

Robust Least Squares Multidimensional Scaling

Jan de Leeuw

November 1, 2024

Combining different loss functions with linear models and minimizing loss with iteratively reweighted least squares (IRLS) has a long history in robust statistics. In this paper we use an IRLS version of the smacof algorithm to minimize various robust multidimensional scaling loss functions. Our results use a general theorem on sharp quadratic majorization of De Leeuw and Lange (2009). We relate this theorem to earlier results in robust statistics, location theory, and sparse recovery. Code in R is included.

Table of contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 5 |
| 2 | Majorization | 7 |
| 3 | Majorizing Strife | 8 |
| 3.1 | Algorithm | 9 |
| 3.2 | Zero Residuals | 9 |
| 3.3 | ℓ_0 loss | 11 |
| 4 | Generalizing Strife | 12 |
| 4.1 | Sharp Quadratic Majorization | 12 |
| 4.2 | Two Support Points | 15 |
| 4.3 | Sufficient Conditions | 18 |
| 4.4 | Algorithm | 21 |

| | | |
|-----------|--|-----------|
| 5 | Power Smoothers | 22 |
| 5.1 | Charbonnier | 22 |
| 5.2 | Generalized Charbonnier | 23 |
| 5.3 | Barron | 24 |
| 6 | Convolution Smoothers | 26 |
| 6.1 | Huber | 26 |
| 6.2 | Gaussian | 28 |
| 7 | A Bouquet of Loss Functions | 29 |
| 7.1 | Andrews | 30 |
| 7.2 | Tukey | 31 |
| 7.3 | Hinich | 31 |
| 7.4 | Cauchy | 32 |
| 7.5 | Welsch | 33 |
| 7.6 | Logistic | 34 |
| 7.7 | Fair | 34 |
| 8 | Examples | 36 |
| 8.1 | Gruijter | 36 |
| 8.1.1 | Least Squares | 37 |
| 8.1.2 | Least Absolute Value | 39 |
| 8.1.3 | Huber | 41 |
| 8.1.4 | Tukey | 43 |
| 8.2 | Rothkopf | 45 |
| 8.2.1 | Least Squares | 45 |
| 8.2.2 | Least Absolute Value | 47 |
| 8.2.3 | Huber | 49 |
| 8.2.4 | Tukey | 51 |
| 9 | Literature | 54 |
| 9.1 | Robust Statistics | 54 |
| 9.2 | Location Analysis | 54 |
| 9.3 | Sparse Recovery | 58 |
| 9.4 | Multivariate Analysis | 58 |
| 10 | Discussion | 60 |
| 10.1 | Bounding the Second Derivative | 60 |
| 10.2 | Fixed Weights | 60 |
| 10.3 | Residual Definition | 61 |
| 10.4 | Robust Nonmetric MDS | 61 |

| | |
|-------------------------------|-----------|
| 10.5 Practicalities | 62 |
| 11 Code | 63 |
| References | 69 |

List of Figures

| | | |
|----|---|----|
| 1 | A Quartic, Majorized at ± 0.5 | 20 |
| 2 | A Quartic, Majorized at ± 1.25 | 20 |
| 3 | Charbonnier Loss | 23 |
| 4 | Generalized Charbonnier Loss | 24 |
| 5 | Barron Loss | 25 |
| 6 | Huber Loss | 27 |
| 7 | Gaussian Convolution Loss | 29 |
| 8 | Andrews Loss | 30 |
| 9 | Tukey Loss | 31 |
| 10 | Hinich Loss | 32 |
| 11 | Cauchy Loss | 33 |
| 12 | Welsch Loss | 34 |
| 13 | Logistic Loss | 35 |
| 14 | Fair Loss | 35 |
| 15 | Gruijter Configuration Least Squares | 37 |
| 16 | Gruijter Shepard Plot Least Squares | 38 |
| 17 | Gruijter Histogram Least Squares Residuals | 38 |
| 18 | Gruijter Configuration Least Absolute Value | 39 |
| 19 | Gruijter Shepard Plot Least Absolute Value | 40 |
| 20 | Gruijter Histogram Least Absolute Value Residuals | 41 |
| 21 | Gruijter Configuration Huber $c = 1$ | 42 |
| 22 | Gruijter Shepard Plot Huber $c = 1$ | 42 |
| 23 | Gruijter Histogram Huber Residuals | 43 |
| 24 | Gruijter Configuration Tukey $c = 2$ | 44 |
| 25 | Gruijter Shepard Plot Tukey $c = 2$ | 44 |
| 26 | Gruijter Histogram Tukey Residuals | 45 |
| 27 | Rothkopf Configuration Least Squares | 46 |
| 28 | Rothkopf Shepard Plot Least Squares | 46 |
| 29 | Rothkopf Histogram Least Squares Residuals | 47 |
| 30 | Rothkopf Configuration Least Absolute Value | 48 |
| 31 | Rothkopf Shepard Plot Least Absolute Value | 48 |
| 32 | Rothkopf Histogram Least Absolute Value Residuals | 49 |
| 33 | Rothkopf Configuration Huber $c = 1$ | 50 |
| 34 | Rothkopf Shepard Plot Huber $c = 1$ | 50 |
| 35 | Rothkopf Histogram Huber Residuals | 51 |
| 36 | Rothkopf Configuration Tukey $c = 1$ | 52 |
| 37 | Rothkopf Shepard Plot Tukey $c = 1$ | 52 |
| 38 | Rothkopf Histogram Tukey Residuals | 53 |

1 Introduction

The title of this paper is somewhat surprising. Least squares estimation is typically not robust, it is sensitive to outliers and pays too much attention to minimizing the largest residuals. What we mean by robust least squares multidimensional scaling (MDS), however, is to use the smacof machinery designed to minimize least squares loss functions of the form

$$\sigma(X) := \sum \omega_k (\delta_k - d_k(X))^2, \quad (1)$$

to minimize robust loss functions. In Equation 1 the ω_k are positive *weights*, the δ_k are positive *dissimilarities*, and the $d_k(X)$ are Euclidean *distances* between two of the rows of the $n \times p$ *configuration* matrix X .

Some terminological conventions we use throughout the paper. By *positive* we mean smaller than or equal to zero. Smaller than zero is *strictly positive*. Some for *negative* and *strictly negative*, for *increasing* and *strictly increasing*, and for *decreasing* and *strictly decreasing*.

The prototypical robust loss function is least absolute value loss

$$\sigma(X) := \sum \omega_k |\delta_k - d_k(X)|, \quad (2)$$

which we will call *strife*, since the names *stress*, *sstress*, and *strain* are already taken. Note we are overloading the symbol σ , because we will use it for all of the loss functions in this paper.

Strife is not differentiable at configurations X for which there is at least one k for which either $d_k(X) = \delta_k$ or $d_k(X) = 0$ (or both). This lack of differentiability complicates the minimization problem. Moreover experience with one-dimensional and city block MDS suggests that having areas where the loss function is not differentiable leads to (many) additional local minima.

In this paper we will discuss (and implement) various variations of strife from Equation 2. They can be interpreted in two different ways. On the one hand they are smoothers of the absolute value function, and consequently of strife. We want to eliminate the problems with differentiability, at least the ones caused by $\delta_k = d_k(X)$. If this is our main goal, then we want to choose the smoother in such a way that it is close to the absolute value function. This is similar to the distance smoothing used by Pliner (1996) and Groenen, Heiser, and Meulman (1999) in the global minimization of stress from Equation 1, except that we do not smooth the distance function but the strife residual.

On the other hand our modified loss functions can be interpreted as more robust versions of the least squares loss function, and consequently of stress. Our goal then is to combine the robustness of the absolute value function with the efficiency and computational ease of least

squares. If robustness is our main goal then there is no reason to stay close to the absolute value function.

Our robust or smooth loss functions are all of the form

$$\sigma(X) := \sum \omega_k f(r_k(X)), \quad (3)$$

where $r_k(X)$ is the *residual*, i.e.

$$r_k(X) := \delta_k - d_k(X). \quad (4)$$

The function f is assumed to be even (i.e. symmetric around zero) and attains its minimum, which is equal to zero, at zero. For now, note that loss in Equation 1 is the special case with $f(x) = x^2$ and loss in Equation 2 is the special case with $f(x) = |x|$.

2 Majorization

Our loss functions will be minimized by using majorization algorithms (these days more commonly known as MM algorithms). This paper discusses a general way to construct majorization algorithms for robust loss functions.

For completeness we give a short introduction to majorization, without aiming for maximum generality. For the general theory of majorization algorithms we refer to their introduction in De Leeuw (1994) and to the excellent and comprehensive book by Lange (2016). As mentioned in Section 1 the robust loss functions in Equation 3 are real-valued functions of a single real variable, defined on the whole real line, they are even, and they attain a minimum equal to zero at zero. Thus they are all symmetric bowls anchored at zero.

Definition 2.1. A function g majorizes a function f at a point y if $g(x) \geq f(x)$ for all x and $g(y) = f(y)$. The point y is a *support point* of the majorization. Majorization of f at y is *strict* if $g(x) > f(x)$ for all $x \neq y$.

Majorizers have at least one support point, but they can have many. The function $g(x) = x^2 + \sin^2(x)$ majorizes $f(x) = x^2$ at every integer multiple of π , and strictly majorizes f at none of its support points. Also any function f majorizes itself at all points of the real line.

Definition 2.2. If \mathfrak{H} is a family of functions that all majorize f at y then $h \in \mathfrak{H}$ is a *sharp majorization* in \mathfrak{H} if $h(x) \leq g(x)$ for all $g \in \mathfrak{H}$. The sharp majorization, if it exists, is by definition unique.

Theorem 2.1. Suppose g majorizes f at y .

- $f(y) = g(y)$.
- If f and g are differentiable at y then $f'(y) = g'(y)$.
- If f and g are twice-differentiable at y then $f''(y) \leq g''(y)$.
- If $f''(y) < g''(y)$ then g strictly majorizes f in a neighborhood of y .

Proof. $h = g - f$ is negative and has a minimum equal to zero at y . Thus the derivative of h vanishes at y and the second derivative is negative at y . If the second derivative is strictly negative then we use the sufficient condition for a local minimum. \square

A majorization algorithm to minimize f is iterative. We update $x^{(\nu)}$, the approximation of the minimizer in iteration ν , by

$$x^{(\nu+1)} \in \underset{x}{\operatorname{argmin}} g_\nu(x), \quad (5)$$

where g_ν majorizes f at $x^{(\nu)}$. If $x^{(\nu)} \in \operatorname{argmin}_x g_\nu(x)$ we stop. Otherwise we find a new majorizer $g_{\nu+1}$, which majorizes f at $x^{(\nu+1)}$, and start a new iteration.

Convergence of majorization algorithms follows from the *sandwich inequality*

$$f(x^{(\nu+1)}) \leq g_\nu(x^{(\nu+1)}) \leq g_\nu(x^{(\nu)}) = f(x^{(\nu)}). \quad (6)$$

This we the algorithm produces a decreasing sequence of loss function values, which converges if loss is bounded below. In Equation 6 the first inequality (from the left) follows from majorization. If the majorization is strict, then so is the inequality. The second inequality follows from minimization of g . It is strict if g has a unique minimizer, for example if g is strictly convex. The final equality in Equation 6 comes from majorization at $x^{(\nu)}$.

If $x^{(\nu)}$ minimizes g_ν we stop and we have finite convergence. If the algorithm generates an infinite sequence, which is the usual case, the second inequality in Equation 6 is always strict, and the algorithm generates a strictly decreasing sequence of loss function values. Note that for the validity of the sandwich inequality it suffices to decrease g_ν in every iteration, and not necessarily to minimize it. Different ways to decrease g_ν correspond with different step-size procedures in gradient methods.

Of course convergence of loss function values does not guarantee convergence of the $x^{(k)}$. For the additional continuity, compactness, and identification conditions that are needed we refer to the majorization and MM literature.

3 Majorizing Strife

The idea of minimizing a least absolute value (LAV) to obtain parameter estimates dates back to the work of Boskovitch in the middle of the eighteenth century. Until fairly recently it has been applied mainly to fit location parameters and more general linear models.

The pioneering work in strife minimization using smacof is Heiser (1988), which builds on earlier work of Heiser (1987). It is based on a creative use of the Arithmetic Mean-Geometric Mean (AM/GM) inequality to find a majorizer of the absolute value function.

The AM/GM inequality says that for all non-negative x and y we have

$$|x||y| = \sqrt{x^2 y^2} \leq \frac{1}{2}(x^2 + y^2), \quad (7)$$

with equality if and only if $|x| = |y|$. If $|y| > 0$ we can write Equation 7 as

$$|x| \leq \frac{1}{2} \frac{1}{|y|} (x^2 + y^2), \quad (8)$$

and this provides a quadratic majorization of $|x|$ at y . There is no quadratic majorization of $|x|$ at $y = 0$, which is a problem we will have to deal with at some point.

Using the majorization Equation 8, and assuming $\delta_k \neq d_k(Y)$ for all k , we define

$$\omega_k(X) := \omega_k \frac{1}{r_k(X)}, \quad (9)$$

and, for a fixed Y ,

$$\eta(X) := \frac{1}{2} \sum \omega_k(Y) (r_k^2(X) + r_k^2(Y)). \quad (10)$$

Now $\sigma(X) \leq \eta(X)$ for all X and $\sigma(Y) = \eta(Y)$, and thus η majorizes σ at Y .

3.1 Algorithm

Reweighted smacof to minimize strife computes $X^{(\nu+1)}$ by minimizing or decreasing

$$\sum \omega_k(X^{(\nu)}) (\delta_k - d_k(X^{(\nu)}))^2, \quad (11)$$

using a standard smacof step. It then computes the new weights $\omega_k(X^{(\nu+1)})$ from Equation 9 and uses them in the next smacof step to update $X^{(\nu+1)}$. And so on, until convergence.

A straightforward variation of the algorithm does a number of smacof steps between upgrading the weights. This still leads to a monotone, and thus convergent, algorithm. How many smacof steps we have to take in these inner iterations is something that needs further study. It is likely to depend on the fit of the data, on the shape of the function near the local minimum, and on how far the current iteration is from the local minimum.

3.2 Zero Residuals

It may happen that for some k we have $d_k(X^{(\nu)}) = \delta_k$ while iterating. There have been various proposals to deal with such an unfortunate event, and we will discuss some of them below. Even more importantly we will see that that the minimizer of the absolute value loss usually satisfies $d_k(X) = \delta_k$ for quite a few elements, which means that near convergence the algorithm may become unstable because the weights from Equation 9 become very large.

A large number of somewhat ad-hoc solutions have been proposed to deal with the problem of zero residuals, both in location analysis and in the statistical literature. We tend to agree with the assessment of Aftab and Hartley (2015).

.. attempts to analyze this difficulty [caused by infinite weights of IRLS for the ℓ_p -loss] have a long history of proofs and counterexamples to incorrect claims.

Schlossmacher (1973) is the first discussion of the majorization method in the statistical literature (for LAV linear regression). His proposal is to simply set a weight equal to zero if the corresponding residual is less than some small positive value ϵ . A similar approach, also used in location analysis, is to cap the weights at some large positive value. In Heiser (1988) all residuals smaller than this epsilon get a weight equal to the weighted average of all these small residuals. Phillips (2002) assumes double-exponential errors in LAV regression and then concludes that the EM algorithm gives the majorization method we have discussed. He uses (??) throughout if all residuals are larger than ϵ . If one or more residuals are smaller than epsilon then the weight for those residuals is set equal to one, while for the remaining residuals the weight is set to epsilon divided by the absolute value of the residual. Often we get the assurance in these papers that the problem is not really important in practice, because it is very rare, and by just wiggling we will get to the unique solution anyway. But both in location analysis and in LAV regression the loss function is convex, however, which guarantees a unique minimum. This is certainly not the case in robust MDS. In this paper we try to follow a more systematic approach that uses smooth parametric approximations to the absolute value function, where the parameter can be used to make the approximation as precise as necessary.

In the case of stress the directional derivatives can be used to prove that if $\omega_k \delta_k > 0$ for all k then stress is differentiable at each local minimum (De Leeuw (1984)). For strife to be differentiable we would have to prove that at a local minimum both $d_k(X) > 0$ and $(d_k(X) - \delta_k) \neq 0$ for all k with $\omega_k > 0$. But this is impossible by the following argument.

In the one-dimensional case we can partition \mathbb{R}^n into $n!$ polyhedral convex cones corresponding with the permutations of x . Within each cone the distances are a linear function of x . Each cone can be partitioned by intersecting it with the polyhedra defined by the linear inequalities $\delta_k - d_k(x) \geq 0$ or $\delta_k - d_k(x) \leq 0$. Some of these intersections can and will obviously be empty. Within each of these non-empty polyhedral regions strife is a linear function of x . Thus it attains its minimum for the region at a vertex, which is a solution for which some distances are zero and/or some residuals are zero. There can be no minima, local or global, in the interior of one of these polyhedral regions. We have thus shown that in one dimension strife is not differentiable at a local minimum, and that there is presumably a large number of them. Even for moderate n the number of regions is of course too large to actually compute or draw.

In the multidimensional case linearity goes out the window. The set of configurations $d_k(X) = \delta_k$ is an ellipsoid and $d_k(X) = 0$ defines a hyperplane. Strife is not differentiable at all intersections of these ellipsoids and hyperplanes. The partitioning of \mathbb{R}^n by these ellipsoids and hyperplanes is not simple to describe. It has convex and non-convex cells, and within each cell strife is the difference of two weighted sums of distances. Anything can happen.

3.3 ℓ_0 loss

A somewhat extreme special case of Equation 3 has

$$f(x) = \begin{cases} 0 & \text{if } x = 0, \\ 1 & \text{otherwise.} \end{cases} \quad (12)$$

This is ℓ_0 loss. Minimizing ℓ_0 loss means maximizing the number of cases with perfect fit, i.e. with $\delta_k = d_k(X)$. The reason we mention it here is that the work of Donoho and Elad (2003) and Candes and Tao (2005) suggests that the minimizer of ℓ_1 loss, i.e. absolute value loss, gives a good approximation to the minimizer of ℓ_0 loss, at least in a number of special cases. In MDS we do not have linearity or convexity, but nevertheless the theoretical results in simpler cases are suggestive. We have seen that at least in the one-dimensional MDS case a number of residuals will indeed be zero at the optimum LAV solution.

There is an excellent review of the use of ℓ_1 in various sparse recovery fields in Candes, Wakin, and Boyd (2008). In that paper they also propose an iteratively reweighted LAV algorithm, which solves ℓ_1 problems between weight updates. Maybe because of that they go so far as calling ℓ_1 “the modern least squares”. But let’s not get carried away, in actual ease and frequency of use ℓ_1 still has a long way to go if it wants to replace ℓ_2 .

4 Generalizing Strife

We have seen that Heiser (1988) applied majorization to minimize strife, using the AM/GM inequality. We now generalize this approach so that it can deal with other robust loss functions. A great number of different loss functions will be discussed. My intention is certainly not to confuse the reader and potential user by presenting a large number of alternatives with rather limited information. We show all these loss functions as examples of a general principle of algorithm construction and as examples of loss functions that have been used in statistics, location analysis, image analysis, and engineering over the years. They are all implemented in the function `smacofRobust()`, written in R (R Core Team (2024)), and listed in the appendix of this paper.

4.1 Sharp Quadratic Majorization

The AM/GM inequality was used in Section 3 to construct a quadratic majorization of strife. In this paper we are specifically interested in sharp quadratic majorization, in which \mathfrak{H} is the set of all (not necessarily convex) quadratics that majorize f at y . This case has been studied in detail (in the case of real-valued functions on the line) in De Leeuw and Lange (2009), and much of this section is taken from their paper. We added some minor extensions and reformulations.

For the loss functions we study in this paper there are two problems that have to be solved. First, we want a general procedure to construct quadratic majorizers. Second, we want to show that some of our majorizers are sharp.

If f is differentiable at y , then all quadratics g that majorize f at y are of the form

$$g(x) := f(y) + f'(y)(x - y) + \frac{1}{2}a(x - y)^2 \quad (13)$$

for some $a = g''(x)$, not necessarily positive. If f is twice differentiable at y then $g''(y) \geq f''(y)$ by Theorem 2.1, and thus we have the necessary condition $a \geq f''(y)$. Note that not all functions have quadratic majorizations. If f is a non-trivial cubic and g is quadratic, then $h = g - f$ is a non-trivial cubic, and consequently we cannot have $h \geq 0$ on the whole real line.

We now look more closely at a in Equation 13, initially concentrating on necessary conditions for majorization. For $x \neq y$ define

$$\alpha(x) := \frac{f(x) - f(y) - f'(y)(x - y)}{\frac{1}{2}(x - y)^2}. \quad (14)$$

Of course α is a different function of x for each y , but since we are dealing with majorization at one fixed y we suppress this dependence.

If f is two times differentiable at y then, by the definition of the second derivative,

$$\lim_{x \rightarrow y} \alpha(x) = f''(y), \quad (15)$$

and thus can we define $\alpha(y) = f''(y)$ to make α continuous at y .

If f is convex, then α is the ratio of two positive convex functions of x , and is thus positive. If f is concave then α is negative. If f is two times differentiable then there is a z between x and y such that $\alpha(x) = f''(z)$. Thus if $f''(x) \leq K$ for all x , then $\alpha(x) \leq K$ as well.

Quadratic majorization of f by g from Equation 13 at y is equivalent to $\alpha(x) \leq a$ for all x . Thus g majorizes f at y if and only if α is bounded above by a . If α is unbounded above there is no quadratic majorizer at y . We can also define

$$a_+ := \sup_x \alpha(x), \quad (16)$$

and say that sharp quadratic majorization at y is possible if $a_+ < +\infty$. We have majorization at y by g if $a \geq a_+$, we have sharp majorization if $a = a_+$. It follows that sharp quadratic majorizations exist if f'' is bounded above.

We illustrate these concepts by applying them to low-degree polynomials. First a cubic. Expand the cubic around y as

$$f(x) = f(y) + f'(y)(x - y) + \frac{1}{2}f''(y)(x - y)^2 + \frac{1}{6}f'''(y)(x - y)^3. \quad (17)$$

Thus

$$\alpha(x) = f''(y) + \frac{1}{3}f'''(y)(x - y). \quad (18)$$

Since α is always unbounded above, no quadratic majorizer exists at any y .

