

1 Likelihood restriction

This section explicates main restrictions of the likelihood function. They are essential for the proposed algorithm of change point detection. Further assume that log-likelihood function $L(\theta) = L(Y, \theta)$, $Y = (Y_1, \dots, Y_n)$, has rather precise approximation by its quadratic part in local region $\Theta_0(r)$ of θ^* , $\Theta_0(r) \subseteq \mathbb{R}^p$, where

$$\theta^* = \operatorname{argmax}_{\theta} \mathbb{E}L(\theta), \quad \hat{\theta} = \operatorname{argmax}_{\theta} L(\theta)$$

and $\Theta_0(r) = \{\|D(\theta - \theta^*)\| < r\}$. [?] provides required conditions for justified quadratic approximation and parameter concentration in the local region. Approximation error involves the next variables for its estimation:

$$\alpha(\theta, \theta_0) = L(\theta) - L(\theta_0) - (\theta - \theta_0)^T \nabla L(\theta_0) + \frac{1}{2} \|D(\theta - \theta_0)\|^2,$$

$$\chi(\theta, \theta_0) = D^{-1} \nabla \alpha(\theta, \theta_0) = D^{-1} (\nabla L(\theta) - \nabla L(\theta_0)) + D(\theta - \theta_0).$$

Let in region $\Theta_0(r)$ with probability $1 - e^{-x}$:

$$\frac{|\alpha(\theta, \theta^*)|}{\|D(\theta - \theta^*)\|} \leq \diamond(r, x), \quad \|\chi(\theta, \theta^*)\| \leq \diamond(r, x), \quad (\text{A})$$

where $\diamond(r, x) = (\delta(r) + 6v_0 z_H(x)\omega)r$,

$$D^2(\theta) = -\nabla^2 \mathbb{E}L(\theta), \quad D = D(\theta^*), \quad (\text{D})$$

$$\|D^{-1} D^2(\theta) D^{-1} - I_p\| \leq \delta(r), \quad (\text{dD})$$

$$\forall \lambda \leq g, \gamma_1 \gamma_2 \in \mathbb{R}^p : \quad \log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \frac{\gamma_1^T \nabla^2 \zeta(\theta) \gamma_2}{\|D\gamma_2\| \|D\gamma_2\|} \right\} \leq \frac{\nu_0^2 \lambda^2}{2}, \quad (\text{ED2})$$

$$z_H(x) = \sqrt{H} + \sqrt{2x} + \frac{g^{-2}x + 1}{g} H, \quad H = 6p,$$

$$\omega = \omega(n) \sim \frac{1}{\sqrt{n}}, \quad \delta(r) \sim \frac{r}{\sqrt{n}}, \quad r \sim \sqrt{p}.$$

Condition (dD) ensures quadratic approximation of $\mathbb{E}L(\theta)$ and (ED2) ensures linear approximation of centered likelihood $\zeta(\theta) = L(\theta) - \mathbb{E}L(\theta)$.

2 LRT Statistic

Provide simple non-strict explanation of what kind of distribution the main statistic T_h is supposed to have. Review T_h as

$$T_h = L(\hat{\theta}) - L(\hat{\theta}_{H_0}), \quad L(\theta) = L_1(\theta_1) + L_2(\theta_2), \quad L_1 = L(Y[1 : h]), \quad L_2 = L(Y[h : 2h]),$$

where $\hat{\theta}_{H_0}$ is argmax of L under condition $H_0 : \theta_1 = \theta_2$. Then under quadratic approximation assumption T_h could be presented in Tailor equation with point $\hat{\theta}$:

$$T_h \approx \frac{1}{2} \|D(\hat{\theta} - \hat{\theta}_{H_0})\|^2.$$

If $\hat{\theta}$ and $\hat{\theta}_{H_0}$ tend to be Normal and H_0 is true then their difference are close to a centered Normal variable. If H_0 is false – the Normal variable will have mean that is equal to $\theta^* - \theta_{H_0}^*$.

The next two theorems includes more formal properties of T_h statistic.

Theorem 1. Assume condition (L*) and quadratic Laplace approximation (A) of L_1 and L_2 are fulfilled with probability $1 - 2e^{-x}$, additionally with probability $1 - 2e^{-x}$

$$\|\xi_i\| \leq z(x), \quad z^2(x) = \max_i p_{B_i} + 6\lambda_{B_i}x,$$

$$B_i = D_i^{-1} \text{Var}(\nabla L_i(\theta^*)) D_i^{-1}, \quad p_B = \text{tr}(B), \quad \lambda_B = \lambda_{\max}(B). \quad (\text{B})$$

Then in the local region with probability $1 - 8e^{-x}$

$$2T_h = \|\xi_{12} + \theta_{12}^*\|^2 + O(\{r + z(x)\} \diamond(r, x)),$$

where

$$\xi_{12} = \Sigma(D_2^{-1}\xi_2 - D_1^{-1}\xi_1), \quad \theta_{12}^* = \Sigma(\theta_2^* - \theta_1^*).$$

Remark. In increasing sample size $n \rightarrow \infty$ the stochastic component tends to Normal distribution:

$$\xi_{12} \rightarrow \mathcal{N}(0, B_1 + B_2).$$

Remark. For the condition $\hat{\theta} \in \Theta_1(r) \cap \Theta_2(r)$ the restriction of the parameter variability θ^* is required

$$\|D(\theta_1^* - \theta_2^*)\| \leq r. \quad (\text{L}^*)$$

Theorem 2. Assume condition (L*) and quadratic Laplace approximation (A) with probability $1 - 2e^{-x}$ are fulfilled. Then with probability $1 - 4e^{-x}$ in the local region $\Theta_1(r) \cap \Theta_2(r)$ took place

$$\left| \sqrt{2T_h} - \|\xi_{12} + \theta_{12}^*\| \right| \leq 10 \diamond(r, x).$$

where ξ_{12} and θ_{12}^* are defined in theorem 1.

