# Information - computation gap in High-Dimensional clustering.

Bertrand Even[1], Christophe Giraud[1], Nicolas Verzelen[2]
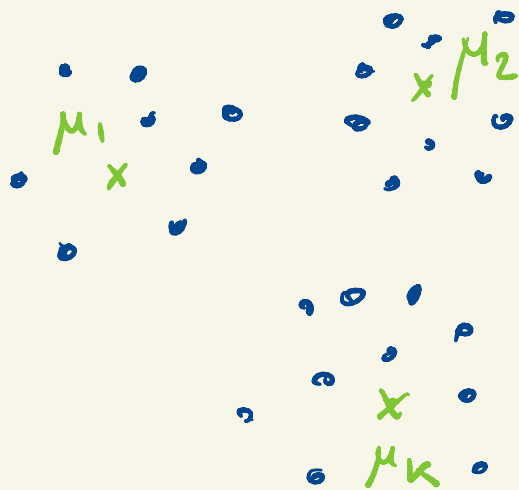
① Labo de maths d'Orsay, université Paris Saclay
② MISTEA, INRAE Montpellier

- clustering = partitioning a set of points into k groups

- Model: we observe $X_1, \ldots, X_m \in \mathbb{R}^d$

  $\exists\, G^*$ partition of $\{1, \ldots, m\}$ in $k$ groups:

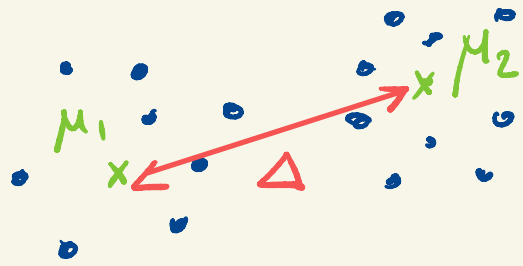  $$X_i \sim \mathcal{N}(\mu_k, \sigma^2 I_d) \quad \forall i \in G^*_k$$

  w.l.o.g

- Goal: recover $G^*$

  ↗ exactly: $\hat{G} = G^*$

  ↠ partially: $err(\hat{G}) :=$ proportion of points well classified
  $$\geqslant c > 0.$$

- Separation: $\Delta^2 = \min_{k \neq \ell} \| \mu_k - \mu_\ell \|^2$

- Assumption: $|G_k^*| \asymp \dfrac{m}{K}$

- Our focus:
  - → high dimension: $d \geq m$
  - → condition on $\Delta$ to recover $\begin{matrix} \text{exactly} \\ \text{partially} \end{matrix}$ $G^*$ ↗ without computational constraints
    
    ↘ with computational constraints

- Plan
  ① Detour on High-Dimensional classification
  ② Information - Computation gap in HD clustering
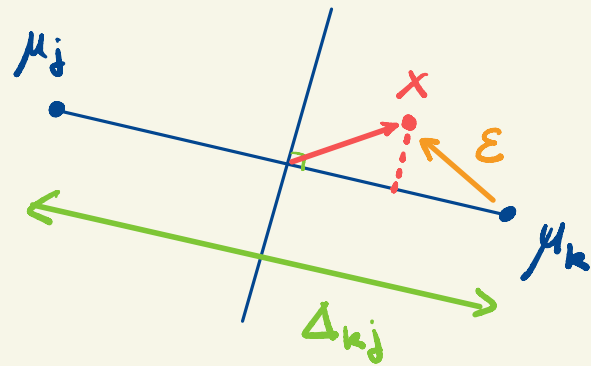  ③ Proving computational barriers

① High-dimensional classification

a/ with $\mu_1, ..., \mu_K$ known



• consider 2 means $\mu_j, \mu_k$ and $x \in \mathbb{R}^d$

$$S_{kj}(x) = \left\langle x - \frac{\mu_j + \mu_k}{2}, \mu_k - \mu_j \right\rangle$$

if $x = \mu_k + \varepsilon$    $\overset{\mathcal{N}(0, I_d)}{\underset{\nearrow}{=}}$    $\left\langle \frac{\mu_k - \mu_j}{2} + \varepsilon, \mu_k - \mu_j \right\rangle = \frac{1}{2} \Delta_{kj}^2 + \Delta_{kj} \mathcal{N}(0,1)$

so    $\mathbb{P}_{X \sim \mathcal{N}(\mu_k, I_d)} \left[ \exists j : S_{kj}(x) < 0 \right] = \sum_j \mathbb{P} \left[ \mathcal{N}(0,1) \leqslant -\frac{1}{2} \Delta_{kj} \right]$

$$\leqslant K e^{-\Delta^2/8}$$

- Setting $\hat{k}(x) = \underset{k'}{\text{argmax}} \ \underset{j : j \neq k'}{\text{min}} \ S_{k'j}(x)$

