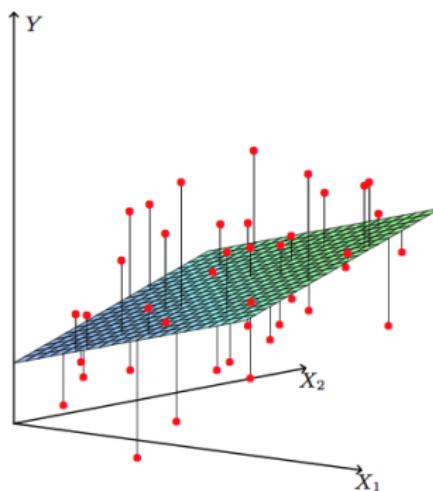


# Statistical Methods for Data Mining

Kuangnan Fang

Xiamen University *Email:* xmufkn@xmu.edu.cn

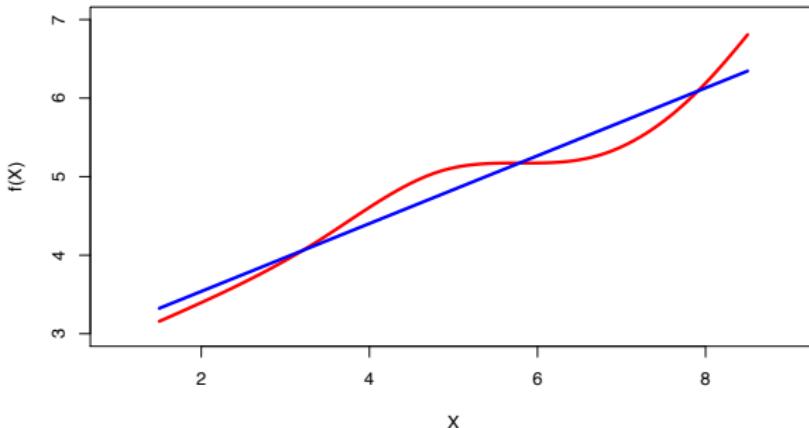


## Linear regression

- Linear regression is a simple approach to supervised learning. It assumes that the dependence of  $Y$  on  $X_1, X_2, \dots, X_p$  is linear.

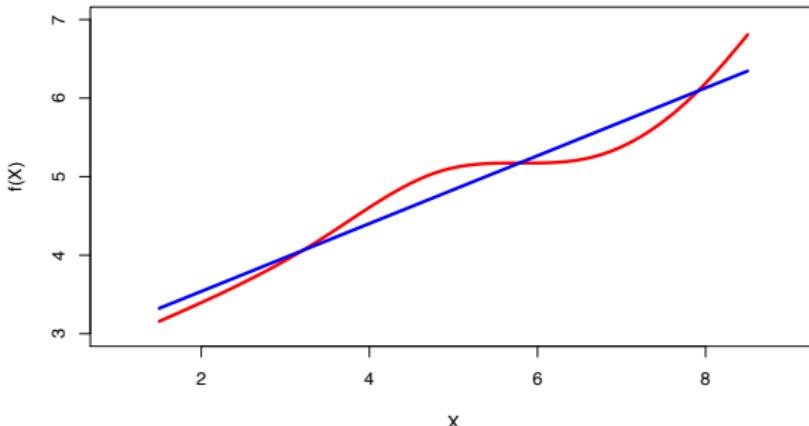
# Linear regression

- Linear regression is a simple approach to supervised learning. It assumes that the dependence of  $Y$  on  $X_1, X_2, \dots, X_p$  is linear.
- True regression functions are never linear!



## Linear regression

- Linear regression is a simple approach to supervised learning. It assumes that the dependence of  $Y$  on  $X_1, X_2, \dots, X_p$  is linear.
- True regression functions are never linear!



- although it may seem overly simplistic, linear regression is extremely useful both conceptually and practically.

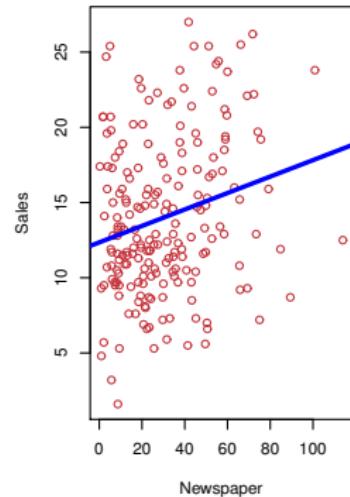
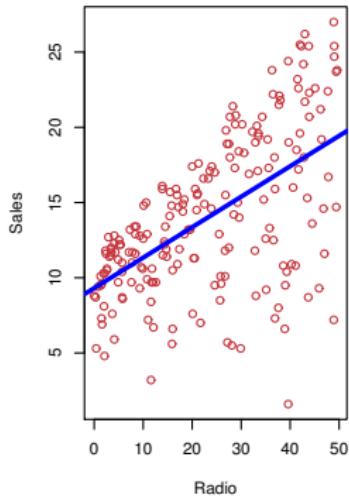
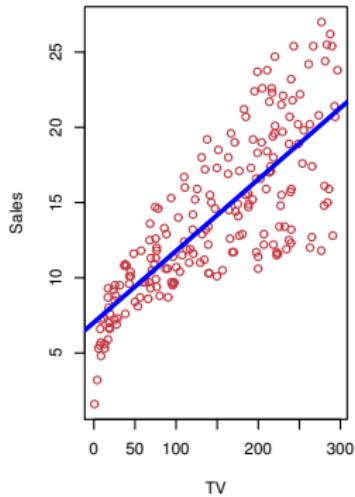
## Linear regression for the advertising data

Consider the advertising data shown on the next slide.

Questions we might ask:

- Is there a relationship between advertising budget and sales?
- How strong is the relationship between advertising budget and sales?
- Which media contribute to sales?
- How accurately can we predict future sales?
- Is the relationship linear?
- Is there synergy among the advertising media?

# Advertising data



# Simple linear regression using a single predictor $X$ .

- We assume a model

$$Y = \beta_0 + \beta_1 X + \epsilon,$$

where  $\beta_0$  and  $\beta_1$  are two unknown constants that represent the *intercept* and *slope*, also known as *coefficients* or *parameters*, and  $\epsilon$  is the error term.

- Given some estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$  for the model coefficients, we predict future sales using

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x,$$

where  $\hat{y}$  indicates a prediction of  $Y$  on the basis of  $X = x$ . The *hat* symbol denotes an estimated value.

## Estimation of the parameters by least squares

- Let  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  be the prediction for  $Y$  based on the  $i$ th value of  $X$ . Then  $e_i = y_i - \hat{y}_i$  represents the  $i$ th *residual*

## Estimation of the parameters by least squares

- Let  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  be the prediction for  $Y$  based on the  $i$ th value of  $X$ . Then  $e_i = y_i - \hat{y}_i$  represents the  $i$ th *residual*
- We define the *residual sum of squares* (RSS) as

$$\text{RSS} = e_1^2 + e_2^2 + \cdots + e_n^2,$$

or equivalently as

$$\text{RSS} = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2.$$

## Estimation of the parameters by least squares

- Let  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  be the prediction for  $Y$  based on the  $i$ th value of  $X$ . Then  $e_i = y_i - \hat{y}_i$  represents the  $i$ th *residual*
- We define the *residual sum of squares* (RSS) as

$$\text{RSS} = e_1^2 + e_2^2 + \cdots + e_n^2,$$

or equivalently as

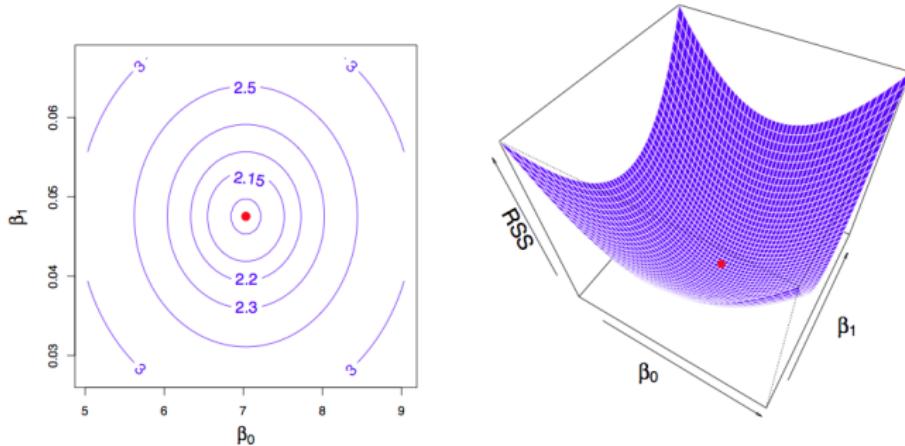
$$\text{RSS} = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2.$$

