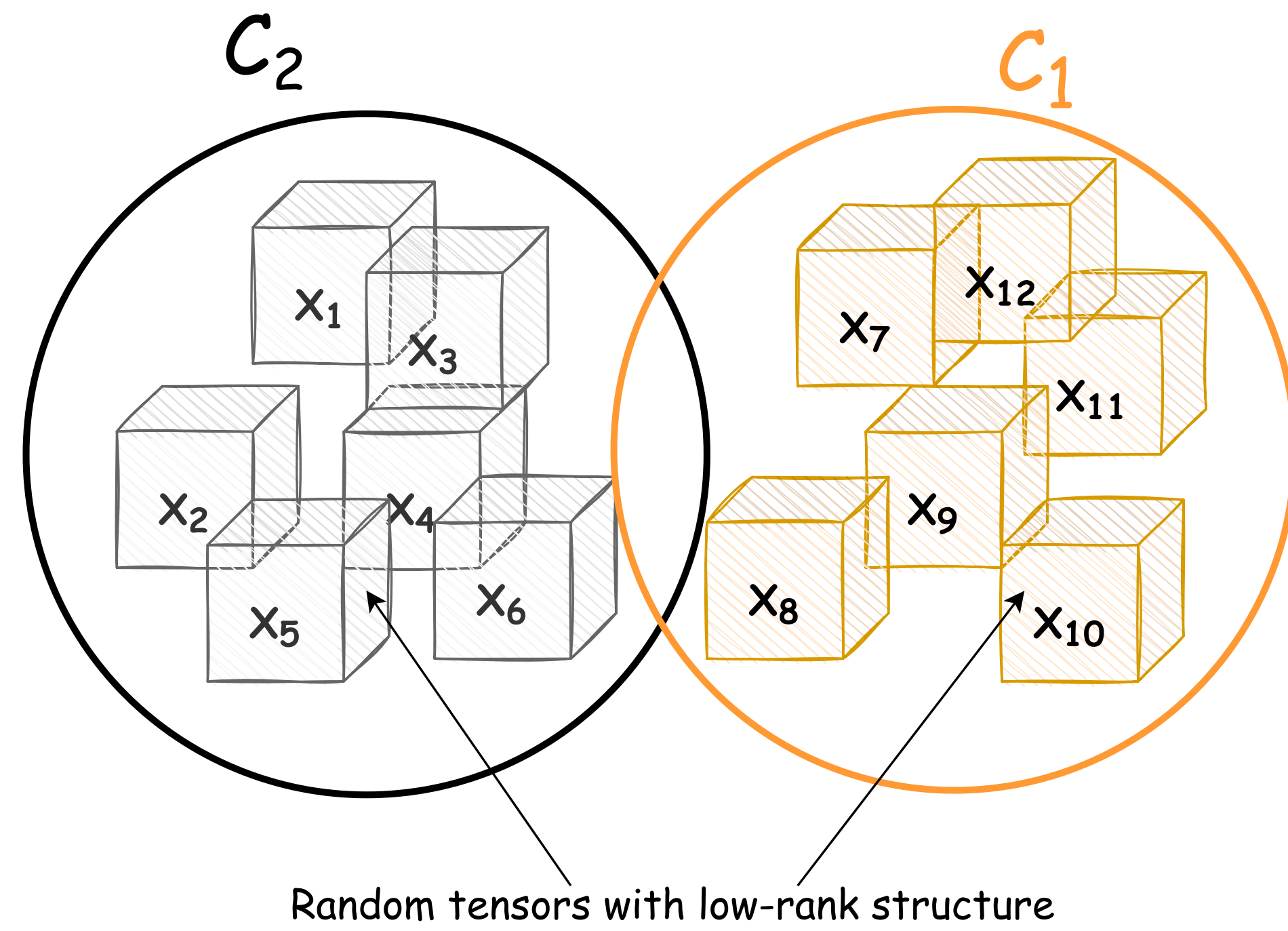


Abstract



- Theoretical analysis of learning from data with hidden **low-rank tensor** structure.
- Quantification of **performance gain** between considering the low-rank **tensor structure** versus treating data as **vectors**.

