

A Random Matrix Analysis of Data Stream Clustering: Coping With Limited Memory Resources

Hugo Lebeau¹ Romain Couillet¹ Florent Chatelain²

¹Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP, LIG, 38000 Grenoble, France

²Univ. Grenoble Alpes, CNRS, Grenoble INP, GIPSA-lab, 38000 Grenoble, France

39th International Conference on Machine Learning
July 17th – 23rd, 2022
Baltimore, Maryland, USA

Random matrix framework for data stream clustering

- Observed data: $\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_t; \dots \in \mathbb{R}^p$

$$\mathbf{x}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\boldsymbol{\mu}; \boldsymbol{\Sigma}_p)$$

Random matrix framework for data stream clustering

- Observed data: $\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_t; \dots \in \mathbb{R}^p$
- Memory

$$\mathbf{x}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\boldsymbol{\mu}, \Sigma_p)$$

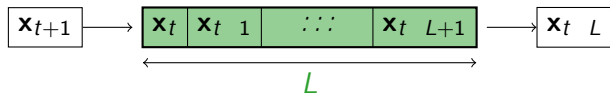


Random matrix framework for data stream clustering

- Observed data: $\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_t; \dots \in \mathbb{R}^p$

$\mathbf{x}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \Sigma)$

- Memory

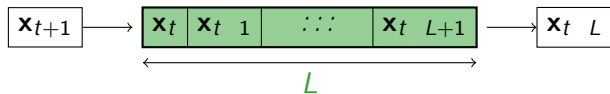


Random matrix framework for data stream clustering

- Observed data: $\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_t; \dots \in \mathbb{R}^p$

\mathbf{x}_t i.i.d. $N(\mu, \Sigma)$

- Memory



- Clustering?

Random matrix framework for data stream clustering

- Observed data: $\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_t; \dots \in \mathbb{R}^p$

$\mathbf{x}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \Sigma)$

- Memory



- Clustering** on the $n > L$ previous points

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_{t-n+1} & \mathbf{x}_{t-n+2} & \dots & \mathbf{x}_t \end{bmatrix} \in \mathbb{R}^{n \times p}$$

Random matrix framework for data stream clustering

- Observed data: $\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_t; \dots \in \mathbb{R}^p$

\mathbf{x}_t i.i.d. $N(\mu; \Sigma)$

- Memory



- Clustering** on the $n > L$ previous points

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_{t-n+1} & \mathbf{x}_{t-n+2} & \dots & \mathbf{x}_t \end{bmatrix} \in \mathbb{R}^{n \times p}$$

$$\begin{matrix} n; p; L & / & +1 \\ p=n & / & c \\ (2L-1)=n & / & " \end{matrix}$$

Random matrix framework for data stream clustering

- Observed data: $\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_t; \dots \in \mathbb{R}^p$

$\mathbf{x}_t \stackrel{\text{i.i.d.}}{\sim} N(\mu; \Sigma_p)$

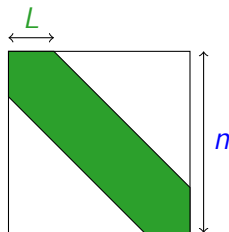
- Memory



- Clustering** on the $n > L$ previous points

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_{t-n+1} & \mathbf{x}_{t-n+2} & \dots & \mathbf{x}_t \end{bmatrix} \in \mathbb{R}^{n \times p} \quad (2L-1) \times n \quad \begin{matrix} +1 \\ c \\ " \end{matrix}$$

- Kernel matrix



Random matrix framework for data stream clustering

- Observed data: $\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_t; \dots \in \mathbb{R}^p$

\mathbf{x}_t i.i.d. $N(\mu; \Sigma)$

- Memory

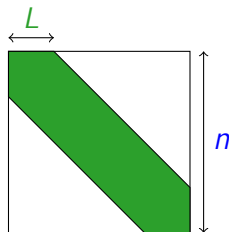


- Clustering** on the $n > L$ previous points

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_{t-n+1} & \mathbf{x}_{t-n+2} & \dots & \mathbf{x}_t \end{bmatrix} \in \mathbb{R}^{n \times p}$$

$$\begin{matrix} n; p; L & / & +1 \\ p=n & / & c \\ (2L-1)=n & / & \end{matrix}$$

- Kernel matrix



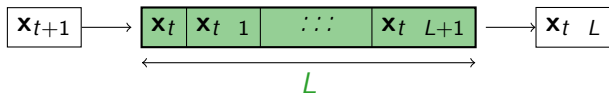
$$\mathbf{K} = \frac{\mathbf{X}^T \mathbf{X}}{p}$$

