

ANOVA:

Analysis of variance.

main idea:

figure out how much of the total variance comes from:

The variance between groupswithin groups.

linear regression assumption:

- error is mean 为 0
- 对所有 X , Σ 方差为 constant
↓
规则
- 误差项 ε 是一个服从正态分布的 R.V.
且相互独立, $\varepsilon \sim N(0, \sigma^2)$

线性: 因变量与自变量之间为线性关系.

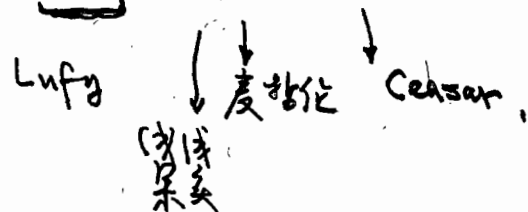
正态: 因变量为正态性.

残差: 独立同分布.

正交: 残差与自变量不相关.
期望为 0

0 同残独立 M

变量线性关系 Master.



Wikipedia 上答案:

1. Weak exogeneity:
 i.e. predictor variable x can be treated as fixed values, rather than r.v.
2. Linearity:
 response variable is a linear combination of the parameters (regression coefficients)
3. constant variance,
4. error 之间互相独立.
5. predictor 之间 No multicollinearity

The common tests arranged:

1 sample proportions	1 sample t	
2 sample proportions	2 Sample t	
	2 Sample t paired t	Correlation/ Regression
	One-way ANOVA	

categorical data
分类即可.

Quantitative data.

Q2 How many samples do you have?

↳ or how many groups do you have.

↳ historically, 你通常有两个.

test group: 做了实验 test 的

Control group: 用于比较, 没做实验 test 的.

stats

3

161017

One sample:

↳ Usually I means I make a comparison against a historical or global value.

Two samples:

↳ Traditional control group vs. test group in an experiment.

↳ comparing one group to the other group.

stats.

461017

1

Confidence intervals.

$\bar{X} \pm t \frac{s}{\sqrt{n}}$ - 总体方差未知, 用 t -statistics.

$\bar{X} \pm z \frac{\sigma}{\sqrt{n}}$ 总体方差已知, 用 z -statistics.

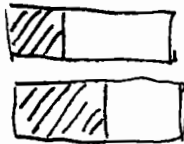
choosing which stat to use:

means, proportions, relationships

stats
16/01/7
2.

Test for proportion

Difference of two proportions



chi-sq test for independence

Difference of two means
(independent samples)

Regression analysis

Test for a mean

Difference of two means
paired

~~Less~~ Quantitative research,
most likely, you perform some hypothesis test in the end.

→ different test? ~~what~~ collected is data first.

Many possibilities:

- Estimate population proportion
- Estimate — mean
- One sample proportion
- Two —

- One sample t (mean)
- Two sample t (mean)
- paired t
- correlation/regression analysis.

- One way ANOVA
- Two way ANOVA
- Chi square test
- One sample variance
- Two sample —

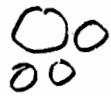
One sample: against hypothesized value.

- mean 250g
- 82% export quality.



Stats
161017
3.

Two samples:



A



B

compared with each other:

- Heavier?
- No difference?
- More export quality?

One sample,
two measurements.



- circumference

- weight.

or: / \
color grade.

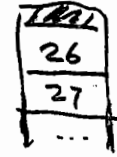
Stats
161017
4

Test
for proportion



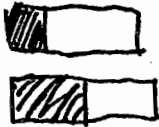
One Sample

Test for a mean



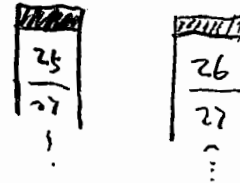
全体的平均値...

Difference of
two proportions.



Two samples.

Differences of two



Two sets of information
on the sample

Chi-sq test
for independence

Regression analysis.

Difference of
two means,
paired.

to how many nuts chocolate sold vs. weather

or male white chocolate
female milk chocolate

お砂糖の量とpreferの割合

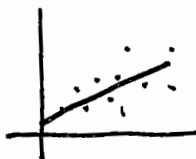
chi-sq & regression

都是用来说 look for ~~regression~~ relationships.

	D	M	W
men			
women			

数据是个 table

regression, 数据是个散点.
scatter plot.



purpose:

- test against a hypothesis value
- comparing two statistics.
- looking for a Relationship.

stats
5

$$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$$

但是

$\text{Cov}(X, Y)$ 不好用.

