

# On depth and deep points: a calculus

Ivan Mizera

**ABSTRACT.** For a general definition of depth in data analysis a differential-like calculus is constructed in which the location case (the framework of Tukey’s median) plays a fundamental role similar to that of linear functions in the mathematical analysis. As an application, a lower bound for maximal regression depth is proved in the general multidimensional case—as conjectured by Rousseeuw and Hubert, and others. This lower bound is demonstrated to have an impact on the breakdown point of the maximum depth estimator.

## 1. INTRODUCTION AND OUTLINE

**1.1. Introduction.** The notion of depth provides a route to possible analogs of the sample median and quantiles—beyond the univariate location model. Unlike other approaches, depth substantially elaborates on the order structure of the sample space.

In the univariate location case, the first relevant observations can be traced in Hotelling (1929) and Chamberlin (1933). For multivariate location, the proposal of Tukey (1975) was developed by Donoho and Gasko (1992); the germ of the idea appeared already in Hodges (1955). A breakthrough broadening the understanding of depth was the invention of regression depth by Rousseeuw and Hubert (1999a), see also Rousseeuw and Hubert (1999b), Hubert, Rousseeuw and Van Aelst (1999); the precursors here were Edgeworth (1888) and Daniels (1954). The interrelations between location and regression were indicated by Hill (1960) and Carrizosa (1996).

In pursuit of the regression version, Rousseeuw and Hubert (1999a) isolated the crucial general essence of halfspace depth: its connection to admissibility in a certain data-analytic sense. In the present paper, we further develop their idea: a general definition of depth is formalized in the framework of vector optimization. Several examples aim at convincing the reader that this way of thinking about depth opens a transparent route to depth-based analysis in various statistical models. Furthermore, vector-optimization approach leads to a sort of differential calculus, where the simplest, multivariate location depth plays the role of the prototype—similarly as linear

---

1991 *Mathematics Subject Classification.* Primary 62H05, Secondary 52A40, 54C60, 55M25, 60B10, 90C29.

*Key words and phrases.* Centerpoint, Compactification, Degree of mapping, Depth, Halfspace, Kronecker index, Median, Multivariate location, Regression, Set-valued analysis, Vector optimization, Weak convergence.

Research supported by Slovak VEGA grants 1/4196/97, 1/7295/20, and by National Scientific and Engineering Research Council of Canada.

structures do in the analysis of nonlinear ones. To illustrate the power of these techniques, we show how they can be successfully used for dealing with the so-called centerpoint problem; in particular, we settle certain standing conjectures from the literature. Finally, we give some statistical implications of centerpoint considerations: results about the bias and breakdown properties of maximum depth estimators in a general setting.

The range of questions opened is broad and difficult to cover in one paper. At the moment, our objective was to lay necessary theoretical and mathematical foundations; the detailed study of depth and depth-based procedures in concrete statistical models is left for the forthcoming work. For applications and aspects of depth, the reader may also consult Liu, Parelius and Singh (1999), Caplin and Nalebuff (1988, 1991b,a), and other references given or not given in this paper. The relevant (largely asymptotic) theory can be found in He and Wang (1997), Nolan (1992, 1998), He and Portnoy (1998), and Bai and He (1999). Our approach to depth is limited to what is called halfspace or Tukey's depth; for other brands of depth in multivariate location see Liu, Parelius and Singh (1999), Liu and Singh (1993), in linear regression Rousseeuw and Hubert (1999a).

**1.2. Outline of the paper.** In Section 2, the definitions of global, local and tangent depth are introduced, based on vector optimization formalism—with weak optimality as the basic notion. General inequalities between, as well as convexity-based criteria for the equality of sorts of depth are given. A few statistical models (covering nonetheless all instances studied in the literature so far) are introduced as examples.

Section 3 extends tangent depth to the general measure-theoretic setting—covering both finite-sample and population distributions. The equality to intuitive definitions, given in literature for several models, is demonstrated on examples. The centerpoint questions are expounded in Subsection 3.3; the most important result is Theorem 3.3. Its less general version Theorem 3.4 admits a shorter and more intuitive proof and extends also to multivariate regression.

Statistical implications of centerpoint considerations are treated in Section 4. Particularly, in 4.1 we show that the bias sets of maximum depth estimators are contained in upper level sets of depth—the fact giving another statistical interpretation for depth contours. In 4.2, we give the lower bounds for the total variation and contamination breakdown points and apply them to the regression depth, with the help of Theorem 5.25—the latter also implies the existence of a fit with maximal regression depth in the general situation.

Mathematical details are collected in Section 5, which, besides the proofs of the results contains also a wealth of the supporting material and adds a lot of illuminating details to the—albeit in principle self-contained—presentation in the Sections 2, 3 and 4. The proofs of all theorems and propositions from Section 2 can be found in Subsection 5.1, in their order of appearance. Theorem 3.2 is proved at the beginning of 5.4, followed later by Proposition 3.1; Theorem 3.4 in 5.7; Subsection 5.8 contains proofs of all theorems from Section 4; and finally, Theorem 5.25, whose direct consequence is Theorem 3.5, is formulated and proved in Subsection 5.9. The notation and techniques developed in Section 5 are also used in the Appendix, devoted to the full proof of Theorem 3.3.

Some of the results have been announced in Mizera (1998) and Portnoy and Mizera (1999).

## 2. DEPTH IN DATA ANALYSIS

**2.1. Preliminaries.** We denote by  $\mathbb{R}^p$  the  $p$ -dimensional Euclidean space, by  $\|\cdot\|$  its Euclidean norm, and by  $0$  the point with zero norm. If  $\vartheta, u \in \mathbb{R}^p$ ,  $u \neq 0$ , we write  $H_{\vartheta, u}$  for the set  $\{x \in \mathbb{R}^p: u^T(x - \vartheta) \geq 0\}$ , the closed halfspace whose boundary contains  $\vartheta$  and is orthogonal to the vector  $u$  pointing into the halfspace. Obviously,  $H_{\vartheta, u} = H_{\vartheta, cu}$  for any  $c > 0$ ; hence we often pick  $u$  from the sphere  $\mathbb{S}^{p-1} = \{x \in \mathbb{R}^p: \|x\| = 1\}$ . If  $\vartheta = 0$ , we abbreviate  $H_{0, u}$  to  $H_u$ .

We denote by  $A^c$  the complement of  $A$ —with respect to a basic set which should be clear from the context. For instance, if  $A$  is a subset of the indexing set  $\{1, 2, \dots, n\}$ ,  $A^c$  denotes the set  $\{1, 2, \dots, n\} \setminus A$ .

**2.2. Generalities.** A typical data-analytic model consists of a collection  $\mathcal{Z}$  of observations  $z_1, z_2, \dots, z_n$ , whose values lie in a *sample space*  $Z$ . For these data, a fit is sought: an element  $\vartheta$  of a *parameter space*  $\Theta$ .

To give each observation its impact on the result, a criterial function  $F_i$  is attached to every observation  $z_i$ ; the lower the value of  $F_i$  at  $\vartheta$ , the better  $\vartheta$  fits  $z_i$ . Fits that yield optimal values for all  $F_i$ —that is, uniformly best solutions—occur only for rare data configurations. Classical strategies thus consider trade-offs, giving each point its share via a compound criterial function—the sum of all  $F_i$ , say.

The approach involving depth could be characterized as elaborating on a “degree of data-analytic admissibility”. The general definition of depth, given by Rousseeuw and Hubert (1999a), says that “the depth of  $\vartheta$  is the smallest number of observations that would need to be removed to make  $\vartheta$  a nonfit”. The word “nonfit” means a parameter value inadmissible from a data-analytic view, a parameter value with zero depth. In what follows, we develop this idea in the framework of vector, multi-objective optimization, in the setting which employs criterial functions. Such an approach provides a guide how to define depth in various statistical models—via a natural transition from classical techniques. It also lays a firm technical foundation for the study of depth and related notions by methods of differential calculus.

**2.3. Global depth.** Let  $n$  stand for  $\text{card } \mathcal{Z}$  and  $N$  for the index set  $\{1, 2, \dots, n\}$ ; if  $A \subseteq N$ , we write  $\#A$  for  $\text{card } A$  divided by  $n$ .

Suppose that a criterial function  $F_i = F_{z_i}$  acting from  $\Theta$  to  $[0, \infty)$  is attached to every observation  $z_i$  from  $\mathcal{Z}$ ; note that the notation implies that criterial functions for two equal observations coincide. A parameter value  $\vartheta \in \tilde{\Theta} \subseteq \Theta$  will be called **weakly optimal** in  $\tilde{\Theta}$  with respect to  $A \subseteq N$ , if  $A \neq \emptyset$  and there is no  $\tilde{\vartheta} \in \tilde{\Theta}$  such that  $F_{z_i}(\tilde{\vartheta}) < F_{z_i}(\vartheta)$  for all  $i \in A$ . We define the **global depth** of  $\vartheta \in \Theta$  to be

$$(1) \quad d_G(\vartheta, \mathcal{Z}) = \min \# \{A \subseteq N: \vartheta \text{ is not weakly optimal in } \Theta \text{ with respect to } A^c\}.$$

If  $\vartheta$  is not weakly optimal with respect to the full collection of  $F_i$ ’s, then the minimal set  $A$  in (1) is empty and  $d_G(\vartheta, \mathcal{Z}) = 0$ . On the other hand, if  $\vartheta$  is weakly optimal with respect to any

subset of  $N$ , then  $d_G(\vartheta, \mathcal{Z}) = 1$ ; by definition,  $\vartheta$  is not weakly optimal with respect to the empty set of criterial functions.

If we choose certain other criterial functions  $\tilde{F}_{z_i}$  to be attached to any data point  $z_i$ , then such a choice may lead to the same depth—as long as  $F_{z_i}(\vartheta) < F_{z_i}(\tilde{\vartheta})$  if and only if  $\tilde{F}_{z_i}(\vartheta) < \tilde{F}_{z_i}(\tilde{\vartheta})$ . Obviously, depth depends essentially only on the order induced by the criterial functions; nevertheless, specific form of criterial functions often allows for better technical handling than abstract order notions.

To achieve immediate compatibility with the population case, we define depth as minimal *proportion*—and not *number* of observations (as common in the literature). Our depth thus has values  $0, 1/n, 2/n, \dots, 1$  instead of  $0, 1, 2, \dots, n$ —a minor detail ignored in the sequel, where we may speak about the equality of depths even if it may actually hold up to a multiplication by  $n$ .

The word “admissibility” we used in the introduction suggests the relationship to those notions from statistical theory related to loss and risk functions under a probabilistic model for the data. While from the decision-theoretic aspect the relationship is really close, we have to raise several cautions. First, “admissibility” considered here is a different, data-analytic one: “residual admissibility”. An example: any multivariate location estimator contained in the convex hull of the observations—the sample mean, for instance—is weakly optimal (and thus “residual admissible”). Compare this with the well-known fact from the statistical theory: the sample mean in dimensions beyond three is inadmissible with respect to quadratic loss function under normal sampling distribution.

Second, “admissibility” in decision theory means what vector optimization literature calls *Pareto optimality*: there is no solution which performs *strictly better in one* criterion and *better or equally well in others*. But, for the definition of depth, we do not use Pareto, but *weak optimality* (also known as “weak Pareto optimality”, “Slater optimality”, “weak efficiency”): there is no solution performing *strictly better in all* criteria. The reason is that in order to have a consistent definition of depth, an omission of a criterial function should not create optimality, only possibly destroy it. This is a property of weak, but not Pareto optimality.

To prevent misunderstandings in the sequel, we therefore avoid the word “admissibility”; we also speak about “criterial” rather than “loss” functions. On the other hand, we do not adopt the word “nonfit” either; as given in Rousseeuw and Hubert (1999a), we consider it semantically on a more general and, to an extent, intuitive level than our rigorous “not weakly optimal fit”.

## 2.4. Examples. We introduce now several statistical models to be analyzed in the sequel.

**EXAMPLE 1** (Multivariate location). In the location model,  $\mathcal{Z} = \Theta = \mathbb{R}^p$ . Natural criterial functions are those based on the distance  $\|\vartheta - z\|$  of  $\vartheta$  from  $z$ ; in fact, any increasing function of this distance results in the same depth. For technical convenience, we choose  $F_z(\vartheta) = \frac{1}{2}\|\vartheta - z\|^2$ .

It is not hard to see that  $\vartheta$  is weakly optimal if and only if it lies in the convex hull of the data points. Applying the definition of the global depth to this model, we obtain the definition of *location depth* used by Rousseeuw and Hubert (1999a):  $d_G(\vartheta, \mathcal{Z})$  is the minimal proportion of

points whose removal makes  $\vartheta$  lying outside the convex hull of the remaining ones. This definition is equivalent (see 2.6) to the halfspace definition of Tukey (1975) and Donoho and Gasko (1992).

**EXAMPLE 2** (Linear regression). In the regression model, an observation  $z_i = (x_i^T, y_i)^T$  consists of a response  $y_i$  and a vector of covariates  $x_i$ . The sample space is  $Z = X \times \mathbb{R}$ , where  $X$  is a subset of  $\mathbb{R}^p$ ; the parameter space is  $\Theta = \mathbb{R}^p$ . Most models have  $x_i = (1, w_i)^T$  and  $X = \{1\} \times \mathbb{R}^{p-1}$  accordingly—regression *with intercept*. When  $w_i$  runs over  $\mathbb{R}$  we speak about the *simple linear regression*; when it runs over  $\mathbb{R}^{p-1}$ , about *multiple linear regression*. The general model covers also cases when  $x_i = g(w_i)$  and  $w_i$  is a lesser-dimensional covariate; for instance, when  $x_i = (1, w_i, w_i^2)^T$  (*quadratic regression*).

Natural criterial functions in regression depend on the residuals. Any increasing function of the absolute residuals leads to the same result (see 2.6), *regression depth* of Rousseeuw and Hubert (1999a). We choose  $F_z(\vartheta) = \frac{1}{2}(y - \vartheta^T x)^2$ ; another possibility could be  $F_z(\vartheta) = |y - \vartheta^T x|$ .

**EXAMPLE 3** (General, nonlinear regression). Nonlinear regression is a generalization of the linear one; however, the usual notation is slightly different. The observations are  $z_i = (w_i^T, y_i)^T$ , drawn from  $Z = W \times \mathbb{R}$ ; the functional dependence is given by a *regression function*  $f(\vartheta, w)$ , the regression being linear if  $f(\vartheta, w) = \vartheta^T g(w)$ .

Similarly to linear regression, we choose  $F_z(\vartheta) = \frac{1}{2}(y - f(\vartheta, w))^2$ . Simple nonlinear (linearizable) models were considered in Rousseeuw and Hubert (1999a) and Van Aelst, Rousseeuw, Hubert and Struyf (2000).

**EXAMPLE 4** (Multivariate linear regression). Multivariate regression is a generalization of Example 2 in another direction. The functional dependence remains linear, but the response  $y_i$  is allowed to be multi-dimensional. An observation  $z_i = (x_i^T, y_i^T)^T$  belongs to  $Z = X \times \mathbb{R}^m$ , where  $X$  is a subset of  $\mathbb{R}^k$ . For notational convenience, we consider the parameter  $\Theta$  to be a  $k \times m$  matrix lying in  $\Theta = \mathbb{R}^p = \mathbb{R}^{km}$  (interpreting it also as a vector, if necessary). The criterial functions are  $F_z(\Theta) = \frac{1}{2}\|y - \Theta^T x\|^2$ .

Our attention to this model was turned by the work of Bern and Eppstein (2000), who gave a geometric definition of multivariate regression depth.

**EXAMPLE 5** (Orthogonal regression). In orthogonal regression, the observations  $z_i$  are points from  $\mathbb{R}^p$  and fits are  $k$ -dimensional affine subspaces in  $\mathbb{R}^p$ . We consider only the simple (traditional) case when  $k = p - 1$  and fits are *hyperplanes*, affine subspaces of codimension 1; this case parallels the classical regression, where one of the variables is interpreted as the response and other as covariates. A hyperplane is parametrized by  $\vartheta = (s, \beta^T)^T$ , where  $\beta \in \mathbb{S}^{p-1}$  is a unit vector orthogonal to the hyperplane and  $s\beta$  is the intersection of the hyperplane with the linear space generated by  $\beta$ . The resulting parameter space is  $\mathbb{R} \times \mathbb{S}^{p-1}$ . We found this parametrization convenient, despite the lack of identification:  $(s, \beta^T)$  and  $(-s, -\beta)$  represent the same hyperplane.

The fact that of the regression is orthogonal is expressed by the choice of the criterial functions. They are based on the orthogonal distances of observations to the fitted hyperplane:  $F_z(\vartheta) = \frac{1}{2}(\beta^T z - s)^2$ .

**2.5. Local depth.** Weak optimality can be effectively studied in a way akin to the classical approach of the differential calculus to extrema. The first step is a transition to local notions.

We define the **local depth** of  $\vartheta \in \Theta$  to be

$$d_{\text{loc}}(\vartheta, \mathcal{Z}) = \min \# \{A \subseteq N : \vartheta \text{ is weakly optimal in no neighborhood of } \vartheta \text{ w.r.t. } A^c\},$$

where “in no neighborhood of  $\vartheta$ ” means “in no open  $\tilde{\Theta} \subseteq \Theta$  containing  $\vartheta$ ”. Since  $\Theta$  is itself a neighborhood of  $\vartheta$ , the global depth never exceeds the local one; the following theorem gives a sufficient condition for their equality. For the definition of quasi-convexity, see 5.1.

**THEOREM 2.1.** *For any  $\vartheta \in \Theta$ ,  $d_G(\vartheta, \mathcal{Z}) \leq d_{\text{loc}}(\vartheta, \mathcal{Z})$ . If  $\Theta$  is an open convex subset of  $\mathbb{R}^p$  and all  $F_i = F_{z_i}$ , attached to the data points  $z_i \in \mathcal{Z}$ , are strictly quasi-convex (in particular, convex), then  $d_G(\vartheta, \mathcal{Z}) = d_{\text{loc}}(\vartheta, \mathcal{Z})$ .*

**2.6. Tangent depth.** In the calculus methodology of handling extremes, the crucial second step is the use of derivatives. Suppose that  $\Theta$  is a  $p$ -dimensional manifold. Given a function  $F$  from  $\Theta$  to  $\mathbb{R}$ , we denote by  $\nabla F(\vartheta)$  the derivative (gradient) of  $F$  at  $\vartheta$ : a linear functional from the tangent space of  $\Theta$  at  $\vartheta$  to  $\mathbb{R}$ , representing the local linear approximation of  $F$  at  $\vartheta$ . Since all tangent spaces are isomorphic to  $\mathbb{R}^p$ , we identify in the usual fashion  $\nabla F(\vartheta)$  with a vector in  $\mathbb{R}^p$ , the vector of partial derivatives of  $F$  at  $\vartheta$ .

(The reader not comfortable with this language may think about derivatives in a less sophisticated way: in most models, the parameter space is actually  $\mathbb{R}^p$  and gradients are the vectors of the partial derivatives taken in the elementary way. Example 5, however, shows the need for the advanced formalism.)

Differential approach to vector optimization dates back to Frisch (1966). As in the classical calculus, we deal here with first-order necessary and second-order sufficient conditions; for Pareto optima those were developed by Smale (1973, 1975b,a) and Wan (1975, 1978). Their first-order, necessary condition for a Pareto optimum turns out to be actually the same as the one required for a weak optimum; see 5.1 for more details. It provides also a sufficient condition when the criterial functions possess a certain degree of convexity; and since this is all we need for the development of depth theory, we do not introduce any second-order sufficient conditions for weak optima in this paper.

If  $S$  is a subset  $S$  of  $\mathbb{R}^p$ , we say that  $S$  **surrounds**  $\vartheta \in \mathbb{R}^p$  ( $\vartheta$  **is surrounded** by  $S$ ) whenever  $\vartheta$  lies in the convex hull of  $S$ .

**PROPOSITION 2.2.** *If  $S \subset \mathbb{R}^p$  surrounds  $\vartheta \in \mathbb{R}^p$ , then it contains a finite subset with at most  $p + 1$  elements that also surrounds  $\vartheta$ . The following are equivalent:*

- (i)  $\vartheta \in \mathbb{R}^p$  is surrounded by  $S$ ;
- (ii) there is no open halfspace in  $\mathbb{R}^p$  which contains  $S$  and has  $\vartheta$  on its boundary;
- (iii) if  $S$  is finite and  $S = \{\eta_1, \eta_2, \dots, \eta_k\}$ , there are nonnegative  $\lambda_1, \lambda_2, \dots, \lambda_k$ , not all equal to zero and such that  $\sum \lambda_i \eta_i = \vartheta$ .

It is not hard to see that if  $\vartheta$  is locally weakly optimal, then the origin  $0$  must be surrounded by  $\nabla F_i(\vartheta)$ ; in other words, once all  $\nabla F_i(\vartheta)$  are contained in an open halfspace with  $0$  on its boundary, then  $\vartheta$  is not locally weakly optimal. Hence, we have a necessary condition for weak optimality which leads to the following definition. Suppose that  $\Phi$  is a function from  $\Theta \times \mathcal{Z} \rightarrow \mathbb{R}^p$ .

Writing its values for  $\vartheta \in \Theta$  and  $z \in \mathcal{Z}$  as  $\Phi_\vartheta(z)$ , we define the **tangent depth** of  $\vartheta \in \Theta$  to be

$$(2) \quad d_T^\Phi(\vartheta, \mathcal{Z}) = \min_{u \neq 0} \# \{i: u^\top \Phi_\vartheta(z_i) \geq 0\} = \min_{\|u\|=1} \# \{i: \Phi_\vartheta(z_i) \in H_u\}.$$

The minimum is taken over all closed halfspaces with 0 on the boundary—we may therefore take it over  $u \neq 0$  or over  $\|u\| = 1$ , whichever is more convenient. The connection between local and tangent depth is established by setting  $\Phi_\vartheta(z) = \nabla F_z(\vartheta)$ ; nevertheless, the general definition opens a room for general  $\Phi$ , not necessarily coming from the vector-optimization problem (an analogy could be general estimating equations not necessarily arising from the maximization of a likelihood). In our examples, we omit the superscript  $\Phi$  whenever its form is clear from the context.

A parameter  $\vartheta$  surrounded by gradients may not yet be a weak local optimum. Nevertheless, it often is—for instance, when all  $F_i$  are convex. For the definition of pseudo-convexity, see 5.1.

**THEOREM 2.3.** *Suppose that all  $F_i = F_{z_i}$ , attached to the data points  $z_i \in \mathcal{Z}$ , are differentiable and let  $\Phi_\vartheta(z_i) = \nabla F_{z_i}(\vartheta)$ . For any  $\vartheta \in \Theta$ ,  $d_{\text{loc}}(\vartheta, \mathcal{Z}) \leq d_T^\Phi(\vartheta, \mathcal{Z})$ . If  $\Theta$  is an open convex subset of  $\mathbb{R}^p$  and all  $F_i$  are pseudo-convex (in particular, convex), then  $d_T^\Phi(\vartheta, \mathcal{Z}) = d_{\text{loc}}(\vartheta, \mathcal{Z}) = d_G(\vartheta, \mathcal{Z})$ .*

The proofs of Theorem 2.1 and 2.3 given in 5.1 reveal that neither strict quasi-convexity nor pseudo-convexity are minimal sufficient conditions required for the equality of depths.

**2.7. Examples of tangent depth.** Now we are ready to reconsider our examples and show that they cover all instances of halfspace depth studied in the literature so far.

**EXAMPLE 1.** The criterial functions  $F_z$  are convex and

$$(3) \quad \Phi_\vartheta(z) = \nabla F_z(\vartheta) = \vartheta - z.$$

Therefore,

$$\begin{aligned} d_G(\vartheta, \mathcal{Z}) &= d_T(\vartheta, \mathcal{Z}) = \min_{u \neq 0} \# \{i: u^\top (\vartheta - z_i) \geq 0\} \\ &= \min_{u \neq 0} \# \{i: u^\top (z_i - \vartheta) \geq 0\} = \min_{\|u\|=1} \# \{i: z_i \in H_{\vartheta, u}\}. \end{aligned}$$

This shows the equality to the original halfspace definition of Tukey (1975) and Donoho and Gasko (1992).

**EXAMPLE 4.** Preserving the matrix dimension of  $\Theta$ , we obtain that

$$(4) \quad \Phi_\vartheta(z) = \nabla F_z(\Theta) = -x(y - \Theta^\top x)^\top = xx^\top \Theta - xy^\top.$$

We denote by  $U \cdot V$  the inner product of two matrices considered as vectors:  $U \cdot V = \text{tr}(U^\top V) = \text{tr}(UV^\top) = \text{tr}(V^\top U) = \text{tr}(VU^\top)$ . Since  $F_z$  is convex for any  $z$ , we have

$$\begin{aligned} (5) \quad d_G(\Theta, \mathcal{Z}) &= d_{\text{loc}}(\Theta, \mathcal{Z}) = d_T^\Phi(\Theta, \mathcal{Z}) = \min_{U \neq 0} \# \{i: U \cdot (x_i x_i^\top \Theta - x_i y_i^\top) \geq 0\} \\ &= \min_{U \neq 0} \# \{i: \text{tr}((x_i x_i^\top \Theta - x_i y_i^\top)^\top U) \geq 0\} = \min_{U \neq 0} \# \{i: -\text{tr}((y_i - \Theta^\top x_i) x_i^\top U) \geq 0\} \\ &= \min_{U \neq 0} \# \{i: -(x_i^\top U)(y_i - \Theta^\top x_i) \geq 0\}, \end{aligned}$$

with  $U$  running over all  $k \times m$  matrices not identically equal to zero.

A question arises about the relationship of not weakly optimal  $\Theta$  and nonfits as defined by Bern and Eppstein (2000). If the two notions are equivalent, that is, pick up identical sets of  $\Theta$ , then the definition of multivariate regression depth given above is equivalent to that of Bern and Eppstein (2000). Despite certain positive evidence (David Eppstein, personal communication), we now believe that the implication holds only in one direction: our definition minimizes over the larger set, thus the multivariate regression depth defined here does not exceed that of Bern and Eppstein (2000); if  $\Theta$  is weakly optimal, it cannot be a nonfit of Bern and Eppstein (2000), but there are nonfits which are not weakly optimal. Note that the inequality has the right direction to ensure that our centerpoint-lower bounds on the multivariate regression depth are applicable also to that of Bern and Eppstein (2000).

