# Statistical inference with optimal transport

WI AS

Leibniz
Association

*Vladimir Spokoiny*

# The team and possible cooperation

*Stochastic*: Vladimir Spokoiny, Andrey Sobolevskiy, Alexey Kroshnin, Alexandra Suvorikova



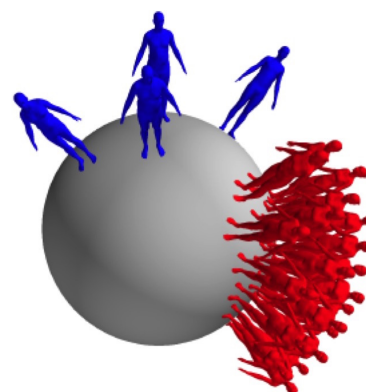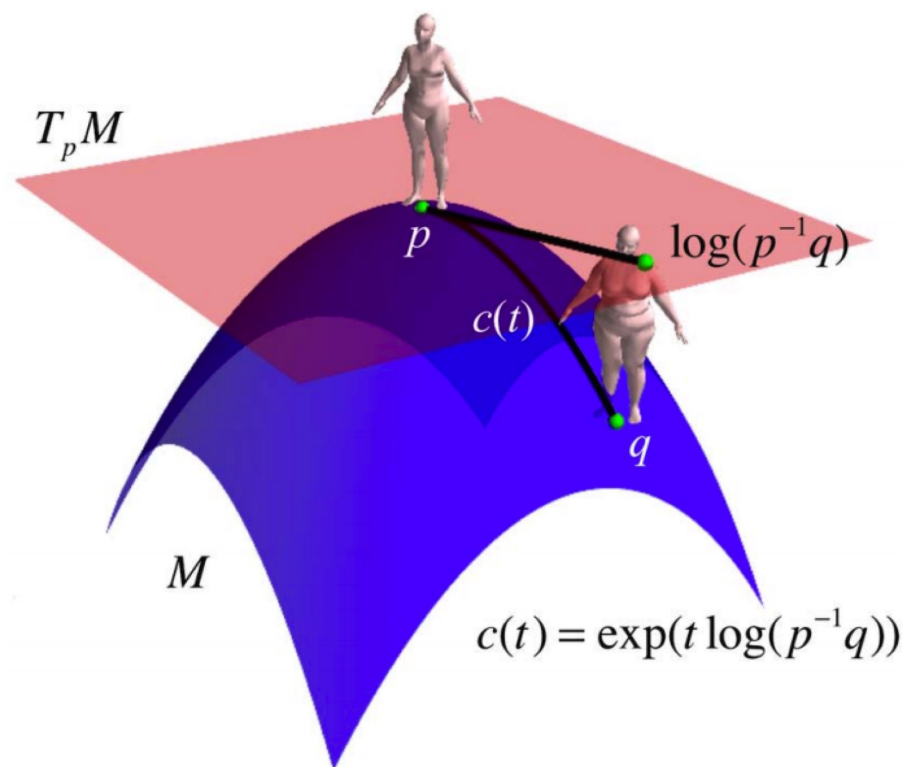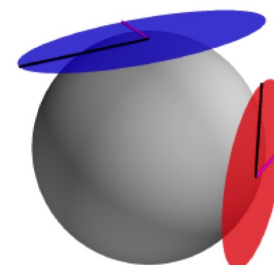*Optimization*: Alexander Gasnikov, Pavel Dvurechensky



*Cooperation*:

