linear kernel model (2) such that the true parameter θ^* is supported on a subset S with cardinality s . In order to state the result, we introduce the convenient shorthand $\Pi_{\perp(X\alpha^T)} = I_n - X\alpha_S^T (\alpha_S X^T X \alpha_S^T)^{-1} \alpha_S X^T$ a type of orthogonal projection matrix.

Theorem 1

Consider an S -sparse linear kernel regression model for which the design matrix satisfies conditions (5) and (6). Then for any choice of regularization parameter such that

$$\lambda_n \geq \frac{2}{1-\rho} \left\| \alpha_{S^c} X^T \Pi_{\perp(X\alpha^T)} \frac{w}{n} \right\|_{\infty} \tag{7}$$

the Lagrangian Kernel Lasso (3) has the following properties:

- (a) Uniqueness: There is a unique optimal solution $\hat{\theta}$
- (b) No false inclusion: The solution has its support set S^* contained within the true support set S .
- (c) l_{∞} -bounds: The error $\hat{\theta} - \theta^*$ satisfies

$$\|\hat{\theta}_S - \theta_S^*\|_{\infty} \leq \underbrace{\left\| (\alpha_S X^T X \alpha_S^T)^{-1} \alpha_S X^T \frac{w}{n} \right\|_{\infty}}_{B(\lambda_n; X\alpha^T)} + \left\| \left(\frac{\alpha_S X^T X \alpha_S^T}{n} \right)^{-1} \right\|_{\infty} \lambda_n \tag{8}$$

(d) No false exclusion: The solution's support set S^* includes all indices $i \in S$ such that $|\theta_i^*| > B(\lambda_n; X^T)$, and hence is sample selection consistent if $\min_{i \in S} |\theta_i^*| > B(\lambda_n; X\alpha^T)$,

The prove of this result utilizes a technique called a primal-dual witness method, which constructs a solution and then the solution was shown be optimal and unique, see Appendix. Here let us try to interpret its main

claims. First the uniqueness claim in part(a) is not trivial. Although the Lagrange objective is convex, it can never be strictly convex when $m > \min\{d, n\}$. Based on the uniqueness claim, we can talk unambiguously about the support of the Kernel Lasso estimate $\hat{\theta}^*$. Part (b) guarantees that the kernel lasso does not falsely include linear kernels that are not in the support of θ^* , or equivalently that $\theta^{*S^c} = 0$. Part (d) is a consequence of the l_∞ -norm bound from part (c): as long as the minimum value of $|\theta_i^*|$ over indices $i \in S$ is not too small, then the Kernel Lasso is kernel selection consistent in the full sense.

Further more, if we make specific assumptions about the noise vector w , more concrete results can be obtained. Before giving the results, we give the definition of the Sub-Gaussian variable.

Definition 1

A random variable w with mean $\mu = E(w)$ is sub-Gaussian if there is a positive number σ such that

$$E[e^{\lambda(w-\mu)}] \leq e^{\sigma^2 \lambda^2 / 2} \quad \text{For all } \lambda \in \mathbb{R}.$$

With this in mind, we give a third assumption, that's (A3) The observation noise $w_i, i \in \{1, \dots, n\}$ is i.i.d. σ -sub-Gaussian variables. Before given the concrete corollary about Linear Kernel Regression, a Lemma about lower bounds for sub-Gaussian variables is needed.

Lemma 2: Let $\{Z_i\}_{i=1}^n$ be a sequence of zero-mean random variables, each sub Gaussian with parameter σ . (No independence assumptions are needed). Consider the lower bound of $Z = \max_{i=1, \dots, n} Z_i$, we have

$$P[Z \geq t] \leq 2ne^{-\frac{t^2}{2\sigma^2}} \quad \text{Valid for all } n \geq 2. \quad (9)$$

Corollary 3: Consider the S -sparse linear kernel regression model based on a noise vector w that assumption (A3) holds, and a deterministic design matrix X that satisfies assumptions (A3) and (A4). And the character Kernel satisfies the C-Column normalization condition: $\max_{k=1, \dots, m} \|X\alpha_k\|_2 / \sqrt{n} \leq C$. Suppose that we solve the Regular Linear Kernel Regression (3) with regularization parameter

$$\lambda_n = \frac{2c_\sigma}{1-p} \left\{ \sqrt{\frac{2 \log(m-s)}{n}} + \delta \right\} \quad (10)$$

for some $\delta > 0$. Then the optimal solution $\hat{\theta}^*$ is unique with its support contained within S , and satisfies the l_∞ -error bound.

