

Consistent kernel change-point detection

Sylvain Arlot¹ (joint works with Alain Celisse², Damien Garreau³ & Zaïd Harchaoui⁴)

¹Université Paris-Sud

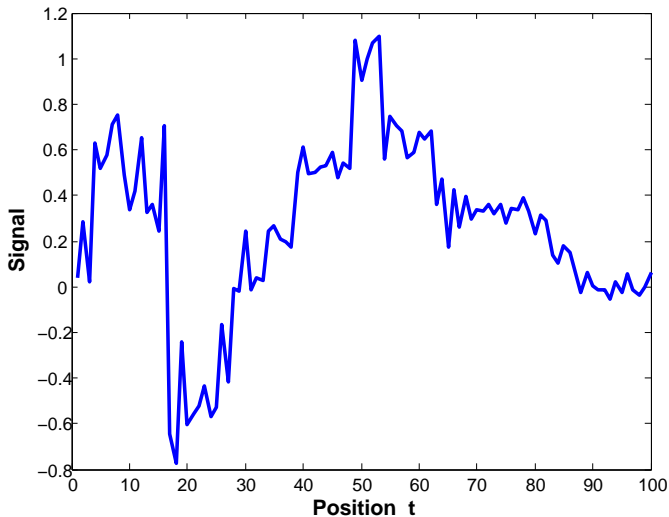
²Université Lille 1

³Max Planck Institute for Intelligent Systems, Tübingen

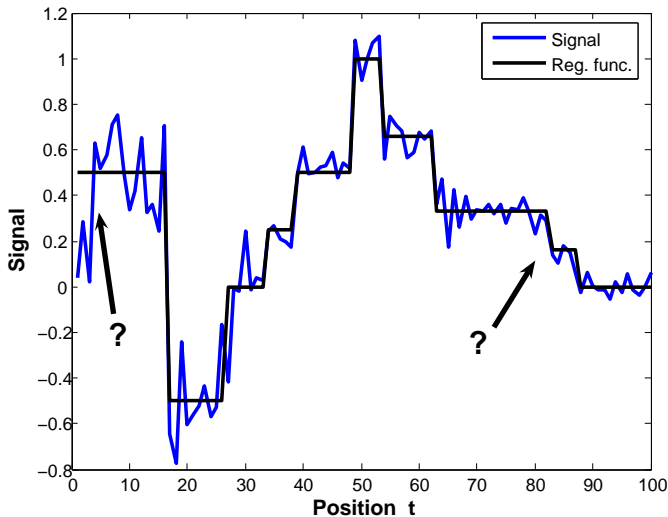
⁴University of Washington

Séminaire MODAL'X, Nanterre, 9 Novembre 2017

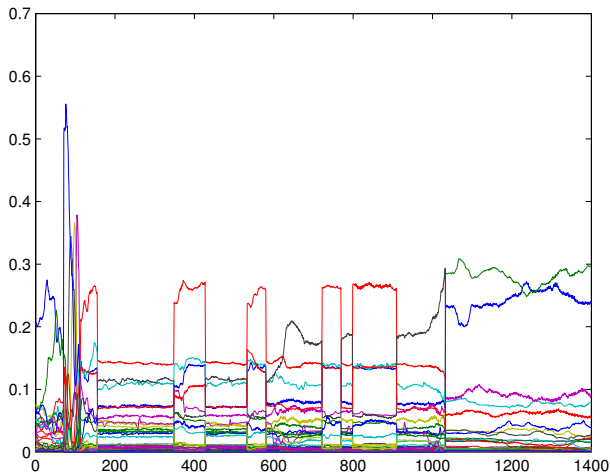
Example 1: 1-D signal



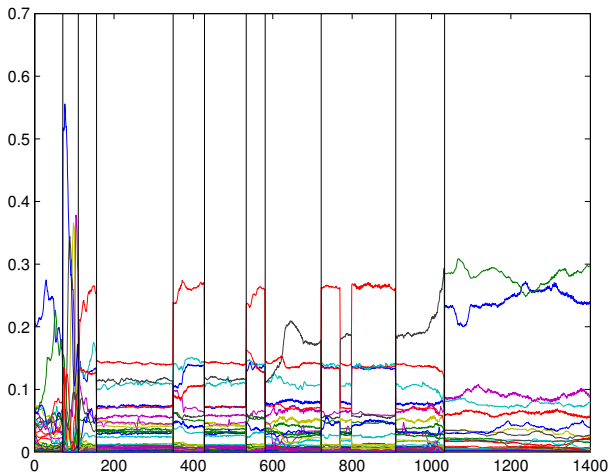
Example 1: 1-D signal: Find abrupt changes in the mean



