

Kernel Regression Utilizing External Information as Constraints

Chi-Shian Dai
Advisor: Jun Shao

- 1 Motivation
- 2 CKR for Summary Level External Data
- 3 CKR for Individual Level External Data
- 4 Simulation

Motivation

- Internal Data: $\{Y_i, \mathbf{U}_i\}_{i=1,\dots,n}$, $\mathbf{U}_i = (\mathbf{X}_i, \mathbf{Z}_i) \in \mathbb{R}^p$, $\mathbf{X} \in \mathbb{R}^q$.
- External Data: Sample size is $m \gg n$. Provide information for $Y \sim \mathbf{X}$.
 - Sources: Population-based census, Past studies...
 - Summary Level Information: Least square estimator $\hat{\beta}$ via $Y \sim \mathbf{X}$.
 - Individual Level information: $\{Y_i, \mathbf{X}_i\}_{i=n+1,\dots,n+m}$.
- Goal: Estimate $E[Y|\mathbf{U} = \mathbf{u}] := \mu(\mathbf{u})$

- Inspiring by [?], they observe that the link between “internal” and “external” can be formulated as constraints. Hence, we consider a constrained kernel regression to estimate μ .
- Example: Let $Y = \beta^\top \mathbf{X} + \gamma^\top \mathbf{Z} + \epsilon$, $\mathbf{X} \perp \mathbf{Z}$, then there is a naive constraints

$$\beta = \hat{\beta}.$$

So, with the help of external data, we only have to focus on estimating γ .

- We generalize this idea to kernel regression. The proposed method call constrained kernel regression (CKR).

- 1 Motivation
- 2 CKR for Summary Level External Data
- 3 CKR for Individual Level External Data
- 4 Simulation

Optimization Form of Kernel Regression

- Given a kernel κ and bandwidths l , and b .
- Kernel Regression estimate for $\mu(\mathbf{u})$:

$$\hat{\mu}_K(\mathbf{u}) = \arg \min_{\mu} \sum_{j=1}^n \kappa_l(\mathbf{u} - \mathbf{U}_j)(Y_j - \mu)^2, \quad (1)$$

where $\kappa_l(\mathbf{u} - \mathbf{U}_j) = l^{-p} \kappa \{ l^{-1}(\mathbf{u} - \mathbf{U}_j) \}$.

- Kernel Regression estimate for $\boldsymbol{\mu} := (\mu_1, \dots, \mu_n)$, $\mu_i = \mu(\mathbf{U}_i)$:

$$\hat{\boldsymbol{\mu}}_K = \arg \min_{\mu_1, \dots, \mu_n} \sum_{i=1}^n \sum_{j=1}^n \kappa_l(\mathbf{U}_i - \mathbf{U}_j)(Y_j - \mu_i)^2$$

Optimization Form of Kernel Regression

- There is another equivalent form.

$$\hat{\mu}_K = \arg \min_{\mu_1, \dots, \mu_n} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \frac{\kappa_l(\mathbf{u}_i - \mathbf{u}_j)}{\sum_{k=1}^n \kappa_l(\mathbf{u}_i - \mathbf{u}_k)} (Y_j - \mu_i)^2 \quad (2)$$

- We prefer (2) since

$$\sum_{j=1}^n \frac{\kappa_l(\mathbf{u}_i - \mathbf{u}_j)}{\sum_{k=1}^n \kappa_l(\mathbf{u}_i - \mathbf{u}_k)} (Y_j - \mu_i)^2 \approx E[(Y - \mu(\mathbf{u}))^2 | \mathbf{u} = \mathbf{u}_i],$$

and

$$\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \frac{\kappa_l(\mathbf{u}_i - \mathbf{u}_j)}{\sum_{k=1}^n \kappa_l(\mathbf{u}_i - \mathbf{u}_k)} (Y_j - \mu_i)^2 \approx E[(Y - \mu(\mathbf{u}))^2]$$

