

# ROBUST REGRESSION 25

## CHAPTER CONTENTS

Nonrobustness of $\ell_2$ -Loss Minimization .....	279
$\ell_1$ -Loss Minimization .....	280
Huber Loss Minimization .....	282
Definition .....	282
Stochastic Gradient Algorithm .....	283
Iteratively Reweighted LS .....	283
$\ell_1$ -Constrained Huber Loss Minimization .....	286
Tukey Loss Minimization .....	290

Although LS is useful in various practical applications, it is sensitive to *outliers*, samples containing large noise. In this chapter, alternative learning methods that possess high robustness against outliers are introduced.

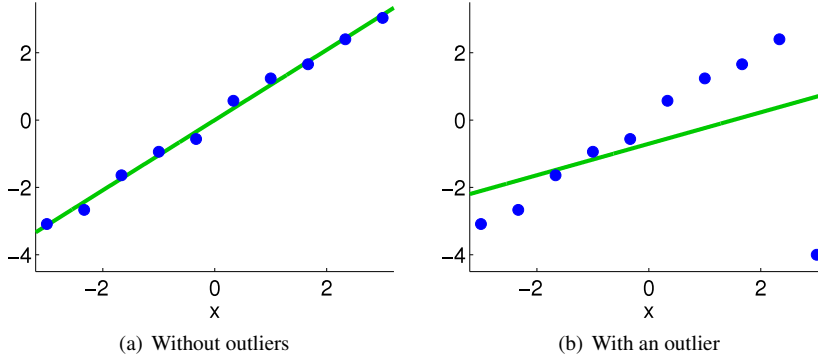
## 25.1 NONROBUSTNESS OF $\ell_2$ -LOSS MINIMIZATION

A LS solution for straight-line model  $f_{\theta}(x) = \theta_1 + \theta_2 x$  obtained from 10 training samples is illustrated in Fig. 25.1. If there is no outlier as in Fig. 25.1(a), a good solution can be obtained by LS. However, if there exists an outlier as in Fig. 25.1(b), the solution is strongly affected by the outlier.

When a large number of training samples are handled, it would be natural to consider that more or less outliers are included. In such a situation, LS is less reliable, as illustrated in Fig. 25.1. One way to cope with outliers is to remove them in advance, which will be discussed in Chapter 38. Another approach is to perform learning so that the solution is less sensitive to outliers. The ability to be less sensitive to outliers is referred to as *robustness*. In this chapter, robust learning methods are explored.

In the LS method, the goodness of fit to training samples is measured by the  $\ell_2$ -loss function:

$$J_{\text{LS}}(\theta) = \frac{1}{2} \sum_{i=1}^n r_i^2,$$

**FIGURE 25.1**

LS solution for straight-line model  $f_{\theta}(x) = \theta_1 + \theta_2 x$ , which is strongly affected by an outlier.

where  $r_i$  denotes the *residual* for the  $i$ th training sample  $(\mathbf{x}_i, y_i)$ :

$$r_i = y_i - f_{\theta}(\mathbf{x}_i).$$

Since the  $\ell_2$ -loss squares the error, large error tends to be magnified significantly. For example, in Fig. 25.1 the output of the rightmost training sample is moved from  $y \approx 3$  to  $y \approx -4$ . Then the residual 7 has the “power” of  $7^2 = 49$  to pull the regression function down, which causes significant change in the solution.

## 25.2 $\ell_1$ -LOSS MINIMIZATION

To cope with the nonrobustness of the LS method, let us employ the  $\ell_1$ -loss or the *absolute loss* (Fig. 25.2):

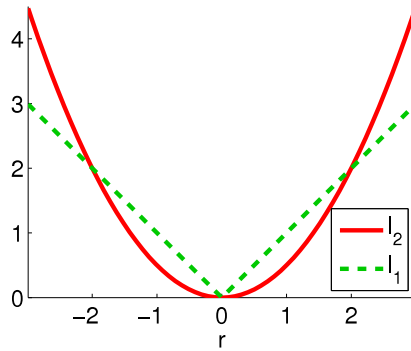
$$\hat{\theta}_{\text{LA}} = \underset{\theta}{\operatorname{argmin}} J_{\text{LA}}(\theta), \quad \text{where } J_{\text{LA}}(\theta) = \sum_{i=1}^n |r_i|.$$

This learning method is called  $\ell_1$ -loss minimization or *least absolute deviations regression*.

