

Consistent Change-point Detection with Kernels

Damien Garreau¹ Sylvain Arlot²

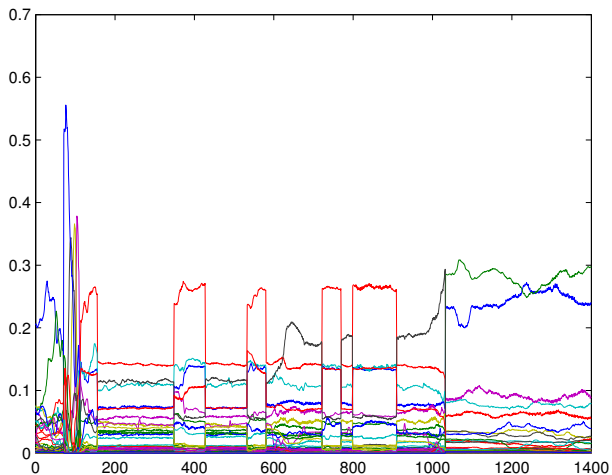
¹Inria, DI ENS

²Université Paris-Sud, Laboratoire de Mathématiques d'Orsay

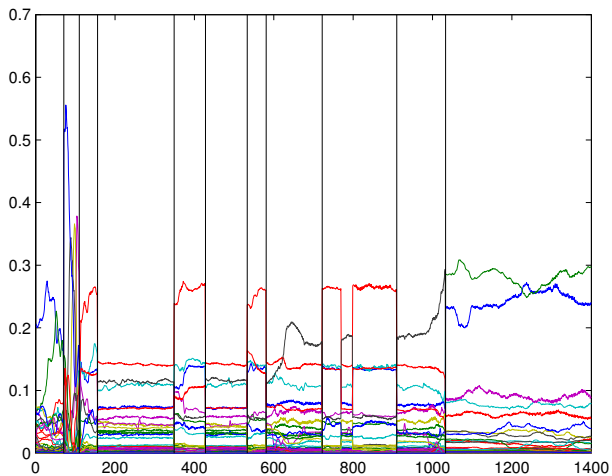
April 6, 2016



An example: shot detection in a movie



An example: shot detection, cont.



Introduction

Overview

The change-point problem

Algorithm

Kernel change-point algorithm

Experimental results

Theoretical results

Hypothesis

Dimension selection

Localization of the change points

Conclusion

We want to:

- ▶ detect abrupt changes in the *distribution* of the data
- ▶ deal with *interesting* (structured) data: each point is a curve, a graph, a histogram, a persistence diagram...

The change-point problem

- ▶ \mathcal{X} arbitrary (measurable) set, $n < +\infty$, and $X_1, \dots, X_n \in \mathcal{X}$ sequence of independent random variables.
- ▶ $\forall i \in \{1, \dots, n\}$, P_{X_i} the distribution of X_i .

The change-point problem can be formalized as follows:

- ▶ Given $(X_i)_{1 \leq i \leq n}$, we want to find the locations of the abrupt changes in the sequence P_{X_1}, \dots, P_{X_n} .

Notations

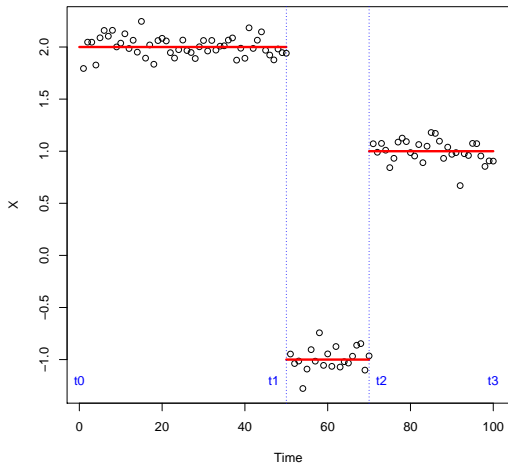
- ▶ Take any $D \in \{1, \dots, n + 1\}$, the set of sequences of $D - 1$ change-points is defined by

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

- ▶ $\tau_1, \dots, \tau_{D-1}$ are the *change-points*, τ is a *segmentation* of $\{1, \dots, n\}$ into D_τ segments.
- ▶ τ^* the true segmentation, $D^* = D_{\tau^*}$ the true number of change-points.

