

Kernel change-point detection

ALAIN CELISSE

¹UMR 8524 CNRS - Université Lille 1

²MODAL INRIA team-project

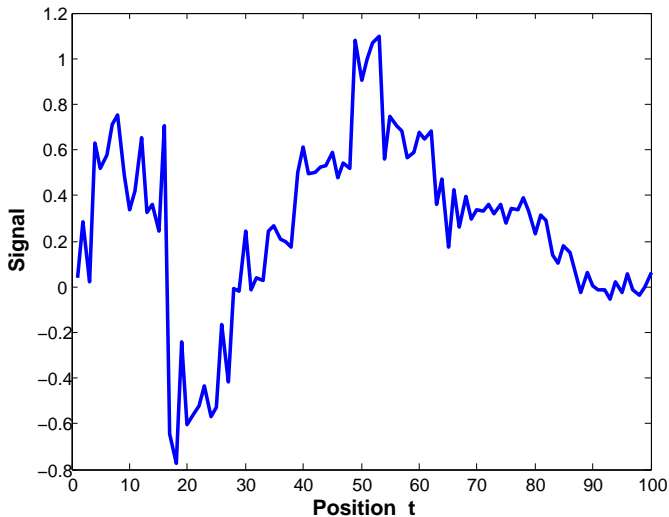
³SSB Group, Paris

joint work with Sylvain Arlot and Zaïd Harchaoui

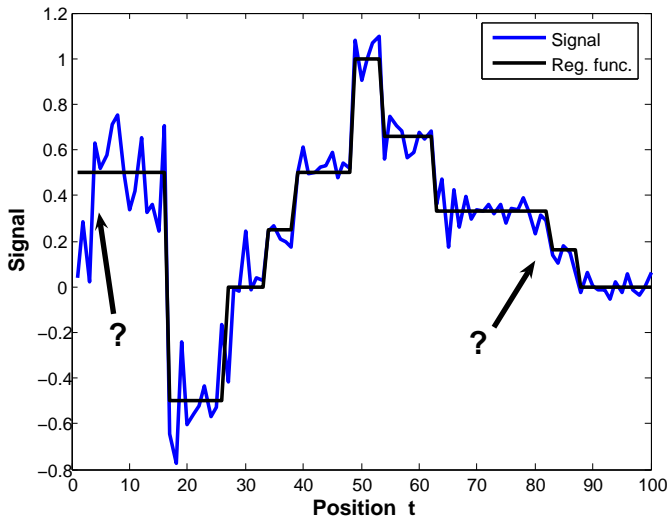
Workshop: "Recent advances in changepoints analysis"

Warwick University, March 28, 2012

1-D signal (example)



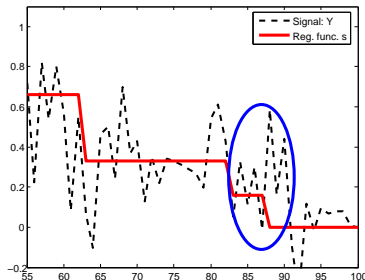
1-D signal (example): Find abrupt changes in the mean



Estimation rather than identification

Fact:

With finite sample, it is impossible to recover change-point in noisy regions.



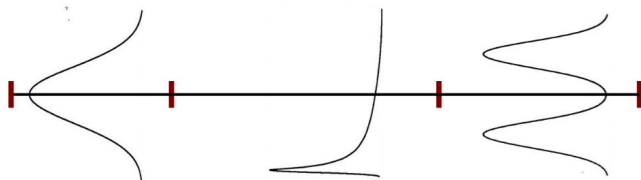
Purpose:

Estimate the regression function.

Idea:

→ Without too strong noise, recover true change-points.

Example 1: Changes in the distribution

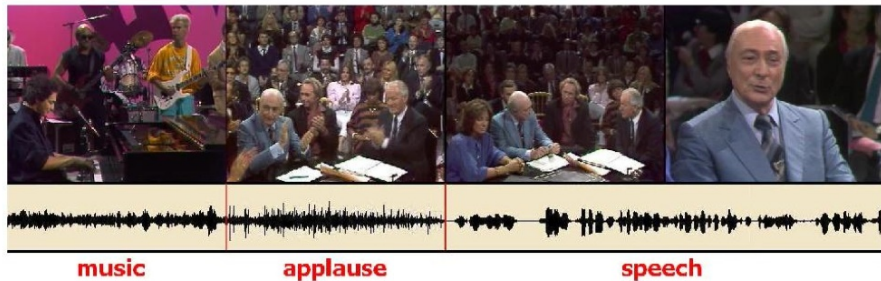


Description:

- Observations generated along the time.
- Observation distribution is piecewise constant along the time.
- Observations have the same mean and variance.