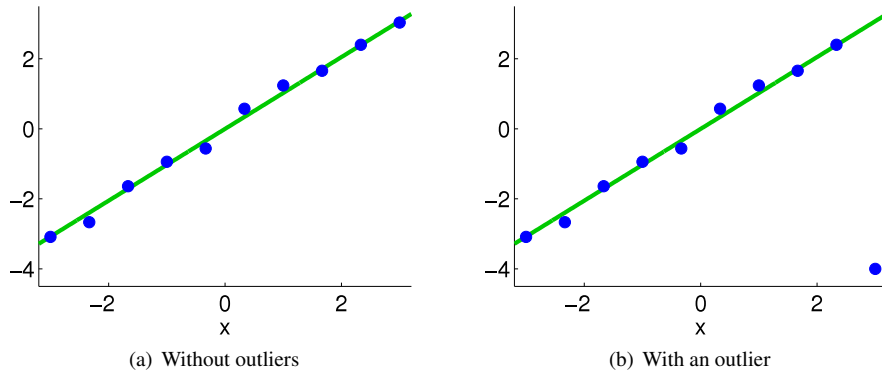
An example of  $\ell_1$ -loss minimization is shown in Fig. 25.3, which is obtained in the same setup as in Fig. 25.1 (how to compute the solution of least absolute deviations will be explained later in Section 25.3). This shows that  $\ell_1$ -loss minimization is much more robust against outliers than  $\ell_2$ -loss minimization, and  $\ell_1$ -loss minimization gives almost the same solution as  $\ell_2$ -loss minimization if there are no outliers.

For a constant model  $f_{\theta}(\mathbf{x}) = \theta$ , the LS method gives the *mean* of training output samples  $\{y_i\}_{i=1}^n$ :

$$\hat{\theta}_{\text{LS}} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n (\theta - y_i)^2 = \operatorname{mean}(\{y_i\}_{i=1}^n).$$

**FIGURE 25.2**

$\ell_2$ -loss and  $\ell_1$ -loss. The  $\ell_2$ -loss magnifies large residuals.

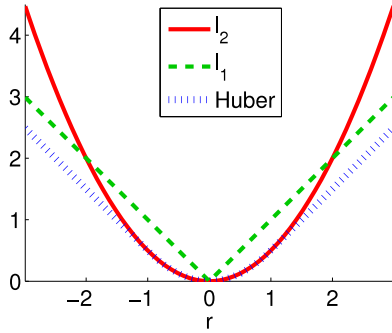
**FIGURE 25.3**

Solution of least absolute deviations for straight-line model  $f_{\theta}(x) = \theta_1 + \theta_2 x$  for the same training samples as Fig. 25.1. Least absolute deviations give a much more robust solution than LS.

On the other hand, least absolute deviations gives the *median* (see Section 2.4.1) of the training output samples  $\{y_i\}_{i=1}^n$ :

$$\hat{\theta}_{\text{LA}} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^n |\theta - y_i| = \operatorname{median}(\{y_i\}_{i=1}^n).$$

The mean of  $\{y_i\}_{i=1}^n$  is affected if one of the sample values is changed, while the median is not affected by changing sample values as long as the order of samples is

**FIGURE 25.4**

Huber loss, with threshold  $\eta = 1$ .

unchanged. The robustness of least absolute deviations is brought by such a property of the  $\ell_1$ -loss.

## 25.3 HUBER LOSS MINIMIZATION

The  $\ell_1$ -loss gives a robust solution, but high robustness implies that learning is performed without fully utilizing the information brought by training samples. In the extreme case, the estimator that does not learn anything from training samples but just always gives zero is the most robust approach, which is clearly meaningless. This extreme example implies that enhancing the robustness too much can degrade the performance of learning from inlier samples. The notion of *efficiency* discussed in Section 13.3 concerns the variance of an estimator, which can be interpreted as how much information can be learned from given training samples. Thus, achieving robustness and efficiency in a balanced way is crucial in practice. In this section, the *Huber loss* [58] is introduced, which can control the balance between robustness and efficiency.

