# PART I: ADVERSARIAL CLUSTERS

**Impossible!**



cluster

**①** Restrict clustering to $U$

**②** Replace TV with earth-mover (EMD)

Domain → metric space

**③** Queries

**④** CLUSTER-REJECT

not allowed!

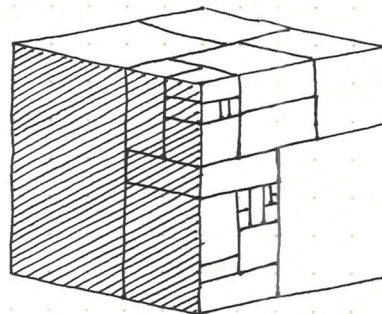$U$: Universe

$G \subseteq U$: "Good" clusterings

E.g.

Convex

depends on distribution

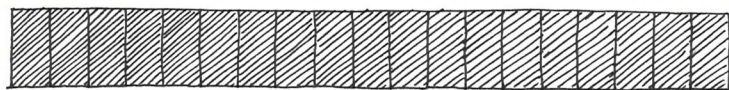High-probability of low diameter boxes, decision trees

RESULT: learning cells is not necessary

# PART II: RANDOM CLUSTERS

**Motivations:**

- Environmental randomness
- Randomized classifier training

**We study:** Testing uniformity,

$\underbrace{\text{cluster}}$

$\overset{\uparrow}{\rho} := \mathbb{P}[\text{separator}]$

$\underbrace{\qquad}_{\text{resolution}}$

**Standard:**

$X$: histogram

Analyze $X^{\mathsf{T}} \underset{\uparrow}{I} X - \|X\|_1$

**Naïve (known clusters):**

$$O\left( \frac{\sqrt{n}}{\rho^{3/2} \varepsilon^2} \right)$$

**Without queries,** $\left( \rho \geq \widetilde{\Omega}\left( n^{-1/5} \varepsilon^{-4/5} \right) \right)$

$$\widetilde{O}\left( \frac{\sqrt{n}}{\rho^{3/2} \varepsilon^2} \right)$$

**With queries**

$$O\left( \frac{\sqrt{n}}{\rho \, \varepsilon^2} \right) \quad \text{Use} \quad [VV'17]$$

**Now:**

Analyze $X^{\mathsf{T}} \Phi X - \|X\|_1$