On algorithms for computation of the Tukey depth

Rainer Dyckerhoff^a Xiaohui Liu^b Karl Mosler^a Pavlo Mozharovskyi^c

^aInstitute of Econometrics and Statistics, University of Cologne ^bSchool of Statistics, Jiangxi University of Finance and Economics; Research Center of Applied Statistics, Jiangxi University of Finance and Economics ^cLTCI, Telecom Paris, Institut Polytechnique de Paris

LTCI Data Science seminar

Paris, June 6, 2019

Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Multivariate data



/		Weight	Age
	Subject 1	1350	32
	Subject 2	1500	32
	Subject 162	1320	28
١,	Subject 161	1150	27

Babies with low birth weight



Babies with low birth weight



Contents

Data depth Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Statistical data depth

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

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

that is:

- affine invariant: D(Ax + b|AX + b) = D(x|X);
- ▶ vanishing at infinity: $\lim_{||\mathbf{x}|| \to \infty} D(\mathbf{x}|X) = 0$;
- ▶ monotone w.r.t. the deepest point: for any $\mathbf{x}^* \in \operatorname{argmax}_{\mathbf{x} \in \mathbb{R}^d} D(\mathbf{x}|X)$, any $\mathbf{x} \in \mathbb{R}^d$, and any $0 \le \alpha \le 1$ it holds: $D(\mathbf{x}|X) \le D(\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*)|X)$;
- upper semicontinuous in **x**: the upper-level sets $D_{\tau}(X) = \{ \mathbf{x} \in \mathbb{R}^d : D(\mathbf{x}|X) \ge \tau \}$ are closed for all τ ;
- (quasiconcave in x): the upper-level sets are convex for all τ .

Tukey (1975) — "Mathematics and the picturing of data"

Tukey depth of $\mathbf{x} \in \mathbb{R}^d$ w.r.t. a *d*-variate random vector X distributed as P is defined as the smallest probability mass of a closed halfspace containing \mathbf{x} :

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

and w.r.t. a data set $\boldsymbol{X} = \{\boldsymbol{x}_1,...,\boldsymbol{x}_n\} \subset \mathbb{R}^d$:

$$D^{T(n)}(\boldsymbol{x}|\boldsymbol{X}) = \frac{1}{n} \min_{\boldsymbol{u} \in \mathbb{S}^{d-1}} \sharp\{i : \boldsymbol{u}' \boldsymbol{x}_i \geq \boldsymbol{u}' \boldsymbol{x}\}.$$

Other depth notions: Mahalanobis ('36), projection (Stahel, '81; Donoho, '82), simplicial volume (Oja, '83), simplicial (Liu, '90), zonoid (Koshevoy, Mosler, '97), spatial (Vardi, Zhang, '00; Serfling, '02) depth.

Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight





Contents

Data depth Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Tukey-trimmed regions

Tukey depth defines a family of (depth-)trimmed (central) regions $D_{\tau}^{T}(X)$, the upper-level sets of the depth function:

$$D_{ au}^{T}(X) = ig\{ oldsymbol{x} \in \mathbb{R}^{d} \, : \, D^{T}(oldsymbol{x}|X) \geq au ig\}.$$

Properties:

Depth:

- Affine invariant;
- Vanishing at infinity;
- Monotone w.r.t. deepest point;
- Upper-semicontinuous;
- Quasiconcave.

Regions:

Affine equivariant;

Bounded;

Nested;

Closed;

Convex.

Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight


Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight





Babies with low birth weight



Babies with low birth weight

0000 00 Age, in weeks - 5

Babies with low birth weight



Babies with low birth weight



Babies with low birth weight



Babies with low birth weight

Babies with low birth weight



Babies with low birth weight



Tukey (=halfspace, location) data depth





Tukey (=halfspace, location) depth region: $\tau = 2/161$



Tukey (=halfspace, location) depth region: $\tau = 5/161$



Tukey (=halfspace, location) depth region: $\tau = 9/161$



Tukey (=halfspace, location) depth region: $\tau = 13/161$



Tukey (=halfspace, location) depth region: $\tau = 17/161$



Tukey (=halfspace, location) depth region: $\tau = 25/161$



•

Tukey (=halfspace, location) depth region: $\tau = 33/161$



•

Tukey (=halfspace, location) depth region: $\tau = 41/161$





Tukey (=halfspace, location) depth region: $\tau = 49/161$



.

