

ロバスト回帰 (6章)

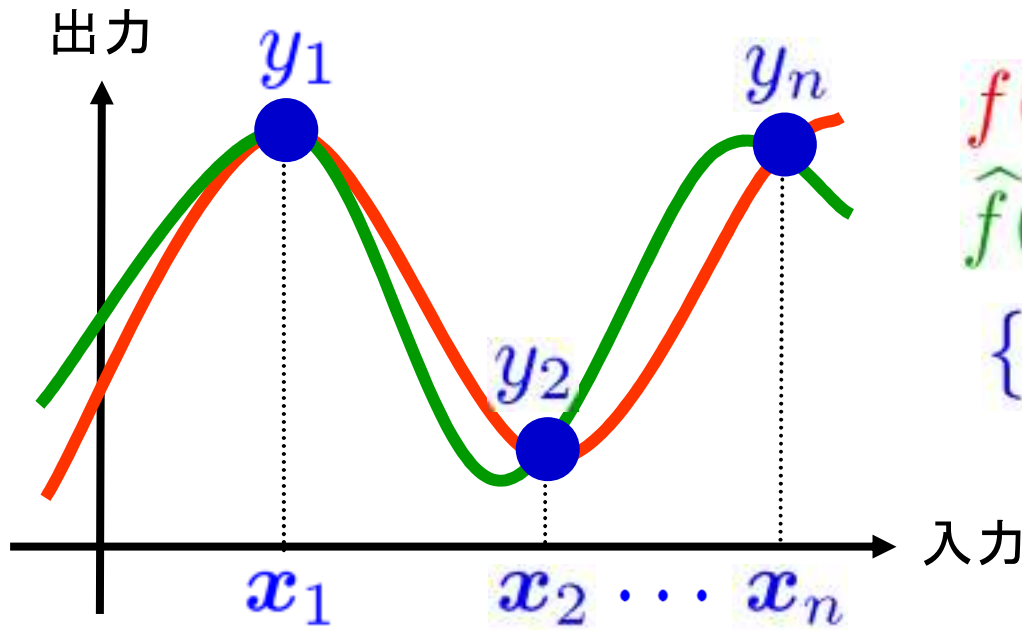
杉山将・本多淳也

sugi@k.u-tokyo.ac.jp, jhonda@k.u-tokyo.ac.jp

<http://www.ms.k.u-tokyo.ac.jp>

回帰 = 関数近似

2



$f(x)$: 学習したい真の関数

$\hat{f}(x)$: 学習結果の関数

$\{(x_i, y_i)\}_{i=1}^n$: 訓練標本

$y_i = f(x_i) (+\text{noise})$

訓練標本から真の関数にできるだけ近い関数を求める

パラメータに関する線形モデル 3

■ 線形モデル:
$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^b \theta_j \phi_j(\mathbf{x})$$

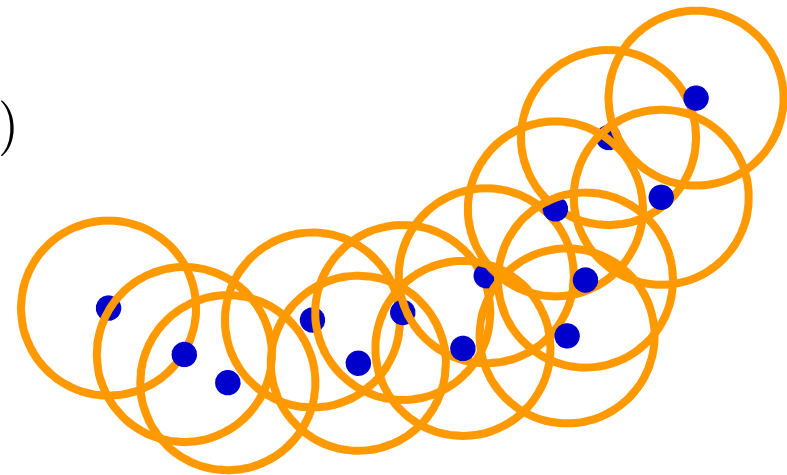
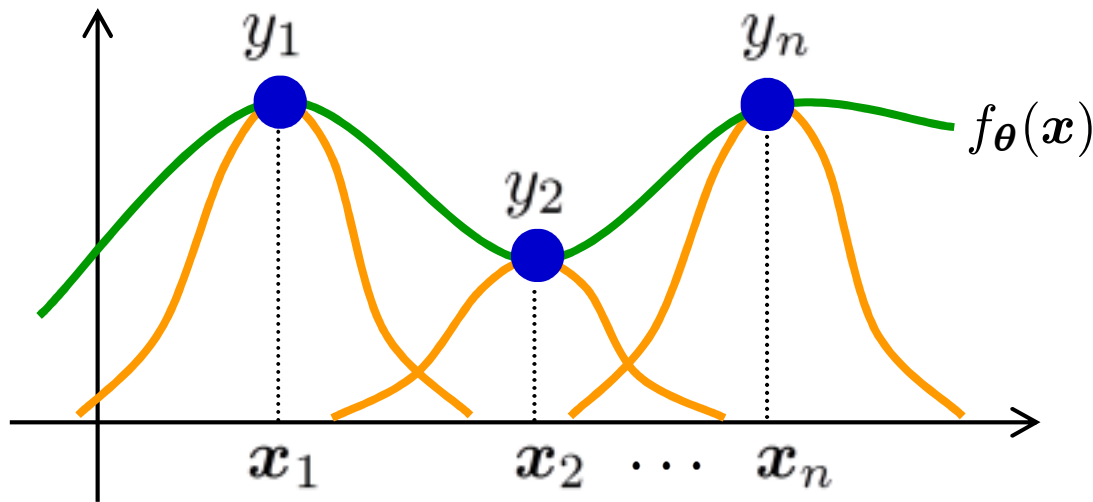
$\{\phi_j(\mathbf{x})\}_{j=1}^b$
: 基底関数

■ カーネルモデル:

ガウスカーネル

$$f_{\theta}(\mathbf{x}) = \sum_{j=1}^n \theta_j K(\mathbf{x}, \mathbf{x}_j)$$

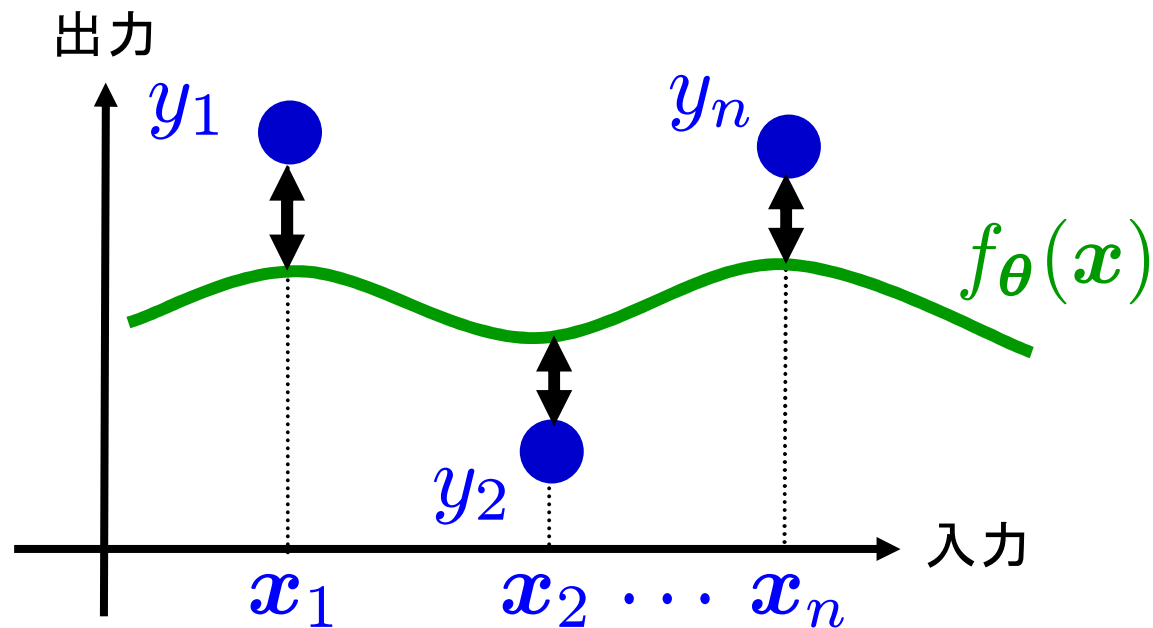
$$K(\mathbf{x}, \mathbf{c}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}\|^2}{2h^2}\right)$$



