

# A Random CSP with Connections to Discrepancy Theory and Randomized Trials

**Eren C. Kızıldağ** (University of Illinois Urbana-Champaign)

International Symposium on Information Theory

July 12, 2024

# Overview of Model

Fix  $d \in \mathbb{N}$  and  $\mathbf{c} = (c_1, \dots, c_d) \in \mathbb{R}_+^d$ . Generate **iid**  $X_1, \dots, X_d \sim \mathcal{N}(0, I_n)$ . Define

$$\mathcal{F}(\mathbf{c}) = \left\{ \boldsymbol{\sigma} \in \{-1, 1\}^n : |\langle \boldsymbol{\sigma}, X_i \rangle| \leq \sqrt{n} 2^{-c_i n}, \forall i \right\}.$$

**Focus:** Non-proportional regime,  $n \rightarrow \infty, d = O(1)$ . **Random CSP**

**Questions:** When is  $\mathcal{F}(\mathbf{c}) \neq \emptyset$ ? How does its 'geometry' look like?

**Today:** Sharp Phase Transition for  $\{\mathcal{F}(\mathbf{c}) \neq \emptyset\}$ . Landscape of  $\mathcal{F}(\mathbf{c})$

# Motivation

## Discrepancy Theory

- Given  $M \in \mathbb{R}^{d \times n}$ , compute/bound its **discrepancy**  $\mathcal{D}(M) := \min_{\sigma \in \{\pm 1\}^n} \|M\sigma\|_\infty$ .
- Note that  $\mathcal{D}(M) \leq B \iff \exists \sigma \in \{\pm 1\}^n : |\langle \sigma, X_i \rangle| \leq B, \forall i$  ( $X_i$  are rows of  $M$ ).
- **Worst-case & random  $M$ . Existential & Algorithmic** results.  
[Spencer 85, Karmarkar-Karp-Lueker-Odlyzko 86, Matousek 99, Chazelle 00, Bansal 10, Lovett-Meka 15]

**Proportional Regime:** For  $d = \Theta(n)$ , and  $M \in \mathbb{R}^{d \times n}$  with  $\mathcal{N}(0, 1)$  entries, whp

$$\mathcal{D}(M) = f(\alpha)\sqrt{n}(1 + o_n(1)), \quad \text{where } \alpha = d/n$$

for explicit  $f(\cdot)$ . [Perkins-Xu 21, Abbe-Li-Sly 21]

# Motivation: Discrepancy and Random CSPs

## Symmetric Binary Perceptron

Fix  $\kappa > 0$  and consider random  $M \in \mathbb{R}^{\alpha n \times n}$ . What is the largest  $\alpha > 0$  for which a

$$\sigma \in \{\pm 1\}^n : \|M\sigma\|_\infty \leq \kappa\sqrt{n}$$

exists whp? When do efficient search algs work?

- **Perceptron:** Model for pattern storage. Popular in **probability, stat phys, statistics** communities [Cover 65, Hopfield 82, Krauth-Mézard 89, Talagrand 99, 10, Franz-Parisi 16, Candès-Sur 20, Perkins-Xu 21, Abbe-Li-Sly 21, 22, Montanari-Zhong-Zhou 21, Sah-Sawhney 23, Nakajima-Sun 23, Barbier-El Alaoui-Krzakala-Zdeborová 23, Gamarnik-K.-Perkins-Xu 22, 23, K.-Wakhare 23]
- **Dual of discrepancy:** Fix  $\alpha > 0$ , seek smallest  $\kappa > 0$  s.t.  $\sigma \in \{\pm 1\}^n : \|M\sigma\|_\infty \leq \kappa\sqrt{n}$  exists

# Motivation: Our Model

$M \in \mathbb{R}^{d \times n}$  with iid rows  $X_1, \dots, X_d \sim \mathcal{N}(0, I_n)$ ,  $\mathbf{c} = (c_1, \dots, c_d) \in \mathbb{R}_+^d$ .

$$\mathcal{F}(\mathbf{c}) = \{\boldsymbol{\sigma} \in \{-1, 1\}^n : |\langle \boldsymbol{\sigma}, X_i \rangle| \leq \sqrt{n}2^{-c_i n}, \forall i\}.$$

$$\boldsymbol{\sigma} \in \mathcal{F}(\mathbf{c}) \iff |M\boldsymbol{\sigma}| \leq \begin{pmatrix} \sqrt{n}2^{-c_1 n} \\ \vdots \\ \sqrt{n}2^{-c_d n} \end{pmatrix}.$$

