

On the Spectrum of Random Features Maps of High Dimensional Data

ICML 2018, Stockholm, Sweden

Zhenyu Liao, Romain Couillet

L2S, CentraleSupélec, Université Paris-Saclay, France

GSTATS IDEX DataScience Chair, GIPSA-lab, Université Grenoble-Alpes, France.



1 Problem Statement

2 Main Results

3 Summary

Problem Setup

Random projection/random feature maps for feature extraction:

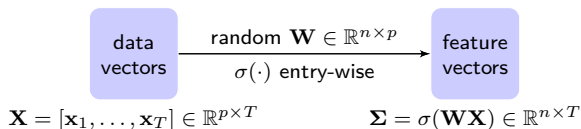


Figure: Illustration of random feature maps

Objective

Gram matrix of random features $\mathbf{G} \equiv \frac{1}{n} \mathbf{\Sigma}^T \mathbf{\Sigma}$ (sample covariance matrix in **feature** space):

- what kind of data **information** are extracted?
- what is the impact of different **nonlinearities**?
- how to perform clustering with \mathbf{G} , what do its **eigenvectors** “look like”?

With RMT: for large n, p, T , eigenspectrum of \mathbf{G} is determined **only** by¹

- the average **kernel** matrix $\Phi_{i,j} \equiv \mathbb{E}_{\mathbf{w}} \mathbf{G}_{i,j} = \mathbb{E}_{\mathbf{w}} \sigma(\mathbf{w}^T \mathbf{x}_i) \sigma(\mathbf{w}^T \mathbf{x}_j)$ (function of \mathbf{X})
- the ratios between n, p, T .

¹Louart Cosme, **Zhenyu Liao**, and Romain Couillet. “A Random Matrix Approach to Neural Networks.” The Annals of Applied Probability 28, no. 2 (2018): 1190-1248.

Some Known Facts

Objective: spectral characterization of Φ , with $\Phi_{i,j} = \mathbb{E}_{\mathbf{w}} \sigma(\mathbf{w}^T \mathbf{x}_i) \sigma(\mathbf{w}^T \mathbf{x}_j)$:
 For standard Gaussian $\mathbf{W} \Rightarrow$ integral calculus on \mathbb{R}^P .

Table: $\Phi_{i,j}$ for commonly used $\sigma(\cdot)$, $\angle \equiv \frac{\mathbf{x}_i^T \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}$.

$\sigma(t)$	$\Phi_{i,j}$
t	$\mathbf{x}_i^T \mathbf{x}_j$
$\max(t, 0)$	$\frac{1}{2\pi} \ \mathbf{x}_i\ \ \mathbf{x}_j\ \left(\angle \arccos(-\angle) + \sqrt{1 - \angle^2} \right)$
$ t $	$\frac{2}{\pi} \ \mathbf{x}_i\ \ \mathbf{x}_j\ \left(\angle \arcsin(\angle) + \sqrt{1 - \angle^2} \right)$
$\varsigma_+ \max(t, 0) + \varsigma_- \max(-t, 0)$	$\frac{1}{2} (\varsigma_+^2 + \varsigma_-^2) \mathbf{x}_i^T \mathbf{x}_j + \frac{\ \mathbf{x}_i\ \ \mathbf{x}_j\ }{2\pi} (\varsigma_+ + \varsigma_-)^2 \left(\sqrt{1 - \angle^2} - \angle \cdot \arccos(\angle) \right)$
$1_{t>0}$	$\frac{1}{2} - \frac{1}{2\pi} \arccos(\angle)$
$\text{sign}(t)$	$\frac{2}{\pi} \arcsin(\angle)$
$\varsigma_2 t^2 + \varsigma_1 t + \varsigma_0$	$\varsigma_2^2 \left(2 \left(\mathbf{x}_i^T \mathbf{x}_j \right)^2 + \ \mathbf{x}_i\ ^2 \ \mathbf{x}_j\ ^2 \right) + \varsigma_1^2 \mathbf{x}_i^T \mathbf{x}_j + \varsigma_2 \varsigma_0 \left(\ \mathbf{x}_i\ ^2 + \ \mathbf{x}_j\ ^2 \right) + \varsigma_0^2$
$\cos(t)$	$\exp\left(-\frac{1}{2} \left(\ \mathbf{x}_i\ ^2 + \ \mathbf{x}_j\ ^2 \right)\right) \cosh(\mathbf{x}_i^T \mathbf{x}_j)$
$\sin(t)$	$\exp\left(-\frac{1}{2} \left(\ \mathbf{x}_i\ ^2 + \ \mathbf{x}_j\ ^2 \right)\right) \sinh(\mathbf{x}_i^T \mathbf{x}_j)$
$\text{erf}(t)$	$\frac{2}{\pi} \arcsin\left(\frac{2\mathbf{x}_i^T \mathbf{x}_j}{\sqrt{(1+2\ \mathbf{x}_i\ ^2)(1+2\ \mathbf{x}_j\ ^2)}}\right)$
$\exp\left(-\frac{t^2}{2}\right)$	$\frac{1}{\sqrt{(1+\ \mathbf{x}_i\ ^2)(1+\ \mathbf{x}_j\ ^2) - (\mathbf{x}_i^T \mathbf{x}_j)^2}}$