# Analysis of data on manifolds
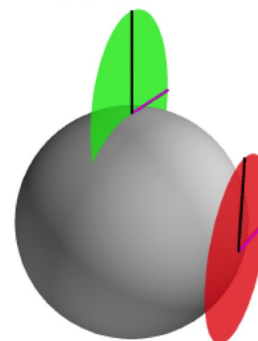
Modern data lives in manifolds:

- underlying geometry
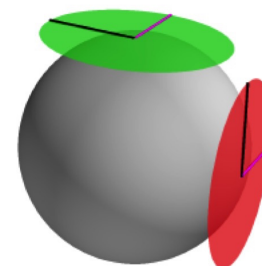- non-linearity of the space



$$c(t) = \exp(t \log(p^{-1}q))$$

(a) Data on a manifold

(b) Data covariances

(c) Linear translation

(d) Covariance Transport

Source: ps.is.tuebingen.mpg.de

# I. Motivation

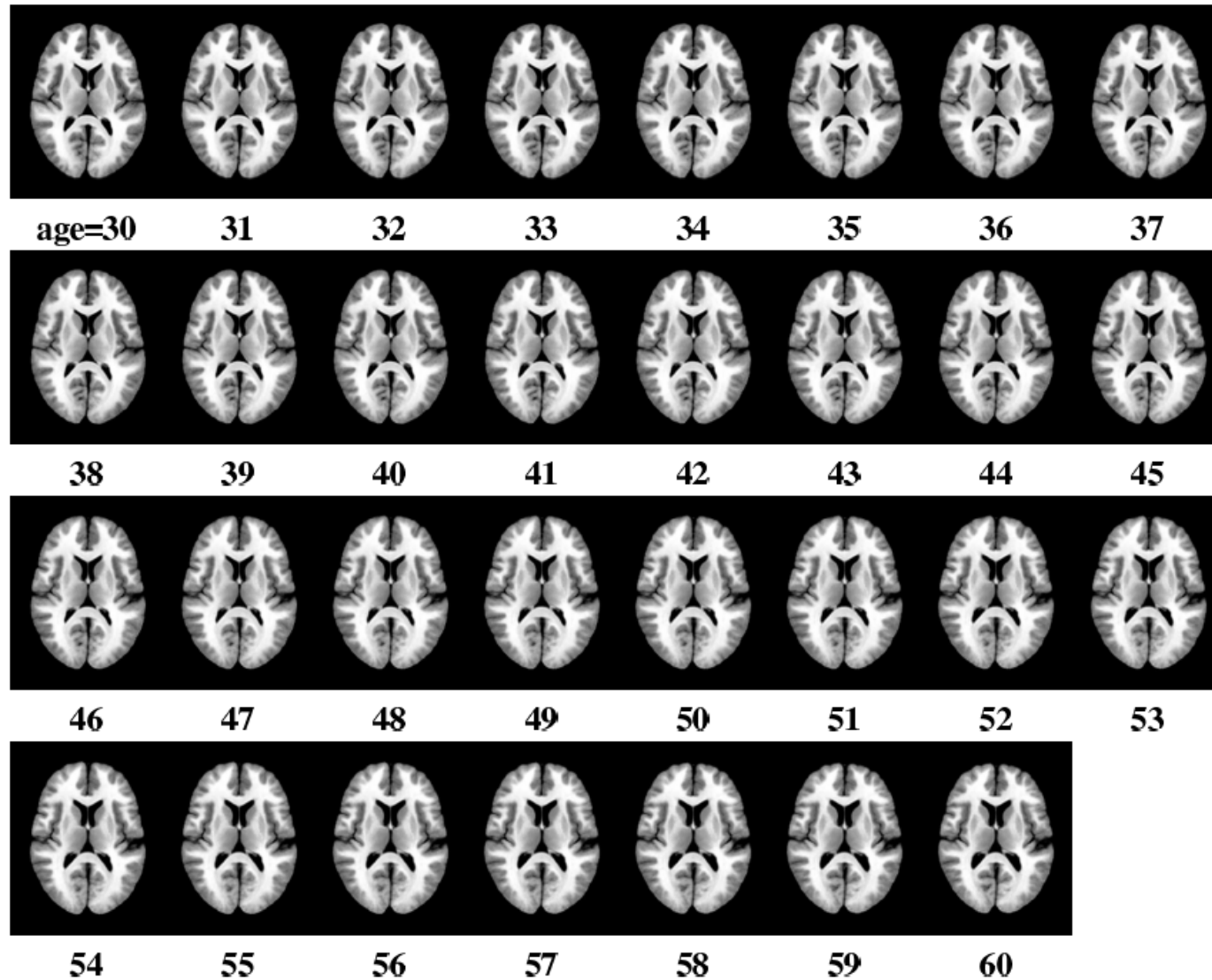Stat Inference with OT

# Example 1: pattern extraction by averaging



Healthy brains, m + f, age of 21 to 72

age=30 31 32 33 34 35 36 37

38 39 40 41 42 43 44 45

46 47 48 49 50 51 52 53

54 55 56 57 58 59 60

Typical ageing pattern, f

# Example 1: pattern extraction by averaging



RECOVERY OF THE TEMPLATE OBJECT

TEMPLATE

OBSERVED SAMPLE

RECOVERED OBJECT

Given $(\mathcal{X}, d)$ and $I\!\!P$ on $\mathcal{X}$ and iid $Y_1, ..., Y_n$, $Y_i \backsim I\!\!P$ :
Template:

$$X^* \overset{\text{def}}{=} \underset{X \in \mathcal{X}}{\operatorname{argmin}} \int_{\mathcal{X}} d^2(X, Y) I\!\!P(dY)$$

Recovered object:

$$X_n \overset{\text{def}}{=} \underset{X \in \mathcal{X}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} d^2(X, Y_i).$$

- how to chose $d$
- confidence sets around $X_n$

# Example 2: stem cell differentiation



Mesenchymal stem cells, $\mu_0$



Chondrogenesis, $t \in [0,1]$

Detect time $t$ when a cell specifies its "type".

Stat Inference with OT