因为 X, Y 是有 units 的.

1	-1×10^6	} 取值一样. 但系数 可以差很多!!
$(X - \mu_X)$	$(Y - \mu_Y)$	
100	-1	
$(X - \mu_X)$	$(Y - \mu_Y)$	

所以人更常用是 correlation,

normalize it, to make it unitless.

$$\text{Corr}(X, Y) := \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}}$$

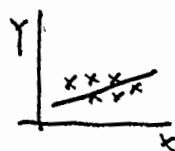


散点图说明 X, Y 正相关. $\text{Cov}(X, Y) > 0$

于是 if $(X > \mu_X)$, 那么得证 $\text{Cov}(X, Y) > 0$.

须 $Y > \mu_Y$

正相关 $\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$



$X > \mu_X$	+	+	Cov +
$X < \mu_X$	-	-	+

X, Y 同涨或同跌 ($\text{Cov} > 0$)

负相关

$\text{Cov}(X, Y) = E[(X - \mu_X)(Y - \mu_Y)]$



$X > \mu_X$	+	-	Cov -
$X < \mu_X$	-	+	-

X 涨, Y 跌 ($\text{Cov} < 0$)

Confidence intervals are based on sample data,
give a range of plausible value for a parameter.

We never construct confidence intervals
for statistics (like \bar{x})

for parameters.

→ $\frac{p}{2}$ common & choice.

We may be 95% confident that μ lies
in the interval $(-0.2, 3.1)$

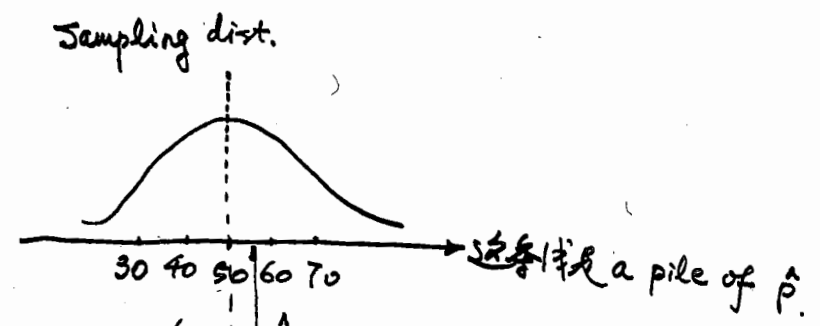
50% students like iphone.
you sample 25.

$$\mu_{\hat{p}} = p$$

mean of the sample distribution \hat{p} .

$$\sigma_{\hat{p}} = \sqrt{\frac{pq}{n}} = \sqrt{\frac{.5(.5)}{25}} = .10$$

std of the sample dist.



你得到一个 \hat{p}
往上, 往下两个 σ
得到一个置信区间.

I don't know
what my true p is,
all I know is my \hat{p}

from sampling

真实的 p 落在此区间
中的概率是 95%

\hat{p} 是 statistics.
由 sampling 得到的
sampling 取值的
 \hat{p} 分布满足 sampling
distribution.
是个正态分布.

cheby:

描述的是: ~~差异性~~.

差或远离均值, 可能性减小.

↓ 这个可能性由方差衡量

↓ 衡量时远离均值的程度.

$$P(|X - \mu| \geq a) \leq \frac{\text{Var}(X)}{a^2}$$

什么叫 Inference?

有 10,000 个苹果, 求平均重量.

你送了 100 个
Sampling

抽样

总体, 即称为 Inference.

由样本推总体.

s_1 / s_2 / s_3
150g 130g 140g

—— 总是有 sampling error 的!

- confidence interval depends on 什么?

- Variance

- sample size.

Small samples
vary more from each other,
and have less info.

↓ leads to
wider confidence intervals.

Population with low
variation
↓
samples w. low
variation
↓
narrow confidence
interval.

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

→ $\sigma_{\bar{x}}$
↓
n.v.

由中心极限定理可得!!

Confidence Interval for mean

$$\bar{X} \pm t \frac{s}{\sqrt{n}}$$

\bar{X} : sample mean
 t : 由置信度或显著性水平 α 所写之值, 查表可得
 s : sample std.
 n : sample size

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

	Sample size				
t	10	15	20	30	100
90%	1.833	1.761	1.729		
95%	2.262	2.145			
99%					

$\frac{s}{\sqrt{n}}$ gives us standard error.

