

# Linear Regression-3

Prof. Asim Tewari  
IIT Bombay

# Multiple Linear Regression

Multiple Linear Regression assumes

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p + \epsilon$$

The model can be expressed as

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \cdots + \hat{\beta}_p x_p$$

with its coefficients being derived by minimizing  
RSS

$$\begin{aligned} \text{RSS} &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \hat{\beta}_2 x_{i2} - \cdots - \hat{\beta}_p x_{ip})^2 \end{aligned}$$

*Handwritten red text:*  $\sum_{i=1}^n |y_i - \hat{y}_i|$

# Multiple Linear Regression

If  $X$  is a vector  $\begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^p \end{bmatrix}$

Data :  $n$ -data points

$(x_1^T, y_1), (x_2^T, y_2) \dots (x_n^T, y_n)$

$\begin{bmatrix} x_{1,1} \\ x_{1,2} \\ x_{1,3} \\ \vdots \\ x_{1,p} \end{bmatrix}$

$\begin{bmatrix} x_{n,1} \\ x_{n,2} \\ \vdots \\ x_{n,p} \end{bmatrix}$

# Multiple Linear Regression

Data:  $(x_1, y_1), \dots, (x_i, y_i), \dots, (x_n, y_n)$

$$x_i = (x_i^1, x_i^2, \dots, x_i^p)^T$$

$p \times 1$

$$\beta = (\beta_0, \beta_1, \beta_2, \dots, \beta_p)^T$$

$(p+1) \times 1$

$$x_i = (1, x_i^1, x_i^2, x_i^3, \dots, x_i^p)$$

$(p+1) \times 1$

# Multiple Linear Regression

$$\begin{array}{c}
 \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}_{n \times (p+1)} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}_{(p+1) \times 1} \\
 \uparrow \quad \quad \quad \uparrow \\
 n \times (p+1) \quad (p+1) \times 1 \\
 \underbrace{\hspace{10em}}_{n \times 1} \\
 \begin{bmatrix} 1 \\ x_{1,1} \\ x_{1,2} \\ \vdots \\ x_{1,p} \end{bmatrix}
 \end{array}
 \quad
 \left\{
 \begin{array}{l}
 Y = X \beta \\
 \begin{array}{ccc}
 \uparrow & \uparrow & \uparrow \\
 n \times 1 & n \times (p+1) & (p+1) \times 1
 \end{array}
 \end{array}
 \right.
 \quad
 \left.
 \begin{array}{l}
 p \rightarrow \# \text{ of features} \\
 n = \# \text{ of samples}
 \end{array}
 \right.$$

# Multiple Linear Regression

$$RSS = \left\| \hat{Y} - Y \right\|^2$$

$\beta^*$  is the value of  $\beta$  which min RSS  
 $\therefore (n > p+1)$

$$RSS(\beta)$$

Local min  $\beta^*$  would be sol of

$$\nabla RSS(\beta^*) = 0$$

# Multiple Linear Regression

$$\begin{aligned} \text{RSS}(\beta) &= \|X\beta - y\|^2 \\ &= (X\beta - y)^T (X\beta - y) \\ &= (X\beta)^T X\beta - (X\beta)^T y - y^T X\beta + y^T y \\ &= \beta^T X^T X \beta - 2 \beta^T X^T y + y^T y \end{aligned}$$

$$\boxed{\nabla_x (a^T x) = a \quad \text{and} \quad \nabla_x (x^T A x) = (A + A^T)x}$$

# Multiple Linear Regression

$$\nabla_x (a^T x) = a$$

$$RSS = \beta^T \underbrace{x^T x}_{a} \beta - 2 \underbrace{\beta^T x^T y}_a + y^T y \quad \nabla_x (x^T A x) = (A + A^T)x$$

$$\nabla_{\beta} RSS(\beta) = \left( \underset{\substack{\uparrow \\ (p+1) \times 1}}{x^T x + x^T x} \right) \beta - 2 x^T y$$

$$\nabla_{\beta} RSS(\beta) = 2 x^T x \beta - 2 x^T y$$

Find  $\nabla_{\beta} RSS(\beta^*) = 0$



# Multiple Linear Regression

$$\nabla_{\beta} \text{RSS}(\beta) = 2X^T X \beta - 2X^T y$$

Find  $\nabla_{\beta} \text{RSS}(\beta^*) = 0$

$$2X^T X \beta^* - 2X^T y = 0$$

$$X^T X \beta^* = X^T y$$

$$\underbrace{(X^T X)^{-1}}_I (X^T X) \beta^* = (X^T X)^{-1} X^T y$$

$$\beta^* = (X^T X)^{-1} X^T y$$

# Multiple Linear Regression

$$\min \text{RSS}(\beta)$$

$$\Rightarrow \beta^* = (X^T X)^{-1} X^T y$$

$\hookrightarrow \min$  (global min?)

$$y = X\beta$$

$$\nabla_{\beta} \text{RSS}(\beta) = 2 X^T X \beta - 2 X^T y$$

$$\nabla^2 \text{RSS}(\beta) = 2 \underbrace{X^T X}_{\substack{(P+1) \times n \\ (P+1) (P+1)}} \quad \begin{matrix} \uparrow \\ n \times (P+1) \end{matrix}$$

# Multiple Linear Regression

$$\underbrace{\nabla^2 \text{RSS}(\beta)}_{\text{Hessian of RSS}} = 2 \underbrace{X^T X}_{\substack{(P+1) \times n \\ (P+1) \quad (P+1)}} \quad \begin{matrix} \uparrow \\ n \times (P+1) \end{matrix}$$

$$\forall \beta \quad \beta^T (2 X^T X) \beta = 2 (X \beta)^T X \beta = 2 \|X \beta\|^2 \geq 0$$

$\beta^*$  is now global minima.