3 Multiplier Bootstrap

Likelihood function in bootstrap case is a zipped sum with i.i.d weights (u_1, \dots, u_n) :

$$L^b(\theta) = \sum_{i=1}^n u_i l_i(\theta).$$

Each weight element has $\text{Var}_b u_i = 1$ and $\mathbb{E}_b u_i = 1$, that results in $(\{l_i(\theta)\})$ assumed to be independent)

$$\mathbb{E}_b L^b(\theta) = L(\theta), \quad \text{Var}_b \nabla L^b(\theta) = \sum_{i=1}^n \nabla l_i(\theta) \nabla l_i(\theta)^T \approx \text{Var} \nabla L(\theta) + \sum_{i=1}^n \mathbb{E} \nabla l_i(\theta) \mathbb{E} \nabla l_i(\theta)^T.$$

Variable $\zeta = \sum_i \zeta_i$ denotes stochastic part of the likelihood ($\zeta_i(\theta) = l_i(\theta) - \mathbb{E} l_i(\theta)$).

Lemma (Deviations of empirical process norm). Given exponential moment restriction with $\|\gamma_1\| = \|\gamma_2\| = 1$

$$\log \mathbb{E} \exp \left\{ \frac{\lambda}{\omega} \gamma_1 D^{-1} \nabla^2 \zeta_i(\theta) D^{-1} \gamma_2 \right\} \leq \frac{\lambda^2 \nu_i^2}{2}, \quad \sum_i \nu_i^2 = \nu_0^2, \quad (\text{EDi})$$

in local region $\Theta_0(r)$ of radius r with probability $1 - e^{-x}$ the next statement is fulfilled for all θ, θ_0

$$\|D_0^{-1}(\nabla \zeta(\theta) - \nabla \zeta(\theta_0))\| \leq 12\nu_0 z(x) \omega(n)r.$$

Paper [?] contains proof for this lemma.

Function $\alpha^b(\theta, \theta_0)$ determines quadratic approximation error for the weighted likelihood function.

$$\alpha^b(\theta, \theta_0) = L^b(\theta) - L^b(\theta_0) - (\theta - \theta_0)^T \nabla L^b(\theta_0) + \frac{1}{2} \|D(\theta - \theta_0)\|^2$$

Mean and deviation of the approximation error are

$$\|D^{-1} \nabla \mathbb{E}_b \alpha^b(\theta, \theta_0)\| = \|D^{-1} \nabla \alpha(\theta, \theta_0)\| \leq 2\Diamond(r, x).$$

$$S^b(\theta, \theta_0) = D^{-1} \{\nabla \alpha^b(\theta, \theta_0) - \mathbb{E}_b \nabla \alpha^b(\theta, \theta_0)\} = \sum_{i=1}^n D^{-1} \{\nabla l_i(\theta) - \nabla l_i(\theta_0)\} (u_i - 1). \quad (\text{Sb})$$

Centered function $\overset{o}{S}_b = S^b(\theta, \theta_0) - \mathbb{E}_Y S^b(\theta, \theta_0)$ with probability $1 - e^{-x}$ fulfills

$$\|\overset{o}{S}_b\| \leq 12\nu_0 z(x) \omega r \sqrt{\sum_{i=1}^n (\nu_i^2 / \nu_0^2) (u_i - 1)^2}.$$

Let u_i has restricted exponential moment:

$$\log \mathbb{E}_b e^{(u_i - 1)^2} \leq 1, \quad (\text{Eu})$$

then with probability $1 - e^{-t}$

$$\sqrt{\sum_{i=1}^n (\nu_i^2 / \nu_0^2) (u_i - 1)^2} \leq \sqrt{1 + t}$$

Consequently

$$\|\overset{o}{S}_b\| \leq 12\nu_0 z(x) \omega r \sqrt{1 + x}.$$

One could get restriction for $\mathbb{E} S(\theta_2, \theta_1)$ from Hoeffding inequality 8 restricting its components norms:

$$\|D^{-1} \nabla^2 \mathbb{E} l_i(\theta) D^{-1}\| \leq \frac{C_i(r)}{n}, \quad (\text{dDi})$$

$$\|\mathbb{E} S^b(\theta, \theta_0)\| \leq \frac{1 + \sqrt{2x}}{2n} C(r) r c_u,$$

where $C(r) = \sqrt{\sum_i C_i^2(r)}$, $c_u = \max(u_i - 1)$.

Theorem 3 (Boot, Wilks). Under conditions (EDi), (dDi), (A) and restrictions for weights (Eu), in local r -region it holds with probability $1 - 4e^{-x}$

$$|\alpha^b(\theta, \theta_0)| \leq \|D(\theta - \theta_0)\| \left(2\Diamond(r, x) + 12\nu_0 z(x) \omega r \sqrt{1 + x} + \frac{1 + \sqrt{2x}}{2n} C(r) r c_u \right) = \quad (\text{Ab})$$

$$= \|D(\theta - \theta_0)\| 2\Diamond^b(r, x).$$

Fisher theorem could be proved from condition

$$\chi^b(\theta, \theta_0) = D^{-1} (\nabla L^b(\theta) - \nabla L^b(\theta_0)) + D(\theta - \theta_0),$$

$$\chi^b(\theta, \theta_0) = D^{-1} \nabla \alpha^b(\theta, \theta_0).$$

With analog steps to 3 theorem proof one get an upper bound for $\chi^b(\theta, \theta_0)$.

Theorem 4 (Boot, Fisher). Under conditions (EDi), (dDi), (A) and restrictions for weights (Eu), in local r -region it holds with probability $1 - 4e^{-x}$

$$\|\chi^b(\hat{\theta}^b, \hat{\theta})\| = \|D(\hat{\theta}^b - \hat{\theta}) - D^{-1} \nabla L^b(\hat{\theta})\| \leq 2\Diamond^b(r, x),$$

where $\hat{\theta}^b, \hat{\theta}$ – MLE of weighted and non-weighted likelihood functions.