EXAMPLE 2. This model is a special case of Example 4; using (5) yields

$$(6) \quad d_G(\vartheta, \mathcal{Z}) = d_{\text{loc}}(\vartheta, \mathcal{Z}) = d_T(\vartheta, \mathcal{Z}) = \min_{u \neq 0} \# \{i: u^T(x_i x_i^T \vartheta - x_i y_i) \geq 0\}$$

Obviously, the criterial functions remain convex. We may reexpress (6) in many ways:

$$(7) \quad d_T(\vartheta, \mathcal{Z}) = \min_{u \neq 0} \# \{i: -(u^T x_i) \operatorname{sgn}(y_i - \vartheta^T x_i) \geq 0\}$$

$$(8) \quad = \min_{u \neq 0} \# \{i: \operatorname{sgn}(u^T x_i) \operatorname{sgn}(y_i - \vartheta^T x_i) \geq 0\}$$

—in all these formulas, we are free to include or drop the minus sign and also to restrict the domain of minimization from  $u \neq 0$  to  $\|u\| = 1$ . Incidentally, (7) can be interpreted as arising from the criterial functions of  $|y - \vartheta^T x|$  (but not all reexpressions should arise in this way).

To see that equations (6)–(8) yield the regression depth whose geometric definition was given by Rousseeuw and Hubert (1999a), note first that the minimized function is piecewise constant; hence the equivalent result is obtained by a minimization over a dense set of directions  $u$ . Fix this set to be the set  $S$  of all  $u \in \mathbb{S}^{p-1}$  such that  $u^T x_i \neq 0$  for all  $x_i$ . When  $x_i = (1, w_i)^T$ , the complement of  $S$  contains only finitely many hyperplanes of lower dimension; the general case is treated similarly below. The expression (8) then equals the minimal proportion over  $S$  of observations such that

$$(9) \quad \text{either } (u^T x_i > 0 \text{ and } \operatorname{sgn}(y_i - \vartheta^T x_i) \geq 0) \quad \text{or} \quad (u^T x_i < 0 \text{ and } \operatorname{sgn}(y_i - \vartheta^T x_i) \leq 0).$$

It is perhaps interesting to mention an equivalent possibility that recently appeared in Adrover, Maronna and Yohai (2000): take the minimal proportion of observations such that

$$\operatorname{sgn} \left( \frac{y_i - \vartheta^T x_i}{u^T x_i} \right) \geq 0.$$

The set of observations satisfying (9) corresponds to the shaded area at Fig. 1. All observations lying on the solid line are included, no observation is expected to lie on the vertical line. The set given by (9) contains all observations that are met by a line during its rotation to the vertical position, in the sense indicated at Fig. 1.

The regression without the intercept (when  $x_i = w_i$ ) behaves similarly; we only have to pay a separate attention to points with  $x = 0$ . They have no influence on the resulting fit, but increase



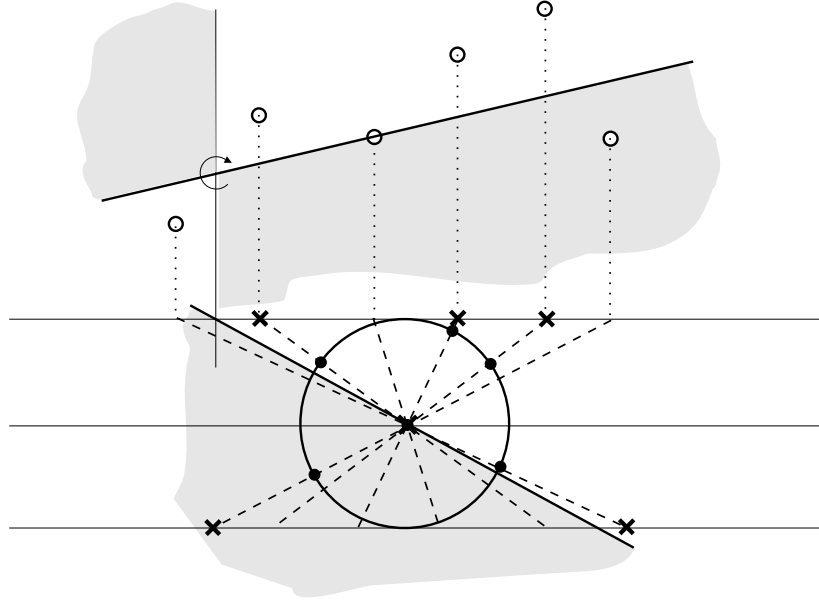


FIGURE 1. Regression depth as location depth.

uniformly the depth of all fits; in fact, we may remove them all without destroying the weak optimality of any fit.

Fig. 1 illustrates how tangent depth may be evaluated in simple regression (and it is not hard to see how the same method works in general). We evaluate the depth of the (solid) regression line  $\vartheta = (\alpha, \beta)^T$ , with respect to the (small empty circles) observations  $(w_i, y_i)$ ; recall that  $x_i = (1, w_i)^T$ . We start by projecting the observations onto the covariate space  $\{1\} \times \mathbb{R}^{p-1}$  (indicated by dotted lines). Then we project them further (along dashed lines) to the unit sphere—that is, we plot their normalized directions (solid circles), in the positive or negative halfspace according to the sign of the residual; if the residual is zero, as happens for one point at Fig 1, the corresponding projection goes to the origin (slightly obscured). The desired result is obtained as the location depth of the origin with respect to these projections (a minimizing halfspace is indicated by the shaded area). Points represented by crosses show how the same method would equivalently work for  $F_z(\vartheta) = |y - \vartheta^T x|$ , when (7) replaces (6).

**EXAMPLE 3.** All previous examples involved convex criterial functions. General nonlinear regression models provide exceptions—yet in many cases local and tangent, or global and local depth are equal. The more thorough analysis of various nonlinear regression models is beyond the scope of this paper.

**EXAMPLE 5.** Also in this example, we limit our analysis for now to the simple two-dimensional case, where the observations are  $z_i = (x_i, y_i)^T$ . For  $\vartheta \in \mathbb{R} \times \mathbb{S}^1$ , consider local coordinates  $(s, t)^T = \varphi(\vartheta)$  with the inverse  $\varphi^{-1}(s, t) = (s, \beta(t)^T)^T = (s, (-\sin t, \cos t)^T)^T$ . Representing

$\nabla F_z(\vartheta)$  by the partial derivatives in these local coordinates, we obtain

$$\nabla F_z(\vartheta) = -(\beta^T z - s)(1, z^T \Gamma \beta)^T,$$

where  $\Gamma$  denotes the matrix of the clockwise rotation by 90 degrees (in  $\mathbb{R}^2$ ); note that

$$-\Gamma \beta = (-\cos t, -\sin t)^T = \frac{\partial}{\partial t} \beta(t).$$

Let  $\vartheta_0$  be the line represented by local coordinates  $(0, 0)$ , the line  $y = 0$  in  $\mathbb{R}^2$ . Since  $\varphi^{-1}(0, 0) = (0, \beta(0)^T)^T = (0, (0, 1)^T)^T$ , we have, for  $z = (x, y)^T$ ,

$$(10) \quad \nabla F_z(\vartheta_0) = -y(1, x)^T = (-y, -yx)^T.$$

The evaluation of the tangent depth of any  $\vartheta$  can be reduced to this canonical case by the appropriate rotation and translation (which leave all orthogonal distances unchanged). In the simple linear regression model, the same line is represented as  $y = 0x + 0$ ; equation (4) gives that  $\nabla F_z(\vartheta) = -y(1, w)^T$ . Recall that  $w$  in simple linear regression corresponds in this example to  $x$ ; in other words, we obtained the same expression in both cases. Therefore, the orthogonal regression tangent depth of  $\vartheta$  is the regression depth of the line  $y = 0$  after a Euclidean change of coordinates that carries the line represented by  $\vartheta$  to the line  $y = 0$ . We obtained a geometric definition of the orthogonal regression tangent depth (Fig. 2 left): 1. given a line orthogonal to line represented by  $\vartheta$ , take the minimum of the proportion of observations in the shaded area, and of those in its complement (counting those lying on the line represented by  $\vartheta$  in both cases); 2. minimize over all orthogonal lines that do not contain any of observations.

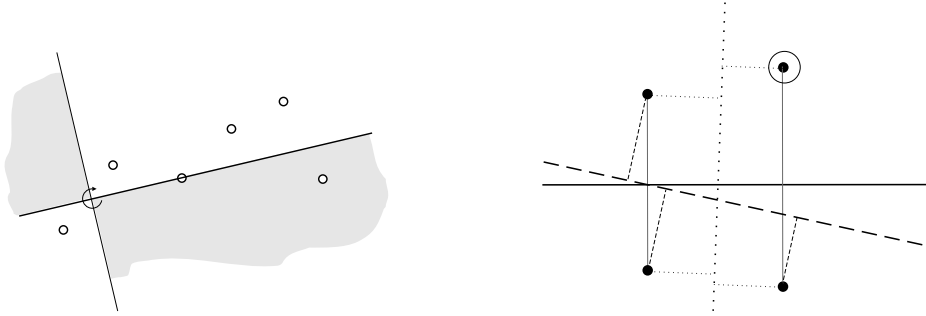


FIGURE 2. Depth in orthogonal regression.

Fig. 2 right shows that the orthogonal regression model is interesting also from another aspect: each of the depth notions introduced above may result in a different value. Consider the fit represented by the solid line. It is easily seen that its tangent depth is 2. After removing the circled point, the orthogonal residuals of the dashed line uniformly supersede those of the solid one—therefore the local depth of the solid line is at most 1; on the other hand, the fit represented by the solid line is locally weakly optimal—hence its local depth is at least 1. Finally, the orthogonal residuals of the dotted line uniformly supersede those of the solid one; hence the fit is not globally weakly optimal and its global depth is 0.

## 3. TANGENT DEPTH IN MEASURE-THEORETIC SETTING

**3.1. Preliminaries.** Under “a measure on  $\mathbb{X}$ ” we always understand a measure defined on the appropriate Borel  $\sigma$ -algebra of a separable metrizable space  $\mathbb{X}$ ; examples of  $\mathbb{X}$  are  $\mathbb{R}^p$ ,  $\mathbb{S}^{p-1}$ , the  $p$ -dimensional unit ball  $\mathbb{D}^p = \{x \in \mathbb{R}^p : \|x\| \leq 1\}$  or the projective plane  $\mathbb{RP}^{p-1}$  (arising from  $\mathbb{S}^{p-1}$  by identifying the antipodal points). Halfspaces and similar geometric constructions are thus measurable; we implicitly assume the measurability of any set or function under consideration. With the only exception of Lebesgue measure, all measures we work with are “subprobability” measures—bounded measures with total mass not exceeding 1.

We denote by  $f \circ g$  the composition mapping assigning the value  $f(g(x))$  to  $x$ . Particularly, if  $P$  is a measure and  $g$  a function on  $\mathbb{X}$ , then  $P \circ g^{-1}$  is another measure with the same total mass, assigning  $P\{x : g(x) \in E\}$  to any  $E$ .

**3.2. Tangent depth in probability fields.** For statistical considerations, we need an analog of depth defined for population distributions—we need a model for sampling, or target for asymptotics, or both. The following general definition of the tangent depth extends that given in Section 2. It stems from the fact that formula (2) can be easily rewritten to depend on the empirical probability supported by the gradient points. The analogy (“plug-in”) principle suggests then to replace the empirical probability by a general one.

Let  $\mathbb{X}$  be any space where halfspaces can be reasonably defined:  $\mathbb{R}^p$ , but also  $\mathbb{S}^{p-1}$  and  $\mathbb{D}^p$  whose halfspaces are formed by the intersection of those from  $\mathbb{R}^p$  with  $\mathbb{S}^{p-1}$  or  $\mathbb{D}^p$ , respectively. We will later add also the cosmic space  $\bar{\mathbb{R}}^p$  (see 5.3). For any measure  $Q$  on  $\mathbb{X}$ , we define

$$d(Q) = \inf_{u \neq 0} Q(H_u) = \inf_{\|u\|=1} Q(H_u).$$

Note that a minimizing halfspace may not exist, therefore the use of  $\inf$  instead of  $\min$  is essential. The infimum, on the other hand, should not be necessarily taken over the full set of directions.

**PROPOSITION 3.1.** *Let  $\mathbb{X}$  be  $\mathbb{R}^p$ ,  $\mathbb{S}^{p-1}$ ,  $\mathbb{D}^p$  or  $\bar{\mathbb{R}}^p$ . For any measure  $Q$  on  $\mathbb{X}$ ,  $d(Q) = \inf_{u \in S} Q(H_u)$  whenever  $S$  is a dense subset of  $\mathbb{S}^{p-1}$ .*

An important instance of a set  $S$  appearing in Proposition 3.1 is given by the following theorem.

**THEOREM 3.2.** *Let  $\mathbb{X}$  be  $\mathbb{R}^p$ ,  $\mathbb{S}^{p-1}$ ,  $\mathbb{D}^p$  or  $\bar{\mathbb{R}}^p$ . For any probability measure  $Q$  on  $\mathbb{X}$ ,  $\mu(S(Q)) = 1$ , where  $\mu$  is the uniform distribution on  $\mathbb{S}^{p-1}$  and  $S(Q)$  is:*

- for  $\mathbb{X} = \mathbb{R}^p$ ,  $\mathbb{D}^p$  or  $\bar{\mathbb{R}}^p$ , the set of all  $u \in \mathbb{S}^{p-1}$  such that  $Q(\partial H_u \setminus \{0\}) = 0$ ;*
- for  $\mathbb{X} = \mathbb{S}^{p-1}$ , the set of all  $u \in \mathbb{S}^{p-1}$  such that  $Q(\partial H_u) = 0$ .*

Let  $\Theta$  be a metrizable topological space. We define a **probability field** on  $\mathbb{X}$  indexed by  $\Theta$  to be a function  $\{P_\vartheta : \vartheta \in \Theta\}$  assigning a probability measure  $P_\vartheta$  on  $\mathbb{X}$  to each  $\vartheta \in \Theta$ . An important instance is when  $\Theta$  is a  $p$ -dimensional manifold,  $\mathbb{X} = \mathbb{R}^p$  is isomorphic to the tangent space of  $\Theta$  at each its point  $\vartheta$ , and  $\Phi_\vartheta(Z) = \nabla F_Z(\vartheta)$ ; in this case we speak about a *gradient probability field*. The motivation for the name is: if  $F$  is a differentiable function on  $\Theta$ , then its gradient  $\nabla F(\vartheta)$  defines a vector field, a mapping that assigns a vector from the tangent space of  $\Theta$  at  $\vartheta$  to each  $\vartheta \in \Theta$ . A collection of  $n$  differentiable functions produces  $n$  gradients; these  $n$  points support an empirical probability.

Given a probability field, then  $d(P_\vartheta)$  naturally yields a depth function on  $\Theta$ . To proceed more in analogy with the definition of tangent depth in 2.6, let us consider again a function  $\Phi$  from  $\Theta \times \mathcal{Z} \rightarrow \mathbb{X}$ , where  $\mathbb{X}$  is now any space like  $\mathbb{R}^p$  or  $\mathbb{S}^{p-1}$ , where the notion of halfspace makes sense. Let  $\Phi_\vartheta^{-1}$  stand for the preimage set-valued mapping corresponding to the function  $\Phi_\vartheta(\cdot)$ :  $\Phi_\vartheta^{-1}(E) = \{z: \Phi_\vartheta(z) \in E\}$ . If  $P$  is a probability on  $\mathcal{Z}$ ,  $\Phi$  gives birth to a probability field  $P_\vartheta = P \circ \Phi_\vartheta^{-1}$ . We define the **tangent depth** of  $\vartheta$  to be

$$d_T^\Phi(\vartheta, P) = d(P_\vartheta) = d(P \circ \Phi_\vartheta^{-1}) = \inf_{u \neq 0} P(\Phi_\vartheta^{-1}(\mathbf{H}_u)).$$

Again, the definition depends on  $\Phi$  and  $P$ . When  $P$  is a distribution of a random variable  $Z$ , we abuse the notation and write  $d_T^\Phi(\vartheta, Z)$  instead of precise  $d_T^\Phi(\vartheta, \mathcal{L}(Z))$ , having  $P_\vartheta(E) = \mathbb{P}[\Phi(Z, \vartheta) \in E]$  for any  $E$  and

$$d_T^\Phi(\vartheta, Z) = \inf_{u \neq 0} \mathbb{P}[Z \in \Phi_\vartheta^{-1}(\mathbf{H}_u)] = \inf_{u \neq 0} \mathbb{P}[\Phi_\vartheta(Z) \in \mathbf{H}_u].$$

The symbol  $d_T$  is now used in several formally different ways, which nonetheless all express the same concept:  $d_T^\Phi(\vartheta, \mathcal{Z})$  is equal to  $d_T^\Phi(\vartheta, Z) = d_T^\Phi(\vartheta, \mathcal{L}(Z))$ , whenever  $\mathcal{L}(Z)$  is the empirical distribution of  $\mathcal{Z}$ .

All general definition of depth that occurred in the literature so far are special instances of our general definition of tangent depth. We do not introduce analogous extensions for global and local depth in this paper. Such a development might be relatively straightforward, but formal subtleties needed for analogs of Theorems 2.1 and 2.3 would digress us from the main theme.

**EXAMPLE 1** (General version). Recall (3): in the multivariate location setting,  $\Phi_\vartheta(z) = \nabla F_z(\vartheta) = \vartheta - z$ . If  $Z$  is a  $\mathcal{Z}$ -valued random variable with distribution  $P$ , then

$$d_T^\Phi(\vartheta, Z) = \inf_{u \neq 0} \mathbb{P}[u^T(\vartheta - Z) \geq 0] = \inf_{u \neq 0} \mathbb{P}[Z \in \mathbf{H}_{\vartheta, u}] = \inf_{u \neq 0} P_\vartheta(\mathbf{H}_u) = \inf_{u \neq 0} P(\mathbf{H}_{\vartheta, u}),$$

where  $P_\vartheta = \mathcal{L}(\vartheta - Z)$ . The analogy with the finite-sample case is straightforward.

**EXAMPLE 2** (General version). In linear regression setting,  $\Phi_\vartheta(z) = \Phi_\vartheta(x, y) = \nabla F_z(\vartheta) = xx^T\vartheta - xy$ . Thus,  $P_\vartheta = \mathcal{L}(XX^T\vartheta - XY)$ . For a random element  $Z = (X, Y)$ , we obtain straightforward analogs of (6), (7), and (8):

$$\begin{aligned} d_T^\Phi(\vartheta, Z) &= \inf_{u \neq 0} P_\vartheta(\mathbf{H}_u) = \inf_{u \neq 0} \mathbb{P}[u^T(XX^T\vartheta - XY) \geq 0] \\ &= \inf_{\|u\|=1} \mathbb{P}[-u^T X \operatorname{sgn}(Y - X^T\vartheta) \geq 0] \\ &= \inf_{u \in \mathbb{S}^{p-1}} \mathbb{P}[\operatorname{sgn}(u^T X) \operatorname{sgn}(Y - X^T\vartheta) \geq 0]. \end{aligned}$$

To establish the equivalence to the geometric definition of Rousseeuw and Hubert (1999a), we have to invoke Proposition 3.1. Consider the set  $S$  of directions  $u$  such that  $\mathbb{P}[X^T u = 0]$ . Theorem 3.2 says that such a set  $S$  is dense in  $\mathbb{S}^{p-1}$ , if  $\mathbb{P}[X = 0] = 0$ . Once this holds,  $d_T(\vartheta, Z)$  is equal to the infimum of

$$\mathbb{P}[u^T X > 0 \text{ and } \operatorname{sgn}(Y - \vartheta^T X) \geq 0 \quad \text{or} \quad u^T X < 0 \text{ and } \operatorname{sgn}(Y - \vartheta^T X) \leq 0],$$

over  $u$  satisfying  $\mathbb{P}[u^\top X] = 0$ , as in the finite-sample case. When the regression is with intercept, the first coordinate of  $X$  is equal to 1, hence  $\mathbb{P}[X \neq 0] = 1$ . The regression without intercept can be treated similarly: we just decompose the distribution of  $Z$  to the part with  $X = 0$  and the rest (see Section 5 for a more formal treatment).

EXAMPLE 4 (General version). Also for this example, the general version is a straightforward extension of the finite-sample case:  $\Phi_\vartheta(z) = \nabla F_z(\Theta) = xx^\top \Theta - xy^\top$ ,  $P_\vartheta = \mathcal{L}(XX^\top \Theta - XY^\top)$  and

$$\begin{aligned} d_T^\Phi(\Theta, \mathcal{Z}) &= \inf_{U \neq 0} \mathbb{P}[U \cdot (XX^\top \Theta - XY^\top) \geq 0] \\ &= \inf_{U \neq 0} \mathbb{P}[\text{tr}((XX^\top \Theta - XY^\top)^\top U) \geq 0] \\ &= \inf_{U \neq 0} \mathbb{P}[-(X^\top U)(Y - \Theta^\top X) \geq 0], \end{aligned}$$

a direct analog of (5).

**3.3. Existence of centerpoints.** If the parametric space is a manifold, then we may speak about its dimension (equal to the dimension of its tangent space); from the statistical point of view, it is the number of independent parameters in the model. For instance, when  $\Theta = \mathbb{R}^p$ , then its dimension is simply  $p$ . A *centerpoint* is a parameter whose (tangent) depth is not less than  $1/(1 + \dim \Theta)$ .

EXAMPLE 1. In the multivariate location model, the existence of a centerpoint for any probability  $P$  on  $\mathbb{R}^p$  is a well-known mathematical result, proved by Neumann (1945) for  $p = 2$  and by Rado (1946) in general. Birch (1959) gave an alternative proof, similar to that by Donoho and Gasko (1992), who rediscovered the result for statistics. A very short and elementary proof for finite-sample point configurations can be found on page 66 of Edelsbrunner (1987).

EXAMPLE 2. The centerpoint problem in the linear regression setting was raised by Rousseeuw and Hubert (1999a), who also established the important special case: they proved the existence of a centerpoint for  $p = 2$ , for simple regression. Their ingenious geometric construction based on the ham-and-sandwich theorem seems to be, however, not extendable to a higher-dimensional case, where they conjectured the same general lower bound  $1/(p+1)$  for the maximal depth: separately for all finite-sample data—Conjecture 1(a) in Rousseeuw and Hubert (1999a), Conjecture 1 in Rousseeuw and Hubert (1999b)—and for any absolutely continuous population distribution—Conjecture 1(b) in Rousseeuw and Hubert (1999a). The following theorem establishes the lower bound for maximal depth for *any* random variable  $Z = (X, Y)$ , settling thus both conjectures of Rousseeuw and Hubert (1999b).

THEOREM 3.3. *Let  $\Phi_\vartheta(z) = \Phi_\vartheta(x, y) = -x(y - x^\top \vartheta)$  for any  $\vartheta \in \Theta = \mathbb{R}^p$  and any  $(x, y)$  from  $\mathbf{X} \times \mathbb{R}$ , where  $\mathbf{X} \subseteq \mathbb{R}^p$ . For every probability  $P$  on  $\mathbf{X} \times \mathbb{R}$ ,*

$$(11) \quad \sup_{\vartheta \in \Theta} d_T^\Phi(\vartheta, P) \geq \frac{1}{1 + \dim \Theta}.$$

The proof of Theorem 3.3 in its full generality is long and employs technical constructions specific for the univariate regression setting. An easier, restricted, but also more general version

of Theorem 3.3 is obtained under an additional assumption involving the *identification index*

$$\Delta(X) = \sup_{\vartheta \neq 0} \mathbb{P}[X^\top \vartheta = 0],$$

the quantity well-known from breakdown and consistency considerations in linear regression. Note that the definition of  $\Delta(X)$  does not depend on whether  $\vartheta$  runs over  $\mathbb{R}^k$  or  $\bar{\mathbb{R}}^k$ .

EXAMPLE 4 (Centerpoints). Amenta, Bern, Eppstein and Teng (2000) gave another proof of Conjecture 1(a), a more geometric one. Later, Bern and Eppstein (2000) introduced (in the finite-sample setting) a multivariate generalization of regression depth, and conjectured a lower bound on the maximal depth for their definition. Motivated by this development, we formulated Theorem 3.4 in the revised version to cover the multivariate case as well. Theorem 3.4 relies on the general Theorem 5.19 and Proposition 5.21.

THEOREM 3.4. *Let  $\Phi_\vartheta(z) = \Phi_\vartheta(x, y) = -x(y^\top - x^\top \vartheta)$  for any  $\vartheta \in \Theta = \mathbb{R}^p = \mathbb{R}^{mk}$  and any  $(x, y)$  from  $\mathbf{X} \times \mathbb{R}^m$ , where  $\mathbf{X} \subseteq \mathbb{R}^p$ . For every probability  $P$  on  $\mathbf{X} \times \mathbb{R}$  the following holds: if  $\Delta(X) < 1/(1+p)$  for any random variable  $Z = (X, Y)$  whose distribution is  $P$ , then there is  $\dot{\Theta} \in \mathbb{R}^p$  such that*

$$d_T^\Phi(\dot{\Theta}, P) = \sup_{\Theta \in \Theta} d_T^\Phi(\Theta, P) \geq \frac{1}{1 + \dim \Theta}.$$

The additional assumption introduced by Theorem 3.4 is probably more restrictive from the mathematical than statistical aspect: it says that the regressors are in (fairly) general position, a condition satisfied in many real situations. In simple linear regression, for instance, more than one third of observations should lie in a covariate point to make the assumption invalid. Note also that if the distribution of  $X$  is absolutely continuous, then  $\Delta(X) = 0$ ; thus Theorem 3.4 immediately proves the original Conjecture 1(b) of Rousseeuw and Hubert (1999a).

The only missing link between Theorem 3.4 and Theorem 3.3 is our inability to handle the smoothing approximation successfully. For the details, see Section 5 and the Appendix; here we only mention that the whole question can be reduced to a problem whether a regression probability field generated by  $P$  can be approximated by a smooth enough (absolute continuity is enough but is not necessary) regression probability fields generated by  $P_\nu$ , in a way that

$$(12) \quad \limsup_{\nu \rightarrow \infty} \sup_{\vartheta \in \mathbb{R}^p} d_T(\vartheta, P_\nu) \leq \sup_{\vartheta \in \mathbb{R}^p} d_T(\vartheta, P).$$

If yes, then the conclusion of Theorem 3.4 holds for  $P$  even if the assumption on  $\Delta(X)$  is not satisfied. Thus, if the existence of such an approximation for general  $P$  would yield a theorem more general than Theorem 3.3. We do not know the answer to this question, though a positive one seems plausible; we would like to remark only we are aware of counterexamples showing that mere weak convergence  $P_\nu$  to  $P$  does not imply (12).

On the other hand, however, it seems that the proper treatment of the smoothing trick would not illuminate any new feature of the problem and is only a technical necessity. To bypass it, we had to undertake a painstaking way expounded in Section 6, where the proof of Theorem 3.3 follows basically the same scheme as that of Theorem 3.4, but on a considerably higher technical level. Realizing the importance of linear regression in statistics, we believe that it is essential to

have Theorem 3.3 rigorously proved in the maximal possible generality—in fact, Theorem 3.4 does not cover even all possible finite-sample cases. Also, our Theorem 3.3 may be of independent mathematical interest. Unfortunately, the technique does not extend to multivariate regression depth, where the complete solution of the centerpoint problem remains open.