# Multiple Linear Regression

$$\beta^* = (X^T X)^{-1} X^T y$$

$$X = \begin{bmatrix} 1 & x_1^1 & x_1^2 & \dots & x_1^p \\ 1 & x_2^1 & x_2^2 & \dots & x_2^p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n^1 & x_n^2 & \dots & x_n^p \end{bmatrix}$$

$$n \times (p+1)$$

$$\begin{matrix} (p+1) \times (p+1) \\ \left( \begin{matrix} X^T & X \end{matrix} \right)^{-1} \\ \begin{matrix} \uparrow & \uparrow \\ (p+1) \times n & n \times (p+1) \end{matrix} \end{matrix}$$

}

# Gauss-Markov Theorem

- The Gauss–Markov theorem states that if we have any other linear estimator  $\tilde{\theta} = \mathbf{c}^T \mathbf{y}$  that is unbiased for  $\boldsymbol{\alpha}^T \boldsymbol{\beta}$ , that is,  $E(\mathbf{c}^T \mathbf{y}) = \boldsymbol{\alpha}^T \boldsymbol{\beta}$ , then

$$\text{Var}(\mathbf{a}^T \hat{\boldsymbol{\beta}}) \leq \text{Var}(\mathbf{c}^T \mathbf{y}).$$

# Gauss-Markov Theorem

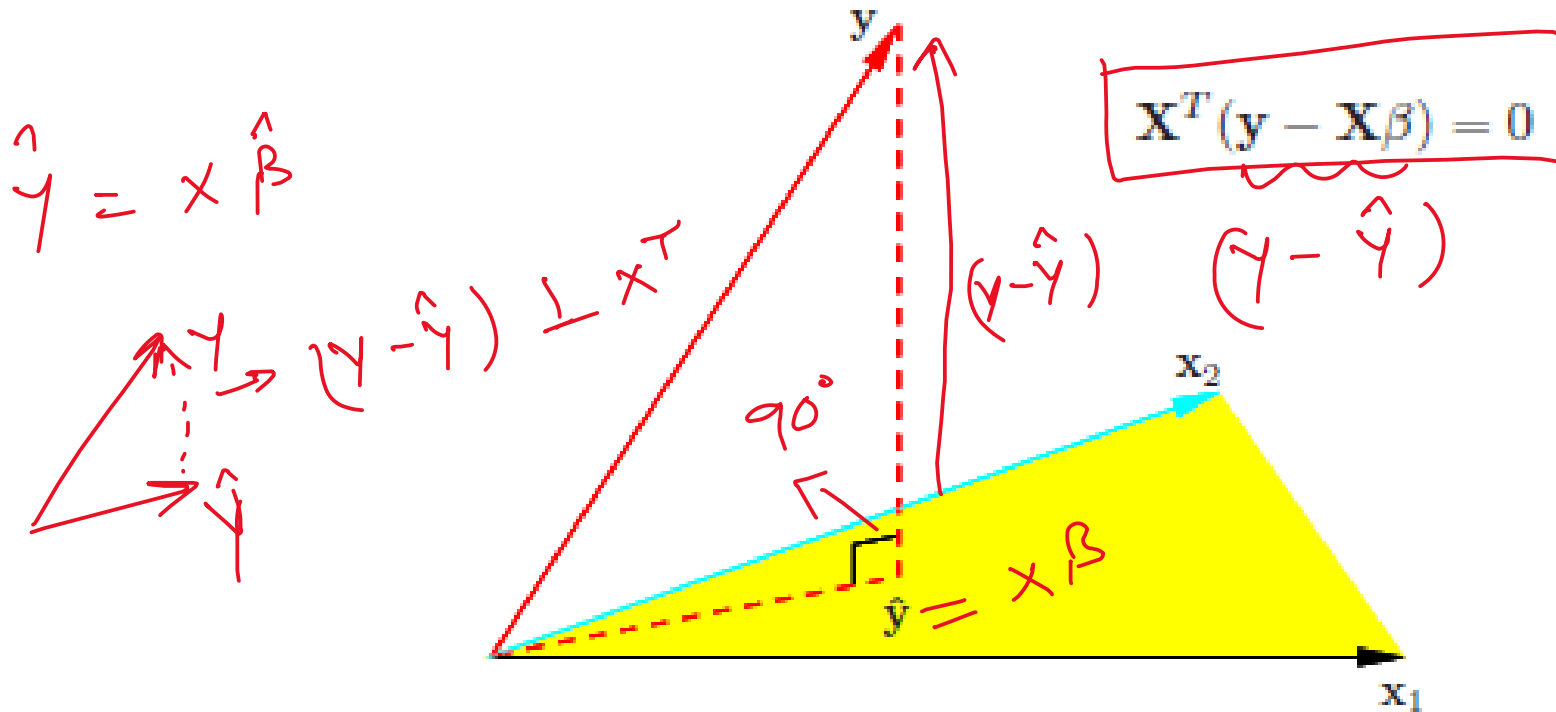
- The least squares estimate of  $\alpha^T \beta$  is

$$\hat{\theta} = a^T \hat{\beta} = a^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

- Considering  $\mathbf{X}$  to be fixed, this is a linear function  $c_0^T \mathbf{y}$  of the response vector  $\mathbf{y}$ .
- If we assume that the linear model is correct,  $\alpha^T \beta$  is unbiased since

$$\begin{aligned} \mathbb{E}(a^T \hat{\beta}) &= \mathbb{E}(a^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}) \\ &= a^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta \\ &= a^T \beta. \end{aligned}$$

# Linear Regression



The N-dimensional geometry of least squares regression with two predictors. The outcome vector  $y$  is orthogonally projected onto the hyperplane spanned by the input vectors  $x_1$  and  $x_2$ . The projection  $\hat{y}$  represents the vector of the least squares predictions

# Multiple Linear Regression

In multiple linear regression, we usually are interested in answering a few important questions.