→ Detecting changes in the mean and variance is useless.

Example 2: Working with non-vectorial objects



Description:

- Video sequences from “Le grand échiquier”, 70s-80s French talk show.
- At each time, one observes an image (high-dimensional).
- Each image is summarized by a histogram.

Example 2: Working with non-vectorial objects

- Preprocessing images (patches in yellow).
- Each histogram bin corresponds to a patch.



Non-vectorial object:

Histograms with D bins belong to

$$\left\{ (p_1, \dots, p_D) \in [0, 1]^D, \sum_{i=1}^D p_i = 1 \right\}.$$

→ Algorithms for vectorial data are useless.

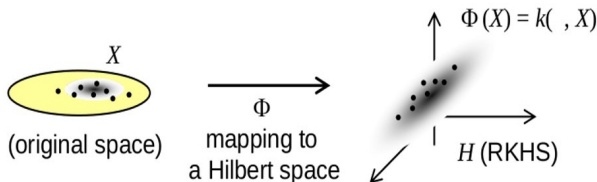
Detect abrupt changes. . .

General purposes:

- 1 Detect **changes in the whole distribution** (not only in the mean)
- 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

Kernel and Reproducing Kernel Hilbert Space (RKHS)

- $X_1, \dots, X_n \in \mathcal{X}$: initial observations.
- $k(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$: reproducing kernel (\mathcal{H} : RKHS).
- $\phi(\cdot) : \mathcal{X} \rightarrow \mathcal{H}$ s.t. $\phi(x) = k(x, \cdot)$: canonical feature map.
- $\langle \cdot, \cdot \rangle_{\mathcal{H}}$: inner-product in \mathcal{H} .



Asset:

Enables to work with **high-dimensional heterogeneous data**.

Model

Mapping of the initial data

$$\forall 1 \leq i \leq n, \quad Y_i = \phi(X_i) \in \mathcal{H} .$$

$\longrightarrow (t_1, Y_1), \dots, (t_n, Y_n) \in [0, 1] \times \mathcal{H} : \text{independent} .$

Model

Mapping of the initial data

$$\forall 1 \leq i \leq n, \quad Y_i = \phi(X_i) \in \mathcal{H} .$$

$\longrightarrow (t_1, Y_1), \dots, (t_n, Y_n) \in [0, 1] \times \mathcal{H} : \quad \text{independent} .$

Mean element

The **mean element of P_{X_i}** (distribution of X_i) is μ_i^* :

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

Fact:

With characteristic kernels,

$$P_{X_i} \neq P_{X_j} \quad \Rightarrow \quad \mu_i^* \neq \mu_j^* .$$

Model

$$\forall 1 \leq i \leq n, \quad Y_i = \mu_i^* + \varepsilon_i \in \mathcal{H},$$

where

- $\mu_i^* \in \mathcal{H}$: mean element of P_{X_i} (distribution of X_i)
- $\forall i, \quad \varepsilon_i := Y_i - \mu_i^*, \quad \text{with} \quad \mathbb{E}\varepsilon_i = 0, \quad v_i := \mathbb{E} \left[\|\varepsilon_i\|_{\mathcal{H}}^2 \right].$

Model

$$\forall 1 \leq i \leq n, \quad Y_i = \mu_i^* + \varepsilon_i \in \mathcal{H},$$

where

- $\mu_i^* \in \mathcal{H}$: **mean element of P_{X_i}** (distribution of X_i)
- $\forall i, \quad \varepsilon_i := Y_i - \mu_i^*, \quad \text{with} \quad \mathbb{E}\varepsilon_i = 0, \quad v_i := \mathbb{E} \left[\|\varepsilon_i\|_{\mathcal{H}}^2 \right].$

Assumptions

- ① $\max_i \|Y_i\|_{\mathcal{H}} \leq M \quad \text{a.s.} \quad (\mathbf{Db}).$
- ② $\max_i v_i \leq v_{\max} \quad (\mathbf{Vmax}).$
- ③ $\mu^* = (\mu_1^*, \dots, \mu_n^*)' \in \mathcal{H}^n$: **piecewise constant.**

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

Goal: \longrightarrow Estimate μ^* to recover change-points.