- The least squares approach chooses  $\hat{\beta}_0$  and  $\hat{\beta}_1$  to minimize the RSS. The minimizing values can be shown to be

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2},$$

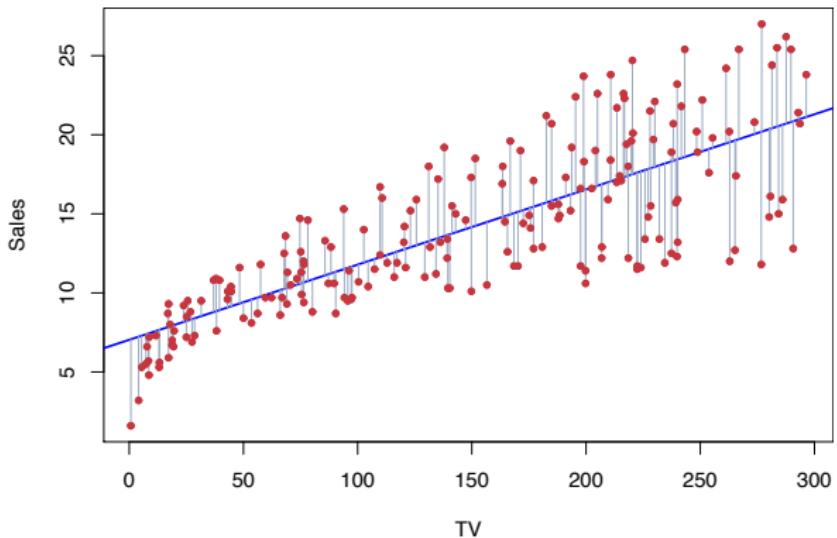
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$

where  $\bar{y} \equiv \frac{1}{n} \sum_{i=1}^n y_i$  and  $\bar{x} \equiv \frac{1}{n} \sum_{i=1}^n x_i$  are the sample means.



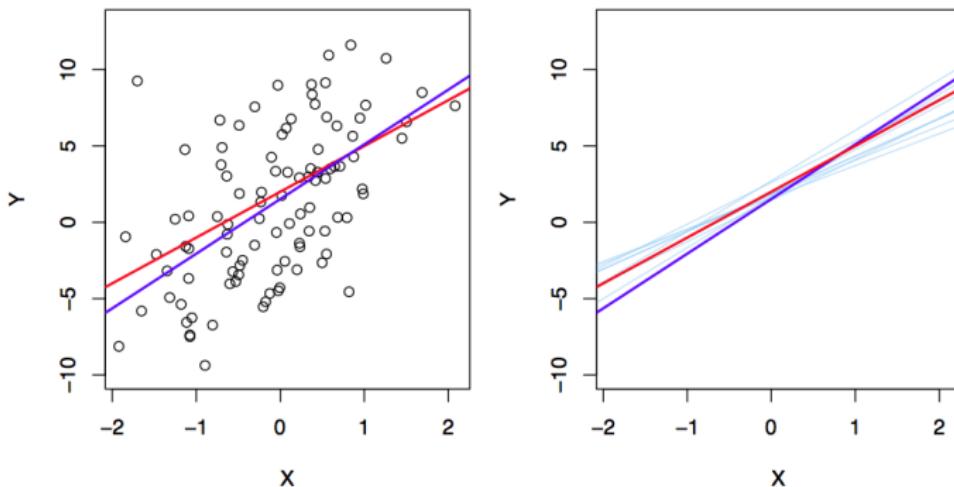
**FIGURE 3.2.** Contour and three-dimensional plots of the RSS on the Advertising data, using **sales** as the response and **TV** as the predictor. The red dots correspond to the least squares estimates  $\hat{\beta}_0$  and  $\hat{\beta}_1$ , given by (3.4).

## Example: advertising data



The least squares fit for the regression of **sales** onto **TV**.  
In this case a linear fit captures the essence of the relationship,  
although it is somewhat deficient in the left of the plot.

# Assessing the Accuracy of the Coefficients Estimates



**FIGURE 3.3.** A simulated data set. Left: The red line represents the true relationship,  $f(X) = 2 + 3X$ , which is known as the population regression line. The blue line is the least squares line; it is the least squares estimate for  $f(X)$  based on the observed data, shown in black. Right: The population regression line is again shown in red, and the least squares line in dark blue. In light blue, 10 least squares lines are shown, each computed on the basis of a separate random set of observations. Each least squares line is different, but on average, the least squares lines are quite close to the population regression line.

## Assessing the Accuracy of the Coefficient Estimates

- The standard error of an estimator reflects how it varies under repeated sampling. We have

$$\text{SE}(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \text{SE}(\hat{\beta}_0)^2 = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right],$$

where  $\sigma^2 = \text{Var}(\epsilon)$

# Assessing the Accuracy of the Coefficient Estimates

- The standard error of an estimator reflects how it varies under repeated sampling. We have

$$\text{SE}(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \text{SE}(\hat{\beta}_0)^2 = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right],$$

where  $\sigma^2 = \text{Var}(\epsilon)$

- These standard errors can be used to compute *confidence intervals*. A 95% confidence interval is defined as a range of values such that with 95% probability, the range will contain the true unknown value of the parameter. It has the form

$$\hat{\beta}_1 \pm 2 \cdot \text{SE}(\hat{\beta}_1).$$

## Confidence intervals — continued

That is, there is approximately a 95% chance that the interval

$$\left[ \hat{\beta}_1 - 2 \cdot \text{SE}(\hat{\beta}_1), \hat{\beta}_1 + 2 \cdot \text{SE}(\hat{\beta}_1) \right]$$

will contain the true value of  $\beta_1$  (under a scenario where we got repeated samples like the present sample)

## Confidence intervals — continued

That is, there is approximately a 95% chance that the interval

$$\left[ \hat{\beta}_1 - 2 \cdot \text{SE}(\hat{\beta}_1), \hat{\beta}_1 + 2 \cdot \text{SE}(\hat{\beta}_1) \right]$$

will contain the true value of  $\beta_1$  (under a scenario where we got repeated samples like the present sample)

For the advertising data, the 95% confidence interval for  $\beta_1$  is  
[0.042, 0.053]

## Hypothesis testing

- Standard errors can also be used to perform *hypothesis tests* on the coefficients. The most common hypothesis test involves testing the *null hypothesis* of

$H_0$  : There is no relationship between  $X$  and  $Y$

versus the *alternative hypothesis*

$H_A$  : There is some relationship between  $X$  and  $Y$ .

## Hypothesis testing

- Standard errors can also be used to perform *hypothesis tests* on the coefficients. The most common hypothesis test involves testing the *null hypothesis* of

$H_0$  : There is no relationship between  $X$  and  $Y$

versus the *alternative hypothesis*

$H_A$  : There is some relationship between  $X$  and  $Y$ .

- Mathematically, this corresponds to testing

$$H_0 : \beta_1 = 0$$

versus

$$H_A : \beta_1 \neq 0,$$

since if  $\beta_1 = 0$  then the model reduces to  $Y = \beta_0 + \epsilon$ , and  $X$  is not associated with  $Y$ .

## Hypothesis testing — continued

- To test the null hypothesis, we compute a *t-statistic*, given by

$$t = \frac{\hat{\beta}_1 - 0}{\text{SE}(\hat{\beta}_1)},$$

- This will have a *t*-distribution with  $n - 2$  degrees of freedom, assuming  $\beta_1 = 0$ .
- Using statistical software, it is easy to compute the probability of observing any value equal to  $|t|$  or larger. We call this probability the *p-value*.

## Results for the advertising data



The table displays statistical results for a regression model. The columns are labeled: Coefficient, Std. Error, t-statistic, and p-value. The Intercept row shows values 7.0325, 0.4578, 15.36, and < 0.0001 respectively. The TV row shows values 0.0475, 0.0027, 17.67, and < 0.0001 respectively.

	Coefficient	Std. Error	t-statistic	p-value
Intercept	7.0325	0.4578	15.36	< 0.0001
TV	0.0475	0.0027	17.67	< 0.0001

## Assessing the Overall Accuracy of the Model

- We compute the *Residual Standard Error*

$$\text{RSE} = \sqrt{\frac{1}{n-2} \text{RSS}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2},$$

where the *residual sum-of-squares* is  $\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ .

## Assessing the Overall Accuracy of the Model

- We compute the *Residual Standard Error*

$$y = \beta_0 + \beta_1 x$$

2.7 教  
n-2

$$\text{RSE} = \sqrt{\frac{1}{n-2} \text{RSS}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2},$$

where the *residual sum-of-squares* is  $\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ .

