

BS-MS Thesis

STATISTICAL DEPTH FUNCTIONS

Satvik Saha (*19MS154*)

Supervised by

Dr. Anirvan Chakraborty^{*}



*Department of Mathematics and Statistics
Indian Institute of Science Education and Research, Kolkata*

April 2024

^{*}Assistant Professor, Department of Mathematics and Statistics, Indian Institute of Science Education and Research, Kolkata

ABSTRACT

DEDICATION

DECLARATION

ACKNOWLEDGMENTS

CONTENTS

1	INTRODUCTION	6
1.1	Centrality vs Density	6
1.2	Nonparametric procedures	6
2	MULTIVARIATE DATA	7
2.1	Multivariate depth functions	8
2.1.1	The projection property	11
2.1.2	Continuity properties	11
2.2	Depth contours	13
2.3	Depth-Depth plots	14
2.4	Testing	15
2.5	Classification	20
2.6	Clustering	23
3	FUNCTIONAL DATA	25
3.1	Functional depth functions	27
3.1.1	Summary depths	27
3.1.2	Band depths	27
3.1.3	Half-region depths	28
3.2	Classification	29
3.2.1	Outlyingness matrices	29
3.2.2	Random projections	30
3.3	Outlier detection	31
3.3.1	Outliergrams	32
3.3.2	Centrality-Stability diagrams	33
3.3.3	MO-VO diagrams	35
3.4	Partially observed functional data	35
4	LOCAL DEPTH	39
4.1	Local depth regions	39
4.2	Regression based on local depth	40
5	CONCLUSION	42

Chapter 1

INTRODUCTION

1.1 Centrality vs Density

1.2 Nonparametric procedures

Chapter 2

MULTIVARIATE DATA

It is desirable for a depth function $D: \mathbb{R}^d \times \mathcal{F} \rightarrow \mathbb{R}$ to satisfy the following properties, described by Zuo and Serfling (2000).

P1. *Affine invariance.* For any random vector \mathbf{X} in \mathbb{R}^d , any $d \times d$ nonsingular matrix A , and any d -vector \mathbf{b} ,

$$D(A\mathbf{x} + \mathbf{b}, F_{A\mathbf{X}+\mathbf{b}}) = D(\mathbf{x}, F_{\mathbf{X}}). \quad (2.0.1)$$

This makes $D(\mathbf{x}, F_{\mathbf{X}})$ independent of the choice of coordinate system.

P2. *Maximality at center.* For any $F \in \mathcal{F}$ having ‘center’ $\boldsymbol{\theta}$,

$$D(\boldsymbol{\theta}, F) = \sup_{\mathbf{x} \in \mathbb{R}^d} D(\mathbf{x}, F). \quad (2.0.2)$$

This means that the deepest point coincides with some center of symmetry of the distribution F .

P3. *Monotonicity relative to deepest point.* For any $F \in \mathcal{F}$ having deepest point $\boldsymbol{\theta}$ and for $\alpha \in [0, 1]$,

$$D(\mathbf{x}, F) \leq D(\boldsymbol{\theta} + \alpha(\mathbf{x} - \boldsymbol{\theta}), F). \quad (2.0.3)$$

Thus, $D(\cdot, F)$ monotonically decreases along any ray pointing away from the deepest point.

P4. *Vanishing at infinity.* For any $F \in \mathcal{F}$,

$$D(\mathbf{x}, F) \rightarrow 0 \text{ as } \|\mathbf{x}\| \rightarrow \infty. \quad (2.0.4)$$

Furthermore, we demand that D be non-negative and bounded. Thus, we may assume hereon that D only takes values in $[0, 1]$.

The notion of a ‘center’ of a distribution in **P2** is typically described in terms of symmetry. We say that a random vector \mathbf{X} is *centrally symmetric* about $\boldsymbol{\theta} \in \mathbb{R}^d$

if $\mathbf{X} - \boldsymbol{\theta} \stackrel{d}{=} \boldsymbol{\theta} - \mathbf{X}$. Similarly, we say that \mathbf{X} is *angularly symmetric* about $\boldsymbol{\theta}$ if $(\mathbf{X} - \boldsymbol{\theta})/\|\mathbf{X} - \boldsymbol{\theta}\|$ is centrally symmetric about $\mathbf{0}$. An even more restrictive notion of symmetry is *spherical symmetry*, where we demand that $U(\mathbf{X} - \boldsymbol{\theta}) \stackrel{d}{=} \mathbf{X} - \boldsymbol{\theta}$ for every orthonormal matrix U . *Elliptical symmetry* requires that $V\mathbf{X}$ is spherically symmetric about $\boldsymbol{\theta}$ for some nonsingular matrix V . Finally, the weakest notions of symmetry discussed here is *halfspace symmetry*, where we impose $P(\mathbf{X} \in H) \geq 1/2$ for every closed halfspace in \mathbb{R}^d containing $\boldsymbol{\theta}$. Thus, the symmetries in decreasing order of strength are $S > E > C > A > H$.

Mosler and Mozharovskiy (2022, Table 2) provides a detailed summary of the properties satisfied by the depth functions discussed in the following section.

2.1 Multivariate depth functions

Definition 2.1.1 (Halfspace/Tukey depth). Denote the collection of all closed halfspaces in \mathbb{R}^d containing \mathbf{x} by $\mathcal{H}_{\mathbf{x}}$. The halfspace depth, or Tukey depth, is defined as

$$D_H(\mathbf{x}, F) = \inf_{H \in \mathcal{H}_{\mathbf{x}}} P_F(H). \quad (2.1.1)$$

Remark. If $F \in \mathcal{F}$ is supported on a convex region $C \subset \mathbb{R}^d$, then $D(\cdot, F)$ vanishes outside C .

Proposition 2.1.2. *The halfspace depth can be formulated as*

$$D_H(\mathbf{x}, F) = \inf_{\mathbf{v} \in S^{d-1}} P_{\mathbf{X} \sim F}(\langle \mathbf{v}, \mathbf{X} \rangle \leq \langle \mathbf{v}, \mathbf{x} \rangle). \quad (2.1.2)$$

Remark. When $d = 1$, the halfspace depth reduces to

$$D_H(x, F) = \min\{P_F(-\infty, x], P_F[x, \infty)\}. \quad (2.1.3)$$

Struyf and Rousseeuw (1999) showed that halfspace depth fully characterizes empirical distributions.

Theorem 2.1.3 (Struyf and Rousseeuw, 1999). *The empirical distribution of any dataset $\{\mathbf{X}_i\}_{i=1}^n \subset \mathbb{R}^d$ is uniquely determined by its halfspace depth function $D(\cdot, \hat{F}_n)$.*

Nagy (2020) offers a comprehensive overview of this characterization problem. Indeed, Nagy (2021) supplies examples of distinct probability distributions F_1, F_2 such that $D_H(\cdot, F_1) = D_H(\cdot, F_2)$. It can be shown that for an α -symmetric distribution F that $D_H(\mathbf{x}, F) = G(-\|\mathbf{x}\|_{\alpha^*})$, where G is the marginal distribution of the first component of $\mathbf{X} \sim F$ and $1/\alpha + 1/\alpha^* = 1$.

Definition 2.1.4 (Mahalanobis depth). Let $\mathbf{X} \sim F$ have mean $\boldsymbol{\mu}$ and covariance matrix Σ . The Mahalanobis depth is defined as

$$D_M(\mathbf{x}, F) = \left(1 + (\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)^{-1}. \quad (2.1.4)$$



Figure 2.1: Depth contours with respect to purple points. Darker contours have higher depth.

Remark. The mean and covariance in the above definition may be replaced with more robust estimates $\boldsymbol{\mu}^*$ and Σ^* , for instance using the minimum covariance determinant (MCD) method. The corresponding depth function is called the robust Mahalanobis depth.

Definition 2.1.5 (Spatial depth). The spatial depth is defined as

$$D_{Sp}(\mathbf{x}, F) = 1 - \left\| \mathbb{E}_{\mathbf{X} \sim F} \left[\frac{\mathbf{x} - \mathbf{X}}{\|\mathbf{x} - \mathbf{X}\|} \right] \right\|. \quad (2.1.5)$$

We use the convention $\mathbf{0}/0 = \mathbf{0}$.

Remark. Spatial depth defined as above does not obey **P1**. Indeed, spatial depth is only invariant under spherical transformations of the form $U\mathbf{X} + \mathbf{b}$ for orthonormal U . We may define an affine invariant version of spatial depth as

$$D_{AISP}(\mathbf{x}, F) = 1 - \left\| \mathbb{E}_{\mathbf{X} \sim F} \left[\frac{\Sigma^{-1/2}(\mathbf{x} - \mathbf{X})}{\sqrt{(\mathbf{x} - \mathbf{X})^\top \Sigma^{-1}(\mathbf{x} - \mathbf{X})}} \right] \right\|. \quad (2.1.6)$$

Remark. Nagy (2017) showed that spatial depth does not obey **P3**.

Lemma 2.1.6. *Spatial depth obeys P4, i.e. $D_{Sp}(\mathbf{x}, F) \rightarrow 0$ as $\|\mathbf{x}\| \rightarrow \infty$.*

Proof. Let $\epsilon > 0$, and let $M > 0$ such that $P_{\mathbf{X} \sim F}(\|\mathbf{X}\| > M) = \epsilon$. Denote $\mathbf{Y} = (\mathbf{x} - \mathbf{X})/\|\mathbf{x} - \mathbf{X}\|$, and observe that

$$\mathbb{E}_{\mathbf{X} \sim F}[\mathbf{Y}] = (1 - \epsilon) \mathbb{E}[\mathbf{Y} | \|\mathbf{X}\| \leq M] + \epsilon \mathbb{E}[\mathbf{Y} | \|\mathbf{X}\| > M] \quad (2.1.7)$$

Thus, using $\|\mathbf{Y}\| = 1$ and the reverse triangle inequality,

$$\|\mathbb{E}[\mathbf{Y}]\| \geq (1 - \epsilon) \|\mathbb{E}[\mathbf{Y} | \|\mathbf{X}\| \leq M]\| - \epsilon. \quad (2.1.8)$$

Let $\alpha = \arccos((1 - 2\epsilon)/(1 - \epsilon))$, and let $r_\alpha = M/\sin \alpha$. It follows that the ball $\{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| \leq M\}$ subtends an angle of at most 2α from any point \mathbf{x} such that $\|\mathbf{x}\| > r_\alpha$. This gives $\|\mathbb{E}[\mathbf{Y} | \|\mathbf{X}\| \leq M]\| \geq \cos \alpha$. Thus, for $\|\mathbf{x}\| > r_\alpha$,

$$\|\mathbb{E}[\mathbf{Y}]\| \geq (1 - \epsilon) \cos \alpha - \epsilon = 1 - 3\epsilon, \quad (2.1.9)$$

whence $D_{Sp}(\mathbf{x}, F) \leq 3\epsilon$. \square

Definition 2.1.7 (Projection depth). The projection depth is defined as

$$D_P(\mathbf{x}, F) = \left(1 + \sup_{\mathbf{v} \in S^{d-1}} \frac{|\langle \mathbf{v}, \mathbf{x} \rangle - \text{med}(\langle \mathbf{v}, \mathbf{X} \rangle)|}{\text{MAD}(\langle \mathbf{v}, \mathbf{X} \rangle)} \right)^{-1}, \quad \mathbf{X} \sim F. \quad (2.1.10)$$

Definition 2.1.8 (Simplicial depth). The simplicial depth is defined as

$$D_{Sim}(\mathbf{x}, F) = P_{\mathbf{X}_i \stackrel{\text{iid}}{\sim} F}(\mathbf{x} \in \text{conv}(\mathbf{X}_1, \dots, \mathbf{X}_{d+1})), \quad (2.1.11)$$

where $\text{conv}(\mathbf{x}_1, \dots, \mathbf{x}_{d+1})$ denotes the convex hull of $\{\mathbf{x}_1, \dots, \mathbf{x}_{d+1}\}$.

Definition 2.1.9 (Oja depth). The simplicial volume depth, or Oja depth, is defined as

$$D_{Oja}(\mathbf{X}, F) = \left(1 + \mathbb{E}_{\mathbf{X}_i \sim F} [\text{vol}(\text{conv}(\mathbf{x}, \mathbf{X}_1, \dots, \mathbf{X}_d))]\right)^{-1}. \quad (2.1.12)$$

Remark. Oja depth does not obey **P1**, since

$$\text{vol}(\text{conv}(A\mathbf{x}_1 + \mathbf{b}, \dots, A\mathbf{x}_{d+1} + \mathbf{b})) = |\det(A)| \text{vol}(\text{conv}(\mathbf{x}_1, \dots, \mathbf{x}_{d+1})). \quad (2.1.13)$$

Instead, we may define an affine invariant version of Oja depth as

$$D_{AIOja}(\mathbf{X}, F) = \left(1 + \mathbb{E}_{\mathbf{X}_i \sim F} \left[\frac{\text{vol}(\text{conv}(\mathbf{x}, \mathbf{X}_1, \dots, \mathbf{X}_d))}{\sqrt{\det(\Sigma)}} \right]\right)^{-1}, \quad (2.1.14)$$

where Σ is the covariance matrix of F .

2.1.1 The projection property

Definition 2.1.10 (Projection property). We say that a depth function D has the projection property if

$$D(\mathbf{x}, F_{\mathbf{X}}) = \inf_{\mathbf{v} \in S^{d-1}} D(\langle \mathbf{v}, \mathbf{x} \rangle, F_{\langle \mathbf{v}, \mathbf{X} \rangle}). \quad (2.1.15)$$

Depths which have this property can be approximated by calculating the univariate depths of the projected data along many directions \mathbf{v} .

Lemma 2.1.11. *The halfspace depth, Mahalanobis depth, and projection depth have the projection property.*

The halfspace depth in particular is often computationally challenging. Thus, the property motivates the definition of the random Tukey depth (Cuesta-Albertos & Nieto-Reyes, 2008a).

Definition 2.1.12 (Random Tukey depth). Let $\mathbf{v}_1, \dots, \mathbf{v}_n$ be a realization of an iid sample from $\mathcal{U}(S^{d-1})$. The random Tukey depth is defined as

$$D_{RT}(\mathbf{x}, F_{\mathbf{X}}) = \min_{1 \leq i \leq n} D_H(\langle \mathbf{v}_i, \mathbf{x} \rangle, F_{\langle \mathbf{v}_i, \mathbf{X} \rangle}). \quad (2.1.16)$$

2.1.2 Continuity properties

It is also desirable for a depth function to obey some notions of continuity.