Model selection

Models:

- $\mathcal{M}_n = \{m, \text{segmentation of } \{1, \dots, n\}\}, \quad D_m = \text{Card}(m).$
- $S_m = \{\mu : (t_1, \dots, t_n) \rightarrow \mathcal{H}, \text{ piecewise const. on } (I_\lambda)_{\lambda \in m}\},$
 $(I_\lambda)_{\lambda \in m}: I_1 = [0, t_{\lambda_1}], I_2 =]t_{\lambda_1}, t_{\lambda_2}], \dots, I_{D_m} =]t_{\lambda_{D_m-1}}, 1].$

Strategy:

$$(S_m)_{m \in \mathcal{M}_n} \longrightarrow (\hat{\mu}_m)_{m \in \mathcal{M}_n} \longrightarrow \hat{\mu}_{\hat{m}} \quad ???$$

Model selection

Models:

- $\mathcal{M}_n = \{m, \text{segmentation of } \{1, \dots, n\}\}, \quad D_m = \text{Card}(m).$
- $S_m = \{\mu : (t_1, \dots, t_n) \rightarrow \mathcal{H}, \text{ piecewise const. on } (I_\lambda)_{\lambda \in m}\},$
 $(I_\lambda)_{\lambda \in m}: I_1 = [0, t_{\lambda_1}], I_2 =]t_{\lambda_1}, t_{\lambda_2}], \dots, I_{D_m} =]t_{\lambda_{D_m-1}}, 1].$

Strategy:

$$(S_m)_{m \in \mathcal{M}_n} \longrightarrow (\hat{\mu}_m)_{m \in \mathcal{M}_n} \longrightarrow \hat{\mu}_{\hat{m}} \quad ???$$

Goal:

Oracle inequality (in expectation, or with large probability):

$$\|\mu^* - \hat{\mu}_{\hat{m}}\|^2 \leq C \inf_{m \in \mathcal{M}_n} \left\{ \|\mu^* - \hat{\mu}_m\|^2 \right\} + r_n .$$

Least-squares estimator

- Empirical risk minimizer over S_m (= model):

$$\hat{\mu}_m \in \arg \min_{u \in S_m} \frac{1}{n} \sum_{i=1}^n \|u(t_i) - Y_i\|_{\mathcal{H}}^2 \left(=: \arg \min_{u \in S_m} P_n \gamma(u) \right).$$

- **Regressogram:**

$$\hat{\mu}_m = \sum_{\lambda \in \mathcal{M}} \hat{\beta}_\lambda \mathbf{1}_{I_\lambda} \quad \hat{\beta}_\lambda = \frac{1}{\text{Card} \{t_i \in I_\lambda\}} \sum_{t_i \in I_\lambda} Y_i.$$

Choose $D - 1$ change-points. . .

Assumption: (Harchaoui, Cappé (2007))

The number $D - 1$ of change-points is known.

Strategy:

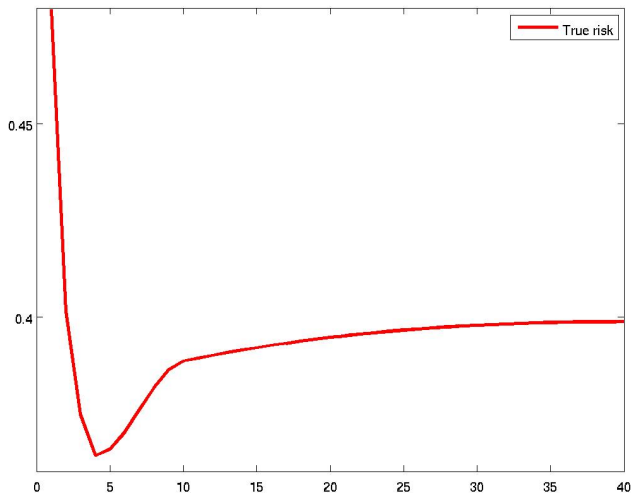
Choose $\hat{m}(D)$ among $\{m \in \mathcal{M}_n, D_m = D\}$.

ERM algorithm:

$$\hat{m}(D) = \hat{m}_{\text{ERM}}(D) = \text{Argmin}_{m|D_m=D} \|Y - \hat{\mu}_m\|^2.$$

(dynamic programming).