# Example 2: stem cell differentiation



$$\begin{cases} H_0 : \text{ data is } \textcolor{red}{\text{homogeneous}} \iff t \text{ is not a change point} \\ H_1 : \text{ data is } \textcolor{red}{\text{non-homogeneous}} \iff t \text{ is a change point} \end{cases}$$

1. compute some cumulative statistics, e.g. means $\mu_l(t)$ and $\mu_r(t)$
2. compare them, e.g. $\text{dist}\big(\mu_l(t), \mu_r(t)\big)$
3. if $\text{dist}\big(\mu_l(t), \mu_r(t)\big) \geq \mathfrak{z}_\alpha(t)$, then $H_0$ is rejected

Goal: non-asymptotic data driven rejection level $\mathfrak{z}_\alpha$

# Common features

All above mentioned examples have a common problem: data sets possess inner geometry

$\mathcal{Q}$ :  what is a good way to define a distance between such objects?

# Common features

All above mentioned examples have a common problem: data sets possess inner geometry

$\mathcal{A}$ : $2$-Wasserstein distance might be a solution



Image source: Marco Cuturi's OT tutorial

# II. Introduction to OT

# What is Optimal Transport?

The natural geometry for probability measures

Monge    Kantorovich   Koopmans    Dantzig    Brenier    Otto    McCann    Villani

Nobel '75                      Fields '10

Image source: Marco Cuturi's OT tutorial

AVERAGING OF IMAGES

EUCLIDEAN MEAN

2-WASSERSTEIN MEAN

# Wasserstein distance and optimal transport

- $W$-*distance* $-$ minimum amount of work that is necessary to convert $\mu$ into $\nu$
- $D(x, y)$ $-$ cost of transportation of unit mass from $x$ to $y$