- R-squared* or fraction of variance explained is

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

- 有問題  
- 有問題

where  $\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2$  is the *total sum of squares*.

## Assessing the Overall Accuracy of the Model

- We compute the *Residual Standard Error*

$$\text{RSE} = \sqrt{\frac{1}{n-2} \text{RSS}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2},$$

where the *residual sum-of-squares* is  $\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ .

- *R-squared* or fraction of variance explained is

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

where  $\text{TSS} = \sum_{i=1}^n (y_i - \bar{y})^2$  is the *total sum of squares*.

- It can be shown that in this simple linear regression setting that  $R^2 = r^2$ , where  $r$  is the correlation between  $X$  and  $Y$ :

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}.$$

## Advertising data results

Quantity	Value
Residual Standard Error	3.26
$R^2$	0.612
F-statistic	312.1

# Multiple Linear Regression

- Here our model is

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

- We interpret  $\beta_j$  as the *average* effect on  $Y$  of a one unit increase in  $X_j$ , *holding all other predictors fixed*. In the advertising example, the model becomes

$$\text{sales} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times \text{newspaper} + \epsilon.$$

# Interpreting regression coefficients

- The ideal scenario is when the predictors are uncorrelated
  - a *balanced design*:  
多元回归 = 遵守平衡原则的回归
  - Each coefficient can be estimated and tested separately.
  - Interpretations such as “*a unit change in  $X_j$  is associated with a  $\beta_j$  change in  $Y$ , while all the other variables stay fixed*”, are possible.
- Correlations amongst predictors cause problems:  
*Var 1*
  - The variance of all coefficients tends to increase, sometimes dramatically
  - Interpretations become hazardous — when  $X_j$  changes, everything else changes. 不再为因果关系。
- *Claims of causality* should be avoided for observational data.

# The woes of (interpreting) regression coefficients

*“Data Analysis and Regression” Mosteller and Tukey 1977*

- a regression coefficient  $\beta_j$  estimates the expected change in  $Y$  per unit change in  $X_j$ , *with all other predictors held fixed*. But **predictors usually change together!**

# The woes of (interpreting) regression coefficients

*“Data Analysis and Regression” Mosteller and Tukey 1977*

- a regression coefficient  $\beta_j$  estimates the expected change in  $Y$  per unit change in  $X_j$ , *with all other predictors held fixed*. But predictors usually change together!
- Example:  $Y$  total amount of change in your pocket;  
 $X_1 = \#$  of coins;  $X_2 = \#$  of pennies, nickels and dimes. By itself, regression coefficient of  $Y$  on  $X_2$  will be  $> 0$ . But how about with  $X_1$  in model?

# The woes of (interpreting) regression coefficients

*“Data Analysis and Regression” Mosteller and Tukey 1977*

- a regression coefficient  $\beta_j$  estimates the expected change in  $Y$  per unit change in  $X_j$ , *with all other predictors held fixed*. But predictors usually change together!
- Example:  $Y$  total amount of change in your pocket;  $X_1 = \#$  of coins;  $X_2 = \#$  of pennies, nickels and dimes. By itself, regression coefficient of  $Y$  on  $X_2$  will be  $> 0$ . But how about with  $X_1$  in model?
- $Y$  = number of tackles by a football player in a season;  $W$  and  $H$  are his weight and height. Fitted regression model is  $\hat{Y} = b_0 + .50W - .10H$ . How do we interpret  $\hat{\beta}_2 < 0$ ?

多元变量将存在相关性

$\hat{\beta}$  焦点时... 可能是负向

## Two quotes by famous Statisticians

*“Essentially, all models are wrong, but some are useful”*

George Box

## Two quotes by famous Statisticians

*“Essentially, all models are wrong, but some are useful”*

George Box

*“The only way to find out what will happen when a complex system is disturbed is to disturb the system, not merely to observe it passively”*

Fred Mosteller and John Tukey, paraphrasing George Box

## Estimation and Prediction for Multiple Regression

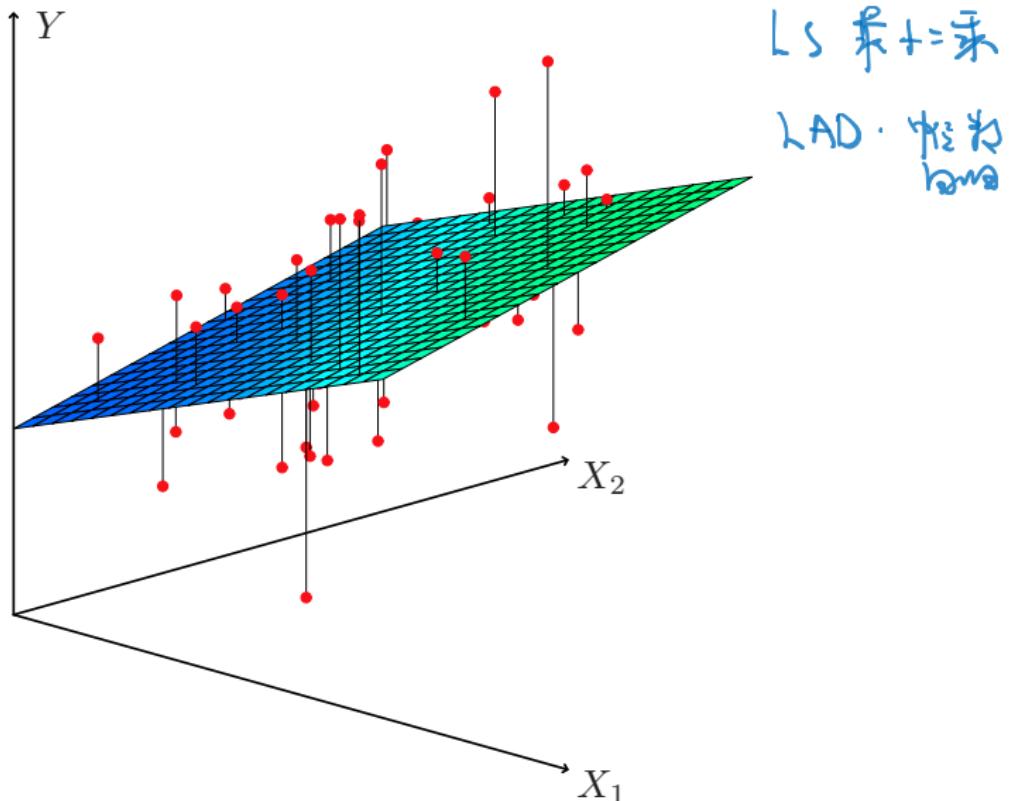
- Given estimates  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ , we can make predictions using the formula

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

- We estimate  $\beta_0, \beta_1, \dots, \beta_p$  as the values that minimize the sum of squared residuals

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

This is done using standard statistical software. The values  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$  that minimize RSS are the multiple least squares regression coefficient estimates.



## Results for advertising data

	Coefficient	Std. Error	t-statistic	p-value
Intercept	2.939	0.3119	9.42	< 0.0001
TV	0.046	0.0014	32.81	< 0.0001
radio	0.189	0.0086	21.89	< 0.0001
newspaper	-0.001	0.0059	-0.18	0.8599

	Correlations:			
	TV	radio	newspaper	sales
TV	1.0000	0.0548	0.0567	0.7822
radio		1.0000	0.3541	0.5762
newspaper			1.0000	0.2283
sales				1.0000

## Some important questions

1. *Is at least one of the predictors  $X_1, X_2, \dots, X_p$  useful in predicting the response?*

## Some 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?*

## Some 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?*

## Some 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?*

## Is at least one predictor useful?

For the first question, we can use the F-statistic

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

Quantity	Value
Residual Standard Error	1.69
$R^2$	0.897
F-statistic	570

## Deciding on the important variables

- The most direct approach is called *all subsets* or *best subsets* regression: we compute the least squares fit for all possible subsets and then choose between them based on some criterion that balances training error with model size.

## Deciding on the important variables

- The most direct approach is called *all subsets* or *best subsets* regression: we compute the least squares fit for all possible subsets and then choose between them based on some criterion that balances training error with model size.
- However we often can't examine all possible models, since they are  $2^p$  of them; for example when  $p = 40$  there are over a billion models!

Instead we need an automated approach that searches through a subset of them. We discuss two commonly used approaches next.

## Forward selection

- Begin with the *null model* — a model that contains an intercept but no predictors.
- Fit  $p$  simple linear regressions and add to the null model the variable that results in the lowest RSS.
- Add to that model the variable that results in the lowest RSS amongst all two-variable models.
- Continue until some stopping rule is satisfied, for example when all remaining variables have a p-value above some threshold.