Setting & Data Model

We consider n data points: $(\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3)_{ijk} = x_{1i}x_{2j}x_{3k}$

$$\mathbf{X}_i \in \mathcal{C}_a \Leftrightarrow \mathbf{X}_i = (-1)^a \boldsymbol{\mu}_1 \otimes \cdots \otimes \boldsymbol{\mu}_k + \mathbf{Z}_i \in \mathbb{R}^{p_1 \times \cdots \times p_k}$$

where $[\mathbf{Z}_i]_{i_1 \dots i_k} \sim \mathcal{N}(0, 1)$ i.i.d. and denote $\mathbf{M} = \boldsymbol{\mu}_1 \otimes \cdots \otimes \boldsymbol{\mu}_k$.

- Generalizes the classical model ($k = 1$), i.e. $\mathbf{x}_i = (-1)^a \boldsymbol{\mu}_1 + \mathbf{z}_i$.
- Even for $k \geq 2$, the standard approach consists in **flattening** the data.
- What is the **optimal** classifier? Theoretical misclassification?

Supervised Setting

Given $\mathbf{X} = [\mathbf{X}_1, \dots, \mathbf{X}_n] \in \mathbb{R}^{p_1 \times \cdots \times p_k \times n}$ and $\mathbf{y} = [y_1, \dots, y_n] \in \{-1, 1\}^n$

Denote $\mathbf{X} = \mathbf{X}_{(k+1)} \in \mathbb{R}^{n \times P}$ with $P = \prod_{i=1}^k p_i$ and consider the *Ridge* classifier:

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \gamma \|\mathbf{w}\|^2 \Leftrightarrow \mathbf{w}^* = (\mathbf{X}^\top \mathbf{X} + \gamma \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

For some $\gamma \gg \|\mathbf{X}^\top \mathbf{X}\|$ (optimal for the above data model):

$$\mathbf{w} = \frac{1}{\sqrt{np}} \mathbf{X}^\top \mathbf{y}$$

where $p = \sum_{i=1}^k p_i$. In tensor notations, the decision function is:

$$f_{\text{R}}(\tilde{\mathbf{X}}_i) = \langle \mathbf{W}, \tilde{\mathbf{X}}_i \rangle \stackrel{\mathcal{C}_1}{\leq} 0 \quad \mathbf{W} \equiv \frac{1}{\sqrt{np}} \mathbf{X} \times_{k+1} \mathbf{y}$$

with $\tilde{\mathbf{X}}_i$ a test datum independent of \mathbf{X} .

Assumption: $p_i = \mathcal{O}(n)$ and $\|\mathbf{M}\| = \mathcal{O}(1)$.

Data Flattening Performance

Theorem: For $\tilde{\mathbf{X}}_i$ independent of \mathbf{X} :

$$\frac{1}{\sigma} (f_{\text{R}}(\tilde{\mathbf{X}}_i) - m_a) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1) \Rightarrow \varepsilon = Q\left(\frac{|m_a|}{\sigma}\right)$$

where $m_a = (-1)^a \|\mathbf{M}\|^2 \sqrt{\frac{n}{p}}$ and $\sigma = \sqrt{\frac{n}{p} \|\mathbf{M}\|^2 + \frac{P}{p}}$.

Tensor-based Classification

The weight tensor \mathbf{W} is a **spiked random tensor**:

$$\mathbf{W} = \sqrt{\frac{n}{p}} \bigotimes_{i=1}^k \boldsymbol{\mu}_i + \frac{1}{\sqrt{p}} \mathbf{Z}$$

with $\mathbf{Z} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{y}_i \mathbf{Z}_i$ (Universality with CLT).