$$\|\hat{\theta}_s - \theta_s^*\|_\infty \leq \frac{\sigma}{\sqrt{c_{\min}}} \left\{ \sqrt{\frac{2 \log s}{n}} + \delta \right\} + \left\| \left(\frac{\alpha_s X^T X \alpha_s^T}{n} \right)^{-1} \right\|_\infty \lambda_n \quad (11)$$

With probability at least $1 - 4e^{-\frac{n\delta^2}{2}}$

Proof: First, let us show that the given choice of $\lambda_n(10)$ of the regularization parameter satisfies the bound(7) with high probability. It suffices to bound the maximum absolute value of the random variables

$$Z_k := \alpha_k^T X^T \left[I_n - X \alpha_s^T (\alpha_s X^T X \alpha_s^T)^{-1} \alpha_s X^T \right] \left(\frac{w}{n} \right) \quad \text{For } k \in S^c.$$

Since $\Pi_{S^c}(X)$ is an orthogonal projection matrix, we have

$$\|\Pi_{S^c} X \alpha_k\|_2 \leq \|X \alpha_k\|_2 \leq C\sqrt{n}$$

where the last inequality follows from the C-Column normalization condition. Therefore, each variable Z_k is sub-Gaussian with parameter at most $C^2 \sigma^2 / n$. By Lemma 2, it holds that,

$$P \left[\max_{j \in S^c} Z_j \geq t \right] \leq 2(m-s)e^{-\frac{nt^2}{2C^2\sigma^2}}$$

Transforming this to the l_∞ -norm form, we have

$$P \left[\frac{2}{1-p} \left\| \alpha_s X^T \Pi_{S^c}(X \alpha^T) \frac{w}{n} \right\|_\infty \geq \lambda_n \right] \leq 2(m-s)e^{-\frac{n(1-p)^2 \lambda_n^2}{8C^2\sigma^2}}$$

Substitute the choice of λ_n (10), we have

$$\lambda_n \geq \frac{2}{1-p} \left\| \alpha_s X^T \Pi_{S^c}(X \alpha^T) \frac{w}{n} \right\|_\infty \quad \text{with probability at least } 1 - 2e^{-\frac{n\delta^2}{2}}.$$

It remains to simplify the l_∞ -bound (8). Consider the random variables $\tilde{Z}_i := e_i^T \left(\frac{1}{n} \alpha_s X^T X \alpha_s^T \right)^{-1} X \alpha_s^T w / n$

Since the elements of the vector w are i.i.d. σ -sub-Gaussian, the variable \tilde{Z}_i is zero-mean and sub-Gaussian with parameter at most

$$\frac{\sigma^2}{n} \left\| \left(\frac{1}{n} \alpha_s X^T X \alpha_s^T \right)^{-1} \right\|_2 \leq \frac{\sigma^2}{c_{\min} n}$$

where we have used the eigenvalue condition (5). Consequently, for any $\delta > 0$, we have

$$P \left[\max_{i=1, \dots, s} |\tilde{Z}_i| > \frac{\sigma}{\sqrt{c_{\min}}} \left\{ \sqrt{\frac{2 \log s}{n}} \right\} \right] \leq 2e^{-\frac{n\delta^2}{2}}$$

Then the claims follows from $(1 - 2e^{-\frac{n\delta^2}{2}})^2 \geq 1 - 4e^{-\frac{n\delta^2}{2}}$

Upper to here, we have given the Kernel selection consistence conditions. By careful analysis of the original optimization problem (3), we find that the values of $(\theta^*)_{k \in S^c}$ gradually diminish, as the penalize parameter λ_n is large enough. From condition (7), if the given parameter λ_n is too small, (7) condition is violated, we have no hope for excluding all fake Kernels. From part (d) of Theorem 1, the $\hat{\theta}_i, i \in S^c$, which cannot be zero, may be miscalculated as zero. In real world data detection, the choice of λ shall be carefully checked. Besides, Cross-validation methods are usually utilized for the decision of λ_n , or other improved methods [17].