Quality of the segmentations



Elementary calculations

Expectations

$$(v_\lambda = \frac{1}{\text{Card}(\lambda)} \sum_{i \in \lambda} v_i)$$

$$\mathbb{E} \left[\|\mu^* - \hat{\mu}_m\|^2 \right] = \|\mu^* - \Pi_m \mu^*\|^2 + \sum_{\lambda \in m} v_\lambda ,$$

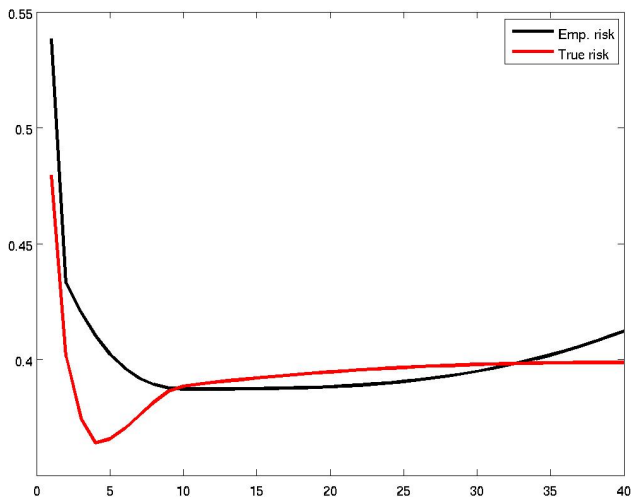
$$\mathbb{E} \left[\|Y - \hat{\mu}_m\|^2 \right] = \|\mu^* - \Pi_m \mu^*\|^2 - \sum_{\lambda \in m} v_\lambda + \text{Cste} ,$$

(Π_m : orthog. proj. operator onto S_m).

Conclusion:

→ ERM prefers models with large $\sum_{\lambda \in m} v_\lambda$ (overfitting).

Overfitting illustration



Choose the number of change-points

From $\{\hat{\mu}_{\hat{m}_D}\}_D$, choose D amounts to choose the “best model”.

Ideal penalty:

$$\begin{aligned} m^* &\in \operatorname{Argmin}_{m \in \mathcal{M}} \|\mu^* - \hat{\mu}_m\|^2 \quad (\text{oracle segmentation}) \\ &= \operatorname{Argmin}_{m \in \mathcal{M}} \left\{ \|Y - \hat{\mu}_m\|^2 + \operatorname{pen}_{\text{id}}(m) \right\}, \end{aligned}$$

with $\operatorname{pen}_{\text{id}}(m) := 2 \|\Pi_m \varepsilon\|^2 - 2 \langle (I - \Pi_m) \mu^*, \varepsilon \rangle$.

Strategy

- ① Concentration inequalities for linear and quadratic terms.
- ② Derive a tight upper bound $\operatorname{pen} \geq \operatorname{pen}_{\text{id}}$ with high probability.

Previous work:

Birgé, Massart (2001): Gaussian assump. + real valued functions.
 → cannot be extended to Hilbert framework.

Concentration of the linear term

Theorem (Linear term)

Let us consider a segmentation m and assume **(Db)** – **(Vmax)** hold true. Then for every $x > 0$ with probability at least $1 - 2e^{-x}$,

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

for every $\theta > 0$.

Concentration of the quadratic term

Theorem (Quadratic term)

Assuming **(Db)**-**(Vmax)**, and

$$\exists \kappa \geq 1, \quad 0 < \frac{M^2}{\kappa} \leq \min_i v_i \quad \mathbf{(Vmin)}$$

for every $m \in \mathcal{M}_n$, $x > 0$, and $\theta \in (0, 1]$,

$$\left| \|\Pi_m \varepsilon\|^2 - \mathbb{E} \left[\|\Pi_m \varepsilon\|^2 \right] \right| \leq \theta \mathbb{E} \left[\|\Pi_m \mu^* - \hat{\mu}_m\|^2 \right] + \theta^{-1} L(\kappa) v_{\max} x,$$

with probability at least $1 - 2e^{-x}$, where $L(\kappa)$ is a constant.