Note that Theorem 3.3 guarantees only a bound on the maximal depth. The existence of a centerpoint then follows from the existence of a point with maximal depth (in  $\mathbb{R}^p$ ). This fact is not obvious, though in many cases comes trivially: for instance, if  $P$  is the empirical distribution of a finite sample. It also follows from Theorem 3.4 under the assumption  $\Delta(X) < 1/(p+1)$ . For (univariate) linear regression, the general existence of the deepest parameter is established in the following theorem.

**THEOREM 3.5.** *Let  $\Phi_{\vartheta}(z) = \Phi_{\vartheta}(x, y) = -x(y - x^T \vartheta)$  for any  $\vartheta \in \Theta = \mathbb{R}^p$  and any  $(x, y)$  from  $\mathbb{X} \times \mathbb{R}$ , where  $\mathbb{X} \subseteq \mathbb{R}^p$ . For any probability  $P$  on  $\mathbb{R} \times \mathbb{X}$ , there is  $\dot{\vartheta} \in \mathbb{R}^p$  such that*

$$d_T^{\Phi}(\dot{\vartheta}, P) = \sup_{\vartheta \in \mathbb{R}^p} d_T^{\Phi}(\vartheta, P).$$

The proof is a direct consequence of Theorem 5.25, formulated and proved in Section 5. Again, we do not know whether the same general result—although plausible—holds for the multivariate regression depth.

#### 4. BIAS, BREAKDOWN POINT, AND MAXIMUM DEPTH ESTIMATORS

**4.1. Bias sets of maximum depth estimators and depth contours.** One of the reasons for the study of centerpoints are their implications on bias and breakdown of maximum depth estimators. For location and regression, these aspects were illustrated by Donoho and Gasko (1992), Rousseeuw and Hubert (1999a), and Van Aelst and Rousseeuw (2000). Here we give a view from more general perspective.

We formalize *an estimator* as a mapping  $\mathcal{T}$  assigning one or more parameters from  $\Theta$  to any probability  $P$  with values in  $\mathbb{Z}$  (many robust estimators yield non-unique results for certain datasets, thus we have to work at this level of generality). The setting in which data are represented by probabilities encompasses a variety of situations and is also perfectly relevant for finite-sample data (see below). If  $\mathcal{E}$  is a set of probabilities, then  $\mathcal{T}(\mathcal{E})$  denotes the set of all possible values of  $\mathcal{T}$  under  $P \in \mathcal{E}$ . We are interested in the behavior of  $\mathcal{T}(\mathcal{E})$  when  $\mathcal{E}$  is a neighborhood of  $P$ .

Neighborhoods are constructed with the help of a distance  $\pi$  defined on the space of probabilities on  $\mathbb{Z}$ . As a rule, this distance depends only on the laws of random variables under consideration. It may be a metric, but not necessarily; all we need is that for any random variable  $Z$ , the balls  $\mathbf{B}_{\pi}(P, \varepsilon) = \{\tilde{P} : \pi(P, \tilde{P}) < \varepsilon\}$  decrease with  $\varepsilon$  and shrink to  $\{P\}$  for  $\varepsilon = 0$ . A frequent choice for  $\pi$  is the total variation metric:  $\pi(P, \tilde{P}) = \inf |P(A) - \tilde{P}(A)|$ , inf taken over all (measurable)  $A$ . Another popular (and non metric) choice is the contamination distance  $\gamma$ :  $\gamma(P, \tilde{P}) \leq \varepsilon$  if  $\tilde{P} = (1 - \varepsilon)P + \varepsilon Q$  for some probability  $Q$  on  $\mathbb{Z}$ . The inequality  $\gamma(P, \tilde{P}) \leq v(P, \tilde{P})$

implies that  $\mathbf{B}_\gamma(P, \varepsilon) \subseteq \mathbf{B}_v(P, \varepsilon)$  for all  $P$  and any  $\varepsilon > 0$ . There may be good reasons for adopting other distances, as well as to reject certain other ones—for a thorough discussion, see Davies (1993).

Maximum depth estimators, defined in particular statistical models, can be viewed as generalizations of the sample median. Given a function  $\Phi$  from  $\Theta \times \mathcal{Z}$  to  $\mathbb{R}^p$ , we call  $\mathcal{T}$  a **maximum depth estimator** if

$$\mathcal{T}(P) = \{\dot{\vartheta} \in \Theta: d_T^\Phi(\dot{\vartheta}, P) \geq d_T^\Phi(\vartheta, P) \text{ for all } \vartheta \in \Theta\}.$$

The following theorem shows that bias sets of any maximum depth estimator are closely related to the upper level sets of depth. This gives another interpretation for depth contours in data analysis—in addition to those found in Liu, Parelius and Singh (1999).

**THEOREM 4.1.** *Let  $\Phi_\vartheta$  be, for any  $\vartheta \in \Theta$ , a (measurable) function from  $\mathcal{Z}$  to  $\mathbb{R}^p$  and let  $\mathcal{T}$  be a maximum depth estimator. For any  $\varepsilon \geq 0$  and any  $P$ , the inclusion*

$$(13) \quad \mathcal{T}(\mathbf{B}_\pi(P, \varepsilon)) \subseteq \{\vartheta: d_T^\Phi(\vartheta, P) \geq \delta\}$$

*holds with*

- (i)  $\delta = \eta - \varepsilon$ , if  $d_T(\dot{\vartheta}, \tilde{P}) \geq \eta$  for all  $\dot{\vartheta} \in \mathcal{T}(\tilde{P})$  and all  $\tilde{P}$ ;
- (ii)  $\delta = \eta - 2\varepsilon$ , if  $d_T(\dot{\vartheta}, P) \geq \eta$  for all  $\dot{\vartheta} \in \mathcal{T}(P)$  and  $\pi = v$ ;
- (iii)  $\delta = \eta(1 - \varepsilon) - \varepsilon$ , if  $d_T(\dot{\vartheta}, P) \geq \eta$  for all  $\dot{\vartheta} \in \mathcal{T}(P)$  and  $\pi = \gamma$ .

We believe that in many concrete cases the inclusion (13) is actually the equality. The attainable uniform bound required by (i) may be lower than the depth of the maximum depth estimator—it means that (ii) or (iii) often give sharper bounds.

**4.2. Breakdown points of maximum depth estimators.** An important indicator of the bias behavior is the breakdown point, which says at which  $\varepsilon$  the estimator “breaks down”: the bias sets start to be unacceptably rich, that is, unbounded (or containing all parameter points or at least some unacceptable ones). In this paper, we limit our analysis to parameter spaces equal to  $\mathbb{R}^p$  or similar—in other words, endowed with a structure of “boundedness”. Knowing once what “bounded” means (in  $\mathbb{R}^p$ : a bounded set is contained in a ball with finite perimeter), we define the **breakdown point** of  $\mathcal{T}$  at  $P$  to be

$$\varepsilon_\pi^*(\mathcal{T}, P) = \inf\{\varepsilon > 0: \mathcal{T}(\mathbf{B}_\pi(P, \varepsilon)) \text{ is not bounded in } \Theta\}.$$

Note that when we restrict our attention to those  $P$  whose laws are empirical distributions of the  $n$ -tuples—provided that the same restriction is applied to all probabilities appearing in the definition of  $\mathbf{B}_v(P, \varepsilon)$ —then for  $\pi = v$  we obtain nothing but the popular finite-sample Donoho-Huber replacement breakdown point. (Thus, it is not “another” breakdown theory which we study here, but the one perfectly relevant also for the finite-sample case. Rigorously, it is not hard to see that the finite-sample Donoho-Huber replacement breakdown point is always bounded from below by  $\varepsilon_v^*(P)$  evaluated at  $P$  whose law is the empirical distribution of a dataset  $\mathcal{Z}$ ; the equality actually holds in all but some artificial cases.)



Let  $\mathcal{T}$  be a maximum depth estimator. A quantity influencing the breakdown point of  $\mathcal{T}$  is “the depth at breakdown”

$$d_T(\infty, P) = \inf_{A \text{ bounded}} \sup\{d_T^\Phi(\vartheta, P) : \vartheta \in A^c\},$$

where  $\inf$  is taken over all bounded subsets  $A$  of  $\Theta$ . It is not hard to see that when  $\Theta = \mathbb{R}^p$ ,  $d_T^\Phi(\infty, P)$  is equal to the supremum of  $\limsup_{\nu \rightarrow \infty} d_T^\Phi(\vartheta_\nu, P)$  taken over all sequences such that  $\|\vartheta_\nu\| \rightarrow \infty$ .

**THEOREM 4.2.** *Let  $\Phi_\vartheta$  be, for any  $\vartheta \in \Theta$ , a (measurable) function from  $Z$  to  $\mathbb{R}^p$  and let  $\mathcal{T}$  be a maximum depth estimator.*

(i) *If  $d_T(\dot{\vartheta}) \geq \eta$  for all  $\dot{\vartheta} \in \mathcal{T}(\tilde{P})$  and all  $\tilde{P}$ , then  $\varepsilon_\pi^*(\mathcal{T}, P) \geq \eta - d_T^\Phi(\infty, P)$  (for both  $\pi$  equal to  $v$  or  $\gamma$ ).*

(ii) *If  $d_T(\dot{\vartheta}) \geq \eta$  for all  $\dot{\vartheta} \in \mathcal{T}(P)$ , then*

$$\varepsilon_v^*(\mathcal{T}, P) \geq \frac{1}{2}(\eta - d_T^\Phi(\infty, P))$$

and

$$\varepsilon_\gamma^*(\mathcal{T}, P) \geq \frac{\eta - d_T^\Phi(\infty, P)}{1 + \eta}.$$

**EXAMPLE 1** (Breakdown point). In the location setting,  $d_T(\infty, P) = 0$  for any  $P$ . Thus, Theorem 4.2(i) gives that both  $\varepsilon_v^*(\mathcal{T}, P)$  and  $\varepsilon_\gamma^*(\mathcal{T}, P)$  are bounded by  $1/(p+1)$  from below,  $p = \dim \Theta$ . If the depth of the Tukey median  $\mathcal{T}$  is  $1/2$ , as often happens for symmetric distributions, then by Theorem 4.2(ii),  $\varepsilon_v^*(\mathcal{T}, P) \geq 1/4$  and  $\varepsilon_\gamma^*(\mathcal{T}, P) \geq 1/3$ , regardless of the dimension.

**EXAMPLE 2** (Breakdown point). Let  $Z = (X, Y)$  be a random variable whose distribution is  $P$ . The theory developed in Section 5 yields upper bounds for  $d_T(\infty, P)$ . The cruder bound  $d_T(\infty, P) \leq \Delta(X)$  follows from Proposition 5.21(ii) and Proposition 5.12(i); Theorem 4.2(i) then yields the general lower bound  $1/(1+p) - \Delta(X)$  for the breakdown point of the maximum depth estimator. For instance, if  $Z$  corresponds to the finite-sample data  $n$ -tuple in general position, then  $\Delta(X) = (p-1)/n$  and the finite sample breakdown point is not less than  $n/(p+1) - p+1$ .

Theorem 4.2(ii) yields the dimension-free bound for the maximum depth estimator: if its depth is  $1/2$  then its breakdown point is  $(1-2\Delta(X))/4$  for the total variation and  $(1-2\Delta(X))/3$  for the contamination breakdown point. Van Aelst and Rousseeuw (2000) and Van Aelst, Rousseeuw, Hubert and Struyf (2000) considered the situation when “the model holds”: there is a parameter  $\vartheta_0$  such that the conditional distribution of  $Z$  given  $X$  is symmetric about  $\ell(\vartheta_0)$ . In such a case, the depth of  $\vartheta_0$  is  $\delta_0 + (1 - \delta_0)/2$  with  $\delta_0 = \mathbb{P}[X^\tau \vartheta_0 = Y]$ . Theorem 5.25 then yields the bound

$$d_T(\infty, P) \leq \Delta_0 + \frac{1}{2}(\Delta(X) - \Delta_0),$$

where  $\Delta_0 = \sup_{\vartheta \neq 0} \mathbb{P}[X^\tau \vartheta = Y \text{ and } X^\tau \vartheta = 0] \leq \delta_0$ ; by Theorem 4.2(ii), we obtain

$$\varepsilon_v^*(\mathcal{T}, P) \geq \frac{1}{4}(1 + \delta_0 - \Delta(X) - \Delta_0) \geq \frac{1}{4}(1 - \Delta(X)).$$

For the contamination breakdown, Theorem 4.2(iii) gives that

$$\varepsilon_\gamma^*(\mathcal{T}, P) \geq \frac{1 + \delta_0 - \Delta(X) - \Delta_0}{3 + \delta_0}.$$

We believe that these bounds are in general sharp; note that if  $Z$  has an absolutely continuous distribution, they reduce to  $1/4$  and  $1/3$ , respectively.

## 5. MATHEMATICAL DETAILS

**5.1. Weak optimality, convexity.** In this subsection, we prove all results from Section 2. Let  $F$  be a real function  $F$  defined on a convex domain  $\Theta$ . Following the terminology of Ponstein (1967), we call  $F$  **quasi-convex** if all upper-level sets are convex ( $F(x) \leq \min\{F(a), F(b)\}$  for any  $x$  lying on the line connecting  $a$  and  $b$ ). We call  $F$  **strictly quasi-convex** if the following holds: given any  $a, b$  with  $F(a) < F(b)$ ,  $F(ta + (1-t)b) < F(b)$  for every  $t \in (0, 1)$ . We call  $F$  **pseudo-convex** if given any  $a, b$  with  $F(a) < F(b)$ , there exists  $c > 0$  and  $\tau \in (0, 1]$  such that  $F(ta + (1-t)b) \leq F(b) - ct$  for every  $t \in [0, \tau]$ .

PROPOSITION 5.1. (i) *If  $F$  is convex, then it is pseudo-convex.*

(ii) *If  $F$  is pseudo-convex and continuous, then it is strictly quasi-convex.*

(iii) *If  $F$  is pseudo-convex, differentiable and  $\nabla F(\vartheta) = 0$ , then  $\vartheta$  is its point of global minimum.*

(iv) *If  $F$  is strictly quasi-convex or pseudo-convex, then it is quasi-convex.*

PROOF. See Ortega and Rheinboldt (1970), pages 102–105 (in a slightly different terminology there).  $\square$

PROOF OF THEOREM 2.1. Suppose that we removed enough observations to destroy local weak optimality. Then  $\vartheta$  cannot be globally weakly optimal; hence we would need not more observations to destroy the global optimality. This proves the general inequality.

To prove the equality under the additional hypothesis, note that if  $\vartheta$  is not a weak optimum in  $\Theta$  with respect to  $A^c$ , then there is  $\tilde{\vartheta} \in \Theta$  such that  $F_i(\tilde{\vartheta}) < F_i(\vartheta)$  for all  $i \in A^c$ . The strict quasi-convexity of all  $F_i$  implies then that  $\vartheta$  cannot be weakly optimal in any open neighborhood of  $\vartheta$ , with respect to the same  $A^c$ ; this proves the converse inequality and hence equality.  $\square$

PROOF OF PROPOSITION 2.2. The first part of the proposition is just a slightly different statement of the Carathéodory theorem. The equivalent conditions follow from direct combinations of known convexity properties: note that (ii) is the consequence of the Minkowski separation theorem.  $\square$

The next theorem gives a sufficient condition for weak optimality when criterial functions are quasi-convex; compare with the necessary and sufficient conditions for Pareto optimality given on page 71 of Smale (1975a).

THEOREM 5.2. *Suppose that  $F_1, F_2, \dots, F_k$  are differentiable at  $\vartheta$ . If  $\vartheta$  is locally weakly optimal with respect to  $F_1, F_2, \dots, F_k$ , then 0 is surrounded by  $\{\nabla F_1(\vartheta), \nabla F_2(\vartheta), \dots, \nabla F_k(\vartheta)\}$ , that is,*

$$(S) \sum \lambda_i \nabla F_i(\vartheta) = 0 \text{ for some nonnegative and not all equal to zero constants } \lambda_1, \lambda_2, \dots, \lambda_k.$$

*Conversely,  $\vartheta$  is a weak local optimum, if (S) holds, all  $F_i$  are quasi-convex, and*

$$(M) \text{ any } F_i \text{ with nonzero } \lambda_i \text{ in (S) is locally nondecreasing in any direction } u \text{ such that } u^T \nabla F_i(\vartheta) = 0.$$

PROOF. If (S) does not hold, then all  $\nabla F_i(\vartheta)$  are contained in an open halfspace with 0 on its boundary, by Proposition 2.2. Hence, there is a direction  $u$  such that  $u^T \nabla F_i(\vartheta) < 0$  for all  $i$ ; therefore, all  $F_i$  are locally decreasing in the direction of  $u$  and hence  $\vartheta$  cannot be weakly optimal.

To prove the converse, we will show that under (S) and (M) weak optimality cannot be violated in any direction—this is sufficient in view of quasi-convexity of criterial functions. Fix a direction  $u$ . If there is  $F_i$  such that  $u^T \nabla F_i(\vartheta) < 0$ , then this  $F_i$  is locally increasing in the direction of  $u$  and the desired conclusion holds. If there is no such  $F_i$ , then all  $\nabla F_i(\vartheta)$  lie in a subspace  $\{x: x^T u = 0\}$ . By (S), there is at least one nonzero  $\lambda_i$ ; by (M), the corresponding  $F_i$  is locally nondecreasing in the direction  $u$  and we obtain the same conclusion again.  $\square$

PROPOSITION 5.3. *Suppose that  $F$  is differentiable at  $\vartheta$ . If  $F$  is pseudo-convex, then  $F$  is nondecreasing in any direction  $u$  such that  $u^T \nabla F(\vartheta) = 0$ .*

PROOF. If  $u^T \nabla F(\vartheta) = 0$ , then the directional derivative of  $F$  in the direction  $u$  is 0. The proposition follows from Proposition 5.1(iii) and the fact that the restriction of a pseudo-convex  $F$  on the line in direction of  $u$  passing through  $\vartheta$  is again pseudo-convex.  $\square$

PROOF OF THEOREM 2.3. If there is a halfspace  $H_u$  containing a proportion  $d_T(\vartheta, \mathcal{Z})$  of the gradients, then after removing the observations corresponding to the gradients from  $H_u$  the remaining ones do not surround 0, due to Proposition 2.2(ii). Theorem 5.2 then implies that  $\vartheta$  cannot be locally weakly optimal. This proves the inequality in the general situation.

To prove the converse inequality under the additional hypothesis, suppose that we have removed a proportion  $d_{\text{loc}}(\vartheta, \mathcal{Z})$  of the observations so that  $\vartheta$  is not anymore a local weak optimum with respect to the remaining ones. By Theorem 5.2, this means that then either (S) or (M) must be violated. It must be (S), since (M) is implied for pseudo-convex  $F_i$  by Proposition 5.3. The remaining observations therefore do not surround 0; but then they lie in an open halfspace, due to Proposition 2.2(ii). Hence  $d_T(\vartheta, \mathcal{Z}) \leq d_{\text{loc}}(\vartheta, \mathcal{Z})$ .  $\square$

**5.2. Atomic decomposition of measures.** For the proof of Theorem 3.2, we need a combinatorial technique generalizing the principle that a bounded measure may possess only a countable number of atoms. The same technique is used for the proof of Lemma 6.3 in Section 6; unfortunately, we do not know about any reference in the literature which would allow us to skip the following technical development.

We introduce the following notation:  $\mathbb{X}$  is the basic set (a measurable space, for instance,  $\mathbb{S}^{p-1}$ ,  $\mathbb{R}^p$  or any other; we do not require any topological assumptions) and  $\mathcal{A}$  is an *atomic system*: a set of (measurable) subsets of  $\mathbb{X}$  that can be written as a union

$$(14) \quad \mathcal{A} = \mathcal{A}^0 \cup \mathcal{A}^1 \cup \mathcal{A}^2 \cup \dots \cup \mathcal{A}^r$$

such that  $\mathcal{A}^0$  consists of a single element (often  $\emptyset$ , but not always) and the following property holds: if  $A \in \mathcal{A}^i$ ,  $E \in \mathcal{A}^j$  and  $A \subseteq E$ , then  $i \leq j$ ; and if  $A \subset E$ , then  $i < j$  (that is, equality of  $i$  and  $j$  implies the equality of  $A$  and  $E$ ). It is not hard to see that any atomic system has the property of the intersection system introduced by Balek and Mizera (1997): the intersection of any two distinct elements from  $\mathcal{A}^i$  belongs to  $\mathcal{A}^j$  with  $j < i$ .

A canonical example is  $\mathcal{A} = \mathcal{A}^1$  consisting of all singletons; here  $\mathcal{A}^0 = \{\emptyset\}$ . For this example, our theory reduces to the well-known considerations concerning atoms. This example could be iterated to higher cardinalities, but we use two other instances instead: the system of linear subspaces and the system of affine subspaces of  $\mathbb{R}^p$  (or  $\bar{\mathbb{R}}^p$ ). In both of them,  $\mathcal{A}^i$  consist of subspaces with dimension equal to  $i$ ; for linear subspaces  $\mathcal{A}^0 = \{0\}$ , while for affine ones  $\mathcal{A}^0 = \{\emptyset\}$ .

If  $Q$  is a measure and  $E$  is a (measurable) set, we denote by  $Q \restriction E$  the restriction (“trace”) of  $Q$  to  $E$ : the measure satisfying  $(Q \restriction E)(A) = Q(E \cap A)$ . We say that  $Q$  is supported by  $A$  if  $Q \restriction A = Q$ ; that is,  $Q(E) = Q(E \cap A)$  for every  $E$ .

**PROPOSITION 5.4.** *Suppose that  $\mathcal{A}$  is an atomic system on  $\mathbb{X}$ . Any finite measure  $Q$  can be written as the sum*

$$(15) \quad Q = \sum_{i=0}^r \sum_{A \in \mathcal{A}_Q^i} Q^A + Q^\omega,$$

where sets  $\mathcal{A}_Q^i$  are at most countable (they may be empty), every measure  $Q^A$  is supported by  $A$  and the measure  $Q^j$  given by

$$(16) \quad Q^j = \sum_{i=0}^j \sum_{A \in \mathcal{A}_Q^i} Q^A$$

satisfies that  $Q(A) - Q^j(A) = 0$  for any  $A \in \mathcal{A}^j$ .

**PROOF.** In the proof, we use the following corollary of the Hahn-Saks-Vitali theorem—see Doob (1993), Theorems III.10 and IX.10, pages 30 and 155: if  $Q_\nu$  is a countable system of measures such that the sum  $\sum_\nu Q_\nu(E)$  converges for all (measurable)  $E$ , then this sum defines a measure  $Q(E)$ . We write  $Q_1 \leq Q_2$ , if  $Q_1(E) \leq Q_2(E)$  for all  $E$ .

We define the decomposition (15) inductively. In the initial step, take the only element  $A$  of  $\mathcal{A}^0$ ; if  $Q(A) > 0$ , then set  $\mathcal{A}_Q^0 = \{A\}$  and  $Q^A = Q \restriction A$ . Otherwise (for instance, if  $A = \emptyset$ ), set  $\mathcal{A}_Q^0 = \emptyset$ . In each case, we have that  $Q^0(E) \leq Q(E)$  for any  $E$ .

Suppose now that for all  $i < j$ , we have already defined the sets  $\mathcal{A}_Q^i$ , as well as measures  $Q^A$  for every  $A \in \mathcal{A}_Q^i$ ; the sets  $\mathcal{A}_Q^i$  are at most countable, the measure  $Q^{j-1}$ , given by (16), satisfies the inequality  $Q^{j-1} \leq Q$ , and finally  $Q(E) - Q^{j-1}(E) = 0$  for any  $A \in \mathcal{A}^i$  and any  $i < j$ . Since  $Q^{j-1} \leq Q$ ,  $Q - Q^{j-1}$  is a measure; denote it by  $\tilde{Q}$ . Let  $\mathcal{A}_Q^j$  be the (possibly empty) collection of all  $A \in \mathcal{A}^j$  such that  $\tilde{Q}(A) > 0$ . For every  $A \in \mathcal{A}_Q^j$ , let  $Q^A = (\tilde{Q} \restriction A)$ . The intersection of any two distinct elements  $\mathcal{A}^j$  belongs to  $\mathcal{A}^i$  with  $i < j$ ; therefore this intersection has measure  $\tilde{Q}$  equal to 0. It follows that  $\mathcal{A}_Q^j$  is at most countable and for any  $E$ ,

$$(17) \quad \sum_{A \in \mathcal{A}_Q^j} Q^A(E) = \sum_{A \in \mathcal{A}_Q^j} \tilde{Q}(E \cap A) = \tilde{Q}\left(E \cap \bigcup_{A \in \mathcal{A}_Q^j} A\right) \leq \tilde{Q}(E) = Q(E) - Q^{j-1}(E);$$

hence  $Q^j$ , defined by (16), satisfies the inequality  $Q^j \leq Q$ . The way how we constructed all measures  $Q^A$  so far implies that

$$\tilde{Q} - \sum_{A \in \mathcal{A}_Q^j} Q^A = Q - Q^{j-1} - \sum_{A \in \mathcal{A}_Q^j} Q^A = Q - Q^j.$$

If  $i < j$  and  $E$  belongs to  $\mathcal{A}^i$ , then  $\tilde{Q}(E)$  is equal to 0 as well as all  $Q^A(E)$ , due to the inequality  $Q^A \leq Q$  holding for any  $A \in \mathcal{A}_Q^j$ . The same holds if  $E$  is from  $\mathcal{A}_Q^j \setminus \mathcal{A}^j$ . Finally, if  $E$  is from  $\mathcal{A}_Q^j$ , then

$$(18) \quad \tilde{Q}(E) = Q^E(E) \leq \sum_{A \in \mathcal{A}_Q^j} Q^A(E).$$

This yields, together with (17), that  $Q(E) - Q^j(E) = 0$ , concluding the induction step.

After repeating the construction for  $j = 1, 2, \dots, r$ , we finally set  $Q^\omega = Q - Q^r$ .  $\square$

We call the system of measures appearing in (15) an *atomic decomposition* of a (finite) measure  $Q$  with respect to an atomic system  $\mathcal{A}$ . We denote by  $\mathcal{A}_Q$  the union of sets  $\mathcal{A}_Q^i$  for  $i = 0, 1, 2, \dots, r$ —the set of *atoms*. Finally, the measure  $Q^\omega$  is a *nonatomic part of  $Q$  with respect to  $\mathcal{A}$* , since  $Q^\omega(A) = 0$  for any  $A \in \mathcal{A}$ .