Now apply the same reasoning to a non-trivial quartic. We find

$$\alpha(x) = f''(y) + \frac{1}{3}f'''(y)(x - y) + \frac{1}{12}f^{iv}(y)(x - y)^2, \quad (19)$$

a quadratic in x . Of course for a quartic $f^{iv}(y)$ is the same for every y and we may as well write f^{iv} . If f^{iv} is strictly positive the quadratic in Equation 19 is unbounded above, and no quadratic majorization exists at any y . If f^{iv} is strictly negative then α has a maximum at $x = y - 2f'''(y)/f^{iv}$, and a sharp quadratic majorization exists at any y .

We can get some more information about a_+ by differentiating α .

Theorem 4.1. *Suppose there is an x where α is two times differentiable and $a_+ = \alpha(x)$. Then*

$$\frac{f(x) - f(y)}{x - y} = \frac{1}{2}(f'(x) + f'(y)), \quad (20a)$$

$$f''(x) \leq \frac{f'(x) - f'(y)}{x - y}. \quad (20b)$$

Moreover

$$a_+ = \frac{f'(x) - f'(y)}{x - y}. \quad (20c)$$

If (20a) and the inequality in (20b) is strict then α has a local maximum at x .

Proof. After some manipulation (20a) and (20b) are the necessary conditions $\alpha'(x) = 0$ and $\alpha''(x) \leq 0$ for a local maximum. If $\alpha'(x) = 0$ and $\alpha''(x) > 0$ the conditions are sufficient. If x satisfies (20a) then substitution in Equation 14 gives (20c). \square

Note that we have not shown that α always attains its maximum. De Leeuw and Lange (2009) give the example of the differentiable function

$$f(x) = \begin{cases} x^2 & \text{if } x \leq 1, \\ 2x - 1 & \text{otherwise,} \end{cases} \quad (21)$$

which has $\alpha(x) = 0$ for $x > 1$ and $\alpha(x) < 2$ for $x \leq 1$, so that $a_+ = \sup_{x \leq 1} \alpha(x) = 2$ and the maximum does not exist.

Also, the conditions of Theorem 4.1 cannot possibly show that α has a *global* maximum at x , and that consequently g of Equation 13 with a given by (20c) is a sharp quadratic majorizer.

Differentiation of α gives the following result.

Corollary 4.1. *Suppose α is differentiable. Then it is strictly increasing if and only if for all x*

$$\frac{f(x) - f(y)}{x - y} < \frac{1}{2}(f'(x) + f'(y)), \quad (22)$$

and strictly decreasing if and only if for all x

$$\frac{f(x) - f(y)}{x - y} < \frac{1}{2}(f'(x) + f'(y)). \quad (23)$$

Proof. These are the conditions $\alpha'(x) > 0$ and $\alpha'(x) < 0$ for all x . \square

If Equation 22 or Equation 23 is true then α does not have a maximum, and the supremum (possibly infinite) is the limit of α to $-\infty$ or $+\infty$.

Because of our robust loss functions we are especially interested in the case that f is even. Setting $x = -y$ makes both sides of (20a) equal to zero. Thus α has a stationary point at $x = -y$. If in addition

$$f''(y) < \frac{f'(y)}{y} \quad (24)$$

then α has a local maximum at $x = -y$.

Theorem 4.2. *Suppose f is even and twice-differentiable. If $f'(x)/x$ is decreasing on the positive real line then α has a local maximum at $x = -y$.*

Proof. If $f'(x)/x$ is strictly decreasing its derivative is strictly negative. Thus $xf''(x) - f'(x) < 0$ for $x > 0$. For an even function $f'(x)/x$ is strictly decreasing on the positive real line if and only if it is strictly increasing on the negative real line. Thus $xf''(x) - f'(x) > 0$ for $x < 0$. In both cases Equation 24 follows. \square

4.2 Two Support Points

We can say more if it is known that the quadratic majorization has more than one support point.

Definition 4.1. Functions f and g have the *two-point property* at y and z if $f(y) = g(y)$ and $f'(y) = f'(z)$.

If g majorizes f at y and z then by Theorem 2.1 g and f have the two-point property at y and z .

Lemma 4.1. *Suppose f is differentiable at y and z . Then there is a quadratic g such that f and g have the two-point property at y and z if and only if*

$$\frac{f(y) - f(z)}{y - z} = \frac{1}{2}(f'(y) + f'(z)). \quad (25)$$

In that case $g(x) = c + bx + \frac{1}{2}ax^2$, with

$$a = \frac{f'(y) - f'(z)}{y - z}, \quad (26)$$

$$b = \frac{yf'(z) - zf'(y)}{y - z}, \quad (27)$$

and

$$c = f(y) - by - \frac{1}{2}ay^2 = f(z) - bz - \frac{1}{2}az^2. \quad (28)$$

Proof. This is elementary computation, but we write it down anyway. Function f and the quadratic $g = c + bx + \frac{1}{2}ax^2$ have the two-point property at y and z if and only if a , b , and c solve the four linear equations.

$$c + by + \frac{1}{2}ay^2 = f(y), \quad (29a)$$

$$c + bz + \frac{1}{2}az^2 = f(z), \quad (29b)$$

$$b + ay = f'(y), \quad (29c)$$

$$b + az = f'(z). \quad (29d)$$

Solve (29c) and (29d) gives the solution for a and b in Equation 26 and Equation 27. We can substitute these a and b in (29a) and (29b) to get two solutions for c . The system (29a)-(29d) is consistent, and has a unique solution, if and only if these two values for c are the same, which is the case if and only if Equation 25 is true. \square

Corollary 4.2. *Suppose the quadratic g majorizes f at y and at $z \neq y$. Then*

$$g(x) = f(y) + f'(y)(x - y) + \frac{1}{2} \frac{f'(z) - f'(y)}{z - y} (x - y)^2 = \quad (30)$$

$$= f(z) + f'(z)(x - z) + \frac{1}{2} \frac{f'(z) - f'(y)}{z - y} (x - z)^2. \quad (31)$$

Proof. From Theorem 2.1 and Lemma 4.1. \square

Again, we have not shown that g with a from Corollary 4.2 majorizes f at y and z . Only the reverse implication, which is that if g majorizes f at y and z then g is uniquely determined by Corollary 4.2. In practice, even if one knows the support points y and z , one still has to prove majorization. This is precisely how Heiser (1988), Verboon and Heiser (1994), and Groenen, Giaquinto, and Kiers (2003) establish their majorizations.

In his master's thesis Van Ruitenburg (2005) takes us one step further down the necessary conditions road by showing that quadratic majorization at two points implies strict quadratic majorization. To present his main result we need two lemmas.

Lemma 4.2. *If different quadratics g and h majorize f at y then either g strictly majorizes h or h strictly majorizes g .*

Proof. We have

$$g(x) = f(y) + f'(y)(x - y) + \frac{1}{2}a_1(x - y)^2, \quad (32a)$$

$$h(x) = f(y) + f'(y)(x - y) + \frac{1}{2}a_2(x - y)^2. \quad (32b)$$

Thus $g(x) - h(x) = \frac{1}{2}(a_1 - a_2)(x - y)^2$, which is either strictly positive or strictly negative for $x \neq y$. \square

Lemma 4.3. *Suppose quadratics g and $h \neq g$ majorize f at y . Suppose, in addition, that g majorizes f at $z \neq y$. Then h strictly majorizes g at y .*

Proof. By Lemma 4.2 one of g and h has to strictly majorize the other. If g strictly majorizes h then $h(z) < g(z) = f(z)$ and thus h does not majorize f , contrary to assumption. It follows that h strictly majorizes g at y . \square

Theorem 4.3. *If the quadratic g majorizes f at y and at $z \neq y$, then g is a sharp majorizer of f at y .*

Proof. Directly from Lemma 4.3. \square

Again our results simplify if the function f is even.

Theorem 4.4. *If f is even and the quadratic g majorizes f at y and $-y$, where $y \neq 0$, then g is the even quadratic given by*

$$g(x) = f(y) + \frac{1}{2} \frac{f'(y)}{y} (x^2 - y^2). \quad (33)$$

Moreover g is the sharp quadratic majorization of f at y and $-y$.

Proof. For even f Equation 26 becomes

$$a = \frac{f'(y)}{y}, \quad (34)$$

while b from [eq-bfunc](#) becomes zero. Moreover because g majorizes f at y

$$g(x) = f(y) + f'(y)(x - y) + \frac{1}{2} \frac{f''(y)}{y} (x - y)^2, \quad (35a)$$

and because g majorizes f at $-y$

$$g(x) = f(y) - f'(y)(x + y) + \frac{1}{2} \frac{f''(y)}{y} (x + y)^2. \quad (35b)$$

Averaging the two equations [\(35a\)](#) and [\(35b\)](#) for g , and simplifying, gives the required result. That the majorization is sharp follows from [Theorem 4.3](#). \square

4.3 Sufficient Conditions

[De Leeuw and Lange \(2009\)](#) gives a way to construct quadratic majorizers of a differentiable function on the real line. We give a slightly more general version of their basic theorem that provides a convenient way to deal with loss functions that are not differentiable everywhere. Our proofs are different.

Theorem 4.5. *Suppose there is a function h , concave on \mathbb{R}^+ , such that $f(x) = h(x^2)$. Then*

$$g(x) = f(y) + \frac{1}{2} \frac{f''(y)}{y} (x^2 - y^2). \quad (36)$$

is a sharp quadratic majorizer of f at y

Proof. By concavity for all $u \in \mathbb{R}^+$ and for all $v \in \mathbb{R}^+$ for which $h'(v)$ exists

$$h(u) \leq h(v) + h'(v)(u - v). \quad (37)$$

Substituting $u = x^2$ and $v = y^2$ gives the quadratic majorization

$$f(x) \leq f(y) + h'(y^2)(x^2 - y^2), \quad (38)$$

which is true for all real x and y , provided h is differentiable at y^2 . If it is then using $f'(y) = 2yh'(y^2)$ in [Equation 38](#) gives

$$f(x) \leq f(y) + \frac{1}{2} \frac{f''(y)}{y} (x^2 - y^2). \quad (39)$$

\square

If we want to get rid of the differentiability assumption altogether we can use the majorization

$$f(x) \leq f(y) + a(y)(x^2 - y^2) \quad (40)$$

where $a(y)$ is any element of $\partial h(y^2)$, the superdifferential of h at y^2 . The superdifferential, which is the analog of the subdifferential for concave functions (Border (2018)), is non-empty on \mathbb{R}^+ , except possibly at zero. Note, by the way, that $f(x) = h(x^2)$ implies that f is even.

Theorem 4.6. $f(x) = h(x^2)$ for concave h on \mathbb{R}^+ if and only if $f(\sqrt{x})$ is concave in x , and, for differentiable f , if and only if $f'(x)/x$ is decreasing on \mathbb{R}^+ .

Proof. The first part of the theorem merely states the obvious. Second part ... □

Some quick examples of the use of Theorem 4.5 before we go to the robust loss functions.

Example 4.1. Suppose h is the square root, so that $f(x) = \sqrt{|x^2|} = |x|$ and for $y \neq 0$ we have $f'(y)/y = 1/|y|$. From Theorem 4.5 we have the quadratic majorization

$$|x| \leq |y| + \frac{1}{2} \frac{1}{|y|} (x^2 - y^2) = \frac{1}{2} \frac{1}{|y|} (x^2 + y^2),$$

which is the AM/GM inequality. So AM/GM majorization of the absolute value is sharp.

Example 4.2. Suppose $h(x) = \frac{1}{2}x(1 - \frac{1}{2}x)$. Thus f is the even quartic

$$f(x) = \frac{1}{2}x^2(1 - \frac{1}{2}x^2).$$

We have $f'(y)/y = 1 - y^2$, and thus the majorizer

$$g(x) = f(y) + \frac{1}{2}(1 - y^2)(x^2 - y^2).$$

For all $y^2 < 1$ the majorizer g has its minimum at zero. If $y^2 > 1$ the majorizer is a concave quadratic, which has no minimum and is unbounded below. If $y^2 = 1$ the majorizer is the horizontal line $x = 0.25$ and the majorization method stops at a local maximum. This is illustrated in Figure 1 and Figure 2.

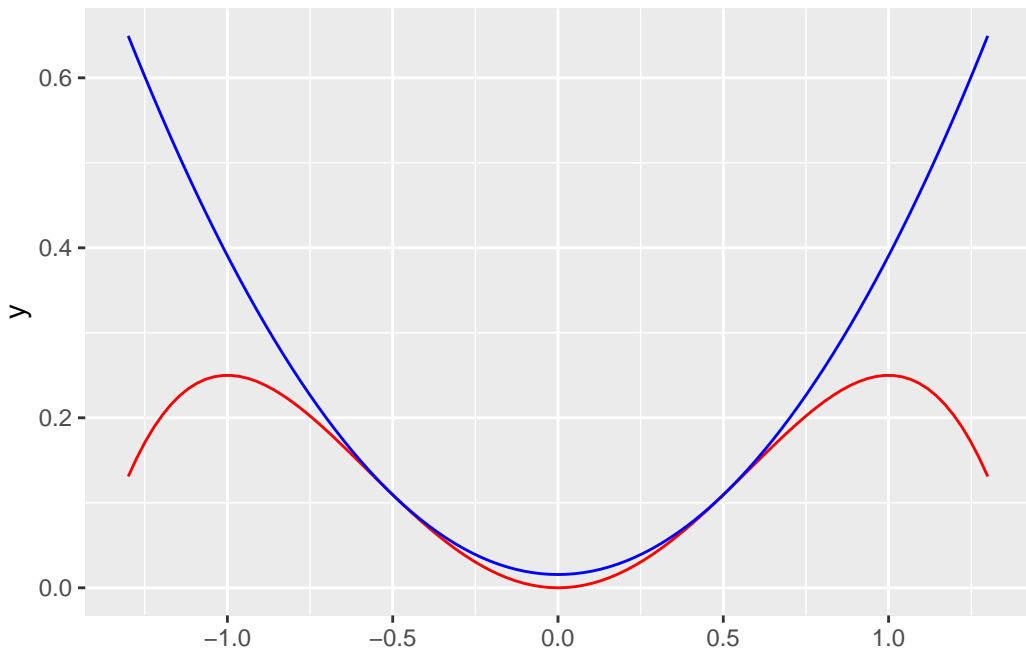


Figure 1: A Quartic, Majorized at ± 0.5

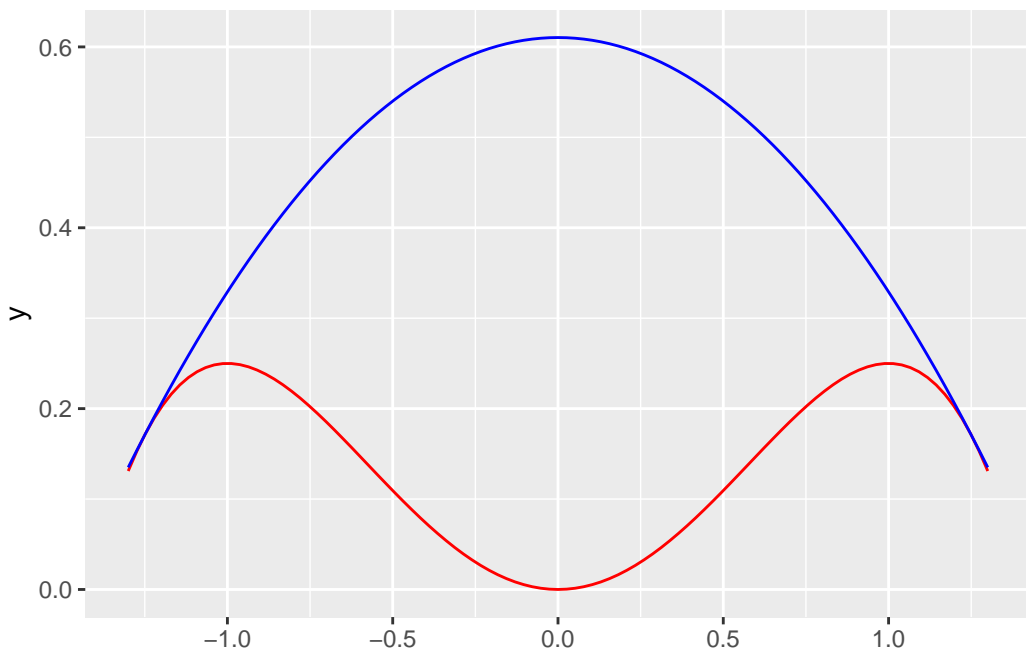


Figure 2: A Quartic, Majorized at ± 1.25

4.4 Algorithm

We now apply Theorem 4.5 to functions of the form

$$\sigma(X) := \sum \omega_k f(\delta_k - d_k(X)), \quad (41)$$

where f satisfies the conditions in the theorem. If

$$\eta(X) := \sum \omega_k \frac{f'(\delta_k - d_k(Y))}{2(\delta_k - d_k(Y))} \{(\delta_k - d_k(X))^2 - (\delta_k - d_k(Y))^2\} + f(\delta_k - d_k(Y)), \quad (42)$$

then η is a sharp quadratic majorization at Y .

In iteration ν the robust smacof algorithm does a smacof step towards minimization of η over X . We can ignore the parts of Equation 42 that only depend on Y , and minimize

$$\sum \omega_k(X^{(\nu)}) (\delta_k - d_k(X))^2, \quad (43)$$

with

$$\omega_k(X^{(\nu)}) := \omega_k \frac{f'(\delta_k - d_k(X^{(\nu)}))}{2(\delta_k - d_k(Y))}. \quad (44)$$

We then recompute the weights $\omega_k(X^{(\nu+1)})$ and go to the smacof step again. This can be thought of as iteratively reweighted least squares (IRLS), and also as nested majorization, with the smacof majorization based on the Cauchy-Schwartz inequality within the sharp quadratic majorization of the loss function based on Theorem 4.5.

5 Power Smoothers

We first discuss a class of smoothers of the absolute value function that maintain most of its structure. They have a shift parameter c that takes care of the non-differentiability. Although different smoothers have different scales and interpretations for c , we will use the same symbol throughout. Also some smoothers have a power parameter q that determines the shape of the loss function bowl.

5.1 Charbonnier

The first, and perhaps most obvious, choice for smoothing the absolute value function is

$$f(x) = \sqrt{x^2 + c^2}. \quad (45)$$

In the engineering literature Equation 45 is known as Charbonnier loss, after Charbonnier et al. (1994), who were possibly the first researchers to use it in image restoration. Ramirez et al. (2014) count the number of computer operations and conclude that Equation 45 is also the “most computationally efficient smooth approximation to $|x|$ ”.

For $c > 0$ we have $f_c(x) > |x|$. If $c \rightarrow 0$ then $f_c(x)$ decreases monotonically to $|x|$. Also $\max_x |f_c(x) - |x|| = c$ attained at $x = 0$, which implies uniform convergence of f_c to $|x|$.

By l’Hôpital

$$\lim_{x \rightarrow 0} \frac{\sqrt{x^2 + c^2} - c}{\frac{1}{2}x^2} = 1. \quad (46a)$$

Of course also

$$\lim_{x \rightarrow \infty} \frac{\sqrt{x^2 + c^2}}{|x|} = 1 \quad (46b)$$

and

$$\lim_{x \rightarrow \pm\infty} \sqrt{x^2 + c^2} - |x| = 0 \quad (46c)$$

Thus if x is much smaller than c then loss is approximately a quadratic in x , and if x is much larger than c then loss is approximately the absolute value.

Loss function Equation 45 is infinitely many times differentiable. Its first derivative is

$$f'_c(x) = \frac{1}{\sqrt{x^2 + c^2}}x, \quad (47)$$

which converges, again in the sup-norm and uniformly, to the sign function if $c \rightarrow 0$. The IRLS weights are

$$w_c(x) = \frac{1}{\sqrt{x^2 + c^2}} \quad (48)$$

which is clearly a decreasing function of x on \mathbb{R}^+ .

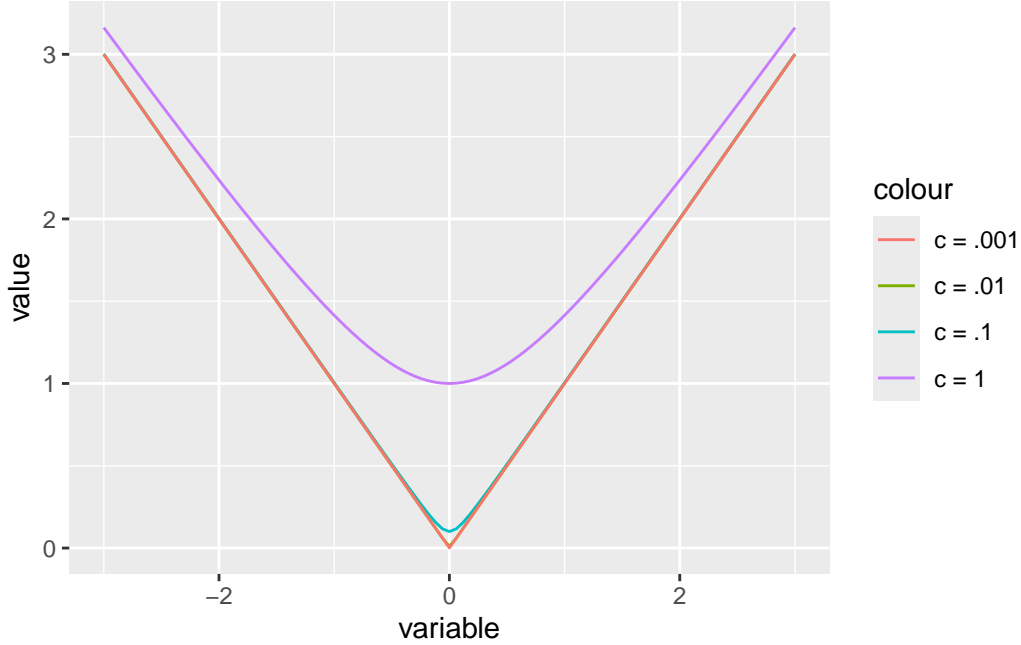


Figure 3: Charbonnier Loss

5.2 Generalized Charbonnier

The loss function $(x^2 + c^2)^{\frac{1}{2}}$ smoothes $|x|$. In the same way generalized Charbonnier loss smoothes ℓ_p loss $|x|^q$. We have a two-parameter family of loss functions in this case.

$$f_{c,q}(x) := (x^2 + c^2)^{\frac{1}{2}q} \quad (49)$$

$$w_{c,q}(x) = q(x^2 + c^2)^{\frac{1}{2}q-1} \quad (50)$$

which is non-increasing for $q \leq 2$. Note that we do not assume that $q > 0$, and consequently generalized Charbonnier loss provides us with more flexibility than Charbonnier loss from Equation 45. Of course if $q < 0$ “loss” becomes “gain”, with a maximum at zero instead of a minimum. To get a proper loss function, take the negative. Figure 4 plots generalized Charbonnier loss for some negative values of q . We see that for $\alpha \rightarrow -\infty$ generalized Charbonnier loss approximates ℓ_0 loss.