- Dual of discrepancy. **Non-proportional** regime,  $d = O_n(1)$ . **Non-uniform** constraints
- $2^{-n}$  scaling:  $\min_{\boldsymbol{\sigma} \in \{\pm 1\}^n} \|M\boldsymbol{\sigma}\|_\infty = \sqrt{n}2^{-\Omega(n/d)} = \sqrt{n}2^{-\Omega(n)}$  for  $d = O_n(1)$   
[Karmarkar-Karp-Lueker-Odlyzko 86, Costello 09, Turner-Meka-Rigollet 20]

# Motivation: Randomized Controlled Trials

- Gold standard for clinical trials (drug/vaccine)
- $n$  individuals, covariates  $Y_1, \dots, Y_n \in \mathbb{R}^d$  (columns of  $M \in \mathbb{R}^{d \times n}$ ).  $n \gg d$ .
- Split into balanced **treatment** & **control**: for thresholds  $t_1, \dots, t_d$  and features  $j \in \{1, \dots, d\}$

$$D_j := \left| \sum_{i:\sigma(i)=+1} Y_i(j) - \sum_{i:\sigma(i)=-1} Y_i(j) \right| \leq t_j, \quad \forall j \in [d].$$

- Any solution to **CSP** gives a valid design:  $\sigma \in \mathcal{F}(\mathbf{c}) \Leftrightarrow D_j \leq t_j$ , for  $t_j = \sqrt{n}2^{-c_j n}$ .

# Main Results: A Sharp Phase Transition

For  $d \in \mathbb{N}$ ,  $\mathbf{c} = (c_1, \dots, c_d) \in \mathbb{R}_+^d$  and iid  $X_1, \dots, X_d \sim \mathcal{N}(0, I_n)$

$$\mathcal{F}(\mathbf{c}) = \{\boldsymbol{\sigma} \in \{-1, 1\}^n : |\langle \boldsymbol{\sigma}, X_i \rangle| \leq \sqrt{n} 2^{-c_i n}, \forall i\}$$

Theorem (K., 2024)

$$\lim_{n \rightarrow \infty} \mathbb{P}[\mathcal{F}(\mathbf{c}) \neq \emptyset] = \begin{cases} 0, & \text{if } \|\mathbf{c}\|_1 > 1 \\ 1, & \text{if } \|\mathbf{c}\|_1 < 1. \end{cases}$$

- [Costello 09]:  $\min_{\boldsymbol{\sigma} \in \{\pm 1\}^n} \|M\boldsymbol{\sigma}\|_\infty \sim \sqrt{n} 2^{-n/d}$  for  $M \in \mathbb{R}^{d \times n}$ ,  $d = O(1)$ . ( $c_i \sim 1/d$ )
- Proof based on the **first moment method** and the **second moment method**

## Proof Sketch: $\|\mathbf{c}\|_1 > 1$

Let  $T = |\mathcal{F}(\mathbf{c})|$ . Our proof is based on the **moment method**.

### First Moment Method for $\|\mathbf{c}\|_1 > 1$

Observe that  $n^{-1/2}\langle \boldsymbol{\sigma}, \mathbf{X}_i \rangle \sim \mathcal{N}(0, 1)$ ,  $1 \leq i \leq d$  are **iid**. So,

$$\mathbb{P}[\boldsymbol{\sigma} \in \mathcal{F}(\mathbf{c})] = \prod_{1 \leq i \leq d} \mathbb{P}[|\mathcal{N}(0, 1)| \leq 2^{-c_i n}] \sim 2^{-\|\mathbf{c}\|_1 n}.$$

Using **Markov's inequality**, we have that for  $\|\mathbf{c}\|_1 > 1$

$$\mathbb{P}[T \geq 1] \leq \mathbb{E}[T] \sim 2^{n(1-\|\mathbf{c}\|_1)} = 2^{-\Theta(n)}.$$

Hence,  $\mathcal{F}(\mathbf{c}) = \emptyset$  whp for  $\|\mathbf{c}\|_1 > 1$ .



**Proof Idea:**  $\|\mathbf{c}\|_1 < 1$

## Paley-Zygmund Inequality (Second Moment Method)

Let  $T \geq 0$  be a rv and  $\theta \in [0, 1]$ . Then,

$$\mathbb{P}[T > \theta \mathbb{E}[T]] \geq (1 - \theta)^2 \frac{\mathbb{E}[T]^2}{\mathbb{E}[T^2]}.$$

Suppose  $T \in \mathbb{Z}$  and  $\mathbb{E}[T^2] = (1 + o(1))\mathbb{E}[T]^2$ . Taking  $\theta = 0$  yields  $T \geq 1$  whp.