\Rightarrow (still) highly **nonlinear** functions of the data \mathbf{x} !

Dig Deeper into the Average Kernel Φ

Data Model

Consider data from a K -class Gaussian mixture model: $\mathbf{x}_i \in \mathcal{C}_a \Leftrightarrow \mathbf{x}_i = \boldsymbol{\mu}_a/\sqrt{p} + \boldsymbol{\omega}_i$, with $\boldsymbol{\omega}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_a/p)$, $a = 1, \dots, K$ of statistical **mean** $\boldsymbol{\mu}_a$ and **covariance** \mathbf{C}_a .

Non-trivial Classification [Neyman-Pearson Minimal]

For p large, we have $\|\boldsymbol{\mu}_a - \boldsymbol{\mu}_b\| = O(1)$, $\|\mathbf{C}_a\| = O(1)$ and $\text{tr}(\mathbf{C}_a - \mathbf{C}_b)/\sqrt{p} = O(1)$.

As a consequence,

$$\|\mathbf{x}_i\|^2 = \underbrace{\|\boldsymbol{\omega}_i\|^2}_{O(1)} + \underbrace{\|\boldsymbol{\mu}_a\|^2/p + 2\boldsymbol{\mu}_a^\top \boldsymbol{\omega}_i/\sqrt{p}}_{O(p^{-1})} = \underbrace{\text{tr} \mathbf{C}_a/p}_{O(1)} + \underbrace{\|\boldsymbol{\omega}_i\|^2 - \text{tr} \mathbf{C}_a/p}_{O(p^{-1/2})} + \underbrace{\|\boldsymbol{\mu}_a\|^2/p + 2\boldsymbol{\mu}_a^\top \boldsymbol{\omega}_i/\sqrt{p}}_{O(p^{-1})}$$

- if relaxed, classification **too easy**: it suffices to compare the norm $\|\mathbf{x}_i\|^2$ and $\|\mathbf{x}_j\|^2$!
- in fact reveals a more **intrinsic** property of **high dimensional data**:

Curse of dimensionality: **little difference** in Euclidean distance between pairs!

Denote $\mathbf{C}^\circ = \sum_{i=1}^K \frac{T_i}{T} \mathbf{C}_a$ and $\mathbf{C}_a = \mathbf{C}_a^\circ + \mathbf{C}^\circ$ for $a = 1, \dots, K$.

Then $\|\mathbf{x}_i\|^2 = \tau + O(p^{-1/2})$ with $\tau \equiv \text{tr}(\mathbf{C}^\circ)/p$, $\|\mathbf{x}_i - \mathbf{x}_j\|^2 = \|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2 - \mathbf{x}_i^\top \mathbf{x}_j \approx 2\tau$:

\Rightarrow Almost **constant** distance no matter from the **same** or **different** classes!

Dig Deeper into the Average Kernel Φ

Why things are still working? \Rightarrow statistical information are **hidden** in smaller order terms!

$$\Rightarrow \|\mathbf{x}_i - \mathbf{x}_j\|^2 = \|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2 - \mathbf{x}_i^\top \mathbf{x}_j \approx 2\tau + \underbrace{\omega_i^\top \omega_j}_{O(p^{-1/2})} + \underbrace{\mu_a^\top \mu_b / p + \mu_a^\top \omega_j / \sqrt{p} + \mu_b^\top \omega_i / \sqrt{p}}_{O(p^{-1})}$$

Small **entry-wise** \neq small in **matrix form** (in operator norm): **repeated** in $p \times p$ large matrix
 \Rightarrow spectral clustering works! 😊