为什么是 \sqrt{n} ?

reflects the more ^{sample} information we have, the less information we get from new samples.

We're 90% confident



95%



99%



$t \frac{s}{\sqrt{n}}$: margin error.

要求越精确, Confidence interval 越大.

总体是一个随机变量 $X \sim F(x)$

样本: (x_1, \dots, x_n)

实际上是几维随机变量

→ n 维之联合分布.
若为正态, 则分量必为正态.

→ 加减操作亦为正态.

- X, Y 各自正态 $\xrightarrow{+}$ (X, Y) 正态
可以保证.

- X, Y 独立正态 $\longrightarrow X+Y$ 正态
独立很重要, 加了独立, 怎么折腾都行.
→ (X, Y) 正态 $N(\dots, 0)$
= 独立正态.
→ $\rho=0$! 这是独立无相关性证明.

- 简单随机样本.
→ i.i.d.

X : 总体.
 x_1, x_2, \dots, x_n

① X_i 独立
② $X_i \sim F(x)$

独立同分布 (i.i.d.)

样本之分布函数

$$F(x_1, x_2, \dots, x_n) = P(x_1 \leq x_1, x_2 \leq x_2, \dots, x_n \leq x_n) \quad \text{独立}$$

$$= P(x_1 \leq x_1) \cdots P(x_n \leq x_n)$$

↓

几个随机变量

$$= F(x_1) \cdots F(x_n) = \prod_{i=1}^n F(x_i)$$

样本之分布函数

即把总体之分布函数

拿来, x 换成 x_i , 连乘起来.

类似的,

样本之概率密度

$$X \sim f(x)$$

$$(x_1, x_2, \dots, x_n) \sim \prod_{i=1}^n f(x_i)$$

简单随机样本

样本之概率密度函数

样本之数字特征

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{X})^2$$

↓ 定义样本之方差.
为 $\frac{1}{n-1}$

样本之均值是几个随机变量求平均

$$E\bar{X} = \frac{1}{n} \sum_{i=1}^n E x_i = \mu$$

→ 就是总体之期望.

$$E S^2$$

$$D S^2$$

$$D\bar{X} =$$

$\bar{X}, \bar{S}^2, E\bar{X}, D\bar{X}, E\bar{S}^2, D\bar{S}^2$ 统计量.

e.g. $f(x) = \begin{cases} 2x & 0 < x < 1 \\ 0 & \text{其他} \end{cases}$

$X_{(4)} = \max(X_1, X_2, X_3, X_4)$ X_1, \dots, X_4 独立.

求 $f_{X_{(4)}}(x)$

$P(X_1 \leq x) P(X_2 \leq x) P(X_3 \leq x) P(X_4 \leq x)$
|||

Soln.

* $F_{X_{(4)}}(x) = P(X_4 \leq x) = P(\max(X_1, X_2, X_3, X_4) \leq x)$
 $= [F(x)]^4$

于是先求 $F(x)$.

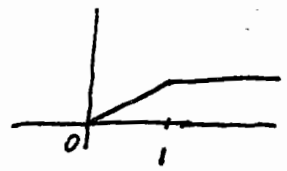
$F(x) = \begin{cases} 0 & x < 0 \\ x^2 & 0 \leq x < 1 \\ 1 & 1 \leq x \end{cases}$

$F_{X_{(4)}}(x) = F^4(x) = \begin{cases} 0 & x < 0 \\ x^8 & 0 \leq x < 1 \\ 1 & 1 \leq x \end{cases}$

$f_{X_{(4)}}(x) = \begin{cases} 8x^7 & 0 \leq x < 1 \\ 0 & \text{其他} \end{cases}$

→ $X \sim U(0,1)$

分布函数:



$\bar{X}, S^2, E\bar{X}, D\bar{X}, ES^2, DS^2$

→ χ^2, t, F

1). Chi

$\chi^2 = X_1^2 + X_2^2 + \dots + X_n^2$

$\sim \chi^2(n)$

$f(x) =$
不必记.

$E(\chi^2) = n,$

$D(\chi^2) = 2n.$

X_i 均服从.

① $X_i \sim N(0,1)$ 标准正态

② X_i 独立

当然可以 $\int x f(x)$

Expectation 期望.

$$\begin{aligned} E\chi^2 &= E(X_1^2 + \dots + X_n^2) = EX_1^2 + \dots + EX_n^2 \\ &= DX_1 + (EX_1)^2 + \dots \\ &= n \end{aligned}$$

↑
典型模式

X_i 标准正态.

