

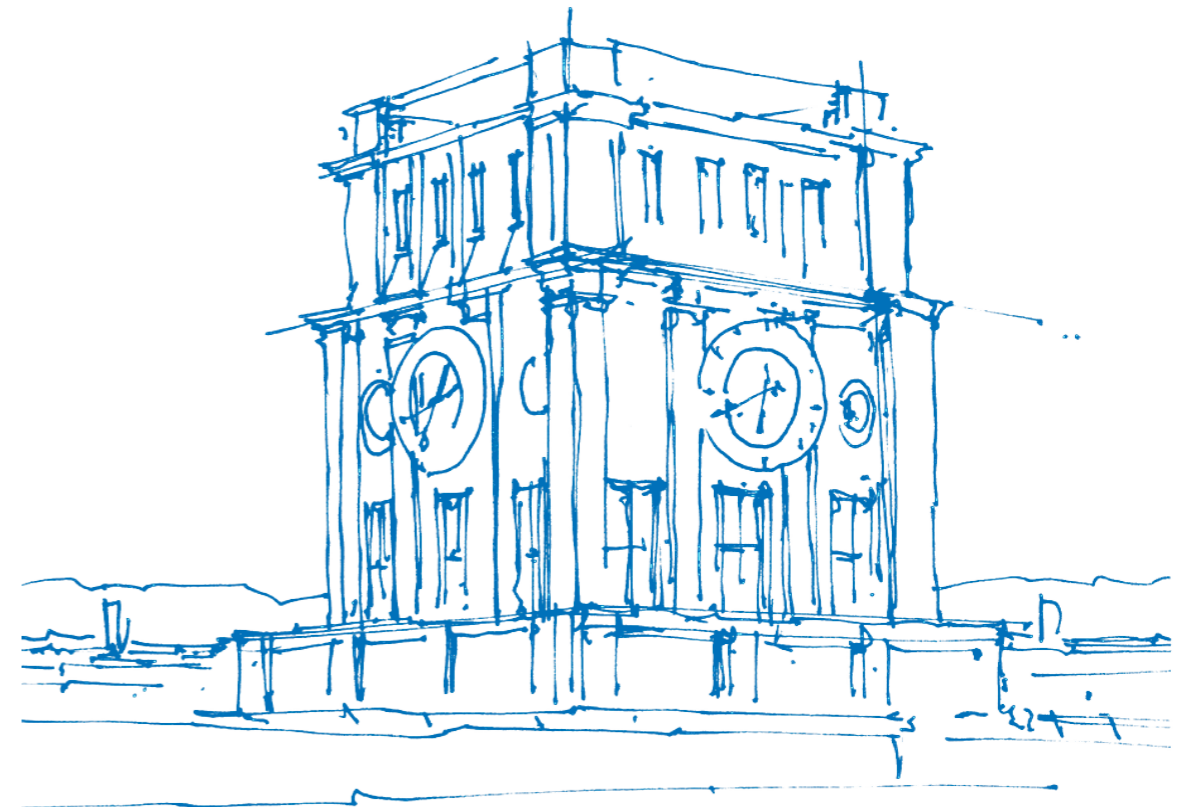
# On universally consistent and fully distribution-free rank tests of vector independence

Hongjian Shi

Technical University of Munich

ETH-UCPH-TUM Workshop

12 October 2022



*TUM Uhrenturm*

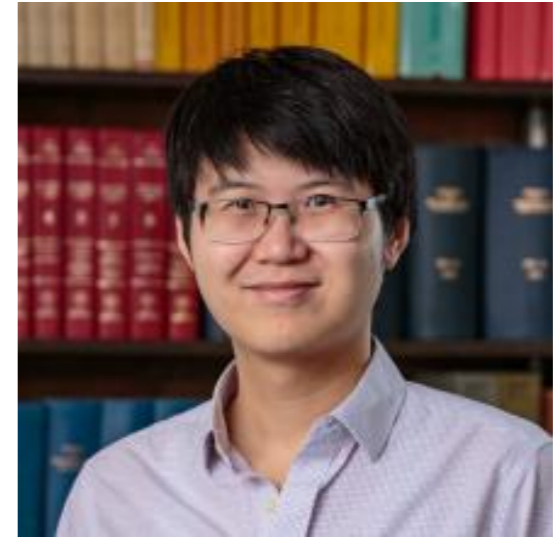
# Collaborators



Marc Hallin



Mathias Drton



Fang Han

# Problem

We consider testing independence of two random vectors with absolutely continuous distributions based on finite observations.

$$\mathbf{X} = \underbrace{(X_1, X_2, \dots, X_p)^\top}_{p \text{ covariates}} \quad \mathbf{Y} = \underbrace{(Y_1, Y_2, \dots, Y_q)^\top}_{q \text{ covariates}}$$

$$H_0 : \mathbf{X} \text{ is independent of } \mathbf{Y}$$

---

Data:

$$\begin{pmatrix} (\mathbf{X}_1, \mathbf{Y}_1) \\ (\mathbf{X}_2, \mathbf{Y}_2) \\ \vdots \\ (\mathbf{X}_n, \mathbf{Y}_n) \end{pmatrix}$$

$n$  independent copies of  $(\mathbf{X}, \mathbf{Y})$

# Paradigm

Criteria:

- The test should be **distribution-free**, and directly implementable without the need of permutation.
- The test should be **consistent** in a certain sense.

a long-standing problem

- The test should be **optimal** under certain standard.

rate-optimality

- ~~■ The dimension  $p$  should be allowed to be **much larger**  
**than the sample size  $n$ .**~~

future work

# Bivariate case

- **Data:**  $\{(X_i \in \mathbb{R}, Y_i \in \mathbb{R}), i \in [n]\}$  i.i.d.
  - **Aim:** testing if “ $H_0 : X \perp\!\!\!\perp Y$ ” is true.
- 

The test should be **distribution-free**

The ranks of  $X_1, \dots, X_n$  are **uniformly distributed** on the set of all permutations of  $[n]$ .

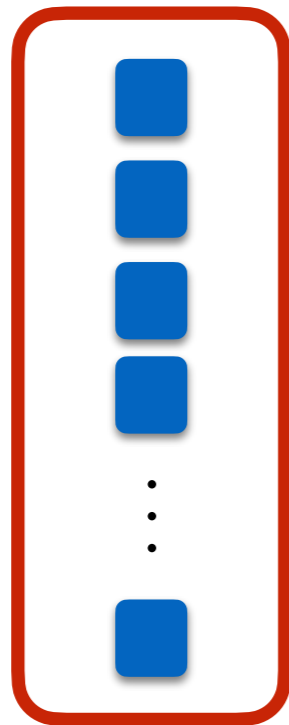
Under  $H_0$ , the marginal ranks of  $\{X_i, i \in [n]\}$  and  $\{Y_i, i \in [n]\}$  are **independent**.

For any test statistic based on ranks, its null distribution is both **determined** and **independent** of  $P_{X,Y}$ , i.e., the test is **distribution-free**.

# (Marginal) rank tests?

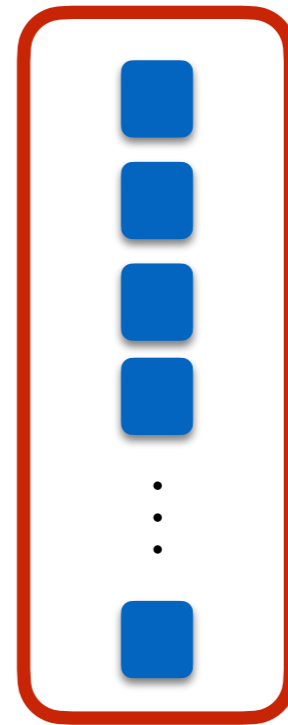
- **Data:**  $\{(\mathbf{X}_i \in \mathbb{R}^p, \mathbf{Y}_i \in \mathbb{R}^q), i \in [n]\}$  i.i.d.
  - **Aim:** testing if “ $H_0 : \mathbf{X} \perp\!\!\!\perp \mathbf{Y}$ ” is true.
- 

