

# Optimal transport in statistics and Pitman efficient multivariate distribution-free testing

Nabarun Deb  
Department of Statistics, Columbia University

Kantorovich Initiative Retreat  
March 18, 2022

# Multivariate distribution-free nonparametric testing

Consider the following **nonparametric hypothesis testing** problem:

Testing for equality of distributions (two-sample goodness-of-fit (GoF))

- **Data:**  $fX_i g_{i=1}^m$  iid  $P$  on  $\mathbb{R}^d$ ;  $fY_j g_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$
- Test if the **two-samples** came from the **same distribution**, i.e.,

$$H_0 : P = Q \quad \text{versus} \quad H_1 : P \neq Q$$

# Multivariate distribution-free nonparametric testing

Consider the following **nonparametric hypothesis testing** problem:

Testing for equality of distributions (two-sample goodness-of-fit (GoF))

- **Data:**  $fX_i g_{i=1}^m$  iid  $P$  on  $\mathbb{R}^d$ ;  $fY_j g_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$
- Test if the **two-samples** came from the **same distribution**, i.e.,

$$H_0 : P = Q \quad \text{versus} \quad H_1 : P \neq Q$$

- When  $d = 1$ : **Student's  $t$ -test (1908)**, **Wilcoxon rank-sum (1947)**, **Cramér-von Mises (1928)**, **Wald and Wolfowitz (1940)**, **Mann and Whitney (1947)**, **Kolmogorov-Smirnov (1939)**
- When  $d > 1$ : **Hotelling's  $T^2$ -statistic (1931)**, **Weiss (1960)**, **Anderson (1962)**, **Friedman and Raksy (1979)**, **Schilling (1986)**, **Rosenbaum (2005)**, **Gretton et al. (2012)**, **Szekely and Rizzo (2013)**, **Biswas et al. (2014)**, **Li and Yuan (2019)**

# Multivariate distribution-free nonparametric testing

Consider the following **nonparametric hypothesis testing** problem:

Testing for equality of distributions (two-sample goodness-of-fit (GoF))

- **Data:**  $fX_i g_{i=1}^m$  iid  $P$  on  $\mathbb{R}^d$ ;  $fY_j g_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$
- Test if the **two-samples** came from the **same distribution**, i.e.,

$$H_0 : P = Q \quad \text{versus} \quad H_1 : P \neq Q$$

- When  $d = 1$ : Student's  $t$ -test (1908), Wilcoxon rank-sum (1947), Cramér-von Mises (1928), Wald and Wolfowitz (1940), Mann and Whitney (1947), Kolmogorov-Smirnov (1939)
- When  $d > 1$ : Hotelling's  $T^2$ -statistic (1931), Weiss (1960), Anderson (1962), Friedman and Raksy (1979), Schilling (1986), Rosenbaum (2005), Gretton et al. (2012), Szekely and Rizzo (2013), Biswas et al. (2014), Li and Yuan (2019)

When  $d = 1$

- **Two-sample  $t$ -test:** Compares  $\bar{X}_m$  and  $\bar{Y}_n$

## When $d = 1$

- **Two-sample  $t$ -test:** Compares  $\bar{X}_m$  and  $\bar{Y}_n$
- Reject if the statistic is larger than the  $(1 - \alpha)$ -th quantile of  $t_{m+n-2}$  (or use a permutation test)
- **Approximate** (not valid for small sample sizes) level  $\alpha$  test, requires additional moment assumptions

# When $d = 1$

- **Two-sample  $t$ -test:** Compares  $\bar{X}_m$  and  $\bar{Y}_n$
  - Reject if the statistic is larger than the  $(1 - \alpha)$ -th quantile of  $t_{m+n-2}$  (or use a permutation test)
  - **Approximate** (not valid for small sample sizes) level  $\alpha$  test, requires additional moment assumptions
- 
- **Distribution-free tests:** Null distribution of the test statistic  $T_n$  is **universal**, i.e.,  $P(T_n \leq c_n) = \alpha$  where  $c_n$ , the deterministic rejection **threshold** can be obtained **before observing** the data

# When $d = 1$

- **Two-sample  $t$ -test:** Compares  $\bar{X}_m$  and  $\bar{Y}_n$
  - Reject if the statistic is larger than the  $(1 - \alpha)$ -th quantile of  $t_{m+n-2}$  (or use a permutation test)
  - **Approximate** (not valid for small sample sizes) level  $\alpha$  test, requires additional moment assumptions
- 
- **Distribution-free tests:** Null distribution of the test statistic  $T_n$  is **universal**, i.e.,  $P(T_n \leq c_n) = \alpha$  where  $c_n$ , the deterministic rejection **threshold** can be obtained **before observing** the data
  - They are **exact** tests and valid for **all sample sizes**



# When $d = 1$

- **Two-sample  $t$ -test:** Compares  $\bar{X}_m$  and  $\bar{Y}_n$
- Reject if the statistic is larger than the  $(1 - \alpha)$ -th quantile of  $t_{m+n-2}$  (or use a permutation test)
- **Approximate** (not valid for small sample sizes) level  $\alpha$  test, requires additional moment assumptions

- **Distribution-free tests:** Null distribution of the test statistic  $T_n$  is **universal**, i.e.,  $P(T_n \leq c_n) = \alpha$  where  $c_n$ , the deterministic rejection **threshold** can be obtained **before observing** the data
- They are **exact** tests and valid for **all sample sizes**
- Based on **univariate ranks** — advent of **classical nonparametrics**

## Comparison of Wilcoxon rank-sum (WRS) test with two-sample $t$ -test

Pool  $(X_1, \dots, X_m; Y_1, \dots, Y_n)$ : (scaled) ranks  $R_{m;n}(X_i)$ 's and  $R_{m;n}(Y_j)$ 's

$$\frac{1}{n} \sum_{j=1}^n R_{m;n}(Y_j) \quad \frac{1}{m} \sum_{i=1}^m R_{m;n}(X_i)$$

- WRS test is **distribution-free** and **exact** for all  $F$  continuous

## Comparison of Wilcoxon rank-sum (WRS) test with two-sample $t$ -test

Pool  $(X_1, \dots, X_m; Y_1, \dots, Y_n)$ : (scaled) ranks  $R_{m;n}(X_i)$ 's and  $R_{m;n}(Y_j)$ 's

$$\frac{1}{n} \sum_{j=1}^n R_{m;n}(Y_j) \quad \frac{1}{m} \sum_{i=1}^m R_{m;n}(X_i)$$

- WRS test is **distribution-free** and **exact** for all  $F$  **continuous**
- WRS test has **0.95 Pitman efficiency** w.r.t.  $t$ -test when  $F$  is **Gaussian**

## Comparison of Wilcoxon rank-sum (WRS) test with two-sample $t$ -test

Pool  $(X_1, \dots, X_m; Y_1, \dots, Y_n)$ : (scaled) ranks  $R_{m;n}(X_i)$ 's and  $R_{m;n}(Y_j)$ 's

$$\frac{1}{n} \sum_{j=1}^n R_{m;n}(Y_j) \quad \frac{1}{m} \sum_{i=1}^m R_{m;n}(X_i)$$

- WRS test is **distribution-free** and **exact** for all  $F$  continuous
- WRS test has **0.95 Pitman efficiency** w.r.t.  $t$ -test when  $F$  is Gaussian
- Non-trivial efficiency **lower bound** of **0.864** w.r.t.  $t$ -test [Hodges and Lehmann (1956)]; efficiency can be **+1** (for heavy-tailed dist.)

## Comparison of Wilcoxon rank-sum (WRS) test with two-sample $t$ -test

Pool  $(X_1, \dots, X_m; Y_1, \dots, Y_n)$ : (scaled) ranks  $R_{m;n}(X_i)$ 's and  $R_{m;n}(Y_j)$ 's

$$\frac{1}{n} \sum_{j=1}^n R_{m;n}(Y_j) \quad \frac{1}{m} \sum_{i=1}^m R_{m;n}(X_i)$$

- WRS test is **distribution-free** and **exact** for all  $F$  **continuous**
- WRS test has **0.95 Pitman efficiency** w.r.t.  $t$ -test when  $F$  is **Gaussian**
- Non-trivial efficiency **lower bound** of **0.864** w.r.t.  $t$ -test [Hodges and Lehmann (1956)]; efficiency can be **+7** (for heavy-tailed dist.)
- Non-trivial efficiency **lower bound** of **1** w.r.t.  $t$ -test [Chernoff and Savage (1958)] when the following revised statistic is used:

$$\frac{1}{n} \sum_{j=1}^n 1(R_{m;n}(Y_j)) \quad \frac{1}{m} \sum_{i=1}^m 1(R_{m;n}(X_i))$$

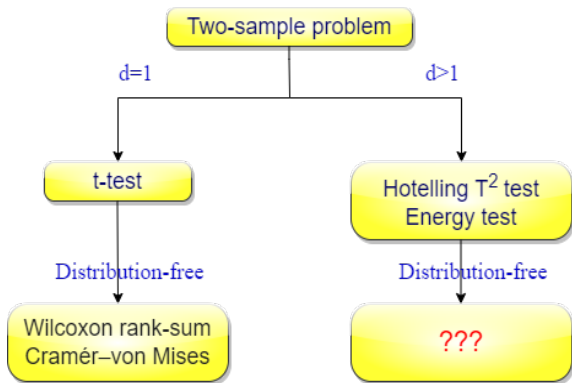
**Generalize** distribution-freeness, efficiency to **multivariate** data

## Question

Can we construct multivariate nonparametric distribution-free tests?

## Question

Can we construct multivariate nonparametric distribution-free tests?



Good news

Tests based on "*ranks*" are distribution-free



Bad news

Tests based on "*ranks*" are distribution-free

## Bad news

### Tests based on "ranks" are distribution-free

How do we define multivariate ranks which lead to distribution-free tests?

What about their statistical efficiency?

## Bad news

Tests based on "ranks" are distribution-free

How do we define multivariate ranks which lead to distribution-free tests?

What about their statistical efficiency?

Optimal transport!

- 1 A (very) brief introduction to optimal transport
- 2 Multivariate ranks using optimal transport
- 3 Multivariate distribution-free tests using optimal transport
  - Rank Hotelling  $T^2$  test and Pitman efficiency
  - Pitman efficiency, comparison with Hotelling  $T^2$

- 1 A (very) brief introduction to optimal transport
- 2 Multivariate ranks using optimal transport
- 3 Multivariate distribution-free tests using optimal transport
  - Rank Hotelling  $T^2$  test and Pitman efficiency
  - Pitman efficiency, comparison with Hotelling  $T^2$

# Optimal (measure) transportation

$$\text{KL}(P \parallel Q) = \sum_j p_j \log \frac{p_j}{q_j} \quad p_j = 1$$
$$\text{TV}(P; Q) = \frac{1}{2} \sum_j |p_j - q_j| \quad p_j, q_j = 1$$

# Optimal (measure) transportation

$$\begin{aligned} \text{KL}(\mathbf{P}; \mathbf{Q}) &= \sum_j^R \log \frac{p_j}{q_j} \quad p_j = 1 = \text{KL}(\mathbf{P}; \mathbf{R}) \\ \text{TV}(\mathbf{P}; \mathbf{Q}) &= \frac{1}{2} \sum_j^R |p_j - q_j| \quad q_j = 1 = \text{TV}(\mathbf{P}; \mathbf{R}) \end{aligned}$$

# Optimal (measure) transportation

$$KL(P \parallel Q) = \int \log \frac{p}{q} p = 1 = KL(P \parallel R)$$

$$TV(P; Q) = \frac{1}{2} \int |p - q| = 1 = TV(P; R)$$

Need a notion of distance that is sensitive to **geometry**

Monge's approach (1781): Given probability measures  $P, Q$  on  $\mathbb{R}^d$ , find an "optimal" map  $T_0 : \mathbb{R}^d \rightarrow \mathbb{R}^d$  satisfying

$$\min_{T \# P = Q} \int \|x - T(x)\|^2 dP(x); \quad T \# P = Q, \quad X \sim P; T(X) \sim Q$$



# Optimal (measure) transportation

$$KL(P \parallel Q) = \sum_j p_j \log \frac{p_j}{q_j} \quad p_j = 1 = KL(P \parallel R)$$

$$TV(P; Q) = \frac{1}{2} \sum_j |p_j - q_j| \quad \sum_j p_j = \sum_j q_j = 1$$

Need a notion of distance that is sensitive to **geometry**

Monge's approach (1781): Given probability measures  $P, Q$  on  $\mathbb{R}^d$ , find an "optimal" map  $T_0 : \mathbb{R}^d \rightarrow \mathbb{R}^d$  satisfying

$$W_2^2(P; Q) = \min_{T \# P = Q} \int \sum_x \|T(x) - x\|^2 dP(x); \quad T \# P = Q, \quad X \sim P; T(X) \sim Q$$

Call **optimizer**  $T_0^{P;Q}$   $T_0$  (if it exists) | **optimal transport (OT) map**

$W_2^2(P; Q)$  | squared **Wasserstein** distance

# Optimal (Measure) Transportation

$$W_2^2(\mathbf{P}; \mathbf{Q}) = kb \quad ak^2, \quad W_2^2(\mathbf{P}; \mathbf{R}) = kc \quad ak^2$$

# Optimal (Measure) Transportation

$$W_2^2(\mathbf{P}; \mathbf{Q}) = k\mathbf{b} \quad \mathbf{a}k^2, \quad W_2^2(\mathbf{P}; \mathbf{R}) = k\mathbf{c} \quad \mathbf{a}k^2$$

$$T_0^{\mathbf{P}; \mathbf{Q}}(x) = x + \mathbf{b} \quad \mathbf{a}, \quad T_0^{\mathbf{P}; \mathbf{R}}(x) = x + \mathbf{c} \quad \mathbf{a}$$



How to estimate the optimal transport map?

