

Imputation with data depth

Pavlo Mozharovskyi^{1,2,3}

joint work with

Julie Josse^{2,4} and François Husson^{2,3}

¹ Centre Henri Lebesgue, Rennes

² Agrocampus Ouest, Rennes

³ Institut de Recherche Mathématique de Rennes (IRMAR)

⁴ Institut National de Recherche en Informatique et en Automatique
(INRIA), Orsay

48ème Journées de Statistique de la SFdS
Montpellier, May 31, 2016

Contents

Motivation

Depth and depth lift

Proposal

Experiments

Conclusions and outlook

Contents

Motivation

Depth and depth lift

Proposal

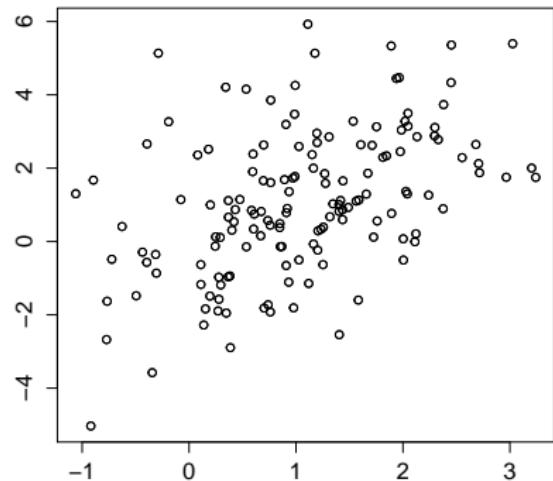
Experiments

Conclusions and outlook

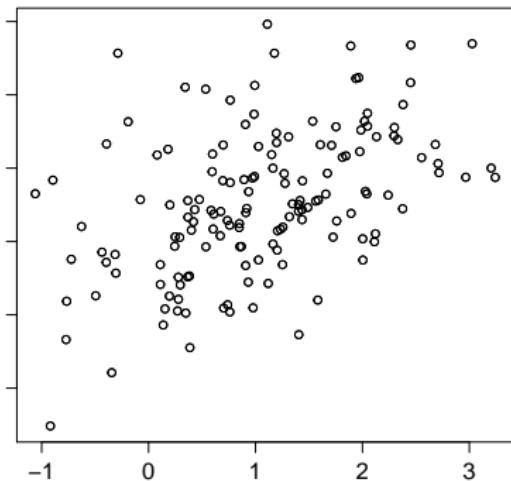
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

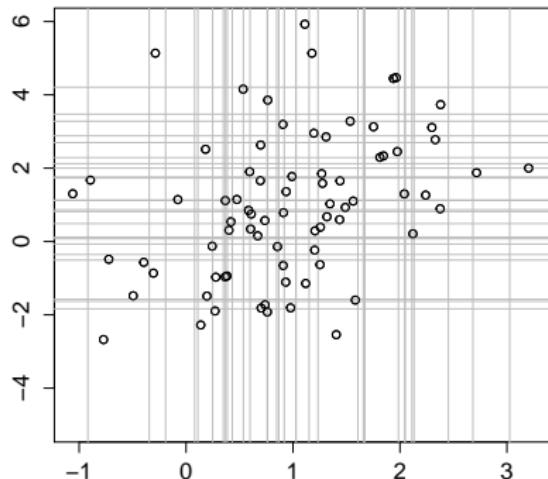
● – zonoid depth

▲ – random forest

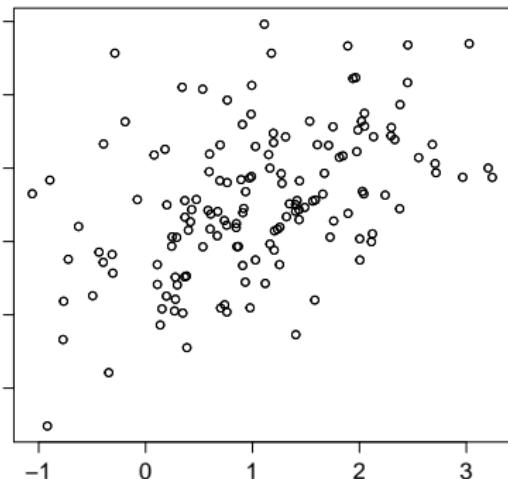
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

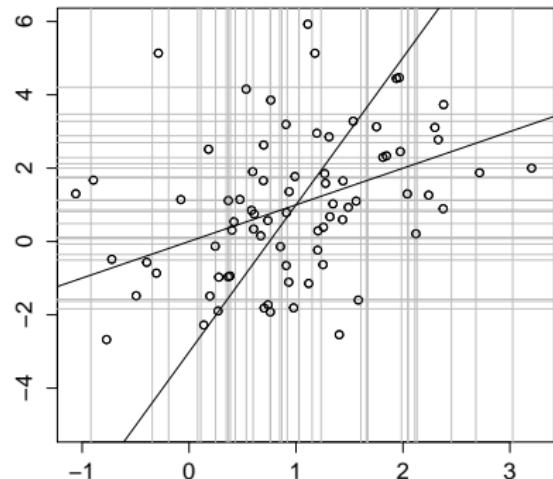
● – zonoid depth

▲ – random forest

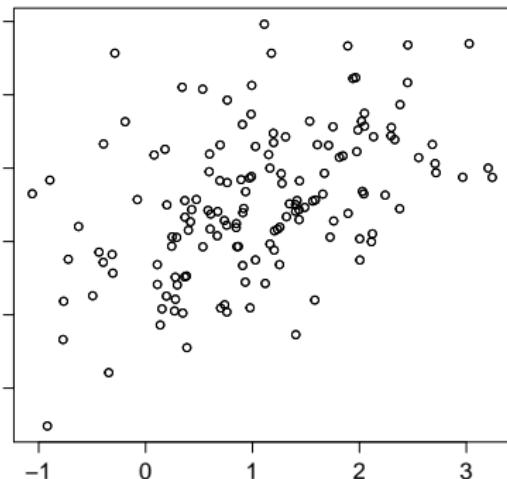
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

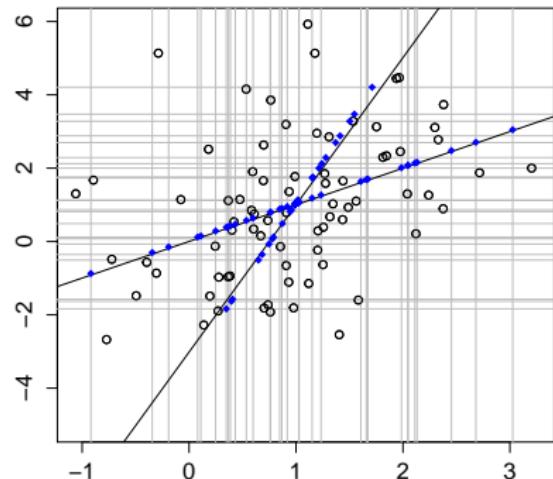
● – zonoid depth

▲ – random forest

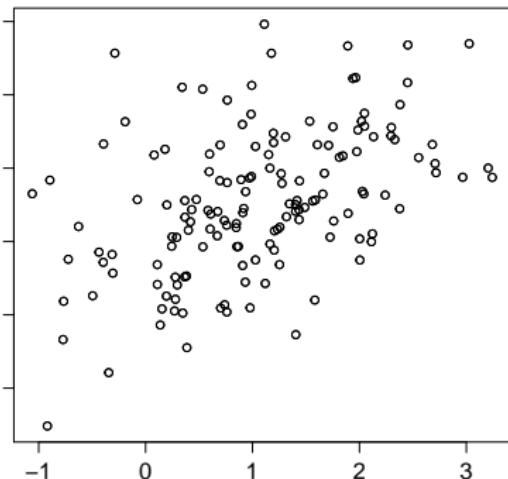
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

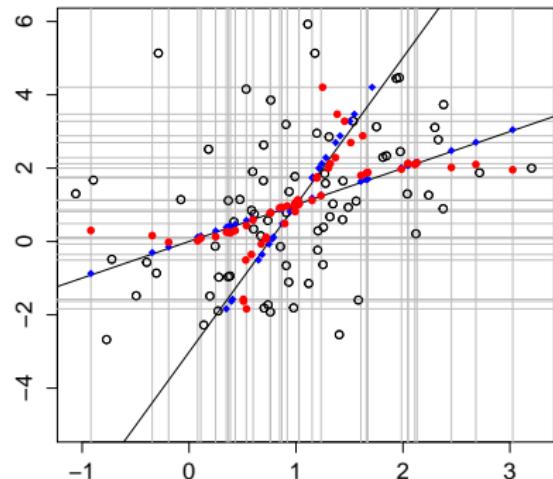
● – zonoid depth

▲ – random forest

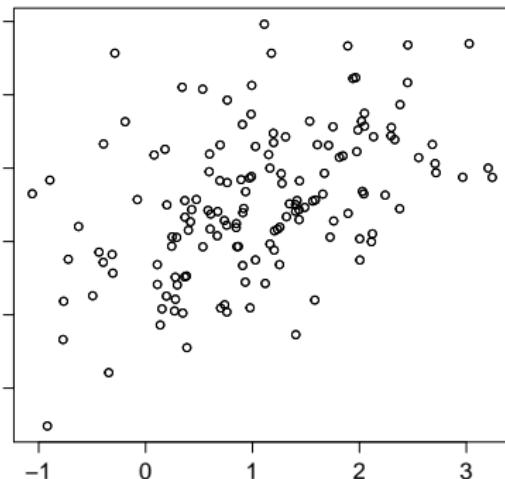
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