$$\mathbb{E}[T^2] = \underbrace{\sum_{(\sigma, \sigma') \in \mathcal{T}_1} \mathbb{P}[\sigma, \sigma' \in \mathcal{F}(\mathbf{c})]}_{:= \Sigma_1} + \underbrace{\sum_{(\sigma, \sigma') \in \mathcal{T}_2} \mathbb{P}[\sigma, \sigma' \in \mathcal{F}(\mathbf{c})]}_{:= \Sigma_2},$$

where for an arbitrary  $\epsilon > 0$ ,

$$\mathcal{T}_1 = \left\{ (\sigma, \sigma') : \frac{1}{n} \langle \sigma, \sigma' \rangle \in [-\epsilon, \epsilon] \right\} \quad \text{and} \quad \mathcal{T}_2 = \left\{ (\sigma, \sigma') : \frac{1}{n} |\langle \sigma, \sigma' \rangle| > \epsilon \right\}$$

## Proof Sketch: $\|\mathbf{c}\|_1 < 1$

- Show  $\Sigma_2 \leq \mathbb{E}[T]^2 e^{-\Theta(n)}$ . For  $\Sigma_1$ , take  $(\boldsymbol{\sigma}, \boldsymbol{\sigma}') \in \mathcal{T}_1$ . Then,

$$\left( \frac{1}{\sqrt{n}} \langle \boldsymbol{\sigma}, \mathbf{X}_i \rangle, \frac{1}{\sqrt{n}} \langle \boldsymbol{\sigma}', \mathbf{X}_i \rangle \right) \sim \mathcal{N} \left( \mathbf{0}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right), \quad \text{where } \rho = \frac{1}{n} \langle \boldsymbol{\sigma}, \boldsymbol{\sigma}' \rangle \in [-\epsilon, \epsilon]$$

- Ignoring absolute constants (not depending on  $\epsilon$ )

$$\mathbb{P}[\boldsymbol{\sigma}, \boldsymbol{\sigma}' \in \mathcal{F}(\mathbf{c})] \leq (1 - \epsilon^2)^{-\frac{d}{2}} 2^{-2\|\mathbf{c}\|_1 n}$$

- This gives  $\Sigma_1 \leq \mathbb{E}[T]^2 (1 - \epsilon^2)^{-\frac{d}{2}}$ . Using  $d = O(1)$ ,

$$\liminf_{n \rightarrow \infty} \mathbb{P}[T \geq 1] \geq \liminf_{n \rightarrow \infty} \frac{\mathbb{E}[T]^2}{\mathbb{E}[T^2]} \geq (1 - \epsilon^2)^{\frac{d}{2}}.$$

- Conclude by  $\epsilon \rightarrow 0$

# Main Results: Solution Space Geometry

$$\mathcal{F}(\mathbf{c}) = \{\boldsymbol{\sigma} \in \{-1, 1\}^n : |\langle \boldsymbol{\sigma}, \mathbf{X}_i \rangle| \leq \sqrt{n} 2^{-c_i n}, \forall i \in [d]\}.$$

## Theorem (K., 2024)

Let  $\|\mathbf{c}\|_1 > \frac{1}{2}$ . There exists  $\beta^* \in (0, 1)$  such that whp

$$\min_{\boldsymbol{\sigma}, \boldsymbol{\sigma}' \in \mathcal{F}(\mathbf{c}), \boldsymbol{\sigma} \neq \boldsymbol{\sigma}'} d_H(\boldsymbol{\sigma}, \boldsymbol{\sigma}') \geq \beta^* n.$$

- **Solutions are Isolated:** If  $\|\mathbf{c}\|_1 > \frac{1}{2}$  then any  $(\boldsymbol{\sigma}, \boldsymbol{\sigma}')$  are  $\Omega(n)$  apart.
- Proof via **first moment method:** let  $T$  count  $\#(\boldsymbol{\sigma}, \boldsymbol{\sigma}') : d_H(\boldsymbol{\sigma}, \boldsymbol{\sigma}') \leq \beta^* n$ , show  $\mathbb{E}[T] = o(1)$ .
- Suggests **algorithmic hardness**

[Achlioptas-Ricci Tersenghi 06, Achlioptas-Coja Oghlan 08, Gamarnik-Sudan 14, 17, Gamarnik-K. 21]

# Main Results: Solution Space Geometry

## Theorem (K., 2024)

Let  $\|\mathbf{c}\|_1 < \frac{1}{2}$  and  $\beta \in (0, 1)$  be arbitrary. Then,

$$\mathbb{E}[N_\beta] = e^{\Theta(n)}, \quad \text{where } N_\beta = |(\sigma, \sigma') : \sigma, \sigma' \in \mathcal{F}(\mathbf{c}), 1 \leq d_H(\sigma, \sigma') \leq \beta n|.$$

- **First moment evidence** that  $\exists(\sigma, \sigma')$  at arbitrarily small distances for  $\|\mathbf{c}\|_1 < \frac{1}{2}$ .
- Matching **second moment bound**: Can we show

$$\mathbb{E}[N_\beta^2] = (1 + o(1))\mathbb{E}[N_\beta]^2$$

and get  $N_\beta \geq 1$  via **Paley-Zygmund**?