Idea of the proof:

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

Oracle inequality

Theorem

Let us assume **(Db)**-**(Vmin)**-**(Vmax)** and for every $x > 0$, define

$$\hat{m} \in \operatorname{Argmin}_m \left\{ \|Y - \hat{\mu}_m\|^2 + \operatorname{pen}(m) \right\} ,$$

where $\operatorname{pen}(m) = v_{\max} D_m \left[C_1 \ln \left(\frac{n}{D_m} \right) + C_2 \right]$ for constants $C_1, C_2 > 0$. Then, there exists an event of probability larger than $1 - 2e^{-x}$ on which

$$\|\mu^* - \hat{\mu}_{\hat{m}}\|^2 \leq \Delta_1 \inf_m \left\{ \|\mu^* - \hat{\mu}_m\|^2 + \operatorname{pen}(m) \right\} + \Delta_2 ,$$

where $\Delta_1 \geq 1$ and $\Delta_2 > 0$ is a remainder term.

Rk:

In Birgé, Massart (2001), $\operatorname{pen}(m) = \sigma^2 D_m \left[c_1 \ln \left(\frac{n}{D_m} \right) + c_2 \right]$.

Model selection procedure

Penalty:

$$\text{pen}(m) = v_{\max} D_m \left[C_1 \ln \left(\frac{n}{D_m} \right) + C_2 \right] = \text{pen}(D_m) .$$

Algorithm

- 1 For every $1 \leq D \leq D_{\max}$,

$$\hat{m}_D \in \text{Argmin}_{m, D_m=D} \left\{ \|Y - \hat{\mu}_m\|^2 \right\} ,$$

- 2 Define

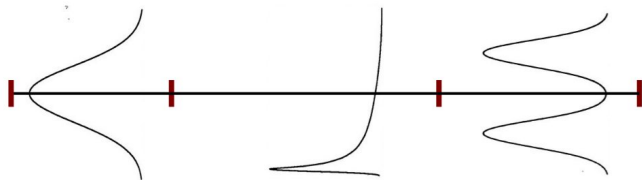
$$\hat{D} = \text{Argmin}_D \left\{ \|Y - \hat{\mu}_{\hat{m}_D}\|^2 + v_{\max} D \left[C_1 \ln \left(\frac{n}{D} \right) + C_2 \right] \right\} .$$

where C_1, C_2 : computed by simulation experiments.

- 3 Final estimator:

$$\hat{\mu}_{\hat{m}} := \hat{\mu}_{\hat{m}_{\hat{D}}} .$$

Changes in the distribution (synthetic data)

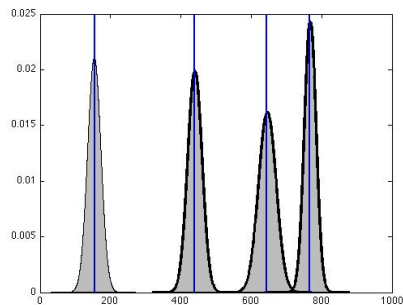
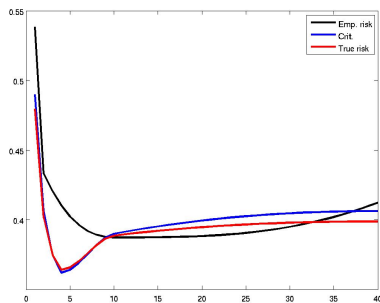


Description:

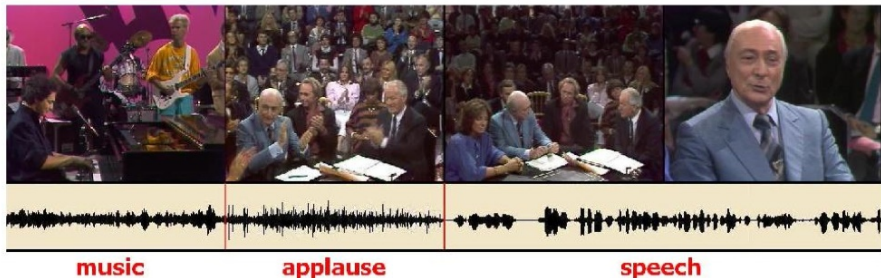
- 1 $n = 1000$, $D^* = 4$, $N_{rep} = 100$.
- 2 In each segment, observations generated according to one distribution within a pool of 10 distributions with same mean and variance.
- 3 Our kernel based approach enables to distinguish them (higher order moments)

Changes in the distribution (synthetic data) (Cont'.)