Prediction error or risk assessment

Another interesting twist is to do a risk assessment for all areas. Consider the Linear Kernel Regression Model (2). Given an estimation of $\hat{\theta}$, the prediction of outbreaks \hat{y} is the superposition summation of the Kernels at the position of x , that's

$$\hat{y} = \sum_{k \in [m]} \hat{\theta}_k < \alpha_{k,x} >$$

As this prediction \hat{y} is the linear combination of the Kernel functions, radiation effects of the true high-risk areas, at the point x , not just the number of the local outbreaks. We see \hat{y} as the risk assessment of the area x . Given a S - Sparse Linear Kernel Regression model, we devoted ourselves to bound on the mean-squared prediction error

$$\frac{\|X \alpha^T (\hat{\theta} - \theta^*)\|_2^2}{n} = \frac{1}{n} \sum_{i=1}^m \left[\sum_{k=1}^m < \alpha_{k,x_i} > (\hat{\theta}_k - \theta_k^*) \right]^2 \quad (12)$$

Before giving the results, some more constraints about the design matrix, namely Restricted Eigenvalue (RE) condition shall be advertised here.

Definition 2

The matrix X satisfies the Restricted Eigenvalue (RE) condition over S with parameters (κ, γ) if

$$\frac{1}{n} \|X \Delta\|_2^2 \geq \kappa \|\Delta\|_2^2 \quad \text{for all } \Delta \in \mathcal{C}_\gamma(S)$$

Here the subset $\mathcal{C}_\gamma(S) = \{\Delta \in \mathbb{R}^d \mid \|\Delta_{S^c}\|_1 \leq \gamma \|\Delta_S\|_1\}$, is a collection of vectors. whose l_1 -norm off the support is dominated by the l_1 -norm on the support.

The Restricted Eigenvalue (RE) condition constructs a link between the recovering error $\|\hat{\theta} - \theta\|_2$ and the prediction error $\|X \alpha^T (\hat{\theta} - \theta^*)\|_2$. With these set-ups, we give the bounds on prediction errors as follows.

Theorem 4

(Prediction error bounds) Consider the Regular Linear Kernel Regression Optimization Problem(3) with a strictly positive regularization

parameter $\lambda n \geq 2\|\alpha X^T w/n\|_\infty$. (a) Any optimal solution $\hat{\theta}$ satisfies the bound

$$\|X\alpha^T(\hat{\theta} - \theta^*)\|_2 \leq 12\|\theta^*\|_2 \lambda_n$$

(b) If θ^* is supported on a subset S of cardinality s , and the design matrix $X\alpha^T$ satisfies the $(\kappa; 3)$ -RE condition over S , then the optimal solution satisfies the bound

$$\frac{\|X\alpha^T(\hat{\theta} - \theta^*)\|_2}{n} \leq \frac{9}{k} s \lambda_n \frac{2}{n}$$

The proof of this prediction bounds can be seen in Appendix. Before that, let's see what we are informed from the bounds. The two upper bounds are all related to λn . In order to reduce the forecast errors, a critical step is to reduce the regularization parameter λn .

Simulation studies

Consider a S -sparse Linear Kernel Regression model (2). Suppose each row x_i of the design matrix $X \in \mathbb{R}^{n \times d}$ is independently draw from a Gaussian distribution $N(0, I_d)$. And each α_j of the candidates kernels $\alpha \in \mathbb{R}^{m \times d}$ is also i.i.d of Gaussian distribution $N(0, I_d)$. The true support S is a subset of $\{1, \dots, m\}$ uniformly distributed in $[m]$. The real parameter of θ_j^* , $j \in S$ is uniformly distributed on $[-1, -0.1] \cup [0.1, 1]$, and $\theta_j^* = 0$ for $j \in S^c$. Let the entries of the noise vector w follow i.i.d of a Gaussian distribution of $N(0, 0.01)$.

Then the response variable Y is calculated by

$$Y = X\alpha^T\theta + w.$$

Applied by the regular method of l_1 -norm, we are interested in the false-inclusion rate and consistency. False-inclusion rate is defined as,

$$\text{False-inclusion} = \frac{\hat{S} / S}{\hat{S}} \tag{13}$$

Where \hat{S} stands for the support of estimation $\hat{\theta}$. And the consistency between \hat{S} and S , is given by

$$\text{Consistency} = \frac{\hat{S} \cap S}{\hat{S} \cup S} \tag{14}$$

(F) and consistency (Con) under given levels of λn in Table 1. From Table 1, we know that the best result is $F=0$, which means no false inclusion and $\text{Con} = 1$, which means the estimation support \hat{S} is identical with the true kernel support S . For λn that's not too small, the no false inclusion holds. And the consistence grows up and down, as λn increases. The consistence holds for a long time as we increase the value of λn , which means that for a long period, the support of $\hat{\theta}$ makes no changes. And this information will be useful for empirical selection of the regularization parameter λn .