误差量小  
但不降低系数

## Backward selection

- Start with all variables in the model.
- Remove the variable with the largest p-value — that is, the variable that is the least statistically significant.
- The new  $(p - 1)$ -variable model is fit, and the variable with the largest p-value is removed.
- Continue until a stopping rule is reached. For instance, we may stop when all remaining variables have a significant p-value defined by some significance threshold.

## Model selection — continued

- Later we discuss more systematic criteria for choosing an “optimal” member in the path of models produced by forward or backward stepwise selection.
- These include *Mallow's  $C_p$* , *Akaike information criterion (AIC)*, *Bayesian information criterion (BIC)*, *adjusted  $R^2$*  and *Cross-validation (CV)*.

变量数不同的模型模型

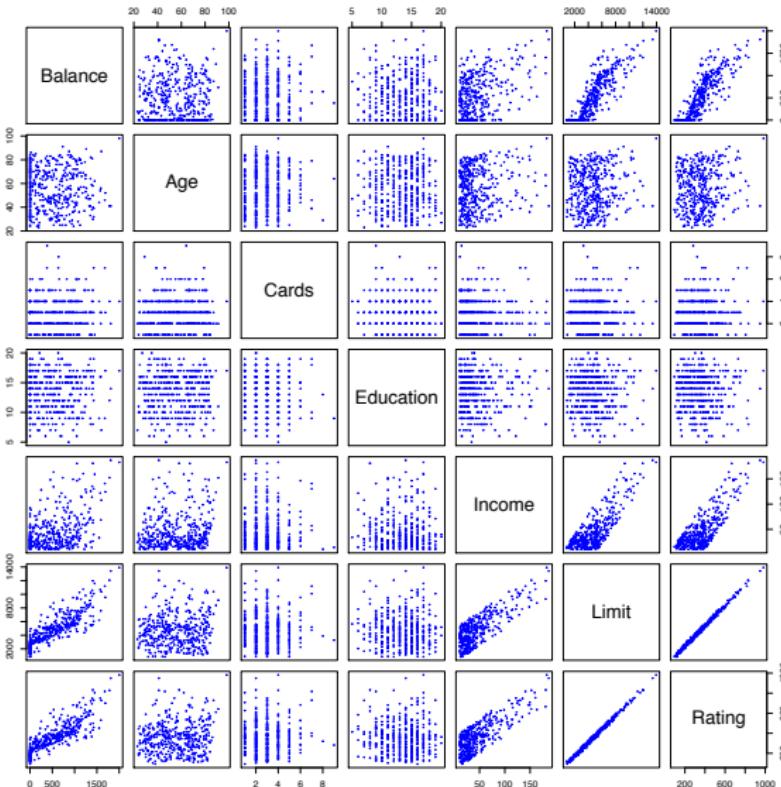
# Other Considerations in the Regression Model

## *Qualitative Predictors*

- Some predictors are not *quantitative* but are *qualitative*, taking a discrete set of values.
- These are also called *categorical* predictors or *factor variables*.
- See for example the scatterplot matrix of the credit card data in the next slide.

In addition to the 7 quantitative variables shown, there are four qualitative variables: **gender**, **student** (student status), **status** (marital status), and **ethnicity** (Caucasian, African American (AA) or Asian).

# Credit Card Data



## Qualitative Predictors — continued

Example: investigate differences in credit card balance between males and females, ignoring the other variables. We create a new variable

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

Resulting model:

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

Annotations in red:

- Top right:  $y_{\text{female}} \text{ 平均值}$
- Bottom right:  $y_{\text{male}} \text{ 平均值}$
- Right side:  $\beta_0 + \beta_1$  (with a red circle around it) is labeled "男女平均值差异" (β<sub>1</sub> 差异, 即).
- Bottom right: ANOVA

Intrepretation?

what if code as 1 and -1?

## Credit card data — continued

Results for gender model:

	Coefficient	Std. Error	t-statistic	p-value
$\beta_0$ Intercept	509.80	33.13	15.389	< 0.0001
$\beta_1$ gender [Female]	19.73	46.05	0.429	0.6690

Annotations in red:

- A red arrow points from the text "男性平均值" to the coefficient value 509.80.
- A red bracket groups the coefficient 19.73 and the p-value 0.6690, with a red arrow pointing from the text "女性平均值" to the bracket.
- Handwritten text "不显著" is written next to the p-value 0.6690.
- Handwritten text "女男相" is written next to the coefficient 19.73.
- Handwritten text "女男相" is also written next to the p-value 0.6690.
- Handwritten text "女生平均值" is written next to the coefficient 19.73.
- Handwritten text "女生平均值" is also written next to the p-value 0.6690.

## Qualitative predictors with more than two levels

- With more than two levels, we create additional dummy variables. For example, for the **ethnicity** variable we create two dummy variables. The first could be

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

and the second could be

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

## Qualitative predictors with more than two levels — continued.

- Then both of these variables can be used in the regression equation, in order to obtain the model

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

- There will always be one fewer dummy variable than the number of levels. The level with no dummy variable — African American in this example — is known as the

*baseline.* *看 design matrix*

$$\left( \begin{array}{cccc} | & | & | & | \\ | & | & | & | \\ | & | & | & | \end{array} \right)$$

三个水平的 dummy 变量  
其中一个被选为 baseline  
即为 3 - 2 = 1 个

无类别  
非白人  
可用 3 个  
dummy  
(水平数相同)

## Results for ethnicity

	Coefficient	Std. Error	t-statistic	p-value
Intercept	531.00	46.32	11.464	< 0.0001
ethnicity[Asian]	-18.69	65.02	-0.287	0.7740
ethnicity[Caucasian]	-12.50	56.68	-0.221	0.8260

## Extensions of the Linear Model

交互項.

Removing the additive assumption: *interactions* and *nonlinearity*

### *Interactions:*

- In our previous analysis of the **Advertising** data, we assumed that the effect on **sales** of increasing one advertising medium is independent of the amount spent on the other media.
- For example, the linear model

$$\widehat{\text{sales}} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times \text{newspaper}$$

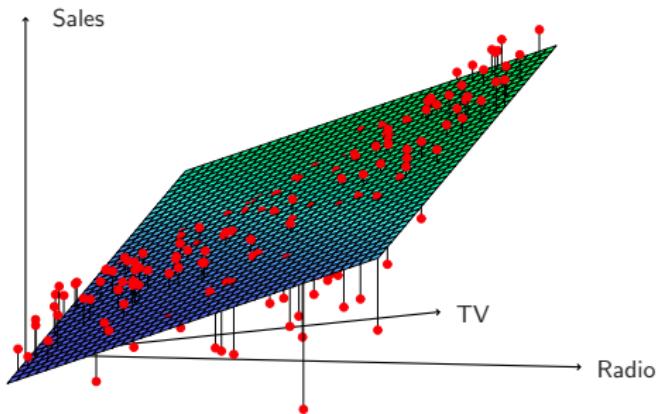
$\beta_1$  是 TV 的回歸係數。  
 $\beta_2$  不回歸

states that the average effect on **sales** of a one-unit increase in **TV** is always  $\beta_1$ , regardless of the amount spent on **radio**.

## Interactions — continued

- But suppose that spending money on radio advertising actually increases the effectiveness of TV advertising, so that the slope term for **TV** should increase as **radio** increases.
- In this situation, given a fixed budget of \$100,000, spending half on **radio** and half on **TV** may increase **sales** more than allocating the entire amount to either **TV** or to **radio**.  
*投一半钱在radio - 销售量 sales 高  
✓ 加倍效果.*
- In marketing, this is known as a *synergy* effect, and in statistics it is referred to as an *interaction* effect.

# Interaction in the Advertising data?



When levels of either **TV** or **radio** are low, then the true **sales** are lower than predicted by the linear model.  
But when advertising is split between the two media, then the model tends to underestimate **sales**.

# Modelling interactions — Advertising data

Model takes the form

$$\begin{aligned}\text{sales} &= \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times (\text{radio} \times \text{TV}) + \epsilon \\ &= \beta_0 + (\beta_1 + \beta_3 \times \text{radio}) \times \text{TV} + \beta_2 \times \text{radio} + \epsilon.\end{aligned}$$

交互项  
↑

TV的系数应该成为radio的系数。

Results:

重估。

$$= \beta_0 + \beta_1 \text{TV} + (\beta_2 + \beta_3 \text{TV}) \text{radio} + \epsilon.$$

同理

	Coefficient	Std. Error	t-statistic	p-value
Intercept	6.7502	0.248	27.23	< 0.0001
TV	0.0191	0.002	12.70	< 0.0001
radio	0.0289	0.009	3.24	0.0014
TV×radio	$\beta_3$	0.0011	20.73	< 0.0001

→ radio越高，  
TV的回归效果越好。  
→ TV越高，  
radio的回归效果越好。  
TV越高，  
Radio的回归效果越好。  
模型存在交互效应。

## Interpretation

- The results in this table suggests that interactions are important.
- The p-value for the interaction term  $\text{TV} \times \text{radio}$  is extremely low, indicating that there is strong evidence for  $H_A : \beta_3 \neq 0$ .
- The  $R^2$  for the interaction model is 96.8%, compared to only 89.7% for the model that predicts **sales** using **TV** and **radio** without an interaction term.

