

Permutation Mixtures and Empirical Bayes

Yanjun Han (NYU Math and Data Science)

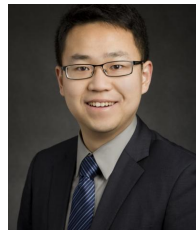
Joint work with:



Jonathan Niles-Weed
(NYU)



Yandi Shen
(CMU)



Yihong Wu
(Yale)

Statistics Seminar, Rutgers
April 23, 2025

Setup

Let P_1, \dots, P_n be n probability distributions over the same space.

A permutation mixture \mathbb{P}_n :

- draw independent $Z_1 \sim P_1, \dots, Z_n \sim P_n$;
- draw a uniformly random permutation $\pi \sim \text{Unif}(S_n)$;
- \mathbb{P}_n is the joint distribution of (X_1, \dots, X_n) with $X_i = Z_{\pi(i)}$;
- in mathematical terms:

$$(X_1, \dots, X_n) \sim \mathbb{E}_{\pi \sim \text{Unif}(S_n)} \left[\bigotimes_{i=1}^n P_{\pi(i)} \right] \quad \text{under } \mathbb{P}_n.$$

Setup

Let P_1, \dots, P_n be n probability distributions over the same space.

A permutation mixture \mathbb{P}_n :

- draw independent $Z_1 \sim P_1, \dots, Z_n \sim P_n$;
- draw a uniformly random permutation $\pi \sim \text{Unif}(S_n)$;
- \mathbb{P}_n is the joint distribution of (X_1, \dots, X_n) with $X_i = Z_{\pi(i)}$;
- in mathematical terms:

$$(X_1, \dots, X_n) \sim \mathbb{E}_{\pi \sim \text{Unif}(S_n)} \left[\bigotimes_{i=1}^n P_{\pi(i)} \right] \quad \text{under } \mathbb{P}_n.$$

An i.i.d. (mean-field) approximation \mathbb{Q}_n :

$$(X_1, \dots, X_n) \sim \left(\frac{1}{n} \sum_{i=1}^n P_i \right)^{\otimes n} \quad \text{under } \mathbb{Q}_n.$$

Setup

Let P_1, \dots, P_n be n probability distributions over the same space.

A permutation mixture \mathbb{P}_n :

- draw independent $Z_1 \sim P_1, \dots, Z_n \sim P_n$;
- draw a uniformly random permutation $\pi \sim \text{Unif}(S_n)$;
- \mathbb{P}_n is the joint distribution of (X_1, \dots, X_n) with $X_i = Z_{\pi(i)}$;
- in mathematical terms:

$$(X_1, \dots, X_n) \sim \mathbb{E}_{\pi \sim \text{Unif}(S_n)} \left[\bigotimes_{i=1}^n P_{\pi(i)} \right] \quad \text{under } \mathbb{P}_n.$$

An i.i.d. (mean-field) approximation \mathbb{Q}_n :

$$(X_1, \dots, X_n) \sim \left(\frac{1}{n} \sum_{i=1}^n P_i \right)^{\otimes n} \quad \text{under } \mathbb{Q}_n.$$

Target of this work

Show that the i.i.d. approximation \mathbb{Q}_n to \mathbb{P}_n is accurate, i.e. the information divergence (or statistical distance) between \mathbb{P}_n and \mathbb{Q}_n is small (and ideally, **independent** of n)

Later in the talk:

- statistics: permutation prior
- information theory: permutation channel
- probability: de Finetti-style theorems
- indirect application (second half): compound decisions and empirical Bayes

Later in the talk:

- statistics: permutation prior
- information theory: permutation channel
- probability: de Finetti-style theorems
- indirect application (second half): compound decisions and empirical Bayes

Bigger picture:

- general mean-field approximation
- information geometry of high-dimensional mixtures

Failure of existing approaches in a toy example

Let $P_1 = \cdots = P_{n/2} = \mathcal{N}(\mu, 1)$ and $P_{n/2+1} = \cdots = P_n = \mathcal{N}(-\mu, 1)$

- $\mathbb{P}_n = \nu_{\mathbb{P}} \star \mathcal{N}(0, I_n)$, where $\nu_{\mathbb{P}}$ is the distribution of n uniformly random draws from the multiset $\{-\mu, \dots, -\mu, \mu, \dots, \mu\}$ **without replacement**;
- $\mathbb{Q}_n = \nu_{\mathbb{Q}} \star \mathcal{N}(0, I_n)$, where $\nu_{\mathbb{Q}}$ is the counterpart **with replacement**;

$$\chi^2(P\|Q) := \sum_x \frac{(p_x - q_x)^2}{q_x}$$

Failure of existing approaches in a toy example

Let $P_1 = \dots = P_{n/2} = \mathcal{N}(\mu, 1)$ and $P_{n/2+1} = \dots = P_n = \mathcal{N}(-\mu, 1)$

- $\mathbb{P}_n = \nu_{\mathbb{P}} \star \mathcal{N}(0, I_n)$, where $\nu_{\mathbb{P}}$ is the distribution of n uniformly random draws from the multiset $\{-\mu, \dots, -\mu, \mu, \dots, \mu\}$ **without replacement**;
- $\mathbb{Q}_n = \nu_{\mathbb{Q}} \star \mathcal{N}(0, I_n)$, where $\nu_{\mathbb{Q}}$ is the counterpart **with replacement**;

Our result

$$\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) = \begin{cases} O(\mu^4) & \text{if } \mu \leq 1, \\ O(\exp(\mu^2)) & \text{if } \mu > 1. \end{cases}$$

- χ^2 -divergence independent of dimension n
- smaller than the one-dimensional divergence $\chi^2(\mathcal{N}(\mu, 1) \| \mathcal{N}(-\mu, 1))$
- existing approaches fail even for this toy example

$$\chi^2(P \| Q) := \sum_x \frac{(p_x - q_x)^2}{q_x}$$

Failed approach I: reduction to two simple distributions

Apply convexity to reduce to the divergence between two simple distributions:

$$\begin{aligned}\text{KL}(\mathbb{P}_n \| \mathbb{Q}_n) &= \text{KL}(\mathbb{E}_{\vartheta \sim \nu_{\mathbb{P}}} [\mathcal{N}(\vartheta, I_n)] \| \mathbb{E}_{\vartheta' \sim \nu_{\mathbb{Q}}} [\mathcal{N}(\vartheta', I_n)]) \\ &\leq \min_{\rho \in \Pi(\nu_{\mathbb{P}}, \nu_{\mathbb{Q}})} \mathbb{E}_{(\vartheta, \vartheta') \sim \rho} [\text{KL}(\mathcal{N}(\vartheta, I_n) \| \mathcal{N}(\vartheta', I_n))] \\ &= \frac{W_2^2(\nu_{\mathbb{P}}, \nu_{\mathbb{Q}})}{2} \asymp \sqrt{n} \mu^2\end{aligned}$$

$$\text{KL}(P \| Q) := \sum_x p(x) \log \frac{p(x)}{q(x)}$$

Failed approach I: reduction to two simple distributions

Apply convexity to reduce to the divergence between two simple distributions:

$$\begin{aligned}\text{KL}(\mathbb{P}_n \| \mathbb{Q}_n) &= \text{KL}(\mathbb{E}_{\vartheta \sim \nu_{\mathbb{P}}} [\mathcal{N}(\vartheta, I_n)] \| \mathbb{E}_{\vartheta' \sim \nu_{\mathbb{Q}}} [\mathcal{N}(\vartheta', I_n)]) \\ &\leq \min_{\rho \in \Pi(\nu_{\mathbb{P}}, \nu_{\mathbb{Q}})} \mathbb{E}_{(\vartheta, \vartheta') \sim \rho} [\text{KL}(\mathcal{N}(\vartheta, I_n) \| \mathcal{N}(\vartheta', I_n))] \\ &= \frac{W_2^2(\nu_{\mathbb{P}}, \nu_{\mathbb{Q}})}{2} \asymp \sqrt{n} \mu^2\end{aligned}$$

→ grows with the dimension n

→ wrong dependence on μ

$$\text{KL}(P \| Q) := \sum_x p(x) \log \frac{p(x)}{q(x)}$$

Failed approach II: reduction to one simple distribution

A more careful coupling:

$$\text{KL}(\mathbb{P}_n \| \mathbb{Q}_n) \leq \min_{\{\nu_{\theta'}\}_{\theta' \in \{\pm\mu\}^n}} \mathbb{E}_{\vartheta' \sim \nu_{\mathbb{Q}}} [\text{KL}(\mathbb{E}_{\vartheta \sim \nu_{\vartheta'}} [\mathcal{N}(\vartheta, I_n)] \| \mathcal{N}(\vartheta', I_n))],$$

where the minimization is over all possible families of distributions $\{\nu_{\theta'}\}_{\theta' \in \{\pm\mu\}^n}$ such that $\mathbb{E}_{\vartheta' \sim \nu_{\mathbb{Q}}}[\nu_{\vartheta'}] = \nu_{\mathbb{P}}$.

Failed approach II: reduction to one simple distribution

A more careful coupling:

$$\text{KL}(\mathbb{P}_n \| \mathbb{Q}_n) \leq \min_{\{\nu_{\theta'}\}_{\theta' \in \{\pm\mu\}^n}} \mathbb{E}_{\vartheta' \sim \nu_{\mathbb{Q}}} [\text{KL}(\mathbb{E}_{\vartheta \sim \nu_{\vartheta'}} [\mathcal{N}(\vartheta, I_n)] \| \mathcal{N}(\vartheta', I_n))],$$

where the minimization is over all possible families of distributions $\{\nu_{\theta'}\}_{\theta' \in \{\pm\mu\}^n}$ such that $\mathbb{E}_{\vartheta' \sim \nu_{\mathbb{Q}}}[\nu_{\vartheta'}] = \nu_{\mathbb{P}}$.