Further, we find that the sample dimension has an important relationship with the consistency. We calculated the average consistency under different dimensions of $d \in [40, 100, 200, 300, \text{ and } 500]$. The average is the result of 100 replicates, where the λn is chosen according to maximizing consistence. From the results in Table 2, it's easy to find that as the dimension increases, the average consistency increases steadily. As the number of candidate Kernels increase, picking out exactly the true supports becomes more difficult. From here, we get the revelation that increasing the dimension of the data is of great significance to correctly pick out the Kernels. And reducing the number of candidate kernels is also useful for reducing false inclusion rates and improving consistency.

Finally, we care about the changing law of the prediction error, by adding tests check. We draw the test data-set by choosing a test design matrix $X^{\text{test}} \in \mathbb{R}^{T \times d}$, with each row x_i^{test} follows a Gaussian distribution $N(0, I^d)$. And the test response variable Y^{test} is calculated by

$$Y^{\text{test}} = X^{\text{test}}\alpha^T\theta + w$$

The prediction error is calculated by equation (12). The curves of prediction error along λn were plotted in Figure 1 under given levels of $s \in [5, 10, 15, 20]$, which stands for the number of true Kernels. Here $n=1000$, $m=1000$, $d=500$, $T=500$, where T stands for the size of test data-set. From the results of the Figure 1, the prediction error grows as λn goes beyond some fixed point. And for the different levels of s , the prediction error grows proportional to s , which is coincident with the result from part (b) of Theorem 4.

Table 1. The False-inclusion rate and consistency of given levels of λn , where $n=m=1000$, $d=500$.

| s | λn | F | Con | s | λn | F | Con |
|------|-------------|-------|-------|------|---------------|-------|-------|
| s=5 | 0 ~ 0.03 | 0.992 | 0.007 | s=15 | 0 0.05 | 0.97 | 0.029 |
| | 0.03 ~ 55.5 | 0 | 1 | | 0.05 ~ 68.7 | 0 | 1 |
| | 55.5 ~ 95.9 | 0 | 0.8 | | 68.7 ~ 72.7 | 0 | 0.933 |
| | 95.9 ~ 242 | 0 | 0.6 | | 72.7 ~ 103 | 0 | 0.867 |
| s=10 | 0 ~ 0.10 | 0.988 | 0.011 | s=20 | 0 ~ 0.05 | 0.964 | 0.035 |
| | 0.10 ~ 103 | 0 | 1 | | 0.05 ~ 82.8 | 0 | 1 |
| | 103 ~ 254 | 0 | 0.9 | | 82.8 ~ 105.5 | 0 | 0.95 |
| | 254 ~ 263.6 | 0 | 0.8 | | 105.5 ~ 128.3 | 0 | 0.9 |

Table 2. The average False-inclusion rate and consistency of dimensions of $d \in [40, 100, 200, 300, 500]$. Here $n=1000$, $s=20$ and the average is calculated over 100 replicates.

| m | d | F | Con | m | d | F | Con |
|-------|-------|-------|-------|--------|-------|-------|-------|
| m=40 | d=40 | 0.205 | 0.772 | m=300 | d=40 | 0.603 | 0.234 |
| | d=100 | 0.043 | 0.956 | | d=100 | 0.33 | 0.593 |
| | d=200 | 0.004 | 0.995 | | d=200 | 0.074 | 0.923 |
| | d=300 | 0.001 | 0.998 | | d=300 | 0.014 | 0.984 |
| | d=500 | 0 | 1 | | d=500 | 0.001 | 0.999 |
| m=100 | d=40 | 0.462 | 0.43 | m=1000 | d=40 | 0.731 | 0.118 |
| | d=100 | 0.158 | 0.83 | | d=100 | 0.475 | 0.366 |
| | d=200 | 0.019 | 0.98 | | d=200 | 0.169 | 0.8 |
| | d=300 | 0.002 | 0.997 | | d=300 | 0.049 | 0.948 |
| | d=500 | 0 | 1 | | d=500 | 0.002 | 0.997 |

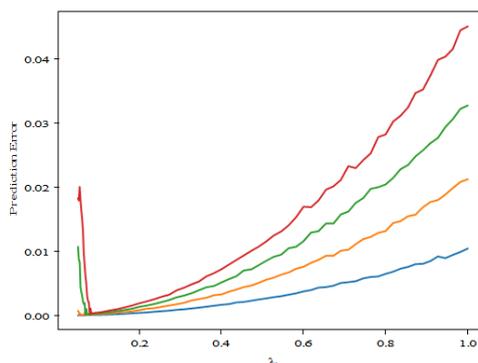


Figure 1. The Prediction error of Regular Linear Kernel Regression. Note: (—) S=5; (—) S=10; (—) S=15; (—) S=20.