- 1. Is at least one of the predictors  $X_1, X_2, \dots, X_p$  useful in predicting the response?*
- 2. Do all the predictors help to explain  $Y$ , or is only a subset of the predictors useful?*
- 3. How well does the model fit the data?*
- 4. Given a set of predictor values, what response value should we predict, and how accurate is our prediction?*



# Extensions of Linear Models

- Removing the Additive Assumption

Introduce the interactive term

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon.$$

$$\begin{aligned} Y &= \beta_0 + (\beta_1 + \beta_3 X_2) X_1 + \beta_2 X_2 + \epsilon \\ &= \beta_0 + \tilde{\beta}_1 X_1 + \beta_2 X_2 + \epsilon \end{aligned}$$

- Where  $\tilde{\beta}_1 = \beta_1 + \beta_3 X_2$ .

- Non-linear Relationships

$$\text{mpg} = \beta_0 + \beta_1 \times \text{horsepower} + \beta_2 \times \text{horsepower}^2 + \epsilon$$

# Nonlinear regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots \dots$$

$$y = \beta_0 + \beta_1 x + \beta_2 x^2$$

$$\equiv y = \beta_0 + \beta_1 x + \beta_2 x^2$$

$$V \propto S$$

# Basis function regression

$$y = \underline{\beta_0} + \underline{\beta_1} f_1(x) + \underline{\beta_2} f_2(x) + \dots + \epsilon$$

eg. (1)

$x$

$x^2$

(2)

$$f_1(x) = x^{1.5}, \quad f_2(x) = \left(1 + \frac{x}{3+x^2}\right)$$

# Hypothesis testing in multi linear regression

$$y = \overbrace{f(x)} + \epsilon \rightarrow f(x) ?$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \epsilon$$

t statistic  $\Rightarrow$  p-value

Intercept	$\beta_0$	p-value
$x_1$	$\beta_1$	
$x_2$	$\beta_2$	
	$\vdots$	

# F-Statistic

$$y = \beta_0 + \beta_1 \overset{\checkmark x}{X_1} + \beta_2 \overset{\checkmark x}{X_2} + \dots + \beta_{100} \overset{\checkmark x}{X_{100}}$$

for 95% C  $\Rightarrow$   $\sum X_i$  ✓

t-statistic  $\Rightarrow$  p-value for each variable  $x_i$

$$F\text{-statistic} = \frac{(TSS - RSS) / p}{RSS / (n - p - 1)}$$

$\left\{ \begin{array}{l} \text{If } F \approx 1 \quad H_0 \quad \checkmark \\ F \gg 1 \quad H_a \quad \checkmark \end{array} \right. \rightarrow n, p$

# F-Statistic

$$y = \beta_0 + \overset{\checkmark x}{\beta_1 X_1} + \overset{\checkmark x}{\beta_2 X_2} + \dots + \overset{\checkmark x}{\beta_{100} X_{100}} + \epsilon$$

for 95% C  $\Rightarrow$   $\sum X_i$  ✓

If  $H_0$  is true and  $\epsilon$  is  $N(0, \sigma)$

then F-statistic follows a F-dist.

$\left\{ \begin{array}{ll} \text{If } F \approx 1 & H_0 \quad \checkmark \\ F \gg 1 & H_a \quad \checkmark \end{array} \right. \rightarrow \eta, p$

# Is There a Relationship Between the Response and Predictors?

We test the null hypothesis,

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$$

*Handwritten notes:* }  $\rightarrow$  If this is true  
 $(TSS - RSS)/p \rightarrow \sigma^2$

versus the alternative

$$H_a : \text{at least one } \beta_j \text{ is non-zero.}$$

This hypothesis test is performed by computing the F-statistic,

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)}$$

*Handwritten notes:*  
If  $H_0$  is true,  $F \hat{=} 1$   
If  $H_0$  is not true,  $F > 1$   
Linear model is correct  $\sigma^2$

Value of F-statistic close to 1 when null hypothesis true  $\leftarrow$

Value of F-statistic greater than 1 when alternative hypothesis true

# Hypothesis testing in multi linear regression

- F is very close to **one** we cannot reject the null hypothesis (thus, in a sense we accepted the **null hypothesis**)
- If F is **very large** we reject the null hypothesis (thus, in a sense we accepted the **alternate hypothesis**)

## How large is large enough?

- This depends upon the values of  $n$  and  $p$ .
- If  $n$  is very large a small value above 1 is also a compelling evidence against the null hypothesis; however if  $n$  is a small then  $F$  has to be very large for us to reject the null hypothesis.
- When the null hypothesis is true and the error follows a Gaussian distribution, then it can be shown that F-statistic follows F-distribution



# Hypothesis testing in multi linear regression

Why do we need F-statistic when t-statistic already exists?

(Given these individual p-values for each variable, why do we need to look at the overall F-statistic? After all, it seems likely that if any one of the p-values for the individual variables is very small, then at least one of the predictors is related to the response.)

- However, the above logic is flawed, especially when the number of predictors  $p$  is large.
- For instance, consider an example in which  $p = 100$  and  $H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$  is true, so no variable is truly associated with the response. In this situation, about 5% of the p-values associated with each variable will be below 0.05 by chance. In other words, we expect to see approximately five small p-values even in the absence of any true association between the predictors and the response.