Moreover, “concentration” brings simplifications: for $\Phi_{i,j} = \mathbb{E}_{\mathbf{w}} \sigma(\mathbf{w}^\top \mathbf{x}_i) \sigma(\mathbf{w}^\top \mathbf{x}_j)$ and ReLU,

$$\Phi_{i,j} = \frac{1}{2\pi} \|\mathbf{x}_i\| \|\mathbf{x}_j\| \left(\angle \arccos(-\angle) + \sqrt{1 - \angle^2} \right)$$

with $\angle \equiv \frac{\mathbf{x}_i^\top \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}$. “**Concentration**”: $\angle = \frac{0}{\tau^2} +$ information terms (μ_a, C_a) !

“Blessing” of Dimensionality

High dimensional “concentration” \Rightarrow Taylor expansion to **linearize** Φ !

Main Results

Asymptotic Equivalent of Φ

For all $\sigma(\cdot)$ listed in the table above, we have, as $n \sim p \sim T \rightarrow \infty$,

$$\|\Phi - \tilde{\Phi}\| \rightarrow 0$$

almost surely, with

$$\tilde{\Phi} \equiv d_1 \left(\Omega + \mathbf{M} \frac{\mathbf{J}^\top}{\sqrt{p}} \right)^\top \left(\Omega + \mathbf{M} \frac{\mathbf{J}^\top}{\sqrt{p}} \right) + d_2 \mathbf{U} \mathbf{B} \mathbf{U}^\top + d_0 \mathbf{I}_T$$

$$\text{and } \mathbf{U} \equiv \left[\frac{\mathbf{J}}{\sqrt{p}}, \phi \right], \quad \mathbf{B} \equiv \begin{bmatrix} \mathbf{t} \mathbf{t}^\top + 2\mathbf{S} & \mathbf{t} \\ \mathbf{t}^\top & 1 \end{bmatrix}.$$

Table: Coefficients d_i in $\tilde{\Phi}$ for different $\sigma(\cdot)$.

$\sigma(t)$	d_1	d_2
t	1	0
$\max(t, 0)$	$\frac{1}{4}$	$\frac{1}{8\pi\tau}$
$ t $	0	$\frac{1}{2\pi\tau}$
$\varsigma_+ \max(t, 0) + \varsigma_- \max(-t, 0)$	$\frac{1}{4}(\varsigma_+ - \varsigma_-)^2$	$\frac{1}{8\tau\pi}(\varsigma_+ + \varsigma_-)^2$
$1_{t>0}$	$\frac{1}{2\pi\tau}$	0
$\text{sign}(t)$	$\frac{2}{\pi\tau}$	0
$\varsigma_2 t^2 + \varsigma_1 t + \varsigma_0$	ς_1^2	ς_2^2
$\cos(t)$	0	$\frac{e^{-\tau}}{4}$
$\sin(t)$	$e^{-\tau}$	0
$\text{erf}(t)$	$\frac{4}{\pi} \frac{1}{2\tau+1}$	0
$\exp(-\frac{t^2}{2})$	0	$\frac{1}{4(\tau+1)^3}$

With $\mathbf{J} \equiv [\mathbf{j}_1, \dots, \mathbf{j}_K]$, \mathbf{j}_a canonical vector of \mathcal{C}_a : $(\mathbf{j}_a)_i = \delta_{\mathbf{x}_i \in \mathcal{C}_a}$ (for clustering), weighted by

- Ω , ϕ random fluctuations of data.
- $\mathbf{M} \equiv [\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K]$, $\mathbf{t} \equiv \left\{ \text{tr} \mathbf{C}_a^\circ / \sqrt{p} \right\}_{a=1}^K$, $\mathbf{S} \equiv \left\{ \text{tr}(\mathbf{C}_a \mathbf{C}_b) / p \right\}_{a,b=1}^K$ statistical information from data distribution.

Consequence

Table: Coefficients d_i in $\tilde{\Phi}$ for different $\sigma(\cdot)$.