**PROPOSITION 5.5.** *Suppose that  $\mathcal{A}$  is an atomic system on  $\mathbb{X}$  and  $Q$  is a finite measure. If  $E \in \mathcal{A}$  and  $Q(E) > 0$ , then  $E$  contains a member from  $\mathcal{A}_Q$ .*

**PROOF.** Since  $Q^\omega(E) = 0$ , Proposition 5.4 implies that there is  $\tilde{A} \in \mathcal{A}_Q$  such that  $Q^{\tilde{A}}(E) > 0$ ; we may pick it from  $\mathcal{A}_Q^j$  with smallest possible  $j$ , so that for all  $A \in \mathcal{A}_Q^i$  with  $i < j$ ,  $Q^A(E) = 0$ . Now,  $0 < Q^{\tilde{A}}(E) = Q^{\tilde{A}}(\tilde{A} \cap E)$ ; thus  $\tilde{A} \cap E$  must belong to  $\mathcal{A}_Q^i$  with  $i \geq j$ . Since  $\tilde{A} \cap E \subseteq \tilde{A}$ , it follows that  $\tilde{A} \cap E = \tilde{A}$ ; that is,  $\tilde{A} \subseteq E$ .  $\square$

**5.3. Cosmic spaces.** Hereafter, the letter  $\nu$  is reserved for passages to infinity along positive integers: all limits involving  $\nu$  are with respect to  $\nu \rightarrow \infty$ .

A **cosmic space**  $\bar{\mathbb{R}}^p$  is a compactification of  $\mathbb{R}^p$  (a compact topological space containing  $\mathbb{R}^p$  as an open dense subset), which is formed by adding a point with infinite distance from the origin in each direction. The name is due to Rockafellar and Wets (1998); we refer there for more details and background. Rigorously spoken,  $\bar{\mathbb{R}}^p$  is the union of  $\mathbb{R}^p$  and a homeomorphic copy  $\partial\mathbb{R}^p$  of  $\mathbb{S}^{p-1}$ : it is convenient to write elements of  $\partial\mathbb{R}^p$  as  $\omega(u)$ , where  $u \in \mathbb{S}^{p-1}$  and  $\omega$  is a one-to-one mapping of  $\mathbb{S}^{p-1}$  to  $\partial\mathbb{R}^p$ .

The space  $\bar{\mathbb{R}}^p$  can be conveniently described in polar coordinates. The *norm*  $\|x\|$  extends the Euclidean norm on  $\mathbb{R}^p$ ; it is equal to  $\infty$  for  $x \in \partial\mathbb{R}^p$ . We define the *direction* of  $x$  to be

$$((x)) = \begin{cases} 0 & \text{if } x = 0, \\ x/\|x\| & \text{if } x \neq 0 \text{ and } x \in \mathbb{R}^p, \\ u & \text{if } x \in \partial\mathbb{R}^p, x = \omega(u) \text{ for } u \in \mathbb{S}^{p-1}. \end{cases}$$

The direction defined in this way satisfies the following equalities for all  $x, y \in \mathbb{R}^p$ :

$$(19) \quad \text{sgn}(x^\top y) = \text{sgn}((x)^\top y) = \text{sgn}(x^\top (y)) = \text{sgn}((x)^\top (y)).$$

In the vein of (19), we define for any  $x, y \in \bar{\mathbb{R}}^p$ ,

$$(20) \quad \text{sgn}(x \cdot y) = \text{sgn}((x)^\top (y)).$$

In what follows, we actually never need the full definition of the inner product, just its sign; nevertheless, for better readability we often write  $x \cdot y \geq 0$  instead of  $\text{sgn}(x \cdot y) \geq 0$ .

Norm and direction characterize the convergence in  $\bar{\mathbb{R}}^p$ :  $x_\nu \rightarrow x$  if and only if  $\|x_\nu\| \rightarrow \|x\|$ , and either  $x = 0$  or  $((x_\nu)) \rightarrow ((x))$  (note that the latter may not hold if  $x = 0$ ). It is not hard to see that the convergence defined in this way extends the standard one in  $\mathbb{R}^p$  and is consistent with the topological properties of  $\bar{\mathbb{R}}^p$ . Particularly,  $\partial\mathbb{R}^p$  is homeomorphic to  $\mathbb{S}^{p-1}$ ;  $\mathbb{R}^p$  is open and dense in  $\bar{\mathbb{R}}^p$ ; and  $\bar{\mathbb{R}}^p$  is compact and homeomorphic to the ball  $\mathbb{D}^p$ .

For later convenience, we extend the mapping  $\omega$  from  $\mathbb{S}^{p-1}$  to  $\bar{\mathbb{R}}^p$ . We set  $\omega(0) = 0$ ; for  $x \neq 0$ , let  $\omega(x) = \omega((x))$ , where the right side refers to the original definition of  $\omega$ . Albeit the extended  $\omega$  is no longer one-one, it satisfies the following identity for all  $x \in \bar{\mathbb{R}}^p$ :

$$(21) \quad ((\omega(x))) = ((x)).$$

We write  $x = -y$  if and only if  $((x)) = -((y))$ . For all  $x, y \in \bar{\mathbb{R}}^p$ ,

$$(22) \quad \text{sgn}((-x) \cdot y) = \text{sgn}(x \cdot (-y)) = -\text{sgn}(x \cdot y).$$

**PROPOSITION 5.6.** *Let  $x_\nu, y_\nu \in \mathbb{R}^p$ . If  $x_\nu \rightarrow x \in \partial\mathbb{R}^p$  and  $\|y_\nu\| < \infty$ , then  $x_\nu + y_\nu \rightarrow x$ .*

**PROOF.** A straightforward verification: just note that  $\|x_\nu + y_\nu\| \geq \|x_\nu\| - \|y_\nu\| \rightarrow \infty$  and then that

$$((x_\nu + y_\nu)) = \frac{x_\nu + y_\nu}{\|x_\nu + y_\nu\|} = \frac{\frac{x_\nu}{\|x_\nu\|} + \frac{y_\nu}{\|x_\nu\|}}{\left\| \frac{x_\nu}{\|x_\nu\|} + \frac{y_\nu}{\|x_\nu\|} \right\|} \rightarrow \frac{((x))}{\|((x))\|} = ((x));$$

these two facts together give the desired convergence.  $\square$

Obviously, any measure on  $\mathbb{R}^p$  may be extended to  $\bar{\mathbb{R}}^p$  and, since  $\mathbb{R}^p \subset \bar{\mathbb{R}}^p$ , any theorem about measures on  $\bar{\mathbb{R}}^p$  holds for measures on  $\mathbb{R}^p$ .

**5.4. General properties of depth.** In this subsection,  $\mathbb{X}$  may be  $\mathbb{R}^p$ ,  $\mathbb{D}^p$ ,  $\mathbb{S}^{p-1}$ , or  $\bar{\mathbb{R}}^p$ ; we do all the proofs for  $\mathbb{X} = \bar{\mathbb{R}}^p$ , the other cases being entirely similar. If  $Q$  is a measure on  $\mathbb{X}$ , we write  $S(Q)$  for the set of all  $u \in \mathbb{S}^{p-1}$  such that  $Q(\partial H_u \setminus \{0\}) = 0$ , and  $h_Q$  for the function from  $\mathbb{S}^{p-1}$  to  $\mathbb{R}$  such that  $h_Q(u) = Q(H_u)$ .

**PROOF OF THEOREM 3.2.** Let, for  $j = 0, 1, 2, \dots, p-1$ , the sets  $\mathcal{A}^j$  consist of all subsets of  $\mathbb{X}$  of the form

$$\{x \in \bar{\mathbb{R}}^p : x^\top u_1 = 0, x^\top u_2 = 0, \dots, x^\top u_{p-i} = 0\} \setminus \{0\},$$

where  $u_1, u_2, \dots, u_{p-i}$  are linearly independent vectors from  $\mathbb{S}^{p-1}$ , for  $i \leq j$ . The union  $\mathcal{A}$  of all  $\mathcal{A}^j$  is an atomic system consisting of all linear subspaces of  $\bar{\mathbb{R}}^p$  with codimension at least 1 with the origin  $\{0\}$  removed. We apply Proposition 5.5 to the measure  $Q \llcorner (\bar{\mathbb{R}}^p \setminus \{0\})$ : any  $E \in \mathcal{A}$  contains  $A$  from a countable system  $\mathcal{A}_Q$ .

For any  $A \in \mathcal{A}$ , let  $A^\perp$  denote the set of all  $u \in \mathbb{S}^{p-1}$  such that  $u \cdot x = 0$  for all  $x$  in  $A$ . Any such  $A^\perp$  has codimension at least 1, hence its dimension is less than the full dimension of  $\mathbb{S}^{p-1}$ —and therefore its measure  $\mu$  is 0, as well as the measure  $\mu$  of any countable union of such sets, in particular of the union  $\mathcal{A}_Q^\perp$  of all  $A^\perp$  such that  $A \in \mathcal{A}_Q$ . If  $A \subseteq E$ , then  $A^\perp \supseteq E^\perp$ ; thus by Proposition 5.5,  $\mathcal{A}_Q^\perp$  contains all  $u$  such that  $Q(\{x \in \bar{\mathbb{R}}^p : x \cdot u = 0\} \setminus \{0\}) > 0$ . Since  $\mu(\mathcal{A}_Q^\perp) = 0$ , we obtain that  $\mu(S(Q)) = 1$ ; therefore  $S(Q)$  is dense in  $\mathbb{S}^{p-1}$ .  $\square$

Weak convergence for subprobability measures is understood in the sense of Billingsley (1968, 1971): we write  $Q_\nu \rightarrow Q$ , if the integrals of all bounded continuous functions with respect to  $Q_\nu$  converge to those with respect to  $Q$ . Weak convergence implies that the total mass of  $Q_\nu$  converges to  $Q$ ; our sequences are always formed by measures with the same total mass, hence it is not hard to see that all relevant theorems hold (trivially if the total mass is zero, and by division of the total mass otherwise—switching to “conditional probabilities”).

**PROPOSITION 5.7.** *Suppose that  $Q_\nu, Q$  are measures on  $\bar{\mathbb{R}}^p$  and  $Q_\nu \rightarrow Q$ . If  $u_\nu, u \in \mathbb{S}^{p-1}$  and  $u_\nu \rightarrow u$ , then*

$$(23) \quad \limsup_{\nu \rightarrow \infty} Q_\nu(H_{u_\nu}) \leq Q(H_u).$$

*If  $Q(\partial H_u) = 0$ , then  $Q_\nu(H_{u_\nu}) \rightarrow Q(H_u)$ .*

**PROOF.** Note first the following. Let  $u_\nu, u \in \mathbb{S}^{p-1}$ ,  $u_\nu \rightarrow u$ , and  $x_\nu, x \in \bar{\mathbb{R}}^p$ ,  $x_\nu \rightarrow x$ . If  $u_\nu^T((x_\nu)) \geq 0$  for infinitely many  $\nu$ , then  $u^T((x)) \geq 0$ ; this holds trivially if  $x = 0$  and follows from the fact that  $((x_\nu)) \rightarrow ((x))$  otherwise. Conversely, if  $u^T((x)) > 0$ , then  $u_\nu^T((x_\nu)) > 0$  for all but a finite number of  $\nu$ , since in this case  $x \neq 0$  and thus  $((x_\nu)) \rightarrow ((x))$ .

We start the proof by invoking the Skorokhod representation: it yields the existence of random vectors  $\xi_\nu, \xi$ , defined on the common probability space  $(\Omega, \mathcal{S}, \mathbb{P})$ , such that  $\mathcal{L}(\xi_\nu) = Q_\nu$ ,  $\mathcal{L}(\xi) = Q$ , and  $\xi_\nu \rightarrow \xi$  almost surely. With the help of (20), we obtain

$$(24) \quad \begin{aligned} \limsup_{\nu \rightarrow \infty} Q_\nu(H_{u_\nu}) &= \limsup_{\nu \rightarrow \infty} \mathbb{P}[u_\nu \cdot \xi_\nu \geq 0] = \limsup_{\nu \rightarrow \infty} \mathbb{P}[u_\nu^T((\xi_\nu)) \geq 0] \\ &\leq \mathbb{P}[u^T((\xi)) \geq 0] = \mathbb{P}[u \cdot \xi \geq 0] = Q(H_u); \end{aligned}$$

this proves (23). If  $Q(\partial H_u) = 0$ , then

$$(25) \quad \begin{aligned} Q(H_u) &= \mathbb{P}[u \cdot \xi > 0] = \mathbb{P}[u^T((\xi)) > 0] \leq \liminf_{\nu \rightarrow \infty} \mathbb{P}[u_\nu^T((\xi_\nu)) > 0] \\ &\leq \liminf_{\nu \rightarrow \infty} \mathbb{P}[u_\nu^T((\xi_\nu)) \geq 0] = \liminf_{\nu \rightarrow \infty} Q_\nu(H_{u_\nu}). \end{aligned}$$

The convergence follows from the combination of (24) and (25).  $\square$

**PROPOSITION 5.8.** *Let  $\Phi_\vartheta(z)$  be a function from  $\Theta \times Z$  to  $\mathbb{R}^p$ ; let  $Q_\nu, Q$  be measures on  $Z$ .*

*(i) If, for  $\vartheta_\nu \rightarrow \vartheta$ ,*

$$(26) \quad ((\Phi_{\vartheta_\nu}(z_\nu))) \rightarrow ((\Phi_\vartheta(z))) \quad \text{whenever } z_\nu \rightarrow z$$

for almost all  $z$  with respect to  $Q$  (in particular, when  $\Phi$  is jointly continuous in  $\vartheta$  and  $z$ ), then  $Q_\nu \rightarrow Q$  and  $\vartheta_\nu \rightarrow \vartheta$  imply that

$$(27) \quad \limsup_{\nu \rightarrow \infty} Q_\nu \circ \Phi_{\vartheta_\nu}^{-1}(\mathbf{H}_u) \leq Q \circ \Phi_\vartheta^{-1}(\mathbf{H}_u).$$

If (27) holds only with  $z_\nu = z$  for all  $\nu$  and  $Q$ -almost all  $z$  (that holds, for instance, if  $\Phi$  is continuous only in  $\vartheta$ ), then (27) still holds with  $Q_\nu = Q$  for all  $\nu$ .

(ii) If  $\Phi_\vartheta(z)$ , as a function of  $\vartheta$ , is continuous at  $\vartheta$  for  $Q$ -almost all  $z$ , then  $Q \circ \Phi_{\vartheta_\nu}^{-1} \rightarrow Q \circ \Phi_\vartheta^{-1}$ .

PROOF. The proof is completely analogous to that of Proposition 5.7; the second part is an application of the continuous mapping theorem from the theory of weak convergence.  $\square$

PROPOSITION 5.9. For any measure  $Q$  on  $\bar{\mathbb{R}}^p$ , a function  $h_Q(u) = Q(\mathbf{H}_u)$  is upper semicontinuous on  $\mathbb{S}^{p-1}$ , and continuous at every  $u \in S(Q)$ , the set of all  $u \in \mathbb{S}^{p-1}$  such that  $Q(\partial \mathbf{H}_u \setminus \{0\}) = 0$ .

PROOF. Upper semicontinuity follows from Proposition 5.7, setting  $Q_\nu = Q$  for all  $\nu$ . To prove continuity, we apply the same proposition to measures defined by  $Q_\nu(E) = Q(E) = Q(E \setminus \{0\})$  for all  $E$  and  $\nu$ .  $\square$

PROOF OF PROPOSITION 3.1. The proposition directly follows from Proposition 5.9. In a more detailed way, if  $h_Q(u) < d(Q) + \varepsilon$ , then there is a sequence  $u_\nu \in S$  such that  $u_\nu \rightarrow u$ ; upper semicontinuity yields that

$$\inf_{u \in S} h_Q(u) \leq \limsup_{\nu \rightarrow \infty} h_Q(u_\nu) \leq h_Q(u) \leq d(Q) + \varepsilon$$

—and since  $\varepsilon$  was arbitrary, (ii) follows (the converse inequality is obvious).  $\square$

PROPOSITION 5.10. For any measure  $Q$  on  $\bar{\mathbb{R}}^p$ ,  $Q(\{0\}) \leq d(Q) \leq (Q(\{0\}) + Q(\bar{\mathbb{R}}^p))/2$ .

PROOF. The first inequality is obvious. The second one follows from Proposition 3.1 via Theorem 3.2, which implies that the set of those  $u$  for which  $Q(\mathbf{H}_u) + Q(\mathbf{H}_{-u}) = 1 + Q(\{0\})$  is dense in  $\mathbb{S}^{p-1}$ .  $\square$

PROPOSITION 5.11. Let  $Q_\nu, Q$  be measures on  $\bar{\mathbb{R}}^p$  such that  $Q_\nu \rightarrow Q$ . If  $Q(\partial \mathbf{H}_u) = 0$  for all  $u \in \mathbb{S}^{p-1}$ , then  $Q_\nu(\mathbf{H}_u) \rightarrow Q(\mathbf{H}_u)$  uniformly in  $u \in \mathbb{S}^{p-1}$ .

PROOF. We topologize the space of all closed proper halfspaces in a natural way:  $\mathbf{H}_{u_\nu} \rightarrow \mathbf{H}_u$  if and only if  $u_\nu \rightarrow u$ . Clearly, this topology makes the space of all halfspaces homeomorphic to  $\mathbb{S}^{p-1}$ , hence compact and metrizable. If  $\mathbf{H}_{u_\nu} \rightarrow \mathbf{H}_u$  and  $Q_\nu \rightarrow Q$ , then  $Q_\nu(\mathbf{H}_{u_\nu}) \rightarrow Q(\mathbf{H}_u)$ , by Proposition 5.7. Finally, we have also that  $Q(\mathbf{H}_{u_\nu}) \rightarrow Q(\mathbf{H}_u)$ , again by Proposition 5.7. We verified all assumptions of Proposition 2.2 from page 6 of Bickel and Millar (1992); the proposition follows.  $\square$

PROPOSITION 5.12. For any measures  $Q_\nu, Q$  on  $\bar{\mathbb{R}}^p$ , the following implications hold:

- (i) if  $\limsup_{\nu \rightarrow \infty} Q_\nu(\mathbf{H}_u) \leq Q(\mathbf{H}_u)$  for any  $u \in \mathbb{S}^{p-1}$ , then  $\limsup_{\nu \rightarrow \infty} d(Q_\nu) \leq d(Q)$ ;
- (ii) if  $Q_\nu \rightarrow Q$  and  $Q_\nu(\mathbf{H}_u) \rightarrow Q(\mathbf{H}_u)$  uniformly in  $u \in \mathbb{S}^{p-1}$ , then  $d(Q_\nu) \rightarrow d(Q)$ .



PROOF. In (i), the assumption asserts that for fixed  $u$ ,  $Q(H_u)$  as a function of  $Q$  is upper semicontinuous (with respect to the weak topology). Depth function is thus the infimum of upper semicontinuous functions and therefore also upper semicontinuous. The assumption in (ii) says that the functions  $h_{Q_\nu}(u) = Q_\nu(H_u)$  converge uniformly to the function  $h_Q(u) = Q(H_u)$ ; the convergence of infima follows.  $\square$

**5.5. Measures concentrated in closed halfspaces.** Hereafter, we denote by  $G_u$  the open halfspace that is the interior of  $H_u$ . Taking complements with respect to  $\mathbb{X}$  (the correct choice of  $\mathbb{X}$  being hopefully clear from the context), we have that  $G_u = H_{-u}^c$ .

We say that  $w$  lies on the *arc* connecting  $u$  and  $v$  from  $\mathbb{S}^{p-1}$  if one of the following three possibilities holds: either  $w = u = v$ ; or  $u = -v$  (and  $w$  is arbitrary); or  $w$  lies between  $u$  and  $v$  on the shorter part of the circumference created by the intersection of  $\mathbb{S}^{p-1}$  and the linear subspace generated by  $u$  and  $v$ .

PROPOSITION 5.13. *Suppose that  $Q$  is a measure on  $\bar{\mathbb{R}}^p$  such that  $Q(G_u) = 0$  for some  $u \in \mathbb{S}^{p-1}$ . For every  $v_1, v_2 \in \mathbb{S}^{p-1}$ , if  $v_1 \neq u$  lies on the arc connecting  $u$  and  $v_2$ , then  $Q(H_{v_1}) \leq Q(H_{v_2})$ .*

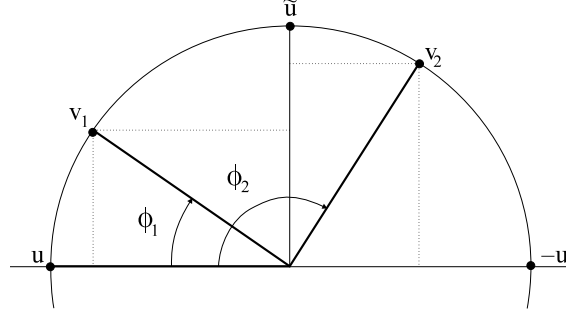


FIGURE 3. An illustration for the proof of Proposition 5.13.

PROOF. Note first that the assumptions imply that  $v_2 \neq u$ ; if  $v_2 = -u$ , then  $Q(H_{v_2}) = Q(\bar{\mathbb{R}}^p)$  and the proposition holds trivially—hence we will assume that  $v_2 \neq -u$  as well. In such a case,  $u$  and  $v_2$  are not collinear and

$$\tilde{u} = \frac{v_2 - u(u^T v_2)}{\|v_2 - u(u^T v_2)\|}$$

is well defined. It is straightforward to verify that (see Figure 3):  $u^T \tilde{u} = 0$ , hence  $u$  is orthogonal to  $\tilde{u}$  and lies in  $H_{\tilde{u}}$ ;  $\tilde{u}$  lies in the linear space spanned by  $u$  and  $v_2$  and therefore on the same circumference as  $u$ ,  $v_2$ , and thus also  $v_1$ ; finally,  $v_2^T \tilde{u} \geq 0$ , hence  $v_2$  lies in the halfspace  $H_{\tilde{u}}$ . Since  $u_1$  lies on the arc connecting  $u$  and  $v_2$ , it also lies in  $H_{\tilde{u}}$ . We express  $v_1$  and  $v_2$  in the polar coordinates: for  $i = 1, 2$ ,

$$v_i = u(u^T v_i) + \tilde{u}(\tilde{u}^T v_i) = u \cos \varphi_i + \tilde{u} \sin \varphi_i,$$

where the angles  $\varphi_i$  are from  $(0, \pi)$ . The condition that  $v_1$  lies on the arc connecting  $u$  and  $v_2$  means that  $\varphi_1 < \varphi_2$  and hence  $\cot \varphi_1 > \cot \varphi_2$ .

Now it is sufficient to show that  $H_{v_1} \setminus H_{v_2} \subseteq H_{-u}^c$ , since then  $Q(H_{v_1} \setminus H_{v_2}) \leq Q(H_{-u}^c) = 0$  and by subadditivity,

$$Q(H_{v_1}) \leq Q(H_{v_1} \setminus H_{v_2}) + Q(H_{v_2}).$$

Suppose that  $x \in H_{v_1} \setminus H_{v_2}$ . Then

$$(28) \quad ((x))^T v_1 = ((x))^T u \cos \varphi_1 + ((x))^T \tilde{u} \sin \varphi_1 \geq 0,$$

$$(29) \quad ((x))^T v_2 = ((x))^T u \cos \varphi_2 + ((x))^T \tilde{u} \sin \varphi_2 < 0.$$

Dividing (28) and (29) by  $\sin \varphi_1$  and  $\sin \varphi_2$ , respectively (we excluded the possibility that any one of them is zero at the beginning), and combining both inequalities, we obtain

$$((x))^T u (\cot \varphi_1 - \cot \varphi_2) > 0.$$

Hence  $((x))^T u > 0$ , that is,  $x \in H_{-u}^c$ .  $\square$

**PROPOSITION 5.14.** *If  $Q$  is a measure on  $\bar{\mathbb{R}}^p$  such that  $Q(G_u) = 0$  for some  $u \in \mathbb{S}^{p-1}$ , then  $d(Q) = d(Q \llcorner \partial H_u)$ .*

**PROOF.** We have to prove that  $d(Q) \leq d(Q \llcorner \partial H_u)$ ; the converse inequality is trivial. Fix  $\varepsilon > 0$ . Consider  $v \in \partial H_u \cap \mathbb{S}^{p-1}$ . For any  $w \in \mathbb{S}^{p-1}$ ,

$$Q(H_w) = Q(H_w \cap G_u) + Q(H_w \cap \partial H_u) + Q(H_w \cap G_{-u}) = Q(H_w \cap \partial H_u) + Q(H_w \cap G_{-u}),$$

since  $Q(G_u) = 0$ . If  $w$  approaches  $u$  along the arc connecting  $u$  and  $v$ , then  $Q(H_w \cap G_{-u}) \rightarrow 0$ . Therefore, there is  $w \in \mathbb{S}^{p-1}$  such that

$$Q(H_w) \leq Q(H_w \cap \partial H_u) + \varepsilon = Q(H_v \cap \partial H_u) + \varepsilon = (Q \llcorner \partial H_u)(H_v) + \varepsilon$$

It follows that  $d(Q) \leq d(Q \llcorner \partial H_u) + \varepsilon$ . Since  $\varepsilon$  was arbitrary, the desired inequality follows.  $\square$

**5.6. Skeletons, centrality, topology.** Throughout this subsection,  $Q$  is a measure on  $\bar{\mathbb{R}}^p$  or  $\mathbb{S}^{p-1}$ ,  $h_Q(u) = Q(H_u)$  is a function on  $\mathbb{S}^{p-1}$ ,  $\lambda$  a nonnegative constant, and  $\mu$  denotes the uniform distribution on  $\mathbb{S}^{p-1}$ . For any  $\lambda$ , we define a  $\lambda$ -skeleton of  $Q$  to be

$$s_\lambda(Q) = \int_{\mathbb{S}^{p-1}} u (h_Q(u) \wedge \lambda) d\mu(u).$$

Clearly, any skeleton assigns to  $Q$  a vector in a unit ball  $\mathbb{D}^p$ . Skeletons give vector-field approximations to probability fields; the inspiration for the notion came from Birch (1959).

**PROPOSITION 5.15.** *If  $\lambda = \lambda(Q) \geq 0$  depends continuously on  $Q$ , with respect to the weak topology, and  $Q(\{0\}) = 0$ , then  $s_{\lambda(\cdot)}(\cdot)$  is continuous at  $Q$ .*

**PROOF.** Suppose that  $Q_\nu \rightharpoonup Q$  and  $Q(\{0\}) = 0$ . By Proposition 5.7 and Theorem 3.2,  $h_{Q_\nu}(u) \rightarrow h_Q(u)$  for  $\mu$ -almost all  $u$ . If  $\lambda(Q_\nu) \rightarrow \lambda(Q)$ , then

$$u (h_{Q_\nu}(u) \wedge \lambda(Q_\nu)) \rightarrow u (h_Q(u) \wedge \lambda(Q))$$

for  $\mu$ -almost all  $u$ . The proof is concluded by the application of the Lebesgue dominated convergence theorem.  $\square$

We call  $Q$   $\lambda$ -**central**, if the level set  $\{u \in \mathbb{S}^{p-1} : h_Q(u) \leq \lambda\}$  surrounds 0. Note that  $\lambda$ -centrality of  $Q$  means also its  $\tilde{\lambda}$ -centrality whenever  $\lambda \leq \tilde{\lambda}$ . The following proposition expresses the basic Helly-Carathéodory principle.