C1. *Continuity in \mathbf{x} .*

$$D(\mathbf{x}_n, F) \rightarrow D(\mathbf{x}, F) \quad \text{when } \mathbf{x}_n \rightarrow \mathbf{x}. \quad (2.1.17)$$

C2. *Continuity in F .*

$$D(\mathbf{x}, F_n) \rightarrow D(\mathbf{x}, F) \text{ when } F_n \xrightarrow{d} F. \quad (2.1.18)$$

C3. *Uniform continuity.*

$$\sup_{\mathbf{x} \in G} |D(\mathbf{x}, F_n) - D(\mathbf{x}, F)| \rightarrow 0 \text{ when } F_n \xrightarrow{d} F. \quad (2.1.19)$$

Property **C1** is rarely satisfied without imposing some regularity conditions on F , such as absolute continuity. Property **C2** helps bridge the gap between the population and empirical versions of depth. Property **C3** becomes relevant when dealing with the convergence of depth contours.

The Mahalanobis depth is trivially continuous in \mathbf{x} , i.e. obeys **C1**. Furthermore, it also satisfies **C2** as long as F has a regular covariance matrix (Mosler & Mozharovskiy, 2022).

The halfspace depth also enjoys all three notions of continuity, under mild restrictions on F .

Theorem 2.1.13 (Mizera and Volau, 2002). *Let $F \in \mathcal{F}$ be such that the probability of every hyperplane in \mathbb{R}^d is zero, i.e. for all $\alpha \in \mathbb{R}$ and $\mathbf{v} \in S^{d-1}$,*

$$P_{\mathbf{X} \sim F}(\langle \mathbf{v}, \mathbf{X} \rangle = \alpha) = 0. \quad (2.1.20)$$

Then for $\mathbf{x}_n \rightarrow \mathbf{x}$ and $F_n \xrightarrow{d} F$, we have $D_H(\mathbf{x}_n, F_n) \rightarrow D_H(\mathbf{x}, F)$.

Remark. Equation 2.1.20 is satisfied whenever F is absolutely continuous.

Remark. It follows that if $F \in \mathcal{F}$ satisfies Equation 2.1.20, then the map $D_H(\cdot, F)$ is continuous.

Corollary 2.1.14. *Let $F \in \mathcal{F}$ satisfy Equation 2.1.20. Then, for $F_n \xrightarrow{d} F$, and compact $K \subset \mathbb{R}^d$,*

$$\sup_{\mathbf{x} \in K} |D_H(\mathbf{x}, F_n) - D_H(\mathbf{x}, F)| \rightarrow 0. \quad (2.1.21)$$

Proof. Denoting $g = D(\cdot, F)$, $g_n = D_H(\cdot, F_n)$, we have the continuity of g along with $g_n(\mathbf{x}_n) \rightarrow g(\mathbf{x})$ whenever $\mathbf{x}_n \rightarrow \mathbf{x}$ in K . If the given conclusion is false, we may pass to a subsequence of g_n and find $\epsilon > 0$ such that each $\sup_{\mathbf{x} \in K} |g_n(\mathbf{x}) - g(\mathbf{x})| \geq \epsilon$. Using the compactness of K , we pass to a further subsequence and find $\mathbf{x} \in K$ such that $\mathbf{x}_n \rightarrow \mathbf{x}$. This contradicts $|g_n(\mathbf{x}_n) - g(\mathbf{x}_n)| \geq \epsilon$. \square

Theorem 2.1.15. *Let $F \in \mathcal{F}$ satisfy Equation 2.1.20. Then, for $F_n \xrightarrow{d} F$,*

$$\sup_{\mathbf{x} \in \mathbb{R}^d} |D_H(\mathbf{x}, F_n) - D_H(\mathbf{x}, F)| \rightarrow 0. \quad (2.1.22)$$

Proof. Let $K_r = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| \leq r\}$ be a continuity set of F . Observe that $D_H(\mathbf{y}, F) \leq P_F(K_r^c)$ for $\mathbf{y} \in K_r^c$, hence

$$\sup_{\mathbf{y} \in K_r^c} |D_H(\mathbf{y}, F_n) - D_H(\mathbf{y}, F)| \leq P_{F_n}(K_r^c) + P_F(K_r^c). \quad (2.1.23)$$

As $n \rightarrow \infty$, we have $P_{F_n}(K_r^c) \rightarrow P_F(K_r^c) = p_r$ (say). Thus, denoting $\delta_n(X) = \sup_{\mathbf{x} \in X} |D_H(\mathbf{x}, F_n) - D_H(\mathbf{x}, F)|$, we have

$$\limsup_{n \rightarrow \infty} \delta_n(\mathbb{R}^d) \leq \lim_{n \rightarrow \infty} \delta_n(K_r) + \limsup_{n \rightarrow \infty} \delta_n(K_r^c) \quad (2.1.24)$$

$$\leq 0 + 2p_r. \quad (2.1.25)$$

Using $p_r \rightarrow 0$ as $r \rightarrow \infty$ completes the proof. \square

The spatial depth is similarly well behaved.

Theorem 2.1.16. *Spatial depth obeys C1 when F is non-atomic, as well as C2.*

Proof. Consider the spatial map

$$S_F: \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad \mathbf{x} \mapsto \mathbb{E}_{\mathbf{X} \sim F} \left[\frac{\mathbf{x} - \mathbf{X}}{\|\mathbf{x} - \mathbf{X}\|} \right]. \quad (2.1.26)$$

The Dominated Convergence Theorem guarantees the continuity of S_F , hence of $D_{Sp}(\cdot, F) = 1 - \|S_F(\cdot)\|$. Furthermore, if $F_n \xrightarrow{d} F$, we have $S_{F_n}(\mathbf{x}) \rightarrow S_F(\mathbf{x})$ by the Portmanteau Lemma for all $\mathbf{x} \in \mathbb{R}^d$. \square

Theorem 2.1.17 (Serfling, 2002). *For $F \in \mathcal{F}$ and compact $K \subset \mathbb{R}^d$,*

$$\sup_{\mathbf{x} \in K} |D_{Sp}(\mathbf{x}, \hat{F}_n) - D_{Sp}(\mathbf{x}, F)| \xrightarrow{a.s.} 0. \quad (2.1.27)$$

Remark. This result can be generalized from compact subsets K to the whole of \mathbb{R}^d , using Lemma 2.1.6 and arguments similar to the proof of Theorem 2.1.15.

2.2 Depth contours

The following definitions are adapted from Liu et al., 1999.

Definition 2.2.1. The contour of depth t is the set $\{\mathbf{x} \in \mathbb{R}^d : D(\mathbf{x}, F) = t\}$.

Definition 2.2.2. The region enclosed by the contour of depth t is the set

$$R_F(t) = \{\mathbf{x} \in \mathbb{R}^d : D(\mathbf{x}, F) > t\}. \quad (2.2.1)$$

Definition 2.2.3. The p -th central region is the set

$$C_F(p) = \bigcap_t \{R_F(t) : P_F(R_F(t)) \geq p\}. \quad (2.2.2)$$

Definition 2.2.4. The p -th level contour, or center-outward contour surface, is the set $Q_F(p) = \partial C_F(p)$.

Example 2.2.5. Consider $\mathcal{U}(B^d)$, i.e. the uniform distribution on the unit ball in \mathbb{R}^d . While there are no proper density contours to speak of, halfspace depth contours are concentric spheres centered at the origin, the deepest point. This illustrates how depth contours are more suited to indicating centrality than density contours.

Definition 2.2.6. Let $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{iid}}{\sim} F$. We introduce depth based order statistics $\mathbf{X}_{[1]}, \dots, \mathbf{X}_{[n]}$, which are a reordering of the sample in decreasing order of depth, i.e. $D(\mathbf{X}_{[1]}, F) \geq \dots \geq D(\mathbf{X}_{[n]}, F)$.

With this, given $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{iid}}{\sim} F$, the sample p -th central region is given by

$$C_{\hat{F}_n}(p) = \text{conv}(\mathbf{X}_{[1]}, \dots, \mathbf{X}_{[np]}). \quad (2.2.3)$$

2.3 Depth-Depth plots

Definition 2.3.1 (DD plot). Let F, G be two distributions on \mathbb{R}^d , and let D be a depth function. The Depth-Depth plot, also known as the DD plot, of F and G is given by

$$\text{DD}(F, G) = \{(D(\mathbf{z}, F), D(\mathbf{z}, G)) : \mathbf{z} \in \mathbb{R}^d\}. \quad (2.3.1)$$

Remark. The above definition generalizes naturally to involve more than two distributions on \mathbb{R}^d .

When the depth function D only takes values in $[0, 1]$, the DD plot is a subset of $[0, 1]^2$ and hence easily visualized. Clearly when $F = G$, the corresponding DD plot is confined to the diagonal $\{(t, t) : t \in [0, 1]\}$. However, when $d \geq 2$ and F, G are absolutely continuous, $\text{DD}(F, G)$ has non-zero area (Lebesgue measure) when $F \neq G$. Assuming that D is affine invariant, Liu et al. (1999) propose this area as an affine invariant measure of the discrepancy between F and G .

If the distributions F, G are unknown, we may use data samples $\mathcal{D}_F = \{\mathbf{X}_i\}$ and $\mathcal{D}_G = \{\mathbf{Y}_j\}$ where $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{iid}}{\sim} F$ and $\mathbf{Y}_1, \dots, \mathbf{Y}_m \stackrel{\text{iid}}{\sim} G$, then construct empirical distributions \hat{F}_n, \hat{G}_m . With this, we may examine the empirical DD plot

$$\text{DD}(\hat{F}_n, \hat{G}_m) = \{(D(\mathbf{z}, \hat{F}_n), D(\mathbf{z}, \hat{G}_m)) : \mathbf{z} \in \mathcal{D}_F \cup \mathcal{D}_G\}. \quad (2.3.2)$$

DD plots can be used as a diagnostic tool to detect differences in location and scale between two multivariate distributions.

1. If $F = G$, the points in $\text{DD}(\hat{F}_n, \hat{G}_m)$ stay close to the diagonal. See Figure 2.2.



Figure 2.2: Empirical DD plot using spatial depth, where both underlying distributions (bivariate normal) are identical. Observe how the points in the DD plot stay close to the diagonal black line.

2. If the same point \mathbf{z}_0 achieves maximum depths with respect to both distributions F and G , this indicates that \mathbf{z}_0 is their common center. See Figure 2.3.
3. Suppose that F and G have the same center. If the points in $\text{DD}(\hat{F}_n, \hat{G}_m)$ arch above the diagonal, i.e. the bulk of points are deeper in G than in F , this indicates that F has a greater spread than G . See Figure 2.4a.

Liu et al. (1999) also demonstrate the use of DD plots to detect differences in skewness and kurtosis. This tool is especially convenient since the DD plot is always two dimensional regardless of the dimension d of the sample points.

2.4 Testing

We are mainly interested in the two sample homogeneity test. Given samples from F and G , we wish to test the null hypothesis $H_0 : F = G$ against an alternate hypothesis that F and G differ in location or scale.

When F, G are distributions on \mathbb{R} , rank based tests such as the Wilcoxon rank-sum test or the Siegel-Tukey test are readily available. A very useful tool in this setting is the probability integral transform.

Proposition 2.4.1. *Let $\mathbf{X} \sim F$, and let the distribution F be continuous. Then, $F(\mathbf{X}) \sim \mathcal{U}[0, 1]$.*

Since $F(\mathbf{X}_j)$ has the same rank within $\{F(\mathbf{X}_i)\}$ as does \mathbf{X}_j within $\{\mathbf{X}_i\}$, the above result is the key towards establishing many distribution-free tests and procedures.

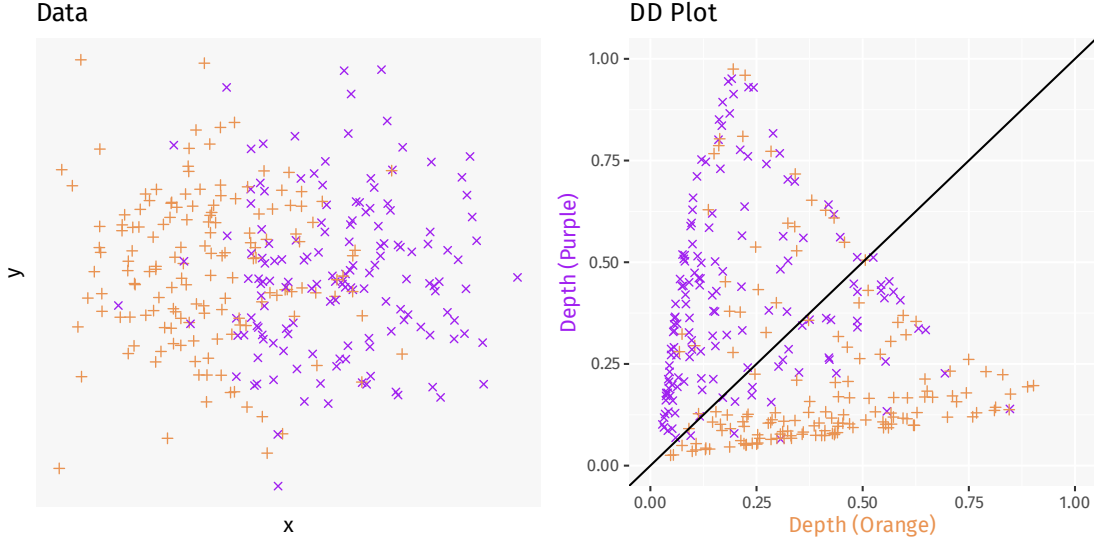


Figure 2.3: Empirical DD plot using spatial depth, where both underlying distributions (bivariate normal) differ only in location. Observe how most of the orange points fall in the lower triangle, while the purple ones fall in the upper triangle. The deepest point with respect to the orange distribution has fairly low depth with respect to the purple one, and vice versa.

In the multivariate setting, Liu and Singh (1993) use the following depth based analogue.

Definition 2.4.2. Denote

$$R(\mathbf{z}, F) = P_{\mathbf{X} \sim F}(D(\mathbf{X}, F) \leq D(\mathbf{z}, F)). \quad (2.4.1)$$

Note that in the empirical setting, $R(\mathbf{z}, \hat{F}_n)$ is simply the proportion of sample points $\{\mathbf{X}_i\}$ which are deeper in F than \mathbf{z} .