### 25.3.1 DEFINITION

The Huber loss is defined with the  $\ell_2$ -loss and the  $\ell_1$ -loss as follows (Fig. 25.4):

$$\rho_{\text{Huber}}(r) = \begin{cases} r^2/2 & (|r| \leq \eta), \\ \eta|r| - \eta^2/2 & (|r| > \eta). \end{cases}$$

This means that the Huber loss is reduced to the  $\ell_2$ -loss if the absolute residual  $|r|$  is less than or equal to a threshold  $\eta$  (i.e., that sample may be an inlier), while it is reduced to the  $\ell_1$ -loss if the absolute residual  $|r|$  is larger than a threshold  $\eta$  (i.e., that

sample may be an outlier). Note that the  $\ell_1$ -loss is modified as  $\eta|r| - \eta^2/2$  so that it is smoothly connected to the  $\ell_2$ -loss  $r^2/2$ . The method of Huber loss minimization is given by

$$\min_{\theta} J(\theta), \quad \text{where } J(\theta) = \sum_{i=1}^n \rho_{\text{Huber}}(r_i). \quad (25.1)$$

### 25.3.2 STOCHASTIC GRADIENT ALGORITHM

Let us consider the linear-in-parameter model:

$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^b \theta_j \phi_j(\mathbf{x}) = \boldsymbol{\theta}^{\top} \boldsymbol{\phi}(\mathbf{x}).$$

Since the Huber loss is once differentiable as

$$\rho'_{\text{Huber}}(r) = \begin{cases} r & (|r| \leq \eta), \\ -\eta & (r < -\eta), \\ \eta & (r > \eta), \end{cases}$$

the solution of Huber loss minimization may be obtained by the *stochastic gradient* algorithm (see Section 15.3) as follows:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \varepsilon \frac{\partial \rho_{\text{Huber}}}{\partial \boldsymbol{\theta}} \bigg|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}},$$

where  $(\mathbf{x}_i, y_i)$  is a randomly chosen training sample and  $\varepsilon > 0$  is the step size.

### 25.3.3 ITERATIVELY REWEIGHTED LS

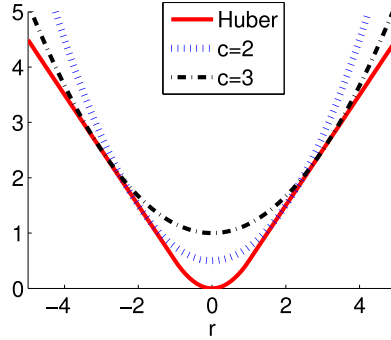
Another approach is *iteratively reweighted LS*, which considers the following quadratic upper bound of the absolute-value part of the Huber loss:

$$\eta|r| - \frac{\eta^2}{2} \leq \frac{\eta}{2c} r^2 + \frac{\eta c}{2} - \frac{\eta^2}{2} \quad \text{for } c > 0.$$

As illustrated in Fig. 25.5, this quadratic upper bound touches the Huber loss at  $r = \pm c$ .

Let us consider an iterative optimization procedure and construct an upper bound for  $c_i = |\tilde{r}_i|$ :

$$\eta|r_i| - \frac{\eta^2}{2} \leq \frac{\eta}{2|\tilde{r}_i|} r_i^2 + \frac{\eta|\tilde{r}_i|}{2} - \frac{\eta^2}{2},$$

**FIGURE 25.5**

Quadratic upper bound  $\frac{\eta r^2}{2c} + \frac{\eta c}{2} - \frac{\eta^2}{2}$  of Huber loss  $\rho_{\text{Huber}}(r)$  for  $c > 0$ , which touches each other at  $r = \pm c$ .

where  $\tilde{r}_i$  denotes the residual of the current solution for the  $i$ th training sample  $(\mathbf{x}_i, y_i)$ . If the absolute-value part of the Huber loss is upper-bounded by this, Huber loss minimization problem Eq. (25.1) yields the following *weighted LS* problem:

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \tilde{J}(\boldsymbol{\theta}), \quad \text{where } \tilde{J}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n \tilde{w}_i r_i^2 + C,$$

where  $C = \sum_{i: |\tilde{r}_i| > \eta} (\eta |\tilde{r}_i|/2 - \eta^2/2)$  is a constant that is independent of  $\boldsymbol{\theta}$ , and weight  $\tilde{w}_i$  is defined as follows (Fig. 25.6):

$$\tilde{w}_i = \begin{cases} 1 & (|\tilde{r}_i| \leq \eta), \\ \eta/|\tilde{r}_i| & (|\tilde{r}_i| > \eta). \end{cases}$$