  $\mathbb{P}[\text{1 point misclassified}] \leq K e^{-\Delta^2/8}$

  $\mathbb{P}[\text{at least 1 out of } n \text{ points misclassified}] \leq n K e^{-\Delta^2/8}$

- So if $\mu_1, \ldots, \mu_K$ known, we need

$$\Delta^2 \gtrsim \begin{matrix} \log(K) & \text{for partial recovery} \\ \log(n) & \text{for exact recovery.} \end{matrix}$$

b/ with $\mu_1, \ldots, \mu_K$ unknown :

we rely on estimators $\hat{\mu}_1, \ldots, \hat{\mu}_K$ computed with sample size $m = \frac{M}{K}$:

$$\hat{\mu}_k = \mu_k + \frac{1}{\sqrt{m}} \, \zeta_k \quad \leftarrow \mathcal{N}(0, I_d)$$

$$\rightsquigarrow \hat{S}_{kj}(x) = \left\langle x - \frac{\hat{\mu}_k + \hat{\mu}_j}{2}, \; \hat{\mu}_k - \hat{\mu}_j \right\rangle$$

$$\hat{S}_{kj}(\mu_k + \varepsilon) = \left\langle \frac{\mu_k - \mu_j}{2} + \varepsilon - \frac{\zeta_k + \zeta_j}{\sqrt{m}}, \; \mu_k - \mu_j + \frac{\zeta_k - \zeta_j}{\sqrt{m}} \right\rangle$$

$$= \frac{1}{2} \Delta_{kj}^2 + \left(1 + O\left(\frac{1}{\sqrt{m}}\right)\right) \left[ \Delta_{kj} \, \mathcal{N}(0,1) + \frac{\langle \varepsilon, \zeta_k - \zeta_j \rangle}{\sqrt{m}} \right]$$

$$\underset{\uparrow}{\geqslant} 0$$

if $\mathcal{N}(0,1) \ni -\Delta_{jk}$ and $\mathcal{N}'(0,1) \ni -\sqrt{\frac{dK}{n}} \, \Delta_{kj}^2$

$$\simeq \sqrt{\frac{d}{m}} \, \mathcal{N}'(0,1)$$

So $\quad \mathbb{P}[\,1 \text{ point misclassified}] \leq K \exp\left(-c\, \Delta^2 \wedge \frac{m\Delta^4}{Kd}\right)$

$\mathbb{P}[\text{at least 1 out of } m \text{ points misclass.}] \leq mK \exp\left(-c\, \Delta^2 \wedge \frac{m\Delta^4}{Kd}\right)$

so, with estimated means we need

$$\Delta^2 \overset{(*)}{\gtrsim} \log(\square) \vee \sqrt{\frac{dK}{m} \log(\square)}$$

$\underbrace{\phantom{\sqrt{\frac{dK}{m} \log(\square)}}}$

with $\square = \dfrac{K}{m}$ for $\begin{array}{l}\text{partial}\\ \text{exact}\end{array}$ recovery

$\left|\begin{array}{l}\text{curse of dimensionality}\\ \text{for } d \gtrsim \frac{m}{K} \log(\square)\end{array}\right.$

. Is $(*)$ the minimal separation for clustering?

## ② Information-Computation gap in HD clustering

### a) Without computational constraints

**Theorem : EGV '24**

Partial / exact recovery minimax impossible if

$$\Delta^2 \leq \log(\square) \vee \sqrt{\frac{dK}{n} \log(\square)}$$

and possible with exact Kmeans if

$$\Delta^2 \overset{(*)}{\gtrsim} \log(\square) \vee \sqrt{\frac{dK}{n} \log(\square)}$$

**Remarks:**
→ $\Delta^2 \gtrsim \log(K)$ already known for d small (Kwon and Caramis CoLT 2020)
→ for K = 2, tight rate for exact recovery in Ndaoud (AOS 2022)

# b) with computational constraints

. Is clustering possible in polynomial time when (*) holds?

