

## Chapter 7

# Fitting in other Norms

In this chapter we discuss methods for parameter estimation in cases, where the definition of “best” is different from the (weighted) least squares, cf Section 5.1. More precisely, we are given a vector function  $\mathbf{r} : \mathbb{R}^n \mapsto \mathbb{R}^m$  and we want to find  $\hat{\mathbf{x}}$  that minimizes some measure of  $\mathbf{r}(\mathbf{x})$ , for instance  $\|\mathbf{r}(\mathbf{x})\|_1$  or  $\|\mathbf{r}(\mathbf{x})\|_\infty$ . These two cases are treated in the first two sections, and in the rest of the chapter we discuss a method which can be thought of as a hybrid between the least squares and the least absolute deviation definitions of “best”.

The function  $\mathbf{r}$  may depend nonlinearly on  $\mathbf{x}$ , and we present algorithms of Gauss–Newton type, cf Section 6.1, based on successive approximations by first order Taylor expansions,

$$\begin{aligned} r_i(\mathbf{x}+\mathbf{h}) &\simeq \ell_i(\mathbf{h}) \equiv r_i(\mathbf{x}) + \nabla r_i(\mathbf{x})^T \mathbf{h}, \quad i = 1, \dots, m, \\ \Downarrow \\ \mathbf{r}(\mathbf{x}+\mathbf{h}) &\simeq \boldsymbol{\ell}(\mathbf{h}) \equiv \mathbf{r}(\mathbf{x}) + \mathbf{J}(\mathbf{x})\mathbf{h}, \end{aligned} \tag{7.1}$$

where  $\nabla r_i \in \mathbb{R}^n$  is the gradient of  $r_i$  and  $\mathbf{J} \in \mathbb{R}^{m \times n}$  is the Jacobian of  $\mathbf{r}$ . It has the rows  $\mathbf{J}_{i:} = (\nabla r_i)^T$ . In the first two sections we combine the Gauss–Newton method with a trust region approach, cf Section 2.4. The generic algorithm is

### Algorithm 7.2. Generic algorithm

Given starting point  $\mathbf{x}$  and trust region radius  $\Delta$

**while not STOP**

$\hat{\mathbf{h}} = \operatorname{argmin}_{\|\mathbf{h}\|_\infty \leq \Delta} L(\mathbf{h})$

$\varrho = ((f(\mathbf{x}) - f(\mathbf{x} + \hat{\mathbf{h}})) / (L(\mathbf{0}) - L(\hat{\mathbf{h}})))$

**if**  $\varrho \geq 0.75$  **then**  $\Delta := \Delta * 2$

**if**  $\varrho \leq 0.25$  **then**  $\Delta := \Delta / 3$

**if**  $\varrho > 0$  **then**  $\mathbf{x} := \mathbf{x} + \hat{\mathbf{h}}$

**end**

The functions  $f$  and  $L$  are defined by  $f(\mathbf{x}) = \|\mathbf{r}(\mathbf{x})\|_p$  and  $L(\mathbf{h}) = \|\ell(\mathbf{h})\|_p$ , respectively. Combining these definitions with (7.1) we see that  $L(\mathbf{0}) = f(\mathbf{x})$ .

### 7.1. Fitting in the $L_\infty$ norm

We seek a minimizer  $\hat{\mathbf{x}}$  of the function

$$f(\mathbf{x}) = \|\mathbf{r}(\mathbf{x})\|_\infty = \max_i |r_i(\mathbf{x})| .$$

Hald and Madsen [25] proposed to use Algorithm 7.2. The linearized problem

$$\hat{\mathbf{h}} = \underset{\|\mathbf{h}\|_\infty \leq \Delta}{\operatorname{argmin}} \{ L(\mathbf{h}) \equiv \|\mathbf{r}(\mathbf{x}) + \mathbf{J}(\mathbf{x}) \mathbf{h}\|_\infty \}$$

is solved by an LP algorithm, cf Appendix A.9. To get an LP formulation of the problem, we introduce the variable  $h_{n+1} = L(\mathbf{h})$  and the extended vector

$$\tilde{\mathbf{h}} = \begin{pmatrix} \mathbf{h} \\ h_{n+1} \end{pmatrix} \in \mathbb{R}^{n+1} .$$

Now  $\hat{\mathbf{h}}$  can be found by solving the problem

$$\begin{array}{ll} \text{minimize} & h_{n+1} \\ \text{subject to} & \left. \begin{array}{l} -\Delta \leq h_i \leq \Delta \\ -h_{n+1} \leq r_i(\mathbf{x}) + \nabla r_i(\mathbf{x})^T \mathbf{h} \leq h_{n+1} \end{array} \right\} \quad i = 1, \dots, n \end{array}$$

See [25] for further information. This simple linear programming formulation is enabled by the change of the trust region definition from the 2-norm discussed in Section 2.4 to the  $\infty$ -norm used here.