Constraints for Summary Level External Data

- $\hat{\beta}$ is a consistent estimate for $\beta_0 := E[\mathbf{X}\mathbf{X}^\top]^{-1}E[\mathbf{X}Y]$, which satisfy

$$E\{(Y - \mathbf{X}^\top \beta_0)\mathbf{X}\} = 0.$$

Hence, the constrained optimization can be

$$\begin{aligned} \hat{\mu} = \arg \min_{\mu_1, \dots, \mu_n} & \sum_{i=1}^n \sum_{j=1}^n \frac{\kappa_l(\mathbf{U}_i - \mathbf{U}_j)}{\sum_{k=1}^n \kappa_l(\mathbf{U}_i - \mathbf{U}_k)} (Y_j - \mu_i)^2 \\ \text{subject to} & \sum_{i=1}^n (\mu_i - \mathbf{x}_i^\top \hat{\beta}) \mathbf{x}_i = 0. \end{aligned} \quad (3)$$

- (3) is a quadratic programming. Hence, it can be solved by Lagrange multiplier.
- For arbitrary $\mathbf{u} \in \mathcal{U}$, we apply additional kernel regression by replacing Y with $\hat{\mu}$.

$$\hat{\mu}_{CK}(\mathbf{u}) = \sum_{i=1}^n \hat{\mu}_i \kappa_b(\mathbf{u} - \mathbf{u}_i) / \sum_{i=1}^n \kappa_b(\mathbf{u} - \mathbf{u}_i). \quad (4)$$

- Kernel Regression (KR): $\hat{\mu}_K(\mathbf{u})$
- Double Kernel Regression (DKR): Consider CKR without applying any constraints. Use notation $\hat{\mu}_{DK}(\mathbf{u})$.
 - Step 1: Estimate $(\mu(\mathbf{U}_1), \dots, \mu(\mathbf{U}_n))$ by Kernel Regression
 - Step 2: Estimate $\mu(\mathbf{u})$ by additional kernel regression replacing Y with the results in first step.

Assumption for Normality

- (A1) The response Y has a finite $E|Y|^s$ with $s > 2 + p/2$. The covariate vector \mathbf{U} has compact support \mathcal{U} , and has a positive definite covariance matrix. The density of \mathbf{U} is bounded away from infinity and zero, and has bounded second-order derivatives.
- (A2) Functions $\mu(\mathbf{u}) = E(Y|\mathbf{U} = \mathbf{u})$ and $\sigma^2(\mathbf{u}) = E[\{Y - \mu(\mathbf{U})\}^2|\mathbf{U} = \mathbf{u}]$ are Lipschitz continuous. The function $\mu(\mathbf{u})$ has bounded third-order derivatives.
- (A3) The kernel κ is a positive bounded density with mean zero and finite sixth moments. Furthermore, κ is Lipschitz continuous.
- (A4) The bandwidths b in (1) and l in (3) are polynomial rate of n , satisfying $b \rightarrow 0$, $l \rightarrow 0$, $l/b \rightarrow \gamma \in (0, \infty)$, and $nb^{4+p} \rightarrow c \in [0, \infty)$ as the internal sample size $n \rightarrow \infty$.
- (A5) The external sample size m satisfies $n = O(m)$, i.e., n/m is bounded by a fixed constant.

Theorem

Theorem 1

Assume conditions (A1)-(A5). Then, as $n \rightarrow \infty$,

$$\sqrt{nb^p}\{\hat{\mu}_t(\mathbf{u}) - \mu(\mathbf{u})\} \rightarrow N(B_t(\mathbf{u}), V_t(\mathbf{u})) \text{ in distribution,} \quad (5)$$

where $t = DK$ or CK ,

$$B_{DK}(\mathbf{u}) = c^{1/2}(1 + \gamma^2)A(\mathbf{u}),$$

$$B_{CK}(\mathbf{u}) = c^{1/2}[(1 + \gamma^2)A(\mathbf{u}) - \gamma^2 \mathbf{x}^\top \Sigma_X^{-1} E\{\mathbf{X}A(\mathbf{U})\}],$$

$$V_{DK}(\mathbf{u}) = \frac{\sigma^2(\mathbf{u})}{f_U(\mathbf{u})} \int \left\{ \int \kappa(\mathbf{w} - \mathbf{v}\gamma) \kappa(\mathbf{v}) d\mathbf{v} \right\}^2 d\mathbf{w},$$

$$V_{CK}(\mathbf{u}) = V_{DK}(\mathbf{u}),$$

$$A(\mathbf{u}) = \int \kappa(\mathbf{v}) \left\{ \frac{1}{2} \mathbf{v}^\top \nabla^2 \mu(\mathbf{u}) \mathbf{v} + \nabla \mu(\mathbf{u})^\top \mathbf{v} \mathbf{v}^\top \nabla \log f_U(\mathbf{u}) \right\} d\mathbf{v}, \quad (6)$$

and f_U is the density of \mathbf{U} .