Proposition 2.4.3 (Liu and Singh, 1993). *Let $\mathbf{X} \sim F$, and let the distribution of $D(\mathbf{X}, F)$ be continuous. Then, $R(\mathbf{X}, F) \sim \mathcal{U}[0, 1]$.*

Definition 2.4.4. Denote the quality index

$$Q(F, G) = P(D(\mathbf{X}, F) \leq D(\mathbf{Y}, F) \mid \mathbf{X} \sim F, \mathbf{Y} \sim G). \quad (2.4.2)$$

Note that $Q(F, G)$ and $Q(G, F)$ are not necessarily the same. We may also write

$$Q(F, G) = \mathbb{E}_{\mathbf{Y} \sim G}[R(\mathbf{Y}, F)]. \quad (2.4.3)$$

It is clear that $Q(F, G) = 1/2$ when $F = G$. It can be shown under special circumstances that $Q(F, G) < 1/2$ if F, G differ in terms of location or scale. This will form the basis of our testing scheme, with $H_0 : F = G$ versus $H_A : Q(F, G) < 1/2$.

Here, we restrict our attention to elliptical distributions on \mathbb{R}^d .



Figure 2.4: Empirical DD plot using spatial depth, where both underlying distributions (bivariate normal) differ only in scale. In (a), observe how the points remain in the upper triangle in the DD plot. In (b), observe how there are more orange points in the lower triangle, and more purple points in the upper triangle in the DD plot, especially in the region close to the origin.

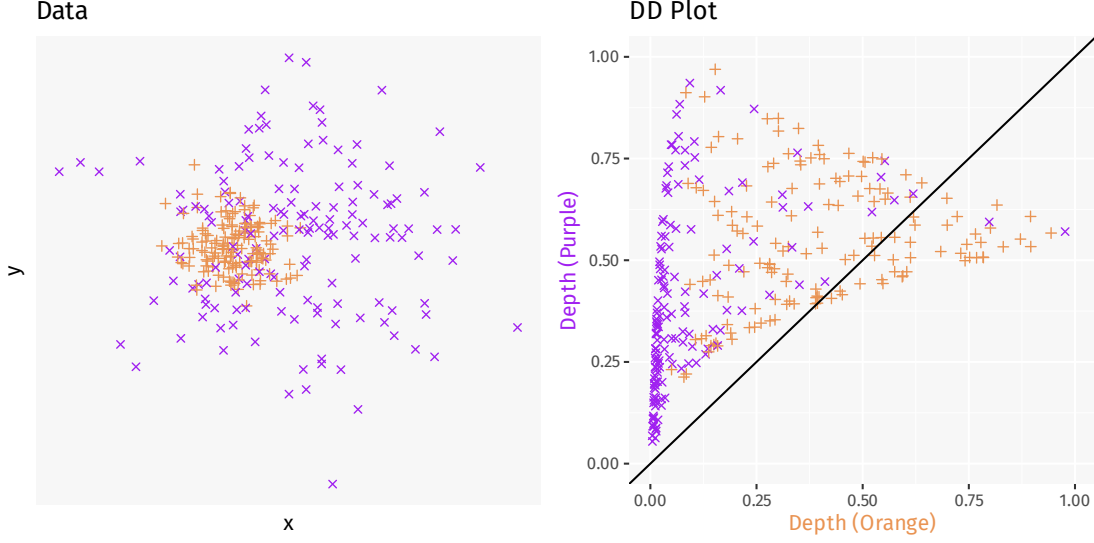


Figure 2.5: Empirical DD plot using spatial depth, where both underlying distributions (bivariate normal) differ in both location and scale. Observe that there is a clear separation between the orange and purple points in the DD plot, although not about the diagonal line.

Definition 2.4.5 (Elliptical distributions). We say that a distribution is elliptical if it has a density of the form

$$f(\mathbf{x}) = c |\Sigma|^{-1/2} h((\mathbf{x} - \boldsymbol{\mu})^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})) \quad (2.4.4)$$

for some non-increasing function h . This is denoted by $\text{Ell}(h; \boldsymbol{\mu}, \Sigma)$.

The quality index obeys the following properties.

Proposition 2.4.6 (Liu and Singh, 1993). Let $F \sim \text{Ell}(h; \boldsymbol{\mu}_1, \Sigma_1)$ and $G \sim \text{Ell}(h; \boldsymbol{\mu}_2, \Sigma_2)$ where $\Sigma_1 - \Sigma_2$ is positive definite. Further suppose that $D(\cdot, F)$ has the affine invariance and monotonicity properties. Then, $Q(F, G) \leq 1/2$ decreases monotonically as $\boldsymbol{\mu}_2$ is moved away from $\boldsymbol{\mu}_1$ along any line.

Proposition 2.4.7 (Liu and Singh, 1993). Let $F \sim \text{Ell}(h; \boldsymbol{\mu}, \Sigma_1)$ and $G \sim \text{Ell}(h; \boldsymbol{\mu}, \Sigma_2)$ where $\Sigma_1 - \Sigma_2$ is positive definite. Consider Huber's contamination of the form

$$G_\alpha = (1 - \alpha)F + \alpha G \quad (2.4.5)$$

where $0 \leq \alpha \leq 1$. Then, $Q(F, G_\alpha)$ decreases monotonically as α increases.

This motivates a modified Wilcoxon rank-sum test in the multivariate setting, using the quality index $Q(F, G)$. Let $\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{\text{iid}}{\sim} F$, and $\mathbf{Y}_1, \dots, \mathbf{Y}_m \stackrel{\text{iid}}{\sim} G$. Since $R(\cdot, F), Q(F, \cdot)$ depend on $D(\cdot, F)$, the latter has to be approximated using $D(\cdot, \hat{F}_{n_0})$, where \hat{F}_{n_0} is based on a (fairly large) additional sample $\mathbf{Z}_1, \dots, \mathbf{Z}_{n_0} \stackrel{\text{iid}}{\sim} F$, with $n_0 \gg n, m$. With this, we compute

$$R(\cdot, \hat{F}_{n_0}) = \frac{1}{n_0} \sum_{i=1}^{n_0} \mathbf{1}(D(\mathbf{Z}_i, \hat{F}_{n_0}) \leq D(\cdot, \hat{F}_{n_0})). \quad (2.4.6)$$

Assign ranks $1, \dots, n+m$ to the arranged values $R(\mathbf{X}_i, \hat{F}_{n_0}), R(\mathbf{Y}_j, \hat{F}_{n_0})$ (ascending order), and define W to be the sum of ranks of the $R(\mathbf{Y}_j, \hat{F}_{n_0})$. If necessary, break ties at random. Under the null hypothesis $F = G$, it is clear that W has the same distribution as the sum of m numbers drawn without replacement from $\{1, \dots, n+m\}$. Under the alternate hypothesis $Q(F, G) < 1/2$, the ranks of $R(\mathbf{Y}_j, \hat{F}_{n_0})$ will tend to be lower on average, making W smaller.

Theorem 2.4.8 (Liu and Singh, 1993). *Let $H_{n,m}$ be the distribution of the sum of m numbers drawn randomly without replacement from $\{1, \dots, n+m\}$. Suppose that F admits a density function f . Under the null hypothesis $F = G$, we have $W \sim H_{n,m}$.*

It is also possible to approximate $Q(F, G)$ more directly via $Q(\hat{F}_n, \hat{G}_m)$ and perform our test this way. This sidesteps the need for the ‘reference’ sample $\mathbf{Z}_1, \dots, \mathbf{Z}_{n_0} \stackrel{\text{iid}}{\sim} F$. Note that

$$Q(\hat{F}_n, \hat{G}_m) = \frac{1}{m} \sum_{j=1}^m R(\mathbf{Y}_j, \hat{F}_n) = \frac{1}{nm} \sum_{i,j} \mathbf{1}(D(\mathbf{X}_i, \hat{F}_n) \leq D(\mathbf{Y}_j, \hat{F}_n)). \quad (2.4.7)$$

This estimate is indeed consistent under mild assumptions.

Theorem 2.4.9 (Liu and Singh, 1993). *Suppose that the distribution of $D(\mathbf{Y}, F)$ is continuous where $\mathbf{Y} \sim G$, and that*

$$\sup_{\mathbf{z} \in \mathbb{R}^d} |D(\mathbf{z}, \hat{F}_n) - D(\mathbf{z}, F)| \xrightarrow{a.s.} 0. \quad (2.4.8)$$

Then, $Q(\hat{F}_n, \hat{G}_n) \xrightarrow{a.s.} Q(F, G)$ as $\min\{n, m\} \rightarrow \infty$.

This allows us to determine the asymptotic null distribution of $Q(\hat{F}_n, \hat{G}_m)$.

Theorem 2.4.10 (Liu and Singh, 1993). *Let F be absolutely continuous, such that $\mathbb{E}_{\mathbf{X} \sim F} \|\mathbf{X}\|^4 < \infty$. Using Mahalanobis depth to define Q , we have*

$$S(\hat{F}_n, \hat{G}_m) = \left[\frac{1}{12} \left(\frac{1}{n} + \frac{1}{m} \right) \right]^{-1/2} \left[Q(\hat{F}_n, \hat{G}_m) - \frac{1}{2} \right] \xrightarrow{d} \mathcal{N}(0, 1) \quad (2.4.9)$$

as $\min\{n, m\} \rightarrow \infty$, under the null hypothesis $F = G$.

Later, Zuo and He (2006) show that under certain mild regularity conditions, the above asymptotic convergence can be extended to a broader class of depth functions, without the assumption $F = G$. They demonstrate that

$$\left[\frac{\sigma_{GF}^2}{n} + \frac{\sigma_{FG}^2}{m} \right]^{-1/2} \left[Q(\hat{F}_n, \hat{G}_m) - Q(F, G) \right] \xrightarrow{d} \mathcal{N}(0, 1), \quad (2.4.10)$$

where

$$\sigma_{FG}^2 = \int P_{\mathbf{X} \sim F}^2(D(\mathbf{X}, F) \leq D(\mathbf{y}, F)) dG(\mathbf{y}) - Q^2(F, G), \quad (2.4.11)$$

$$\sigma_{GF}^2 = \int P_{\mathbf{Y} \sim G}^2(D(\mathbf{x}, F) \leq D(\mathbf{Y}, F)) dF(\mathbf{x}) - Q^2(F, G). \quad (2.4.12)$$

Observe that given two samples, we have a choice between using $Q(\hat{F}_n, \hat{G}_m)$ or $Q(\hat{G}_m, \hat{F}_n)$. It may be advantageous to use the sample with a greater number of observations as the reference distribution. Shi et al. (2023) propose a weighted combination of the form

$$W_{n,m}^\alpha = \alpha S(\hat{F}_n, \hat{G}_m)^2 + (1 - \alpha) S(\hat{G}_m, \hat{F}_n)^2 \quad (2.4.13)$$

for $\alpha \in [0, 1]$, or a maximum

$$M_{n,m} = \max\{S(\hat{F}_n, \hat{G}_m)^2, S(\hat{G}_m, \hat{F}_n)^2\}. \quad (2.4.14)$$

Under similar assumptions, they show that both $W_{n,m}^\alpha \xrightarrow{d} \chi_1^2$ and $M_{n,m} \xrightarrow{d} \chi_1^2$ as $\min\{n, m\} \rightarrow \infty$ and n/m converges to a positive constant, under the null hypothesis $F = G$.

2.5 Classification

The k -class classification task involves assigning an observation \mathbf{x} to one of k populations, described by distributions F_i for $1 \leq i \leq k$. The populations may also be associated with prior probabilities π_i .

Definition 2.5.1 (Classifier). A classifier is a map $\hat{l}: \mathbb{R}^d \rightarrow \{1, \dots, k\}$.

Example 2.5.2 (Bayes classifier). Suppose that the population densities f_i for each $1 \leq i \leq k$ are known. The Bayes classifier assigns \mathbf{x} to the \hat{l}_B -th population where

$$\hat{l}_B(\mathbf{x}) = \arg \max_{1 \leq i \leq k} \pi_i f_i(\mathbf{x}). \quad (2.5.1)$$

One way of measuring the performance of a classifier (given the population distributions and their priors) is by measuring its average misclassification rate.

Definition 2.5.3 (Average misclassification rate). The average misclassification rate of a classifier \hat{l} is given by

$$\Delta(\hat{l}) = \sum_{i=1}^k \pi_i P_{\mathbf{X} \sim F_i}(\hat{l}(\mathbf{X}) \neq i). \quad (2.5.2)$$

Proposition 2.5.4. *The Bayes classifier has the lowest possible average misclassification rate. This is known as the optimal Bayes risk, denoted Δ_B .*

The simplest depth based classifier is the maximum depth classifier (Ghosh & Chaudhuri, 2005).

Example 2.5.5 (Maximum depth classifier). Suppose that the prior probabilities π_i are equal. The maximum depth classifier \hat{l}_D for a choice of depth function D is described by

$$\hat{l}_D(\mathbf{x}) = \arg \max_{1 \leq i \leq k} D(\mathbf{x}, F_i). \quad (2.5.3)$$

In practice, instead of having direct access to the population distributions F_i , we have typically deal with labeled training data

$$\mathcal{D} = \{(\mathbf{x}_{ij}, i)\} \subset \mathbb{R}^d \times \{1, \dots, k\}, \quad (2.5.4)$$

where $\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i}$ is an instance of an iid sample from F_i for each $1 \leq i \leq k$. The empirical maximum depth classifier simply replaces the population distributions F_i with their empirical counterparts \hat{F}_i determined by $\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i}$. Thus, it is given by

$$\hat{l}_D(\mathbf{x}) = \arg \max_{1 \leq i \leq k} D(\mathbf{x}, \hat{F}_i). \quad (2.5.5)$$

Under certain restrictions, this classifier becomes asymptotically optimal in the following sense.

Theorem 2.5.6 (Ghosh and Chaudhuri, 2005). *Suppose that the population density functions f_i are elliptically symmetric, with $f_i(\mathbf{x}) = g(\mathbf{x} - \boldsymbol{\mu}_i)$ for parameters $\boldsymbol{\mu}_i$ and a density function g such that $g(k\mathbf{x}) \leq g(\mathbf{x})$ for every \mathbf{x} and $k > 1$. Further suppose that the priors on the populations are equal, and the depth function D is one of HD, SD, MJD, PD. Then, $\Delta(\hat{l}_D) \rightarrow \Delta_B$ as $\min\{n_1, \dots, n_k\} \rightarrow \infty$.*

Note that this result deals with elliptic population densities differing only in location. Relax this assumption, and instead suppose that $f_i \sim \text{Ell}(h_i; \boldsymbol{\mu}_i, \Sigma)$, i.e.

$$f_i(\mathbf{x}) = c_i |\Sigma|^{-1/2} h_i((\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)) \quad (2.5.6)$$

for strictly decreasing h_i , and that the depths can be expressed as $D(\cdot, F_i) = l_i(f_i(\cdot))$ for strictly increasing functions l_i . It follows that the Bayes decision rule can be reformulated as

$$\pi_i f_i(\mathbf{x}) > \pi_j f_j(\mathbf{x}) \iff D(\mathbf{x}, F_i) > r_{ij}(D(\mathbf{x}, F_j)) \quad (2.5.7)$$

for some real increasing function r_{ij} . Using this observation, the DD classifier (Li et al., 2012) picks separating functions r_{ij} which best classify the training data \mathcal{D} .