- a judicious choice [Ding'22] leads to an upper bound $O(\mu^2)$ for small μ
- however, can show that any such upper bound must be $\Omega(\mu^2)$

Failed approach III: method of moments

A powerful approach to upper bound the statistical difference between two mixtures distributions, with many recent applications [Cai and Low'11, Hardt and Price'15, Wu and Yang'20, Han et al.'20, ...]

$$\text{TV}(P, Q) := \frac{1}{2} \sum_x |p_x - q_x|$$

Failed approach III: method of moments

A powerful approach to upper bound the statistical difference between two mixtures distributions, with many recent applications [Cai and Low'11, Hardt and Price'15, Wu and Yang'20, Han et al.'20, ...]

Idea: express the Gaussian likelihood ratio in terms of **Hermite polynomials**

$$\frac{\varphi(x - \theta)}{\varphi(x)} = \sum_{k=0}^{\infty} \frac{H_k(x)}{k!} \theta^k,$$

$$\text{TV}(P, Q) := \frac{1}{2} \sum_x |p_x - q_x|$$

Failed approach III: method of moments

A powerful approach to upper bound the statistical difference between two mixtures distributions, with many recent applications [Cai and Low'11, Hardt and Price'15, Wu and Yang'20, Han et al.'20, ...]

Idea: express the Gaussian likelihood ratio in terms of **Hermite polynomials**

$$\frac{\varphi(x - \theta)}{\varphi(x)} = \sum_{k=0}^{\infty} \frac{H_k(x)}{k!} \theta^k,$$

so that

$$\begin{aligned} \text{TV}(\mu \star \mathcal{N}(0, 1), \nu \star \mathcal{N}(0, 1))^2 &= \frac{1}{4} \left(\mathbb{E}_{Z \sim \mathcal{N}(0, 1)} \left| \mathbb{E}_{U \sim \mu} \left[\frac{\varphi(Z - U)}{\varphi(Z)} \right] - \mathbb{E}_{V \sim \nu} \left[\frac{\varphi(Z - V)}{\varphi(Z)} \right] \right| \right)^2 \\ &= \frac{1}{4} \left(\mathbb{E}_{Z \sim \mathcal{N}(0, 1)} \left| \sum_{k=0}^{\infty} \frac{H_k(Z)}{k!} (\mathbb{E}_{U \sim \mu}[U^k] - \mathbb{E}_{V \sim \nu}[V^k]) \right| \right)^2 \\ &\stackrel{\text{C-S}}{\leq} \frac{1}{4} \mathbb{E}_{Z \sim \mathcal{N}(0, 1)} \left(\sum_{k=0}^{\infty} \frac{H_k(Z)}{k!} (\mathbb{E}_{U \sim \mu}[U^k] - \mathbb{E}_{V \sim \nu}[V^k]) \right)^2 \\ &= \frac{1}{4} \sum_{k=0}^{\infty} \frac{(\mathbb{E}_{U \sim \mu}[U^k] - \mathbb{E}_{V \sim \nu}[V^k])^2}{k!} \end{aligned}$$

$$\text{TV}(P, Q) := \frac{1}{2} \sum_x |p_x - q_x|$$

Failed approach III: method of moments (cont'd)

In general dimensions:

$$\text{TV}(\nu_{\mathbb{P}} \star \mathcal{N}(0, I_n) \| \nu_{\mathbb{Q}} \star \mathcal{N}(0, I_n))^2 \leq \frac{1}{4} \sum_{\vec{\alpha} \in \mathbb{N}^n} \frac{(m_{\vec{\alpha}}(\nu_{\mathbb{P}}) - m_{\vec{\alpha}}(\nu_{\mathbb{Q}}))^2}{\vec{\alpha}!}$$

→ $\vec{\alpha} = (\alpha_1, \dots, \alpha_n)$ is a multi-index, with $\vec{\alpha}! := \alpha_1! \cdots \alpha_n!$

→ $m_{\vec{\alpha}}(\mu) := \mathbb{E}_{\vartheta \sim \mu}[\vartheta_1^{\alpha_1} \cdots \vartheta_n^{\alpha_n}]$ denotes the joint moment

Failed approach III: method of moments (cont'd)

In general dimensions:

$$\text{TV}(\nu_{\mathbb{P}} \star \mathcal{N}(0, I_n) \| \nu_{\mathbb{Q}} \star \mathcal{N}(0, I_n))^2 \leq \frac{1}{4} \sum_{\vec{\alpha} \in \mathbb{N}^n} \frac{(m_{\vec{\alpha}}(\nu_{\mathbb{P}}) - m_{\vec{\alpha}}(\nu_{\mathbb{Q}}))^2}{\vec{\alpha}!}$$

- $\vec{\alpha} = (\alpha_1, \dots, \alpha_n)$ is a multi-index, with $\vec{\alpha}! := \alpha_1! \cdots \alpha_n!$
- $m_{\vec{\alpha}}(\mu) := \mathbb{E}_{\vartheta \sim \mu}[\vartheta_1^{\alpha_1} \cdots \vartheta_n^{\alpha_n}]$ denotes the joint moment

Application to our toy example:

- non-zero moment difference starting from $|\vec{\alpha}| = 2$, suggesting an $O(\mu^4)$ dependence
- however, too many cross terms in high dimensions: the total contributions of $|\vec{\alpha}| = 2\ell$ are at least $\Omega_{\ell}(\mu^{4\ell} n^{\ell-1})$, which is growing with n for $\ell \geq 2$

Failed approach IV: method of cumulants

A recent development based on cumulants [Schramm and Wein'22]:

$$\chi^2(\nu_{\mathbb{P}} \star \mathcal{N}(0, I_n) \| \nu_{\mathbb{Q}} \star \mathcal{N}(0, I_n)) \leq \sum_{\vec{\alpha} \in \mathbb{N}^d} \frac{\kappa_{\vec{\alpha}}^2}{\vec{\alpha}!},$$

where $\kappa_{\vec{\alpha}}$ is the joint cumulant

$$\kappa_{\vec{\alpha}} = \kappa_{\nu_{\mathbb{Q}}} \left(\frac{d\nu_{\mathbb{P}}}{d\nu_{\mathbb{Q}}}, \vartheta_1, \dots, \vartheta_1, \vartheta_2, \dots, \vartheta_2, \dots, \vartheta_n \right).$$

$$\kappa(X_1, \dots, X_n) := \frac{\partial^n}{\partial t_1 \dots \partial t_n} \Big|_{t_1 = \dots = t_n = 0} \log \mathbb{E} [\exp (\sum_{i=1}^n t_i X_i)]$$

Failed approach IV: method of cumulants

A recent development based on cumulants [Schramm and Wein'22]:

$$\chi^2(\nu_{\mathbb{P}} \star \mathcal{N}(0, I_n) \| \nu_{\mathbb{Q}} \star \mathcal{N}(0, I_n)) \leq \sum_{\vec{\alpha} \in \mathbb{N}^d} \frac{\kappa_{\vec{\alpha}}^2}{\vec{\alpha}!},$$

where $\kappa_{\vec{\alpha}}$ is the joint cumulant

$$\kappa_{\vec{\alpha}} = \kappa_{\nu_{\mathbb{Q}}} \left(\frac{d\nu_{\mathbb{P}}}{d\nu_{\mathbb{Q}}}, \vartheta_1, \dots, \vartheta_1, \vartheta_2, \dots, \vartheta_2, \dots, \vartheta_n \right).$$

- a better behavior for certain cross terms
- however, can show that $\kappa_{(1,\ell,0,\dots,0)} \asymp C^\ell \ell!$ for odd ℓ , so summing along this subsequence gives a diverging result

$$\kappa(X_1, \dots, X_n) := \frac{\partial^n}{\partial t_1 \dots \partial t_n} \Big|_{t_1=\dots=t_n=0} \log \mathbb{E} [\exp(\sum_{i=1}^n t_i X_i)]$$

Main result

Let $P_1, \dots, P_n \in \mathcal{P}$. Define the following **dimension-independent** quantities:

Definition (Quantities of \mathcal{P})

- χ^2 channel capacity: $C_{\chi^2}(\mathcal{P}) = \sup_{\rho \in \Delta(\mathcal{P})} I_{\chi^2}(P; X)$, with $P \sim \rho$ and $X \sim P$
- χ^2 diameter: $D_{\chi^2}(\mathcal{P}) = \sup_{P_1, P_2 \in \mathcal{P}} \chi^2(P_1 \| P_2)$

$$I_{\chi^2}(X; Y) := \chi^2(P_{XY} \| P_X P_Y)$$

Main result

Let $P_1, \dots, P_n \in \mathcal{P}$. Define the following **dimension-independent** quantities:

Definition (Quantities of \mathcal{P})

- χ^2 channel capacity: $C_{\chi^2}(\mathcal{P}) = \sup_{\rho \in \Delta(\mathcal{P})} I_{\chi^2}(P; X)$, with $P \sim \rho$ and $X \sim P$
- χ^2 diameter: $D_{\chi^2}(\mathcal{P}) = \sup_{P_1, P_2 \in \mathcal{P}} \chi^2(P_1 \| P_2)$

Theorem (H., Niles-Weed'24)

$$\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) \leq \min \left\{ 10 \sum_{\ell=2}^n C_{\chi^2}(\mathcal{P})^\ell, (1 + D_{\chi^2}(\mathcal{P}))^{1+C_{\chi^2}(\mathcal{P})} - 1 \right\}$$