Discussion

High-risk areas selection and risk assessment of COVID-19

In this section, we try to explore high-risk areas of COVID-19 with the outbreaks of COVID-19, provided by <https://fy.onesight.com/>. This data includes the number of confirmed cases in each country since the outbreak of COVID-19. As the data collection deadline is September 2021, we had to collect the relevant covariates of each country up to that point in time. First, we collected the coordinates of the capital of each country, from the latitude β_t and longitude β_g , by

$$\{x=\cos(\beta_t)\cos(\beta_g), y=\cos(\beta_t) \sin(\beta_g), z=\sin(\beta_t)\}$$

As mentioned earlier, the COVID-19 now spreads mainly in the form of person-to-person spread. So in the second part of data collection, the variables of population transmission were collected, including population, population density per square kilometer, population birth rate and Gross Domestic Product (GDP), which concerns the activity of each country.

And in the third part of the data collection, the history of death toll of COVID-19 by weeks is also collected from <https://www.kaggle.com/tarunkr/covid-19-case-study-analysis-viz-comparisons> from January 2020 to September 2021. In summary, we collected outbreaks and associated covariates of 114 dimensions for 189 countries.

The first concern is the selection of high-risk areas. We make all 189 countries or regions be candidates of Kernels. Before the regular Linear Kernel Regression model being implemented, a normalize scale of data set is employed, as follows

$$x_{scale} = \frac{x - \bar{x}}{\sqrt{S^2(x)}}$$

Here, \bar{x} stands for sample mean of x and $S^2(x)$ stands for sample covariates, $S^2(x)=\frac{1}{n} \sum (x_i - \bar{x})^2$. Given several levels of λ , the Kernel areas selected were given in Table 3. In the table, high-risk areas are countries whose coefficients are positive. Purification areas are countries, whose coefficients are negative.

Another meaningful task is to carry out a risk assessment of all candidate areas. Given $\lambda=0.5$, we made a risk assessment for all 189 countries or areas. The risk is calculated by

$$Y_{risk} = \sum_{k=1}^m \hat{\theta}_k \langle x, \alpha_k \rangle$$

As all variables of covariates were scaled with zero-mean and standard deviation $S^2(x)=1$, including the number of outbreaks. So the risk Y_{risk} of each area takes values around 0. For high-risk areas, the response Y_{risk} is also large. In particular, the assessment results are presented in Table 4, 5. From the risk assessment table, we see that Brazil, United Kingdom and Russia are of high risk in risk assessment program. While in purification areas selection, they are viewed as purification areas. The reason lies in the those areas are close in "distance" with high-risk areas and actually they are victims. They are surrounded by the true COVID-19 emission sources. The relation between high-risk selections with risk assessment is that high-risk selection is finding the causes or emission sources and risk assessment is the presentation of spreading of COVID-19. For some areas that are judged to be of great risk assessment, they are likely to be innocent and may be purification areas. With Kernel selection of high-risk areas, only then can we decide whether an area is a true source of transmission, not just by risk assessment.

Table 3. The high-risk areas and purification areas selected by regular linear kernel regression.

| | High – risk area | Purification area |
|-----------------|--|---|
| $n=2$ | United States, India, Malaysia, Indonesia | Brazil, Russia, Mexico, Peru, China |
| $\lambda n=1$ | United States, India, Malaysia, Indonesia, | Brazil, United Kingdom, Russia, Mexico, Peru, Iraq, Czechia, Chile, Israel, Portugal, Jordan, United Arab Emirates, Oman, Bahrain, Equatorial Guinea, Saudi Arabia, China, Italy, Monaco. |
| $\lambda n=0.5$ | United States, India, Malaysia, Indonesia, Turkey, Japan | Brazil, United Kingdom, Russia, Mexico, Peru, Iraq, Czechia, Chile, Israel, Portugal, Jordan, United Arab Emirates, Oman, Bahrain, Equatorial Guinea, Saudi Arabia, China, Italy, Monaco. Poland, Ukraine, Netherlands, Czechia, Libya, Egypt, Zambia, Democratic Republic of Congo |