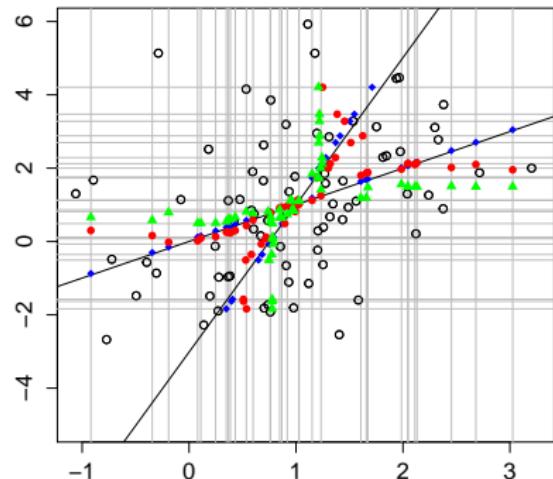
● – zonoid depth

▲ – random forest

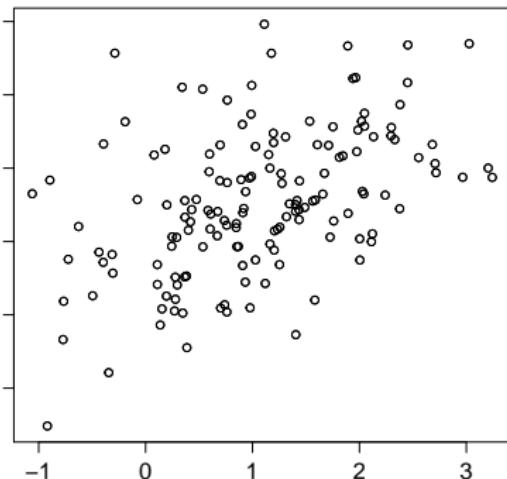
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

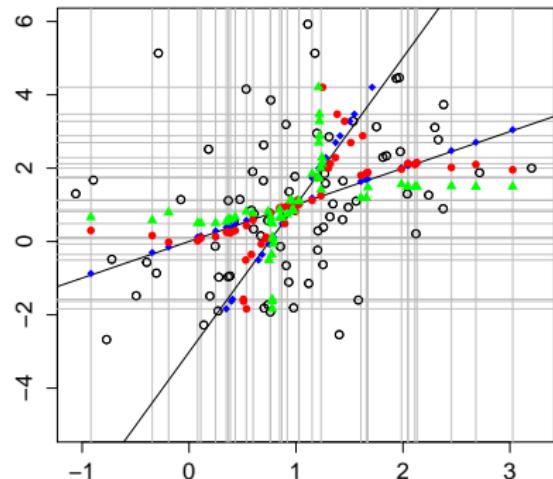
● – zonoid depth

▲ – random forest

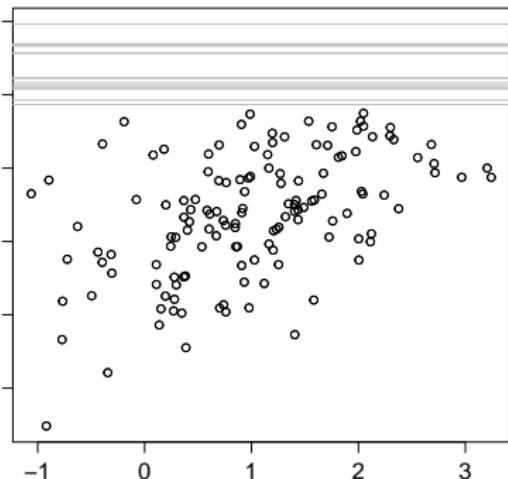
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

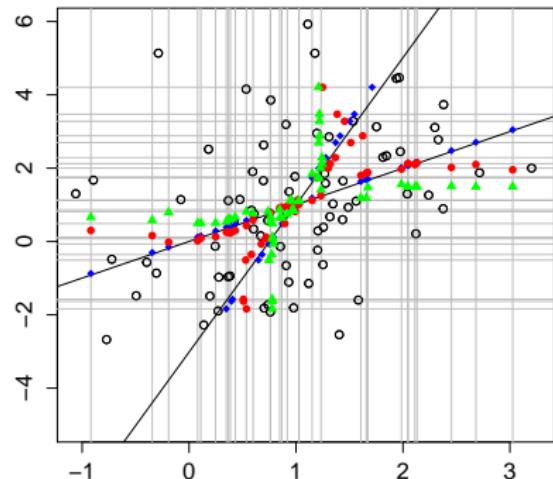
● – zonoid depth

▲ – random forest

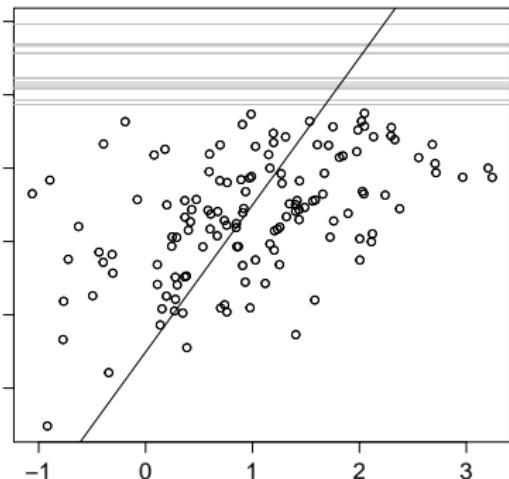
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

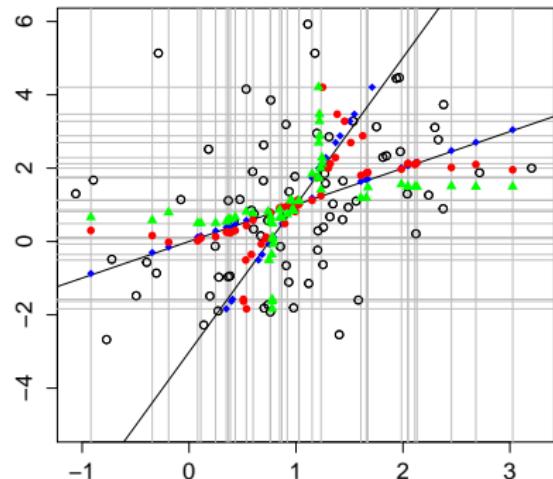
● – zonoid depth

▲ – random forest

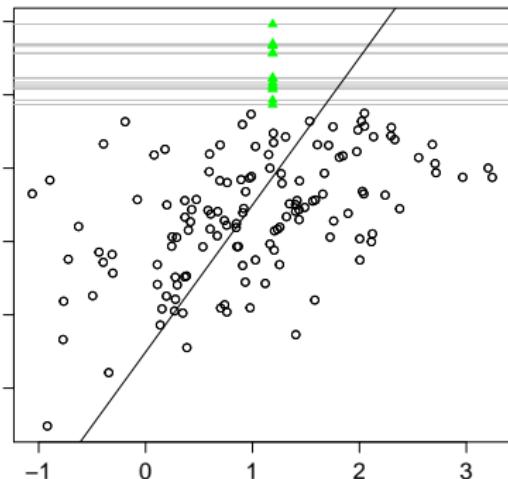
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

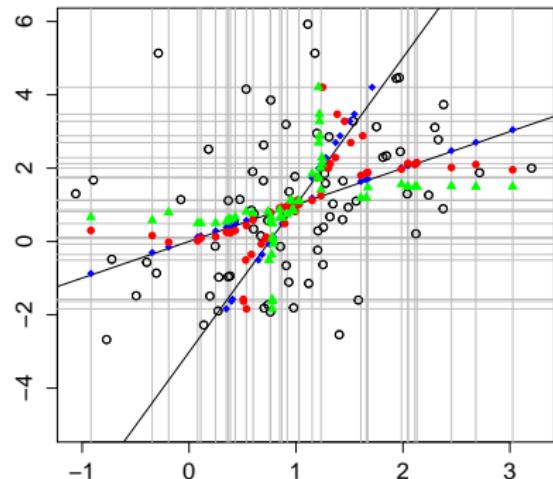
● – zonoid depth

▲ – random forest

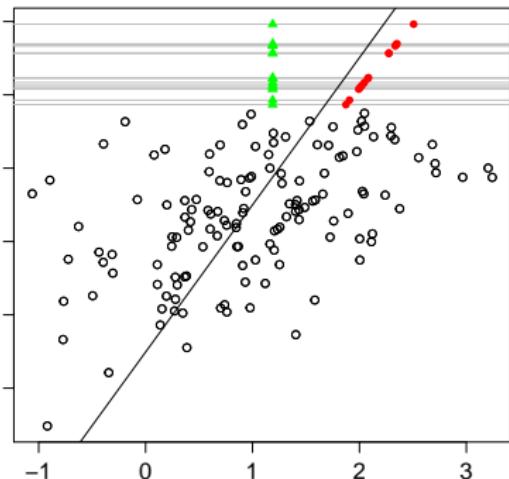
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

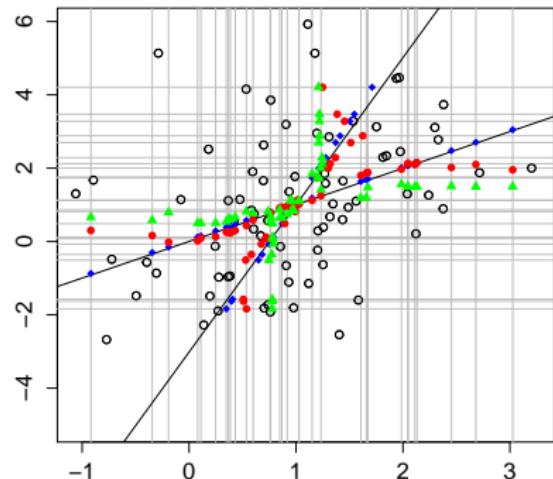
● – zonoid depth

▲ – random forest

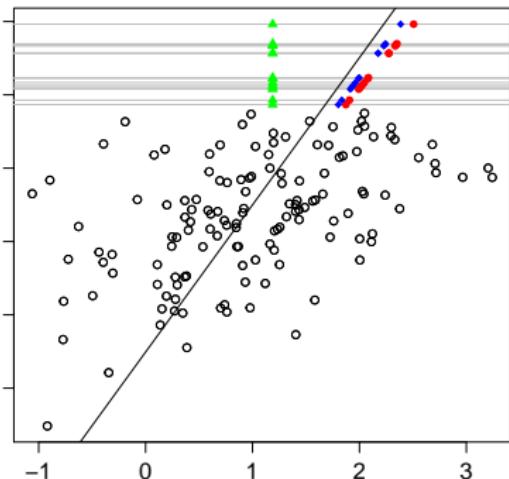
Imputation close to data and extrapolation

150 observations of $X \sim Normal\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Missing completely at random