Tensor-Ridge classifier is defined as:

$$f_{\text{TR}}(\tilde{\mathbf{X}}_i) = \left\langle \lambda^* \bigotimes_{i=1}^k \mathbf{u}_i^*, \tilde{\mathbf{X}}_i \right\rangle \stackrel{\mathcal{C}_1}{\leq} 0 \stackrel{\mathcal{C}_2}{\geq} 0$$

where (best rank-one approximation of \mathbf{W}):

$$(\lambda^*, \{\mathbf{u}_i^*\}_{i=1}^k) = \arg \min_{\lambda \in \mathbb{R}^+, \mathbf{u}_i \in \mathbb{S}^{p_i-1}} \left\| \mathbf{W} - \lambda \bigotimes_{i=1}^k \mathbf{u}_i \right\|_{\text{F}}^2$$

Remark: The above MLE is **NP-hard** but feasible if $\|\mathbf{M}\| \geq \mathcal{O}(P^{1/4}/p^{1/2})$.

Tensor-based Performance

Theorem: For $\tilde{\mathbf{X}}_i$ independent of \mathbf{X} :

$$\frac{1}{\sigma} (f_{\text{TR}}(\tilde{\mathbf{X}}_i) - m_a) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1) \Rightarrow \varepsilon = Q\left(\frac{|m_a|}{\sigma}\right)$$

where $m_a = (-1)^a \sigma \|\mathbf{M}\| \prod_{j=1}^k q_j(\sigma)$ and $f\left(\sigma, \|\mathbf{M}\| \sqrt{\frac{n}{p}}\right) = 0$ with q_j and f in [1].

Unsupervised Setting

- Linear clustering:** compute the **left singular vector** of:

$$\mathbf{X} = \mathbf{X}_{(k+1)} = \mathbf{y} \otimes \text{flatten}(\mathbf{M}) + \mathbf{Z} \in \mathbb{R}^{n \times P} \rightarrow \hat{\mathbf{y}}$$

- Tensor-based clustering:** compute the **best rank-one approximation** of:

$$\mathbf{X} = \mathbf{M} \otimes \mathbf{y} + \mathbf{Z} \in \mathbb{R}^{p_1 \times \cdots \times p_k \times n} \rightarrow \hat{\mathbf{y}}$$

Theorem: The estimated class for \mathbf{X}_i is given by $\text{sign}(\hat{y}_i)$:

$$\frac{1}{\sqrt{1-\alpha^2}} (\sqrt{n} \hat{y}_i - \alpha y_i) \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1) \Rightarrow \varepsilon = Q\left(\frac{\alpha}{\sqrt{1-\alpha^2}}\right)$$

- Linear:** $\alpha = \kappa \left(\|\mathbf{M}\| \sqrt{\frac{n}{P+n}, \frac{n}{P+n}} \right)^{-1}$ with κ in [1].
- Tensor:** $\alpha = q_{k+1}(\lambda^*)$ with $f\left(\lambda^*, \|\mathbf{M}\| \sqrt{\frac{n}{P+n}}\right) = 0$ (λ^* spectral norm of \mathbf{X}).