Table 4 The risk assessment of all 189 countries or areas.

| Index | Area | Risk | Index | Area | Risk |
|-------|----------------|-------|-------|------------------------------|-------|
| 1 | United States | 17.62 | 47 | Egypt | -1.29 |
| 2 | India | 13.6 | 48 | Czechia | -1.29 |
| 3 | Brazil | 7.31 | 49 | Georgia | -1.31 |
| 4 | Russia | 2.59 | 50 | Portugal | -1.31 |
| 5 | United Kingdom | 2.49 | 51 | Israel | -1.31 |
| 6 | France | 1.92 | 52 | Austria | -1.32 |
| 7 | Argentina | 1.39 | 53 | Saudi Arabia | -1.33 |
| 8 | Colombia | 1.15 | 54 | China | -1.33 |
| 9 | Mexico | 0.98 | 55 | Croatia | -1.34 |
| 10 | Spain | 0.98 | 56 | Kenya | -1.34 |
| 11 | Peru | 0.87 | 57 | Bosnia and Herzegovina | -1.35 |
| 12 | Italy | 0.69 | 58 | Cameroon | -1.36 |
| 13 | Germany | 0.66 | 59 | Serbia | -1.36 |
| 14 | Iran | 0.44 | 60 | Afghanistan | -1.37 |
| 15 | Indonesia | 0.37 | 61 | Kazakhstan | -1.37 |
| 16 | Turkey | 0.3 | 62 | Nigeria | -1.37 |
| 17 | South Africa | -0.18 | 63 | Oman | -1.38 |
| 18 | Poland | -0.18 | 64 | Chile | -1.38 |
| 19 | Ukraine | -0.25 | 65 | Uganda | -1.38 |
| 20 | Philippines | -0.54 | 66 | Algeria | -1.39 |
| 21 | Belgium | -0.55 | 67 | Jordan | -1.39 |
| 22 | Romania | -0.66 | 68 | Zimbabwe | -1.39 |
| 23 | Canada | -0.67 | 69 | Zambia | -1.39 |
| 24 | Netherlands | -0.9 | 70 | Ireland | -1.4 |
| 25 | Malaysia | -0.91 | 71 | Azerbaijan | -1.4 |
| 26 | Japan | -0.96 | 72 | Namibia | -1.4 |
| 27 | Vietnam | -0.99 | 73 | Tanzania | -1.41 |
| 28 | Bangladesh | -1.06 | 74 | Democratic Republic of Congo | -1.41 |
| 29 | Tunisia | -1.09 | 75 | Angola | -1.41 |
| 30 | Pakistan | -1.12 | 76 | Bahrain | -1.41 |
| 31 | Hungary | -1.14 | 77 | Libya | -1.42 |
| 32 | Nepal | -1.15 | 78 | Sudan | -1.42 |
| 33 | Sweden | -1.15 | 79 | United Arab Emirates | -1.42 |
| 34 | Ecuador | -1.16 | 80 | Mozambique | -1.42 |
| 35 | Bulgaria | -1.17 | 81 | Lebanon | -1.42 |
| 36 | Thailand | -1.17 | 82 | North Macedonia | -1.43 |
| 37 | Bolivia | -1.19 | 83 | Moldova | -1.43 |
| 38 | Greece | -1.19 | 84 | Slovakia | -1.43 |
| 39 | Iraq | -1.2 | 85 | Malawi | -1.43 |
| 40 | Pakistan | -1.21 | 86 | Rwanda | -1.44 |
| 41 | Paraguay | -1.21 | 87 | Somalia | -1.44 |
| 42 | Sri Lanka | -1.23 | 88 | Kuwait | -1.44 |
| 43 | Morocco | -1.23 | 89 | Lithuania | -1.44 |
| 44 | Switzerland | -1.24 | 90 | Armenia | -1.44 |
| 45 | Myanmar | -1.26 | 91 | Syria | -1.44 |
| 46 | Ethiopia | -1.28 | 92 | Botswana | -1.45 |

Table 5. The risk assessment of all 189 countries or areas.