最小二乗回帰

- 訓練出力との二乗誤差を最小にする:

$$\min_{\theta} \sum_{i=1}^n \left(f_{\theta}(x_i) - y_i \right)^2$$

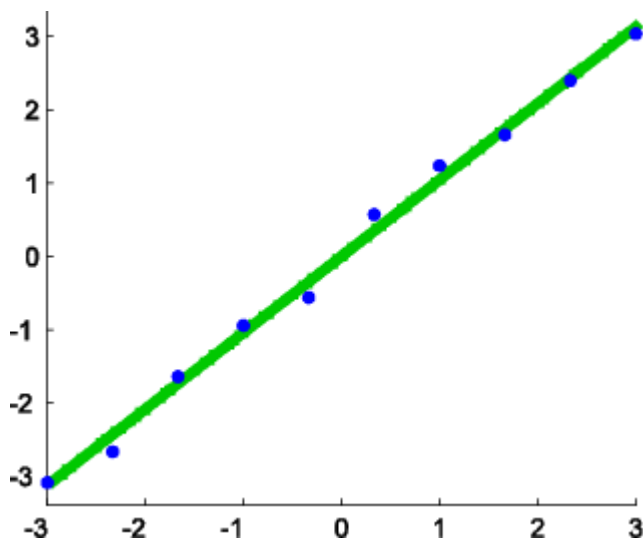


最小二乗回帰の問題点

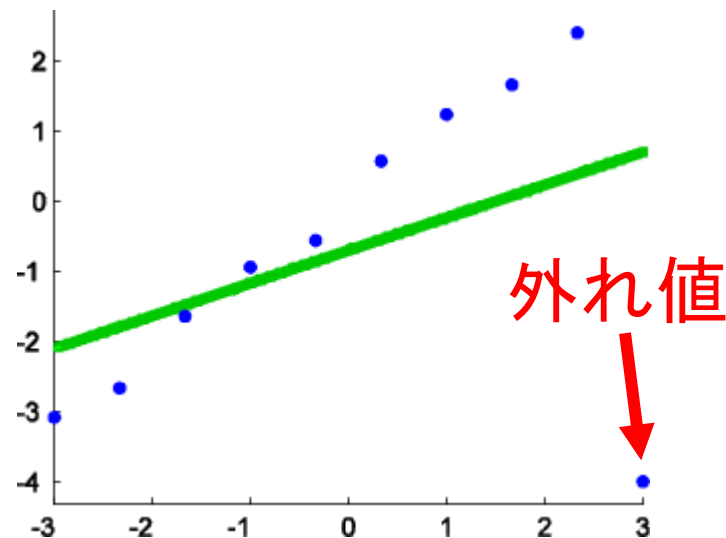
5

- たった一つの**外れ値**が，学習結果を大きく変えてしまう！

$$f_{\theta}(x) = \theta_1 + \theta_2 x$$



最小二乗回帰
(外れ値なし)

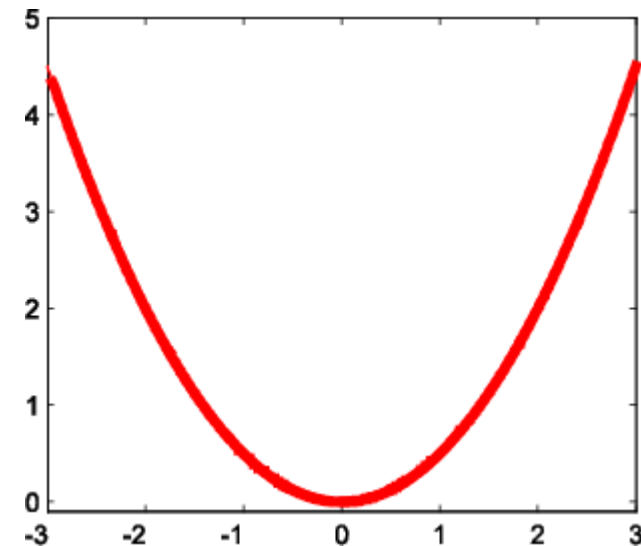


最小二乗回帰
(外れ値あり)

ℓ_2 -損失関数

$$\sum_{i=1}^n \left(f_{\theta}(x_i) - y_i \right)^2$$

- **最小二乗回帰**: 訓練出力との適合のよさを ℓ_2 -損失関数で測る
- 外れ値は“**二乗の強さ**”で学習結果に影響を及ぼす
- 学習結果は安定させるには外れ値の影響を小さくする必要がある