Definition 2.5.7 (Empirical misclassification rate). The empirical misclassification rate of a classifier \hat{l} , with respect to data \mathcal{D} , is given by

$$\hat{\Delta}(\hat{l}) = \sum_{i=1}^k \frac{\pi_i}{n_i} \sum_{j=1}^{n_i} \mathbf{1}(\hat{l}(\mathbf{x}_{ij}) \neq i). \quad (2.5.8)$$

Definition 2.5.8 (DD classifier). Suppose that $k = 2$, that D is a depth function, and that $r: [0, 1] \rightarrow [0, 1]$ is an increasing function. The DD classifier $\hat{l}_{D,r}$ is given by

$$\hat{l}_{D,r}(\mathbf{x}) = \begin{cases} 1, & \text{if } D(\mathbf{x}, F_2) \leq r(D(\mathbf{x}, F_1)), \\ 2, & \text{if } D(\mathbf{x}, F_2) > r(D(\mathbf{x}, F_1)). \end{cases} \quad (2.5.9)$$

The empirical DD classifier $\hat{l}_{D,\hat{r}}$ replaces F_i by their empirical counterparts \hat{F}_i . Here, the separating curve \hat{r} is chosen from a family Γ so as to minimize the empirical misclassification rate, i.e.

$$\hat{r} = \arg \min_{r \in \Gamma} \hat{\Delta}(\hat{l}_{D,r}). \quad (2.5.10)$$

Remark. The maximum depth classifier $\hat{\iota}_D$ is simply the DD classifier $\hat{\iota}_{D,\text{id}}$, where $\text{id}(x) = x$. Figure 2.5 clearly illustrates how this choice of separating function may not always be appropriate.

Li et al. (2012) show that under certain restrictions, the empirical DD classifier is asymptotically equivalent to the Bayes rule. We give one such instance below.

Lemma 2.5.9. *Suppose that the following conditions hold.*

1. Γ is the class of polynomial functions on $[0, 1]$.
2. The depth functions $D(\cdot, F_i)$ are continuous.
3. As $\min\{n_1, n_2\} \rightarrow \infty$, we have for each $i \in \{1, 2\}$,

$$\sup_{\mathbf{z} \in \mathbb{R}^d} |D(\mathbf{z}, \hat{F}_i) - D(\mathbf{z}, F_i)| \xrightarrow{a.s.} 0. \quad (2.5.11)$$

4. The distributions F_i are elliptical and satisfy for all $\delta \in \mathbb{R}$

$$P_{\mathbf{Z} \sim F_i}(D(\mathbf{Z}, F_i) = \delta) = 0. \quad (2.5.12)$$

Then, $\Delta(\hat{\iota}_{D,\hat{r}}) \rightarrow \Delta_B$ as $\min\{n_1, n_2\} \rightarrow \infty$.

In all the depth based classifiers we have seen so far, the classification rule depends on the observation \mathbf{x} only through the depths $D(\mathbf{x}, F_i)$. Thus, we are motivated to define the following transformation from \mathbb{R}^d to a depth feature space.

Definition 2.5.10. The depth feature vector \mathbf{x}^D of an observation \mathbf{x} , with respect to the population distributions F_i and a choice of depth function D , is defined as

$$\mathbf{x}^D = (D(\mathbf{x}, F_1), \dots, D(\mathbf{x}, F_k)). \quad (2.5.13)$$

Remark. The graph

$$\text{DD}(F_1, \dots, F_k) = \{\mathbf{x}^D : \mathbf{x} \in \mathbb{R}^d\} \quad (2.5.14)$$

is the analogue of the **DD plot**, with k distributions.

Assuming that the depth function D only takes values in $[0, 1]$, the map $\mathbf{x} \mapsto \mathbf{x}^D$ takes values in $[0, 1]^k$, regardless of the dimensionality of the original vector \mathbf{x} . With this, the maximum depth classification rule can be expressed as

$$\hat{\iota}_D(\mathbf{x}) = i \iff \mathbf{x}^D \in R_i^D = \{\mathbf{y} \in [0, 1]^k : y_i = \max_j y_j\}. \quad (2.5.15)$$

Indeed, any partition of the unit cube $[0, 1]^k$ into k decision regions R_i^D gives rise to a depth based classifier. The DD classifier achieves this by using an increasing separating function r to partition $[0, 1]^2$. Furthermore, $r \in \Gamma$ is chosen so as to best separate the training data \mathcal{D} transformed into the depth feature space. However, we can in principle use the transformed training data

$$\mathcal{D}^D = \{(\mathbf{x}_{ij}^D, i)\} \subset [0, 1]^k \times \{1, \dots, k\} \quad (2.5.16)$$

along with any multivariate classification algorithm (LDA, QDA, k NN, GLM, etc) to devise suitable decision regions. This is the basis of the DD^G classifier (Cuesta-Albertos et al., 2017).

2.6 Clustering

The unsupervised clustering grouping a collection of observations, such that points within the same group are more similar to each other than those from different groups.

Definition 2.6.1 (Clustering). Given observations $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$, a clustering assignment is a choice of a partition I_1, \dots, I_K of $\{1, \dots, N\}$.

With this notation, the k -th cluster consists of the points $\{\mathbf{x}_i\}_{i \in I_k}$. A good cluster assignment is one that maximizes similarity within clusters, as well as dissimilarity between clusters. Thus, the problem of clustering can be framed as the optimization of some objective function which combines these notions of similarity and dissimilarity. A simple algorithm such as the K -means clustering seeks to minimize

$$\{I_1, \dots, I_K\} \mapsto \frac{1}{N} \sum_{k=1}^K \sum_{i \in I_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2, \quad (2.6.1)$$

the average sum of square distances between each point and its cluster mean

$$\boldsymbol{\mu}_k = \frac{1}{|I_k|} \sum_{i \in I_k} \mathbf{x}_i = \arg \min_{\boldsymbol{\mu} \in \mathbb{R}^d} \sum_{i \in I_k} \|\mathbf{x}_i - \boldsymbol{\mu}\|^2. \quad (2.6.2)$$

Jörnsten (2004) proposes a depth based approach to this problem, by examining the depth of a point within its cluster, relative to its depth within the best competing cluster.

In this section, we will abbreviate $D_k(\mathbf{x}) = D(\mathbf{x}, \hat{F}_{I_k})$, i.e. the empirical depth of \mathbf{x} with respect to the points in the k -th cluster. Jörnsten (2004) chooses L_1 depth, the empirical version of spatial depth.

Definition 2.6.2. The within cluster depth of \mathbf{x}_i is $D_i^w = D_k(\mathbf{x}_i)$, where $i \in I_k$.

To deal with dissimilarity between clusters, we represent each cluster by its L_1 -median.

Definition 2.6.3 (L_1 -median). The L_1 -median of the k -th cluster is given by

$$\boldsymbol{\theta}_k = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \sum_{i \in I_k} \|\mathbf{x}_i - \boldsymbol{\theta}\|. \quad (2.6.3)$$

Definition 2.6.4. The between cluster depth of \mathbf{x}_i is $D_i^b = D_\ell(\mathbf{x}_i)$, where

$$\ell = \arg \min_{k: i \notin I_k} \|\mathbf{x}_i - \boldsymbol{\theta}_k\|. \quad (2.6.4)$$

In other words, the between cluster depth of \mathbf{x}_i is its depth within the best competing cluster.

Definition 2.6.5 (Relative depth). The relative depth of \mathbf{x}_i is $\text{ReD}_i = D_i^w - D_i^b$.

A point \mathbf{x}_i is *well clustered* if ReD_i is very high, i.e. it is deep within its own cluster, and has low depth with respect to its next best competing cluster. Thus, to obtain a good clustering, we may choose to maximize the objective function

$$\{I_1, \dots, I_K\} \mapsto \frac{1}{N} \sum_{k=1}^K \sum_{i \in I_k} \text{ReD}_i, \quad (2.6.5)$$

which is simply the average relative depth. This maximization can be achieved iteratively, starting with a random cluster assignment and reassigning a subset of observations with low ReD_i to their nearest competing clusters. The reassignment is accepted if the objective function increases, and the process is repeated. Jörnsten (2004) also suggests the use of simulated annealing to overcome the problem of getting trapped in local maxima. Here, the reassignment is accepted with some probability $P(\beta, \delta)$ where δ is the change in the objective function value, even if the objective function decreases at that step. $P(\beta, \delta)$ is chosen to decrease with increasing β and δ . The tuning parameter β can be increased every iteration so that the probability of accepting poorer clustering assignments drops to zero eventually.

Another notion of similarity and dissimilarity involves *silhouette width*.

Definition 2.6.6 (Silhouette width). Denote the average distance of \mathbf{z} from points in the k -th cluster not equal to \mathbf{z} by

$$\bar{d}_k(\mathbf{z}) = \frac{1}{|\{i \in I_k : \mathbf{x}_i \neq \mathbf{z}\}|} \sum_{\substack{i \in I_k \\ \mathbf{x}_i \neq \mathbf{z}}} \|\mathbf{x}_i - \mathbf{z}\|. \quad (2.6.6)$$

The silhouette width of \mathbf{x}_i where $i \in I_k$ is given by

$$\text{Sil}_i = \frac{b_i - a_i}{\max\{a_i, b_i\}}, \quad a_i = \bar{d}_k(\mathbf{x}_i), \quad b_i = \min_{\ell \neq k} \bar{d}_\ell(\mathbf{x}_i). \quad (2.6.7)$$

It has been observed that the silhouette width is greatly affected by differences in scale between clusters, while the relative depth is not. An objective function of the form

$$\{I_1, \dots, I_K\} \mapsto \frac{1}{N} \sum_{k=1}^K \sum_{i \in I_k} (1 - \lambda) \text{Sil}_i + \lambda \text{ReD}_i \quad (2.6.8)$$

may be used to combine both notions. Here, $\lambda \in [0, 1]$ controls the influence of the relative depth. It seems that small values of λ encourages equal scale clusters, while large values of λ allows unequal scale clusters. Thus, λ may be tuned accordingly to favour these different kinds of clustering assignments.

Chapter 3

FUNCTIONAL DATA

Consider a class of functions \mathcal{X} of the form $\mathbf{x}: [0, 1] \rightarrow \mathbb{R}^d$, equipped with a norm, and let \mathcal{F} be a suitable class of distributions on \mathcal{X} . Typically, we choose \mathcal{X} to be either $L_2[0, 1]$ or $\mathcal{C}[0, 1]$.

It is desirable for a depth function $D: \mathcal{X} \times \mathcal{F} \rightarrow \mathbb{R}$ to satisfy the following properties (Gijbels & Nagy, 2017).

P0. *Non-degeneracy.* For $F \in \mathcal{F}$,

$$\inf_{\mathbf{x} \in \mathcal{X}} D(\mathbf{x}, F) < \sup_{\mathbf{x} \in \mathcal{X}} D(\mathbf{x}, F). \quad (3.0.1)$$

The property **P0** requires some care; the natural functional analogue to half-space/Tukey depth when \mathcal{X} is a Banach space,

$$D_H(\mathbf{x}, F) = \inf_{\mathbf{v}^* \in \mathcal{X}^*} P_{\mathbf{X} \sim F}(\mathbf{v}^* \mathbf{X} \leq \mathbf{v}^* \mathbf{x}), \quad (3.0.2)$$

turns out to be degenerate for a wide class of distributions \mathcal{F} (Chakraborty & Chaudhuri, 2014a). For example, when $\mathcal{X} = C[0, 1]$ with the supremum norm and \mathbf{X} is a Gaussian process with a positive definite covariance kernel, we have $D_H(\cdot, F_{\mathbf{X}}) = 0$ almost surely. A similar result holds for the analogue to the projection depth. However, neither the functional random Tukey depth nor the functional spatial depth

$$D_S(\mathbf{x}, F_{\mathbf{X}}) = 1 - \left\| \mathbb{E}_{\mathbf{X} \sim F} \left[\frac{\mathbf{x} - \mathbf{X}}{\|\mathbf{x} - \mathbf{X}\|_2} \right] \right\|_2 \quad (3.0.3)$$

suffer this deficiency when $\mathcal{X} = L_2$ (Cuesta-Albertos & Nieto-Reyes, 2008a; Gijbels & Nagy, 2017).

The notion of affine invariance in **P1** can be generalized in many ways; Gijbels and Nagy (2017) recommend the following.

P1S. *Scalar-affine invariance.* For $a, b \in \mathbb{R}$ with a non-zero and $\mathbf{x} \in \mathcal{X}$,

$$D(a\mathbf{x} + b, F_{a\mathbf{X}+b}) = D(\mathbf{x}, F_{\mathbf{X}}). \quad (3.0.4)$$

P1F. *Function-affine invariance.* For $a, b, x \in \mathcal{X}$ with $ax \in \mathcal{X}$,

$$D(ax + b, F_{aX+b}) = D(x, F_X). \quad (3.0.5)$$

When generalizing **P2**, we must first define a notion of symmetry of $F \in \mathcal{F}$. To this end, we say that F_X is symmetric about $\theta \in \mathcal{X}$ if for all $\varphi \in \mathcal{X}^*$, we have $\varphi(\mathbf{X})$ is symmetric about $\varphi(\theta)$. Again, we are free to choose our notion of univariate symmetry for $\varphi(\mathbf{X})$. Gijbels and Nagy (2017) consider central and halfspace symmetry.

P2C. *Maximality at center of central symmetry.* Any centrally symmetric $F \in \mathcal{F}$ is symmetric about $\theta \in \mathcal{X}$ if and only if $D(\theta, F) = \sup_{\mathbf{x} \in \mathcal{X}} D(\mathbf{x}, F)$.

P2H. *Maximality at center of halfspace symmetry.* Any halfspace symmetric $F \in \mathcal{F}$ is symmetric about $\theta \in \mathcal{X}$ if and only if $D(\theta, F) = \sup_{\mathbf{x} \in \mathcal{X}} D(\mathbf{x}, F)$.

Earlier, Nieto-Reyes and Battey (2016) proposed the following variant of **P2**.