Tests built on marginal ranks are no longer distribution-free:



ranks of  $\{X_{1,j}, \dots, X_{n,j}\}$

possibly correlated with



ranks of  $\{X_{1,k}, \dots, X_{n,k}\}$

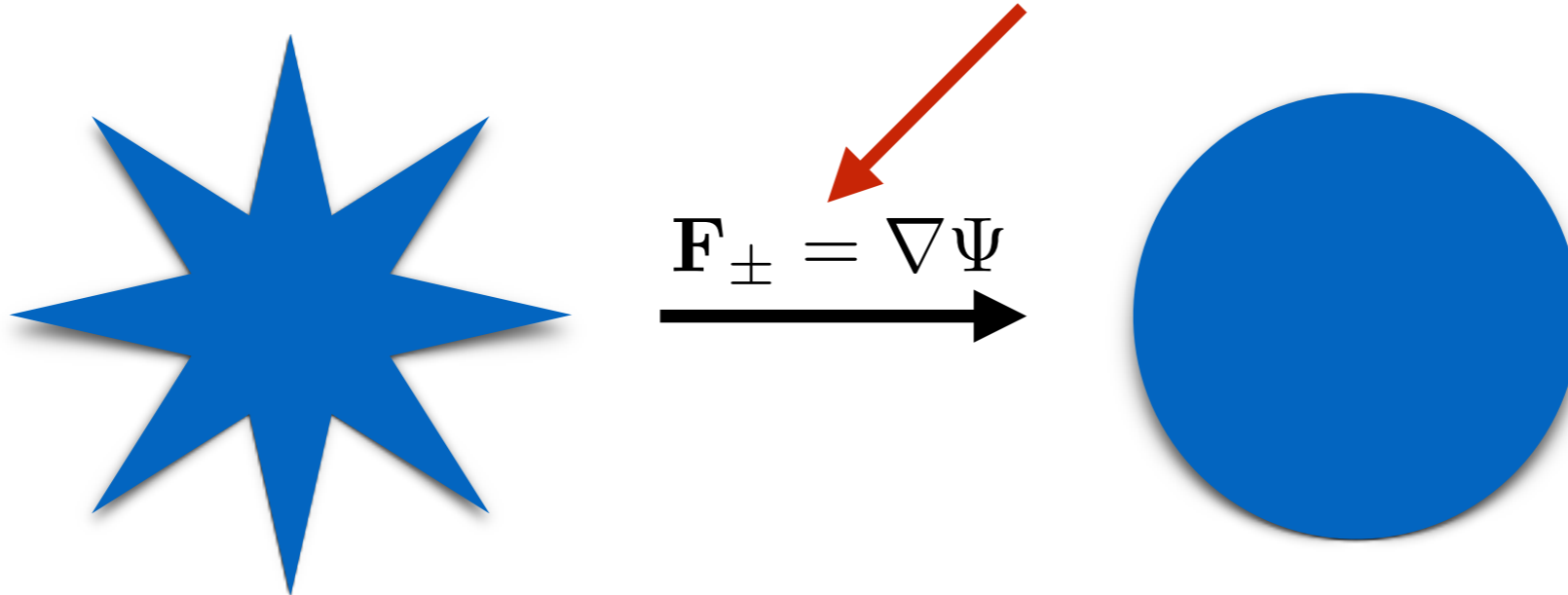
# Outline

- **Center-outward multivariate rank**
- **The proposed test**
- **Discussion**

# Solution

- Center-outward **CDF** in general dimension

Center-outward population distribution function



$P_d$

a general absolutely continuous probability measure in dimension  $d$

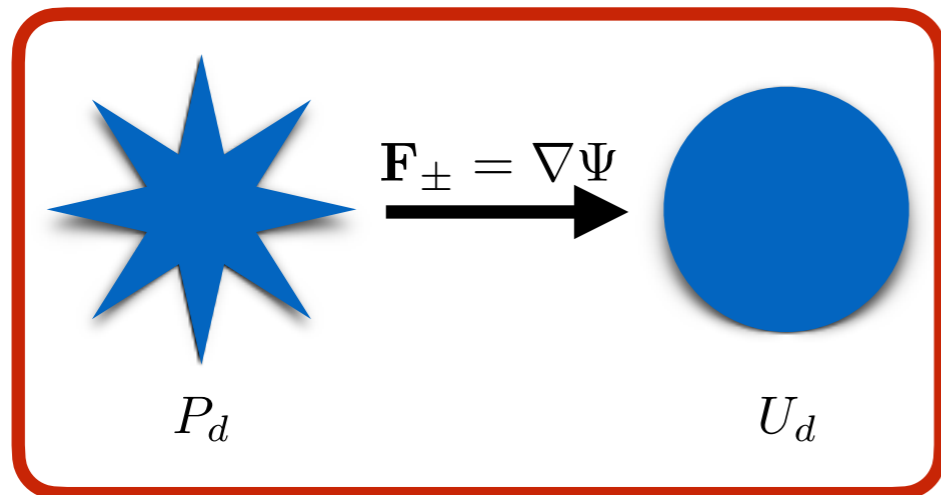
$U_d$

spherical uniform measure over  $d$ -dimensional unit ball



# Solution

- Center-outward **CDF** in general dimension



$$\inf_T \int_{\mathbb{R}^d} \left\| T(\mathbf{x}) - \mathbf{x} \right\|_2^2 dP_d$$

$$\text{subject to } T_\# P_d = U_d$$

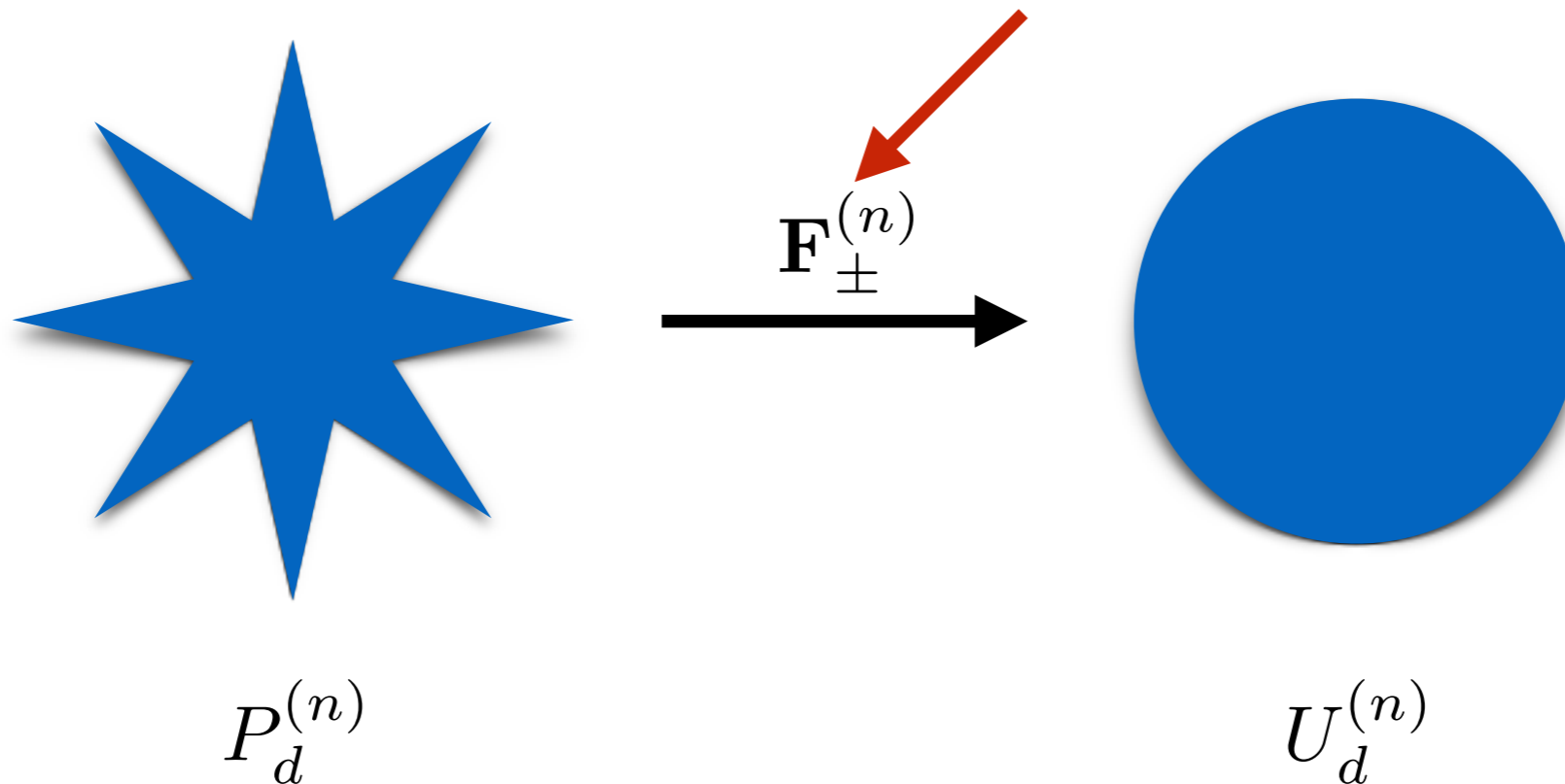
optimal transport problem

- Existence and uniqueness: Main Theorem in [McCann \(1995\)](#)