# Estimation | a plug-in approach

$$T_0 = \arg \min_{T \# P = Q} \int k(x - T(x))^2 dP(x); \quad T \# P = Q, \quad X \sim P; \quad T(X) \sim Q$$

$T_0$  is unique  $P$  a.s.) if  $P; Q$  are absolutely continuous ([McCann \(1995\)](#))

# Estimation | a plug-in approach

$$T_0 = \arg \min_{T \# P = Q} \int k(x, T(x))^2 dP(x); \quad T \# P = Q, \quad X \sim P; \quad T(X) \sim Q$$

$T_0$  is unique  $P$  a.s.) if  $P; Q$  are absolutely continuous ([McCann \(1995\)](#))

**Data:**  $X_1; X_2; \dots; X_m$  iid  $P$  (unknown, absolutely continuous) and  $Y_1; \dots; Y_n$  iid  $Q$  (unknown, absolutely continuous)

# Estimation | a plug-in approach

$$T_0 = \arg \min_{T \# P = Q} \int k(x) T(x)^2 dP(x); \quad T \# P = Q, \quad X \sim P; \quad T(X) \sim Q$$

$T_0$  is unique  $P$  a.s.) if  $P; Q$  are absolutely continuous (McCann (1995))

**Data:**  $X_1; X_2; \dots; X_m$  iid  $P$  (unknown, absolutely continuous) and  $Y_1; \dots; Y_n$  iid  $Q$  (unknown, absolutely continuous)

**Empirical distributions:**  $P_m := \frac{1}{m} \sum_{i=1}^m \delta_{X_i}; \quad Q_n := \frac{1}{n} \sum_{j=1}^n \delta_{Y_j}$



# Estimation | a plug-in approach

$$T_0 = \arg \min_{T \# P = Q} \int k(x) T(x) k^2 dP(x); \quad T \# P = Q, \quad X \sim P; \quad T(X) \sim Q$$

$T_0$  is unique  $P$  a.s.) if  $P; Q$  are absolutely continuous (McCann (1995))

**Data:**  $X_1; X_2; \dots; X_m$  iid  $P$  (unknown, absolutely continuous) and  $Y_1; \dots; Y_n$  iid  $Q$  (unknown, absolutely continuous)

**Empirical distributions:**  $P_m := \frac{1}{m} \sum_{i=1}^m \delta_{X_i}; \quad Q_n := \frac{1}{n} \sum_{j=1}^n \delta_{Y_j}$

When  $m = n$

$$\hat{T} := \arg \min_{T \# P_n = Q_n} \int k(x) T(x) k^2 dP_n(x)$$

# Estimation | a plug-in approach

$$T_0 = \arg \min_{T \# P = Q} \int_{\mathcal{X}} k(x) |T(x) - Q(x)|^2 dP(x); \quad T \# P = Q, \quad \mathcal{X} \subseteq \mathcal{P}; \quad T(\mathcal{X}) \subseteq \mathcal{Q}$$

$T_0$  is unique  $P$  a.s.) if  $P; Q$  are absolutely continuous (McCann (1995))

**Data:**  $X_1; X_2; \dots; X_m$  iid  $P$  (unknown, absolutely continuous) and  $Y_1; \dots; Y_n$  iid  $Q$  (unknown, absolutely continuous)

**Empirical distributions:**  $P_m := \frac{1}{m} \sum_{i=1}^m \delta_{X_i}; \quad Q_n := \frac{1}{n} \sum_{j=1}^n \delta_{Y_j}$

When  $m = n$

$$\hat{T} := \arg \min_{T \# P_n = Q_n} \int_{\mathcal{X}} k(x) |T(x) - Q_n(x)|^2 dP_n(x) = \arg \min_{T \# P_n = Q_n} \frac{1}{n} \sum_{i=1}^n k(X_i) |T(X_i) - Q_n(X_i)|^2$$

Recall  $T \# P_n = Q_n$  means if  $X \sim P_n$ , then  $T(X) \sim Q_n$

# Estimation | a plug-in approach

$$T_0 = \arg \min_{T \# P = Q} \int k(x) T(x) k^2 dP(x); \quad T \# P = Q, \quad X \sim P; \quad T(X) \sim Q$$

$T_0$  is unique  $P$  a.s.) if  $P; Q$  are absolutely continuous (McCann (1995))

**Data:**  $X_1; X_2; \dots; X_m$  iid  $P$  (unknown, absolutely continuous) and  $Y_1; \dots; Y_n$  iid  $Q$  (unknown, absolutely continuous)

**Empirical distributions:**  $P_m := \frac{1}{m} \sum_{i=1}^m \delta_{X_i}; \quad Q_n := \frac{1}{n} \sum_{j=1}^n \delta_{Y_j}$

When  $m = n$

$$\hat{T} := \arg \min_{T \# P_n = Q_n} \int k(x) T(x) k^2 dP_n(x) = \arg \min_{T \# P_n = Q_n} \frac{1}{n} \sum_{i=1}^n k(X_i) T(X_i) k^2$$

Recall  $T \# P_n = Q_n$  means if  $X \sim P_n$ , then  $T(X) \sim Q_n$

$T \# P_n = Q_n$ :  $(T(X_1); \dots; T(X_n))$  is some **permutation** of  $(Y_1; \dots; Y_n)$

# Estimation | a plug-in approach

$$T_0 = \arg \min_{T \# P = Q} \int k(x) T(x) k^2 dP(x); \quad T \# P = Q, \quad X \sim P; \quad T(X) \sim Q$$

$T_0$  is unique (P a.s.) if  $P; Q$  are absolutely continuous (McCann (1995))

**Data:**  $X_1; X_2; \dots; X_m$  iid  $P$  (unknown, absolutely continuous) and  $Y_1; \dots; Y_n$  iid  $Q$  (unknown, absolutely continuous)

**Empirical distributions:**  $P_m := \frac{1}{m} \sum_{i=1}^m \delta_{X_i}; \quad Q_n := \frac{1}{n} \sum_{j=1}^n \delta_{Y_j}$

When  $m = n$

$$\hat{T} := \arg \min_{T \# P_n = Q_n} \int k(x) T(x) k^2 dP_n(x) = \arg \min_{T \# P_n = Q_n} \frac{1}{n} \sum_{i=1}^n k(X_i) T(X_i) k^2$$

Recall  $T \# P_n = Q_n$  means if  $X \sim P_n$ , then  $T(X) \sim Q_n$

**Assignment** problem (linear program { exact algorithm with complexity  $O(n^3)$ ; parallel computing { Date and Nagi (2016)})

# Estimation | a plug-in approach (Continued)

What happens when  $m < n$ ?

Can we still define

$$\hat{T} := \arg \min_{T \in \mathcal{P}_m = \mathcal{Q}_n} \int_{\mathcal{Z}} k(x, T(x))^2 dP_m(x)??$$

# Estimation | a plug-in approach (Continued)

What happens when  $m < n$ ?

Can we still define

$$\hat{T} := \arg \min_{T \# P_m = Q_n} \int k(x) |T(x) - Q_n(x)|^2 dP_m(x)??$$

**NOT FEASIBLE!**

There is **no function**  $T$  such that  $T \# P_m = Q_n$

(Kantorovic **relaxation**)

Let  $(P; Q)$  be the set of probability measures  
(**coupling**) on  $\mathbb{R}^d \times \mathbb{R}^d$ ,  
with marginals  $P; Q$ .

Then

$$W_2^2(P; Q) = \inf_{Z \in \Pi(P; Q)} \int k(x, y)^2 dZ(x, y)$$

**Examples:**  $P, Q$ ,  
 $(x; y) \mapsto \mathbb{1}(y = T_0(x))$

(Kantorovic **relaxation**)

Let  $(P; Q)$  be the set of probability measures (**coupling**) on  $\mathbb{R}^d \times \mathbb{R}^d$ , with marginals  $P; Q$ .

Then

$$W_2^2(P; Q) = \inf_{\gamma \in \Pi(P; Q)} \int_{\mathbb{R}^d \times \mathbb{R}^d} \|x - y\|^2 d\gamma(x; y)$$

**Examples:**  $P, Q$ ,  
 $d(x; y) = |y - T_0(x)|$

Always **has a minimizer**  
which matches  $T_0$  if  $P$  is absolutely continuous



Solve

$$\min_{x \in P_m, y \in Q_n} \|x - y\|_2^2$$

via a **linear program**

Solve

$$z \in \arg \min_{z \in \mathbb{R}^d} \int_{\mathcal{P}_m; \mathcal{Q}_n} \|x - y\|^2 d(x; y)$$

via a **linear program**

D., Ghosal, and Sen (NeurIPS, 2021): Define our estimator (**barycentric projection**) as

$$\hat{T}(x) = E_b[Y | X = x] = \int_{\mathcal{R}} y \frac{db(x; y)}{\int_y db(x; y)}$$

Both definitions coincide when  $m = n$

What is the rate of convergence of  $\bar{d}_t$  to  $T_0$ ?

What is the rate of convergence  $\hat{T}$  to  $T_0$ ?

When  $m = n \dots$

Empirical OT map:

$$\hat{T} := \arg \min_{T \# P_n = Q_n} \int_Z k(x, T(x)) k^2 dP_n(x)$$

Population OT map:

$$T_0 := \arg \min_{T \# P = Q} \int_Z k(x, T(x)) k^2 dP(x)$$

What is the rate of convergence  $\hat{T}$  to  $T_0$ ?

When  $m = n \dots$

Empirical OT map:

$$\hat{T} := \arg \min_{T \# P_n = Q_n} \int_Z k(x, T(x)) k^2 dP_n(x)$$

Population OT map:

$$T_0 := \arg \min_{T \# P = Q} \int_Z k(x, T(x)) k^2 dP(x)$$

Different parameter spaces

What is the rate of convergence of  $\hat{T}_n$  to  $T_0$ ?

Rate of convergence (D., Ghosal, and Sen (NeurIPS, 2021))

Assume that  $T_0$  is Lipschitz, and both  $P$  and  $Q$  are compactly supported (can be relaxed). Then, for  $d \geq 4$ ,

$$\frac{1}{m} \sum_{i=1}^m \mathbb{E} \| \hat{T}_n(X_i) - T_0(X_i) \|^2 = O\left(m^{-\frac{2}{d}} + n^{-\frac{2}{d}}\right)$$

What is the rate of convergence  $\hat{T}_m$  to  $T_0$ ?

Rate of convergence (D., Ghosal, and Sen (NeurIPS, 2021))

Assume that  $T_0$  is Lipschitz, and both  $P$  and  $Q$  are compactly supported (can be relaxed). Then, for  $d \geq 4$ ,

$$\frac{1}{m} \sum_{i=1}^m \mathbb{E} \| \hat{T}(X_i) - T_0(X_i) \|^2 \leq m^{-\frac{2}{d}} + n^{-\frac{2}{d}}.$$

The proof requires **convex analysis**, **chaining** and **Talagrand's concentration** arguments

Minimax optimal for  $d \geq 4$  (Hutter and Rigollet (2019))

What is the rate of convergence of  $\hat{T}$  to  $T_0$ ?

Rate of convergence (D., Ghosal, and Sen (NeurIPS, 2021))

Assume that  $T_0$  is Lipschitz, and both  $P$  and  $Q$  are compactly supported (**can be relaxed**). Then, for  $d \geq 4$ ,

$$\frac{1}{m} \sum_{i=1}^m \mathbb{E} \| \hat{T}(X_i) - T_0(X_i) \|^2 \leq m^{-\frac{2}{d}} + n^{-\frac{2}{d}}:$$

For  $d = 1; 2; 3$ ,  $m = n$ , rates are  $n^{-4/5}$ ,  $n^{-2/3}$ ,  $n^{-4/7}$  (ongoing work)

The proof requires **convex analysis**, **chaining** and **Talagrand's concentration** arguments

**Minimax optimal** for  $d \geq 4$  (**Hutter and Rigollet (2019)**)

These are the **fast** rates for a **practically computable** estimator of the OT map  $T_0$  (note  $\hat{T}$  requires **no tuning**) ▶ skip



## Question

Can we construct multivariate distribution-free tests?

- 1 A (very) brief introduction to optimal transport
- 2 **Multivariate ranks using optimal transport**
- 3 Multivariate distribution-free tests using optimal transport
  - Rank Hotelling  $T^2$  test and Pitman efficiency
  - Pitman efficiency, comparison with Hotelling  $T^2$

# Ranks: Where $d = 1$

Rank map  $\mathbb{P}^n$  assigns  $X_1; X_2; \dots; X_n$  to elements of  $\frac{1}{n}; \frac{2}{n}; \dots; \frac{n}{n}g$

$$\text{Define } \mu_n := \frac{1}{n} \sum_{i=1}^n x_i \text{ and } \nu_n := \frac{1}{n} \sum_{j=1}^n \frac{j}{n}$$

# Ranks: Where $d = 1$

Rank map  $\mathbb{P}^n$  assigns  $X_1; X_2; \dots; X_n$  to elements of  $\frac{1}{n}; \frac{2}{n}; \dots; \frac{n}{n}g$

$$\text{Define } \mathbb{P}^n := \frac{1}{n} \sum_{i=1}^n X_i \text{ and } \mathbb{P}^n := \frac{1}{n} \sum_{j=1}^n \frac{j}{n}$$

# Ranks: Where $d = 1$

Rank map  $\mathbb{R}_n$  assigns  $X_1; X_2; \dots; X_n$  to elements of  $\frac{1}{n}; \frac{2}{n}; \dots; \frac{n}{n}$

$$\text{Define } \bar{x}_n := \frac{1}{n} \sum_{i=1}^n X_i \text{ and } \bar{y}_n := \frac{1}{n} \sum_{j=1}^n Y_j$$

$$\mathbb{R}_n := \arg \max_{T: T \# \bar{x}_n = \bar{y}_n} \frac{1}{n} \sum_{i=1}^n X_i \cdot T(X_i) = \arg \min_{T: T \# \bar{x}_n = \bar{y}_n} \frac{1}{n} \sum_{i=1}^n |X_i - T(X_i)|^2$$

# Ranks: When $d = 1$

Rank map  $\mathbb{R}_n$  assigns  $(X_1; X_2; \dots; X_n)$  to elements of  $\{ \frac{1}{n}; \frac{2}{n}; \dots; \frac{n}{n} \}$

Define  $\mu_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$  and  $\nu_n := \frac{1}{n} \sum_{j=1}^n \delta_{\frac{j}{n}}$

$$\mathbb{R}_n := \arg \max_{T: T\# \mu_n = \nu_n} \frac{1}{n} \sum_{i=1}^n X_i T(X_i) = \arg \min_{T: T\# \mu_n = \nu_n} \frac{1}{n} \sum_{i=1}^n j X_i T(X_i)^2$$

$\mathbb{R}_n$  is the empirical OT map from  $\mu_n$  to  $\nu_n$

# Multivariate ranks $\phi$ (1)

Empirical rank map assigns  $X_1, \dots, X_n$  to  $c_1, \dots, c_n$  in  $\mathbb{R}^d$  |  
grid of "reference" points (e.g., a random sample from  $\text{Unif}([0, 1]^d)$ ,  
 $N(0; I_d)$  distribution, deterministic quasi-Monte Carlo sequences)

# Multivariate ranks (1)

**Empirical rank map** assigns  $(X_1, \dots, X_n) \in \mathbb{R}^d$  to a grid of "reference" points (e.g., a random sample from  $\text{Unif}([0, 1]^d)$ ,  $N(0, I_d)$  distribution, deterministic quasi-Monte Carlo sequences)

**Sample rank map** (Hallin (2017)) is defined as the **empirical OT map** :

$$\mathbb{R}_n := \arg \min_{T: T\# \mu_n = \nu_n} \frac{1}{n} \sum_{i=1}^n \|X_i - T(X_i)\|^2$$

where  $T$  transports  $\mu_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$  to  $\nu_n := \frac{1}{n} \sum_{j=1}^n \delta_{c_j}$



# Multivariate ranks (1)

**Empirical rank map** assigns  $(X_1, \dots, X_n) \in \mathbb{R}^d$  to a grid of "reference" points (e.g., a random sample from  $\text{Unif}([0, 1]^d)$ ,  $N(0; I_d)$  distribution, deterministic quasi-Monte Carlo sequences)

**Sample rank map** (Hallin (2017)) is defined as the **empirical OT map** :

$$\hat{R}_n := \arg \min_{T: T \# \mu_n = \mu} \frac{1}{n} \sum_{i=1}^n \|X_i - T(X_i)\|^2$$

where  $T$  transports  $\mu_n := \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$  to  $\mu := \frac{1}{n} \sum_{j=1}^n \delta_{c_j}$

# Multivariate rank function as an OT map

$X$  ; is a probability measure on  $\mathbb{R}^d$  (abs. cont.)