    $\rightarrow$ in low dimension: (kind of) yes   [Liu and Li 2022]

    $\rightarrow$ in high-dimension?

**Theorem** (informal) EGV '24

  For $d \geq m$, under computational constraints (specified later)

$$\Delta^2 \gtrsim \log^{\beta} \sqrt{\frac{d\,K^2}{m}} \wedge \sqrt{d}$$

  is required and enough for partial recovery

$\rightsquigarrow$ information - computation gap

# Remarks

$\rightarrow$ when $\Delta^2 \gtrsim \overset{\log}{\sqrt{d}}$ : recovery possible with hierarchical clustering
                                                                with single linkage

$\rightarrow$ when $\Delta^2 \gtrsim \overset{\log}{\sqrt{\frac{dK^2}{n}}}$ : recovery possible with SDP relaxation of Kmeans
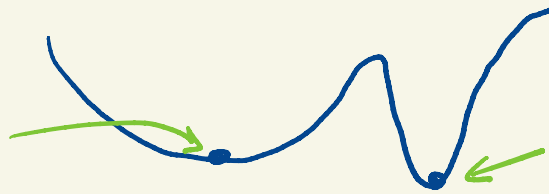                                                                    (G.V. '19)

$\rightarrow$ computational gap conjectured in Lesieur et al. (2016)
   based on the computation of fixed points of State Evolution of AMP
   $\leftrightarrow$ local minima of Bethe free energy
   $\hookrightarrow$ replica theory predicts that multiple minima $\leftrightarrow$ gap

local minima achieved
by "local" search with non
informative initialization

global optima
hard to achieve

③ **Proving computational barriers**

→ **Worst case complexity:**
proving that a problem is NP-hard
(e.g. minimizing Kmeans exactly)
↳ not our case, as we consider some
    random instances with separation

→ **Reduction:** to a problem that we
believe to be hard
  (e.g. planted clique)

→ **Computation model:** prove that
some classes of algorithms fail.
Ex: - SQ algorithms
     - SoS algorithms
     - local algorithms ($\cap$ $c$ $\cap$ $c$)
       (landscape analysis)
     - low degree polynomials
       ↳ our choice here

Rmk: there are connections between
these notions and also with tools from
statistical physics (replica symmetry
and cavity method).

# Recipe 1 : from clustering to estimation

combinatorial ⇝ continuous

- Partnership matrix:

  $\Pi_{ij}^{G} := \mathbb{1}_{i \overset{G}{\sim} j} \quad \in \{0,1\}^{m \times m}$

  $\Pi^{*} := \Pi^{G^{*}}$

- Estimation error

  $R(\hat{\Pi}) := \frac{1}{m(m-1)} \sum_{i \neq j} (\hat{\Pi}_{ij} - \Pi_{ij})^2$

- Relation to clustering

  $R(\Pi^{\hat{G}}) \leq 2 \, \text{err}(\hat{G})$

proportion of misclustered points

So

$$\inf_{\hat{\Pi} \text{ poly-time}} R(\hat{\Pi}) \leq 2 \inf_{\hat{G} \text{ poly-time}} \text{err}(\hat{G})$$

# Recipe 2 : introducing a generative model

$\rightarrow k_1, \ldots, k_m \overset{iid}{\sim} \mathcal{U}(\{1, -, K\})$

and $G_k^{*} = \{i : k_i = k\}$

$\rightarrow \mu_1, \ldots, \mu_K \overset{iid}{\sim} \mathcal{U}\left[\left\{-\frac{\Delta}{\sqrt{d}}, +\frac{\Delta}{\sqrt{d}}\right\}^d\right]$

$\Rightarrow \|\mu_j - \mu_k\|^2 \asymp \Delta^2$

We can investigate $\mathbb{E}[R(\hat{\Pi})]$

expectation relative to prior + data generation

# Low degree polynomials : (Schramm & Wein 2022)

We restrict to $\hat{M}$ s.t.

$$\hat{M}_{ij} = f_{ij}(x) \quad \text{with} \quad f_{ij} \in \mathbb{R}_D[x]$$

$$D = O(\log(m))$$

$\rightsquigarrow$ approximate spectral

AMP

etc...