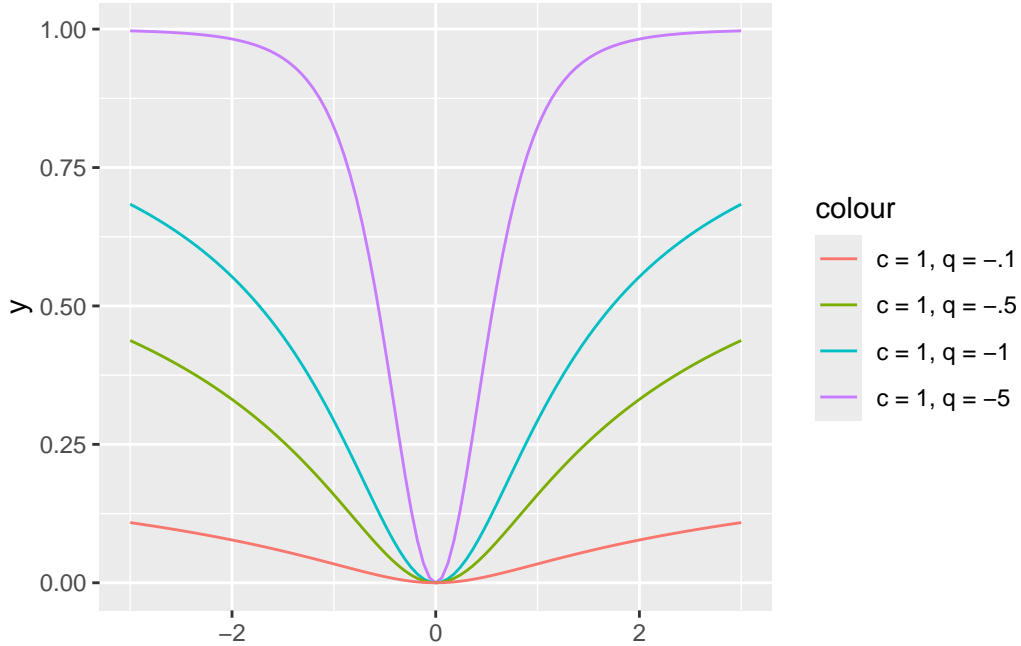


Figure 4: Generalized Charbonnier Loss

5.3 Barron

There are a fair number of generalizations of the power smoother loss functions in the engineering literature. We will discuss one nice generalization from Barron (2019).

$$f_{\alpha,c}(x) = \frac{|\alpha - 2|}{\alpha} \left(\left(\frac{(x/c)^2}{|\alpha - 2|} + 1 \right)^{\alpha/2} - 1 \right). \quad (51)$$

To quote Barron

Here $\alpha \in \mathbb{R}$ is a shape parameter that controls the robustness of the loss and $c > 0$ is a scale parameter that controls the size of the loss's quadratic bowl near $x = 0$.

A number of interesting special cases of Equation 51 are obtained by selecting various values of the α parameters. For $\alpha = 1$ it becomes Charbonnier loss, and for $\alpha = -2$ it is Geman-McClure loss. There are also some limiting cases. For $\alpha \rightarrow 2$ Barron loss becomes squared error loss, for $\alpha \rightarrow 0$ it becomes Cauchy loss, and for $\alpha \rightarrow -\infty$ it becomes Welsch loss.

Accordingly

$$f'_{\alpha,c}(x) = \begin{cases} \frac{x}{c^2} & \text{if } \alpha = 2, \\ \frac{2x}{x^2+2c^2} & \text{if } \alpha = 0, \\ \frac{x}{c^2} \exp\left(-\frac{1}{2}(x/c)^2\right) & \text{if } \alpha \rightarrow -\infty, \\ \frac{x}{c^2} \left(\frac{(x/c)^2}{|\alpha-2|} + 1\right)^{\frac{1}{2}\alpha-1} & \text{otherwise.} \end{cases} \quad (52)$$

and thus

$$h_{\alpha,c}(x) = \begin{cases} \frac{1}{c^2} & \text{if } \alpha = 2, \\ \frac{2}{x^2+2c^2} & \text{if } \alpha = 0, \\ \frac{1}{c^2} \exp\left(-\frac{1}{2}(x/c)^2\right) & \text{if } \alpha \rightarrow -\infty, \\ \frac{1}{c^2} \left(\frac{(x/c)^2}{|\alpha-2|} + 1\right)^{\frac{1}{2}\alpha-1} & \text{otherwise.} \end{cases} \quad (53)$$

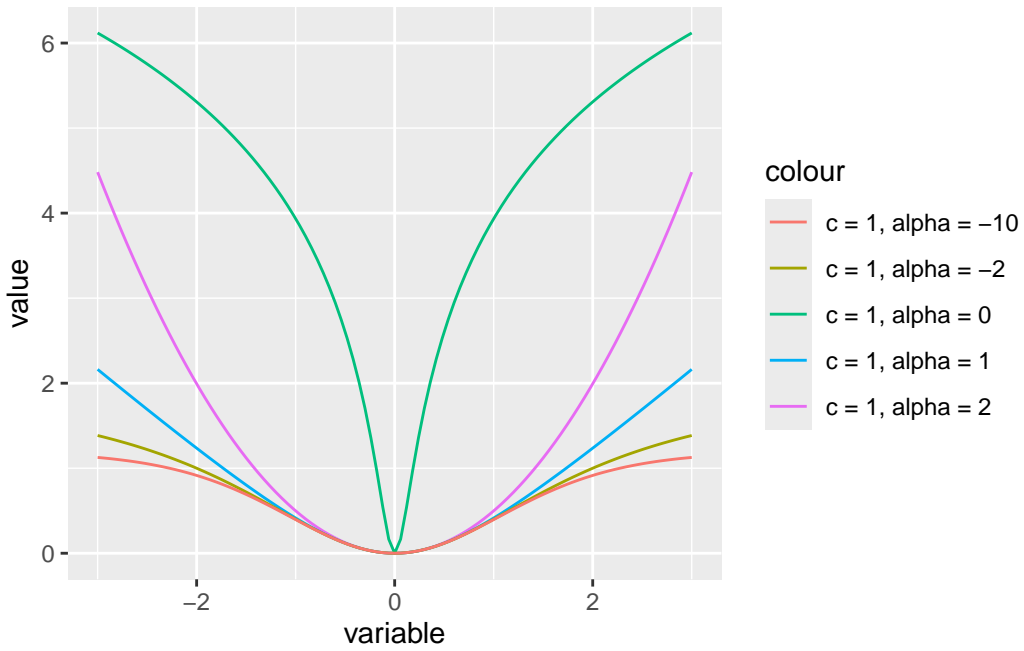


Figure 5: Barron Loss

6 Convolution Smoothers

Suppose π is a probability density, symmetric around zero, with finite or infinite support, expectation zero, and variance one. Define the convolution

$$f_c(x) := \frac{1}{c} \int_{-\infty}^{+\infty} |x - y| \pi\left(\frac{y}{c}\right) dy.$$

Now $c^{-1}\pi(y/c)$ is still a symmetric probability density integrating to one, with expectation zero, but it now has variance c^2 . Thus if $c \rightarrow 0$ it becomes more and more like the Dirac delta function and $f_c(x)$ converges to the absolute value function.

It is clear that we can use any scale family of probability densities to define convolution smoothers. There is an infinite number of possible choices, with finite or infinite support, smooth or nonsmooth, using splines or wavelets, and so on. We give two quite different examples.

6.1 Huber

Take

$$\pi(x) = \begin{cases} \frac{1}{2} & \text{if } |x| \leq 1, \\ 0 & \text{otherwise.} \end{cases}$$

Then

$$f_c(x) = \frac{1}{2c} \int_{-c}^{+c} |x - y| dy = \begin{cases} \frac{1}{2c}(x^2 + c^2) & \text{if } |x| \leq c, \\ |x| & \text{otherwise.} \end{cases} \quad (54)$$

The Huber function (Huber (1964)) is traditionally transformed linearly so that it is zero for $x = 0$. This gives

$$f_c(x) = \begin{cases} \frac{1}{2}x^2 & \text{if } |x| < c, \\ c|x| - \frac{1}{2}c^2 & \text{otherwise.} \end{cases} \quad (55)$$

For robust estimation and IRLS it does not matter if we use Equation 54 or Equation 55. Our discussion in the introduction suggests that if we just want a smoother of the absolute value function, then Equation 54 is the natural choice, if we want a robust loss function that combines the advantages of least squares and least absolute value then that leads us to Equation 55.

Because Charbonnier loss behaves the same way as Huber loss, as absolute value loss for large x and as squared loss for small x , it is also known as Pseudo-Huber loss.

The Huber function is differentiable, although not twice differentiable. Its derivative is

$$f'(x) = \begin{cases} c & \text{if } x \geq c, \\ x & \text{if } |x| \leq c, \\ -c & \text{if } x \leq -c. \end{cases}$$

$$\omega(x) = \begin{cases} \frac{c}{x} & \text{if } x \geq c, \\ 1 & \text{if } |x| \leq c, \\ -\frac{c}{x} & \text{if } x \leq -c. \end{cases}$$

The Huber function is even and differentiable. Moreover $f'(x)/x$ decreases from. Thus Theorem 4.5 applies.

The MDS majorization algorithm for the Huber loss is to update Y by minimizing (or by performing one smacof step to decrease)

$$\sum \omega_k(Y)(\delta_k - d_k(X))^2$$

where

$$\omega_k(Y) = \begin{cases} \omega_k & \text{if } |\delta_k - d_k(Y)| < c, \\ \frac{c\omega_k}{|\delta_k - d_k(Y)|} & \text{otherwise.} \end{cases}$$

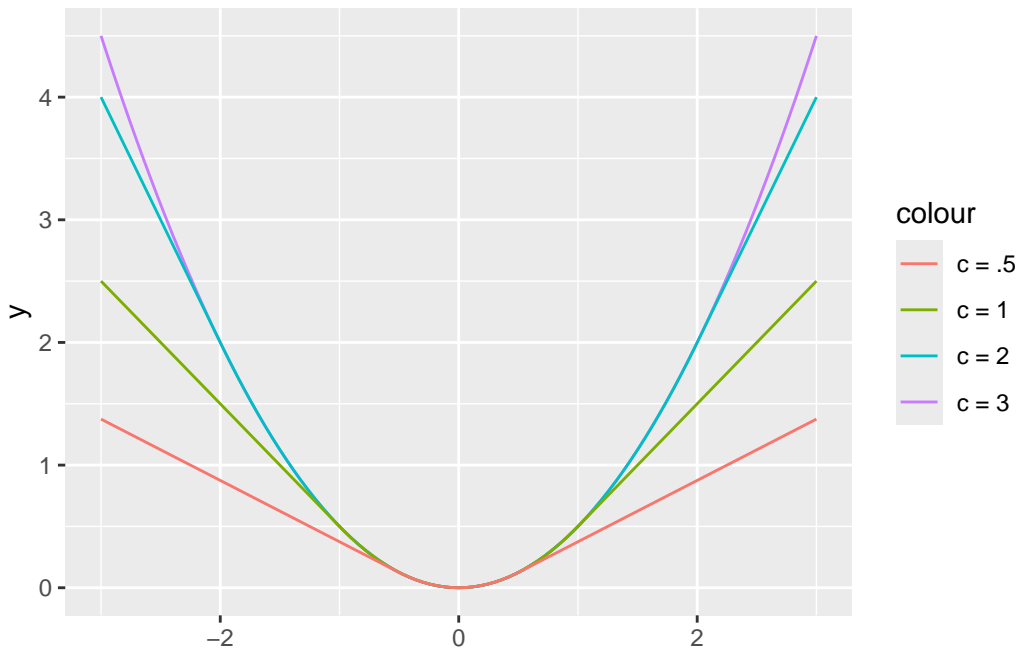


Figure 6: Huber Loss

6.2 Gaussian

In De Leeuw (2018) we also discussed the convolution smoother proposed by Voronin, Ozkaya, and Yoshida (2014). The idea is to use the convolution of the absolute value function and a Gaussian pdf.

$$f(x) = \frac{1}{c\sqrt{2\pi}} \int_{-\infty}^{+\infty} |x - y| \exp \left\{ -\frac{1}{2} \left(\frac{y}{c} \right)^2 \right\} dy$$

Carrying out the integration gives

$$f_c(x) = x \{ 2\Phi(x/c) - 1 \} + 2c\phi(x/c).$$

The derivative is

$$f'_c(x) = 2\Phi(x/c) - 1$$

It may not be immediately obvious in this case that the weight function $f'(x)/x$ is non-increasing on \mathbb{R}^+ . We prove that its derivative is negative on $(0, +\infty)$. The derivative of $f'(x)/x$ has the sign of $xf''(x) - f'(x)$, which is $z\phi(z) - \Phi(z) + 1/2$, with $z = x/c$. It remains to show that $\Phi(z) - z\phi(z) \geq \frac{1}{2}$, or equivalently that $\int_0^z \phi(x) dx - z\phi(z) \geq 0$. Now if $0 \leq x \leq z$ then $\phi(x) \geq \phi(z)$ and thus $\int_0^z \phi(x) dx \geq \phi(z) \int_0^z dx = z\phi(z)$, which completes the proof.

$$\omega_k(Y) = \frac{\Phi((\delta_k - d_k(Y))/c) - \frac{1}{2}}{\delta_k - d_k(Y)}$$

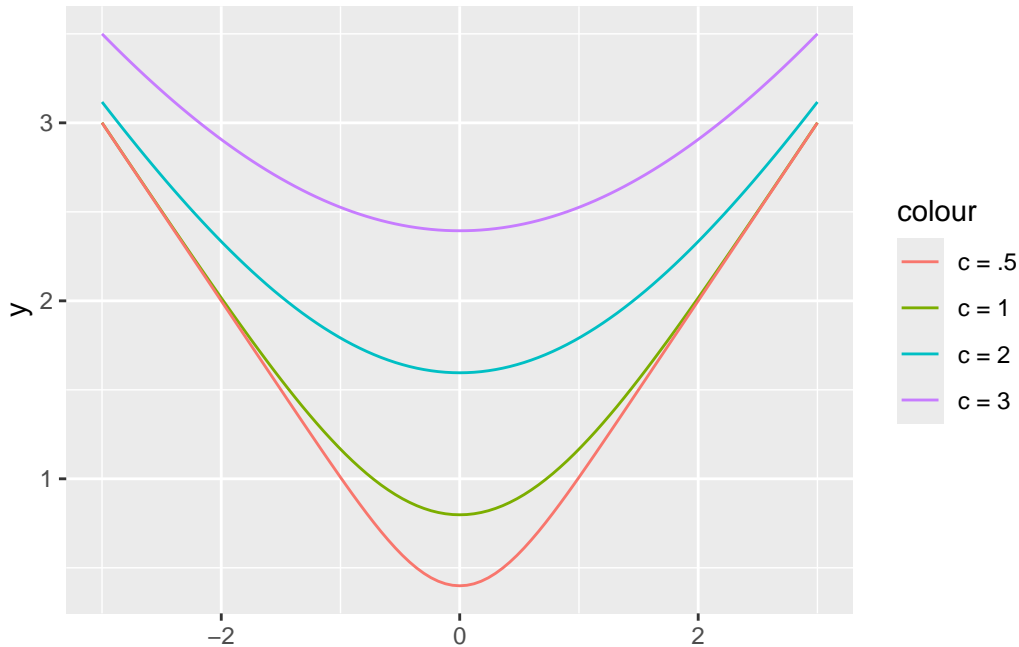


Figure 7: Gaussian Convolution Loss

7 A Bouquet of Loss Functions

In the early seventies, after the pioneering mostly theoretical work in robust statistics of Huber, Hampel, and Tukey, the mainframe computer allowed statisticians to make large-scale comparisons of many robust loss functions. The most impressive of such comparisons was the Princeton Robustness Study (Andrews et al. (1972)).

In Holland and Welsch (1977) the computer package ROSEPACK was introduced that made it relatively easy to compute robust estimators using several different loss functions. Eight different weight functions were implemented as options. Somewhat later Coleman et al. (1980) made an more modern computer implementation available, using the same eight weight functions, which was not limited to mainframes.

We have implemented the same eight weight functions in smacofRobust. Below we give formulas for the loss function, the influence function, and the weight function. One of the eight is Huber loss, which we already discussed in the convolution section. We graph the remaining seven loss functions for selected values of the “tuning constants” c .

Holland and Welsch (1977), following Andrews et al. (1972), distinguish between “hard redescenders” that have an influence function f' equal to zero if x is large enough (Andrews, Tukey, and Hinich loss), “soft redescenders” with influence functions asymptotic to zero for

large x (Cauchy, Welsch loss), and loss functions with a monotone influence function (Huber, Logistic, Fair loss)

7.1 Andrews

The first loss function in this section is taken from Andrews et al. (1972).

$$f(x) = \begin{cases} c^2(1 - \cos(x/c)) & \text{if } |x| \leq \pi c, \\ 2c^2 & \text{otherwise.} \end{cases} \quad (56)$$

$$f'(x) = \begin{cases} c \sin(x/c) & \text{if } |x| \leq \pi c, \\ 0 & \text{otherwise.} \end{cases} \quad (57)$$

$$\omega(x) = \begin{cases} (x/c)^{-1} \sin(x/c) & \text{if } |x| \leq \pi c, \\ 0 & \text{otherwise.} \end{cases} \quad (58)$$

Because \cos is even and $\sin(x)/x$ decreases on $[0, \pi]$ Theorem 4.5 applies.

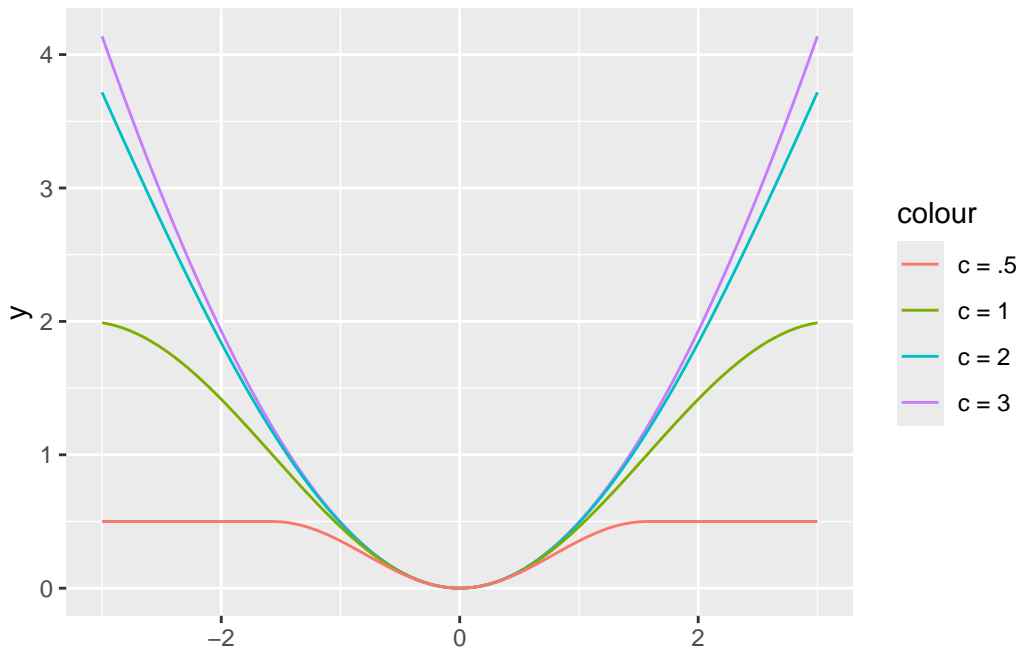


Figure 8: Andrews Loss

7.2 Tukey

The usual reference for Tukey loss is Beaton and Tukey (1974), although closely related hard redescenders are also in Andrews et al. (1972).

$$f(x) = \begin{cases} \frac{c^2}{6} (1 - (1 - (x/c)^2)^3) & \text{if } |x| \leq c, \\ \frac{c^2}{6} & \text{otherwise.} \end{cases} \quad (59)$$

$$f'(x) = \begin{cases} x (1 - (1 - (x/c)^2)^2) & \text{if } |x| \leq c, \\ 0 & \text{otherwise.} \end{cases} \quad (60)$$

$$\omega(x) = \begin{cases} (1 - (1 - (x/c)^2)^2) & \text{if } |x| \leq c, \\ 0 & \text{otherwise.} \end{cases} \quad (61)$$

The conditions of Theorem 4.5 are clearly satisfied.

