

Multivariable Linear Regression

1.1 Examples.

• Polynomial Regression

$$E[Y|X] = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \dots + \beta_k x_1^k = XB$$

$$X = [1, x_1, x_1^2, \dots, x_1^k] \quad B = [\beta_0 \ \beta_1 \ \dots \ \beta_k]^T$$

• Multiple Regression

$$E[Y|X] = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k = XB$$

$$X = [1, x_1, x_2, \dots, x_k] \quad B = [\beta_0 \ \beta_1 \ \dots \ \beta_k]^T$$

• Model with Interaction

$$E[Y|X] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 = XB$$

$$X = [1 \ x_1 \ x_2 \ x_1 x_2] \quad B = [\beta_0 \ \beta_1 \ \beta_2 \ \beta_{12}]^T$$

• Model Transformation

$$E[Y|X] = \beta_0 + \beta_1 \log(x_1) + \beta_2 \sin(x_2) = XB$$

$$X = [1 \ \log(x_1) \ \sin(x_2)] \quad B = [\beta_0 \ \beta_1 \ \beta_2]^T$$

Model Setup

In all cases, we use a linear model

$$E[Y|X] = XB \text{ where } X = [1 \ x_1 \ x_2 \ \dots \ x_k] \in \mathbb{R}^P$$

$B = [\beta_0 \ \beta_1 \ \dots \ \beta_k]$ is a $P \times 1$ column vector

$$P=1+k \quad (P=2 \text{ for SLR})$$

Assumption

- ① There are k predictors $x_1, x_2 \dots x_k$
We make no assumption about their distribution
They may or may not be dependent. X without subscript refers to $X = [x_1 \ x_2 \ \dots \ x_k]$
- ② There is a single response y
(multiple responses - multivariate)
- ③ Linear model
 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon = X\beta + \varepsilon$
for some constant $\beta_0, \beta_1 \dots \beta_k$
- ④ The noise has $E(\varepsilon) = E(\varepsilon|X) = 0$ (zero mean)
and $\text{Var}(\varepsilon) = \text{Var}(\varepsilon|X) = \sigma^2$ (constant variance)

If we have n observations $(y_i, \underbrace{x_{i1} \ x_{i2} \ \dots \ x_{ik}}_{x_i}) \ i=1 \dots n$
 $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}$
 $= x_i \beta + \varepsilon_i$
 $\quad (\beta_0, \beta_1 \dots \beta_k)$

In matrix form

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & & & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

$$Y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \varepsilon_{n \times 1} \quad (p=k+1)$$

where $E(\varepsilon|X) = 0_{n \times 1}$
 $\text{Var}(\varepsilon|X) = \sigma^2 I_n$ diag : $\text{Var}(\varepsilon_i)$
others : $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0$
assumption

$$\Rightarrow E(Y|X) = X\beta$$

$$\text{Var}(Y|X) = \sigma^2 I_n$$

1.3 Parameter Interpretation

σ^2 : the variance of the noise around the true regression line.

β_0 : the expected value of y when $x_1 = 0, x_2 = 0 \dots x_k = 0$

$$E(Y | x_1 = 0 \dots x_k = 0) = \bar{X}\beta$$

$$= \beta_0 + 0 \cdot \beta_1 + \dots + 0 \cdot \beta_k$$

$$= \beta_0$$

Note: In some case, $E(Y | x=0) = \beta_0$ might not have a meaning when $x=0$.

β_i for $i = 1 \dots k$. Take β_1 as an example.

$$\begin{aligned} & E(Y | x_1 = x_1 + 1, x_2 = x_2, \dots x_k = x_k) - E(Y | x_1 = x_1, \dots x_k = x_k) \\ &= (\beta_0 + \underline{\beta_1(x_1+1)} + \dots + \underline{\beta_k x_k}) - (\beta_0 + \underline{\beta_1 x_1} + \dots + \underline{\beta_k x_k}) \\ &= \beta_1 \end{aligned}$$

If we select two sets of cases from the distribution of the data, where x_1 differs by 1. We expect y to differ by β_1 on average

Note: $\beta_1 \neq Y | (x_1 = x_1 + 1, \dots x_k = x_k) - Y | (x_1 = x_1, \dots x_k = x_k)$

1.4 Least Squares Estimation

Estimates Provided with observed data $\{(y_1, x_1), \dots, (y_n, x_n)\}$ we estimate β using LS.

$$\begin{aligned} \text{We minimize } S(\beta) &= \sum_{i=1}^n (y_i - x_i \beta)^2 \\ &= (Y - X\beta)^T (X\beta - Y) \\ &= \|Y - X\beta\|^2 \end{aligned}$$

$$\frac{\partial S}{\partial \beta_j} = -2 \sum_{i=1}^n (y_i - x_i \beta_j) x_{ij} = 0 \quad j = 0, 1 \dots k$$

We have $p = k+1$ least squares normal equations.

which can also written in the matrix form

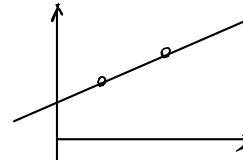
$$\nabla_{\beta} S(\beta) = -2X^T(Y - X\beta) = 0_p \Rightarrow X^T X \beta = X^T Y$$