As explained in Section 22.1, the solution of weighted LS is given analytically as

$$\hat{\boldsymbol{\theta}} = (\boldsymbol{\Phi}^\top \tilde{\mathbf{W}} \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \tilde{\mathbf{W}} \mathbf{y}, \quad (25.2)$$

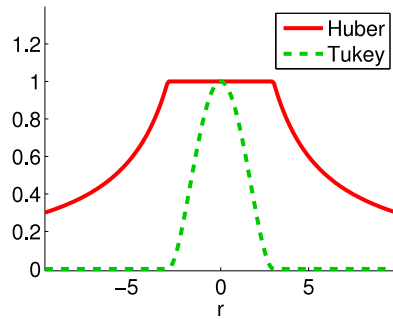
where  $\tilde{\mathbf{W}}$  denotes the diagonal matrix with diagonal elements  $\tilde{w}_1, \dots, \tilde{w}_n$ .

Since the upper bound  $\tilde{J}$  touches the original objective function  $J$  at current solution  $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}$ ,

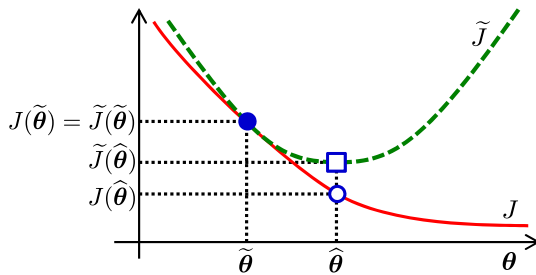
$$J(\hat{\boldsymbol{\theta}}) = \tilde{J}(\hat{\boldsymbol{\theta}})$$

holds. Since  $\hat{\boldsymbol{\theta}}$  is the minimizer of  $\tilde{J}$ ,

$$\tilde{J}(\hat{\boldsymbol{\theta}}) \geq \tilde{J}(\tilde{\boldsymbol{\theta}}).$$

**FIGURE 25.6**

Weight functions for Huber loss minimization and Tukey loss minimization.

**FIGURE 25.7**

Updated solution  $\hat{\theta}$  is no worse than current solution  $\tilde{\theta}$ .

Furthermore,  $\tilde{J}$  upper-bounds  $J$  and thus

$$\tilde{J}(\hat{\theta}) \geq J(\hat{\theta})$$

holds. Summarizing the above relations,

$$J(\tilde{\theta}) = \tilde{J}(\tilde{\theta}) \geq \tilde{J}(\hat{\theta}) \geq J(\hat{\theta})$$

holds, meaning that the updated solution  $\hat{\theta}$  is no worse than the current solution  $\tilde{\theta}$  in terms of the objective value  $J$  (Fig. 25.7). By iterating this upper bound minimization, the solution of Huber loss minimization can be obtained. The pseudocode of this algorithm, called *iteratively reweighted LS*, is provided in Fig. 25.8. If *singular value*

1. Initialize parameter  $\theta$ , e.g., by ordinary LS as

$$\theta \leftarrow (\Phi^T \Phi)^{-1} \Phi^T y.$$

2. Compute the weight matrix  $W$  from the current solution  $\theta$  as

$$W = \text{diag}(w_1, \dots, w_n), \quad \text{where } w_i = \begin{cases} 1 & (|r_i| \leq \eta), \\ \eta/|r_i| & (|r_i| > \eta), \end{cases}$$

and  $r_i = y_i - f_\theta(x_i)$  is the residual.

3. Compute the solution  $\theta$  based on the weight matrix  $W$  as

$$\theta \leftarrow (\Phi^T W \Phi)^{-1} \Phi^T W y.$$

4. Iterate 2–3 until convergence.

**FIGURE 25.8**

Iteratively reweighted LS for Huber loss minimization.

*decomposition* (see Fig. 22.2) of design matrix,

$$\Phi = \sum_{j=1}^b \kappa_j \psi_j \varphi_j^T,$$

is computed in advance, update of  $\theta$  can be performed more efficiently as

$$\theta \leftarrow \sum_{j=1}^b \frac{\psi_j^T W y}{\kappa_j \psi_j^T W \psi_j} \varphi_j.$$

A MATLAB code of iteratively reweighted LS for Huber loss minimization is provided in Fig. 25.9, and its behavior is illustrated in Fig. 25.10. This shows that only two iterations give almost the final solution. The iteration converges after four iterations and a robust solution can be obtained.

