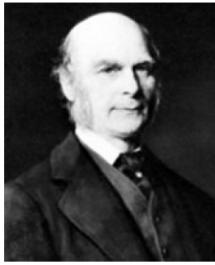


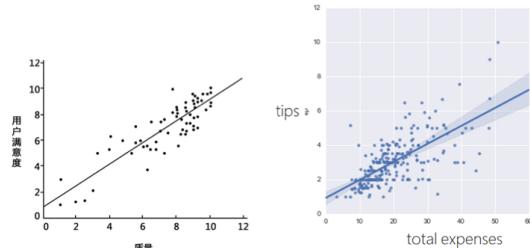
# Regression

- Proposed by Francis Galton (left) and Karl Pearson (right), in the publication "Regression towards mediocrity in hereditary"
- The characteristics (e.g., height) in the offspring regress towards a mediocre point (mean) of that of their parents
- Generalization : predict the dependent variables  $y$  from the independent variables  $\mathbf{x}$  :  $y = f(\mathbf{x})$  or  $y = E[y|\mathbf{x}]$



## Applications

- Predict medical expenses from the individual profiles of the patients
- Predict the scores on Douban from the quality of the movies
- Predict the tips from the total expenses

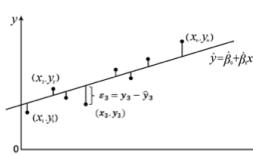


回归是一神有监督学习方法

## Univariate Linear Model 一元线性模型

- Linear model :  $y = w_0 + w_1 x + \epsilon$ , where  $w_0$  and  $w_1$  are regression coefficients,  $\epsilon$  is the error or noise
- Assume  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ , where  $\sigma^2$  is a fixed but unknown variance; then  $y|x \sim \mathcal{N}(w_0 + w_1 x, \sigma^2)$
- Assume the samples  $\{(x_i, y_i)\}_{i=1}^n$  are generated from this conditional distribution, i.e.,  $y_i|x_i \sim \mathcal{N}(w_0 + w_1 x_i, \sigma^2)$
- Intuitively, find the best straight line ( $w_0$  and  $w_1$ ) such that the sample points fit it well, i.e., the residuals are minimized,

$$(\hat{w}_0, \hat{w}_1) = \arg \min_{w_0, w_1} \sum_{i=1}^n (y_i - w_0 - w_1 x_i)^2$$



样本点大致分布在直线附近  
采样会有误差、噪音。  
噪音  $\epsilon$  一般满足  $N(0, \sigma^2)$ ,  $\sigma^2$  确定但未知  
因此  $y|x \sim N(w_0 + w_1 x, \sigma^2)$ . 并假定所有样本符合这个分布。

用最小二乘法能找到  $\hat{w}_0$  和  $\hat{w}_1$ .

使误差在平均意义上不大，使残差平方和最小。

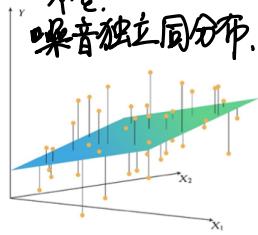
# Multivariate Linear Model 多元线性模型

- Linear model:  $y = f(\mathbf{x}) + \epsilon = w_0 + w_1 x_1 + \dots + w_p x_p + \epsilon$ , where  $w_0, w_1, \dots, w_p$  are regression coefficients,  $\mathbf{x} = (x_1, \dots, x_p)^T$  is the input vector whose components are independent variables or attribute values,  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  is the noise
- For the size  $n$  samples  $\{(x_i, y_i)\}_{i=1}^n$ , let  $\mathbf{y} = (y_1, \dots, y_n)^T$  be the response or dependent variables,  $\mathbf{w} = (w_0, w_1, \dots, w_p)^T$ ,  $\mathbf{X} = [1, (x_1, \dots, x_n)^T] \in \mathbb{R}^{n \times (p+1)}$ , and  $\epsilon = (\epsilon_1, \dots, \epsilon_n)^T \sim \mathcal{N}(0, \sigma^2 I_n)$ .

模型输出  $\hat{\mathbf{y}} = \mathbf{X} \vec{\mathbf{w}}$  有  $n$  个样本就有  $n$  个单位矩阵

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon$$

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix} = \begin{pmatrix} \mathbf{1} \\ \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{pmatrix}$$



同样有  $y | \mathbf{x} \sim \mathcal{N}(w_0 + w_1 x_1 + \dots + w_p x_p, \sigma^2)$

样本  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , 同样有最小二乘法.

$$\min_{w_0, \dots, w_p} \sum_{i=1}^n (y_i - w_0 - w_1 x_{i,1} - w_2 x_{i,2} - \dots - w_d x_{i,d})^2$$

$$||\mathbf{y}_i - (w_0, w_1, w_2, \dots, w_d) \begin{pmatrix} 1 \\ x_{i,1} \\ x_{i,2} \\ \vdots \\ x_{i,d} \end{pmatrix}||^2$$

$$\Rightarrow \min_{\vec{\mathbf{w}}} \sum_{i=1}^n (y_i - \vec{\mathbf{w}}^T \vec{\mathbf{x}}_i)^2$$

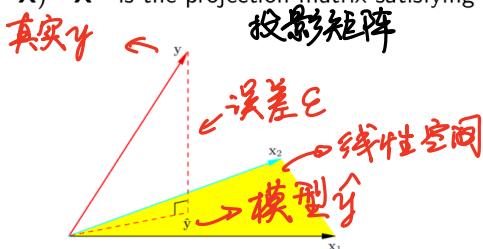
$$\Rightarrow \min_{\vec{\mathbf{w}}} \|\vec{\mathbf{y}} - \vec{\mathbf{X}} \vec{\mathbf{w}}\|_2^2$$

$$\vec{\mathbf{y}} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \vec{\mathbf{X}} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}, \vec{\mathbf{w}} = \begin{pmatrix} w_0 \\ \vdots \\ w_p \end{pmatrix}$$

误差、投影、极大似然估计都能解释为何要最小化残差平方和.

## Least Square (LS)

- Minimize the total residual sum-of-squares :  
 $RSS(\mathbf{w}) = \sum_{i=1}^n (y_i - w_0 - w_1 x_1 - \dots - w_p x_p)^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$
- When  $\mathbf{X}^T \mathbf{X}$  is invertible, the minimizer  $\hat{\mathbf{w}}$  satisfies  
 $\nabla_{\mathbf{w}} RSS(\hat{\mathbf{w}}) = 0 \Rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- The prediction  $\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{P}\mathbf{y}$  is a projection of  $\mathbf{y}$  onto the linear space spanned by the column vectors of  $\mathbf{X}$ ;  
 $\mathbf{P} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  is the projection matrix satisfying  $\mathbf{P}^2 = \mathbf{P}$



$$RSS(\vec{\mathbf{w}}) = \|\vec{\mathbf{y}} - \vec{\mathbf{X}} \vec{\mathbf{w}}\|_2^2$$

$$\Rightarrow \frac{\partial RSS(\vec{\mathbf{w}})}{\partial \vec{\mathbf{w}}} = -2 \vec{\mathbf{X}}^T (\vec{\mathbf{y}} - \vec{\mathbf{X}} \vec{\mathbf{w}}) = \vec{0}$$

$$\Rightarrow \hat{\vec{\mathbf{w}}} = (\vec{\mathbf{X}}^T \vec{\mathbf{X}})^{-1} \vec{\mathbf{X}}^T \vec{\mathbf{y}}$$

即将不在线性空间上的  $\mathbf{y}$  投影到线性空间上, 能使得误差最小.