| Index | Area | Risk | Index | Area | Risk |
|-------|-------------|-------|-------|--------------|-------|
| 93 | Belarus | -1.45 | 142 | Liberia | -1.52 |
| 94 | Yemen | -1.45 | 143 | Senegal | -1.52 |
| 95 | Burundi | -1.45 | 144 | Vatican | -1.52 |
| 96 | Comoros | -1.45 | 145 | Honduras | -1.52 |
| 97 | Maldives | -1.46 | 146 | Sierra Leone | -1.52 |
| 98 | South Korea | -1.46 | 147 | Gibraltar | -1.53 |
| 99 | Uruguay | -1.46 | 148 | Mauritania | -1.53 |
| 100 | Singapore | -1.46 | 149 | Gambia | -1.53 |

| | | | | | |
|-----|--------------------------|-------|-----|--------------------------|-------|
| 101 | Monaco | -1.46 | 150 | Luxembourg | -1.53 |
| 102 | Qatar | -1.46 | 151 | Guinea-Bissau | -1.54 |
| 103 | Denmark | -1.46 | 152 | San Marino | -1.54 |
| 104 | Slovenia | -1.46 | 153 | Brunei | -1.54 |
| 105 | South Sudan | -1.46 | 154 | Estonia | -1.54 |
| 106 | Eritrea | -1.46 | 155 | Venezuela | -1.54 |
| 107 | Seychelles | -1.46 | 156 | Liechtenstein | -1.54 |
| 108 | Uzbekistan | -1.46 | 157 | Andorra | -1.55 |
| 109 | Australia | -1.46 | 158 | Martinique | -1.55 |
| 110 | Mauritius | -1.46 | 159 | Cuba | -1.55 |
| 111 | Djibouti | -1.47 | 160 | Timor | -1.56 |
| 112 | Cambodia | -1.47 | 161 | Cape Verde | -1.56 |
| 113 | Chad | -1.47 | 162 | Dominican Republic | -1.58 |
| 114 | Lesotho | -1.47 | 163 | Panama | -1.58 |
| 115 | Norway | -1.47 | 164 | Trinidad and Tobago | -1.59 |
| 116 | Congo | -1.47 | 165 | Papua New Guinea | -1.59 |
| 117 | Niger | -1.47 | 166 | Iceland | -1.6 |
| 118 | Equatorial Guinea | -1.48 | 167 | Suriname | -1.61 |
| 119 | Central African Republic | -1.48 | 168 | Guyana | -1.63 |
| 120 | Gabon | -1.48 | 169 | New Zealand | -1.64 |
| 121 | Ghana | -1.49 | 170 | New Caledonia | -1.64 |
| 122 | Sao Tome and Principe | -1.49 | 171 | Jamaica | -1.65 |
| 123 | Cyprus | -1.49 | 172 | Barbados | -1.65 |
| 124 | Burkina Faso | -1.49 | 173 | Saint Lucia | -1.65 |
| 125 | Cote d'Ivoire | -1.49 | 174 | Grenada | -1.65 |
| 126 | Albania | -1.5 | 175 | Antigua and Barbuda | -1.65 |
| 127 | Benin | -1.5 | 176 | Dominica | -1.66 |
| 128 | Togo | -1.5 | 177 | Saint Kitts and Nevis | -1.66 |
| 129 | Mali | -1.5 | 178 | Anguilla | -1.66 |
| 130 | Tajikistan | -1.5 | 179 | El Salvador | -1.66 |
| 131 | Latvia | -1.5 | 180 | Haiti | -1.66 |
| 132 | Guatemala | -1.5 | 181 | Montserrat | -1.66 |
| 133 | Finland | -1.5 | 182 | Bahamas | -1.67 |
| 134 | Montenegro | -1.5 | 183 | Aruba | -1.67 |
| 135 | Costa Rica | -1.51 | 184 | Turks and Caicos Islands | -1.67 |
| 136 | Kyrgyzstan | -1.51 | 185 | Fiji | -1.67 |
| 137 | Mongolia | -1.51 | 186 | Belize | -1.69 |
| 138 | Malta | -1.51 | 187 | Nicaragua | -1.69 |
| 139 | Guinea | -1.51 | 188 | Samoa | -1.7 |
| 140 | Laos | -1.52 | 189 | French Polynesia | -1.72 |
| 141 | Bhutan | -1.52 | | | |

Conclusion

In this paper, we considered the condition of Kernel selection consistence, in Regular Linear Kernel Regression Model with l_1 -norm. And bounds on the prediction error of RLKR are also given with mild conditions. Then we applied the RLKR to high-risk selection and risk assessment for COVID-19, in a view point of area-to-area spreading mode. In the simulation of random design matrix, we explored the F-also-inclusion and Consistency under different levels of λ_n and the number m of candidate kernels. The prediction error bounds were also checked with a no-wise result. In real world data of COVID-19, the risk assessments of 189 numbers of countries or areas were given. We believe this will provide valid advices on traveling.