分布:

$$T = \frac{X}{\sqrt{Y/n}} \sim t(n)$$

$$\begin{cases} X: \text{标准正态.} \\ Y \sim \chi^2(n) \\ X, Y \text{ 独立.} \end{cases}$$

分布最重要, 性质不必掌握

↓ 但是个偏态. α 分位点对称.

$$F \text{ 分布: } F = \frac{X/n_1}{Y/n_2} \sim F(n_1, n_2)$$

$$\begin{cases} \textcircled{1} X \sim \chi^2(n_1) \\ \textcircled{2} Y \sim \chi^2(n_2) \\ \textcircled{3} X, Y \text{ 独立.} \end{cases}$$

$\frac{1}{F}$ 亦为 F 分布

$$\frac{1}{F} \sim F(n_2, n_1)$$

χ^2

$$X_1^2 + X_2^2 + \dots + X_n^2$$

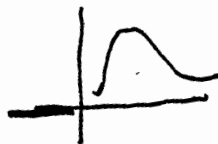
$$X_i \sim N(0, 1)$$

独立

$$EX^2 = n$$

$$DX^2 = 2n$$

分布和
取值范围



t

$$\frac{X}{\sqrt{Y/n}}$$

$$X \sim N(0, 1)$$

$$Y \sim \chi^2(n)$$

X, Y 独立

$f(x)$

偏态



F

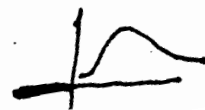
$$\frac{X/n_1}{Y/n_2}$$

$$X \sim \chi^2(n_1)$$

$$Y \sim \chi^2(n_2)$$

X, Y 独立.

$$\frac{1}{F} \sim F(n_2, n_1)$$



正态总体. $X \sim N(\mu, \sigma^2)$ 总体均值

样本 x_1, x_2, \dots, x_n , i.i.d $\sim N(\mu, \sigma^2)$

联合起来是几维正态.

4.

$$\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$$

sample mean 不是 n 维正态了!

2. \bar{X} 与 σ^2 独立.

$$3. T = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim t(n-1)$$

两个正态总体. 类似.

$$\sum_{i=1}^n (x_i - \bar{x})^2 \text{ 自由度只有 } n-1 \text{ 个!!}$$

$$\text{因为 } \sum_{i=1}^n (x_i - \bar{x}) = 0$$

有一个约束! 看似 n 个, 实际是 $n-1$ 个自由度

$$\frac{(n-1)\sigma^2}{\sigma^2} \sim \chi^2(n-1)$$

$$(n-1)\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\frac{1}{n-1}} = \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\frac{(n-1)\sigma^2}{\sigma^2} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2}$$

x_i : 正态.

$\bar{x} = \frac{1}{n}(x_1 + \dots + x_n)$ 正态.

σ : 常数

$$E\left(\frac{x_i - \bar{x}}{\sigma}\right) = \frac{E x_i - E \bar{x}}{\sigma} = 0$$

$$D\left(\frac{(n-1)\sigma^2}{\sigma^2}\right) = 2(n-1)$$

常数提出来, 要平方.

$$\frac{(n-1)^2}{\sigma^4} D(\sigma^2) = 2(n-1) \Rightarrow D\sigma^2 = \frac{2\sigma^4}{n-1}$$

$$X \sim B(1, p)$$

X 服从 0, 1 分布.

即一次 Bernoulli 试验

一个简单随机样本, X_1, X_2, \dots, X_n .

求 X_1, X_2, \dots, X_n 的分布律.

X	0	1
	q	p

Soln: 难点在于如何把离散分布列写成分布函数
分布列不太好看.

$$P(X=k) = C_1^k p^k q^{1-k}$$

$$k=0, 1 \quad C_1^k = C_1^0 = 1$$

$$C_1^1 = 1$$

$$= p^k q^{1-k}$$

分布函数.

$$P(x) = \begin{cases} p^x q^{1-x} & x=0, 1 \\ 0 & \text{其他} \end{cases}$$

独立同分布 (把 x 换成 x_i)

$$P(X_1, \dots, X_n) = \begin{cases} p^{x_1 + \dots + x_n} q^{n - x_1 - \dots - x_n} & x_i = 0, 1 \\ 0 & \text{其他} \end{cases}$$

$$= \begin{cases} p^{n\bar{x}} q^{n(1-\bar{x})} & x_i = 0, 1 \\ 0 & \text{其他} \end{cases}$$

2. $X \sim P(\lambda)$ \bar{X} 分布律.

