# Tomorrow's large dimensional AI:
# renewed intuitions and new mathematics?
### *Workshop MACS COMET-SCA on "Automatics and AI"*

**Romain COUILLET**

CentraleSupélec, L2S, University of ParisSaclay, France
GSTATS IDEX DataScience Chair, GIPSA-lab, University Grenoble–Alpes, France.

June 2, 2021

**Clustering setting in (not so) large** $n, p$:

**Clustering setting in (not so) large $n, p$:**

▶ GMM setting: $x_1^{(a)}, \ldots, x_{n_a}^{(a)} \sim \mathcal{N}(\mu_a, C_a)$, $a = 1, \ldots, k$

**Clustering setting in (not so) large $n, p$:**

▶ GMM setting: $x_1^{(a)}, \ldots, x_{n_a}^{(a)} \sim \mathcal{N}(\mu_a, C_a)$, $a = 1, \ldots, k$

▶ Non-trivial task:

$$\|\mu_a - \mu_b\| = O(1), \quad \text{tr}\,(C_a - C_b) = O(\sqrt{p}), \quad \text{tr}\,[(C_a - C_b)^2] = O(p)$$

# The curse of dimensionality and its consequences

**Clustering setting in (not so) large $n, p$:**

▶ GMM setting: $x_1^{(a)}, \ldots, x_{n_a}^{(a)} \sim \mathcal{N}(\mu_a, C_a)$, $a = 1, \ldots, k$

▶ Non-trivial task:

$$\|\mu_a - \mu_b\| = O(1), \quad \text{tr}\,(C_a - C_b) = O(\sqrt{p}), \quad \text{tr}\,[(C_a - C_b)^2] = O(p)$$

**Classical method: spectral clustering**

# The curse of dimensionality and its consequences

**Clustering setting in (not so) large $n, p$:**

▶ GMM setting: $x_1^{(a)}, \ldots, x_{n_a}^{(a)} \sim \mathcal{N}(\mu_a, C_a)$, $a = 1, \ldots, k$

▶ Non-trivial task:

$$\|\mu_a - \mu_b\| = O(1), \quad \mathsf{tr}\,(C_a - C_b) = O(\sqrt{p}), \quad \mathsf{tr}\,[(C_a - C_b)^2] = O(p)$$

**Classical method: spectral clustering**

▶ Extract and cluster the dominant eigenvectors of

$$K = \{\kappa(x_i, x_j)\}_{i,j=1}^n$$

# The curse of dimensionality and its consequences

**Clustering setting in (not so) large $n, p$:**

▶ GMM setting: $x_1^{(a)}, \ldots, x_{n_a}^{(a)} \sim \mathcal{N}(\mu_a, C_a)$, $a = 1, \ldots, k$

▶ Non-trivial task:

$$\|\mu_a - \mu_b\| = O(1), \quad \operatorname{tr}(C_a - C_b) = O(\sqrt{p}), \quad \operatorname{tr}[(C_a - C_b)^2] = O(p)$$

**Classical method: spectral clustering**

▶ Extract and cluster the dominant eigenvectors of

$$K = \{\kappa(x_i, x_j)\}_{i,j=1}^n, \quad \kappa(x_i, x_j) = f\left(\frac{1}{p}\|x_i - x_j\|^2\right).$$

**Clustering setting in (not so) large $n, p$:**

▶ GMM setting: $x_1^{(a)}, \ldots, x_{n_a}^{(a)} \sim \mathcal{N}(\mu_a, C_a)$, $a = 1, \ldots, k$

▶ Non-trivial task:

$$\|\mu_a - \mu_b\| = O(1), \quad \mathsf{tr}\,(C_a - C_b) = O(\sqrt{p}), \quad \mathsf{tr}\,[(C_a - C_b)^2] = O(p)$$

**Classical method: spectral clustering**

▶ Extract and cluster the dominant eigenvectors of

$$K = \{\kappa(x_i, x_j)\}_{i,j=1}^n, \quad \kappa(x_i, x_j) = f\left(\frac{1}{p}\|x_i - x_j\|^2\right).$$