## Maximal Likelihood Estimate (MLE)

- A probabilistic viewpoint :  
 $y | \mathbf{x} \sim \mathcal{N}(w_0 + w_1 x_1 + \dots + w_p x_p, \sigma^2)$
- Likelihood function :  $\rightarrow$  独立  $\Rightarrow$  每一个的乘积.  
 $L(\mathbf{w}; \mathbf{X}, \mathbf{y}) = P(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{i=1}^n P(y_i | x_i, \mathbf{w})$  with  
 $P(y_i | x_i, \mathbf{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - w_0 - w_1 x_{i,1} - \dots - w_p x_{i,p})^2}{2\sigma^2}}$
- Maximal likelihood estimate : given the samples from some unknown parametric distribution, find the parameters such that the samples the most probably seem to be drawn from that distribution, i.e.,  $\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} L(\mathbf{w}; \mathbf{X}, \mathbf{y})$
- Equivalent to maximize the log-likelihood function  
 $I(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \log L(\mathbf{w}; \mathbf{X}, \mathbf{y}) \rightarrow$  最大化  
 $-n \log(\sqrt{2\pi\sigma^2}) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w_0 - w_1 x_{i,1} - \dots - w_p x_{i,p})^2 \rightarrow$  最小化  
 $\rightarrow$  即最小二乘法.
- The same minimizer as LS :  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

$y_i$  独立同分布.

对  $\mathbf{w}$  叫似然函数, 对  $\mathbf{y}$  叫联合概率.

$\rightarrow$  即最小二乘法.

# Projection by Orthogonalization

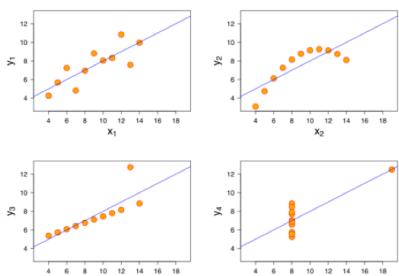
- Another useful formulation : let  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ ,  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ , then OLS can be formulated by using the centralized data  $\{\tilde{x}_i, \tilde{y}_i\}_{i=1}^n = \{x_i - \bar{x}, y_i - \bar{y}\}_{i=1}^n$ ,  $RSS(\tilde{w}) = \sum_{i=1}^n (\tilde{y}_i - w_1 \tilde{x}_{i1} - \dots - w_p \tilde{x}_{ip})^2 = \|\tilde{y} - \tilde{X} \tilde{w}\|_2^2$ , with  $\hat{w}_0 = \bar{y} - \tilde{w}^T \tilde{x}$
- Ordinary least square (OLS) prediction  $\hat{y} = \mathbf{P}y$  is the projection of  $y$  on the linear space spanned by the columns of  $\mathbf{X}$ , i.e.,  $\mathcal{X} = \text{Span}\{x_{\cdot,0}, x_{\cdot,1}, \dots, x_{\cdot,p}\}$ , recall that  $x_{\cdot,0} = \mathbf{1}_n$
- If  $\{x_{\cdot,0}, x_{\cdot,1}, \dots, x_{\cdot,p}\}$  forms a set of orthonormal basis, then  $\hat{y} = \sum_{i=0}^p \langle y, x_{\cdot,i} \rangle x_{\cdot,i}$  → 直接投影
- If not, we can first do orthogonalization by Gram-Schmidt procedure for the set  $\{x_{\cdot,0}, x_{\cdot,1}, \dots, x_{\cdot,p}\}$
- Similar orthogonalization procedures can be done by QR decomposition or SVD of the matrix  $\mathbf{X}^T \mathbf{X}$  (classic topics in numerical linear algebra)

# Regression by Successive Orthogonalization

- The expansion of  $y$  on the standard orthonormal basis after Gram-Schmidt procedure can be summarised in the following algorithm :
  - Initialize  $z_0 = x_0 = \mathbf{1}_n$
  - For  $j = 1, \dots, p$  :
    - Regress  $x_j$  on  $\{z_0, \dots, z_{j-1}\}$  to produce coefficients  $\hat{\gamma}_{lj} = \langle z_l, x_j \rangle / \langle z_l, z_l \rangle$  with  $l = 0, \dots, j-1$  and residual vectors  $z_j = x_j - \sum_{k=0}^{j-1} \hat{\gamma}_{kj} z_k$
    - Regress  $y$  on the residual  $z_p$  to give the estimate  $\hat{w}_p$
- If  $x_p$  is highly correlated with some of the other  $x_k$ 's, the residual vector  $z_p$  will be close to zero ; in such situation, the coefficient  $\hat{w}_p$  with small Z-score  $\frac{\hat{w}_p}{\hat{\sigma}_p}$  could be thrown out, where  $\hat{\sigma}_p^2 = \frac{\hat{\delta}_p^2}{\|z_p\|_2^2}$  is an estimate of  $\text{Var}(\hat{w}_p) = \frac{\sigma^2}{\|z_p\|_2^2}$

# Shortcomings of Fitting Nonlinear Data

- Evaluating the model by Coefficient of Determination  $R^2$  :  $R^2 := 1 - \frac{SS_{res}}{SS_{tot}}$  ( $= \frac{SS_{reg}}{SS_{tot}}$  for linear regression), where  $SS_{tot} = \sum_{i=1}^n (y_i - \bar{y})^2$  is the total sum of squares,  $SS_{reg} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$  is the regression sum of squares, and  $SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$  is the residual sum of squares.
- The larger the  $R^2$ , the better the model



可用决定系数  $R^2$  评判效果

→ 总平方和

→ 回归平方和

→ 残差平方和

$$S_{tot} = S_{res} + S_{reg}$$

$R^2 \in [0, 1]$ , 模型越好,  $R^2$  越大.

因为用线性模型拟合非线性数据本来不合理,  
所以拟合效果一般较差

$$\begin{aligned} S_{tot} &= \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) \\ &= S_{res} + S_{reg} + 2(\vec{y} - \vec{\hat{y}})^T (\vec{\hat{y}} - \vec{\bar{y}}) \\ &= S_{res} + S_{reg} + 2(\vec{y} - \vec{\hat{y}})^T (\vec{X} \vec{w} - \vec{\bar{y}}) \\ &= S_{res} + S_{reg} + 2(\vec{y} - \vec{\hat{y}})^T \vec{X} \vec{w} - 2\vec{y}^T (\vec{y} - \vec{\hat{y}})^T \\ \because (\vec{y} - \vec{\hat{y}})^T \vec{X} &= \vec{y}^T \vec{X} - \vec{\hat{y}}^T \vec{X} = \vec{0}, (\vec{y} - \vec{\hat{y}})^T \vec{I} = 0 \\ \therefore S_{tot} &= S_{res} + S_{reg} \end{aligned}$$

# Multicollinearity 多重共线性

- If the columns of  $\mathbf{X}$  are almost linearly dependent, i.e., multicollinearity, then  $\det(\mathbf{X}^T \mathbf{X}) \approx 0$ , the diagonal entries in  $(\mathbf{X}^T \mathbf{X})^{-1}$  is quite large. This implies the variances of  $\hat{\mathbf{w}}$  get large, and the estimate is not accurate
- Eg : 10 samples are drawn from the true model  
 $y = 10 + 2x_1 + 3x_2 + \epsilon$ ; the LS estimator is  $\hat{w}_0 = 11.292$ ,  $\hat{w}_1 = 11.307$ ,  $\hat{w}_2 = -6.591$ , far from the true coefficients; correlation coefficient is  $r_{12} = 0.986$
- Remedies : ridge regression, principal component regression, partial least squares regression, etc.