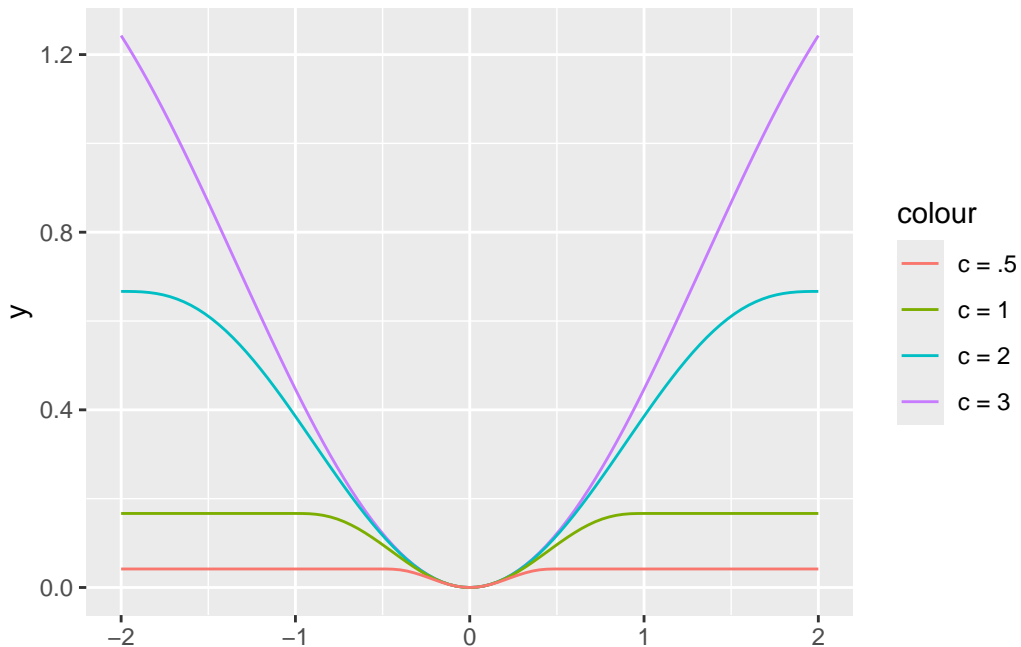


Figure 9: Tukey Loss

7.3 Hinich

Hinich loss, from Hinich and Talwar (1975), is somewhat special because it is not differentiable at c . For $x \neq c$ and $x > 0$ the function $f'(x)/x$ is discontinuous, but non-increasing on

$[0, +\infty)$.

$$f(x) = \begin{cases} \frac{1}{2}x^2 & \text{if } |x| \leq c, \\ \frac{1}{2}c^2 & \text{otherwise.} \end{cases} \quad (62)$$

$$g(x) = \begin{cases} x & \text{if } |x| \leq c, \\ 0 & \text{otherwise.} \end{cases} \quad (63)$$

$$h(x) = \begin{cases} 1 & \text{if } |x| \leq c, \\ 0 & \text{otherwise.} \end{cases} \quad (64)$$

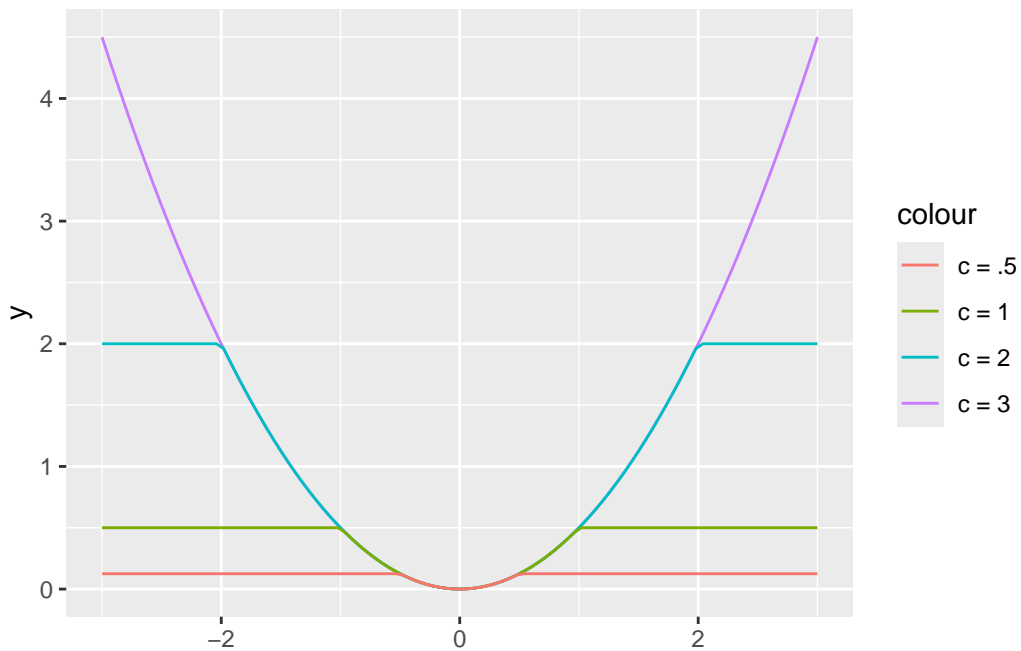


Figure 10: Hinich Loss

7.4 Cauchy

Cauchy loss seems to have many names. Black and Anandan (1996) call it Lorentzian loss, and Holland and Welsch (1977) call it t-likelihood loss. It is related to the Cauchy distribution, which is Student's t distribution with one degree of freedom.

Mlotshwa, Van Deventer, and Sergeevna Bosman (2023)

$$f(x) = \frac{1}{2}c^2 \log(1 + \{\frac{x}{c}\}^2), \quad (65)$$

$$f'(x) = x \frac{1}{\{1 + \frac{x}{c}\}^2}, \quad (66)$$

$$\omega(x) = \frac{1}{\{1 + \frac{x}{c}\}^2} \quad (67)$$

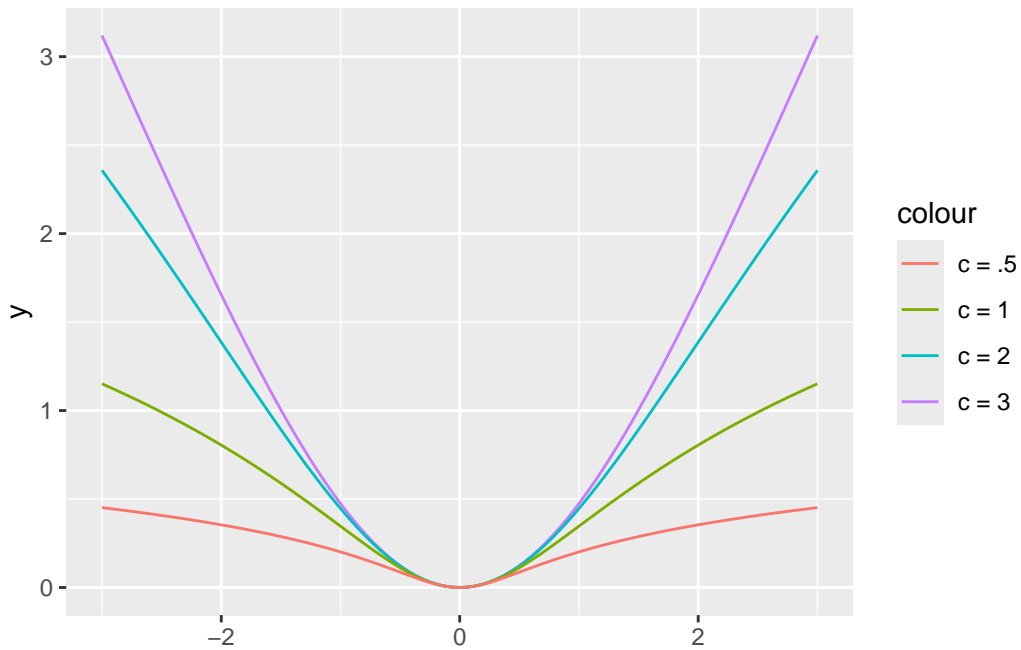


Figure 11: Cauchy Loss

7.5 Welsch

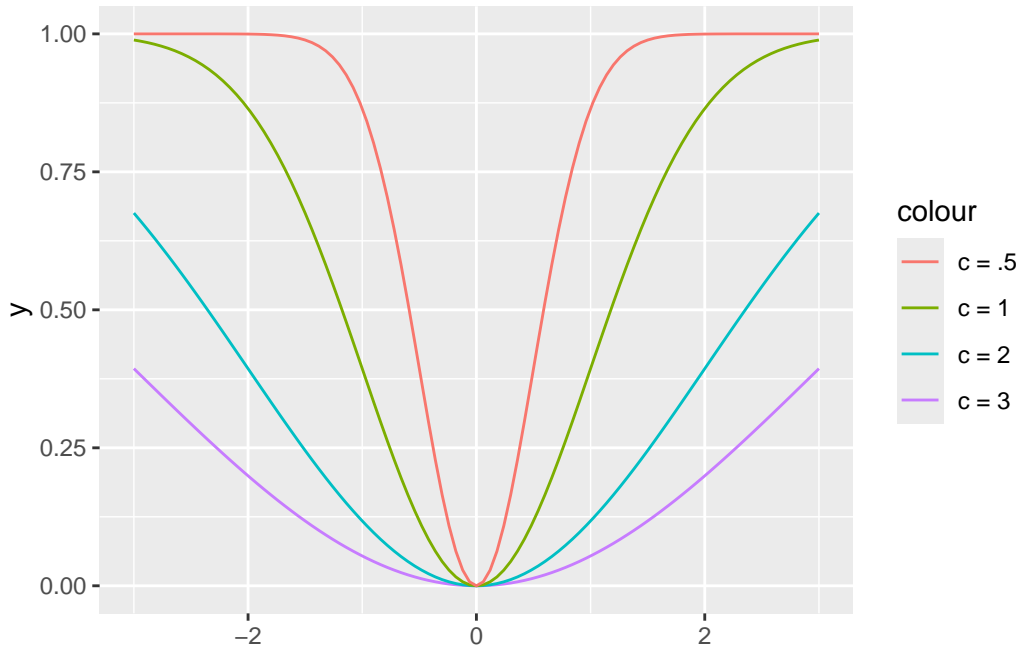
Dennis Jr and Welsch (1978)

Leclerc loss

$$f(x) = \frac{1}{2}c^2[1 - \exp(-\{\frac{x}{c}\}^2)], \quad (68)$$

$$f'(x) = x \exp(-\{\frac{x}{c}\}^2), \quad (69)$$

$$\omega(x) = \exp(-\{\frac{x}{c}\}^2), \quad (70)$$



7.6 Logistic

$$f(x) = c^2 [\log(\cosh(x/c))], \quad (71)$$

$$f'(x) = c \tanh(x/c), \quad (72)$$

$$\omega(x) = (x/c)^{-1} \tanh(x/c). \quad (73)$$

7.7 Fair

$$f(x) = c^2 \{|x|/c - \log(1 + |x|/c)\}, \quad (74)$$

$$f'(x) = x(1 + (|x|/c))^{-1}, \quad (75)$$

$$\omega(x) = (1 + (|x|/c))^{-1}. \quad (76)$$

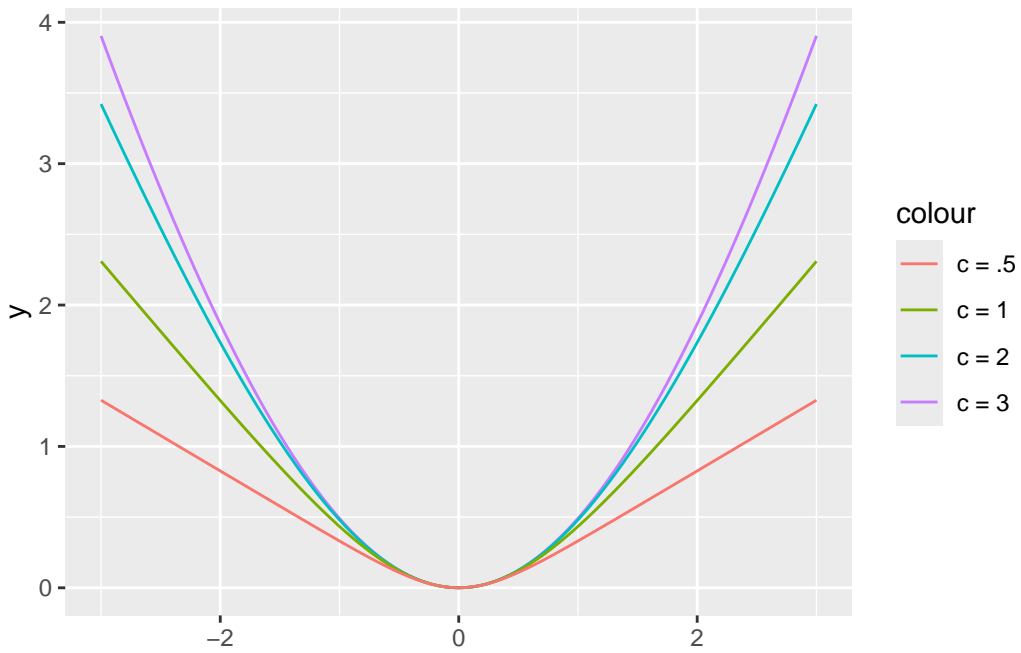


Figure 13: Logistic Loss

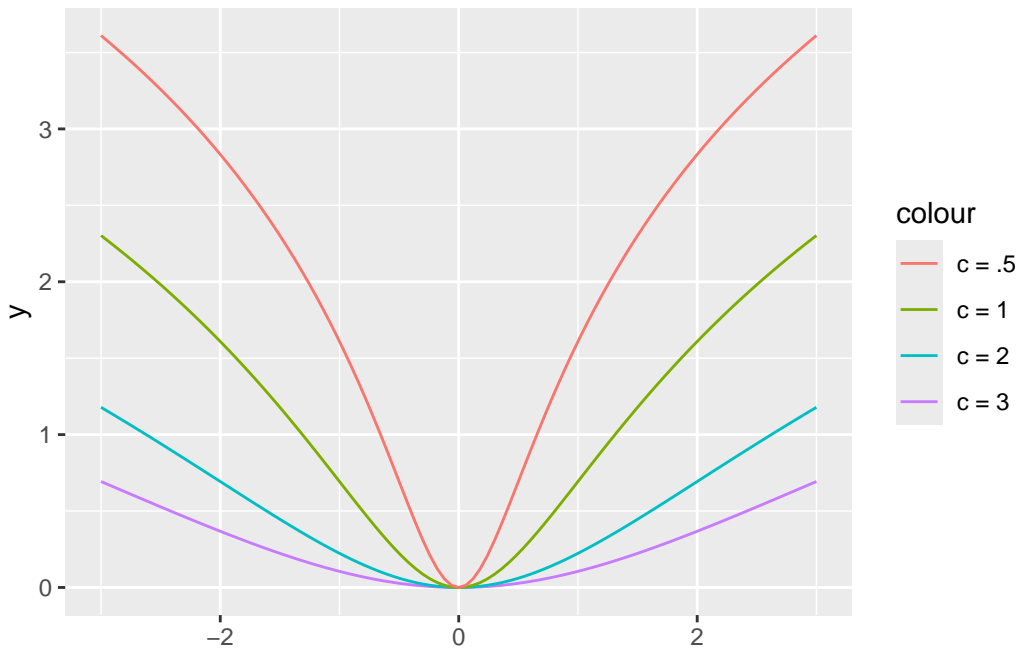


Figure 14: Fair Loss

8 Examples

8.1 Gruijter

The example we use are dissimilarities between nine Dutch political parties, collected by De Gruijter (1967). They are averages over a politically heterogenous group of 100 introductory psychology students, and consequently they regress to the mean. Any reasonable MDS analysis of these data would at least allow for an additive constant.

Some background on Dutch politics around that time may be useful.

- CPN - Communists.
- PSP - Pacifists, left-wing.
- PvdA - Labour, Democratic Socialists.
- D'66 - Pragmatists, nether left-wing nor right-wing, brand new in 1967.
- KVP - Christian Democrats, catholic.
- ARP - Christian Democrats, protestant.
- CHU - Christian Democrats, protestant.
- VVD - Liberals, European flavour, conservative.
- BP - Farmers, protest party, right-wing.

The dissimilarities are in the table below.

| | KVP | PvdA | VVD | ARP | CHU | CPN | PSP | BP | D66 |
|------|------|------|------|------|------|------|------|------|------|
| KVP | 0.00 | 5.63 | 5.27 | 4.60 | 4.80 | 7.54 | 6.73 | 7.18 | 6.17 |
| PvdA | 5.63 | 0.00 | 6.72 | 5.64 | 6.22 | 5.12 | 4.59 | 7.22 | 5.47 |
| VVD | 5.27 | 6.72 | 0.00 | 5.46 | 4.97 | 8.13 | 7.55 | 6.90 | 4.67 |
| ARP | 4.60 | 5.64 | 5.46 | 0.00 | 3.20 | 7.84 | 6.73 | 7.28 | 6.13 |
| CHU | 4.80 | 6.22 | 4.97 | 3.20 | 0.00 | 7.80 | 7.08 | 6.96 | 6.04 |
| CPN | 7.54 | 5.12 | 8.13 | 7.84 | 7.80 | 0.00 | 4.08 | 6.34 | 7.42 |
| PSP | 6.73 | 4.59 | 7.55 | 6.73 | 7.08 | 4.08 | 0.00 | 6.88 | 6.36 |
| BP | 7.18 | 7.22 | 6.90 | 7.28 | 6.96 | 6.34 | 6.88 | 0.00 | 7.36 |
| D66 | 6.17 | 5.47 | 4.67 | 6.13 | 6.04 | 7.42 | 6.36 | 7.36 | 0.00 |

The reason we have chosen this example is partly because CPN and BP are outliers, and we can expect the robust loss functions to handle outlying dissimilarities differently from the bulk of the data.

Unless otherwise indicated we run `smacofRobust()` with a maximum of 10,000 iterations, and we decide that we have convergence if the difference between consecutive stress values is less than 10^{-15} . We perform one single `smacof` iteration between the updates of the weights. For

each analysis we show the configuration plot, the Shepard plot, and a histogram of the absolute values of the residuals. In the Shepard plot points corresponding to the eight CPN-dissimilarities are labeled “C”, while BP-dissimilarities are “B”.

8.1.1 Least Squares

We start with a least squares analysis, actually with Huber loss with $c = 10$, which for these data is equivalent to least squares. The process converges in 859 iterations.

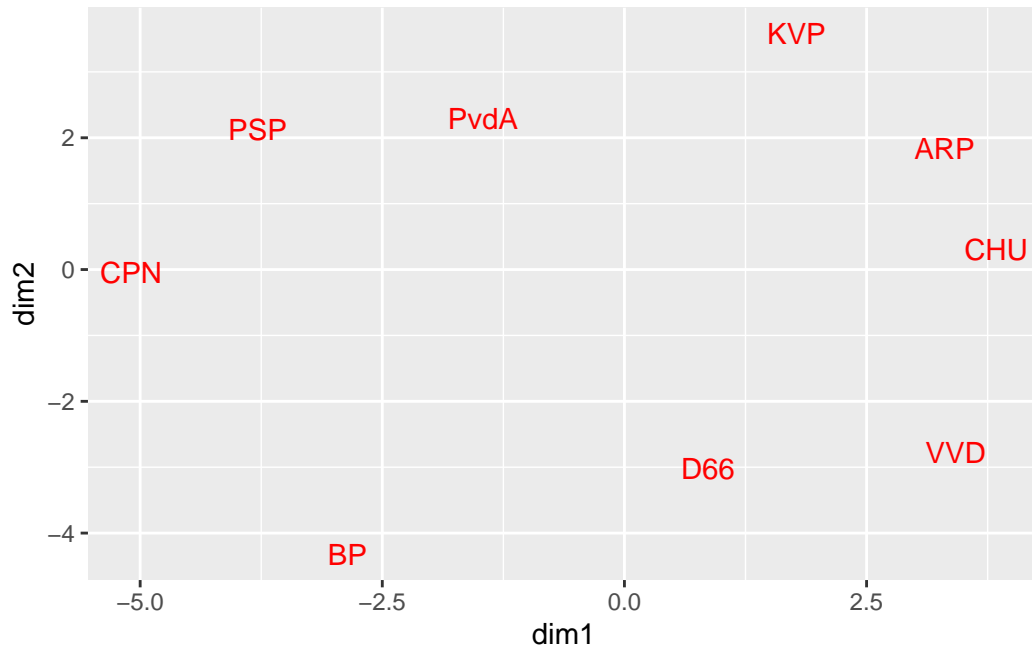


Figure 15: Gruijter Configuration Least Squares

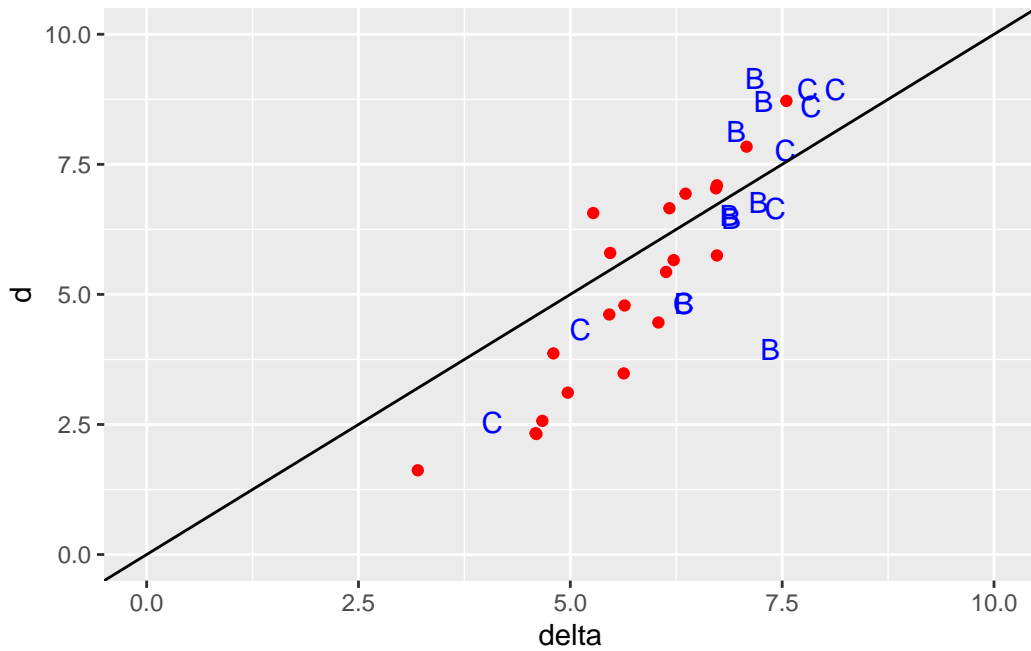


Figure 16: Grubjer Shepard Plot Least Squares

The Shepard plot clearly shows why an additive constant would be very beneficial in this case.

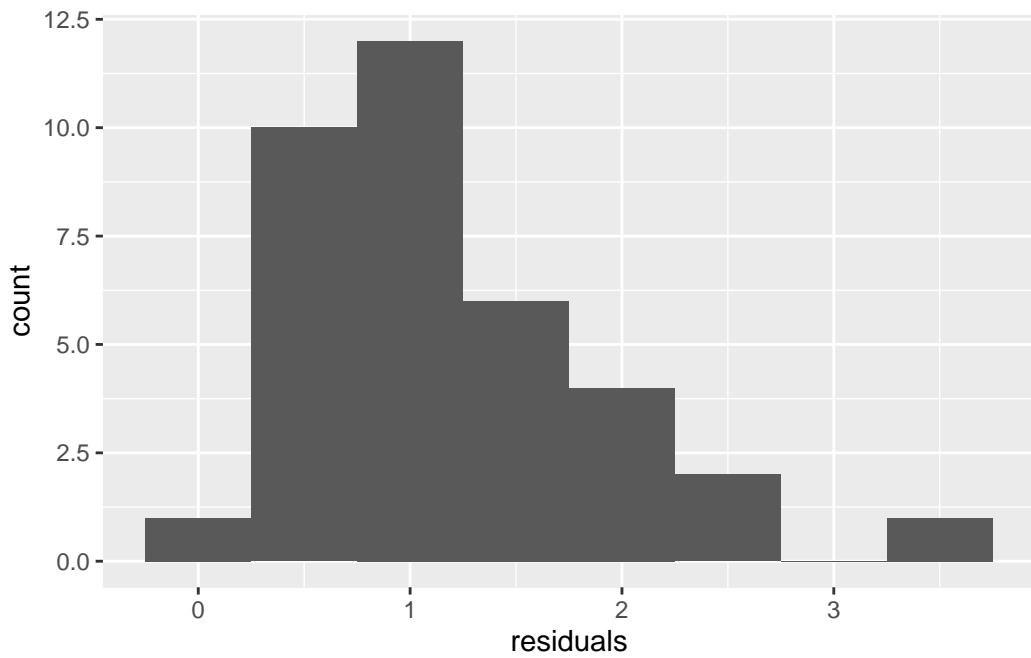


Figure 17: Grubjer Histogram Least Squares Residuals

8.1.2 Least Absolute Value

For our LAV smacof we use engine smacofCharbonnier with $c = .001$. We have convergence in 637 iterations.