- \mathbb{P}_n is contiguous to \mathbb{Q}_n : $\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) = \mathcal{O}_{\mathcal{P}}(1)$ if $D_{\chi^2}(\mathcal{P}) < \infty$
- high-probability events under the simpler product measure \mathbb{Q}_n translate to high-probability events under the mixture \mathbb{P}_n

$$I_{\chi^2}(X; Y) := \chi^2(P_{XY} \| P_X P_Y)$$

Theorem (H., Niles-Weed'24)

$$\chi^2(\mathbb{P}_n \parallel \mathbb{Q}_n) \leq \min \left\{ 10 \sum_{\ell=2}^n C_{\chi^2}(\mathcal{P})^\ell, (1 + D_{\chi^2}(\mathcal{P}))^{1+C_{\chi^2}(\mathcal{P})} - 1 \right\}$$

Examples

Theorem (H., Niles-Weed'24)

$$\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) \leq \min \left\{ 10 \sum_{\ell=2}^n C_{\chi^2}(\mathcal{P})^\ell, (1 + D_{\chi^2}(\mathcal{P}))^{1+C_{\chi^2}(\mathcal{P})} - 1 \right\}$$

Example I (Two-component Gaussian)

$\mathcal{P} = \{\mathcal{N}(\mu, 1), \mathcal{N}(-\mu, 1)\}$: $C_{\chi^2}(\mathcal{P}) \leq 1 - e^{-\mu^2}$, so

$$\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) = \begin{cases} O(\mu^4) & \text{if } \mu \leq 1, \\ O(\exp(\mu^2)) & \text{if } \mu > 1. \end{cases}$$

Theorem (H., Niles-Weed'24)

$$\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) \leq \min \left\{ 10 \sum_{\ell=2}^n C_{\chi^2}(\mathcal{P})^\ell, (1 + D_{\chi^2}(\mathcal{P}))^{1+C_{\chi^2}(\mathcal{P})} - 1 \right\}$$

Example I (Two-component Gaussian)

$\mathcal{P} = \{\mathcal{N}(\mu, 1), \mathcal{N}(-\mu, 1)\}$: $C_{\chi^2}(\mathcal{P}) \leq 1 - e^{-\mu^2}$, so

$$\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) = \begin{cases} O(\mu^4) & \text{if } \mu \leq 1, \\ O(\exp(\mu^2)) & \text{if } \mu > 1. \end{cases}$$

Example II (Bounded Gaussian)

$\mathcal{P} = \{\mathcal{N}(\theta, 1) : |\theta| \leq \mu\}$: $C_{\chi^2}(\mathcal{P}) = O(\mu \wedge \mu^2)$, $D_{\chi^2}(\mathcal{P}) = \exp(O(\mu^2))$, so

$$\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) = \begin{cases} O(\mu^4) & \text{if } \mu \leq 1, \\ \exp(O(\mu^3))^a & \text{if } \mu > 1. \end{cases}$$

^aWith Y. Liang, recently improved to $\exp(O(\mu^2))$ by higher-order Cheeger inequality

Applications

Sequence model in statistics: observe $X_i \sim P_{\theta_i}$ with unknown $\theta = (\theta_1, \dots, \theta_n)$

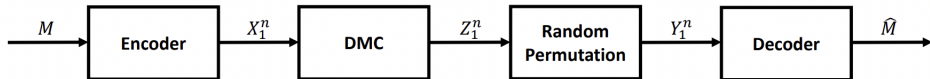
- a common “permutation prior”: $\theta = (v_{\pi(1)}, \dots, v_{\pi(n)})$ for a known vector v and a random permutation π
- a quantity of interest: mutual information $I(\theta; X^n)$

Our result: can pretend as if the coordinates $\theta_i \sim \frac{1}{n} \sum_{j=1}^n \delta_{v_j}$ are i.i.d.

Mutual information under a permutation prior

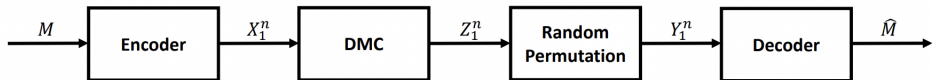
$$I_{\mathbb{Q}_n}(\theta; X^n) - \mathcal{O}_{\mathcal{P}}(1) \leq I_{\mathbb{P}_n}(\theta; X^n) \leq I_{\mathbb{Q}_n}(\theta; X^n)$$

Information theory: permutation channel



The noisy permutation channel [Makur'20]

Information theory: permutation channel

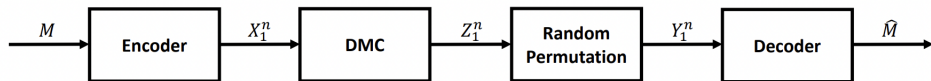


The noisy permutation channel [Makur'20]

- target: find the channel capacity $C_n(\mathcal{P}) = \max_{p(x^n)} I(X^n; Y^n)$
- known achievability [Makur'20] and converse [Tang and Polyanskiy'23]:

$$C_n(\mathcal{P}) \sim \frac{\text{rank}(P_{Z|X}) - 1}{2} \log n \quad \text{for discrete } \mathcal{P}.$$

Information theory: permutation channel



The noisy permutation channel [Makur'20]

- target: find the channel capacity $C_n(\mathcal{P}) = \max_{p(x^n)} I(X^n; Y^n)$
- known achievability [Makur'20] and converse [Tang and Polyanskiy'23]:

$$C_n(\mathcal{P}) \sim \frac{\text{rank}(P_{Z|X}) - 1}{2} \log n \quad \text{for discrete } \mathcal{P}.$$

Our result: for general \mathcal{P} , can pretend as if Y^n have independent coordinates

Converse for general permutation channels

$$C_n(\mathcal{P}) \leq \text{Red}(\text{conv}(\mathcal{P})^{\otimes n}) + \mathcal{O}_{\mathcal{P}}(1)$$

Theorem (de Finetti)

Any exchangeable distribution P_{X^∞} can be written as an i.i.d. mixture:

$$P_{X^\infty}(x^\infty) = \mathbb{E}_\theta \left[\prod_{i=1}^{\infty} Q_\theta(x_i) \right].$$

The joint distribution of (X_1, \dots, X_n) is exchangeable if $(X_1, \dots, X_n) \stackrel{d}{=} (X_{\pi(1)}, \dots, X_{\pi(n)})$

Theorem (de Finetti)

Any exchangeable distribution P_{X^∞} can be written as an i.i.d. mixture:

$$P_{X^\infty}(x^\infty) = \mathbb{E}_\theta \left[\prod_{i=1}^{\infty} Q_\theta(x_i) \right].$$

Approximately holds for exchangeable distribution P_{X^n} with finite n :

- $\text{KL}(P_{X^k} \| \mathbb{E}_\theta[Q_\theta^{\otimes k}]) \lesssim \frac{k^2}{n}$ [Diaconis and Freedman'80]
- for small $|\mathcal{X}|$, $\text{KL}(P_{X^k} \| \mathbb{E}_\theta[Q_\theta^{\otimes k}]) \lesssim \frac{|\mathcal{X}|k^2}{n(n+1-k)}$ [Stam'78]
- more recent refinements [Gavalakis and Kontoyiannis'21; Johnson, Gavalakis, and Kontoyiannis'24]

The joint distribution of (X_1, \dots, X_n) is exchangeable if $(X_1, \dots, X_n) \stackrel{d}{=} (X_{\pi(1)}, \dots, X_{\pi(n)})$

Using the first upper bound and $C_{\chi^2}(\mathcal{P}) \leq |\mathcal{X}|$:

χ^2 -type finite de Finetti

For exchangeable distribution P_{X^n} and $k \leq n$:

$$\chi^2 \left(P_{X^k} \| \mathbb{E}_\theta [Q_\theta^{\otimes k}] \right) \lesssim \frac{k^2 |\mathcal{X}|^2}{n^2} \quad \text{if } k < \frac{n}{|\mathcal{X}|}.$$

Our extensions

Using the first upper bound and $C_{\chi^2}(\mathcal{P}) \leq |\mathcal{X}|$:

χ^2 -type finite de Finetti

For exchangeable distribution P_{X^n} and $k \leq n$:

$$\chi^2 \left(P_{X^k} \| \mathbb{E}_\theta [Q_\theta^{\otimes k}] \right) \lesssim \frac{k^2 |\mathcal{X}|^2}{n^2} \quad \text{if } k < \frac{n}{|\mathcal{X}|}.$$

Using the second upper bound:

Noisy de Finetti

Let P_{Y^n} be the output distribution with an input exchangeable distribution P_{X^n} and a channel \mathcal{P} . Then for $k \leq n$:

$$\chi^2 \left(P_{Y^k} \| \mathbb{E}_\theta [Q_\theta^{\otimes k}] \right) = \mathcal{O}_{\mathcal{P}} \left(\frac{k^2}{n^2} \right) \quad \text{if } D_{\chi^2}(\mathcal{P}) < \infty.$$

Sketch of the first upper bound

Toy example: a different basis

→ Hermite basis:

$$\frac{\varphi(x - \theta)}{\varphi(x)} = \sum_{k=0}^{\infty} \frac{H_k(x)}{k!} \theta^k$$

where φ is the density of $\mathcal{N}(0, 1)$.

$$(\theta_1, \dots, \theta_n) = (\mu, \dots, \mu, -\mu, \dots, -\mu).$$

Toy example: a different basis

→ Hermite basis:

$$\frac{\varphi(x - \theta)}{\varphi(x)} = \sum_{k=0}^{\infty} \frac{H_k(x)}{k!} \theta^k$$

where φ is the density of $\mathcal{N}(0, 1)$.