Asymptotic Mean Integrated Square Error

- $\text{AMISE}(\hat{\mu}_t) = E[\{B_t(\mathbf{U})\}^2 + V_t(\mathbf{U})]$, $t = CK$ or DK ,
- Observe that $\boldsymbol{\xi} := \Sigma_X^{-1} E\{\mathbf{X}A(\mathbf{U})\}$ is a linear coefficient of fitting $A(\mathbf{U})$ by \mathbf{X} . Hence,

$$E[\{A(\mathbf{U}) - \mathbf{X}^\top \boldsymbol{\xi}\} \mathbf{X}] = 0.$$

- From this we can derive

$$\begin{aligned} E\{B_{CK}(\mathbf{U})\}^2 &= c(1 + \gamma^2)^2 E\{A(\mathbf{U}) - \mathbf{X}^\top \boldsymbol{\xi}\}^2 + cE\{\mathbf{X}^\top \boldsymbol{\xi}\}^2, \\ E\{B_{DK}(\mathbf{U})\}^2 &= c(1 + \gamma^2)^2 E\{A(\mathbf{U})\}^2 \\ &= c(1 + \gamma^2)^2 E\{A(\mathbf{U}) - \mathbf{X}^\top \boldsymbol{\xi}\}^2 + c(1 + \gamma^2)^2 E\{\mathbf{X}^\top \boldsymbol{\xi}\}^2. \end{aligned}$$

- $B_K(\mathbf{u}) = c^{1/2}A(\mathbf{u})$ and $V_K(\mathbf{u}) = \frac{\sigma^2(\mathbf{u})}{f_U(\mathbf{u})} \int \{\kappa(\mathbf{v})\}^2 d\mathbf{v}$.
- $B_{CK}(\mathbf{u}) = c^{1/2}[(1 + \gamma^2)A(\mathbf{u}) - \gamma^2 \mathbf{x}^\top \Sigma_X^{-1} E\{\mathbf{X}A(\mathbf{U})\}]$,
- $V_{CK}(\mathbf{u}) = \frac{\sigma^2(\mathbf{u})}{f_U(\mathbf{u})} \int \left\{ \int \kappa(\mathbf{w} - \mathbf{v}\gamma) \kappa(\mathbf{v}) d\mathbf{v} \right\}^2 d\mathbf{w}$,
- If $\gamma = 0$, $B_{CK} = B_K$ and $V_{CK} = V_K$. So,
 $\inf_{\gamma} \text{AMISE}(\hat{\mu}_{CK})(\gamma) \leq \text{AMISE}(\hat{\mu}_K)$
- If γ increases, V_{CK} decreases and B_{CK}^2 increases.

Theorem 2

Under the conditions in Theorem 1 and an additional condition that $\int \nabla^2 \kappa(\mathbf{u}) \kappa(\mathbf{u}) d\mathbf{u}$ being strictly negative definite, $\text{AMISE}(\hat{\mu}_{CK}) < \text{AMISE}(\hat{\mu}_K)$ for c and γ in a neighborhood of 0.

- 1 Motivation
- 2 CKR for Summary Level External Data
- 3 CKR for Individual Level External Data**
- 4 Simulation

CKR for Individual Level External Data

- First, estimate $h = E[Y|X]$ via kernel regression.
- Second, observe that for all real function g

$$E\{Y - h(\mathbf{X})\}g(\mathbf{X}) = 0.$$

- Consider the corresponding constrained optimization.

$$\begin{aligned} \hat{\mu} = \arg \min_{\mu_1, \dots, \mu_n} & \sum_{i=1}^n \sum_{j=1}^n \frac{\kappa_l(\mathbf{U}_i - \mathbf{U}_j)}{\sum_{k=1}^n \kappa_l(\mathbf{U}_i - \mathbf{U}_k)} (Y_j - \mu_i)^2 \\ \text{subject to} & \sum_i^n \{\mu_i - \hat{h}(\mathbf{X}_i)\} g(\mathbf{X}_i) = 0. \end{aligned} \quad (7)$$

- Question: How to choose g ?

