Can Random Matrices Change the Future of Machine Learning? Digicosme "Junior Conference on Wireless and Optical Communications 2020"

Romain COUILLET

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

December 3, 2020



























Basics of Random Matrix Theory Motivation: Large Sample Covariance Matrices Spiked Models

Outline

Basics of Random Matrix Theory

Motivation: Large Sample Covariance Matrices Spiked Models

Basics of Random Matrix Theory Motivation: Large Sample Covariance Matrices Spiked Models

Baseline scenario: $y_1, \ldots, y_n \in \mathbb{C}^p$ (or \mathbb{R}^p) i.i.d. with $E[y_1] = 0$, $E[y_1y_1^*] = C_p$:

Baseline scenario: $y_1, \ldots, y_n \in \mathbb{C}^p$ (or \mathbb{R}^p) i.i.d. with $E[y_1] = 0$, $E[y_1y_1^*] = C_p$: If $y_1 \sim \mathcal{N}(0, C_p)$, ML estimator for C_p is the sample covariance matrix (SCM)

$$\hat{C}_p = \frac{1}{n} Y_p Y_p^* = \frac{1}{n} \sum_{i=1}^n y_i y_i^*$$

 $(Y_p = [y_1, \ldots, y_n] \in \mathbb{C}^{p \times n}).$

Baseline scenario: $y_1, \ldots, y_n \in \mathbb{C}^p$ (or \mathbb{R}^p) i.i.d. with $E[y_1] = 0$, $E[y_1y_1^*] = C_p$: If $y_1 \sim \mathcal{N}(0, C_p)$, ML estimator for C_p is the sample covariance matrix (SCM)

$$\hat{C}_p = \frac{1}{n} Y_p Y_p^* = \frac{1}{n} \sum_{i=1}^n y_i y_i^*$$

 $(Y_p = [y_1, \dots, y_n] \in \mathbb{C}^{p \times n}).$ If $n \to \infty$, then, strong law of large numbers

$$\hat{C}_p \xrightarrow{\text{a.s.}} C_p.$$

or equivalently, in spectral norm

$$\left\| \hat{C}_p - C_p \right\| \xrightarrow{\text{a.s.}} 0.$$

Baseline scenario: $y_1, \ldots, y_n \in \mathbb{C}^p$ (or \mathbb{R}^p) i.i.d. with $E[y_1] = 0$, $E[y_1y_1^*] = C_p$: If $y_1 \sim \mathcal{N}(0, C_p)$, ML estimator for C_p is the sample covariance matrix (SCM)

$$\hat{C}_p = \frac{1}{n} Y_p Y_p^* = \frac{1}{n} \sum_{i=1}^n y_i y_i^*$$

 $(Y_p = [y_1, \dots, y_n] \in \mathbb{C}^{p \times n}).$ If $n \to \infty$, then, strong law of large numbers

$$\hat{C}_p \xrightarrow{\text{a.s.}} C_p.$$

or equivalently, in spectral norm

$$\left\| \hat{C}_p - C_p \right\| \xrightarrow{\text{a.s.}} 0.$$

Random Matrix Regime

▶ No longer valid if $p, n \to \infty$ with $p/n \to c \in (0, \infty)$,

$$\left\| \hat{C}_p - C_p \right\| \not\to 0.$$

Baseline scenario: $y_1, \ldots, y_n \in \mathbb{C}^p$ (or \mathbb{R}^p) i.i.d. with $E[y_1] = 0$, $E[y_1y_1^*] = C_p$: If $y_1 \sim \mathcal{N}(0, C_p)$, ML estimator for C_p is the sample covariance matrix (SCM)

$$\hat{C}_p = \frac{1}{n} Y_p Y_p^* = \frac{1}{n} \sum_{i=1}^n y_i y_i^*$$

 $(Y_p = [y_1, \dots, y_n] \in \mathbb{C}^{p \times n}).$ If $n \to \infty$, then, strong law of large numbers

$$\hat{C}_p \xrightarrow{\text{a.s.}} C_p.$$

or equivalently, in spectral norm

$$\left\| \hat{C}_p - C_p \right\| \xrightarrow{\text{a.s.}} 0.$$

Random Matrix Regime

 $\blacktriangleright \text{ No longer valid if } p,n \to \infty \text{ with } p/n \to c \in (0,\infty),$

$$\left\| \hat{C}_p - C_p \right\| \not\to 0.$$

For practical p, n with $p \simeq n$, leads to dramatically wrong conclusions

Baseline scenario: $y_1, \ldots, y_n \in \mathbb{C}^p$ (or \mathbb{R}^p) i.i.d. with $E[y_1] = 0$, $E[y_1y_1^*] = C_p$: If $y_1 \sim \mathcal{N}(0, C_p)$, ML estimator for C_p is the sample covariance matrix (SCM)

$$\hat{C}_p = \frac{1}{n} Y_p Y_p^* = \frac{1}{n} \sum_{i=1}^n y_i y_i^*$$

 $(Y_p = [y_1, \dots, y_n] \in \mathbb{C}^{p \times n}).$ If $n \to \infty$, then, strong law of large numbers

$$\hat{C}_p \xrightarrow{\text{a.s.}} C_p.$$

or equivalently, in spectral norm

$$\left\| \hat{C}_p - C_p \right\| \xrightarrow{\text{a.s.}} 0.$$

Random Matrix Regime

 $\blacktriangleright \ \, \text{No longer valid if } p,n \to \infty \text{ with } p/n \to c \in (0,\infty),$

$$\left\| \hat{C}_p - C_p \right\| \not\to 0.$$

For practical p, n with p ≃ n, leads to dramatically wrong conclusions
Even for p = n/100.