# Space of images $(\mathcal{P}(I\!\!R^d), W_2)$

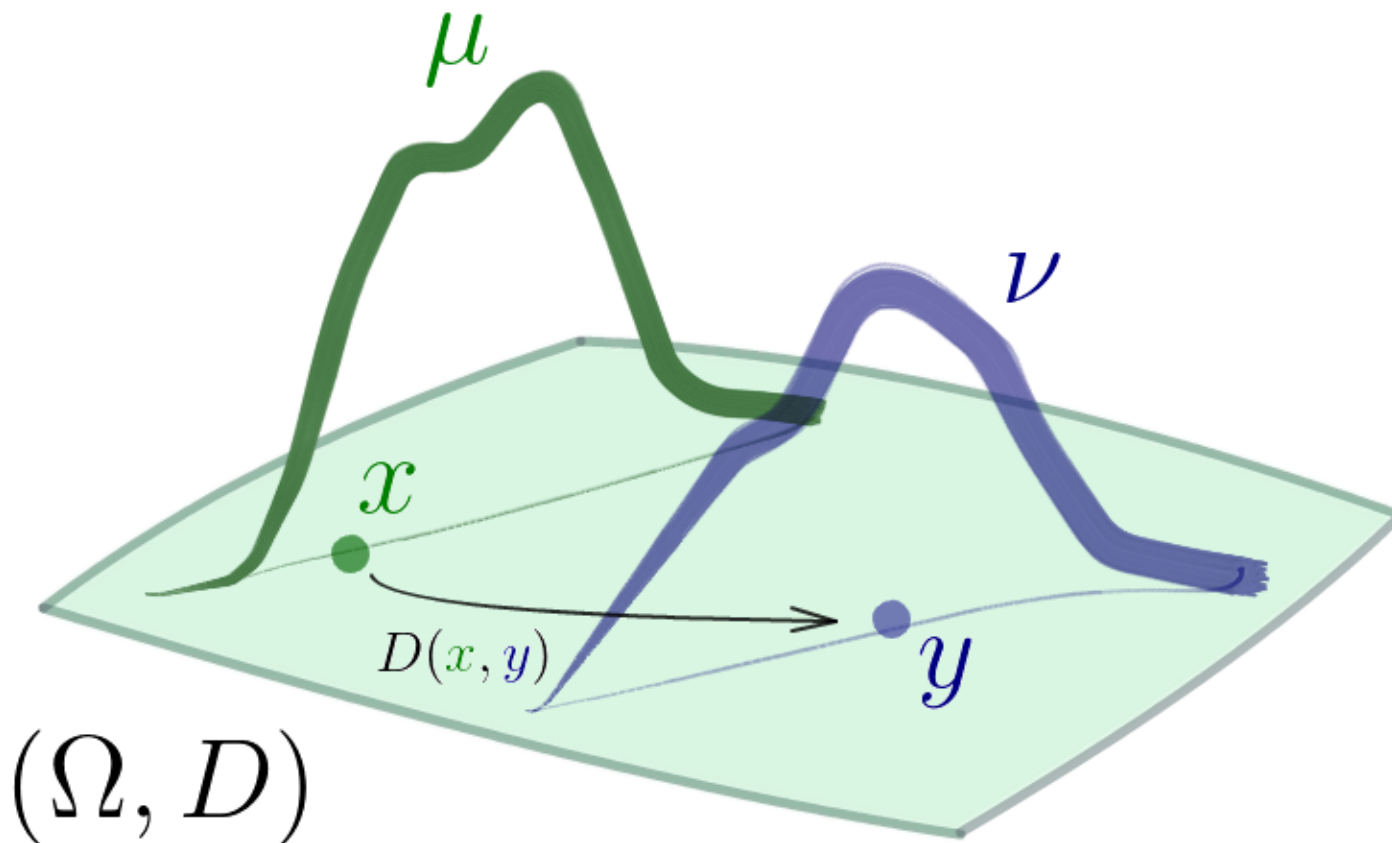$\mathcal{P}_2(I\!\!R^d) = \{$Probability measures $\nu$ on $I\!\!R^d\}$,
metrizied by $2$-W distance:

$$W_2^2(\mu, \nu) \overset{\text{def}}{=} \inf_{\pi \in \Pi(\mu, \nu)} \int_{I\!\!R^d \times I\!\!R^d} \|x - y\|^2 d\pi(x, y)$$

$\Pi(\mu, \nu)$ – set of all prob. measures $\pi$ on $I\!\!R^d \times I\!\!R^d$ with marginals $\mu$ and $\nu$

$$\begin{cases} \int_{\Omega \times \Omega} \pi(x, y)dy = \mu(x) \\ \int_{\Omega \times \Omega} \pi(x, y)dx = \nu(y) \end{cases}$$

# Recommended literature

- Villani, C. (2008). Optimal transport: old and new (Vol. 338). Springer Science & Business Media.