The goal: lower bound

$$\text{MMSE}_D := \inf_{f_{ij} \in \mathbb{R}_D[x]} \mathbb{E}\left[ R(f(x)) \right]$$

for $D \simeq \log(m)$.

---

## Theorem EGV'24

if $\Delta^2 \overset{\log^6}{\lesssim} \sqrt{\frac{dK^2}{m}} \wedge \sqrt{d}$ then

$$\text{MMSE}_{O(\log(m))} = \frac{1}{K} - \frac{1+o(1)}{K^2}$$

Remark : $\tilde{M}_{ij} = \frac{1}{K}$ for $i \neq j$

fulfills

$$\mathbb{E}[R(\tilde{M})] = \frac{1}{K} - \frac{1}{K^2}$$

## Sketch of proof :

① focusing on a single entry

② relating MMSE to cumulants $\left(\begin{array}{c}\text{Schramm}\\ \text{& Wein 22}\end{array}\right)$

③ bounding cumulants (technical)

① Since the $\text{MMSE}_D$ optimisation problem is separable

$$\text{MMSE}_D = \inf_{g \in \mathbb{R}_D[x]} \mathbb{E}\left[ \left( g(x) - \underbrace{\Pi^*_{12}}_{} \right)^2 \right] \underset{=:m}{}$$

② Relating MMSE to cumulants

(Schramm & Wein '22)

$$\text{MMSE}_D = \| m - P_D\, m \|^2_{L^2} = \| m \|^2_{L^2} - \| P_D m \|^2_{L^2}$$

$$= \| m \|^2_{L^2} - \underbrace{\left[ \sup_{g \in \mathbb{R}_D(x)} \frac{\langle m, g(x) \rangle_{L^2}^2}{\| g(x) \|_{L^2}} \right]^2}_{=: \text{Corr}^2_D}$$

$$= \frac{1}{K} \qquad =: \text{Corr}^2_D$$

Below $\quad X = Z + E$

$m \times d \quad \uparrow \quad \nwarrow$

$\mathbb{E}(x) \quad \varepsilon_{ij} \overset{iid}{\sim} \mathcal{N}(0,1)$

---

Lemma (SW'22, translated in our setting )

$$\text{Corr}^2_D \leq \sum_{\substack{\alpha \in \mathbb{N}^{m \times d} \\ |\alpha| \leq D}} \frac{\mathcal{K}^2_\alpha}{\alpha!} \qquad (**)$$

where $\mathcal{K}_\alpha = \text{cumulant}(m, \dots, \underbrace{z_{ij}, \dots, z_{ij}, \dots}_{d_{ij} \text{ times}})$

$\cdot\ \alpha! = \prod_{ij} \alpha_{ij}!$

Proof:

· Inequality from Jensen

$$\mathbb{E}\left[ g(x)^2 \right] \geq \mathbb{E}_E\left[ \mathbb{E}_Z[ g(z+E)]^2 \right]$$

· expansion on Hermite polynomials

· Linear algebra

· recognize recursion of cumulants

□

How can we exploit (**) ?

(i) exploit the property
$$X \perp\!\!\!\perp Y \implies \text{cumulant}(X,Y)=0$$
to detect the $\kappa_\alpha = 0$ and
to prune $\sum_\alpha \dfrac{\kappa_\alpha^2}{\alpha!}$
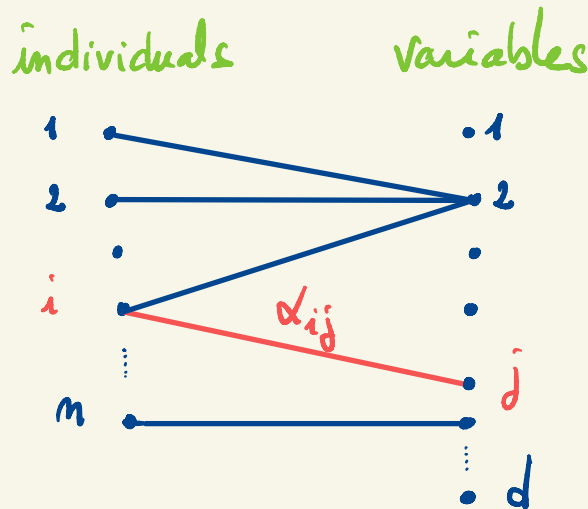
(ii) relate cumulants to moments
$$\kappa_\alpha \overset{(***)}{=} \mathbb{E}[m\, Z^\alpha] - \sum_{\beta \lneq \alpha} \kappa_\beta \binom{\alpha}{\beta} \mathbb{E}[Z^{\alpha-\beta}]$$