When the threshold  $\eta$  is taken to be sufficiently small, the Huber loss may be regarded as a smooth approximation to the  $\ell_1$ -loss. Thus, the solution of  $\ell_1$ -loss minimization can be approximately obtained by the above iteratively weighted LS algorithm.

### 25.3.4 $\ell_1$ -CONSTRAINED HUBER LOSS MINIMIZATION

As explained in Chapter 24, the  $\ell_1$ -constraint tends to give a sparse solution. Here,  $\ell_1$ -constrained Huber loss minimization is considered:



```

n=10; N=1000; x=linspace(-3,3,n)'; X=linspace(-4,4,N)';
y=x+0.2*randn(n,1); y(n)=-4;

p(:,1)=ones(n,1); p(:,2)=x; t=p\y; e=1;
for o=1:1000
    r=abs(p*t-y); w=ones(n,1); w(r>e)=e./r(r>e);
    t0=(p'*( repmat(w,1,2) .*p))\ (p'*(w.*y));
    if norm(t-t0)<0.001, break, end
    t=t0;
end
P(:,1)=ones(N,1); P(:,2)=X; F=P*t;
figure(1); clf; hold on; axis([-4 4 -4.5 3.5]);
plot(X,F,'g-'); plot(x,y,'bo');

```

**FIGURE 25.9**

MATLAB code of iteratively reweighted LS for Huber loss minimization. Straight-line model  $f_\theta(x) = \theta_1 + \theta_2 x$  is used, with threshold  $\eta = 1$ .

$$\min_{\theta} \sum_{i=1}^n \rho_{\text{Huber}}(r_i) \quad \text{subject to } \|\theta\|_1 \leq R^2,$$

which produces a robust and sparse solution.

As shown above, the optimization problem of Huber loss minimization can be solved by iteratively reweighted LS, which was obtained by a quadratic upper bounding technique. Similarly,  $\ell_1$ -regularized LS, whose solution was obtained by ADMM in Section 24.2, can also be solved by iteratively reweighted LS. More specifically, the absolute value of parameter  $\theta$  can be upper-bounded by a quadratic function as

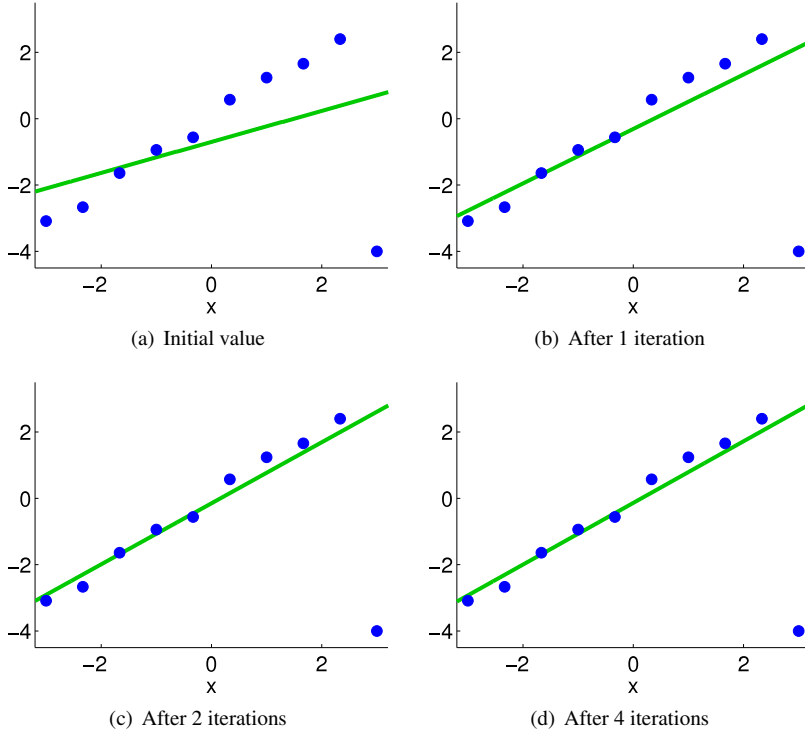
$$|\theta| \leq \frac{\theta^2}{2c} + \frac{c}{2} \quad \text{for } c > 0,$$

which touches the absolute value at  $\theta = \pm c$  (see Fig. 25.11).