This gives the solution for the least squares estimator

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

Note: $(X^T X)^{-1}$ exists if

- $n \geq p+1$
- The predictors must be linear independent, i.e. No columns of X is a linear combination of the other columns. (To check column rank of $X = k+1$)



at least you can have
a model (fit a
regression)

Fitted value The fitted value of the regression model

corresponding to $x_i = (1, x_{i1}, x_{i2}, \dots, x_{ik})$ is

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_k x_{ik}, \quad i=1 \dots n$$

In the matrix format, it is $\hat{Y} = X \hat{\beta} = X (X^T X)^{-1} X^T Y$

Statistical Properties of Least-squares Estimators

For estimator $\hat{\beta}$, if X is non-random, then

$$E(\hat{\beta}) = E(\hat{\beta}|X) = \beta$$

$$\text{Var}(\hat{\beta}|X) = \sigma^2 \frac{(X^T X)^{-1}}{C_{(k+1) \times (k+1)}} = C \cdot \sigma^2$$

Estimation of σ^2 - Estimate the variation of the data

The residual sum of squares for MLR is

$$\begin{aligned} SS_{\text{Res}} &= (Y - \hat{Y})^T (Y - \hat{Y}) \\ &= (Y - X \hat{\beta})^T (Y - X \hat{\beta}) \\ &= \|Y - X \hat{\beta}\|_2^2 \end{aligned}$$

For the estimator of σ^2 , we have

$$\begin{aligned} \hat{\sigma}^2 &= MS_{\text{Res}} = \frac{SS_{\text{Res}}}{n-p} = \frac{1}{n-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \frac{\|Y - X \hat{\beta}\|_2^2}{n-p} \end{aligned}$$

In R, residual standard error $\hat{\sigma} = \sqrt{\hat{\sigma}^2}$

We can show that $\hat{\sigma}^2$ is an unbiased estimator.

$$E(\hat{\sigma}^2 | X) = \sigma^2$$

Student's t-test

Requires the additional AN-SLR assumptions,
we further assume that

$$\begin{cases} \varepsilon \sim N(0, \sigma^2 I_n) \\ \varepsilon \text{ is independent of } X \end{cases}$$

From these assumptions conditional on X ,

$Y|X$ has a M.V (multivariate) Gaussian distribution

$$Y|X \sim N(X\beta, \sigma^2 I_n)$$

We can show that, $\hat{\beta}$ has a sample distribution

$$\hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$$

$$\hat{\beta}_j \sim N(\beta_j, \underbrace{[\sigma^2 (X^T X)^{-1}]_{(j+1)(j+1)}}_{\text{diagonal entry}}) \quad j=0 \dots k$$

To test $\begin{cases} H_0: \beta_j = 0 \\ H_1: \beta_j \neq 0 \end{cases}$ where $j = 0 \dots k$

Equivalently, $\begin{cases} H_0: Y = \beta_0 + \beta_1 X_1 + \dots + \beta_{j-1} X_{j-1} + \beta_{j+1} X_{j+1} + \dots + \beta_k X_k + \varepsilon \\ H_1: Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \varepsilon \end{cases}$

If H_0 is rejected, it indicates that the predictor X_j is likely to be meaningful addition to the model

We use the statistic $T_j = \frac{\hat{\beta}_j - 0}{\text{se}(\hat{\beta}_j)}$

where $\text{se}(\hat{\beta}_j) = \sqrt{C_{jj}} = \sqrt{\sigma^2 [(\mathbf{X}^T \mathbf{X})^{-1}]_{(j+1)(j+1)}}$
 $\text{MSRES} = \frac{\text{SSRES}}{n-p}$

If under the AN-SLR, $\varepsilon \sim N(0, \sigma^2 I_n)$, and under the $H_0: \beta_j = 0$

$$T_j \sim t_{n-p}$$

We reject H_0 if

$$|T_j| > k \equiv t_{\frac{\alpha}{2}, n-p}$$

or use the P-value

$$P(T > |T_j|) = \text{P-value} < \alpha.$$

F-test

Consider the regression $Y = X\beta + \varepsilon = [X_{(1)}, X_{(2)}] \begin{bmatrix} \beta_{(1)} \\ \beta_{(2)} \end{bmatrix} \xrightarrow{X} \beta_{(1)} + \beta_{(2)} + \varepsilon$

$$= X_{(1)}\beta_{(1)} + X_{(2)}\beta_{(2)} + \varepsilon$$

$\downarrow n \times p-r \quad \downarrow p-r x_1 \quad \downarrow n-r \quad \downarrow x_2$

We would like to test

$$H_0: \beta_{(2)} = 0 \text{ } r_{x_2}$$

$$H_1: \beta_{(2)} \neq 0 \text{ } r_{x_2}$$

Equivalently,

$$H_0: Y = X_{(1)}\beta_{(1)} + \varepsilon \text{ (reduced model)}$$

$$H_1: Y = X_{(1)}\beta_{(1)} + X_{(2)}\beta_{(2)} + \varepsilon \text{ (full model)}$$