- Santambrogio, F. (2015). Optimal transport for applied mathematicians. Birkaeuser, NY, 99-102.

- Ambrosio, L. (2003). Lecture notes on optimal transport problems. In Mathematical aspects of evolving interfaces (pp. 1-52). Springer Berlin Heidelberg.

and many more...

# III. Some obtained results



Stat Inference with OT

# III.a Non-asymptotic confidence sets

Stat Inference with OT

$$\mathcal{P}_2(I\!\!R^d), W_2 \text{ -- metric space and } I\!\!P \in \mathcal{I\!\!P}\big(\mathcal{P}_2(I\!\!R^d)\big)$$

RECOVERY OF THE TEMPLATE OBJECT



TEMPLATE

OBSERVED SAMPLE

RECOVERED OBJECT

Wasserstein *population* barycenter

$$\mu^* \subseteq \operatorname*{argmin}_{\mu \in \mathcal{P}_2(I\!\!R^d)} \int_{\mathcal{P}_2(I\!\!R^d)} W_2^2(\mu, \nu) I\!\!P(d\nu)$$

$\nu_1, ..., \nu_n$ – observed random iid sample, $\nu_i \backsim I\!\!P$
Wasserstein *empirical* barycenter

$$\mu_n \subseteq \operatorname*{argmin}_{\mu \in \mathcal{P}_2(I\!\!R^d)} \frac{1}{n} \sum_{i=1}^{n} W_2^2(\mu, \nu_i).$$

# Confidence set around $\mu_n$

### Real world

$$T_n \overset{\text{def}}{=} \sqrt{n}W_2(\mu^*, \mu_n)$$

we <u>do not</u> know $\mathfrak{z}(\alpha)$ :

$$\mathfrak{z}(\alpha) \overset{\text{def}}{=} \underset{\mathfrak{z}>0}{\text{argmin}}\big\{ I\!\!P\big(T_n \geq \mathfrak{z}\big) \leq \alpha \big\}$$

### Bootstrap world (mimics $T_n$)

$$T_n^\flat \overset{\text{def}}{=} \sqrt{n}W_2(\mu_n, \mu_n^\flat)$$

we know $\mathfrak{z}^\flat(\alpha)$ :

$$\mathfrak{z}^\flat(\alpha) \overset{\text{def}}{=} \underset{\mathfrak{z}>0}{\text{argmin}}\big\{ I\!\!P^\flat(T_n^\flat \geq \mathfrak{z}) \leq \alpha \big\}$$

Goal: replace $\mathfrak{z}(\alpha)$ with $\mathfrak{z}^\flat(\alpha)$

Result: Under some technical assumptions it holds with h.p. $I\!\!P$, $I\!\!P^\flat$

$$\big| I\!\!P\big(T_n \geq \mathfrak{z}^\flat(\alpha)\big) - \alpha \big| \leq \mathtt{C}/\sqrt{n}.$$

# Multiplier bootstrap

## Real world

Observed sample $\nu_i \overset{\text{iid}}{\curvearrowright} \mathbb{P}$

$W$ -*population* barycenter:

$$\mu^* = \operatorname*{argmin}_{\mu} \int W_2^2(\mu, \nu) \mathbb{P}(d\nu),$$

$W$ -*empirical* barycenter:

$$\mu_n = \operatorname*{argmin}_{\mu} \frac{1}{n} \sum W_2^2(\mu, \nu_i),$$

$$T_n = \sqrt{n} W_2(\mu^*, \mu_n)$$

## Bootstrap world, $u_i \curvearrowright \text{Po}(1)$

Training sample $\nu_i' \overset{\text{iid}}{\curvearrowright} \mathbb{P}$

$W^\flat$ -*population* barycenter:

$$\mu_n = \operatorname*{argmin}_{\mu} \frac{1}{n} \sum W_2^2(\mu, \nu_i'),$$

$W^\flat$ -*empirical* barycenter:

$$\mu_n^\flat = \operatorname*{argmin}_{\mu} \frac{1}{n} \sum W_2^2(\mu, \nu_i') u_i,$$

$$T_n^\flat = \sqrt{n} W_2(\mu_n, \mu_n^\flat)$$

# III.b Non-parametric 2-sample test

# Two-sample testing

*T-statistics:* William S. Gosset (1908), economical way to monitor the quality of stout for Guinness