# Solution

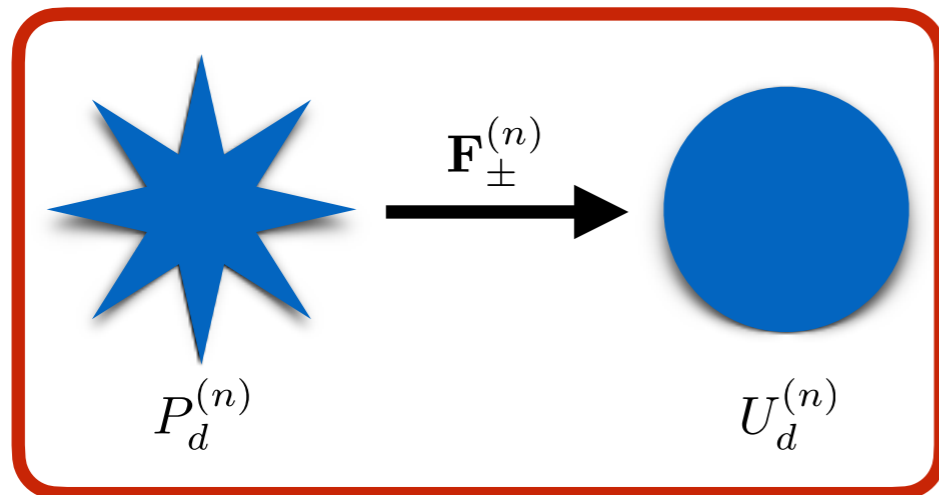
- Center-outward **Empirical CDF** in general dimension

Center-outward empirical distribution function



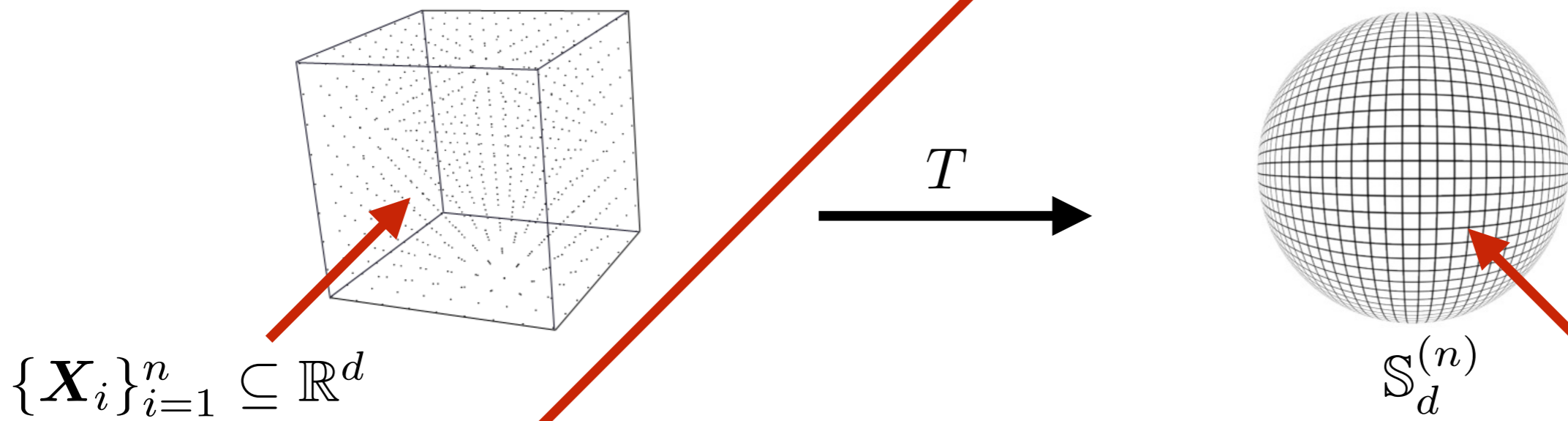
# Solution

- Center-outward **Empirical CDF** in general dimension



$$\mathbf{F}_{\pm}^{(n)} := \operatorname{argmin}_{T \in \mathcal{T}} \sum_{i=1}^n \left\| \mathbf{X}_i - T(\mathbf{X}_i) \right\|_2^2$$

assignment problem

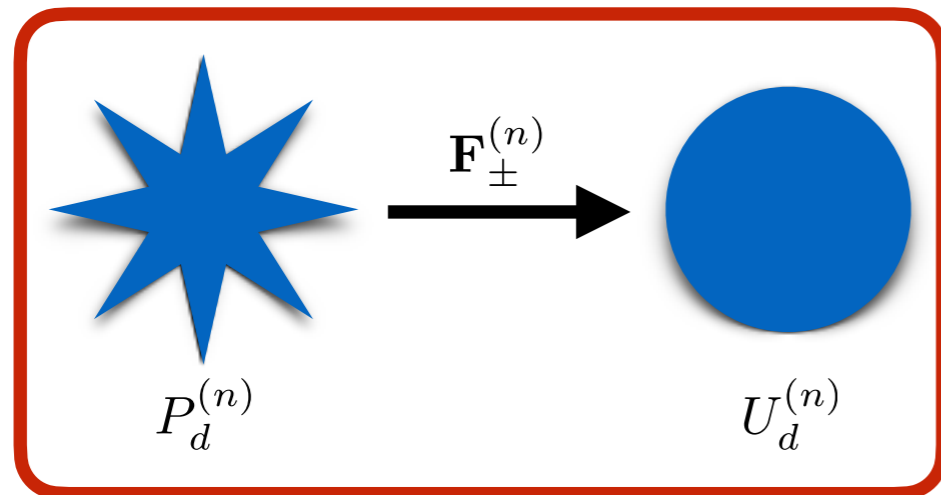


the collection of all bijective mappings between  $\{\mathbf{X}_i\}_{i=1}^n$  and  $\{\mathbf{u}_i\}_{i=1}^n$

consisting of  $n$  points  $\{\mathbf{u}_i\}_{i=1}^n \subseteq \mathbb{R}^d$  that approximate the spherical uniform measure over the unit ball

# Solution

- Center-outward **Empirical CDF** in general dimension



$$\mathbf{F}_{\pm}^{(n)} := \operatorname{argmin}_{T \in \mathcal{T}} \sum_{i=1}^n \left\| \mathbf{X}_i - T(\mathbf{X}_i) \right\|_2^2$$

assignment problem

Center-outward multi-rank is **strongly consistent** and **distribution-free**:

**Hallin et al. (2021)**: Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be i.i.d. with absolutely continuous distribution  $P_d$ . Then

$$\left\| \mathbf{F}_{\pm}^{(n)}(\mathbf{X}_i) - \mathbf{F}_{\pm}(\mathbf{X}_i) \right\|_2 \xrightarrow{\text{a.s.}} 0 \quad \text{as } n \rightarrow \infty,$$

and  $(\mathbf{F}_{\pm}^{(n)}(\mathbf{X}_1), \dots, \mathbf{F}_{\pm}^{(n)}(\mathbf{X}_n))$  is **uniformly distributed over all permutations of  $\mathcal{S}_d^{(n)}$** .

# Outline

- Center-outward multivariate rank
- **The proposed test**
- Discussion

# Distance covariance

- **SRB's insight:**

The squared distance covariance between two random vectors  $\mathbf{X} \in \mathbb{R}^p$  and  $\mathbf{Y} \in \mathbb{R}^q$  with finite first moments, introduced by Székely, Rizzo and Bakirov (2007), is defined as

$$\text{dCov}^2(\mathbf{X}, \mathbf{Y}) := \frac{1}{c_p c_q} \int_{\mathbb{R}^{p+q}} \frac{\|\phi_{\mathbf{X}, \mathbf{Y}}(\mathbf{t}, \mathbf{s}) - \phi_{\mathbf{X}}(\mathbf{t})\phi_{\mathbf{Y}}(\mathbf{s})\|_2^2}{\|\mathbf{t}\|_2^{1+p} \|\mathbf{s}\|_2^{1+q}} d\mathbf{t} d\mathbf{s},$$

where  $\phi_{\mathbf{X}}$ ,  $\phi_{\mathbf{Y}}$  and  $\phi_{\mathbf{X}, \mathbf{Y}}$  are the individual and joint characteristic functions of  $X$  and  $Y$  respectively.