↓ 每一组中  $x_{1,i}$  与  $x_{2,i}$  的取值相近,  $x_1$  与  $x_2$  线性相关, 几乎不独立.

No.	1	2	3	4	5	6	7	8	9	10
$x_1$	1.1	1.4	1.7	1.7	1.8	1.8	1.9	2.0	2.3	2.4
$x_2$	1.1	1.5	1.8	1.7	1.9	1.8	1.8	2.1	2.4	2.5
$\epsilon_i$	0.8	-0.5	0.4	-0.5	0.2	1.9	1.9	0.6	-1.5	-1.5
$y_i$	16.3	16.8	19.2	18.0	19.5	20.9	21.1	20.9	20.3	22.0

若有多元线性变量, 会出现多重共线性.

可能有两个坐标分量是线性相关的.  $\det(\mathbf{X}^T \mathbf{X}) \approx 0$ .  
 此时  $(\mathbf{X}^T \mathbf{X})^{-1}$  的值会较大.

$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  的方差会变大.

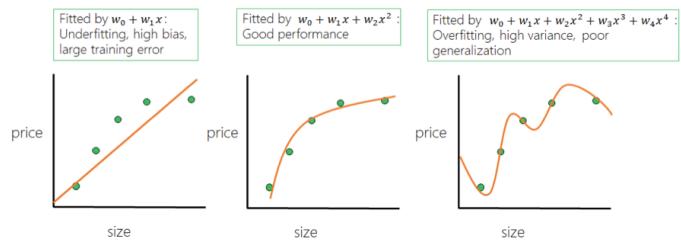
可用岭回归、主成分回归、偏最小二乘回归优化.

线性相关, 几乎不独立. (主成分分析)  
 可踢掉  $x_1$  或  $x_2$ , 或通过  $x_1$  和  $x_2$  生成  $\vec{x}_3$

本质上都是为了消除多重共线性.

## Overfitting

- Easily to be overfitted when introducing more variables, e.g., regress housing price with housing size
- The high degree model also fits the noises in the training data, so generalizes poorly to new data
- Remedy : regularization 正则化



# Bias-Variance Decomposition 偏差-方差分解

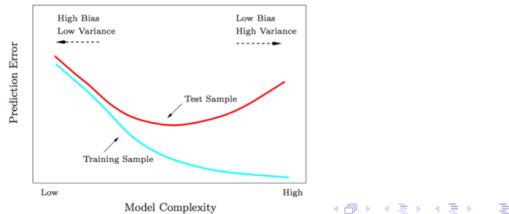
- Bias-variance decomposition of generalization error in  $L^2$  loss :

$$E_{train} R_{exp}(\hat{f}(x)) = E_{train} E_P[(y - \hat{f}(x))^2 | x] = \underbrace{\text{Var}(\hat{f}(x))}_{\text{variance}} + \underbrace{\text{Bias}^2(\hat{f}(x))}_{\text{bias}} + \underbrace{\sigma^2}_{\text{noise}}$$

train表示测试数据的联合分布

where  $P = P(y|x)$  is the conditional probability of  $y$  given  $x$

- Bias :  $\text{Bias}(\hat{f}(x)) = E_{train} \hat{f}(x) - f(x)$  is the average accuracy of prediction for the model (deviation from the truth)
- Variance :  $\text{Var}(\hat{f}(x)) = E_{train} (\hat{f}(x) - E_{train} \hat{f}(x))^2$  is the variability of the model prediction due to different data set (stability)



Model  $y = f(x) + \epsilon$ , with  $E(\epsilon) = 0$  and  $\text{Var}(\epsilon) = \sigma^2$  (system error)

$$\begin{aligned} E_{train} R_{exp}(\hat{f}(x)) &= E_P[(y - f(x))^2 | x] + E_{train}[(f(x) - \hat{f}(x))^2] \\ &\quad + 2 E_{train} E_P[(y - f(x))(f(x) - \hat{f}(x)) | x] \\ &\quad \text{vanishes since } E_P(y - f(x))|x=0 \rightarrow E(\epsilon)=0 \\ &= \sigma^2 + E_{train}[(f(x) - E_{train} \hat{f}(x))^2] + E_{train}[(E_{train} \hat{f}(x) - \hat{f}(x))^2] \\ &\quad + 2 E_{train}[(f(x) - E_{train} \hat{f}(x))(E_{train} \hat{f}(x) - \hat{f}(x))] \\ &\quad \text{vanishes since } E_{train}[E_{train} \hat{f}(x) - \hat{f}(x)] = 0 \\ &= \sigma^2 + \text{Bias}^2(\hat{f}(x)) + \text{Var}(\hat{f}(x)) \end{aligned}$$

对训练样本

有一个大的泛化误差就不会小

The more complicated the model, the lower the bias, but the higher the variance.

## Bias-Variance Decomposition = kNN Regression

- kNN can be used to do regression if the mode (majority vote) is replaced by mean :  $\hat{f}(x) = \frac{1}{k} \sum_{x(i) \in N_k(x)} y(i)$  → 最近的  $k$  个点的平均预测值就是当前点的值.
- Generalization error of kNN regression is

$$\begin{aligned} E_{train} R_{exp}(\hat{f}(x)) &= \sigma^2 + (f(x) - \frac{1}{k} \sum_{x(i) \in N_k(x)} f(x(i)))^2 \\ &\quad \text{bias}^2 \\ &\quad + E_{train} \left[ \frac{1}{k} \sum_{x(i) \in N_k(x)} (y(i) - f(x(i))) \right]^2 \\ &\quad \text{variance}^2 \\ &\quad \underbrace{\frac{1}{k} \sigma^2}_{\text{noise}} \end{aligned}$$

where we have used the fact that  $E_{train} y_i = f(x_i)$  and  $\text{Var}(y_i) = \sigma^2$ .

- For small  $k$ , overfitting, bias ↘, variance ↗
- For large  $k$ , underfitting, bias ↗, variance ↘

$$\begin{aligned} \hat{f} &= \hat{\omega} = (X^T X)^{-1} X^T y \quad (\text{线性回归}) \\ R_{exp}(\hat{f}(x)) &= E_{P(y|x)} (y - \hat{f}(x))^2 \quad (\text{泛化误差, 给定 } x) \\ &\quad E_{P(y|x)} (y - f(x))^2 = 0 \\ &\quad -f(x) + \hat{f}(x) \\ \text{Generalization error} &= E_{\{(x_i, y_i)\} \sim P} R_{exp}(\hat{f}(x)) = E_{\{(x_i, y_i)\} \sim P} E_{P(y|x)} (y - \hat{f}(x))^2 \\ &= E_{\{(x_i, y_i)\} \sim P} E_{P(y|x)} [(y - f(x))^2 + (f(x) - \hat{f}(x))^2 + 2(y - f(x))(f(x) - \hat{f}(x))] \\ &\quad \text{noise}^2 \quad -E_{\{(x_i, y_i)\}} \hat{f}(x) + E_{\{(x_i, y_i)\}} \hat{f}(x) \\ &= E_{\{(x_i, y_i)\} \sim P} E_{P(y|x)} [(y - f(x))^2 + (f(x) - E_{x,y} \hat{f}(x))^2 \rightarrow \text{bias}^2] \\ &\quad + (E_{x,y} \hat{f}(x) - \hat{f}(x))^2 \rightarrow \text{Variance} \\ &\quad + 2(f(x) - E_{x,y} \hat{f}(x))(E_{x,y} \hat{f}(x) - \hat{f}(x)) \\ &= 0 \end{aligned}$$