Acknowledgement

Our thanks go to Professor Bo Li for stimulating this research. Our

research is supported by NSF of China (Grant No. 61877023) and the Fundamental Research Funds for the Central Universities (Grant No. CCNU19TD009).

References

1. Fauci Anthony S, H. Clifford Lane and Robert R. Redfield. "Covid-19 - Navigating the Uncharted". *N Engl J Med* 382(2020):1268-1269.
2. Yuki Koichi, Miho Fujiogi and Sophia Koutsogiannaki. "Covid-19 Pathophysiology: A Review". *Clin Immunol* 215(2020):108427.
3. John H Beigel, Kay M. Tomashek, Lori E. Dodd and Aneesh K Mehta, et al. "Remdesivir for the Treatment of Covid-19 - Final Report". *N Engl J Med* 383(2020):992-994.
4. Felsenstein Susanna, Jenny A. Herbert, Paul S. McNamara and Christian M. Hedrich. "Covid-19: Immunology and Treatment Options". *Clin Immunol* 215(2020):108448

5. Lotfi Melika, Michael R.Hamblin and NimaRezaei. "Covid-19: Transmission, Prevention, and Potential Therapeutic Opportunities". *Cli. Chim Acta* 508(2020):254-266.
6. Ndairou Faical, Ivan Area, Juan J.Nieto and Delfim F.M.Torres. "Mathematical Modeling of Covid-19 Transmission Dynamics with a Case Study of Wuhan". *Chaos Soli* 135 (2020):109846.
7. Gostic Katelyn, Ana CR Gomez, Riley O Mummah and Adam J Kucharski et al. "Estimated Effectiveness of Symptom and Risk Screening to Prevent the Spread of Covid-19". *Elife* 9(2020):55570.
8. Bleier Benjamin S, Murugappan Ramanathan and Andrew P. Lane. "Covid-19 Vaccines may not Prevent Nasal Sars-Cov-2 Infection and Asymptomatic Transmission". *Otolaryngology. J OTOLARYNGOL-HEAD N* 164(2021):305-307.
9. Osborne Michael R, Brett Presnell and Berwin A. Turlach. "On the Lasso and its Dual". *J Comput Graph Stat* 9(2000):319-337.
10. Zhao Peng and Bin Yu. "On Model Selection Consistency of Lasso". *J Mach Learn Res* 7(2006):2541-2563.
11. Hastie Trevor, Andrea Montanari, Saharon Rosset and Ryan J. Tibshirani. "Surprises in High-Dimensional Ridgeless Least Squares Interpolation".1-6(2020).
12. Bartlett Peter L, Philip M. Long, Gabor Lugosi and Alexander Tsigler. "Benign Overfitting in Linear regression". *Proceedings of the National Academy of Sciences* 117(2020):30063-30070.
13. Tuo Rui, Yan Wang and C. F. Jeff Wu. "On the Improved Rates of Convergence for Matérn-type Kernel Ridge Regression, with Application to Calibration of Computer Models". *Int J Uncertain Quantif* 8(2020):1522-1547.
14. Fanghui Liu, Zhenyu Liao and Johan A.K. Suykens. "Kernel Regression in High Dimensions: Refined Analysis beyond Double Descent". *Proceedings of Machine Learning Research*. 130(2021):649-657.
15. Lorin Crawford, Kris C. Wood, Xiang Zhou and Sayan Mukherjee. "Bayesian Approximate Kernel Regression with Variable Selection". *J Am Stat Assoc* 113(2018):1710-1721.
16. Liang Tengyuan and Alexander Rakhlin. "Just Interpolate: Kernel "Ridgeless" Regression can Generalize". *Ann Stat* 48(2020):1329-1347.
17. Braga Igor and Maria Carolina Monard. "Improving the Kernel Regularized Least Squares Method for Small-Sample Regression". *Neurocomputing* 163(2015):106-114.

How to cite this article: Lu xan, Ba lin "Consistence Condition of Kernel Selection In Regular Linear Kernel Regression And Its Application In Covid-19 High-risk Areas Exploration." *Virol Curr Res.S4* (2023): 003.