The sample squared distance covariance is defined as

$$\text{dCov}_n^2 \left( (\mathbf{X}_i)_{i=1}^n, (\mathbf{Y}_i)_{i=1}^n \right) := \binom{n}{4}^{-1} \sum_{i_1 \neq \dots \neq i_4} \frac{1}{4 \cdot 4!} g(\mathbf{X}_{i_1}, \mathbf{X}_{i_2}, \mathbf{X}_{i_3}, \mathbf{X}_{i_4}) g(\mathbf{Y}_{i_1}, \mathbf{Y}_{i_2}, \mathbf{Y}_{i_3}, \mathbf{Y}_{i_4}),$$

where  $g(\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3, \mathbf{z}_4)$

$$:= \|\mathbf{z}_1 - \mathbf{z}_2\|_2 + \|\mathbf{z}_3 - \mathbf{z}_4\|_2 - \|\mathbf{z}_1 - \mathbf{z}_3\|_2 - \|\mathbf{z}_2 - \mathbf{z}_4\|_2.$$

# Test

- **Data:**  $\{(\mathbf{X}_i \in \mathbb{R}^p, \mathbf{Y}_i \in \mathbb{R}^q), i \in [n]\}$  i.i.d. distributed with absolutely continuous probability measures  $P_{\mathbf{X}}, P_{\mathbf{Y}}$ .
  - **Aim:** testing if “ $H_0 : \mathbf{X} \perp\!\!\!\perp \mathbf{Y}$ ” is true.
- 

Proposal:

- Calculate **center-outward ranks**  $\mathbf{F}_{\mathbf{X},\pm}^{(n)}(\mathbf{X}_1), \dots, \mathbf{F}_{\mathbf{X},\pm}^{(n)}(\mathbf{X}_n)$  and  $\mathbf{F}_{\mathbf{Y},\pm}^{(n)}(\mathbf{Y}_1), \dots, \mathbf{F}_{\mathbf{Y},\pm}^{(n)}(\mathbf{Y}_n)$ ;
- Combine **center-outward ranks** with **distance covariance**, obtaining the test statistic

$$\widehat{M}_n := n \cdot \text{dCov}_n^2\left(\left(\mathbf{F}_{\mathbf{X},\pm}^{(n)}(\mathbf{X}_i)\right)_{i=1}^n, \left(\mathbf{F}_{\mathbf{Y},\pm}^{(n)}(\mathbf{Y}_i)\right)_{i=1}^n\right);$$

- Reject  $H_0$  if  $\widehat{M}_n$  is large enough.

# Theory

- **Multivariate Hájek asymptotic representation**
  - Consider the “**population**” center-outward signed-ranks  $\mathbf{F}_{\mathbf{X},\pm}(\mathbf{X}_i) \sim U_p$  and  $\mathbf{F}_{\mathbf{Y},\pm}(\mathbf{Y}_i) \sim U_q$ ;
  - Lead to the “**oracle**” test statistic:

$$\widetilde{M}_n := n \cdot \text{dCov}_n^2\left(\left(\mathbf{F}_{\mathbf{X},\pm}(\mathbf{X}_i)\right)_{i=1}^n, \left(\mathbf{F}_{\mathbf{Y},\pm}(\mathbf{Y}_i)\right)_{i=1}^n\right);$$

- Standard exercise (e.g. [Jakobsen \(2017, Theorem 5.10\)](#)) gives, under  $H_0$ ,

$$\widetilde{M}_n \xrightarrow{d} \sum_{k=1}^{\infty} \lambda_k (\xi_k^2 - 1),$$

where  $\lambda_1, \lambda_2, \dots$  are positive constants depending only on  $p, q$ .



# Theory

- **Multivariate Hájek asymptotic representation**

**Main Theorem (SHDH 2022).** Let  $(\mathbf{X}_1, \mathbf{Y}_1), \dots, (\mathbf{X}_n, \mathbf{Y}_n)$  be independent copies of  $(\mathbf{X}, \mathbf{Y})$  with absolutely continuous probability measures  $P_{\mathbf{X}}, P_{\mathbf{Y}}$  and  $\mathbf{X}$  and  $\mathbf{Y}$  are **independent**. Then it holds that

$$\widehat{M}_n - \widetilde{M}_n = o_{\mathbb{P}}(1),$$

where

$$\begin{aligned}\widehat{M}_n &:= n \cdot \text{dCov}_n^2 \left( \left( \mathbf{F}_{\mathbf{X}, \pm}^{(n)}(\mathbf{X}_i) \right)_{i=1}^n, \left( \mathbf{F}_{\mathbf{Y}, \pm}^{(n)}(\mathbf{Y}_i) \right)_{i=1}^n \right) \\ \widetilde{M}_n &:= n \cdot \text{dCov}_n^2 \left( \left( \mathbf{F}_{\mathbf{X}, \pm}(\mathbf{X}_i) \right)_{i=1}^n, \left( \mathbf{F}_{\mathbf{Y}, \pm}(\mathbf{Y}_i) \right)_{i=1}^n \right).\end{aligned}$$

As an immediate corollary,

$$\widehat{M}_n \xrightarrow{d} \sum_{k=1}^{\infty} \lambda_k (\xi_k^2 - 1).$$

# Theory

## ■ Consistency

The test based on  $\widehat{M}_n$  takes the form

$$T_\alpha := \mathbb{1}(\widehat{M}_n > q_{1-\alpha}),$$

where  $q_{1-\alpha}$  is the  $(1 - \alpha)$ -quantile of  $\sum_{k=1}^{\infty} \lambda_k (\xi_k^2 - 1)$ .

---

**Proposition (Uniform validity and consistency).** The test  $T_\alpha$  is uniformly valid in the sense that

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}_p^{\text{ac}} \otimes \mathcal{P}_q^{\text{ac}}} P(T_\alpha = 1) = \alpha.$$

Moreover, for fixed alternative such that  $\mathbf{X} \in \mathcal{P}_p^{\text{ac}}$  and  $\mathbf{Y} \in \mathcal{P}_q^{\text{ac}}$  are dependent, it holds that

$$\lim_{n \rightarrow \infty} P(T_\alpha = 1) = 1.$$

# Theory

- Local power

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} := \begin{pmatrix} \mathbf{I}_p & \delta \mathbf{M} \\ \delta \mathbf{M}' & \mathbf{I}_q \end{pmatrix} \begin{pmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{pmatrix} = \mathbf{A}_\delta \begin{pmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{pmatrix}$$

Sequence of local alternatives:

$$H_{1,n}(\delta_0) : \delta = \delta_n, \quad \text{where } \delta_n := n^{-1/2} \delta_0$$

---

**Proposition (SHDH 2022).** If some regularity assumption holds, we have that for any number  $\beta > 0$  satisfying  $\alpha + \beta < 1$  there exists a constant  $c_\beta > 0$  only depending on  $\beta$  such that as long as  $|\delta_0| < c_\beta$

$$\inf_{\bar{T}_\alpha \in \mathcal{T}_\alpha} \mathbb{P}\{\bar{T}_\alpha = 0 \mid H_{1,n}(\delta_0)\} \geq 1 - \alpha - \beta$$

for all sufficiently large  $n$ . Here the infimum is taken over all size- $\alpha$  tests.

# Theory

- Local power

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} := \begin{pmatrix} \mathbf{I}_p & \delta \mathbf{M} \\ \delta \mathbf{M}' & \mathbf{I}_q \end{pmatrix} \begin{pmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{pmatrix} = \mathbf{A}_\delta \begin{pmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{pmatrix}$$

Sequence of local alternatives:

This is the **detection boundary!**

$$H_{1,n}(\delta_0) : \delta = \delta_n, \quad \text{where } \delta_n := n^{-1/2} \delta_0$$

---

**Theorem (SHDH 2022).** If some regularity assumption holds, then for any number  $\beta > 0$ , there exists some sufficiently large constant  $C_\beta > 0$  only depending on  $\beta$  such that, as long as  $|\delta_0| > C_\beta$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}\{\mathsf{T}_\alpha = 1 \mid H_{1,n}(\delta_0)\} \geq 1 - \beta.$$

multivariate Hájek asymptotic representation and Le Cam's third lemma!