Figure: Histogram of the eigenvalues of \hat{C}_p for c = 1/4, $C_p = I_p$.



Figure: Histogram of the eigenvalues of \hat{C}_p for c = 1/4, $C_p = I_p$.



Figure: Histogram of the eigenvalues of \hat{C}_p for c = 1/4, $C_p = I_p$.



Figure: Histogram of the eigenvalues of \hat{C}_p for c = 1/4, $C_p = I_p$.



Figure: Histogram of the eigenvalues of \hat{C}_p for c = 1/4, $C_p = I_p$.



Figure: Histogram of the eigenvalues of \hat{C}_p for c = 1/4, $C_p = I_p$.

Definition (Empirical Spectral Density)

Empirical spectral density (e.s.d.) μ_p of Hermitian matrix $A_p \in \mathbb{C}^{p imes p}$ is

$$\mu_p = \frac{1}{p} \sum_{i=1}^p \delta_{\lambda_i(A_p)}.$$

Definition (Empirical Spectral Density)

Empirical spectral density (e.s.d.) μ_p of Hermitian matrix $A_p \in \mathbb{C}^{p imes p}$ is

$$\mu_p = \frac{1}{p} \sum_{i=1}^p \delta_{\lambda_i(A_p)}.$$

Theorem (Marčenko–Pastur Law [Marčenko,Pastur'67]) $X_p \in \mathbb{C}^{p \times n}$ with i.i.d. zero mean, unit variance entries. As $p, n \to \infty$ with $p/n \to c \in (0, \infty)$, e.s.d. μ_p of $\frac{1}{n}X_pX_p^*$ satisfies

$$\mu_p \xrightarrow{\text{a.s.}} \mu_c$$

weakly, where

•
$$\mu_c(\{0\}) = \max\{0, 1 - c^{-1}\}$$

Definition (Empirical Spectral Density)

Empirical spectral density (e.s.d.) μ_p of Hermitian matrix $A_p \in \mathbb{C}^{p imes p}$ is

$$\mu_p = \frac{1}{p} \sum_{i=1}^p \delta_{\lambda_i(A_p)}$$

Theorem (Marčenko–Pastur Law [Marčenko,Pastur'67]) $X_p \in \mathbb{C}^{p \times n}$ with i.i.d. zero mean, unit variance entries. As $p, n \to \infty$ with $p/n \to c \in (0, \infty)$, e.s.d. μ_p of $\frac{1}{n}X_pX_p^*$ satisfies

$$\mu_p \xrightarrow{\mathrm{a.s.}} \mu_c$$

weakly, where

•
$$\mu_c(\{0\}) = \max\{0, 1 - c^{-1}\}$$

• on $(0,\infty)$, μ_c has continuous density f_c supported on $[(1-\sqrt{c})^2,(1+\sqrt{c})^2]$

$$f_c(x) = \frac{1}{2\pi cx} \sqrt{(x - (1 - \sqrt{c})^2)((1 + \sqrt{c})^2 - x)}.$$



Figure: Marčenko-Pastur law for different limit ratios $c = \lim_{p \to \infty} p/n$.



Figure: Marčenko-Pastur law for different limit ratios $c = \lim_{p \to \infty} p/n$.



Figure: Marčenko-Pastur law for different limit ratios $c = \lim_{p \to \infty} p/n$.

Outline

Basics of Random Matrix Theory Motivation: Large Sample Covariance Matrices Spiked Models

Spiked Models

Small rank perturbation: $C_p = I_p + P$, P of low rank.



Spiked Models

Small rank perturbation: $C_p = I_p + P$, P of low rank.


Small rank perturbation: $C_p = I_p + P$, P of low rank.



Small rank perturbation: $C_p = I_p + P$, P of low rank.



Theorem (Eigenvalues [Baik,Silverstein'06]) Let $Y_p = C_p^{\frac{1}{2}} X_p$, with $\searrow X_p$ with i.i.d. zero mean, unit variance, $E[|X_p|_{ij}^4] < \infty$.

•
$$C_p = I_p + P$$
, $P = U\Omega U^*$, where, for K fixed,

$$\Omega = \operatorname{diag} \left(\omega_1, \dots, \omega_K \right) \in \mathbb{R}^{K \times K}, \text{ with } \omega_1 \ge \dots \ge \omega_K > 0.$$

Theorem (Eigenvalues [Baik,Silverstein'06]) Let $Y_p = C_p^{\frac{1}{2}} X_p$, with X_p with i.i.d. zero mean, unit variance, $E[|X_p|_{ij}^4] < \infty$. $C_p = I_p + P$, $P = U\Omega U^*$, where, for K fixed, $\Omega = \text{diag}(\omega_1, \dots, \omega_K) \in \mathbb{R}^{K \times K}$, with $\omega_1 \ge \dots \ge \omega_K > 0$.