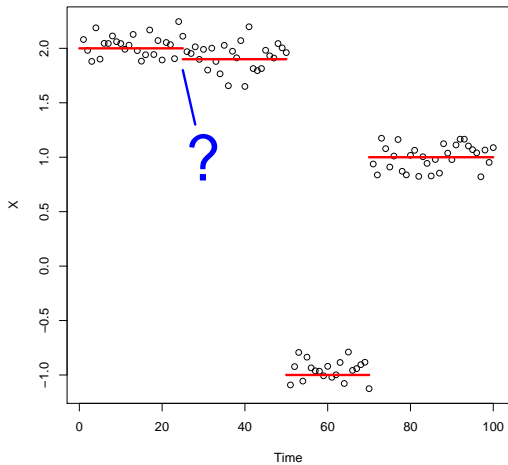
In pictures

Here, $\mathcal{X} = \mathbb{R}$, $D^* = 3$, and $\tau^* = (0, 50, 70, 100)$.



In pictures, cont.

It is not an easy question:



- ▶ With finite sample size, it is *not easy* to recover the true change-points in presence of noise.
- ▶ When $\mathcal{X} = \mathbb{R}^d$ and the changes occur in the first moments of the distribution, the problem has already received considerable attention, cf. [Basseville and Nikiforov, 1993].
- ▶ Kernel change-point detection can tackle more subtle changes / less conventional data.

Plan

Introduction

Overview

The change-point problem

Algorithm

Kernel change-point algorithm

Experimental results

Theoretical results

Hypothesis

Dimension selection

Localization of the change points

Conclusion

Kernels: a quick reminder

- ▶ Let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a semidefinite kernel.
- ▶ k is a measurable function s.t. $\forall x_1, \dots, x_m \in \mathcal{X}$, the matrix $(k(x_i, x_j))_{1 \leq i, j \leq m}$ is *positive semi-definite*. Think *inner product*.
- ▶ Examples include
 - ▶ the *linear* kernel $k(x, y) = \langle x, y \rangle$,
 - ▶ the *Gaussian* kernel $k(x, y) = \exp(-\|x - y\|^2 / (2h^2))$,
 - ▶ the *histogram* kernel $k(x, y) = \sum_{k=1}^p \min(x_k, y_k)$,
 - ▶ ...

The kernel least-squares criterion

- ▶ Intuition: least-squares criterion

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

- ▶ Define

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

- ▶ This is just a kernelized version, the two definitions coincide when $\mathcal{X} = \mathbb{R}$ and $k(x, y) = xy$.

Most important slide of the talk

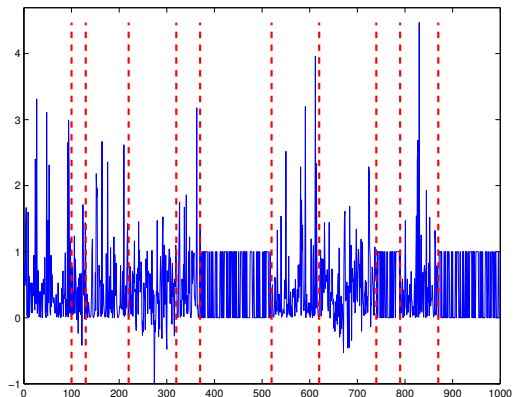
We investigate the properties of

$$\hat{\tau} \in \arg \min_{\tau \in \mathcal{T}_n} \left\{ \overbrace{\hat{\mathcal{R}}_n(\tau)}^{\text{least-squares criterion}} + \underbrace{\text{pen}(\tau)}_{\text{penalty function}} \right\},$$

where pen is a function increasing with D_τ .

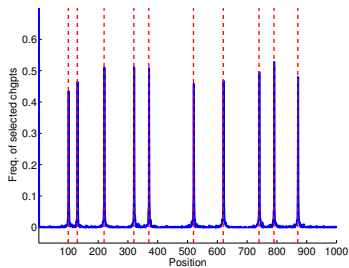
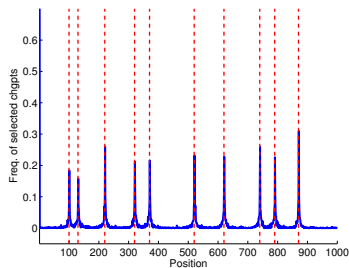
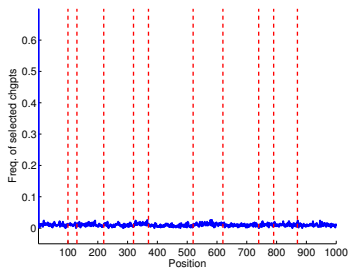
Constant mean and variance

Constant mean and variance: the distribution of X_i is chosen among $\mathcal{B}(0.5)$, $\mathcal{N}(0.5, 0.25)$ and $\Gamma(1, 0.5)$.