Missing (not) at random



◆ – regression

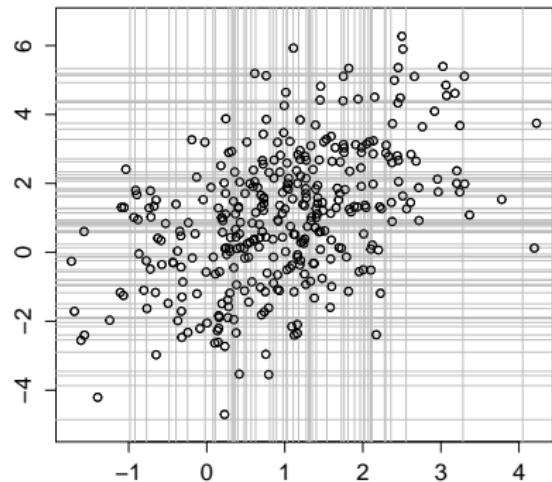
● – zonoid depth

▲ – random forest

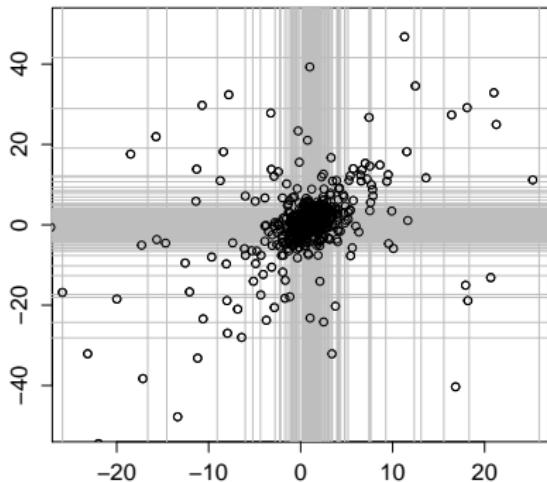
Robust imputation

15% (left) and 100% (right) from $\text{Cauchy}\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Robust to outliers (MCAR)



Robust to distribution (MCAR)



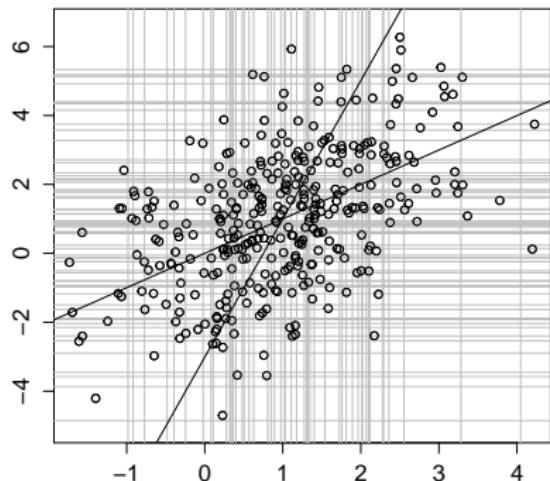
◆ – regression

● – Tukey depth

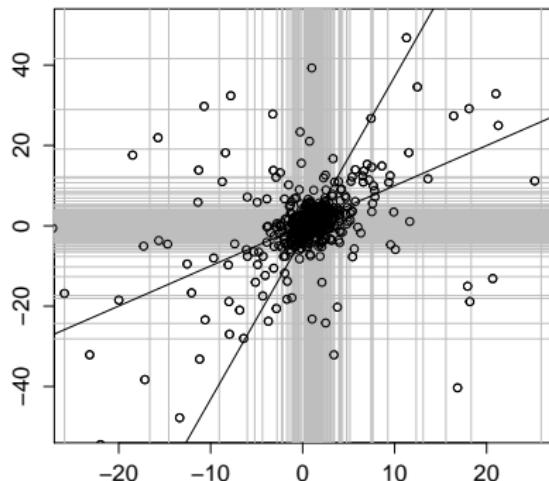
Robust imputation

15% (left) and 100% (right) from $\text{Cauchy}\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Robust to outliers (MCAR)



Robust to distribution (MCAR)



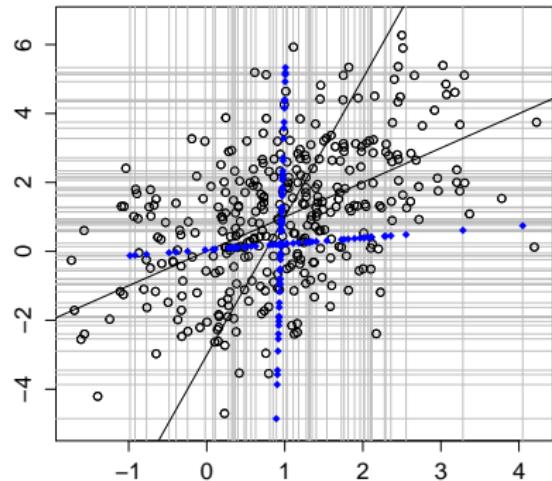
◆ – regression

● – Tukey depth

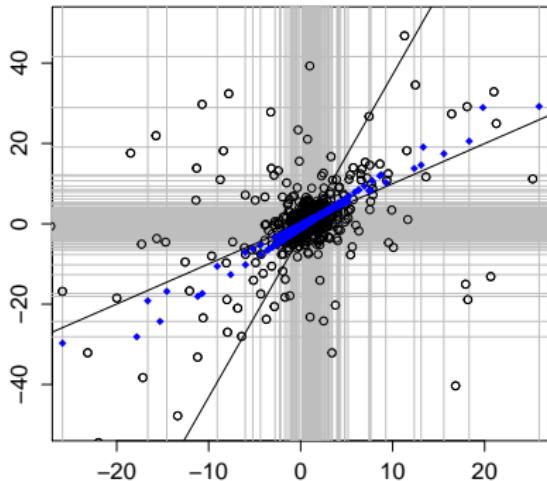
Robust imputation

15% (left) and 100% (right) from $\text{Cauchy}\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Robust to outliers (MCAR)



Robust to distribution (MCAR)



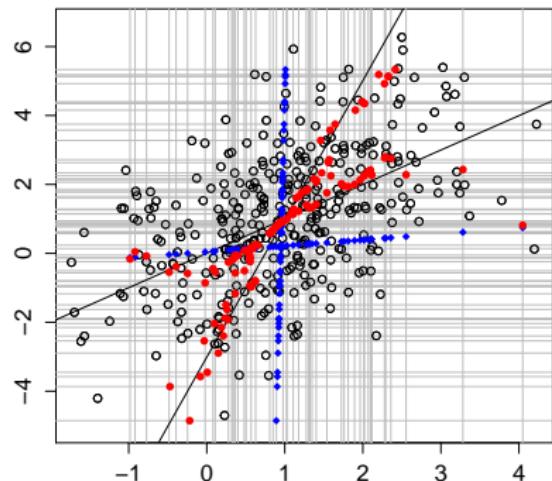
◆ – regression

● – Tukey depth

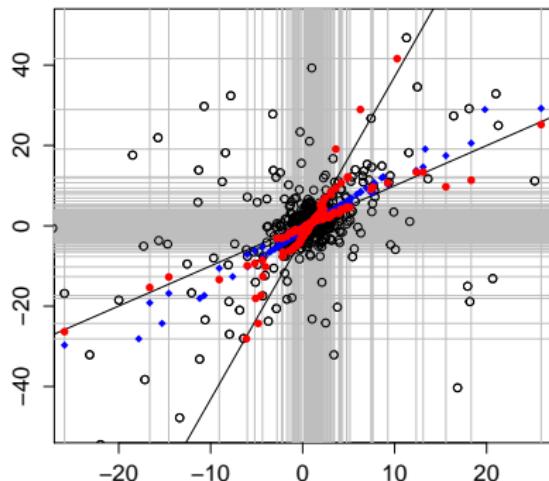
Robust imputation

15% (left) and 100% (right) from $\text{Cauchy}\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 1 \\ 1 & 4 \end{pmatrix}\right)$

Robust to outliers (MCAR)



Robust to distribution (MCAR)



◆ – regression

● – Tukey depth

Contents

Motivation

Depth and depth lift

Proposal

Experiments

Conclusions and outlook

Statistical data depth

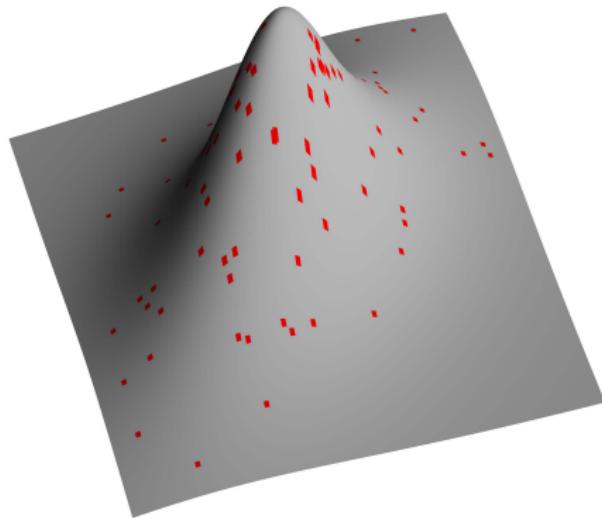
A **data depth** measures, how “close” a given point is located to the “center” of a distribution. For $\mathbf{x} \in \mathbb{R}^d$ and a d -variate random vector X distributed as $P \in \mathcal{M}$, a data depth is a function

$$D : \mathbb{R}^d \times \mathcal{M} \rightarrow [0, 1], (\mathbf{x}, P) \mapsto D(\mathbf{x}|P)$$

that is:

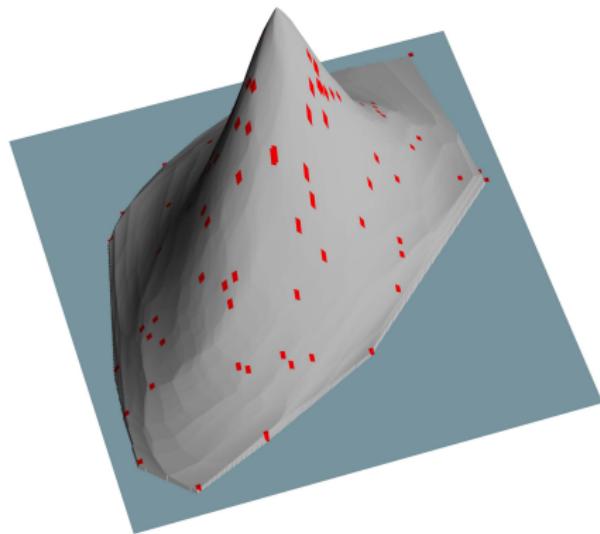
- ▶ **affine invariant**: $D(A\mathbf{x} + b|AX + b) = D(\mathbf{x}|X)$;
- ▶ **vanishing at infinity**: $\lim_{||\mathbf{x}|| \rightarrow \infty} D(\mathbf{x}|X) = 0$;
- ▶ **monotone w.r.t. the deepest point**: for any $\mathbf{x}^* \in \max_{\mathbf{x} \in \mathbb{R}^d} D(\mathbf{x}|X)$, any $\mathbf{x} \in \mathbb{R}^d$, and any $0 \leq \alpha \leq 1$ it holds: $D(\mathbf{x}|X) \leq D(\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*)|X)$;
- ▶ **upper semicontinuous in \mathbf{x}** : the upper-level sets, **depth regions**, $D_\alpha(X) = \{\mathbf{x} \in \mathbb{R}^d : D(\mathbf{x}|X) \geq \alpha\}$ are closed for all α ;
- ▶ **(quasiconcave in \mathbf{x})**: the upper-level sets are convex for all α .