Let us consider an iterative optimization procedure and construct an upper bound for  $c_j = \tilde{\theta}_j$ :

$$|\theta_j| \leq \frac{|\tilde{\theta}_j|^\dagger}{2} \theta_j^2 + \frac{|\tilde{\theta}_j|}{2},$$

where  $\tilde{\theta}_j$  is the current solution and  $^\dagger$  denotes the *generalized inverse* (see Fig. 22.1). If the  $\ell_1$ -norm of parameter vector  $\theta = (\theta_1, \dots, \theta_b)^\top$  is upper-bounded by this, the

**FIGURE 25.10**

Examples of iteratively reweighted LS for Huber loss minimization. Straight-line model  $f_{\theta}(x) = \theta_1 + \theta_2 x$  is used, with threshold  $\eta = 1$ .

$\ell_1$ -regularized LS problem,

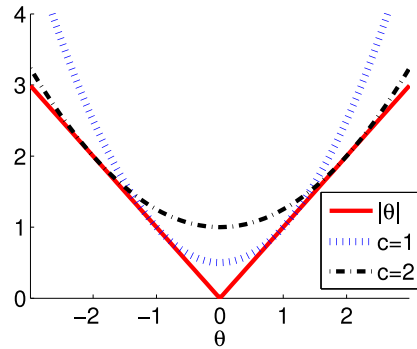
$$\min_{\theta} \left[ \frac{1}{2} \|\mathbf{y} - \Phi\theta\|^2 + \lambda \|\theta\|_1 \right],$$

yields the following *generalized  $\ell_2$ -regularized LS*:

$$\hat{\theta} = \operatorname{argmin}_{\theta} \tilde{J}(\theta), \quad \text{where } \tilde{J}(\theta) = J_{\text{LS}}(\theta) + \frac{\lambda}{2} \theta^{\top} \tilde{\Theta}^{\dagger} \theta + C,$$

where  $C = \sum_{j=1}^b |\tilde{\theta}_j|/2$  is a constant that is independent of  $\theta$ . As explained in Section 23.2, the solution of generalized  $\ell_2$ -regularized LS is given analytically as

$$\hat{\theta} = (\Phi^{\top} \Phi + \lambda \tilde{\Theta}^{\dagger})^{-1} \Phi^{\top} \mathbf{y}.$$

**FIGURE 25.11**

Quadratic upper bound  $\frac{\theta^2}{2c} + \frac{c}{2}$  of absolute value  $|\theta|$  for  $c > 0$ , which touches each other at  $\theta = \pm c$ .

1. Initialize parameter  $\theta$ , e.g., randomly or by ordinary LS as

$$\theta \leftarrow \Phi^\dagger \mathbf{y}.$$

2. Compute the weight matrix  $\mathbf{W}$  and regularization matrix  $\Theta$  from the current solution  $\theta$  as

$$\mathbf{W} = \text{diag}(w_1, \dots, w_n) \quad \text{and} \quad \Theta = \text{diag}(|\theta_1|, \dots, |\theta_b|),$$

where weight  $w_i$  is defined using residual  $r_i = y_i - f_\theta(\mathbf{x}_i)$  as

$$w_i = \begin{cases} 1 & (|r_i| \leq \eta), \\ \eta/|r_i| & (|r_i| > \eta). \end{cases}$$

3. Compute the solution  $\theta$  based on the weight matrix  $\mathbf{W}$  regularization matrix  $\Theta$  as

$$\theta \leftarrow (\Phi^\top \mathbf{W} \Phi + \lambda \Theta^\dagger)^{-1} \Phi^\top \mathbf{W} \mathbf{y}.$$