▶ Why? **Finite-dimensional intuition**

$$K = \begin{pmatrix} \begin{array}{|c|c|c|} \hline \substack{\kappa(x_i,x_j) \\ \gg 1} & \substack{\kappa(x_i,x_j) \\ \ll 1} & \substack{\kappa(x_i,x_j) \\ \ll 1} \\ \hline \substack{\kappa(x_i,x_j) \\ \ll 1} & \substack{\kappa(x_i,x_j) \\ \gg 1} & \substack{\kappa(x_i,x_j) \\ \ll 1} \\ \hline \substack{\kappa(x_i,x_j) \\ \ll 1} & \substack{\kappa(x_i,x_j) \\ \ll 1} & \substack{\kappa(x_i,x_j) \\ \gg 1} \\ \hline \end{array} \end{pmatrix} \begin{array}{l} \updownarrow \mathcal{C}_1 \\ \updownarrow \mathcal{C}_2 \\ \updownarrow \mathcal{C}_3 \end{array}$$
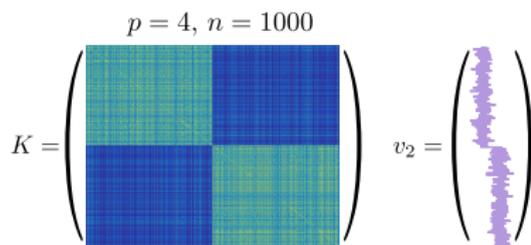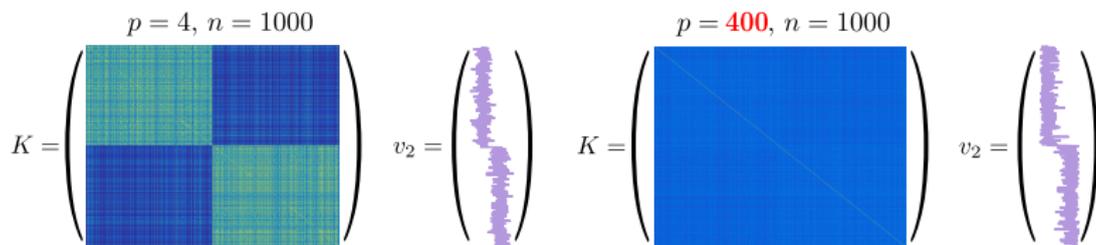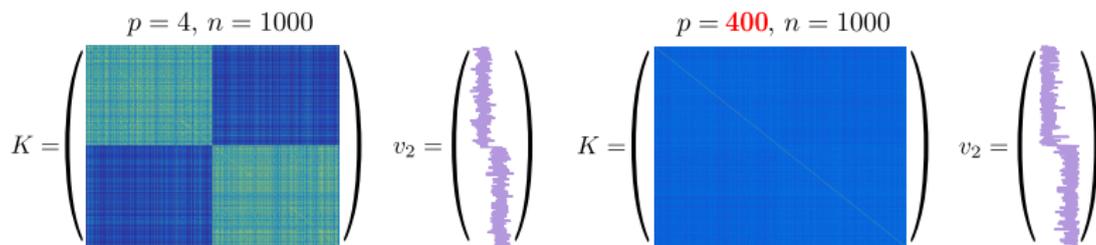
**In reality, here is what happens...**

Kernel $K_{ij} = \exp(-\frac{1}{2p}\|x_i - x_j\|^2)$ and second eigenvector $v_2$

$(x_i \sim \mathcal{N}(\pm\mu, I_p),\ \mu = (2, 0, \ldots, 0)^\mathsf{T} \in \mathbb{R}^p)$.