~~解~~

Soln:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

↪ X_i 样本! 必独立.

$X_1 + X_2 \sim P(2\lambda)$ 对于泊松分布, 有可加性.

$$\sum_{i=1}^n X_i \sim P(n\lambda)$$

$$\text{即 } P\left(\sum_{i=1}^n X_i = k\right) = \frac{(n\lambda)^k}{k!} e^{-n\lambda} \quad k=0, 1, 2, \dots$$

于是

$$P\left(\frac{1}{n} \sum_{i=1}^n X_i = \frac{k}{n}\right) = P\left(\bar{X} = \frac{k}{n}\right) = \frac{\lambda^k}{k!} e^{-\lambda}$$

不再是泊松分布

$X_1, X_2, X_3, X_4 \sim N(0, 2^2)$ 样本
(简单随机)

$$\chi^2 = a(X_1 - 2X_2)^2 + b(3X_3 - 4X_4)^2 \sim \chi^2(n)$$

Soln:

- 随机变量函数分布. 用什么? 定义法可以但巨麻烦

- χ^2 用典型模式!!

$$X_1 - 2X_2 \sim N(0, 20)$$

$$\begin{aligned} D(X_1 - 2X_2) &= D(X_1) + 4D(X_2) = 56^2 \\ &\quad \downarrow \\ &\text{非标准正态, 于是标准化.} \end{aligned}$$

$$\frac{X_1 - 2X_2}{\sqrt{20}} \sim N(0, 1)$$

于是令 $a = \frac{1}{20}$

$$\left(\frac{X_1 - 2X_2}{\sqrt{20}} \right)^2 \text{ 是一个标准正态之平方.}$$

于是 $n = 2$

可能 a 或 $b = 0$. 于是 $n = 1$.

4. $T \sim t(n)$

则 T^2 服从何分布?

Soln:

$T \sim t(n)$

$T = \frac{X}{\sqrt{Y/n}} \quad , \quad T^2 = \frac{X^2}{Y/n} \Rightarrow F(1, n)$

$\nearrow \chi^2(1)$

$\nwarrow \chi^2(n)$

看 χ^2, t, F , 不要用反义法!

→ 因为不记得 $f(x)$.

$X_i \sim N(0, 1)$ 来自标准正态分布的简单随机样本

问 1. $n\bar{X} \sim N(0, 1)$?

1. $\left\{ \begin{array}{l} - n\bar{X} = \sum_{i=1}^n X_i \sim N(0, n) \quad \leftarrow \text{方差相加} \\ - \bar{X} \sim N(0, \frac{1}{n}) \\ n\bar{X} \sim N(0, n) \end{array} \right.$

\rightarrow 方差相加, 要乘以 n^2

2. $nS^2 \sim \chi^2(n)$

2. $S^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$

$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1) \quad \checkmark$

$nS^2 \sim \chi^2(n) \quad \times$

\downarrow S^2 中隐含 n .

3. $\frac{(n-1)\bar{X}}{S} \sim t(n-1)$

3. $\sigma^2 = 1$

$(n-1)S^2 \sim \chi^2(n-1)$

$\frac{\bar{X}}{\sqrt{Y/n}}$

$\frac{(n-1)^2 \bar{X}}{(n-1)S^2} = \frac{\bar{X}}{\frac{(n-1)S^2}{(n-1)^2}}$

不对.

$X \sim N(0, \sigma^2)$,
 X_1, \dots, X_9
 9个样本.

求使

$P(1 < \bar{X} < 3)$ 最大.

soln:

$$\bar{X} = \frac{1}{n} \sum X_i \sim N\left(0, \frac{\sigma^2}{9}\right)$$

$\downarrow n=9$

于是标准化

$$\frac{\bar{X} - 0}{\sigma/3} = \frac{3\bar{X}}{\sigma}$$

$$P\left(\frac{3}{\sigma} < \frac{3\bar{X}}{\sigma} \leq \frac{9}{\sigma}\right) = \Phi\left(\frac{9}{\sigma}\right) - \Phi\left(\frac{3}{\sigma}\right) = P(\xi)$$

\downarrow 求导!