(iii) upper-bound the moments
$$\mathbb{E}[m Z^\alpha] \quad \text{and} \quad \mathbb{E}[Z^{\alpha-\beta}]$$

(iv) bound $\kappa_\alpha$ by induction from (***)

(i) represent $\alpha \in \mathbb{N}^{m \times d}$ as weighted bipartite graph $G_\alpha$



Lemma: If $\kappa_\alpha \neq 0$ then
- $G_\alpha^+$ connex
- individuals 1 and 2 $\in G_\alpha^+$
- each variable $j \in G_\alpha^+$ connected to at least 2 individuals.

(iii) Bounds on the moments ─────────────────

Define: $C_\alpha :=$ # connected components of $G_\alpha^+$

$l_\alpha =$ # nodes of $G_\alpha^+$

Then $\mathbb{E}\left[ \textcolor{red}{m} Z^\alpha \right] \leqslant \left( \dfrac{\Delta}{\sqrt{d}} \right)^{|\alpha|_1} \left( \dfrac{|\alpha|_1^{|\alpha|_1}}{K^{l_\alpha - \frac{1}{2}|\alpha|_1 - C_\alpha}} \wedge \dfrac{1}{\textcolor{red}{K}} \right)$

Idea: reminder: $k_1, \ldots, k_m \overset{iid}{\sim} \mathcal{U}\{1, \ldots, K\}$

$\mu_1, \ldots, \mu_K \overset{iid}{\sim} \mathcal{U}\left\{ -\dfrac{\Delta}{\sqrt{d}}, \dfrac{\Delta}{\sqrt{d}} \right\}^d$

and $Z^\alpha = \prod_{i,j} \mu_{k_i d}^{\alpha_{ij}} = \prod_{k,j} \mu_{kj}^{\sum_{i \in G_u} \alpha_{ij}}$

So $\mathbb{E}\left[ Z^\alpha \mid k_1 \ldots k_m \right] = \left( \dfrac{\Delta}{\sqrt{d}} \right)^{|\alpha|_1} \times \begin{cases} 1 & \text{if } \sum_{i \in G_u} \alpha_{ij} \text{ even } \forall k, j \\ 0 & \text{otherwise} \end{cases}$

Hence $\mathbb{E}[Z^\alpha] = \left(\dfrac{\Delta}{\sqrt{d}}\right)^{|\alpha|_1} \underset{G}{\mathbb{P}}\left[\underset{i \in G_k}{\sum} \alpha_{ij} \text{ even } \forall k,j\right]$

<span style="color:red">$\underbrace{\hspace{5cm}}$</span>

<span style="color:red">where delicate combinatrics kicks in ...</span>

## (iv) Bound $K_\alpha$ by induction

· $k_0 = \mathbb{E}[m] = \dfrac{1}{K}$

· induction: from (***)

$$K_\alpha \leqslant \left(\frac{\Delta}{\sqrt{d}}\right)^{|\alpha|_1} (1+|\alpha|_1)^{|\alpha|_1} \left[\frac{|\alpha|_1^{|\alpha|_1}}{K^{l_\alpha - \frac{1}{2}|\alpha|_1 - 1}} \wedge \frac{1}{K}\right]$$

## (v) Conclusion

$$\text{Cor}_D^2 \leqslant \underset{|\alpha|_1 \leqslant D}{\sum} \frac{K_\alpha^2}{\alpha!} \leqslant \frac{1+o(1)}{K^2}$$

<span style="color:green">$\uparrow$</span>

<span style="color:green">for $\Delta^2 \overset{\log \beta}{\leqslant} \sqrt{\dfrac{dk^2}{m}} \wedge \sqrt{d}$ and $D = O(\log(n))$</span>

# Take Home Message :

- Low degree polynomials are handy for proving the existence of computational barriers, at the price of spurious log factors.