### 7.2. Fitting in the $L_1$ norm

We seek a minimizer  $\hat{\mathbf{x}}$  of the function

$$f(\mathbf{x}) = \|\mathbf{r}(\mathbf{x})\|_1 = \sum_{i=1}^m |r_i(\mathbf{x})| .$$

In some applications the solution is referred to as the *least absolute deviation* fit.

Hald and Madsen [26] proposed to use Algorithm 7.2. The linearized problem

$$\hat{\mathbf{h}} = \underset{\|\mathbf{h}\|_\infty \leq \Delta}{\operatorname{argmin}} \left\{ L(\mathbf{h}) \equiv \|\mathbf{r}(\mathbf{x}) + \mathbf{J}(\mathbf{x}) \mathbf{h}\|_1 \right\}$$

is solved by an LP algorithm: We introduce auxiliary variables  $h_{n+i}$ ,  $i = 1, \dots, m$ , and the extended vector

$$\tilde{\mathbf{h}} = \begin{pmatrix} \mathbf{h} \\ h_{n+1} \\ \vdots \\ h_{n+m} \end{pmatrix} \in \mathbb{R}^{n+m}.$$

Now  $\hat{\mathbf{h}}$  can be found by solving the problem

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m h_{n+i} \\ & \text{subject to} && \left. \begin{aligned} & -\Delta \leq h_i \leq \Delta \\ & -h_{n+i} \leq r_i(\mathbf{x}) + \nabla r_i(\mathbf{x})^T \mathbf{h} \leq h_{n+i} \end{aligned} \right\} \quad i = 1, \dots, m \end{aligned}$$

See [26] for further information.

**Example 7.1.** When  $\mathbf{r}$  is an affine function,  $\mathbf{r}(\mathbf{x}) = \mathbf{y} - \mathbf{F}\mathbf{x}$ , the gradients are constant, and the LP formulation is

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m x_{n+i} \\ & \text{subject to} && -x_{n+i} \leq y_i - \mathbf{F}_{i,:} \mathbf{x} \leq x_{n+i}, \quad i = 1, \dots, m. \end{aligned}$$

This formulation is the background for the *Barrodale–Roberts* algorithm [2], which is often used for robust parameter estimation with linear fitting models, cf Examples 5.5 and 7.3. ■

### 7.3. Huber estimation

This approach combines the smoothness of the least squares estimator with the robustness of the  $L_1$ -estimator, cf Example 5.5.

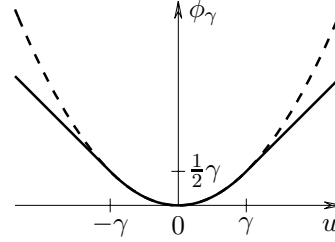
For a function  $\mathbf{r} : \mathbb{R}^n \mapsto \mathbb{R}^m$  the *Huber estimator*  $\mathbf{x}_\gamma$  is defined by<sup>1)</sup>

$$\mathbf{x}_\gamma = \underset{\mathbf{x}}{\operatorname{argmin}} \{ f_\gamma(\mathbf{x}) \equiv \sum_{i=1}^m \phi_\gamma(r_i(\mathbf{x})) \} , \quad (7.3)$$

where  $\phi_\gamma$  is the *Huber function*,

$$\phi_\gamma(u) = \begin{cases} \frac{1}{2\gamma} u^2 & \text{if } |u| \leq \gamma , \\ |u| - \frac{1}{2} \gamma & \text{if } |u| > \gamma . \end{cases} \quad (7.4)$$

**Figure 7.1.** *Huber function (full line) and the scaled L<sub>2</sub> function  $\frac{1}{2\gamma}u^2$  (dotted line).*



The threshold  $\gamma$  is used to distinguish between “small” and “large” function values (residuals). Based on the values of the  $r_i(\mathbf{x})$  we define a generalized sign vector  $\mathbf{s} = \mathbf{s}_\gamma(\mathbf{x})$  and an “activity matrix”  $\mathbf{W} = \mathbf{W}_\gamma(\mathbf{x}) = \operatorname{diag}(w_1, \dots, w_m)$ . Note that  $w_i = 1 - s_i^2$ .