$$P'(\xi) = \varphi\left(\frac{9}{\sigma}\right) \left(-\frac{9}{\sigma^2}\right) - \varphi\left(\frac{3}{\sigma}\right) \left(-\frac{3}{\sigma^2}\right) = 0$$

$$\varphi\left(\frac{9}{\sigma}\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{81}{2\sigma^2}}$$

$$\therefore \sigma = \frac{6}{\sqrt{\ln 3}}$$

X_1, X_2, X_3 独立 $N(0, \sigma^2)$

证明: $\sqrt{\frac{2}{3}} \frac{X_1 + X_2 + X_3}{|X_2 - X_3|} \sim t(1)$

Soln: 典型模式 $T = \frac{X}{\sqrt{Y/1}}$

- ① $X \sim N(0, 1)$
- ② $Y \sim \chi^2(1)$
- ③ X, Y 独立

$\frac{1}{\sqrt{3}}(X_1 + X_2 + X_3) \sim N(0, \sigma^2)$

$\frac{1}{\sqrt{3}}(X_1 + X_2 + X_3) \sim N(0, \sigma^2)$

标准化-T

$\frac{X_1 + X_2 + X_3}{\sqrt{3}\sigma} \sim N(0, 1)$

独立证.

$X_1 + X_2 + X_3$ 与 $X_2 - X_3$ 独立.

即证.

$X_2 + X_3$ 与 $X_2 - X_3$ 独立.

X_1, X_2, X_3 独立, 则

$X_2 + X_3$ 与 $X_2 - X_3$ 均正交.

独立正交情况下, 独立 \Rightarrow 相关系数=0

$\text{COV}(X_2 + X_3, X_2 - X_3)$

$= \text{COV}(X_2, X_2) - \text{COV}(X_3, X_2) + \text{COV}(X_2, X_3) - \text{COV}(X_3, X_3)$

$Y = \frac{|X_2 - X_3|}{\sqrt{2}\sigma}$

$Y^2 = \frac{(X_2 - X_3)^2}{2\sigma^2}$

$= \left(\frac{X_2 - X_3}{\sqrt{2}\sigma} \right)^2$

↓ 标准化

$X_2 - X_3 \sim N(0, 2\sigma^2)$

$\frac{X_2 - X_3}{\sqrt{2}\sigma} \sim N(0, 1)$

$Y^2 \sim \chi^2(1) \Rightarrow T = \frac{X}{\sqrt{Y}} \sim t(1)$

$$E S^2$$

$$X \sim N(\mu, \sigma^2)$$

正态总体.

样本: X_1, X_2, \dots, X_{2n} .

$$\bar{X} = \frac{1}{2n} \sum_{i=1}^{2n} X_i$$

$$Y = \sum_{i=1}^n (X_i + X_{n+i} - 2\bar{X})^2 \text{ 求 } EY.$$

Soln.

$$EY = \sum_{i=1}^n (EX_i + EX_{n+i} - 2\bar{X})^2$$

新样本 $(X_1 + X_{n+1}), (X_2 + X_{n+2}), \dots, (X_n + X_{2n})$

n 个分量.

这个样本取自 $N(2\mu, 2\sigma^2)$

$$\text{均值: } \frac{1}{n} \sum_{i=1}^n (X_i + X_{n+i})$$

$$= \frac{1}{n} \sum_{i=1}^{2n} X_i = 2\bar{X}$$

$$\text{新样本 方差: } \frac{1}{n-1} \sum_{i=1}^n (X_i + X_{n+i} - 2\bar{X})^2 = \frac{Y}{n-1}$$

$$E\left(\frac{Y}{n-1}\right) = 2\sigma^2$$

→ 新总体之方差.

$$EY = 2(n-1)\sigma^2.$$

参数估计.

→ 关键是L如何写.

Stats
15

- 矩估计, 最大似然估计.

估计方法:

无偏估计, 一致性, 有效性.

↓ 何为无偏?

你得到一个 estimator, $\hat{\theta}$, 求期望

$$E\hat{\theta} = \theta.$$

- 估计量的求法.