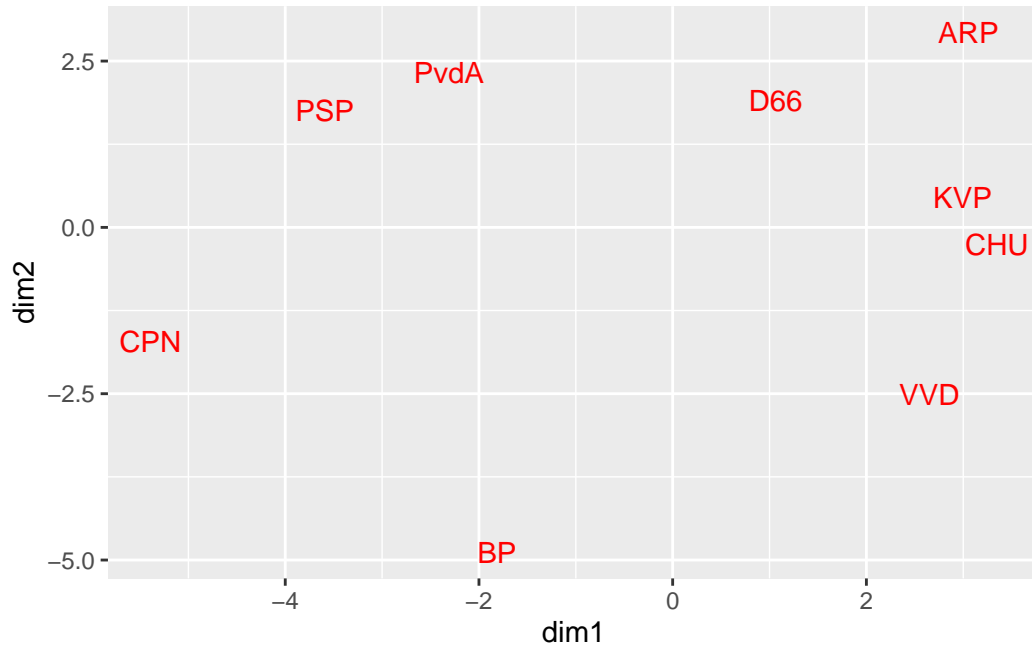


Figure 18: Gruijter Configuration Least Absolute Value

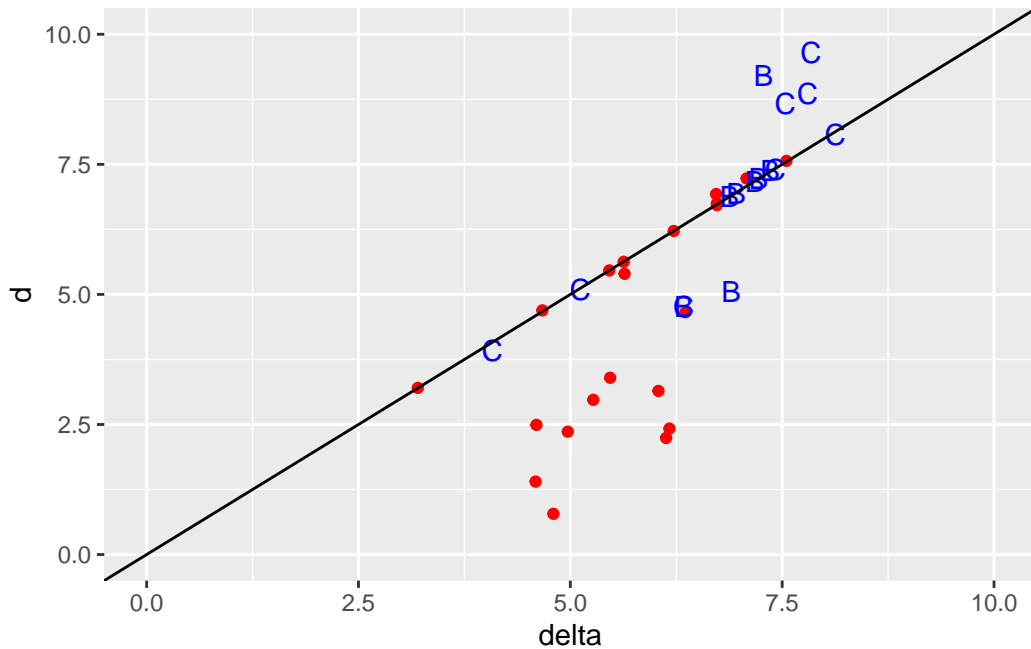


Figure 19: Gruijter Shepard Plot Least Absolute Value

In the Shepard plot we see that there are a number of dissimilarities which are fitted exactly. If we count them there are about 15-20. Note that configurations in two dimensions have $(n - 1) + (n - 2) = 2n - 3$ degrees of freedom, which is 15 in this case. Thus if we take the 15 dissimilarities which are fitted exactly, give them weight one, give all other 21 dissimilarities weight zero, and do a regular non-robust smacof analysis using these weights, then we will have perfect fit in two dimensions, and the solution will be the LAV solution. All this is easier said than done, because it presumes that we use Charbonnier loss with $c = 0$ and that we are able to decide which residuals are exactly equal to zero. The LAV analysis also suggests the possibility of a huge number of local minima, because there are so many ways to pick 15 out of 36 dissimilarities.

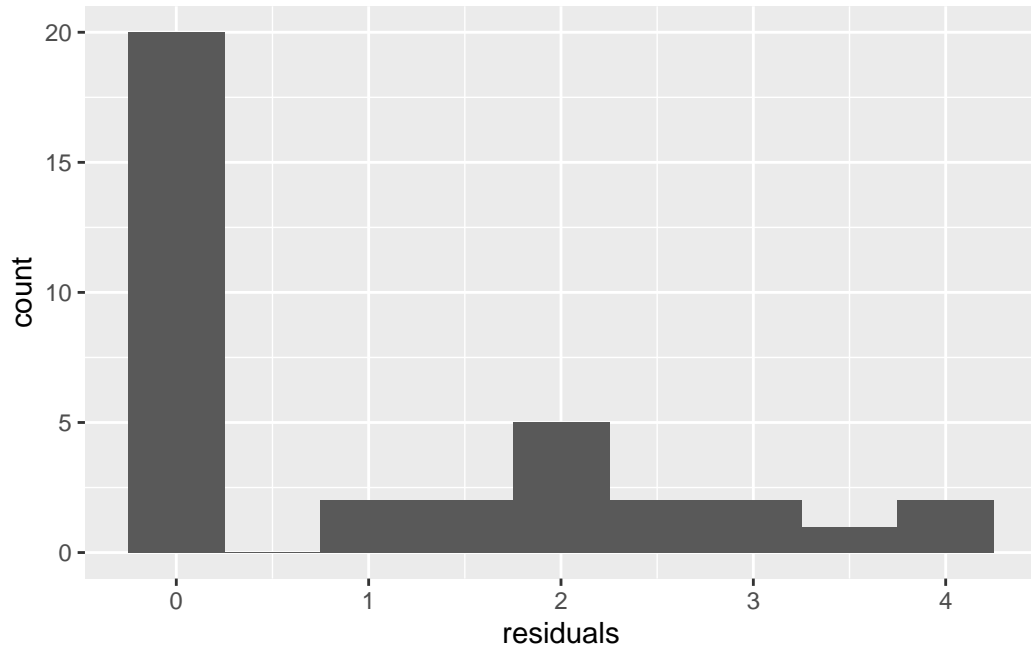


Figure 20: Gruijter Histogram Least Absolute Value Residuals

8.1.3 Huber

smacofHuber with $c = 1$ converges in 165 iterations.

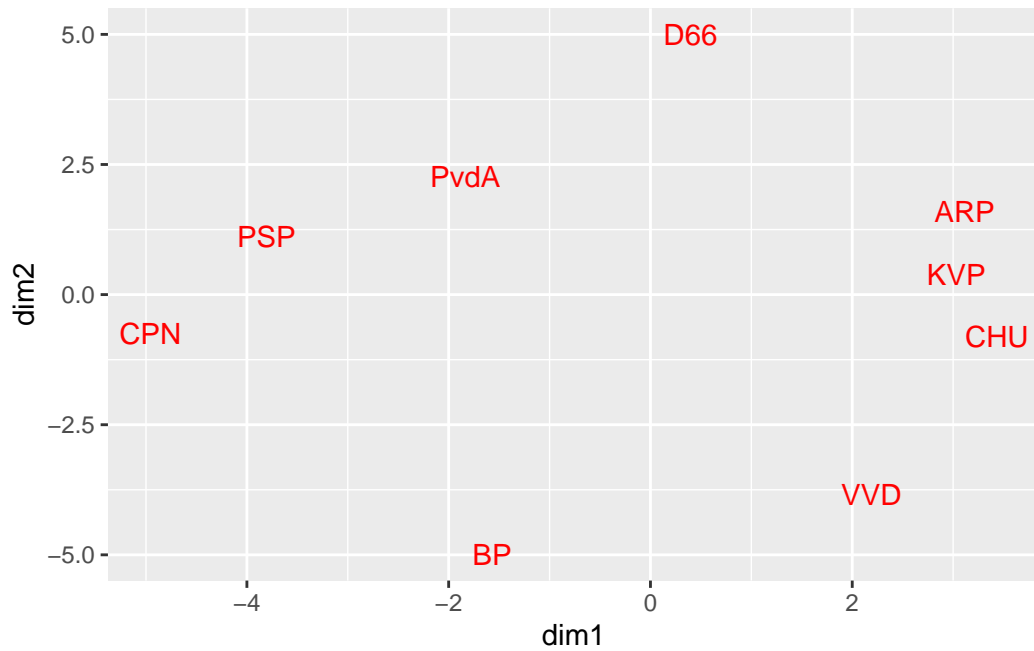


Figure 21: Gruijter Configuration Huber $c = 1$

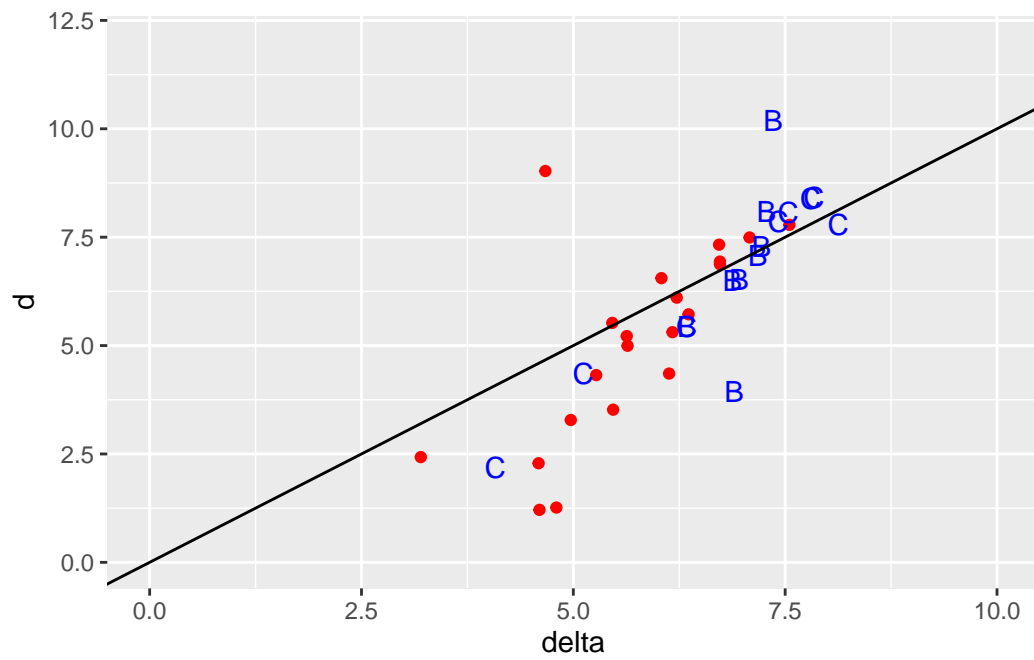


Figure 22: Gruijter Shepard Plot Huber $c = 1$

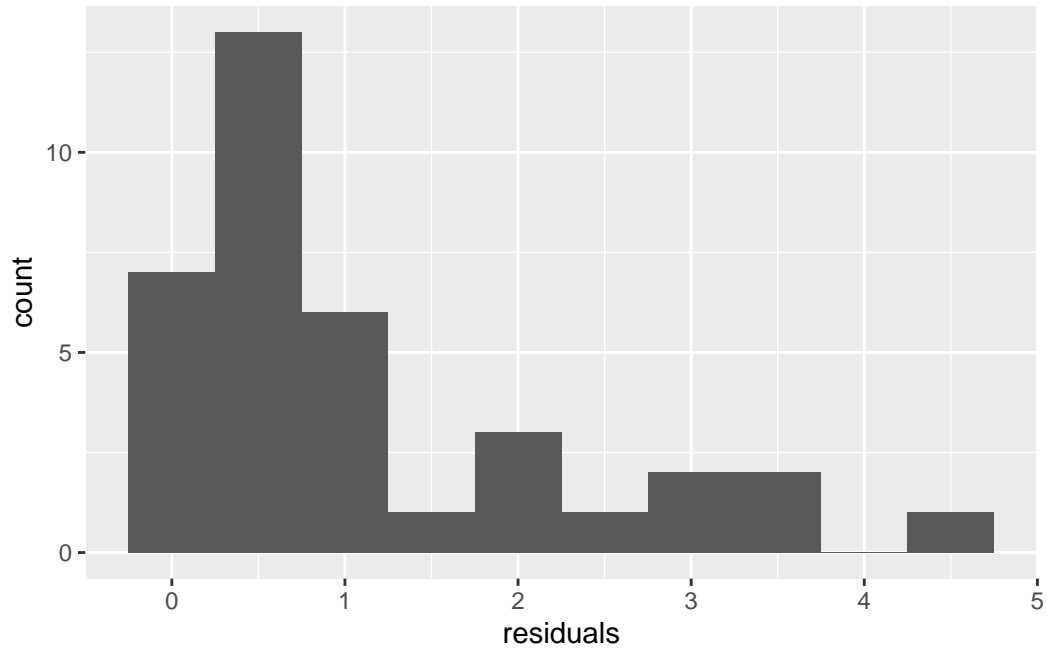


Figure 23: Gruijter Histogram Huber Residuals

8.1.4 Tukey

smacofTukey with $c = 2$ converges in 180 iterations.

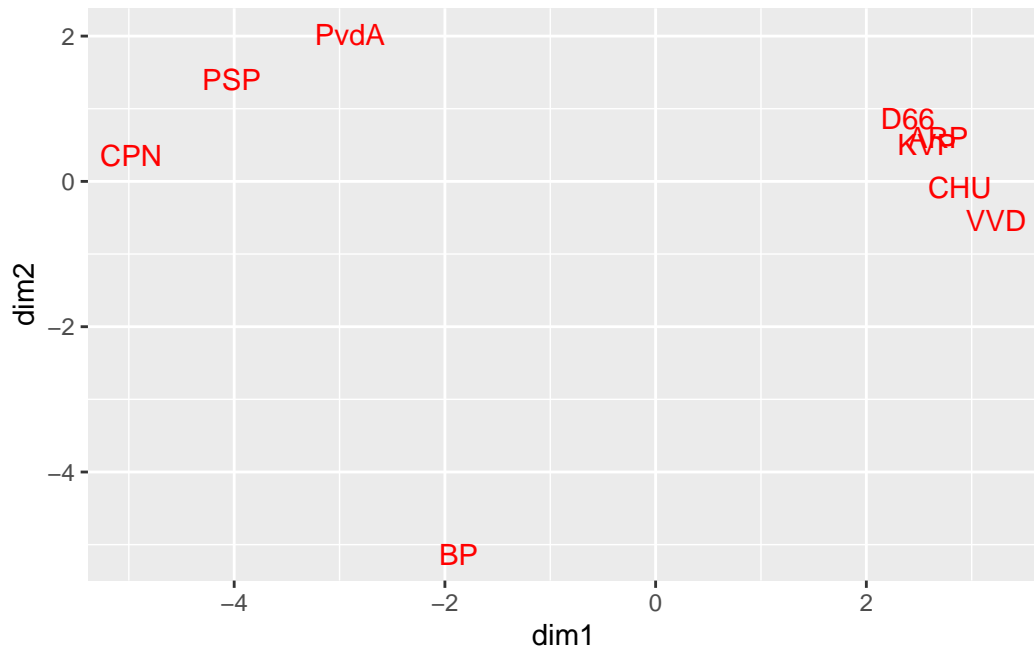


Figure 24: Gruijter Configuration Tukey $c = 2$

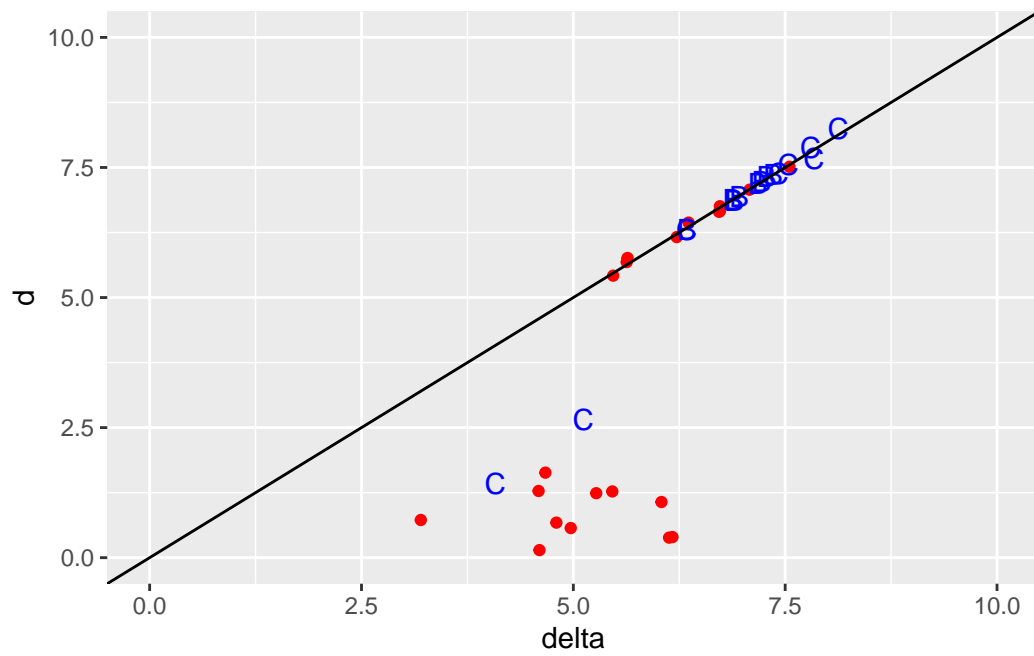


Figure 25: Gruijter Shepard Plot Tukey $c = 2$

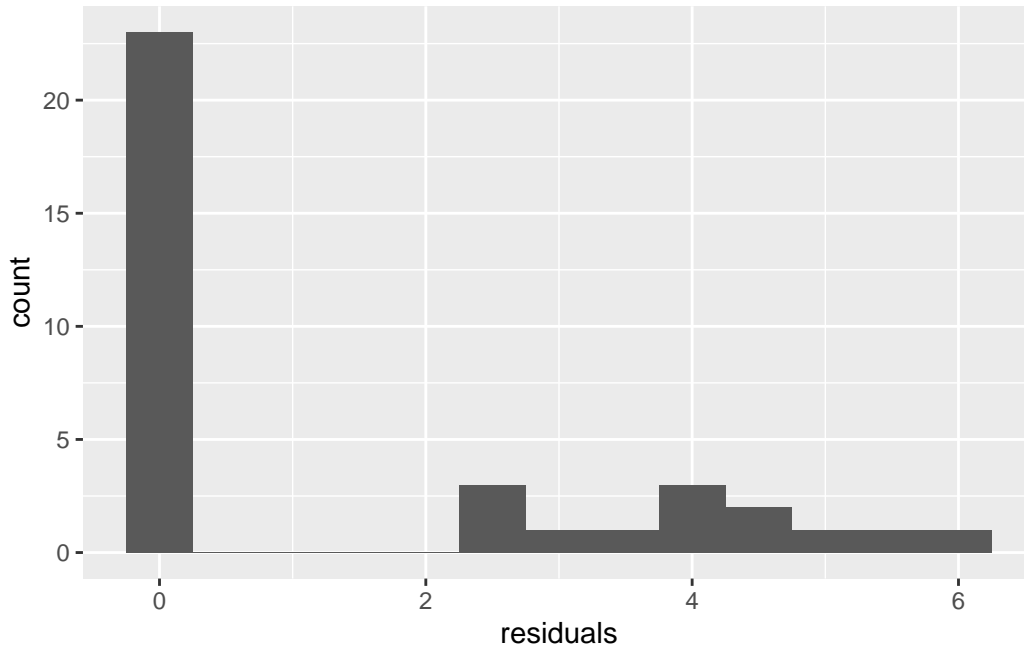


Figure 26: Gruijter Histogram Tukey Residuals

8.2 Rothkopf

Our second example are the Rothkopf Morse data (Rothkopf (1957)), which have a better fit and have fewer outliers than the Gruijter data. We used the asymmetric confusion matrix from the smacof package (De Leeuw and Mair (2009)) and defined dissimilarities by the Shepard-Luce formula

$$\delta_{ij} = -\log \frac{p_{ij}p_{ji}}{p_{ii}p_{jj}}.$$

8.2.1 Least Squares

For least squares we use the smacofHuber engine with $c = 25$, well outside the range of the residuals. We have convergence in 213 iterations.