4 Uniform bands

Quadratic approximation of $\alpha(\theta_0, \theta)$ is correct with high probability for each window position t , but it takes place uniformly for all the windows with higher probability than independent events because of correlations. Corresponding probability of the uniform approximation could be obtained by grouping window positions into high correlation regions. After split the whole data range of size T into such regions of size r_t the uniform approximations are to be true for all the regions with probability $1 - e^{-x_t}$, $x_t = x - \log(T/r_t)$. The first grouping approach uses bootstrap related equations.

$$L(\theta, t) = \sum_{i=1}^n l_i(\theta) u_i(t), \quad 1 \leq t \leq n = h + 2r_t,$$

where $u_i(t)$ denotes kernel function that selects elements around coordinate t with distance less than $h/2$:

$$u_i(t) = \mathbb{I}\{|t - i| \leq h/2\}, \quad \mathbb{E}_t u_i(t) = \frac{h}{n}.$$

Weights $u_i(t)$ have probabilistic nature since they are dependently generated by sampling position $t \in [1, n]$ uniformly. The expectation comes from all window positions of circle closed data ($Y_{n+i} = Y_i$).

$$\alpha(t, \theta, \theta_0) = L(\theta, t) - L(\theta_0, t) - (\theta - \theta_0) \nabla L(\theta_0, t) + \frac{1}{2} \|D_h(\theta - \theta_0)\|^2$$

Precise quadratic approximation requires small value of $\nabla \alpha(t, \theta, \theta_0)$ for all From (A) and weights mean follows

$$\|\mathbb{E}_t D_h^{-1} \nabla \alpha(t)\| \leq \sqrt{\frac{h}{n}} \cdot 2 \diamond(r(\theta^*), x, n) \sim \sqrt{\frac{h}{n}} \sqrt{\frac{p}{n}}.$$

Defining variable similar to (Sb) one obtains its mean and deviation bounds in separate steps.

$$S(t) = D_h^{-1} \{\nabla \alpha(t) - \mathbb{E}_t \nabla \alpha(t)\}, \quad \|S(t) - \mathbb{E}_Y S(t)\| \leq \sqrt{1 - \frac{h}{n}} \cdot 12\nu_0 z(x, p) \omega(h) \max_i r(\theta_i^*).$$

Here $r(\theta_i^*)$ denotes max distance from $\arg\max \mathbb{E} l_i(\theta)$ to (θ, θ_0) . From asymptotic assumptions for $\omega(h) \sim 1/\sqrt{h}$ and $z(x) \sim \sqrt{p}$ one get

$$\|S(t) - \mathbb{E}_Y S(t)\| \sim \sqrt{\frac{p}{h}} \bar{r}.$$

The mean part requires additional condition (dDi) for divergence of the likelihood components

$$\|\mathbb{E}_Y S(t)\| \leq \max_i \frac{C(r_i)}{h} r_i \sum_{j=1}^n |u_j - \mathbb{E}_t u_j| = 2 \max_i C(r_i) r_i \left(1 - \frac{h}{n}\right),$$

where

$$\sum_{j=1}^n |u_j - \mathbb{E}_t u_j| = (n - h) \left|0 - \frac{h}{n}\right| + h \left|1 - \frac{h}{n}\right| = 2h \left(1 - \frac{h}{n}\right).$$

Finally, the uniform estimation error for region n is

$$12\nu_0 z(x, p) \omega(h) \max_i r(\theta_i^*) + 2 \max_i C(r_i) r_i \left(1 - \frac{h}{n}\right).$$

From which with $C(r_i) \sim r_i$ one get $1/\sqrt{h} \sim (1 - h/n)$, consequently $n - h \sim \sqrt{h}$.

The second grouping approach treats window position t as a part of parameter θ . Define variable y as difference of stochastic likelihood part gradients:

$$y(\theta, t) = D_h^{-1} (\nabla \zeta(\theta, t) - \nabla \zeta(\theta_0, t_0)).$$

Define new parameters u, τ aiming for restriction $\|u\| \leq r_0, \|\tau\| \leq r_0$:

$$u = D_h(\theta - \theta_0), \quad \tau = \frac{r_0}{r_t}(t - t_0), \quad |t - t_0| \leq r_t.$$

Estimation for r_t comes from equating parts of y gradients from each parameter.

$$\nabla_u y = D_h^{-1} \nabla^2 \zeta(\theta, t) D_h^{-1}, \quad \nabla_\tau y = \frac{r_t}{r_0} D_h^{-1} \nabla_t \nabla \zeta(\theta, t)$$

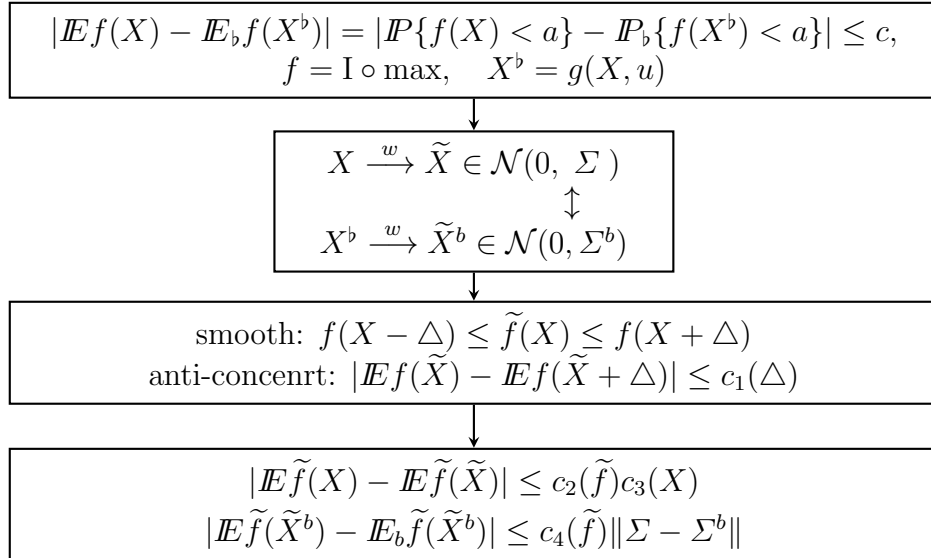
$$\|D_h^{-1} \nabla_t \nabla \zeta(\theta, t)\| \leq 2\|D_h^{-1} \nabla l_i(\theta_i^*, t)\| + \|D_h^{-1} \nabla^2 \zeta_i(\cdot, t) D_h^{-1}\| r_0 \sim \sqrt{\frac{p}{h}}.$$

$$\|D_h^{-1} \nabla^2 \zeta(\theta, t) D_h^{-1}\| \sim \omega(h) z(x) \sim \sqrt{\frac{p}{h}}.$$

Finally, this approach results in $n - h = r_t \sim r_0$.