**In reality, here is what happens...**

Kernel $K_{ij} = \exp(-\frac{1}{2p}\|x_i - x_j\|^2)$ and second eigenvector $v_2$
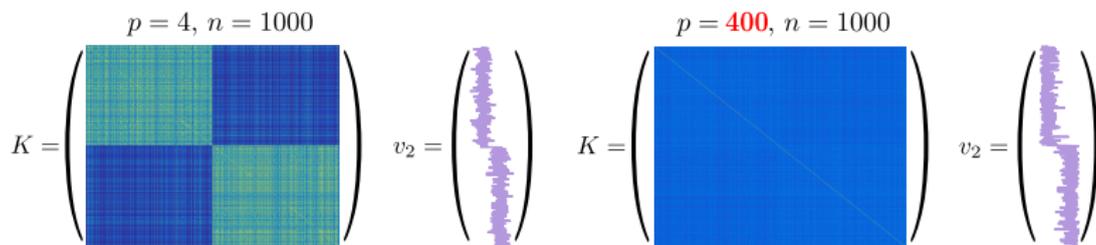($x_i \sim \mathcal{N}(\pm\mu, I_p)$, $\mu = (2, 0, \ldots, 0)^{\mathsf{T}} \in \mathbb{R}^p$).

$p = 4$, $n = 1000$



$$K = \begin{pmatrix} & \\ & \end{pmatrix} \quad v_2 = \begin{pmatrix} \\ \end{pmatrix}$$

**In reality, here is what happens...**

Kernel $K_{ij} = \exp(-\frac{1}{2p}\|x_i - x_j\|^2)$ and second eigenvector $v_2$
($x_i \sim \mathcal{N}(\pm\mu, I_p)$, $\mu = (2, 0, \ldots, 0)^\mathsf{T} \in \mathbb{R}^p$).



$p = 4$, $n = 1000$

$p = 400$, $n = 1000$

**In reality, here is what happens...**

Kernel $K_{ij} = \exp(-\frac{1}{2p}\|x_i - x_j\|^2)$ and second eigenvector $v_2$
($x_i \sim \mathcal{N}(\pm\mu, I_p)$, $\mu = (2, 0, \ldots, 0)^\mathsf{T} \in \mathbb{R}^p$).



$p = 4$, $n = 1000$      $p = 400$, $n = 1000$

**Key observation**: Under growth rate assumptions,

$$\max_{1 \le i \ne j \le n} \left\{ \left| \frac{1}{p} \|x_i - x_j\|^2 - \tau \right| \right\} \xrightarrow{\text{a.s.}} 0, \quad \tau = \frac{2}{p} \sum_{i=1}^{k} \operatorname{tr} \frac{n_a}{n} C_a.$$
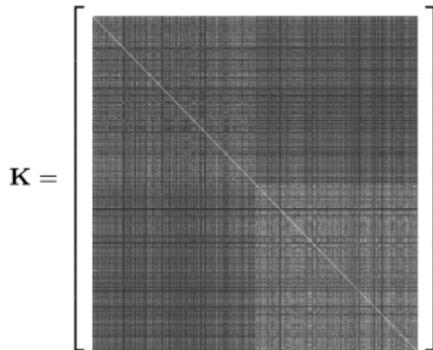
**In reality, here is what happens...**

Kernel $K_{ij} = \exp(-\frac{1}{2p}\|x_i - x_j\|^2)$ and second eigenvector $v_2$

$(x_i \sim \mathcal{N}(\pm\mu, I_p),\ \mu = (2, 0, \ldots, 0)^\mathsf{T} \in \mathbb{R}^p).$



$p = 4,\ n = 1000$     $p = 400,\ n = 1000$

**Key observation**: Under growth rate assumptions,

$$\max_{1 \le i \ne j \le n} \left\{ \left| \frac{1}{p}\|x_i - x_j\|^2 - \tau \right| \right\} \xrightarrow{\text{a.s.}} 0, \quad \tau = \frac{2}{p} \sum_{i=1}^{k} \text{tr}\, \frac{n_a}{n} C_a.$$

▶ this suggests $K \simeq f(\tau)1_n 1_n^\mathsf{T}$!