**PROPOSITION 5.16.** *If  $\lambda > d(Q)$  and  $s_{\lambda}(Q) = 0$ , then  $Q$  is  $\lambda$ -central and  $\lambda \geq 1/(p+1)$ .*

**PROOF.** Since  $\lambda > d(Q)$ , Proposition 3.1 and Theorem 3.2 imply the existence of  $u \in S(Q)$  such that  $h_Q(u) < \lambda$ . By Proposition 5.9, the level set  $E_{\lambda} = \{u \in \mathbb{S}^{p-1} : h_Q(u) \leq \lambda\}$  contains some neighborhood of  $u$ ; hence  $\mu(E_{\lambda}) > 0$ . Since  $s_{\lambda}(Q) = 0$ , the set  $E_{\lambda}$  cannot be contained in any open halfspace  $\mathbf{H}_u^c$ ; hence  $E_{\lambda}$  surrounds 0, by Proposition 2.2(ii), and  $Q$  is  $\lambda$ -central.

Proposition 2.2(iii) then yields the existence of a set  $V \subseteq E_{\lambda}$  surrounding 0 with cardinality at most  $m+1$ . If a set  $V$  surrounds 0, then the system  $\{\mathbf{H}_v\}_{v \in V}$  covers  $\bar{\mathbb{R}}^p$  and

$$1 \leq \sum_{v \in V} P_{\partial}(\mathbf{H}_v) \leq (p+1)\lambda. \quad \square$$

The following proposition specifies how skeletons behave on the compactifying boundary.

**PROPOSITION 5.17.** *Let  $Q$  be a probability measure on  $\bar{\mathbb{R}}^p$  such that  $Q(\mathbf{H}_{-u}) = 1$  for some  $u \in \mathbb{S}^{p-1}$ . For any  $\lambda \geq 0$ ,  $s_{\lambda}(Q) \in \mathbf{H}_{-u}$ .*

**PROOF.** Given  $u \in \mathbb{S}^{p-1}$ , we write  $v^*$  for the reflection of  $v \in \mathbb{S}^{p-1}$  about  $\partial\mathbf{H}_u$ :  $v^* = v - 2u(u^T v)$ . This reflection places points from the open hemispheres  $\mathbf{G}_u \cap \mathbb{S}^{p-1}$  and  $\mathbf{G}_{-u} \cap \mathbb{S}^{p-1}$  into one-one correspondence.

Suppose that  $v \in \mathbf{G}_u \cap \mathbb{S}^{p-1}$ ,  $v \neq u$ . Let  $\tilde{v} = v - u(u^T v)$ ; note that  $\tilde{v}$  is orthogonal to  $u$  and thus lies in  $\partial\mathbf{H}_u$ . We will decompose both  $v$  and  $v^*$  to a sum of vectors collinear with  $u$  and  $\tilde{v}$ . First,  $u^T v^* = -u^T v$ , hence the corresponding components in the direction of  $u$  are of the same magnitude. Second, since

$$\tilde{v}^T v^* = \tilde{v}^T v = 1 - (u^T v)^2 \geq 0,$$

$v$  and  $v^*$  lie in  $\mathbf{H}_{\tilde{v}}$ , and  $v$  lies on the arc connecting  $u$  and  $v^*$ : all the three are collinear and if  $v \in \mathbf{G}_u$ , that is,  $u^T v > 0$ , then  $u^T v^* < 0$ . Proposition 5.13 implies that  $Q(\mathbf{H}_v) \leq Q(\mathbf{H}_{v^*})$ ; therefore,

$$(30) \quad h_Q(v) \wedge \lambda \leq h_Q(v^*) \wedge \lambda$$

for all  $\lambda \geq 0$ . Computing the  $\lambda$ -skeleton, we obtain

$$(31) \quad \begin{aligned} s_{\lambda}(Q) &= \int_{\mathbb{S}^{p-1}} v(h_Q(v) \wedge \lambda) d\mu(v) \\ &= \int_{\mathbf{G}_u \cup \mathbf{G}_{-u}} v(h_Q(v) \wedge \lambda) d\mu(v) + \int_{\partial\mathbf{H}_u} v(h_Q(v) \wedge \lambda) d\mu(v) \end{aligned}$$

The second integral results in a vector lying in  $\partial H_u$ . The first integral in (31) can be further orthogonally decomposed to

$$(32) \quad \begin{aligned} & u \int_{G_u} (u^T v) ((h_Q(v) \wedge \lambda) - (h_Q(v^*) \wedge \lambda)) d\mu(v) \\ & + \int_{G_u} \tilde{v} \frac{\tilde{v}^T v}{\|\tilde{v}\|^2} ((h_Q(v) \wedge \lambda) + (h_Q(v^*) \wedge \lambda)) d\mu(v). \end{aligned}$$

The second integral in (32) results again in  $\partial H_u$ . The first one in  $H_{-u}$ , due to (30); since  $\partial H_u \subseteq H_{-u}$ , we obtain by (31) and (32) that  $s_\lambda(Q) \in H_{-u}$ .  $\square$

A point where the vector field is equal to 0 is its critical point. Critical points of continuous vector fields are effectively hunted by the Kronecker index theory. Let  $S, T$  be topological spaces. Continuous mappings  $f, g$  from  $S$  to  $T$  are called *homotopic*, if one can be continuously transformed to another: there is a continuous mapping  $\pi$  from  $S \times [0, 1]$  to  $T$  such that  $\pi(x, 0) = f(x)$  and  $\pi(x, 1) = g(x)$ . A continuous mapping from  $\mathbb{S}^{p-1}$  to  $\mathbb{S}^{p-1}$  is homotopic to a constant mapping if and only if it can be extended to a continuous mapping of  $\mathbb{D}^p$  to  $\mathbb{S}^{p-1}$ . Homotopy is an equivalence relation: if  $f$  is homotopic to a constant and  $f$  and  $g$  are homotopic, then  $g$  is homotopic to a constant. A topological space  $T$  is called *contractible* if the identity mapping from  $T$  to  $T$  is homotopic to a constant. More thorough picture of homotopic equivalence of continuous mappings of spheres gives the theory of degree: a mapping is homotopic to a constant if and only if its topological degree is 0. On the sphere, neither the identity mapping  $f(x) = x$  (degree 1) nor the antipodal mapping  $f(x) = -x$  (degree  $-1$ ) are homotopic to a constant. If a restriction to the boundary  $\partial U$  of a region  $U$  of a vector field is nonsingular, its degree is called an *index* of a vector field over  $U$ . The following proposition is a reformulation of the well-known principle from the theory of vector fields, saying that a continuous vector field on  $\partial U$  with a nonzero index possesses a critical point inside  $U$ .

**PROPOSITION 5.18.** *Let  $\xi$  be a continuous vector field on  $\mathbb{R}^p$ . If  $\mathbb{R}^p$  contains a subset  $S$  homeomorphic to  $\mathbb{S}^{p-1}$  and such that the restriction of  $\xi$  to  $S$  is not homotopic to a constant, then  $\xi$  has a critical point in the region bounded by  $S$  (possibly lying on its boundary).*

**PROOF.** See Dodson and Parker (1997), page 25. Instead of  $\mathbb{R}^p$ , we might take any contractible topological space  $\Theta$ , due to the pathwise connectedness of  $\mathbb{S}^{p-1}$ ; see, for instance, Exercise 19(ii), page 26 of Rotman (1988).  $\square$

The following theorem is the core of our general approach to centerpoint hunting, in the vein of Birch (1959) and Donoho and Gasko (1992). The general theorem given below covers the simplest case when parameter space is  $\Theta = \mathbb{R}^p$  (Examples 1, 2, and 4 showed the practical importance of this case). More sophisticated parametric spaces can be analyzed similarly, only the adequate topological engine would be different.

**THEOREM 5.19.** *If a probability field  $\{P_\vartheta: \vartheta \in \bar{\mathbb{R}}^p\}$  on  $\bar{\mathbb{R}}^p$  is*

- (i) *continuous in the weak topology at all  $\vartheta \in \bar{\mathbb{R}}^p$ :  $P_{\vartheta_\nu} \rightharpoonup P_\vartheta$  whenever  $\vartheta_\nu \rightarrow \vartheta$ ,*
- (ii)  *$P_\vartheta(H_{-\vartheta}) = 0$  for all  $\vartheta \in \partial \bar{\mathbb{R}}^p$ ,*

(iii)  $P_\vartheta(\{0\}) = 0$  for all  $\vartheta \in \mathbb{R}^p$ ,

then there exists a point  $\dot{\vartheta} \in \mathbb{R}^p$  such that  $d(P_{\dot{\vartheta}}) \geq 1/(p+1)$ .

PROOF. For fixed  $\lambda > 0$ , the continuity of the probability field and the fact that  $P_\vartheta(\{0\}) = 0$  for all  $\vartheta$  imply, via Proposition 5.15, the continuity of the  $\lambda$ -skeleton  $s_\lambda(P_\vartheta)$  as a vector field dependent on  $\vartheta$ . Assumption (ii) implies, via Proposition 5.17, that  $s_\lambda(P_\vartheta) \in H_\vartheta$  for any  $\vartheta \in \partial\mathbb{R}^p$ . The Poincaré-Bohl theorem—see Dodson and Parker (1997), page 19—says that if  $f, g$  are two continuous functions from  $\mathbb{S}^{p-1}$  to  $\mathbb{S}^{p-1}$  not pointing to opposite directions at any point, then  $f$  and  $g$  are homotopic. This holds in our case; the skeleton  $s_\lambda(P_\vartheta)$  never points in the direction opposite to hence  $\vartheta$ , hence its degree on  $\partial\mathbb{R}^p$  is equal to the degree of the identity mapping on  $\mathbb{S}^{p-1}$ . The latter is equal to 1; therefore the  $\lambda$ -skeleton cannot be homotopic to a constant. By Proposition 5.18, such a  $\lambda$ -skeleton has a critical point in  $\mathbb{R}^p$ . By Proposition 5.16, if  $\lambda > \sup\{d(P_\vartheta) : \vartheta \in \mathbb{R}^p\}$ , then  $\lambda$  cannot be less than  $1/(p+1)$ . Thus,  $\sup\{d(P_\vartheta) : \vartheta \in \mathbb{R}^p\} \geq 1/(p+1)$ . Proposition 5.12 and compactness of  $\mathbb{R}^p$  yield a deepest parameter, a point  $\dot{\vartheta} \in \mathbb{R}^p$  such that  $d(P_{\dot{\vartheta}}) = \sup\{d(P_\vartheta) : \vartheta \in \mathbb{R}^p\}$ . Apparently,  $d(P_{\dot{\vartheta}}) \geq 1/(p+1)$ ; finally,  $\dot{\vartheta}$  cannot lie in  $\partial\mathbb{R}^p$ , since from (ii) follows that  $d(P_\vartheta) = 0$  for any  $\vartheta \in \partial\mathbb{R}^p$ .  $\square$

**5.7. Applications: regression probability fields.** Two principal techniques to overcome technical difficulties in application of Theorem 5.19 are compactification and smoothing. Recall that  $\bar{\mathbb{X}}$  is a compactification of  $\mathbb{X}$  if it is compact and  $\mathbb{X}$  is dense in  $\bar{\mathbb{X}}$ ; we denote  $\bar{\mathbb{X}} \setminus \mathbb{X}$  by  $\partial\bar{\mathbb{X}}$ . We say that a probability field  $\{\bar{P}_\vartheta : \vartheta \in \bar{\Theta}\}$  on  $\bar{\mathbb{X}}$  is a **compactification** of a probability field  $\{P_\vartheta : \vartheta \in \Theta\}$  on  $\mathbb{X}$ , if  $\bar{\mathbb{X}}$  is a compactification of  $\mathbb{X}$ , the indexing set  $\bar{\Theta}$  is a compactification of  $\Theta$ , and  $\bar{P}_\vartheta = P_\vartheta$  for all  $\vartheta \in \Theta$ , in the sense that  $\bar{P}_\vartheta \llcorner \mathbb{X} = P_\vartheta$  and  $\bar{P}_\vartheta(\partial\bar{\mathbb{X}}) = 0$ .

For models with  $\Theta = \mathbb{R}^p$ , the useful compactifications are those with  $\bar{\Theta} = \bar{\mathbb{R}}^p$ —for these, Theorem 5.19 yields the existence of a centerpoint—in the original  $\Theta = \mathbb{R}^p$  (so that we may eventually forget about the compactification trick). The appropriate compactified probability field can be found by a natural process of continuous extension. As far as the assumptions of Theorem 5.19 are concerned, note first that (ii) implies (iii) also for all  $x \in \partial\mathbb{R}^p$ ; for verifying (i), the following scheme may be helpful.

PROPOSITION 5.20. *Suppose that  $\Theta$  is an open, dense subset in  $\bar{\Theta}$ . The probability field  $\{\bar{P}_\vartheta : \vartheta \in \bar{\Theta}\}$  is continuous with respect to weak topology, if  $\vartheta_\nu \rightarrow \vartheta$  implies weak convergence of  $\bar{P}_{\vartheta_\nu}$  to  $\bar{P}_\vartheta$  in the following three cases:*

- (a)  $\vartheta_\nu \in \Theta, \vartheta \in \Theta$ ;
- (b)  $\vartheta_\nu \in \Theta, \vartheta \in \bar{\Theta} \setminus \Theta$ ;
- (c)  $\vartheta_\nu \in \bar{\Theta} \setminus \Theta, \vartheta \in \bar{\Theta} \setminus \Theta$ .

PROOF. A straightforward manipulation with subsequences.  $\square$

To check (a) in Proposition 5.20, we verify the continuity of the original probability field—setting  $\bar{P}_\vartheta = P_\vartheta$  for  $\vartheta \in \Theta$  and extending the measures  $\bar{P}_\vartheta$ , if needed. If we define  $\bar{P}_\vartheta$  for  $\vartheta \in \partial\mathbb{R}^p$  by continuous extension, then (b) comes automatically. Finally, checking (c) is often easy.

Another problem in Theorem 5.19 is its assumption (iii), usually not holding when the underlying distribution has atoms. A technique to overcome this difficulty is smoothing: we add

a small perturbation with absolutely continuous distribution. However, it may be not that easy to backtrack the smoothing approximation successfully; unless we are able to control the depth on the compactifying boundary, for instance.

Let us illustrate all this technology on the simplest example.

**EXAMPLE 1** (Centerpoints). Let  $P$  be a probability on  $\mathbb{R}^p$ ; let  $Z$  be a random variable whose distribution is  $P$ . The function  $\Psi_\vartheta(z) = \vartheta - z$  is continuous, jointly in  $\vartheta$  and  $z$ , hence the gradient probability field  $P_\vartheta = \mathcal{L}(\vartheta - Z)$  is continuous at any  $\vartheta \in \mathbb{R}^p$ , by Proposition 5.8—we checked assumption (a) of Proposition 5.20. If  $\vartheta \in \partial\mathbb{R}^p$ , then Proposition 5.6 yields that  $\vartheta_\nu - z$  converges to  $\vartheta$ ; we extend the definition of  $\Phi$  setting  $\Phi_\vartheta(z) = \vartheta$  for  $\vartheta \in \partial\mathbb{R}^p$ . By Proposition 5.6, this extension preserve the continuity of  $\Phi$ , jointly in  $\vartheta$  and  $z$ ; therefore, assumptions (b) and (c) are checked as well and Proposition 5.20 yields assumption (i) of Theorem 5.19.

Note that for  $\vartheta \in \partial\mathbb{R}^p$ , the probability  $P_\vartheta$  is concentrated in  $\vartheta$ , therefore  $P_\vartheta(\mathbf{H}_{-\vartheta}) = 0$ ; this checks assumption (ii) of Theorem 5.19. At this point, (ii) holds without any continuity assumption on  $Z$ ; this implies, in particular, that  $d(P_\vartheta) = 0$  for any  $\vartheta \in \partial\mathbb{R}^p$ .

The only remaining assumption of Theorem 5.19 is (iii). If it holds—for instance, when the distribution of  $Z$  is absolutely continuous—then we are able to prove the existence of a centerpoint in  $\mathbb{R}^p$ . However, assumption (iii) does not hold when the distribution of  $Z$  has atoms with nonzero probability. To overcome this difficulty, we consider a sequence of compactified probability fields arising from random variables  $Z + Z_\nu$ , with the following properties:  $Z_\nu \rightarrow 0$  and hence the distributions  $P_\nu$  of  $Z + Z_\nu$  converge weakly to  $P$ ; the distributions of all  $Z_\nu$ , and hence of  $Z + Z_\nu$ , are absolutely continuous. All the perturbed probability fields satisfy (i), (ii), and (iii); Theorem 5.19 yields a sequence of centerpoints  $\dot{\vartheta}_\nu \in \mathbb{R}^p$  such that  $d_T(\dot{\vartheta}_\nu, Z + Z_\nu) \geq 1/(p+1)$ .

Such a sequence has a subsequence converging to  $\dot{\vartheta} \in \bar{\mathbb{R}}^p$ . The continuity of  $\Phi$  implies, via Proposition 5.8, that  $\limsup_{\nu \rightarrow \infty} P_\nu \circ \Phi_{\dot{\vartheta}_\nu}^{-1}(\mathbf{H}_u) \leq P \circ \Phi_{\dot{\vartheta}}^{-1}(\mathbf{H}_u)$  (we abuse the subsequence notation in the usual way). Proposition 5.12 then yields that  $d_T(\dot{\vartheta}, P) \geq 1/(p+1)$ . Finally,  $\dot{\vartheta}$  cannot lie in  $\partial\mathbb{R}^p$  since all  $\vartheta \in \partial\mathbb{R}^p$  have  $d_T(\vartheta, Z) = 0$  (recall that this was shown without the help of the absolute continuity assumption).

In regression models, we follow essentially the same scheme. The required technical facts are summarized in the following proposition.

**PROPOSITION 5.21.** *Let  $Z = (X, Y)$  be a random variable with values in  $\mathbf{X} \times \mathbb{R}^m$  where  $\mathbf{X} \subseteq \mathbb{R}^k$ . If  $\{P_\Theta: \Theta \in \mathbb{R}^p\}$  is a probability field such that  $P_\Theta = \mathcal{L}(XX^\top\Theta - XY^\top)$  for any  $\Theta \in \mathbb{R}^p = \mathbb{R}^{km}$ , then it has a compactification, a probability field  $\{P_\Theta: \Theta \in \bar{\mathbb{R}}^p\}$  on  $\bar{\mathbb{R}}^p$ , with the following properties:*

- (i) *the probability field  $\{P_\Theta: \Theta \in \bar{\mathbb{R}}^p\}$  is continuous (in weak topology) at every  $\Theta \in \mathbb{R}^p$  and at every  $\Theta \in \partial\mathbb{R}^p$  such that  $\mathbb{P}[X \cdot \Theta = 0]$ .*
- (ii)  *$\limsup_{\nu \rightarrow \infty} P_{\Theta_\nu} \leq P_\Theta$ , whenever  $\Theta_\nu \in \mathbb{R}^p$  and  $\Theta_\nu \rightarrow \Theta \in \bar{\mathbb{R}}^p$ ;*
- (iii) *for  $\Theta \in \partial\mathbb{R}^p$ ,  $P(\mathbf{H}_\Theta) = 1$  and  $d(P_\Theta) = P_\Theta(\{0\}) = P_\Theta(\mathbf{H}_{-\Theta}) = \mathbb{P}[X \cdot \Theta = 0]$ .*

*If, moreover,  $\Delta(X) = 0$ , then*

- (iv)  *$\{P_\Theta: \Theta \in \bar{\mathbb{R}}^p\}$  is continuous (in weak topology) at all  $\Theta \in \bar{\mathbb{R}}^p$ ;*

(v) for every  $\Theta \in \partial\mathbb{R}^p$ ,  $d(P_\Theta) = P(\mathbf{H}_{-\Theta}) = 0$ .

The values of the compactification at  $\vartheta \in \partial\mathbb{R}^p$  are given by  $P_\Theta = \mathcal{L}(\omega(XX^T((\Theta))))$ .

Before proving Proposition 5.21, we show how it can be applied in the proof of Theorem 3.4. Note that Proposition 5.21 and Theorem 5.19 imply Conjecture 1(b) of Rousseeuw and Hubert (1999a) immediately. If the distribution of  $Z$  is absolutely continuous (with respect to the appropriate Lebesgue measure), then  $\Delta(X) = 0$ , assumption (iii) of Theorem 5.19 holds, and Proposition 5.21(iv) and (v) implies assumptions (i) and (ii), respectively.

**PROOF OF THEOREM 3.4.** If  $\Delta(X) > 0$ , then the compactified probability field may not be continuous with respect to weak topology at points from  $\partial\mathbb{R}^p$ . We consider a sequence of compactified probability fields generated by random variables  $Z + Z_\nu = (X + X_\nu, Y + Y_\nu)$  where  $Z_\nu \rightarrow 0$ . All  $Z_\nu$ , and thus all  $Z + Z_\nu$ , have absolutely continuous distributions. Hence,  $\Delta(X + X_\nu) = 0$  and all probability fields generated by perturbed random variables satisfy all assumptions of Theorem 5.19. We obtain a sequence of centerpoints  $\dot{\vartheta}_\nu \in \mathbb{R}^p$ ,  $d_T(\dot{\vartheta}_\nu, Z + Z_\nu) \geq 1/(p+1)$ ; take its subsequence with a limit  $\dot{\vartheta} \in \bar{\mathbb{R}}^p$ . From Proposition 5.21(ii) follows, via Proposition 5.12(i), that  $d_T(\dot{\vartheta}, Z) \geq 1/(p+1)$ . The final step is to show that  $\dot{\vartheta}$  cannot lie in  $\partial\mathbb{R}^p$ . By Proposition 5.21(iii),  $d_T(\vartheta, Z) = \mathbb{P}[X \cdot \vartheta = 0] \leq \Delta(X)$  for  $\vartheta \in \partial\mathbb{R}^p$ . Therefore, the desired result follows—if  $\Delta(X) < 1/(p+1)$ .

By Proposition 5.21(ii), 5.21(i), and 5.12, tangent depth is upper semicontinuous on  $\bar{\mathbb{R}}^p$ , which is compact; this proves the existence of the deepest parameter  $\dot{\vartheta}$  in  $\bar{\mathbb{R}}^p$ . Proposition 5.21(iii) then implies that  $\dot{\vartheta}$  cannot lie in  $\partial\mathbb{R}^p$ .  $\square$

**PROOF OF PROPOSITION 5.21.** The continuity of the probability field at all  $\Theta \in \mathbb{R}^p$  follows, via Proposition 5.8(ii), from the joint continuity (in  $x$ ,  $y$ , and  $\Theta$ ) of the function  $\Phi_\Theta(x, y) = xx^T\Theta - xy^T$ . In the same way, (ii) follows from Proposition 5.8(ii) for  $\Theta \in \mathbb{R}^p$ .

Suppose now that  $\Theta_\nu \rightarrow \Theta$ ,  $\Theta_\nu \in \mathbb{R}^p$ ,  $\Theta \in \partial\mathbb{R}^p$ . Since  $\|\Theta_\nu\| \rightarrow \infty$ , we may assume that  $\Theta_\nu \neq 0$ ; in view of the fact that also  $((\Theta_\nu)) \rightarrow ((\Theta))$ , we obtain that

$$(33) \quad ((x_\nu x_\nu^T \Theta_\nu - x_\nu y_\nu^T)) = ((x_\nu x_\nu^T ((\Theta_\nu)) - x_\nu y_\nu^T \|\Theta_\nu\|^{-1})) \rightarrow ((xx^T((\Theta))))$$

This suggests the form of the compactified probability field:  $P_\Theta = \mathcal{L}(\omega(XX^T((\Theta))))$ . Under this definition, Proposition 5.8(i), in view of (33), completes the proof of (ii). To finish also part (i), we need first to show that

$$(34) \quad xx^T \Theta_\nu - xy^T \rightarrow \omega(xx^T((\Theta)))$$

for almost all  $(x, y)$  with respect to  $P = \mathcal{L}(X, Y)$ . Suppose that  $\mathbb{P}[X^T((\Theta)) = 0] = 0$ ; then  $X \neq 0$  almost surely, and therefore  $XX^T((\Theta)) \neq 0$  almost surely. Hence, for  $P$ -almost all  $x$  (and  $y$ ),  $\|xx^T((\Theta_\nu))\| \rightarrow \|xx^T((\Theta))\| > 0$  and therefore

$$(35) \quad \|xx^T \Theta_\nu\| = \|\Theta_\nu\| \|xx^T((\Theta_\nu))\| \rightarrow \infty$$

and

$$(36) \quad ((xx^T \Theta_\nu)) = \frac{xx^T((\Theta_\nu))}{\|xx^T((\Theta_\nu))\|} \rightarrow \frac{xx^T((\Theta))}{\|xx^T((\Theta))\|} = ((xx^T((\Theta)))) = ((\omega(xx^T((\Theta))))).$$

By Proposition 5.6, (35) and (36) imply (34). Finally, if  $\Theta_\nu, \Theta \in \partial\mathbb{R}^p$  and  $\Theta_\nu \rightarrow \Theta$ , then

$$(37) \quad ((\omega(xx^T((\Theta_\nu)))) = ((xx^T((\Theta_\nu)))) \rightarrow ((xx^T((\Theta)))) = ((\omega(xx^T((\Theta))))),$$

finishing the proof of (iii) completely.

Let  $\Theta \in \partial\mathbb{R}^p$ . By (21), (20) and (19),

$$(38) \quad \begin{aligned} P_\Theta(\mathbf{H}_\Theta) &= \mathbb{P}[\Theta \cdot \omega(XX^T((\Theta))) \geq 0] = \mathbb{P}[(\Theta) \cdot (XX^T((\Theta))) \geq 0] \\ &= \mathbb{P}[(\Theta) \cdot (XX^T((\Theta))) \geq 0] = \mathbb{P}[\text{tr}((\Theta))^T XX^T((\Theta)) \geq 0] \\ &= \mathbb{P}[X^T((\Theta))(\Theta)^T X \geq 0] = \mathbb{P}[\|(\Theta)^T X\|^2 \geq 0] = 1, \end{aligned}$$

Similarly, for all  $U \in \mathbb{S}^{p-1}$ ,

$$(39) \quad \begin{aligned} P_\Theta(\mathbf{H}_U) &= \mathbb{P}[U \cdot \omega(XX^T((\Theta))) \geq 0] = \mathbb{P}[\text{tr}(U^T XX^T((\Theta))) \geq 0] \\ &= \mathbb{P}[X^T((\Theta))U^T X \geq 0] \geq \mathbb{P}[X^T((\Theta)) = 0] = \mathbb{P}[X \cdot \Theta = 0] \end{aligned}$$

—that is,  $X^T \vartheta_j = 0$  for  $j = 1, 2, \dots, m$ , where  $\vartheta_j$  is the  $j$ -th column of  $\Theta$ . By (22),

$$(40) \quad \begin{aligned} P_\Theta(\mathbf{H}_{-\Theta}) &= \mathbb{P}[-(\Theta) \cdot (XX^T((\Theta))) \geq 0] \\ &= \mathbb{P}[-\text{tr}((\Theta))^T XX^T((\Theta)) \geq 0] = \mathbb{P}[-X^T((\Theta))(\Theta)^T X \geq 0] \\ &= \mathbb{P}[\|(\Theta)^T X\|^2 \leq 0] = \mathbb{P}[(\Theta)^T X = 0] = \mathbb{P}[X \cdot \Theta = 0]. \end{aligned}$$

Since (39) and (40) hold for any  $\Theta \in \partial\mathbb{R}^p$ , we obtain the proof of (iii):

$$(41) \quad \begin{aligned} d(P_\Theta) &= P_\Theta(\mathbf{H}_{-\Theta}) = \mathbb{P}[(\Theta)^T X = 0] = \mathbb{P}[(\Theta)^T XX^T = 0] \\ &= \mathbb{P}[(\Theta)^T XX^T = 0] = \mathbb{P}[XX^T((\Theta)) = 0] = P_\Theta(\{0\}). \end{aligned}$$

The rest of the proposition is easy: just note that if  $\Delta(X) = 0$ , then  $\mathbb{P}[X \cdot \Theta = 0] = 0$  for all  $\Theta \in \bar{\mathbb{R}}^p$ .  $\square$

**5.8. Bias and breakdown.** In this subsection, the function  $\Phi$  involved in the definition of tangent depth is supposed to be an arbitrary function from  $\Theta \times \mathbf{Z}$  to  $\mathbb{R}^p$ , measurable as a function from  $\mathbf{Z}$  to  $\mathbb{R}^p$  for any  $\vartheta \in \Theta$  and we suppress the dependence on it in the notation. As in Section 4,  $v$  stands for the variation metric and  $\gamma$  for the contamination distance.