5 LRT Critical values calibration

Scheme of \mathbb{P} measure and variable $f(X)$ (f is non random) approximation by corresponded bootstrap (empirical) measure \mathbb{P}_b and variable $f(X^b)$ could be done in three steps: approximate X, X^b by Normal distributions, approximate $f(\cdot)$ by smooth function, present measures distance as product of f properties and its arguments distribution parameters distance.



Weighted LRT statistic T_h^b from multiplier bootstrap has analog of theorem 2 due to 4 and 3.

Theorem 5. Let (??) and (??) are true. Then with probability $1 - 16e^{-x}$ in local region $\Theta_1(r) \cap \Theta_2(r)$

$$\left| \sqrt{2T_h^b(\theta_1^b, \theta_2^b)} - \|\xi_{12}^b + \hat{\theta}_{12}\| \right| \leq 10\Diamond^b(r, x),$$

where

$$\xi_{12}^b = \Sigma(D_2^{-1}\xi_2^b - D_1^{-1}\xi_1^b), \quad \hat{\theta}_{12} = \Sigma(\hat{\theta}_2 - \hat{\theta}_1), \quad \theta^b = \operatorname{argmax}_{\theta} L^b(\theta).$$

Remark. Generator of ξ_{12}^b use statistic $T_h^b = \sqrt{2T_h(\theta_1^b, \theta_1^b + \widehat{\theta}_{12})}$, for which 5 theorem is also applicable with $\widehat{\theta}_{12} = 0$,

$$|T_h^b - \|\xi_{12}^b\|| \leq 10\Diamond^b(r, x).$$

Precision of distribution estimation $\xi_{12} + m$ using bootstrap (it's analog $\xi_{12}^b + m^b$, m , m^b non random) could be obtained by normal approximation for both variables. The sequence of comparison is

1. $\xi_{12}^b(\widehat{\theta}_1, \widehat{\theta}_2) \approx \xi_{12}^b(\theta_1^*, \theta_2^*)$;
2. $\|\xi_{12}^b(\theta_1^*, \theta_2^*) + m^b\| \approx \|\widetilde{\xi}^b\|$, $\widetilde{\xi}^b \sim \mathcal{N}(m^b, \Sigma^b)$;
3. $\|\xi_{12}(\theta_1^*, \theta_2^*) + m\| \approx \|\widetilde{\xi}\|$, $\widetilde{\xi} \sim \mathcal{N}(m, \Sigma)$;
4. $\|\widetilde{\xi}\| \approx \|\widetilde{\xi}^b\|$.

1) Compare $\xi_{12}^b(\widehat{\theta}_1, \widehat{\theta}_2)$ and $\xi_{12}^b(\theta_1^*, \theta_2^*)$. Due to

$$\xi_i^b(\theta) = D^{-1} \sum_{i=1}^n \nabla l_i(\theta)(u_i - 1),$$

and estimation for (Sb), the next inequality holds with probability $1 - 3e^x$, $i = \{1, 2\}$

$$\|\xi_i^b(\theta_2) - \xi_i^b(\theta_1)\| \leq \left(12\nu_0 z(x) \omega r \sqrt{1+x} + \frac{1 + \sqrt{2x}}{2n} C(r) r c_u \right) = \Diamond_\xi^b(r, x).$$

By definition

$$\xi_{12}^b = A \begin{pmatrix} \xi_1^b \\ \xi_2^b \end{pmatrix}, \quad A = \Sigma \begin{pmatrix} -D_1^{-1} & D_2^{-1} \end{pmatrix},$$

which leads to statement

Lemma.

$$\|\xi_{12}^b(\widehat{\theta}_1, \widehat{\theta}_2) - \xi_{12}^b(\theta_1^*, \theta_2^*)\| \leq 2\Diamond_{\xi, F}^b(r, x),$$

where $F = A^T A$, $H_2 = H_2(F) + 4p$,

$$\Diamond_{\xi, F}^b(r, x) = \left(12\nu_0 z_F(x) \omega r \sqrt{1+x} + \frac{1 + \sqrt{2x}}{2n} C(r) r c_u \right),$$

$$z_F(x) = \dots$$

2),3) This step is weak approximation for $\|\xi_{12}(\theta_1^*, \theta_2^*) + m\| \approx \|\widetilde{\xi}\|$, $\widetilde{\xi} \sim \mathcal{N}(m, \Sigma)$. Namely, one need to compare functions of these norms (pointwise argument comparison for each $\|\xi_{12}(\theta_1^*, \theta_2^*) + m\|(t)$ has worse precision for its aggregations):

$$\max_{\tau} \sum_{t \in \tau} P(t) \|\xi_{12}(\theta_1^*, \theta_2^*) + m\|(t). \quad (1)$$

For non-limiting simplification assume that $m^b = m = 0$. Norm of a vector ξ could be presented as $\max_{\|\gamma\|=1} \gamma^T \xi$, and max as smooth max function:

$$h_\beta(x) = \beta^{-1} \log \left(\sum_i e^{\beta x_i} \right), \quad x = \{x_i\}.$$

After smoothing the considered expression (1) is

$$h_\beta \left(\sum_{t \in \tau} P(t) \gamma_t^T \xi_{12}(t) \right) = h_\beta \left(\sum_{t \in \tau} P(t) \gamma_t^T \sum_{i=1}^n w_i(t) \varepsilon_i \right) = h_\beta \left(\sum_{i=1}^n c(\tau)^T \varepsilon_i \right),$$

where $w_i(t)$ are kernel weights for window position t , $P(t)$ pattern values, ε_i – components of ξ_{12} . Function h_β corresponds to max over iterable arguments (τ, γ_τ) .