 $\textit{Then, as } p,n \to \infty, \ p/n \to c \in (0,\infty), \textit{ denoting } \lambda_m = \lambda_m (\tfrac{1}{n} Y_p Y_p^*) \ (\lambda_m > \lambda_{m+1}),$

$$\lambda_m \xrightarrow{\text{a.s.}} \begin{cases} 1 + \omega_m + c \frac{1 + \omega_m}{\omega_m} > (1 + \sqrt{c})^2 &, \ \omega_m > \sqrt{c} \\ (1 + \sqrt{c})^2 &, \ \omega_m \in (0, \sqrt{c}]. \end{cases}$$

Theorem (Eigenvectors [Paul'07]) Let $Y_p = C_p^{\frac{1}{2}} X_p$, with

▶ X_p with i.i.d. zero mean, unit variance, $E[|X_p|_{ij}^4] < \infty$.

•
$$C_p = I_p + P$$
, $P = U\Omega U^* = \sum_{i=1}^K \omega_i u_i u_i^*$, $\omega_1 > \ldots > \omega_M > 0$.

Theorem (Eigenvectors [Paul'07]) Let $Y_p = C_p^{\frac{1}{2}} X_p$, with $\searrow X_p$ with *i.i.d.* zero mean, unit variance, $E[|X_p|_{ij}^4] < \infty$. $\searrow C_p = I_p + P$, $P = U\Omega U^* = \sum_{i=1}^K \omega_i u_i u_i^*$, $\omega_1 > \ldots > \omega_M > 0$.

Then, as $p, n \to \infty$, $p/n \to c \in (0, \infty)$, for $a, b \in \mathbb{C}^p$ deterministic and \hat{u}_i eigenvector of $\lambda_i(\frac{1}{n}Y_pY_p^*)$,

$$a^*\hat{u}_i\hat{u}_i^*b - \frac{1 - c\omega_i^{-2}}{1 + c\omega_i^{-1}}a^*u_iu_i^*b \cdot \mathbf{1}_{\omega_i > \sqrt{c}} \xrightarrow{\text{a.s.}} 0$$

In particular,

$$|\hat{u}_i^* u_i|^2 \xrightarrow{\text{a.s.}} \frac{1 - c\omega_i^{-2}}{1 + c\omega_i^{-1}} \cdot 1_{\omega_i > \sqrt{c}}.$$



Population spike ω_1

Figure: Simulated versus limiting $|\hat{u}_1^{\mathsf{T}}u_1|^2$ for $Y_p = C_p^{\frac{1}{2}}X_p$, $C_p = I_p + \omega_1 u_1 u_1^{\mathsf{T}}$, p/n = 1/3, varying ω_1 .



Population spike ω_1

Figure: Simulated versus limiting $|\hat{u}_1^{\mathsf{T}}u_1|^2$ for $Y_p = C_p^{\frac{1}{2}}X_p$, $C_p = I_p + \omega_1 u_1 u_1^{\mathsf{T}}$, p/n = 1/3, varying ω_1 .



Population spike ω_1

Figure: Simulated versus limiting $|\hat{u}_1^{\mathsf{T}}u_1|^2$ for $Y_p = C_p^{\frac{1}{2}}X_p$, $C_p = I_p + \omega_1 u_1 u_1^{\mathsf{T}}$, p/n = 1/3, varying ω_1 .



Population spike ω_1

Figure: Simulated versus limiting $|\hat{u}_1^{\mathsf{T}}u_1|^2$ for $Y_p = C_p^{\frac{1}{2}}X_p$, $C_p = I_p + \omega_1 u_1 u_1^{\mathsf{T}}$, p/n = 1/3, varying ω_1 .

Similar results for multiple matrix models:

▶
$$Y_p = \frac{1}{n}(I+P)^{\frac{1}{2}}X_pX_p^*(I+P)^{\frac{1}{2}}$$

▶ $Y_p = \frac{1}{n}X_pX_p^* + P$
▶ $Y_p = \frac{1}{n}X_p^*(I+P)X$
▶ $Y_p = \frac{1}{n}(X_p+P)^*(X_p+P)$
▶ etc.

Basics of Random Matrix Theory Motivation: Large Sample Covariance Matrices Spiked Models

Application to Machine Learning

Takeaway Message 1

"RMT Explains Why Machine Learning Intuitions Collapse in Large Dimensions"

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:

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

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

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

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

Classical method: spectral clustering

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:

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

Classical method: spectral clustering

Extract and cluster the dominant eigenvectors of

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

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:

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

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:

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

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

$$K = \begin{pmatrix} \kappa(x,x) & \kappa(x,x) \\ \gg 1 & \ll 1 & \ll 1 \\ \kappa(x,x) & \kappa(x,x) & \kappa(x,x) \\ \approx 1 & \gg 1 & \ll 1 \\ \hline \kappa(x,x) & \kappa(x,x) & \kappa(x,x) \\ \kappa(x,x) & \kappa($$

In reality, here is what happens...

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

In reality, here is what happens...

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



In reality, here is what happens...

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



In reality, here is what happens...

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



Key observation: Under growth rate assumptions,

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

In reality, here is what happens...

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



Key observation: Under growth rate assumptions,

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

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



(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 $p, n \rightarrow \infty$,

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

(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 $p, n \rightarrow \infty$,

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

(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 $p, n \rightarrow \infty$,

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

with $J = [j_1, \dots, j_k] \in \mathbb{R}^{n \times k}$, $j_a = (0, 1_{n_a}, 0)^{\mathsf{T}}$ (the clusters!)