Random matrix framework for data stream clustering

- Observed data: $\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_t; \dots \in \mathbb{R}^p$

$\mathbf{x}_t \stackrel{\text{i.i.d.}}{\sim} N(\mu, \Sigma)$

- Memory

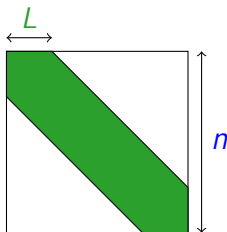


- Clustering** on the $n > L$ previous points

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_{t-n+1} & \mathbf{x}_{t-n+2} & \dots & \mathbf{x}_t \end{bmatrix} \in \mathbb{R}^{n \times p}$$

$$\begin{matrix} n > L & / & +1 \\ p & = & n & / & c \\ (2L-1) & = & n & / & \end{matrix}$$

- Kernel matrix



$$\mathbf{K}_L = \frac{\mathbf{X}^T \mathbf{X}}{p}$$

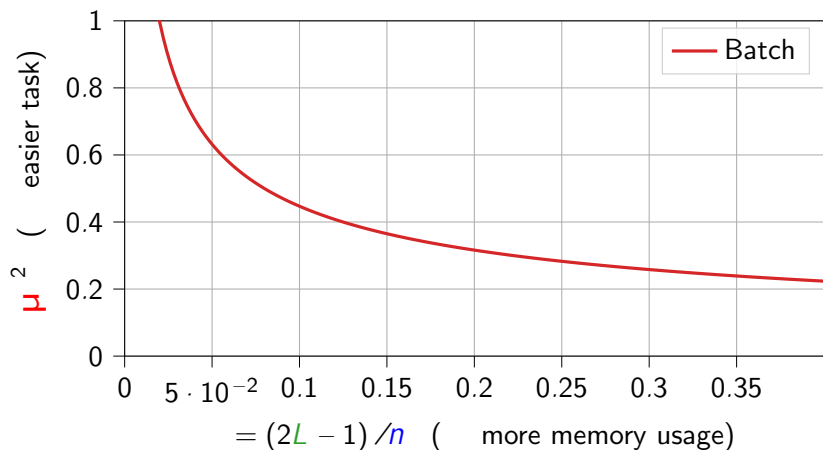
$$\begin{matrix} 2 & 1 & \dots & 1 & & 0 & 3 \\ 6 & \dots & \dots & \dots & \dots & 7 & \\ 6 & \dots & \dots & \dots & \dots & 7 & \\ 6 & 1 & \dots & \dots & \dots & 1 & 5 \\ & \dots & \dots & \dots & \dots & \dots & \\ | & 0 & \dots & 1 & \dots & 1 & \\ & & & \mathbf{T} & & & \end{matrix}$$

Data stream clustering: improving over batch clustering

- Size- L memory *batch clustering* vs. $\mathbf{K}_L = \frac{1}{p} \mathbf{X}^T \mathbf{X}$ \mathbf{T}

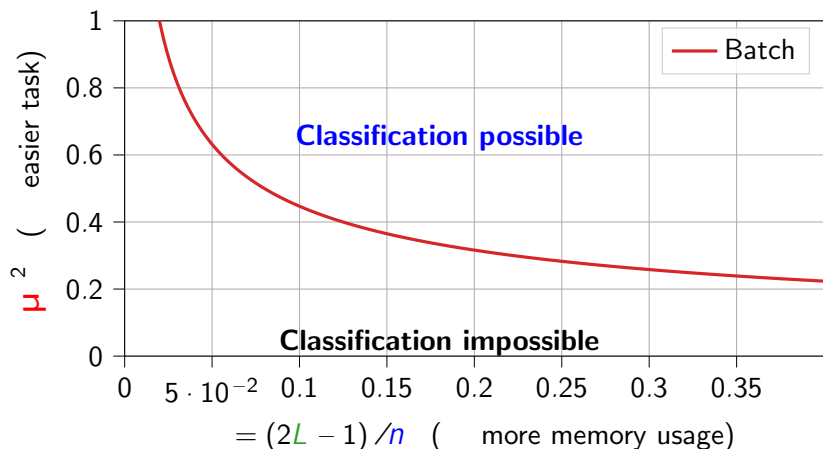
Data stream clustering: improving over batch clustering

- Size- L memory *batch clustering* vs. $\mathbf{K}_L = \frac{1}{p} \mathbf{X}^T \mathbf{X}$ \mathbf{T}
- **Spectral clustering phase transition** ($n=p=100$ () $c=0.01$)