Example 2: shot detection in a movie



Example 2: shot detection in a movie



The change-point problem

- **Observation:** $X_1, \dots, X_n \in \mathcal{X}$ independent random variables (\mathcal{X} : arbitrary measurable set).
 - P_{X_i} : distribution of X_i .
- ⇒ find where are the **abrupt changes in the sequence**
 P_{X_1}, \dots, P_{X_n} ?

Notation:

$$\tau \in \mathcal{T}_n^D := \{(\tau_0, \dots, \tau_D) \in \mathbb{N}^{D+1}, 0 = \tau_0 < \tau_1 < \dots < \tau_D = n\}$$

segmentation (of $\{1, \dots, n\}$) into $D_\tau = D \in \{1, \dots, n\}$ segments.

Challenges for (multiple) change-point detection

- ① Detect **changes in the whole distribution** (not only the mean)
 - Mean:
 - homoscedastic: Birgé & Massart (2001), Comte & Rozenholc (2002, 2004), Baraud, Giraud & Huet (2010)...
 - heteroscedastic: A. & Celisse (2011)
 - Mean and variance: Picard et al. (2007), Fryzlewicz and Subba Rao (2014)
 - Full distribution: Zou et al. (2014) in \mathbb{R} , Matteson & James (2014) in \mathbb{R}^d

Challenges for (multiple) change-point detection

- ① Detect **changes in the whole distribution** (not only the mean)
 - Mean:
 - homoscedastic: Birgé & Massart (2001), Comte & Rozenholc (2002, 2004), Baraud, Giraud & Huet (2010)...
 - heteroscedastic: A. & Celisse (2011)
 - Mean and variance: Picard et al. (2007), Fryzlewicz and Subba Rao (2014)
 - Full distribution: Zou et al. (2014) in \mathbb{R} , Matteson & James (2014) in \mathbb{R}^d
- ② **High-dimensional data** of different nature:
 - Vectorial: measures in \mathbb{R}^d , curves (sound recordings, ...)
 - Non vectorial: phenotypic data, graphs, DNA sequence, ...
 - Both vectorial and non vectorial data.

Challenges for (multiple) change-point detection

- ① Detect **changes in the whole distribution** (not only the mean)
 - Mean:
 - homoscedastic: Birgé & Massart (2001), Comte & Rozenholc (2002, 2004), Baraud, Giraud & Huet (2010)...
 - heteroscedastic: A. & Celisse (2011)
 - Mean and variance: Picard et al. (2007), Fryzlewicz and Subba Rao (2014)
 - Full distribution: Zou et al. (2014) in \mathbb{R} , Matteson & James (2014) in \mathbb{R}^d
- ② **High-dimensional data** of different nature:
 - Vectorial: measures in \mathbb{R}^d , curves (sound recordings, ...)
 - Non vectorial: phenotypic data, graphs, DNA sequence, ...
 - Both vectorial and non vectorial data.
- ③ **Efficient algorithm** allowing to deal with large data sets

Kernels: a quick reminder

- $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ measurable is a **positive semidefinite kernel** if $\forall x_1, \dots, x_m \in \mathcal{X}$, the matrix $(k(x_i, x_j))_{1 \leq i, j \leq m}$ is **positive semidefinite**.
- Examples:
 - **linear** kernel: $k(x, y) = \langle x, y \rangle$,
 - **polynomial** kernel: $k(x, y) = (1 + \langle x, y \rangle)^p$,
 - **Gaussian** kernel: $k(x, y) = \exp(-\|x - y\|^2 / (2h^2))$,
 - χ^2 kernel on Δ^d : $k(x, y) = \exp\left(-\frac{1}{h \cdot d} \sum_{i=1}^d \frac{(x_i - y_i)^2}{x_i + y_i}\right)$
 - ...