講義の流れ



1. ℓ_1 -損失

1. 中央値との関係
2. ロバスト性と推定精度

2. フーバー損失

3. テューキー損失

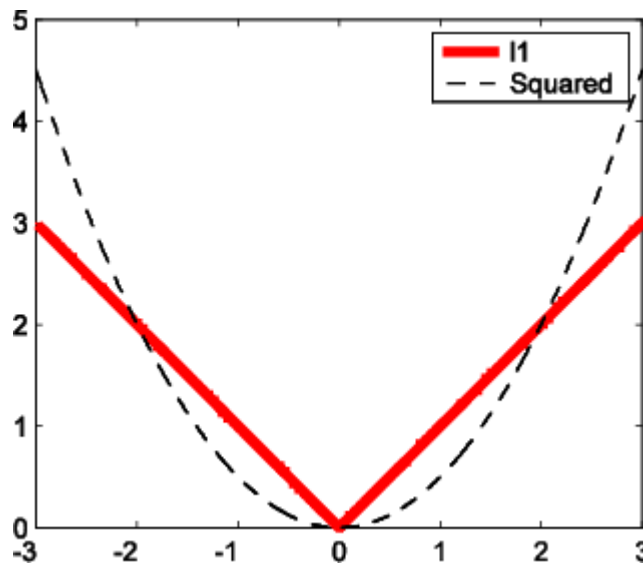
ℓ_1 -損失を用いたロバスト回帰

8

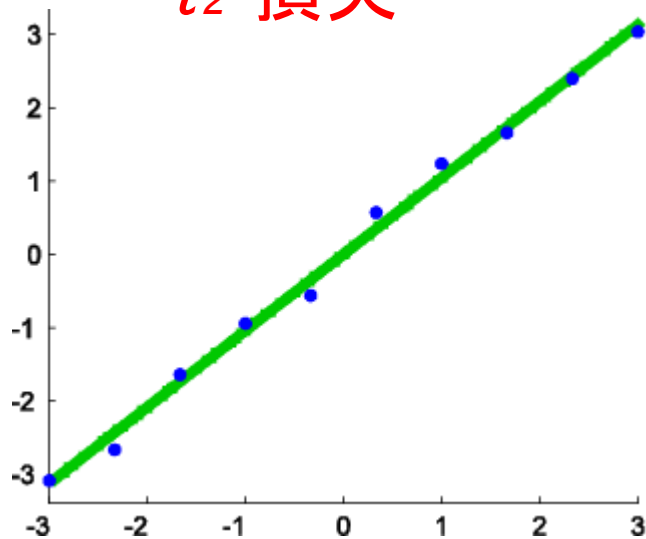
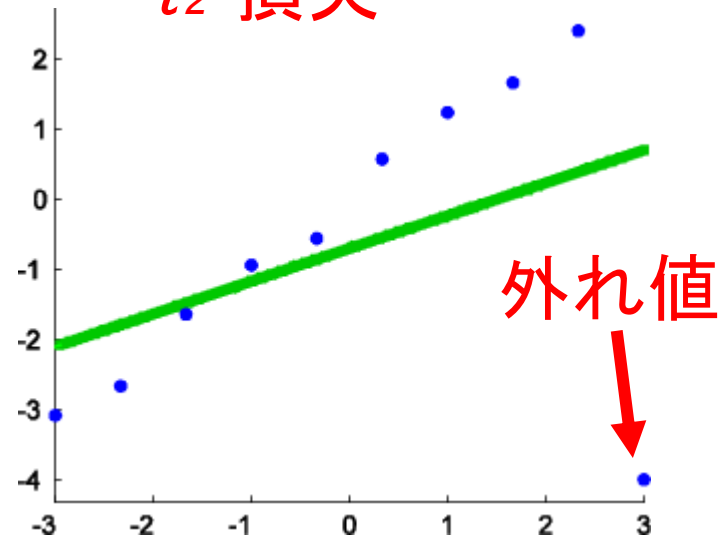
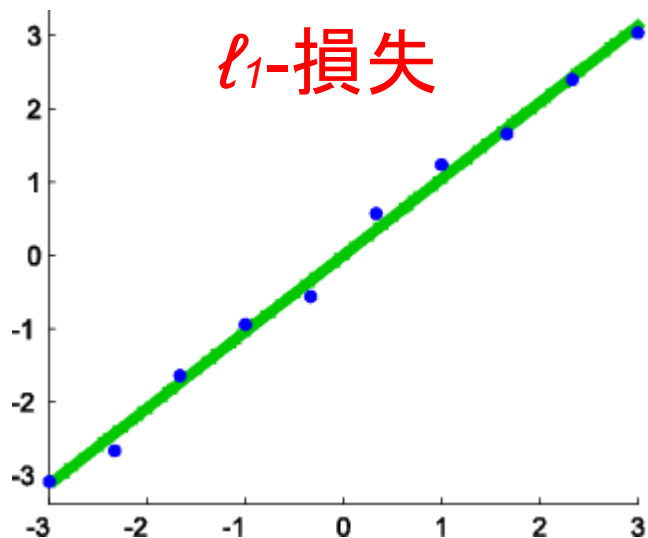
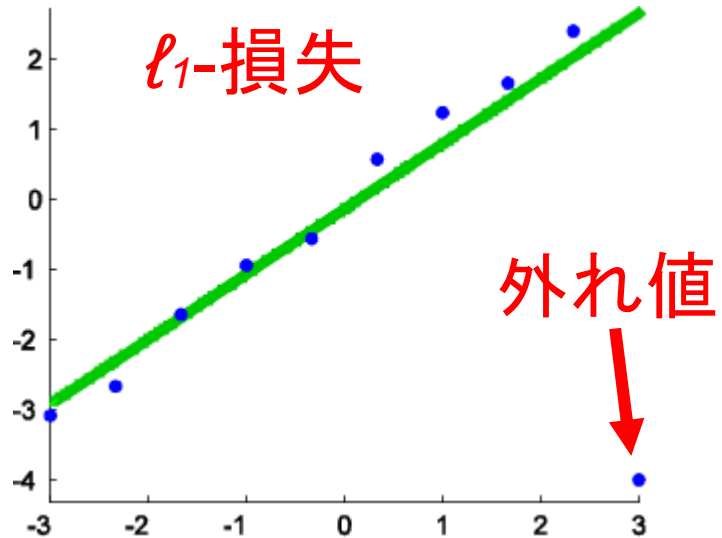
- ℓ_1 -損失関数で訓練出力との適合のよさを測る:

$$\min_{\theta} \sum_{i=1}^n \left| f_{\theta}(x_i) - y_i \right|$$

- 外れ値の影響は線形にしか効いてこない.



例

 ℓ_2 -損失 ℓ_2 -損失 ℓ_1 -損失 ℓ_1 -損失

解の求め方は後ほど

講義の流れ



1. ℓ_1 -損失

1. 中央値との関係
2. ロバスト性と推定精度

2. フーバー損失

3. テューキー損失

累積分布関数

- 連続型の確率変数が x 以下の値をとる確率

$$P(x) = \text{Prob}(X \leq x) = \int_{-\infty}^x p(u) du$$

$P(x)$: 累積分布関数(cumulative distribution function)

$$P'(x) = \frac{dP(x)}{dx} = p(x)$$

- 広義単調増加:

$$x_1 < x_2 \implies P(x_1) \leq P(x_2)$$

- 範囲:

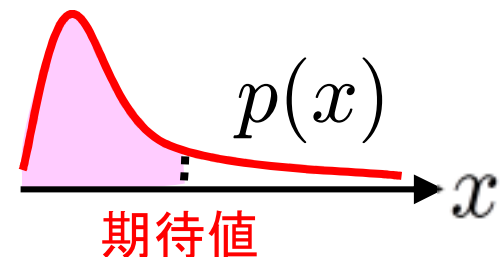
$$x \rightarrow -\infty \implies P(x) \rightarrow 0$$

$$x \rightarrow \infty \implies P(x) \rightarrow 1$$

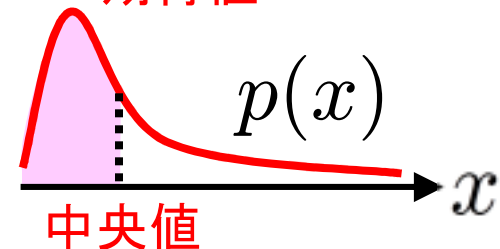
期待値と中央値

12

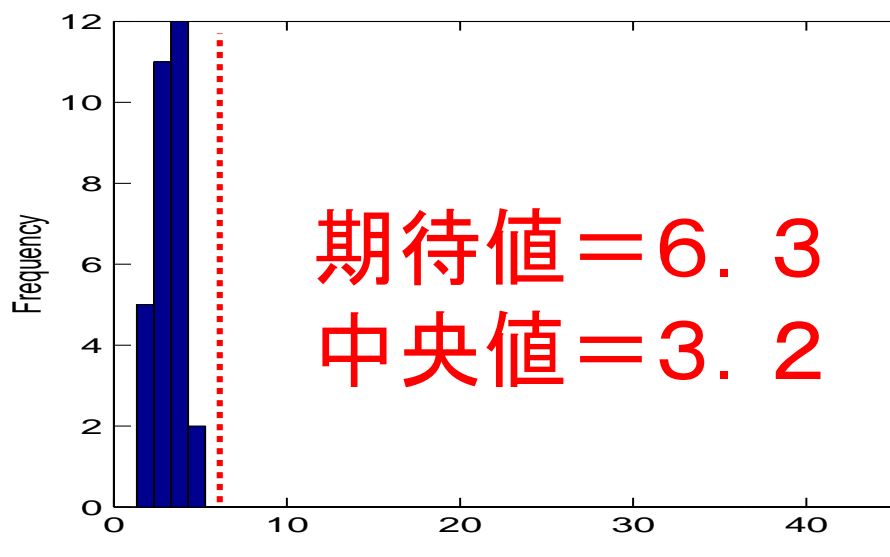
■ 期待値: $E[X] = \int xp(x)dx$



■ 中央値: $\text{Prob}(X \leq x) = \frac{1}{2}$ を満たす x



■ 期待値は, 外れ値 (outlier) があるときに直感と合わない値になることがある



期待値 = 6.3
中央値 = 3.2

例: 収入分布. 一人超大金持ち (外れ値) がいると, その人以外全員が期待値以下になってしまう

外れ値

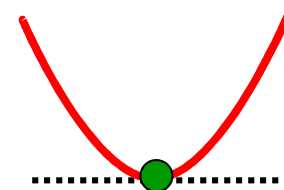


- $[a, b]$ 上に定義された確率密度関数 $p(y)$ を考える
- 次の二乗誤差 $J_2(y)$ を最小にする θ_2 は, y の期待値であることを示せ

$$\theta_2 = \operatorname{argmin}_{\theta} J_2(\theta)$$

$$J_2(\theta) = \int_a^b (y - \theta)^2 p(y) dy$$

- ヒント: 最小点での微分はゼロ



- $[a, b]$ 上の確率密度関数 $p(y)$ を考える
- **絶対値誤差** $J_1(\theta)$ を最小にする θ_1 は y の **中央値** である($P(\theta_1) = 1/2$ となる)ことを示せ

$$\theta_1 = \operatorname{argmin}_{\theta} J_1(\theta)$$

$$J_1(\theta) = \int_a^b |y - \theta| p(y) dy$$

- ヒント: 累積分布の微分 $P'(y) = p(y)$ について部分積分を適用

$$\int_a^b f(y) g'(y) dy = \left[f(y) g(y) \right]_a^b - \int_a^b f'(y) g(y) dy$$

■ 観測値＝真値＋雑音: $\{y_i \mid y_i = \mu^* + \epsilon_i\}_{i=1}^n$

■ ℓ_2 -損失関数: 観測値の期待値の推定に対応

$$\operatorname{argmin}_{\theta} \mathbb{E} \left[(y - \theta)^2 \right] = \mu^* + \operatorname{mean}(\epsilon)$$