模型越复杂，偏差越小、方差越大。

模型复杂后，越适应于当前训练样本 → 偏差小  
但引入测试集后，更容易造成模型变化  
→ 方差大。

当  $\text{Var}(x) + \text{Bias}^2(x)$  最小时最好

$k$  越小，方差越大，偏差越小，过拟合

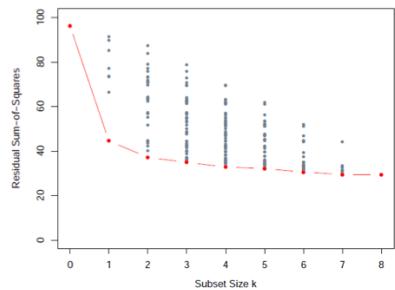
# Regulation by Subset Selection

- In high dimensions, the more the input attributes, the larger the variance
- Shrinking some coefficients or setting them to zero can reduce the overfitting
- Using less input variables also help interpretation with the most important variables
- Subset selection: retaining only a subset of the variables, while eliminating the rest variables from the model
- Best-subset selection : find for each  $k \in \{0, 1, \dots, p\}$  the subset  $S_k \subset \{1, \dots, p\}$  of size  $k$  that gives the smallest  $RSS(\mathbf{w}) = \sum_{i=1}^n (y_i - w_0 - \sum_{j \in S_k} w_j x_{ij})^2$

每次只选一部分最相关的变量

## Best-Subset Selection

- The best subset of size  $k + 1$  may not include the the variables in the best subset of size  $k$
- The RSS of the best subset of size  $k$  is not necessarily decreasing with  $k$
- Choose  $k$  based on bias-variance tradeoff, usually by AIC and BIC, or practically by cross-validation



## Forward (Backward) Stepwise Selection

- Forward-stepwise selection : start with the intercept  $\bar{y}$ , then sequentially add into the model the variables that improve the fit most (reduce RSS most)
- QR factorization helps search the candidate variables to add
- Greedy algorithm : the solution could be sub-optimal
- Computationally more efficient than best-subset selection ; statistically the constrained search enjoys lower variance than best-subset selection
- Backward-stepwise selection : start with the full model, then sequentially delete from the model the variables that has the least impact on the fit most (the candidate for dropping is the variable with the smallest Z-score) ; can only be used when  $n > p$  in order to fit the full model by OLS

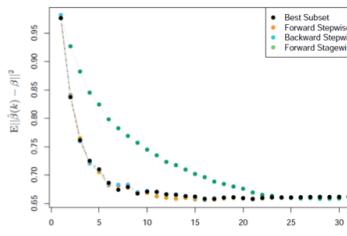
正向每次加被模型增强的变量

贪心算法

反过来一步步从模型中减少变量.

## FS Selection

- Starts with the intercept and centered variables with 0 coefficients
- At each step, identify the variables (among all variables) most correlated with the current residual, then regress the residual on this chosen variable and increment the current coefficient with the new regression coefficient
- Ends when no variables are correlated with the residual (arrive at the OLS fit when  $n > p$ )
- Slower than forward-stepwise : the other variables and their coefficients are not changed at each step except the chosen variable

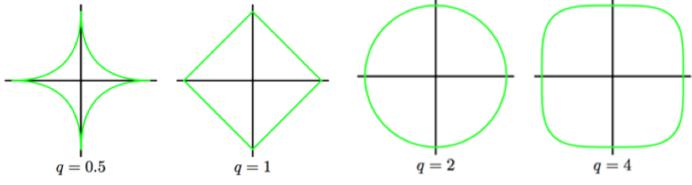


# Regularization by Penalties

- Add a penalty term, in general  $l_q$ -norm

$$\begin{aligned} & \sum_{i=1}^n (y_i - w_0 - w_1 x_1 - \cdots - w_p x_p)^2 + \lambda \|\mathbf{w}\|_q^q \\ & = \|\mathbf{y} - \mathbf{Xw}\|_2^2 + \lambda \|\mathbf{w}\|_q^q \end{aligned}$$

- $q = 2$  : ridge regression
- $q = 1$  : LASSO regression



## Ridge Regression

- The optimization problem turns to be

$$\begin{aligned} \hat{\mathbf{w}} &= \arg \min_{\mathbf{w}} \sum_{i=1}^n (y_i - w_0 - w_1 x_1 - \cdots - w_p x_p)^2 + \lambda \|\mathbf{w}\|_2^2 \\ &= \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{Xw}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \end{aligned}$$

- $\lambda \geq 0$  is a fixed parameter which has to be tuned by cross-validation

- Equivalent to the constraint minimization problem :

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{Xw}\|_2^2, \quad \text{subject to} \quad \|\mathbf{w}\|_2 \leq \mu,$$

where  $\mu \geq 0$  is a prescribed threshold (tuning parameter)

- The large  $\lambda$  corresponds to the small  $\mu$ .

$$\lambda \rightarrow \mu$$

## Solving Ridge Regression

- Easy to show that  $\hat{\mathbf{w}}^{ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{p+1})^{-1} \mathbf{X}^T \mathbf{y}$
- The estimator is also a projection of  $\mathbf{y}$  :  

$$\hat{\mathbf{y}}^{ridge} = \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{p+1})^{-1} \mathbf{X}^T \mathbf{y}$$
- $\mathbf{X}$  can be diagonalized by SVD :  $\mathbf{X} = \mathbf{P} \mathbf{D} \mathbf{Q}^T$  with  
 $\mathbf{D} = \text{diag}(\nu_1, \dots, \nu_{p+1})$ , and  $\mathbf{P} \in \mathbb{R}^{n \times (p+1)}$ ,  $\mathbf{Q} \in \mathbb{R}^{(p+1) \times (p+1)}$   
 being orthogonal matrices ( $\mathbf{P}^T \mathbf{P} = \mathbf{I}_{p+1}$ )
- $\hat{\mathbf{y}}^{ridge} = \mathbf{P} \text{diag}(\frac{\nu_1^2}{\nu_1^2 + \lambda}, \dots, \frac{\nu_{p+1}^2}{\nu_{p+1}^2 + \lambda}) \mathbf{P}^T \mathbf{y}$ , while  $\hat{\mathbf{y}}^{OLS} = \mathbf{P} \mathbf{P}^T \mathbf{y}$
- In the spectral space, the ridge regression estimator is a shrinkage of the OLS estimator ( $\lambda = 0$ )

在线性回归中，多重共线性导致  $\mathbf{X}$  比较奇异， $\det(\mathbf{X}^T \mathbf{X}) \approx 0$ 。但是由于  $\mathbf{X}^T \mathbf{X}$  是半正定的， $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{p+1}$  必为正定的，能算出较稳定的逆。因此岭回归的效果会比一般的线性回归效果好。

$$J(\mathbf{w}) = \|\mathbf{y} - \mathbf{Xw}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

$$\nabla_{\mathbf{w}} J = 2 \mathbf{X}^T (\mathbf{Xw} - \mathbf{y}) + 2\lambda \mathbf{w} = 0$$

$$\Rightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\Rightarrow \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\mathbf{w}} = \mathbf{X} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\mathbf{X} = \mathbf{P} \Lambda \mathbf{Q}^T \text{ for } \mathbf{P} \in \mathbb{R}^{n \times (d+1)}, \mathbf{Q} \in \mathbb{R}^{(d+1) \times (d+1)}$$

$$\begin{array}{ll} \text{奇异值分解} & \mathbf{P}^T \mathbf{P} = \mathbf{I}_{d+1} \quad \mathbf{Q} \mathbf{Q}^T = \mathbf{Q}^T \mathbf{Q} = \mathbf{I}_{d+1} \\ (\text{SVD}) & \end{array}$$