(courtesy of [Arlot et al., 2012])

Constant mean and variance, cont.



Linear, Hermite, and Gaussian kernels (courtesy of [Arlot et al., 2012]).

Plan

Introduction

Overview

The change-point problem

Algorithm

Kernel change-point algorithm

Experimental results

Theoretical results

Hypothesis

Dimension selection

Localization of the change points

Conclusion

More notations

- ▶ Along with the kernel k comes a reproducing kernel Hilbert space \mathcal{H} endowed with $\langle \cdot, \cdot \rangle_{\mathcal{H}}$.
- ▶ There exists a mapping $\Phi : \mathcal{X} \rightarrow \mathcal{H}$ such that, for any $x, y \in \mathcal{X}$, $k(x, y) = \langle \Phi(x), \Phi(y) \rangle_{\mathcal{H}}$.
- ▶ The algorithm is looking for breaks in the “mean” of $Y_i := \Phi(X_i) \in \mathcal{H}$.
- ▶ Whenever possible, define μ_i^* the mean of Y_i ; it satisfies

$$\forall g \in \mathcal{H}, \quad \langle \mu_i^*, g \rangle_{\mathcal{H}} = \mathbb{E} [g(X_i)] = \mathbb{E} [\langle Y_i, g \rangle_{\mathcal{H}}].$$

- ▶ We write $Y_i = \mu_i^* + \varepsilon_i$.

Hypothesis

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

$$\exists M \in (0, +\infty), \quad \forall 1 \leq i \leq n, \quad k(X_i, X_i) \leq M^2. \quad (\text{Db})$$

- ▶ Finite variance:

$$\forall 1 \leq i \leq n, \quad v_i := \mathbb{E} \left[\|\varepsilon_i\|_{\mathcal{H}}^2 \right] \leq V < +\infty. \quad (\text{V})$$

Under **(Db)**, an oracle inequality has been proven.

→ See [Arlot et al., 2012] for the result.

Dimension selection, light version

- ▶ Assume that (Db) holds true;
- ▶ Suppose that $\text{pen}(\cdot)$ is “large enough”;
- ▶ Suppose that $\underline{\Delta}^2 \times \underline{\Gamma}$ is “large enough”, where $\underline{\Delta} := \inf_{\mu_i^* \neq \mu_{i+1}^*} \|\mu_i^* - \mu_{i+1}^*\|_{\mathcal{H}}$ is the size of the smallest jump in \mathcal{H} , and $\underline{\Gamma}$ depends only on the geometry of τ^* ;
- ▶ Then, with high probability, $D_{\hat{\tau}} = D^*$.

If k is characteristic, we recover *all* the changes in P_{X_i} .

Theorem

Let y be a positive number. Assume that **(Db)** holds true and that

$$\forall \tau \in \mathcal{T}_n, \quad \text{pen}(\tau) = \frac{CD_\tau M^2}{n} \left(1 + \sqrt{2 \left(4 + y + \log \frac{n}{D_\tau} \right)} \right)^2,$$

with $C \geq (2D^* + 1)(5 + y + \log D^*)$. Suppose that

$$\underline{\Delta}^2 \times \underline{\Gamma} \gtrsim \frac{CD^* M^2}{n} \left(y + \log \frac{n}{D^*} \right).$$

Then $\mathbb{P}(D_{\hat{\tau}} = D^*) \geq 1 - e^{-y}$.

Distance between segmentations

- ▶ We consider only segmentation with the same number of segments D^* .
- ▶ Several possibilities, equivalent under assumptions regarding

$$\underline{\Lambda}_\tau := \frac{1}{n} \min_{\lambda \in \tau} |\lambda|.$$

- ▶ We focus on

$$d_\infty(\tau^1, \tau^2) := \max_{1 \leq i \leq D^* - 1} \left| \tau_i^1 - \tau_i^2 \right|.$$

Localization of the change-points, light version

- ▶ Assume that D^* is known and that (V) holds true.
- ▶ Take $\delta_n > 0$, and choose $\hat{\tau}$ in

$$\mathcal{T}_n^{D^*}(\delta_n) := \{\tau \in \mathcal{T}_n, \underline{\Lambda}_\tau \geq \delta_n, D_\tau = D^*\}.$$

- ▶ Then, for any $0 < x < \underline{\Lambda}_{\tau^*}$,

$$\mathbb{P}\left(\frac{1}{n}d_\infty(\hat{\tau}, \tau^*) \geq x\right) \lesssim \frac{V}{nx} \left(\frac{1}{\delta_n} + \frac{1}{x}\right).$$

- ▶ This goes to 0 whenever $\delta_n \rightarrow 0$ and $n\delta_n \rightarrow +\infty$.