## Interpretation — continued

未被解释部分有69%可用交互项解释。

- This means that  $(96.8 - 89.7)/(100 - 89.7) = 69\%$  of the variability in **sales** that remains after fitting the additive model has been explained by the interaction term.
- The coefficient estimates in the table suggest that an increase in TV advertising of \$1,000 is associated with increased sales of
$$(\hat{\beta}_1 + \hat{\beta}_3 \times \text{radio}) \times 1000 = 19 + 1.1 \times \text{radio}$$
 units.
- An increase in radio advertising of \$1,000 will be associated with an increase in sales of
$$(\hat{\beta}_2 + \hat{\beta}_3 \times \text{TV}) \times 1000 = 29 + 1.1 \times \text{TV}$$
 units.

# Hierarchy

交互项显著但用于建立的原变量不显著  
main effect

- Sometimes it is the case that an interaction term has a very small p-value, but the associated main effects (in this case, **TV** and **radio**) do not.
- The *hierarchy principle*: *原则：既然 main effect 7. 显著，那么交互项显著，main effect 也要加进来*  
*If we include an interaction in a model, we should also include the main effects, even if the p-values associated with their coefficients are not significant.*

## Hierarchy — continued

原因：无 main effect 的交互项没有意义。

- The rationale for this principle is that interactions are hard to interpret in a model without main effects — their meaning is changed.
- Specifically, the interaction terms also contain main effects, if the model has no main effect terms.

## Interactions between qualitative and quantitative variables

Consider the **Credit** data set, and suppose that we wish to predict **balance** using **income** (quantitative) and **student** (qualitative).

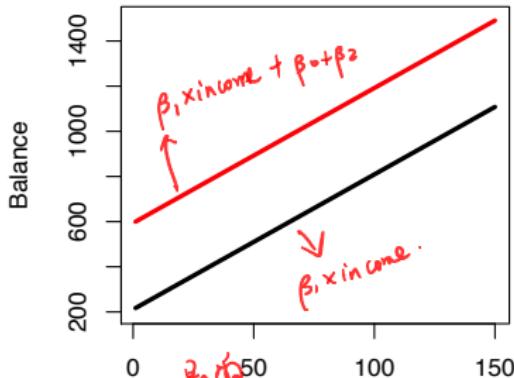
Without an interaction term, the model takes the form

$$\begin{aligned}\text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 & \text{if } i\text{th person is a student} \\ 0 & \text{if } i\text{th person is not a student} \end{cases} \\ &= \beta_1 \times \text{income}_i + \begin{cases} \beta_0 + \beta_2 & \text{if } i\text{th person is a student} \\ \beta_0 & \text{if } i\text{th person is not a student.} \end{cases}\end{aligned}$$

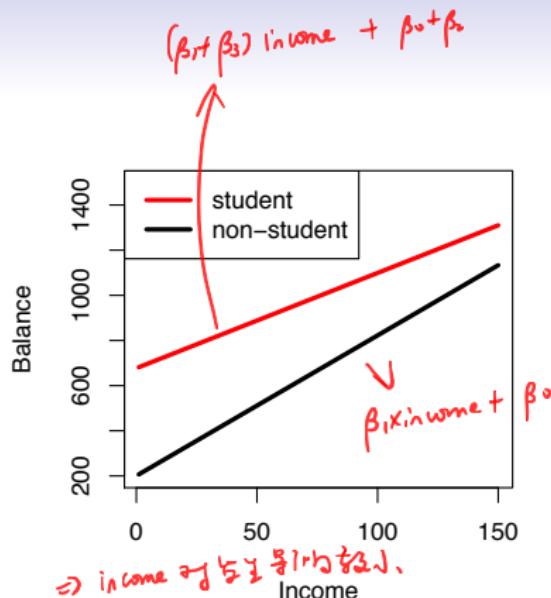
$$\beta_0 + \beta_1 \times \text{income} + (\beta_2)_0 + (\beta_2)_1 \times \text{income}$$

With interactions, it takes the form

$$\begin{aligned}\text{balance}_i &\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 + \beta_3 \times \text{income}_i & \text{if student} \\ 0 & \text{if not student} \end{cases} \\ &= \begin{cases} (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \times \text{income}_i & \text{if student} \\ \beta_0 + \beta_1 \times \text{income}_i & \text{if not student} \end{cases}\end{aligned}$$



$\Rightarrow$  由高收入者比低收入者  
大學生的存款額也較高



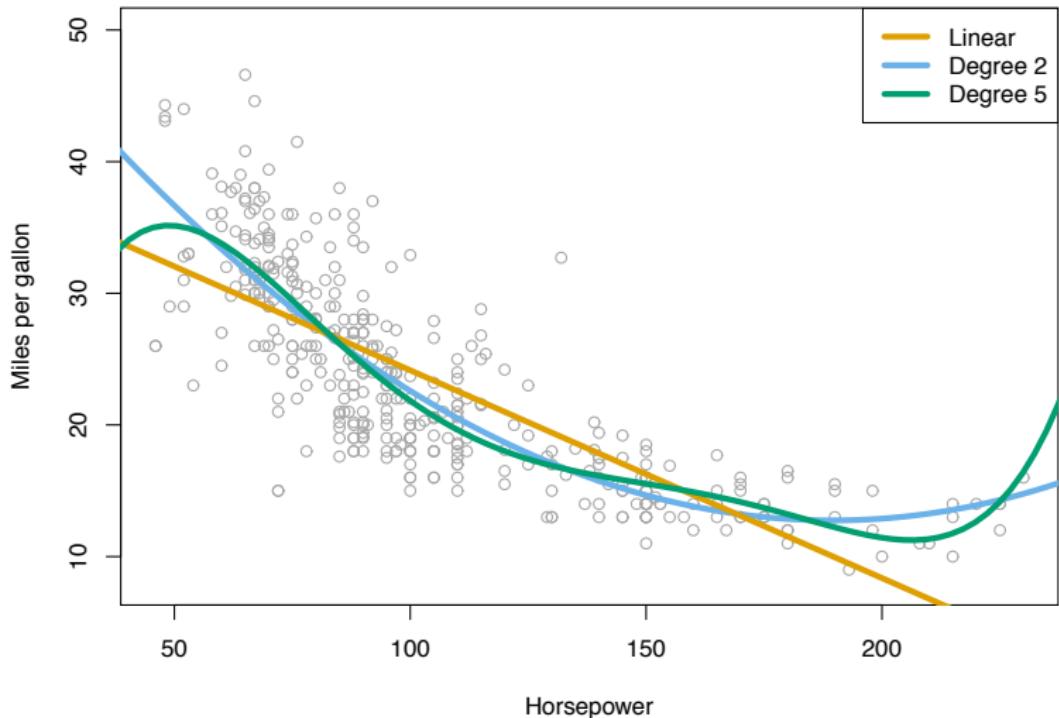
$\Rightarrow$  income 與 student 有交互作用。

Credit data; Left: no interaction between **income** and **student**.  
 Right: with an interaction term between **income** and **student**.

## Non-linear effects of predictors

多项式回归  
polynomial regression on Auto data

拟合用CV选



The figure suggests that

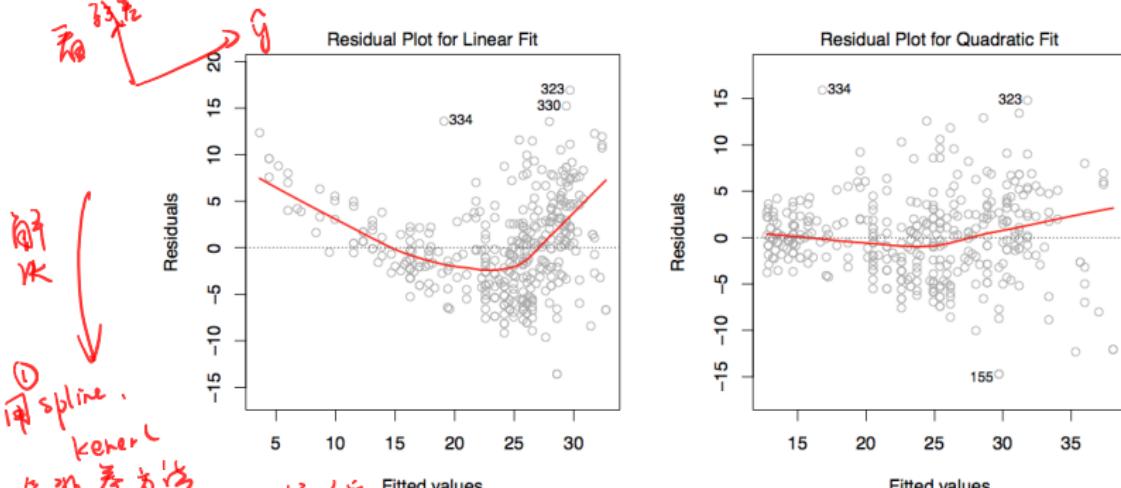
$$\text{mpg} = \beta_0 + \beta_1 \times \text{horsepower} + \beta_2 \times \text{horsepower}^2 + \epsilon$$

may provide a better fit.