Mahalanobis depth (Mahalanobis, 1936)



- ▶ Mahalanobis distance:
$$(d^M)^2(\mathbf{x}|X) = (\mathbf{x} - \boldsymbol{\mu}_X)' \boldsymbol{\Sigma}_X^{-1} (\mathbf{x} - \boldsymbol{\mu}_X)$$
- ▶ **Mahalanobis depth:**
$$D^M(\mathbf{x}|X) = \frac{1}{1 + (d^M)^2(\mathbf{x}|X)}$$
- ▶ - moment estimates for $\boldsymbol{\mu}_X$ and $\boldsymbol{\Sigma}_X$,
or robust estimates:
 - minimum volume ellipsoid,
 - minimum covariance determinant.

Zonoid depth (Koshevoy, Mosler, 1997)



- ▶ Define the **zonoid trimmed region**:

$$D_{\alpha}^z(X) = \left\{ \int_{\mathbb{R}^d} \mathbf{x} g(\mathbf{x}) dP : \begin{array}{l} g : \mathbb{R}^d \rightarrow [0, \frac{1}{\alpha}] \\ \text{measurable and} \\ \int_{\mathbb{R}^d} g(\mathbf{x}) dP = 1 \end{array} \right\}$$

for $\alpha \in (0, 1]$,

$$D_0^z(X) = \text{cl}\left(\cup_{\alpha \in (0, 1]} D_{\alpha}^z(X)\right)$$

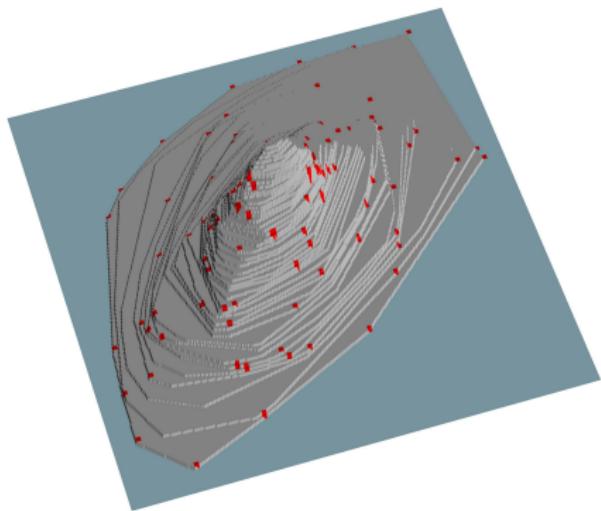
for $\alpha = 0$.

- ▶ **Zonoid depth** is then:

$$D^z(\mathbf{x}|X) = \sup\{\alpha : \mathbf{x} \in D_{\alpha}^z(X)\},$$

or 0 if no such $D_{\alpha}^z(X)$ exists for $\alpha \in [0, 1]$.

Tukey depth (Tukey, 1975)



- ▶ The smallest probability of X in a closed halfspace containing \mathbf{x} :

$$D^T(\mathbf{x}|X) = \inf\{P_X(H) : \mathbf{x} \in H, H \text{ a closed halfspace}\}$$

- ▶ Empirical version of the **Tukey depth**:

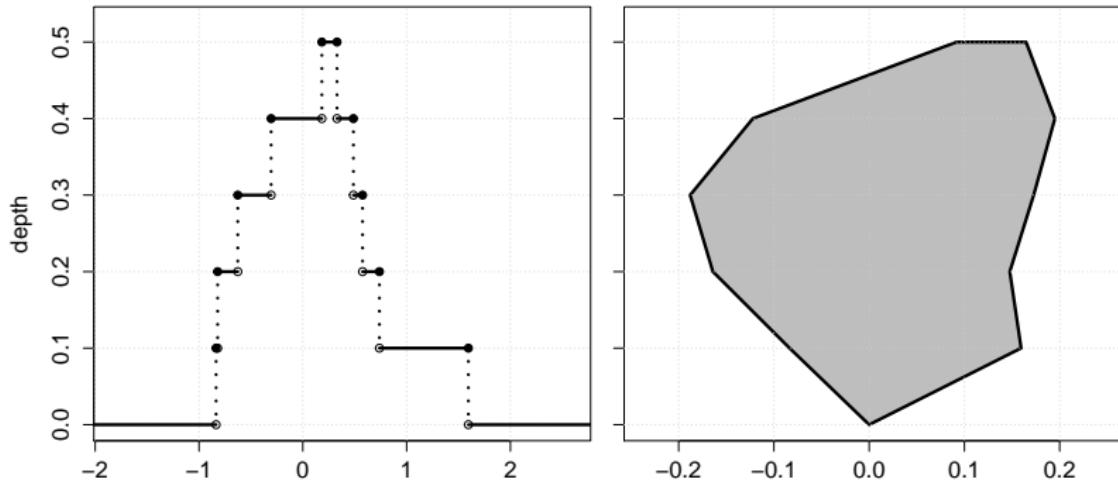
$$D^{T(n)}(\mathbf{x}|X) = \frac{1}{n} \min_{\mathbf{r} \in S^{d-1}} \#\{i : \mathbf{x}'_i \mathbf{r} \geq \mathbf{x}' \mathbf{r}\}$$

Depth lift

By adding a real dimension to the trimmed regions $D_\alpha(X)$, $\alpha \in [0, \alpha_{max}]$, define the **depth lift**:

$$D(X) = \{(\alpha, \mathbf{y}) \in [0, \alpha_{max}] \times \mathbb{R}^d : \mathbf{y} = \alpha \mathbf{x}, \mathbf{x} \in D_\alpha(X), \alpha \in [0, \alpha_{max}]\}$$

Example: Tukey depth lift (10 points)

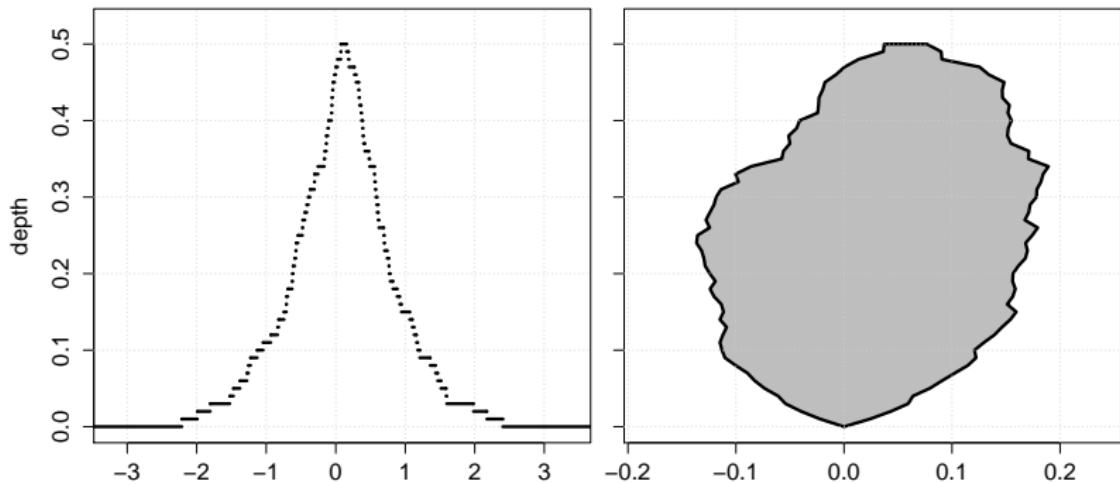


Depth lift

By adding a real dimension to the trimmed regions
 $D_\alpha(X)$, $\alpha \in [0, \alpha_{max}]$, define the **depth lift**:

$$D(X) = \{(\alpha, \mathbf{y}) \in [0, \alpha_{max}] \times \mathbb{R}^d : \mathbf{y} = \alpha \mathbf{x}, \mathbf{x} \in D_\alpha(X), \alpha \in [0, \alpha_{max}]\}$$

Example: Tukey depth lift (100 points)

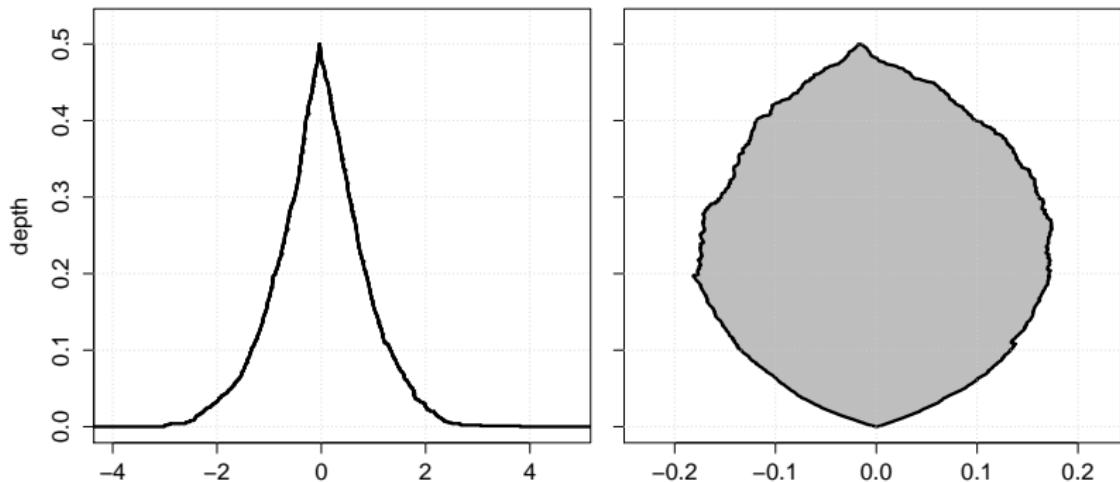


Depth lift

By adding a real dimension to the trimmed regions $D_\alpha(X)$, $\alpha \in [0, \alpha_{max}]$, define the **depth lift**:

$$D(X) = \{(\alpha, \mathbf{y}) \in [0, \alpha_{max}] \times \mathbb{R}^d : \mathbf{y} = \alpha \mathbf{x}, \mathbf{x} \in D_\alpha(X), \alpha \in [0, \alpha_{max}]\}$$