Reference dist.  $U$  on  $S \subset \mathbb{R}^d$  ( =  $\text{Unif}([0; 1]^d)$ ,  $N(0; I_d)$ )

Find OT map  $T$  s.t.  $T(X) \stackrel{d}{=} U$  ( abs. cont.)

# Multivariate rank function as an OT map

$X$  ; is a probability measure on  $\mathbb{R}^d$  (abs. cont.)

Reference dist.  $U$  on  $S \subseteq \mathbb{R}^d$  ( $= \text{Unif}([0; 1]^d), N(0; I_d)$ )

Find OT map  $T$  s.t.  $T(X) \stackrel{d}{=} U$  (abs. cont.)

Population rank function (a.k.a OT map) [Chernozhukov et al. (2017)]

If  $E \|X\|^2 < \infty$ , rank fn.  $R : \mathbb{R}^d \rightarrow S$  is the population transport map

$$R := \arg \min_{T: \mathbb{R}^d \rightarrow S} E \|X - T(X)\|^2$$

# Multivariate rank function as an OT map

$X$  ; is a probability measure on  $\mathbb{R}^d$  (abs. cont.)

Reference dist.  $U$  on  $S \subset \mathbb{R}^d$  ( $= \text{Unif}([0; 1]^d), N(0; I_d)$ )

Find OT map  $T$  s.t.  $T(X) \stackrel{d}{=} U$  (abs. cont.)

Population rank function (a.k.a OT map) [Chernozhukov et al. (2017)]

If  $E \|X\|^2 < \infty$ , rank fn.  $R : \mathbb{R}^d \rightarrow S$  is the population transport map

$$R := \arg \min_{T: T\#X = U} E \|X - T(X)\|^2$$

Properties of population rank function [Brenier (1991), McCann (1995)]

$R(\cdot)$  characterizes distribution  $R_1(x) = R_2(x) \iff \int \phi(x) dP_1 = \int \phi(x) dP_2$

# Multivariate rank function as an OT map

$X$  ; is a probability measure on  $\mathbb{R}^d$  (abs. cont.)

Reference dist.  $U$  on  $S \subset \mathbb{R}^d$  ( $= \text{Unif}([0; 1]^d), N(0; I_d)$ )

Find OT map  $T$  s.t.  $T(X) \stackrel{d}{=} U$  (abs. cont.)

Population rank function (a.k.a OT map) [Chernozhukov et al. (2017)]

If  $E \|X\|^2 < \infty$ , rank fn.  $R : \mathbb{R}^d \rightarrow S$  is the population transport map

$$R := \arg \min_{T: T\#X = U} E \|X - T(X)\|^2$$

Properties of population rank function [Brenier (1991), McCann (1995)]

$R(\cdot)$  characterizes distribution  $R_1(x) = R_2(x)$   $\Leftrightarrow \int \phi(x) dP_1 = \int \phi(x) dP_2$

$R(\cdot)$  is the gradient of a convex function and smoothly invertible

# Multivariate rank function as an OT map

$X$  ; is a probability measure on  $\mathbb{R}^d$  (abs. cont.)

Reference dist.  $U$  on  $S \subset \mathbb{R}^d$  ( $= \text{Unif}([0; 1]^d), N(0; I_d)$ )

Find OT map  $T$  s.t.  $T(X) \stackrel{d}{=} U$  (abs. cont.)

Population rank function (a.k.a OT map) [Chernozhukov et al. (2017)]

If  $E \|X\|^2 < \infty$ , rank fn.  $R : \mathbb{R}^d \rightarrow S$  is the population transport map

$$R := \arg \min_{T: T\#X = U} E \|X - T(X)\|^2$$

Properties of population rank function [Brenier (1991), McCann (1995)]

$R(\cdot)$  characterizes distribution  $R_1(x) = R_2(x) \iff \int \mathbb{1}_A(x) dP_1 = \int \mathbb{1}_A(x) dP_2$

$R(\cdot)$  is the gradient of a convex function and smoothly invertible

When  $d = 1$ ,  $R(\cdot)$  is the CDF of  $X$ , when  $U = \text{Unif}([0; 1])$

**Distribution-freeness:** If  $F$  is absolutely continuous, then

$$(F_n(X_1); \dots; F_n(X_n))$$

is uniformly distributed over the  $n!$  permutations of  $c_1; \dots; c_n$

**Distribution-freeness:** If  $\mu$  is absolutely continuous, then

$$(\mu_n(X_1); \dots; \mu_n(X_n))$$

is uniformly distributed over the  $n!$  permutations of  $c_1; \dots; c_n$

**Consistency:** If  $\mu_n := \frac{1}{n} \sum_{j=1}^n \mu_j$  (abs. cont.), then

$$\frac{1}{n} \sum_{i=1}^n \|\mu_n(X_i) - R(X_i)\|_2^2 \rightarrow 0 \quad \text{as } n \rightarrow \infty;$$

where  $R$  is the unique OT map from  $\mu$  to  $\mu$ .

**No moment assumptions** needed on the model



- 1 A (very) brief introduction to optimal transport
- 2 Multivariate ranks using optimal transport
- 3 **Multivariate distribution-free tests using optimal transport**
  - Rank Hotelling  $T^2$  test and Pitman efficiency
  - Pitman efficiency, comparison with Hotelling  $T^2$

- 1 A (very) brief introduction to optimal transport
- 2 Multivariate ranks using optimal transport
- 3 **Multivariate distribution-free tests using optimal transport**  
Rank Hotelling  $T^2$  test and Pitman efficiency  
Pitman efficiency, comparison with Hotelling  $T^2$

## Testing for equality of two multivariate distributions

Data:  $\{X_i\}_{i=1}^m$  iid  $P$  on  $\mathbb{R}^d$ ;  $\{Y_j\}_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$

Test if the **two samples** came from the **same distribution**, i.e.,

$$H_0 : P = Q \quad \text{versus} \quad H_1 : P \neq Q$$

## Testing for equality of two multivariate distributions

Data:  $\{X_i\}_{i=1}^m$  iid  $P$  on  $\mathbb{R}^d$ ;  $\{Y_j\}_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$

Test if the **two samples** came from the **same distribution**, i.e.,

$$H_0 : P = Q \quad \text{versus} \quad H_1 : P \neq Q$$

Let  $N = m + n$  and assume  $\frac{m}{N} \rightarrow \lambda \in (0, 1)$

## Testing for equality of two multivariate distributions

Data:  $f X_i g_{i=1}^m$  iid  $P$  on  $\mathbb{R}^d$ ;  $f Y_j g_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$

Test if the **two samples** came from the **same distribution**, i.e.,

$$H_0 : P = Q \quad \text{versus} \quad H_1 : P \neq Q$$

Let  $N = m + n$  and assume  $\frac{m}{N} \rightarrow \lambda \in (0, 1)$

Hotelling  $T^2$  statistic [**Hotelling (1931)**]: The **multivariate analogue** of Student's **t-statistic**, given by

$$T_{m;n}^2 := \frac{mn}{m+n} (X - Y)^T S_{m;n}^{-1} (X - Y);$$

where  $S_{m;n}$  is **pooled covariance matrix**

Reject  $H_0$  if  $T_{m;n}^2 > c$  [**asympt. cut-off**  $c : (1 - \alpha)$  quantile of  $\chi_d^2$ ]

## Testing for equality of two multivariate distributions

Data:  $\{X_i\}_{i=1}^m$  iid  $P$  on  $\mathbb{R}^d$ ;  $\{Y_j\}_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$

Test if the **two samples** came from the **same distribution**, i.e.,

$$H_0 : P = Q \quad \text{versus} \quad H_1 : P \neq Q$$

Let  $N = m + n$  and assume  $\frac{m}{N} \rightarrow \lambda \in (0, 1)$

Hotelling  $T^2$  statistic [**Hotelling (1931)**]: The **multivariate analogue** of Student's  $t$ -statistic, given by

$$T_{m;n}^2 := \frac{mn}{m+n} \bar{X} - \bar{Y} > S_{m;n}^{-1} (\bar{X} - \bar{Y})$$

where  $S_{m;n}$  is **pooled covariance matrix**

Reject  $H_0$  if  $T_{m;n}^2 > c$  [**asympt. cut-off**  $c : (1 - \alpha)$  quantile of  $\chi^2_d$ ]

**Consistency:**  $P(T_{m;n}^2 > c) \rightarrow 1$  when  $E[X_1] \neq E[Y_1]$

Data:  $\{X_i\}_{i=1}^m$  iid  $P$  (abs. cont.),  $\{Y_j\}_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$

Reference dist.:  $\mathcal{S}$  on  $\mathbb{R}^d$  (abs. cont.; e.g.,  $= \text{Unif}([0; 1]^d)$ )

Proposed tests [D. & Sen (JASA, 2020); D., Bhattacharya & Sen (2020)]

Joint rank map: The sample ranks of the pooled observations:

$$\hat{R}_{m;n} : \{X_1, \dots, X_m, Y_1, \dots, Y_n\} \rightarrow \{c_1, \dots, c_{m+n}\} \subseteq \mathcal{S}$$

Rank Hotelling:  $RT_{m;n}^2 := T_{m;n}^2(\hat{R}_{m;n}(X_i), \hat{R}_{m;n}(Y_j))$

Data:  $\{X_i\}_{i=1}^m$  iid  $P$  (abs. cont.),  $\{Y_j\}_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$

Reference dist.:  $\pi$  on  $S \times \mathbb{R}^d$  (abs. cont.; e.g.,  $\pi = \text{Unif}([0; 1]^d)$ )

Proposed tests [D. & Sen (JASA, 2020); D., Bhattacharya & Sen (2020)]

**Joint rank map:** The sample ranks of the **pooled** observations:

$$\hat{R}_{m;n} : \{X_1, \dots, X_m, Y_1, \dots, Y_n\} \rightarrow \{c_1, \dots, c_{m+n}\} \subseteq S$$

Rank Hotelling:  $RT_{m;n}^2 := T_{m;n}^2(\hat{R}_{m;n}(X_i), \hat{R}_{m;n}(Y_j))$

**General principle:** Start with a "good" test and **replace** the  $X_i$ 's and  $Y_j$ 's with their **pooled multivariate ranks**



Data:  $\{X_i\}_{i=1}^m$  iid  $P$  (abs. cont.),  $\{Y_j\}_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$

Reference dist.:  $\pi$  on  $S \times \mathbb{R}^d$  (abs. cont.; e.g.,  $\pi = \text{Unif}([0; 1]^d)$ )

Proposed tests [D. & Sen (JASA, 2020); D., Bhattacharya & Sen (2020)]

**Joint rank map:** The sample ranks of the **pooled** observations:

$$\hat{R}_{m;n} : \{X_1, \dots, X_m, Y_1, \dots, Y_n\} \rightarrow \{c_1, \dots, c_{m+n}\} \subseteq S$$

Rank Hotelling:  $RT_{m;n}^2 := T_{m;n}^2 \left( \hat{R}_{m;n}(X_i), \hat{R}_{m;n}(Y_j) \right)$

**General principle:** Start with a "good" test and **replace** the  $X_i$ 's and  $Y_j$ 's with their **pooled multivariate ranks**

This yields the **Wilcoxon rank-sum** test when applied to **the test**.  
Therefore  $RT_{m;n}^2$  is equivalent to **Wilcoxon rank-sum** when  $d=1$

Data:  $\{X_i\}_{i=1}^m$  iid  $P$  (abs. cont.),  $\{Y_j\}_{j=1}^n$  iid  $Q$  on  $\mathbb{R}^d$ ,  $d \geq 1$

Reference dist.:  $\pi$  on  $S \times \mathbb{R}^d$  (abs. cont.; e.g.,  $\pi = \text{Unif}([0; 1]^d)$ )

Proposed tests [D. & Sen (JASA, 2020); D., Bhattacharya & Sen (2020)]

**Joint rank map:** The sample ranks of the **pooled** observations:

$$\hat{R}_{m;n} : \{X_1, \dots, X_m, Y_1, \dots, Y_n\} \rightarrow \{c_1, \dots, c_{m+n}\} \subseteq S$$

Rank Hotelling:  $RT_{m;n}^2 := T_{m;n}^2 \left( \hat{R}_{m;n}(X_i), \hat{R}_{m;n}(Y_j) \right)$

**General principle:** Start with a "good" test and **replace** the  $X_i$ 's and  $Y_j$ 's with their **pooled multivariate ranks**

This yields the **Wilcoxon rank-sum** test when applied to **the test**.  
Therefore  $RT_{m;n}^2$  is equivalent to **Wilcoxon rank-sum** when  $d=1$

Distribution-freeness [D. & Sen (JASA, 2020)]

Under  $H_0$ , distributions of  $RT_{m;n}^2$  are **free** of  $P \times Q$

Rank Hotelling test :  $m;n$   $1f RT_{m;n}^2 > c^{(m;n)}g$  | **distribution-free**

$c^{(m;n)}$  depends on  $\sigma_i$ 's,  $m$ ,  $n$ , and  $d$

Rank Hotelling test :  $m; n$   $1f RT_{m;n}^2 > c^{(m;n)}g$  | **distribution-free**

$c^{(m;n)}$  depends on  $\sigma_i$ 's,  $m$ ,  $n$ , and  $d$

Power (D., Bhattacharya, and Sen, 2021)

Under **location shift** alternatives, we have

$$\lim_{m;n \rightarrow \infty} E_{H_1}[m;n] = 1 :$$

Rank Hotelling test :  $m; n \quad 1f RT_{m;n}^2 > c^{(m;n)}g \mid$  **distribution-free**

$c^{(m;n)}$  depends on  $c_j$ 's,  $m$ ,  $n$ , and  $d$

Power (D., Bhattacharya, and Sen, 2021)

Under **location shift** alternatives, we have

$$\lim_{m;n \rightarrow 1} E_{H_1}[ \quad ] = 1 :$$

Asymptotic null distribution (D., Bhattacharya, and Sen, 2021)

Under  $H_0$ , if  $N := \frac{1}{N} \prod_{j=1}^N c_j \neq 1^d$ , then

$$RT_{m;n}^2 \rightarrow \frac{2}{d} :$$

Rank Hotelling test :  $m; n$   $1f RT_{m;n}^2 > c^{(m;n)} g$  | **distribution-free**

$c^{(m;n)}$  depends on  $c_i$ 's,  $m$ ,  $n$ , and  $d$

Power (D., Bhattacharya, and Sen, 2021)

Under **location shift** alternatives, we have

$$\lim_{m;n \rightarrow \infty} E_{H_1}[RT_{m;n}^2] = 1 :$$

Asymptotic null distribution (D., Bhattacharya, and Sen, 2021)

Under  $H_0$ , if  $N := \frac{1}{N} \sum_{j=1}^N c_j \rightarrow d$ , then

$$RT_{m;n}^2 \xrightarrow{d} \chi^2_{\frac{2}{d}}$$

**Goal**

How does rank Hotelling test compare with Hotelling  $T^2$  test?