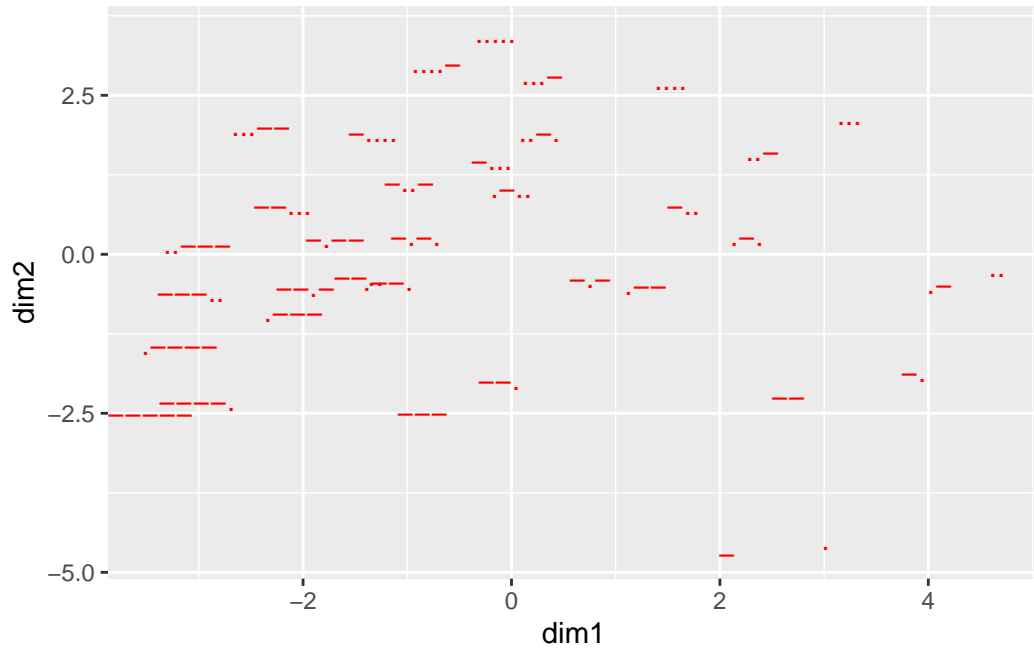


Figure 27: Rothkopf Configuration Least Squares

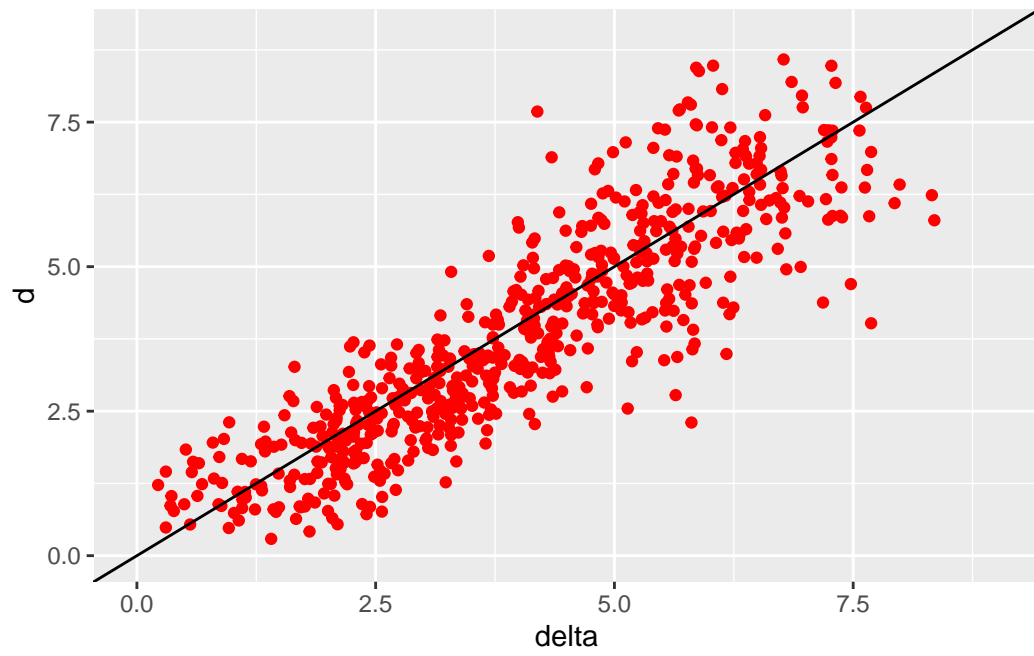


Figure 28: Rothkopf Shepard Plot Least Squares

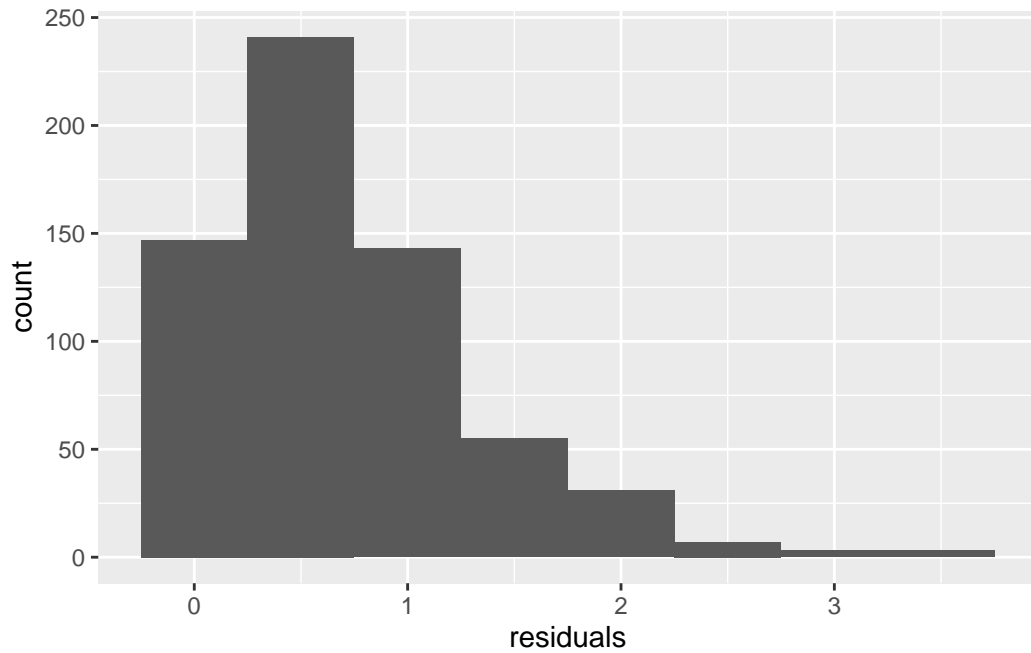


Figure 29: Rothkopf Histogram Least Squares Residuals

8.2.2 Least Absolute Value

For least absolute value we use Chardonnier loss with $c = .001$. We have convergence in 2291 iterations.

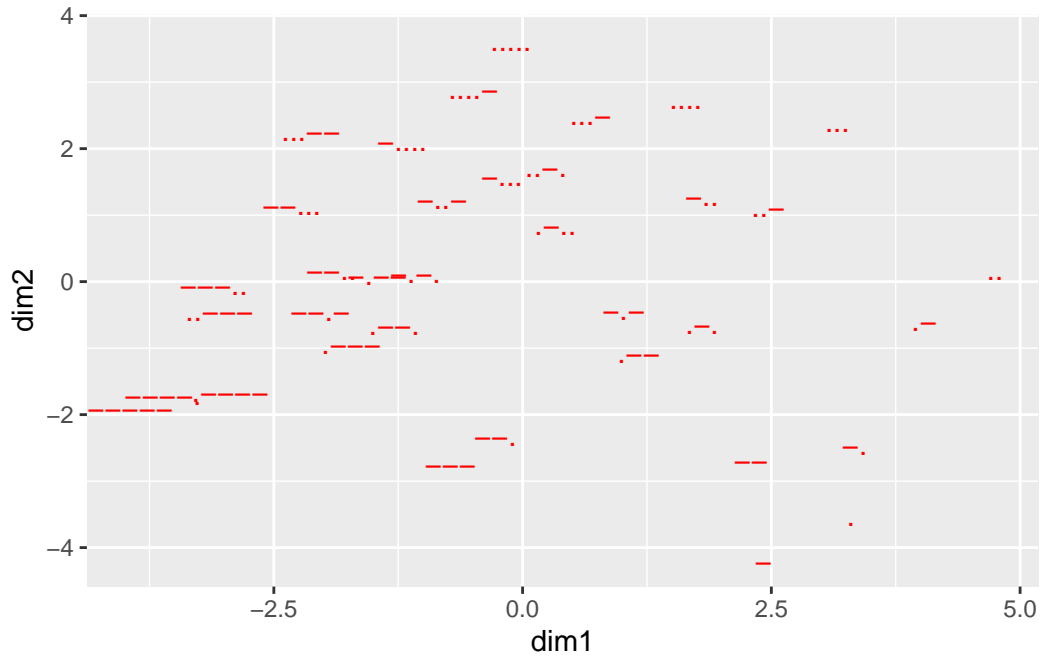


Figure 30: Rothkopf Configuration Least Absolute Value

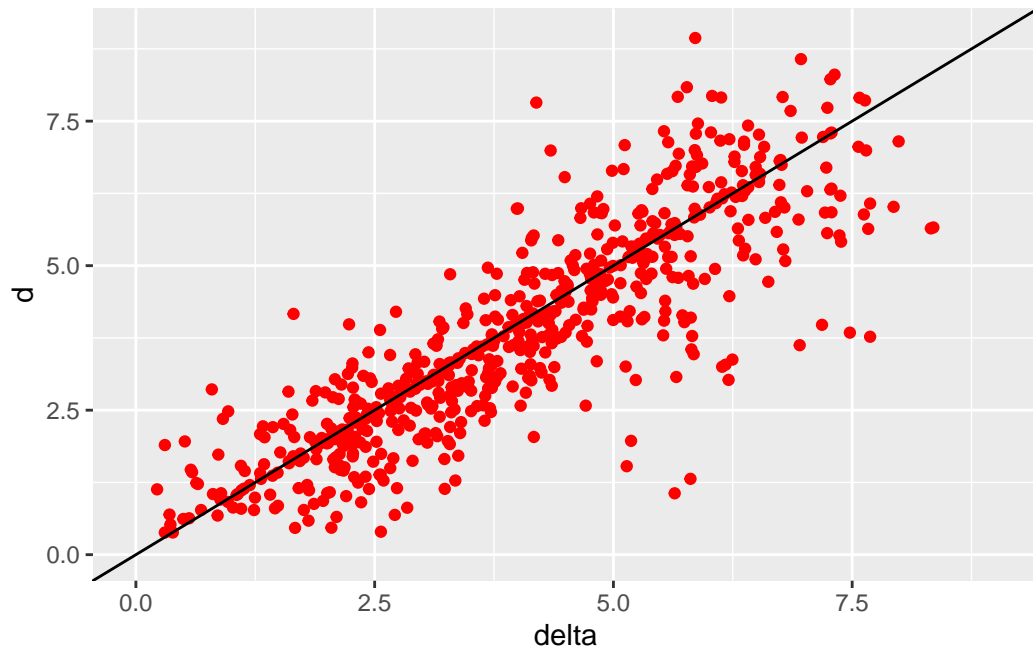


Figure 31: Rothkopf Shepard Plot Least Absolute Value

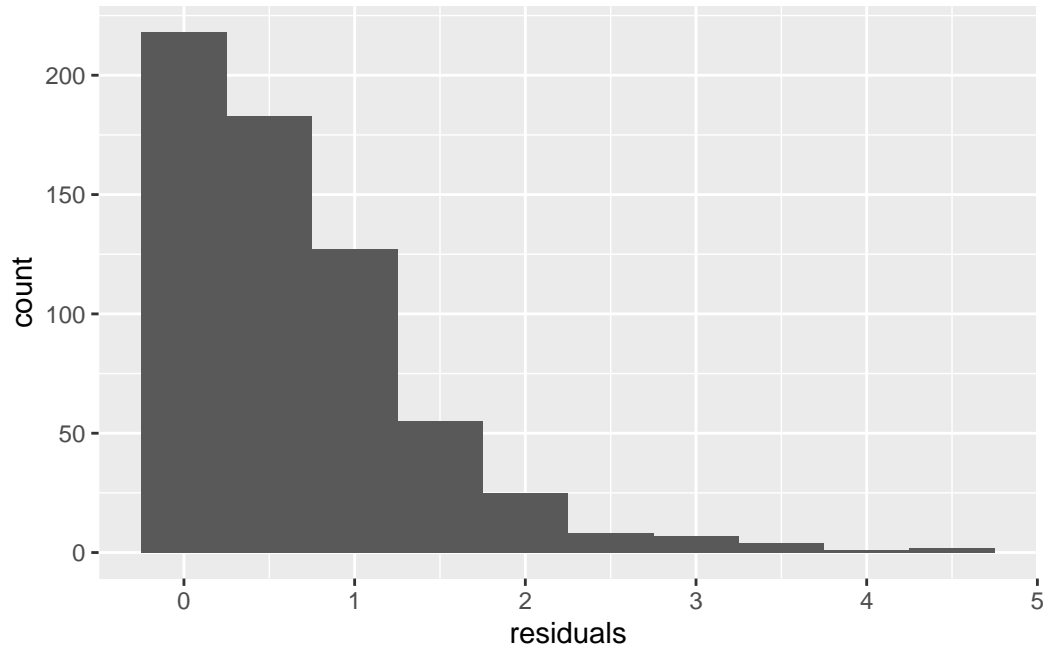


Figure 32: Rothkopf Histogram Least Absolute Value Residuals

8.2.3 Huber

smacofHuber with $c = 1$ converges in 680 iterations.

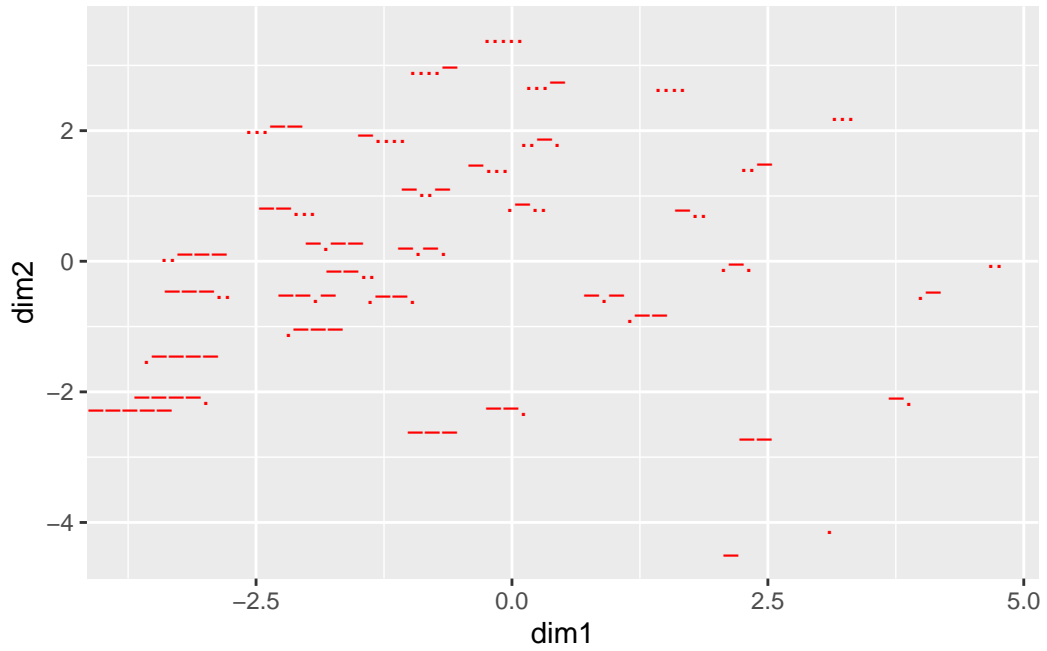


Figure 33: Rothkopf Configuration Huber $c = 1$

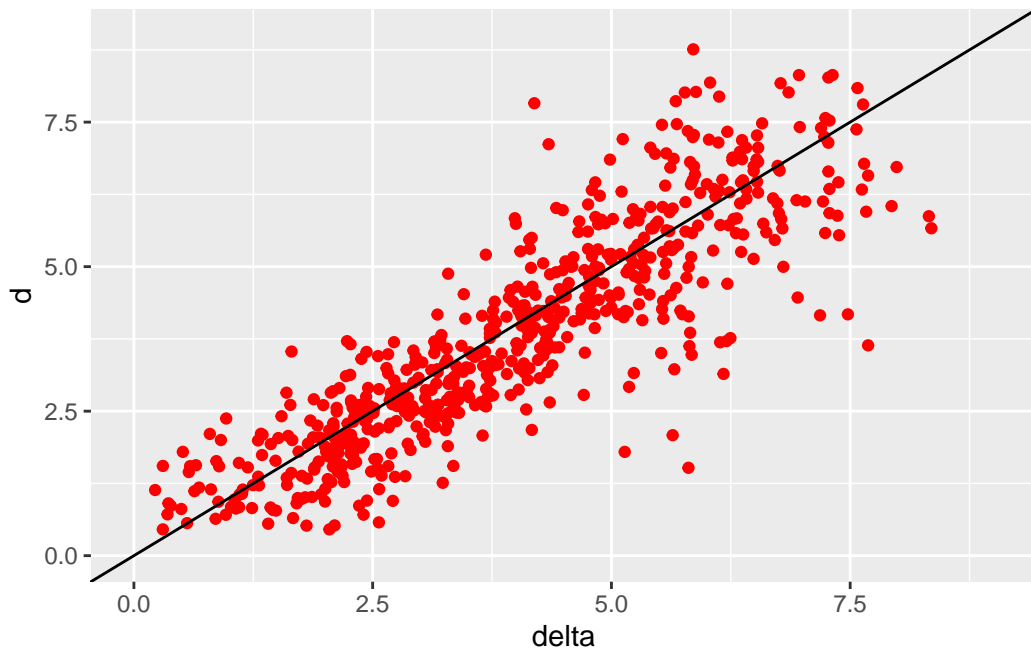


Figure 34: Rothkopf Shepard Plot Huber $c = 1$

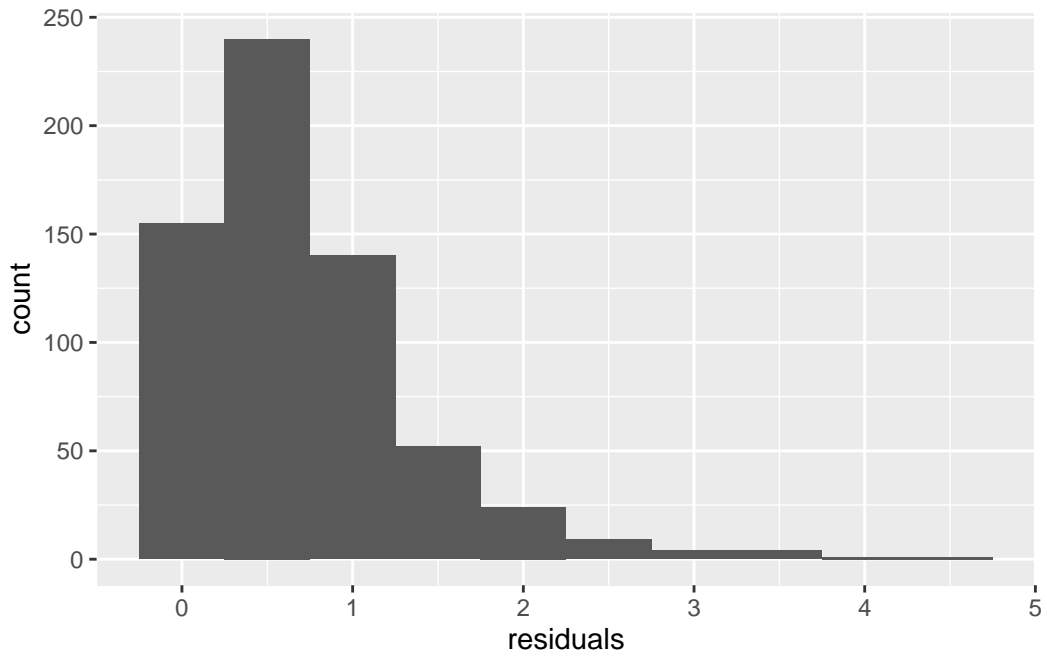


Figure 35: Rothkopf Histogram Huber Residuals

8.2.4 Tukey

Tukey with $c = 1$ converges in 812 iterations.