**PROPOSITION 5.22.** *Suppose that  $\pi(P, \tilde{P}) \leq \varepsilon$ . Then*

- (i)  $|d_T(\vartheta, P) - d_T(\vartheta, \tilde{P})| \leq \varepsilon$  if  $\pi = v$ ;
- (ii)  $(1 - \varepsilon)d_T(\vartheta, P) \leq d_T(\vartheta, \tilde{P}) \leq d_T(\vartheta, P) + \varepsilon$  if  $\pi = \gamma$ .

**PROOF.** The assumption  $v(P, \tilde{P}) \leq \varepsilon$  implies that

$$\left| P(\Phi_\vartheta^{-1}(\mathbf{H}_u)) - \tilde{P}(\Phi_\vartheta^{-1}(\mathbf{H}_u)) \right| \leq \varepsilon$$

for any  $u \neq 0$  and any  $\vartheta \in \Theta$ . This implies (i). The second inequality in (ii) follows from (i) and the inequality  $v(P, \tilde{P}) \leq \gamma(P, \tilde{P}) \leq \varepsilon$ . To see the first inequality in (ii), just note that

$$(1 - \varepsilon)P(\Phi_\vartheta^{-1}(\mathbf{H}_u)) \leq \tilde{P}(\Phi_\vartheta^{-1}(\mathbf{H}_u))$$

for any  $u \neq 0$  and any  $\vartheta \in \Theta$ .  $\square$



PROOF OF THEOREM 4.1. To see (i), note first that  $\tilde{\vartheta}$  can be from  $\mathcal{T}(\tilde{P})$  only if  $d_T(\tilde{\vartheta}, \tilde{P}) \geq \eta$ ; but then, by Proposition 5.22,  $d_T(\tilde{\vartheta}, P) \geq \eta - \varepsilon$ , proving (i). Let  $\dot{\vartheta}$  be from  $\mathcal{T}(P)$ . Any  $\tilde{\vartheta}$  can be from  $\mathcal{T}(\tilde{P})$  only if  $d_T(\dot{\vartheta}, \tilde{P}) \leq d_T(\tilde{\vartheta}, \tilde{P})$ . We obtain that

$$(42) \quad \eta \leq d_T(\dot{\vartheta}, P) \leq d_T(\dot{\vartheta}, \tilde{P}) + \varepsilon \leq d_T(\tilde{\vartheta}, \tilde{P}) + \varepsilon \leq d_T(\tilde{\vartheta}, P) + 2\varepsilon,$$

by the repeated application of Proposition 5.22(i). This proves (ii). Finally, the proof of (iii) is analogous, only now Proposition 5.22(ii) gives

$$(43) \quad \eta(1 - \varepsilon) \leq (1 - \varepsilon)d_T(\dot{\vartheta}, P) \leq d_T(\dot{\vartheta}, \tilde{P}) \leq d_T(\tilde{\vartheta}, \tilde{P}) \leq d_T(\tilde{\vartheta}, P) + \varepsilon$$

instead of (42).  $\square$

Suppose now that  $\Theta = \mathbb{R}^p$  (all reasonings are valid for any locally compact space).

PROPOSITION 5.23. *If  $\vartheta_n \in \Theta$ ,  $\vartheta_n \rightarrow \infty$ , then  $\limsup_{n \rightarrow \infty} d_T(\vartheta_n, P) \leq d_T(\infty, P)$ .*

PROOF. A straightforward consequence of the definition of  $d_T(\infty, P)$ .  $\square$

PROOF OF THEOREM 4.2. Fix  $\varepsilon > \varepsilon_\pi^*(\mathcal{T}, P)$ . Then  $\mathcal{T}(\mathcal{B}_\pi(P, \varepsilon))$  is not bounded and hence there is a sequence  $\tilde{\vartheta}_n \in \mathcal{T}(\tilde{P}_n)$  such that  $d_T(\tilde{\vartheta}_n, \tilde{P}_n) \geq \eta$ . By Propositions 5.22 and 5.23,  $\varepsilon \geq \eta - d_T(\infty, P)$ . Since  $\varepsilon$  was arbitrary, (i) follows.

The proof of (ii) starts in the same way and we arrive to the observation similar to (42), with  $\tilde{\vartheta}$  replaced by  $\tilde{\vartheta}_n$ . Proposition 5.23 then yields

$$\eta \leq \limsup_{n \rightarrow \infty} d_T(\tilde{\vartheta}_n, P) + 2\varepsilon \leq d_T(\infty, P) + 2\varepsilon$$

which proves (ii). Finally, (iii) comes also in an analogous way, only (43) is used instead of (42), giving

$$(1 - \varepsilon)\eta \leq \limsup_{n \rightarrow \infty} d_T(\tilde{\vartheta}_n, P) + \varepsilon \leq d_T(\infty, P) + \varepsilon,$$

implying (iii).  $\square$

**5.9. A refined analysis of the regression model.** In this subsection,  $Z = (X, Y)$  is a random variable with values in  $\mathbf{X} \times \mathbb{R}$  and  $\{P_\vartheta: \vartheta \in \mathbb{R}^p\}$  is a (linear, not multivariate) regression probability field:  $P_\vartheta = \mathcal{L}(XX^T\vartheta - XY^T)$  for  $\vartheta \in \mathbb{R}^p$ .

The geometric definition of the regression depth shows that it is independent of the parametrization (this principle holds in greater generality, but we will not develop this theme here). The following notion comes from the theory of nonlinear regression. Given a regression function  $f(x, \vartheta)$ , the *solution locus* of  $\vartheta \in \Theta$  is defined to be

$$\ell(\vartheta) = \{(y, x): y = f(x, \vartheta), x \in \mathbf{X}\}.$$

In linear regression,  $\ell(\vartheta) = \{(y, x): y = x^T\vartheta, x \in \mathbf{X}\}$ . A closer investigation of graphs of fitted surfaces reveals how this definition can be extended to  $\vartheta$  from  $\partial\mathbb{R}^p$ : exploring the limits of sets  $\ell(\vartheta)$  in terms of set convergence, we arrive to the definition

$$\ell(\vartheta) = \{(y, x): x \bullet \vartheta = 0, x \in \mathbf{X}\}.$$

This indicates a statistical interpretation for the parameter points from the cosmic extension: they correspond to “vertical” regression fits. Since  $x \bullet \vartheta = 0$  if and only if  $x \bullet (-\vartheta) = 0$ , we have that  $\ell(\vartheta) = \ell(-\vartheta)$ ; that is, each vertical fit may be interpreted as occurring once in upwards and once in downwards orientation. Therefore, a projective plane might be more appropriate compactification than the cosmic one from the statistical point of view; however, the latter has more convenient topological properties. Note that by Proposition 5.21(iii),  $d(P_{-\vartheta}) = \mathbb{P}[X \bullet (-\vartheta) = 0] = \mathbb{P}[X \bullet \vartheta = 0] = d(P_{\vartheta})$  for any  $\vartheta \in \partial\mathbb{R}^p$ , hence the cosmic compactification is consistent from the depth point of view.

For the proof Theorem 5.25, the following technical result is required.

PROPOSITION 5.24. *If  $\vartheta_\nu \rightarrow \vartheta \in \partial\mathbb{R}^p$ , then  $P_{\vartheta_\nu}(\mathbf{G}_{-\vartheta}) \rightarrow 0$ .*

PROOF. Just note that, for any  $\vartheta \in \mathbb{R}^p$  and any  $u \in \bar{\mathbb{R}}^p$ ,

$$(44) \quad \begin{aligned} \mathbb{P}[u \bullet (XX^T \vartheta - XY) = 0] &= \mathbb{P}[-((u))^T (XX^T \vartheta - XY) = 0] \\ &\geq \mathbb{P}[(u)^T X = 0] = \mathbb{P}[X \bullet u = 0] \end{aligned}$$

and therefore

$$\begin{aligned} &\limsup_{\nu \rightarrow \infty} \mathbb{P}[(-\vartheta) \bullet (XX^T \vartheta_\nu - XY) > 0] \\ &= \limsup_{\nu \rightarrow \infty} (\mathbb{P}[(-\vartheta) \bullet (XX^T \vartheta_\nu - XY) \geq 0] - \mathbb{P}[(-\vartheta) \bullet (XX^T \vartheta_\nu - XY) = 0]) \\ &\leq \limsup_{\nu \rightarrow \infty} \mathbb{P}[(-\vartheta) \bullet (XX^T \vartheta_\nu - XY) \geq 0] - \mathbb{P}[X \bullet (-\vartheta) = 0] \\ &\leq \mathbb{P}[X \bullet \vartheta = 0] - \mathbb{P}[X \bullet \vartheta = 0] = 0. \quad \square \end{aligned}$$

Assume the same setting as in the definition of tangent depth. For any  $E \subseteq \mathbf{Z}$ , we define

$$d_T^\Phi(\vartheta, P \llcorner E) = d((P \llcorner E) \circ \Phi_\vartheta^{-1}) = \inf_{u \neq 0} P(\Phi_\vartheta^{-1}(\mathbf{H}_u) \cap E) = d(P_\vartheta \llcorner \Phi(E, \vartheta)),$$

where  $\Phi(E, \vartheta) = \{\Phi(z, \vartheta) : z \in E\}$  and  $P_\vartheta = P \circ \Phi_\vartheta^{-1}$ . If  $P$  is the distribution of a random variable  $Z$ , then we write  $d_T^\Phi(\vartheta, Z \llcorner E)$  for  $d_T^\Phi(\vartheta, \mathcal{L}(Z) \llcorner E)$ ; in such a case,

$$d_T^\Phi(\vartheta, Z \llcorner E) = \inf_{u \neq 0} \mathbb{P}[Z \in \Phi_\vartheta^{-1}(\mathbf{H}_u) \cap E] = \inf_{u \neq 0} \mathbb{P}[\Phi_\vartheta(Z) \in \mathbf{H}_u \text{ and } Z \in E].$$

THEOREM 5.25. *Suppose that  $Z = (X, Y)$  is a random variable with values in  $\mathbf{X} \times \mathbb{R} \subseteq \mathbb{R}^{p+1}$ , and  $\Phi_\vartheta(z) = \Phi_\vartheta(y, x) = xx^T \vartheta - xy$  for any  $\vartheta \in \Theta = \mathbb{R}^p$ . If  $\vartheta_\nu$  is a sequence from  $\mathbb{R}^p$  such that  $\vartheta_\nu \rightarrow \vartheta \in \partial\mathbb{R}^p$  for  $\nu \rightarrow \infty$ , then*

$$(45) \quad \limsup_{\nu \rightarrow \infty} d_T^\Phi(\vartheta_\nu, Z) \leq \sup_{\eta \in \mathbb{R}^p} d_T^\Phi(\eta, Z \llcorner \ell(\vartheta)).$$

PROOF. We prove first that

$$(46) \quad \limsup_{\nu \rightarrow \infty} d(P_{\vartheta_\nu}) \leq \limsup_{\nu \rightarrow \infty} d(P_{\vartheta_\nu} \llcorner \partial\mathbf{H}_\vartheta)$$

(which actually means the equality since the converse is obvious) and then the equality

$$(47) \quad \limsup_{\nu \rightarrow \infty} d(P_{\vartheta_\nu} \llcorner \partial\mathbf{H}_\vartheta) = \limsup_{\nu \rightarrow \infty} d_T(\vartheta_\nu, Z \llcorner \ell(\vartheta)).$$

Fix  $\varepsilon > 0$ . By Proposition 5.24, there is  $\nu_0$  such that  $P_{\vartheta_\nu}(\mathbf{G}_{-\vartheta}) < \varepsilon$  for  $\nu > \nu_0$ ; thus by Proposition 5.14,

$$d(P_{\vartheta_\nu}) \leq d(P_{\vartheta_\nu} \sqcup \mathbf{H}_\vartheta) + \varepsilon = d(P_{\vartheta_\nu} \sqcup \partial \mathbf{H}_\vartheta) + \varepsilon.$$

Since  $\varepsilon$  is arbitrary, (46) follows. The verification of (47) is rather straightforward. Note first that

$$d(P_{\vartheta_\nu} \sqcup \partial \mathbf{H}_\vartheta) = \inf_{u \neq 0} \mathbb{P}[u^T (XX^T \vartheta_\nu - XY) \geq 0 \text{ and } ((\vartheta))^\top (XX^T \vartheta_\nu - XY) = 0]$$

and

$$d_T(\vartheta_\nu, Z \sqcup \ell(\vartheta)) = \inf_{u \neq 0} \mathbb{P}[u^T (XX^T \vartheta_\nu - XY) \geq 0 \text{ and } ((\vartheta))^\top X = 0].$$

While in the general case only the obvious inequality holds given by (44) holds, the equality (47) holds when  $\vartheta_\nu \rightarrow \vartheta \in \partial \mathbb{R}^p$ . Let  $\zeta_\nu = XX^T \vartheta_\nu - XY$ . We obtain that

$$\begin{aligned} \limsup_{\nu \rightarrow \infty} d(P_{\vartheta_\nu} \sqcup \partial \mathbf{H}_\vartheta) &= \limsup_{\nu \rightarrow \infty} \inf_{u \neq 0} \mathbb{P}[u^T \zeta_\nu \geq 0 \text{ and } ((\vartheta))^\top (XX^T \vartheta_\nu - XY) = 0] \\ &= \limsup_{\nu \rightarrow \infty} \inf_{u \neq 0} \left( \mathbb{P}[u^T \zeta_\nu \geq 0 \text{ and } ((\vartheta))^\top X = 0] \right. \\ &\quad \left. + \mathbb{P}[u^T \zeta_\nu \geq 0 \text{ and } ((\vartheta))^\top X \neq 0 \text{ and } X^T \vartheta_\nu - Y = 0] \right) \\ &\leq \limsup_{\nu \rightarrow \infty} \left( \inf_{u \neq 0} \mathbb{P}[u^T \zeta_\nu \geq 0 \text{ and } ((\vartheta))^\top X = 0] + \mathbb{P}[((\vartheta))^\top X \neq 0 \text{ and } X^T \vartheta_\nu = Y] \right) \\ &\leq \limsup_{\nu \rightarrow \infty} d_T(\vartheta_\nu, Z \sqcup \ell(\vartheta)) + \limsup_{\nu \rightarrow \infty} \mathbb{P}[((\vartheta))^\top X \neq 0 \text{ and } X^T((\vartheta_\nu)) = Y \| \vartheta_\nu \|^{-1}] \\ &\leq \limsup_{\nu \rightarrow \infty} d_T(\vartheta_\nu, Z \sqcup \ell(\vartheta)) + \mathbb{P}[((\vartheta))^\top X \neq 0 \text{ and } X^T((\vartheta)) = 0], \end{aligned}$$

yielding the desired result, since the last probability is equal to zero.  $\square$

In the linear regression with  $\Theta = \mathbb{R}^p$ , Theorem 5.25 gives that

$$(48) \quad d_T(\infty, Z) \leq \sup_{\vartheta \in \mathbb{R}^p} \sup_{\eta \in \mathbb{R}^p} d_T^\Phi(\eta, Z \sqcup \ell(\vartheta)).$$

As already mentioned, the bound  $d_T(\infty, Z) \leq \Delta(X)$  follows from Propositions 5.21(ii) and 5.12(i). If “model holds”, that is, there is a parameter  $\vartheta_0$  such that the conditional distribution of  $Z$  given  $X$  is symmetric about  $\ell(\vartheta_0)$ , then the depth of  $\vartheta_0$  is  $\delta_0 + (1 - \delta_0)/2$  with  $\delta_0 = \mathbb{P}[X^T \vartheta_0 = Y]$ ; by Theorem 5.25,

$$d_T(\infty, Z) \leq \Delta_0 + \frac{1}{2}(\Delta(X) - \Delta_0),$$

where  $\Delta_0 = \sup_{\vartheta \neq 0} \mathbb{P}[X^T \vartheta = Y \text{ and } X^T \vartheta = 0] \leq \delta_0$ .

#### APPENDIX: PROOF OF THEOREM 3.3 IN THE FULL GENERALITY

**Additional facts from measure theory.** We write  $\text{cl } E$  for the topological closure of a set  $E$ ; all other notation introduced throughout the paper is kept. In what follows,  $Q_\nu$ ,  $Q_\nu^i$ ,  $Q^i$  and  $Q$  are all (bounded) measures on a separable metric space  $\mathbb{X}$ .

LEMMA 6.1. *Suppose that  $Q_\nu \rightarrow Q$ . If  $E_1, E_2$  are closed sets such that  $Q_\nu(E_1 \cup E_2) = Q(E_1 \cup E_2)$  and  $Q_\nu(E_1 \cap E_2) = Q(E_1 \cap E_2) = 0$  for all  $\nu$ , then  $Q_\nu(E_i) \rightarrow Q(E_i)$  for  $i = 1, 2$ .*

PROOF. Since  $E_i$  are closed, we have for  $i = 1, 2$ ,

$$(49) \quad \limsup_{\nu \rightarrow \infty} Q_\nu(E_i) \leq Q(E_i).$$

Using the fact that  $Q_\nu(E_2) = Q_\nu(E_1 \cup E_2) - Q_\nu(E_1)$  for all  $\nu$  and also for  $Q$  in place of  $Q_\nu$ , we obtain, in view of (49) for  $E_2$ , that

$$Q(E_1) \leq \liminf_{\nu \rightarrow \infty} Q_\nu(E_1)$$

—together with (49) for  $E_1$  this yields the convergence for  $E_1$ . The proof for  $E_2$  is symmetric.  $\square$

LEMMA 6.2. *Suppose that  $Q_\nu = \sum_{i=1}^{\infty} Q_\nu^i$  and  $Q = \sum_{i=1}^{\infty} Q^i$ .*

(i) *If  $Q_\nu^i \rightarrow Q^i$  for all  $i$ , then  $Q_\nu \rightarrow Q$ .*

(ii) *If  $\mathcal{E}$  is a system of sets such that  $Q_\nu^i(E) \rightarrow Q^i(E)$  uniformly for all  $E \in \mathcal{E}$ , then  $Q_\nu(E) \rightarrow Q(E)$  uniformly for all  $E$ .*

PROOF. We prove (ii) first. Fix  $\varepsilon > 0$ . Choose  $k$  such that

$$(50) \quad \sum_{i=k+1}^{\infty} Q^i(\mathbb{X}) \leq \frac{\varepsilon}{4}.$$

Choose  $\nu_0 = \nu(\varepsilon)$  such that for  $\nu > \nu_0$ ,

$$(51) \quad |Q_\nu^i(\mathbb{E}) - Q^i(\mathbb{E})| \leq \frac{\varepsilon}{4k}$$

for all  $i = 1, 2, \dots, k$ , all  $E \in \mathcal{E}$ , and also  $E = \mathcal{X}$ . Note that

$$(52) \quad \begin{aligned} \sum_{i=1}^k Q_\nu^i(\mathbb{X}) &\geq \sum_{i=1}^k Q^i(\mathbb{X}) - \sum_{i=1}^k |Q_\nu^i(\mathbb{X}) - Q^i(\mathbb{X})| \\ &\geq 1 - \frac{\varepsilon}{4} - \frac{k\varepsilon}{4k} = 1 - \frac{\varepsilon}{2}. \end{aligned}$$

Combining (50), (51), and (52) yields for the given  $\varepsilon > 0$ : there is  $\nu_0$  such that for all  $\nu > \nu_0$  and all  $E \in \mathcal{E}$ ,

$$\begin{aligned} |Q_\nu(E) - Q(E)| &\leq \left| \sum_{i=k+1}^{\infty} Q_\nu^i(E) \right| + \left| \sum_{i=1}^k Q_\nu^i(E) - Q^i(E) \right| + \left| \sum_{i=k+1}^{\infty} Q^i(E) \right| \\ &\leq \sum_{i=k+1}^{\infty} Q_\nu^i(\mathbb{X}) + \sum_{i=1}^k |Q_\nu^i(E) - Q^i(E)| + \sum_{i=k+1}^{\infty} Q^i(\mathbb{X}) \leq \frac{\varepsilon}{4} + \frac{k\varepsilon}{4k} + \frac{\varepsilon}{2} = \varepsilon. \end{aligned}$$

Since  $\varepsilon$  was arbitrary, the uniform continuity is proved.

The proof of (i) goes along the same lines, only in (52) we take  $E$  to be just a single  $Q$ -continuity set; it is clear that then it is also a  $Q^i$ -continuity set for all  $i$ .  $\square$

**Additional geometric facts.** We need a lemma similar to Theorem 3.2, but specific for regression setting.

LEMMA 6.3. *For any probability  $Q$  on  $\mathbf{X} \times \mathbb{R}$ , the set  $\{\vartheta: Q(\{(x, y): y = \vartheta^\top x\}) > 0\}$  has zero Lebesgue measure in  $\Theta = \mathbb{R}^p$ .*

PROOF. Consider an atomic system  $\mathcal{A}$  in  $\mathbf{Z} = \mathbb{R} \times \mathbf{X}$ , the union of  $\mathcal{A}^j$  for  $j = 0, 1, 2, \dots, k$ , defined in the following way. The single element of  $\mathcal{A}^0$  is  $\emptyset$ ;  $\mathcal{A}^1$  contains all points  $(x, y) \in \mathbf{Z}$ , that is, all sets consisting of all  $(x, x^\top \vartheta) \in \mathbf{Z}$  such that  $\vartheta$  belong to a codimension 1 affine subspace of  $\mathbb{R}^p$ ;  $\mathcal{A}^2$  contains all sets consisting of all  $(x, x^\top \vartheta) \in \mathbf{Z}$  such that  $\vartheta$  belong to a codimension 2 affine subspace of  $\mathbb{R}^p$  (lines in the simple linear regression); generally,  $\mathcal{A}^j$  contains all sets composed from all or  $(x, x^\top \vartheta) \in \mathbf{Z}$  such that  $\vartheta$  belongs to a codimension  $j$  affine subspace of  $\mathbb{R}^p$ . It is straightforward to verify that  $\mathcal{A}$  is an atomic system. By Propositions 5.4 and 5.5, any  $E \in \mathcal{A}$  such that  $Q(A) > 0$  contains an  $A$  from a countable system  $\mathcal{A}_Q$ ; there is a similar dual structure, the set of all  $\vartheta$  such that  $(x, x^\top \vartheta) \in E$  is contained in the set of all  $\vartheta$  with  $(x, x^\top \vartheta) \in A$ ; and any such set has codimension at least 1, thus has zero Lebesgue measure in  $\mathbb{R}^p$ , as does any countable union of them.  $\square$

The rest of this subsection is oriented towards the proof of Lemma 6.6, which gives a special condition for the uniform convergence of  $Q_\nu$  to  $Q$  on halfspaces, applicable to regression probability fields. We keep the same halfspace notation  $\mathbf{H}_u$  and  $\mathbf{G}_u$  for the corresponding hemispheres  $\mathbf{H}_u \cap \mathbb{S}^{p-1}$  and  $\mathbf{G}_u \cap \mathbb{S}^{p-1}$ . The collection of all finite intersections of closed hemispheres in  $\mathbb{S}^{p-1}$  is denoted by  $\mathcal{H}^{p-1}$ . Recall that a system of sets is called a *convergence-determining class* if  $Q_\nu(E) \rightarrow Q(E)$  for all its elements  $E$  implies that  $Q_\nu \rightarrow Q$ .

LEMMA 6.4. *For any  $p$ ,  $\mathcal{H}^{p-1}$  is a convergence-determining class on  $\mathbb{S}^{p-1}$ .*

PROOF. The lemma follows from Corollary 1, page 14 of Billingsley (1968):  $\mathcal{H}^{p-1}$  is closed under the formation of finite intersections and given any  $x \in \mathbb{S}^{p-1}$  and any  $\varepsilon > 0$ , there is a finite intersection  $E$  of hemispheres such that  $x$  lies in the interior of  $E$  and  $E$  is contained in the ball with center  $x$  and radius  $\varepsilon$ .  $\square$

Suppose that  $\mathcal{A}$  is an atomic system composed from all linear (proper) subspaces of  $\mathbb{R}^p$  intersected with  $\mathbb{S}^{p-1}$ . That is,  $\mathcal{A}^0$  contains only  $\emptyset$ ,  $\mathcal{A}^1$  only antipodal pairs of points,  $\mathcal{A}^2$  circles, and so forth.

LEMMA 6.5. *Let  $H$  be from  $\mathcal{H}^{p-1}$ . If  $A \in \mathcal{A}^j$ , then there is a closed set  $A_H \subseteq A$  such that  $(A \cap H) \cup A_H = A$ , and  $(A \cap H) \cap A_H = H \cap A_H$  is contained in a (possibly empty) union of finite number of elements from  $\mathcal{A}^i$ ,  $i < j$ .*

PROOF. We define  $A_H$  to be  $\text{cl}(A \cap H^c)$ . Since  $A$  is closed,  $A_H \subseteq A_H$ ; thus  $(A \cap H) \cup A_H \subseteq A$ . Conversely, the properties of closure give that

$$(A \cap H) \cup A_H \supseteq (A \cap H) \cup (A \cap H^c) = A \cap (H \cup H^c) = A;$$

hence  $(A \cap H) \cup A_H = A$ .