The kernel least-squares criterion

- **Least-squares criterion** (when $\mathcal{X} = \mathbb{R}$): $\forall \tau \in \mathcal{T}_n := \bigcup_{D \geq 1} \mathcal{T}_n^D$,

$$\widehat{\mathcal{R}}_n(\tau) := \frac{1}{n} \sum_{\ell=1}^D \sum_{i=\tau_{\ell-1}+1}^{\tau_{\ell}} (X_i - \bar{X}_{\tau_{\ell-1}+1, \tau_{\ell}})^2.$$

- **Kernel least-squares criterion:**

$$\widehat{\mathcal{R}}_n(\tau) := \frac{1}{n} \sum_{i=1}^n k(X_i, X_i) - \frac{1}{n} \sum_{\ell=1}^D \left[\frac{1}{\tau_{\ell} - \tau_{\ell-1}} \sum_{i=\tau_{\ell-1}+1}^{\tau_{\ell}} \sum_{j=\tau_{\ell-1}+1}^{\tau_{\ell}} k(X_i, X_j) \right].$$

- The two definitions coincide when $\mathcal{X} = \mathbb{R}$ and $k(x, y) = xy$.

Kernel change-point detection (KCP)

$$\hat{\tau} \in \underset{\tau \in \mathcal{T}_n}{\operatorname{argmin}} \left\{ \overbrace{\hat{\mathcal{R}}_n(\tau)}^{\substack{\text{kernel} \\ \text{least-squares} \\ \text{criterion}}} + \underbrace{\operatorname{pen}(\tau)}_{\substack{\text{penalty} \\ \text{function}}} \right\} \quad (\text{A., Celisse \& Harchaoui, 2012})$$

where pen is a function increasing with D_{τ} , such as:

$$\operatorname{pen}(\tau) = \frac{1}{n} \left(c_1 \log \left(\frac{n-1}{D_{\tau}-1} \right) + c_2 D_{\tau} \right)$$

$$\operatorname{pen}(\tau) = \frac{D_{\tau}}{n} \left(c_1 \log \left(\frac{n}{D_{\tau}} \right) + c_2 \right)$$

$$\operatorname{pen}(\tau) = \frac{c_1 D_{\tau}}{n}.$$

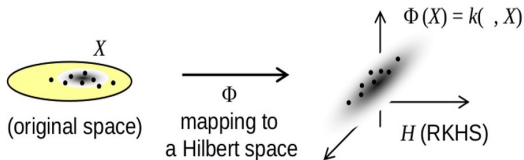
For $\mathcal{X} = \mathbb{R}$, linear kernel, Birgé & Massart (2001) and Lebarbier (2005) take $\operatorname{pen}(\tau) = \frac{\sigma^2 D_{\tau}}{n} \left[c_1 \log \left(\frac{n}{D_{\tau}} \right) + c_2 \right]$.

(Abstract) intuition on KCP

- KCP \Leftrightarrow kernelized version of (penalized) least-squares change-point detection

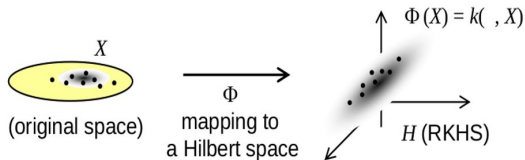
(Abstract) intuition on KCP

- KCP \Leftrightarrow kernelized version of (penalized) least-squares change-point detection
- Canonical feature map $\Phi : x \in \mathcal{X} \mapsto k(x, \cdot) \in \mathcal{H}$ reproducing kernel Hilbert space (RKHS)
- $Y_i = \Phi(X_i) \in \mathcal{H}$ are independent \mathcal{H} -valued r.v.



(Abstract) intuition on KCP

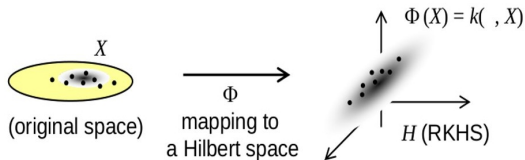
- KCP \Leftrightarrow kernelized version of (penalized) least-squares change-point detection
- Canonical feature map $\Phi : x \in \mathcal{X} \mapsto k(x, \cdot) \in \mathcal{H}$ reproducing kernel Hilbert space (RKHS)
- $Y_i = \Phi(X_i) \in \mathcal{H}$ are independent \mathcal{H} -valued r.v.