Then the null hypothesis $\beta_{(2)} = 0$ can be tested by the statistic

$$F_0 = \frac{\overline{SSR}(\beta_{(2)} | \beta_{(1)}) / r}{\overline{SSRes} / (n-p)} = \frac{\overline{SSR}(\beta_{(2)} | \beta_{(1)}) / r}{MS_{Res}}$$

(r) degree of freedom

$$\begin{aligned} \overline{SSR}(\beta_{(2)} | \beta_{(1)}) &= \overline{SSR}(\hat{\beta}) - \overline{SSR}(\hat{\beta}_{(1)}) \\ &= \hat{\beta}^T X^T Y - \hat{\beta}_{(1)}^T X_{(1)}^T Y \\ &\quad \cdot \cdot \cdot \\ &\quad (X^T X)^{-1} X^T Y \quad (X_{(1)}^T X_{(1)})^{-1} X_{(1)}^T Y \end{aligned}$$

Extra sums-of-squares due to $\beta_{(2)}$ when $\beta_{(1)}$ is in the model.

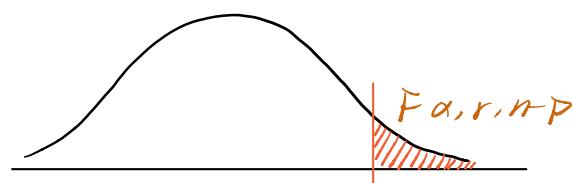
i.e. Extra contribution in \overline{SSR} due to the predictors in

$X_{(2)}$ under $H_0: \beta_{(2)} = 0$.

$$MS_{Res} = \frac{SS_{Res}}{n-p} = \frac{Y^T Y - \hat{\beta}^T Y}{n-p}$$

Under $H_0: \beta_{(2)} = 0$, under AN-SLR, the F_0 follows F distribution with $df_1 = r$, $df_2 = n-p$, we reject H_0 if

$$F_0 > F_{\alpha, r, n-p}$$



left threshold of p-value

Sequential F-test

• Test 1

$$H_0: Y = \beta_0 + \varepsilon$$

$$H_1: Y = \beta_0 + \beta_1 x_1 + \varepsilon$$

• Test 2

$$H_0: Y = \beta_0 + \beta_1 x_1 + \varepsilon$$

$$H_1: Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

drawback: Order matters!

• Test 3

$$H_0: Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

$$H_1: Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon$$

To compute F-statistic, we need to compute $\overline{SSR}(\beta_1 | \beta_0)$ for test 1, $\overline{SSR}(\beta_2 | \beta_1, \beta_0)$ for test 2 and $\overline{SSR}(\beta_3 | \beta_2, \beta_1, \beta_0)$ for test 3

$$F_1 = \frac{\overline{SSR}(\beta_1 | \beta_0) / 1}{\overline{SSRes}(\beta_1, \beta_2, \beta_1, \beta_0) / (n-4)} \quad \text{ending full model}$$

↳ β_2, β_3 will be considered later.

$$F_1 = \frac{\overline{SSR}(\beta_1 | \beta_0) / 1}{\overline{SSRes}(\beta_1, \beta_0) / (n-2)} \Rightarrow \text{results are equivalent when } n \rightarrow \infty$$

$$F_2 = \frac{\overline{SSR}(\beta_2 | \beta_1, \beta_0) / 1}{\overline{SSRes}(\beta_0, \beta_1, \beta_2, \beta_3) / (n-4)}$$

$$F_3 = \frac{\overline{SSR}(\beta_3 | \beta_2, \beta_1, \beta_0) / 1}{\overline{SSRes}(\beta_0, \beta_1, \beta_2, \beta_3) / (n-4)}$$

Note: The denominators are always the same.

$$\begin{aligned} & \overline{SSR}(\beta_3, \beta_2, \beta_1 | \beta_0) \\ &= \frac{\overline{SSR}(\beta_1, \beta_0) - \overline{SSR}(\beta_0)}{\overline{SSR}(\beta_1 | \beta_0)} + \frac{\overline{SSR}(\beta_2, \beta_1, \beta_0) - \overline{SSR}(\beta_1, \beta_0)}{\overline{SSR}(\beta_2 | \beta_1, \beta_0)} \\ & \quad + \frac{\overline{SSR}(\beta_3, \beta_2, \beta_1, \beta_0) - \overline{SSR}(\beta_2, \beta_1, \beta_0)}{\overline{SSR}(\beta_3 | \beta_2, \beta_1, \beta_0)} \\ &= \frac{\overline{SSR}(\beta_1 | \beta_0) + \overline{SSR}(\beta_2 | \beta_1, \beta_0) + \overline{SSR}(\beta_3 | \beta_2, \beta_1, \beta_0)}{\overline{SSR}(\beta_3, \beta_2, \beta_1, \beta_0) - \overline{SSR}(\beta_0)} \end{aligned}$$

Some technical details

Originally, we have the sum of squares decomposition

$$SST = SS_{\text{Res}} + SS_R$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

Now, we consider a different sum of squares.