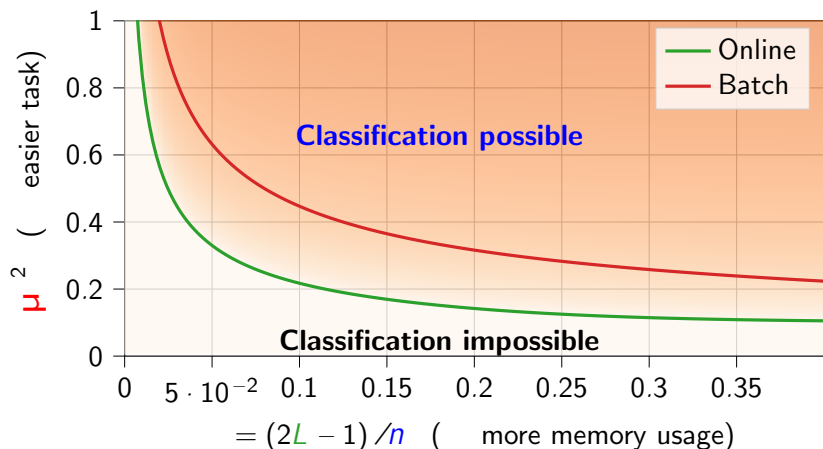
Data stream clustering: improving over batch clustering

- Size- L memory *batch clustering* vs. $\mathbf{K}_L = \frac{1}{p} \mathbf{X}^T \mathbf{X} \mathbf{T}$
- **Spectral clustering phase transition** ($n=p=100$) ($c=0.01$)



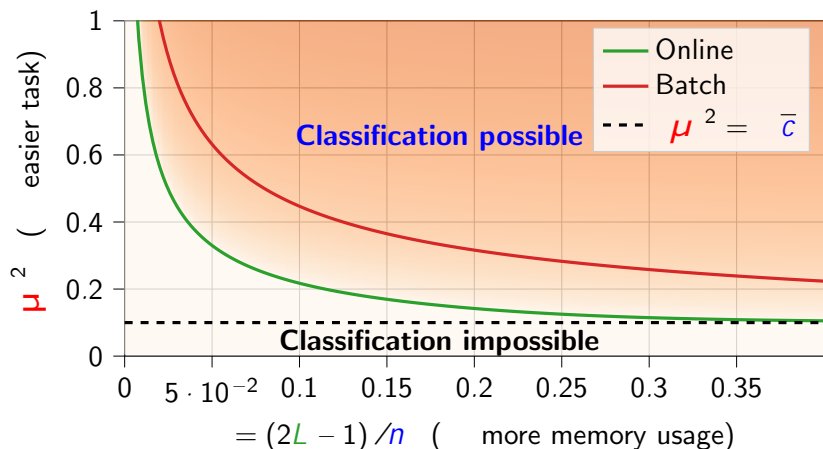
Data stream clustering: improving over batch clustering

- Size- L memory *batch clustering* vs. $\mathbf{K}_L = \frac{1}{p} \mathbf{X}^T \mathbf{X} \mathbf{T}$
- **Spectral clustering phase transition** ($n=p=100$) ($c=0.01$)



Data stream clustering: improving over batch clustering

- Size- L memory *batch clustering* vs. $\mathbf{K}_L = \frac{1}{p} \mathbf{X}^T \mathbf{X}$ \mathbf{T}
- **Spectral clustering phase transition** ($n=p=100$ (\cdot) $c=0.01$)



From theory to practice: online kernel spectral clustering

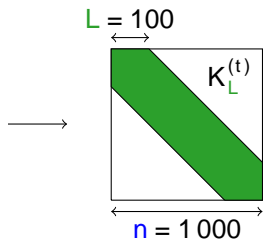
- **BigGAN images** (VGG features, $p = 4\,096$)

$T = 20\,000$

From theory to practice: online kernel spectral clustering

BigGAN images (VGG features $p = 4096$)

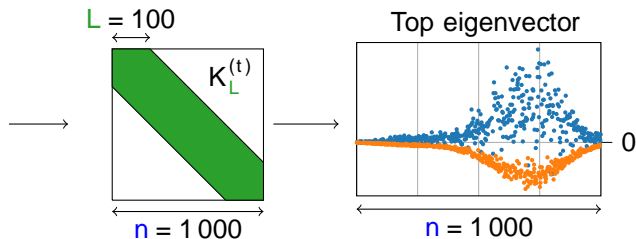
$T = 20000$



From theory to practice: online kernel spectral clustering

BigGAN images (VGG features $p = 4096$)