$$\Rightarrow \hat{\mathbf{y}} = \mathbf{P} \Lambda \mathbf{Q}^T (\mathbf{Q} \Lambda \mathbf{P}^T \mathbf{P} \Lambda \mathbf{Q}^T + \lambda \mathbf{I})^{-1} \mathbf{Q} \Lambda \mathbf{P}^T \mathbf{y}$$

$$= \mathbf{P} \Lambda \mathbf{Q}^T (\mathbf{Q} (\lambda^2 + \lambda \mathbf{I}) \mathbf{Q}^T)^{-1} \mathbf{Q} \Lambda \mathbf{P}^T \mathbf{y}$$

$$= \mathbf{P} \Lambda \mathbf{Q}^T \mathbf{Q} (\lambda^2 + \lambda \mathbf{I})^{-1} \mathbf{Q}^T \mathbf{Q} \Lambda \mathbf{P}^T \mathbf{y}$$

$$= \mathbf{P} \Lambda (\lambda^2 + \lambda \mathbf{I})^{-1} \Lambda \mathbf{P}^T \mathbf{y}$$

$$= \mathbf{P} \text{diag}(\frac{\nu_1^2}{\nu_1^2 + \lambda}, \dots, \frac{\nu_{d+1}^2}{\nu_{d+1}^2 + \lambda}) \mathbf{P}^T \mathbf{y}$$

通过入将特征值做了收缩  
 入控制了泛化误差，需取一个不大不小的数。

# Bayesian Viewpoint of Ridge Regression

- Given  $\mathbf{X}$  and  $\mathbf{w}$ , the conditional distribution of  $\mathbf{y}$  is  
 $P(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}) \propto \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w})\right)$
- In addition, assume  $\mathbf{w}$  has a prior distribution  
 $P(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0) \propto \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_0)^T \boldsymbol{\Lambda}_0^{-1}(\mathbf{w} - \boldsymbol{\mu}_0)\right)$
- By Bayes theorem, the posterior distribution of  $\mathbf{w}$  given the data  $\mathbf{X}$  and  $\mathbf{y}$  is

$$P(\mathbf{w}|\mathbf{X}, \mathbf{y}) \propto P(\mathbf{y}|\mathbf{X}, \mathbf{w})P(\mathbf{w})$$

$$\begin{aligned} &\propto \exp\left(-\frac{1}{2\sigma^2}(\mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} - 2\mathbf{y}^T \mathbf{X}\mathbf{w})\right) \\ &\quad - \frac{1}{2}(\mathbf{w}^T \boldsymbol{\Lambda}_0^{-1} \mathbf{w} - 2\boldsymbol{\mu}_0^T \boldsymbol{\Lambda}_0^{-1} \mathbf{w}) \\ &\propto \exp\left(-\frac{1}{2}(\mathbf{w} - \boldsymbol{\mu}_m)^T \boldsymbol{\Lambda}_m^{-1}(\mathbf{w} - \boldsymbol{\mu}_m)\right) \end{aligned}$$

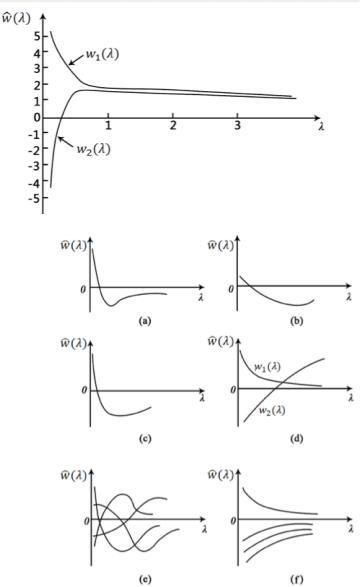
where  $\boldsymbol{\Lambda}_m = (\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \boldsymbol{\Lambda}_0^{-1})^{-1}$  and  $\boldsymbol{\mu}_m = \boldsymbol{\Lambda}_m(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{y} + \boldsymbol{\Lambda}_0^{-1} \boldsymbol{\mu}_0)$

- If  $\boldsymbol{\mu}_0 = 0$  and  $\boldsymbol{\Lambda}_0 = \frac{\sigma^2}{\lambda} \mathbf{I}_{p+1}$ , then  $\hat{\mathbf{w}} = \boldsymbol{\mu}_m = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{p+1})^{-1} \mathbf{X}^T \mathbf{y}$  maximizes the posterior probability  $P(\mathbf{w}|\mathbf{X}, \mathbf{y})$

## Ridge Trace 谷迹

- The functional plot of  $\hat{\mathbf{w}}^{ridge}(\lambda)$  with  $\lambda$  is called ridge trace
- The large variations in ridge trace indicate the multicollinearity in variables
- When  $\lambda \in (0, 0.5)$ , the ridge traces have large variations, it suggests to choose  $\lambda = 1$
- Before plot ridge trace, do scaling for the variables
- The coefficients with stable trace and small absolute values should have little influence on  $y$ , as in (a)
- The coefficients with large stable absolute values should have great impact on  $y$ , as in (b) and (c)
- The ridge traces of the coefficients of two variables are not stable, but the sum of the coefficients is stable. This implies the multicollinearity as in (d)
- The stable ridge traces of all variables suggest good performance using OLS as in (f)

$\lambda$	0	0.1	0.15	0.2	0.3	0.4	0.5	1.0	1.5	2.0	3.0
$\hat{w}_1^{ridge}(\lambda)$	11.31	3.48	2.99	2.71	2.39	2.20	2.06	1.66	1.43	1.27	1.03
$\hat{w}_2^{ridge}(\lambda)$	-6.59	0.63	1.02	1.21	1.39	1.46	1.49	1.41	1.28	1.17	0.98



线性回归中的  $\mathbf{w}$  也会有预测的分布  
在线性回归中， $\mathbf{w}$  的预测分布是高斯分布  
对于不同回归， $\mathbf{w}$  的分布意义不同。

一次项和常数项都在括号里

$\lambda$  越大， $\hat{\mathbf{w}}$  的每个分量都趋向于稳定，有效避免了多重共线性的影响。

$\lambda$  越大，模型越简单，压缩更厉害。

不可控制过拟合 / 欠拟合

(d)(e) 随入增大不能稳定 (或入要很大)

(f) 随入增大趋于稳定

# LASSO Regression

- Proposed by R. Tibshirani, short for "Least Absolute Shrinkage and Selection Operator"
- Can be used to estimate the coefficients and select the important variables simultaneously
- Reduce the model complexity, avoid overfitting, and improve the generalization ability
- Also improve the model interpretability
- The optimization problem

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} E(\mathbf{w})$$

$$E(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1$$

误差 模型复杂度

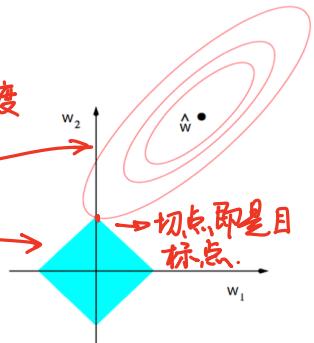
- Equivalent to the constraint minimization problem :

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2,$$

$$\text{subject to } \|\mathbf{w}\|_1 \leq \mu,$$

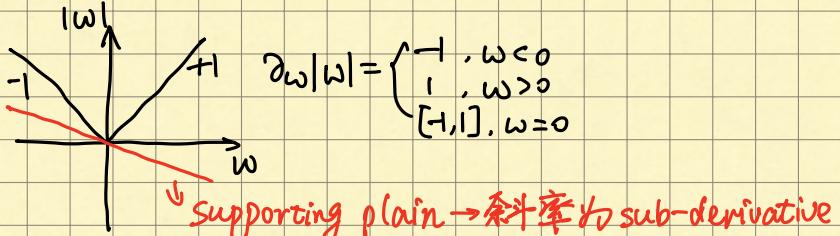
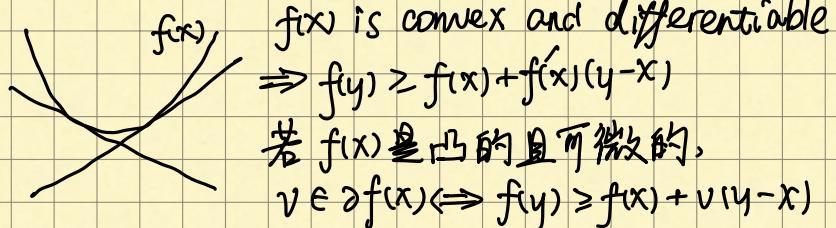
- The large  $\lambda$  corresponds to the small  $\mu$ .