- 1 A (very) brief introduction to optimal transport
- 2 Multivariate ranks using optimal transport
- 3 **Multivariate distribution-free tests using optimal transport**  
Rank Hotelling  $T^2$  test and Pitman efficiency  
Pitman efficiency, comparison with Hotelling  $T^2$

Question: How to compare two consistent tests  $S_N$  and  $T_N$ ?

Asymptotic relative (Pitman) efficiency (ARE) [Pitman (1948), Serfling (1980), Lehmann & Romano (2005), van der Vaart (1998)]



Question: How to compare two consistent tests  $S_N$  and  $T_N$ ?

Asymptotic relative (Pitman) efficiency (ARE) [Pitman (1948), Serfling (1980), Lehmann & Romano (2005), van der Vaart (1998)]

$X_1, \dots, X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$ ;  $\frac{m}{N} \rightarrow 2 \in (0, 1)$   
f P g 2  $\mathbb{R}^p$ : "smooth" (satisfies DQM) parametric family

Question: How to compare two consistent tests  $S_N$  and  $T_N$ ?

Asymptotic relative (Pitman) efficiency (ARE) [Pitman (1948), Serfling (1980), Lehmann & Romano (2005), van der Vaart (1998)]

$X_1, \dots, X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$ ;  $\frac{m}{N} \rightarrow 2 \in (0, 1)$   
f.g.  $\mathcal{P}_2$  on  $\mathbb{R}^2$ : "smooth" (satisfies DQM) parametric family

Test  $H_0: \theta_2 = \theta_1$  vs.  $H_1: \theta_2 = \theta_1 + \delta$ ;  $\delta \neq 0$

Question: How to compare two consistent tests  $S_N$  and  $T_N$ ?

Asymptotic relative (Pitman) efficiency (ARE) [Pitman (1948), Serfling (1980), Lehmann & Romano (2005), van der Vaart (1998)]

$X_1, \dots, X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$ ;  $\frac{m}{N} \rightarrow \lambda \in (0, 1)$   
 f P g  $\mathbb{R}^p$ : "smooth" (satisfies DQM) parametric family

Test  $H_0: \theta = \theta_0$  vs.  $H_1: \theta = \theta_1$ ;  $\theta_1 \neq \theta_0$

Fix  $\alpha \in (0, 1)$  (level) and  $1 - \beta \in (0, 1)$  (power)

Question: How to compare two consistent tests  $S_N$  and  $T_N$ ?

Asymptotic relative (Pitman) efficiency (ARE) [Pitman (1948), Serfling (1980), Lehmann & Romano (2005), van der Vaart (1998)]

$X_1, \dots, X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$ ;  $\frac{m}{N} \rightarrow \lambda \in (0, 1)$   
 f P g  $\mathbb{R}^p$ : "smooth" (satisfies DQM) parametric family

Test  $H_0: \theta = \theta_0$  vs.  $H_1: \theta = \theta_1$ ;  $\theta_1 \neq \theta_0$

Fix  $\alpha \in (0, 1)$  (level) and  $1 - \beta \in (0, 1)$  (power)

Question: How to compare two consistent tests  $S_N$  and  $T_N$ ?

Asymptotic relative (Pitman) efficiency (ARE) [Pitman (1948), Serfling (1980), Lehmann & Romano (2005), van der Vaart (1998)]

$X_1, \dots, X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$ ;  $\frac{m}{N} \rightarrow \lambda \in (0, 1)$   
 f P g  $\mathbb{R}^p$ : "smooth" (satisfies DQM) parametric family

Test  $H_0: \theta = \theta_0$  vs.  $H_1: \theta = \theta_1$ ;  $\theta_1 \neq \theta_0$

Fix  $\alpha \in (0, 1)$  (level) and  $1 - \beta \in (0, 1)$  (power)

Let  $N(\alpha, \beta)$  denote the minimum number of samples s.t.:

$$E_{H_0}[T_N] \leq N \quad \text{and} \quad E_{H_1}[T_N] \leq N$$

Question: How to compare two **consistent** tests  $S_N$  and  $T_N$ ?

Asymptotic relative (Pitman) efficiency (ARE) [Pitman (1948), Serfling (1980), Lehmann & Romano (2005), van der Vaart (1998)]

$X_1, \dots, X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$ ;  $\frac{m}{N} \rightarrow \alpha \in (0, 1)$   
 f P g 2  $\mathbb{R}^p$ : "smooth" (satisfies DQM) parametric family

Test  $H_0: \theta = \theta_0$  vs.  $H_1: \theta = \theta_1$ ;  $\theta_1 \neq \theta_0$

Fix  $\alpha \in (0, 1)$  (level) and  $1 - \beta \in (0, 1)$  (power)

Let  $N(\alpha, \beta)$  denote the **minimum** number of **samples** s.t.:

$$E_{H_0}[T_N] \leq N(\alpha, \beta) \quad \text{and} \quad E_{H_1}[T_N] \leq N(\alpha, \beta)$$

The **asymptotic (Pitman) efficiency** of  $S_N$  w.r.t.  $T_N$  is given by

$$\text{ARE}(S_N; T_N) := \lim_{\alpha \rightarrow 0} \frac{N(\alpha, \beta)}{N(\alpha, \beta)}$$

Question: How to compare two **consistent** tests  $S_N$  and  $T_N$ ?

Asymptotic relative (Pitman) efficiency (ARE) [Pitman (1948), Serfling (1980), Lehmann & Romano (2005), van der Vaart (1998)]

$X_1, \dots, X_m \stackrel{iid}{\sim} P_1$  &  $Y_1, \dots, Y_n \stackrel{iid}{\sim} P_2$ ;  $N = m + n$ ;  $\frac{m}{N} \rightarrow \lambda \in (0, 1)$   
 $f \in \mathcal{P}_g$   $\mathbb{R}^p$ : "smooth" (satisfies DQM) parametric family

Test  $H_0: \theta = \theta_0$  vs.  $H_1: \theta = \theta_1$ ;  $\theta_1 \neq \theta_0$

Fix  $\lambda \in (0, 1)$  (level) and  $1 - \beta \in (0, 1)$  (power)

Let  $N(\lambda, \beta)$  denote the **minimum** number of **samples** s.t.:

$$E_{H_0}[T_N] \leq N \quad \text{and} \quad E_{H_1}[T_N] \leq N$$

The **asymptotic (Pitman) efficiency** of  $S_N$  w.r.t.  $T_N$  is given by

$$ARE(S_N; T_N) := \lim_{\lambda \rightarrow 0} \frac{N(\lambda, \beta)}{N(\lambda, \beta)}$$

In principle,  $ARE(S_N; T_N)$  can depend on  $\lambda$  and  $\beta$ , but in many interesting cases **they don't**

$X_1, \dots, X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$

$f \in \mathcal{P}(\mathbb{R}^p)$ : "smooth" (satisfies DQM) parametric family

Consider  $H_0: \mu = \mu_1$  vs.  $H_1: \mu = \mu_1 + hN^{-1/2}$ ;  $h \in \mathbb{R}^p$

ARE ( $T_{m;n}^2; T_{m;n}^2$ ) can be derived from the distribution of both test statistics under above alternatives



$X_1, \dots, X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$

$f, g \in \mathcal{P}(\mathbb{R}^p)$ : "smooth" (satisfies DQM) parametric family

Consider  $H_0: \theta = \theta_1$  vs.  $H_1: \theta = \theta_1 + hN^{-1/2}$ ;  $h \in \mathbb{R}^p$

ARE ( $RT_{m;n}^2; T_{m;n}^2$ ) can be derived from the distribution of both test statistics under above alternatives

### Some observations

Expression of ARE ( $RT_{m;n}^2; T_{m;n}^2$ ) does not depend on  $\theta$  and

Asymp. dist. of  $RT_{m;n}^2$  can depend on the choice of

$X_1, \dots, X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$

$f \in \mathcal{P}(\mathbb{R}^p)$ : "smooth" (satisfies DQM) parametric family

Consider  $H_0: \mu_2 = \mu_1$  vs.  $H_1: \mu_2 = \mu_1 + hN^{-1/2}$ ;  $h \in \mathbb{R}^p$

ARE ( $RT_{m;n}^2; T_{m;n}^2$ ) can be derived from the distribution of both test statistics under above alternatives

### Some observations

Expression of ARE ( $RT_{m;n}^2; T_{m;n}^2$ ) does not depend on  $\mu$  and  $\Sigma$

Asymp. dist. of  $RT_{m;n}^2$  can depend on the choice of  $\mu$  and  $\Sigma$

Can we lower bound ARE for sub-classes of multivariate dists., i.e.,

$$\min_F \text{ARE} (RT_{m;n}^2; T_{m;n}^2) = ??$$

$X_1; \dots; X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1; \dots; Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$

### Independent coordinates case

$F_{\text{ind}} = f P g_2$  has density  $p(z_1; \dots; z_d) = \prod_{i=1}^d f_i(z_i)$ ,  $z_i \in \mathbb{R}^d$

### Theorem (D., Bhattacharya, and Sen (2021))

Suppose  $\frac{m}{N} \rightarrow \alpha \in (0, 1)$ . If  $\mathbf{P}_N := \frac{1}{N} \sum_{j=1}^N \delta_{z_j} \stackrel{d}{\rightarrow} \text{Unif}([0, 1]^d)$ , then

$$\min_{F_{\text{ind}}} \text{ARE}(RT_{m;n}^2; T_{m;n}^2) = 0.864$$

$X_1; \dots; X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1; \dots; Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$

### Independent coordinates case

$F_{\text{ind}} = f P g_2$  has density  $p(z_1; \dots; z_d) = \prod_{i=1}^d f_i(z_i)$ ,  $z_i \in \mathbb{R}^d$

### Theorem (D., Bhattacharya, and Sen (2021))

Suppose  $\frac{m}{N} \rightarrow \alpha \in (0, 1)$ . If  $N := \frac{1}{N} \sum_{j=1}^N \mathbf{q}_j \stackrel{\text{iid}}{\sim} \text{Unif}([0, 1]^d)$ , then

$$\min_{F_{\text{ind}}} \text{ARE}(RT_{m;n}^2; T_{m;n}^2) = 0.864$$

If  $N \stackrel{\text{iid}}{\sim} N(0; I_d)$ , then

$$\min_{F_{\text{ind}}} \text{ARE}(RT_{m;n}^2; T_{m;n}^2) = 1$$

$X_1; \dots; X_m \stackrel{\text{iid}}{\sim} P_1$  &  $Y_1; \dots; Y_n \stackrel{\text{iid}}{\sim} P_2$ ;  $N = m + n$

### Independent coordinates case

$F_{\text{ind}} = f P g_2$  has density  $p(z_1; \dots; z_d) = \prod_{i=1}^d f_i(z_i)$ ,  $z_i \in \mathbb{R}^d$

### Theorem (D., Bhattacharya, and Sen (2021))

Suppose  $\frac{m}{N} \rightarrow \alpha \in (0, 1)$ . If  $N := \frac{1}{N} \sum_{j=1}^N \epsilon_j \stackrel{\text{iid}}{\sim} \text{Unif}([0; 1]^d)$ , then

$$\min_{F_{\text{ind}}} \text{ARE}(RT_{m;n}^2; T_{m;n}^2) = 0.864$$

If  $N \stackrel{\text{iid}}{\sim} N(0; I_d)$ , then

$$\min_{F_{\text{ind}}} \text{ARE}(RT_{m;n}^2; T_{m;n}^2) = 1$$

Generalizes Hodges & Lehmann (1956), Chernoff & Savage (1958)

ARE can be arbitrarily large (can tend to  $\infty$ ) for heavy tailed dists.