**MNIST**
raw
$p = 784$, $n = 500$
↓

**ImageNet**
VGG-features
$p = 3084$, $n = 500$
↓

**20NewsGroup**
BERT embedding
$p = 300$, $n = 500$
↓

$\mathbf{K} =$

$\mathbf{v}_2 =$

(ici, classes "5" et "0")

(ici, classes "bird" et "plane")

(ici, classes "sports" et "sales")

**(Major) consequences**:

- ▶ Most **machine learning intuitions collapse**

**(Major) consequences**:

▶ Most **machine learning intuitions collapse**

▶ **But luckily**, concentration of distances allows for Taylor expansion, linearization...

**(Major) consequences**:
- ▶ Most **machine learning intuitions collapse**
- ▶ **But luckily**, concentration of distances allows for Taylor expansion, linearization...

Theorem (**[C-Benaych'16]** Asymptotic Kernel Behavior)

*Under growth rate assumptions, as $p, n \to \infty$,*

$$\left\| K - \hat{K} \right\| \xrightarrow{\text{a.s.}} 0, \quad \hat{K} \simeq \underbrace{f(\tau) 1_n 1_n^{\mathsf{T}}}_{O_{\|\cdot\|}(n)}$$

**(Major) consequences**:

- ▶ Most **machine learning intuitions collapse**
- ▶ **But luckily**, concentration of distances allows for Taylor expansion, linearization...

**Theorem ([C-Benaych'16]** Asymptotic Kernel Behavior)

*Under growth rate assumptions, as $p, n \to \infty$,*

$$\left\| K - \hat{K} \right\| \xrightarrow{\text{a.s.}} 0, \quad \hat{K} \simeq \underbrace{f(\boldsymbol{\tau}) 1_n 1_n^\mathsf{T}}_{O_{\|\cdot\|}(n)} + f'(\boldsymbol{\tau}) \frac{1}{p} Z Z^\mathsf{T} + J A J^\mathsf{T} + *$$

**(Major) consequences**:

▶ Most **machine learning intuitions collapse**

▶ **But luckily**, concentration of distances allows for Taylor expansion, linearization...

### Theorem ([C-Benaych'16] Asymptotic Kernel Behavior)

*Under growth rate assumptions, as $p, n \to \infty$,*

$$\left\| K - \hat{K} \right\| \xrightarrow{\text{a.s.}} 0, \quad \hat{K} \simeq \underbrace{f(\boldsymbol{\tau}) 1_n 1_n^{\mathsf{T}}}_{O_{\|\cdot\|}(n)} + f'(\boldsymbol{\tau}) \frac{1}{p} Z Z^{\mathsf{T}} + J A J^{\mathsf{T}} + *$$
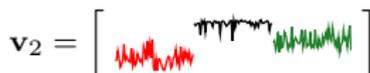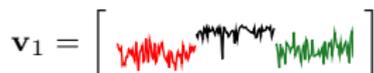
*with $J = [j_1, \ldots, j_k] \in \mathbb{R}^{n \times k}$, $j_a = (0, 1_{n_a}, 0)^{\mathsf{T}}$ (the clusters!)*

## The curse of dimensionality and its consequences (4)

**(Major) consequences**:

▶ Most **machine learning intuitions collapse**

▶ **But luckily**, concentration of distances allows for Taylor expansion, linearization…

Theorem (**[C-Benaych'16]** Asymptotic Kernel Behavior)

*Under growth rate assumptions, as $p, n \to \infty$,*

$$\left\| K - \hat{K} \right\| \xrightarrow{\text{a.s.}} 0, \quad \hat{K} \simeq \underbrace{f(\tau) 1_n 1_n^\mathsf{T}}_{O_{\|\cdot\|}(n)} + f'(\tau) \frac{1}{p} ZZ^\mathsf{T} + JAJ^\mathsf{T} + *$$
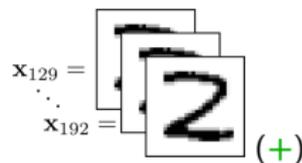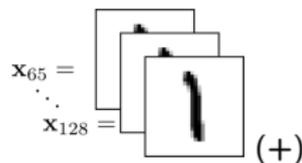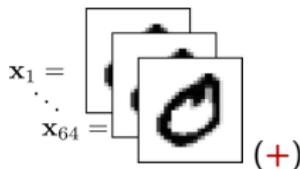
*with $J = [j_1, \ldots, j_k] \in \mathbb{R}^{n \times k}$, $j_a = (0, 1_{n_a}, 0)^\mathsf{T}$ (the clusters!) and $A \in \mathbb{R}^{k \times k}$ function of:*

▶ $f(\tau)$, $f'(\tau)$, $f''(\tau)$

▶ $\|\mu_a - \mu_b\|$, $tr(C_a - C_b)$, $tr((C_a - C_b)^2)$, for $a, b \in \{1, \ldots, k\}$.