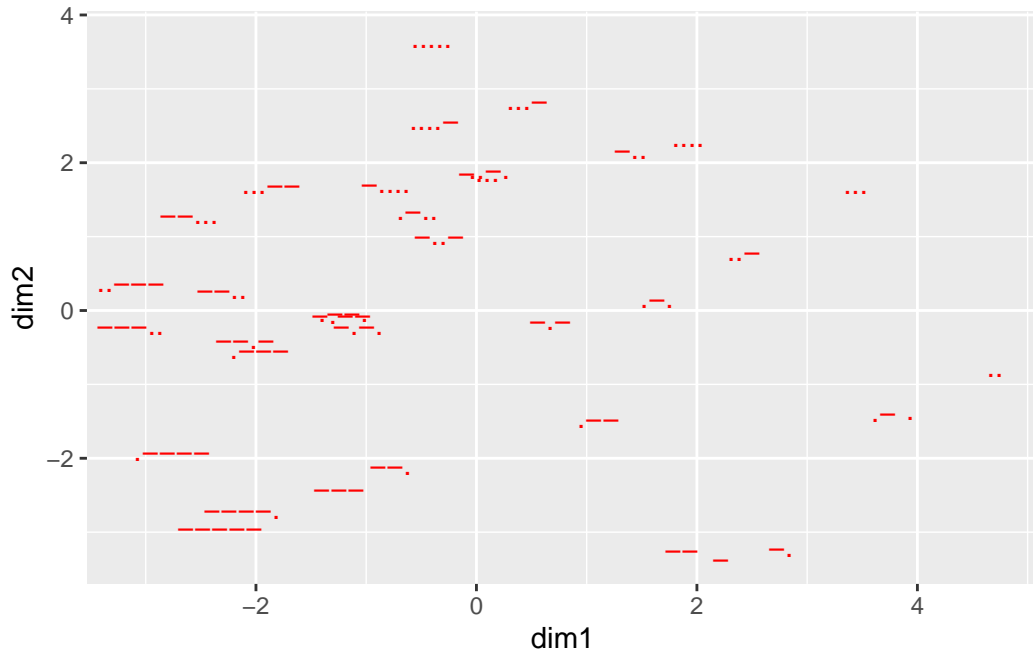


Figure 36: Rothkopf Configuration Tukey c = 1

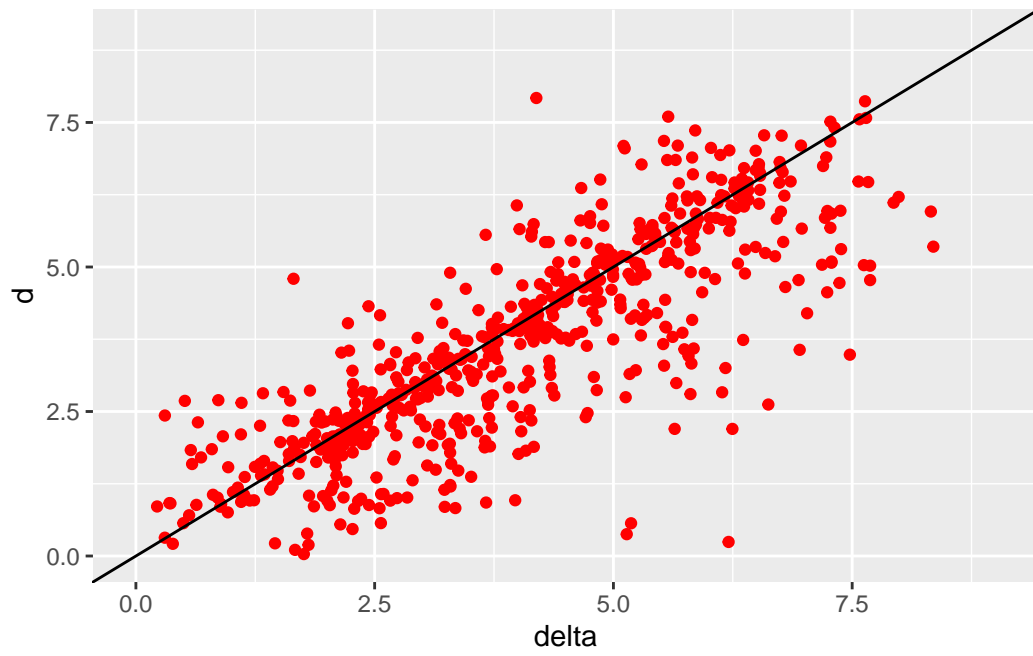


Figure 37: Rothkopf Shepard Plot Tukey c = 1

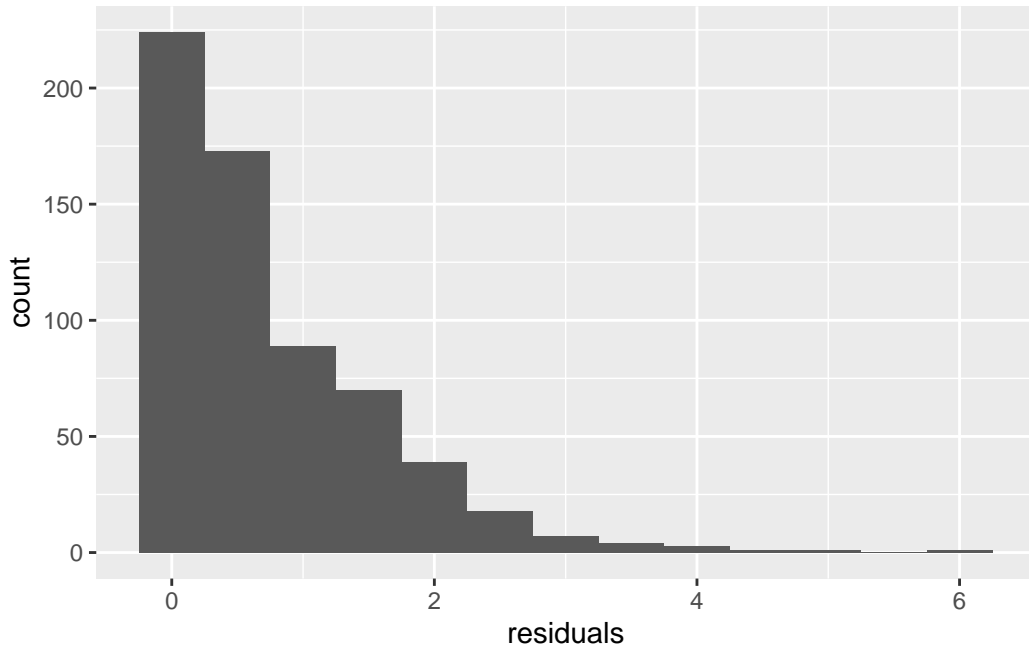


Figure 38: Rothkopf Histogram Tukey Residuals

9 Literature

The literature on results like Theorem 4.5 and Theorem 4.6 is difficult to review. There are various reasons for that. Relevant results have been published in robust statistics, computational statistics, optimization, location analysis, image restoration, sparse recovery. As is often the case, there are not many references between fields, almost everything is within. Even the names of the loss functions differ between fields. Much of it is hard to find in conference proceedings. Also, in most cases, the authors have specific applications in mind, which they then embed in a likelihood, Bayesian, linear regression, logistic regression, facility location, or EM framework and language.

De Leeuw and Lange (2009) give some references to previous work on results like Theorem 4.5, notably Groenen, Giaquinto, and Kiers (2003), Jaakkola and Jordan (2000), and Hunter and Li (2005). In these earlier papers we do not find Theorem 4.5 in its full generality. In Groenen, Giaquinto, and Kiers (2003) majorization of the log logistic function is considered. Besides requiring equality of the function and the majorizing quadratic at the support point y they also require equality at $-y$ and then check that the resulting quadratic is indeed a majorizer. In Jaakkola and Jordan (2000) also consider a symmetrized version of the log logistic function. They note that the resulting function is a convex function of x^2 , and use a linear majorizer at x^2 to obtain a quadratic majorization. Hunter and Li (2005) come closest to Theorem 4.5. In their proposition 3.1 they approximate the general penalty function they use for variable selection at y by a quadratic with coefficient $f'(y)/2y$, and then show that it provides a quadratic majorization. In neither of the three papers there is a notion of sharp quadratic majorization.

I will discuss some of the literature under the headings “robust statistics”, “location analysis”, and “sparse recovery”. Since I most definitely am not an expert in either of these three fields the literature reviews will be biased and incomplete. A final section, where I am somewhat more sure-footed, is “multivariate analysis”.

9.1 Robust Statistics

In robust statistics it has been known for a long time that iterative reweighted least squares (IRLS) with weights $f'(x)/x$ gives a quadratic majorization algorithm. This result, and the corresponding IRLS algorithm, is often attributed to Beaton and Tukey (1974).

9.2 Location Analysis

In location analysis the first majorization/IRLS method is generally attributed to a 16-year old Hungarian mathematics prodigy (Weiszfeld (1937)). His algorithm *avant-la-lettre* was

intended to find the minimum of a function of the form

$$\sigma(x) = \sum_k w_k d_k(x), \quad (77)$$

where $d_k(x) = \|x - y_k\|$, over x in \mathbb{R}^n . The y_k are known locations, called *anchors* in the literature, and the norm is Euclidean. Actually Weiszfeld (1937) did not use weights w_k and worked in three-dimensional space. There is an English translation of Weiszfeld's paper, with bibliography and comments, in Weiszfeld and Plastria (2009).

The problem of minimizing Equation 77 is known under various names, usually consisting of one of the seven different non-empty selections from the triple (Torricelli, Fermat, Weber). The history of the problem is discussed, for example, in Plastria (2011).

Weiszfeld first acknowledges that Sturm (1884) has already established the existence and uniqueness of the minimum point. He then proceeds to give three new proofs. We are interested in the first one, described in his first theorem. It defines the iterative sequence

$$x^{(\nu+1)} = \frac{\sum \omega_k(x^{(\nu)}) y_k}{\sum \omega_k(x^{(\nu)})} \quad (78)$$

with weights

$$\omega_\nu(x) = \frac{1}{d_k(x^{(\nu)})}. \quad (79)$$

The proof then consists of showing that the sequence converges to the unique minimum point of Equation 77. It would have been nice if we were told where Equation 78 came from, but it is simply taken as the starting point.

We can guess what suggested this particular sequence. Observe first that the corresponding problem with $d_k(x)$ replaced by $d_k^2(x)$ is easy to solve. The solution is simply the weighted mean of the y_k . This suggests the rewrite

$$\sigma(x) = \sum_k w_k d_k(x) = \sum_k w_k \frac{1}{d_k(x)} d_k^2(x), \quad (80)$$

which in turn suggests Equation 78. Alternatively, differentiate the loss in Equation 77 and set the partials equal to zero. This gives

$$\sum_k w_k \frac{1}{d_k(x)} (x - y_k) = 0, \quad (81)$$

or

$$x = \frac{\sum \omega_k(x) y_k}{\sum \omega_k(x)}. \quad (82)$$

Sturm (1884) mentions that he derived his existence and uniqueness theorem without using differentiation, but he mentions a paper by Lorenz Lindelöf from 1866 which derives his necessary conditions using differentiation (I have not been able to find a copy of that paper). The discussion is a clear example of the nineteenth century tension between using synthetic (geometric) methods or analytic (calculus) methods.

Even more synthetic were the methods in a paper by Lamé and Clapeyron of 1829. (I have not been able to find a copy of that paper either). There is a partial translation in Franksen and Grattan-Guinness (1989). Lamé and Clapeyron suggest solving their “moindre distances” problems, which generalize the single facility location problem in various ways, by ingenious systems of pulleys. There is a translation of the general principles section of the Lamé and Clapeyron paper in Franksen and Grattan-Guinness (1989). From that translation we read in General Principle 9 their defense of the mechanical method they propose.

But if one considers that the proposed problem is totally insoluble, in its entire generality, by the current means of analysis and geometry; that it is only in very special and very simple cases that one may obtain a complete graphical solution; that finally in the applications the data themselves are only approximate; one will be forced to admit that in the state of imperfection in which algebraic analysis is still found today, the manner of solution in question here is the only one which obtains for the proposed problem.

In General Principle 21 they suggest an iterative method of “trial and error” to solve the weighted location problem, where the weights in each iteration are adjusted by multiplying them with the inverse of the distances in the previous iteration. Thus they propose the Weiszfeld algorithm, albeit in a version using pulleys.

Over the years the location problem has been generalized in numerous directions, to multiple locations, to using different norms, to unknown anchors, to nonlinear manifolds, and to obnoxious anchors you want to be far from. A good recent overview is Beck and Sabach (2015). A paper close in spirit to our paper is Aftab, Hartley, and Trumpf (2015), which has generalizations to ℓ_q norms and to Riemannian manifolds of non-negative curvature.

There is a huge literature on the convergence of the Weiszfeld algorithm. As in our Section 3.1 we can simply use AM/GM inequality. Thus

$$\|x - y_\nu\| \leq \frac{1}{2} \frac{1}{\|x^{(\nu)} - y_\nu\|} (\|x - y_\nu\|^2 + \|x^{(\nu)} - y_\nu\|^2), \quad (83)$$

which immediately gives Equation 78. Convergence follows from the general majorization of MM theory.

Unlike robust smacof the Toricelli-Fermat-Weber problem is convex, and consequently has no problems with non-global local minima. A most elegant proof of convergence using the tools

of modern convex analysis is in Mordukhovich and Nam (2019). Older proofs sometimes have difficulty dealing with cases in which the iterates coincide with one of the anchors or in which the solution is actually one of the anchors. This creates problems similar to the problems in our Section 3.2, but in this simple case the problem is can be completely resolved using convexity and has no serious algorithmic consequences.

In a straightforward generalization of the Toricelli-Fermat-Weber problem , which is particularly relevant for the developments in our paper, Katz (1969) proposes to minimize

$$\sigma(x) = \sum_{k=1}^m w_k d_k^q(x), \quad (84)$$

and even

$$\sigma(x) = \sum_{k=1}^m w_k f_k(d_k(x)) \quad (85)$$

with Euclidean distances and f_k functions defined on the non-negative reals. Note that if f is convex and increasing then σ of Equation 85 is convex. If $q \geq 1$ then σ of Equation 84 is convex.

The algorithm Katz suggests for minimizing σ of Equation 84 generalizes the decomposition in Equation 80 to

$$\sigma(x) = \sum_k w_k d_k(x) = \sum_k w_k \frac{1}{d_k^{2-q}(x)} d_k^2(x), \quad (86)$$

which leads to

$$x^{(\nu+1)} = \frac{\sum_{k=1}^m w_k \frac{1}{d_k^{2-q}(x^{(\nu)})} y_k}{\sum_{k=1}^m w_k \frac{1}{d_k^{2-q}(x^{(\nu)})}}. \quad (87)$$

For Equation 85, analogous with Equation 81, we set the derivative of Equation 85 equal to zero. Thus we solve

$$\sum_{k=1}^m w_k \frac{f'_k(d_k(x))}{d_k(x)} (x - y_k) = 0, \quad (88)$$

and the iteration becomes

$$x^{(\nu+1)} = \frac{\sum_{k=1}^m w_k \frac{f'(d_k(x^{(\nu)}))}{d_k(x^{(\nu)})} y_k}{\sum_{k=1}^m w_k \frac{f'(d_k(x^{(\nu)}))}{d_k(x^{(\nu)})}}. \quad (89)$$

The conditions in Katz (1969) on f needed for the convergence proof, and also the proof itself, are rather complicated. We can simply use the conditions of Theorem 4.5 that $f'(x)/x$ is non-increasing on the non-negatives reals to construct a quadratic majorization algorithm in which we minimize

$$\sum_{k=1}^m \sum_{k=1}^m w_k \frac{f'(d_k(x^{(\nu)}))}{d_k(x^{(\nu)})} d_k^2(x) \quad (90)$$

to find x^{k+1} . This gives directly the update Equation 87.

9.3 Sparse Recovery

This is a field which is difficult to delineate. A somewhat ad-hoc definition is recovering complete information from incomplete information, often in the context of specific engineering problems. There is overlap with signal detection, image analysis, matrix completion, ... But “sparse recovery” scientific activities “sparse recovery” could be extended far beyond these boundaries. Since classical statistics infers properties of the population from those of a sample it is a form of sparse recovery. Since science infers properties of the real world from outcomes of experiments it is sparse recovery too.

9.4 Multivariate Analysis

The smacof majorization method for multidimensional scaling was first presented at the *US-Japan Seminar on Theory, Methods and Applications of Multidimensional Scaling and Related Techniques* at UCSD in La Jolla, August 1975. Shortly after that I read the basic EM paper by Dempster, Laird, and Rubin (1977), and shortly after that I realized that smacof and EM were both special cases of a general minimization strategy, which I called majorization at the time. In June 1978 both Nan Laird and I attended the *Fifth International Symposium on Multivariate Analysis* at the University of Pittsburgh. I remember mentioning majorization, excitedly, to Nan on the conference bus.

The smacof majorization method was fully discussed in De Leeuw (1977), De Leeuw and Heiser (1977), and De Leeuw and Heiser (1980). The familiar picture illustrating two steps of the general majorization algorithm first appears in De Leeuw (1988a). But unlike EM, which took off as a rocket in 1977, the general idea of majorization remained unpublished, until De Leeuw (1994) and Heiser (1995). Majorization was used regularly in the Gifi project. The book Gifi (1990), which is a version of 1981 lecture notes, mentions majorization only once, but since then a stream of papers and dissertations from the Data Theory department in Leiden using majorization appeared. Heiser (1995) mentions most of them. In De Leeuw (1988b) another large majorization subfield, the *aspect approach* to multivariate analysis, was developed. In section 7 of that paper the general majorization/minorization approach to optimization is outlined, possibly for the first time in print.

Robust versions of low rank matrix approximation, a.k.a. principal component analysis, were first considered by Gabriel and Odoroff (1984). They start by discussing the alternating least squares algorithm for least squares weighted matrix approximation of Gabriel and Zamir (1979). The alternating is to compute new row scores for currently fixed column scores by linear regression, and then computing new column scores corresponding with the new row

scores, again by linear regression. Gabriel and Odoroff (1984) suggest to replace the linear least squares weighted averages in each of the two stages by medians or trimmed means to get a robust PCA. There is no sign of a convergence proof, but there is the suggestion to use alternating least absolute value methods to minimize the sum of absolute residuals of the matrix approximation. This suggestion was taken up by Verboon and Heiser (1994) using the majorization approach and the Huber and Tukey robust loss functions. Their robust PCA method is very similar to our robust MDS method, but the presentation of their method has some magical elements. The Huber and Tukey majorization functions are presented without any discussion where they came from, and it is then verified that they are indeed majorizations. There is clearly nothing wrong with this, but using our Theorem 4.5 gives a more general and more direct approach.

Heiser (1986) was the first to connect the Weiszfeld problem with correspondence analysis and multidimensional scaling, emphasizing the majorization aspects. As we have seen in Heiser (1987) and Heiser (1988) he constructed majorization algorithms for multidimensional scaling and correspondence analysis.

The IRLS approach to robustifying multivariate matrix approximation techniques could easily lead to a large and varied number of publications. There are some excellent examples making their way through the usual publication channels. I will just give two recent examples, with good bibliographies. They are Huber Principal Component Analysis (He et al. (2023)) and Cauchy Factor Analysis (Li (2024)).

10 Discussion

10.1 Bounding the Second Derivative

In some cases our basic theorems may not apply, but there may be an alternative way to majorize loss. In fact, this is classic quadratic bounding as in Vosz and Eckhardt (1980) or Böhning and Lindsay (1988). As before, we want to minimize $\sum \omega_k f(\delta_k - d_k(X))$, but now we suppose that there is a $K > 0$ such that $f''(x) \leq K$. We then have the majorization

$$f(\delta_k - d_k(X)) \leq f(\delta_k - d_k(Y)) + f'(\delta_k - d_k(Y))(d_k(Y) - d_k(X)) + \frac{1}{2}K(d_k(Y) - d_k(X))^2 \quad (91)$$

and in iteration k we minimize, or at least decrease,

$$\sum \omega_k [d_k(X) - \{d_k(X^{(\nu)}) - K^{-1}f'(\delta_k - d_k(X^{(\nu)}))\}]^2 \quad (92)$$

Note that in this algorithm the weights do not change. Instead of fitting a fixed target with moving weights, we fit a moving target with fixed weights.

We can apply bounding the second derivative, for example, to Charbonnier loss, using the inequality

$$f_c''(x) = (x^2 + c^2)^{-\frac{1}{2}} - x^2(x^2 + c^2)^{-\frac{3}{2}} \leq (x^2 + c^2)^{-\frac{1}{2}} \leq c^{-1}, \quad (93)$$

Of course this method requires that the second derivative exists at x . Although I have not done any comparisons it will probably require more iterations and take longer than the method in Section 5.1.

The paper by Vosz and Eckhardt (1980) deserves some special mention here.

$$\mathcal{D}^2\sigma(x) = \sum \frac{1}{d_i(x)} \left\{ I - \frac{(x - y_i)(x - y_i)'}{d_i^2(x)} \right\}$$

10.2 Fixed Weights

One could also consider using the fixed weights in regular non-robust smacof to achieve some form of robustness. Redefine stress as

$$\sigma(X) := \sum_k \omega_k f(\delta_k)(\delta_k - d_k(X))^2 \quad (94)$$

For example, we can choose a negative power for f , so that it downweights the large dissimilarities. If the dissimilarities is large, then it should have less influence on the fit, and thus on the solution X . This type of fixed power-weighting is used in various places (De Leeuw and Heiser (1980), Groenen and Van de Velden (2016)) to approximate loss functions such the one with logarithmic residuals in Ramsay (1977).

But we have to keep in mind that downweighting large dissimilarities is not the same thing as downweighting large residuals. The residuals depend on X , and it is perfectly possible that some small dissimilarities have large residuals. On the other hand emphasizing small dissimilarities in the loss function means that we want small dissimilarities to be fitted relatively well, which means that on average we want small dissimilarities to have small residuals. The Shepard plot will tend to fan out at the high end.

Despite these reservations, it will be useful to study if and how fixed weights can be used to improve robustness of smacof. If only because fixed weights correspond with a simpler and presumably more efficient algorithm.

10.3 Residual Definition

In our examples and in our code we use the residuals $\delta_k - d_k(X)$ are arguments of our loss functions. From the statistical point of view we have to remember, however, that most of these loss functions were designed for the robust estimation of a location parameter or a linear regression function. The error distributions were explicitly or implicitly assumed to be symmetric around zero, and defined on the whole real line, which was reflected in the fact that loss functions were even and had infinite support. In MDS, however, distances and dissimilarities are non-negative and reasonable error functions are not symmetric. One could follow the example of Ramsay (1977) and measure residuals as $\log \delta_{ij} - \log d_{ij}(X)$. This does not have any effect on the majorization of the loss functions, but it means that in the smacof step to find $X^{(\nu+1)}$ we have to minimize

$$\sigma(X) = \sum \omega_k(X^{(\nu)}) (\log \delta_{ij} - \log d_{ij}(X))^2,$$

which is considerably more complicated (De Leeuw, Groenen, and Mair (2016)).

10.4 Robust Nonmetric MDS

Our discussion and our software is all about metric MDS. It seems easy to extend the discussion to non-linear and non-metric MDS by adding an alternating least squares step optimally scaling the dissimilarities. This would take place between two majorizations of the robust loss function,

so one or more transformation and smacof steps can be taken between updating the weights. But this paradigm does not work for robust smacof.

Consider any hard redescender, such as Tukey or Hinich. At iteration ν , for current weights, first improve the configuration, then compute the optimal transformation of the dissimilarities, and then compute new weights. This is a recipe for disaster. At some point we minimize

$$\sigma(\hat{d}) = \sum \omega_k(X^{(\nu)})(\hat{d}_k - d_k(X^{(\nu+1)}))^2$$

over the disparities \hat{d} , which must be monotone with the dissimilarities. Because of the hard redescending some of the weights, for current absolute residuals larger than c , will be zero. The monotone regression is done for the observations with non-zero weights, and the disparities corresponding with zero weights are only determined by the order they are required to have. Thus they can be freely chosen in an interval between two disparities obtained from the monotone regression. That interval can be large, in fact if one of the zero weights corresponds with the largest dissimilarity it can be infinite. What we choose in the interval will determine the new residual and thus the next set of weights.