- The optimal solution is sparse with  $\hat{w}_2 = 0$



$$E(\vec{w}) = \|\vec{y} - \vec{X}\vec{w}\|_2^2 + \lambda \|\vec{w}\|_1$$

$$\begin{aligned} \partial_{\vec{w}} E(\vec{w}) &= 2\vec{X}^T(\vec{y} - \vec{X}\vec{w}) + \lambda \nabla_{\vec{w}} \|\vec{w}\|_1, \\ &= 2(\vec{X}^T\vec{y} - \vec{X}^T\vec{X}\vec{w}) + \lambda \nabla_{\vec{w}} \|\vec{w}\|_1, \\ &= 2(\hat{\vec{w}}^{OLS} - \vec{w}) + \lambda \partial_{\vec{w}} \|\vec{w}\|_1 \\ &= \begin{cases} 2(\hat{\vec{w}} - \vec{w}) + \lambda, & \vec{w} > 0 \Rightarrow \vec{w} = \hat{\vec{w}} + \frac{1}{2}\lambda \\ 2(\hat{\vec{w}} - \vec{w}) - \lambda, & \vec{w} < 0 \Rightarrow \vec{w} = \hat{\vec{w}} - \frac{1}{2}\lambda \\ 2\hat{\vec{w}} + [-\lambda, \lambda], & \vec{w} = 0 \Rightarrow \hat{\vec{w}} \in \frac{1}{2}[-\lambda, \lambda] \end{cases} \end{aligned}$$

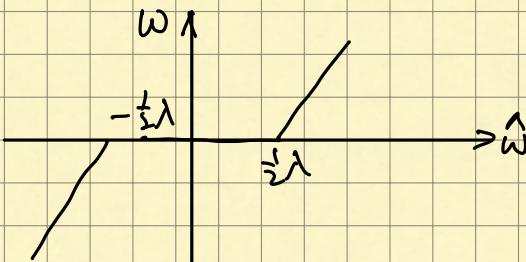


数据量不够、维数很大时也能算出

LASSO 回归没有解析解.

在  $\|\mathbf{w}\|_1$  有界下最小化  $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$

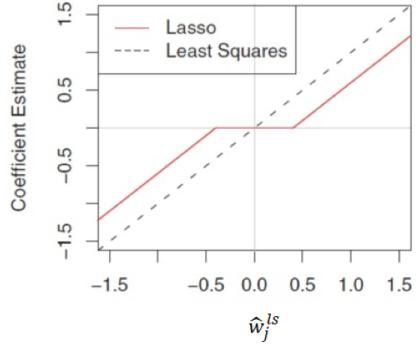
LASSO 能使交点落在坐标轴上，从而降低更大的回归系数入有更小的惩罚项  $\mu$ .



$f(x)$  is differentiable :  $\min f(x) \Rightarrow f'(x) = 0$   
 $f(x)$  is convex :  $\min f(x) \Rightarrow \partial f(x) = 0$ .

# Shrinkage and Selection Property of LASSO

- Assume  $\mathbf{X}^T \mathbf{X} = \mathbf{I}_{p+1}$ , then  $\hat{\mathbf{w}}^{OLS} = \mathbf{X}^T \mathbf{y}$
  - $\partial_{\mathbf{w}} E(\mathbf{w}) = \mathbf{w} - \mathbf{X}^T \mathbf{y} + \lambda (\partial|w_0| \times \dots \times \partial|w_p|)$
  - $\mathbf{0} \in \partial_{\mathbf{w}} E(\hat{\mathbf{w}}^{LASSO})$  implies  $0 \in \hat{w}_i^{LASSO} - \hat{w}_i^{OLS} + \lambda \partial |\hat{w}_i^{LASSO}|$
  - If  $\hat{w}_i^{LASSO} > 0$ ,  $\partial |\hat{w}_i^{LASSO}| = \{1\}$ , and  $\hat{w}_i^{LASSO} = \hat{w}_i^{OLS} - \lambda$  with  $\hat{w}_i^{OLS} > \lambda$
  - If  $\hat{w}_i^{LASSO} < 0$ ,  $\partial |\hat{w}_i^{LASSO}| = \{-1\}$ , and  $\hat{w}_i^{LASSO} = \hat{w}_i^{OLS} + \lambda$  with  $\hat{w}_i^{OLS} < -\lambda$
  - If  $\hat{w}_i^{LASSO} = 0$ ,  $\partial |\hat{w}_i^{LASSO}| = [-1, 1]$ , and  $\hat{w}_i^{OLS} \in [-\lambda, \lambda]$
  - In summary,  $\hat{w}_i^{LASSO} = (|\hat{w}_i^{OLS}| - \lambda)_+ \text{sign}(\hat{w}_i^{OLS})$
- $\hat{w}_i^{LASSO} = (|\hat{w}_i^{OLS}| - \lambda)_+ \text{sign}(\hat{w}_i^{OLS})$  is called soft thresholding of  $\hat{w}_i^{OLS}$ , where  $(a)_+ = \max(a, 0)$  is the positive part of  $a$



## Maximum A Posteriori (MAP) Estimation 极大后验估计

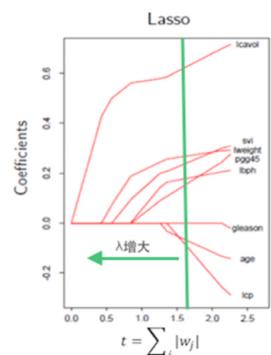
- Given  $\theta$ , the conditional distribution of  $\mathbf{y}$  is  $P(\mathbf{y}|\theta)$
- In addition, assume the parameter  $\theta$  has a prior distribution  $P(\theta)$
- The posterior distribution of  $\theta$  given the data  $\mathbf{y}$  is  $P(\theta|\mathbf{y}) \propto P(\mathbf{y}|\theta)P(\theta)$
- MAP choose the point of maximal posterior probability :

$$\hat{\theta}^{MAP} = \arg \max_{\theta} P(\theta|\mathbf{y}) = \arg \max_{\theta} (\log P(\mathbf{y}|\theta) + \log P(\theta))$$

- If  $\theta = \mathbf{w}$ , and we choose the log-prior proportional to  $\lambda \|\mathbf{w}\|_2^2$  (i.e., the normal prior  $\mathcal{N}(0, \frac{\sigma^2}{\lambda} \mathbf{I})$ ), we recover the ridge regression
- If the log-prior is proportional to  $\lambda \|\mathbf{w}\|_1$ , i.e., the prior is the tensor product of Laplace (or double exponential) distribution  $\text{Laplace}(0, \frac{2\sigma^2}{\lambda})$
- Different log-prior lead to different penalties (regularization), but this is not the case in general : some penalties may not be the logarithms of probability distributions, some other penalties depend on the data (prior is independent of the data)

## LASSO Path

- When  $\lambda$  varies, the values of the coefficients form paths (regularization paths)
- The paths are piecewise linear with the same change points, may cross the x-axis many times
- In practice, choose  $\lambda$  by cross-validation

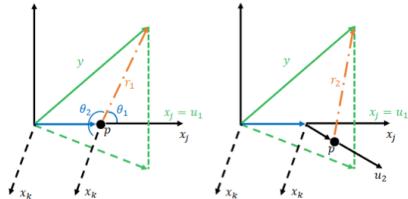


可通过交叉验证找到  $\lambda$ .