$\sigma(t)$	d_1	d_2
t	1	0
$\max(t, 0)$	$\frac{1}{4}$	$\frac{1}{8\pi\tau}$
$ t $	0	$\frac{1}{2\pi\tau}$
$\varsigma_+ \max(t, 0) + \varsigma_- \max(-t, 0)$	$\frac{1}{4}(\varsigma_+ - \varsigma_-)^2$	$\frac{1}{8\tau\pi}(\varsigma_+ + \varsigma_-)^2$
$1_{t>0}$	$\frac{1}{2\pi\tau}$	0
$\text{sign}(t)$	$\frac{2}{\pi\tau}$	0
$\varsigma_2 t^2 + \varsigma_1 t + \varsigma_0$	ς_1^2	ς_2^2
$\cos(t)$	0	$\frac{e^{-\tau}}{4}$
$\sin(t)$	$e^{-\tau}$	0
$\text{erf}(t)$	$\frac{4}{\pi} \frac{1}{2\tau+1}$	0
$\exp(-\frac{t^2}{2})$	0	$\frac{1}{4(\tau+1)^3}$

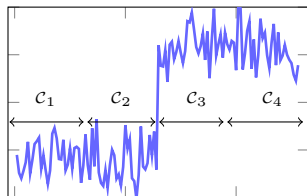
A natural classification of $\sigma(\cdot)$:

- **mean-oriented**, $d_1 \neq 0, d_2 = 0$:
 $t, 1_{t>0}, \text{sign}(t), \sin(t)$ and $\text{erf}(t)$
 \Rightarrow separate with difference in means \mathbf{M} ;
- **covariance-oriented**, $d_1 = 0, d_2 \neq 0$:
 $|t|, \cos(t)$ and $\exp(-t^2/2)$
 \Rightarrow track differences in covariances \mathbf{t}, \mathbf{S} ;
- **"balanced"**, both $d_1, d_2 \neq 0$:
 - ▶ ReLU function $\max(t, 0)$,
 - ▶ Leaky ReLU function
 $\varsigma_+ \max(t, 0) + \varsigma_- \max(-t, 0)$,
 - ▶ quadratic function $\varsigma_2 t^2 + \varsigma_1 t + \varsigma_0$. \Rightarrow make use of **both** statistics!

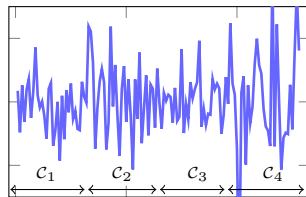
Numerical Validations: Gaussian Data

Example: Gaussian mixture data of four classes: $\mathcal{N}(\mu_1, \mathbf{C}_1)$, $\mathcal{N}(\mu_1, \mathbf{C}_2)$, $\mathcal{N}(\mu_2, \mathbf{C}_1)$ and $\mathcal{N}(\mu_2, \mathbf{C}_2)$ with Leaky ReLU function $\varsigma_+ \max(t, 0) + \varsigma_- \max(-t, 0)$.

Case 1: $\varsigma_+ = \varsigma_- = 1$ (equivalent to linear map $\sigma(t) = t$)

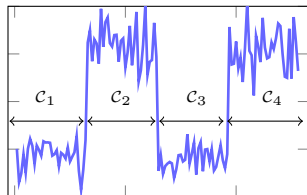


Eigenvector 1

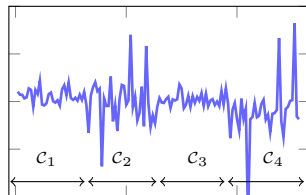


Eigenvector 2

Case 2: $\varsigma_+ = -\varsigma_- = 1$ (equivalent to $\sigma(t) = |t|$)



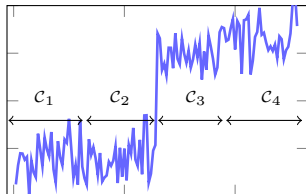
Eigenvector 1



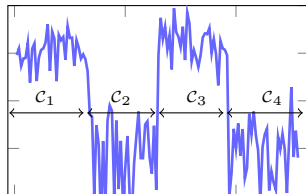
Eigenvector 2

Numerical Validations: Gaussian Data

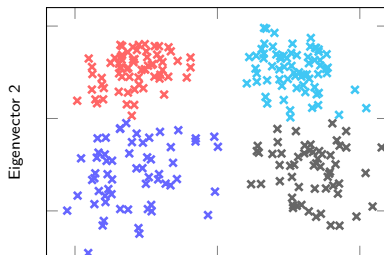
Case 3: $\zeta_+ = 1, \zeta_- = 0$ (the ReLU function)



Eigenvector 1



Eigenvector 2



Eigenvector 1

Numerical Validations: Real Datasets

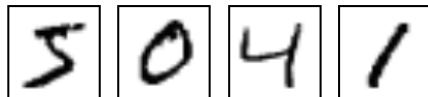


Figure: The MNIST image database.