# Multiple Linear Regression

- Predictors with Only Two Levels
- Define new variables

$$x_i = \begin{cases} 1 & \text{if } i\text{th person is female} \\ 0 & \text{if } i\text{th person is male,} \end{cases}$$

The model then takes the form

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i\text{th person is female} \\ \beta_0 + \epsilon_i & \text{if } i\text{th person is male.} \end{cases}$$

# Multiple Linear Regression

- Predictors with more than Two Levels
- Define new variables

$$x_{i1} = \begin{cases} 1 & \text{if } i\text{th person is Asian} \\ 0 & \text{if } i\text{th person is not Asian,} \end{cases}$$

$$x_{i2} = \begin{cases} 1 & \text{if } i\text{th person is Caucasian} \\ 0 & \text{if } i\text{th person is not Caucasian.} \end{cases}$$

The model then takes the form

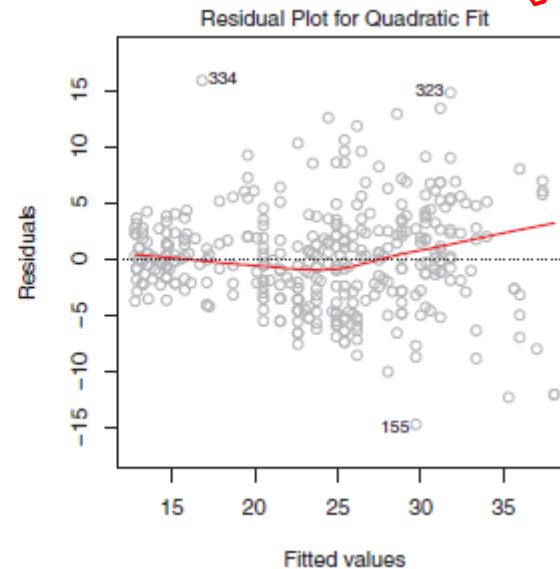
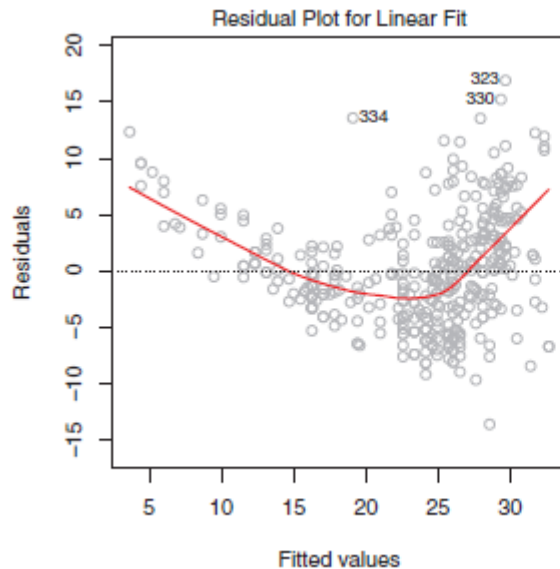
$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i = \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i\text{th person is Asian} \\ \beta_0 + \beta_2 + \epsilon_i & \text{if } i\text{th person is Caucasian} \\ \beta_0 + \epsilon_i & \text{if } i\text{th person is African American.} \end{cases}$$

# Potential problems of Linear Regression

1. *Non-linearity of the response-predictor relationships.*
2. *Correlation of error terms.*
3. *Non-constant variance of error terms.*
4. *Outliers.*
5. *High-leverage points.*
6. *Collinearity.*