(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 $p, n \rightarrow \infty$,

$$\left\| K - \hat{K} \right\| \xrightarrow{\text{a.s.}} 0, \quad \hat{K} \simeq \underbrace{f(\tau) \mathbf{1}_n \mathbf{1}_n^{\mathsf{T}}}_{\mathcal{O}_{\|\cdot\|}(n)} + \frac{1}{p} Z Z^{\mathsf{T}} + J A J^{\mathsf{T}} + \ast$$

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:

- $\blacktriangleright f(\tau), f'(\tau), f''(\tau)$
- $\|\mu_a \mu_b\|$, $tr(C_a C_b)$, $tr((C_a C_b)^2)$, for $a, b \in \{1, \dots, k\}$.

(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 $p, n \rightarrow \infty$,

$$\left\| K - \hat{K} \right\| \xrightarrow{\text{a.s.}} 0, \quad \hat{K} \simeq \underbrace{f(\tau) \mathbf{1}_n \mathbf{1}_n^{\mathsf{T}}}_{\mathcal{O}_{\|\cdot\|}(n)} + \frac{1}{p} Z Z^{\mathsf{T}} + J A J^{\mathsf{T}} + \ast$$

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:

► $f(\tau), f'(\tau), f''(\tau)$ ► $\|\mu_a - \mu_b\|, tr(C_a - C_b), tr((C_a - C_b)^2), \text{ for } a, b \in \{1, ..., k\}.$

→ This is a spiked model! We can study it fully!

Performance prediction: spectral clustering

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



$$\mathbf{v}_1 = \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix} \quad \mathbf{v}_2 = \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \mathbf{w}_3 \end{bmatrix} \quad \mathbf{v}_3 = \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \mathbf{w}_3 \end{bmatrix} \begin{bmatrix} \mathbf{w}_2 \\ \mathbf{w}_3 \\ \mathbf{w}_4 \end{bmatrix} \begin{bmatrix} \mathbf{w}_3 \\ \mathbf{w}_4 \end{bmatrix} \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \mathbf{w}_3 \end{bmatrix} \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \mathbf$$

Performance prediction: spectral clustering

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



$$\mathbf{v}_1 = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix} \quad \mathbf{v}_2 = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \end{bmatrix} \quad \mathbf{v}_3 = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \quad \mathbf{v}_3 = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_3 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_4 \end{bmatrix} \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_4 \end{bmatrix} \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_4 \end{bmatrix} \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_4 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \mathbf{v}_4 \end{bmatrix} \end{bmatrix}$$



Performance prediction: spectral clustering

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



 \mathbf{v}_3




Takeaway Message 2

"RMT Reassesses and Improves Data Processing"

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

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

• Going further than ([Kammoun,Couillet'17]), if $f'(\tau) = 0$,

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

• Going further than ([Kammoun,Couillet'17]), if $f'(\tau) = 0$,

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

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

 \rightarrow <u>Remark</u>: highly counter-intuitive kernel!





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



25 / 47

Semi-supervised learning: a great idea that never worked!

Semi-supervised learning: a great idea that never worked!

Setting: assume now

$$\begin{array}{l} \bullet \hspace{0.1 cm} x_{1}^{(a)}, \ldots, x_{n_{a,}[l]}^{(a)} \hspace{0.1 cm} \text{already labelled (few)}, \\ \bullet \hspace{0.1 cm} x_{n_{a,}[l]+1}^{(a)}, \ldots, x_{n_{a}}^{(a)} \hspace{0.1 cm} \text{unlabelled (a lot)}. \end{array}$$

Semi-supervised learning: a great idea that never worked!

- Setting: assume now
 - $\begin{array}{l} \blacktriangleright \ x_1^{(a)},\ldots,x_{n_a,[l]}^{(a)} \text{ already labelled (few),} \\ \\ \blacktriangleright \ x_{n_a,[l]+1}^{(a)},\ldots,x_{n_a}^{(a)} \text{ unlabelled (a lot).} \end{array}$

• Machine Learning original idea: find "scores" F_{ia} for x_i to belong to class a

$$F = \operatorname{argmin}_{F \in \mathbb{R}^{n \times k}} \sum_{a=1}^{k} \sum_{i,j} K_{ij} \left(F_{ia} - F_{ja} \right)^{2}, \quad F_{ia}^{[l]} = \delta_{\{x_{i} \in \mathcal{C}_{a}\}}.$$

Semi-supervised learning: a great idea that never worked!

- Setting: assume now
 - $\begin{array}{l} \blacktriangleright \ x_1^{(a)},\ldots,x_{n_a,[l]}^{(a)} \text{ already labelled (few),} \\ \\ \blacktriangleright \ x_{n_a,[l]+1}^{(a)},\ldots,x_{n_a}^{(a)} \text{ unlabelled (a lot).} \end{array}$

• Machine Learning original idea: find "scores" F_{ia} for x_i to belong to class a

$$F = \operatorname{argmin}_{F \in \mathbb{R}^{n \times k}} \sum_{a=1}^{k} \sum_{i,j} K_{ij} \left(F_{ia} D_{ii}^{\alpha} - F_{ja} D_{jj}^{\alpha} \right)^{2}, \quad F_{ia}^{[l]} = \delta_{\{x_{i} \in \mathcal{C}_{a}\}}.$$

Semi-supervised learning: a great idea that never worked!