■ ℓ_1 -損失関数: 観測値の中央値の推定に対応

$$\operatorname{argmin}_{\theta} \mathbb{E} [|y - \theta|] = \mu^* + \operatorname{median}(\epsilon)$$

■ より一般に出力が入力に依存するモデル

$y_i = \mu^*(\mathbf{x}_i) + \epsilon_i(\mathbf{x}_i)$ の場合は, $\mu^*(\mathbf{x}) + \epsilon(\mathbf{x})$ の各点での期待値あるいは中央値の推定に対応

講義の流れ



1. ℓ_1 -損失

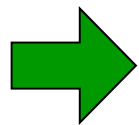
1. 中央値との関係
2. ロバスト性と推定精度

2. フーバー損失

3. テューキー損失

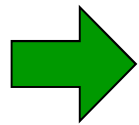
■ **破綻点** : 標本のいくつかを**無限**に飛ばしたとき、学習結果が**有限**に留まる(破綻しない)標本数の割合の上限

- ℓ_2 -損失関数 : **0%**

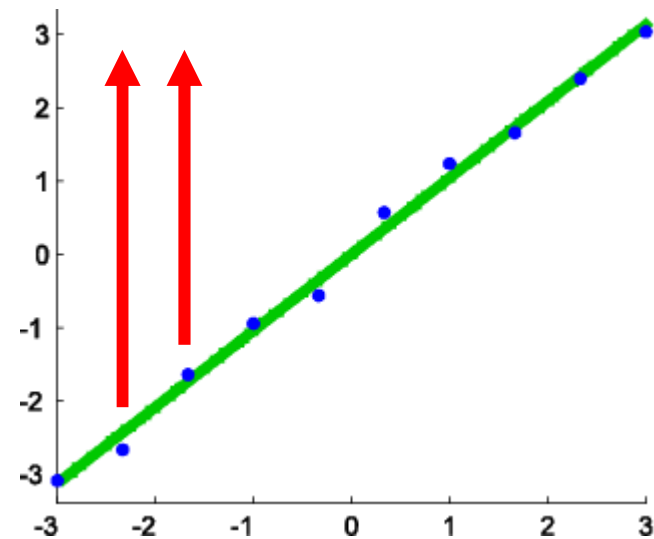


全くロバストでない

- ℓ_1 -損失関数 : **50%**



非常にロバスト



■ ただし, ℓ_1 -損失関数はガウス雑音に対して**有効性を満たさない**(分散が大きい)

- ロバスト性が高い = 訓練標本をきちんと見ない
 - 常に0を出力する無意味な学習が最もロバスト
- 現実的な要請：
 - 外れ値がない場合には最小二乗法に近い
 - 外れ値が多い or 大きい場合にロバスト

講義の流れ

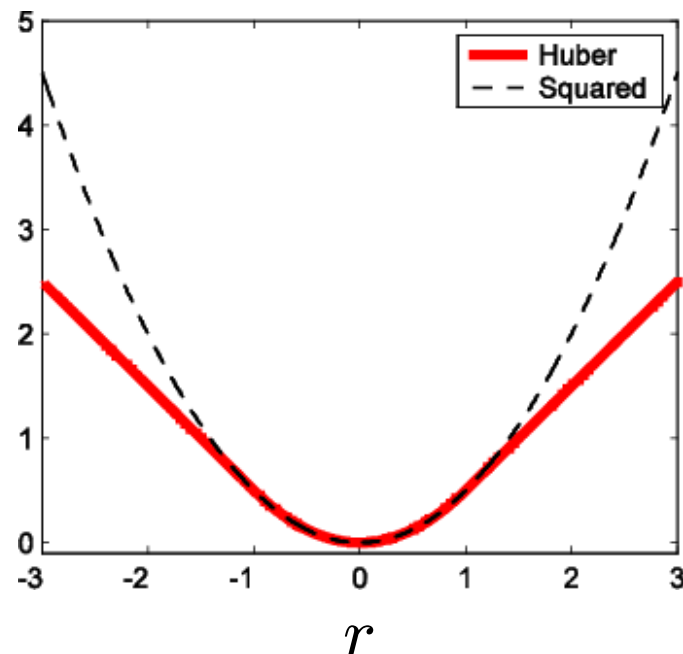


1. ℓ_1 -損失
2. フーバー損失
3. テューキー損失

フーバー損失

- ℓ_1 -損失と ℓ_2 -損失の折衷案
- 小さな誤差に対しては二乗
- 大きな誤差に対しては絶対値

$$\min_{\theta} \sum_{i=1}^n \rho(f_{\theta}(\mathbf{x}_i) - y_i)$$



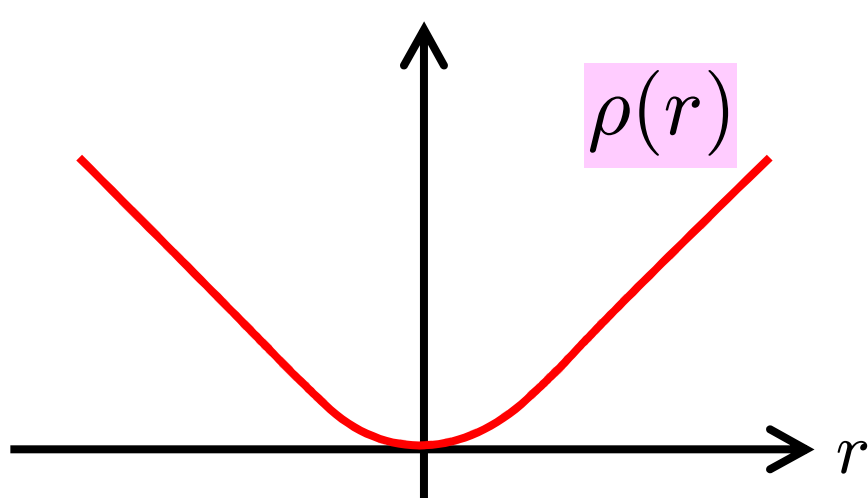
$$\rho(r) = \begin{cases} r^2/2 & (|r| \leq \eta) \\ \eta|r| - \eta^2/2 & (|r| > \eta) \end{cases}$$

$$\eta \geq 0$$

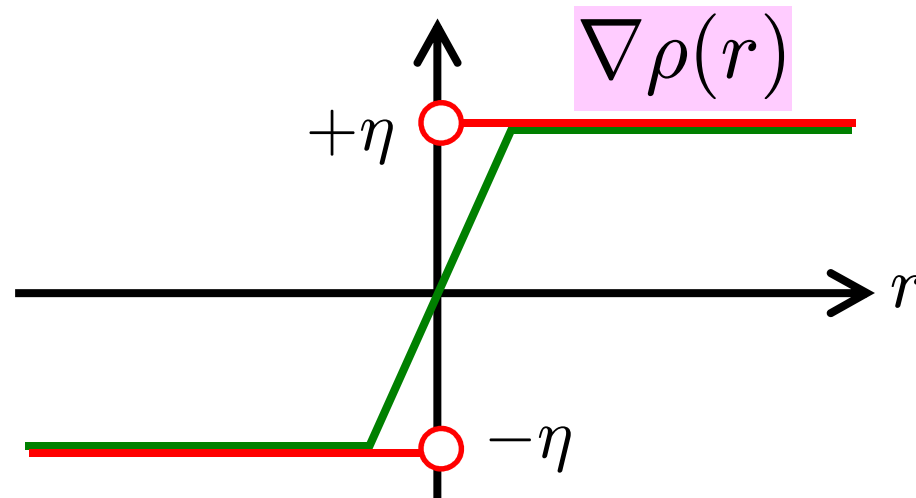
- パラメータ η は誤差が外れ値由来である可能性が現れ始める点としてユーザが設計

解の求め方1

■ フーバー損失は連続的微分可能



$$\rho(r) = \begin{cases} r^2/2 & (|r| \leq \eta) \\ \eta|r| - \eta^2/2 & (|r| > \eta) \end{cases}$$

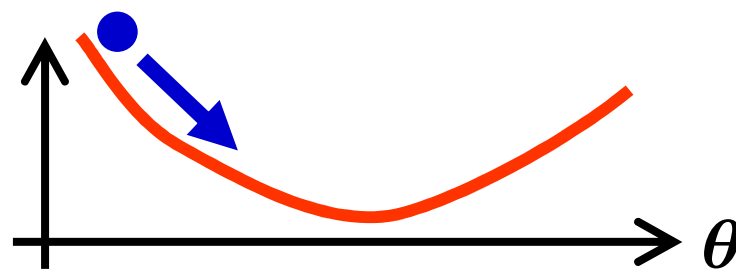


$$\rho'(r) = \begin{cases} r & (|r| \leq \eta) \\ \text{sign}(r)\eta & (|r| > \eta) \end{cases}$$