LASSO 回归 对应 Laplace 分布的极大后验估计  
岭回归 对应 高斯分布的极大后验估计

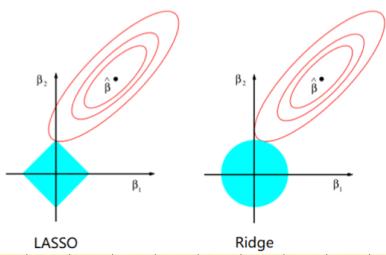
# Solving LASSO by LARS (Hastie and Efron)

1. Start with all coefficients  $w_i$  equal to zero
2. Find the predictor  $x_i$  most correlated with  $y$
3. Increase the coefficient  $w_i$  in the direction of the sign of its correlation with  $y$ . Take residuals  $r = y - \hat{y}$  along the way. Stop when some other predictor  $x_k$  has as much correlation with  $r$  as  $x_i$  has
4. Increase  $(w_i, w_k)$  in their joint least squares direction, until some other predictor  $x_m$  has as much correlation with the residual  $r$
5. Continue until all predictors are in the model



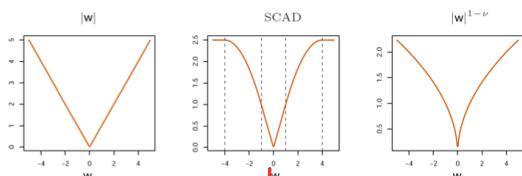
## Other Solvers

- "glmnet" by Friedman, Hastie and Tibshirani, implemented by coordinate descent, can be used in linear regression, logistic regression, etc., with LASSO ( $\ell_1$ ), ridge ( $\ell_2$ ) and elastic net ( $\ell_1 + \ell_2$ ) regularization terms
- Why LASSO seeks the sparse solution in comparison with ridge?



## Related Regulation Models

- Elastic net :  $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda_1 \|\mathbf{w}\|_2^2 + \lambda_2 \|\mathbf{w}\|_1$
- Group LASSO :  $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \sum_{g=1}^G \lambda_g \|\mathbf{w}_g\|_2$ , where  $\mathbf{w} = (\mathbf{w}_1, \dots, \mathbf{w}_G)$  is the group partition of  $\mathbf{w}$
- Dantzig Selector :  $\min_{\mathbf{w}} \|\mathbf{w}\|_1$ , subject to  $\|\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w})\|_\infty \leq \mu$
- Smoothly clipped absolute deviation (SCAD) penalty by Fan and Li (2005) : replace the penalty  $\lambda \sum_{i=0}^p |w_i|$  by  $\sum_{i=0}^p J_\alpha(w_i, \lambda)$ , where  $J_\alpha(x, \lambda)$  satisfies (for  $a \geq 2$ ) :  $\frac{dJ_\alpha}{dx} = \lambda \text{sign}(x) I(|x| \leq \lambda) + \frac{(a\lambda - |x|)_+}{(a-1)\lambda} I(|x| > \lambda)$
- Adaptive LASSO : weighted penalty  $\sum_{i=0}^p \mu_i |w_i|$  where  $\mu_i = \frac{1}{|\hat{w}_i^{OLS}|^\nu}$  with  $\nu > 0$ , as an approximation to  $|w_i|^{1-\nu}$ , non-convex penalty



希望原点处  
不导数，为模拟  
LASSO.

Ridge : SVD . 先验估计

LASSO: soft threshold . Laplace

LASSO能使惩罚项降为零,而岭回归只能使系数  
趋近于零而不能降为零.

Elastic net:

$$\begin{aligned} & \| \vec{y} - \vec{X} \vec{w} \|_2^2 + \lambda_1 \| \vec{w} \|_2^2 + \lambda_2 \| \vec{w} \|_1 \\ &= \| \vec{y} \|_2^2 - 2 \vec{y}^T \vec{X} \vec{w} + \vec{w}^T \vec{X}^T \vec{X} \vec{w} + \lambda_2 \vec{w}^T \vec{w} + \lambda_1 \| \vec{w} \|_1 \end{aligned}$$

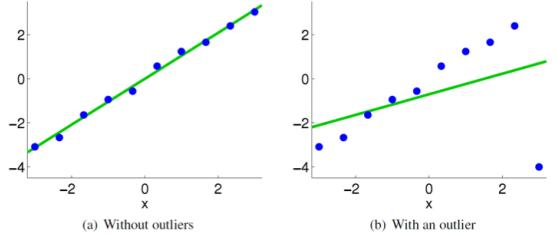
$$= \vec{w}^T (\vec{X}^T \vec{X} + \lambda_2 I)$$

## Problem with Outliers — Robustness

- $l_2$ -loss minimization is sensitive to outliers (non-robust)
- It penalizes greatly for the large residual, probably at outliers
- Consider  $l_2$ -regression towards a constant model  $f_\theta(x) = \theta$ :

$$\hat{\theta}_{LS} = \arg \min_{\theta} \sum_{i=1}^n (y_i - \theta)^2 = \frac{1}{n} \sum_{i=1}^n y_i = \text{Mean}(\{y_i\})$$

- If  $l_1$ -loss is used :  $\hat{\theta}_{LS} = \arg \min_{\theta} \sum_{i=1}^n |y_i - \theta| = \text{Median}(\{y_i\})$



$l_2$ -loss 对离群值敏感

常数回帰

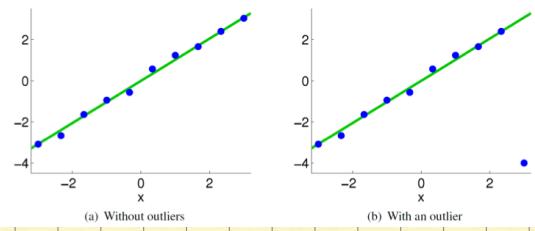
## Least Absolute Deviations Regression

- With  $l_1$ -loss :  $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_1$
- Recast to linear programming (Boyd & Vandenberghe, 2004)

$$\begin{aligned} \min_{\mathbf{w}} \sum_{i=1}^n r_i \\ \text{subject to } -r_i \leq y_i - \mathbf{w}^T \mathbf{x}_i \leq r_i, \quad i = 1, \dots, n \end{aligned}$$

← 线性规划  
等价替换

$l_2$ -loss 不鲁棒,  $l_1$ -loss 无解析解

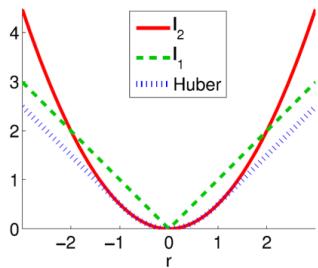


## Huber Loss Minimization

- Combine  $l_2$ -loss with  $l_1$ -loss :  $L(y, f) = \rho_{\text{Huber}}(y - f)$ , where

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

- Huber loss minimization :  $\min_{\mathbf{w}} J(\mathbf{w}) \triangleq \sum_{i=1}^n \rho_{\text{Huber}}(y_i - \mathbf{w}^T \mathbf{x}_i)$



局部使用  $l_2$ -loss, 使得整体解光滑.