- Setting: assume now
 - $\begin{array}{l} \blacktriangleright \ x_1^{(a)},\ldots,x_{n_a,[l]}^{(a)} \text{ already labelled (few),} \\ \\ \blacktriangleright \ x_{n_a,[l]+1}^{(a)},\ldots,x_{n_a}^{(a)} \text{ unlabelled (a lot).} \end{array}$

• Machine Learning original idea: find "scores" F_{ia} for x_i to belong to class a

$$F = \operatorname{argmin}_{F \in \mathbb{R}^{n \times k}} \sum_{a=1}^{k} \sum_{i,j} K_{ij} \left(F_{ia} D_{ii}^{\alpha} - F_{ja} D_{jj}^{\alpha} \right)^{2}, \quad F_{ia}^{[l]} = \delta_{\{x_{i} \in \mathcal{C}_{a}\}}.$$

Explicit solution:

$$F^{[u]} = \left(I_{n_{[u]}} - D_{[u]}^{-1-\alpha} K_{[uu]} D^{\alpha}{}_{[u]}\right)^{-1} D_{[u]}^{-1-\alpha} K_{[ul]} D^{\alpha}{}_{[l]} F^{[l]}$$

where $D = \text{diag}(K1_n)$ (degree matrix) and [ul], [uu], ... blocks of labeled/unlabeled data.

The finite-dimensional case: What we expect



Figure: Outcome F of Laplacian algorithms ($\alpha = -1$) for $\mathcal{N}(\pm \mu, I_p)$ with p = 1.

The finite-dimensional case: What we expect



Figure: Outcome F of Laplacian algorithms ($\alpha = -1$) for $\mathcal{N}(\pm \mu, I_p)$ with p = 1.

The reality: What we see!



Figure: Outcome F of Laplacian algorithms ($\alpha = -1$) for $\mathcal{N}(\pm \mu, I_p)$ with p = 80.

The reality: What we see!



Figure: Outcome F of Laplacian algorithms ($\alpha = -1$) for $\mathcal{N}(\pm \mu, I_p)$ with p = 80.

The reality: What we see! (on MNIST)



Figure: Vectors $[F^{(u)}]_{\cdot,a}$, a = 1, 2, 3, for 3-class MNIST data (zeros, ones, twos), n = 192, p = 784, $n_l/n = 1/16$, Gaussian kernel.

The reality: What we see! (on MNIST)



Figure: Vectors $[F^{(u)}]_{\cdot,a}$, a = 1, 2, 3, for 3-class MNIST data (zeros, ones, twos), n = 192, p = 784, $n_l/n = 1/16$, Gaussian kernel.

The reality: What we see! (on MNIST)



Figure: Vectors $[F^{(u)}]_{\cdot,a}$, a = 1, 2, 3, for 3-class MNIST data (zeros, ones, twos), n = 192, p = 784, $n_l/n = 1/16$, Gaussian kernel.

Consequences of the finite-dimensional "mismatch"

A priori, the algorithm should not work

- A priori, the algorithm should not work
- Indeed "in general" it does not!

- A priori, the algorithm should not work
- Indeed "in general" it does not!
- ▶ But, luckily, after some (not clearly motivated) renormalization (e.g., $\alpha = -1$, $F_{i.} \leftarrow F_{i.}/n_{[l],i}$), it works again...

- A priori, the algorithm should not work
- Indeed "in general" it does not!
- But, luckily, after some (not clearly motivated) renormalization (e.g., $\alpha = -1$, $F_{i.} \leftarrow F_{i.}/n_{[l],i}$), it works again...
- BUT it does not use efficiently unlabelled data!

Consequences of the finite-dimensional "mismatch"

- A priori, the algorithm should not work
- Indeed "in general" it does not!
- But, luckily, after some (not clearly motivated) renormalization (e.g., $\alpha = -1$, $F_{i.} \leftarrow F_{i.}/n_{[l],i}$), it works again...
- BUT it does not use efficiently unlabelled data!

Chapelle, Schölkopf, Zien, "Semi-Supervised Learning", Chapter 4, 2009.

Our concern is this: it is frequently the case that we would be better off just discarding the unlabeled data and employing a supervised method, rather than taking a semi-supervised route. Thus we worry about the embarrassing situation where the addition of unlabeled data degrades the performance of a classifier.

Asymptotic Performance Analysis

Theorem ([Mai,C'18] Asymptotic Performance of SSL) For $x_i \in C_b$ unlabelled, score vector $F_{i,\cdot} \in \mathbb{R}^k$ satisfies:

 $F_{i,\cdot} - G_b \to 0, \ G_b \sim \mathcal{N}(m_b, \Sigma_b)$

with $m_b \in \mathbb{R}^k$, $\Sigma_b \in \mathbb{R}^{k \times k}$ function of

• $f(\tau), f'(\tau), f''(\tau), \mu_1, \dots, \mu_k, C_1, \dots, C_k$

• only n_l .