Results



“Le grand échiquier”, 70s-80s French talk show



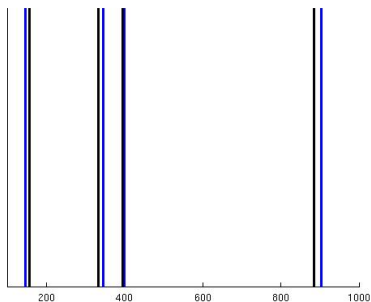
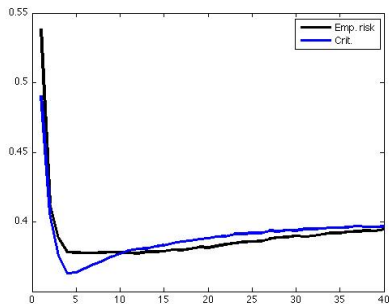
- Audio and video recordings.
- Audio: different situations can be distinguished from sound recordings (music, applause, speech, ...).
- Video: different video scenes can be distinguished by their backgrounds or specific actions of people (clapping hands, discussing, ...).

Video sequence

Description:

- $n = 10\,000$, $D^* = 4$.
- Each image summarized by a histogram with 1 024 bins.
- χ^2 kernel:
$$k_d(x, y) = \sum_{i=1}^d \frac{(x_i - y_i)^2}{x_i + y_i}.$$

Results:



Concluding remarks

Open questions:

- 1 Relax the assumption on the variance.
- 2 Use resampling strategies (heteroscedasticity).
- 3 Influence of the choice of kernel.

Preprint:

- ArXiv
- <http://www.math.univ-lille1.fr/~celisse/>

Concluding remarks

Open questions:

- 1 Relax the assumption on the variance.
- 2 Use resampling strategies (heteroscedasticity).
- 3 Influence of the choice of kernel.

Preprint:

- ArXiv
- <http://www.math.univ-lille1.fr/~celisse/>

Thank you!

Intro.
○○○○

Framework
○○○○○

Which change-points? (D known)
○○○○

How many change-points?
○○○○○

Empirical assessment
○○○○○

Sketch of proof

$$\textcircled{1} \quad \|\Pi_m \varepsilon\|^2 = \sum_{\lambda \in m} \frac{1}{n_\lambda} \left\| \sum_{i \in \lambda} \varepsilon_i \right\|_{\mathcal{H}}^2 = \sum_{\lambda \in m} T_\lambda.$$

$$\textcircled{2} \quad \left\{ \left\| \sum_{i \in \lambda} \varepsilon_i \right\|_{\mathcal{H}}^2 \right\}_{\lambda \in m} \text{ are independent r.v. .}$$

$$\textcircled{3} \quad \text{Bernstein's inequality to } \|\Pi_m \varepsilon\|^2 \quad (\star).$$

$$\textcircled{4} \quad \text{For every } q \geq 2, \text{ upper bound of } \mathbb{E} [T_\lambda^q].$$

$$\textcircled{5} \quad \text{Pinelis-Sakhanenko's inequality on } \left\| \sum_{i \in \lambda} \varepsilon_i \right\|_{\mathcal{H}}:$$

$$\forall x > 0, \quad \mathbb{P} \left[\left\| \sum_{i \in \lambda} \varepsilon_i \right\|_{\mathcal{H}} > x \right] \leq 2 \exp \left[- \frac{x^2}{2(\sigma_\lambda^2 + b_\lambda x)} \right],$$

$$\text{with } b_\lambda = 2M/3 \text{ and } \sigma_\lambda^2 = \sum_{i \in \lambda} v_i.$$

Bernstein rather than Talagrand

Talagrand's inequality

$$\|\Pi_m \varepsilon\| = \sup_{f \in B_n} \langle f, \Pi_m \varepsilon \rangle = \sup_{f \in B_n} \sum_{i=1}^n \langle f_i, (\Pi_m \varepsilon)_i \rangle_{\mathcal{H}}$$

$$\mathbb{P} \left[\|\Pi_m \varepsilon\| \leq \mathbb{E} [\|\Pi_m \varepsilon\|] + \sqrt{2\nu x} + \frac{b}{3}x \right],$$

with $\nu = \sum_{i=1}^n \sup_f \mathbb{E} (\langle f_i, (\Pi_m \varepsilon)_i \rangle_{\mathcal{H}}^2) + 16b\mathbb{E} [\|\Pi_m \varepsilon\|]$.

Bernstein's inequality

$$\sigma^2 = \sup_f \sum_{i=1}^n \mathbb{E} (\langle f_i, (\Pi_m \varepsilon)_i \rangle_{\mathcal{H}}^2) = \mathbb{E} [\|\Pi_m \varepsilon\|^2].$$