In the unfortunate situation that the current absolute residuals are all larger than c , even after choosing the optimal \hat{d} , the next weights will all be zero and the algorithm stops with zero stress.

10.5 Practicalities

Recommending one particular loss function from the many we have discussed is not easy. In some cases, for example for Cauchy loss, one can justify the choice of a loss function by assuming a particular error distribution and using the maximum likelihood principle. But in general perhaps the best way to proceed for a given MDS problem is to take what we could call a *trajectory approach*. Choose one particular parametric family, for example the Huber one, and compute the robust smacof solution $X(c)$ for a number of increasing positive c values. For small c we start with a close approximation of the LAV solution, increasing c will eventually take us to the LS solution. The starting point for computing each solution will be the solution for the previous c . We can plot the trajectory of the points in the configurations $X(c)$, and even make an animation. It seems that the Huber family is a good candidate for such a study, with the generalized Charbonnier a good second. If the main concern is to suppress the influence of outliers then trying some of the hard redescenders, such as the Tukey family, makes sense. Studying trajectories for some of robust loss functions is clearly interesting, but it is not something we can or will explore in this paper.

11 Code

The function `smacofRobust` has a parameter “engine”, which can be equal to `smacofCharbonnier`, `smacofGeneralizedCharbonnier`, `smacofBarron`, `smacofHuber`, `smacofTukey`, `smacofHinnich`, `smacofCauchy`, `smacofFair`, `smacofAndrews`, `smacofLogistic`, `smacofWelsch`, or `smacofGaussian`. These thirteen small modules compute the respective loss function values and weights for the IRLS procedure. This makes it easy for interested parties to add additional robust loss functions.

```
smacofRobust <- function(delta,
                          weights = 1 - diag(nrow(delta)),
                          ndim = 2,
                          xold = smacofTorgerson(delta, ndim),
                          engine = smacofAV,
                          cons = 0,
                          itmax = 1000,
                          eps = 1e-15,
                          verbose = TRUE) {
  nobj <- nrow(delta)
  wmax <- max(weights)
  dold <- as.matrix(dist(xold))
  h <- engine(nobj, weights, delta, dold, cons)
  rold <- h$resi
  wold <- h$wght
  sold <- h$strs
  itel <- 1
  repeat {
    vmat <- -wold
    diag(vmat) <- -rowSums(vmat)
    vinv <- solve(vmat + (1 / nobj)) - (1 / nobj)
    bmat <- -wold * delta / (dold + diag(nobj))
    diag(bmat) <- -rowSums(bmat)
    xnew <- vinv %*% (bmat %*% xold)
    dnew <- as.matrix(dist(xnew))
    h <- engine(nobj, weights, delta, dnew, cons)
    rnew <- h$resi
    wnew <- h$wght
    snew <- h$strs
    if (verbose) {
      cat(
```

```

    "itel ",
    formatC(itel, width = 4, format = "d"),
    "sold ",
    formatC(sold, digits = 10, format = "f"),
    "snew ",
    formatC(snew, digits = 10, format = "f"),
    "\n"
  )
}
if ((itel == itmax) || ((sold - snew) < eps)) {
  break
}
xold <- xnew
dold <- dnew
sold <- snew
wold <- wnew
rold <- rnew
itel <- itel + 1
}
return(list(
  x = xnew,
  s = snew,
  d = dnew,
  r = rnew,
  itel = itel
))
}

smacofTorgerson <- function(delta, ndim) {
  dd <- delta^2
  rd <- apply(dd, 1, mean)
  md <- mean(dd)
  sd <- -.5 * (dd - outer(rd, rd, "+") + md)
  ed <- eigen(sd)
  return(ed$vectors[, 1:ndim] %*% diag(sqrt(ed$values[1:ndim])))
}

smacofCharbonnier <- function(nobj, wmat, delta, dmat, cons) {
  resi <- sqrt((delta - dmat)^2 + cons)
  resi <- ifelse(resi < 1e-10, 2 * max(wmat), resi)
}

```



```

rmin <- sqrt(cons)
wght <- wmat / (resi + diag(nobj))
strs <- sum(wmat * resi) - rmin * sum(wmat)
return(list(
  resi = resi,
  wght = wght,
  strs = strs
))
}

smacofGeneralizedCharbonnier <- function(nobj, wmat, delta, dmat, cons) {
  resi <- ((delta - dmat) ^ 2 + cons[1]) ^ cons[2]
  rmin <- cons[1] ^ cons[2]
  wght <- wmat * ((delta - dmat) ^ 2 + cons[1] + diag(nobj)) ^ (cons[2] - 1)
  strs <- sum(wmat * resi) - rmin * sum(wmat)
  return(list(
    resi = resi,
    wght = wght,
    strs = strs
  ))
}

smacofBarron <- function(nobj, wmat, delta, dmat, cons) {
  f1 <- abs(cons[2] - 2) / cons[2]
  f2 <- (((delta - dmat) / cons[1]) ^ 2) / abs(cons[2] - 2) + 1)
  resi <- f1 * (f2 ^ (cons[2] / 2) - 1)
  wght <- wmat * f2 ^ (cons[2] / 2 - 1)
  strs <- sum(wmat * resi)
  return(list(
    resi = resi,
    wght = wght,
    strs = strs
  ))
}

smacofGauss <- function(nobj, wmat, delta, dmat, cons) {
  difi <- delta - dmat
  resi <- difi * (2 * pnorm(difi / cons) - 1) + 2 * cons * dnorm(difi / cons)
  rmin <- 2 * cons * dnorm(0)
  wght <- wmat * (pnorm(difi / cons) - 0.5) / (difi + diag(nobj))
}

```

```

    strs <- sum(wmat * resi) - rmin * sum(wmat)
    return(list(
      resi = resi,
      wght = wght,
      strs = strs
    ))
  }

smacofHuber <- function(nobj, wmat, delta, dmat, cons) {
  difi <- delta - dmat
  resi <- ifelse(abs(difi) < cons, (difi ^ 2) / 2, cons * abs(difi) - ((cons ^ 2) / 2))
  wght <- ifelse(abs(difi) < cons,
    wmat,
    wmat * sign(difi - cons) * cons / (difi + diag(nobj)))
  strs <- sum(wmat * resi)
  return(list(
    resi = resi,
    wght = wght,
    strs = strs
  ))
}

smacofTukey <- function(nobj, wmat, delta, dmat, cons) {
  cans <- (cons ^ 2) / 6
  difi <- delta - dmat
  resi <- ifelse(abs(difi) < cons, cans * (1 - (1 - (difi / cons)^2)^3), cans)
  wght <- wmat * ifelse(abs(difi) < cons, (1 - (difi / cons)^2)^2, 0)
  strs <- sum(wmat * resi)
  return(list(
    resi = resi,
    wght = wght,
    strs = strs
  ))
}

smacofCauchy <- function(nobj, wmat, delta, dmat, cons) {
  difi <- delta - dmat
  resi <- log(((difi / cons)^2 + 1))
  wght <- wmat * (1 / ((difi / cons)^2 + 1))
  strs <- sum(wmat * resi)
}

```

```

return(list(
  resi = resi,
  wght = wght,
  strs = strs
))
}

smacofWelsch <- function(nobj, wmat, delta, dmat, cons) {
  difi <- delta - dmat
  resi <- 1 - exp(-(difi / cons)^2)
  wght <- wmat * exp(-(difi / cons)^2)
  strs <- sum(wmat * resi)
  return(list(
    resi = resi,
    wght = wght,
    strs = strs
  ))
}

smacofAndrews <- function(nobj, wmat, delta, dmat, cons) {
  difi <- delta - dmat
  resi <- ifelse(abs(difi) < pi * cons,
                (cons ^ 2) * (1 - cos(x / cons)),
                2 * (cons^2))
  wght <- wmat * ifelse(abs(difi) < pi * cons, sin(x / cons) / (x / cons), 0)
  strs <- sum(wmat * resi)
  return(list(
    resi = resi,
    wght = wght,
    strs = strs
  ))
}

smacofHinich <- function(nobj, wmat, delta, dmat, cons) {
  difi <- delta - dmat
  resi <- ifelse(abs(difi) < cons, (difi^2) / 2, (cons^2) / 2)
  wght <- wmat * ifelse(abs(difi) < cons, 1, 0)
  strs <- sum(wmat * resi)
  return(list(
    resi = resi,

```

```

    wght = wght,
    strs = strs
  ))
}

smacofLogistic <- function(nobj, wmat, delta, dmat, cons) {
  difi <- delta - dmat
  resi <- (cons ^ 2) * log(cosh(x / cons))
  wght <- wmat * tanh(x / cons) / (x / cons)
  strs <- sum(wmat * resi)
  return(list(
    resi = resi,
    wght = wght,
    strs = strs
  ))
}

smacofFair <- function(nobj, wmat, delta, dmat, cons) {
  difi <- delta - dmat
  resi <- log((difi / cons)^2 + 1)
  wght <- wmat * (1 / ((difi / cons) ^ 2 + 1))
  strs <- sum(wmat * resi)
  return(list(
    resi = resi,
    wght = wght,
    strs = strs
  ))
}

```

References

- Aftab, K., and R. Hartley. 2015. “Convergence of Iteratively Re-Weighted Least Squares to Robust m-Estimators.” In *2015 IEEE Winter Conference on Applications of Computer Vision*, 480–87. <https://doi.org/10.1109/WACV.2015.70>.
- Aftab, K., R. Hartley, and J. Trumpf. 2015. “Generalized Weiszfeld Algorithms for Lq Optimization.” *IEEE Transactions on Pattern Analysis and Machine Intelligence* 37 (4): 728–44. <https://doi.org/10.1109/TPAMI.2014.2353625>.
- Andrews, D. F., P. J. Bickel, F. R. Hampel, P. J. Huber, W. H. Rogers, and J. W. Tukey. 1972. *Robust Estimators of Location: Survey and Advances*. Princeton University Press.
- Barron, J. T. 2019. “A General and Adaptive Robust Loss Function.” In *Proceedings 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4331–39. https://openaccess.thecvf.com/content_CVPR_2019/papers/Barron_A_General_and_Adaptive_Robust_Loss_Function_CVPR_2019_paper.pdf.
- Beaton, A. E., and W. Tukey J. 1974. “The Fitting of Power Series, Meaning Polynomials, Illustrated on Band-Spectroscopic Data.” *Technometrics* 16 (147–185).
- Beck, A., and S. Sabach. 2015. “Weiszfeld’s Method: Old and New Results.” *Journal of Optimization Theory and Applications* 164: 1–40.
- Black, M. J., and P. Anandan. 1996. “The Robust Estimation of Multiple Motions: Parametric and Piecewise-Smooth Flow Fields.” *Computer Vision and Image Understanding* 63 (1): 75–104.
- Böhning, D., and B. G. Lindsay. 1988. “Monotonicity of Quadratic-approximation Algorithms.” *Annals of the Institute of Statistical Mathematics* 40 (4): 641–63.
- Border, K. C. 2018. “Supergradients.” <https://healy.econ.ohio-state.edu/kcb/Notes/Supergrad.pdf>.
- Candes, E. J., and T. Tao. 2005. “Decoding by Linear Programming.” *IEEE Transactions on Information Theory* 51 (12): 4203–15.
- Candes, E. J., M. B. Wakin, and S. P. Boyd. 2008. “Enhancing Sparsity by Reweighted l_1 Minimization.” *Journal of Fourier Analysis and Applications* 14: 877–905. <https://doi.org/10.1007/s00041-008-9045-x>.
- Charbonnier, P., L. Blanc-Feraud, G. Aubert, and M. Barlaud. 1994. “Two deterministic half-quadratic regularization algorithms for computed imaging.” *Proceedings of 1st International Conference on Image Processing 2*: 168–72. <https://doi.org/10.1109/icip.1994.413553>.
- Coleman, D., P. Holland, N. Kaden, V. Klema, and S. C. Peters. 1980. “A System of Subroutines for Iteratively Reweighted Least Squares Computations.” *ACM Transactions on Mathematical Software* 6 (3): 327–36.
- De Gruijter, D. N. M. 1967. “The Cognitive Structure of Dutch Political Parties in 1966.” Report E019-67. Psychological Institute, University of Leiden.
- De Leeuw, J. 1977. “Applications of Convex Analysis to Multidimensional Scaling.” In

- Recent Developments in Statistics*, edited by J. R. Barra, F. Brodeau, G. Romier, and B. Van Cutsem, 133–45. Amsterdam, The Netherlands: North Holland Publishing Company.
- . 1984. “Differentiability of Kruskal’s Stress at a Local Minimum.” *Psychometrika* 49: 111–13.
- . 1988a. “Convergence of the Majorization Method for Multidimensional Scaling.” *Journal of Classification* 5: 163–80.
- . 1988b. “Multivariate Analysis with Optimal Scaling.” In *Proceedings of the International Conference on Advances in Multivariate Statistical Analysis*, edited by S. Das Gupta and J. K. Ghosh, 127–60. Calcutta, India: Indian Statistical Institute.
- . 1994. “Block Relaxation Algorithms in Statistics.” In *Information Systems and Data Analysis*, edited by H. H. Bock, W. Lenski, and M. M. Richter, 308–24. Berlin: Springer Verlag. <https://jansweb.netlify.app/publication/deleeuw-c-94-c/deleeuw-c-94-c.pdf>.
- . 2018. “MM Algorithms for Smoothed Absolute Values.” 2018. <https://jansweb.netlify.app/publication/deleeuw-e-18-f/deleeuw-e-18-f.pdf>.
- De Leeuw, J., P. Groenen, and P. Mair. 2016. “Minimizing rStress Using Majorization.” 2016. <https://jansweb.netlify.app/publication/deleeuw-groenen-mair-e-16-a/deleeuw-groenen-mair-e-16-a.pdf>.
- De Leeuw, J., and W. J. Heiser. 1977. “Convergence of Correction Matrix Algorithms for Multidimensional Scaling.” In *Geometric Representations of Relational Data*, edited by J. C. Lingoes, 735–53. Ann Arbor, Michigan: Mathesis Press.
- . 1980. “Multidimensional Scaling with Restrictions on the Configuration.” In *Multivariate Analysis, Volume V*, edited by P. R. Krishnaiah, 501–22. Amsterdam, The Netherlands: North Holland Publishing Company.
- De Leeuw, J., and K. Lange. 2009. “Sharp Quadratic Majorization in One Dimension.” *Computational Statistics and Data Analysis* 53: 2471–84.
- De Leeuw, J., and P. Mair. 2009. “Multidimensional Scaling Using Majorization: SMACOF in R.” *Journal of Statistical Software* 31 (3): 1–30. <https://www.jstatsoft.org/article/view/v031i03>.
- Dempster, A. P., N. M. Laird, and D. B. Rubin. 1977. “Maximum Likelihood for Incomplete Data via the EM Algorithm.” *Journal of the Royal Statistical Society* B39: 1–38.
- Dennis Jr, J. E., and R. E. Welsch. 1978. “Techniques for Nonlinear Least Squares and Robust Regression.” *Communications in Statistics - Simulation and Computation* 7 (4): 345–59.
- Donoho, D. L., and M. Elad. 2003. “Optimally Sparse Representation in General (Nonorthogonal) Dictionaries via ℓ_1 Minimization.” *Proceedings of the National Academy of Sciences* 100 (5): 2197–2202.
- Franksen, O. I., and I. Grattan-Guinness. 1989. “The Earliest Contribution to Location Theory ? Spatio-Economic Equilibrium with Lamé and Clapeyron.” *Mathematics and Computers in Simulation* 3w1: 195–220.
- Gabriel, K. R., and Ch. L. Odoroff. 1984. “Resistant Lower Rank Approximation of Matrices.” In *Data Analysis and Informatics*, edited by E. Diday, M. Jambu, L. Lebart, J. Pages, and R. Tomassone, 3:23–30. North Holland Publishing Company.

- Gabriel, K. R., and S. Zamir. 1979. “Lower Rank Approximation of Matrices by Least Squares with Any Choize of Weights.” *Technometrics* 21 (4): 489–98.
- Gifi, A. 1990. *Nonlinear Multivariate Analysis*. New York, N.Y.: Wiley.
- Groenen, P. J. F., P. Giaquinto, and H. A. L. Kiers. 2003. “Weighted Majorization Algorithms for Weighted Least Squares Decomposition Models.” Econometric Institute Report EI 2003-09. Econometric Institute, Erasmus University Rotterdam. <https://repub.eur.nl/pub/1700>.
- Groenen, P. J. F., W. J. Heiser, and J. J. Meulman. 1999. “Global Optimization in Least-Squares Multidimensional Scaling by Distance Smoothing.” *Journal of Classification* 16: 225–54.
- Groenen, P. J. F., and M. Van de Velden. 2016. “Multidimensional Scaling by Majorization: A Review.” *Journal of Statistical Software* 73 (8): 1–26. <https://www.jstatsoft.org/index.php/jss/article/view/v073i08>.
- He, Y., L. Li, D. Liu, and W.-Z. Zhou. 2023. “Huber Principal Component Analysis for Large-Dimensional Factor Models.” <https://arxiv.org/abs/2303.02817>.
- Heiser, W. J. 1986. “A Majorization Algorithm for the Reciprocal Location Problem.” RR-86-12. Department of Data Theory, University of Leiden.
- . 1987. “Correspondence Analysis with Least Absolute Residuals.” *Computational Statistics and Data Analysis* 5: 337–56.
- . 1988. “Multidimensional Scaling with Least Absolute Residuals.” In *Classification and Related Methods of Data Analysis*, edited by H. H. Bock, 455–62. North-Holland Publishing Co.
- . 1995. “Convergent Computing by Iterative Majorization: Theory and Applications in Multidimensional Data Analysis.” In *Recent Advantages in Descriptive Multivariate Analysis*, edited by W. J. Krzanowski, 157–89. Oxford: Clarendon Press.
- Hinich, M. J., and P. P. Talwar. 1975. “A Simple Method for Robust Regression.” *Journal of the American Statistical Association* 70: 113–19.
- Holland, P. W., and R. E. Welsch. 1977. “Robust Regression Using Iteratively Reweighted Least-Squares.” *Communications in Statistics - Theory and Methods* 6 (9): 813–27. <https://doi.org/10.1080/03610927708827533>.
- Huber, P. J. 1964. “Robust Estimation of a Location Parameter.” *Annals of Mathematical Statistics* 35 (1): 73–101.
- Hunter, D. R., and R. Li. 2005. “Variable Selection Using MM Algorithms.” *The Annals of Statistics* 33: 1617–42.
- Jaakkola, T. S., and M. I. Jordan. 2000. “Bayesian Parameter Estimation via Variational Methods.” *Statistics and Computing* 10: 25–37.
- Katz, I. N. 1969. “On the Convergence of a Numerical Scheme for Solving Some Locational Equilibrium Problems.” *SIAM Journal on Applied Mathematics* 17 (6): 1224–31.
- Lange, K. 2016. *MM Optimization Algorithms*. SIAM.
- Li, J. 2024. “Robust Matrix Factor Analysis Method with Adaptive Parameter Adjustment Using Cauchy Weighting.” *Computational Statistics*. <https://doi.org/10.1007/s00180-024-01548-4>.
- Mlotshwa, T., H. Van Deventer, and Sergeevna Bosman A. 2023. “Cauchy Loss Function:

- Robustness Under Gaussian and Cauchy Noise.” <https://arxiv.org/abs/2302.07238>.
- Mordukhovich, B. S., and N. M. Nam. 2019. “The Fermat-Torricelli Problem and Weiszfeld’s Algorithm in the Light of Convex Analysis.” *Journal of Applied Numerical Optimization* 1 (3): 205–19. <https://doi.org/https://doi.org/10.23952/jano.1.2019.3.02>.
- Phillips, R. F. 2002. “Least Absolute Deviations Estimation via the EM Algorithm.” *Statistics and Computing* 12: 281–85.
- Plastria, F. 2011. “The Weiszfeld Algorithm: Proof, Amendments and Extensions.” In *Foundations of Location Analysis*, edited by H. A. Eiselt and V. Marianov, 155:357–89. International Series in Operations Research and Management Science. Springer.
- Pliner, V. 1996. “Metric Unidimensional Scaling and Global Optimization.” *Journal of Classification* 13: 3–18.
- R Core Team. 2024. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Ramirez, C., R. Sanchez, V. Kreinovich, and M. Arguez. 2014. “ $\sqrt{x^2 + \mu}$ is the Most Computationally Efficient Smooth Approximation to $|x|$.” *Journal of Uncertain Systems* 8: 205–10.
- Ramsay, J. O. 1977. “Maximum Likelihood Estimation in Multidimensional Scaling.” *Psychometrika* 42: 241–66.
- Rothkopf, E. Z. 1957. “A Measure of Stimulus Similarity and Errors in some Paired-associate Learning.” *Journal of Experimental Psychology* 53: 94–101.
- Schlossmacher, E. J. 1973. “An Iterative Technique for Absolute Deviations Curve Fitting.” *Journal of the American Statistical Association* 68: 857–59.
- Sturm, R. 1884. “Ueber Den Punkt Kleinster Entfernungssumme von Gegebenen Punkten.” *Journal für Die Reine Und Angewandte Mathematik* 97: 49–61.
- Van Ruitenburg, J. 2005. “Algorithms for Parameter Estimation in the Rasch Model.” Measurement and Research Department Reports 2005-04. Arnhem, Netherlands: CITO. https://www.researchgate.net/publication/355568984_Algorithms_for_parameter_estimation_in_the_Rasch_model_Measurement_and_Research_Report_05-04#fullTextFileContent.
- Verboon, P., and W. J. Heiser. 1994. “Resistant Lower Rank Approximation of Matrices by Iterative Majorization.” *Computational Statistics and Data Analysis* 18: 457–67.
- Voronin, S., G. Ozkaya, and Y. Yoshida. 2014. “Convolution Based Smooth Approximations to the Absolute Value Function with Application to Non-smooth Regularization.” 2014. <https://arxiv.org/abs/1408.6795>.
- Vosz, H., and U. Eckhardt. 1980. “Linear Convergence of Generalized Weiszfeld’s Method.” *Computing* 25: 243–51.
- Weiszfeld, E. 1937. “Sur le Point par lequel la Somme des Distances de n Points Donnés est Minimum.” *Tohoku Mathematics Journal* 43: 355–86.
- Weiszfeld, E., and F. Plastria. 2009. “On the Point for Which the Sum of the Distances to n Given Points Is Minimum.” *Annals of Operations Research* 167: 7–41.