P2G. *Maximality at Gaussian process mean.* For a zero-mean, stationary, almost surely continuous Gaussian process $F \in \mathcal{F}$, we have $D(\theta, F) = \sup_{\mathbf{x} \in \mathcal{X}} D(\mathbf{x}, F)$ where θ is the zero mean function.

The above notions of maximality at the center are **P2H** > **P2C** > **P2G** in order of strength.

The properties **P3** and **P4** have straightforward generalization.

P3D. *Monotonicity relative to deepest point.* For $F \in \mathcal{F}$ such that $D(\theta, F) = \sup_{\mathbf{x} \in \mathcal{X}} D(\mathbf{x}, F)$, we have for $\alpha \in [0, 1]$,

$$D(\mathbf{x}, F) \leq D(\theta + \alpha(\mathbf{x} - \theta), F). \quad (3.0.6)$$

P4V. *Vanishing at infinity.* For any $F \in \mathcal{F}$,

$$D(\mathbf{x}, F) \rightarrow 0 \text{ as } \|\mathbf{x}\| \rightarrow \infty. \quad (3.0.7)$$

Nieto-Reyes and Battey (2016) and Gijbels and Nagy (2017) also deal with the notions of continuity in F . Let $d_{\mathcal{F}}$ metrize the topology of weak convergence in \mathcal{F} .

C2W. *Weak continuity in F .* For all $\epsilon > 0$ and $F \in \mathcal{F}$, there exists $\delta > 0$ such that for all $G \in \mathcal{F}$ such that $d_{\mathcal{F}}(F, G) < \delta$, we have $|D(\mathbf{x}, F) - D(\mathbf{x}, G)| < \epsilon$, F -almost surely.

C2U. *Uniform continuity in F .* For all $\epsilon > 0$ and $F \in \mathcal{F}$, there exists $\delta > 0$ such that for all $G \in \mathcal{F}$ such that $d_{\mathcal{F}}(F, G) < \delta$, we have $\sup_{\mathbf{x} \in \mathcal{X}} |D(\mathbf{x}, F) - D(\mathbf{x}, G)| < \epsilon$.

Gijbels and Nagy (2017, Table 1) provides a detailed summary of which of these properties are satisfied by the depth functions discussed in the following section.

3.1 Functional depth functions

3.1.1 Summary depths

Let D be a multivariate depth function. We can use this to define the depth of a curve \mathbf{x} by first computing the multivariate D -depth of each time slice $\mathbf{x}(t)$, then ‘summarizing’ these depths over all $t \in [0, 1]$. One possibility is to take a simple or weighted time average as follows.

Definition 3.1.1 (Fraiman-Muniz depth). The integrated depth, or Fraiman-Muniz depth, is defined as

$$D_F(\mathbf{x}, F_{\mathbf{X}}) = \int_{[0,1]} D(\mathbf{x}(t), F_{\mathbf{X}(t)}) w(t) dt. \quad (3.1.1)$$

Here, w is a weight function.

Alternatively, we may choose the lowest or ‘worst’ depth over time. This way, low depth values over small portions of time, which indicate a deviation from centrality, are better reflected in the summary.

Definition 3.1.2 (Infimal depth). The infimal depth is defined as

$$D_{Inf}(\mathbf{x}, F_{\mathbf{X}}) = \inf_{t \in [0,1]} D(\mathbf{x}(t), F_{\mathbf{X}(t)}). \quad (3.1.2)$$

3.1.2 Band depths

López-Pintado and Romo (2009) later introduced the notion of band depth for univariate functional data.

Definition 3.1.3 (Band depth). The band depth, for some index $J \geq 2$, is defined as

$$D_B^J(\mathbf{x}, F_{\mathbf{X}}) = \sum_{j=2}^J P_{\mathbf{X}_i \sim F_{\mathbf{X}}}^{\text{iid}}(\mathbf{x} \in \text{conv}(\mathbf{X}_1, \dots, \mathbf{X}_j)). \quad (3.1.3)$$

The empirical version of band depth is defined as

$$D_B^J(\mathbf{x}, \hat{F}_n) = \sum_{j=2}^J \binom{n}{j}^{-1} \sum_{1 \leq i_1 < \dots < i_j \leq n} \mathbf{1}(\mathbf{x} \in \text{conv}(\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_j})). \quad (3.1.4)$$

This is simply the proportion of j -tuples of curves (for $2 \leq j \leq J$) which envelope \mathbf{x} . Note that if two curves intersect at a point, a third curve is enveloped by them only when it passes through the point of intersection, i.e. with probability zero. Thus, we generally use $J = 3$.

Definition 3.1.4 (Modified band depth). Define the enveloping time

$$\text{ET}(\mathbf{x}; \mathbf{x}_1, \dots, \mathbf{x}_j) = m_1(\{t \in [0, 1] : \mathbf{x} \in \text{conv}(\mathbf{x}_1, \dots, \mathbf{x}_j)\}), \quad (3.1.5)$$

where m_1 is the Lebesgue measure on \mathbb{R} . The modified band depth is defined as

$$D_{MB}^J(\mathbf{x}, F_{\mathbf{X}}) = \sum_{j=2}^J \mathbb{E}_{\mathbf{X}_i \stackrel{\text{iid}}{\sim} F_{\mathbf{X}}} [\text{ET}(\mathbf{x}; \mathbf{X}_1, \dots, \mathbf{X}_j)]. \quad (3.1.6)$$

The empirical version of modified band depth is defined as

$$D_{MB}^J(\mathbf{x}, \hat{F}_n) = \sum_{j=2}^J \binom{n}{j}^{-1} \sum_{1 \leq i_1 < \dots < i_j \leq n} \text{ET}(\mathbf{x}; \mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_j}). \quad (3.1.7)$$

We generally use $J = 2$ for ease of computation, and denote the corresponding modified band depth simply as $D_{MB}(\cdot, \cdot)$, dropping the superscript.

$$D_{MB}(\mathbf{x}, \hat{F}_n) = \binom{n}{2}^{-1} \sum_{i=1}^n \sum_{j=i+1}^n \text{ET}(\mathbf{x}; \mathbf{x}_i, \mathbf{x}_j). \quad (3.1.8)$$

3.1.3 Half-region depths

Later, López-Pintado and Romo (2011) introduced the half-region depth.

Definition 3.1.5. We say that \mathbf{y} is in the hypograph of \mathbf{x} , denoted, $\mathbf{y} \in H_{\mathbf{x}}$, if $\mathbf{y}(t) \leq \mathbf{x}(t)$ for all $t \in [0, 1]$. Similarly, we say that \mathbf{y} is in the epigraph of \mathbf{x} , denoted, $\mathbf{y} \in E_{\mathbf{x}}$, if $\mathbf{y}(t) \geq \mathbf{x}(t)$ for all $t \in [0, 1]$.

Definition 3.1.6 (Half-region depth). The half-region depth is defined as

$$D_{HR}(\mathbf{x}, F) = \min\{P_F(H_{\mathbf{x}}), P_F(E_{\mathbf{x}})\}. \quad (3.1.9)$$

The quantity $P_F(E_{\mathbf{x}})$ is called the epigraph index, which measures the proportion of curves that lie entirely above \mathbf{x} .

Definition 3.1.7 (Modified half-region depth). Denote the superior modified hypograph (MHI) and epigraph (MEI) indices

$$\text{MHI}_F(\mathbf{x}) = \mathbb{E}_{\mathbf{X} \sim F}[m_1(\{t \in [0, 1] : \mathbf{x}(t) \geq \mathbf{X}(t)\})], \quad (3.1.10)$$

$$\text{MEI}_F(\mathbf{x}) = \mathbb{E}_{\mathbf{X} \sim F}[m_1(\{t \in [0, 1] : \mathbf{x}(t) \leq \mathbf{X}(t)\})]. \quad (3.1.11)$$

The modified half-region depth is defined as

$$D_{MHR}(\mathbf{x}, F) = \min\{\text{MHI}_F(\mathbf{x}), \text{MEI}_F(\mathbf{x})\}. \quad (3.1.12)$$

3.2 Classification

Observe that the classification procedures for multivariate data described in Section 2.5 (the maximum depth classifier, the DD classifier, and the DD^G classifier) only depend on the data through the depth feature vectors

$$\mathbf{x}^D = (D(\mathbf{x}, F_1), \dots, D(\mathbf{x}, F_k)) \in \mathbb{R}^k. \quad (3.2.1)$$

By simply choosing an appropriate functional data depth D , all of these classification procedures naturally generalize to the functional setting.

3.2.1 Outlyingness matrices

Dai and Genton (2018) proposed a method which measures the outlyingness of \mathbf{x} with respect to a population via depth as follows.

Definition 3.2.1. Let \mathbf{X} be a d -variate stochastic process of continuous functions. At each time point $t \in [0, 1]$, the directional outlyingness is defined as

$$\mathbf{O}(t) = \mathbf{O}(\mathbf{X}(t), F_{\mathbf{X}(t)}) = \left(\frac{1}{D(\mathbf{X}(t), F_{\mathbf{X}(t)})} - 1 \right) \mathbf{v}(t), \quad (3.2.2)$$

where $\mathbf{v}(t)$ is the unit vector pointing from the median of $F_{\mathbf{X}(t)}$ to $\mathbf{X}(t)$.

Definition 3.2.2. The functional directional outlyingness is defined as

$$\text{FO}(\mathbf{X}, F_{\mathbf{X}}) = \int_{[0,1]} \|\mathbf{O}(t)\|^2 w(t) dt. \quad (3.2.3)$$

Definition 3.2.3. The mean directional outlyingness is defined as

$$\mathbf{MO}(\mathbf{X}, F_{\mathbf{X}}) = \int_{[0,1]} \mathbf{O}(t) w(t) dt. \quad (3.2.4)$$

Definition 3.2.4. The variation of directional outlyingness is defined as

$$\text{VO}(\mathbf{X}, F_{\mathbf{X}}) = \int_{[0,1]} \|\mathbf{O}(t) - \mathbf{MO}(t)\|^2 w(t) dt. \quad (3.2.5)$$

Here, w is a weight function on $[0, 1]$. In our discussion, we set $w = 1$.

It is easily verified that

$$\text{FO}^2 = \|\mathbf{MO}\|^2 + \text{VO}. \quad (3.2.6)$$

Dai and Genton (2018) propose using the $(d + 1)$ -variate feature vectors

$$\mathbf{Y}(\mathbf{X}, F_{\mathbf{X}}) = (\mathbf{MO}^\top, \text{VO})^\top \quad (3.2.7)$$

corresponding to the curve \mathbf{X} for the purposes of classification. For instance, one may define the classifier

$$\hat{l}(\mathbf{X}) = \arg \max_{1 \leq i \leq k} D'(\mathbf{Y}(\mathbf{X}, F_i), F_{\mathbf{Y}(\mathbf{X}, F_i)}), \quad (3.2.8)$$

where D' is a multivariate depth function. This is simply a maximum depth classifier applied on the feature vectors \mathbf{Y} . When D' is chosen to be the robust Mahalanobis depth, we have the classifier

$$\hat{l}_{RM}(\mathbf{X}) = \arg \max_{1 \leq i \leq k} D_{RM}(\mathbf{Y}(\mathbf{X}, F_i), F_{\mathbf{Y}(\mathbf{X}, F_i)}). \quad (3.2.9)$$

Definition 3.2.5. The functional directional outlyingness matrix is defined as

$$\text{FOM}(\mathbf{X}, F_{\mathbf{X}}) = \int_{[0,1]} \mathbf{O}(t) \mathbf{O}(t)^\top w(t) dt. \quad (3.2.10)$$

Definition 3.2.6. The functional directional outlyingness matrix is defined as

$$\text{VOM}(\mathbf{X}, F_{\mathbf{X}}) = \int_{[0,1]} (\mathbf{O}(t) - \mathbf{MO}(t)) (\mathbf{O}(t) - \mathbf{MO}(t))^\top w(t) dt. \quad (3.2.11)$$

Again, it is easily verified that

$$\text{FOM} = \mathbf{MO} \mathbf{MO}^\top + \text{VOM}, \quad (3.2.12)$$

and that

$$\text{FO} = \text{trace}(\text{FOM}), \quad \text{VO} = \text{trace}(\text{VOM}). \quad (3.2.13)$$

We may also use the feature matrix VOM, or its matrix norm $\|\text{VOM}\|_F$ corresponding to the curve \mathbf{X} for the purposes of classification. Here, $\|\cdot\|_F$ denotes the Frobenius norm. For instance, a VOM based classifier may be defined as

$$\hat{l}_{\text{VOM}}(\mathbf{X}) = \arg \min_{1 \leq i \leq k} \|\text{VOM}(\mathbf{X}, F_i)\|_F. \quad (3.2.14)$$

3.2.2 Random projections

Another approach is to use a feature vector consisting of multiple projections of \mathbf{X} . Given functions $\mathbf{v}_1, \dots, \mathbf{v}_d$ chosen at random, we examine the d -variate feature vectors

$$\mathbf{V}(\mathbf{X}, F_{\mathbf{X}}) = (\langle \mathbf{v}_1, \mathbf{X} \rangle, \dots, \langle \mathbf{v}_d, \mathbf{X} \rangle) \quad (3.2.15)$$

and apply a depth based multivariate classifier. For instance, given a multivariate depth function D' , we may define a classifier

$$\hat{l}_{D'}^d(\mathbf{X}) = \arg \max_{1 \leq i \leq k} D'(\mathbf{V}(\mathbf{X}, F_i), F_{\mathbf{V}(\mathbf{X}, F_i)}). \quad (3.2.16)$$

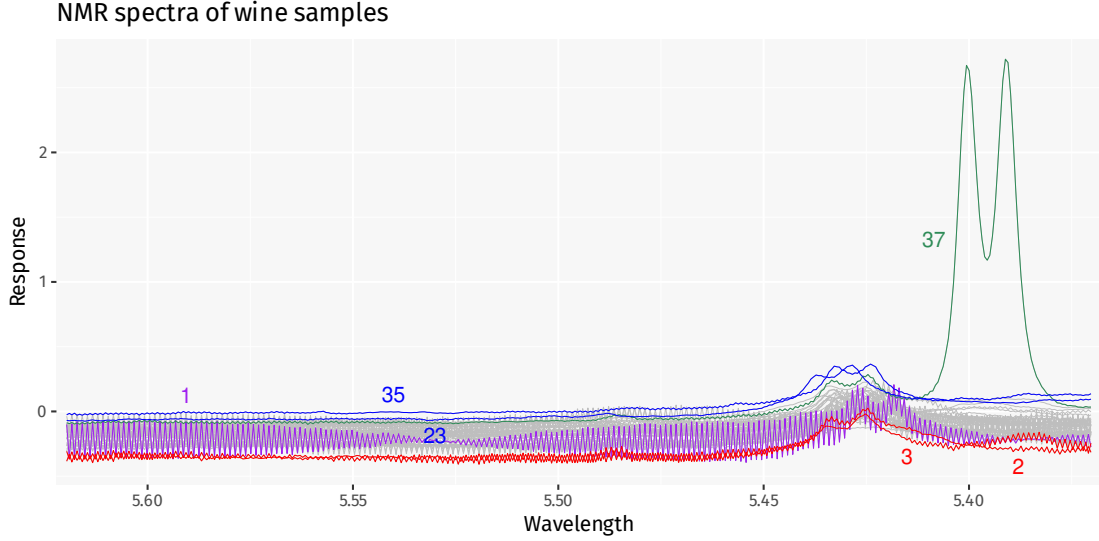


Figure 3.1: NMR spectra of 40 wine samples, from the R package `speaq`, with some curves showing outlying behaviour highlighted. The green curve #37 is an isolated outlier, the blue and red curves are shift outliers, and the purple curve is a shape outlier.