# Independent Instances

$\mathcal{F}(\mathbf{c})$  defined before. For  $\mathbf{c}' = (c'_1, \dots, c'_d) \in \mathbb{R}_+^d$  and **iid**  $X'_1, \dots, X'_d \sim \mathcal{N}(0, I_n)$ , let

$$\mathcal{F}'(\mathbf{c}') = \{\boldsymbol{\sigma} \in \{-1, 1\}^n : |\langle \boldsymbol{\sigma}, X'_i \rangle| \leq \sqrt{n} 2^{-c'_i n}, \forall i \in [d]\}.$$

When is  $\mathcal{F}(\mathbf{c}) \cap \mathcal{F}'(\mathbf{c}') \neq \emptyset$ ? If  $\cap$  is empty, how far  $\mathcal{F}(\mathbf{c})$  and  $\mathcal{F}'(\mathbf{c}')$  are?

## Motivation

- RCT  $\boldsymbol{\sigma} \in \mathcal{F}(\mathbf{c})$ . Design a new RCT  $\boldsymbol{\sigma}'$  involving a new population & different constraints  $\mathbf{c}'$ .
- Repeat similar **RCT** at different regions or many years later: populations **do not** overlap

Can the same RCT  $\boldsymbol{\sigma}$  be used as is? If not, how many changes are needed?

# Solutions Spaces of Independent Instances

When is  $\mathcal{F}(\mathbf{c}) \cap \mathcal{F}(\mathbf{c}') \neq \emptyset$ ?

Consider  $\bar{\mathbf{c}} = (\mathbf{c}, \mathbf{c}') \in \mathbb{R}_+^{2d}$  and **iid**  $X_1, \dots, X_d, X'_1, \dots, X'_d \sim \mathcal{N}(0, I_n)$ . We immediately obtain

Corollary (to Theorem 1)

$\mathcal{F}(\mathbf{c}) \cap \mathcal{F}'(\mathbf{c}') \neq \emptyset$  *whp* if  $\|\mathbf{c}\|_1 + \|\mathbf{c}'\|_1 < 1$  and  $\mathcal{F}(\mathbf{c}) \cap \mathcal{F}'(\mathbf{c}') = \emptyset$  *whp* if  $\|\mathbf{c}\|_1 + \|\mathbf{c}'\|_1 > 1$ .

Suppose  $\|\mathbf{c}\|_1 + \|\mathbf{c}'\|_1 > 1$ . How far  $\mathcal{F}(\mathbf{c})$  and  $\mathcal{F}'(\mathbf{c}')$  are?

# Main Results: Distance between Independent Instances

Let  $\|\mathbf{c}\|_1 + \|\mathbf{c}'\|_1 > 1$  and  $d(\mathbf{c}, \mathbf{c}') := \min_{\sigma \in \mathcal{F}(\mathbf{c}), \sigma' \in \mathcal{F}'(\mathbf{c}')} \frac{d_H(\sigma, \sigma')}{n}$ , where  $d_H$  is Hamming distance.

## Theorem (K., 2024)

Suppose  $\gamma^* \in (0, \frac{1}{2})$  is the unique value such that

$$h(\gamma^*) = \|\mathbf{c}\|_1 + \|\mathbf{c}'\|_1 - 1, \quad \text{where } h(p) = -p \log_2 p - (1-p) \log_2(1-p).$$

Then, for any  $\epsilon > 0$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}[|d(\mathbf{c}, \mathbf{c}') - \gamma^*| \leq \epsilon] = 1$ .

- $d(\mathbf{c}, \mathbf{c}') \xrightarrow{\text{i.p.}} \gamma^*$ , which is well-defined as  $h : [0, \frac{1}{2}] \rightarrow [0, 1]$  is a bijection.
- $\mathcal{F}(\mathbf{c})$  and  $\mathcal{F}'(\mathbf{c}')$  are  $\Omega(n)$  apart.

# Future Work

- **Universality**: Is **Gaussianity** necessary?
- Second moment calculation for  $N_\beta$
- **Algorithmic** guarantees: Can we find a  $\sigma \in \mathcal{F}(c)$  in poly time?
- What are the **fundamental limits** of algs? **Overlap Gap Property**  
[Gamarnik-Sudan 14, 17, Gamarnik-Jagannath-Wein 20, Gamarnik-K. 21, Gamarnik-K.-Perkins-Xu 22,23]



# Thank you!