Theorem 3

Assume (A1)-(A5) in Theorem 1 and

- (A1') The matrix $\Sigma_g = E\{g(\mathbf{X})g(\mathbf{X})^\top\}$ is positive definite.
- (A2') The functions $h(\mathbf{x})$ and $g(\mathbf{x})$ are Lipschitz continuous. The function $h(\mathbf{x})$ has bounded third-order derivatives.
- (A3') The kernel function used in the kernel regression based on the external dataset satisfies condition (A3).
- (A4') The bandwidth used in the kernel regression based on the external dataset is of the order $m^{-1/(4+q)}$ as $m \rightarrow \infty$.

Then, as $n \rightarrow \infty$, CKR with constraints (7) have following properties.

$$\sqrt{nb^p}\{\hat{\mu}_{CK}(\mathbf{u}) - \mu(\mathbf{u})\} \rightarrow N(B_{CK}(\mathbf{u}), V_{CK}(\mathbf{u})) \quad \text{in distribution,}$$

where

$$B_{CK}(\mathbf{u}) = c^{1/2}[(1 + \gamma^2)A(\mathbf{u}) - \gamma^2 g(\mathbf{x})^\top \Sigma_g^{-1} E\{g(\mathbf{X})A(\mathbf{U})\}]$$

and $V_{CK}(\mathbf{u})$ and $A(\mathbf{u})$ are the same as those in Theorem 1.

Choose of g

- $E\{B_{CK}(\mathbf{U})\}^2 = c\{(1+\gamma^2)^2 - 1\}E\{A(\mathbf{U}) - g(\mathbf{X})^\top \boldsymbol{\xi}_g\}^2 + cE\{A(\mathbf{U})\}^2$,
where $\boldsymbol{\xi}_g = \boldsymbol{\Sigma}_g^{-1}E\{g(\mathbf{X})A(\mathbf{U})\}$ is a linear coefficient of $A(\mathbf{U}) \sim g(\mathbf{X})$.
- Find g to minimize $E\{A(\mathbf{U}) - g(\mathbf{X})^\top \boldsymbol{\xi}_g\}^2$.
- The best one is $g^* = E[A(\mathbf{U})|\mathbf{X}]$.
- $A(\mathbf{u}) = \int \kappa(\mathbf{v}) \left\{ \frac{1}{2} \mathbf{v}^\top \nabla^2 \mu(\mathbf{u}) \mathbf{v} + \nabla \mu(\mathbf{u})^\top \mathbf{v} \mathbf{v}^\top \nabla \log f_U(\mathbf{u}) \right\} d\mathbf{v}$,
is estimable.
- g^* is also estimable.

Theorem 4

Assume the conditions in Theorem 3 and the following additional conditions.

- (C1) The kernel κ in (A3) satisfies $\int u_k^2 \kappa(\mathbf{u}) d\mathbf{u} = 1$ and $\int u_k u_j \kappa(\mathbf{u}) d\mathbf{u} = 0$ when $k \neq j$. The kernel $\tilde{\kappa}$ in the estimators $\hat{\nu}_k$ and $\nabla_{kk}^2 \hat{f}_U$, $k = 1, \dots, p$, has finite second-order moments, bounded $\nabla_{kk}^2 \tilde{\kappa}$, finite $\int |\nabla_{kk}^2 \tilde{\kappa}(\mathbf{u})| d\mathbf{u}$, and bounded $\sup_{\mathbf{u}} \lambda^{-2} |\tilde{\kappa}(\mathbf{u}/\lambda)|$ and $\sup_{\mathbf{u}} \lambda^{-3} |\nabla_k \tilde{\kappa}(\mathbf{u}/\lambda)|$ as $\lambda \rightarrow 0$, $k = 1, \dots, p$.
- (C2) The bandwidth λ_1 for $\hat{\nu}_0$ and \hat{f}_U has order $n^{-1/(p+4)}$, the bandwidth λ_2 for $\hat{\nu}_k$ and $\nabla_{kk}^2 \hat{f}_U$ has order $n^{-1/(p+8)}$, and the bandwidth δ in estimating g^* has order $n^{-1/(q+4)}$.

Then, the result in Theorem 3 with $g = g^*$ holds for $\hat{\mu}_{CK}$ using the estimated constraints \hat{g}^* .