# Paradigm

Goals to reach:

- The test should be **distribution-free**, and directly implementable without the need of permutation.

Yes, by distribution-freeness of multivariate ranks and Hájek asymptotic representation!



- The test should be **consistent** in a certain sense.

Yes, by consistency of distance covariance and P-a.s. invertibility of center-outward distribution function!



- The test should be **optimal** under certain standard.

Yes, by multivariate Hájek asymptotic representation and Le Cam's third lemma!



# Outline

- Center-outward multivariate rank
- The proposed test (cont.)
- Discussion

# Test

- **Data:**  $\{(\mathbf{X}_i \in \mathbb{R}^p, \mathbf{Y}_i \in \mathbb{R}^q), i \in [n]\}$  i.i.d.
- **Aim:** testing if “ $H_0 : \mathbf{X} \perp\!\!\!\perp \mathbf{Y}$ ” is true.

Proposal:

- Calculate **center-outward ranks**  $\mathbf{F}_{\mathbf{X},\pm}^{(n)}(\mathbf{X}_1), \dots, \mathbf{F}_{\mathbf{X},\pm}^{(n)}(\mathbf{X}_n)$  and  $\mathbf{F}_{\mathbf{Y},\pm}^{(n)}(\mathbf{Y}_1), \dots, \mathbf{F}_{\mathbf{Y},\pm}^{(n)}(\mathbf{Y}_n)$ ;

- Compute  $\widehat{\mathbf{W}}_J^{(n)} := \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathbf{J}_{\mathbf{X}} \left( \mathbf{F}_{\mathbf{X};\pm}^{(n)}(\mathbf{X}_i) \right) \mathbf{J}_{\mathbf{Y}} \left( \mathbf{F}_{\mathbf{Y};\pm}^{(n)}(\mathbf{Y}_i) \right),$

where  $\mathbf{J}(\mathbf{u}) := J(\|\mathbf{u}\|) \frac{\mathbf{u}}{\|\mathbf{u}\|} \mathbf{1}_{[\|\mathbf{u}\| \neq 0]}$ ; e.g.  $J(u) = 1$  (sign);  
 $J(u) = u$  (Wilcoxon);

- Reject  $H_0$  if  $\|\widehat{\mathbf{W}}_J^{(n)}\|_{\mathbb{F}}^2$  is large enough.  $J(u) = (F_{\chi_d^2}^{-1}(u))^{1/2}$  (vdW).

# Pitman efficiency

- Pitman efficiency

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} := \begin{pmatrix} \mathbf{I}_p & \delta \mathbf{M} \\ \delta \mathbf{M}' & \mathbf{I}_q \end{pmatrix} \begin{pmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{pmatrix} = \mathbf{A}_\delta \begin{pmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{pmatrix}$$

Sequence of local alternatives:

$$H_{1,n}(\delta_0) : \delta = \delta_n, \quad \text{where } \delta_n := n^{-1/2} \delta_0$$

---

**Theorem (SDHH 2021).** If  $\mathbf{X}^*$  and  $\mathbf{Y}^*$  are **elliptically symmetric distributions** satisfying some regularity assumption. Then, the **Pitman asymptotic relative efficiency (ARE)** of the center-outward test based on the **van der Waerden score functions**  $\left(F_{\chi_d^2}^{-1}(\cdot)\right)^{1/2}$  with respect to Wilks' test is larger than or equal to 1.



# Simulation

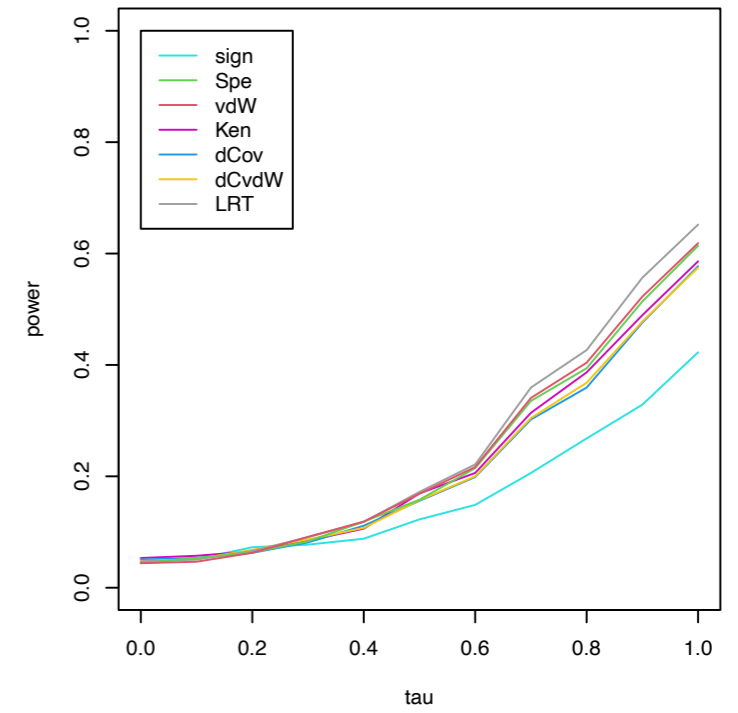
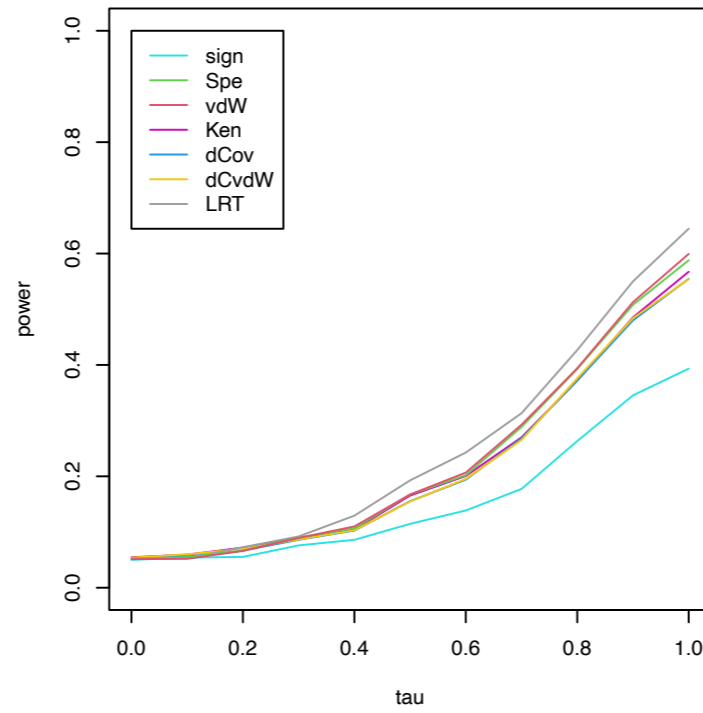
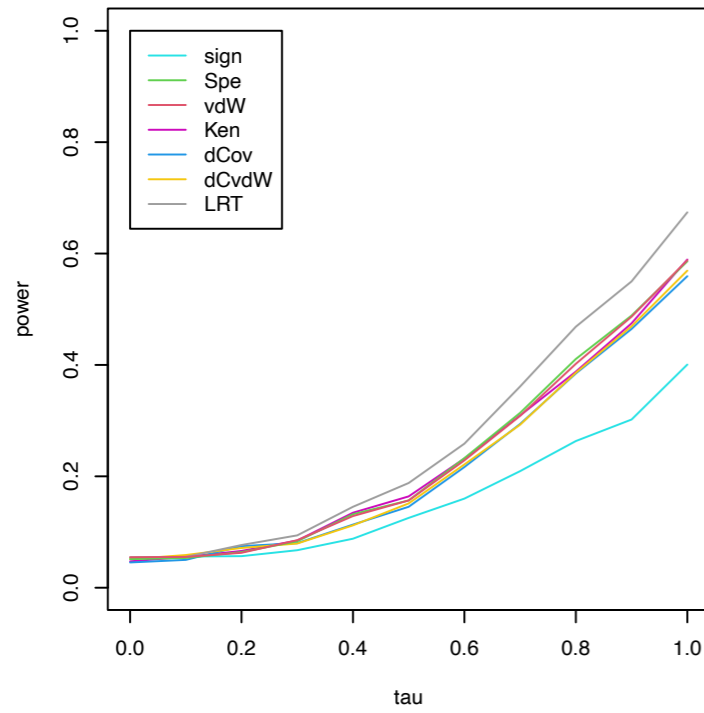
## ■ Gaussian