Asymptotic Performance Analysis

Theorem ([Mai,C'18] Asymptotic Performance of SSL) For $x_i \in C_b$ unlabelled, score vector $F_{i,\cdot} \in \mathbb{R}^k$ satisfies:

$$F_{i,\cdot} - G_b \to 0, \ G_b \sim \mathcal{N}(m_b, \Sigma_b)$$

with $m_b \in \mathbb{R}^k$, $\Sigma_b \in \mathbb{R}^{k \times k}$ function of

•
$$f(\tau), f'(\tau), f''(\tau), \mu_1, \dots, \mu_k, C_1, \dots, C_k$$

• only n_k



Figure: Accuracy as a function of $n_{[u]}/p$ with $n_{[l]}/p = 2$, $c_1 = c_2$, p = 100, $-\mu_1 = \mu_2 = [1; \mathbf{0}_{p-1}], \{\mathbf{C}\}_{i,j} = .1^{|i-j|}$. Graph constructed with $K_{ij} = e^{-||x_i - x_j||^2/p}$.

Asymptotic Performance Analysis

Theorem ([Mai,C'18] Asymptotic Performance of SSL) For $x_i \in C_b$ unlabelled, score vector $F_{i,\cdot} \in \mathbb{R}^k$ satisfies:

$$F_{i,\cdot} - G_b \to 0, \ G_b \sim \mathcal{N}(m_b, \Sigma_b)$$

with $m_b \in \mathbb{R}^k$, $\Sigma_b \in \mathbb{R}^{k \times k}$ function of

•
$$f(\tau), f'(\tau), f''(\tau), \mu_1, \dots, \mu_k, C_1, \dots, C_k$$

• only n_k .



Figure: Accuracy as a function of $n_{[u]}/p$ with $n_{[l]}/p = 2$, $c_1 = c_2$, p = 100, $-\mu_1 = \mu_2 = [1; \mathbf{0}_{p-1}], \{\mathbf{C}\}_{i,j} = .1^{|i-j|}$. Graph constructed with $K_{ij} = e^{-||x_i - x_j||^2/p}$.

Solution: From RMT calculus (but not from ML intuition!), solution is to replace K by

$$\tilde{K} \equiv PKP, \quad P = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^{\mathsf{T}}.$$

Solution: From RMT calculus (but not from ML intuition!), solution is to replace K by

$$\tilde{K} \equiv PKP, \quad P = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^{\mathsf{T}}.$$

Theorem ([Mai,C'19] Asymptotic Performance of Improved SSL) For $x_i \in C_b$ unlabelled, score vector $\tilde{F}_{i,\cdot} \in \mathbb{R}^k$ satisfies:

$$\tilde{F}_{i,\cdot} - \tilde{G}_b \to 0, \ \tilde{G}_b \sim \mathcal{N}(\tilde{m}_b, \tilde{\Sigma}_b)$$

with $\tilde{m}_b \in \mathbb{R}^k$, $\tilde{\Sigma}_b \in \mathbb{R}^{k \times k}$ function of

• $f(\tau), f'(\tau), f''(\tau), \mu_1, \dots, \mu_k, C_1, \dots, C_k$

 \triangleright n_l and n_u .

Solution: From RMT calculus (but not from ML intuition!), solution is to replace K by

$$\tilde{K} \equiv PKP, \quad P = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^{\mathsf{T}}.$$

Theorem ([Mai,C'19] Asymptotic Performance of Improved SSL) For $x_i \in C_b$ unlabelled, score vector $\tilde{F}_{i,\cdot} \in \mathbb{R}^k$ satisfies:

$$\tilde{F}_{i,\cdot} - \tilde{G}_b \to 0, \ \tilde{G}_b \sim \mathcal{N}(\tilde{m}_b, \tilde{\Sigma}_b)$$

with $\tilde{m}_b \in \mathbb{R}^k$, $\tilde{\Sigma}_b \in \mathbb{R}^{k \times k}$ function of

• $f(\tau), f'(\tau), f''(\tau), \mu_1, \dots, \mu_k, C_1, \dots, C_k$

 \triangleright n_l and n_u .



 $n_{[u]}/p$

Solution: From RMT calculus (but not from ML intuition!), solution is to replace K by

$$\tilde{K} \equiv PKP, \quad P = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^{\mathsf{T}}.$$

Theorem ([Mai,C'19] Asymptotic Performance of Improved SSL) For $x_i \in C_b$ unlabelled, score vector $\tilde{F}_{i,\cdot} \in \mathbb{R}^k$ satisfies:

$$\tilde{F}_{i,\cdot} - \tilde{G}_b \to 0, \ \tilde{G}_b \sim \mathcal{N}(\tilde{m}_b, \tilde{\Sigma}_b)$$

with $\tilde{m}_b \in \mathbb{R}^k$, $\tilde{\Sigma}_b \in \mathbb{R}^{k \times k}$ function of

• $f(\tau), f'(\tau), f''(\tau), \mu_1, \dots, \mu_k, C_1, \dots, C_k$

 \triangleright n_l and n_u .



 $n_{[u]}/p$



Figure: Top: distribution of normalized pairwise distances for noisy MNIST data (8,9). Bottom: average accuracy as a function of $n_{[u]}$ with $n_{[l]} = 10$, computed over 1000 random realizations.



Figure: Top: distribution of normalized pairwise distances for noisy MNIST data (8,9). Bottom: average accuracy as a function of $n_{[u]}$ with $n_{[l]} = 10$, computed over 1000 random realizations.

What about real data?



Figure: Top: distribution of normalized pairwise distances for noisy MNIST data (8,9). Bottom: average accuracy as a function of $n_{|u|}$ with $n_{|l|} = 10$, computed over 1000 random realizations.

What about real data?