**(Major) consequences**:

▶ Most **machine learning intuitions collapse**
▶ **But luckily**, concentration of distances allows for Taylor expansion, linearization...

Theorem (**[C-Benaych'16]** Asymptotic Kernel Behavior)

*Under growth rate assumptions, as $p, n \to \infty$,*

$$\left\| K - \hat{K} \right\| \xrightarrow{\text{a.s.}} 0, \quad \hat{K} \simeq \underbrace{f(\boldsymbol{\tau}) 1_n 1_n^{\mathsf{T}}}_{O_{\|\cdot\|}(n)} + f'(\boldsymbol{\tau}) \frac{1}{p} ZZ^{\mathsf{T}} + JAJ^{\mathsf{T}} + *$$

*with $J = [j_1, \ldots, j_k] \in \mathbb{R}^{n \times k}$, $j_a = (0, 1_{n_a}, 0)^{\mathsf{T}}$ (the clusters!) and $A \in \mathbb{R}^{k \times k}$ function of:*

▶ $f(\boldsymbol{\tau})$, $f'(\boldsymbol{\tau})$, $f''(\boldsymbol{\tau})$
▶ $\|\mu_a - \mu_b\|$, $\text{tr}(C_a - C_b)$, $\text{tr}((C_a - C_b)^2)$, for $a, b \in \{1, \ldots, k\}$.
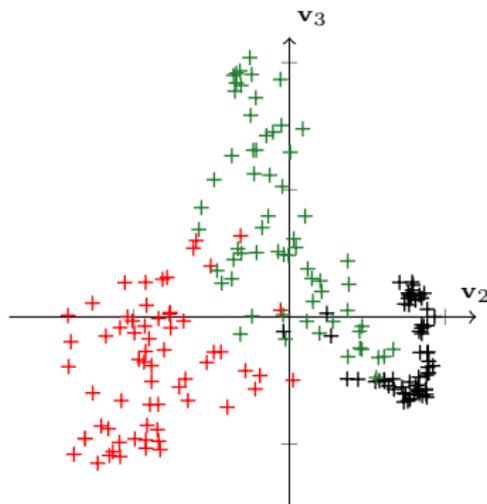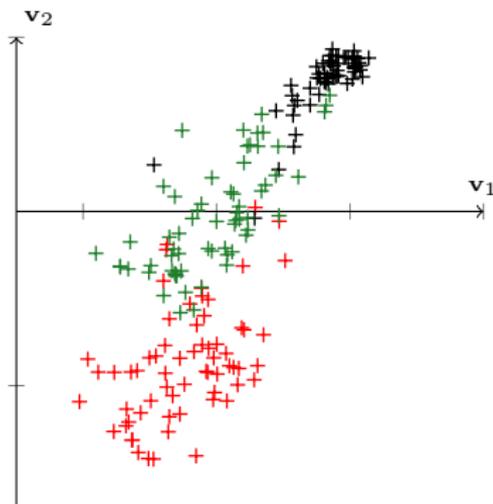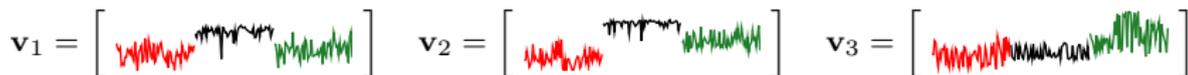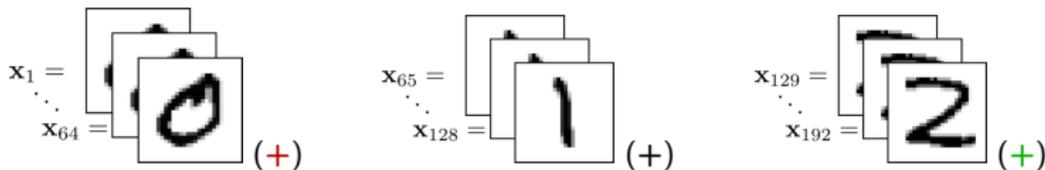
↪ **This is a spiked model! We can study it fully!**