- $\mathbb{E}[\sqrt{k(X_i, X_i)}] < \infty \Rightarrow$ can define $\mu_i^* \in \mathcal{H}$ the “mean” of Y_i
- \Rightarrow KCP detects jumps of the “mean” μ_i^* of Y_i

(Abstract) intuition on KCP

- KCP \Leftrightarrow kernelized version of (penalized) least-squares change-point detection
- Canonical feature map $\Phi : x \in \mathcal{X} \mapsto k(x, \cdot) \in \mathcal{H}$ reproducing kernel Hilbert space (RKHS)
- $Y_i = \Phi(X_i) \in \mathcal{H}$ are independent \mathcal{H} -valued r.v.



- $\mathbb{E}[\sqrt{k(X_i, X_i)}] < \infty \Rightarrow$ can define $\mu_i^* \in \mathcal{H}$ the “mean” of Y_i
- \Rightarrow KCP detects jumps of the “mean” μ_i^* of Y_i
- Remark: if k is characteristic (eg, Gaussian kernel), μ_i^* characterizes P_{X_i} .

KCP for fixed D (Harchaoui & Cappé, 2007)

$$\hat{\tau}(D) \in \operatorname{argmin}_{\tau \in \mathcal{T}_n^D} \{\hat{\mathcal{R}}_n(\tau)\}$$

- Dynamic programming algorithm
- **No computation in \mathcal{H}** , only needs to compute the $k(X_i, X_j)$ (cost \mathcal{C}_k)
- Complexity of computing $(\hat{\tau}(D))_{1 \leq D \leq D_{\max}}$:

time $\mathcal{O}((\mathcal{C}_k + D_{\max})n^2)$ and space $\mathcal{O}(D_{\max}n)$

(Celisse et al., 2016).

Main assumptions

- \mathcal{H} separable
- Bounded kernel/data:

$$\exists M < +\infty, \forall i \in \{1, \dots, n\}, \quad k(X_i, X_i) \leq M^2 \text{ a.s.} \quad (\text{Db})$$

⇒ always satisfied for Gaussian and χ^2 kernel.

$D = D_{\tau^*}$ known: notations

- True segmentation τ^* :

$$\mu_1^* = \dots = \mu_{\tau_1^*}^* \neq \mu_{\tau_1^*+1}^* = \dots = \mu_{\tau_2^*}^* \neq \dots \neq \mu_{\tau_{D_{\tau^*}-1}^*+1}^* = \dots = \mu_n^*.$$

- Smallest jump size: $\underline{\Delta} := \min_{i / \mu_i^* \neq \mu_{i+1}^*} \|\mu_i^* - \mu_{i+1}^*\|_{\mathcal{H}}$
(MMD, Gretton et al. 2006).
- Smallest segment length: $\underline{\Lambda}_{\tau} := \frac{1}{n} \min_{1 \leq \ell \leq D_{\tau}} |\tau_{\ell} - \tau_{\ell-1}|$.
- Loss between segmentations $\tau^1, \tau^2 \in \mathcal{T}_n$:

$$\begin{aligned} d_{\infty, n}(\tau^1, \tau^2) &:= \frac{1}{n} \max_{1 \leq i \leq D_{\tau^1}-1} \left\{ \min_{1 \leq j \leq D_{\tau^2}-1} \left| \tau_i^1 - \tau_j^2 \right| \right\} \\ &= \frac{1}{n} \max_{1 \leq i \leq D_{\tau^1}-1} \left| \tau_i^1 - \tau_i^2 \right| \quad \text{if } D_{\tau^1} = D_{\tau^2} \text{ and } \tau^1, \tau^2 \text{ "close"} \end{aligned}$$

$D = D_{\tau^*}$ known: estimation of change-points locations

Theorem (A. & Garreau, 2016)

Assume: \mathcal{H} separable, (\mathbf{Db}) , $y > 0$ and

$$\underline{\Delta}_{\tau^*} > v_n(y) := \frac{148 D_{\tau^*} M^2}{\underline{\Delta}^2} \cdot \frac{y + \log n + 1}{n}.$$

Then, with probability $1 - e^{-y}$,

$$\forall \hat{\tau}(D_{\tau^*}) \in \operatorname{argmin}_{\tau \in \mathcal{T}_n^{D_{\tau^*}}} \{\hat{\mathcal{R}}_n(\tau)\}, \quad d_{\infty, n}(\tau^*, \hat{\tau}(D_{\tau^*})) \leq v_n(y).$$

$D = D_{\tau^*}$ known: estimation of change-points locations (2)

Corollary (A. & Garreau, 2016, simplified result)

Assume: \mathcal{H} separable, **(Db)** and $\frac{\Delta^2}{M^2} \gtrsim \frac{D_{\tau^*}}{\underline{\Delta}_{\tau^*}} \cdot \frac{\log n}{n}$.