The rest of the proposition is proved by the induction with respect to number of hemispheres whose intersection gives  $H$ . Suppose first that  $H$  is itself a closed hemisphere. If  $A \subseteq H$ , then  $A \cap H^c = \emptyset$  and thus also  $A_H = H \cap A_H = \emptyset$ . Otherwise, again by the properties of closure,

$$H \cap A_H \subseteq H \cap A \cap \text{cl}(H^c) = A \cap H \cap H^c = A \cap \partial H;$$

the properties of the atomic system imply that the last intersection belongs to the  $\mathcal{A}^i$  with  $i < j$  ( $A \cap \partial H$  is a subset of  $H$ , but it is not equal to  $H$ , since  $A$  is not a subset of  $H$ ).

Suppose now that the proposition holds for  $H_1$  and  $H_2$ ; let  $H = H_1 \cap H_2$ . By the properties of closure,

$$\begin{aligned} (53) \quad H \cap A_H &= H \cap \text{cl}(A \cap (H_1 \cap H_2)^c) = H \cap (\text{cl}(A \cap H_1^c) \cup \text{cl}(A \cap H_2^c)) \\ &= (H \cap \text{cl}(A \cap H_1^c)) \cup (H \cap \text{cl}(A \cap H_2^c)) \\ &= (H_1 \cap \text{cl}(A \cap H_1^c)) \cup (H_2 \cap \text{cl}(A \cap H_2^c)) \end{aligned}$$

The induction assumption implies that any of the two last terms in (53) are contained in the union of a finite number of elements of  $\mathcal{A}^i$  with  $i < j$ . Therefore the same is true for  $A \cap A_H$ .  $\square$

We call a set  $E \subseteq \mathbb{S}^{p-1}$  *symmetric* if  $E = -E = \{-x : x \in E\}$ . Let  $\pi$  be the factorization mapping from  $\mathbb{S}^{p-1}$  to the projective plane  $\mathbb{RP}^{p-1}$ , the mapping identifying antipodal points:  $\pi(-x) = \pi(x)$  for all  $x \in \mathbb{S}^{p-1}$ . Any symmetric set is a preimage of a subset of in  $\mathbb{RP}^{p-1}$  under  $\pi$ .

LEMMA 6.6. *Let  $Q_\nu$  and  $Q$  be probability measures on  $\mathbb{S}^{p-1}$  such that  $Q_\nu(E) = Q(E)$  for all symmetric sets  $E$  and all  $\nu$ . If  $Q_\nu \rightarrow Q$ , then  $Q_\nu(H_u) \rightarrow Q(H_u)$  uniformly in  $u \in \mathbb{S}^{p-1}$ .*

PROOF. Consider the atomic decomposition of  $Q_\nu$  and  $Q$ , with the atomic system  $\mathcal{A}$  used in Lemma 6.5 above. Since all sets in  $\mathcal{A}$  are symmetric, it follows that  $\mathcal{A}_{Q_\nu}^j = \mathcal{A}_Q^j$  for all  $\nu$  and all  $j$ ; only measures  $Q_\nu^A$  may differ.

We show first that  $Q_\nu^A \rightarrow Q^A$  for all  $A$ . In view of Lemma 6.5, it is sufficient to prove that  $Q_\nu^A(E) \rightarrow Q^A(E)$  for all  $E$  which are intersections of finitely many closed hemispheres. Suppose that  $A \in \mathcal{A}_Q^j$  and let  $E \in \mathcal{H}^{p-1}$ . The set  $A \cap E$  is closed, and by Lemma 6.5 there is a closed set  $A_E \subseteq A$  such that  $(A \cap E) \cup A_E = A$ . The intersection  $(A \cap E) \cap A_E = E \cap A_E$  lies in the union of finitely many sets from  $\mathcal{A}^i$ ,  $i < j$ ; hence  $Q^A((A \cap E) \cap A_E) = 0$ , as well as  $Q_\nu^A((A \cap E) \cap A_E) = 0$  for all  $\nu$ . The application of Lemma 6.1 then gives the desired convergence. The weak convergence of  $Q_\nu^\omega$  to  $Q^\omega$  follows from Lemma 6.2(i).

Fix  $A \in \mathcal{A}^j$ . Let  $H_u$  be a closed hemisphere in  $\mathbb{S}^{p-1}$ . If  $A \subseteq H_u$ , then

$$(54) \quad Q_\nu^A(H_u) = Q_\nu^A(A) = Q^A(A) = Q^A(H_u),$$

thus the convergence is clearly uniform in this case. If  $H_u$  does not contain  $A$ , then  $A \cap \partial H_u$  lies in  $\mathcal{A}^i$ ,  $i < j$ , and therefore  $Q^A(\partial H_u) = 0$ . The intersections with  $A$  of all hemispheres not containing  $A$  form a set of hemispheres which is isomorphic to the set of all hemispheres in  $\mathbb{S}^{j-1}$ . By the properties of the atomic decomposition given by Proposition 5.4, all these hemispheres are  $Q^A$ -continuity sets. Thus, we may apply Proposition 5.11 and obtain the uniform convergence for these hemispheres. Combining this convergence with that implied with (54), we have that

$Q_\nu^A(H_u) \rightarrow Q^A(H_u)$  uniformly for all  $u \in \mathbb{S}^{p-1}$ . Since  $Q^\omega(\partial H_u) = 0$  for all  $u$ , Proposition 5.11 gives the same conclusion for the sequence  $Q_\nu^\omega$ . The proof of the lemma is then finished by the application of Lemma 6.2(ii).  $\square$

Lemma 6.6 helps us to prove convergence of depth for sequences of measures on  $\mathbb{S}^{p-1}$  with a constant marginal on  $\mathbb{RP}^{p-1}$ . Given a probability  $Q$  on  $\mathbb{X}$  and a mapping  $f$  whose range is a subset  $\mathbb{X}$ , let  $f^{-1}(Q)$  denote the set of all probabilities of the form  $Q \circ f^{-1}$ . For any probability  $Q$  on  $\mathbb{S}^{p-1}$ , the definitions  $h_Q$ ,  $S(Q)$ ,  $d(Q)$  and  $s_\lambda(Q)$  are extended in the natural way: hemispheres stand for halfspaces.

LEMMA 6.7. *For any probability  $\bar{Q}$  on  $\mathbb{RP}^{p-1}$ ,  $\pi^{-1}(\bar{Q})$  is a compact convex subset of the space of all probabilities on  $\mathbb{S}^{p-1}$  (with respect to the weak topology). The function  $d(\cdot)$ , as well as  $s_{d(\cdot)+\varepsilon}(\cdot)$  for any  $\varepsilon > 0$ , is continuous on  $\pi^{-1}(\bar{Q})$ .*

PROOF. From the continuity of  $\pi$  follows that  $\pi^{-1}(\bar{Q})$  is closed; since it is a subset of the compact space of the probabilities on  $\mathbb{S}^{p-1}$ , it is itself compact. Convexity follows by the straightforward verification involving symmetric sets.

Suppose that  $Q_\nu \rightarrow Q$  and all  $Q_\nu$  and  $Q$  belong to  $\pi^{-1}(\bar{Q})$ . Lemma 6.6 gives that  $h_{Q_\nu}(u) \rightarrow h_Q(u)$  uniformly for all  $u \in \mathbb{S}^{p-1}$ ; the convergence of  $d(Q_\nu)$  to  $d(Q)$  follows. The continuity of the skeleton function follows then from Proposition 5.15.  $\square$

LEMMA 6.8. *Let  $Q$  be a probability on  $\mathbb{R}^p$  such that  $Q(H_{-u}) = 1$  for some  $u \in \mathbb{S}^{p-1}$ . If  $\lambda > d(Q)$  and  $Q(G_{-u}) > 0$ , then  $s_\lambda(Q) \in G_{-u}$ ; in particular,  $s_\lambda(Q) \neq 0$ .*

PROOF. The proof is a follow-up of that of Proposition 5.17 and uses the same notation. To prove that  $s_\lambda(Q) \in G_{-u}$ , we have to show only that the first integral in (32) is nonzero. This is true if sharp inequality in (30) holds for some set of  $v$ 's with positive measure  $\mu$ .

Let  $u$  satisfy the assumptions of the lemma. We will prove that there is a neighborhood  $\mathcal{U}(u, \varepsilon)$  of this  $u$  such that

$$(55) \quad h_Q(v) < h_Q(v^*) \quad \text{for every } v \in \mathcal{U}(u, \varepsilon).$$

(recall that  $v^* = v - 2u(u^T v)$  is the reflection of  $v$  about  $\partial H_u$ ). Given  $\delta > 0$ , we write  $H_u^\delta$  for the  $\delta$ -fattening of  $H_u$ :  $H_u^\delta = \{x \in \mathbb{R}^p : ((x))^T u \geq -\delta\}$ . The proof of Proposition 5.17 showed, via Proposition 5.13, that

$$Q(H_v) = Q(H_v \setminus H_{v^*}) + Q(H_v \cap H_{v^*}) = Q(H_v \cap H_{v^*}).$$

Since, on the other hand,

$$Q(H_{v^*}) = Q(H_{v^*} \setminus H_v) + Q(H_v \cap H_{v^*}),$$

we obtain that

$$(56) \quad Q(H_{v^*}) - Q(H_v) = Q(H_{v^*} \setminus H_v).$$

For arbitrary  $v \in \mathbb{S}^{p-1}$  the following holds: if  $x \in H_v$ , then

$$(57) \quad ((x))^T u = ((x))^T (u - v) + ((x))^T v \geq ((x))^T (u - v) \geq -\|u - v\|$$

and therefore  $x \in H_u^{\|u-v\|}$ . Hence,  $H_v \subseteq H_u^{\|u-v\|}$  and thus

$$(58) \quad H_v^c \supseteq (H_u^{\|u-v\|})^c.$$

The reflection about  $\partial H_u$  carries  $v$  to  $v^*$  and  $u$  to  $-u$ ; it is a Euclidean transformation, hence  $\|u - v\| = \|(-u) - v^*\|$ . Suppose that  $x \in H_{v^*}^c$ . Then  $x^T v^* < 0$ ; similarly as in (57), we obtain that

$$\begin{aligned} ((x))^T(-u) &= ((x))^T((-u) - v^*) + ((x))^T v^* \leq ((x))^T((-u) - v^*) \\ &\leq \|(-u) - v^*\| = \|u - v\|; \end{aligned}$$

hence  $x \in H_u^{\|u-v\|}$ . Thus,

$$(59) \quad H_{v^*}^c \supseteq (H_u^{\|u-v\|})^c.$$

Combining (58) and (59), we obtain from (56) that

$$Q(H_{v^*}) - Q(H_v) \geq Q((H_u^{\|u-v\|})^c).$$

For  $\|u - v\| \rightarrow 0$ , the set  $(H_u^{\|u-v\|})^c$  monotonically increases to  $H_u^c = G_{-u}$ ; since  $Q(G_u) > 0$ , we may conclude that there is  $\varepsilon > 0$  such that  $Q(H_{v^*}) - Q(H_v) > 0$  whenever  $\|u - v\| < \varepsilon$ . This proves (55).

Suppose now that  $\lambda > d(Q)$ . There is  $v_2$  such that  $h_Q(v_2) < \lambda$ . Take the arc connecting  $u$  and  $v_2$ ; according to Proposition 5.13,  $h_Q$  is nonincreasing along this arc when  $v$  approaches  $u$ . Therefore, we may choose  $v_1 \neq u$  in  $\mathcal{U}(u, \varepsilon)$  such that  $h_Q(v_1) < \lambda$ . By Proposition 5.9,  $h_Q$  is upper semicontinuous, hence there is a neighborhood  $\mathcal{U}(v_1, \eta)$  such that  $h_Q(v) < \lambda$  for all  $v \in \mathcal{U}(v_1, \eta)$ ; we may choose this neighborhood to be contained in  $\mathcal{U}(u, \varepsilon)$ . In view of (55), we have proved that  $h_Q(v) \wedge \lambda < h_Q(v^*) \wedge \lambda$  for all  $v \in \mathcal{U}(v_1, \eta)$ . This concludes the proof, since  $\mu(\mathcal{U}(v_1, \eta)) > 0$ .  $\square$

Note that Proposition 3.1 implies that the set  $\{v \in \mathbb{S}^{p-1} : h_Q(v) \leq \lambda\} \cap S$  is nonempty for any  $S$  dense in  $\mathbb{S}^{p-1}$  and any  $\lambda > d(Q)$ .

**Set-valued topology and spheric closures.** The proof of Theorem 3.3 uses the theory of set-valued vector fields; for more background on this mathematical apparatus, see Rockafellar and Wets (1998) or Klein and Thompson (1984). A *set-valued mapping*  $\mathcal{F}$  from a metric space  $S$  to a metric space  $T$  is any mapping from  $S$  to the set  $2^T$  of all subsets of  $T$ . Any set-valued mapping has the *graph*  $\text{Gr}(\mathcal{F}) = \{(s, t) \in S \times T : t \in \mathcal{F}(s)\}$ . Conversely, any subset of  $S \times T$  defines a set-valued mapping from  $S$  to  $T$ . The *domain* of a set-valued mapping  $\mathcal{F}$  is the  $\text{Dom } \mathcal{F} = \{s : \mathcal{F}(s) \neq \emptyset\}$ , the *range* is the set  $\text{Rng } \mathcal{F} = \{t : t \in \mathcal{F}(s) \text{ for some } s \in S\}$ . The *composition* of set-valued mappings  $\mathcal{F}_1$  from  $S_1$  to  $S_2$  and  $\mathcal{F}_2$  from  $S_2$  to  $S_3$  is the set-valued mapping  $\mathcal{F}$  from  $S_1$  to  $S_3$  such that

$$\mathcal{F}(s_3) = \{s_3 : s_3 \in \mathcal{F}_2(s_2) \text{ and } s_2 \in \mathcal{F}_1(s_1)\}.$$

A set-valued mapping  $\mathcal{F}$  is called *strongly outer semicontinuous* at  $s$ , if for any open set  $E \supseteq \mathcal{F}(s)$ , any sequence  $t_n \in \mathcal{F}(s_n)$  with  $s_n \rightarrow s$  is eventually in  $E$ . The more widespread terminology is “upper semicontinuous”, but we would like to avoid confusion with the upper



semicontinuity used in the previous sections. A set-valued mapping  $\mathcal{F}$  is *single-valued* at  $s$ , if  $\mathcal{F}(s)$  is a singleton; in such a case, it is strongly outer semicontinuous if and only if it is continuous as the ordinary mapping in the usual sense. A *closure* of  $\mathcal{F}$  is the set-valued function  $\bar{\mathcal{F}}$  such that  $\text{Gr}(\bar{\mathcal{F}}) = \text{cl}(\text{Gr}(\mathcal{F}))$ . A function with closed graph ( $\bar{\mathcal{F}} = \mathcal{F}$ ) is called *closed*—“outer semicontinuous” by Rockafellar and Wets (1998). A set-valued function  $\mathcal{F}$  is *closed-valued* (*compact-valued*, *convex-valued*) if  $\mathcal{F}(s)$  is closed (compact, convex) for all  $s$ , respectively. Note that each closed function is closed-valued, but the converse is not necessarily true.

LEMMA 6.9. *Let  $S, T$  be metrizable spaces. If  $T$  is compact, then a set-valued mapping  $\mathcal{F}$  from  $S$  to  $T$  is closed if and only if it is strongly outer semicontinuous and closed-valued.*

PROOF. See Klein and Thompson (1984), Theorems 7.1.15 and 7.1.16, page 78, or Nikaido (1968), Lemma 4.4, page 66.  $\square$

Let  $S \subseteq \bar{S}$  and  $\mathcal{F}, \bar{\mathcal{F}}$  be set-valued functions from  $S, \bar{S}$ , respectively, to  $T$ ;  $\bar{\mathcal{F}}$  is called an *extension* of  $\mathcal{F}$ , if  $\bar{\mathcal{F}}(s) = \mathcal{F}(s)$  for all  $s \in S$ .

LEMMA 6.10. *Let  $S, \bar{S}, T$  be metrizable spaces. Suppose that  $T$  is compact and  $S$  is a dense subset of  $\bar{S}$ .*

(i) *If  $\mathcal{F}$  is closed, with respect to the relative topology on  $S \times T$  (particularly, if  $\mathcal{F}$  is a single-valued mapping continuous on  $S$ ), then its closure  $\bar{\mathcal{F}}$  in  $S \times T$  is a closed extension of  $\mathcal{F}$  to  $\bar{S}$ .*

(ii) *If  $\mathcal{F}$  is single-valued on  $S$ , with values given by an ordinary (single-valued) mapping  $f$  continuous on  $S$ , then the values of the closed extension  $\bar{\mathcal{F}}$  at  $s \in \bar{S} \setminus S$  are exactly all limit values of all possible sequences  $f(s_\nu)$  with  $s_\nu \rightarrow s$ .*

PROOF. For (i), see Nikaido (1968), Theorem 4.7, page 72. To see (ii), note first that all limit points of sequences  $f(s_\nu)$  with  $s_\nu \rightarrow s$  are in  $\text{Gr}(\bar{\mathcal{F}})$ , since the latter is closed. Suppose that  $s_0 \in \bar{S} \setminus S$  and  $t_0 \in \bar{\mathcal{F}}(s_0)$ . If there is an open neighborhood  $\mathcal{U}$  of  $(s_0, t_0)$  in  $\bar{S} \times T$  containing no point of  $\text{Gr}(\bar{\mathcal{F}})$  with  $s \in S$  then  $\text{Gr}(\bar{\mathcal{F}}) \setminus \mathcal{U}$  is a closed set containing  $\text{Gr}(\mathcal{F})$ ; this is a contradiction with the minimality of the closure. Hence there is a sequence  $s_\nu \rightarrow s_0$  such that  $f(s_\nu) \rightarrow t_0$ .  $\square$

Unlike the closure operation, we define the *closed convex hull* of the set-valued function  $\mathcal{F}$  pointwise:  $(\text{cc } \mathcal{F})(s) = \text{cl}(\text{conv}(\mathcal{F}(s)))$  for all  $s$ .

LEMMA 6.11. *Let  $\mathcal{F}$  be a compact-valued set-valued mapping from  $S$  to  $T$ , a complete metrizable convex subspace of a linear topological space. The strong outer semicontinuity of  $\mathcal{F}$  implies that of  $\text{cc } \mathcal{F}$ .*

PROOF. See Nikaido (1968), Theorem 4.8, page 72, or Borisovich, Gel'man, Myshkis and Obukhovskii (1982), Theorem 1.3.21 on page 138.  $\square$

Recall that the set of all probability measures on  $\mathbb{S}^{p-1}$  with the weak topology is a compact, convex subset of complete metrizable space of all positive continuous linear functionals on the space  $\mathcal{C}(\mathbb{S}^{p-1})$  of continuous functions from  $\mathbb{S}^{p-1}$  to  $\mathbb{R}$ . The latter is the convex cone in the locally convex topological linear space of all continuous linear functionals on  $\mathcal{C}(\mathbb{S}^{p-1})$  with the weak\*-topology—see Tjur (1980), page 19. Recall also the definition of  $P \circ g^{-1}$  for a measure  $P$

and a function  $g$  on  $\mathbb{X}$ . This definition works well even when  $g$  is defined on a subset  $\tilde{\mathbb{X}}$ , resulting in  $(P \llcorner \tilde{\mathbb{X}}) \circ g^{-1}$ . We will use the notation  $P \circ g^{-1}$  in this extended meaning; of course, if  $P$  is supported by  $\tilde{\mathbb{X}}$ , the result is the same.

Let  $\{P_\vartheta : \vartheta \in \Theta\}$  be a probability field on  $\mathbb{R}^p$ . We define a new, set-valued mapping  $\mathcal{P}(\vartheta)$  from  $\text{cl } \dot{\Theta}$  to the set of probability measures on  $\mathbb{S}^{p-1}$ —a set-valued probability field on  $\mathbb{S}^{p-1}$ , indexed by  $\Theta$ . We start with the set  $\dot{\Theta}$  of all  $\vartheta$  such that  $P_\vartheta(\{0\}) = 0$ ; we require that the original probability field  $\{P_\vartheta : \vartheta \in \Theta\}$  is continuous (in the weak topology) on  $\dot{\Theta}$ , and also that  $\dot{\Theta}$  is dense in  $\Theta$ ; these conditions give the construction its sense. The closure of  $\dot{\Theta}$  could be understood in  $\Theta$ ; however, it is more convenient to consider it in a compactification  $\bar{\Theta}$  instead—if  $\dot{\Theta}$  is dense in  $\Theta$ , so it is in  $\bar{\Theta}$ .

We define  $\mathcal{P}$  in three steps. Let  $\kappa$  be the mapping assigning to  $x$  from  $\mathbb{R}^p$  its direction  $((x))$ ; note that the mapping is continuous on  $\mathbb{R}^p \setminus \{0\}$ . In the first step, we set  $\mathcal{P}(\vartheta) = \{P_\vartheta \circ \kappa^{-1}\}$  for any  $\vartheta \in \dot{\Theta}$ ; any such  $P_\vartheta$  is supported by the complement of  $\{0\}$ , hence  $\mathcal{P}$  is single-valued on  $\dot{\Theta}$ .

In the second step, we construct a strongly outer semicontinuous set-valued extension  $\dot{\mathcal{P}}$  of  $\mathcal{P}$  to  $\text{cl } \dot{\Theta}$ . By the continuous mapping theorem,  $P_{\vartheta_\nu} \circ \kappa^{-1}$  converges weakly to  $P_\vartheta \circ \kappa^{-1}$ , whenever  $\vartheta \in \dot{\Theta}$  and  $P_{\vartheta_\nu} \rightharpoonup P_\vartheta$ . Therefore,  $\mathcal{P}$  is continuous on  $\dot{\Theta}$ . This is all we need to apply Lemma 6.10(i): the closure  $\dot{\mathcal{P}}$  of  $\mathcal{P}$  is closed on  $\text{cl } \dot{\Theta}$ ; then Lemma 6.9 yields its strong outer semicontinuity. The closure is also an extension of  $\mathcal{P}$ , that is,  $\dot{\mathcal{P}}(\vartheta)$  is single-valued and equal to  $\mathcal{P}(\vartheta)$  for all  $\vartheta \in \dot{\Theta}$ .

In the third step, we define  $\mathcal{P}(\vartheta)$  to be the convex closure of  $\dot{\mathcal{P}}$ : for all  $\vartheta \in \Theta$ ,  $\mathcal{P}(\vartheta) = \text{cc } \dot{\mathcal{P}}(\vartheta)$ . The convexification leaves the singleton values for  $\vartheta \in \dot{\Theta}$  unchanged, hence our notation remains consistent. Lemma 6.11 implies that  $\mathcal{P}(\vartheta)$  is strongly outer semicontinuous (on the whole  $\Theta$ ); obviously, it is also convex-valued.

We call the resulting set-valued probability field on  $\mathbb{S}^{p-1}$  indexed by  $\bar{\Theta}$  a **spheric closure** of the original field  $\{P_\vartheta : \vartheta \in \mathbb{R}^p\}$ . Spheric closures replace probability fields on  $\mathbb{R}^p$  by those on  $\mathbb{S}^{p-1}$ , bypassing thus assumption (iii) of Theorem 5.19. For  $\vartheta$  outside  $\dot{\Theta}$ , the spheric closure may be thought of as replacing the original probability by the set of probabilities representing all ways how the mass from the origin can be spread along  $\mathbb{S}^{p-1}$ . In the tangent depth setting, such a replacement can never raise depth. Recall that  $h_Q(u) = Q(\mathbf{H}_u)$ , for any measure  $Q$  on  $\mathbb{S}^{p-1}$  (where  $\mathbf{H}_u$  denotes hemispheres) or on  $\mathbb{R}^p$  or  $\mathbb{R}^p$  (where  $\mathbf{H}_u$  denotes halfspaces).

**LEMMA 6.12.** *Let  $\Phi$  be a function from  $\Theta \times \mathbb{Z}$  to  $\mathbb{R}^p$  continuous in the variable  $\vartheta$  for  $Q$ -almost all  $z$ . Let  $P$  be a probability on  $\mathbb{Z}$  and let  $\mathcal{P}$  be a spheric closure of the probability field  $\{P_\vartheta : \vartheta \in \Theta\}$  such that  $P_\vartheta = P \circ \Phi_\vartheta^{-1}$ . For any  $\vartheta \in \Theta$  and any  $Q \in \mathcal{P}(\vartheta)$ ,  $h_Q(u) \leq h_{P_\vartheta}(u)$  for all  $u \in \mathbb{S}^{p-1}$ ; in particular,  $d(Q) \leq d(\vartheta, P)$ .*

**PROOF.** If  $Q$  is the single element of  $\mathcal{P}(\vartheta)$  for  $\vartheta \in \dot{\Theta}$ , then the lemma trivially holds. Fix  $\vartheta \in \Theta \setminus \dot{\Theta}$ . We will show that any  $Q \in \mathcal{P}(\vartheta)$  is of the form

$$(60) \quad Q = \hat{P} \circ \Phi_\vartheta^{-1} \circ \kappa^{-1} + \tilde{Q}$$

where  $\hat{P} = P \llcorner (\mathbb{Z} \setminus \Phi_\vartheta^{-1}(\{0\}))$  and

$$(61) \quad \tilde{Q}(\mathbb{S}^{p-1}) = 1 - \hat{P}(\mathbb{S}^{p-1}) = P(\Phi_\vartheta^{-1}(\{0\})) = P_\vartheta(\{0\})$$

Observe that the set of  $Q$  satisfying (60) is convex and closed under weak convergence; hence it is sufficient to prove (60) for all  $Q \in \dot{\mathcal{P}}(\vartheta)$ .