Tukey (=halfspace, location) depth region: $\tau = 57/161$



.

Tukey (=halfspace, location) depth region: $\tau = 65/161$



Tukey (=halfspace, location) depth region: $\tau = 68/161$



Computation of the Tukey depth: literature

- Rousseeuw, P.J. and Ruts, I. (1996).
 "Algorithm AS 307: Bivariate location depth." Journal of the Royal Statistical Society, Series C, 45, 516-526.
- Ruts, I. and Rousseeuw, P.J. (1996).
 "Computing depth contours of bivariate point clouds." Computational Statistics and Data Analysis, 23, 153-168.
- Rousseeuw, P.J. and Struyf, A. (1998).
 "Computing location depth and regression depth in higher dimensions." Statistics and Computing, 8, 193-203.

 Hallin, M., Paindaveine, D., and Šiman, M. (2010).
 "Multivariate quantiles and multiple-output regression quantiles: From L₁ optimization to halfspace depth." The Annals of Statistics, 38, 635-669.

- Paindaveine, D. and Šiman, M. (2011).
 "On directional multiple-output quantile regression." Journal of Multivariate Analysis, 102, 193–212.
- Liu, X. and Zuo, Y. (2014).
 "Computing halfspace depth and regression depth."
 Communications in Statistics Simulation and Computation, 43, 969-985.
Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Applications of data depth:

- Multivariate data analysis (Liu, Parelius, Singh '99);
- Statistical quality control (Liu, Singh '93);
- Cluster analysis and classification (Mosler, Hoberg '06; Li, Cuesta-Albertos, Liu '12; M., Mosler, Lange '15);
- Tests for multivariate location, scale, symmetry (Liu '92; Dyckerhoff '02; Dyckerhoff, Ley, Paindaveine '15);
- Outlier detection (Hubert, Rousseeuw, Segaert '15);
- Multivariate risk measurement (Cascos, Mochalov '07);
- Robust linear programming (Bazovkin, Mosler '15);
- Missing data imputation (M., Josse, Husson '18);
- etc.

R-package **ddalpha** (Pokotylo, M., Dyckerhoff, Nagy): calculates a number of depths; performs depth-based classification of multivariate and functional data; contains 50 multivariate and 5 functional data sets.

Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth Theoretical background

Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Theoretical background

• **Reminder:** Tukey depth of $x \in \mathbb{R}^d$ w.r.t. a data set X is:

$$D^{\mathcal{T}(n)}(\boldsymbol{x} \mid \boldsymbol{X}) = rac{1}{n} \min_{\boldsymbol{p} \in \mathbb{R}^d \setminus \{\boldsymbol{0}\}} \sharp \{i : \boldsymbol{p}' \boldsymbol{x}_i \geq \boldsymbol{p}' \boldsymbol{x}\}.$$

Due to the affine-invariance property we can rewrite:

$$D^{T(n)}(\boldsymbol{z} | \boldsymbol{x}_1, \ldots, \boldsymbol{x}_n) = D^{T(n)}(\boldsymbol{0} | \boldsymbol{x}_1 - \boldsymbol{z}, \ldots, \boldsymbol{x}_n - \boldsymbol{z}).$$

▶ We call a vector $m{p} \in \mathbb{R}^d ackslash \{ m{0} \}$ optimal for the data set $m{X}$ if

$$D^{\mathcal{T}(n)}(\mathbf{0} \mid \mathbf{X}) = \frac{1}{n} \#\{i \mid \mathbf{p}' \mathbf{x}_i \ge 0\}.$$

For *p* ≠ 0 we define:

$$I_{\boldsymbol{p}}^{+} = \{i \mid \boldsymbol{p}' \boldsymbol{x}_{i} > 0\}, \quad I_{\boldsymbol{p}}^{0} = \{i \mid \boldsymbol{p}' \boldsymbol{x}_{i} = 0\}, \quad I_{\boldsymbol{p}}^{-} = \{i \mid \boldsymbol{p}' \boldsymbol{x}_{i} < 0\}.$$

Proposition

If $\mathbf{p} \neq 0$ is optimal for \mathbf{X} , then $I_{\mathbf{p}}^0 = \emptyset$, i.e., no data points lie on the boundary of the closed halfspace defined by \mathbf{p} .