	$r_i(\mathbf{x}) < -\gamma$	$ r_i(\mathbf{x})  \leq \gamma$	$r_i(\mathbf{x}) > \gamma$
$s_i(\mathbf{x})$	-1	0	1
$w_i(\mathbf{x})$	0	1	0

(7.5)

Now the objective function in (7.3) can be expressed as

$$f_\gamma(\mathbf{x}) = \frac{1}{2\gamma} \mathbf{r}^T \mathbf{W} \mathbf{r} + \mathbf{r}^T \mathbf{s} - \frac{1}{2} \gamma \mathbf{s}^T \mathbf{s} , \quad (7.6)$$

where we have omitted the argument  $(\mathbf{x})$  and index  $\gamma$  on the right-hand side. The gradient is

$$\nabla f_\gamma = \frac{1}{\gamma} \mathbf{J}^T (\mathbf{W} \mathbf{r} + \gamma \mathbf{s}) . \quad (7.7)$$

---

<sup>1)</sup> Strictly speaking, it is misleading to discuss this approach under the heading “other norms”:  $f_\gamma$  in (7.3) is not a norm: the triangle inequality is not satisfied.

At a minimizer the gradient must be zero. Thus, a necessary condition for  $\mathbf{x}_\gamma$  being a minimizer of  $f_\gamma$  is that it satisfies the equation

$$\mathbf{J}(\mathbf{x})^T (\mathbf{W}_\gamma(\mathbf{x})\mathbf{r}(\mathbf{x}) + \gamma \mathbf{s}_\gamma(\mathbf{x})) = \mathbf{0} . \quad (7.8)$$

### 7.3.1. Linear Huber estimation

First consider an affine function

$$\mathbf{r}(\mathbf{x}) = \mathbf{y} - \mathbf{F} \mathbf{x} , \quad (7.9)$$

where  $\mathbf{y} \in \mathbb{R}^m$  and  $\mathbf{F} \in \mathbb{R}^{m \times n}$  are constant. It follows from (7.5) that  $s_i(\mathbf{x}) = 0$  for all  $\mathbf{x}$  between the two hyperplanes  $r_i(\mathbf{x}) = -\gamma$  and  $r_i(\mathbf{x}) = \gamma$ . These  $2m$  hyperplanes divide  $\mathbb{R}^n$  into subregions  $\{D_k\}$ . Inside each subregion  $\mathbf{s}_\gamma(\mathbf{x})$  and  $\mathbf{W}_\gamma(\mathbf{x})$  are constant, so  $f_\gamma$  is a *piecewise quadratic*. The gradient is

$$\nabla f_\gamma(\mathbf{x}) = -\frac{1}{\gamma} \mathbf{F}^T (\mathbf{W}(\mathbf{x})(\mathbf{y} - \mathbf{F} \mathbf{x}) + \gamma \mathbf{s}(\mathbf{x})) . \quad (7.10)$$

It varies continuously across the hyperplanes, while the Hessian

$$\nabla^2 f_\gamma(\mathbf{x}) = \frac{1}{\gamma} \mathbf{F}^T \mathbf{W}(\mathbf{x}) \mathbf{F} \equiv \mathbf{H}(\mathbf{x}) \quad (7.11)$$

is constant in the interior of each  $D_k$  and jumps as  $\mathbf{x}$  crosses one of the dividing hyperplanes.

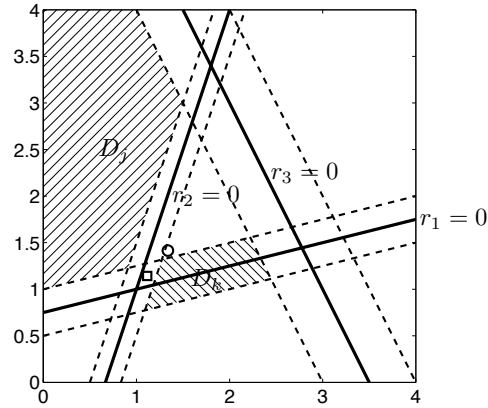
**Example 7.2.** Consider the affine function  $\mathbf{r} : \mathbb{R}^2 \mapsto \mathbb{R}^3$  given by

$$\mathbf{r}(\mathbf{x}) = \begin{pmatrix} 1.5 \\ 2.0 \\ -3.5 \end{pmatrix} - \begin{pmatrix} -0.5 & 2.0 \\ 3.0 & -1.0 \\ -1.0 & 0.5 \end{pmatrix} \mathbf{x} ,$$

and let the threshold  $\gamma = 0.5$ .