■ 勾配法:

$$\theta \leftarrow \theta - \varepsilon \nabla J(\theta)$$

$$J(\theta) = \sum_{i=1}^n \rho(f_{\theta}(\mathbf{x}_i) - y_i)$$

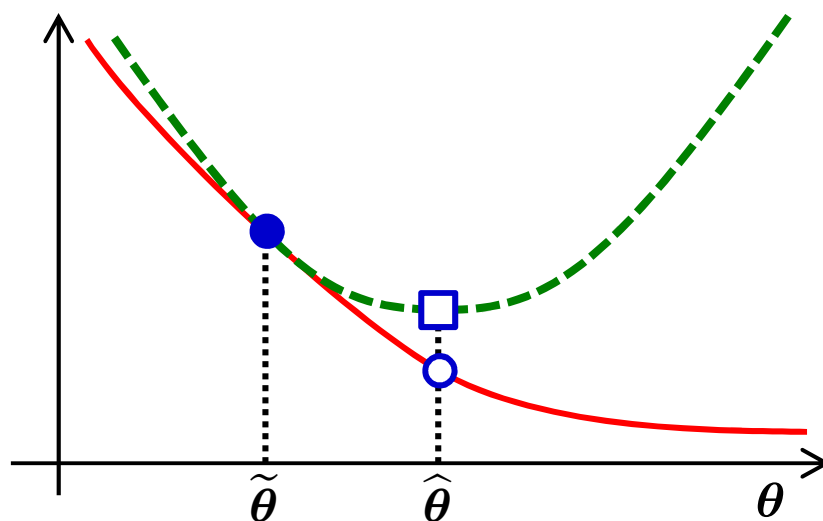


解の求め方2

■ 勾配法はステップ幅の調整が厄介

■ 繰り返し最小二乗アルゴリズム:

- フーバー損失を現在の解で接する二次関数で上から抑える(ニュートン法とは異なる)
- 二次上界を解析的に最小化することにより, 少しずつ良い解を求めていく



■ $|r| > \eta$ の場合のフーバー損失:

$$\eta|r| - \frac{\eta^2}{2}$$

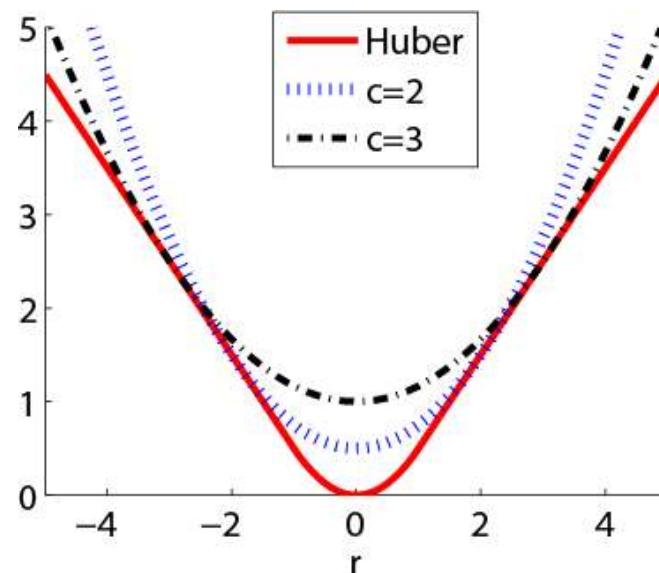
$$\rho(r) = \begin{cases} r^2/2 & (|r| \leq \eta) \\ \eta|r| - \eta^2/2 & (|r| > \eta) \end{cases}$$

に $r = \pm c, c > \eta$ で接する二次関数を求めよ

■ ヒント: $\pm c$ で接する二次関数は対称性より

$$ar^2 + b$$

と表される



二次上界の最小化

$$\rho(r) = \begin{cases} r^2/2 & (|r| \leq \eta) \\ \eta|r| - \eta^2/2 & (|r| > \eta) \end{cases}$$

- 現在の解 $\tilde{\theta}$ に対する残差 $\tilde{r} = f_{\tilde{\theta}}(x) - y$ に対する二次上界 $\tilde{\rho}(r) \geq \rho(r)$:

$$\tilde{\rho}(r) = \begin{cases} r^2/2 & (|\tilde{r}| \leq \eta) \\ \frac{\eta}{2|\tilde{r}|} r^2 + \underbrace{\frac{\eta|\tilde{r}|}{2} - \frac{\eta^2}{2}}_{\text{定数}} & (|\tilde{r}| > \eta) \end{cases}$$

定数

$$= \frac{\tilde{w}}{2} r^2 + \text{const}$$

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

二次上界の最小化(続き)

29

- 元々の最小化したい損失:

$$J(\boldsymbol{\theta}) = \sum_{i=1}^n \rho(f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i)$$

$$\rho(r) = \begin{cases} r^2/2 & (|r| \leq \eta) \\ \eta|r| - \eta^2/2 & (|r| > \eta) \end{cases}$$

- 現在の解 $\tilde{\boldsymbol{\theta}}$ から求めた J の上界 \tilde{J} の最小化:

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

$$\tilde{J}(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n \tilde{w}_i \left(f_{\boldsymbol{\theta}}(\mathbf{x}_i) - y_i \right)^2$$

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

$$\tilde{r}_i = f_{\tilde{\boldsymbol{\theta}}}(\mathbf{x}_i) - y_i$$

二次上界の最小化(続き)

30

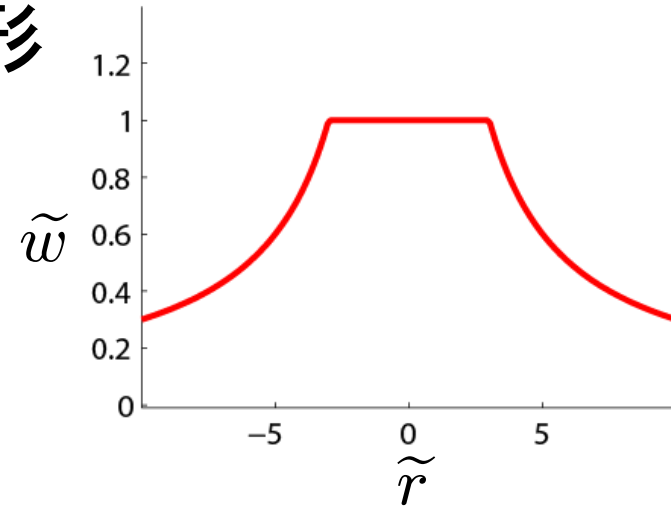
■ 上界は重み付き最小二乗法の形

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n \tilde{w}_i \left(f_{\theta}(\mathbf{x}_i) - y_i \right)^2$$

● 異常値に対する重みを小さく設定

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

$$\tilde{r}_i = f_{\tilde{\theta}}(\mathbf{x}_i) - y_i$$



■ 上界の最小解は解析的に次式で求められる

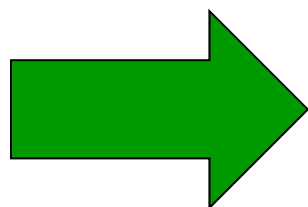
$$\hat{\theta} = (\Phi^{\top} \tilde{\mathbf{W}} \Phi)^{-1} \Phi^{\top} \tilde{\mathbf{W}} \mathbf{y}$$

$$\tilde{\mathbf{W}} = \text{diag}(\tilde{w}_1, \dots, \tilde{w}_n)$$

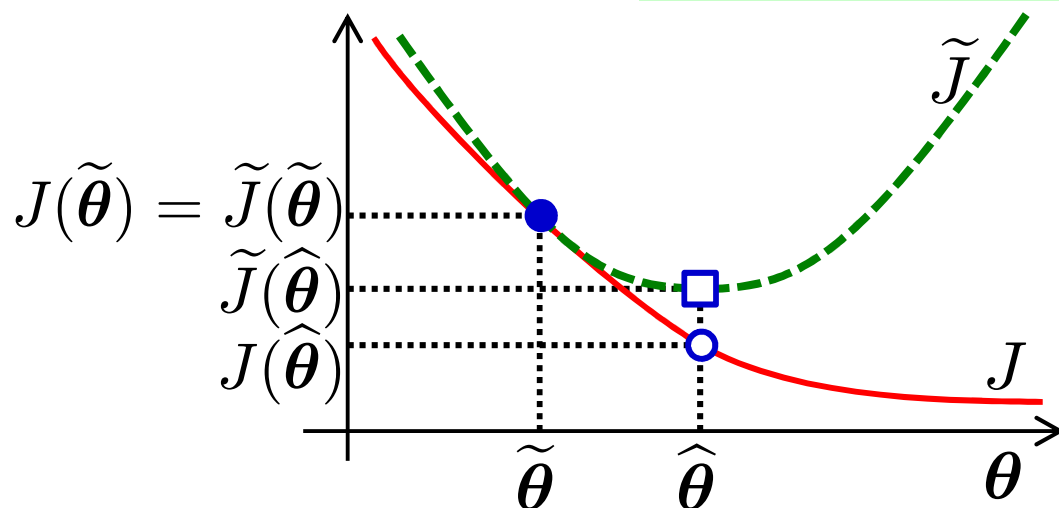
証明は
宿題

二次上界の最小化(続き)