# Potential problems of Linear Regression

## 1. Non-linearity of the response-predictor relationships



$$y = A_0 + A_1 x + A_2 x^2 + \epsilon$$

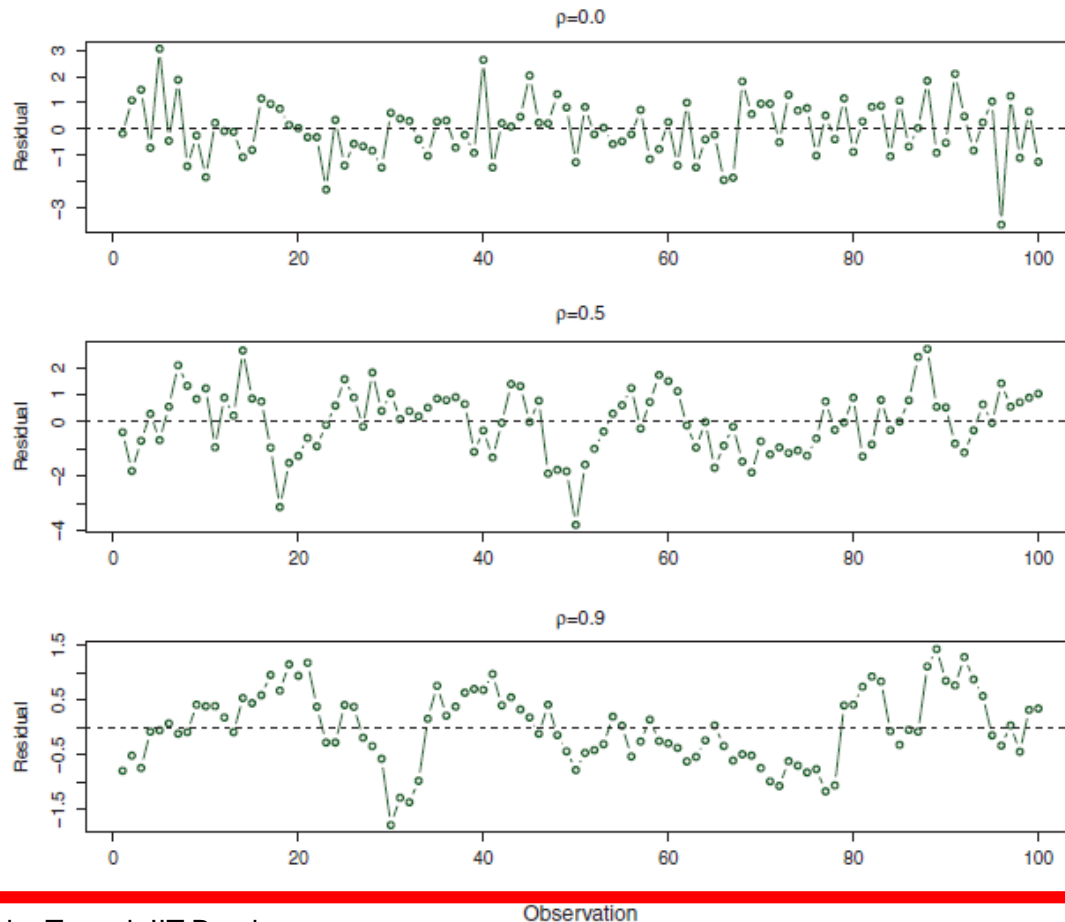
Hyper Parameters

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$$

Plots of residuals versus predicted (or fitted) values for the Auto data set. In each plot, the red line is a smooth fit to the residuals, intended to make it easier to identify a trend. Left: A linear regression of mpg on horsepower. A strong pattern in the residuals indicates non-linearity in the data. Right: A quadratic regression of mpg on horsepower and horsepower<sup>2</sup>. There is little pattern in the residuals.

# Potential problems of Linear Regression

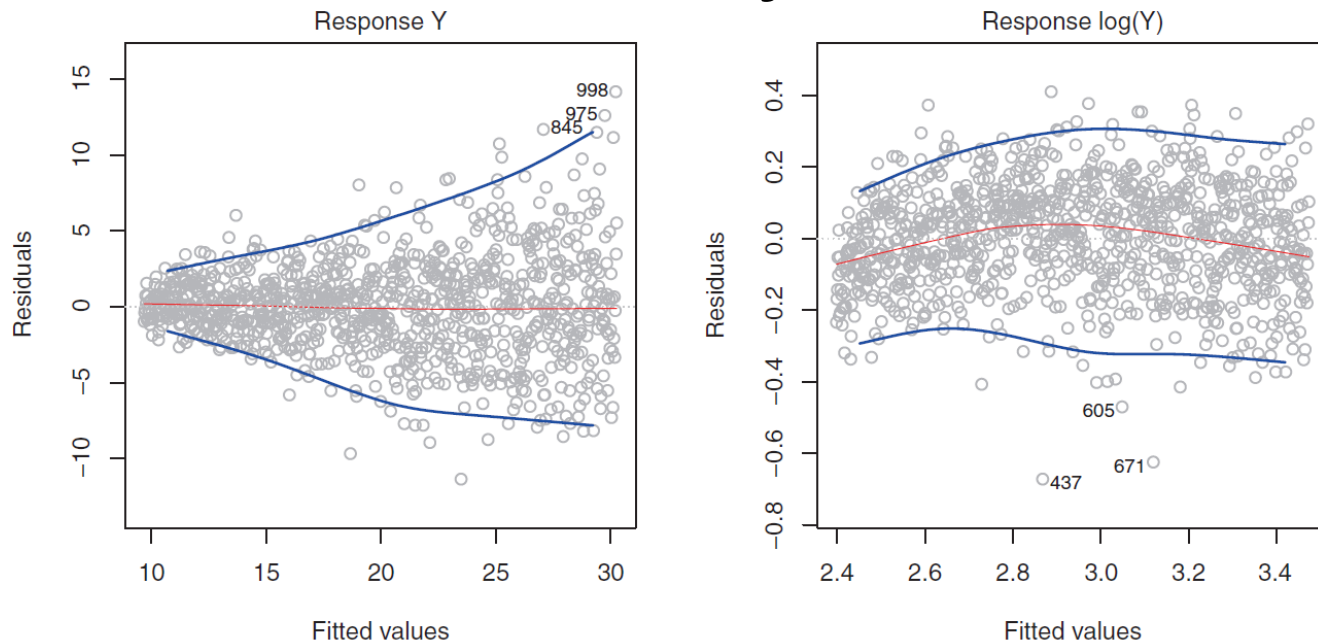
## 2. *Correlation of error terms*



*Plots of residuals from simulated time series data sets generated with differing levels of correlation  $\rho$  between error terms for adjacent time points.*

# Potential problems of Linear Regression

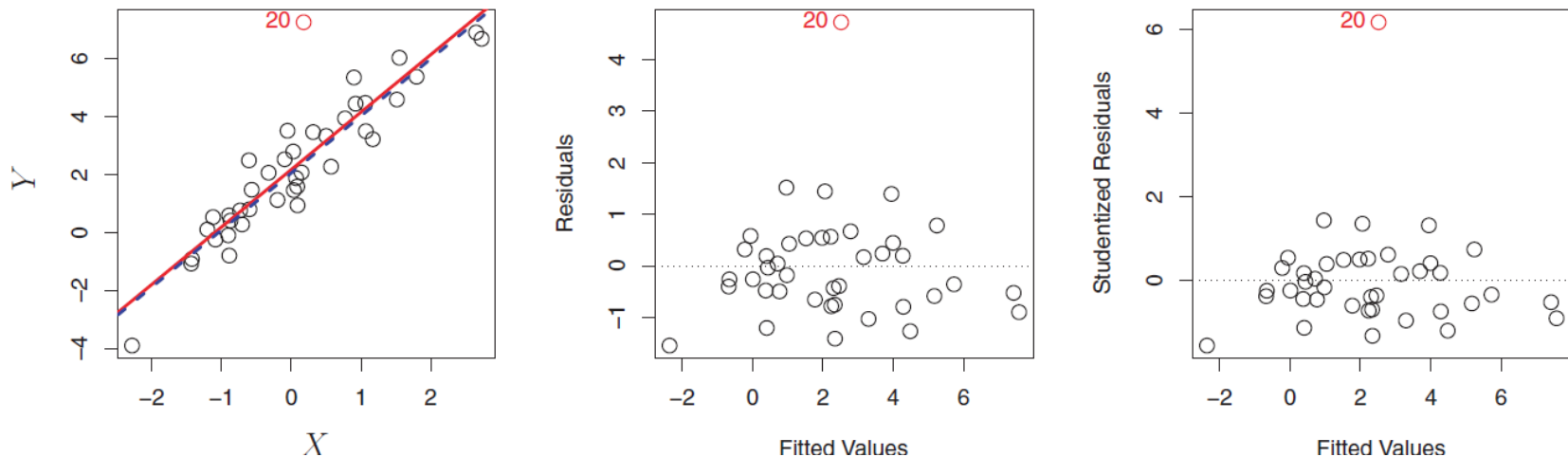
## 3. *Non-constant variance of error terms.*