3.3 Outlier detection

A curve $\mathbf{x}: [0, 1] \rightarrow \mathbb{R}$ may exhibit outlying behaviour with respect to a body of curves in many ways; we use the useful classification as detailed in Hubert et al. (2015). It may deviate significantly over a short interval, in which case we call it an *isolated outlier*. Alternatively, it may deviate over a large, or perhaps even the whole interval, in which case we call it a *persistent outlier*. If this deviation is in terms of shape – for instance, the curve may be rougher or smoother – we call it a *shape outlier*. Otherwise, if the curve has the same shape as the rest but appears above or below them, we call it a *shift outlier*. Another possibility is that the curve differs in scale, in which case we call it an *amplitude outlier*. Some of these behaviours have been illustrated in the dataset in Figure 3.1.

An important consideration when dealing with shape outliers is that each time slice $\mathbf{x}(t)$ may be fairly inconspicuous with respect to the marginal $F_{\mathbf{X}(t)}$. Thus, it seems clear that a tool like the Fraiman-Muniz depth may succeed in identifying shift or amplitude outliers, but fall short against shape outliers. In general, the basic algorithm of iteratively selecting curves with low functional depth as outliers is often insufficient.

A common approach towards examining the shapes of curves in a dataset is to bundle them with their derivatives.

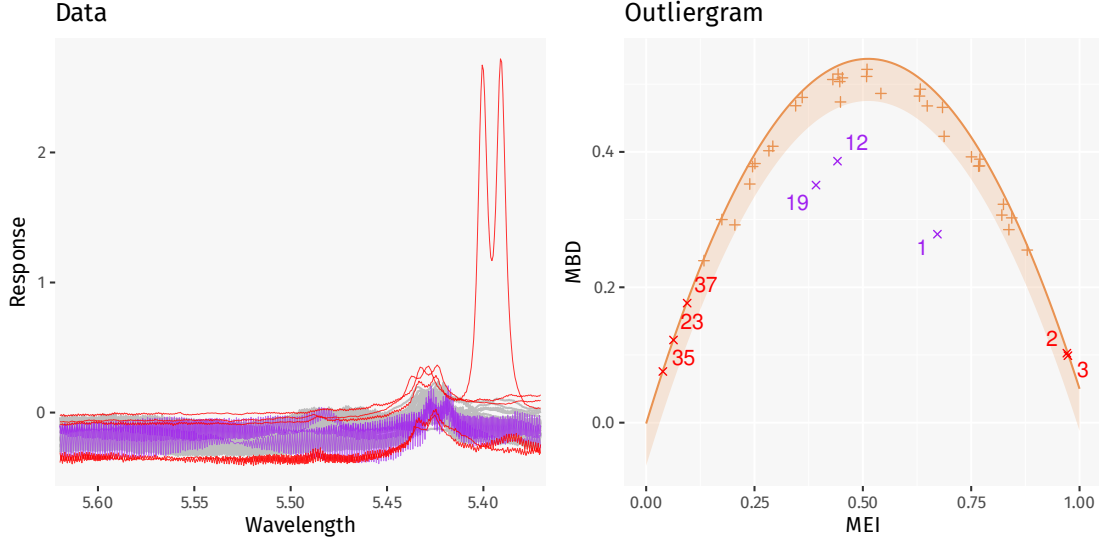


Figure 3.2: Outliergram for the NMR spectra of 40 wine samples. The three purple curves have been identified as shape outliers, as they fall outside the orange ribbon in the outliergram. Although the red curves lie on the orange parabola, they have low MBD and extreme MEI values, indicating that they lie above or below the main mass of curves.

3.3.1 Outliergrams

Arribas-Gil and Romo (2014) combined the notions of the modified epigraph index (MEI) and the modified band depth (MBD), proposing the outliergram as a tool for detecting shape outliers. They show that for a sample $\{\mathbf{x}_i\}_{i=1}^n$, each

$$\text{MBD}(\mathbf{x}_i) = D_{MB}(\mathbf{x}_i, \hat{F}_n) \leq a_0 + a_1 \text{MEI}(\mathbf{x}_i) + a_2 n^2 \text{MEI}(\mathbf{x}_i)^2 \quad (3.3.1)$$

where $a_0 = a_2 = -2/n(n-1)$ and $a_1 = 2(n+1)/(n-1)$. The distance

$$d_i = a_0 + a_1 \text{MEI}(\mathbf{x}_i) + a_2 n^2 \text{MEI}(\mathbf{x}_i) - \text{MBD}(\mathbf{x}_i) \quad (3.3.2)$$

is indicative of the outlyingness of \mathbf{x}_i . Arribas-Gil and Romo (2014) consider shape outlying curves as those for which $d_i \geq d^* = Q_3 + 1.5 \text{IQR}$, where Q_3 and IQR are the third quartile and the interquartile range of $\{d_i\}_{i=1}^n$ respectively.

Definition 3.3.1 (Outliergram). The outliergram for a dataset $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^n$ is the graph

$$\{(\text{MEI}(\mathbf{x}_i, \hat{F}_n), \text{MBD}(\mathbf{x}_i, \hat{F}_n)) : 1 \leq i \leq n\}. \quad (3.3.3)$$

Shape outliers are curves \mathbf{x}_i such that $(\text{MEI}_i, \text{MBD}_i)$ falls outside a ribbon of height d^* under the parabola $a_0 + a_1 \text{MEI} + a_2 n^2 \text{MEI}^2$.

Figure 3.2 illustrates the use of the outliergram. We have also highlighted curves with fairly low or high MEI values as shift outliers; a low MEI value indicates that the curve lies above the main mass of curves, and a high MEI indicates that it lies below.

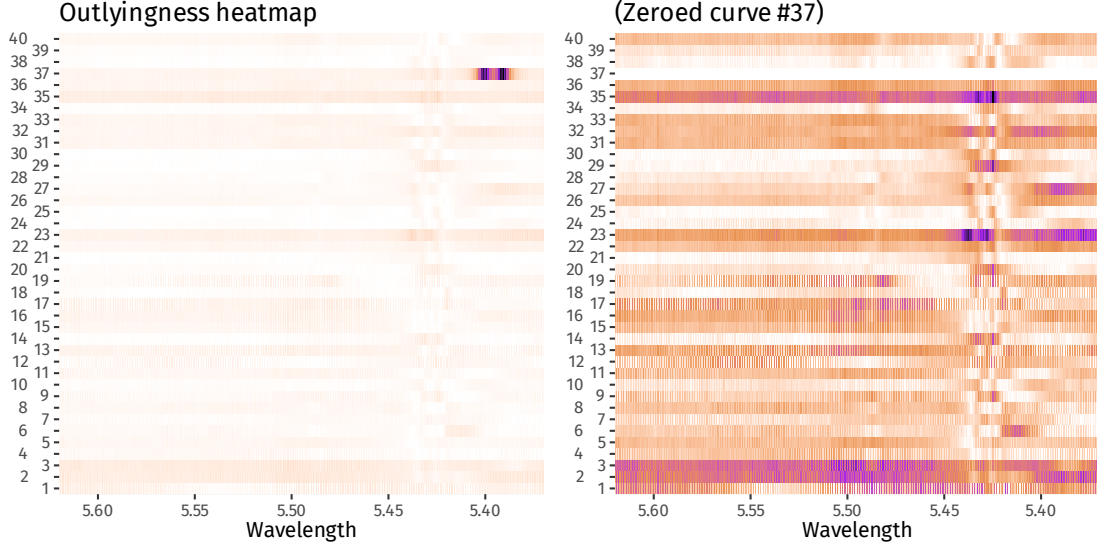


Figure 3.3: Outlyingness heatmap for the NMR spectra of 40 wine samples. The extreme curve #37 has been zeroed out in the second diagram to better illustrate the variation in outlyingness for the remaining curves.

3.3.2 Centrality-Stability diagrams

In their discussion of methods of functional outlier detection, Hubert et al. (2015) proposed the centrality-stability diagram, where both the ‘centrality’ of a curve (measured by depth) and its variability in cross-sectional outlyingness over time are accounted for. A deviation in centrality may point towards a shift outlier, while a deviation in stability may point towards an isolated or shape outlier.

Hubert et al. (2015) begin by choosing a multivariate depth function of the form $D'(\mathbf{x}(t), F_{\mathbf{X}(t)}) = (1 + O(\mathbf{x}(t)))^{-1}$, where $O(\cdot)$ is an outlyingness function. The Mahalanobis, projection, and Oja depths clearly fit this description. Here, for the purposes of computation in the univariate case, we choose

$$O(x(t)) = \frac{|x(t) - \text{med}(X(t))|}{\text{MAD}(X(t))}, \quad (3.3.4)$$

instead of using the skew-adjusted version; the differences are minor enough for us to ignore. The variation in $O(\mathbf{x}(\cdot))$ over time for different curves, in the form of an *outlyingness heatmap*, is quite revealing; Figure 3.3 shows that curves may have large outlyingness for short or long intervals.

Corresponding to D' , we have an integrated Fraiman-Muniz depth

$$D_F(\mathbf{x}, F) = \int_{[0,1]} (1 + O(\mathbf{x}(t)))^{-1} dt. \quad (3.3.5)$$

However, a spike in outlyingness over a short time interval, such as in curve #37 in Figure 3.3, may potentially be ‘washed out’ in this averaging. With this, we seek

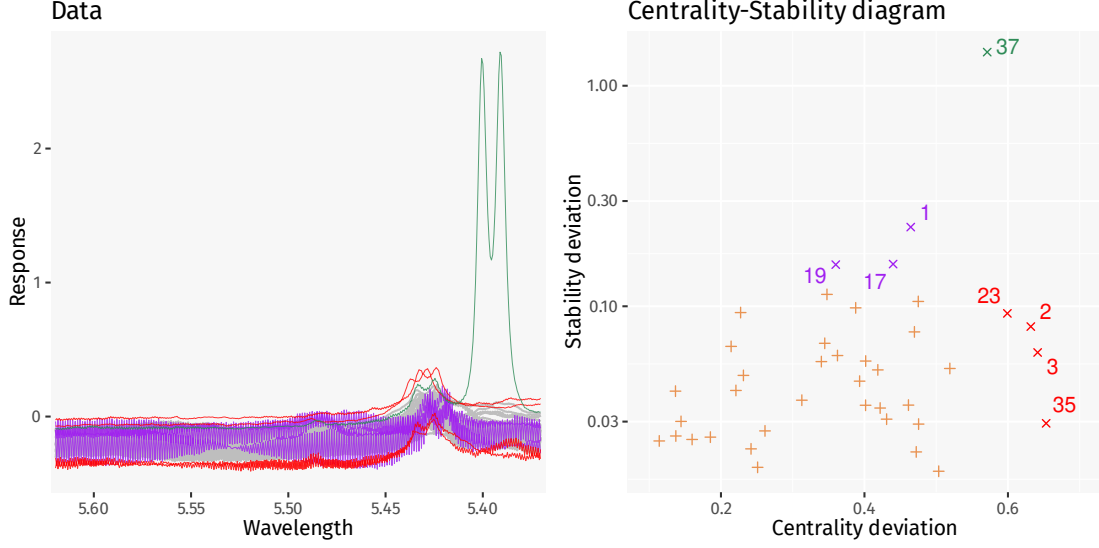


Figure 3.4: Centrality-Stability diagram for the NMR spectra of 40 wine samples. The red curves are seen to deviate in terms of centrality, indicated by the fact that the corresponding points in the centrality-stability diagram fall towards the right. The purple curves deviate in terms of stability, with the green curve showing extreme deviation.

a method of detecting sharp bursts in outlyingness. Note that by setting

$$\widetilde{\text{MO}}(\mathbf{x}, F) = \int_{[0,1]} O(\mathbf{x}(t)) dt, \quad (3.3.6)$$

Cauchy-Schwarz gives us the relation

$$D_F(\mathbf{x}, F) \cdot (1 + \widetilde{\text{MO}}(\mathbf{x}, F)) \geq 1. \quad (3.3.7)$$

Equality is achieved only when $O(\mathbf{x}(\cdot))$ remains constant over time. Any sudden variation in outlyingness over time will be detected by the *stability deviation*

$$\Delta S(\mathbf{x}, F) = (1 + \widetilde{\text{MO}}(\mathbf{x}, F)) - \frac{1}{D_F(\mathbf{x}, F)}, \quad (3.3.8)$$

the difference between the arithmetic and harmonic means of $1 + O(\mathbf{x}(\cdot))$. Defining the *centrality deviation* simply as $\Delta C(\mathbf{x}, F) = 1 - D_F(\mathbf{x}, F)$, we have our *centrality-stability diagram*.

Definition 3.3.2 (Centrality-Stability diagram). The centrality-stability diagram for a dataset $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^n$ is the graph

$$\{(\Delta C(\mathbf{x}_i, \hat{F}_n), \Delta S(\mathbf{x}_i, \hat{F}_n)) : 1 \leq i \leq n\}. \quad (3.3.9)$$

Remark. We make a distinction between **MO** from Definition 3.2.3 and $\widetilde{\text{MO}}$; the outlyingness $O(\cdot)$ used in the latter is real and positive.

Figure 3.4 illustrates the use of the centrality-stability diagram as a summary of the outlyingness heatmap from Figure 3.3. This time, the isolated outlier curve #37 is well separated from the shift and shape outliers, unlike in the outliergram in Figure 3.2.