Main result

- ▶ Denote by L_k the set of all subsets I of order k of {1,..., n} such that the points (x_i)_{i∈I} are linearly independent.
- ▶ For any $I \in \mathcal{L}_k$ with $I = \{i_1, \ldots, i_k\}$, let $P_I = [\mathbf{x}_{i_1}, \ldots, \mathbf{x}_{i_k}]$.
- For any *I* ∈ L_k with *I* = {*i*₁,...,*i*_k}, let *a*₁,..., *a*_{d-k} be a basis of the orthogonal complement of span(*x*_{i1},...,*x*_{ik}) and *A*_l the matrix whose columns are the *a*_i.
- Denote by I^c the complement of a set I.
- For a set *I* = {*i*₁,...,*i_k*} of indices, we denote by *I*^{*} = {*j* | *x_j* ∈ span(*x_{i1}*,...,*x_{ik}*), *j* = 1,...,*n*} the set of all indices *j* such that *x_j* is contained in the linear hull of *x_{i1}*,...,*x_{ik}*.

Theorem

For each k such that $1 \le k < d$ it holds that

$$n \cdot D^{T(n)}(\mathbf{0} \mid \mathbf{X}) = \min_{I \in \mathcal{L}_k} \left[n \cdot D^{T(n)}(\mathbf{0} \mid \mathbf{A}_I' \mathbf{X}_{(I^*)^c}) + n \cdot D^{T(n)}(\mathbf{0} \mid \mathbf{P}_I' \mathbf{X}_{I^*}) \right]$$

Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth Theoretical background

Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Example: combinatorial algorithm with k = d - 2

1: function NHD_COMB2
$$(d, \mathbf{x}_1, \dots, \mathbf{x}_n)$$
 > Halfspace depth of 0
2: if $d = 1$ then return NHD1 $(\mathbf{x}_1, \dots, \mathbf{x}_n)$
3: if $d = 2$ then return NHD2 $(\mathbf{x}_1, \dots, \mathbf{x}_n)$
4: $n_{min} \leftarrow n$
5: for each subset $l \subset \{1, \dots, n\}$ of order $d - 2$ do
6: if $(\mathbf{x}_i)_{i \in l}$ linearly independent then
7: Compute basis $\mathbf{a}_1, \mathbf{a}_2$ of orthogonal complement of $(\mathbf{x}_i)_{i \in l}$
8: $\mathbf{A}_l \leftarrow \text{Matrix}[\mathbf{a}_1, \mathbf{a}_2]$
9: $\mathbf{P}_l \leftarrow \text{Matrix}[\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_{d-2}}]$
10: for all \mathbf{x}_j do
11: if $\mathbf{A}'_l \mathbf{x}_j \neq \mathbf{0}$ then $\mathbf{y}_j \leftarrow \mathbf{A}'_l \mathbf{x}_j$
12: else $\mathbf{z}_j \leftarrow \mathbf{P}'_l \mathbf{x}_j$
13: $l \leftarrow \#\{j : \mathbf{A}'_l \mathbf{x}_j \neq \mathbf{0}\}$
14: $n_{new} \leftarrow \text{NHD2}(\mathbf{y}_{j_1}, \dots, \mathbf{y}_{j_l})$
15: if $n - l > d - 2$ then
16: $n_{new} \leftarrow n_{min}$ + NHD_COMB2 $(d - 2, \mathbf{z}_{j_1}, \dots, \mathbf{z}_{j_{n-l}})$
17: if $n_{new} < n_{min}$ then $n_{min} \leftarrow n_{new}$