Then, with probability $1 - n^{-2}$,

$$\forall \hat{\tau}(D_{\tau^*}) \in \operatorname{argmin}_{\tau \in \mathcal{T}_n^{D_{\tau^*}}} \{\hat{\mathcal{R}}_n(\tau)\}, \quad d_{\infty, n}(\tau^*, \hat{\tau}(D_{\tau^*})) \lesssim \frac{D_{\tau^*} M^2 \log n}{\underline{\Delta}^2 \cdot n}.$$

- $\frac{\Delta^2}{M^2} \approx$ signal-to-noise ratio.
- Matches **minimax lower bound** $\log(n)/n$ (Brunel, 2014).
- Remark: no $\log(n)$ factor in the standard “asymptotic” setting (Korostelev & Tsybakov, 2012).

KCP: data-driven D by model selection

- Notation: $Y = (Y_1, \dots, Y_n) \in \mathcal{H}^n$, $\mu^* = (\mu_1^*, \dots, \mu_n^*) \in \mathcal{H}^n$
 - For any $\tau \in \mathcal{T}_n$, $\Pi_\tau : \mathcal{H}^n \rightarrow \mathcal{H}^n$ orthogonal projection onto $F_\tau = \{(f_1, \dots, f_n) \in \mathcal{H}^n / f_{\tau_{\ell-1}+1} = \dots = f_{\tau_\ell} \forall \ell = 1, \dots, D_\tau\}$
- ⇒ **Least-squares estimator** $\hat{\mu}_\tau = \Pi_\tau Y$
and least-squares criterion:
- $$\hat{\mathcal{R}}_n(\tau) = \frac{1}{n} \|Y - \hat{\mu}_\tau\|^2 = \frac{1}{n} \sum_{i=1}^n \|Y_i - (\hat{\mu}_\tau)_i\|_{\mathcal{H}}^2$$

KCP: data-driven D by model selection

- Notation: $Y = (Y_1, \dots, Y_n) \in \mathcal{H}^n$, $\mu^* = (\mu_1^*, \dots, \mu_n^*) \in \mathcal{H}^n$
- For any $\tau \in \mathcal{T}_n$, $\Pi_\tau : \mathcal{H}^n \rightarrow \mathcal{H}^n$ orthogonal projection onto $F_\tau = \{(f_1, \dots, f_n) \in \mathcal{H}^n / f_{\tau_{\ell-1}+1} = \dots = f_{\tau_\ell} \forall \ell = 1, \dots, D_\tau\}$

⇒ **Least-squares estimator** $\widehat{\mu}_\tau = \Pi_\tau Y$

and least-squares criterion:

$$\widehat{\mathcal{R}}_n(\tau) = \frac{1}{n} \|Y - \widehat{\mu}_\tau\|^2 = \frac{1}{n} \sum_{i=1}^n \|Y_i - (\widehat{\mu}_\tau)_i\|_{\mathcal{H}}^2$$

- **Quadratic risk** of $\mu \in \mathcal{H}^n$:

$$\mathcal{R}(\mu) = \frac{1}{n} \|\mu - \mu^*\|^2 = \frac{1}{n} \sum_{i=1}^n \|\mu_i - \mu_i^*\|_{\mathcal{H}}^2 .$$

- Usual approach for **model selection**: take a penalty such that

$$\forall \tau \in \mathcal{T}_n, \quad \text{pen}(\tau) \geq \text{pen}_{\text{id}}(\tau) := \mathcal{R}(\mu) - \widehat{\mathcal{R}}_n(\tau) + \text{cst} .$$

Oracle inequality for KCP

Theorem (A., Celisse & Harchaoui, 2012–2016)

Assume: \mathcal{H} separable, $(\mathbf{D}\mathbf{b})$, $y > 0$, $C \geq 119$ and

$$\forall \tau \in \mathcal{T}_n, \quad \text{pen}(\tau) \geq \frac{CM^2}{n} \left[\log \binom{n-1}{D_\tau-1} + D_\tau \right].$$