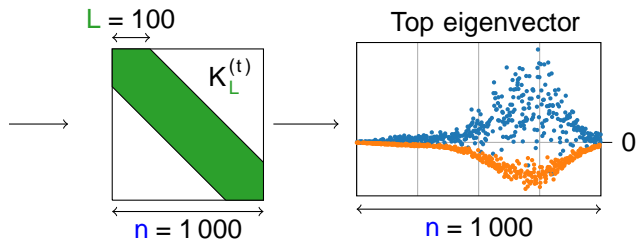
$T = 20\,000$



From theory to practice: online kernel spectral clustering

BigGAN images (VGG features $p = 4096$)

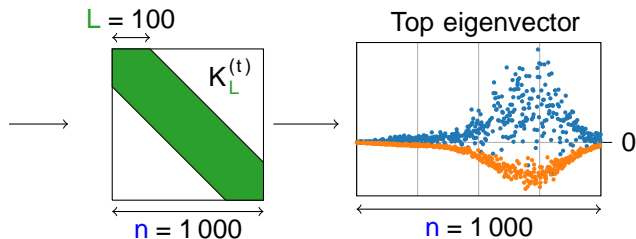
$T = 20\,000$



From theory to practice: online kernel spectral clustering

BigGAN images (VGG features $p = 4096$)

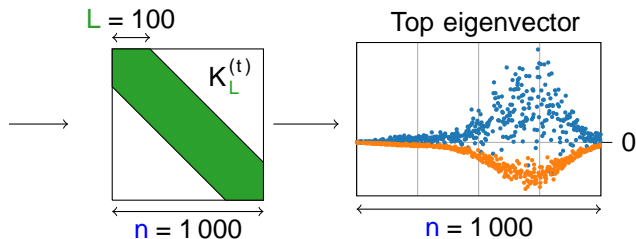
$T = 20\,000$



From theory to practice: online kernel spectral clustering

BigGAN images (VGG features $p = 4096$)

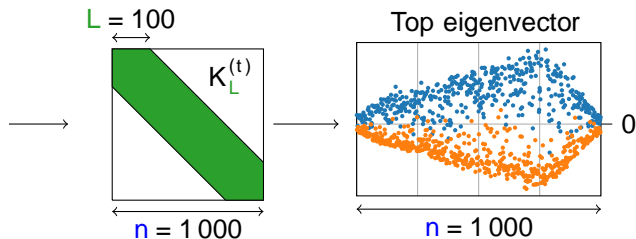
$T = 20\,000$



From theory to practice: online kernel spectral clustering

Fashion-MNIST images (raw, $p = 784$)

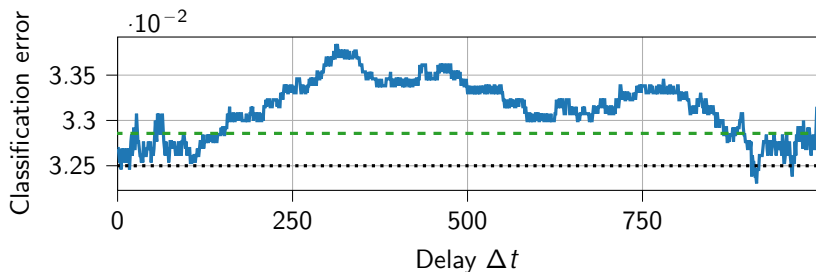
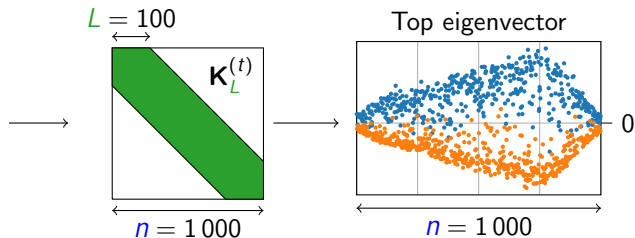
$T = 14\,000$



From theory to practice: online kernel spectral clustering

- Fashion-MNIST images (raw, $p = 784$)

$T = 14\,000$



A Random Matrix Analysis of Data Stream Clustering: Coping With Limited Memory Resources

Hugo Lebeau¹ Romain Couillet¹ Florent Chatelain²

¹Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP, LIG, 38000 Grenoble, France

²Univ. Grenoble Alpes, CNRS, Grenoble INP, GIPSA-lab, 38000 Grenoble, France

39th International Conference on Machine Learning
July 17th – 23rd, 2022
Baltimore, Maryland, USA