	Coefficient	Std. Error	t-statistic	p-value
Intercept	56.9001	1.8004	31.6	< 0.0001
horsepower	-0.4662	0.0311	-15.0	< 0.0001
horsepower <sup>2</sup>	0.0012	0.0001	10.1	< 0.0001

自变量较少  
2.1.1 判斷有無非線性  
若自變量較少時  $y \uparrow$

## Potential Problems-nonlinearity



解法

① spline

kernel

多項式方法  
② 增加 項數 / 增加自由度

**FIGURE 3.9.** 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 linear regression of `mpg` on `horsepower` and `horsepower2`. There is little pattern in the residuals.

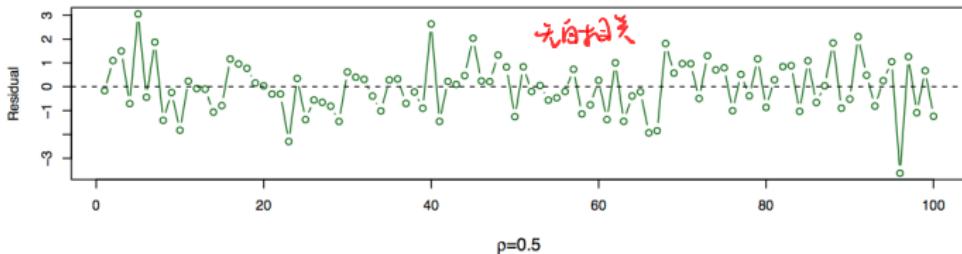
残差自相关

## Potential Problems-correlation of error term

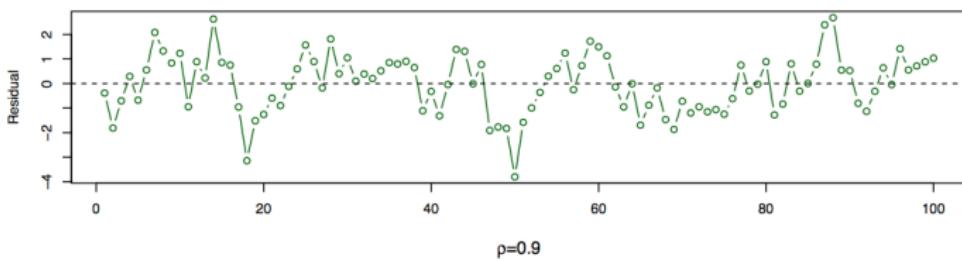
方差恒定  
GARCH < 方差破 HARCH

用时间窗

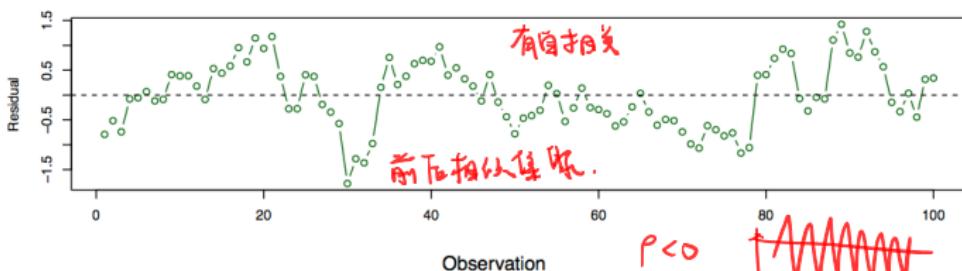
$p=0.0$



$p=0.5$



$p=0.9$



有自相关  
前后相似集聚。

$p<0$



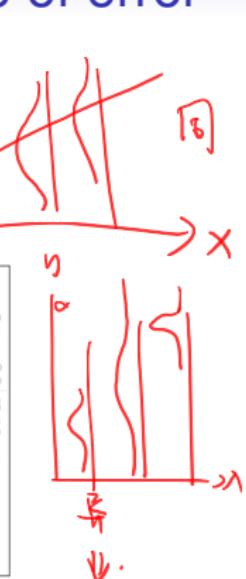
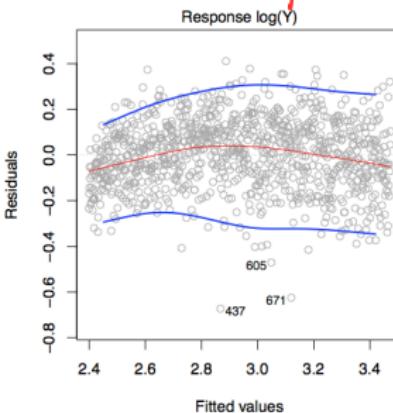
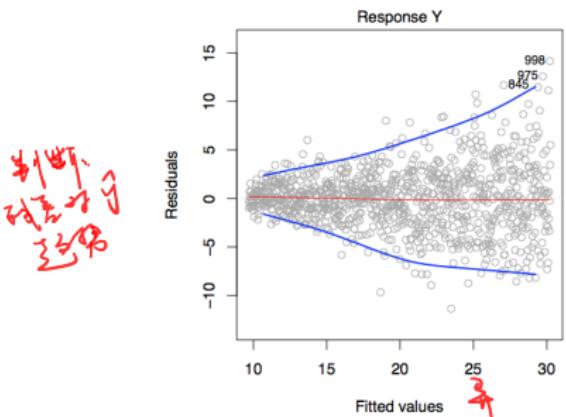
且  $t \uparrow \Rightarrow p \downarrow$   
因果性  
(推验不可信)

估价失准  
过低估  
置信区间  
变窄  
95% 区间  
包含真值  
参数的概  
率  $< 95\%$

## Potential Problems-non-constant variance of error terms

$\frac{1}{f(x_i)}$  反样本 (非先知性权重)

- $Var(\epsilon) \neq \sigma^2$  异方差
- heteroscedasticity – transform or WLS

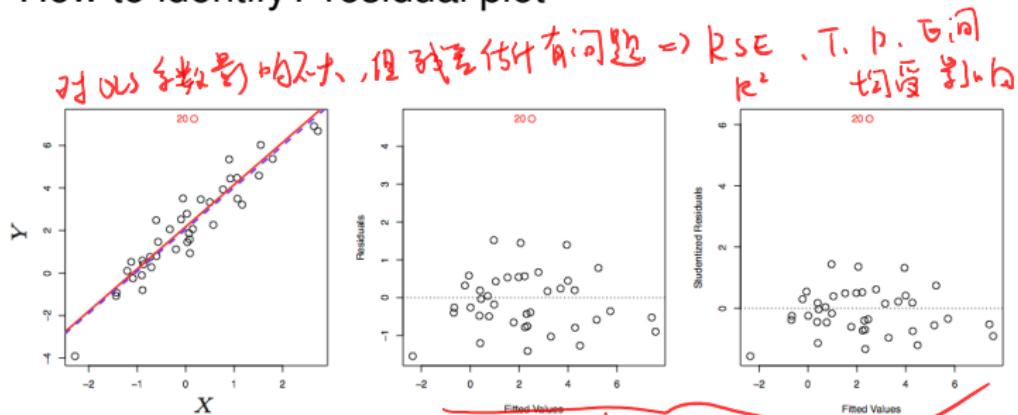


**FIGURE 3.11.** Residual plots. In each plot, the red line is a smooth fit to the residuals, intended to make it easier to identify a trend. The blue lines track the outer quantiles of the residuals, and emphasize patterns. Left: The funnel shape indicates heteroscedasticity. Right: The predictor has been log-transformed, and there is now no evidence of heteroscedasticity.

仍先偏，  
但非异方  
解决区间 T, F, P 问题

# Potential Problems-Outlier

- RSE 1.09 → 0.77 ( $R^2$  89.2% → 80.5%) removed
- How to identify? residual plot



**FIGURE 3.12.** 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  $3$ .

若已知为数据采集错误  $\Rightarrow$  删掉

若真实数据本身存在 outlier  $\Rightarrow$  更多方法

$$\text{学生化残差} < \frac{e}{RSE}$$

一般  $> 3$  有 outlier.