Time of depth calculation (sec.): k = 1, d - 2, d - 1

	<i>n</i> = 40	80	160	320	640	1280	2560	5120	10240	20480	40960	81920
<i>d</i> = 3	0.000 0.002 0.000	0.000 0.002 0.003	0.000 0.003 0.014	0.011 0.016 0.117	0.047 0.063 0.936	0.184 0.250 7.60	0.780 1.03 61.3	3.19 4.22 519	13.2 17.4	54.5 72.2	225 293	933 1210 —
4	0.006 0.005 0.005	0.048 0.038 0.055	0.402 0.302 0.784	3.36 2.50 12.3	28.2 20.4 203	235 166 3290	1960 1360 —					
5	0.205 0.055 0.047	3.62 0.952 1.24	62.5 16.1 35.7	1070 269 1110								
6	7.32 0.506 0.392	275 18.4 21.6	633 1250									
7	257 3.60 2.66	 278 305										
8	21.4 14.3	3550 3470										
9	 107 69.8											
10	439 289											

Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions Existing approaches

The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Directional quantiles

Kong, Mizera (2008)

For a random vector $X \in \mathbb{R}^d$ define the *u*-directional au-th quantile:

$$Q(\tau, \boldsymbol{u}, X) = \inf\{x : P(\boldsymbol{u}'X \leq x) \geq \tau\}.$$

Directional quantile envelope:

$$R_{ au}(X) = \bigcap_{oldsymbol{u}\in\mathbb{S}^{d-1}} Hig(oldsymbol{u}, Q(au,oldsymbol{u},X)ig),$$

where $H(\mathbf{u}, q) = \{x : \mathbf{u}' x \ge q\}$ is the supporting halfspace determined by $\mathbf{u} \in \mathbb{S}^{d-1}$ and $q \in \mathbb{R}$. Then for every $p \in (0; \frac{1}{2}]$ directional quantile envelopes coincide with the Tukey regions:

$${\it R}_{ au}(X)={\it D}_{ au}^{\,{\cal T}}(X)\,$$
 for every $\, au\in(0;rac{1}{2}].$

Ø 0 0 °œ °°











0

Multiple-output regression quantiles

Hallin, Paindaveine, Šiman (2010); Paindaveine, Šiman (2011) Regress $Y \in \mathbb{R}^m$ on $X = (1, W')' \in \mathbb{R}^p$. For a data set $(X, Y) = \{(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}^m; i = 1, ..., n\}$, for $\tau \in (0, 1)$, and for $u \in S^{m-1}$ a (τu) -quantile positive halfspace is any

$$\begin{split} & \mathcal{H}_{\tau \boldsymbol{u}}^{(n)+} := \{ (\boldsymbol{w}', \boldsymbol{y}')' \in \mathbb{R}^{p-1} \times \mathbb{R}^m : \hat{\boldsymbol{b}}_{\tau \boldsymbol{u}}' \boldsymbol{y} - \hat{\boldsymbol{a}}_{\tau \boldsymbol{u}}' (1, \boldsymbol{w}')' \geq 0 \} \text{ with} \\ & (\hat{\boldsymbol{a}}_{\mathsf{HP\check{S}};\tau \boldsymbol{u}}', \hat{\boldsymbol{b}}_{\mathsf{HP\check{S}};\tau \boldsymbol{u}}')' = \operatorname{argmin} \sum_{i=1}^n \rho_\tau (\boldsymbol{b}' \boldsymbol{y}_i - \boldsymbol{a}' \boldsymbol{x}_i) \text{ subject to } \boldsymbol{u}' \boldsymbol{b} = 1, \\ & (\hat{\boldsymbol{a}}_{\mathsf{proj};\tau \boldsymbol{u}}', \hat{\boldsymbol{b}}_{\mathsf{proj};\tau \boldsymbol{u}}')' = \operatorname{argmin} \sum_{i=1}^n \rho_\tau (\boldsymbol{b}' \boldsymbol{y}_i - \boldsymbol{a}' \boldsymbol{x}_i) \text{ subject to } \boldsymbol{u} = \boldsymbol{b}, \end{split}$$

where $ho_{ au}(x) = x (au - I(x < 0))$ is the au-quantile check function.

i=1

Multiple-output regression quantiles (location)