$$\overline{SST} = SS_{\text{Res}} + \overline{SS_R}$$

$$\sum_{i=1}^n y_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n \hat{y}_i^2$$

$$\text{Since } \sum_{i=1}^n y_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n \hat{y}_i^2 + 2 \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i) \hat{y}_i}_0$$

Therefore,

$$\begin{aligned}
 \overline{SS_R} &= \sum_{i=1}^n y_i^2 \\
 &= (X\hat{\beta})^T (X\hat{\beta}) \\
 &= (\hat{\beta} X^T) (X(X^T X)^{-1} X^T Y) \\
 &= \hat{\beta} X^T Y \quad \text{def } P
 \end{aligned}$$

from normal equation

$$\underbrace{\sum_{i=1}^n (y_i - x_i \hat{\beta}) x_i \hat{\beta}}_0$$

Confidence Interval

① Confidence interval for estimates

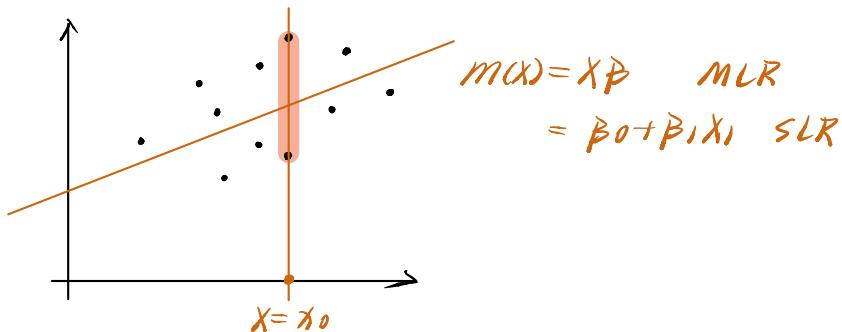
The $100(1-\alpha)\%$ confidence interval for the regression coefficients $\hat{\beta}_i$, $i = 0, 1, 2 \dots k$

$$\hat{\beta}_i \pm k \cdot \sqrt{\hat{\sigma}^2 \cdot C_{ii}} \quad k = t_{\frac{\alpha}{2}, n-p} \quad C = (X^T X)^{-1}$$

for $i = 0, 1 \dots k$

Note: The confidence interval will change using different batch of data.

② Confidence interval for mean response



$$x_0 = [1, x_{01}, x_{02} \dots x_{0k}]$$

The fitted value at this point is

$$\hat{m}(x_0) = x_0 \hat{\beta} = [1, x_0 \dots x_{0k}] \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_k \end{bmatrix}$$

This is the unbiased estimator of $E[Y|X=x_0]$. Since

$$E[\hat{m}(x_0)] = x_0 E[\hat{\beta}] = x_0 \hat{\beta}$$

and variance is

$$\text{Var}[\hat{m}(x_0)] = \hat{\sigma}^2 x_0 (X^T X)^{-1} x_0^T.$$

Therefore, a $100(1-\alpha)\%$ confidence interval for the mean response at $X=x_0 = [1, x_{01}, x_{02} \dots x_{0k}]$

$$CI(m(x_0)) = [\hat{m}(x_0) \pm k \cdot \sqrt{\hat{\sigma}^2 \cdot x_0 (X^T X)^{-1} x_0^T}]$$

③ Prediction interval of new observation

A $100(1-\alpha)\%$ prediction interval for the future observation predicated at $\hat{x} = \hat{x}_0 = [1, x_{01}, x_{02}, \dots, x_{0k}]$

$$PI(\hat{y}_0) = [\hat{m}(\hat{x}_0) \pm k \cdot \sqrt{\hat{\sigma}^2 (1 + \hat{x}_0 (\hat{x}^T \hat{x})^{-1} \hat{x}_0^T)}]$$

MLR

AN-MLR

estimator $\hat{\beta}, \hat{\sigma}^2$ ✓

unbiasedness of variance, E and variance of $\hat{\beta}, \hat{\sigma}^2$ ✓

t -test ✓

F -test ✓

CI($\hat{\beta}_i$) ✓

CI($m(\hat{x}_0)$) ✓

PI(\hat{y}_0) ✓

Model Adequacy Checking

The residuals are

$$e_i = y_i - \hat{y}_i = y_i - x_i \cdot \hat{\beta} \quad i=1 \dots n$$

written in the matrix for $\mathbf{e} = (e_1, e_2 \dots e_n)^T$

$$\begin{aligned} \mathbf{e} &= \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{x}\hat{\beta} = \mathbf{y} - \frac{\mathbf{x}(\mathbf{x}^T\mathbf{x})^{-1}\mathbf{x}^T\mathbf{y}}{H} \\ &= (\mathbf{I}_n - \mathbf{H})\mathbf{y}. \end{aligned}$$

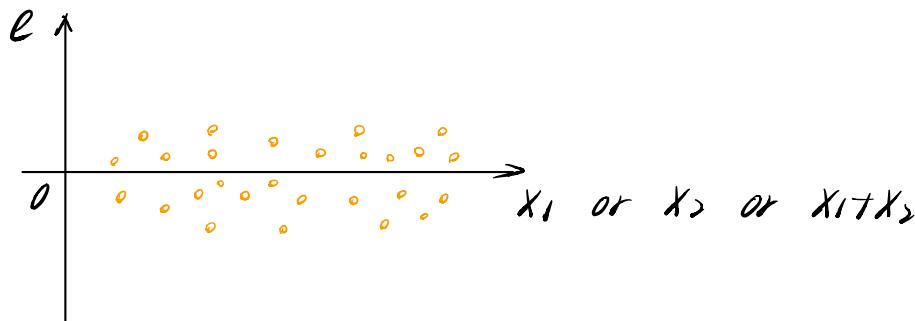
We have the following properties of e .