$n = 216$

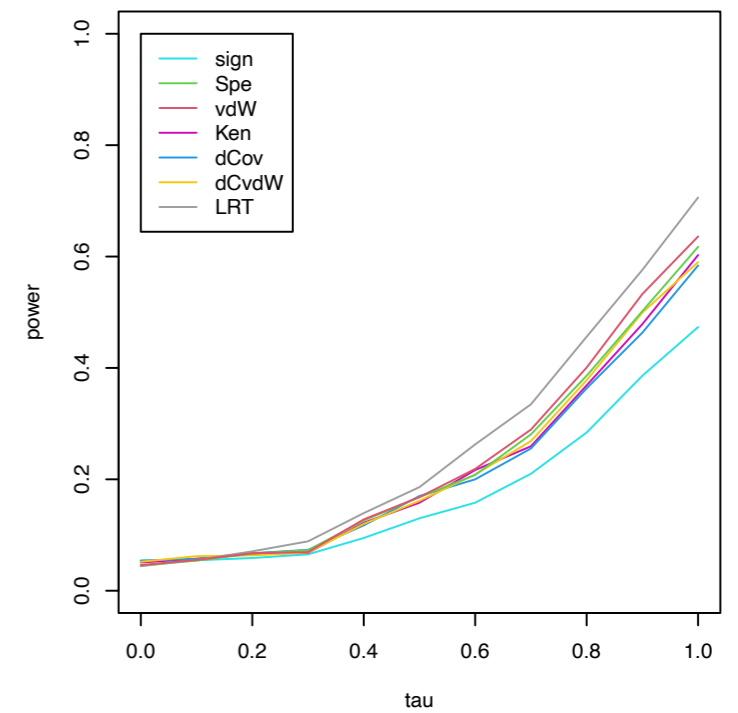
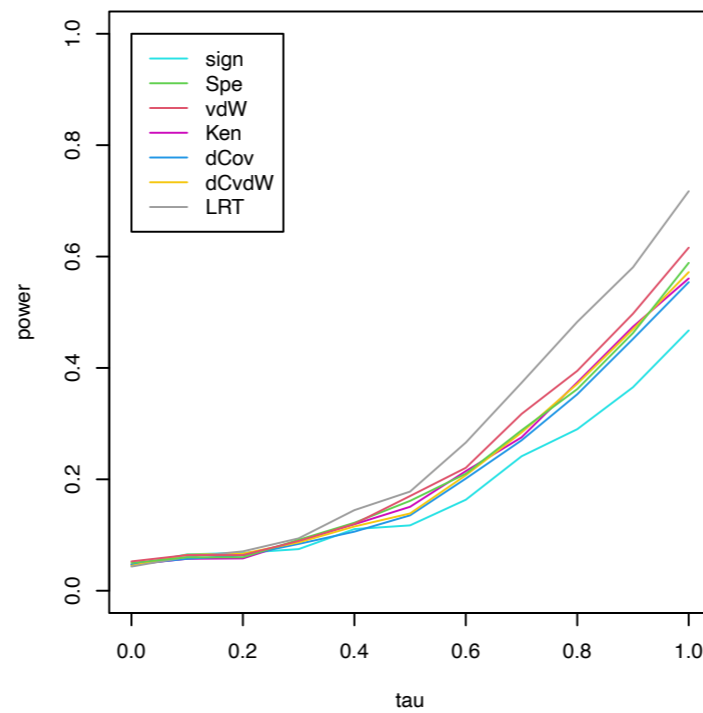
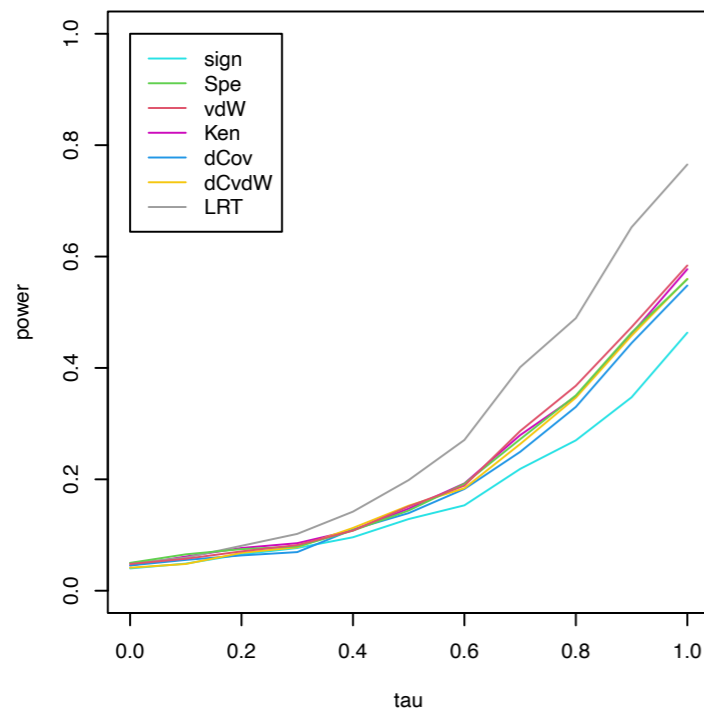
$n = 432$