Following property $h_\beta(x)$, $x \in \mathbb{R}^M$ characterise the error of smooth max approximation

$$\max_i(x_i) \leq h_\beta(x) \leq \max_i(x_i) + \frac{\log(M)}{\beta}.$$

Dimension includes $\gamma_t \in G_\varepsilon$, $t \in \{1, \dots, n\}$ and τ , which results in $M < n|G_\varepsilon|$. The next lemma apply indicator for both parts of the above inequality.

Lemma. For function g_δ which is smooth indicator (g_δ grows from 0 to 1 inside interval of size δ) took place following inequality with $\delta = \beta^{-1} \log(M)$

$$g_\delta h_\beta(x - \delta) \leq \mathbb{I} \left[\max_{0 \leq i \leq M} x_i > 0 \right] \leq g_\delta h_\beta(x + \delta).$$

Consequently, under condition $|\mathbb{E} g_\delta h_\beta(x) - \mathbb{E} g_\delta h_\beta(\tilde{x})| \leq C(\beta, M) \mu_n$,

$$\left| \mathbb{P} \left(\max_{0 \leq i \leq M} x_i > 0 \right) - \mathbb{P} \left(\max_{0 \leq i \leq M} \tilde{x}_i > \pm 2\delta \right) \right| \leq C(\delta, M) \mu_n. \quad (2)$$

The shift of length $\pm 2\delta$ could be moved outside applying anti-concentration lemma for distribution density $\max_i \tilde{x}_i$.

Lemma. Let $x \in \mathcal{N}(m, \Sigma) \in \mathbb{R}^M$, $\sigma_1 \leq \sqrt{\Sigma_{ii}} \leq \sigma_2$, $a_M = \max_i(\tilde{x}_i - m_i)/\sqrt{\Sigma_{ii}}$, then $\forall c$

$$\mathbb{P}(|\max_i \tilde{x}_i - c| \leq \varepsilon) \leq C_{\text{ak}}(M, \Sigma),$$

$$C_{\text{ak}}(M, \Sigma) = \frac{4\varepsilon}{\sigma_1} \left(\frac{\sigma_2}{\sigma_1} a_M + \left(\frac{\sigma_2}{\sigma_1} - 1 \right) \sqrt{2 \log \left(\frac{\sigma_1}{\varepsilon} \right)} + 2 - \frac{\sigma_1}{\sigma_2} \right) \approx 4\varepsilon \frac{\sigma_2}{\sigma_1^2} \sqrt{2 \log \left(\frac{\sigma_1 M}{\varepsilon} \right)}.$$

Combining this lemma and 2 one get result for measure difference of max of vector x and normal vector \tilde{x} .

$$\left| \mathbb{P} \left(\max_{0 \leq i \leq M} x_i > 0 \right) - \mathbb{P} \left(\max_{0 \leq i \leq M} \tilde{x}_i > 0 \right) \right| \leq C(\delta, M) \mu_n + \delta C_{\text{ak}}(M, \Sigma). \quad (3)$$

The next required step is involved upper bound of obtaining

$$|\mathbb{E} g_\delta h_\beta(x) - \mathbb{E} g_\delta h_\beta(\tilde{x})| \leq C(\beta, M) \mu_n.$$

Lemma. Function $f = g_\delta h_\beta$, where g_δ has bounded third order derivations, $\delta = \beta^{-1} \log(M)$, has representation in Tailor form

$$|f(x + d) - f(x) - d^T f'(x) - d^T f''(x) d/2| \leq C(\delta, M) \|d\|_\infty^3,$$

where

$$C(\delta, M) = \frac{1}{6\delta^3} (|g'''| + \log(M)|g''| + \log^2(M)|g'|).$$

Basing on this lemma use that

$$x_{\tau, \gamma} = \sum_{i=1}^n c_i(\tau)^T \varepsilon_i,$$

where $\text{Var} \sum_i c_i(\tau)^T \varepsilon_i = \sum_i c_i(\tau)^T V_i^2 c_i(\tau)$.

The final statement for approximation of ε_i by normal variables $\tilde{\varepsilon}$ with the same mean and variance is

$$\left| \mathbb{E} g_\delta h_\beta \left(\sum_{i=1}^n c(\tau)^T \varepsilon_i \right) - \mathbb{E} g_\delta h_\beta \left(\sum_{i=1}^n c(\tau)^T \tilde{\varepsilon}_i \right) \right| \leq C(\delta, M) \mu_n, \quad (4)$$

where

$$\mu_n = \max_{\tau} \sum_{i=1}^n (\|\varepsilon_i\|_\infty^3 + C \log^{3/2}(M) \sigma_i^3) \|c_i(\tau)\|^3,$$

$$\|c_i(\tau)\| \leq \sum_{t \in \tau} P(t) w_i(t) \|\gamma_t^T\| \sim \mathbb{I}(\tau - 2h \leq i \leq \tau + 2h), \quad \mu_n \sim \frac{1}{\sqrt{h}}.$$

2'), 3') Consider the other type of smoothing for (1):

$$h_\beta \left(\sum_{i \in \tau} P(i) \|\xi_{12}\| (i) \right).$$

One have to estimate derivative

$$(g_\delta h_\beta \{\xi_{12} + td\})_t''' = g_\delta' h_\beta''' + 3g_\delta'' h_\beta' h_\beta'' + g_\delta''' (h_\beta')^3.$$

Assume that approximately

$$\|\nabla_P h_\beta(P_1, \dots, P_n)\|_1 \leq 1, \quad \|\nabla_P^2 h_\beta(P_1, \dots, P_n)\|_1 \leq \beta, \quad \|\nabla_P^3 h_\beta(P_1, \dots, P_n)\|_1 \leq \beta^2,$$

$$|g_\delta'|_\infty \leq \frac{1}{\delta}, \quad |g_\delta''|_\infty \leq \frac{1}{\delta^2}, \quad |g_\delta'''|_\infty \leq \frac{1}{\delta^3}.$$

Also the estimation for $\max_\tau P'_\tau(t)$ is required, where $P_\tau(t) = \sum_{i \in \tau} P(i) \|\xi_{12} + td\| (i)$:

$$\max_\tau P'_\tau(t) \leq \|d\| \max_\tau \sum_{i \in \tau} P(i).$$

Consider pattern weights with sum equal to 1, which leads to

$$(g_\delta h_\beta \{\xi_{12} + td\})_t''' \leq \left(\frac{\beta^2}{\delta} + 3 \frac{\beta}{\delta^2} + \frac{1}{\delta^3} \right) \|d\|^3.$$

4) Compare the same functions $\mathbb{E} g_\delta h_\beta$ from two normal vectors $X = \tilde{\xi}$ and $Y = \tilde{\xi}^b$ with different mean and variance.