① The distribution of e should have a center around 0.

$$E[e|X] = 0_n$$

$$\begin{aligned} \text{Since } E[e|X] &= E[Y - \hat{Y}|X] = E[Y|X] - E[\hat{Y}|X] \\ &= X\beta - E[X\hat{\beta}|X] \\ &= 0 \quad E[\hat{\beta}] = \beta \end{aligned}$$

Plot e against any predictor or the linear combination of any predictor. This plot should have a constant mean equal to zero.



② The covariance between \hat{Y} and e is zero.

$$\text{Cov}(e, \hat{Y}) = 0_{n \times n}$$

$$\begin{aligned} \text{Cov}(e, \hat{Y}) &= \text{Cov}((\mathbf{I}_n - \mathbf{H})\mathbf{y}, \mathbf{H}\mathbf{y}) \\ &= \alpha^2 \underline{\mathbf{H}(\mathbf{I}_n - \mathbf{H})} \\ &= 0 \quad H = HH \text{ idempotent.} \end{aligned}$$

Plot e against \hat{Y} , the plot should be patternless.

MLR

③ The variance of residual e is $\text{Var}(e|X) = \sigma^2(I_n - H)$

$$\begin{aligned}\text{since } \text{Var}(e|X) &= \text{Var}((I_n - H)\gamma|X) \\ &= (I_n - H)^T(I_n - H) \text{Var}(\gamma|X) \\ &= (I_n - H) \sigma^2\end{aligned}$$

This implies that e_1, e_2, \dots, e_n are correlated, not like $\epsilon_1, \dots, \epsilon_n$ which are uncorrelated.

In the scalar form

$$\text{Var}(e_i|X) = \sigma^2(1 - h_{ii}) \quad i=1 \dots n$$

Usually there is no outlier in the data, the h_{ii} will be small. The variance of e should be roughly constant.

Plot e against x_i , the plot should have constant variance.

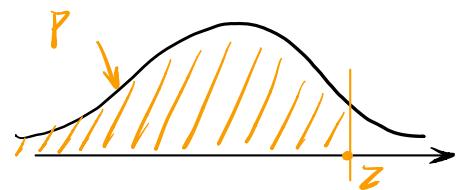
④ only for AN-MLR

If $\gamma|X$ is Gaussian, then $e = \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix}$ is also Gaussian with mean zero and variance $\sigma^2 I_n$. So each term $e_i \sim N(0, \sigma^2)$ for $i=1 \dots n$

- You can make a histogram of the distribution of the residuals, and compare that to a Gaussian
- quantile-quantile plot (Q-Q plot)

For a continuous distribution, its CDF $F(z) = P(Z \leq z)$ has an inverse function.

$$F^{-1}(P) = z \quad \text{s.t. } P = F(z)$$



If the distribution is $N(\mu, \sigma^2)$

$$F^{-1}(P) = \sigma \Phi^{-1}(P) + \mu$$

inverse of CDF function for $N(0,1)$

If we plot $F^{-1}(P)$ against $\Phi^{-1}(P)$, we'll get a straight line.

But in practice, F or F^{-1} is unknown, but we can estimate quantile using data.

We have n residuals e_1, e_2, \dots, e_n

Put them in an increasing order $e_{(1)}, e_{(2)}, \dots, e_{(n)}$, where $e_{(i)}$ is greater to a fraction $\frac{i}{n}$ of the residuals.

e.g. $e_{(2)} \geq e_{(1)}, e_{(1)} \left(\frac{1}{n}\right)$ of n obs.)

$$\Rightarrow \hat{F}^{-1}\left(\frac{i}{n}\right) = e_{(i)}$$

Now make a plot $\hat{F}^{-1}\left(\frac{i}{n}\right) = e_{(i)} \sim \Phi^{-1}\left(\frac{i}{n}\right)$ to approximate $F\left(\frac{i}{n}\right) \sim \Phi^{-1}\left(\frac{i}{n}\right)$ and hope to obtain a straight line if the distribution of e_i 's indeed Gaussian.

Transformation

Variance-stability transformation

We should consider the use of transformation when **constant variance assumption is violated**.

For example, $Y \sim \text{Poi}(\lambda)$. Then $E(Y) = \text{Var}(Y) = \lambda$.

Then if we model $E(Y|X=x) = XB = \text{Var}(Y|X=x)$

\uparrow not constant!

The mean of Y is related to X . To fix this, we can regress $Y' = \sqrt{Y}$ against X . Since the variance of \sqrt{Y} is independent of the mean.