## Elliptically symmetric distributions

$F_{\text{ell}} = \{P \mid g_2\}$  is class of **elliptically symmetric** distributions on  $\mathbb{R}^d$ , i.e.,

$$p(x) / (\det(\Sigma))^{-\frac{1}{2}} f\left(\frac{x - \mu}{\Sigma^{-1/2}}\right) = f\left(\frac{x - \mu}{\Sigma^{-1/2}}\right); \quad \text{for all } x \in \mathbb{R}^d:$$

## Elliptically symmetric distributions

$F_{\text{ell}} = \{f \mid P \text{ g}_2\}$  is class of **elliptically symmetric** distributions on  $\mathbb{R}^d$ , i.e.,

$$p(x) / (\det(\Sigma))^{-\frac{1}{2}} f_{\Sigma}^{-1}(x - \mu) > 1(x - \mu) ; \text{ for all } x \in \mathbb{R}^d:$$

## Theorem (D., Bhattacharya, and Sen (2021))

Suppose: (i)  $N \stackrel{!}{=} N(0; I_d)$ , (ii)  $\frac{m}{N} \stackrel{!}{=} 2(0; 1)$ . Then,

$$\min_{F_{\text{ell}}} \text{ARE}(RT_{m;n}^2; T_{m;n}^2) = 1:$$

This generalizes the famous result of **Cherno and Savage (1958)**

## Model for Independent Component Analysis (ICA)

$F_{ICA} = f_1(x_1) \dots f_d(x_d) : f_1, \dots, f_d \in \mathcal{F}_{\mathbb{R}^d}$  where  $f_1, \dots, f_d$  has the form

$$f_1(x_1; \dots; x_d) = \prod_{i=1}^d f_i \left( \sum_{j=1}^d a_{ij} x_j \right)$$

where  $f_1, f_2, \dots, f_d$  are univariate densities, and  $A = (a_{ij})_{d \times d}$  is an orthogonal matrix (unknown)

Thus,  $f_1$  is the density of  $X_{d-1}$  where

$$X = AW$$

with  $W_{d-1}$  having independent components



## Model for Independent Component Analysis (ICA)

$F_{ICA} = f_{f_1}(\cdot) : f_1 \in \mathcal{F}_{\mathbb{R}^d}$  where  $f_1 \in \mathcal{F}$  has the form

$$f_1(x_1; \dots; x_d) = \prod_{i=1}^d f_i \left( \sum_{j=1}^d a_{ij} x_j \right)$$

where  $f_1; f_2; \dots; f_d$  are univariate densities, and  $A = (a_{ij})_{d \times d}$  is an **orthogonal matrix** (unknown)

Thus,  $f_1$  is the density of  $X_{d-1}$  where

$$X = AW$$

with  $W_{d-1}$  having **independent components**

## Theorem (D., Bhattacharya, and Sen (2021))

Suppose: (i)  $N \stackrel{d}{\sim} N(0; I_d)$ , (ii)  $\frac{m}{N} \rightarrow 2 \in (0; 1)$ . Then,

$$\min_{F_{ICA}} \text{ARE} (RT_{m;n}^2; T_{m;n}^2) = 1$$

# Recap

Exact **distribution-free, nonparametric** test, gives uniformly level test and **high efficiency** compared to Hotelling  $\bar{T}^2$  test

Provides the first comprehensive extension of classical nonparametric testing to the multivariate setting [▶ skip](#)

Using **Gaussian** reference distribution ensures  $ARE(RT_{m;n}^2; T_{m;n}^2) \rightarrow 1$  for many popular subfamilies. Note that the test is **agnostic** to the underlying subfamily [▶ skip](#)

# Recap

Exact **distribution-free, nonparametric** test, gives uniformly level test and **high efficiency** compared to Hotelling  $\bar{T}^2$  test

Provides the first comprehensive extension of classical nonparametric testing to the multivariate setting [▶ skip](#)

Using **Gaussian** reference distribution ensures  $ARE(RT_{m;n}^2; T_{m;n}^2) \rightarrow 1$  for many popular subfamilies. Note that the test is **agnostic** to the underlying subfamily [▶ skip](#)

Robust against outliers and better finite-sample performance under **heavy-tailed** distributions [▶ skip](#)

# Recap

Exact **distribution-free, nonparametric** test, gives uniformly level test and **high efficiency** compared to Hotelling  $\bar{T}^2$  test

Provides the first comprehensive extension of classical nonparametric testing to the multivariate setting [▶▶ skip](#)

Using **Gaussian** reference distribution ensures  $ARE(RT_{m;n}^2; T_{m;n}^2) \rightarrow 1$  for many popular subfamilies. Note that the test is **agnostic** to the underlying subfamily [▶▶ skip](#)

Robust against outliers and better finite-sample performance under **heavy-tailed** distributions [▶▶ skip](#)

Thank you. Questions?

# Multivariate two-sample goodness-of- t test

Recall **general strategy**: Start with a "**good**" test and **replace** the  $X_i$ 's and  $Y_j$ 's with their **pooled multivariate ranks**

# Multivariate two-sample goodness-of- t test

Start with a "good" consistent test, say the energy statistic (Szekely and Rizzo, 2013) and apply our general strategy

# Multivariate two-sample goodness-of-fit test

Start with a "good" consistent test, say the energy statistic (Székely and Rizzo, 2013) and apply our general strategy

Suppose  $X; X^0 \stackrel{\text{iid}}{\sim} P, Y; Y^0 \stackrel{\text{iid}}{\sim} Q, K(s; t) := \|s - t\|_k$ , then energy dist. (or kernel MMD (see Gretton et al., 2008)):

$$E^2(P; Q) := 2 E K(X; Y) - E K(X; X^0) - E K(Y; Y^0) \geq 0$$

# Multivariate two-sample goodness-of-fit test

Start with a "good" consistent test, say the energy statistic (Székely and Rizzo, 2013) and apply our general strategy

Suppose  $X; X^0 \stackrel{\text{iid}}{\sim} P, Y; Y^0 \stackrel{\text{iid}}{\sim} Q, K(s; t) := ks - tk$ , then energy dist. (or kernel MMD (see Gretton et al., 2008)):

$$E^2(P; Q) := 2 E K(X; Y) - E K(X; X^0) - E K(Y; Y^0) = 0$$

Characterizes equality of distributions  $E(P; Q) = 0 \iff P = Q$

E-statistic:  $E_{m;n}^2(f; X_i, g; Y_j) := 2A - B - C$  where

$$A = \frac{1}{mn} \sum_{i,j=1}^{X;n} K(X_i; Y_j); \quad B = \frac{1}{m^2} \sum_{i,j=1}^{X^n} K(X_i; X_j); \quad C = \frac{1}{n^2} \sum_{i,j=1}^{X^n} K(Y_i; Y_j)$$



# Multivariate two-sample goodness-of-fit test

Start with a "good" consistent test, say the energy statistic (Székely and Rizzo, 2013) and apply our general strategy

Suppose  $X; X^0 \stackrel{\text{iid}}{\sim} P, Y; Y^0 \stackrel{\text{iid}}{\sim} Q, K(s; t) := ks - tk$ , then energy dist. (or kernel MMD (see Gretton et al., 2008)):

$$E^2(P; Q) := 2 E K(X; Y) - E K(X; X^0) - E K(Y; Y^0) \geq 0$$

Characterizes equality of distributions  $E(P; Q) = 0 \iff P = Q$

E-statistic:  $E_{m;n}^2(f; X_i, g_{j=1}^m; f; Y_j, g_{j=1}^n) := 2A - B - C$  where

$$A = \frac{1}{mn} \sum_{i,j=1}^{X;n} K(X_i; Y_j); \quad B = \frac{1}{m^2} \sum_{i,j=1}^{X;n} K(X_i; X_j); \quad C = \frac{1}{n^2} \sum_{i,j=1}^{Y;n} K(Y_i; Y_j)$$

Energy test: Reject  $H_0$  if  $E_{m;n}^2(f; X_i, g_{j=1}^m; f; Y_j, g_{j=1}^n) > \tau$  (depends on  $P$ ; we can use permutation test)

# Proposed statistic

Rank energy statistic [D. and Sen (JASA, 2020)]

**Joint rank map:** The sample ranks of the **pooled** observations:

$$\hat{R}_{m;n} : f X_1, \dots, X_m, Y_1, \dots, Y_n g \mapsto f c_1, \dots, c_{m+n} g \quad [0, 1]^d$$

Rank energy:  $RE_{m;n}^2 := E_{m;n}^2 \left( f \hat{R}_{m;n}(X_i) g_{i=1}^m ; f \hat{R}_{m;n}(Y_j) g_{j=1}^n \right)$

# Proposed statistic

Rank energy statistic [D. and Sen (JASA, 2020)]

**Joint rank map:** The sample ranks of the **pooled** observations:

$$\hat{R}_{m;n} : \{X_1, \dots, X_m, Y_1, \dots, Y_n\} \rightarrow \{c_1, \dots, c_{m+n}\} \subset [0, 1]^d$$

Rank energy:  $RE_{m;n}^2 := E_{m;n}^2 \left( \sum_{i=1}^m \hat{R}_{m;n}(X_i) g_{j=1}^m; \sum_{j=1}^n \hat{R}_{m;n}(Y_j) g_{j=1}^n \right)$

Distribution-freeness

Under  $H_0$ , distribution of  $RE_{m;n}$  is **free** of  $P = Q$ , if  $P$  is **abs. cont.**

**Dist. of  $RE_{m;n}$**  just depends on  $c_i$ 's,  $m$ ;  $n$  and  $d$

Rank energy test: Reject  $H_0$  if  $RE_{m;n} > \tau$  (**universal threshold, free of  $P = Q$** )

# Proposed statistic

## Rank energy statistic [D. and Sen (JASA, 2020)]

**Joint rank map:** The sample ranks of the **pooled** observations:

$$\hat{R}_{m;n} : \{X_1, \dots, X_m, Y_1, \dots, Y_n\} \rightarrow \{c_1, \dots, c_{m+n}\} \subset [0, 1]^d$$

Rank energy:  $RE_{m;n}^2 := E_{m;n}^2 \left( \int \hat{R}_{m;n}(X_i) g_{j=1}^m; \int \hat{R}_{m;n}(Y_j) g_{j=1}^n \right)$

## Distribution-freeness

Under  $H_0$ , distribution of  $RE_{m;n}$  is **free** of  $P = Q$ , if  $P$  is **abs. cont.**

**Dist. of  $RE_{m;n}$**  just depends on  $c_j$ 's,  $m$ ,  $n$  and  $d$

Rank energy test: Reject  $H_0$  if  $RE_{m;n} > \tau$  (**universal threshold, free of  $P = Q$** )

## Simplification when $d = 1$

$RE_{m;n}^2$  is exactly **equivalent** to the two-sample **Cramer-von Mises statistic**

## Power [D. and Sen (JASA, 2020)]

Under (ii) and  $P \notin Q$ , if  $\frac{m}{m+n} \rightarrow 2 \in (0; 1)$  then,

$$P(\text{RE}_{m;n} > \alpha^{(m;n)}) \rightarrow 1 \quad \text{as } m; n \rightarrow \infty :$$

Proposed test has **asymptotic power 1**, against **all fixed** alternatives

## Power [D. and Sen (JASA, 2020)]

Under (ii) and  $P \notin Q$ , if  $\frac{m}{m+n} \rightarrow 2 \in (0; 1)$  then,

$$P(\text{RE}_{m;n} > \chi^2_{(m;n)}(1-\alpha)) \rightarrow 1 \quad \text{as } m, n \rightarrow \infty :$$

Proposed test has **asymptotic power 1**, against **all fixed** alternatives

## Limiting distribution under $H_0$ [D. and Sen (JASA, 2020)]

If (i)  $P \in Q$  is **abs. cont.**, and

$$(ii) \frac{1}{N} \sum_{i=1}^N c_i \xrightarrow{d} \text{a.s. } (\mathbb{N} = m + n)$$

Then, under  $H_0$ , there is a **universal** distribution s.t.

$$\frac{mn}{m+n} \text{RE}_{m;n}^2 \xrightarrow{d} \sum_{j=1}^{\infty} Z_j^2 \quad \text{as } m, n \rightarrow \infty \quad \text{where } \sum_{j=1}^{\infty} \lambda_j = 1$$

# Pitman efficiency

## Some observations

Efficiency of  $RE_{m;n}$  w.r.t.  $E_{m;n}$  depends on the type I error and power, which makes it hard to obtain efficiency lower bounds

Existing tests which are both consistent and distribution-free usually do not have Pitman efficiency

## Some observations

Efficiency of  $RE_{m;n}$  w.r.t.  $E_{m;n}$  depends on the type I error and power, which makes it hard to obtain efficiency lower bounds

Existing tests which are both consistent and distribution-free usually do not have Pitman efficiency

Fix a level parameter  $\alpha \in (0, 1)$ . Consider:

$$H_0 : \mu_1 = 0 \quad \text{versus} \quad H_1 : \mu_1 = hN^{-1/2}$$



## Some observations

Efficiency of  $RE_{m;n}$  w.r.t.  $E_{m;n}$  depends on the type I error and power, which makes it hard to obtain efficiency lower bounds

Existing tests which are both consistent and distribution-free usually do not have Pitman efficiency

Fix a level parameter  $\alpha \in (0, 1)$ . Consider:

$$H_0 : \theta_1 = 0 \quad \text{versus} \quad H_1 : \theta_1 = hN^{-1/2}$$

In [Bhattacharya \(2019\)](#), the author showed that for many asymptotic distribution-free tests

$$P_{H_1}(T_{m;n} \text{ rejects } H_0) \rightarrow \alpha \quad (\text{powerless})$$

## Some observations

Efficiency of  $RE_{m;n}$  w.r.t.  $E_{m;n}$  depends on the type I error and power, which makes it hard to obtain efficiency lower bounds

Existing tests which are both consistent and distribution-free usually do not have Pitman efficiency

Fix a level parameter  $\alpha \in (0; 1)$ . Consider:

$$H_0 : \theta = 0 \quad \text{versus} \quad H_1 : \theta = hN^{-1/2}$$

In Bhattacharya (2019), the author showed that for many asymptotic distribution-free tests

$$P_{H_1}(T_{m;n} \text{ rejects } H_0) \rightarrow 0 \quad (\text{powerless})$$

## Rank Energy $RE_{m;n}$ [D., Bhattacharya, and Sen (working paper)]

$$\lim_{m;n \rightarrow \infty} P_{H_1}(RE_{m;n} \text{ rejects } H_0) > 0$$


Only consistent, exactly dist.-free test that can distinguish  $H_0$  &  $H_1$

# Summary

Introduced **optimal transport** and obtained first tuning-free, **minimax optimal estimator** of optimal transport map

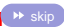
# Summary

Introduced **optimal transport** and obtained first tuning-free, **minimax optimal estimator** of optimal transport map

**Multivariate distribution-free** nonparametric testing procedures with **high efficiency**, constructed using **optimal transport** 

# Summary


Introduced **optimal transport** and obtained **fast tuning-free, minimax optimal estimator** of optimal transport map

**Multivariate distribution-free** nonparametric testing procedures with **high efficiency**, constructed using **optimal transport** 

Proposed a **general framework**, other examples include **independence testing**, testing for **symmetry**, testing **equality of distributions** ...

# Summary

Introduced **optimal transport** and obtained first tuning-free, **minimax optimal estimator** of optimal transport map


**Multivariate distribution-free** nonparametric testing procedures with **high efficiency**, constructed using **optimal transport** 

Proposed a **general framework**, other examples include **independence testing**, testing for **symmetry**, testing **equality of distributions** ...