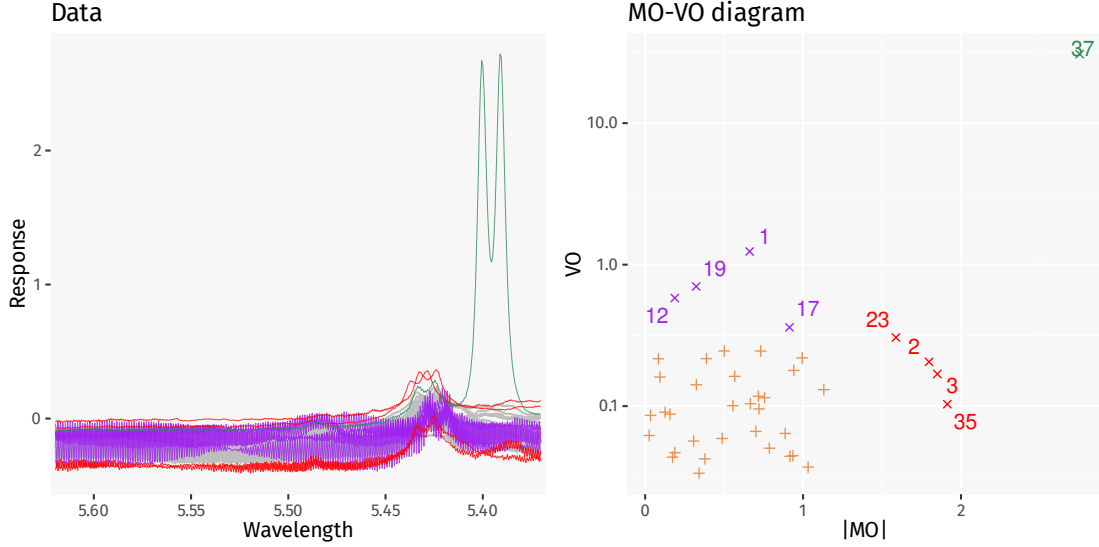


Figure 3.5: MO-VO diagram for the NMR spectra of 40 wine samples. We plot $|MO|$ rather than MO here for better comparison with the centrality stability diagram in Figure 3.4; nevertheless, the signed MO values would reveal whether the shift outliers lie above or below the main mass of curves.

3.3.3 MO-VO diagrams

We observe that the MO-VO diagram from Dai and Genton (2018) neatly falls under a general category of centrality-stability diagrams. The quantity $MO(\mathbf{x}, F)$ may indeed be treated as a measure of deviation from centrality of \mathbf{x} . Again, $VO(\mathbf{x}, F)$ being the variance of $\mathbf{O}(t)$, may be treated as a measure of deviation from stability, since it captures the variability of outlyingness over time and is sensitive to changes over short intervals.

For the purposes of computation in the univariate case, we use the directional outlyingness function

$$O(x(t)) = \frac{x(t) - \text{med}(X(t))}{\text{MAD}(X(t))}. \quad (3.3.10)$$

Figure 3.5 illustrates the use of the MO-VO diagram. Note the similarities with the centrality-stability diagram from Figure 3.4.

We use all of the above diagnostic tools on a different dataset in Figures 3.7 and 3.6.

3.4 Partially observed functional data

Consider the setting where the stochastic process \mathbf{X} of continuous functions is not observed on the entire interval $[0, 1]$, but rather on a random subset $O \subseteq [0, 1]$. Then, a dataset of partially observed curves is of the form $\mathcal{D} = \{(\mathbf{X}_i, O_i)\}_{i=1}^n$,



Figure 3.6: Outliergram, centrality-stability, and MO-VO diagrams for the NIR spectra of 39 gasoline samples, from the R package `rrcov`. The six purple curves #25, 26, 36-39 correspond to samples containing added alcohol. While the outliergram does not clearly identify these outliers, the centrality-stability and MO-VO diagrams show a marked separation from the main curves. Indeed, there is no cutoff d^* defining the lower boundary of the orange ribbon in the outliergram which properly excludes the six outliers.

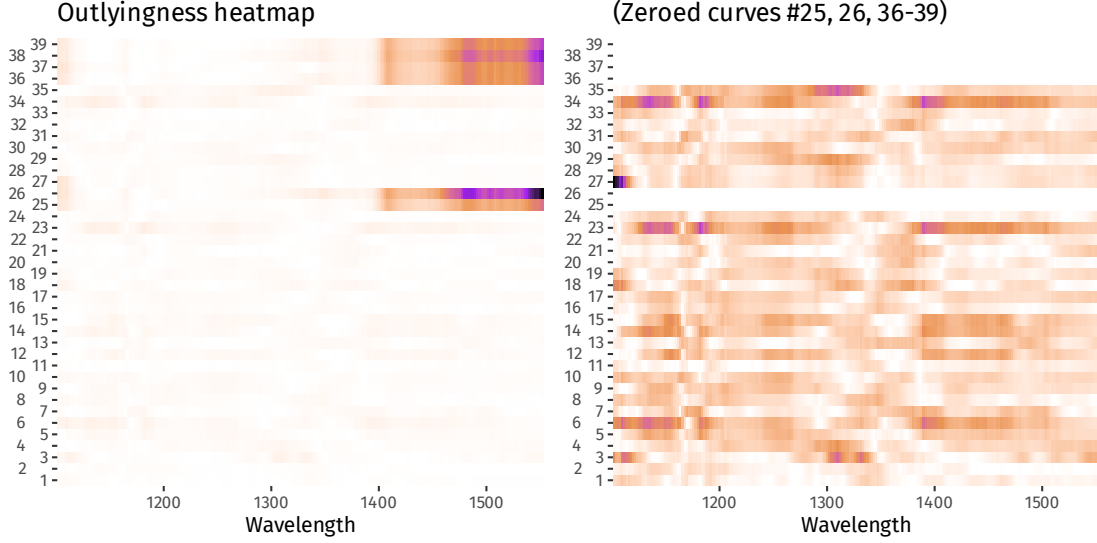


Figure 3.7: Outlyingness heatmap for the NIR spectra of 39 gasoline samples. The outlying curves have been zeroed out in the second diagram.

where $\mathbf{X}_i \stackrel{\text{iid}}{\sim} F_{\mathbf{X}}$, $O_i \stackrel{\text{iid}}{\sim} Q$ where Q generates random compact subsets of $[0, 1]$, independent of \mathbf{X}_i . In other words, $(\mathbf{X}_i, O_i) \stackrel{\text{iid}}{\sim} F_{\mathbf{X}} \times Q$. This setup is known as the ‘Missing Completely at Random’ assumption.

We set $\mathcal{J}(t) = \{j : t \in O_j\}$ to keep track of which curves \mathbf{X}_i have been observed at time t . Furthermore, denote $q(t) = |\mathcal{J}(t)|$ as the number of curves \mathbf{X}_i observed at time t .

Elías, Jiménez, Paganoni, and Sangalli (2023) propose the following.

Definition 3.4.1 (Partially observed integrated functional depth). Let D be a d -variate depth function. The Partially Observed Integrated Functional Depth (POIFD) is defined as

$$D_{POIFD}((\mathbf{x}, o), F_{\mathbf{X}} \times Q) = \int_o D(\mathbf{x}(t), F_{\mathbf{X}(t)}) w_o(t) dt, \quad (3.4.1)$$

where $w_o(t) = q(t) / \int_o q(t) dt$.

We can now proceed with tasks such as classification, outlier detection, etc. on our partially observed dataset, via depth based procedures using POIFD values. Another natural problem is one of curve reconstruction: given a partially observed curve (\mathbf{X}, O) , can we estimate \mathbf{X} on $M = [0, 1] \setminus O$? For instance, we may search for a reconstruction operator $\mathcal{R}: L_2(O) \rightarrow L_2(M)$ that minimizes the mean integrated prediction squared error $\mathbb{E}[\|\mathbf{X}_M - \mathcal{R}(\mathbf{X}_O)\|^2]$. Here, \mathbf{X}_O denotes the curve \mathbf{X} restricted to O , and similarly for \mathbf{X}_M . The best predictor in this sense is the conditional expectation $\mathbb{E}[\mathbf{X}_M | \mathbf{X}_O]$, which is in general a non-linear operator. Thus, Kraus (2015) and Kneip and Liebl (2020) search for continuous linear operators \mathcal{A} , using methods based on estimating terms of the Karhunen-Loève expansion of \mathbf{X} .

Elías, Jiménez, and Shang (2023) offer a depth based solution to the reconstruction problem, adapted from a similar algorithm for time-series forecasting (Elías et al., 2022). The main idea involves selecting a collection of curves, with indices \mathcal{J} , which best envelope (\mathbf{X}, O) , then taking a weighted linear combination. In particular, they suggest

$$\hat{\mathbf{X}}(t) = \frac{\sum_{i \in \mathcal{J}(t)} w_i \mathbf{X}_i(t)}{\sum_{i \in \mathcal{J}(t)} w_i}, \quad w_i = \exp \left(-\theta \frac{\|(\mathbf{X}, O) - (\mathbf{X}_i, O_i)\|}{\delta} \right), \quad (3.4.2)$$

where $\mathcal{J}(t) = \mathcal{J} \cap \mathcal{J}(t) = \{i \in \mathcal{J} : t \in O_i\}$ and $\delta = \min_{i \in \mathcal{J}} \|(\mathbf{X}, O) - (\mathbf{X}_i, O_i)\|$. Here, θ is a tuning parameter, perhaps chosen by minimizing the mean squared error on (\mathbf{X}, O) . Furthermore, we have denoted

$$\|(\mathbf{X}, O) - (\mathbf{X}', O')\| = \frac{1}{m(O \cap O')} \left(\int_{O \cap O'} \|\mathbf{X}(t) - \mathbf{X}'(t)\|^2 dt \right)^{1/2}. \quad (3.4.3)$$

Choosing the best envelope \mathcal{J} involves both depth and distance. Elías, Jiménez, and Shang (2023) use the following three criteria to devise an algorithm that iteratively selects \mathcal{J} .

1. (\mathbf{X}, O) should be as deep as possible in the collection of curves $\{(\mathbf{X}, O)\} \cup \{(\mathbf{X}_i, O_i)\}_{i \in \mathcal{J}}$, in the sense of POIFD.
2. (\mathbf{X}, O) should be enveloped by $\{(\mathbf{X}_i, O_i)\}_{i \in \mathcal{J}}$ as much as possible, i.e. we want to maximize the enveloping time $\text{ET}((\mathbf{X}, O); \{\mathbf{X}_i, O_i\}_{i \in \mathcal{J}})$.
3. $\{(\mathbf{X}_i, O_i)\}_{i \in \mathcal{J}}$ should contain as many near curves to (\mathbf{X}, O) as possible, in the sense of the distance 3.4.3.

Chapter 4

LOCAL DEPTH

4.1 Local depth regions

Given a distribution $F_{\mathbf{X}}$, we may define a symmetrized distribution about a point $\mathbf{x} \in \mathcal{X}$ as

$$F_{\mathbf{X}}^{\mathbf{x}} = \frac{1}{2}F_{\mathbf{X}} + \frac{1}{2}F_{2\mathbf{x}-\mathbf{X}}. \quad (4.1.1)$$

With this, \mathbf{x} becomes the point of central symmetry, hence the deepest point in $F_{\mathbf{X}}^{\mathbf{x}}$ with respect to a depth function D that obeys **P2**. Thus, the β -th central regions of $F_{\mathbf{X}}^{\mathbf{x}}$ behave like neighbourhoods of \mathbf{x} .

Definition 4.1.1 (Paindaveine and Van Bever, 2013). The probability- β depth-based neighbourhood of \mathbf{x} with respect to the distribution F is defined as

$$N_{\beta}^{\mathbf{x}}(F) = C_{F^{\mathbf{x}}}(\beta), \quad (4.1.2)$$

i.e. the β -th central region of F symmetrized about \mathbf{x} .

When working with a sample $\mathcal{D} = \{\mathbf{X}_i\}_{i=1}^n$ from F , we may obtain the β depth-based neighbourhood of \mathbf{x} by first computing the reflected sample $\mathcal{D}' = \{2\mathbf{x} - \mathbf{X}_i\}_{i=1}^n$, then arranging the elements of the symmetrized sample $\mathcal{D}^{\mathbf{x}} = \mathcal{D} \cup \mathcal{D}'$ in descending order by their empirical depth values and choosing the first β proportion of elements. The neighbourhood $N_{\beta}^{\mathbf{x}}(\hat{F}_n)$ is the convex hull of these elements.

Definition 4.1.2 (Paindaveine and Van Bever, 2013). Let D be a depth function, and let $F_{\beta}^{\mathbf{x}}$ denote the distribution F conditional on the neighbourhood $N_{\beta}^{\mathbf{x}}(F)$. The corresponding local depth function at locality level $\beta \in (0, 1]$ is defined as

$$LD_{\beta}(\mathbf{x}, F) = D(\mathbf{x}, F_{\beta}^{\mathbf{x}}) \quad (4.1.3)$$

Again, when working with a sample $\mathcal{D} = \{\mathbf{X}_i\}_{i=1}^n$, we obtain $LD_{\beta}(\mathbf{x}, \hat{F}_n)$ by arranging the elements of \mathcal{D} in descending order by their empirical depth values in the symmetrized sample $\mathcal{D}^{\mathbf{x}}$, choosing the first β proportion of elements, and computing the depth of \mathbf{x} with respect to these elements.

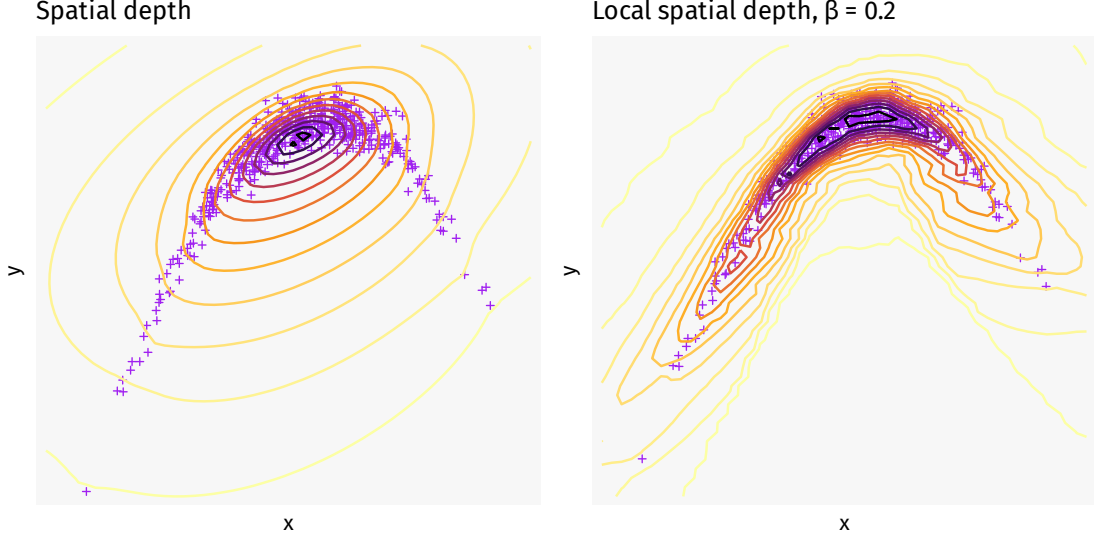


Figure 4.1: Depth contours with respect to a ‘banana-shaped’ distribution. Observe that the spatial depth contours fail to adequately capture the curved shape of the data cloud, in contrast with the local spatial depth (with $\beta = 0.2$) contours.