In  $\mathbb{R}^2$  a hyperplane is a straight line, and Figure 7.2 below shows the lines, along which  $r_i(\mathbf{x}) = 0$  (full lines); the dividing hyperplanes,  $r_i(\mathbf{x}) = \pm\gamma$  (dotted lines) and two of the subregions are indicated by hatching. For  $\mathbf{x} \in D_k$  we have  $\mathbf{s} = (0 \ -1 \ -1)^T$ , while  $\mathbf{s} = (-1 \ 1 \ -1)^T$  for  $\mathbf{x} \in D_j$ .

**Figure 7.2.** Dividing hyperplanes in  $\mathbb{R}^2$ .



The square indicates the Huber solution for  $\gamma = 0.5$ ,  $\mathbf{x}_\gamma = (1.116 \ 1.143)^T$  with  $\mathbf{s}_\gamma = (0 \ 0 \ -1)^T$ . For comparison, the circle marks the least squares solution to  $\mathbf{r}(\mathbf{x}) \simeq \mathbf{0}$ :  $\mathbf{x}_{(2)} = (1.337 \ 1.415)^T$ . ■

With a slight modification of (5.13) we can show that the Hessian (7.11) is positive semidefinite. This implies that if we find a stationary point  $\mathbf{x}_s$ , ie  $\nabla f_\gamma(\mathbf{x}_s) = \mathbf{0}$ , then it is a minimizer for  $f_\gamma$ :  $\mathbf{x}_\gamma = \mathbf{x}_s$ . We use Newton's method with line search to find a stationary point.

**Algorithm 7.12. Linear Huber**

Given starting point  $\mathbf{x}$  and threshold  $\gamma$

**repeat**

Find  $\mathbf{v}$  by solving  $\mathbf{H}(\mathbf{x})\mathbf{v} = -\nabla f_\gamma(\mathbf{x})$

$\mathbf{x} := \text{line\_search}(\mathbf{x}, \mathbf{v})$

**until** STOP

The line search is very simple:  $\varphi(\alpha) = f_\gamma(\mathbf{x} + \alpha\mathbf{v})$  is a piecewise quadratic in the scalar variable  $\alpha$ , so  $\alpha^*$  is a root for  $\varphi'$ , which is piecewise linear; the coefficients change when  $\mathbf{x} + \alpha\mathbf{v}$  passes a dividing hyperplane for increasing  $\alpha$ . The stopping criterion is that  $\mathbf{s}_\gamma(\mathbf{x} + \alpha^*\mathbf{v}) = \mathbf{s}_\gamma(\mathbf{x})$ . See [33] and [38] about details and efficient implementation.

**Example 7.3.** The function (7.9) may be the residual from data fitting with a linear model, cf Chapter 5. *Peter Huber* [29] introduced his M-estimator as a method for reducing the effect of wild points, cf Example 5.5, and from (7.8) we see that this is achieved by replacing a large  $r_i$  by just the threshold  $\gamma$  times the *sign* of this wild point residual. ■

**Example 7.4.** Suppose that  $\gamma$  is changed to  $\delta$  and that  $\mathbf{s}_\delta(\mathbf{x}_\delta) = \mathbf{s}_\gamma(\mathbf{x}_\gamma) = \mathbf{s}$ , and therefore  $\mathbf{W}_\delta(\mathbf{x}_\delta) = \mathbf{W}_\gamma(\mathbf{x}_\gamma) = \mathbf{W}$ . Then it follows from (7.8) and (7.10) that the Huber solution is a linear function of the threshold,

$$\mathbf{x}_\delta = \mathbf{x}_\gamma + (\gamma - \delta) \left( \mathbf{F}^T \mathbf{W} \mathbf{F} \right)^{-1} \mathbf{F}^T \mathbf{s} . \quad (7.13)$$

Let  $\mathbf{x}_{(2)}$  denote the least squares solution to  $\mathbf{r}(\mathbf{x}) \simeq \mathbf{0}$ . This is also the Huber solution (with  $\mathbf{s} = \mathbf{0}$ ) for all  $\gamma \geq \gamma_\infty \equiv \|\mathbf{r}(\mathbf{x}_{(2)})\|_\infty$ . For  $\gamma < \gamma_\infty$  one or more  $r_i$  will be “large”, so  $\mathbf{s}$  changes.