*Residual plots. In each plot, the red line is a smooth fit to the residuals, intended to make it easier to identify a trend. The blue lines track the outer quantiles of the residuals, and emphasize patterns. Left: The funnel shape indicates heteroscedasticity. Right: The response has been log transformed, and there is now no evidence of heteroscedasticity*

# Potential problems of Linear Regression

## 4. *Outliers.*



*Left: The least squares regression line is shown in red, and the regression line after removing the outlier is shown in blue. Center: The residual plot clearly identifies the outlier. Right: The outlier has a studentized residual of 6; typically we expect values between  $-3$  and  $3$ .*



# Robust Regression

- The average quadratic error functional (RSS) is very sensitive to outliers
- Robust error functionals aim to reduce the influence of outliers.
- Linear regression with robust error functionals is called robust linear regression.

# Robust Regression

- One example of a robust error functional is the Huber function where the errors are only squared if they are smaller than a threshold  $\epsilon > 0$ , otherwise they have only a linear impact

$$E_H = \sum_{k=1}^n \begin{cases} e_k^2 & \text{if } |e_k| < \epsilon \\ 2\epsilon \cdot |e_k| - \epsilon^2 & \text{otherwise} \end{cases}$$

$e_k^2$

# Robust Regression

- Another example of a robust error functional is least trimmed squares which sorts the errors so that

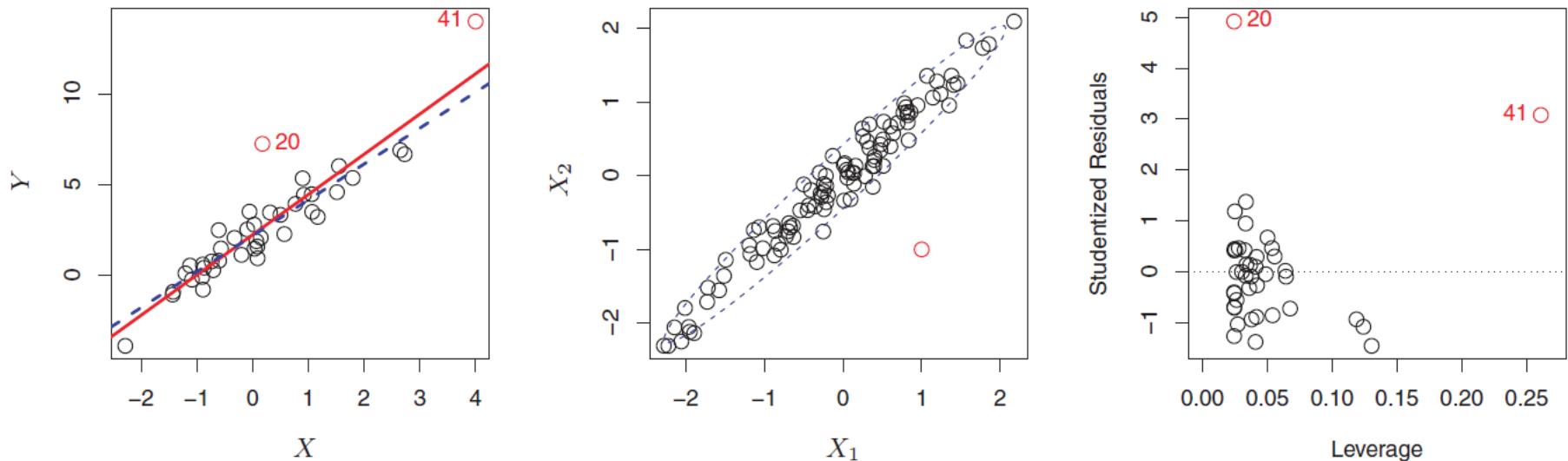
$$e'_1 \leq e'_2 \leq \dots \leq e'_n$$

and only considers the  $m$  smallest errors,  $1 \leq m \leq n$ .

$$E_{LTS} = \sum_{k=1}^m e_k'^2$$

# Potential problems of Linear Regression

## 5. High-leverage points



*Left: Observation 41 is a high leverage point, while 20 is not. The red line is the fit to all the data, and the blue line is the fit with observation 41 removed. Center: The red observation is not unusual in terms of its  $X_1$  value or its  $X_2$  value, but still falls outside the bulk of the data, and hence has high leverage. Right: Observation 41 has a high leverage and a high residual*

# *High-leverage points*

- In order to quantify an observation's leverage, we compute the *leverage statistic*.
- A large value of this statistic indicates an observation with high leverage.
- For a simple linear regression,

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i'=1}^n (x_{i'} - \bar{x})^2}.$$

- The leverage statistic  $h_i$  is always between  $1/n$  and 1.
- The average leverage for all the observations is always equal to  $(p+1)/n$ .