- **Asymptotic analysis of eigenvectors of** $K$: (MNIST, $p = 28 \times 28 (= 784)$)



$\mathbf{x}_1 = $
$\cdots$
$\mathbf{x}_{64} = $ (+)

$\mathbf{x}_{65} = $
$\cdots$
$\mathbf{x}_{128} = $ (+)

$\mathbf{x}_{129} = $
$\cdots$
$\mathbf{x}_{192} = $ (+)

$$\mathbf{v}_1 = \left[\begin{array}{c} \end{array}\right] \qquad \mathbf{v}_2 = \left[\begin{array}{c} \end{array}\right] \qquad \mathbf{v}_3 = \left[\begin{array}{c} \end{array}\right]$$
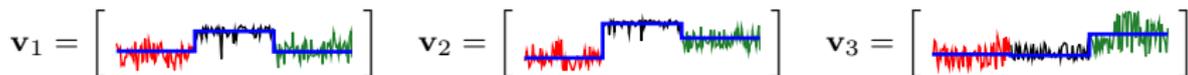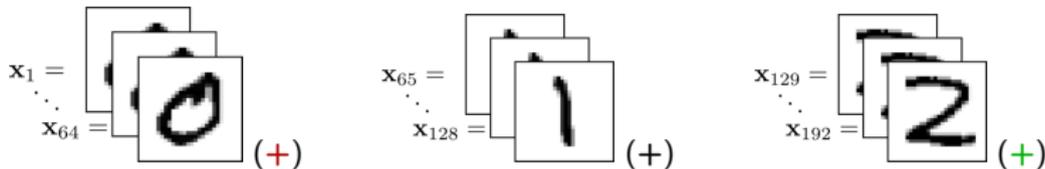
# Performance prediction: spectral clustering

- **Asymptotic analysis of eigenvectors of** $K$: (MNIST, $p = 28 \times 28 (= 784)$)

- **Asymptotic analysis of eigenvectors of $K$:** (MNIST, $p = 28 \times 28(= 784)$)



$\mathbf{x}_1 =$
$\mathbf{x}_{64} =$
$(+)$

$\mathbf{x}_{65} =$
$\mathbf{x}_{128} =$
$(+)$

$\mathbf{x}_{129} =$
$\mathbf{x}_{192} =$
$(+)$

$\mathbf{v}_1 = \begin{bmatrix} \phantom{xx} \end{bmatrix}$
$\mathbf{v}_2 = \begin{bmatrix} \phantom{xx} \end{bmatrix}$
$\mathbf{v}_3 = \begin{bmatrix} \phantom{xx} \end{bmatrix}$

theoretical prediction

**Takeaway Message 2**

"RMT Reassesses and Improves Data Processing"

**Today's menu:**

**1.** New Counter-Intuitive Kernels

**2.** Resurrecting Semi-Supervised Learning

**3.** Making complex ML frameworks simple: Multitask Learning

**4.** Towards cheap "environment-friendly" learning

**5.** Using Random Matrices to Study... Random Tensors!
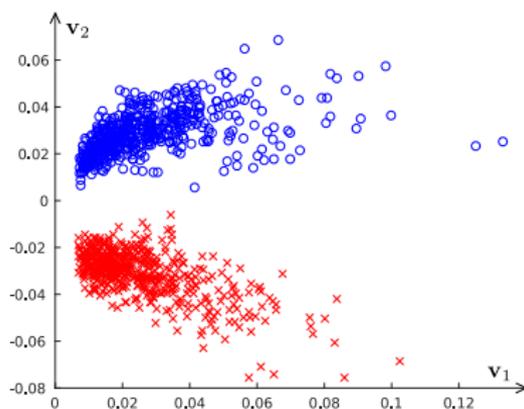
**1.** New Counter-Intuitive Kernels