$n = 864$

$(p, q) = (2, 2)$



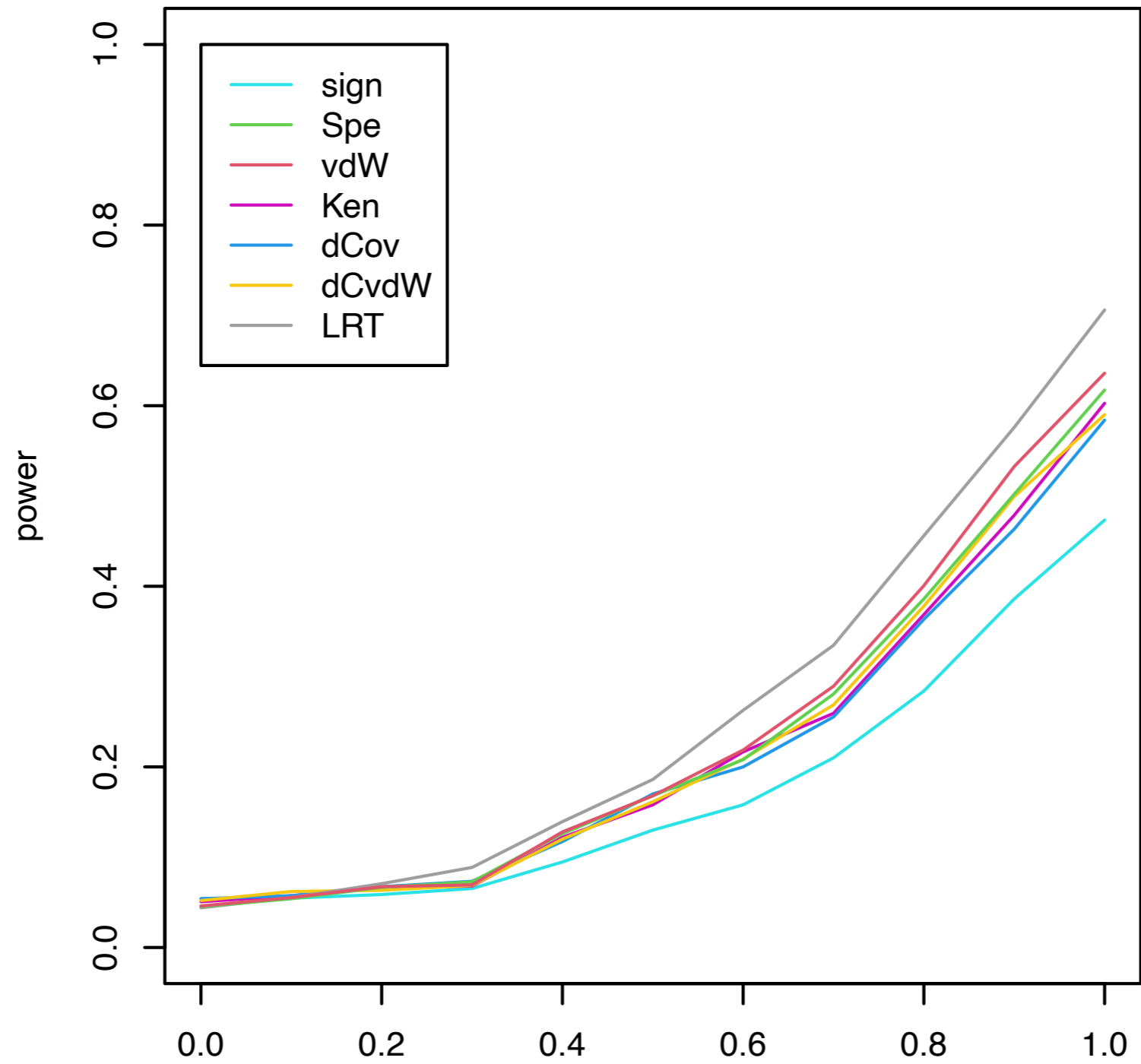
$(p, q) = (3, 3)$



# Simulation

## ■ Gaussian

$$n = 864, (p, q) = (3, 3)$$



# Simulation

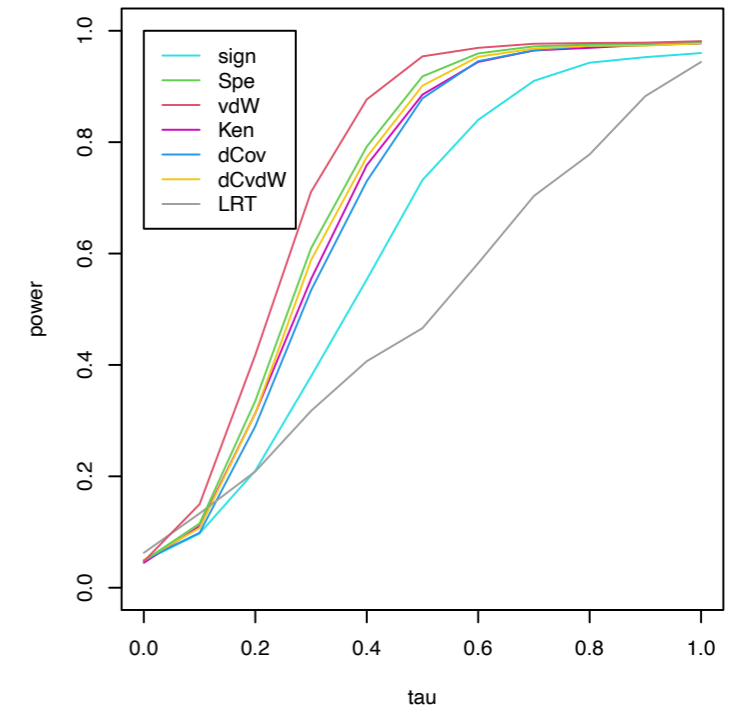
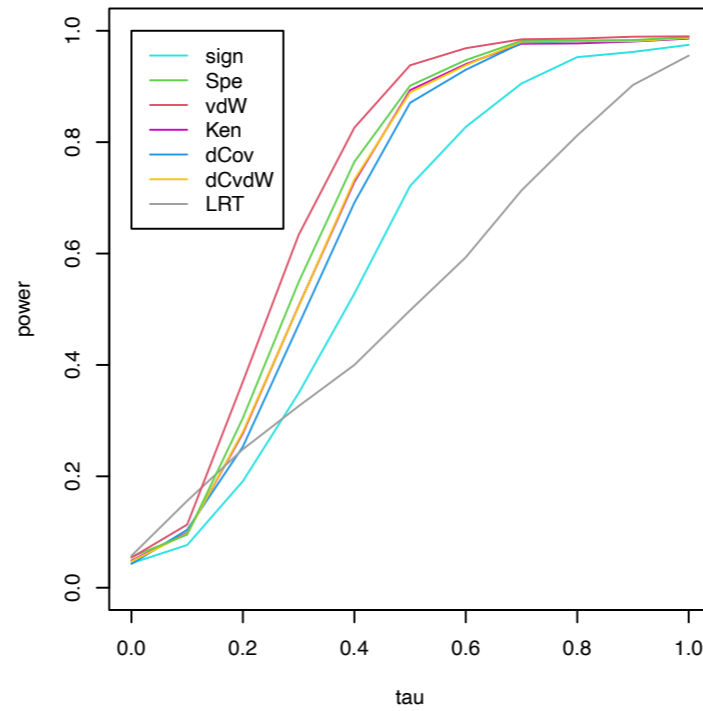
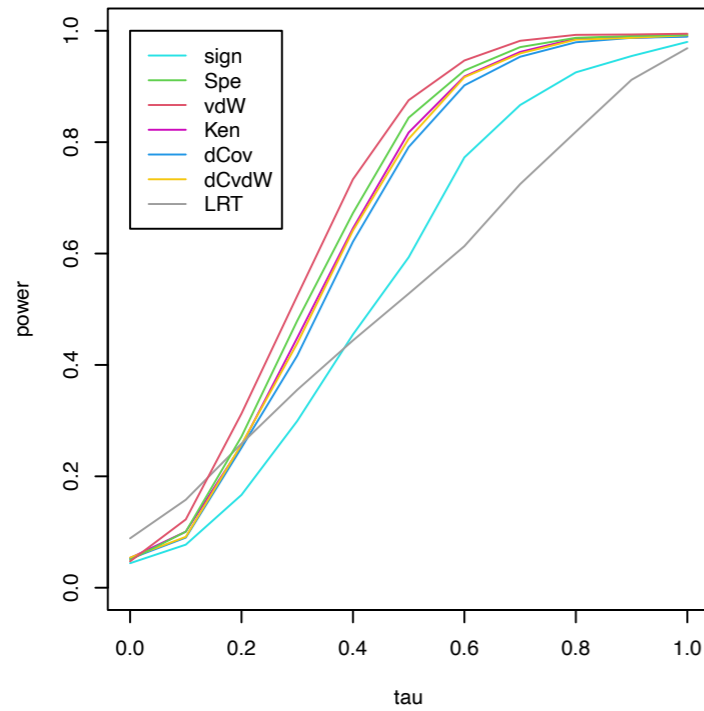
## ■ Cauchy