- 上界が $\tilde{\theta}$ で接することから $J(\tilde{\theta}) = \tilde{J}(\tilde{\theta})$
- $\hat{\theta}$ を上界の最小解とすれば $\tilde{J}(\tilde{\theta}) \geq \tilde{J}(\hat{\theta})$
- \tilde{J} が J の上界であることから $\tilde{J}(\hat{\theta}) \geq J(\hat{\theta})$
- まとめると $J(\tilde{\theta}) = \tilde{J}(\tilde{\theta}) \geq \tilde{J}(\hat{\theta}) \geq J(\hat{\theta})$



解を $\tilde{\theta}$ から $\hat{\theta}$ に更新すれば
 J の値は(一般に)減少する



$$\hat{\theta} = \operatorname{argmin}_{\theta} \tilde{J}(\theta)$$

繰り返し再重み付け 最小二乗アルゴリズム

- θ を適当に初期化する
- 以下を収束するまで繰り返す
 - 現在の解 θ から行列 W を計算する(上界を求める)

$$W = \text{diag}(w_1, \dots, w_n)$$

$$w_i = \begin{cases} 1 & (|f_{\theta}(\mathbf{x}_i) - y_i| \leq \eta) \\ \eta / |f_{\theta}(\mathbf{x}_i) - y_i| & (|f_{\theta}(\mathbf{x}_i) - y_i| > \eta) \end{cases}$$

- 解 θ を更新する(上界を最小化する)

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

実装例

■ 直線モデル $f_{\theta}(x) = \theta_1 + \theta_2 x$ に対する フーバー回帰の繰り返し最小二乗アルゴリズム

```
clear all; rand('state',0); randn('state',0);
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; t0=p\y; e=1;
for o=1:1000
    r=abs(p*t0-y); w=ones(n,1); w(r>e)=e./r(r>e);
    t=(p'*( repmat(w,1,2).*p ))\ (p'*(w.*y));
    if norm(t-t0)<0.001, break, end
    t0=t;
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');
```

η が小さい場合
 ℓ_1 -損失最小化と
ほぼ一致

講義の流れ

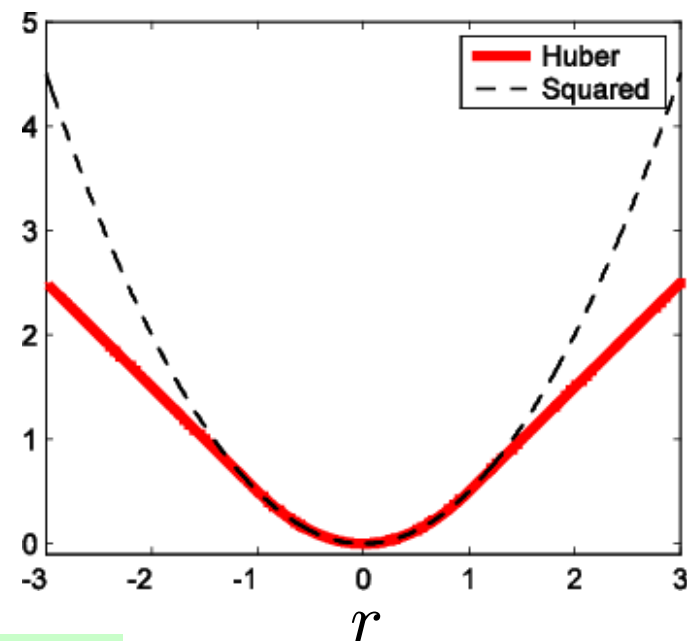


1. ℓ_1 -損失
2. フーバー損失
3. テューキー損失

強い外れ値の影響

- フーバー損失は ℓ_2 -損失と比べてロバスト性が高い
- しかし、損失に上界がないため
非常に強い外れ値の影響は受けてしまう

$$\min_{\theta} \sum_{i=1}^n \rho_{\text{Huber}}(f_{\theta}(\mathbf{x}_i) - y_i)$$

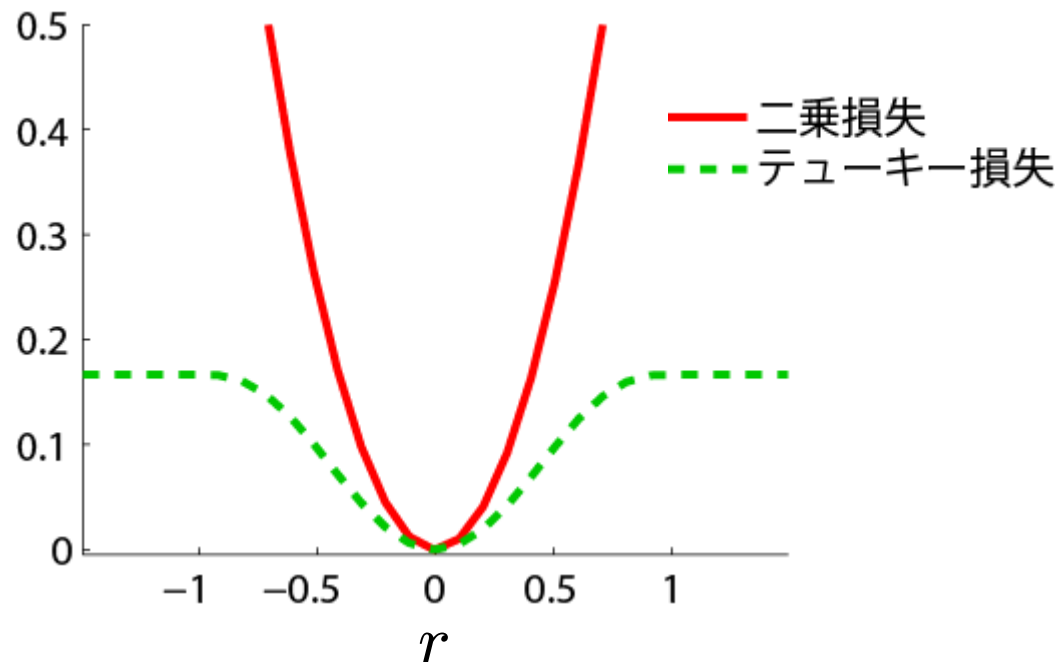


$$\rho_{\text{Huber}}(r) = \begin{cases} r^2/2 & (|r| \leq \eta) \\ \eta|r| - \eta^2/2 & (|r| > \eta) \end{cases}$$

■ 上界のある損失を考える

$$r = f_{\theta}(x) - y$$

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



一般の損失に対する 二次上界の最小化

$$\min_{\theta} \sum_{i=1}^n \rho(f_{\theta}(\mathbf{x}_i) - y_i)$$

- 微分可能で対称な損失 $\rho(r)$ に対して \tilde{r} で接する二次上界は(存在するなら)次式で与えられる:

$$\tilde{\rho}(r) = \frac{\tilde{w}}{2} r^2 + \text{const}$$

$$\tilde{w} = \rho'(\tilde{r})/\tilde{r}$$

証明は
宿題

- 繰り返し最小二乗アルゴリズム:

