post-selection inference (Kuchibhotla, 2022)

Setting: $Q = \{q\}$ all possible selection (model, covariates, transformation, ...)

interested in misspecification-robust target $\theta_q = \arg\min\limits_{\theta \in \Theta_q} \frac{1}{n}\sum\limits_{i=1}^{n} E[l_q(\theta, Z_i)]$

① It's only a loss function, no need to give a param model

stage 1: select $\hat{q}$ (with some data-driven procedure)

② $\theta_q$ can depend on data generating process

stage 2: estimate $\theta_{\hat{q}}$ using $\hat{\theta}_{\hat{q}}$ from data.

$\{Z_i Y_i\}_{i=1}^{n}$

$\hat{q} \longrightarrow \hat{\theta}_{\hat{q}}$ | ?? $\liminf\limits_{n\to\infty} P(\theta_{\hat{q}} \in \widehat{CI}_{\hat{q}}) \geq 1-\alpha$, $\widehat{CI}_{\hat{q}}$ based on $\hat{\theta}_{\hat{q}}$ | example of VIDE.

$q \longrightarrow \hat{\theta}_q$ easy to find $\widehat{CI}_q$ s.t. $\liminf\limits_{n\to\infty} P(\theta_q \in \widehat{CI}_q) \geq 1-\alpha$

Further reading: difference with dimension reduction (Berk 2013)

Solution 1: Data Splitting. $\{Z_i Y_i\}_{i=1}^{m} \perp \{Z_i Y_i\}_{i=m+1}^{n}$ $\qquad P(\hat{\theta}_{\hat{q}} - \theta_{\hat{q}} \in A | \hat{q} = q) = P(\hat{\theta}_q - \theta_q \in A)$

$\downarrow \qquad\qquad \downarrow$

$\hat{q} \longrightarrow \hat{\theta}_{\hat{q}}$

pro: no restriction on selection procedure

conformal inference

$X_i, Y_i \qquad \hat{Y}_i$

cons: ① unacceptable model (?)

② invalid for dependent data.

③ effect brought by splitting (size & randomness).

Note: in this case, the inference being conservative is because: ① it should be (original $\hat{\theta}_{\hat{q}}$ biased)
② smaller sample size.

Solution 2: Simultaneous Inference.

$\forall \hat{q} \in Q. \quad P(\theta_{\hat{q}} \in \widehat{CI}_{\hat{q}}) \geq \underline{P(\bigcap\limits_{q\in Q}\{\theta_q \in \widehat{CI}_q\})}$ $\qquad$ Controlling this needs to construct $\widehat{CI}_q$ for $\forall q \in Q$

$\liminf\limits_{n\to\infty} P(\theta_{\hat{q}} \in \widehat{CI}_{\hat{q}}) \geq \liminf\limits_{n\to\infty} P(\bigcap\limits_{q\in Q}\{\theta_q \in \widehat{CI}_q\}) \geq 1-\alpha$ $\longrightarrow$ ?? have to ~~use~~ do simultaneous inference for arbitrary selection procedure?

Construction procedure: Assumptions: ① uniform asymptotic linear representation, $\exists \{\psi_q(\cdot), \forall q \in Q\}$

$\boxed{\max\limits_{q\in Q}} |\psi_{n,q}^{-\frac{1}{2}}(\hat{\theta}_q - \theta_q - \frac{1}{n}\sum\limits_{i=1}^{n}\psi_q(Z_i))| = o_P(\frac{1}{\sqrt{n}})$ $\qquad \frac{\sqrt{n}(\hat{\theta}_q - \theta_q)}{\psi_{n,q}} \approx \frac{\sqrt{n}(\frac{1}{n}\sum\limits_{i=1}^{n}\psi_q(Z_i))}{\psi_{n,q}}$

② Some conditions guarantee $\frac{\frac{1}{\sqrt{n}}\sum\limits_{i=1}^{n}\psi_q(Z_i)}{\psi_{n,q}} \xrightarrow{d} N(0,1)$.

($\{Z_i Y_i\}$ needs to be independent / weakly dependent).

$\implies (\frac{1}{\sqrt{n}}\psi_{n,q}^{-\frac{1}{2}}(\hat{\theta}_q - \theta_q) : q \in Q) \xrightarrow{d} (G_q : q \in Q) \sim N(0, R)$ $\quad diag(R) = 1.$

$\max |\cdot| \qquad\qquad\qquad \longrightarrow \max\limits_{q} |G_q| \qquad$ correlation between $q$ unknown.