Remark. When $\beta = 1$, the local depth LD_1 reduces to the original global depth D .

Remark. The notions of depth based neighbourhoods and local depth make sense for any distribution F on a space \mathcal{X} as long as the process of symmetrization around $\mathbf{x} \in \mathcal{X}$ can be achieved.

4.2 Regression based on local depth

Definition 4.2.1. Let D be a depth function, and let $\tilde{F}_\beta^\mathbf{x}$ denote the symmetrized distribution $F^\mathbf{x}$ conditional on the neighbourhood $N_\beta^\mathbf{x}(F)$. Given $\mathbf{x} \in \mathcal{X}$, we may define a local depth kernel at locality level β centered at \mathbf{x} as

$$K_\beta^\mathbf{x}: N_\beta^\mathbf{x}(F) \rightarrow \mathbb{R}, \quad \mathbf{z} \mapsto D(\mathbf{z}, \tilde{F}_\beta^\mathbf{x}). \quad (4.2.1)$$

This naturally extends to a map $\mathcal{X} \rightarrow \mathbb{R}$ as $K_\beta^\mathbf{x}(\mathbf{z}) = 0$ for $\mathbf{z} \notin N_\beta^\mathbf{x}(F)$.

Note that $\tilde{F}_\beta^\mathbf{x}$ is angularly symmetric about \mathbf{x} . As a result, $K_\beta^\mathbf{x}$ is maximized at and decreases away from \mathbf{x} , for reasonably well behaved depth functions (**P2** and **P3** for multivariate depth functions).

With this, we propose the (linear) estimator

$$\hat{\mathbf{y}}_\beta(\mathbf{x}) = \sum_i w_i(\mathbf{x}) \mathbf{y}_i, \quad w_i(\mathbf{x}) = \frac{K_\beta^\mathbf{x}(\mathbf{x}_i)}{\sum_j K_\beta^\mathbf{x}(\mathbf{x}_j)}. \quad (4.2.2)$$

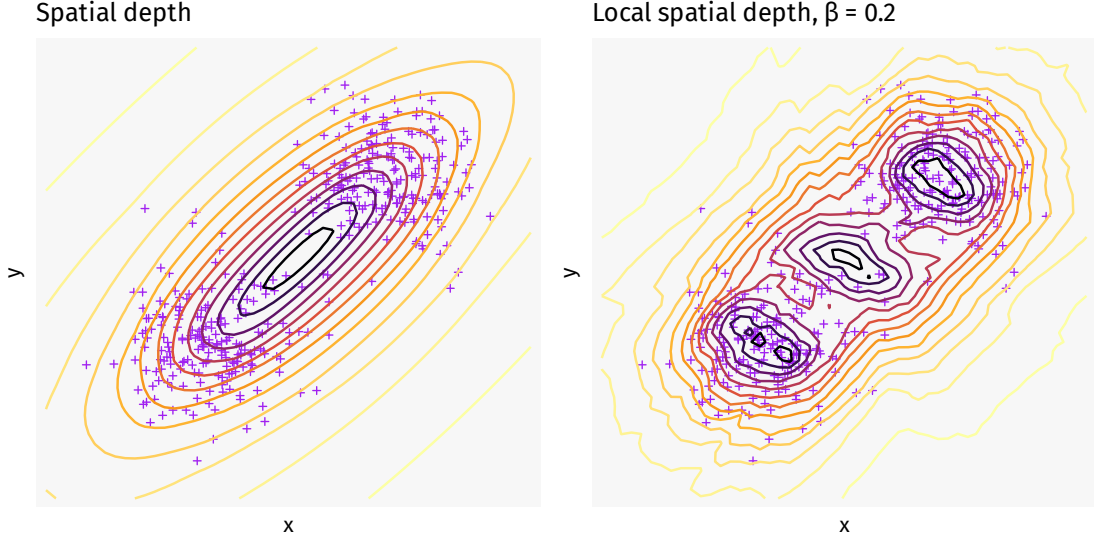


Figure 4.2: Depth contours with respect to a bimodal distribution. Although the local spatial depth contours capture the two modes correctly, it erroneously ascribes high depth values to a region in between them.

The locality level $\beta \in (0, 1]$ is a tuning parameter which may be chosen via methods such as cross-validation.

The kernel function $K_\beta^\mathbf{x}$ is supported on the neighbourhood $N_\beta^\mathbf{x}(F)$, whose shape may vary with changing $\mathbf{x} \in \mathcal{X}$. Indeed, since $N_\beta^\mathbf{x}(\hat{F}_n)$ contains the β proportion of points from $\{\mathbf{x}_i\}$ ‘closest’ to \mathbf{x} (in the sense of being more central in the symmetrized dataset $\mathcal{D}^\mathbf{x}$), this neighbourhood ought to be smaller when \mathbf{x} is more central, and larger when \mathbf{x} has fewer points nearby. Thus, $K_\beta^\mathbf{x}$ behaves somewhat like a variable bandwidth kernel, whose shape adjusts to the dataset. Furthermore, the ‘bandwidth’ of $K_\beta^\mathbf{x}$ is controlled solely by the parameter β regardless of the dimensionality or nature of \mathcal{X} . This stands in contrast with more traditional kernels which often require a selection of multiple bandwidths. For instance, a Gaussian kernel of the form

$$\mathbf{z} \mapsto \exp \left(- \sum_i \frac{(x_i - z_i)^2}{2h_i^2} \right) \quad (4.2.3)$$

needs d parameters $\{h_i\}_{i=1}^d$ to be determined.

Equation 4.2.2 may also be thought of as a weighted KNN estimator, since $N_\beta^\mathbf{x}(\hat{F}_n)$ always captures the same number of points.

When the depth function D is chosen to be affine invariant, the estimator 4.2.2 is also affine invariant, in the sense that it is unchanged by an affine transformation of \mathcal{X} . This is because $N_\beta^{A\mathbf{x}+\mathbf{b}}(F_{A\mathbf{X}+\mathbf{b}})$ will simply be the affine image of $N_\beta^\mathbf{x}(F_\mathbf{X})$.

Chapter 5

CONCLUSION

BIBLIOGRAPHY

- Agostinelli, C., & Romanazzi, M. (2011). Local depth. *Journal of Statistical Planning and Inference*, 141(2), 817–830.
- Arribas-Gil, A., & Romo, J. (2014). Shape outlier detection and visualization for functional data: the outliergram. *Biostatistics*, 15(4), 603–619.
- Chakraborty, A., & Chaudhuri, P. (2014a). On data depth in infinite dimensional spaces. *Annals of the Institute of Statistical Mathematics*, 66(2), 303–324.
- Chakraborty, A., & Chaudhuri, P. (2014b). The spatial distribution in infinite dimensional spaces and related quantiles and depths. *The Annals of Statistics*, 42(3), 1203–1231.
- Chernozhukov, V., Galichon, A., Hallin, M., & Henry, M. (2017). Monge–Kantorovich depth, quantiles, ranks and signs. *The Annals of Statistics*, 45(1), 223–256.
- Cuesta-Albertos, J. A., Febrero-Bande, M., & Oviedo de la Fuente, M. (2017). The DD^G -classifier in the functional setting. *TEST*, 26(1), 119–142.
- Cuesta-Albertos, J. A., & Nieto-Reyes, A. (2008a). The random Tukey depth. *Computational Statistics & Data Analysis*, 52(11), 4979–4988.
- Cuesta-Albertos, J. A., & Nieto-Reyes, A. (2008b). The Tukey and the random Tukey depths characterize discrete distributions. *Journal of Multivariate Analysis*, 99(10), 2304–2311.
- Dai, W., & Genton, M. G. (2018). An outlyingness matrix for multivariate functional data classification. *Statistica Sinica*, 28(4), 2435–2454.
- Donoho, D. L., & Gasko, M. (1992). Breakdown Properties of Location Estimates Based on Halfspace Depth and Projected Outlyingness. *The Annals of Statistics*, 20(4), 1803–1827.
- Elías, A., Jiménez, R., Paganoni, A. M., & Sangalli, L. M. (2023). Integrated depths for partially observed functional data. *Journal of Computational and Graphical Statistics*, 32(2), 341–352.
- Elías, A., Jiménez, R., & Shang, H. L. (2022). On projection methods for functional time series forecasting. *Journal of Multivariate Analysis*, 189, 104890.
- Elías, A., Jiménez, R., & Shang, H. L. (2023). Depth-based reconstruction method for incomplete functional data. *Computational Statistics*, 38(3), 1507–1535.
- Fraiman, R., Liu, R. Y., & Meloche, J. (1997). Multivariate Density Estimation by Probing Depth. *Lecture Notes-Monograph Series*, 31, 415–430.
- Gao, Y. (2003). Data depth based on spatial rank. *Statistics & Probability Letters*, 65(3), 217–225.

- Ghosh, A. K., & Chaudhuri, P. (2005). On maximum depth and related classifiers. *Scandinavian Journal of Statistics*, 32(2), 327–350.
- Gijbels, I., & Nagy, S. (2017). On a General Definition of Depth for Functional Data. *Statistical Science*, 32(4), 630–639.
- He, X., & Wang, G. (1997). Convergence of Depth Contours for Multivariate Datasets. *The Annals of Statistics*, 25(2), 495–504.
- Hubert, M., Rousseeuw, P. J., & Segaeert, P. (2015). Multivariate functional outlier detection. *Statistical Methods & Applications*, 24(2), 177–202.
- Jörnsten, R. (2004). Clustering and classification based on the L_1 data depth [Special Issue on Multivariate Methods in Genomic Data Analysis]. *Journal of Multivariate Analysis*, 90(1), 67–89.
- Kneip, A., & Liebl, D. (2020). On the optimal reconstruction of partially observed functional data. *The Annals of Statistics*, 28(3), 1692–1717.
- Koltchinskii, V. I. (1997). M-Estimation, Convexity and Quantiles. *The Annals of Statistics*, 25(2), 435–477.
- Kong, L., & Zuo, Y. (2010). Smooth depth contours characterize the underlying distribution. *Journal of Multivariate Analysis*, 101(9), 2222–2226.
- Koshevoy, G. A. (2002). The Tukey Depth Characterizes the Atomic Measure. *Journal of Multivariate Analysis*, 83(2), 360–364.
- Kraus, D. (2015). Components and completion of partially observed functional data. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 77(4), 777–801.
- Li, J., Cuesta-Albertos, J. A., & Liu, R. Y. (2012). DD-Classifier: Nonparametric Classification Procedure Based on DD-Plot. *Journal of the American Statistical Association*, 107(498), 737–753.
- Liu, R. Y. (1990). On a Notion of Data Depth Based on Random Simplices. *The Annals of Statistics*, 18(1), 405–414.
- Liu, R. Y., Parelius, J. M., & Singh, K. (1999). Multivariate analysis by data depth: descriptive statistics, graphics and inference. *The Annals of Statistics*, 27(3), 783–858.
- Liu, R. Y., & Singh, K. (1993). A Quality Index Based on Data Depth and Multivariate Rank Tests. *Journal of the American Statistical Association*, 88, 252–260.
- López-Pintado, S., & Romo, J. (2009). On the concept of depth for functional data. *Journal of the American Statistical Association*, 104(486), 718–734.
- López-Pintado, S., & Romo, J. (2011). A half-region depth for functional data. *Computational Statistics & Data Analysis*, 55(4), 1679–1695.
- Mizera, I., & Volauf, M. (2002). Continuity of Halfspace Depth Contours and Maximum Depth Estimators: Diagnostics of Depth-Related Methods. *Journal of Multivariate Analysis*, 83(2), 365–388.
- Mosler, K., & Mozharovskiy, P. (2022). Choosing Among Notions of Multivariate Depth Statistics. *Statistical Science*, 37(3), 348–368.
- Nagy, S. (2017). Monotonicity properties of spatial depth. *Statistics & Probability Letters*, 129, 373–378.

- Nagy, S. (2020). The halfspace depth characterization problem. In M. La Rocca, B. Liseo, & L. Salmaso (Eds.), *Nonparametric statistics* (pp. 379–389). Springer International Publishing.
- Nagy, S. (2021). Halfspace depth does not characterize probability distributions. *Statistical Papers*, 62(3), 1135–1139.
- Nagy, S., Gijbels, I., & Hlubinka, D. (2017). Depth-Based Recognition of Shape Outlying Functions. *Journal of Computational and Graphical Statistics*, 26(4), 883–893.
- Nieto-Reyes, A., & Battey, H. (2016). A Topologically Valid Definition of Depth for Functional Data. *Statistical Science*, 31(1), 61–79.
- Paindaveine, D., & Van Bever, G. (2013). From depth to local depth: A focus on centrality. *Journal of the American Statistical Association*, 108(503), 1105–1119.
- Serfling, R. (2002). A Depth Function and a Scale Curve Based on Spatial Quantiles. In Y. Dodge (Ed.), *Statistical Data Analysis Based on the L1-Norm and Related Methods* (pp. 25–38). Birkhäuser Basel.
- Shi, X., Zhang, Y., & Fu, Y. (2023). Two-sample tests based on data depth. *Entropy*, 25(2).
- Struyf, A., & Rousseeuw, P. J. (1999). Halfspace Depth and Regression Depth Characterize the Empirical Distribution. *Journal of Multivariate Analysis*, 69(1), 135–153.
- Tukey, J. W. (1975). Mathematics and the Picturing of Data. *Proceedings of the International Congress of Mathematicians, Vancouver, 1975*, 2, 523–531.
- Villani, C. (2003). *Topics in optimal transportation*. American Mathematical Society.
- Zuo, Y., & He, X. (2006). On the limiting distributions of multivariate depth-based rank sum statistics and related tests. *The Annals of Statistics*, 34(6), 2879–2896.
- Zuo, Y., & Serfling, R. (2000). General notions of statistical depth function. *The Annals of Statistics*, 28(2), 461–482.