Example: Tukey depth lift (1000 points)

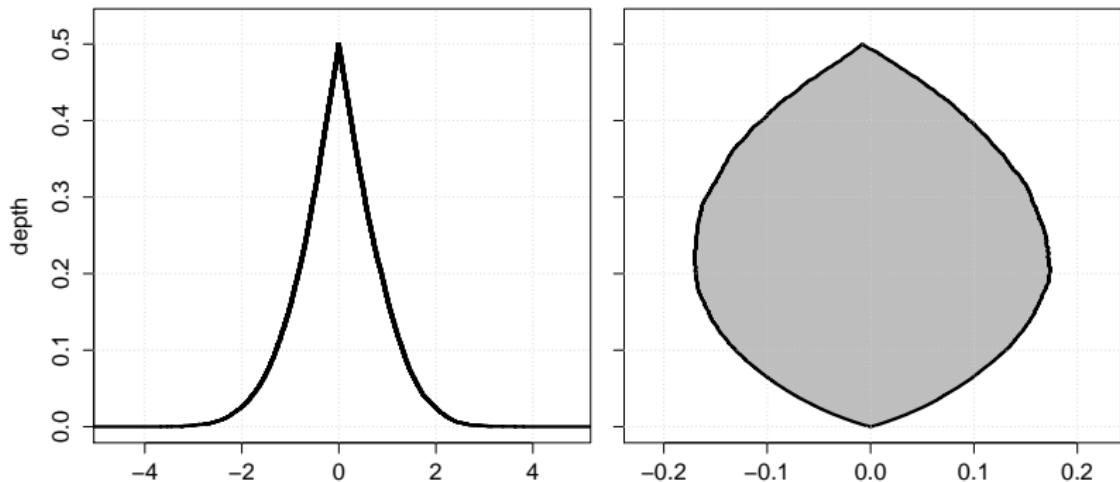


Depth lift

By adding a real dimension to the trimmed regions $D_\alpha(X)$, $\alpha \in [0, \alpha_{max}]$, define the **depth lift**:

$$D(X) = \{(\alpha, \mathbf{y}) \in [0, \alpha_{max}] \times \mathbb{R}^d : \mathbf{y} = \alpha \mathbf{x}, \mathbf{x} \in D_\alpha(X), \alpha \in [0, \alpha_{max}]\}$$

Example: Tukey depth lift (10000 points)



Contents

Motivation

Depth and depth lift

Proposal

Experiments

Conclusions and outlook

The idea

- ▶ A **compromise** between **global** (regression) and **local** (k NN) imputation.
- ▶ Unsureness in normality **away from the center** appeals to **tending to the unconditional mean** rather than the conditional one.
- ▶ Cautious to **introduce new variance** into the data.

Solution: min. **covariance determinant**. Why **not?**

- ▶ take into account **data geometry**,
- ▶ get rid of **moment assumptions**.

Given a data set $\mathbf{X} = (\mathbf{X}_{miss}, \mathbf{X}_{obs})$, impute with \mathbf{Y} :

$$\mathbf{Y} \in \underset{\mathbf{Y}_{obs} = \mathbf{X}_{obs}, \mathbf{Y}_{miss} \in \mathbb{R}^{\#\mathbf{x}_{miss}}}{\operatorname{argmin}} \operatorname{mes}(D^{(n)}(\mathbf{Y})),$$

where “ $\operatorname{mes}(A)$ ” denotes the Lebesgue measure of A .

Properties

Finite-sample extreme-value theorem

Let $\mathbf{X} = (\mathbf{X}_{miss}, \mathbf{X}_{obs})$ be a data set with missing values. Then, for any $0 < M < \infty$,

$$f(\mathbf{Y}_{miss}) = \text{mes}(D^{(n)}(\mathbf{Y} | \mathbf{Y}_{obs} = \mathbf{X}_{obs}))$$

attains its minimum on $[-M, M]^{\#\mathbf{X}_{miss}}$ for **Tukey**, **zonoid**, and **Mahalanobis** (if $\Sigma_{\mathbf{X}}$ is invertable for all $\mathbf{X}_{miss} \in [-M, M]^{\#\mathbf{X}_{miss}}$) depths.

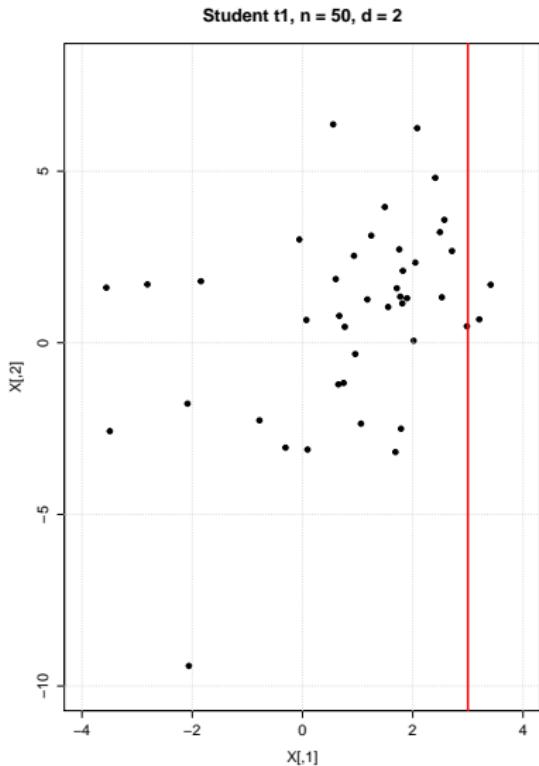
Infinite-observed extreme-value theorem

Further, let $\mathbf{X}_{obs}^{(n)}$ being sampled from absolutely continuous $\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, F)$. Then, for any $0 < M < \infty$, for **Tukey**, **zonoid** (finite 1st moment), and **Mahalanobis** (finite 2nd moment) depths the above holds, and for M large enough in the minimum for all \mathbf{y} with missingness

$$\lim_{n \rightarrow \infty} \mathbf{y}_{miss} \xrightarrow{a.s.} \boldsymbol{\mu}_{miss} + \boldsymbol{\Sigma}_{miss, obs} \boldsymbol{\Sigma}_{obs, obs}^{-1} (\mathbf{y}_{obs} - \boldsymbol{\mu}_{obs}).$$

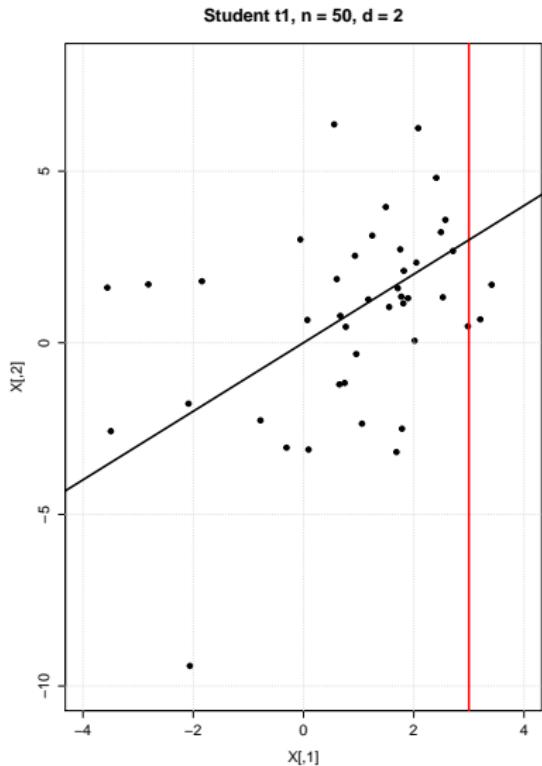
Asymptotics

-1.84206859	1.79279953
0.55599704	6.36432907
1.06121143	-2.35623912
1.75508455	2.72156758
2.01510259	0.05933206
1.71022037	1.58864924
-2.05965874	-9.41484259
3.41324550	1.69117380
1.49315060	3.95687823
2.49331460	3.22377039
2.57342204	3.58135948
.....
0.07041123	0.66430372
0.74614531	-1.17108834
-3.49366370	-2.57489754
1.78620809	-2.50428772
2.04467078	2.33675844
2.98223544	0.48124126
1.17635526	1.26232904
0.93401823	2.53479658
1.82078645	2.09722487
2.71138114	2.67442090
-10.71645912	-12.16363642
3.00000000	NA