Madsen and Nielsen [34] showed that there exists a  $\gamma_0 > 0$  such that  $\mathbf{s}_\gamma(\mathbf{x}_\gamma^*)$  is constant for  $\gamma \leq \gamma_0$  and used this together with (7.13) to develop an efficient method for linear  $L_1$  estimation. This might replace the LP algorithm used to solve the linearized problems in Section 7.2. ■

### 7.3.2. Nonlinear Huber Estimation

Now consider a function  $\mathbf{r}$  that depends nonlinearly on  $\mathbf{x}$ . We must use iteration to find the minimizer  $\mathbf{x}_\gamma$  of  $f_\gamma$ . At the current iterate  $\mathbf{x}$  we use the approximation (7.1) and the corresponding approximation to the objective function

$$\begin{aligned} f_\gamma(\mathbf{x} + \mathbf{h}) &\simeq L(\mathbf{h}) \\ &= \frac{1}{2\gamma} \boldsymbol{\ell}^T \mathbf{W} \boldsymbol{\ell} + \boldsymbol{\ell}^T \mathbf{s} - \frac{1}{2} \gamma \mathbf{s}^T \mathbf{s} , \end{aligned}$$

where  $\boldsymbol{\ell} = \boldsymbol{\ell}(\mathbf{h}) = \mathbf{f}(\mathbf{x}) + \mathbf{J}(\mathbf{x})\mathbf{h}$  and  $\mathbf{s}$  and  $\mathbf{W}$  are given by (7.5) with  $\mathbf{r}(\mathbf{x})$  replaced by  $\boldsymbol{\ell}(\mathbf{h})$ .

As in Sections 7.1 and 7.2 we can combine this Gauss–Newton model with a trust region approach, cf [27]. Instead we shall describe a Levenberg–Marquardt like algorithm, cf 6.18, see Algorithm 7.14.

The linearized problem at  $1^\circ$  is solved by a slight modification of Algorithm 7.12 with starting point  $\mathbf{v} = \mathbf{0}$ . Note that  $L(\mathbf{0}) = f_\gamma(\mathbf{x})$ .

**Algorithm 7.14. Nonlinear Huber**

Given starting point  $\mathbf{x}$ , threshold  $\gamma$  and  $\mu > 0$ .  $\nu = 2$

**while not STOP**

$$\hat{\mathbf{h}} = \operatorname{argmin}_{\mathbf{h}} \{L(\mathbf{h}) + \frac{1}{2}\mu \mathbf{h}^T \mathbf{h}\} \quad \{1^\circ\}$$

$$\varrho := (f_\gamma(\mathbf{x}) - f_\gamma(\mathbf{x} + \hat{\mathbf{h}})) / (L(\mathbf{0}) - L(\hat{\mathbf{h}}))$$

**if**  $\varrho > 0$  {step acceptable}

$$\mathbf{x} := \mathbf{x} + \hat{\mathbf{h}}$$

$$\mu := \mu * \max\{\frac{1}{3}, 1 - (2\varrho - 1)^3\}; \quad \nu := 2$$

**else**

$$\mu := \mu * \nu; \quad \nu := 2 * \nu$$

**end**

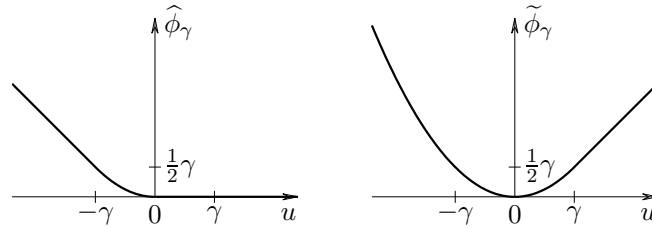
**end**

**end**

Finally, we give two examples of modified Huber functions with reduced influence from positive components of  $\mathbf{r}$ .

$$\hat{\phi}_\gamma(u) = \begin{cases} |u| - \frac{1}{2}\gamma & \text{for } u < -\gamma, \\ \frac{1}{2\gamma} u^2 & \text{for } -\gamma \leq u \leq 0, \\ 0 & \text{for } u > 0. \end{cases}$$

$$\tilde{\phi}_\gamma(u) = \begin{cases} \frac{1}{2\gamma} u^2 & \text{for } u \leq \gamma, \\ u - \frac{1}{2}\gamma & \text{for } u > \gamma. \end{cases}$$



**Figure 7.3.** One-sided Huber functions. one-sided Huber function

See [39] about implementation and some applications of these functions.