$n = 216$

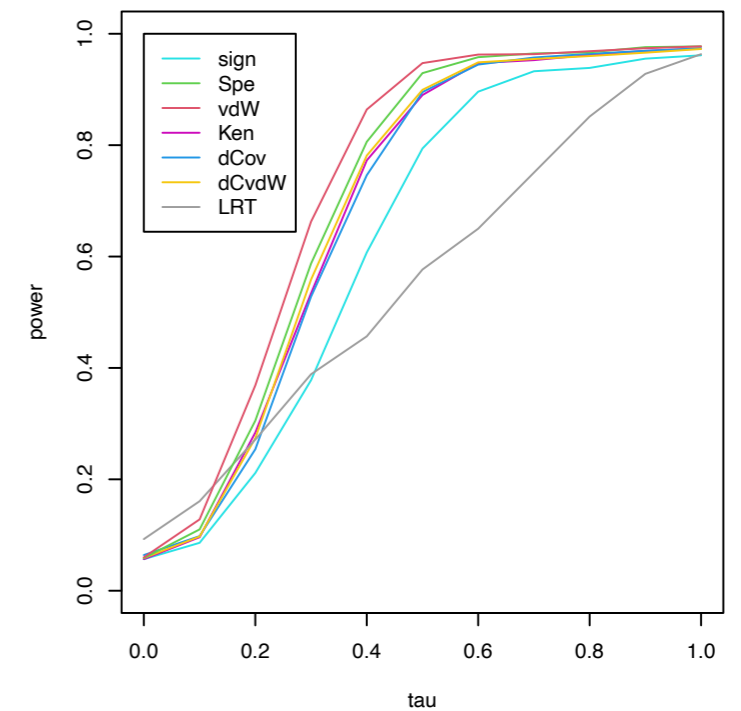
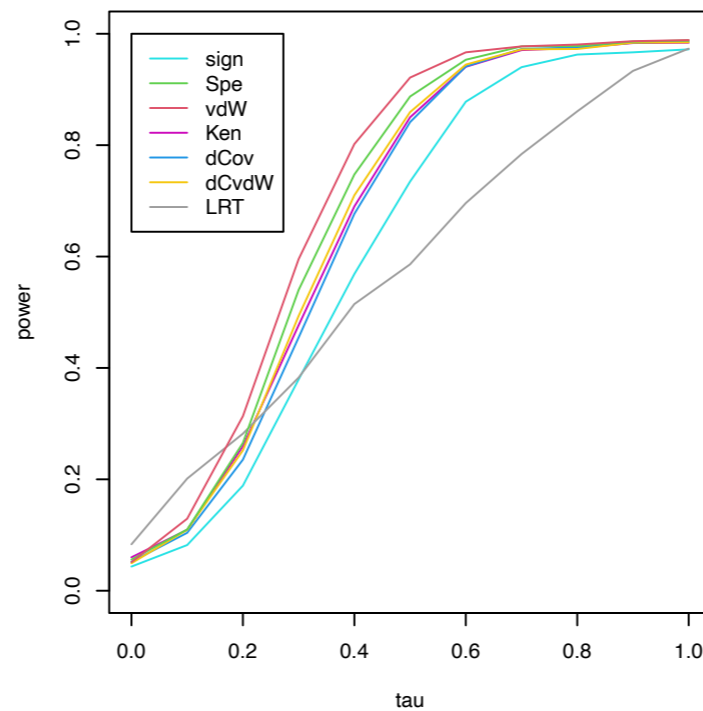
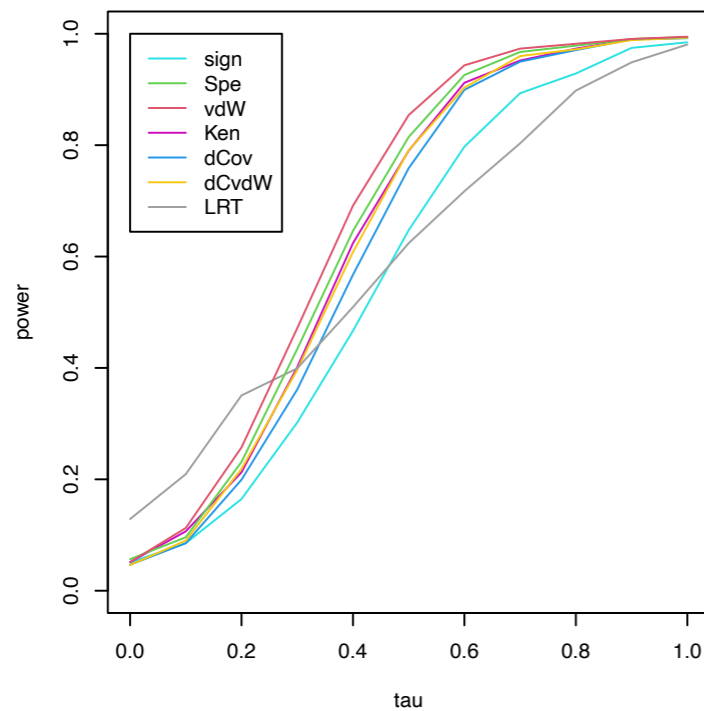
$n = 432$

$n = 864$

$(p, q) = (2, 2)$



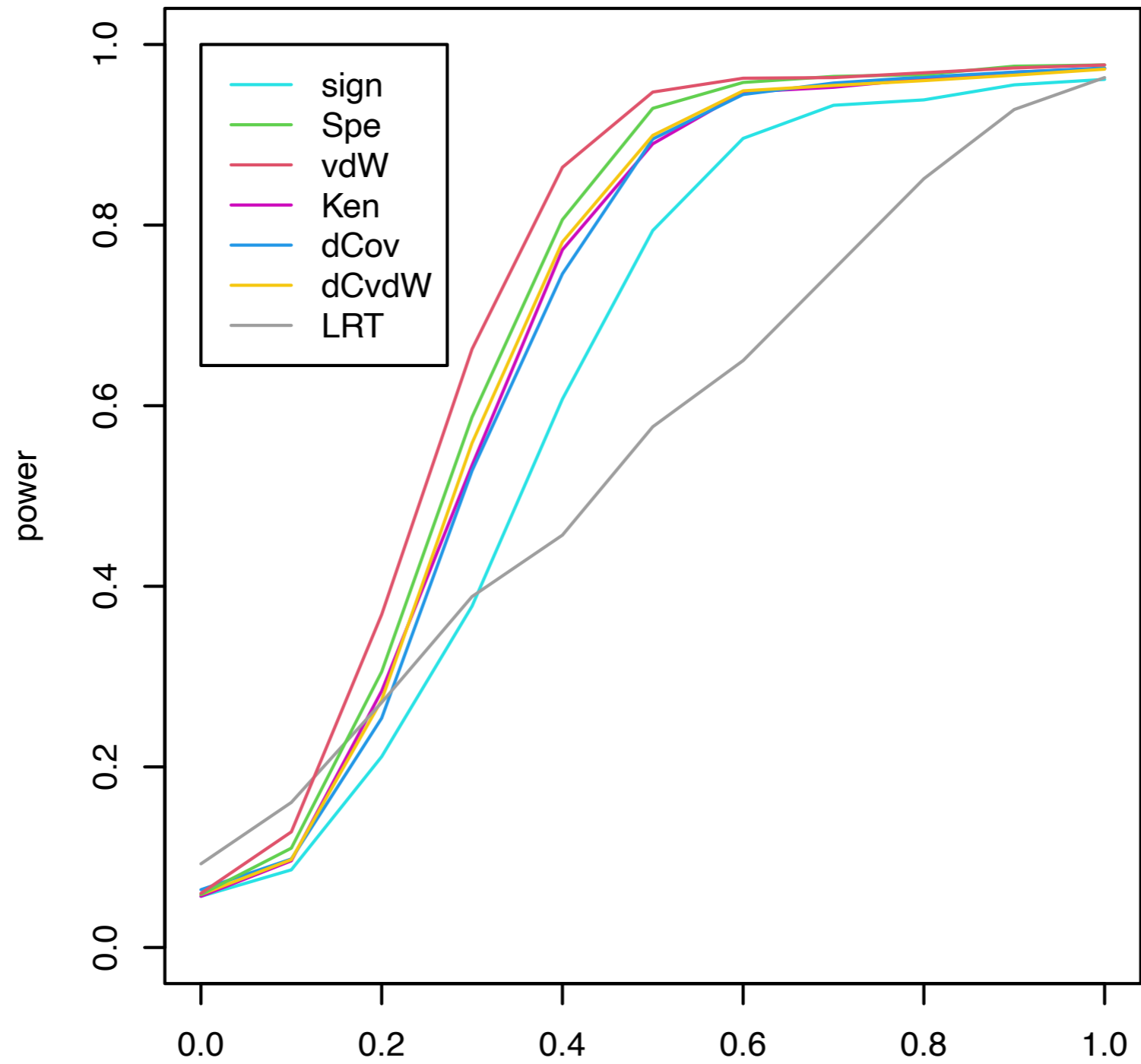
$(p, q) = (3, 3)$



# Simulation

## ■ Cauchy

$$n = 864, (p, q) = (3, 3)$$



# Outline

- Center-outward multivariate rank
- The proposed test
- Discussion

# Discussion

- Computational complexity of optimal transport-based multivariate ranks.

General dimension:  $\tilde{O}(n^{5/2})$  complexity,  
and  $\tilde{O}(n^{3/2})$  complexity if using **fast approximation**;  
if dimension is 2:  $\tilde{O}(n^{3/2+\delta})$  complexity,  
and  $\tilde{O}(n^{5/4})$  if using **fast approximation**.

- Future works: High-dimensional, Conditional independence testing...
- See our papers [SDH 2022 \(JASA, 117:395–410\)](#), [SHDH 2022 \(AoS, 50:1933–1959\)](#), and [SDHH 2021 \(arXiv:2111.15567v1\)](#) for more results.

# Papers

- Shi, H., Drton, M., and Han, F. (2022). Distribution-free consistent independence tests via center-outward ranks and signs. *J. Amer. Statist. Assoc.* 117(537):395–410.
- Shi, H., Hallin, M., Drton, M., and Han, F. (2022). On universally consistent and fully distribution-free rank tests of vector independence. *Ann. Statist.* 50(4):1933–1959.
- Shi, H., Drton, M., Hallin, M., and Han, F. (2021). Center-outward sign-and rank-based quadrant, Spearman, and Kendall tests for multivariate independence. Available at [arXiv:2111.15567v1](https://arxiv.org/abs/2111.15567v1).

**Thanks!**