Variance-mean relationship	Transformation.
$\sigma^2 \propto \text{constant}$	$Y' = Y$
$\sigma^2 \propto E(Y)$	$Y' = \sqrt{Y}$ (Poisson)
$\sigma^2 \propto E(Y)(1-E(Y))$	$Y' = \sin^{-1}(Y)$ (binomial)
$\sigma^2 \propto [E(Y)]^2$	$Y' = \ln(Y)$

Transformations to linearity

If we detect non-linearity using scatterplot or residual plots. In some case, a non-linear function can be linearized using a suitable transformation.

e.g. The Cobb-Douglas production function for observed economic data $i = 1 \dots n$

$$Q_i = e^{\beta_0} L_i^{\beta_1} C_i^{\beta_2} + u_i$$

↑ random error term
 total production labor input capital input

Take natural log, we have that

$$\ln(Q_i) = \beta_0 + \beta_1 \ln(L_i) + \beta_2 \ln(C_i) + \varepsilon_i$$

(ln(Q_i)) (ln(L_i)) (ln(C_i)) random error term

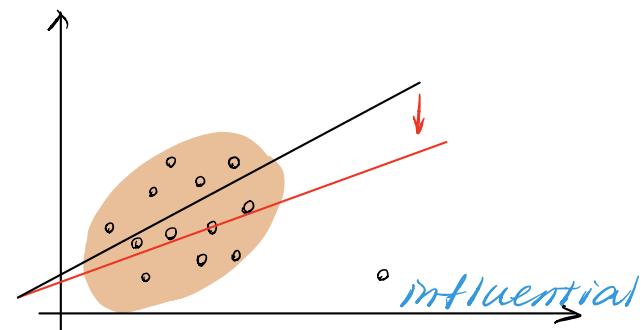
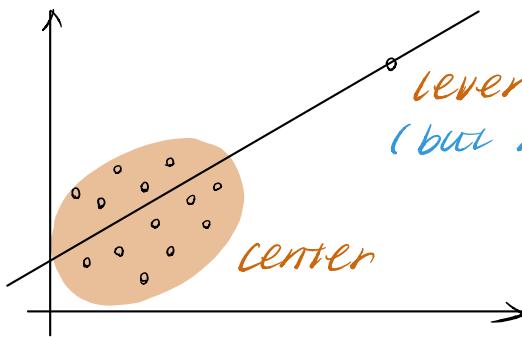
The transformation requires $\varepsilon = \ln(u) \sim N(0, \sigma^2)$

e.g. $y = \beta_0 + \beta_1 (\frac{1}{x}) + \varepsilon$ can be linearized

$$x' = \frac{1}{x}, y = \beta_0 + \beta_1 x' + \varepsilon$$

Original function	Transformation	Linear form
$y = \beta_0 x^{\beta_1}$	$y' = \log y, x' = \log x$	$y' = \log \beta_0 + \beta_1 x'$
$y = \beta_0 e^{\beta_1 x}$	$y' = \log y$	$y' = \log \beta_0 + \beta_1 x$
$y = \beta_0 + \beta_1 \log x$	$x' = \log x$	$y' = \beta_0 + \beta_1 x'$
$y = \frac{x}{\beta_0 x - \beta_1}$	$y' = \frac{1}{y}, x' = \frac{1}{x}$	$y' = \beta_0 - \beta_1 x'$

Leverage and Influence Point



Recall that $\hat{Y} = HY$, $H = X(X^T X)^{-1} X^T$

$$\text{Var}(\hat{Y} | X) = \sigma^2 H$$

$$\Rightarrow \text{Var}(Y - \hat{Y} | X) = \text{Var}(\epsilon | H) = \sigma^2 (I_n - H)$$

The elements of H dictate how much each data point $y_1 \dots y_n$ affects the prediction $\hat{y}_1 \dots \hat{y}_n$. i.e. $\hat{y}_i = \sum_{j=1}^n h_{ij} y_j$
 h_{ij} can be interpreted as the amount of **leverage** exerted by the j th observation y_j on the i th fitted value \hat{y}_i .

We focus on the diagonal elements h_{ii} of H

$$h_{ii} = x_i^T (X^T X)^{-1} x_i \quad \cdot \text{ith row of the matrix } X.$$

If h_{ii} is large, then the i th observation has high leverage.



Influence (Cook's distance)

Each data point may have a different "influence" on the estimation of B

$$D_i = \frac{(\hat{B}_{(i)} - \hat{B})^T X^T X (\hat{B}_{(i)} - \hat{B})}{c} \quad \cdot \text{PMSRes}$$

$\hat{B}_{(i)}$ is the fitted estimate with i th observation removed.

We can also rewrite D_i as

$$D_i = \frac{e_i^2}{\hat{\sigma}^2} \cdot \frac{h_{ii}}{(1-h_{ii})} \cdot \frac{1}{P} \quad \text{K+1 (# predictors including intercept)}$$

(standard residual)² leverage

$\Rightarrow D_i$ is affected by $\begin{cases} \text{residual } e_i \uparrow D_i \uparrow \\ \text{leverage } h_{ii} \uparrow D_i \uparrow \end{cases}$

D_i combines the residual for the i th observation and the location of that point to measure the influence of this point. A point with high D_i , potentially would be **outlier**.