Simulations

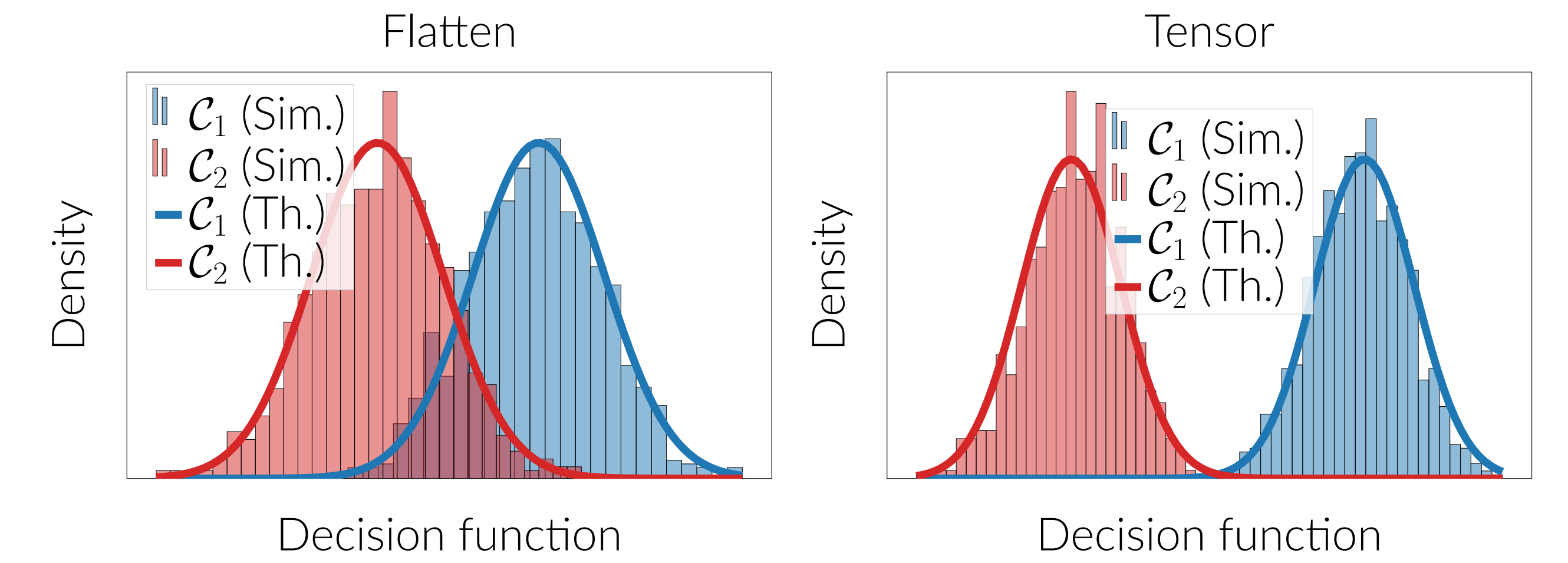
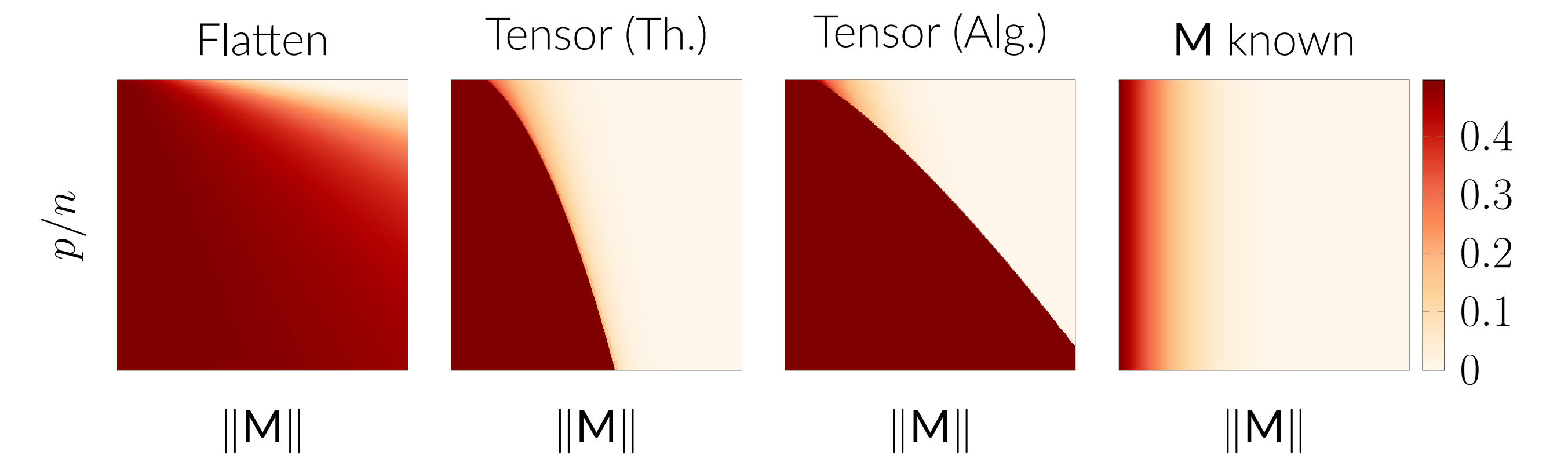


Figure 1. $n = 200$, tensor shape (15, 30, 20) and $\|\mathbf{M}\| = 3$.