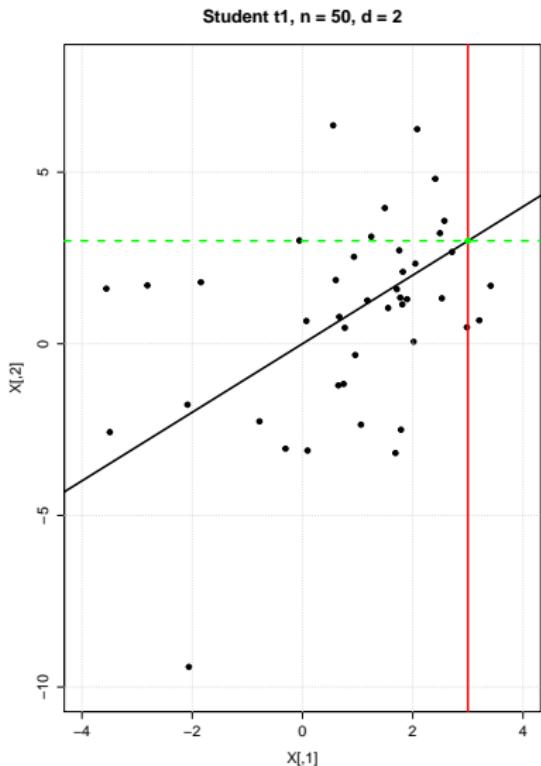
Asymptotics

-1.84206859	1.79279953
0.55599704	6.36432907
1.06121143	-2.35623912
1.75508455	2.72156758
2.01510259	0.05933206
1.71022037	1.58864924
-2.05965874	-9.41484259
3.41324550	1.69117380
1.49315060	3.95687823
2.49331460	3.22377039
2.57342204	3.58135948
.....
0.07041123	0.66430372
0.74614531	-1.17108834
-3.49366370	-2.57489754
1.78620809	-2.50428772
2.04467078	2.33675844
2.98223544	0.48124126
1.17635526	1.26232904
0.93401823	2.53479658
1.82078645	2.09722487
2.71138114	2.67442090
-10.71645912	-12.16363642
3.00000000	NA



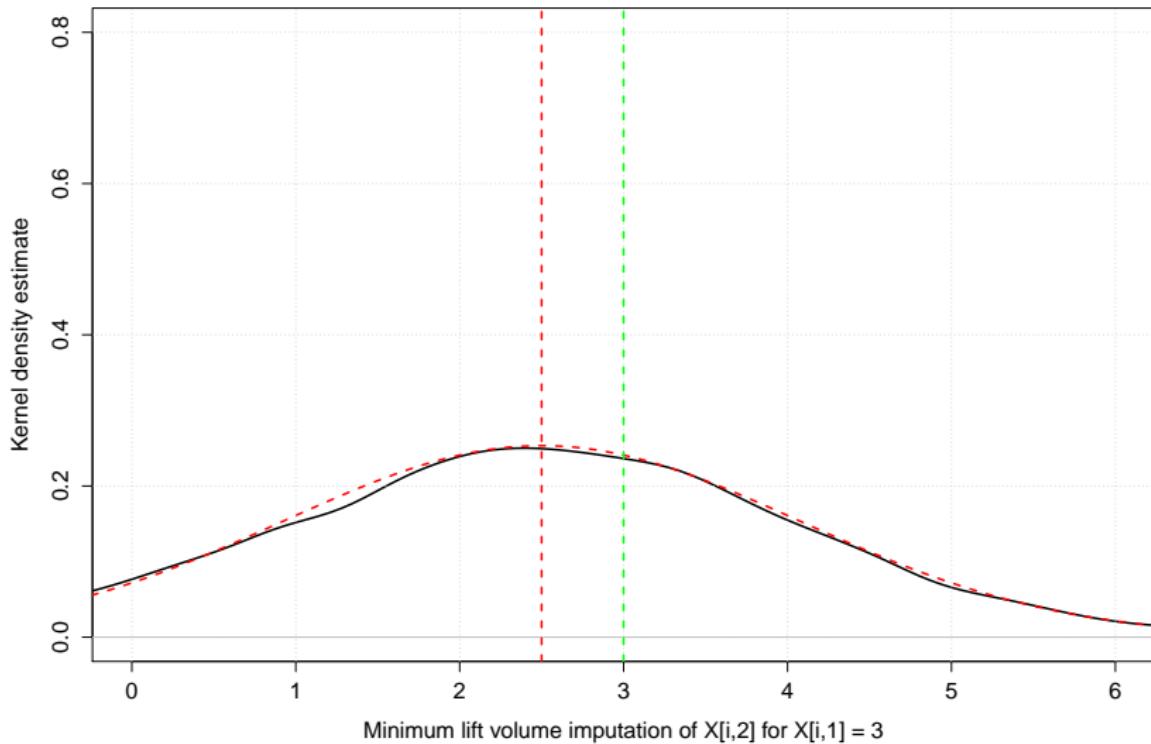
Asymptotics

-1.84206859	1.79279953
0.55599704	6.36432907
1.06121143	-2.35623912
1.75508455	2.72156758
2.01510259	0.05933206
1.71022037	1.58864924
-2.05965874	-9.41484259
3.41324550	1.69117380
1.49315060	3.95687823
2.49331460	3.22377039
2.57342204	3.58135948
.....
0.07041123	0.66430372
0.74614531	-1.17108834
-3.49366370	-2.57489754
1.78620809	-2.50428772
2.04467078	2.33675844
2.98223544	0.48124126
1.17635526	1.26232904
0.93401823	2.53479658
1.82078645	2.09722487
2.71138114	2.67442090
-10.71645912	-12.16363642
3.00000000	3.00000000



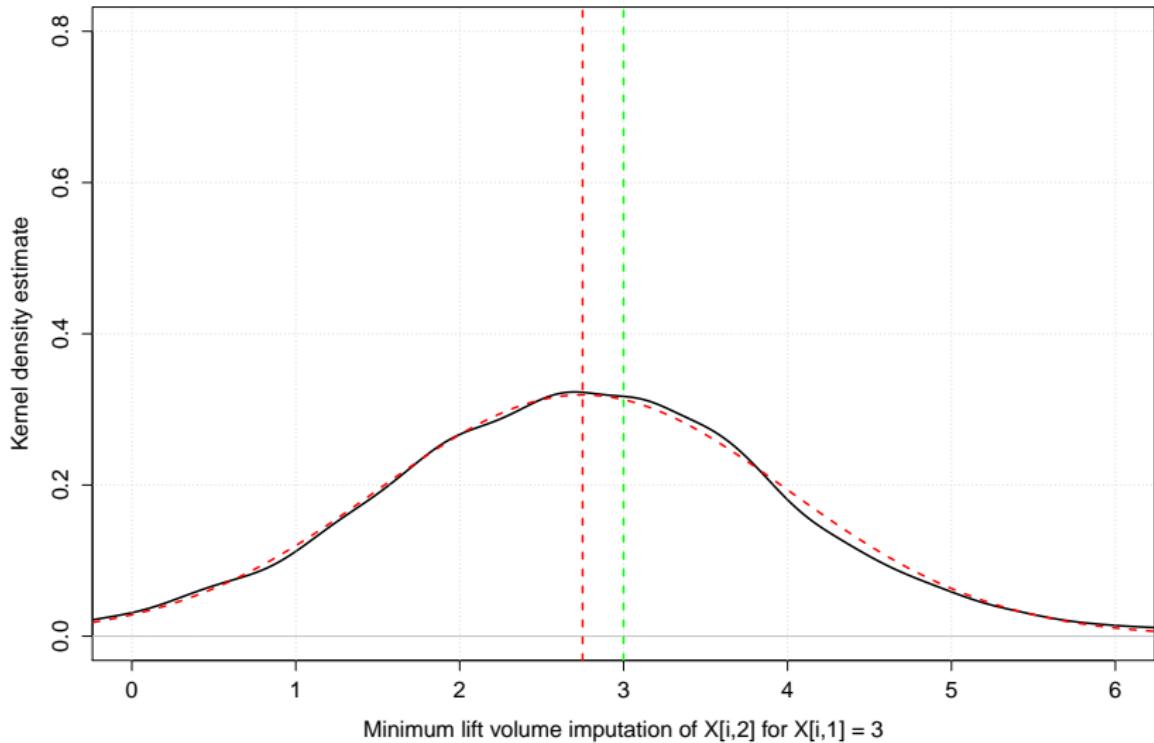
Asymptotics (Tukey depth)

Student t1, n = 50, d = 2, k = 10000



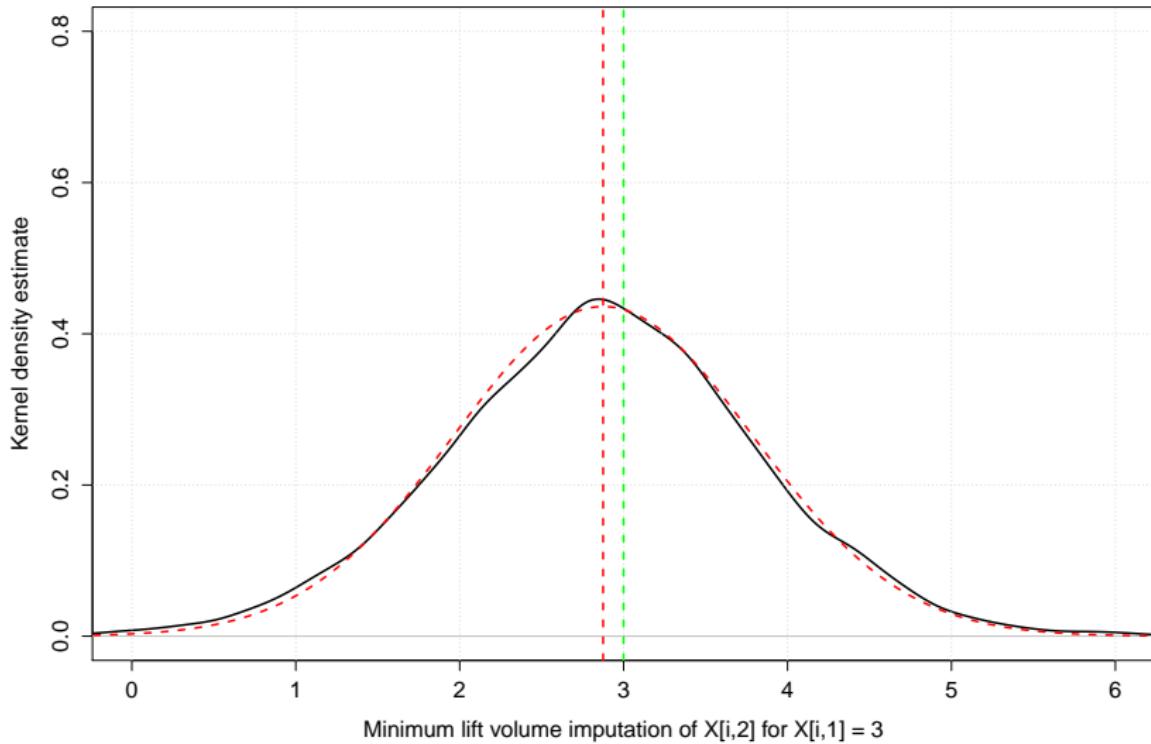
Asymptotics (Tukey depth)