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^n \tilde{w}_i \left(f_{\theta}(\mathbf{x}_i) - y_i \right)^2$$

$$\tilde{w}_i = \rho'(\tilde{r}_i)/\tilde{r}_i$$

$$\tilde{r}_i = f_{\tilde{\theta}}(\mathbf{x}_i) - y_i$$

■チューキー損失

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

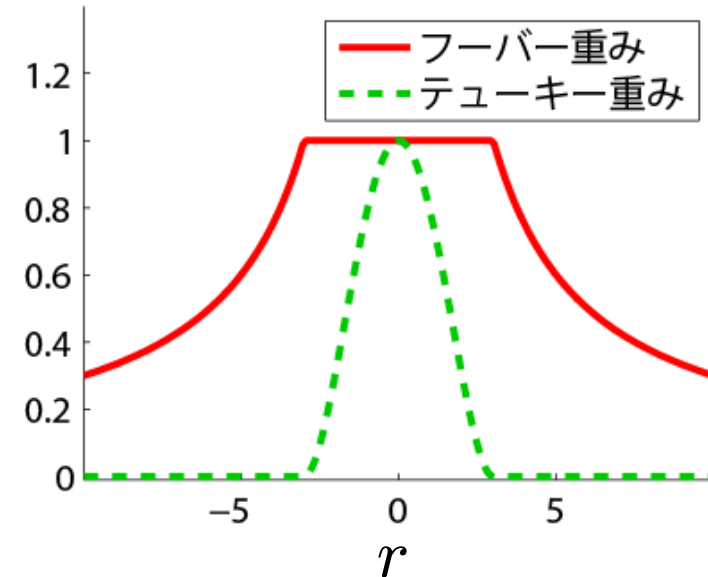
に対する重みは次式で与えられる:

$$\tilde{w} = \rho'(\tilde{r})/\tilde{r}$$

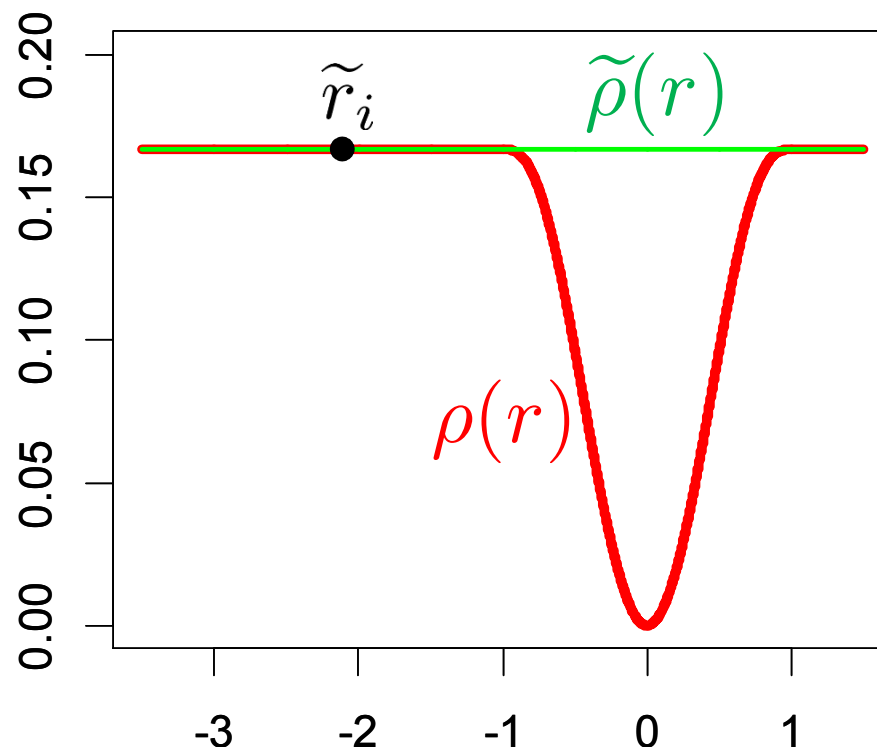
$$w_{\text{Tukey}} = \begin{cases} (1 - r^2/\eta^2)^2 & (|r| \leq \eta) \\ 0 & (|r| > \eta) \end{cases}$$

- 大きな残差に対する重みがゼロ

$$w_{\text{Huber}} = \begin{cases} 1 & (|r| \leq \eta) \\ \eta/|r| & (|r| > \eta) \end{cases}$$

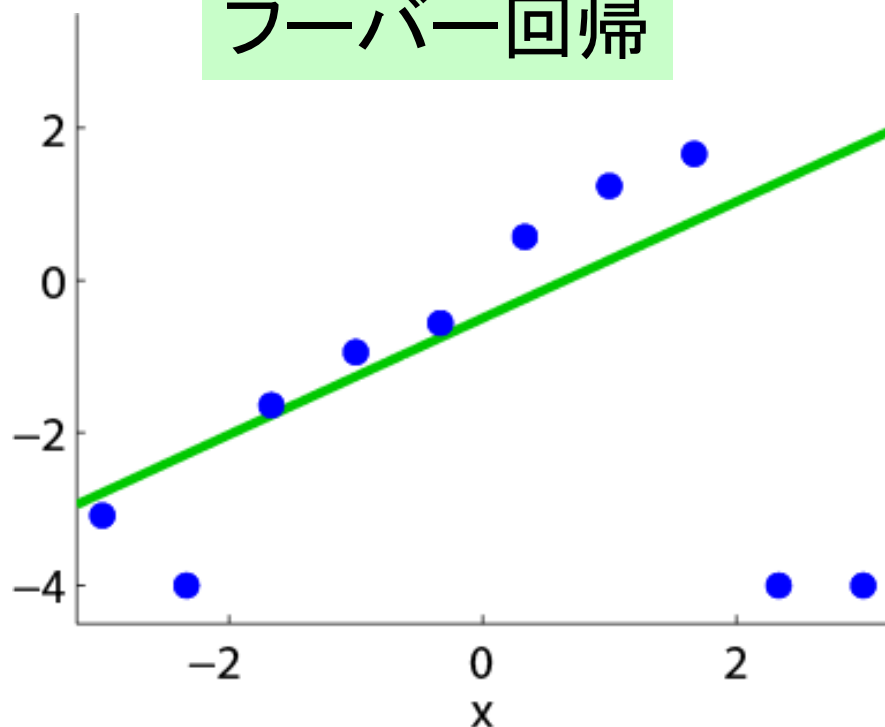


損失関数の非凸性

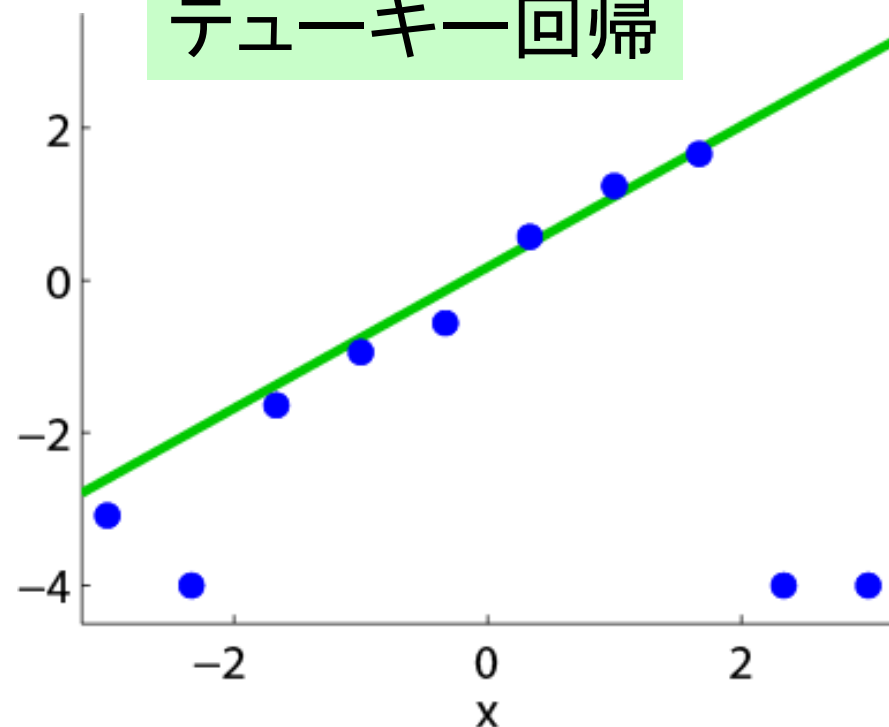


- 非凸な損失関数の二次上界をとった場合は最適解方向に進めるとは限らない

フーバー回帰



テューキー回帰



- テューキー回帰の方が外れ値に強い
- ただし、**非凸最適化**のため、得られる解は初期値の選び方に依存する

- **二乗損失** (平均値) は外れ値に弱い
- **絶対値損失** (中央値) は外れ値に強い
- **フーバー損失** はロバスト性と有効性のバランスがとれている
 - 解は解析的に求められない
- **テューキー損失** を用いるとロバスト性が更に向上
 - ただし非凸最適化になり最適化が困難