→ Hyperbolic basis?

$$\frac{\varphi(x - \theta)}{\varphi_0(x)} = 1 + \tanh(\mu x) \frac{\theta}{\mu}, \quad \theta \in \{\pm\mu\}$$

where $\varphi_0(x) = \frac{\varphi(x-\mu) + \varphi(x+\mu)}{2}$ is the common marginal of \mathbb{P}_n^2 and \mathbb{Q}_n

$$(\theta_1, \dots, \theta_n) = (\mu, \dots, \mu, -\mu, \dots, -\mu).$$

Toy example: a different basis

→ Hermite basis:

$$\frac{\varphi(x - \theta)}{\varphi(x)} = \sum_{k=0}^{\infty} \frac{H_k(x)}{k!} \theta^k$$

where φ is the density of $\mathcal{N}(0, 1)$.

→ $\{H_0(x), H_1(x), \dots\}$ are orthogonal in $L^2(\varphi)$

→ Hyperbolic basis?

$$\frac{\varphi(x - \theta)}{\varphi_0(x)} = 1 + \tanh(\mu x) \frac{\theta}{\mu}, \quad \theta \in \{\pm\mu\}$$

where $\varphi_0(x) = \frac{\varphi(x-\mu) + \varphi(x+\mu)}{2}$ is the common marginal of \mathbb{P}_n^2 and \mathbb{Q}_n

$$(\theta_1, \dots, \theta_n) = (\mu, \dots, \mu, -\mu, \dots, -\mu).$$

Toy example: a different basis

→ Hermite basis:

$$\frac{\varphi(x - \theta)}{\varphi(x)} = \sum_{k=0}^{\infty} \frac{H_k(x)}{k!} \theta^k$$

where φ is the density of $\mathcal{N}(0, 1)$.

→ $\{H_0(x), H_1(x), \dots\}$ are orthogonal in $L^2(\varphi)$

→ Hyperbolic basis?

$$\frac{\varphi(x - \theta)}{\varphi_0(x)} = 1 + \tanh(\mu x) \frac{\theta}{\mu}, \quad \theta \in \{\pm\mu\}$$

where $\varphi_0(x) = \frac{\varphi(x-\mu) + \varphi(x+\mu)}{2}$ is the common marginal of \mathbb{P}_n^2 and \mathbb{Q}_n

→ $\{1, \tanh(\mu x)\}$ are orthogonal in $L^2(\varphi_0)$

$$(\theta_1, \dots, \theta_n) = (\mu, \dots, \mu, -\mu, \dots, -\mu).$$

Toy example: a different basis

→ Hermite basis:

$$\frac{\varphi(x - \theta)}{\varphi(x)} = \sum_{k=0}^{\infty} \frac{H_k(x)}{k!} \theta^k$$

where φ is the density of $\mathcal{N}(0, 1)$.

→ $\{H_0(x), H_1(x), \dots\}$ are orthogonal in $L^2(\varphi)$

→ Hyperbolic basis?

$$\frac{\varphi(x - \theta)}{\varphi_0(x)} = 1 + \tanh(\mu x) \frac{\theta}{\mu}, \quad \theta \in \{\pm\mu\}$$

where $\varphi_0(x) = \frac{\varphi(x-\mu) + \varphi(x+\mu)}{2}$ is the common marginal of \mathbb{P}_n^2 and \mathbb{Q}_n

→ $\{1, \tanh(\mu x)\}$ are orthogonal in $L^2(\varphi_0)$

→ $\mathbb{E}\left[\frac{\theta}{\mu}\right] = 0$ for $\theta \sim \text{Unif}(\{\pm\mu\})$

$$(\theta_1, \dots, \theta_n) = (\mu, \dots, \mu, -\mu, \dots, -\mu).$$

Toy example: a different basis

→ Hermite basis:

$$\frac{\varphi(x - \theta)}{\varphi(x)} = \sum_{k=0}^{\infty} \frac{H_k(x)}{k!} \theta^k$$

where φ is the density of $\mathcal{N}(0, 1)$.

→ $\{H_0(x), H_1(x), \dots\}$ are orthogonal in $L^2(\varphi)$

→ $\mathbb{E}[\theta^k]$ possibly non-zero for $\theta \sim \text{Unif}(\{\pm\mu\})$

→ Hyperbolic basis?

$$\frac{\varphi(x - \theta)}{\varphi_0(x)} = 1 + \tanh(\mu x) \frac{\theta}{\mu}, \quad \theta \in \{\pm\mu\}$$

where $\varphi_0(x) = \frac{\varphi(x-\mu) + \varphi(x+\mu)}{2}$ is the common marginal of \mathbb{P}_n^2 and \mathbb{Q}_n

→ $\{1, \tanh(\mu x)\}$ are orthogonal in $L^2(\varphi_0)$

→ $\mathbb{E}[\frac{\theta}{\mu}] = 0$ for $\theta \sim \text{Unif}(\{\pm\mu\})$

$$(\theta_1, \dots, \theta_n) = (\mu, \dots, \mu, -\mu, \dots, -\mu).$$

Toy example: a different basis

→ Hermite basis:

$$\frac{\varphi(x - \theta)}{\varphi(x)} = \sum_{k=0}^{\infty} \frac{H_k(x)}{k!} \theta^k$$

where φ is the density of $\mathcal{N}(0, 1)$.

→ $\{H_0(x), H_1(x), \dots\}$ are orthogonal in $L^2(\varphi)$

→ $\mathbb{E}[\theta^k]$ possibly non-zero for $\theta \sim \text{Unif}(\{\pm\mu\})$

→ Hyperbolic basis?

$$\frac{\varphi(x - \theta)}{\varphi_0(x)} = 1 + \tanh(\mu x) \frac{\theta}{\mu}, \quad \theta \in \{\pm\mu\}$$

where $\varphi_0(x) = \frac{\varphi(x-\mu) + \varphi(x+\mu)}{2}$ is the common marginal of \mathbb{P}_n^2 and \mathbb{Q}_n

→ $\{1, \tanh(\mu x)\}$ are orthogonal in $L^2(\varphi_0)$

→ $\mathbb{E}[\frac{\theta}{\mu}] = 0$ for $\theta \sim \text{Unif}(\{\pm\mu\})$

Under the new basis:

$$\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}(x^n) = \mathbb{E}_{\pi} \left[\prod_{i=1}^n \frac{\varphi(x_i - \theta_{\pi(i)})}{\varphi_0(x_i)} \right] = \mathbb{E}_{\pi} \left[\prod_{i=1}^n \left(1 + \tanh(\mu x_i) \frac{\theta_{\pi(i)}}{\mu} \right) \right]$$

$$(\theta_1, \dots, \theta_n) = (\mu, \dots, \mu, -\mu, \dots, -\mu).$$

Toy example: a different basis

→ Hermite basis:

$$\frac{\varphi(x - \theta)}{\varphi(x)} = \sum_{k=0}^{\infty} \frac{H_k(x)}{k!} \theta^k$$

where φ is the density of $\mathcal{N}(0, 1)$.

→ $\{H_0(x), H_1(x), \dots\}$ are orthogonal in $L^2(\varphi)$

→ $\mathbb{E}[\theta^k]$ possibly non-zero for $\theta \sim \text{Unif}(\{\pm\mu\})$

→ Hyperbolic basis?

$$\frac{\varphi(x - \theta)}{\varphi_0(x)} = 1 + \tanh(\mu x) \frac{\theta}{\mu}, \quad \theta \in \{\pm\mu\}$$

where $\varphi_0(x) = \frac{\varphi(x-\mu) + \varphi(x+\mu)}{2}$ is the common marginal of \mathbb{P}_n^2 and \mathbb{Q}_n

→ $\{1, \tanh(\mu x)\}$ are orthogonal in $L^2(\varphi_0)$

→ $\mathbb{E}\left[\frac{\theta}{\mu}\right] = 0$ for $\theta \sim \text{Unif}(\{\pm\mu\})$

Under the new basis:

$$\begin{aligned} \frac{d\mathbb{P}_n}{d\mathbb{Q}_n}(x^n) &= \mathbb{E}_{\pi} \left[\prod_{i=1}^n \frac{\varphi(x_i - \theta_{\pi(i)})}{\varphi_0(x_i)} \right] = \mathbb{E}_{\pi} \left[\prod_{i=1}^n \left(1 + \tanh(\mu x_i) \frac{\theta_{\pi(i)}}{\mu} \right) \right] \\ &= \sum_{S \subseteq [n]} \mathbb{E}_{\pi} \left[\prod_{i \in S} \frac{\theta_{\pi(i)}}{\mu} \right] \prod_{i \in S} \tanh(\mu x_i) \end{aligned}$$

$$(\theta_1, \dots, \theta_n) = (\mu, \dots, \mu, -\mu, \dots, -\mu).$$

Toy example: full analysis

$$\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}(x^n) = \sum_{S \subseteq [n]} \mathbb{E}_\pi \left[\prod_{i \in S} \frac{\theta_{\pi(i)}}{\mu} \right] \prod_{i \in S} \tanh(\mu x_i)$$

Toy example: full analysis

$$\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}(x^n) = \sum_{S \subseteq [n]} \mathbb{E}_\pi \left[\prod_{i \in S} \frac{\theta_{\pi(i)}}{\mu} \right] \prod_{i \in S} \tanh(\mu x_i)$$

→ orthogonality of $\{1, \tanh(\mu x)\}$ under $L^2(\varphi_0)$:

$$\mathbb{E}_{\mathbb{Q}_n} \left[\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n} \right)^2 \right] = \sum_{S \subseteq [n]} \left(\mathbb{E}_\pi \left[\prod_{i \in S} \frac{\theta_{\pi(i)}}{\mu} \right] \right)^2 c_{\chi^2(\mathcal{P})}^{|S|}$$

Toy example: full analysis

$$\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}(x^n) = \sum_{S \subseteq [n]} \mathbb{E}_\pi \left[\prod_{i \in S} \frac{\theta_{\pi(i)}}{\mu} \right] \prod_{i \in S} \tanh(\mu x_i)$$