Linear ($\varepsilon = 6.3\%$) Tensor ($\varepsilon = 0.1\%$)

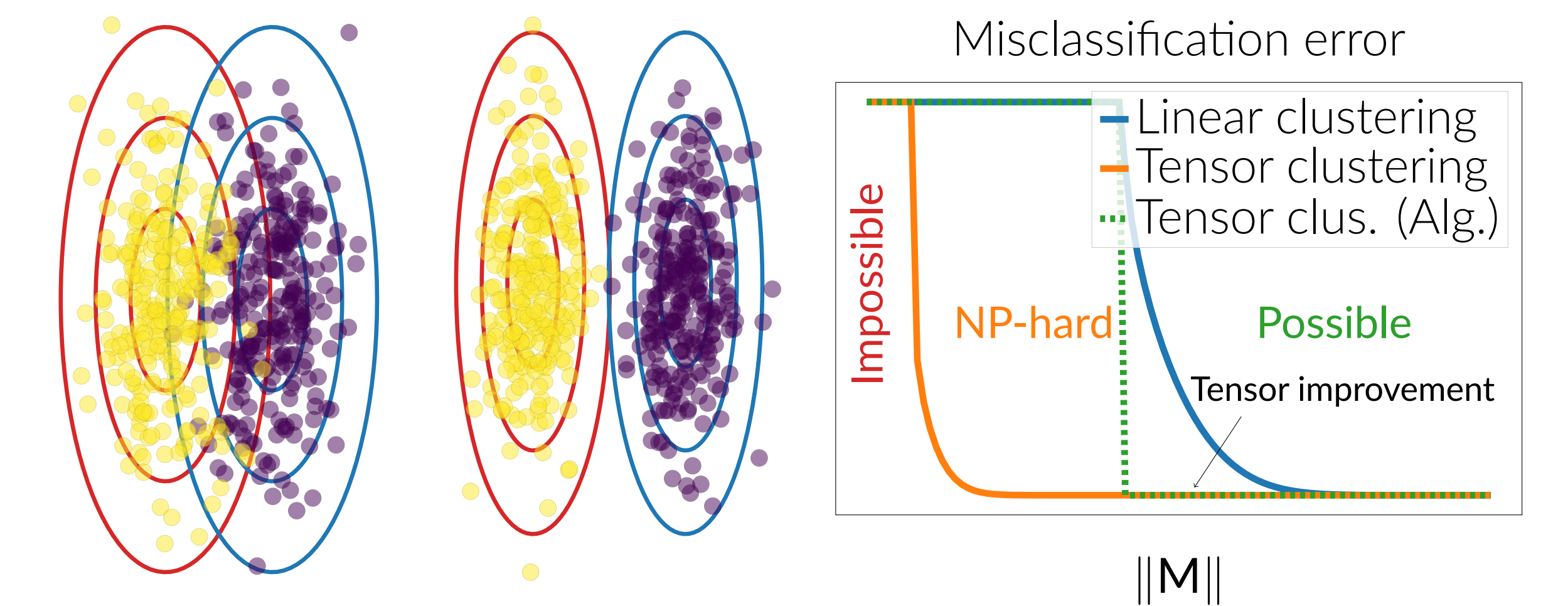


Figure 2. $n = 200$, tensor shape (15, 30, 20) and $\|\mathbf{M}\| = 3$.

Conclusion

- This work analyzes learning from **low-rank tensor** data and shows **performance gains**.
- It applies **random tensor theory** to evaluate simple learning methods.
- This paves the way for **improving machine learning** algorithms for tensor-structured data.

[1] MEA.Seddik, M.Guillaud, R.Couillet, "When Random Tensors meet Random Matrices", *Annals of Applied Probability* 2023.