Figure: Top: distribution of normalized pairwise distances for noisy MNIST data (8,9). Bottom: average accuracy as a function of $n_{|u|}$ with $n_{|l|} = 10$, computed over 1000 random realizations.

What about real data?



Figure: Top: distribution of normalized pairwise distances for noisy MNIST data (8,9). Bottom: average accuracy as a function of $n_{[u]}$ with $n_{[l]} = 10$, computed over 1000 random realizations.
Experimental evidence: MNIST

| 2 | | | | | | | |
|--|--|--|--|--|--|--|--|
| (6,9) | | | | | | | |
| $n_u = 100$ | | | | | | | |
| 85.3±5.9 85.3±5.9 70.0±5.5 81.4±6.8 82.8±6.5 | | | | | | | |
| | | | | | | | |
| 92.6 ± 1.6 92.9 ± 1.4 69.5 ± 3.7 92.0 ± 1.6 91.4 ± 2.0 | | | | | | | |
| | | | | | | | |

Table: Comparison of classification accuracy (%) on MNIST datasets with $n_l = 10$. Computed over 1000 random iterations for $n_u = 100$ and 100 for $n_u = 1000$.

Experimental evidence: Traffic signs (HOG features)

| 6 | | 3 | 0 | 5 | | 30 |
|---|---|----|---|---|----|----|
| | 1 | | 0 | | 30 | |
| | E | 70 | Ø | | 0 | 03 |

| Class ID | (2,7) | (9,10) | (11,18) | | | | | |
|--------------------------------|------------------|-----------------|-----------------|--|--|--|--|--|
| $n_u = 100$ | | | | | | | | |
| Centered kernel (RMT) | 79.0±10.4 | 77.5±9.2 | 78.5±7.1 | | | | | |
| Iterated centered kernel (RMT) | 85.3±5.9 | 89.2±5.6 | 90.1±6.7 | | | | | |
| Laplacian | 73.8±9.8 | 77.3±9.5 | 78.6±7.2 | | | | | |
| Iterated Laplacian | 83.7±7.2 | 88.0±6.8 | 87.1±8.8 | | | | | |
| Manifold | 77.6 ± 8.9 | $81.4{\pm}10.4$ | $82.3{\pm}10.8$ | | | | | |
| $n_u = 1000$ | | | | | | | | |
| Centered kernel (RMT) | 83.6±2.4 | 84.6±2.4 | 88.7±9.4 | | | | | |
| Iterated centered kernel (RMT) | 84.8±3.8 | $88.0{\pm}5.5$ | 96.4±3.0 | | | | | |
| Laplacian | 72.7±4.2 | 88.9±5.7 | 95.8±3.2 | | | | | |
| Iterated Laplacian | $83.0 {\pm} 5.5$ | 88.2±6.0 | $92.7{\pm}6.1$ | | | | | |
| Manifold | 77.7±5.8 | $85.0{\pm}9.0$ | $90.6{\pm}8.1$ | | | | | |

Table: Comparison of classification accuracy (%) on German Traffic Sign datasets with $n_l = 10$. Computed over 1000 random iterations for $n_u = 100$ and 100 for $n_u = 1000$.

- Computation cost reduction: $(p, n \gg 1)$
 - $\rightarrow \ \varepsilon$ -subsampling $K \in \mathbb{R}^{n \varepsilon \times n \varepsilon}$



• Computation cost reduction: $(p, n \gg 1)$

 $\rightarrow \varepsilon$ -subsampling $K \in \mathbb{R}^{n \varepsilon \times n \varepsilon}$





• Computation cost reduction: $(p, n \gg 1)$

 $\rightarrow \ \varepsilon\text{-subsampling } K \in \mathbb{R}^{n\varepsilon \times n\varepsilon}$ $\rightarrow \ K_{\varepsilon} \equiv K \odot B \text{ with } B_{ij} \sim \text{Bern}(\varepsilon) \text{ i.i.d.}$





• Computation cost reduction: $(p, n \gg 1)$

 $\rightarrow \ \varepsilon\text{-subsampling } K \in \mathbb{R}^{n\varepsilon \times n\varepsilon}$ $\rightarrow \ K_{\varepsilon} \equiv K \odot B \text{ with } B_{ij} \sim \text{Bern}(\varepsilon) \text{ i.i.d.}$





• Computation cost reduction: $(p, n \gg 1)$

 $\rightarrow \ \varepsilon\text{-subsampling } K \in \mathbb{R}^{n\varepsilon \times n\varepsilon}$ $\rightarrow \ K_{\varepsilon} \equiv K \odot B \text{ with } B_{ij} \sim \text{Bern}(\varepsilon) \text{ i.i.d.}$





Takeaway Message 3

"RMT Also Grasps 'Real Data' Processing"

From i.i.d. to concentrated random vectors

Beyond Gaussian Mixtures: results still valid for concentrated random vectors.

From i.i.d. to concentrated random vectors

Beyond Gaussian Mixtures: results still valid for concentrated random vectors.

Definition (Concentrated Random Vector)

 $x \in \mathbb{R}^p$ is concentrated if, for all Lipschitz $f : \mathbb{R}^p \to \mathbb{R}$, there exists $m_f \in \mathbb{R}$, such that

 $P\left(|f(x) - m_f| > \varepsilon\right) \le e^{-g(\varepsilon)}, \quad g \text{ increasing function.}$

From i.i.d. to concentrated random vectors

Beyond Gaussian Mixtures: results still valid for concentrated random vectors.