*Two-sample test*

$$X = (X_1, ..., X_n), \; X_i \overset{iid}{\frown} \mu_X,$$

$$Y = (Y_1, ..., Y_m), \; Y_j \overset{iid}{\frown} \mu_Y$$

*Goal*

$$H_0 : \; \mu_X = \mu_Y, \quad H_1 : \; \mu_X \neq \mu_Y$$



A 1904 brand poster (Guinness)

# Underlying idea

Find <u>some</u> transformation $T$ of $\mu_X$, $\mu_Y$ :

- measure-preserving,
- yields universal critical level $\mathfrak{z}_{nm}$, which does not depend on $\mu_X$, $\mu_Y$.

Test:

$$\mathsf{dist}\big(T \# \mu_X^n, T \# \mu_Y^m\big) \geq \mathfrak{z}_{nm} \;\longrightarrow\; H_1$$

1-D case: $T$ is a ranking or quantile map (OT to $\mathcal{U}[0,1]$ )

d-D case: $T$ is a Monge map (optimal transport with $c(x,y) = \|x - y\|^2$ )

$$T := \underset{T'_\# \mu = \nu}{\operatorname{argmin}} \int_{\mathsf{supp}(\mu)} \big\|T'(x) - x\big\|^2 \mu(dx),$$
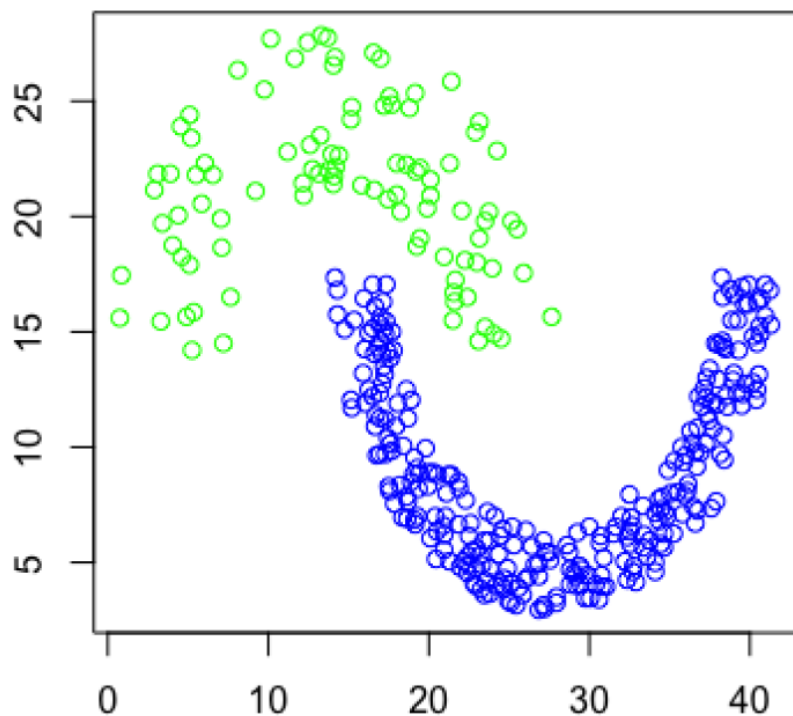
where $\nu$ some fixed reference measure and
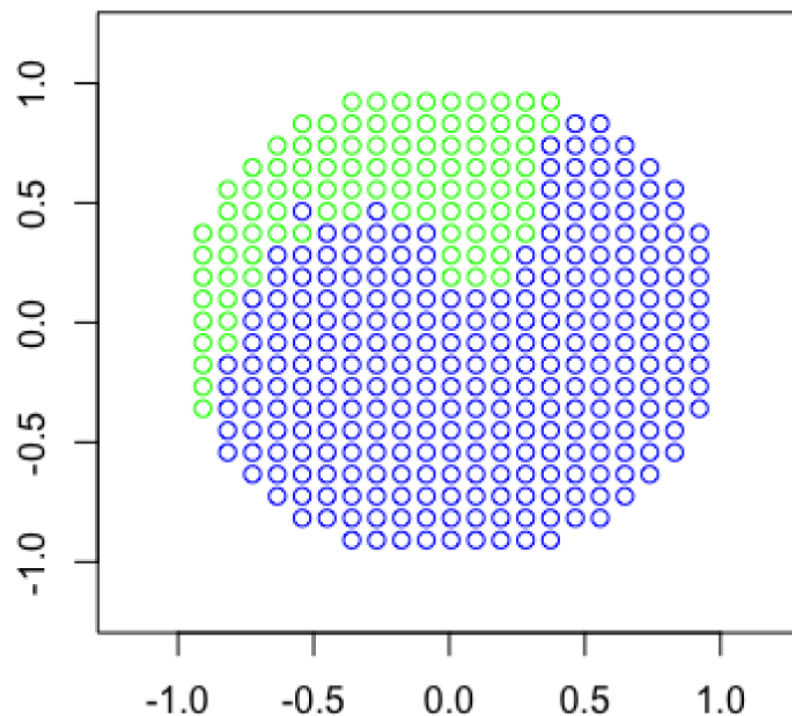
$$\mu := \alpha \mu_X + (1 - \alpha)\mu_Y$$