高杠杆

## Potential Problems-High leverage Points

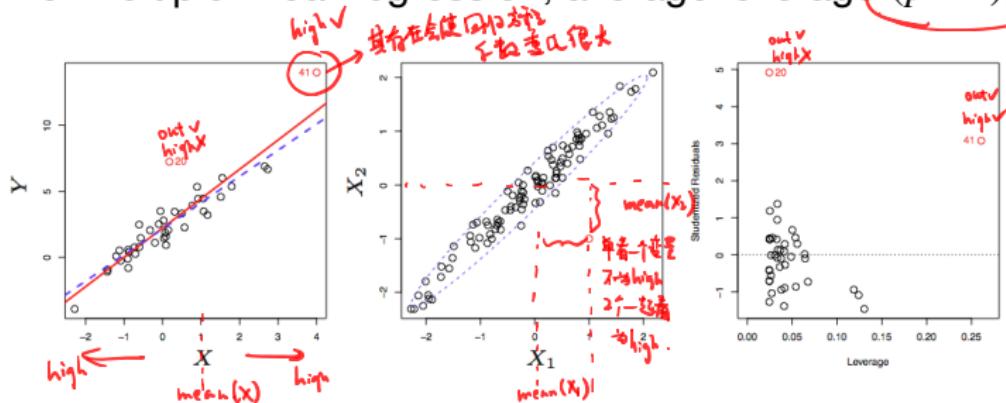
- leverage statistics. For simple linear regression

杠杆统计量  $h_i$

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

点越靠近均值，  
h↓  
正态  
h↑

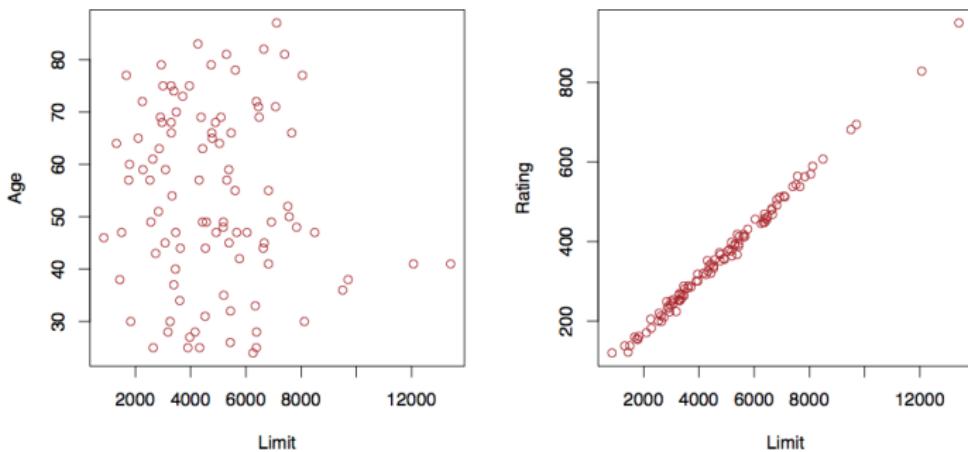
- For multiple linear regression, average leverage  $(p + 1)/n$



**FIGURE 3.13.** 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.

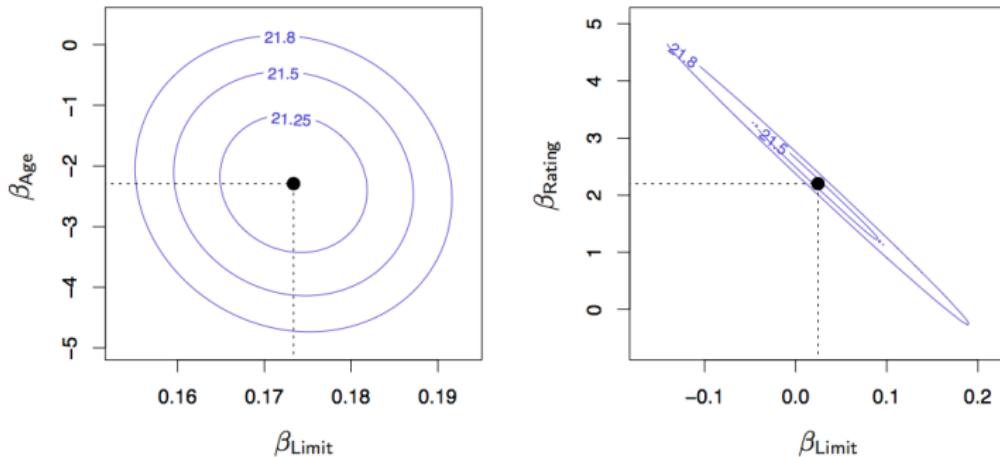
共线性.

## Potential Problems-Collinearity



**FIGURE 3.14.** 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.

## Potential Problems-Collinearity



**FIGURE 3.15.** Contour plots for the RSS values as a function of the parameters  $\beta$  for various regressions involving the Credit data set. In each plot, the black dots represent the coefficient values corresponding to the minimum RSS. Left: A contour plot of RSS for the regression of `balance` onto `age` and `limit`. The minimum value is well defined. Right: A contour plot of RSS for the regression of `balance` onto `rating` and `limit`. Because of the collinearity, there are many pairs  $(\beta_{\text{Limit}}, \beta_{\text{Rating}})$  with a similar value for RSS.

## Potential Problems-Collinearity

- Consequence of collinearity: standard error for  $\hat{\beta}_j$  increase, decline tstatistic, reduce the power of the hypothesis test
- How to detect? Correlation matrix. Variance inflation factor(VIF)

$$VIF(\hat{\beta}_j) = \frac{1}{1 - R^2_{X_j | X_{-j}}}$$

where  $R^2_{X_j | X_{-j}}$  is the  $R^2$  from a regression of  $X_j$  onto all of the other predictors. ( $VIF \geq 10$ )

		Coefficient	Std. Error	t-statistic	p-value
Model 1	Intercept	-173.411	43.828	-3.957	< 0.0001
	age	-2.292	0.672	-3.407	0.0007
	limit	0.173	0.005	34.496	< 0.0001
Model 2	Intercept	-377.537	45.254	-8.343	< 0.0001
	rating	2.202	0.952	2.312	0.0213
	limit	0.025	0.064	0.384	0.7012

## Potential Problems-Endogeneity

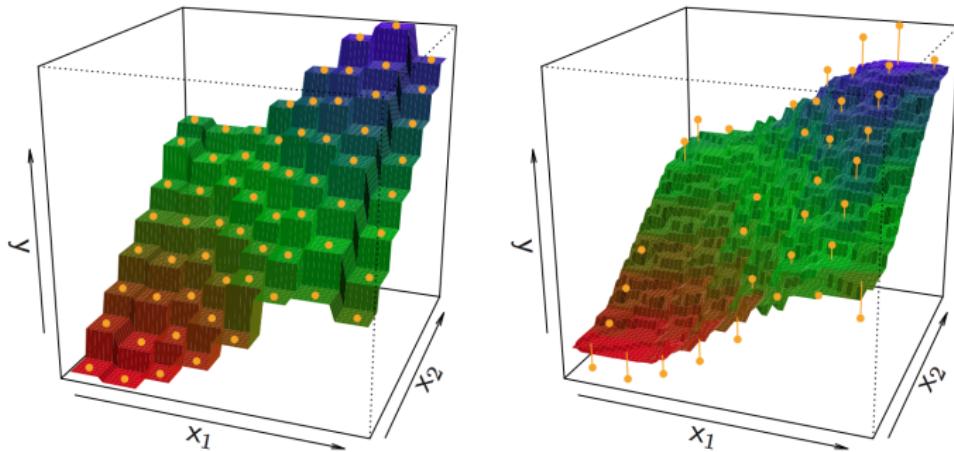
- $E(X\epsilon) \neq 0$ , Endogeneity 内生性
- Inconsistent estimator
- Instrument Variable

## Linear regression VS KNN regression

- linear regression (parametric method) has some advantages: easy to fit, easy to interpret; it also has some advantage: strong assumption about the form of  $f(X)$ , specification error?
- Nonparametric methods don't assume a parametric form for  $f(X)$ .
- The simplest nonparametric method, K-nearest neighbors regression (KNN regression).
- Given a value for  $K$  and a prediction point  $x_0$ , KNN regression first identifies the  $K$  training observations that are closest to  $x_0$ , represented by  $\mathcal{N}_0$ .