There is another formula for D_i :

$$\text{Because } X\hat{\beta}_{(i)} - X\hat{\beta} = \hat{Y}_{(i)} - \hat{Y}$$

$$D_i = \frac{(\hat{Y}_{(i)} - \hat{Y})^T (\hat{Y}_{(i)} - \hat{Y})}{PM_{\text{Res}}} \quad \text{more robust}$$

You can interpret it as the squared distance that the fitted value moves when the i th observation is removed.

$$\left\langle \min_B \|Y - XB\|_2 \right\rangle$$

Special Types of predictors in MLR

Polynomial Term

Suppose the model is

$$\begin{aligned} E[Y|X] &= \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \dots + \beta_k X_1^k \\ &= \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k \end{aligned}$$

$k=3$ cubic function

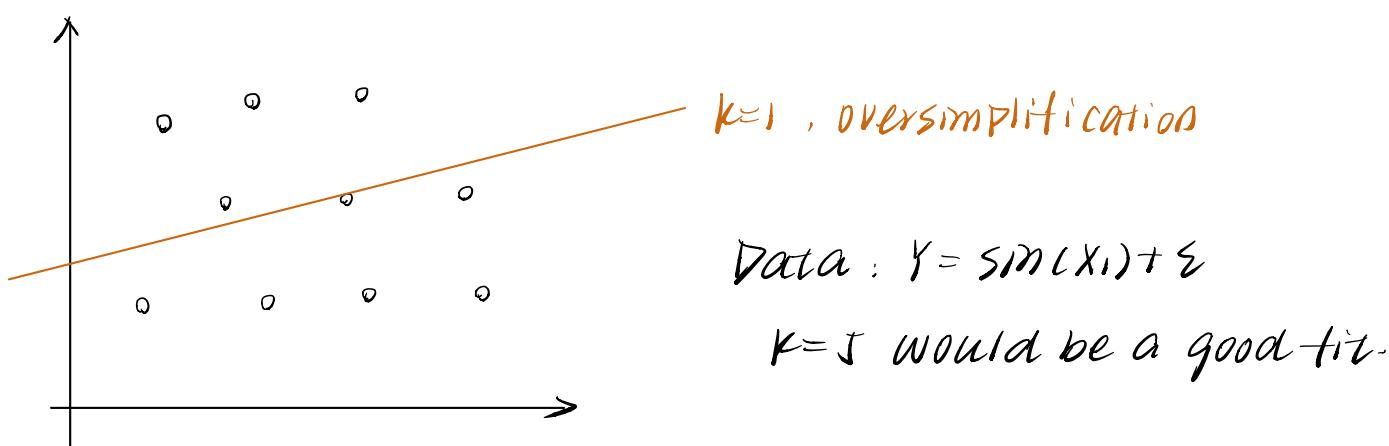
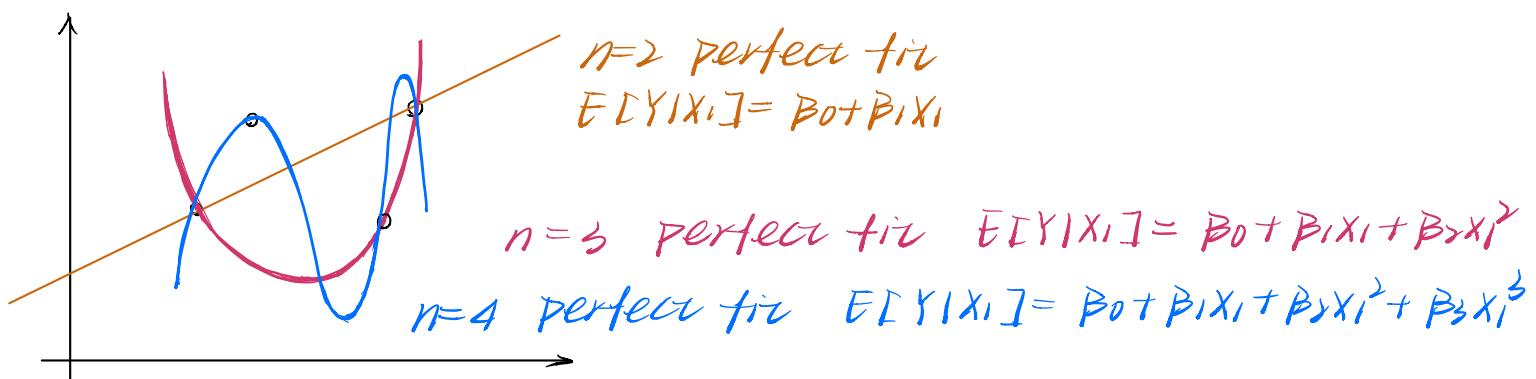


fig 1.7

Penalization

If we fit a polynomial regression with order k .