# Intuition: a kind of $d$-dim quantile map

*Observed distribution:* $\mu = \alpha\mu_X + (1-\alpha)\mu_Y$

*Reference distribution:* $\nu = \mathcal{U}[B_0(1)]\,,\ B_0(1) \subset I\!\!R^d$



Original data

Transformed data

# Idea of the method

*Observed distribution:* $\mu = \alpha \mu_X + (1 - \alpha)\mu_Y$

*Reference distribution:* $\nu = \mathcal{U}[B_0(1)], \; B_0(1) \subset I\!\!R^d$

*Monge map*:

$$T_\# \mu = \nu, \quad \nu = \alpha \nu_X + (1 - \alpha)\nu_Y,$$

$$\nu_X\big(T(B)\big) = \mu_X(B), \quad \nu_Y\big(T(B)\big) = \mu_Y(B), \quad \forall B \in \mathcal{B}.$$

Under $H_0$ :

$$\nu_X = \nu_Y = \mathcal{U}[B_0(1)]$$

Test:

$$D_{nm}^T \stackrel{\text{def}}{=} W_2(\nu_X^n, \nu_Y^m) \geq \mathfrak{z}_{nm} \;\Rightarrow\; H_1$$

# IV. Open problems

Stat Inference with OT

# What we are currently doing

The list of open problems and related literature: http://strlearn.ru/topics/

## Hypothesis testing with Hellinger–Kantorovich distance
Responsible persons: Alexanda Suvorikova, Pavel Dvurechensky, Alexey Kroshnin, Andrey Sobolevskii, Vladimir Spokoiny

## Domain adaptation using optimal transportation
Responsible persons: Alexanda Suvorikova, Pavel Dvurechensky, Alexey Kroshnin, Andrey Sobolevskii, Vladimir Spokoiny

## Bootstrap for empirical barycenters
Responsible persons: Alexanda Suvorikova, Alexey Kroshnin, Andrey Sobolevskii, Vladimir Spokoiny

## Two sample test for high dimensional data using Monge–Kantorovich transform
Responsible persons: Alexanda Suvorikova, Alexey Kroshnin, Andrey Sobolevskii, Vladimir Spokoiny

**References:**
[SAN15] Santambrogio F. Optimal transport for applied mathematicians. Birkäuser, NY, 2015.
[VIL08] Villani, C. Optimal transport: old and new. Springer Science and Business Media, 2008.

# Thank you for your attention!

Stat Inference with OT