# Potential problems of Linear Regression

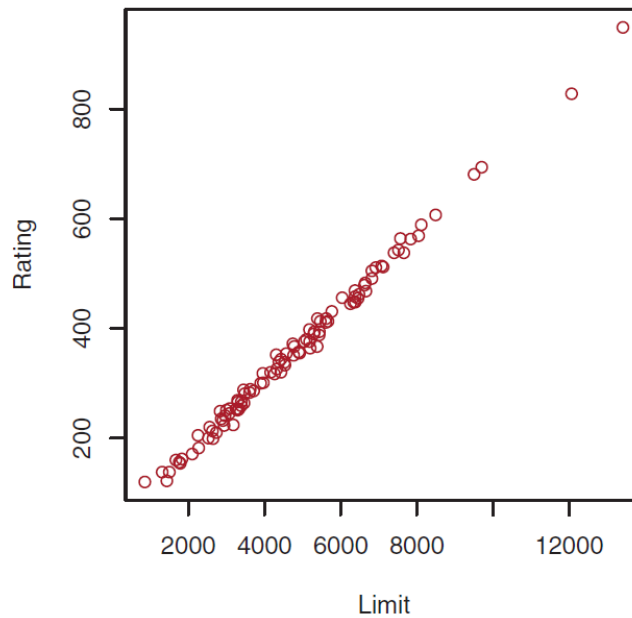
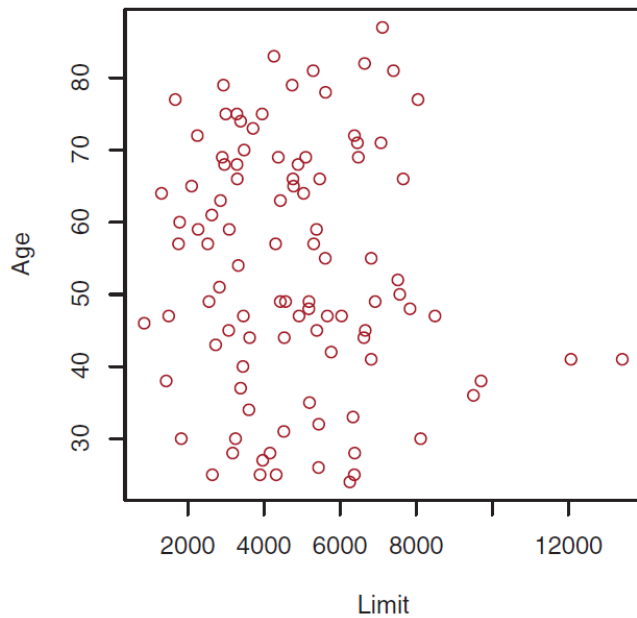
$$y = f(x) + \epsilon$$

↑

## 6. Collinearity

## Regression

$$\beta^* = (X^T X)^{-1} X^T y$$



Scatterplots of the observations from the Credit data set. Left: A plot of age versus limit. These two variables are not collinear. Right: A plot of rating versus limit. There is high collinearity.

# Singular Value Decomposition (SVD)

$$M_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}^*$$

Diagram illustrating the dimensions of the matrices in the SVD decomposition:

- $U$  is  $m \times m$
- $\Sigma$  is  $m \times n$
- $V^*$  is  $n \times n$

$$\begin{bmatrix} \Sigma_{ii} & 0 \\ 0 & 0 \end{bmatrix}$$

Both  $U$  and  $V^*$  are unitary matrices  
(Real or complex)

$\Rightarrow U^*U = I_{m \times m}$   
and  $V^*V = I_{n \times n}$

$\Sigma$  is a diagonal rectangular matrix of non-negative real numbers.

# Singular Value Decomposition (SVD)

Singular value decomposition is a generalization of the eigen-decomposition of a square Matrix to a non-square matrix.

$$M = U \Sigma V^*$$

Where,

$M$  is a real or complex  $m \times n$  matrix

$U$  is an  $m \times m$ , real or complex unitary matrix (conjugate transpose,  $U^*$  is also its inverse)

$\Sigma$  is an  $m \times n$  rectangular diagonal matrix with non-negative real numbers

$V^*$  is an  $n \times n$ , real or complex unitary matrix

If  $M$  is real then  $U$  and  $V$  are real orthogonal matrices

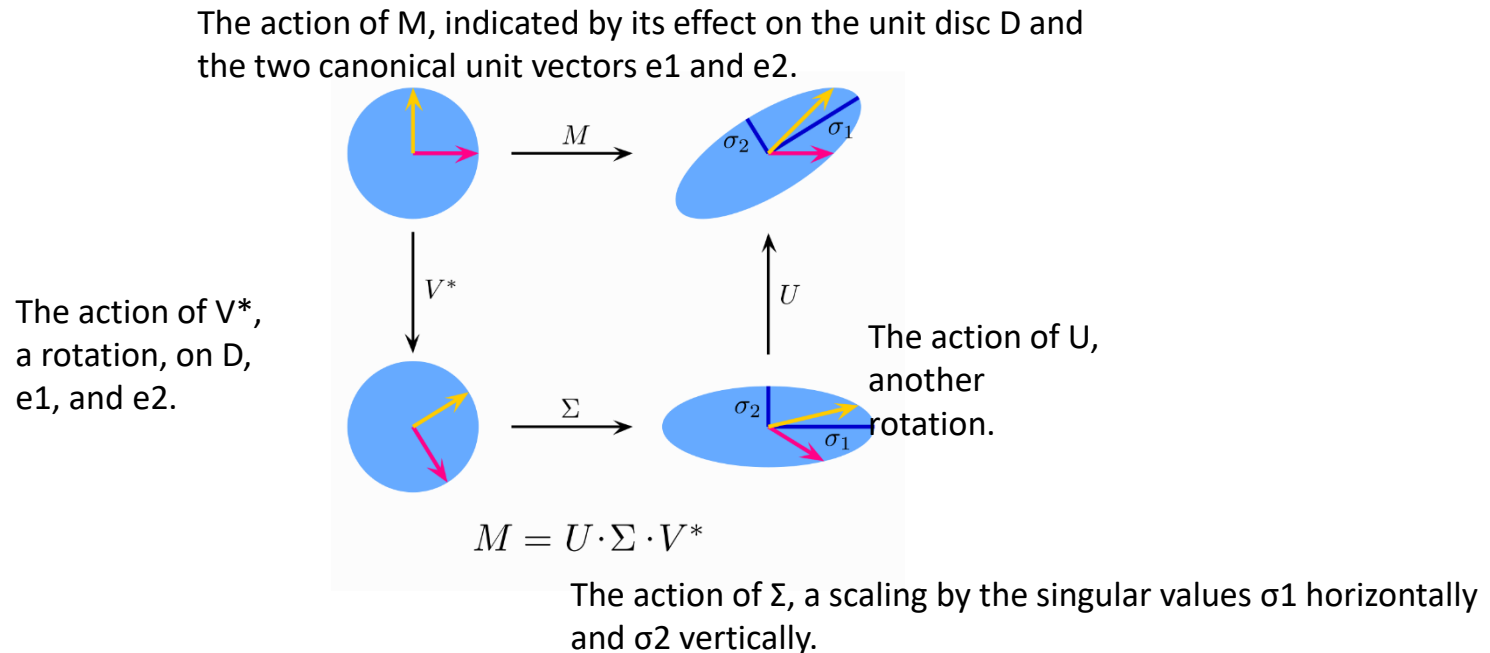
The diagonal values of  $\Sigma$  are known as the singular values. By convention they are written in descending order. In this case  $\Sigma$  (but not always  $U$  and  $V^*$ ) is uniquely determined by  $M$ .