→ orthogonality of $\{1, \tanh(\mu x)\}$ under $L^2(\varphi_0)$:

$$\mathbb{E}_{\mathbb{Q}_n} \left[\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n} \right)^2 \right] = \sum_{S \subseteq [n]} \left(\mathbb{E}_\pi \left[\prod_{i \in S} \frac{\theta_{\pi(i)}}{\mu} \right] \right)^2 \mathbb{C}_{\chi^2(\mathcal{P})}^{|S|}$$

→ the inner expectation: for $|S| = \ell$,

$$\left(\mathbb{E}_\pi \left[\prod_{i \in S} \frac{\theta_{\pi(i)}}{\mu} \right] \right)^2 \leq \frac{\mathbb{1}_{\ell \text{ is even}}}{\binom{n}{\ell}}$$

Toy example: full analysis

$$\frac{d\mathbb{P}_n}{d\mathbb{Q}_n}(x^n) = \sum_{S \subseteq [n]} \mathbb{E}_\pi \left[\prod_{i \in S} \frac{\theta_{\pi(i)}}{\mu} \right] \prod_{i \in S} \tanh(\mu x_i)$$

→ orthogonality of $\{1, \tanh(\mu x)\}$ under $L^2(\varphi_0)$:

$$\mathbb{E}_{\mathbb{Q}_n} \left[\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n} \right)^2 \right] = \sum_{S \subseteq [n]} \left(\mathbb{E}_\pi \left[\prod_{i \in S} \frac{\theta_{\pi(i)}}{\mu} \right] \right)^2 C_{\chi^2(\mathcal{P})}^{|S|}$$

→ the inner expectation: for $|S| = \ell$,

$$\left(\mathbb{E}_\pi \left[\prod_{i \in S} \frac{\theta_{\pi(i)}}{\mu} \right] \right)^2 \leq \frac{\mathbb{1}_{\ell \text{ is even}}}{\binom{n}{\ell}}$$

→ piecing everything together:

$$\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) = \mathbb{E}_{\mathbb{Q}_n} \left[\left(\frac{d\mathbb{P}_n}{d\mathbb{Q}_n} \right)^2 \right] - 1 \leq C_{\chi^2(\mathcal{P})}^2 + C_{\chi^2(\mathcal{P})}^4 + \cdots + C_{\chi^2(\mathcal{P})}^n$$

Importance of zero-mean: a Maclaurin-type inequality

For a vector $x = (x_1, \dots, x_n)$, define the elementary symmetric polynomial

$$e_\ell(x) := \sum_{|S|=\ell} \prod_{i \in S} x_i$$

Importance of zero-mean: a Maclaurin-type inequality

For a vector $x = (x_1, \dots, x_n)$, define the elementary symmetric polynomial

$$e_\ell(x) := \sum_{|S|=\ell} \prod_{i \in S} x_i$$

Theorem (Upper bound on ESPs for centered vector)

Let $\sum_{i=1}^n x_i = 0$ and $\sum_{i=1}^n |x_i|^2 = n$.

→ If $x \in \mathbb{R}^n$, then $|e_\ell(x)|^2 \leq 10 \binom{n}{\ell}$;

→ If $x \in \mathbb{C}^n$, a weaker upper bound holds:

$$|e_\ell(x)|^2 \leq \frac{n^n}{\ell^\ell (n-\ell)^{n-\ell}} < 3\sqrt{\ell+1} \binom{n}{\ell}.$$

Importance of zero-mean: a Maclaurin-type inequality

For a vector $x = (x_1, \dots, x_n)$, define the elementary symmetric polynomial

$$e_\ell(x) := \sum_{|S|=\ell} \prod_{i \in S} x_i$$

Theorem (Upper bound on ESPs for centered vector)

Let $\sum_{i=1}^n x_i = 0$ and $\sum_{i=1}^n |x_i|^2 = n$.

→ If $x \in \mathbb{R}^n$, then $|e_\ell(x)|^2 \leq 10 \binom{n}{\ell}$;

→ If $x \in \mathbb{C}^n$, a weaker upper bound holds:

$$|e_\ell(x)|^2 \leq \frac{n^n}{\ell^\ell (n-\ell)^{n-\ell}} < 3\sqrt{\ell+1} \binom{n}{\ell}.$$

→ similar problems have been recently studied in [\[Gopalan and Yehudayoff'14; Meka, Reingold, and Tal'19; Doron, Hatami, and Hoza'20; Tao'23\]](#)

→ best known bound due to [\[Tao'23\]](#):

$$|e_\ell(x)|^2 \leq \binom{n}{\ell}^2 \left(\frac{\ell-1}{n-1} \right)^\ell \leq e^\ell \binom{n}{\ell}$$

→ we crucially need to improve the base e to the best possible constant 1

Proof of the inequality

For the real case, can argue via the method of Lagrangian multipliers that the maximizer x^* is only supported on two points, i.e. it suffices to consider $x = x^{(k)}$ for some k :

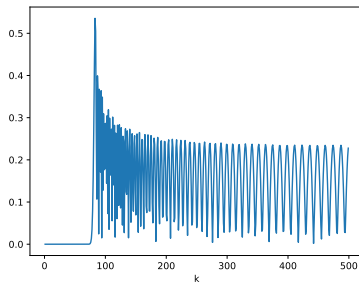
$$x^{(k)} = \left(\underbrace{\sqrt{\frac{k}{n-k}}, \dots, \sqrt{\frac{k}{n-k}}}_{n-k \text{ copies}}, \underbrace{-\sqrt{\frac{n-k}{k}}, \dots, -\sqrt{\frac{n-k}{k}}}_{k \text{ copies}} \right)$$

Proof of the inequality

For the real case, can argue via the method of Lagrangian multipliers that the maximizer x^* is only supported on two points, i.e. it suffices to consider $x = x^{(k)}$ for some k :

$$x^{(k)} = \left(\underbrace{\sqrt{\frac{k}{n-k}}, \dots, \sqrt{\frac{k}{n-k}}}_{n-k \text{ copies}}, \underbrace{-\sqrt{\frac{n-k}{k}}, \dots, -\sqrt{\frac{n-k}{k}}}_{k \text{ copies}} \right)$$

However, upper bounding $|e_\ell(x^{(k)})|$ is still very challenging!!



The quantity $|e_\ell(x^{(k)})|^2 / \binom{n}{\ell}$ vs. k for $n = 1000, \ell = 300$.

Saddle point analysis

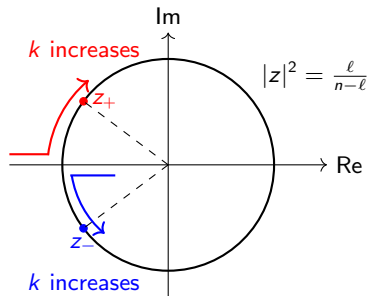
Cauchy's formula :
$$e_{\ell}(x) = \frac{1}{2\pi i} \oint_{|z|=r} \frac{\prod_{i=1}^n (1 + x_i z)}{z^{\ell}} \frac{dz}{z}$$

Saddle point equation :
$$\frac{\ell}{z} = \sum_{i=1}^n \frac{x_i}{1 + x_i z}$$

Saddle point analysis

Cauchy's formula :
$$e_\ell(x) = \frac{1}{2\pi i} \oint_{|z|=r} \frac{\prod_{i=1}^n (1 + x_i z)}{z^\ell} \frac{dz}{z}$$

Saddle point equation :
$$\frac{\ell}{z} = \sum_{i=1}^n \frac{x_i}{1 + x_i z}$$

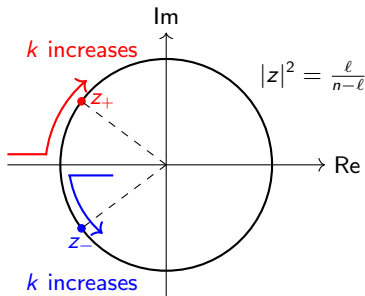


Saddle points for $x = x^{(k)}$

Saddle point analysis

Cauchy's formula :
$$e_{\ell}(x) = \frac{1}{2\pi i} \oint_{|z|=r} \frac{\prod_{i=1}^n (1 + x_i z)}{z^{\ell}} \frac{dz}{z}$$

Saddle point equation :
$$\frac{\ell}{z} = \sum_{i=1}^n \frac{x_i}{1 + x_i z}$$



Saddle points for $x = x^{(k)}$

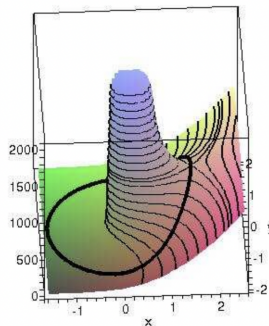


Illustration of saddle point method

Application of saddle point method

Saddle points suggest the contour choice of $\Gamma = \{z : |z| = r\}$ with $r = \sqrt{\frac{\ell}{n-\ell}}$:

$$|e_\ell(x)| = \left| \frac{1}{2\pi i} \oint_\Gamma \frac{\prod_{i=1}^n (1 + x_i z)}{z^\ell} \frac{dz}{z} \right| \leq \max_{|z|=r} \left| \frac{\prod_{i=1}^n (1 + x_i z)}{z^\ell} \right|$$

Application of saddle point method

Saddle points suggest the contour choice of $\Gamma = \{z : |z| = r\}$ with $r = \sqrt{\frac{\ell}{n-\ell}}$:

$$|e_\ell(x)| = \left| \frac{1}{2\pi i} \oint_{\Gamma} \frac{\prod_{i=1}^n (1 + x_i z)}{z^\ell} \frac{dz}{z} \right| \leq \max_{|z|=r} \left| \frac{\prod_{i=1}^n (1 + x_i z)}{z^\ell} \right|$$

Use AM-GM:

$$\begin{aligned} \prod_{i=1}^n |1 + x_i z|^2 &= \prod_{i=1}^n (1 + 2\Re(x_i z) + |x_i|^2 r^2) \\ &\leq \left(\frac{1}{n} \sum_{i=1}^n (1 + 2\Re(x_i z) + |x_i|^2 r^2) \right)^n = (1 + r^2)^n. \end{aligned}$$

This proves the inequality for the complex case.