Fix  $Q \in \dot{\mathcal{P}}(\vartheta)$ . According to Lemma 6.10(ii), there is a sequence  $\vartheta_\nu \rightarrow \vartheta$  such that

$$(62) \quad P_{\vartheta_\nu} \circ \kappa^{-1} = P \circ \Phi_{\vartheta_\nu}^{-1} \circ \kappa^{-1} \rightharpoonup Q.$$

The decomposition

$$P = \widehat{P} + \dot{P} = P \llcorner (Z \setminus \Phi_\vartheta^{-1}(\{0\})) + P \llcorner \Phi_\vartheta^{-1}(\{0\})$$

(the dependence of  $\vartheta$  is suppressed in the notation) leads to the decomposition for  $P_\vartheta$

$$(63) \quad P \circ \Phi_\vartheta^{-1} = \widehat{P} \circ \Phi_\vartheta^{-1} + \dot{P} \circ \Phi_\vartheta^{-1}$$

and also for  $P_{\vartheta_\nu} \circ \kappa^{-1}$ :

$$(64) \quad P \circ \Phi_{\vartheta_\nu}^{-1} \circ \kappa^{-1} = \widehat{P} \circ \Phi_{\vartheta_\nu}^{-1} \circ \kappa^{-1} + \dot{P} \circ \Phi_{\vartheta_\nu}^{-1} \circ \kappa^{-1}.$$

By the continuous mapping theorem,

$$(65) \quad \widehat{P} \circ \Phi_{\vartheta_\nu}^{-1} \circ \kappa^{-1} \rightharpoonup \widehat{P} \circ \Phi_\vartheta^{-1} \circ \kappa^{-1}.$$

From (62) and (65) follows that there is  $\tilde{Q}$  such that

$$(66) \quad \dot{P} \circ \Phi_{\vartheta_\nu}^{-1} \circ \kappa^{-1} \rightharpoonup \tilde{Q}.$$

The decomposition (60) follows from (62), (64), (65) and (66); the identity (61) from (66) and the fact that

$$\dot{P} \circ \Phi_{\vartheta_\nu}^{-1} \circ \kappa^{-1}(\mathbb{S}^{p-1}) = 1 - \widehat{P} \circ \Phi_{\vartheta_\nu}^{-1} \circ \kappa^{-1}(\mathbb{S}^{p-1}) = 1 - (1 - P(\Phi_\vartheta^{-1}(\{0\}))).$$

The decomposition (60) yields that

$$Q(\mathbf{H}_u) = \widehat{P} \circ \Phi_\vartheta^{-1} \circ \kappa^{-1}(\mathbf{H}_u) + \tilde{Q}(\mathbf{H}_u) \leq \widehat{P} \circ \Phi_\vartheta^{-1}(\mathbf{H}_u) + P(\Phi_\vartheta^{-1}(\{0\})) = P \circ \Phi_\vartheta^{-1}(\mathbf{H}_u)$$

where  $\mathbf{H}_u$  stands for a hemisphere or halfspace, whichever is appropriate.  $\square$

The appropriate topological tool for spheric closures uses an approach similar to that of Cellina (1969), with a slightly more powerful approximation theorem.

**LEMMA 6.13.** *Suppose that  $S$  is a compact, convex subset of a convex metrizable subspace of a locally convex topological vector space. Let  $\mathcal{F}$  be a set-valued vector field which is a composition of strongly outer semicontinuous convex-valued mapping  $F$  from  $\mathbb{D}^p$  to  $S$  and a single-valued mapping given by a continuous function  $f$  from  $S$  to  $\mathbb{D}^p$ . Suppose that  $\mathcal{F}(\vartheta) \subseteq \mathbf{G}_{g(\vartheta)}$  for all  $\vartheta \in \mathbb{S}^{p-1}$ , where  $g$  is a continuous function from  $\mathbb{S}^{p-1}$  to  $\mathbb{S}^{p-1}$ . If  $g$  is not homotopic to a constant, then there is a point  $\vartheta$  in the interior of  $\mathbb{D}^p$  such that  $0 \in \mathcal{F}(\vartheta)$ .*

**PROOF.** First, we show that there is  $\eta > 0$  such that

$$(67) \quad u^\top g(\vartheta) \geq \eta \quad \text{whenever } u \in \mathcal{F}(\vartheta) \text{ and } \vartheta \in \mathbb{S}^{p-1}.$$

Assume the contrary. Then there are  $u_\nu, \vartheta_\nu$  such that  $u_\nu \in \mathcal{F}(\vartheta_\nu)$  and  $u_\nu^\top g(\vartheta_\nu) \rightarrow 0$ . From the compactness of  $\mathbb{S}^{p-1}$  and  $\mathbb{D}^p$  follows that we may assume, passing to a subsequence if necessary, that  $\vartheta_\nu \rightarrow \vartheta \in \mathbb{S}^{p-1}$  and  $u_\nu \rightarrow u \in \mathbb{D}^p$ . The composition  $\mathcal{F}$  of the strongly outer semicontinuous

mapping  $F$  and single-valued continuous, hence also strongly outer semicontinuous mapping  $f$  (we abuse slightly the notation) is itself strongly outer semicontinuous; see Klein and Thompson (1984), Theorem 7.3.11 on page 87, or Nikaido (1968), Theorem 4.6 on page 71. We have that  $u_\nu^T g(\vartheta_\nu) \rightarrow u^T g(\vartheta) = 0$ , a contradiction.

Now, we invoke the approximation theorem: Proposition 1.1.10 on page 70 of Borisovich, Gel'man, Myshkis and Obukhovskii (1980). For any  $\varepsilon > 0$ , there exists a single-valued  $\varepsilon$ -approximation of  $F$ : a continuous function  $f_\varepsilon$  from  $\mathbb{D}^p$  to  $S$  with the following property: for every  $\vartheta \in \mathbb{D}^p$  there is  $\vartheta_\varepsilon \in \mathbb{D}^p$  and  $x_\varepsilon \in F(\vartheta_\varepsilon)$  such that  $\|\vartheta - \vartheta_\varepsilon\| < \varepsilon$  and  $d(f_\varepsilon(\vartheta), x_\varepsilon) < \varepsilon$ , where  $d$  is the metric on  $S$  and the values of  $f_\varepsilon$  lie in the convex closure of  $S$  (which is  $S$  itself).

Using the uniform continuity of  $f$  ( $S$  is compact), we can extend this approximation to the composition of  $F$  and  $f$ , and construct a sequence  $f_\nu$  of functions from  $\mathbb{D}^p$  to  $S$  such that for every  $\vartheta \in \mathbb{D}^p$  there is  $\vartheta_\nu \in \mathbb{D}^p$  and  $u_\nu \in f(F(\vartheta)) = \mathcal{F}(\vartheta)$  satisfying

$$(68) \quad \|\vartheta - \vartheta_\nu\| < \frac{1}{\nu}, \quad \text{and} \quad \|f(f_\nu(\vartheta)) - u_\nu\| < \frac{1}{\nu},$$

and, moreover,

$$(69) \quad \|g(\vartheta) - g(\vartheta_\nu)\| < \frac{1}{\nu},$$

using also the uniform continuity of  $g$ .

Suppose that  $\vartheta \in \mathbb{S}^{p-1}$ . We will show that  $u_\nu^T g(\vartheta) > \eta$  for all  $\nu$  sufficiently large, uniformly for  $\vartheta \in \mathbb{S}^{p-1}$ . Assume the contrary. Then there is a sequence  $\vartheta^\nu$  such that (abusing the notation again)  $u_\nu^T g(\vartheta^\nu) \leq \eta$ . We may again suppose (by compactness and passing to a subsequence) that  $u_\nu \rightarrow u$  and  $\vartheta^\nu \rightarrow \vartheta$ , hence  $\vartheta_\nu^\nu \rightarrow \vartheta$  too. Then we obtain that  $u^T g(\vartheta) \leq \eta$ , a contradiction with (67): the strongly outer semicontinuity of  $\mathcal{F}$  yields that  $u \in \mathcal{F}(\vartheta)$ , since  $u_\nu \in \mathcal{F}(\vartheta_\nu^\nu)$ .

Now, for  $\nu$  sufficiently large:  $(g(\vartheta_\nu))^\top u_\nu^\top > \eta$ , and therefore also  $(g(\vartheta_\nu))^\top f(f_\nu(\vartheta)) > \eta$ , uniformly for  $\vartheta \in \mathbb{S}^{p-1}$ , due to (68) and (69). Therefore,  $f(f_\nu(\vartheta))$  and  $g(\vartheta)$  cannot point into opposite directions for any  $\vartheta$  (if  $\nu$  is sufficiently large); this means that  $f \circ f_\nu$  and  $g$  are homotopic. Thus, if  $g$  is not homotopic to a constant, then neither is  $f \circ f_\nu$ ; Proposition 5.18 yields a critical point  $\vartheta^\nu \in \mathbb{D}^p$  such that  $f(f_\nu(\vartheta^\nu)) = 0$ . By the compactness of  $\mathbb{D}^p$ , there is  $\vartheta \in \mathbb{D}^p$  such that (possibly for a subsequence)  $\vartheta^\nu \rightarrow \vartheta$ . Let  $\vartheta_\nu^\nu$  and  $u_\nu$  be the points satisfying (68); we have that  $\vartheta_\nu^\nu \rightarrow \vartheta$  and  $u_\nu \rightarrow 0$ . Since  $u_\nu \in \mathcal{F}(\vartheta_\nu^\nu)$ , the strongly outer semicontinuity of  $\mathcal{F}$  yields that  $0 \in \mathcal{F}(\vartheta)$ . Finally, the possibility that  $\vartheta \in \mathbb{S}^{p-1}$  is ruled out by our assumptions.  $\square$

As can be seen from the proof, the assumption that  $\mathcal{F}(\vartheta) \subseteq \mathbb{G}_{g(\vartheta)}$  is not really necessary. We might formulate Lemma 6.13 in the vein of Proposition 5.18: either there is a critical point on the boundary (and then there is nothing to prove), or  $0 \notin \mathcal{F}(\vartheta)$  for all  $\vartheta \in \mathbb{S}^{p-1}$ ; Lemma 6.13 holds under this assumption as well, only the proof is more tedious.

**Spheric closures of regression fields.** Recall that  $\pi$  denotes the factorization mapping from  $\mathbb{S}^{p-1}$  to the projective plane  $\mathbb{RP}^{p-1}$  and  $\pi^{-1}(\bar{Q})$  stands for the set of all probabilities of the form  $\bar{Q} \circ \pi^{-1}$  for a given probability  $\bar{Q}$  on  $\mathbb{RP}^{p-1}$ .

LEMMA 6.14. Let  $Z = (X, Y)$  be a random variable with values in  $\mathbf{X} \times \mathbb{R}$  where  $\mathbf{X} \subseteq \mathbb{R}^p$ ; let  $\Theta = \mathbb{R}^p$ . If  $\{P_\vartheta: \vartheta \in \Theta\}$  is the regression probability field defined by  $P_\vartheta = \mathcal{L}(XX^\top \vartheta - XY)$  for any  $\vartheta \in \Theta$ , then

- (i)  $\{P_\vartheta: \vartheta \in \Theta\}$  has a spheric closure  $\mathcal{P}(\vartheta)$ , well-defined for all  $\vartheta \in \bar{\Theta} = \bar{\mathbb{R}}^p$ ;
- (ii) there is a probability  $\bar{Q}$  on  $\mathbb{R}P^{p-1}$  such that  $Q \in \pi^{-1}(\bar{Q})$  for any  $Q \in \mathcal{P}(\vartheta)$  and any  $\vartheta \in \bar{\mathbb{R}}^p$ ;
- (iii) for any  $\vartheta \in \partial \mathbb{R}^p$ ,  $Q(\mathbf{H}_\vartheta) = 1$  and  $Q(\mathbf{G}_\vartheta) > 1 - \Delta(X)$  for any  $Q \in \mathcal{P}(\vartheta)$ .

PROOF. The proof of (i) consists merely in noting that in this particular setting, the density of  $\dot{\Theta}$  in  $\Theta$  follows from Lemma 6.3; and the continuity of  $\{P_\vartheta: \vartheta \in \Theta\}$  in any  $\vartheta \in \Theta$ , and thus in  $\dot{\Theta}$ , from Proposition 5.21(i).

For any  $\vartheta \in \dot{\Theta}$ , the measure  $P_\vartheta \circ \kappa^{-1}$  belongs to  $\pi^{-1}(\bar{Q})$  where  $\bar{Q}$  is a fixed probability measure on  $\mathbb{R}P^{p-1}$ ; just note that for any symmetric set  $E$ ,

$$\mathbb{P}[(XX^\top \vartheta - XY) \in E] = \mathbb{P}[(X) \operatorname{sgn}(X^\top \vartheta - XY) \in E] = \mathbb{P}[(X) \in E]$$

—thus  $(P_\vartheta \circ \kappa^{-1})(E)$  is the same for all  $\vartheta \in \dot{\Theta}$ . The rest of (ii) follows from the fact that  $\pi^{-1}(\bar{Q})$  is convex and closed in the weak topology.

In particular, for any  $\vartheta \in \bar{\mathbb{R}}^p$  and any  $Q \in \mathcal{P}(\vartheta)$ ,

$$(70) \quad Q(\partial \mathbf{H}_u) = \mathbb{P}[(X) \in \mathbf{H}_u] = \mathbb{P}[u^\top X = 0].$$

Fix  $\vartheta \in \partial \mathbb{R}^p$ . If  $Q \in \dot{\mathcal{P}}(\vartheta)$ , then Lemma 6.10(ii) yields a sequence  $\vartheta_\nu \in \dot{\Theta}$  such that  $\vartheta_\nu \rightarrow \vartheta$ , and  $Q_\nu = P_{\vartheta_\nu} \circ \kappa^{-1} \rightarrow Q$ . By Lemma 6.6,

$$\begin{aligned} Q(\mathbf{H}_{-(\vartheta)}) &= \lim_{\nu \rightarrow \infty} Q_\nu(\mathbf{H}_{-(\vartheta)}) = \lim_{\nu \rightarrow \infty} P_{\vartheta_\nu}(\mathbf{H}_{-(\vartheta)}) \\ &= \lim_{\nu \rightarrow \infty} \mathbb{P}[-((\vartheta))^\top (XX^\top \vartheta_\nu - XY) \geq 0] \\ &= \lim_{\nu \rightarrow \infty} \mathbb{P}\left[-((\vartheta))^\top \left(XX^\top ((\vartheta_\nu)) - \frac{XY}{\|\vartheta_\nu\|}\right) \geq 0\right] \\ &\leq \mathbb{P}[-((\vartheta))^\top XX^\top ((\vartheta)) \geq 0] = \mathbb{P}[X^\top ((\vartheta)) = 0] \end{aligned}$$

(here  $\mathbf{H}_u$  stands again both for a hemisphere or halfspace). By (70),  $Q(\mathbf{G}_{-(\vartheta)}) = 0$  and hence  $Q(\mathbf{H}_{(\vartheta)}) = 1$ , for any  $Q \in \dot{\mathcal{P}}(\vartheta)$ . This conclusion is preserved under convex combinations and also under taking limits in the weak topology; thus,  $Q(\mathbf{H}_{(\vartheta)}) = 1$  for all  $Q \in \mathcal{P}(\vartheta)$ . By (70) again, we obtain that  $Q(\mathbf{G}_{(\vartheta)}) = 1 - Q(\partial \mathbf{H}_{(\vartheta)}) = 1 - \mathbb{P}[(X)^\top \vartheta = 0] > 1 - \Delta(X)$ .  $\square$

PROOF OF THEOREM 3.3. Let  $Z = (X, Y)$  be a random variable with values in  $\mathbf{X} \times \mathbb{R}$  whose distribution is  $P$ . We proceed by induction with respect to  $p$ . For  $p = 1$ , the model reduces to the univariate location model and the theorem holds: the sample median has depth not less than  $1/2$ .

Suppose that (11) holds for all dimensions less than  $p$ . If  $\Delta(X) = 1$ , the problem is not identifiable; we can reparametrize it to a lesser dimensional one and use the induction assumption. The lower bound obtained in this way is more stringent than  $1/(p+1)$  and all our fits in are extensions of the lower-dimensional ones; therefore, the inequality (11) holds.

Thus, suppose that  $\Delta(X) < 1$  and assume the contrary to the inequality (11): there is  $\varepsilon > 0$  such that

$$(71) \quad \sup_{\vartheta \in \mathbb{R}^p} d(\vartheta, P) + \varepsilon < \frac{1}{p+1}.$$

We know that Lemma 6.14(i) guarantees the existence of the spheric closure of the corresponding regression probability field, a strongly outer semicontinuous and convex-valued set-valued mapping  $\mathcal{P}(\cdot)$  defined on  $\bar{\mathbb{R}}^p$ , whose range is, by Lemma 6.14(ii), a subset of  $\pi^{-1}(\bar{Q})$  for some fixed probability  $\bar{Q}$  on  $\mathbb{R}^{p-1}$ ; note that  $\pi^{-1}(\bar{Q})$  is a convex, closed (and hence compact, with respect to the weak topology) subset of the space of all probability measures on  $\mathbb{S}^{p-1}$ . Lemma 6.7 says that  $s_{d(\cdot)+\varepsilon}(\cdot)$ , a (single-valued) function assigning every  $Q \in \pi^{-1}(\bar{Q})$  a vector in the unit ball  $\mathbb{D}^p$ , is continuous (in the weak topology). Let  $\mathcal{F}$  be the composition of  $\mathcal{P}(\cdot)$  and  $s_{d(\cdot)+\varepsilon}(\cdot)$ : a set-valued mapping from  $\bar{\mathbb{R}}^p$  to  $\mathbb{D}^p$ . Lemmas 6.14 and 6.8 guarantee that for  $\vartheta \in \partial\mathbb{R}^p$ ,  $\mathcal{F}(\vartheta) \subseteq \mathbf{G}_{((\vartheta))}$ . Since  $\bar{\mathbb{R}}^p$  and  $\partial\mathbb{R}^p$  are topologically equivalent to  $\mathbb{D}^p$  and  $\mathbb{S}^{p-1}$ , we may apply Lemma 6.13, which yields an existence of  $\dot{\vartheta} \in \mathbb{R}^p$  such that  $0 \in \mathcal{F}(\dot{\vartheta})$ ; that is, there is  $Q \in \mathcal{P}(\dot{\vartheta})$  such that  $s_{d(Q)+\varepsilon}(Q) = 0$ . By Proposition 5.16 and Lemma 6.12,  $1/(p+1) \leq d(Q) + \varepsilon \leq d(\vartheta, P) + \varepsilon$ . We obtained a contradiction with (71).  $\square$

#### ACKNOWLEDGEMENTS

I am grateful to Yadolah Dodge, Marc Hallin, and Soros Travel Fund, who made possible my participation at the 1997 Neuchâtel L1 conference, where Peter Rousseeuw kindly gave me a preprint of his joint paper with Mia Hubert; he deserves thanks also for valuable discussions. So do David Eppstein, Xuming He, and particularly Roger Koenker and Steve Portnoy for their suggestions, but not only for that: I am indebted to all last three for the great hospitality during my stay at Urbana-Champaign, as well as to other members of the Department of Statistics. Finally, I would like to appreciate the libraries of the University of Illinois, where I spent some of the most exciting moments during the work on this paper—and also in general.

#### REFERENCES

- ADROVER, J., MARONNA, R. and YOHAI, V. (2000). Relationships between maximum depth projection regression estimates. Preprint.
- AMENTA, N., BERN, M., EPPSTEIN, D. and TENG, S.-H. (2000). Regression depth and center points. *Disc. Comp. Geom.* **23** 305–323.
- BAI, Z.-D. and HE, X. (1999). Asymptotic distributions of the maximal depth estimators for regression and multivariate location. *Ann. Statist.* **27** 1616–1637.
- BALEK, V. and MIZERA, I. (1997). The closeness of the range of a probability on a certain system of random events—an elementary proof. *Bull. Belg. Math. Soc.* **4** 621–624.
- BERN, M. and EPPSTEIN, D. (2000). Multivariate regression depth. In *Proceedings of the 16th Annual Symposium on Computational Geometry* 315–321. ACM Press.
- BICKEL, P. J. and MILLAR, P. W. (1992). Uniform convergence of probability measures on classes of functions. *Statistica Sinica* **2** 1–15.

- BILLINGSLEY, P. (1968). *Convergence of probability measures*. Wiley, New York.
- BILLINGSLEY, P. (1971). *Weak convergence of measures: Applications in probability*. Regional Conference Series in Applied Mathematics No. 5, Society for Industrial and Applied Mathematics, Philadelphia.
- BIRCH, B. J. (1959). On 3N points in a plane. *Proc. Cambridge Philos. Soc.* **55** 289–293.
- BORISOVICH, Y. G., GEL'MAN, B. D., MYSHKIS, A. D. and OBUKHOVSKII, V. V. (1980). Topological methods in the fixed-point theory of multi-valued maps. *Uspekhi Mat. Nauk* **35** 59–126 [Russian Math. Surveys 35 (1980), 65–143].
- BORISOVICH, Y. G., GEL'MAN, B. D., MYSHKIS, A. D. and OBUKHOVSKII, V. V. (1982). Multivalued mappings. In *Mathematical analysis* **19** 127–230. Akad. Nauk SSSR, Vsesoyuz. Inst. Nauchn. i Tekhn. Informatsii, Moscow [in Russian].
- CAPLIN, A. and NALEBUFF, B. (1988). On 64%-majority rule. *Econometrica* **56** 787–814.
- CAPLIN, A. and NALEBUFF, B. (1991a). Aggregation and social choice: a mean voter theorem. *Econometrica* **59** 1–23.
- CAPLIN, A. and NALEBUFF, B. (1991b). Aggregation and imperfect competition: on the existence of equilibrium. *Econometrica* **59** 25–59.
- CARRIZOSA, E. (1996). A characterization of halfspace depth. *J. Multivariate Anal.* **58** 21–26.
- CELLINA, A. (1969). Approximation of set-valued functions and fixed point theorems. *Annali di Matematica Pura ed Applicata* **4** 17–24.
- CHAMBERLIN, E. (1933). *The theory of monopolistic competition*. Harvard University Press, Cambridge.
- DANIELS, H. E. (1954). A distribution-free test for regression parameters. *Ann. Math. Statist.* **25** 499–513.
- DAVIES, P. L. (1993). Aspects of robust linear regression. *Ann. Statist.* **21** 1843–1899.
- DODSON, C. T. J. and PARKER, P. E. (1997). *A user's guide to algebraic topology*. Kluwer, Dordrecht.
- DONOHU, D. L. and GASKO, M. (1992). Breakdown properties of location estimates based on halfspace depth and projected outlyingness. *Ann. Statist.* **20** 1803–1827.
- DOOB, J. L. (1993). *Measure theory*. Springer-Verlag, New York.
- EDELSBRUNNER, H. (1987). *Algorithms in Combinatorial Geometry*. Springer-Verlag, Berlin.
- EDGEWORTH, F. Y. (1888). On a new method of reducing observations relating to several quantities. *Philosophical Magazine* **25** 184–191.
- FRISCH, R. (1966). *Maxima and minima: theory and economic applications*. Reidel, Dordrecht.
- HE, X. and PORTNOY, S. (1998). Asymptotics of the deepest line. In *Applied statistical science III: papers in honor of A. K. Md. E. Saleh* (S. E. Ahmed, M. Ahsanullah and B. K. Sinha, eds.) Nova Science Publications, Commack, N.Y.
- HE, X. and WANG, G. (1997). Convergence of depth contours for multivariate datasets. *Ann. Statist.* **25** 495–504.
- HILL, B. M. (1960). A relationship between Hodges' bivariate sign test and a non-parametric test of Daniels'. *Ann. Math. Statist.* **31** 1190–1192.
- HODGES, J. L., JR (1955). A bivariate sign test. *Ann. Math. Statist.* **26** 523–527.

- HOTELLING, H. (1929). Stability in competition. *Econom. J.* **39** 41–57.
- HUBERT, M., ROUSSEEUW, P. J. and VAN AELST, S. (1999). Similarities between location depth and regression depth. In *Statistics in Genetics and in the Environmental Sciences* (L. Fernholz, S. Morgenthaler and W. Stahel, eds.) 159–172. Birkhäuser Verlag, Basel.
- KLEIN, E. and THOMPSON, A. C. (1984). *Theory of correspondences*. Wiley, New York.
- LIU, R. Y. and SINGH, K. (1993). A quality index based on data depth and multivariate rank tests. *J. Amer. Stat. Assoc.* **88** 252–260.
- LIU, R. Y., PARELIUS, J. M. and SINGH, K. (1999). Multivariate analysis by data depth: descriptive statistics, graphics and inference. *Ann. Statist.* **27** 783–840.
- MIZERA, I. (1998). On depth and deep points: a calculus (abstract). *IMS Bull.* **27** 207.
- NEUMANN, B. H. (1945). On an invariant of plane regions and mass distributions. *J. London Math. Soc.* 226–237.
- NIKAIDO, H. (1968). *Convex structures and economic theory*. Academic Press, New York.
- NOLAN, D. (1992). Asymptotics for multivariate trimming. *Stoch. Proc. Appl.* **42** 157–169.
- NOLAN, D. (1998). On min-max majority and deepest points. Preprint.
- ORTEGA, J. M. and RHEINOLDT, W. C. (1970). *Iterative solution of nonlinear equations in several variables*. Academic Press, New York.
- PONSTEIN, J. (1967). Seven kinds of convexity. *Siam Rev.* **9** 115–119.
- PORTNOY, S. and MIZERA, I. (1999). Comment to regression depth by Rousseeuw and Hubert. *J. Amer. Statist. Assoc.* **94** 417–419.
- RADO, R. (1946). A theorem on general measure. *J. London Math. Soc.* **21** 291–300.
- ROCKAFELLAR, R. T. and WETS, R. J.-B. (1998). *Variational analysis*. Springer-Verlag, Berlin.
- ROTMAN, J. J. (1988). *An introduction to algebraic topology*. Graduate Texts in Mathematics 119, Springer-Verlag, New York.
- ROUSSEEUW, P. J. and HUBERT, M. (1999a). Regression depth. *Journal of the American Statistical Association* **94** 388–402.
- ROUSSEEUW, P. J. and HUBERT, M. (1999b). Depth in an arrangement of hyperplanes. *Discrete and Computational Geometry* **22** 167–176.
- SMALE, S. (1973). Global analysis and economics I: Pareto optimum and a generalization of Morse theory. In *Dynamical Systems (Proc. Sympos., Univ. Bahia, Salvador, 1971)* (M. M. Peixoto, ed.) 531–544. Academic Press, New York.
- SMALE, S. (1975a). Sufficient conditions for an optimum. In *Dynamical Systems—Warwick 1974* (A. Manning, ed.) 287–292. Lecture Notes in Mathematics 468, Springer, Berlin.
- SMALE, S. (1975b). Optimizing several functions. In *Manifolds—Tokyo 1973 (Proc. of the Internat. Conf. on Manifolds and Related Topics in Topology, Tokyo, 1973)* (A. Hattori, ed.) 69–75. University of Tokyo Press, Tokyo.
- TJUR, T. (1980). *Probability based on Radon measures*. Wiley, New York.
- TUKEY, J. W. (1975). Mathematics and the picturing of data. In *Proceedings of the International Congress of Mathematicians (Vancouver, B. C., 1974), Vol. 2* 523–531. Canad. Math. Congress, Quebec.



- VAN AELST, S. and ROUSSEEUW, P. J. (2000). Robustness properties of deepest regression. *J. Multiv. Analysis* **73** 82–106.
- VAN AELST, S., ROUSSEEUW, P. J., HUBERT, M. and STRUYF, A. (2000). The deepest regression method. Preprint.
- WAN, Y. H. (1975). On local Pareto optima. *J. Math. Econom.* **2** 35–42.
- WAN, Y. H. (1978). On the structure and stability of local Pareto optima in a pure exchange economy. *J. Math. Econom.* **5** 255–274.

DEPARTMENT OF PROBABILITY AND STATISTICS  
COMENIUS UNIVERSITY IN BRATISLAVA  
MLYNSKÁ DOLINA, SK-84215 BRATISLAVA  
SLOVAKIA

and

DEPARTMENT OF MATHEMATICAL AND STATISTICAL SCIENCES  
UNIVERSITY OF ALBERTA  
EDMONTON, ALBERTA, T6G2G1  
CANADA

*E-mail address:* mizera@stat.ualberta.ca