- 1 Motivation
- 2 CKR for Summary Level External Data
- 3 CKR for Individual Level External Data
- 4 Simulation

- Internal sample size : 200
- External sample size : 1000
- X, Z are normal distribution with variance 1, covariance 0.5.
- $Y = \mu(X, Z) + \epsilon; \epsilon \sim N(0, \sigma^2)$

① Additive Models:

$$\mu = X^3 + Z^2$$

$$\mu = 2^{-1}\cos(2X) + \cos(Z)$$

$$\mu = \cos(X) + \cos(Z)$$

② Non-Additive Models:

$$\mu = X^3 + XZ + Z^2$$

Simulation

- Let $R = 200$ be the number of independently replication.
- Let $L = 121$ be the sample size of test data.
 1. Fixed grid points on $[-1, 1] \times [-1, 1]$
 2. Random sample without replacement from the covariate \mathbf{U}' s of the internal data set.
- Use estimated MISE to evaluate performance.

$$\text{MISE} = \frac{1}{R} \sum_{r=1}^R \frac{1}{L} \sum_{l=1}^L \{\hat{\mu}_r(\mathbf{T}_{r,l}) - \mu(\mathbf{T}_{r,l})\}^2,$$

- Best Bandwidth: Evaluate MISE in a pool of bandwidths and display the one have the best performance.
- 10 folds cross-validation

-

$$\text{Imp}\% = 1 - \frac{\min\{\text{MISE}(\hat{\mu}_{CKR}) \text{ over all CKR methods}\}}{\min\{\text{MISE}(\hat{\mu}_K), \text{MISE}(\hat{\mu}_{DK})\}}.$$

Test Data: Sample

model	σ	test data	b, l	estimator (CKR) with $g =$					estimator			Imp%
				1	$(1, X)$	$(1, \hat{h})$	\hat{g}^*	g^*	(CKR-s)	(KR)	(DKR)	
1	3	sample	best	1.045	0.924	0.912	1.055	0.843	1.152	1.081	1.056	20.17
			CV	1.165	1.148	1.073	1.19	1.176	1.409	1.239	1.181	9.14
2	3	sample	best	0.220	0.225	0.222	0.259	0.210	0.201	0.266	0.260	22.69
			CV	0.297	0.290	0.323	0.341	0.282	0.274	0.335	0.343	18.20
3	3	sample	best	0.338	0.347	0.333	0.437	0.365	0.298	0.443	0.439	32.13
			CV	0.537	0.512	0.581	0.643	0.539	0.47	0.620	0.640	24.19
4	3	sample	best	1.102	1.066	1.018	1.102	1.020	1.437	1.117	1.099	7.37
			CV	1.276	1.294	1.329	1.262	1.410	1.624	1.208	1.270	-4.47

Test Data: Fixed Grid Points on $[-1, 1]^2$

model	σ	test data	b, l	estimator (CKR) with $g =$					estimator			Imp%
				1	$(1, X)$	$(1, \hat{h})$	\hat{g}^*	g^*	(CKR-s)	(KR)	(DKR)	
1	3	grid	best	0.175	0.171	0.181	0.205	0.206	0.243	0.210	0.208	17.78
			CV	0.382	0.365	0.341	0.388	0.355	0.432	0.412	0.384	11.19
2	3	grid	best	0.111	0.108	0.076	0.148	0.132	0.100	0.153	0.153	50.32
			CV	0.141	0.138	0.117	0.164	0.151	0.134	0.160	0.162	26.87
3	3	grid	best	0.102	0.100	0.070	0.129	0.131	0.089	0.135	0.132	46.96
			CV	0.123	0.122	0.099	0.145	0.151	0.121	0.142	0.142	30.28
4	3	grid	best	0.231	0.244	0.229	0.250	0.251	0.338	0.251	0.251	8.76
			CV	0.414	0.426	0.359	0.400	0.369	0.486	0.367	0.397	2.17

- ① We will study more when and why the CKR improves the standard KR, both empirically and theoretically.
- ② With individual-level external data, we only considered kernel estimators of $h(\mathbf{x})$. Different types of estimators of h will be studied—for example, linear models or generalized linear models.
- ③ We mainly considered the AMISE as a criterion in evaluating kernel estimators. We will study some other performance measures.
- ④ We focused on the situation where the underlying distribution of \mathbf{U} is the same in both internal and external data. We will consider extensions in situations where the distributions of \mathbf{U} are different in two datasets.

References