Application of saddle point method

Saddle points suggest the contour choice of $\Gamma = \{z : |z| = r\}$ with $r = \sqrt{\frac{\ell}{n-\ell}}$:

$$|e_\ell(x)| = \left| \frac{1}{2\pi i} \oint_\Gamma \frac{\prod_{i=1}^n (1 + x_i z)}{z^\ell} \frac{dz}{z} \right| \leq \max_{|z|=r} \left| \frac{\prod_{i=1}^n (1 + x_i z)}{z^\ell} \right|$$

Use AM-GM:

$$\begin{aligned} \prod_{i=1}^n |1 + x_i z|^2 &= \prod_{i=1}^n (1 + 2\Re(x_i z) + |x_i|^2 r^2) \\ &\leq \left(\frac{1}{n} \sum_{i=1}^n (1 + 2\Re(x_i z) + |x_i|^2 r^2) \right)^n = (1 + r^2)^n. \end{aligned}$$

This proves the inequality for the complex case.

Real case: a more careful saddle point analysis for $x = x^{(k)}$.

Compound decisions and empirical Bayes

Empirical Bayes

The empirical Bayes (EB) framework [\[Robbins'51; '56\]](#):

- idea: estimate the prior distribution from data
- lots of empirical successes but limited theoretical understanding

Empirical Bayes

The empirical Bayes (EB) framework [Robbins'51; '56]:

- idea: estimate the prior distribution from data
- lots of empirical successes but limited theoretical understanding

A competitive paradigm [Hannan and Robbins'55; Zhang'03; Greenshtein and Ritov'09; Efron'19]:

- compound decision setting: independent $X_i \sim P_{\theta_i}$, aim to estimate $\theta = (\theta_1, \dots, \theta_n)$
- target: find an estimator with a small regret compared with powerful oracles

$$\text{regret}(\hat{\theta}) = \sup_{\theta} \left(\mathbb{E}_{\theta}[L(\theta, \hat{\theta})] - \inf_{\hat{\theta}^{\text{oracle}}} \mathbb{E}_{\theta}[L(\theta, \hat{\theta}^{\text{oracle}})] \right)$$

- simple/separable oracle: best estimator in the form $\hat{\theta}_i^{\text{S}} = f(X_i)$ for a single function f
- permutation invariant oracle: best estimator in the form

$$\hat{\theta}_{\pi(i)}^{\text{PI}}(X_{\pi(1)}, \dots, X_{\pi(n)}) = \hat{\theta}_i^{\text{PI}}(X_1, \dots, X_n)$$

Empirical Bayes

The empirical Bayes (EB) framework [Robbins'51; '56]:

- idea: estimate the prior distribution from data
- lots of empirical successes but limited theoretical understanding

A competitive paradigm [Hannan and Robbins'55; Zhang'03; Greenshtein and Ritov'09; Efron'19]:

- compound decision setting: independent $X_i \sim P_{\theta_i}$, aim to estimate $\theta = (\theta_1, \dots, \theta_n)$
- target: find an estimator with a small regret compared with powerful oracles

$$\text{regret}(\hat{\theta}) = \sup_{\theta} \left(\mathbb{E}_{\theta}[L(\theta, \hat{\theta})] - \inf_{\hat{\theta}^{\text{oracle}}} \mathbb{E}_{\theta}[L(\theta, \hat{\theta}^{\text{oracle}})] \right)$$

- simple/separable oracle: best estimator in the form $\hat{\theta}_i^{\text{S}} = f(X_i)$ for a single function f
- permutation invariant oracle: best estimator in the form

$$\hat{\theta}_{\pi(i)}^{\text{PI}}(X_{\pi(1)}, \dots, X_{\pi(n)}) = \hat{\theta}_i^{\text{PI}}(X_1, \dots, X_n)$$

Question

Can we apply a “mean-field” approximation of the complicated $\hat{\theta}^{\text{PI}}$ by the simple $\hat{\theta}^{\text{S}}$?

The Gaussian case

- observation vector: $X^n \sim \mathcal{N}(\theta^n, I_n)$
- a postulated Bayes model: θ^n is a uniform permutation of a given multiset $\{\theta_1^*, \dots, \theta_n^*\}$
- under the quadratic loss:

$$\hat{\theta}_i^S = \mathbb{E}[\theta_i \mid X_i], \quad \hat{\theta}_i^{\text{PI}} = \mathbb{E}[\theta_i \mid X^n].$$

The Gaussian case

- observation vector: $X^n \sim \mathcal{N}(\theta^n, I_n)$
- a postulated Bayes model: θ^n is a uniform permutation of a given multiset $\{\theta_1^*, \dots, \theta_n^*\}$
- under the quadratic loss:

$$\hat{\theta}_i^S = \mathbb{E}[\theta_i \mid X_i], \quad \hat{\theta}_i^{\text{PI}} = \mathbb{E}[\theta_i \mid X^n].$$

Greenshtein and Ritov (2009)

If $\|\theta^*\|_\infty \leq \mu$ with $\mu \geq 1$,

$$\mathbb{E} \left[\|\hat{\theta}^S - \hat{\theta}^{\text{PI}}\|^2 \right] = e^{O(\mu^2)}.$$

- an $O(1)$ upper bound even if the vectors are n -dimensional
- becomes meaningless when $\mu \gg \sqrt{\log n}$

A tight upper bound

Theorem ([H., Niles-Weed, Shen, Wu'25])

If $\|\theta^*\|_\infty \leq \mu$ with $\mu \geq 1$,

$$\mathbb{E} \left[\|\hat{\theta}^S - \hat{\theta}^{PI}\|^2 \right] = O(\mu \log^2 n).$$

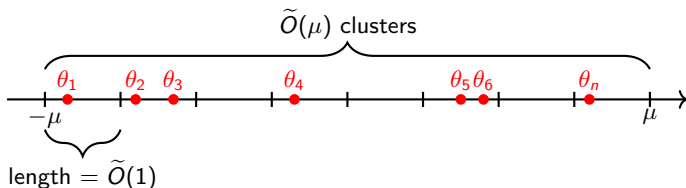
A tight upper bound

Theorem ([H., Niles-Weed, Shen, Wu'25])

If $\|\theta^*\|_\infty \leq \mu$ with $\mu \geq 1$,

$$\mathbb{E} \left[\|\hat{\theta}^S - \hat{\theta}^{PI}\|^2 \right] = O(\mu \log^2 n).$$

Optimal dependence on μ , which is the number of “subproblems”:



- by the concentration of Gaussian, each interval roughly corresponds to an independent subproblem
- overall problem is a “direct sum” of subproblems

Application: Competitive Distribution Estimation

A Poisson sequence model:

$$(N_1, \dots, N_k) \sim \text{Poi}(np_1) \otimes \dots \otimes \text{Poi}(np_k)$$

- n : sample size
- k : support size
- $p = (p_1, \dots, p_k)$: an unknown probability vector

Competitive distribution estimation

A Poisson sequence model:

$$(N_1, \dots, N_k) \sim \text{Poi}(np_1) \otimes \dots \otimes \text{Poi}(np_k)$$

- n : sample size
- k : support size
- $p = (p_1, \dots, p_k)$: an unknown probability vector

Competitive distribution estimation: based on observed counts (N_1, \dots, N_k) , devise an estimator \hat{p} to minimize the **KL regret**:

$$\text{regret}(\hat{p}) = \sup_p \mathbb{E} \left[\text{KL}(p \| \hat{p}) - \text{KL}(p \| \hat{p}^{\text{PI}}) \right],$$

where \hat{p}^{PI} is the best permutation-invariant decision rule which knows the ground truth p

“Why is Good–Turing Good”

Upper bound ([Orlitsky and Suresh'15])

A modified Good–Turing estimator \hat{p}^{MGT} achieves

$$\text{regret}(\hat{p}^{\text{MGT}}) = \tilde{O}\left(\min\left\{\frac{k}{n}, \frac{1}{\sqrt{n}}\right\}\right).$$

“Why is Good–Turing Good”

Upper bound ([Orlitsky and Suresh'15])

A modified Good–Turing estimator \hat{p}^{MGT} achieves

$$\text{regret}(\hat{p}^{\text{MGT}}) = \tilde{O} \left(\min \left\{ \frac{k}{n}, \frac{1}{\sqrt{n}} \right\} \right).$$

The Good–Turing estimator \hat{p}^{GT} [Good'53]: for $N_i = y$,

$$\hat{p}_i^{\text{GT}} = \frac{y+1}{n} \cdot \frac{\sum_{j=1}^k 1(N_j = y+1)}{\sum_{j=1}^k 1(N_j = y)}$$

“Why is Good–Turing Good”

Upper bound ([Orlitsky and Suresh'15])

A modified Good–Turing estimator \hat{p}^{MGT} achieves

$$\text{regret}(\hat{p}^{\text{MGT}}) = \tilde{O} \left(\min \left\{ \frac{k}{n}, \frac{1}{\sqrt{n}} \right\} \right).$$

The Good–Turing estimator \hat{p}^{GT} [Good'53]: for $N_i = y$,

$$\hat{p}_i^{\text{GT}} = \frac{y+1}{n} \cdot \frac{\sum_{j=1}^k 1(N_j = y+1)}{\sum_{j=1}^k 1(N_j = y)}$$