## Gradient descent

$$x_{n+1} = x_n - \alpha \nabla f(x_n)$$

step size

若  $f(x)$  是强凸函数能走到最小值

# Optimizing Huberized Regression

- Gradient descent :

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

求导链式法则使负变正。

then updating  $\mathbf{w} \leftarrow \mathbf{w} + \tau \sum_{i=1}^n \mathbf{x}_i \rho'_{\text{Huber}}(y_i - \mathbf{w}^T \mathbf{x}_i)$  with  $\tau$  being learning rate (or step size)

- Stochastic gradient descent : replace the summation by randomly selecting one sample  $(\mathbf{x}_i, y_i)$  at each step
- Can be used in combination with  $l_1$  regularization (like Lasso)

每次更新梯度随机选取一个样本来算梯度代替原来的加和梯度。在期望上效果同最优一样。

## Iteratively Reweighted Least Square (IRLS)

- Quadratically upper bound the absolute-value part of Huber loss :

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

this upper bound touches Huber loss at  $r = \pm c$  (i.e., " $=$ " holds)

- Take  $c_i^{(k)} = |y_i - (\mathbf{w}^{(k)})^T \mathbf{x}_i|$  to be the absolute residual for the  $i$ -th training sample at  $k$ -th step iteration, instead of directly minimize Huber loss, we can minimize its  $k$ -th upper bound :

$$\mathbf{w}^{(k+1)} = \arg \min_{\mathbf{w}} \tilde{J}(\mathbf{w}) \triangleq \frac{1}{2} \sum_{i=1}^n \alpha_i^{(k)} (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + C^{(k)},$$

$$\text{where } C^{(k)} = \sum_{c_i^{(k)} > \eta} \left( \frac{1}{2} \eta c_i^{(k)} - \frac{1}{2} \eta^2 \right) \text{ and } \alpha_i^{(k)} = \begin{cases} 1, & c_i^{(k)} \leq \eta \\ \eta/c_i^{(k)}, & c_i^{(k)} > \eta \end{cases}$$

- Easy to see that :  $J(\mathbf{w}^{(k)}) = \tilde{J}(\mathbf{w}^{(k)}) \geq \tilde{J}(\mathbf{w}^{(k+1)}) \geq J(\mathbf{w}^{(k+1)})$

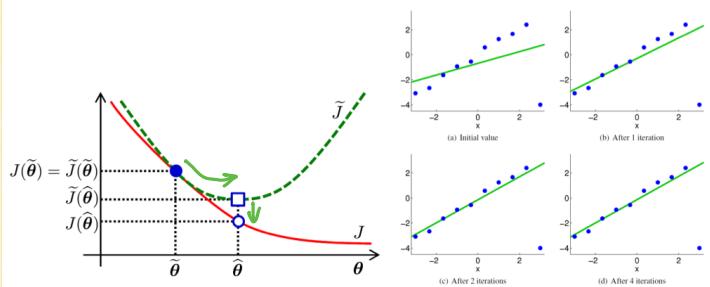
- This is also called majorization-minimization (MM) algorithm
- Continue the iteration until convergence (guaranteed) to global minimizer (since Huber loss is strictly convex)
- In each step, there is an analytical formula for the weighted least square minimization

效率不一定比 gradient decent 高。

在局部用二次函数限制，该二次函数永远在 Huber loss 上面且相切。

通过调整参数来控制残差影响。

Huber Loss 是严格凸的，能得到全局最小值。  
每一步都得到一个加权最小二乘问题。



## Tukey Loss Minimization

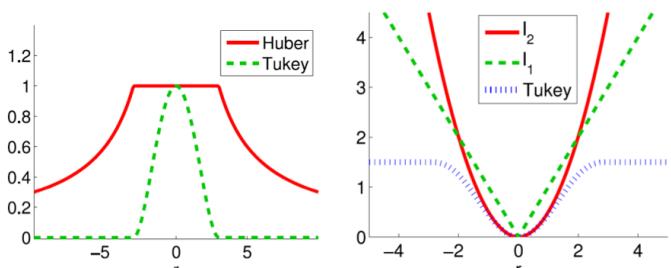
- Further reduce the weights for large residual (left plot)
- Tukey loss (right plot) :

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

- Can also be solved using (IRLS), but not guaranteed to converge to global minimizer (due to the nonconvexity of Tukey loss)

在残差不太大时保留，太大时直接归零。

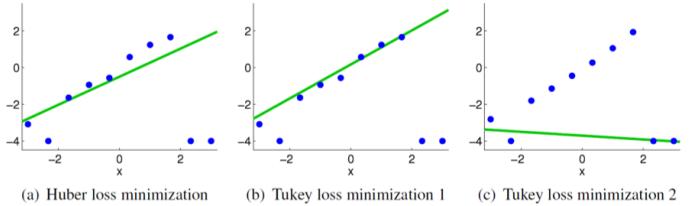
可用 IRLS 解，但不一定收敛于全局最小值。  
因为 Tukey Loss 是非凸的。



# Other Robust Regressions

- In Tukey regression, slightly changing the noise included in training output samples gives another local optimal solution
- Other candidate losses : pinball loss, deadzone-linear loss, Chebyshev (minimax) approximation, Conditional Value-At-Risk (CVaR) measure, etc.

→ 实际上就是 SVR



## Model Assessment

### Errors and $R^2$

- Mean absolute error (MAE) :  $MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$
- Mean square error (MSE) :  $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$
- Root mean square error (RMSE) :  
 $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$
- Coefficient of Determination  $R^2$  :  $R^2 := 1 - \frac{SS_{res}}{SS_{tot}}$ , where  
 $SS_{tot} = \sum_{i=1}^n (y_i - \bar{y})^2$  is the total sum of squares, and  
 $SS_{res} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$  is the residual sum of squares;  
 $R^2 \in [0, 1]$  (might be negative); the larger the  $R^2$ , the smaller the ratio of  $SS_{res}$  to  $SS_{tot}$ , thus the better the model

平均绝对误差 :  $\frac{1}{n} \sum |y_i - \hat{y}_i|$   
 均方误差

均根误差

$$\because S_{tot} = S_{reg} + S_{res} \therefore R^2 = 1 - \frac{S_{res}}{S_{tot}} = \frac{S_{reg}}{S_{tot}}$$

$R^2$  越大，模型越好。 $R^2$  为负数时模型效果很差（说明甚至没随机取效果好）

### Adjusted Coefficient of Determination

- Adjusted coefficient of determination :  $R_{adj}^2 = 1 - \frac{(1-R^2)(n-1)}{n-p-1}$
- $n$  is the number of samples,  $p$  is the dimensionality (or the number of attributes)
- The larger the  $R_{adj}^2$  value, the better performance the model
- When adding important variables into the model,  $R_{adj}^2$  gets larger and  $SS_{res}$  is reduced
- When adding unimportant variables into the model,  $R_{adj}^2$  may get smaller and  $SS_{res}$  may increase
- In fact, one can show that  $1 - R_{adj}^2 = \frac{\hat{\sigma}^2}{S^2}$ , where  
 $\hat{\sigma}^2 = \frac{1}{n-p-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2$  and  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$  with  
 $(n-p-1) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p-1}^2$  and  $(n-1) \frac{S^2}{\sigma^2} \sim \chi_{n-1}^2$  if  $w = 0$ .

高维用调整后的比较好。