- Minimal information separation:

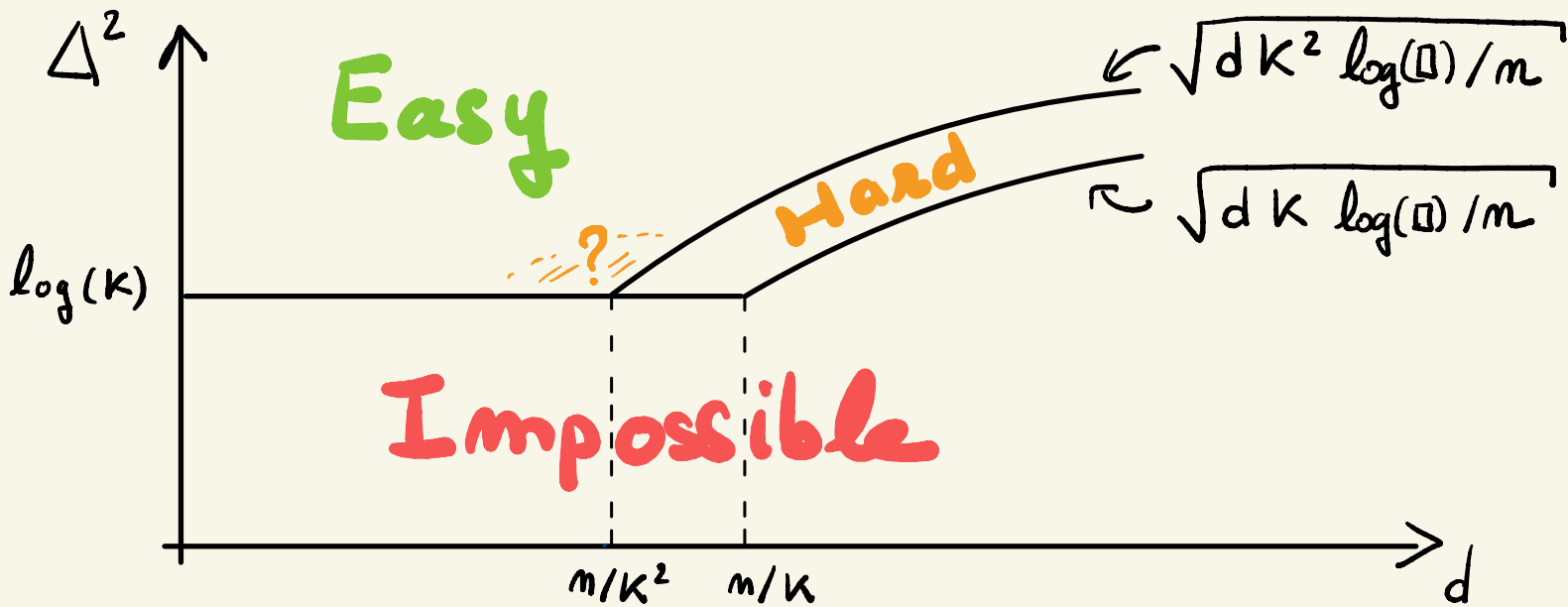$$\Delta_I^2 \asymp \log(\square) \vee \sqrt{\frac{dK \log(\square)}{n}}$$

- Minimal computational separation: (conjectured)

$$\Delta_c^2 \asymp \log(\square) \vee \sqrt{\left(\frac{K^2}{n} \wedge 1\right) d \log(\square)}$$

proved for $d \geq n$
or $d$ small

in progress for
$$\frac{n}{K^2} \leq d \leq n$$

- What is special with $\Delta^2 \overset{log}{\gtrsim} \sqrt{\dfrac{d K^2}{m}}$ ?

↝ related to BBP transition for "isotropic" $\mu_1 \ldots \mu_K$ :

$m \Delta^4 \gtrsim d K^2$   is where   $K$ largest eigenvalues of the Gram matrix escape of the bulk.

. **Remarkable feature**: the Information-Computation gap disappear in an active setting.

active setting: we can sample each point multiple time.

with a total budget of $L$ observations, minimal separation is

$$\Delta_*^2 \asymp \frac{m}{L} \left[ \log(m) \vee \sqrt{\frac{dk}{m} \log(m)} \right] \Big\}$$

Victor Thuot, Alexandra
Carpentier, C.G., Nicolas
Verzelen 2024.

and no computationnal barrier

**why?**
⟶ we can collect localized information