Hallin, Paindaveine, Šiman (2010); Paindaveine, Šiman (2011) Regress $X \in \mathbb{R}^d$ on $1 \in \mathbb{R}$. For a data set $\mathbf{X} = {\mathbf{x}_i \in \mathbb{R}^d; i = 1, ..., n}$, for $\tau \in (0, 1)$, and for $\mathbf{u} \in \mathbb{S}^{d-1}$ a $(\tau \mathbf{u})$ -quantile positive halfspace is any

$$H^{(n)+}_{ auoldsymbol{u}}:=\{oldsymbol{x}'\in\mathbb{R}^d\ :\ \hat{oldsymbol{b}}_{ auoldsymbol{u}}'oldsymbol{x}-\hat{a}_{ auoldsymbol{u}}\geq0\}$$
 with

$$(\hat{a}_{\mathsf{HP\check{S}};\tau\boldsymbol{u}}, \hat{\boldsymbol{b}}_{\mathsf{HP\check{S}};\tau\boldsymbol{u}}')' = \operatorname{argmin} \sum_{i=1}^{n} \rho_{\tau}(\boldsymbol{b}'\boldsymbol{x}_{i} - \boldsymbol{a}) \text{ subject to } \boldsymbol{u}'\boldsymbol{b} = 1,$$

or

$$(\hat{a}_{\text{proj};\tau u}, \hat{b}'_{\text{proj};\tau u})' = \operatorname{argmin} \sum_{i=1}^{n} \rho_{\tau} (b' x_i - a) \text{ subject to } u = b,$$

where $\rho_{\tau}(x) = x(\tau - I(x < 0))$ is the τ -quantile check function. Then for some $\mathcal{L}_{\mathbf{X},\tau}$ with $\#\mathcal{L}_{\mathbf{X},\tau} < \infty$ it holds:

$$D_{\tau}^{T(n)}(\boldsymbol{X}) = R_{\tau}^{(n)}(\boldsymbol{X}) = \bigcap_{\boldsymbol{u} \in \mathbb{S}^{d-1}} \{H_{\tau \boldsymbol{u}}^{(n)+}\} = \bigcap_{\boldsymbol{u} \in \mathcal{L}_{\boldsymbol{X},\tau}} \{H_{\tau \boldsymbol{u}}^{(n)+}\}.$$

Multiple-output regression quantiles (illustration)



Multiple-output regression quantiles (illustration)


















































0



Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Univariate projection quantiles (algorithm)

Input: $X = \{x_1, \dots, x_n\} \in \mathbb{R}^d, \ 2 < d < n < \infty.$

Step 1: Set matrix $\mathcal{A} = (\texttt{false}_n)^{d-1}$, tree $\mathcal{G}_{\tau} = \emptyset$, queue $\mathcal{Q} = \emptyset$.

Step 2: Generate a sufficient initial set of (d-1)-tuples (*e.g.* ridges of the data convex hull), push them into Q and set them true in A.

Step 3: Pop a (d-1)-tuple, say $[i_1, \cdots, i_{d-1}]$, from Q.

Step 4: Find all subscripts $i_d \in \{1, \dots, n\} \setminus \{i_1, \dots, i_{d-1}\}$ s. t. $\{x_{i_1}, \dots, x_{i_{d-1}}, x_{i_d}\}$ determine a τ -th halfspace, store these in \mathcal{T} .

Step 5: For each $j(\text{determining halfspace } \mathbf{g}_j) \in \mathcal{T}$ do:

5.1: If $\mathbf{g}_j \notin \mathcal{G}_{\tau}$, add \mathbf{g}_j to \mathcal{G}_{τ} , else go to the next element in \mathcal{T} . **5.2:** For each ridge of \mathbf{g}_j , except $[i_1, \cdots, i_{d-1}]$, say $[j_1, \dots, j_{d-1}]$, do: **5.2.1:** If $\mathcal{A}_{[j_1, \dots, j_{d-1}]} = \text{false, set } \mathcal{A}_{[j_1, \dots, j_{d-1}]} = \text{true, push}$ $[j_1, \dots, j_{d-1}]$ into \mathcal{Q} .

Step 6: If Q is not empty, go to **Step 3**, else stop.

Step 7: Eliminate redundant halfspaces, compute region's facets. **Output:** \mathcal{G}_{τ} , τ -region's facets.