Student t1, n = 100, d = 2, k = 10000



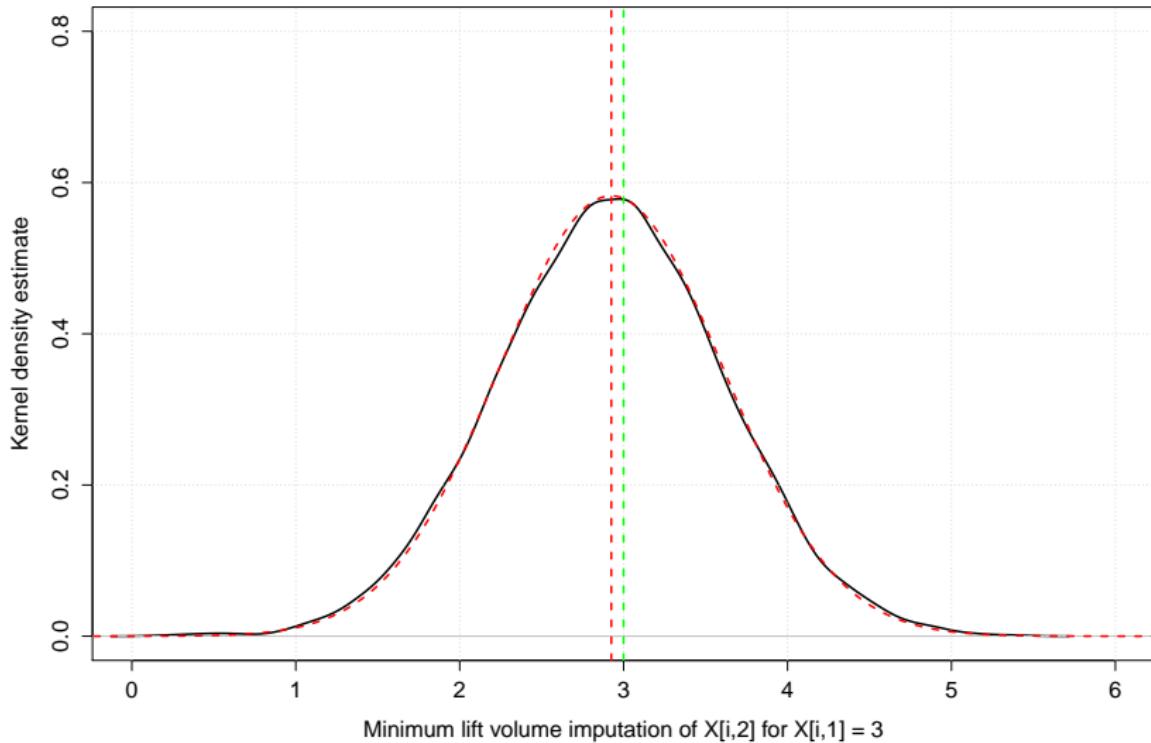
Asymptotics (Tukey depth)

Student t1, n = 200, d = 2, k = 10000



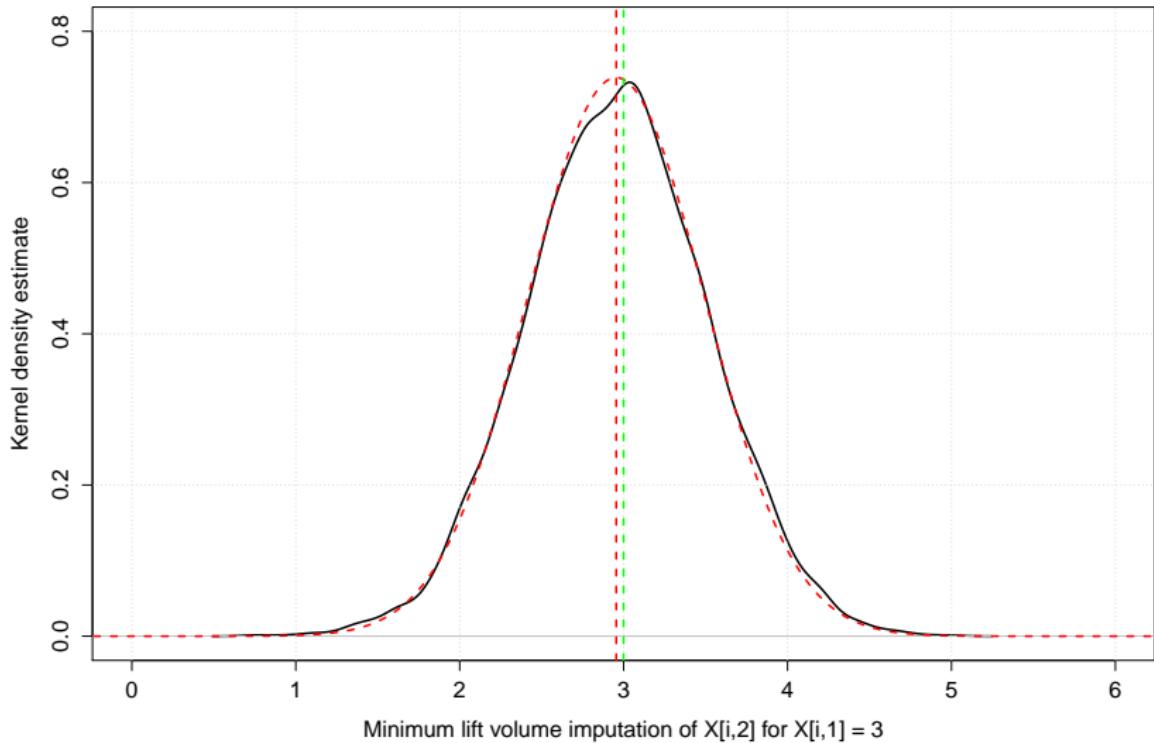
Asymptotics (Tukey depth)

Student t1, n = 500, d = 2, k = 10000



Asymptotics (Tukey depth)

Student t1, n = 1000, d = 2, k = 10000



Special case: Mahalanobis depth

Impute $\mathbf{X} = (\mathbf{X}_{miss}, \mathbf{X}_{obs})$ in \mathbb{R}^d with \mathbf{Y} so that for each \mathbf{y} having missing values holds:

$$\mathbf{y}_{miss} = \operatorname{argmin}_{\mathbf{y}_{obs}=\mathbf{x}_{obs}} \text{mes}(D^M(\mathbf{Y})).$$

Then for each \mathbf{y} having missing values holds as well:

- ▶ \mathbf{y} is imputed with the **conditional mean**:

$$\mathbf{y}_{miss} = \boldsymbol{\mu}(\mathbf{Y})_{miss} + \boldsymbol{\Sigma}(\mathbf{Y})_{miss, obs} \boldsymbol{\Sigma}^{-1}(\mathbf{Y})_{obs, obs} (\mathbf{y}_{obs} - \boldsymbol{\mu}(\mathbf{Y})_{obs}),$$

- ▶ \mathbf{y} is imputed with single-output **regression**,
- ▶ \mathbf{y} is imputed with **regularized PCA** by Josse & Husson (2012) with any $0 < \sigma^2 \leq \lambda_d$:

$$\mathbf{y}_{miss(i)} = \sum_{j=1}^d \sqrt{\frac{\lambda_j - \sigma^2}{\lambda_j}} \mathbf{u}_j \mathbf{v}_{miss(i), j}$$

with $\mathbf{Y} = \mathbf{U} \Lambda \mathbf{V}$ being the centered SVD,

- ▶ \mathbf{y} is imputed with **maximum depth**:

$$\mathbf{y}_{miss} = \operatorname{argmax}_{\mathbf{y}_{obs}=\mathbf{x}_{obs}} D^M(\mathbf{y} | \mathbf{Y}).$$

Computation

Direct minimization requires evaluation of empirical version of $\text{mes}(D^{(n)}(\mathbf{X}))$ for n points in \mathbb{R}^d , namely

- ▶ for Mahalanobis depth:

$$\begin{aligned}\text{mes}(D^{M(n)}(\mathbf{X})) &= \frac{\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2} + 1)} \prod_{i=1}^{\lfloor \frac{d}{2} \rfloor} \frac{2i + 2(d \bmod 2)}{i + \frac{1}{2} + d \bmod 2} \times \\ &\quad \times \left(1 - \left(1 - \frac{\pi}{8}\right) \mathbb{I}(d \bmod 2 \neq 0)\right) \sqrt{\|\Sigma\mathbf{x}\|},\end{aligned}$$

- ▶ for zonoid depth:

$$\text{mes}(D^{Z(n)}(\mathbf{X})) = \frac{1}{n^{d+1}} \sum_{\{i_0, \dots, i_d\} \subset \{1, \dots, n\}} \|((1, \mathbf{x}'_{i_0})', \dots, (1, \mathbf{x}'_{i_d})')\|,$$

- ▶ for Tukey depth:

$$\text{mes}(D^{T(n)}(\mathbf{X})) = \sum_{i=1}^{n\alpha_{\max}} \frac{i^{d+1} - (i-1)^{d+1}}{(d+1)n^{d+1}} \text{mes}(D_{\frac{i}{n}}^{T(n)}(\mathbf{X})).$$

Computation

Denote

$$\alpha^* = \inf_{\alpha \in (0;1)} \{ \alpha \mid \text{int } D_\alpha(\mathbf{X}) \cap \{ \mathbf{z} \mid \mathbf{z} \in \mathbb{R}^d, \mathbf{z}_{obs} = \mathbf{x}_{obs} \} = \emptyset \},$$

and impute \mathbf{x} having missing values with

$$\mathbf{y} = \text{ave} \left(\arg \min_{\mathbf{u} \in \mathbb{R}^d, \mathbf{u}_{obs} = \mathbf{x}_{obs}} \{ \| \mathbf{u} - \mathbf{v} \| \mid \mathbf{v} \in \mathbb{R}^d, \mathbf{v} \in D_{\alpha^*}(\mathbf{X}), \| \cdot \| \} \right).$$

For continuous depth this is equivalent to:

$$\mathbf{y} = \arg \max_{\mathbf{z} \in \mathbb{R}^d, \mathbf{z}_{obs} = \mathbf{x}_{obs}} D(\mathbf{z} | \mathbf{X}).$$

After an arbitrary **initialization** continue imputation for each point **iteratively** until **convergence**.