Lemma. For X, Y independent normal vectors with $m_X, m_Y, \Sigma_X, \Sigma_Y$,

$$\Delta m = m_2 - m_1, \quad \Delta \Sigma = \Sigma_2 - \Sigma_1.$$

it holds

$$|\mathbb{E} g(\delta^{-1} h(X)) - \mathbb{E} g(\delta^{-1} h(Y))| \leq \left(\frac{\beta \|g'\|_\infty}{\delta} + \frac{\|g''\|_\infty}{2\delta^2} \right) \|\Delta \Sigma\|_{\text{op}} + \frac{\|g'\|_\infty}{2\delta} \|\Delta m\|_{\text{op}}.$$

Following Section deals with estimation of $\|\Delta \Sigma\|_{\text{op}}$.

6 Matrix deviations

6.1. Matrix inequality in change point statistic

A sequence of noise variables produced by different window positions.

$$\xi(t) = D_h^{-1}(\nabla L(\theta^*) - \nabla \mathbb{E}L(\theta^*)) = D_h^{-1} \sum_{i=t}^{t+h} \nabla l_i(\theta^*),$$

$$\xi^b(t) = D_h^{-1}(\nabla L(\hat{\theta}) - \nabla \mathbb{E}_b L(\hat{\theta})) = D_h^{-1} \sum_{i=t}^{t+h} \nabla l_i(\hat{\theta})(u_i - 1).$$

Unite all variables involved into LRT statistic.

$$X = \xi_{12}(t)_{t=1}^n, \quad X^b = \xi_{12}^b(t)_{t=1}^n, \quad \xi_{12}(t) = \frac{1}{\sqrt{2}}(\xi(t+1) - \xi(t))$$

Consider difference in variance between X and X^b in case of regression model:

$$Y = \Psi_{[n \times p]}^T \theta + \varepsilon, \quad \varepsilon \in \mathcal{N}(0, \Sigma), \quad \Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$$

Variables X and X^b could be presented as linear transformation of a standart Normal vectors ε and ε^b correspondingly

$$\nabla l(Y_i) = \Psi_i \varepsilon_i$$

$$\nabla l^b(Y_i)(u_i - 1) = \Psi_i(Y - \Pi Y)_i(u_i - 1) \Psi_i \varepsilon_i^b, \quad \Pi = \Psi^T(\Psi \Psi^T)^{-1} \Psi$$

$$X = A H D_h^{-1} D(\Psi) \varepsilon = V \varepsilon, \quad X^b = A H D_h^{-1} D(\Psi) \varepsilon^b = V \varepsilon^b.$$

$$D(\Psi) = \text{diag}(\Psi_1, \dots, \Psi_n),$$

$$H_{np, np} = \begin{pmatrix} I_p & I_p & \cdots & I_p & 0 & 0 \cdots & 0 \\ 0 & I_p & \cdots & I_p & I_p & 0 \cdots & 0 \\ 0 & 0 & \cdots & I_p & I_p & I_p \cdots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots & & \\ 0 & 0 & \dots\dots\dots & 0 & 0 & & \end{pmatrix} \quad A_{np, np} = \begin{pmatrix} I_p & -I_p & \cdots & 0 & \\ 0 & I_p & -I_p & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \\ 0 \cdots & 0 & I_p & -I_p & \end{pmatrix}$$

We are to compare S and S^b by

$$\|S\| \|S^{-1/2}(S^b - S)S^{-1/2}\|$$

$$S = \text{Var}(X) = V \Sigma V^T, \quad S^b = \text{Var}_b(X^b) = V \text{diag}((Y - \Pi Y)(Y - \Pi Y)^T) V^T$$

Involving variable \mathcal{U} the goal estimation $(S^b - S)$ becomes

$$\mathcal{U} = S^{-1/2} V \Sigma^{1/2}, \quad \mathcal{U} \mathcal{U}^T = I_{pn}$$

$$S^{-1/2}(S^b - S)S^{-1/2} = \mathcal{U} \Sigma^{-1/2} \text{diag}((Y - \Pi Y)(Y - \Pi Y)^T) \Sigma^{-1/2} \mathcal{U}^T - I_{np}$$

Variable $Y - \Pi Y$ has decomposition to bias and noise part

$$\Sigma^{-1/2}(Y - \Pi Y) = \Sigma^{-1/2}(f - \Pi f) + \Sigma^{-1/2}(\varepsilon - \Pi \varepsilon) = B + \tilde{\varepsilon}$$

$$\begin{aligned}
S^{-1/2}(S^\flat - S)S^{-1/2} &= \\
&= \mathcal{U} \text{diag}(BB^\top)\mathcal{U}^\top + 2\mathcal{U} \text{diag}(\tilde{\varepsilon}B^\top)\mathcal{U}^\top + \\
&\quad + \mathcal{U} \text{diag}(\tilde{\varepsilon}\tilde{\varepsilon}^\top - \mathbb{E}\tilde{\varepsilon}\tilde{\varepsilon}^\top)\mathcal{U}^\top + \mathcal{U} \text{diag}(\mathbb{E}\tilde{\varepsilon}\tilde{\varepsilon}^\top)\mathcal{U}^\top - I_{np}.
\end{aligned}$$

The last component in this equation is a part of the non-random component, which becomes high in very large model, while B is small in large model.