$\implies \lim\limits_{n\to\infty} P(\max\limits_{q\in Q} |\sqrt{n}\,\psi_{n,q}^{-\frac{1}{2}}(\hat{\theta}_q - \theta_q)| \leq K_\alpha) = 1-\alpha \qquad K_\alpha$: upper $\alpha$ quantile of $\max\limits_{q}|G_q|$

$\implies \widehat{CI}_q = [\hat{\theta}_q - K_\alpha\sqrt{\psi_{n,q}/n},\ \hat{\theta}_q + K_\alpha\sqrt{\psi_{n,q}/n}] \qquad K_\alpha, \psi_{n,q}$ unknown.

need to approximate $K_\alpha$, $\gamma_{n,q}$ using bootstrap.

$$\widehat{CI}_q = \left[ \hat{\theta}_q - \hat{K}_\alpha \frac{\hat{\gamma}_{n,q}^{\frac{1}{2}}}{\sqrt{n}}, \; \hat{\theta}_q + \hat{K}_\alpha \frac{\hat{\gamma}_{n,q}^{\frac{1}{2}}}{\sqrt{n}} \right]$$

$$\widehat{CI}_q^{unadj} = \left[ \hat{\theta}_q - Z_{\frac{\alpha}{2}} \frac{\hat{\gamma}_{n,q}^{\frac{1}{2}}}{\sqrt{n}}, \; \hat{\theta}_q + Z_{\frac{\alpha}{2}} \frac{\hat{\gamma}_{n,q}^{\frac{1}{2}}}{\sqrt{n}} \right]$$

more conservative $\dfrac{\hat{K}_\alpha}{Z_{\frac{\alpha}{2}}} \geq 1$

~~grow~~ ratio grows with dimension. ( dim↑ conservative↑ )

$\hat{K}_\alpha$ depends on correlation.

pros: ① arbitrary selection procedure ( graphical / repeated / multiple model )

② change to block bootstrap: apply to dependent data ( ?? will CLT hold )

③ better selection result on larger sample.

Cons: ① might be too conservative ~~espec~~ ( but not for arbitrary selection method ).

② still only handles weakly dependent data.

③ ~~#~~ need to do calculation on all $q \in Q$.

④ $Q$ needs to be decided before data exploration.

⑤ Too conservative if $P(\hat{q} \in Q_0 \subset Q_u)$ is large & $Q_0$ is small in size.

Assumptions might not be needed ~~so~~ for some Gaussian models.

Solution 3: Conditional Selective Inference.

interested in: $\forall q$, $\displaystyle\liminf_{n\to\infty} P(\theta_q \in \widehat{CI}_q \mid \hat{q} = q) \geq 1 - \alpha$

Idea: with ~~Gaussian assumption~~ (CLT), condition on subset of $\text{span}(\{z_i\})$ that decides the selection $q$.

Assumptions: 1. $\exists$ random vector $D_{n,q} \in \mathbb{R}^{d_D}$ s.t. it decides selection: $\{\hat{q} = q\} \equiv \{D_{n,q} \leq 0\}$

2. nonzero denominator: $\displaystyle\liminf_{n\to\infty} P(\hat{q} = q) = \liminf_{n\to\infty} (D_{n,q} \leq 0) > 0$ ( often fails, but can be dropped ).

3. $\exists$ vec $\mu_{n,q} \in \mathbb{R}^{d_D}$ & cov mat $\Omega_q$ s.t.

$$\begin{bmatrix} \sqrt{n}(\hat{\theta}_q - \theta_q) \\ D_{n,q} - \mu_{n,q} \end{bmatrix} \xrightarrow{d} \begin{bmatrix} G_{\theta q} \\ G_{D q} \end{bmatrix} \sim N(0, \Omega_q)$$

$\Omega_q$ P.D. $\implies$ $\displaystyle\sup_C \left| P\left( \begin{bmatrix} \sqrt{n}(\hat{\theta}_q - \theta_q) \\ D_{n,q} - \mu_{n,q} \end{bmatrix} \in C \right) - P\left( \begin{bmatrix} G_{\theta q} \\ G_{D q} \end{bmatrix} \in C \right) \right| \to 0$

↳ need this uniformity.

4. $\Omega_q$ estimable ( for later use )

$$\hat{\Omega}_q = \begin{pmatrix} \hat{W}_q^2 & \hat{\Omega}_{\theta D} \\ \hat{\Omega}_{D\theta} & \hat{\Omega}_{DD} \end{pmatrix} \xrightarrow{P} \Omega_q = \begin{pmatrix} W_q^2 & \Omega_{\theta D} \\ \Omega_{D\theta} & \Omega_{DD} \end{pmatrix}$$

Algorithm: not copied here.

intuitive understanding: use $\hat{\theta}_q$, $\hat{\Omega}_q$ to construct an empirical conditional distribut and build CI on it.

Improvements: ① Data carving: randomized $D_{n,q} \implies$ bdd expected CI length.

② combine simul... & selective inference $\implies$ shorter length.

pros: ① ~~sle~~ selection on full data.

② computationally easier.

③ result can vary in naive CI (short) to data splitting CI (larger).

Cons: ① $\{\hat{p}=\hat{q}=q\} \equiv \{D_{n,q} \lesssim 0\}$ too strong. (restriction on selection).

② Theoretical analysis $\neq$ in every new implementation.

③ Vanilla version may yields infinite width of CI.

## Uniform Validity

$$\underset{n\to\infty}{\text{Liminf}} \underset{P\in P^{\partial n}}{\inf} P(\theta_{\hat{q}} \in \widehat{CI}_{\hat{q}}) \geq 1-\alpha$$

$$\cdots \cdots P(\theta_{\hat{q}} \in \widehat{CI}_{\hat{q}} \mid \hat{q} = q)$$

uniform for $P \in P^{\partial n}$ htd for all 3 methods with condition on $P^{\partial n}$

this is done on partial model.

Impossibility result for simultaneous & selective inference.

can't get uniform estimation because the inference is done on $\beta_0$ : full model parameter.