# Singular Value Decomposition (SVD)

$$M = U\Sigma V^*$$

Illustration of the singular value decomposition  $U\Sigma V^*$  of a real  $2 \times 2$  matrix  $M$ .



# Compact Singular Value Decomposition

Compact singular value decomposition is similar to SVD with a square diagonal matrix  $\Sigma_d$  .

$$M = U_c \Sigma_d V_c^*$$

Where,

$M$  is a real or complex  $m \times n$  matrix

$U_c$  is an  $m \times r$ , semi-unitary matrix ( $U_c^* U_c = I_{r \times r}$ )

$\Sigma_d$  is an  $r \times r$  square diagonal matrix with positive real numbers

$V_c$  is an  $n \times r$ , semi-unitary matrix ( $V_c^* V_c = I_{r \times r}$ )

$r \leq \min\{n, m\}$  is the rank of  $M$  matrix, and  $\Sigma_d$  has only the non-zero singular values of  $M$ .

$$\text{Thus, } \Sigma = \begin{bmatrix} \Sigma_d & 0 \\ 0 & 0 \end{bmatrix}$$

# Singular Value Decomposition (SVD)

$$M_{m \times n} = U_{m \times m} \sum_{m \times n} V^*_{n \times n}$$

The diagonal value are called  
the Singular values.



# Ridge Regression

$$y = f(x) + \epsilon$$

↑  
Random Error with mean  $\Rightarrow$  Singular values will never be zero if for colinear points but will be close to zero.

# Ridge Regression

If singular values are not zero then

$(X^T X)^{-1}$  will exist.

$$X = U \Sigma V^T \rightarrow \begin{bmatrix} \Sigma & 0 \\ 0 & 0 \end{bmatrix}$$

Diagram illustrating the SVD decomposition of matrix  $X$ . The matrix  $X$  is decomposed into  $U$ ,  $\Sigma$ , and  $V^T$ . A blue arrow points from the  $\Sigma$  term to a large blue matrix structure  $\begin{bmatrix} \Sigma & 0 \\ 0 & 0 \end{bmatrix}$ . A red arrow points from the  $\Sigma$  term to a red matrix structure  $\begin{bmatrix} \Sigma \\ 0 \end{bmatrix}$ .

$X = U \Sigma V^T = U \begin{bmatrix} \Sigma_d & \\ & 0 \end{bmatrix} V^T \rightarrow$  diagonal matrix with singular values  
**Ridge Regression**

$$(X^T X)^{-1} = \left( V \begin{bmatrix} \Sigma_d & 0 \\ 0 & 0 \end{bmatrix} U^T U \begin{bmatrix} \Sigma_d \\ 0 \end{bmatrix} V^T \right)^{-1}$$

$I$

If singular values are very small  $\Rightarrow$  large  $\beta$

$$\begin{aligned}
 &= \left( V \Sigma_d^2 V^T \right)^{-1} \\
 &= (V^T)^{-1} (\Sigma_d^2)^{-1} (V)^{-1} \\
 &= V \Sigma_d^{-2} V^T \rightarrow \text{diagonal with } 1 \text{ singular } 2
 \end{aligned}$$

# Ridge Regression

OLS regression

$$\min_{\beta} \|X\beta - y\|_2^2 \Rightarrow (X^T X)^T X^T y$$

Ridge Regression

$$\min_{\beta} \|X\beta - y\|_2^2 + \lambda \|\beta\|_2^2$$

↑  
+ve fitting param.



# Ridge Regression

$$RSS(\beta) = \|X\beta - y\|_2^2 + \lambda \|\beta\|_2^2$$

$$= \beta^T X^T X \beta - 2 \beta^T X^T y + y^T y + \lambda \beta^T \beta$$

$$\nabla_{\beta} RSS(\beta) = 2 X^T X \beta - 2 X^T y + 2 \lambda \beta$$

$$\nabla RSS(\beta) = 0$$

# Ridge Regression

$$\nabla_{\beta} \text{RSS}(\beta) = 2X^T X \beta - 2X^T Y + 2\lambda \beta$$

$$g) \quad \nabla \text{RSS} = 0$$

$$\Rightarrow 2X^T X \beta - 2X^T Y + 2\lambda \beta = 0$$

$$(2X^T X + 2\lambda I) \beta = 2X^T Y$$
$$\beta = (X^T X + \lambda I)^{-1} X^T Y$$