回帰のまとめ



1. 学習モデル(2章)
2. 最小二乗回帰(3章)
3. 正則化回帰(4章)
4. スパース回帰(5章)
5. ロバスト回帰(6章)

回帰のまとめ

■ 関数を学習するためのモデル:

- 線形モデル
- カーネルモデル
- 非線形モデル

■ 最小二乗回帰:

- 訓練標本との二乗誤差を最小化
- 解は解析的に計算できる

■ オンライン回帰

- データを1つずつ取り出して逐次的に学習
- 大量のデータを扱える

■ ℓ_2 -正則化回帰:

- 最小二乗回帰の過適合を軽減
- 解は解析的に計算できる
- モデル選択には交差確認法を用いる

■ ℓ_1 -正則化回帰(スパース回帰)

- 真のパラメータ値の多くがゼロとなるデータを適切に学習

■ ロバスト回帰

- 異常な値に対するロバスト性を強化

回帰のまとめ(続き)

■ 線形モデル／カーネルモデルに対する回帰法:

<div> <div></div> <div>制約条件</div> </div> <div> <div>損失関数</div> <div></div> </div>		無し	ℓ_2 -制約	ℓ_1 -制約
			正則化	正則化 & スパース
ℓ_2 -損失	有効	解析解	解析解	二次計画
フーバー損失		二次計画	二次計画	二次計画
ℓ_1 -損失	ロバスト	線形計画	二次計画	線形計画

次回の予告

■ 最小二乗分類(7章)



■ 線形モデル

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

 $\{\phi_j(\boldsymbol{x})\}_{j=1}^b$
: 基底関数

に対する重み付き最小二乗法

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{i=1}^n \tilde{w}_i \left(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) - y_i \right)^2$$

の解が次式で与えられることを示せ:

$$\hat{\boldsymbol{\theta}} = (\boldsymbol{\Phi}^\top \widetilde{\boldsymbol{W}} \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \widetilde{\boldsymbol{W}} \boldsymbol{y}$$

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi_1(\boldsymbol{x}_1) & \cdots & \phi_b(\boldsymbol{x}_1) \\ \vdots & \ddots & \vdots \\ \phi_1(\boldsymbol{x}_n) & \cdots & \phi_b(\boldsymbol{x}_n) \end{pmatrix}$$

$$\widetilde{\boldsymbol{W}} = \text{diag}(\tilde{w}_1, \dots, \tilde{w}_n)$$

$$\boldsymbol{y} = (y_1, \dots, y_n)^\top$$

宿題1のヒント

48

- $\Phi^T \widetilde{W} \Phi$ を展開したとき各 \tilde{w}_i にどのような行列が掛かっているか調べる

- 微分可能で対称な損失 $\rho(r)$ に対して \tilde{r} で接する二次上界は(存在するなら)次式で与えられることを示せ

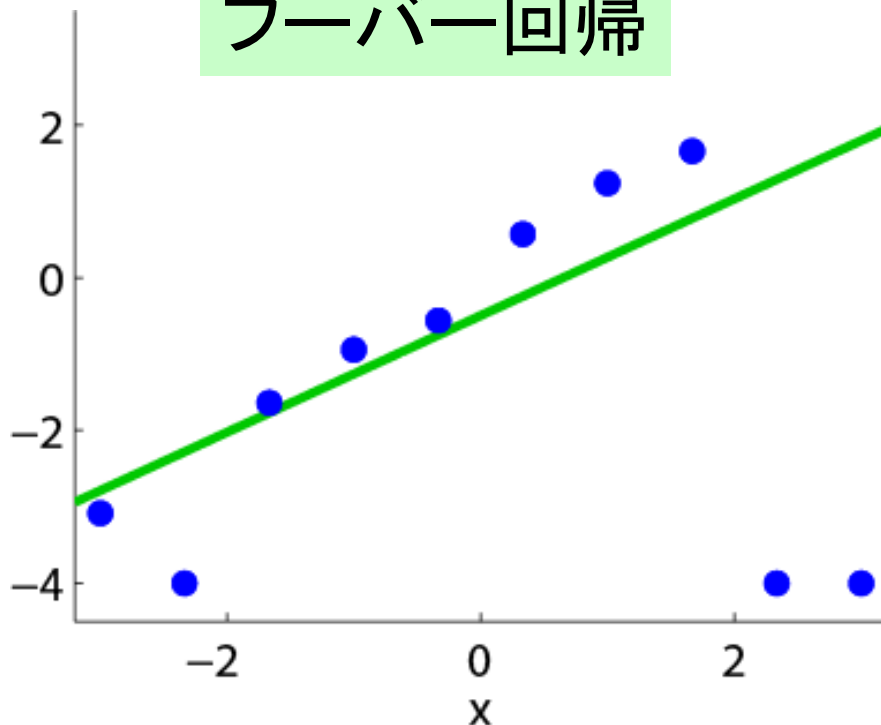
$$\tilde{\rho}(r) = \frac{\tilde{w}}{2} r^2 + \text{const}$$

$$\tilde{w} = \rho'(\tilde{r})/\tilde{r}$$

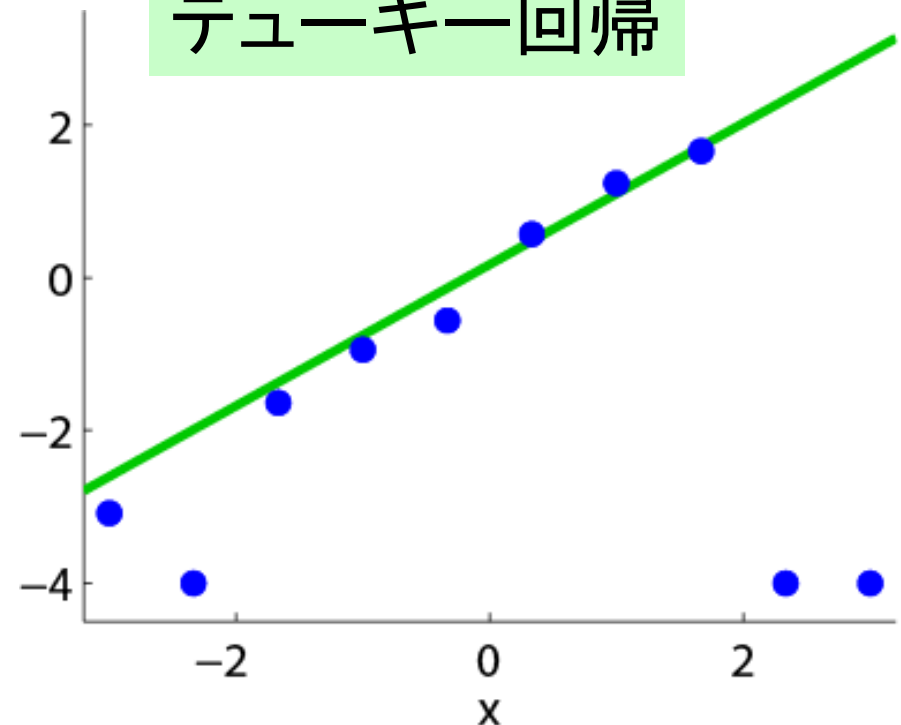
- const : r に依存しない値

- 直線モデル $f_{\theta}(x) = \theta_1 + \theta_2 x$ に対して、
テューキー回帰の繰り返し最小二乗
アルゴリズムを実装せよ（結果は初期値に
依存して変わる場合があることに注意）

フーバー回帰



テューキー回帰



宿題3の続き

- データとしては例えば以下のものを用いてもよい

```
clear all; rand('state',0); randn('state',0);  
n=10; N=1000; x=linspace(-3,3,n)';  
y=x+0.2*randn(n,1); y(n)=-4; y(n-1)=-4; y(2)=-4;
```