•






















0



Complexity

Notation:

- n number of points;
- d dimension;
- \sum_{u} number of relevant directions.

Algorithmic complexity in terms of time:

- Directional quantiles:
- Multiple-output regression quantiles: \sum_{u} is never larger than $O(n^d)$ \sum_{u} is on an average $O(n^{d-1})$
- Univariate projection quantiles: \sum_{u} is never larger than $\binom{n}{d-1}$

 $O(n \sum_{u})$ $O_{i} + O(n \sum_{u})$ $\implies O(n^{d+1})$ $\implies av. O(n^{d})$ $O_{i} + O(n \log n \sum_{u})$ $\implies O(n^{d} \log n)$

Relevant hyperplanes/facets

Logarithmized (average) ratio of the number of facets of the Tukey region to the number of relevant hyperplanes:



- Cases correspond to dimensions: p = 3, 4, 5, 6.
- Lines correspond to depths: $\tau = 0.025, 0.1, 0.2, 0.3$.
- Points correspond to sample sizes:
 n = 40, 80, 160, 320, 640, 1280, 2560, 5120.

Upper bound on the number of facets of $D_{\tau}^{T(n)}$



Upper bound on the number of facets of $D_{\tau}^{\mathcal{T}(n)}$



Upper bound on the number of facets of $D_{\tau}^{T(n)}$



Proposition

For a given data set $\mathbf{X} = {\mathbf{x}_1, ..., \mathbf{x}_n}$ in \mathbb{R}^d , the number of the (non-redundant) facets of the Tukey region $D_{\tau}^{T(n)}(\mathbf{X})$ is bounded from above by $2\binom{n}{d-1}$.

Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Tukey median (algorithm)

Input:
$$X = \{x_1, \dots, x_n\} \in \mathbb{R}^d$$
, $2 < d < n < \infty$.

Step 1: Initialize bounds on τ^* :

1.1: Compute coordinate-wise median
$$\mathbf{x}_0$$
.
1.2: Compute $d_0 = D^T(\mathbf{x}_0 | \mathbf{X})$.
1.3: Set $\tau_{low} = \max\left\{\frac{1}{n} \lceil \frac{n}{p+1} \rceil, d_0\right\}, \tau_{up} = \frac{1}{n} \lfloor \frac{n-p+2}{2} \rfloor + \frac{1}{n}^{**}$.

Step 2: Update bounds:

2.1: Let $\overline{\tau} = \frac{1}{n} \lfloor \frac{n(\tau_{low} + \tau_{up})}{2} \rfloor$, and compute the region $\mathcal{D}(\overline{\tau}) = D_{\overline{\tau}}^{T}(\mathbf{X})$. **2.2:** If $\mathcal{D}(\overline{\tau})$ does not exist then set $\tau_{up} = \overline{\tau}$, otherwise: calculate the barycenter \mathbf{c} of $\mathcal{D}(\overline{\tau})$ and set $\tau_{low} = D^{T}(\mathbf{c}|\mathbf{X})$. **2.3:** If $\tau_{low} < \tau_{up} - \frac{1}{n}$, then repeat **Step 2**, else stop. **Output:** The Tukey median $D_{\tau_{u}}^{T}(\mathbf{X})$.

**Liu, Luo, Zuo (2017). Some results on the computing of Tukey's halfspace median. Statistical Papers. https://doi.org/10.1007/s00362-017-0941-5.

Tukey median (illustration)



Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

The projection property

A depth D satisfies the projection property (Dyckerhoff, '04) if

$$D(\boldsymbol{z}|X) = \inf_{\boldsymbol{p} \in \mathbb{S}^{d-1}} D(\boldsymbol{p}^T \boldsymbol{z} | \boldsymbol{p}^T X).$$

Proposition (Dyckerhoff, '04)

- Let D be a depth satisfying projection property,
- ▶ $\mathbb{S}^{d-1} \to [0,\infty), \, \boldsymbol{\rho} \mapsto D(\boldsymbol{\rho}^T \boldsymbol{z} | \boldsymbol{\rho}^T \boldsymbol{X}) \text{ upper-semicontinuous,}$
- ▶ p₁, p₂, ... a sequence of independent identically uniformly on S^{d-1} distributed random vectors.