- **Going further than ([Kammoun,Couillet'17])**,

$$K \simeq \underbrace{f(\tau)1_n1_n^{\mathsf{T}}}_{O_{\|\cdot\|}(n)} + f'(\tau)\frac{1}{p}ZZ^{\mathsf{T}} + JAJ^{\mathsf{T}}, \text{ avec } A = F\left(\begin{array}{c} f(\tau), f'(\tau), f''(\tau) \\ \|\mu_a - \mu_b\|, \mathrm{tr}\,(C_a - C_b), \ldots \end{array}\right).$$

- **Going further than ([Kammoun,Couillet'17])**, if $\underline{f'(\tau) = 0}$,

$$K \simeq \underbrace{f(\tau)1_n1_n^\mathsf{T}}_{O_{\|\cdot\|}(n)} + f'(\tau)\frac{1}{p}ZZ^\mathsf{T} + JAJ^\mathsf{T}, \text{ avec } A = F \left( \begin{array}{c} f(\tau), f'(\tau), f''(\tau) \\ \|\mu_a - \mu_b\|, \operatorname{tr}(C_a - C_b), \dots \end{array} \right).$$

# Improving Kernel Spectral Clustering

- **Going further than ([Kammoun,Couillet'17])**, if $\underline{f'(\tau) = 0}$,

$$K \simeq \underbrace{f(\tau)1_n 1_n^\mathsf{T}}_{O_{\|\cdot\|}(n)} + \cancel{f'(\tau)\frac{1}{p}ZZ^\mathsf{T}} + JAJ^\mathsf{T}, \text{ avec } A = F\left(\begin{array}{c} f(\tau), f'(\tau), f''(\tau) \\ \cancel{\|\mu_a - \mu_b\|}, \operatorname{tr}(C_a - C_b), \dots \end{array}\right).$$
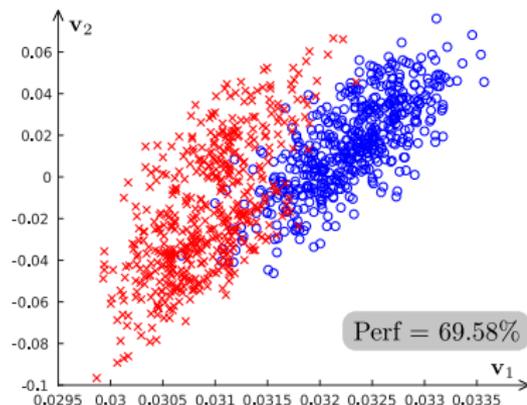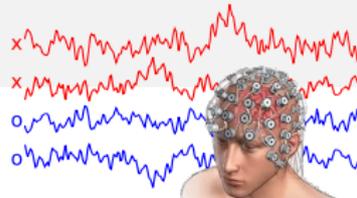
- **Gaussian case**: $\mathcal{N}(0, \mathbf{C}_1)$ vs. $\mathcal{N}(0, \mathbf{C}_2)$



Kernel $K_{ij} = \exp(-\frac{1}{2p}\|x_i - x_j\|^2)$



Kernel $K_{ij} = (\frac{1}{p}\|x_i - x_j\|^2 - \tau)^2$

- **EEG data**: sane vs. epileptic patients



Kernel $K_{ij} = \exp(-\frac{1}{2p}\|x_i - x_j\|^2)$

Kernel $K_{ij} = (\frac{1}{p}\|x_i - x_j\|^2 - \tau)^2$

- **EEG data**: sane vs. epileptic patients



Kernel $K_{ij} = \exp(-\frac{1}{2p}\|x_i - x_j\|^2)$

Kernel $K_{ij} = (\frac{1}{p}\|x_i - x_j\|^2 - \tau)^2$

$\rightarrow$ **_Remark_**: _highly counter-intuitive kernel!_