Then, with probability $1 - e^{-y}$,

$$\forall \hat{\tau} \in \operatorname{argmin}_{\tau \in \mathcal{T}_n} \left\{ \hat{\mathcal{R}}_n(\tau) + \text{pen}(\tau) \right\},$$

$$\mathcal{R}(\hat{\mu}_{\hat{\tau}}) \leq 2 \inf_{\tau \in \mathcal{T}_n} \left\{ \mathcal{R}(\hat{\mu}_\tau) + \text{pen}(\tau) \right\} + \frac{83yM^2}{n}.$$

- applies to $\text{pen}(\tau) = \frac{CM^2 D_\tau}{n}$ if $C \geq 465 \log(n)$.
- $\mathcal{X} = \mathbb{R}$, linear kernel: Birgé & Massart (2001), Lebarbier (2005).

Change-point estimation performance of KCP

Theorem (A. & Garreau, 2016)

Assume: \mathcal{H} separable, (\mathbf{Db}) , $y > 0$ and

$$C_{\min} := \frac{74}{3}(D_{\tau^*} + 1)(y + \log n + 1) < C < C_{\max} := \frac{\Delta^2}{M^2} \frac{\Lambda_{\tau^*}}{6D_{\tau^*}} n.$$

Then, with probability $1 - e^{-y}$:

$$\forall \hat{\tau} \in \operatorname{argmin}_{\tau \in \mathcal{T}_n} \left\{ \hat{\mathcal{R}}_n(\tau) + \frac{CM^2 D_{\tau}}{n} \right\}, \quad D_{\hat{\tau}} = D_{\tau^*}$$

$$\text{and} \quad d_{\infty, n}(\tau^*, \hat{\tau}) \leq v_n(y) := \frac{148 D_{\tau^*} M^2}{\Delta^2} \cdot \frac{y + \log n + 1}{n}.$$

Previous works (Lavielle & Moulines, 2000, among many others):
real case ($\mathcal{H} = \mathbb{R}$) only (with dependent data).

Change-point estimation performance of KCP (2)

Corollary (A. & Garreau, 2016, simplified result)

Assume: \mathcal{H} separable, **(Db)** and

$$D_{\tau^*} \log n \lesssim C \lesssim \frac{\Delta^2}{M^2} \frac{\Lambda_{\tau^*}}{D_{\tau^*}} n.$$

Then, with probability $1 - n^{-2}$:

$$\forall \hat{\tau} \in \operatorname{argmin}_{\tau \in \mathcal{T}_n} \left\{ \hat{\mathcal{R}}_n(\tau) + \frac{CM^2 D_{\tau}}{n} \right\}, \quad D_{\hat{\tau}} = D_{\tau^*}$$

$$\text{and} \quad d_{\infty, n}(\tau^*, \hat{\tau}) \lesssim \frac{D_{\tau^*} M^2}{\Delta^2} \cdot \frac{\log n}{n}.$$

- $\frac{\Delta^2}{M^2} \approx$ signal-to-noise ratio.
- Lower bound on C : $\log(n)$ necessary (Birgé & Massart, 2007).

Oracle inequality: proof ideas

- Notation: $\varepsilon = Y - \mu^* \in \mathcal{H}^n$
- **Ideal penalty:**

$$\begin{aligned} \text{pen}_{\text{id}}(\tau) &:= \mathcal{R}(\mu) - \widehat{\mathcal{R}}_n(\tau) + \frac{1}{n} \|\varepsilon\|^2 \\ &= \frac{2}{n} \underbrace{\langle \Pi_{\tau} \mu^* - \mu^*, \varepsilon \rangle}_{=-L_{\tau} \text{ (linear term)}} + \frac{2}{n} \underbrace{\|\Pi_{\tau} \varepsilon\|^2}_{=Q_{\tau} \text{ (quadratic term)}} \end{aligned}$$

- **Concentration** for L_{τ} and Q_{τ} around their expectations
- \Rightarrow show that **$\text{pen}(\tau) \geq \text{pen}_{\text{id}}(\tau)$ simultaneously for all $\tau \in \mathcal{T}_n$** , with probability $\geq 1 - e^{-\gamma}$.
- Previous work (Birgé & Massart, 2001): Gaussian assumption + real-valued functions \Rightarrow does not apply to RKHS case.