Then with probability one

d

$$\min_{1\leq i\leq m} D(\boldsymbol{p}_i^T\boldsymbol{z}|\boldsymbol{p}_i^T\boldsymbol{X}) \xrightarrow{m\to\infty} D(\boldsymbol{z}|\boldsymbol{X}).$$

Application to the Tukey depth:

$$\underbrace{\underbrace{D^{T}(\boldsymbol{z}|\boldsymbol{X})}_{\text{-variate depth}}}_{\text{exprime}} = \inf_{\boldsymbol{p} \in \mathbb{S}^{d-1}} \underbrace{D^{T}(\boldsymbol{p}^{T}\boldsymbol{z}|\boldsymbol{p}^{T}\boldsymbol{X})}_{\text{p}^{T}\boldsymbol{X}(\boldsymbol{p}^{T}\boldsymbol{z}), 1 - F_{\boldsymbol{p}^{T}\boldsymbol{X}}(\boldsymbol{p}^{T}\boldsymbol{z}^{-}) \}.$$

















Computation of projection depths: Remarks

- A number of depths (Mahalanobis, halfspace, zonoid, ect.) satisfy the assumptions of the preceding proposition.
- The projection property enables the approximate computation of depths even in high dimensions.
- ▶ Due to affine invariance one can restrict *p* to a hemisphere.
- Since a large number of univariate depths have to be computed there is need in efficient algorithms for them.
- Many univariate depths can be exactly computed with time complexity O(n).
- Computing the approximate depth leads to optimization of a (non-differentiable) function on the sphere S^{d-1}. So what does the function

$$\mathbb{S}^{d-1} \to \mathbb{R}, \, \boldsymbol{\rho} \mapsto D(\boldsymbol{\rho}^T \boldsymbol{z} | \boldsymbol{\rho}^T \boldsymbol{X})$$

look like? For example, does it have many local minima?

Illustration: The map $\boldsymbol{p} \mapsto D(\boldsymbol{p}^T \boldsymbol{z} | \boldsymbol{p}^T \boldsymbol{X})$ in \mathbb{R}^3



Contents

Data depth

Tukey depth: definition Tukey trimmed regions Applications

Computation of the Tukey depth

Theoretical background Algorithm and simulations

Computation of Tukey trimmed regions

Existing approaches The proposed algorithm Tukey median

Outlook for approximations

Conclusions

Some conclusions

- Developed a family of exact algorithms for computing Tukey depth in the Euclidean space of any dimension, which are implemented in R-package ddalpha accessible on CRAN.
- Exact algorithm for computing Tukey trimmed regions, implemented in R-package TukeyRegion.
- Ongoing work: development of approximate algorithms for computation of bf depths that satisfy the projection property.

Open challenges:

- Still high complexity and computation time scales exponentially with dimension.
- Precision challenge: too many potentially relevant hyperplanes, precision-constant problems.
- For approximation: adaptation of the algorithm(s) to optimization over the surface of a hypersphere.

Thank you for your attention! Questions?

- Rousseeuw, P.J. and Ruts, I. (1996). Algorithm AS 307: Bivariate location depth. *Journal of the Royal Statistical Society, Series C*, 45, 516-526.
- Ruts, I. and Rousseeuw, P.J. (1996). Computing depth contours of bivariate point clouds. *Computational Statistics and Data Analysis*, 23, 153-168.
- Rousseeuw, P.J. and Struyf, A. (1998). Computing location depth and regression depth in higher dimensions. *Statistics and Computing*, 8, 193-203.
- Hallin, M., Paindaveine, D., and Šiman, M. (2010). Multivariate quantiles and multiple-output regression quantiles: From L₁ optimization to halfspace depth. The Annals of Statistics, 38, 635-669.
- Paindaveine, D. and Šiman, M. (2011). On directional multiple-output quantile regression. *Journal of Multivariate Analysis*, 102, 193–212.
- Liu, X. and Zuo, Y. (2014). Computing halfspace depth and regression depth. Communications in Statistics - Simulation and Computation, 43, 969-985.
- Liu, X. (2017). Fast implementation of the Tukey depth. Computational Statistics, 32, 1395-1410.