**Independence testing**: In **D., Bhattacharya, and Sen (2021)** we obtain multivariate **distribution-free** extensions of **Spearman's correlation** and **kernel tests** of dependence; obtain similar results


# Summary

Introduced **optimal transport** and obtained **fast tuning-free, minimax optimal estimator** of optimal transport map

**Multivariate distribution-free** nonparametric testing procedures with **high efficiency**, constructed using **optimal transport** 

Proposed a **general framework**, other examples include **independence testing**, testing for **symmetry**, testing **equality of distributions** ...

**Independence testing**: In **D., Bhattacharya, and Sen (2021)** we obtain multivariate **distribution-free** extensions of **Spearman's correlation** and **kernel tests** of dependence; obtain similar results

Tuning-free, robust, **computationally feasible** procedures, performs particularly well under heavy-tailed data 

# Summary



- D. and Sen, (2020). <https://arxiv.org/pdf/1909.08733> (JASA, to appear)
- D., Ghosal, and Sen (2021). <https://arxiv.org/abs/2107.01718> (NeurIPS, 2021)
- D., Bhattacharya, and Sen (2021). <https://arxiv.org/abs/2104.01986>
- D., Bhattacharya and Sen (2021+). (working paper)

# Nonparametric association/conditional association

Suppose

$$\begin{matrix} X \\ Y \end{matrix} \sim N \begin{matrix} 2 \\ 2 \end{matrix} \begin{matrix} 0 & 1 \\ 0 & 1 \end{matrix}$$

The correlation measures **association** (linear) between  $X$  and  $Y$

# Nonparametric association/conditional association

Suppose

$$\begin{matrix} X \\ Y \end{matrix} \sim N_2 \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix}$$

The correlation measures **association** (linear) between  $X$  and  $Y$

$= 0$  i  $X$  and  $Y$  are independent

$= 1$  i  $Y$  is an **exact, measurable** (linear) function of  $X$

# Nonparametric association/conditional association

Suppose

$$\begin{matrix} X \\ Y \end{matrix} \sim N \begin{pmatrix} 0 \\ 0 \end{pmatrix} ; \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$$

The correlation measures **association** (linear) between  $X$  and  $Y$

$\rho = 0$  if  $X$  and  $Y$  are independent

$\rho = 1$  if  $Y$  is an **exact, measurable** (linear) function of  $X$

Can we find a **nonparametric measure of association** between random elements on **topological spaces**?

# Nonparametric association/conditional association

Suppose

$$\begin{matrix} X \\ Y \end{matrix} \sim N \begin{pmatrix} 0 \\ 0 \end{pmatrix}; \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

The correlation measures **association** (linear) between  $X$  and  $Y$

$= 0$  if  $X$  and  $Y$  are independent

$= 1$  if  $Y$  is an **exact, measurable** (linear) function of  $X$

Can we find a **nonparametric measure of association** between random elements on **topological spaces**?

**Goal:** Given random elements  $(X_1; Y_1); \dots; (X_n; Y_n) \stackrel{i.i.d.}{\sim} (X; Y)$ , define  $T(X; Y)$  and  $T_n$  (**estimator**), such that:

# Nonparametric association/conditional association

Suppose

$$\begin{matrix} X \\ Y \end{matrix} \sim N \begin{pmatrix} 0 \\ 0 \end{pmatrix}; \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

The correlation measures **association** (linear) between  $X$  and  $Y$

$= 0$  i  $X$  and  $Y$  are independent

$= 1$  i  $Y$  is an **exact, measurable** (linear) function of  $X$

Can we find a **nonparametric measure of association** between random elements on **topological spaces**?

**Goal:** Given random elements  $(X_1; Y_1); \dots; (X_n; Y_n) \stackrel{i.i.d.}{\sim} (X; Y)$ , define  $T(X; Y)$  and  $T_n$  (**estimator**), such that:

$$T_n; T(X; Y) \in [0; 1]$$

$T(X; Y) = 0$  i  $X$  and  $Y$  are **independent**

# Nonparametric association/conditional association

Suppose

$$\begin{matrix} X \\ Y \end{matrix} \sim N \begin{pmatrix} 0 \\ 0 \end{pmatrix}; \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

The correlation measures **association** (linear) between  $X$  and  $Y$

$= 0$  i  $X$  and  $Y$  are independent

$= 1$  i  $Y$  is an **exact, measurable** (linear) function of  $X$

Can we find a **nonparametric measure of association** between random elements on **topological spaces**?

**Goal:** Given random elements  $(X_1; Y_1); \dots; (X_n; Y_n) \stackrel{i.i.d.}{\sim} (X; Y)$ , define  $T(X; Y)$  and  $T_n$  (**estimator**), such that:

$$T_n; T(X; Y) \in [0; 1]$$

$T(X; Y) = 0$  i  $X$  and  $Y$  are **independent**

$T(X; Y) = 1$  i  $Y = f(X)$  for (unknown) meas. function  $f(\cdot)$

# Nonparametric association/conditional association

Suppose

$$\begin{matrix} X \\ Y \end{matrix} \sim N \begin{pmatrix} 0 & 1 \\ 0 & 1 \end{pmatrix}$$

The correlation measures **association** (linear) between  $X$  and  $Y$

$= 0$  if  $X$  and  $Y$  are independent

$= 1$  if  $Y$  is an **exact, measurable** (linear) function of  $X$

Can we find a **nonparametric measure of association** between random elements on **topological spaces**?

**Goal:** Given random elements  $(X_1; Y_1); \dots; (X_n; Y_n) \stackrel{i.i.d.}{\sim} (X; Y)$ , define  $T(X; Y)$  and  $T_n$  (**estimator**), such that:

$$T_n; T(X; Y) \in [0; 1]$$

$T(X; Y) = 0$  if  $X$  and  $Y$  are **independent**

$T(X; Y) = 1$  if  $Y = f(X)$  for (unknown) meas. function  $f(\cdot)$

$$T_n \xrightarrow{p.s.} T(X; Y)$$



# Measure of Association

We answer this question in the **adaptive** by combining ideas from **reproducing kernel Hilbert spaces (RKHS)** and **geometric graphs** (e.g., nearest neighbors, minimum spanning trees), to come up with a large **class** of such measures

Our measures are completely **nonparametric** (unlike

We can also extend this to measuring **conditional association** with applications in **variable selection**, conditional independence testing ...

# Measure of Association

We answer this question in the **a rnative** by combining ideas from **reproducing kernel Hilbert spaces (RKHS)** and **geometric graphs** (e.g., nearest neighbors, minimum spanning trees), to come up with a large **class** of such measures

Our measures are completely **nonparametric** (unlike

We can also extend this to measuring **conditional association** with applications in **variable selection**, conditional independence testing ...

## References:

- D., Ghosal, and Sen (2020). <https://arxiv.org/pdf/2010.01768.pdf> (RR at AoS)
- Huang, D., and Sen (2021). <https://arxiv.org/pdf/2012.14804.pdf> (JMLR, to appear)
- Auddy, D., and Nandy (2021). <https://arxiv.org/pdf/2104.15140.pdf> (RR at Bernoulli)

## Applied probability (Ising type models):

D. & Mukherjee (2020). <https://arxiv.org/pdf/2005.00710.pdf> (AoAP, to appear)

D., Mukherjee, Mukherjee & Yuan (2021). <https://arxiv.org/abs/2012.05784>

## Applied probability (Ising type models):

D. & Mukherjee (2020). <https://arxiv.org/pdf/2005.00710.pdf> (AoAP, to appear)

D., Mukherjee, Mukherjee & Yuan (2021). <https://arxiv.org/abs/2012.05784>

## Measuring (conditional) association on topological data:

D., Ghosal & Sen (2020). <https://arxiv.org/pdf/2010.01768.pdf> (RR at AoS)

Huang, D. & Sen (2021). <https://arxiv.org/pdf/2012.14804.pdf> (JMLR, to appear)

Auddy, D. & Nandy (2021). <https://arxiv.org/pdf/2104.15140.pdf> (RR at Bernoulli)

## Applied probability (Ising type models):

D. & Mukherjee (2020). <https://arxiv.org/pdf/2005.00710.pdf> (AoAP, to appear)

D., Mukherjee, Mukherjee & Yuan (2021). <https://arxiv.org/abs/2012.05784>

## Measuring (conditional) association on topological data:

D., Ghosal & Sen (2020). <https://arxiv.org/pdf/2010.01768.pdf> (RR at AoS)

Huang, D. & Sen (2021). <https://arxiv.org/pdf/2012.14804.pdf> (JMLR, to appear)

Auddy, D. & Nandy (2021). <https://arxiv.org/pdf/2104.15140.pdf> (RR at Bernoulli)

## Gaussian mixture models and multiple testing:

D., Saha, Guntuboyina & Sen (2020).

<https://www.tandfonline.com/doi/full/10.1080/01621459.2021.1888739>

(JASA, published online)

## Causal inference:

Ghosh, D., Karmakar & Sen (2021). <https://arxiv.org/abs/2111.15524>

▶ skip

# Rate of convergence result

$$T_0 = \arg \min_{T \# P=Q} \int_{\mathcal{Z}} k(x) |T(x) - 1|^2 dP(x);$$
$$W_2^2(P; Q) = \min_{T \# P=Q} \int_{\mathcal{Z}} k(x) |T(x) - 1|^2 dP(x)$$

Data:  $X_1; X_2; \dots; X_n$  iid  $P$  and  $Y_1; \dots; Y_n$  iid  $Q$

# Rate of convergence result

$$T_0 = \arg \min_{T \# P=Q} \int_{\mathcal{Z}} k(x) |T(x)|^2 dP(x);$$
$$W_2^2(P; Q) = \min_{T \# P=Q} \int_{\mathcal{Z}} k(x) |T(x)|^2 dP(x)$$

Data:  $X_1; X_2; \dots; X_n$  iid  $P$  and  $Y_1; \dots; Y_n$  iid  $Q$

Estimator:

$$\hat{T} := \arg \min_{T \# P_n=Q_n} \int_{\mathcal{Z}} k(x) |T(x)|^2 dP_n(x)$$

# Rate of convergence result

$$T_0 = \arg \min_{T \# P=Q} \int_{\mathcal{X}} k(x, T(x)) k^2 dP(x);$$
$$W_2^2(P; Q) = \min_{T \# P=Q} \int_{\mathcal{X}} k(x, T(x)) k^2 dP(x)$$

Data:  $X_1; X_2; \dots; X_n$  iid  $P$  and  $Y_1; \dots; Y_n$  iid  $Q$

Estimator:

$$\hat{T} := \arg \min_{T \# P_n=Q_n} \int_{\mathcal{X}} k(x, T(x)) k^2 dP_n(x)$$

Rate of convergence (D., Ghosal, Sen, NeurIPS, 2021)

Assume that  $T_0$  is Lipschitz, and both  $P$  and  $Q$  are compactly supported (**can be relaxed**). Then,

$$\frac{1}{n} \sum_{i=1}^n k(\hat{T}(X_i), T_0(X_i)) k^2 \leq n^{-\frac{2}{d}}$$

for  $d \geq 4$ . For  $d = 1; 2; 3$ , the rates are  $n^{-1}$ ,  $n^{-2/3}$  and  $n^{-4/7}$ .



Nonparametric regression: Say  $Y_i = f_0(X_i) + \epsilon_i$ ,  $i = 1, \dots, n$ , where  $\epsilon_i$ 's are iid  $N(0; 1)$ , and  $f_0 \in F$

$$\hat{f}_n := \arg \min_{f \in F} \sum_{i=1}^n (Y_i - f(X_i))^2$$

Both  $\hat{f}_n$  and  $f_0$  belong to  $F$ , which yields the **basic inequality**:

$$\sum_{i=1}^n (Y_i - \hat{f}_n(X_i))^2 \geq \sum_{i=1}^n (Y_i - f_0(X_i))^2$$

Nonparametric regression: Say  $Y_i = f_0(X_i) + \epsilon_i$ ,  $i = 1, \dots, n$ , where  $\epsilon_i$ 's are iid  $N(0; 1)$ , and  $f_0 \in F$

$$\hat{f}_n := \arg \min_{f \in F} \sum_{i=1}^n (Y_i - f(X_i))^2$$

Both  $\hat{f}_n$  and  $f_0$  belong to  $F$ , which yields the **basic inequality**:

$$\sum_{i=1}^n (f_0(X_i) - \hat{f}_n(X_i))^2 \leq 2 \sum_{i=1}^n (\hat{f}_n(X_i) - f_0(X_i))^2$$

Nonparametric regression: Say  $Y_i = f_0(X_i) + \epsilon_i$ ,  $i = 1, \dots, n$ , where  $\epsilon_i$ 's are iid  $N(0; 1)$ , and  $f_0 \in F$

$$\hat{f}_n := \arg \min_{f \in F} \sum_{i=1}^n (Y_i - f(X_i))^2$$

Both  $\hat{f}_n$  and  $f_0$  belong to  $F$ , which yields the **basic inequality**:

$$\sum_{i=1}^n (Y_i - \hat{f}_n(X_i))^2 \leq \sum_{i=1}^n (Y_i - f_0(X_i))^2$$

Nonparametric regression: Say  $Y_i = f_0(X_i) + \epsilon_i$ ,  $i = 1, \dots, n$ , where  $\epsilon_i$ 's are iid  $N(0, 1)$ , and  $f_0 \in \mathcal{F}$

$$\hat{f}_n := \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n (Y_i - f(X_i))^2$$

Both  $\hat{f}_n$  and  $f_0$  belong to  $\mathcal{F}$ , which yields the **basic inequality**:

$$\sum_{i=1}^n (Y_i - \hat{f}_n(X_i))^2 \leq \sum_{i=1}^n (Y_i - f_0(X_i))^2$$

OT problem:

$$\hat{T} := \arg \min_{T \in \mathcal{P}_n} \int k(x, T(x))^2 dP_n(x)$$

Constraint set:  $T_n := \{T : T \in \mathcal{P}_n\}$ .

$\hat{T}_n \in T_n$  but  $T_0 \notin T_n$

## Dual form

Alternatively,

$$W_2^2(P_n; Q_n) = \min_{T \# P_n = Q_n} \int \kappa(x, T(x)) k^2 dP_n(x) = \min_{f;g} \int f dP_n + \int g dQ_n$$

such that  $f(x) + g(y) \leq \kappa(x, y)k^2$  for all  $x, y \in \mathbb{R}^d$ .