Let $\text{Var}(\tilde{\varepsilon}) = R$ and $\|(S^{-1/2}V)_i\| \leq \delta$, then from the consequence for theorem 13 with probability $1 - e^x$

$$\|\mathcal{U} \text{diag}(\tilde{\varepsilon}\tilde{\varepsilon}^\top - \mathbb{E}\tilde{\varepsilon}\tilde{\varepsilon}^\top)\mathcal{U}^\top\|_{\text{op}} \leq \lambda_{\max}(R)(2\delta\sqrt{\mathbf{x}_{np}} + 2\delta^2\mathbf{x}_{np}),$$

where $\mathbf{x}_{np} = x + \log(np)$,

$$\|\mathcal{U} \text{diag}(\tilde{\varepsilon}B^\top)\mathcal{U}^\top\|_{\text{op}} \leq \sqrt{\lambda_{\max}(R)}\delta^2\|B\|\sqrt{2\mathbf{x}}.$$

The other two components are restricted by means of

Lemma.

$$\|\mathcal{U}A\mathcal{U}^\top\|_{\text{op}} \leq \min(\|A\|_{\text{op}}, \delta^2 \text{tr}\{A\})$$

Proof.

$$\|\mathcal{U}A\mathcal{U}^\top\|_{\text{op}} = \sup_{\|\gamma\|=1} \gamma^\top \mathcal{U}A\mathcal{U}^\top \gamma \leq \sup_{\|\gamma\|=1, \max \gamma \leq \delta} \gamma^\top A \gamma \leq \min(\|A\|_{\text{op}}, \delta^2 \text{tr}\{A\})$$

□

Applying the lemma one get

$$\|\mathcal{U} \text{diag}(BB^\top)\mathcal{U}^\top\|_{\text{op}} \leq \min(\|B\|_{\text{op}}^2, \delta^2 \|B\|^2)$$

and

$$\|\mathcal{U} \text{diag}(\mathbb{E}\tilde{\varepsilon}\tilde{\varepsilon}^\top)\mathcal{U}^\top - I_{np}\|_{\text{op}} = \|\mathcal{U}(I_n - \Pi)^2 - I_n\mathcal{U}^\top\|_{\text{op}} = \|\mathcal{U}\Pi\mathcal{U}^\top\|_{\text{op}} \leq 1.$$

Finally the upper bound for variance difference is

$$\lambda_{\max}(R)(2\delta\sqrt{\mathbf{x}_{np}} + 2\delta^2\mathbf{x}_{np}) + \sqrt{\lambda_{\max}(R)}\delta^2\|B\|\sqrt{2\mathbf{x}} + \min(\|B\|_{\text{op}}^2, \delta^2 \|B\|^2) + 1.$$

The last step is bound estimation for $\|S\|_{\text{op}}$ which is

$$\|A\|_{\text{op}}^2 \|H\|_{\text{op}}^2 \|D_h^{-1}\|_{\text{op}}^2 \|D(\Psi)D(\Psi)^\top\|_{\text{op}}$$

$$\|D(\Psi)D(\Psi)^\top\|_{\text{op}} \leq \max_i \|\Psi_i\Psi_i^\top\|_{\text{op}}$$

Lemma. Spectrum of an upper triangular matrix A is its diagonal.

Proof. $\text{Res}(A)$ set is equal to $\text{Res}(\text{diag}(A))$, because of $\text{rank}(A - \lambda I) = \text{rank}(\text{diag}(A) - \lambda I)$. □

Correspondingly, $\|A\|_{\text{op}}^2 = \|H\|_{\text{op}}^2 = 1$. Bound for $\|D_h^{-1}\|_{\text{op}}^2$ is $O(h^{-1})$.

6.2. Generalization to i.i.d model

Retreat regression model with an i.i.d model. In this case comparing variables would be

$$X = AHD_h^{-1} \begin{pmatrix} \nabla l_1(\theta_1^*) \\ \vdots \\ \nabla l_n(\theta_n^*) \end{pmatrix} = AHD_h^{-1} \nabla, \quad X = AHD_h^{-1} \begin{pmatrix} \nabla l_1(\hat{\theta}_1)(u_1 - 1) \\ \vdots \\ \nabla l_n(\hat{\theta}_n)(u_n - 1) \end{pmatrix} = AHD_h^{-1} \nabla^b.$$

Here instead of noise variance one get block-diagonal one

$$\Sigma = \begin{pmatrix} \text{Var } \nabla l_1(\theta_1^*) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \text{Var } \nabla l_n(\theta_n^*) \end{pmatrix}$$

Analogically divide $\hat{\nabla} = \nabla(\hat{\theta})$ into mean and stochastic parts

$$\Sigma^{-1/2} \hat{\nabla} = \Sigma^{-1/2} \mathbb{E} \hat{\nabla} + \Sigma^{-1/2} (\hat{\nabla} - \mathbb{E} \hat{\nabla}) = B + \zeta.$$

Normalized variance difference term is

$$\begin{aligned} S^{-1/2} (S^b - S) S^{-1/2} &= \\ &= \mathcal{U} \text{blockDiag}(BB^\top) \mathcal{U}^\top + 2\mathcal{U} \text{blockDiag}(\zeta B^\top) \mathcal{U}^\top + \\ &+ \mathcal{U} \text{blockDiag}(\zeta \zeta^\top - \mathbb{E} \zeta \zeta^\top) \mathcal{U}^\top + \mathcal{U} \text{blockDiag}(\mathbb{E} \zeta \zeta^\top) \mathcal{U}^\top - I_{np}. \end{aligned}$$

Where blockDiag determines block diagonal matrix with block shape $[p \times p]$, corresponded to each component in sequence $\nabla l_1, \dots, \nabla l_n$. Start consideration with term $\mathcal{U} \text{blockDiag}(\zeta \zeta^\top - \mathbb{E} \zeta \zeta^\top) \mathcal{U}^\top$.