Contents

Motivation

Depth and depth lift

Proposal

Experiments

Conclusions and outlook

Student t-distribution (MCAR)

$n = 50, \mu = (1, 1, 1)', \Sigma = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 4 & 4 \\ 1 & 4 & 8 \end{pmatrix}, \text{miss} = 25\%.$

MSE	Normal	t_{10}	t_5	t_3	t_2	Cauchy
D.Tuk	2.773 (0.8294)	3.392 (1.129)	4.401 (1.742)	6.461 (3.570)	11.49 (8.006)	135.5 (164.7)
D.Zon	2.690 (0.8251)	3.360 (1.161)	4.451 (1.823)	6.844 (4.176)	13.08 (9.493)	179.8 (215.1)
D.Mah	2.719 (0.8593)	3.355 (1.158)	4.451 (1.905)	6.743 (3.872)	13.15 (9.369)	175.1 (207.1)
EM	2.532 (0.7893)	3.124 (1.046)	4.042 (1.647)	6.169 (3.390)	11.55 (8.114)	162.8 (193.1)
RF	3.143 (0.9778)	3.967 (1.402)	5.096 (2.007)	7.592 (4.170)	12.95 (9.193)	153.8 (184.8)
k NN	3.261 (1.0415)	3.919 (1.462)	5.104 (2.292)	7.266 (3.657)	13.10 (9.024)	156.6 (189.2)
rPCA2	2.719 (0.8576)	3.355 (1.158)	4.444 (1.905)	6.752 (3.894)	13.16 (9.404)	175.8 (211.8)
rPCA1	2.712 (0.8222)	3.299 (1.143)	4.310 (1.713)	6.433 (3.619)	12.24 (8.787)	159.4 (192.4)

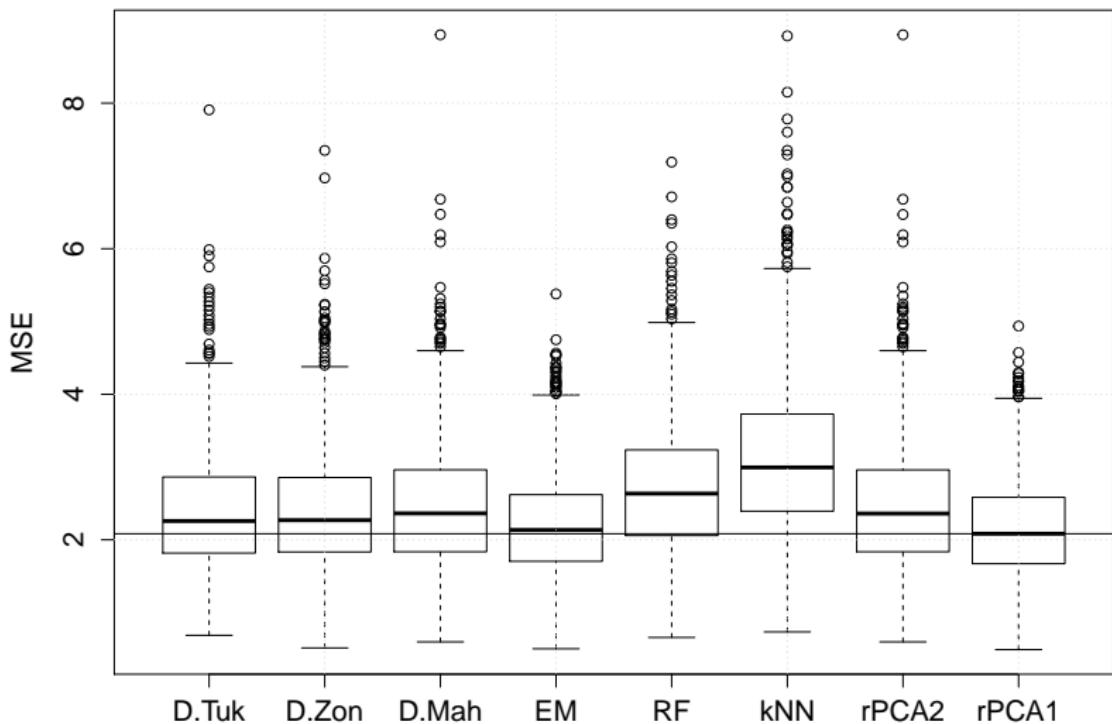
Student t-distribution + Cauchy outliers (MCAR)

$n = 50, \mu = (1, 1, 1)', \Sigma = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 4 & 4 \\ 1 & 4 & 8 \end{pmatrix}, \text{miss} = 25\%, \text{outl} = 15\%.$

MSE	Normal	<i>t</i> 10	<i>t</i> 5	<i>t</i> 3	<i>t</i> 2	Cauchy
D.Tuk	3.025 (1.066)	3.678 (1.374)	4.676 (1.997)	6.555 (3.563)	12.03 (8.692)	125.5 (154.4)
D.Zon	3.480 (1.463)	4.420 (1.895)	5.371 (2.576)	8.115 (4.784)	14.78 (11.61)	162.0 (198.5)
D.Mah	3.787 (1.848)	4.519 (2.136)	5.614 (2.784)	8.150 (5.013)	14.86 (11.53)	169.8 (210.5)
EM	3.516 (1.673)	4.288 (1.950)	5.261 (2.557)	7.554 (4.622)	13.82 (10.67)	153.4 (189.8)
RF	3.520 (1.204)	4.251 (1.602)	5.324 (2.276)	7.930 (4.398)	13.98 (10.30)	139.9 (172.3)
kNN	3.535 (1.227)	4.294 (1.647)	5.313 (2.285)	7.690 (4.299)	13.55 (9.275)	143.5 (174.6)
rPCA2	3.793 (1.846)	4.532 (2.139)	5.609 (2.796)	8.151 (5.013)	14.86 (11.55)	168.8 (210.3)
rPCA1	3.657 (1.774)	4.424 (2.201)	5.520 (2.736)	8.339 (5.283)	14.51 (11.36)	157.9 (197.0)

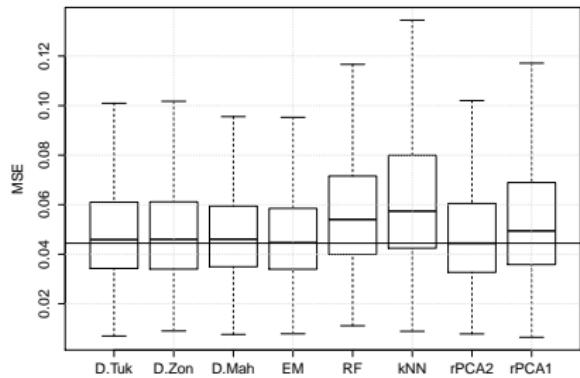
Normal distribution (MAR)

$n = 50, \mu = (1, 1, 1)', \Sigma = \begin{pmatrix} 1 & 1.75 & 2 \\ 1.75 & 4 & 4 \\ 2 & 4 & 8 \end{pmatrix}$, miss = 24%.

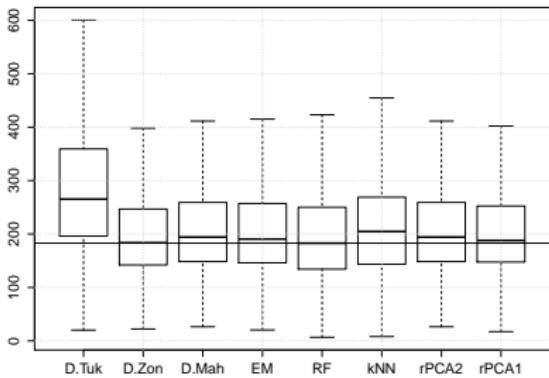


Real data (MCAR, 10%)

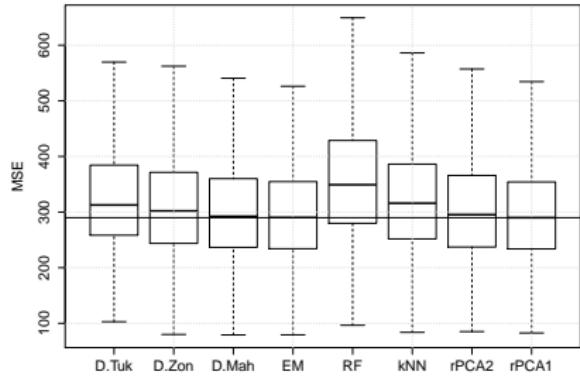
Iris (setosa), $n = 50$, $d = 4$



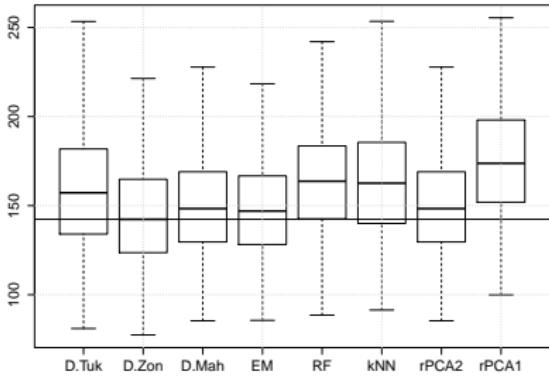
Glass (float), $n = 70$, $d = 3$



Indian diabetes, $n = 200$, $d = 4$



Blood transfusion, $n = 502$, $d = 3$



Contents

Motivation

Depth and depth lift

Proposal

Experiments

Conclusions and outlook

Conclusions and outlook

Suggested is a **framework** for **single imputation** based on **data depth**, which incorporates features of both **global** and **local** approaches, and imputes **close to data** geometry, preserves functionality under **missing at random** (MAR) mechanism, is **robust** both in sense of distribution and outliers.

Future work:

- ▶ Detailed exploration of **theoretical properties**, especially for the family of **elliptically symmetric** distributions.
- ▶ Extension to **higher-dimensional** data, **larger samples**, further **experimental study**.
- ▶ Extension to **multiple imputation**.
- ▶ Implementation in an **R-package**.

Questions? Suggestions? Ideas?

Thank you for your attention!