- 区间估计.

$$X_i \sim N(\mu, \sigma^2)$$

$$\hat{\sigma}^2 = C \sum_{i=1}^{n-1} (X_{i+1} - X_i)^2$$

$$\text{欲使 } E\hat{\sigma}^2 = \sigma^2, \quad C = ?$$

即 $\hat{\sigma}^2$ 是 σ^2 的无偏估计.

$$E \sum_{i=1}^{n-1} (X_{i+1} - X_i)^2 = E \sum_{i=1}^{n-1} (X_{i+1}^2 - 2X_i X_{i+1} + X_i^2)$$

$$E \sum_{i=1}^{n-1} [(X_{i+1} - \mu) - (X_i - \mu)]^2 = E \sum_{i=1}^{n-1} \left(\underbrace{(X_{i+1} - \mu)^2}_{\Downarrow \sigma^2} - 2 \underbrace{(X_{i+1} - \mu)(X_i - \mu)}_{\substack{0 \\ \text{独立}}} + \underbrace{(X_i - \mu)^2}_{\Downarrow \sigma^2} \right)$$

F-test.

$$y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p.$$

We had some sort of population



within this population, we
have some true value (β_1)

所以假设检验的 ultimate
power 来自于什么?

- 中心极限定理

Sample distribution
趋于正态分布.

借助中心极限定理
的强大力量, 构建 t-test

对于单变量 regression 而言,
没问题.

若 test 多个变量, β_0, β_1, \dots

Not feasible within current framework.

我们只有 population 的一部分 (Sample)
利用这有限的样本, 我们要来做 test.
(Not 全集)

test if $\beta_1 = 0$.

- β_1 是 population 中的 true mean.

- 这便是 so-called 假设检验.

H_0 : null hypothesis.

$$\beta_1 = 0.$$

如 Null hypothesis:

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

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

$$\beta_1, \beta_2, \beta_3 = 0$$

p

Full population

$H_0:$

Null hypothesis is $\beta_1, \beta_2, \dots, \beta_p$ are non-significant.

$H_1:$

Alternative hypothesis is any of β_i is significant.

$$\beta_i \neq 0$$

If $\beta_i \neq 0$, it leads to us rejecting the NULL hypothesis.

Unrestricted regression

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \hat{u}$$

$$SSR = \sum_{i=1}^N \hat{u}_i^2 = (u_1^2 + u_2^2 + \dots + u_N^2)$$

estimation of population error.

We call residues.

Form restricted model.

$$y = \alpha$$

$$SSR_R > SSR_{UR}$$

"Statistically significant variables"

Bond=

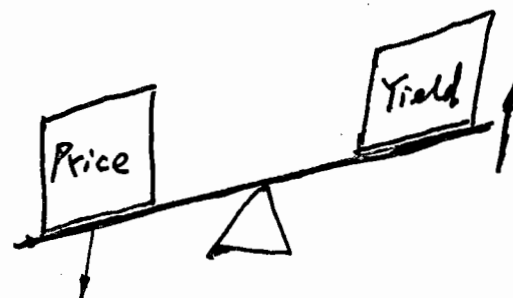
F2N

Price and yield are inverse,

Longer duration = more sensitive to rate moves

Bonds up

→ prices up (yield 其 其 down)



Old bond

Price \$100
Coupon \$5
Yield 5%

所以 old bond
你还欠 100 块，
便没人买，于是
你只好降价。

if price falls to \$83.33, a coupon
payment of \$5, will yield 6%
($5/83.33$)

New bond

Price \$100
Coupon \$6
Yield 6%

Obviously, every investor
will choose the new bond,
interest rate ~~has~~ has gone up,
so get more interest on
their loans.

→ interest rate goes higher,
so lenders demanding higher yield.

为什么只有 F -test $\bigg|$ ANOVA
 t -test

? 因为人们往往只关心 β 阶!

高阶的并不常用!

Volatility $\sim \chi^2$

What is the probability of developing a Surgical site infection (SSI) after a certain type of surgery?

At the 400 patients, who had this surgery, 12 develop an SSI,

$$\hat{p} = \frac{12}{400} = 0.03.$$

- \hat{p} : (sample proportion) is a point ^{estimation} ~~proportion~~ of p (population proportion)
- ^{点估计} \hat{p} When we draw our sample, \hat{p} will get a value.

→ But this is not usually the case, in statistics. Just giving a point estimation, is not good enough.

→ We want to use the data to say as much as we can about the parameter p .

- 两种常见之办法:

- Construct a confidence interval for p .
- We may have a hypothesis value about p to test.

To construct - 个 infer Inference procedures for p ,

We need to know the sampling distribution of the sample \hat{p} .

样本 proportion

(由样本得到的比例)

\hat{p} has a mean of p .