## Dual form

Alternatively,

$$W_2^2(P_n; Q_n) = \min_{T \# P_n = Q_n} \int \|\mathbf{x} - T(\mathbf{x})\|^2 dP_n(\mathbf{x}) = \min_{f; g} \int f dP_n + \int g dQ_n$$

such that  $f(\mathbf{x}) + g(\mathbf{y}) \leq \|\mathbf{x} - \mathbf{y}\|^2$  for all  $\mathbf{x}; \mathbf{y} \in \mathbb{R}^d$ .

Note that the constraints are **not data driven**.

## Basic inequality (D., Ghosal and Sen, 2021)

Suppose  $T_0$  is Lipschitz. Write  $T_0 = r \circ \rho$  and  $Q_n := T_0 \# P_n$ . Then,

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - T_0(\mathbf{x}_i)\|^2 \leq W_2^2(P_n; Q_n) \leq W_2^2(P_n; Q_n) + \int g d(Q_n - Q_n)$$

where  $g(\mathbf{y}) = \rho(\mathbf{y}) - (1-L)\|\mathbf{y}\|^2$ ,  $\rho(\mathbf{y}) := \sup_{\mathbf{x} \in \mathbb{R}^d} (\|\mathbf{x}\| - L\|\mathbf{y} - \mathbf{x}\|)$   
 (Legendre-Fenchel dual of  $\rho(\cdot)$ )

Proof requires arguments from **convex analysis**

## Dual form

Alternatively,

$$W_2^2(P_n; Q_n) = \max_{f;g} \int f dP_n + \int g dQ_n$$

such that  $f(x) + g(y) \leq \|x - y\|^2$  for all  $x, y \in \mathbb{R}^d$ .

Note that the constraints are **not data driven**.

## Basic inequality (D., Ghosal and Sen, 2021)

Suppose  $T_0$  is Lipschitz. Write  $T_0 = r \circ \rho$  and  $Q_n := T_0 \# P_n$ . Then,

$$\frac{1}{n} \sum_{i=1}^n \|T_n(X_i) - T_0(X_i)\|^2 \leq W_2^2(P_n; Q_n) \leq W_2^2(P_n; Q_n) + \int g d(Q_n - Q_n)$$

where  $\rho(y) = \rho(y) = (1/2)\|y\|^2$ ,  $\rho(y) := \sup_{x \in \mathbb{R}^d} (\langle y, x \rangle - \rho(x))$   
(Legendre-Fenchel dual of  $\rho(\cdot)$ )

Proof requires arguments from **convex analysis**

Using the **dual** form of  $W_2^2(\cdot; \cdot)$ , coupled with **chaining** and **Talagrand's concentration** inequality proves the rate of convergence result

Thank you. Questions?



# Properties

$$T_0 \stackrel{???}{:=} \arg \min_{T \# P = Q} \int_Z \|x - T(x)\|^2 dP(x); \quad T \# P = Q, \quad X \subseteq P; T(X) \subseteq Q:$$

Does a solution always exist?

Is the solution unique?

# Properties

$$T_0 \stackrel{???}{:=} \arg \min_{T \# P=Q} \int_{\mathcal{X}} \|T(x) - k\|^2 dP(x); \quad T \# P = Q, \quad \mathcal{X} \subseteq \mathcal{P}; \quad T(\mathcal{X}) \subseteq \mathcal{Q}$$

Does a solution always exist?

No! Take  $P = \delta_0$  and  $Q = \text{Unif}[0; 1]$ .

Is the solution unique?

No! Take  $P = 0.5 \delta_p + 0.5 \delta_{-p}$  and  $Q = 0.5 \delta_q + 0.5 \delta_{-q}$ .

# Properties

$$T_0 \stackrel{???}{:=} \arg \min_{T \# P = Q} \int_{\mathcal{X}} \|T(x) - k\|^2 dP(x); \quad T \# P = Q, \quad \mathcal{X} \subseteq \mathcal{Y}; \quad T(X) \sim Q:$$

Does a solution always exist?

No! Take  $P = \delta_0$  and  $Q = \text{Unif}[0; 1]$ .

Is the solution unique?

No! Take  $P = 0.5 \delta_p + 0.5 \delta_{-p}$  and  $Q = 0.5 \delta_q + 0.5 \delta_{-q}$ .

# Crossmatch test (Rosenbaum 2005)

## Pitman asymptotics for crossmatch test (Rosenbaum 2005)

Consider the testing set-up from before (with additional regularity assumptions). Then, for any  $\eta$ , we have:

$$\lim_{m;n \uparrow} P_{H_1}(T_{m;n} \text{ rejects } H_0) = \eta :$$

# Crossmatch test (Rosenbaum 2005)

## Pitman asymptotics for crossmatch test (Rosenbaum 2005)

Consider the testing set-up from before (with additional regularity assumptions). Then, for any  $\eta$ , we have:

$$\lim_{m;n \uparrow} P_{H_1}(T_{m;n} \text{ rejects } H_0) = \eta$$

Therefore, crossmatch test **does not** distinguish between the null and the alternative at the contiguous scale

The same phenomena happens for many **other graph-based asymptotically distribution-free tests**, see [Bhattacharya \(2019, Theorem 3.1\)](#)

# Power plot with varying sample size

Figure:  $X_1; Y_1$  are i.i.d. Epanechnikov with location parameters  $0$  and  $1$  respectively.  $X_2; X_3 \sim X_1; Y_2; Y_3 \sim Y_1$  and  $X := (X_1; X_2; X_3)$ ,  $Y := (Y_1; Y_2; Y_3)$ . Here

$$e(\text{RankUnif}; \text{Hotelling}) = 0.864$$

and  $e(\text{RankGaussian}; \text{Hotelling}) > 1$  [▶ skip](#)

# Power plot with varying location parameter

Log-normal location problem (slightly heavy-tailed)

Figure:  $U_1; U_2$  are iid standard normal, and  $V_1; V_2$  are normal with variance 1 and varying mean. Define  $X_i := \exp(U_i)$  and  $Y_i := \exp(V_i)$ . Set  $X := (X_1; X_2)$  and  $Y := (Y_1; Y_2)$  | sample size  $n = 200$  [▶ skip](#)

# Power plot with varying location parameter

**Figure:** (Left panel)  $X_1, Y_1$  are i.i.d. normal with mean 0 and  $\sigma^2$  respectively (and unit variance).  $X_2, X_3 \sim X_1; Y_2, Y_3 \sim Y_1$  and  $X := (X_1; X_2; X_3)$ . Similarly define  $Y$ . [▶ skip](#)

(Right panel)  $U := (U_1; U_2; U_3)$  and  $V := (V_1; V_2; V_3)$  where  $U_i = \exp(X_i)$ ,  $V_i = \exp(Y_i)$  and  $X_1; X_2; X_3; Y_1; Y_2; Y_3$  has the same distribution as above.

**Red - Rank energy, Black - Crossmatch, Blue - Energy, Green - HHG.**



## More simulations

	(RB)	(HHG)	(EN)	(REN)
V1	0.13	0.15	0.13	0.34
V2	0.34	0.94	0.94	0.89
V3	0.41	0.34	0.34	0.46
V4	0.34	0.31	0.33	0.32
V5	0.73	0.70	0.56	0.93
V6	0.90	0.88	0.82	0.99
V7	0.13	0.51	0.65	0.63
V8	0.11	0.39	0.35	0.43
V9	0.06	1.00	0.97	1.00
V10	0.28	0.99	1.00	0.59

**Table:** Proportion of times the null hypothesis was rejected across 10 settings. Here  $n = 200$ ,  $d = 3$ . Here RB - Rosenbaum's crossmatch test ([Rosenbaum, 2005](#)), HHG - Heller, Heller and Gor ne ([Heller et al., 2013](#)), En - energy statistic ([Székely and Rizzo, 2013](#)).

# Asymptotic stabilization

	(100)	(300)	(500)	(700)	(900)
0.05	0.39	0.40	0.39	0.40	0.40
0.1	0.36	0.36	0.36	0.36	0.36

Table: Thresholds for  $\alpha = 0.05, 0.1$  and  $n = 100; 300; 500; 700; 900$ ,  $d = 2$ .

	(100)	(300)	(500)	(700)	(900)
0.05	1.37	1.38	1.38	1.38	1.38
0.1	1.34	1.35	1.35	1.35	1.35

Table: Thresholds for  $\alpha = 0.05, 0.1$  and  $n = 100; 300; 500; 700; 900$ ,  $d = 8$ .

# What happens for $\alpha = 1$ ?

$$T_0 \stackrel{???}{:=} \arg \min_{T \# P=Q} \int_{\mathcal{X}} |T(x) - x|^2 dP(x):$$

Assume  $Q = \text{Unif}[0; 1]$  and  $X \sim P$  with cdf  $F$

# What happens for $d = 1$ ?

$$T_0 \stackrel{???}{:=} \arg \min_{T \# P=Q} \int_{\mathbb{R}} |T(x) - x|^2 dP(x):$$

Assume  $Q = \text{Unif}[0; 1]$  and  $X \sim P$  with cdf  $F$

Given  $x_1, x_2 \in \mathbb{R}$ , note that

$$(x_1 - T_0(x_1))^2 + (x_2 - T_0(x_2))^2 = (x_1 - T_0(x_2))^2 + (x_2 - T_0(x_1))^2$$

Expect  $T_0(\cdot)$  to be monotone and  $T_0(X) \stackrel{d}{=} \text{Unif}[0; 1]$ .

$T_0(\cdot)$  is the **distribution function** of  $X$ , say  $F(\cdot)$

# What happens for $\alpha = 1$ ?

$$F = \arg \min_{T \# P=Q} \int_{\mathbb{R}} |T(x) - x|^2 dP(x):$$

Assume  $Q = \text{Unif}[0; 1]$  and  $X \sim P$  with cdf  $F$

Given  $x_1, x_2 \in \mathbb{R}$ , note that

$$(x_1 - T_0(x_1))^2 + (x_2 - T_0(x_2))^2 = (x_1 - T_0(x_2))^2 + (x_2 - T_0(x_1))^2$$

,  $T_0(x_1) < T_0(x_2)$

Expect  $T_0(\cdot)$  to be monotone and  $T_0(X) \stackrel{d}{=} \text{Unif}[0; 1]$ .

$T_0(\cdot)$  is the **distribution function** of  $X$ , say  $F(\cdot)$

Note that increasing functions can be viewed as "**derivatives**" of **convex functions**.

Suppose  $T_0 \in C^{\alpha}$  (Hölder or Sobolev class),  $\alpha > 1$

The minimax rate of convergence is

$$n^{-\frac{2\alpha}{2+2\alpha}} + n^{-1}$$

Suppose  $T_0 \in C^{\alpha}$  (Hölder or Sobolev class),  $\alpha > \frac{1}{2}$

The minimax rate of convergence is

$$n^{-\frac{2\alpha}{2+2\alpha}} + n^{-1}$$

In ongoing work, we can show that using a kernel density based estimator yields the optimal rate (up to log-factors)

$$\frac{1}{n} \sum_{i=1}^n E k(\hat{T}(X_i) - T_0(X_i))^2 \sim n^{-\frac{2\alpha}{2+2\alpha}} + n^{-1}$$

In Manole et al. (2021), Hutter and Rigollet (2019), wavelet based estimators have been used to get optimal rate [▶ skip](#)

$T_0(X) = Y$ ,  $X$  with density  $f$ ,  $Y$  with density  $g$ . Then change of variable formula implies

$$g(T(x)) \det(J_{T_0}(x)) = f(x)$$

Estimating anti-derivative of  $f$ ;  $g$  related to estimating  $T_0$



$T_0(X) = Y, X$  with density  $f$ ,  $Y$  with density  $g$ . Then change of variable formula implies

$$g(T(x)) \det(J_{T_0}(x)) = f(x)$$

Estimating anti-derivative of  $f$ ;  $g$  related to estimating  $T_0$

(Caarelli Regularity, 1992, 1996) |  $T_0 \in C^1$  corresponds to  $f, g \in C^1, \lambda > 1$

$T_0(X) = Y$ ,  $X$  with density  $f$ ,  $Y$  with density  $g$ . Then change of variable formula implies

$$g(T(x)) \det(J_{T_0}(x)) = f(x)$$

Estimating anti-derivative of  $f$ ;  $g$  related to estimating  $T_0$

(Caarelli Regularity, 1992, 1996) |  $T_0 \in C^1$  corresponds to  $f, g \in C^1, \alpha > 1$

Goal is to estimate the anti-derivative of  $C^1$  functions

$T_0(X) = Y, X$  with density  $f, Y$  with density  $g$ . Then change of variable formula implies

$$g(T(x)) \det(J_{T_0}(x)) = f(x)$$

Estimating anti-derivative of  $f; g$  related to estimating  $T_0$

(Caarelli Regularity, 1992, 1996) |  $T_0 \in C^1$  corresponds to  $f, g \in C^1, \alpha > 1$

Goal is to estimate the anti-derivative of  $C^1$  functions

(Muller and Gasser, 1979) | optimal minimax lower bounds for estimating  $k$ -th derivative of  $\alpha$ -smooth functions is

$$n^{-\frac{2(\alpha-k)}{2+\alpha}}$$

$T_0(X) = Y, X$  with density  $f, Y$  with density  $g$ . Then change of variable formula implies

$$g(T(x)) \det(J_{T_0}(x)) = f(x)$$

Estimating anti-derivative of  $f; g$  related to estimating  $T_0$

(Caarelli Regularity, 1992, 1996) |  $T_0 \in C^1$  corresponds to  $f; g \in C^1, \alpha > 1$

Goal is to estimate the anti-derivative of  $C^1$  functions

(Muller and Gasser, 1979) | optimal minimax lower bounds for estimating  $k$ -th derivative of  $\alpha$ -smooth functions is

$$n^{-\frac{2(\alpha - k)}{2 + d}}$$

Use  $\alpha = 1$  and  $k = 1$  (anti-derivative), the lower bound reduces to  $n^{-\frac{2}{2+d}}$  [▶ skip](#)

Choose the **reference** distribution as spherical uniform

Choose the **reference** distribution as spherical uniform

If  $X$  is **spherically symmetric**, then

$$R(X) = \frac{X}{\|X\|} G(\|X\|)$$

where  $G$  is the dist. fn. of  $\|X\|$

Choose the **reference** distribution as spherical uniform

If  $X$  is **spherically symmetric**, then

$$R(X) = \frac{X}{\|X\|} G(\|X\|)$$

where  $G$  is the dist. fn. of  $\|X\|$

**Pooled**:  $\|X_1\|; \dots; \|X_m\|; \|Y_1\|; \dots; \|Y_n\|$ . Let  $G_{m;n}$  be the **empirical cdf** of the pooled data

Choose the **reference** distribution as spherical uniform

If  $X$  is **spherically symmetric**, then

$$R(X) = \frac{X}{\|X\|} G(\|X\|)$$

where  $G$  is the dist. fn. of  $\|X\|$

**Pooled**:  $\|X_1\|; \dots; \|X_m\|; \|Y_1\|; \dots; \|Y_n\|$ . Let  $G_{m;n}$  be the **empirical cdf** of the pooled data

**Modified Rank Hotelling**  $T^2$ :

$$RT_{m;n}^2 := T_{m;n}^2 \left( \frac{X_1}{\|X_1\|} G_{m;n}(\|X_1\|); \dots; \frac{Y_1}{\|Y_1\|} G_{m;n}(\|Y_1\|); \dots \right)$$

Test is **distribution-free** | if  $Y \stackrel{d}{=} X$  then **detection boundary** at  $\frac{d}{d+n}$  [▶ skip](#)



## $d > 1$ | A step in the right direction

### Brenier, '91, McCann '95, Polar Factorization Theorem

Assume that  $P$  is absolutely continuous on  $\mathbb{R}^d$ , Then there exists a **unique** ( $P$  a.s.)  $T_0 : \mathbb{R}^d \rightarrow \mathbb{R}^d$  such that  $T_0(\cdot)$  is the **gradient of a convex function** and

$$T_0\# P = Q:$$

If both  $P$  and  $Q$  have finite second moments, then  $T_0(\cdot)$  solves

$$\min_{T\# P=Q} \int_{\mathbb{R}^d} \|x - T(x)\|^2 dP(x):$$

Existence of  $T_0(\cdot)$  **does not require** any moment assumptions

Uniqueness:  $T_0\# P = Q$  and  $T_0\# R = Q$  will imply  **$P = R$**

# Rank functions as transport maps: Where 1

$X \sim F$  on  $\mathbb{R}$ ,  $F$  abs. cont. c.d.f.