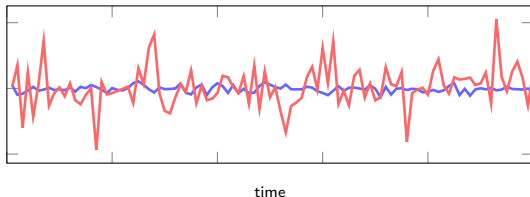


Figure: The epileptic EEG datasets.²

Reproducibility: codes available at <https://github.com/Zhenyu-LIAO/RMT4RFM>.

²<http://www.meb.unibonn.de/epileptologie/science/physik/eegdata.html>.

Numerical Validations: Real Datasets

Table: Empirical estimation of differences in means and covariances of the MNIST and epileptic EEG datasets.

	$\ \mathbf{M}^T\mathbf{M}\ $	$\ \mathbf{t}\mathbf{t}^T + 2\mathbf{S}\ $
MNIST data	172.4	86.0
EEG data	1.2	182.7

Table: Clustering accuracies on MNIST dataset.

	$\sigma(t)$	$T = 64$	$T = 128$
mean-oriented	t	88.94%	87.30%
	$1_{t>0}$	82.94%	85.56%
	$\text{sign}(t)$	83.34%	85.22%
	$\sin(t)$	87.81%	87.50%
	$\text{erf}(t)$	87.28%	86.59%
cov-oriented	$ t $	60.41%	57.81%
	$\cos(t)$	59.56%	57.72%
	$\exp(-\frac{t^2}{2})$	60.44%	58.67%
balanced	$\text{ReLU}(t)$	85.72%	82.27%

Table: Clustering accuracies on EEG dataset.

	$\sigma(t)$	$T = 64$	$T = 128$
mean-oriented	t	70.31%	69.58%
	$1_{t>0}$	65.87%	63.47%
	$\text{sign}(t)$	64.63%	63.03%
	$\sin(t)$	70.34%	68.22%
	$\text{erf}(t)$	70.59%	67.70%
cov-oriented	$ t $	99.69%	99.50%
	$\cos(t)$	99.38%	99.36%
	$\exp(-\frac{t^2}{2})$	99.81%	99.77%
balanced	$\text{ReLU}(t)$	87.91%	90.97%

Numerical Validations: Real Datasets

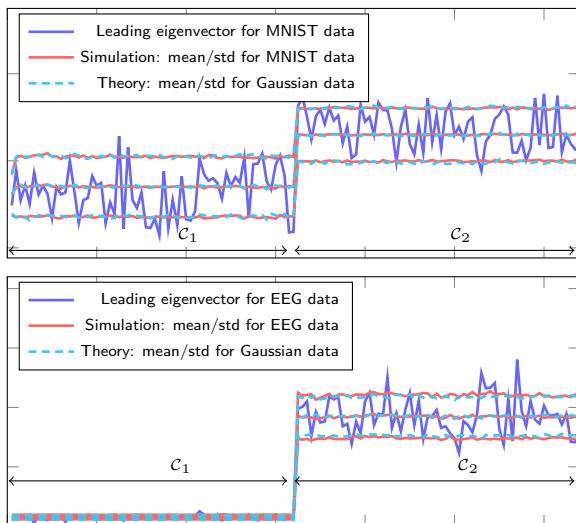


Figure: Leading eigenvector of Φ for the MNIST (top) and EEG (bottom) with Gaussian mixture data (of same statistics) with a width of ± 1 standard deviations.

Take-away message:

- “concentration” of high dimensional data to handle the **nonlinearity**
- different nonlinearities into three attributes: **mean-**, **covariance-**oriented and “**balanced**”
- **optimize** the choice of nonlinearity as a function of **data** (quadratic and LReLU)
- novel insight into **understanding of neural networks** for high dimensional data

Future work:

- study of the eigenvalue distribution \Rightarrow the (asymptotic) behavior of **leading eigenvectors**
- combination of different type of nonlinearities, e.g., $\sin + \cos \Rightarrow$ Gaussian kernel
- directly linking $\sigma(\cdot)$ and the coefficients d_0 , d_1 and d_2

Thank you

Thank you!

Poster # 62