Lemma. Let $\{U_i\}$ be submatrices of matrix \mathcal{U} corresponded to each block, such that $\mathcal{U} = (U_1, \dots, U_n)$. Assume that

$$\|U_i^\top U_i\|_{\text{op}} \leq \delta^2.$$

Than

$$\|(U_i^\top U_i)^{-1} \log \mathbb{E} \exp\{(U_i^\top U_i)^{1/2} (\zeta_i \zeta_i^\top - \mathbb{E} \zeta_i \zeta_i^\top) (U_i^\top U_i)^{1/2}\}\|_{\text{op}} \leq \frac{\nu_i^2}{2} \delta^2.$$

Proof. Denote $\zeta_i \zeta_i^\top - \mathbb{E} \zeta_i \zeta_i^\top$ as A and $(U_i^\top U_i)^{1/2}$ as U . With condition $\mathbb{E} A = 0$:

$$\|\log \mathbb{E} e^{UAU}\|_{\text{op}} = \left\| \left(\frac{1}{2} U \mathbb{E}[AU^2 A] U + O((UA)^3) \right) \right\|_{\text{op}}.$$

With condition $\|UA\|_{\text{op}} < 1/2$ (with high probability):

$$\|U^{-2} \log \mathbb{E} e^{UAU}\|_{\text{op}} \leq \|U^2\|_{\text{op}} \|\mathbb{E} A^2\|_{\text{op}}. \quad (5)$$

$$\|\mathbb{E}(\zeta_i \zeta_i^\top - \mathbb{E} \zeta_i \zeta_i^\top)^2\|_{\text{op}} = \|\mathbb{E} \|\zeta_i\|^2 \zeta_i \zeta_i^\top - (\mathbb{E} \zeta_i \zeta_i^\top)^2\|_{\text{op}},$$

where

$$\|\zeta_i\|^2 \sim p, \quad \|\mathbb{E} \zeta_i \zeta_i^\top\|_{\text{op}} = \left\| \Sigma_i^{-1/2} \Sigma_i(\hat{\theta}_i) \Sigma_i^{-1/2} - I_p + I_p \right\|_{\text{op}} \leq 1 + \delta_\Sigma(r).$$

Finally, $\nu_i^2 \sim 2p(1 + \delta_\Sigma(r))$. □

From proof for theorem 13

$$\mathbb{P} \left(\|\mathcal{U} \text{diag}(\zeta \zeta^\top - \mathbb{E} \zeta \zeta^\top) \mathcal{U}^\top\|_{\text{op}} > t \right) \leq 2np \inf_{\theta} \exp \left\{ -\theta t + \frac{\delta^2 \theta^2 \nu_i^2}{2} \right\},$$

Subsequently, with $t = \sqrt{2x_{np}}\delta\nu_l$, $x_{np} = x + \log(2np)$:

$$\mathbb{P} \left(\|\mathcal{U} \text{diag}(\zeta\zeta^T - \mathbb{E}\zeta\zeta^T)\mathcal{U}^\top\|_{\text{op}} > t \right) \leq e^x.$$

For the other term $\mathcal{U} \text{blockDiag}(\zeta B^\top)\mathcal{U}^\top$ by means of 5 one get with $\|\zeta_i\|^2 \sim p$

$$\|U^{-2} \log \mathbb{E} \exp\{U\zeta_i B_i^T U\}\|_{\text{op}} \leq p \|B_i\|^2 \delta^2.$$

Correspondingly, from proof for theorem 13

$$\mathbb{P} \left(\|\mathcal{U} \text{diag}(\zeta B^T)\mathcal{U}^\top\|_{\text{op}} > t \right) \leq 2np \inf_{\theta} \exp \left\{ -\theta t + p \max_i \|B_i\|^2 \delta^2 \right\},$$

Finally, with $t = 2\sqrt{x_{np}p}\delta \max_i \|B_i\|$, $x_{np} = x + \log(2np)$:

$$\mathbb{P} \left(\|\mathcal{U} \text{diag}(\zeta B^T)\mathcal{U}^\top\|_{\text{op}} > t \right) \leq e^x.$$

7 Max-norm Gaussian

Comparison of $\max_t \langle P, (\|X_1\|, \dots, \|X_n\|) \rangle(t)$ and $\max_t \langle P, (\|X_1^b\|, \dots, \|X_n^b\|) \rangle(t)$ could be done by Slepian inequality and replacing norm with max-product function ($\max_t P \max_{\gamma_1 \dots \gamma_n} \Gamma X$), where

$$\Gamma = \text{diag}(\gamma_1^T \dots \gamma_n^T), \quad P = \text{diag}(P_h^T \dots P_h^T)$$

Lemma.

$$\max_{\gamma \in G_\varepsilon} \gamma^T S \leq \|S\| \leq \frac{1}{1-\varepsilon} \max_{\gamma \in G_\varepsilon} \gamma^T S,$$

where $|G_\varepsilon| \leq (2/\varepsilon)^{p-1}$, $S(t) \in \mathbb{R}^p$, $1 \leq t \leq M$ By anti-concentration lemma one could obtain inequality for measures

$$\left| \mathbb{P}(\max_t \|S(t)\| < q) - \mathbb{P}(\max_t \max_{\gamma \in G_\varepsilon} \gamma^T S(t) < q) \right| \leq 8\sqrt{p}\varepsilon q \frac{\sigma_2}{\sigma_1^2} \sqrt{\log \left(M \frac{2}{\varepsilon} \right)}$$

Since \max_{γ_i} functions are independent the next theorem take place with $|Gh_\varepsilon| \leq h(2/\varepsilon)^{p-1}$

Theorem 6.

$$\left| \mathbb{P}(\max_t \langle P, (\|X_1\|, \dots, \|X_n\|) \rangle(t) < q) - \mathbb{P}(\max_t \max_{\gamma_1 \dots \gamma_n} P \Gamma X < q) \right| \leq 8\sqrt{p}\varepsilon q \frac{\sigma_2}{\sigma_1^2} \sqrt{\log \left(\frac{2nh}{\varepsilon} \right)}$$

Theorem 7.

$$\begin{aligned} & \mathbb{P}(\max_t \max_{\gamma_1 \dots \gamma_n} P \Gamma X < q) - \mathbb{P}(\max_t \max_{\gamma_1 \dots \gamma_n} P \Gamma X^b < q) \leq \\ & \leq \left(\frac{\beta \|g'\|_\infty}{\delta} + \frac{\|g''\|_\infty}{2\delta^2} \right) \|\Delta P \Gamma \Sigma \Gamma^T P^T\|_{\text{op}} + 8\sqrt{p}\delta \frac{\sigma_2}{\sigma_1^2} \sqrt{\log \left(\frac{2nh}{\varepsilon} \right)} \end{aligned}$$

8 Appendix