Lower bound ([Orlitsky and Suresh'15])

$$\inf_{\hat{p}} \text{regret}(\hat{p}) = \Omega \left(\min \left\{ \frac{k}{n}, \frac{1}{n^{2/3}} \right\} \right).$$

Better Good–Turing: NPMLE

Our estimator is similar to [\[Jiang and Zhang'09\]](#) and relies on two cornerstones:

- **EB**: think of $p_1, \dots, p_k \stackrel{\text{i.i.d.}}{\sim} G^*$, with the empirical measure $G^* = \frac{1}{k} \sum_{i=1}^k \delta_{p_i}$
- **Nonparametric MLE (NPMLE)** [\[Kiefer and Wolfowitz'56\]](#): a natural estimator for G^* maximizes the marginal likelihood

$$\hat{G} = \operatorname{argmax}_G \sum_{i=1}^k \log \mathbb{E}_G [\mathbb{P}(\text{Poi}(np) = N_i)]$$

- the final estimator \hat{p}^{NPMLE} is the Bayes rule under the “data-driven prior” \hat{G} :

$$\hat{p}^{\text{NPMLE}} = \text{normalized version of } (\mathbb{E}_{\hat{G}}[p_1 \mid N_1], \dots, \mathbb{E}_{\hat{G}}[p_k \mid N_k])$$

Better Good–Turing: NPMLE

Our estimator is similar to [Jiang and Zhang'09] and relies on two cornerstones:

- **EB**: think of $p_1, \dots, p_k \stackrel{\text{i.i.d.}}{\sim} G^*$, with the empirical measure $G^* = \frac{1}{k} \sum_{i=1}^k \delta_{p_i}$
- **Nonparametric MLE (NPMLE)** [Kiefer and Wolfowitz'56]: a natural estimator for G^* maximizes the marginal likelihood

$$\hat{G} = \operatorname{argmax}_G \sum_{i=1}^k \log \mathbb{E}_G [\mathbb{P}(\text{Poi}(np) = N_i)]$$

- the final estimator \hat{p}^{NPMLE} is the Bayes rule under the “data-driven prior” \hat{G} :

$$\hat{p}^{\text{NPMLE}} = \text{normalized version of } (\mathbb{E}_{\hat{G}}[p_1 \mid N_1], \dots, \mathbb{E}_{\hat{G}}[p_k \mid N_k])$$

Efficient, tuning parameter-free, and optimal competitive guarantee:

Theorem (H., Niles-Weed, Shen, Wu'25)

The above estimator \hat{p}^{NPMLE} achieves

$$\text{regret}(\hat{p}^{\text{NPMLE}}) = \tilde{O} \left(\min \left\{ \frac{k}{n}, \frac{1}{n^{2/3}} \right\} \right).$$

Part I of regret: \hat{p}^{NPMLE} against the separable oracle

$$\hat{p}^{\text{S}} = \text{normalized version of } (\mathbb{E}_{G^*}[p_1 \mid N_1], \dots, \mathbb{E}_{G^*}[p_k \mid N_k])$$

→ use the theory of NPMLE to argue that $\mathbb{E}_{\hat{G}}[p_i \mid N_i] \approx \mathbb{E}_{G^*}[p_i \mid N_i]$

Part I of regret: \hat{p}^{NPMLE} against the separable oracle

$$\hat{p}^S = \text{normalized version of } (\mathbb{E}_{G^*}[p_1 | N_1], \dots, \mathbb{E}_{G^*}[p_k | N_k])$$

→ use the theory of NPMLE to argue that $\mathbb{E}_{\hat{G}}[p_i | N_i] \approx \mathbb{E}_{G^*}[p_i | N_i]$

Part II of regret: separable oracle \hat{p}^S against the PI oracle \hat{p}^{PI}

→ our technique applied to the Poisson case gives

$$\mathbb{E} \left[\text{KL} \left(\hat{p}^{\text{PI}} \| \hat{p}^S \right) \right] = \frac{\tilde{O}(\# \text{ of subproblems in the Poisson model})}{n}$$

→ it turns out that

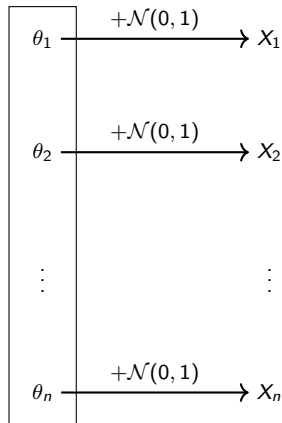
$$\# \text{ of subproblems in the Poisson model} = O \left(\min \left\{ k, n^{1/3} \right\} \right)$$

Proof for the Gaussian case

A (failed) information-theoretic argument

→ Recall that

$$\hat{\theta}_1^S = \mathbb{E}[\theta_1 \mid X_1], \quad \hat{\theta}_1^{\text{PI}} = \mathbb{E}[\theta_1 \mid X^n].$$



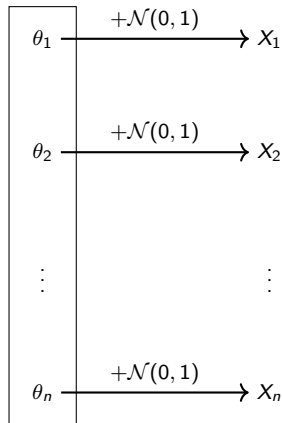
A (failed) information-theoretic argument

→ Recall that

$$\hat{\theta}_1^S = \mathbb{E}[\theta_1 \mid X_1], \quad \hat{\theta}_1^{\text{PI}} = \mathbb{E}[\theta_1 \mid X^n].$$

→ Tao's inequality:

$$\mathbb{E} \left[(\mathbb{E}[\theta_1 \mid X_1] - \mathbb{E}[\theta_1 \mid X^n])^2 \right] = \tilde{O}(1) \cdot I(\theta_1; X_2^n \mid X_1).$$



A (failed) information-theoretic argument

→ Recall that

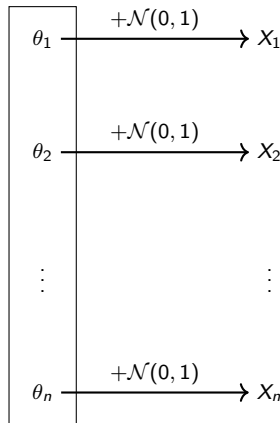
$$\hat{\theta}_1^S = \mathbb{E}[\theta_1 | X_1], \quad \hat{\theta}_1^{\text{PI}} = \mathbb{E}[\theta_1 | X^n].$$

→ Tao's inequality:

$$\mathbb{E} \left[(\mathbb{E}[\theta_1 | X_1] - \mathbb{E}[\theta_1 | X^n])^2 \right] = \tilde{O}(1) \cdot I(\theta_1; X_2^n | X_1).$$

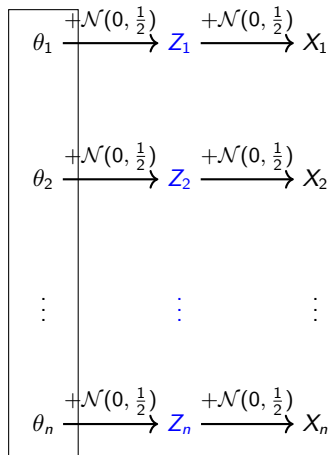
→ A “model-free” upper bound:

$$\begin{aligned} I(\theta_1; X_2^n | X_1) &= H(\theta_1 | X_1) - H(\theta_1 | X^n) \\ &\leq H(\theta_1 | X_1) - \frac{1}{n} H(\theta^n | X^n) \\ &= H(\theta_1) - \frac{H(\theta^n)}{n} - \underbrace{\left(I(\theta_1; X_1) - \frac{I(\theta^n; X^n)}{n} \right)}_{\geq 0 \text{ as } P_{X^n|\theta^n} = \prod_i P_{X_i|\theta_i}} \\ &\leq H(\theta_1) - \frac{H(\theta^n)}{n} = \frac{1}{n} \text{KL}(P_{\theta^n} \| \prod_i P_{\theta_i}) \\ &= \tilde{O} \left(\frac{|\text{supp}(\{\theta_1, \dots, \theta_n\})|}{n} \right) \end{aligned}$$



Improvement via “noisy” θ^n

→ idea: add a noisy Z_i between θ_i and X_i

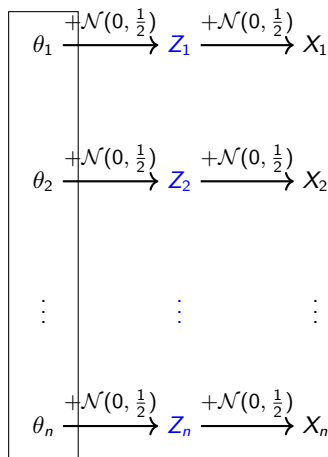


Improvement via “noisy” θ^n

→ idea: add a noisy Z_i between θ_i and X_i

→ key identity:

$$\begin{aligned}\mathbb{E}[\theta_1 \mid X_1] - \mathbb{E}[\theta_1 \mid X^n] \\ = 2 (\mathbb{E}[Z_1 \mid X_1] - \mathbb{E}[Z_1 \mid X^n])\end{aligned}$$



Improvement via “noisy” θ^n

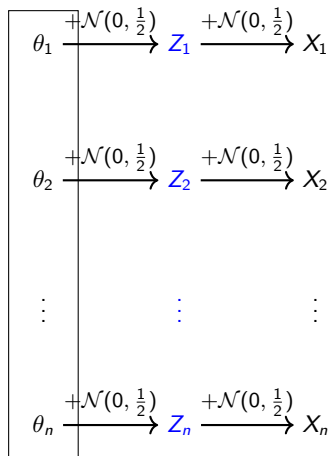
→ idea: add a noisy Z_i between θ_i and X_i