4. Iterate 2–3 until convergence.

**FIGURE 25.12**

Iteratively reweighted LS for  $\ell_1$ -regularized Huber loss minimization.

```

n=50; x=linspace(-3,3,n)'; pix=pi*x;
y=sin(pix)./(pix)+0.1*x+0.2*randn(n,1); y(n/2)=-0.5;
hh=2*0.3^2; l=0.1; e=0.1; t=randn(n,1); x2=x.^2;
k=exp(-( repmat(x2,1,n)+repmat(x2',n,1)-2*x*x')/hh);
for o=1:1000
    r=abs(k*t-y); w=ones(n,1); w(r>e)=e./r(r>e);
    Z=k*( repmat(w,1,n).*k)+l*pinv(diag(abs(t)));
    t0=(Z+0.000001*eye(n))\ (k*(w.*y));
    if norm(t-t0)<0.001, break, end
    t=t0;
end
N=1000; X=linspace(-3,3,N)';
K=exp(-( repmat(X.^2,1,n)+repmat(x2',N,1)-2*X*x')/hh); F=K*t;
figure(1); clf; hold on; axis([-2.8 2.8 -1 1.5]);
plot(X,F,'g-'); plot(x,y,'bo');

```

**FIGURE 25.13**

MATLAB code of iteratively reweighted LS for  $\ell_1$ -regularized Huber loss minimization with Gaussian kernel model.

The iteratively reweighted LS algorithms for  $\ell_1$ -regularized LS and Huber loss minimization can be combined to obtain a sparse and robust solution, as summarized in Fig. 25.12.

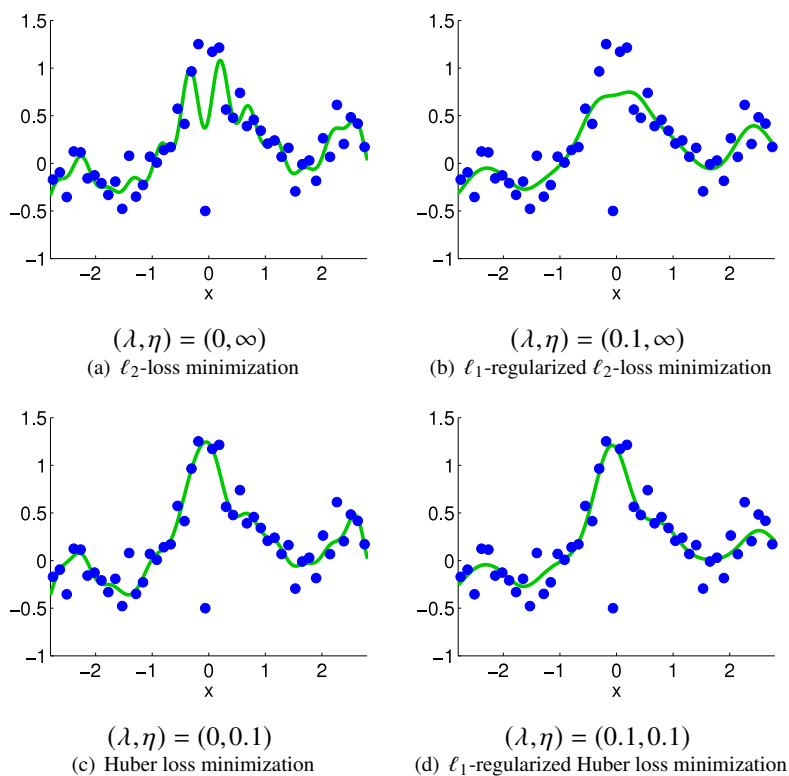
A MATLAB code of iteratively reweighted LS for  $\ell_1$ -regularized Huber loss minimization is provided in Fig. 25.13, where the Gaussian kernel model,

$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^n \theta_j \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_j\|^2}{2h^2}\right),$$

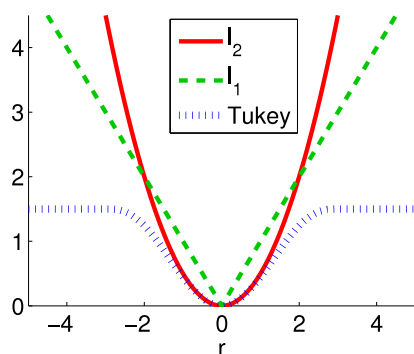
is used and the Gaussian bandwidth is set at  $h = 0.3$ . Since plain LS is numerically unstable, it is stabilized by adding  $10^{-6}$  to the diagonal elements of  $\Phi^T \mathbf{W} \Phi$  (more precisely,  $\mathbf{K} \mathbf{W} \mathbf{K}$  since the Gaussian kernel model is used). The behavior of  $\ell_1$ -regularized Huber loss minimization is illustrated in Fig. 25.14. Regardless of the presence or absence of the  $\ell_1$ -regularizer,  $\ell_2$ -loss minimization is strongly affected by the outlier at around  $x = 0$ , while Huber loss minimization can successfully suppress its influence. With the  $\ell_1$ -regularizer,  $\ell_2$ -loss minimization gives 38 zero parameters out of 50 parameters, and Huber loss minimization gives 36 zero parameters.