Localization of the change-points

Theorem

Assum that (V) holds true. Take $\delta_n > 0$ and choose

$$\hat{\tau} \in \arg \min_{\tau \in \mathcal{T}_n^{D^*}(\delta_n)} \{\hat{\mathcal{R}}_n(\tau)\}.$$

Suppose that $\delta_n \leq \underline{\Lambda}_{\tau^*}$. Then, for any $0 < x \leq \underline{\Lambda}_{\tau^*}$,

$$\mathbb{P}(d_\infty(\hat{\tau}, \tau^*) \geq x) \lesssim \frac{VD^*}{nx\underline{\Delta}^2} \left(\frac{1}{\delta_n} + \frac{(D^*)^3 \overline{\Delta}^2}{x\underline{\Delta}^2} \right).$$

For instance: take $\delta_n = n^{-1/2}$: $d_\infty(\hat{\tau}, \tau^*) = o_P(n^{-1/2})$.

Plan

Introduction

Overview

The change-point problem

Algorithm

Kernel change-point algorithm

Experimental results

Theoretical results

Hypothesis

Dimension selection





Localization of the change points

Conclusion

- ▶ Kernelized version of the change-point detection procedure of [Lebarbier, 2005].
- ▶ Detection of changes in the distribution, not only the first moments.
- ▶ Possible to deal with structured data more efficiently.
- ▶ Under reasonable assumptions and for a class of penalty functions,
 - ▶ we dispose of an oracle inequality
 - ▶ the procedure is consistent
 - ▶ it recovers the true localization of the change-points

- ▶ Exchange the hypothesis and *still* prove our results (in progress);
- ▶ Tackle dependency structures within the X_i s as in [Lavielle and Moulines, 2000];
- ▶ Learn how to choose the kernel;
- ▶ Find interesting data!

Thank you for your attention!

-  Arlot, S., Celisse, A., and Harchaoui, Z. (2012). Kernel change-point detection.
arXiv preprint arXiv:1202.3878.
-  Basseville, M. and Nikiforov, I. (1993). *Detection of abrupt changes: theory and application.* Prentice Hall Englewood Cliffs.
-  Lavielle, M. and Moulines, E. (2000). Least-squares estimation of an unknown number of shifts in a time series.
Journal of time series analysis, 21(1):33–59.
-  Lebarbier, É. (2005). Detecting multiple change-points in the mean of a Gaussian process by model selection.
Signal Proces., 85:717–736.

Bonus: elements of proof (dimension selection)

Main idea: $\forall \tau \in \mathcal{T}_n$ s.t. $D_\tau \neq D^*$, w.h.p.,

$$\mathcal{R}(\tau) + \text{pen}(\tau) > \mathcal{R}(\tau^*) + \text{pen}(\tau^*).$$

Since $\hat{\tau}$ minimizes $\mathcal{R}(\cdot) + \text{pen}(\cdot)$, $\hat{D} = D^*$ w.h.p..

Bonus: elements of proof (dimension selection)

Main idea: $\forall \tau \in \mathcal{T}_n$ s.t. $D_\tau \neq D^*$, w.h.p.,

$$\mathcal{R}(\tau) + \text{pen}(\tau) > \mathcal{R}(\tau^*) + \text{pen}(\tau^*).$$

Since $\hat{\tau}$ minimizes $\mathcal{R}(\cdot) + \text{pen}(\cdot)$, $\hat{D} = D^*$ w.h.p..

$$\mathcal{R}(\tau) = \frac{1}{n} \|\mu^* - \mu_\tau^*\|^2 + \frac{2}{n} \langle \mu^* - \mu_\tau^*, \varepsilon \rangle - \frac{1}{n} \|\Pi_\tau \varepsilon\|^2 + \frac{1}{n} \|\varepsilon\|^2,$$

where μ_τ^* is the projection of μ^* on the subset of \mathcal{H}^n “constant on the segments of τ ”, i.e.,

$$F_\tau := \{f \in \mathcal{H}^n, f_{\tau_{\ell-1}+1} = \dots = f_{\tau_\ell} \forall 1 \leq \ell \leq D_\tau\}.$$