Concentration of the quadratic term

Proposition (A., Celisse & Harchaoui, 2012–2016)

Assume: \mathcal{H} separable and **(Db)**. Then, for every $\tau \in \mathcal{T}_n$, $x > 0$:

$$\|\Pi_{\tau}\varepsilon\|^2 - \mathbb{E} \left[\|\Pi_{\tau}\varepsilon\|^2 \right] \leq \frac{14M^2}{3} (x + 2\sqrt{2x}D_{\tau}) ,$$

with probability at least $1 - e^{-x}$.

Proof ideas:

- Pinelis-Sakhanenko's inequality ($\|\sum_{i \in \lambda} \varepsilon_i\|_{\mathcal{H}}$).
- Bernstein's inequality (upper bounding moments).

Concentration of the linear term

Proposition

Assume: \mathcal{H} separable and **(Db)**. Then, for every $\tau \in \mathcal{T}_n$, $x > 0$, with probability at least $1 - 2e^{-x}$:

$$|\langle \Pi_{\tau} \mu^* - \mu^*, \varepsilon \rangle| \leq \theta \|\Pi_{\tau} \mu^* - \mu^*\|^2 + \left(\frac{1}{2\theta} + \frac{4}{3} \right) M^2 x,$$

for every $\theta > 0$.

Proof: Bernstein's inequality.

Identification of change-points: proof ideas

$$\hat{\tau} \in \operatorname{argmin}_{\tau \in \mathcal{T}_n} \{ \hat{\mathcal{R}}_n(\tau) + \operatorname{pen}(\tau) \}$$

- Empirical risk:

$$\hat{\mathcal{R}}_n(\tau) = \underbrace{\frac{1}{n} \|\mu^* - \Pi_{\tau} \mu^*\|^2}_{=A_{\tau}(\text{approximation})} + \underbrace{\frac{2}{n} \langle \mu^* - \Pi_{\tau} \mu^*, \varepsilon \rangle}_{=L_{\tau}(\text{linear term})} - \underbrace{\frac{1}{n} \|\Pi_{\tau} \varepsilon\|^2}_{=Q_{\tau}(\text{quadratic term})} + \underbrace{\frac{1}{n} \|\varepsilon\|^2}_{(\text{constant})}$$

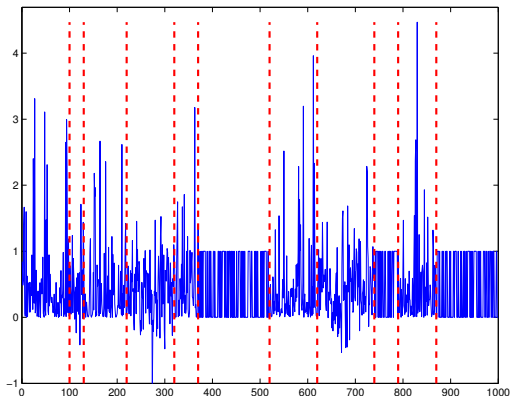
- Previous concentration inequalities for L_{τ} , Q_{τ} .
- Deterministic bounds on A_{τ} :

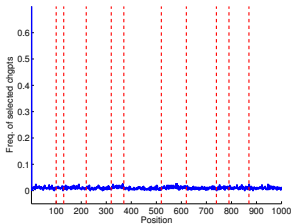
$$D_{\tau} < D_{\tau^*} \Rightarrow \frac{1}{n} A_{\tau} \geq \frac{1}{2} \underline{\Lambda}_{\tau^*} \underline{\Delta}^2 \quad (\text{for showing } D_{\hat{\tau}} \geq D_{\tau^*})$$

$$\frac{1}{n} A_{\tau} \geq \frac{1}{2} \min \left\{ \underline{\Lambda}_{\tau^*}, d_{\infty, n}(\tau^*, \tau) \right\} \underline{\Delta}^2 \quad (\text{for } \hat{\tau}(D_{\tau^*}))$$

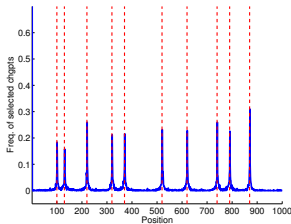
Constant mean and variance: data

Constant mean and variance: the distribution of X_i is chosen among $\mathcal{B}(0.5)$, $\mathcal{N}(0.5, 0.25)$ and $\mathcal{E}(0.5)$.