Rank: The rank of  $x \in \mathbb{R}$  is  $F(x)$  (aka the c.d.f. at  $x$ )

Property:  $F(X) \sim \text{Uniform}([0, 1])$

Thus,  $F$  transports the distribution of  $X$  to  $U \sim \text{Uniform}([0, 1])$

# Rank functions as transport maps: Why?

$X \sim F$  on  $\mathbb{R}$ ,  $F$  abs. cont. c.d.f.

Rank: The **rank** of  $x \in \mathbb{R}$  is  $F(x)$  (aka the **c.d.f.** at  $x$ )

Property:  $F(X) \sim \text{Uniform}([0, 1])$

Thus,  $F$  **transports** the distribution of  $X$  to  $U \sim \text{Uniform}([0, 1])$

In fact, if  $E[X^2] < \infty$ , c.d.f.  $F$  is the **optimal transport map** as

$$F = \arg \min_{T: T(X) \stackrel{d}{=} U} E|X - T(X)|^2$$

# Rank functions as transport maps: Why?

$X \sim F$  on  $\mathbb{R}$ ,  $F$  abs. cont. c.d.f.

Rank: The **rank** of  $x \in \mathbb{R}$  is  $F(x)$  (aka the **c.d.f.** at  $x$ )

Property:  $F(X) \sim \text{Uniform}([0, 1])$

Thus,  $F$  **transports** the distribution of  $X$  to  $U \sim \text{Uniform}([0, 1])$

In fact, if  $E[X^2] < \infty$ , c.d.f.  $F$  is the **optimal transport map** as

$$F = \arg \min_{T: T(X) \stackrel{d}{=} U} E |X - T(X)|^2$$

**Sample rank map** (aka empirical c.d.f.) is also a **transport map**:

$$\hat{R}_n := \arg \min_{T \in \mathcal{S}_n} \frac{1}{n} \sum_{i=1}^n |X_i - T(X_i)|^2 = \arg \min_T \frac{1}{n} \sum_{i=1}^n |X_i - T(X_i)|^2$$

where  $T$  **transports**  $\frac{1}{n} \sum_{i=1}^n \delta_{X_i}$  to  $\frac{1}{n} \sum_{i=1}^n \delta_{\frac{i}{n}}$

# Multivariate rank functions as transport maps

$X$  ; is a probability measure in  $\mathbb{R}^d$  (abs. cont.)

$U$  Uniform( $[0, 1]^d$ )

Goal: Find the "optimal" transport map  $T$  s.t.  $T(X) \stackrel{d}{=} U$

# Multivariate rank functions as transport maps

$X$  ; is a probability measure in  $\mathbb{R}^d$  (abs. cont.)

$U$  Uniform( $[0; 1]^d$ )

Goal: Find the "optimal" transport map  $T$  s.t.  $T(X) \stackrel{d}{=} U$

If  $E\|X\|^2 < \infty$ , the population rank function  $R(\cdot)$  is the transport map s.t.

$$R := \arg \min_{T: T(X) \stackrel{d}{=} U; X} E\|X - T(X)\|^2$$

# Multivariate rank functions as transport maps

$X$  ; is a probability measure on  $\mathbb{R}^d$  (abs. cont.)  
 $U$  Uniform( $[0, 1]^d$ )

Goal: Find the "optimal" transport map  $T$  s.t.  $T(X) \stackrel{d}{=} U$

If  $E\|X\|^2 < \infty$ , the "population rank function"  $R(\cdot)$  is the transport map s.t.

$$R := \arg \min_{T: T(X) \stackrel{d}{=} U; X} E\|X - T(X)\|^2$$

Data:  $X_1, \dots, X_n$  iid (abs. cont.) on  $\mathbb{R}^d$

$c_1, \dots, c_n$   $\mathbb{R}^d$  | grid of "reference" points

Sample multivariate rank map is defined as the transport map s.t.

$$\hat{R}_n = \arg \min_{T \in \mathcal{S}_n} \frac{1}{n} \sum_{i=1}^n \|X_i - T(X_i)\|^2 = \arg \min_T \frac{1}{n} \sum_{i=1}^n \|X_i - T(X_i)\|^2$$

where  $T$  transports  $\frac{1}{n} \sum_{i=1}^n \delta_{X_i}$  to  $\frac{1}{n} \sum_{i=1}^n \delta_{c_i}$

If  $\|k_X\|^2 < 1$ , the **population rank function**  $R(\cdot)$  is defined as

$$R := \arg \min_{T: T(X) \stackrel{d}{=} U; X} \|k_X - T(X)\|^2$$



If  $\|k_X\|^2 < 1$ , the **population rank function**  $R(\cdot)$  is defined as

$$R := \arg \min_{T: T(X) \stackrel{d}{=} U; X} \|k_X - T(X)\|^2$$

Even when  $\|k_X\|^2 = +\infty$ , **population rank function**  $R(\cdot)$  can also be defined [More details](#)

If  $E\|X\|^2 < 1$ , the **population rank function  $\mathbf{R}(\cdot)$**  is defined as

$$\mathbf{R} := \arg \min_{\mathbf{T}: \mathbf{T}(\mathbf{X}) \stackrel{d}{=} U; X} E\|X - \mathbf{T}(X)\|^2$$

Even when  $E\|X\|^2 = +\infty$ , **population rank function  $\mathbf{R}(\cdot)$**  can also be defined [More details](#)

**Sample multivariate rank map  $\hat{\mathbf{R}}_n(\cdot)$**  is defined as

$$\hat{\mathbf{R}}_n = \arg \min_{\mathbf{T}} \frac{1}{n} \sum_{i=1}^n \|X_i - \mathbf{T}(X_i)\|^2$$

where  $\mathbf{T}$  transports  $\frac{1}{n} \sum_{i=1}^n \delta_{X_i}$  to  $\frac{1}{n} \sum_{i=1}^n \delta_{c_i}$

If  $E\|X\|^2 < \infty$ , the **population rank function**  $\mathbf{R}(\cdot)$  is defined as

$$\mathbf{R} := \arg \min_{\mathbf{T}: \mathbf{T}(\mathbf{X}) \stackrel{d}{=} U; X} E\|X - \mathbf{T}(\mathbf{X})\|^2$$

Even when  $E\|X\|^2 = \infty$ , **population rank function**  $\mathbf{R}(\cdot)$  can also be defined [More details](#)

**Sample multivariate rank map**  $\hat{\mathbf{R}}_n(\cdot)$  is defined as

$$\hat{\mathbf{R}}_n = \arg \min_{\mathbf{T}} \frac{1}{n} \sum_{i=1}^n \|X_i - \mathbf{T}(X_i)\|^2$$

where  $\mathbf{T}$  transports  $\frac{1}{n} \sum_{i=1}^n \delta_{X_i}$  to  $\frac{1}{n} \sum_{i=1}^n \delta_{c_i}$

**Regularity:  $L_2$ -convergence [D. and Sen, JASA 2020]**

$X_1, \dots, X_n$  iid (abs. cont.). If  $\frac{1}{n} \sum_{i=1}^n \delta_{c_i} \xrightarrow{w} \text{Unif}([0, 1]^d)$ , then

$$\frac{1}{n} \sum_{i=1}^n \|\hat{\mathbf{R}}_n(X_i) - \mathbf{R}(X_i)\|^2 \xrightarrow{a.s.} 0 \quad \text{as } n \rightarrow \infty$$

Result gives the required **regularity** of the **empirical multivariate rank map**

# Population version

Assume  $m=(m+n) = 2(0;1)$ .

Rank energy distance [D. and Sen, JASA 2020]

- **Joint rank map:** The “pooled” population rank map:

$$R : R(\mathbf{Z}) \sim \text{Uniform}([0;1]^d)$$

where  $\mathbf{Z} = P + (1 - )Q$ .

# Population version

Assume  $m=(m+n) = 2(0;1)$ .

Rank energy distance [D. and Sen, JASA 2020]

- **Joint rank map:** The “pooled” population rank map:

$$R : \mathbb{R}^d \rightarrow \text{Uniform}([0;1]^d)$$

where  $\mathbf{Z} = P + (1 - \alpha)Q$ .

- **Rank energy:**  $RE^2(P; Q) := E^2(R(\mathbf{X}); R(\mathbf{Y}))$ .
- $RE = 0$  iff  $P = Q$  provided  $P, Q$  are absolutely continuous.

# Population version

Assume  $m=(m+n) = 2(0;1)$ .

Rank energy distance [D. and Sen, JASA 2020]

- **Joint rank map:** The “pooled” population rank map:

$$R : R(\mathbf{Z}) \sim \text{Uniform}([0;1]^d)$$

where  $\mathbf{Z} = P + (1 - )Q$ .

- **Rank energy:**  $RE^2(P; Q) := E^2(R(X); R(Y))$ .
- $RE = 0$  iff  $P = Q$  provided  $P, Q$  are absolutely continuous.
- Our **general principle** could have been used with any other procedure for testing equality of distributions, e.g., the **MMD** statistic [Gretton et al. (2008)] which uses ideas from RKHS, ...

# Population version

Assume  $m=(m+n) = 2(0;1)$ .

Rank energy distance [D. and Sen, JASA 2020]

- **Joint rank map:** The “pooled” population rank map:

$$R : R(\mathbf{Z}) \sim \text{Uniform}([0;1]^d)$$

where  $\mathbf{Z} = P + (1 - )Q$ .

- **Rank energy:**  $RE^2(P; Q) := E^2(R(X); R(Y))$ .
- $RE = 0$  iff  $P = Q$  provided  $P, Q$  are absolutely continuous.
- Our **general principle** could have been used with any other procedure for testing equality of distributions, e.g., the **MMD** statistic [Gretton et al. (2008)] which uses ideas from RKHS, ...
- For  $d = 1$ , we prove that  $RE_{m,n}^2$  and  $RE^2$  are exactly equivalent to the sample and population two-sample Cramér-von Mises statistic.

# Pitman efficiency

Consider  $X_1, \dots, X_n \sim P_1$  and  $Y_1, \dots, Y_m \sim P_2$ , with  $m, n \rightarrow \infty$  and  $m/(m+n) \rightarrow \lambda \in (0, 1)$ . We want to test:

$$H_0: \theta_1 = \theta_2 \quad \text{versus} \quad H_1: \theta_1 = h(m+n)^{-1/2} \theta_2;$$



# Pitman efficiency

Consider  $X_1, \dots, X_n \sim P_1$  and  $Y_1, \dots, Y_m \sim P_2$ , with  $m = m + n = 2(0;1)$ . We want to test:

$$H_0: \theta_1 = 0 \quad \text{versus} \quad H_1: \theta_1 = h(m+n)^{-1/2};$$

Fix  $\alpha$  (size) and  $1 - \beta$  (power).

Two test functions  $T_{m;n}$  and  $S_{m;n}$ .

# Pitman efficiency

Consider  $X_1, \dots, X_n \sim P_1$  and  $Y_1, \dots, Y_m \sim P_2$ , with  $m = m + n = 2(0;1)$ . We want to test:

$$H_0: \theta_1 = 0 \quad \text{versus} \quad H_1: \theta_1 = h(m+n)^{-1/2};$$

Fix  $\alpha$  (size) and  $1 - \beta$  (power).

Two test functions  $T_{m;n}$  and  $S_{m;n}$ .

$K(T_{m;n})$  denotes minimum number of samples such that:

$$E_{H_0}(T_{m;n}) \leq K \quad \text{and} \quad E_{H_1}(T_{m;n}) \leq K \cdot \frac{1 - \beta}{\alpha};$$

# Pitman efficiency

Consider  $X_1, \dots, X_n \sim P_1$  and  $Y_1, \dots, Y_m \sim P_2$ , with  $m = m + n = 2(0;1)$ . We want to test:

$$H_0: \theta_1 = 0 \quad \text{versus} \quad H_1: \theta_1 = h(m+n)^{1/2};$$

Fix  $\alpha$  (size) and  $1 - \beta$  (power).

Two test functions  $T_{m;n}$  and  $S_{m;n}$ .

$K(T_{m;n})$  denotes minimum number of samples such that:

$$E_{H_0}(T_{m;n}) \leq \alpha \quad \text{and} \quad E_{H_1}(T_{m;n}) \geq 1 - \beta;$$

The Pitman efficiency of  $S_{m;n}$  with respect to  $T_{m;n}$  is given by

$$\lim_{m+n \rightarrow \infty} \frac{K(T_{m;n})}{K(S_{m;n})}.$$