Definition (Concentrated Random Vector)

 $x \in \mathbb{R}^p$ is concentrated if, for all Lipschitz $f : \mathbb{R}^p \to \mathbb{R}$, there exists $m_f \in \mathbb{R}$, such that

 $P\left(|f(x) - m_f| > \varepsilon\right) \le e^{-g(\varepsilon)}, \quad g \text{ increasing function.}$



Theorem ([Louart,C'18] [Seddik,C'19] Kernel Universality) For $x_i \sim \mathcal{L}(\mu_a, C_a)$ concentrated random vector, under the conditions of [C-Benaych'16],

$$\|K - \hat{K}\| \xrightarrow{\text{a.s.}} 0, \quad \hat{K} = f(\tau) \mathbf{1}_n \mathbf{1}_n^\mathsf{T} + \frac{1}{p} Z Z^\mathsf{T} + J A J^\mathsf{T} + *$$

with A only dependent on $f(\tau), f'(\tau), f''(\tau), \mu_1, \ldots, \mu_k, C_1, \ldots, C_k$.

Theorem ([Louart,C'18] [Seddik,C'19] Kernel Universality) For $x_i \sim \mathcal{L}(\mu_a, C_a)$ concentrated random vector, under the conditions of [C-Benaych'16],

$$\|K - \hat{K}\| \xrightarrow{\text{a.s.}} 0, \quad \hat{K} = f(\tau) \mathbf{1}_n \mathbf{1}_n^{\mathsf{T}} + \frac{1}{p} Z Z^{\mathsf{T}} + J A J^{\mathsf{T}} + *$$

with A only dependent on $f(\tau), f'(\tau), f''(\tau), \mu_1, \ldots, \mu_k, C_1, \ldots, C_k$.

Same result as [C-Benaych'16]... Universality of first two moments!

Key Finding. GAN-generated data are concentrated random vectors!

$\mathsf{Ok}.\,.\,\mathsf{so}\,\,\mathsf{what}?$



Key Finding. GAN-generated data are concentrated random vectors!





Results. [Seddik,C'19]









44 / 47

Our Research Activities:

Random Matrix Theory for Data Processing











Institut Fourier

géométrie





GIPSA

tenseurs







+PhD



GIPSA théorie de l'info

O. Michel GIPSA signal, physique





Trait. signal

M. Seddik

Apprentissage

applis vision

araphes

GIPSA

statistiques

L. Dall'Amico Physique Stats

C. Louart Mathématiques concentration

transfer, SSL

LIG

M. Tiomoko Apprentissage

traitement langage stats, physique

H. Chakroun Mathématiques géométrie

GIPSA

C. Doz Apprentissage RMT et radar

GIPSA

graphes

T. Zarrouk Apprentissage

C. Séjourné RMT structure RMT non convexe

Apprentissage

B. Nabet Finance ML & fi-stats H. Goulart

tenseurs







GIPSA

tenseurs

concentration

(+P.D.



traitement langage stats, physique



+PhD





S. Zozor GIPSA théorie de l'info

O. Michel GIPSA signal, physique





H. Goulart



Institut Fourier

géométrie



araphes

GIPSA

statistiques

M. Seddik Apprentissage applis vision

L. Dall'Amico

C. Louart Physique Stats Mathématiques

LIG

M. Tiomoko Apprentissage transfer, SSL

H. Chakroun Mathématiques géométrie

GIPSA

C. Doz Apprentissage RMT et radar

GIPSA

graphes

T. Zarrouk Apprentissage

C. Séjourné RMT structure RMT non convexe

Apprentissage

B. Nabet Finance ML & fi-stats

Trait. signal tenseurs

Join us I

The End

Thank you!



[C-Benaych'16] R. Couillet, Benaych-Georges, "Kernel Spectral Clustering of Large Dimensional Data", Electronic Journal of Statistics, vol. 10, no. 1, pp. 1393-1454, 2016. [article]



[Mai,C'18] X. Mai, R. Couillet, "A random matrix analysis and improvement of semi-supervised learning for large dimensional data", Journal of Machine Learning Research, vol. 19, no. 79, pp. 1-27, 2018. [article]



[Louart,C'18] C. Louart, Z. Liao, R. Couillet, "A Random Matrix Approach to Neural Networks", The Annals of Applied Probability, vol. 28, no. 2, pp. 1190-1248, 2018. [article]



[Seddik,C'19] M. Seddik, M. Tamaazousti, R. Couillet, "Kernel Random Matrices of Large Concentrated Data: The Example of GAN-Generated Image", IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP'19), Brighton, UK, 2019. [article]



💘 H. Tiomoko Ali, R. Couillet, "Improved spectral community detection in large heterogeneous networks", Journal of Machine Learning Research, vol. 18, no. 225, pp. 1-49, 2018. [article]



R. Couillet, M. Tiomoko, S. Zozor, E. Moisan, "Random matrix-improved estimation of covariance matrix distances". Journal of Multivariate Analysis, vol. 174, pp. 104531, 2019. [article]



🚬 Z. Liao, R. Couillet, "A Large Dimensional Analysis of Least Squares Support Vector Machines", IEEE Transactions on Signal Processing, vol. 67, no.4, pp. 1065-1074, 2018. [article]