$$\hat{f}(x_0) = \frac{1}{K} \sum_{x_i \in \mathcal{N}_0} y_i$$

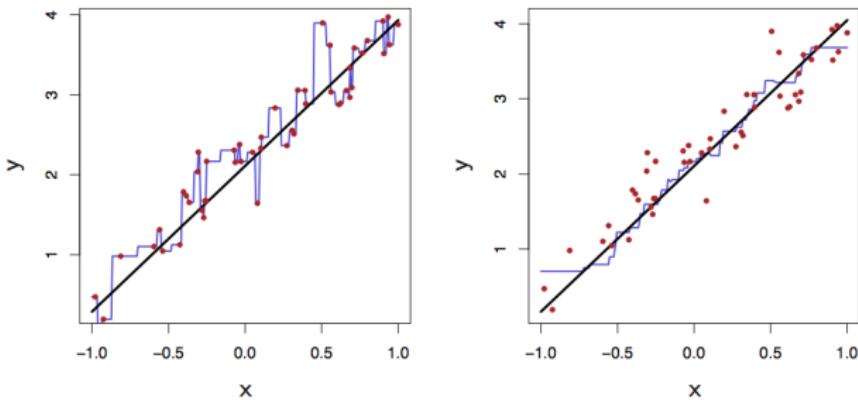
# Linear regression VS KNN regression



**FIGURE 3.16.** Plots of  $\hat{f}(X)$  using KNN regression on a two-dimensional data set with 64 observations (orange dots). Left:  $K = 1$  results in a rough step function fit. Right:  $K = 9$  produces a much smoother fit.

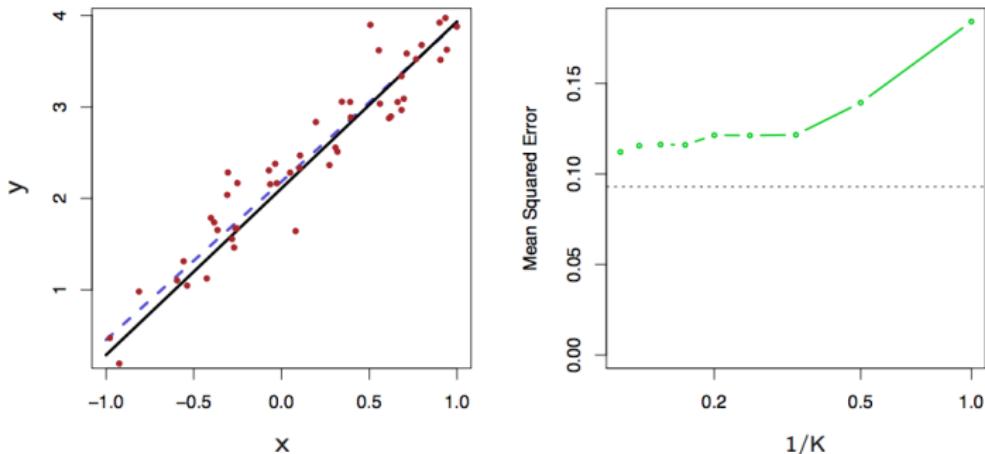
## Linear regression VS KNN regression

- The parametric approach will outperform the non-parametric approach if the parametric form that has been selected is close to the true form of  $f$ .



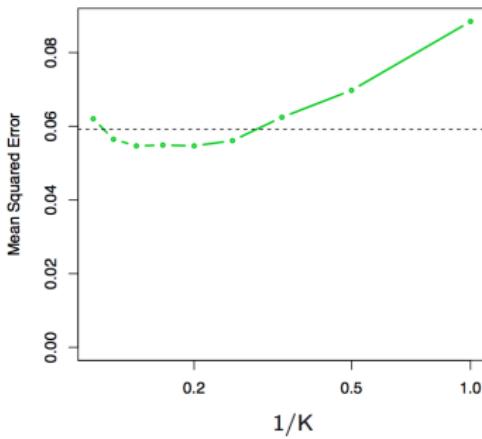
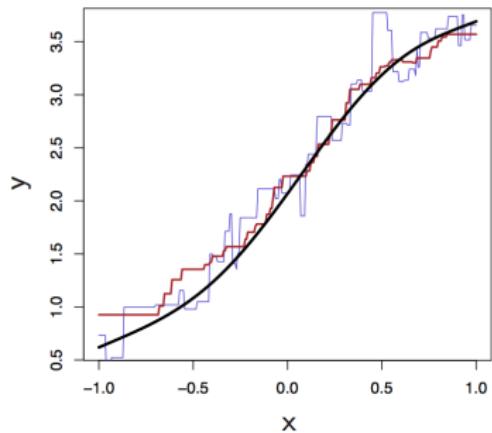
**FIGURE 3.17.** Plots of  $\hat{f}(X)$  using KNN regression on a one-dimensional data set with 100 observations. The true relationship is given by the black solid line. Left: The blue curve corresponds to  $K = 1$  and interpolates (i.e. passes directly through) the training data. Right: The blue curve corresponds to  $K = 9$ , and represents a smoother fit.

## Linear regression VS KNN regression

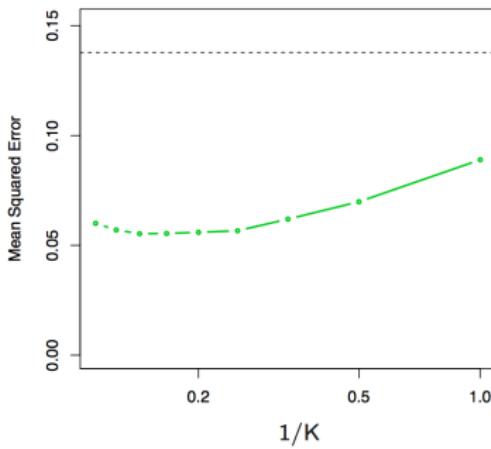
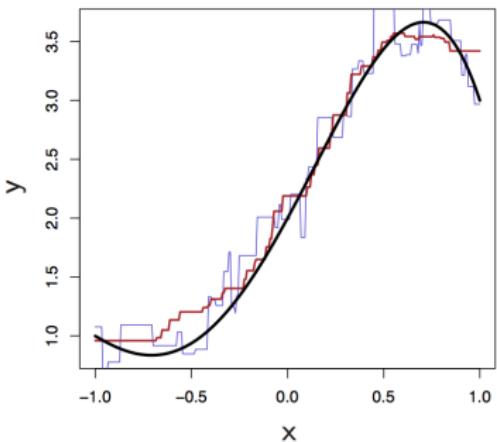


**FIGURE 3.18.** The same data set shown in Figure 3.17 is investigated further. Left: The blue dashed line is the least squares fit to the data. Since  $f(X)$  is in fact linear (displayed as the black line), the least squares regression line provides a very good estimate of  $f(X)$ . Right: The dashed horizontal line represents the least squares test set MSE, while the green solid line corresponds to the MSE for KNN as a function of  $1/K$  (on the log scale). Linear regression achieves a lower test MSE than does KNN regression, since  $f(X)$  is in fact linear. For KNN regression, the best results occur with a very large value of  $K$ , corresponding to a small value of  $1/K$ .

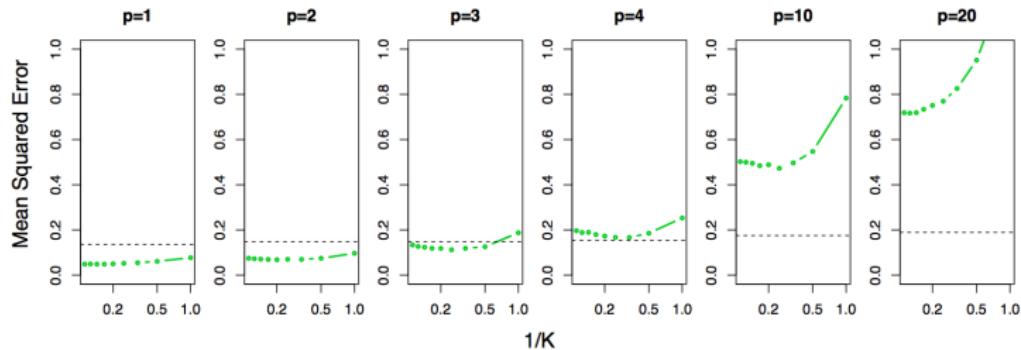
# Linear regression VS KNN regression



# Linear regression VS KNN regression



# Linear regression VS KNN regression



# Generalizations of the Linear Model

In much of the rest of this course, we discuss methods that expand the scope of linear models and how they are fit:

- *Classification problems*: logistic regression, support vector machines
- *Non-linearity*: kernel smoothing, splines and generalized additive models; nearest neighbor methods.
- *Interactions*: Tree-based methods, bagging, random forests and boosting (these also capture non-linearities)
- *Regularized fitting*: Ridge regression and lasso