Bonus: elements of proof (dimension selection)

Main idea: $\forall \tau \in \mathcal{T}_n$ s.t. $D_\tau \neq D^*$, w.h.p.,

$$\mathcal{R}(\tau) + \text{pen}(\tau) > \mathcal{R}(\tau^*) + \text{pen}(\tau^*).$$

Since $\hat{\tau}$ minimizes $\mathcal{R}(\cdot) + \text{pen}(\cdot)$, $\hat{D} = D^*$ w.h.p..

$$\mathcal{R}(\tau) = \frac{1}{n} \|\mu^* - \mu_\tau^*\|^2 + \frac{2}{n} \langle \mu^* - \mu_\tau^*, \varepsilon \rangle - \frac{1}{n} \|\Pi_\tau \varepsilon\|^2 + \frac{1}{n} \|\varepsilon\|^2,$$

where μ_τ^* is the projection of μ^* on the subset of \mathcal{H}^n “constant on the segments of τ ”, i.e.,

$$F_\tau := \{f \in \mathcal{H}^n, f_{\tau_{\ell-1}+1} = \dots = f_{\tau_\ell} \forall 1 \leq \ell \leq D_\tau\}.$$

We are reduced to show that if $D_\tau \neq D^*$, w.h.p.,

$$\begin{aligned} \frac{1}{n} \|\mu^* - \mu_\tau^*\|^2 + \frac{2}{n} \langle \mu^* - \mu_\tau^*, \varepsilon \rangle + \frac{1}{n} \|\Pi_{\tau^*} \varepsilon\|^2 - \frac{1}{n} \|\Pi_\tau \varepsilon\|^2 &> \\ &> \text{pen}(\tau^*) - \text{pen}(\tau). \end{aligned}$$

We control each term $\forall \tau$, with probability $1 - e^{-x}$:

- ▶ the linear term: $|\langle \mu^\star - \mu_\tau^\star, \varepsilon \rangle| \lesssim \theta \|\mu^\star - \mu_\tau^\star\|^2 + \frac{1}{\theta} M^2 x$,
- ▶ the quadratic term:
$$\|\Pi_\tau \varepsilon\|^2 - \mathbb{E} \left[\|\Pi_\tau \varepsilon\|^2 \right] \lesssim (x + \sqrt{x D_\tau}) M^2,$$
- ▶ $\text{pen}(\tau) - \text{pen}(\tau^\star)$ via technical lemmas.

Elements of proof, cont.

Union bound,

$$\mathbb{P} \left(\bigcap_{\tau \in \mathcal{T}_n} \Omega_\tau \right) \geq 1 - 4 \sum_{\tau \in \mathcal{T}_n} e^{-x_\tau}.$$

Recall $\#\mathcal{T}_n^D = \binom{n-1}{D-1} \leq (ne/d)^d$.

$$\begin{aligned} \sum_{\tau \in \mathcal{T}_n} e^{-x_\tau} &\leq \sum_{d=1}^n \exp(d + d \log n - d \log d - 4d - dy - d \log n + d \log d) \\ &= \sum_{d=1}^n \exp(d(-3 - y)) = \sum_{d=1}^n \left(\frac{e^{-y}}{e^3} \right)^d \\ &= e^{-3-y} \cdot \frac{1 - (e^{-3-y})^n}{1 - e^{-3-y}} \leq e^{-y} / 4. \end{aligned}$$

This part *fails* if we do not have bounded assumption.

Bonus: a word about the concentration result

Lemma

For every $x > 0$, with probability $\geq 1 - e^{-x}$,

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

Proof.

Write $\|\Pi_{\tau}\varepsilon\|^2$ as $\sum_{1 \leq \ell \leq D_{\tau}} T_{\ell}$, a sum of independent random variables, where $T_{\ell} := \frac{1}{\tau_{\ell} - \tau_{\ell-1}} \left\| \sum_{j=\tau_{\ell-1}+1}^{\tau_{\ell}} \varepsilon_j \right\|^2$. Apply Bernstein's inequality. The tricky part is to check that the moment conditions for Bernstein are satisfied. Idea: Write $\mathbb{E} [T_{\ell}^q]$ as an integral depending upon $\mathbb{P} \left(\left\| \sum_{j=\tau_{\ell-1}+1}^{\tau_{\ell}} \varepsilon_j \right\| \geq y \right)$ for which Pinelis-Sakhanenko's inequality gives an upper-bound. □