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_1^k + \varepsilon$$

or

$$E[Y|X_1] = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_1^k$$

Given n observations of data $\{y_i, x_{ii}\}_{i=1}^n$. You can estimate the value of $\beta_0, \beta_1, \dots, \beta_k$ by least squares. Denote $\beta = (\beta_0 \dots \beta_k)^T$

$$\min_{\beta} \sum_{i=1}^n \frac{(y_i - (\beta_0 + \beta_1 x_{ii} + \dots + \beta_k x_{ii}^k))^2}{\text{residual}}$$

If k is large, the model becomes too complicated or "*overfitted*". The corresponding fitted curve becomes wiggly.

To fix this issue, we restrict the size of the coefficients β by using penalization, i.e.

$$\min_{\beta} \sum_{i=1}^n \frac{(y_i - (\beta_0 + \beta_1 x_{ii} + \dots + \beta_k x_{ii}^k))^2}{\text{residual}}$$

such that $\|\beta'\|_2 \leq r$ where

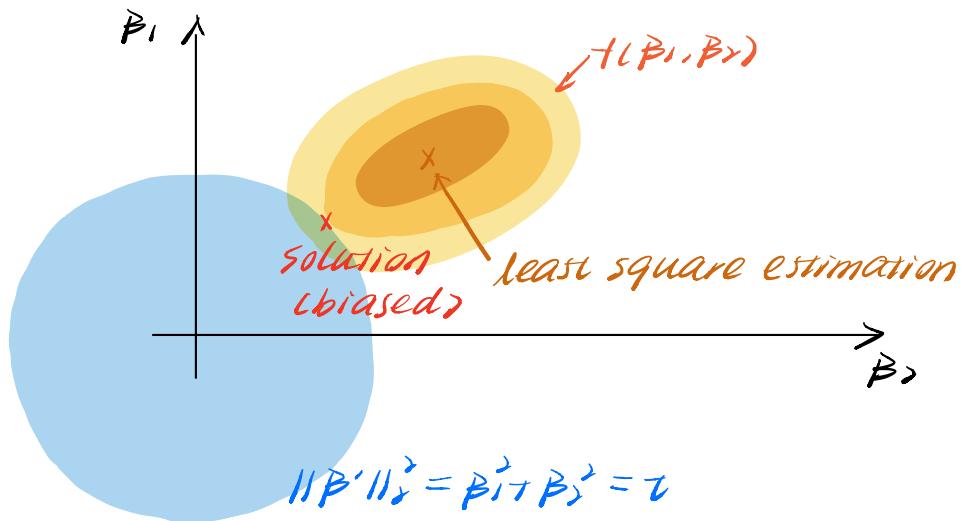
$$\|\beta'\|_2 = \sqrt{\beta_1^2 + \beta_2^2 + \dots + \beta_k^2} \quad \text{usually } \beta_0 \text{ is not penalized}$$

e.g. Consider the model with only $\beta_0, \beta_1, \beta_2$

$$\min_{\beta_0, \beta_1, \beta_2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{ii} - \beta_2 x_{ii}^2)^2 = c$$

$\downarrow f(\beta_1, \beta_2)$

$$\text{s.t. } \|\beta'\|_2 \leq r, \beta' = (\beta_1, \beta_2)^T$$



Mathematically, it is difficult to solve. Consider an alternative form.

$$\min_{\beta} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_k x_{ik})^2 + \frac{\lambda \|\beta'\|_2^2}{\lambda \geq 0}$$

For demonstration purpose, consider a no-intercept model.

$$\begin{aligned} & \min_{\beta} \sum_{i=1}^n (y_i - \beta_1 x_{i1} - \dots - \beta_k x_{ik})^2 + \lambda \|\beta'\|_2^2 \\ & = \min_{\beta} \frac{\|Y - X\beta\|^2}{\text{f}(P)} + \lambda \|\beta'\|_2^2 \quad \text{ridge regression} \end{aligned}$$

where $Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} x_{11} & \dots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nk} \end{pmatrix}$

$$\nabla_{\beta} f(\beta) = -2(Y - X\beta)^T X + 2\lambda\beta = 0$$

We solve β

$$Y^T X - X^T X \beta = \lambda \beta$$

$$(X^T X + \lambda I_P) \beta = Y^T X$$

$$\Rightarrow \hat{\beta} = (X^T X + \lambda I_P)^{-1} Y^T X$$

usage of ridge regression

- ① control model complexity
- ② deal with high collinearity between predictors.

Y	x_1	x_2	x_3
1	2		
2	4		
3	6		
4	8		

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

becomes singular when two columns of x_j, x_i are perfectly correlated.

- ③ When $n < p$, n is the number of obversations, p is the number of predictors.

$$n \left\{ \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \left(\begin{matrix} 1 & x_{11} & \dots & x_{1k} \\ \vdots & & & \\ 1 & x_{n1} & \dots & x_{nk} \end{matrix} \right) \right\}$$

P

MLR could fail in this case,
but you should try ridge regression.

- ④ SVM with overlapped case (linear model)
 ≈ logistic regression with ridge penalization.

$$\text{SVM: } \sum_{i=1}^n (1 - y_i(\beta_0 + \beta^T x_i)) + \lambda \|\beta\|_2$$

$$\text{Logistic: } \sum_{i=1}^n \log(1 + e^{-y_i(\beta_0 + \beta^T x_i)}) + \lambda \|\beta\|_2$$