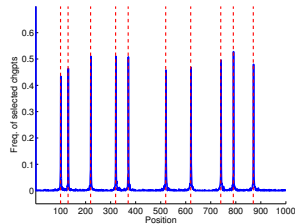


Constant mean and variance: results ($D_{\mathcal{T}^*}$)

Linear

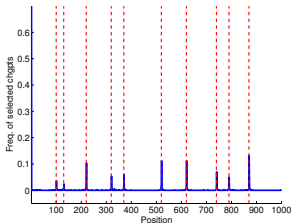


Hermite

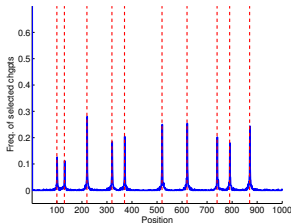


Gaussian

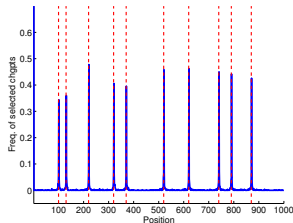
KCP with $D_{\mathcal{T}^*}$ known.

Constant mean and variance: results (\hat{D})

Linear



Hermite

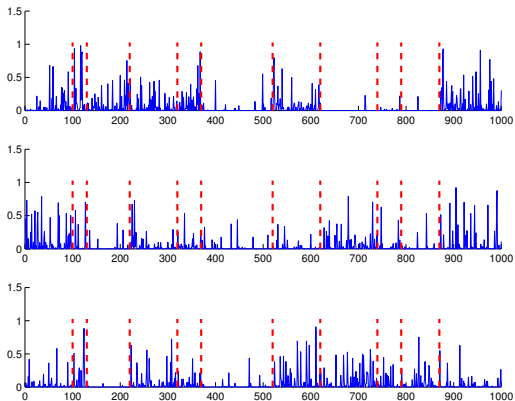


Gaussian

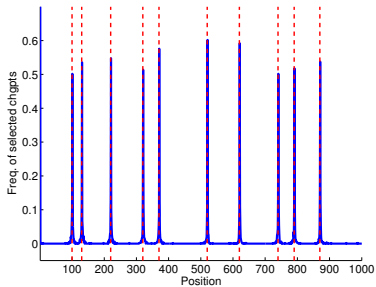
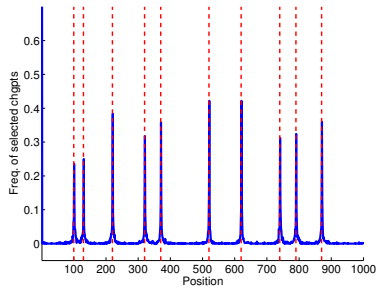
KCP with \hat{D} data-driven.

Histogram-valued data

$X_i \in d$ -dimensional simplex, Dirichlet distribution $(p_1^\ell, \dots, p_d^\ell)$ on the ℓ -th segment, with p_i^ℓ independent $\sim \mathcal{U}([0, 0.2])$.

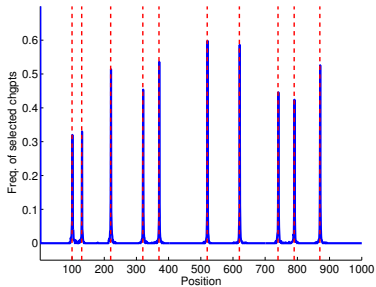
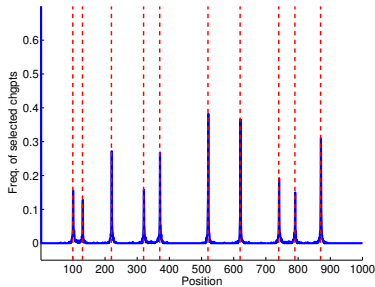


(first three coordinates)

Histogram-valued data: results (D_{T^*}) χ^2 kernel

Gaussian kernel

KCP with D_{T^*} known.

Histogram-valued data: results (\hat{D}) χ^2 kernel

Gaussian kernel

KCP with \hat{D} data-driven.

Conclusion

Take-home message:

- Kernelized version of penalized least-squares change-point detection (eg, Lebarbier, 2005).
- Detection of **changes in the distribution**, not only the first moments.
- Can deal with **structured data**.
- Under reasonable assumptions and for a class of penalty functions:
 - **oracle inequality**;
 - identifies the correct **number of change-points**;
 - estimates at the correct rate the **change-points locations**.

Future work:

- Unbounded data/kernel.
- Dependent data?
- Learn how to choose the kernel.