## 25.4 TUKEY LOSS MINIMIZATION

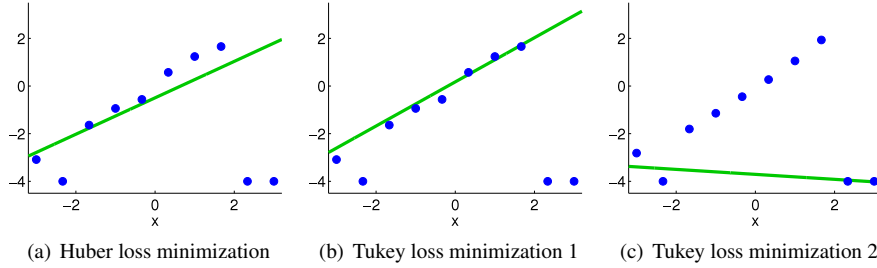
The Huber loss nicely controls the balance between efficiency and robustness by combining the  $\ell_1$ -loss and the  $\ell_2$ -loss. However, as long as the  $\ell_1$ -loss is used, strong

**FIGURE 25.14**

Example of  $\ell_1$ -regularized Huber loss minimization with Gaussian kernel model.

**FIGURE 25.15**

Tukey loss, with threshold  $\eta = 3$ .

**FIGURE 25.16**

Example of Tukey loss minimization. Tukey loss minimization gives more robust solutions than Huber loss minimization, but only a local optimal solution can be obtained.

outliers can still affect the solution significantly. Indeed, as illustrated in Fig. 25.6, the weight  $\tilde{w}_i$  used in Huber loss minimization is not zero for large absolute residuals.

In such a hard circumstance, the *Tukey loss* [73] is useful (Fig. 25.15):

$$\rho_{\text{Tukey}}(r) = \begin{cases} \left(1 - [1 - r^2/\eta^2]^3\right) \eta^2/6 & (|r| \leq \eta), \\ \eta^2/6 & (|r| > \eta). \end{cases}$$

The Tukey loss is constant  $\eta^2/6$  if the absolute residual  $|r|$  is larger than a threshold  $\eta$ . Thus, Tukey loss minimization is expected to be more robust than Huber loss minimization.

However, the Tukey loss is not a *convex function* (see Fig. 8.3). Thus, there may exist multiple local optimal solutions and finding the global one is not straightforward. In practice, *iteratively reweighted LS* with the following weight is used to find a local optimal solution (Fig. 25.6):

$$w = \begin{cases} (1 - r^2/\eta^2)^2 & (|r| \leq \eta), \\ 0 & (|r| > \eta). \end{cases}$$

Since the Tukey weight is zero for  $|r| > \eta$ , Tukey loss minimization is not at all influenced by strong outliers.

A MATLAB code of iteratively reweighted LS for Tukey loss minimization is essentially the same as that for Huber loss minimization shown in Fig. 25.9. Only the weight computation for the Huber loss,

```
w=ones(n,1); w(r>e)=e./r(r>e);
```

is replaced with that for the Tukey loss:

```
w=zeros(n,1); w(r<=e)=(1-r(r<=e).^2/e^2).^2;
```

The behavior of Tukey loss minimization is illustrated in [Fig. 25.16](#), where the threshold is set at  $\eta = 1$ . In this example, Tukey loss minimization gives an even more robust solution than Huber loss minimization. However, since Tukey loss minimization is a nonconvex optimization problem, different solutions may be obtained if the initial solution is changed or training samples are slightly perturbed. Indeed, as illustrated in [Fig. 25.16\(c\)](#), slightly changing the noise included in training output samples gives another local optimal solution, which is significantly different from the original one in [Fig. 25.16\(b\)](#).

For general differentiable symmetric loss  $\rho(r)$ , iteratively reweighted LS with weight  $w_i = \rho'(r_i)/r_i$  gives a (local) optimal solution, where  $\rho'(r)$  denotes the derivative of  $\rho(r)$  [57].