→ key identity:

$$\begin{aligned}\mathbb{E}[\theta_1 \mid X_1] - \mathbb{E}[\theta_1 \mid X^n] \\ = 2 (\mathbb{E}[Z_1 \mid X_1] - \mathbb{E}[Z_1 \mid X^n])\end{aligned}$$

→ the previous “model-free” bound now gives

$$\mathbb{E} \left[(\mathbb{E}[\theta_1 \mid X_1] - \mathbb{E}[\theta_1 \mid X^n])^2 \right] \lesssim \frac{1}{n} \text{KL}(P_{Z^n} \parallel \prod_i P_{Z_i})$$



Improvement via “noisy” θ^n

→ idea: add a noisy Z_i between θ_i and X_i

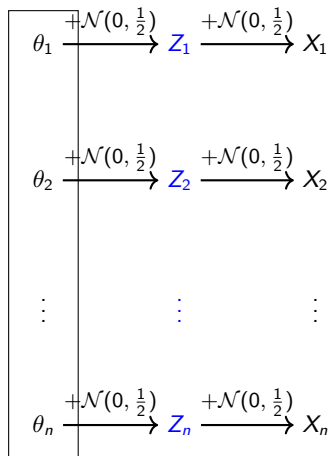
→ key identity:

$$\begin{aligned}\mathbb{E}[\theta_1 \mid X_1] - \mathbb{E}[\theta_1 \mid X^n] \\ = 2 (\mathbb{E}[Z_1 \mid X_1] - \mathbb{E}[Z_1 \mid X^n])\end{aligned}$$

→ the previous “model-free” bound now gives

$$\mathbb{E} \left[(\mathbb{E}[\theta_1 \mid X_1] - \mathbb{E}[\theta_1 \mid X^n])^2 \right] \lesssim \frac{1}{n} \text{KL}(P_{Z^n} \parallel \prod_i P_{Z_i})$$

→ the final quantity $\text{KL}(P_{Z^n} \parallel \prod_i P_{Z_i})$ is now between a Gaussian permutation mixture and its i.i.d. approximation!



Concluding remarks

Take home messages:

- permutations induce weak dependency, quantitatively
- centered basis is preferred in the method of “moments”
- NPMLE + EB outperforms Good–Turing

Concluding remarks

Take home messages:

- permutations induce weak dependency, quantitatively
- centered basis is preferred in the method of “moments”
- NPMLE + EB outperforms Good–Turing

Further questions:

- method of “moments” for two high-dimensional mixtures?
- a better understanding of the noisy Z ? non-divisible distribution?

Concluding remarks

Take home messages:

- permutations induce weak dependency, quantitatively
- centered basis is preferred in the method of “moments”
- NPMLE + EB outperforms Good–Turing

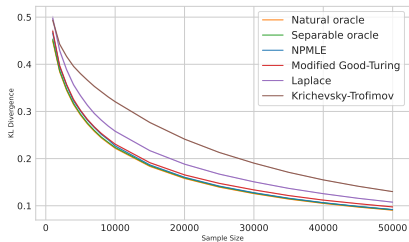
Further questions:

- method of “moments” for two high-dimensional mixtures?
- a better understanding of the noisy Z ? non-divisible distribution?

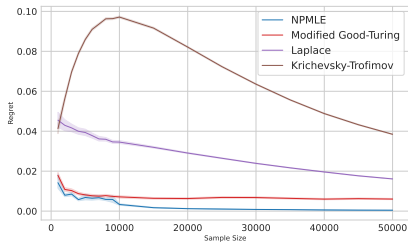
Thank You!

Backup Slides

Experiments on sqrt-Cauchy distribution

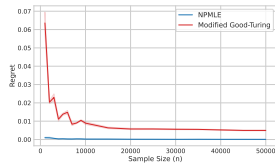


(a) KL risks.

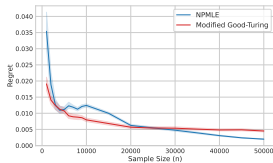


(b) Regret over the separable oracle.

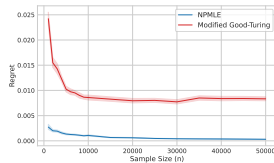
Experiments on more distributions



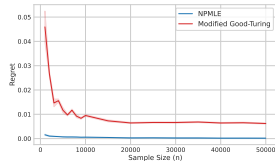
(a) Uniform



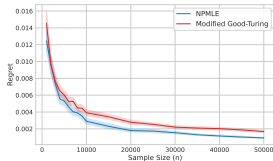
(b) Zipf ($\alpha = 1$).



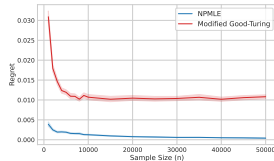
(c) Dirichlet ($c = 1$)



(d) Step.

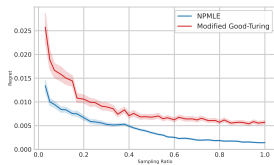


(e) Zipf ($\alpha = 1.5$).

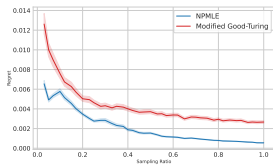


(f) Dirichlet ($c = 0.5$)

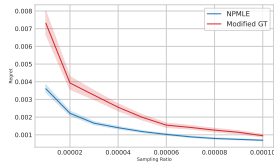
Experiments on real data



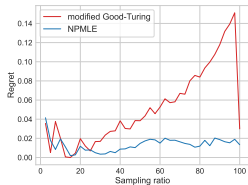
(a) Hamlet (random).



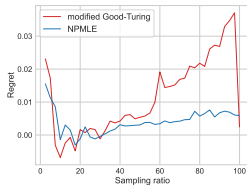
(b) LOTR (random).



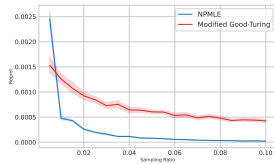
(c) 2020 Census Detailed DHC-A.



(d) Hamlet (consecutive).



(e) LOTR (consecutive).



(f) 2010 Census surname.

An alternative view from matrix permanent

Drawbacks of the first upper bound:

- meaningless when $C_{\chi^2}(\mathcal{P}) \geq 1$
- why loose: Banach's inequality may overlook the benefits from different rows

An observation thanks to permutations:

χ^2 divergence as matrix permanents

$$\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) = \frac{n^n}{n!} \text{Perm}(A) - 1,$$

where $A \in \mathbb{R}^{n \times n}$ is given by $A_{i,j} = \mathbb{E}_{\bar{P}} \left[\frac{dP_i}{d\bar{P}} \frac{dP_j}{d\bar{P}} \right]$.

The famous van der Waerden conjecture (proven in 1980's) states that $\text{Perm}(A) \geq \frac{n!}{n^n}$ for all doubly stochastic matrices, so showing $\chi^2(\mathbb{P}_n \| \mathbb{Q}_n) = O(1)$ essentially means that $\text{Perm}(A)$ is nearly as small as possible

Properties of matrix A

Properties of A

- A is PSD and doubly stochastic;
- $\text{Tr}(A) \leq C_{\chi^2}(\mathcal{P}) + 1$;
- its spectral gap satisfies $1 - \lambda_2(A) \geq \frac{1}{D_{\chi^2}(\mathcal{P}) + 1}$.

Suggests to use the eigendecomposition $A = UDU^\top$ and expand

$$\frac{n^n}{n!} \text{Perm}(UDU^\top) = \sum_{\ell=0}^n S_\ell(\lambda_2, \dots, \lambda_n),$$

with homogeneous polynomials S_ℓ of total degree ℓ

Key idea: express S_ℓ using **complex normal random variables**

Expressing the sum $\sum_{\ell=0}^n S_{\ell}$

Complex normal random variable:

- $z \sim \mathcal{CN}(0, 1)$ iff $z = x + iy$ with independent $x, y \sim \mathcal{N}(0, \frac{1}{2})$
- moment condition: $\mathbb{E}[z^m \bar{z}^n] = n! \mathbb{1}_{m=n}$ for $z \sim \mathcal{CN}(0, 1)$

Fact 1 ([Gurvit'03])

$$\sum_{\ell=0}^n S_{\ell} \propto \mathbb{E} \left[\prod_{i=1}^n \left| \left(UD^{1/2} z \right)_i \right|^2 \right], \quad z_1, \dots, z_n \sim \mathcal{CN}(0, 1).$$

Applying AM-GM to the product gives

$$\sum_{\ell=0}^n S_{\ell} \leq \sum_{\ell_2 + \dots + \ell_n \leq n} \lambda_2^{\ell_2} \dots \lambda_n^{\ell_n} \leq \prod_{i=2}^n \frac{1}{1 - \lambda_i}$$

- the trace and spectral gap properties lead to the second upper bound

Expressing the individual term S_ℓ

Fact II

$$S_\ell \propto \mathbb{E} \left[\left| e_\ell \left((\tilde{U} \tilde{D}^{1/2} z)_1, \dots, (\tilde{U} \tilde{D}^{1/2} z)_n \right) \right|^2 \right], \quad z_1, \dots, z_{n-1} \sim \mathcal{CN}(0, 1),$$

where (\tilde{U}, \tilde{D}) takes out the leading eigenvector/eigenvalue in (U, D) .

- can show that the vector $\tilde{U} \tilde{D}^{1/2} z$ sums into zero
- using our key inequality eventually leads to

$$S_\ell \leq 3\sqrt{\ell+1} \sum_{\ell_2+\dots+\ell_n=\ell} \lambda_2^{\ell_2} \dots \lambda_n^{\ell_n}$$

recall that $e_\ell(x_1, \dots, x_n) = \sum_{|S|=\ell} \prod_{i \in S} x_i$.