- 样本 proportion is an unbiased estimator of the population proportion.

总体

- standard deviation of sampling distribution of \hat{p}

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$$

总体方差是 $p(1-p)$

但是 p 是总体的参数

未知量!!

← We'll be able to walk around that.

- Sampling distribution of \hat{p} is approximately normal if the sample size is large.

Inference for
proportion
2.
stats.

The assumptions of the one-sample inference procedures
for a single proportion

↓
一次抽样.

Inference for
proportion

3

stats.

→ We have a simple random sample from the population of interest.
总体.

→ The sample size is large enough, for the
normal approximation to be reasonable.

→ It's reasonable 用 Z 估计, if:

$$n\hat{p} \geq 15 \quad \& \quad n(1-\hat{p}) \geq 15.$$

即若 \hat{p} 小, 则需要 n sample size 更大.

A $(1-\alpha)$ 100% confidence interval for p is given by:

$$\hat{p} \pm Z_{\alpha/2} \times SE(\hat{p})$$

↓
best estimator
of p .

margin of error.

standard error
of \hat{p} .

→ recall true standard deviation
of sampling distribution of \hat{p}

$$\sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$$

→ 我们不知 p . 于是 do the next best thing:

$$SE(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

假设检验的原理.

Mathematics

stat.

一个班里学生的得分.

假设为 80 分平均分, 5 分标准差.

开始检验

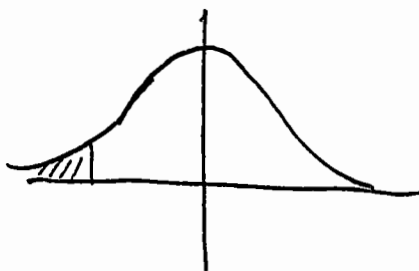
你上来抽一个 50 分. 跟你的 80 分有 66.

其概率是很低的.

于是你大概率上认为
你的假设没有问题.

either 平均分,

抑或是标准差出了问题.



否定域, 占据了面积为 α
的区域, 便否定之.

显著性水平 α .

通常取 5%.

但有时也取 10% (对结果之要求没那么
严格, 判错了也没关系时)

- 从回归方程可知, 计算过程中并不需要知道 Y 与 X 是否具有线性相关之关系.

→ 若不存在, 则求得之回归方程毫无意义.

→ 于是需对回归方程进行假设检验.

- 统计上称 β_1 是 $E(Y)$ 随 X 线性变化之变化率. 若 $\beta_1 = 0$, 则

$E(Y)$ 并不随 X 线性变化.

仅当 $\beta_1 \neq 0$ 时, $E(Y)$ 随 X 线性变化, 此时回归方程才有意义.

$$H_0: \beta_1 = 0, H_1: \beta_1 \neq 0$$

通常采用三种统计量 (构造统计量)

1) t 检验法: H_0 成立时, 统计量

$$T = \frac{\hat{\beta}_1}{\text{sd}(\hat{\beta}_1)} = \frac{\hat{\beta}_1 \sqrt{S_{xx}}}{\hat{\sigma}} \sim t(n-2)$$

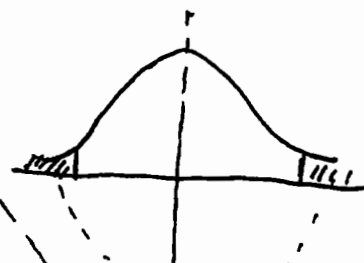
对于给定显著性水平 α , 检验拒绝域为:

$$|T| > t_{\alpha/2}(n-2)$$

符合以 $(n-2)$ 为自由度之 t 分布.
即, 有 20 个样本点,
则服从 $t(18)$

$\hat{\beta}_1, \sqrt{S_{xx}}, \hat{\sigma}$ 均可通过样本来算出, 不包含未知参数.

$$t = \frac{\hat{\beta}_1 \sqrt{S_{xx}}}{\hat{\sigma}} \text{ 即 } R \text{ 中 } \text{lm} \text{ 得 } t \text{ 值}$$



落于此二区域外, 为拒绝域
拒绝 H_0 , 即 $\beta_1 \neq 0$.
即 $E(Y)$ 随 X 线性变化.

t value $Pr(>|t|)$
13.51 9.50×10^{-8}

是构造的统计量 $T = \frac{\hat{\beta}_1 \sqrt{S_{xx}}}{\hat{\sigma}}$ 在当前样本下取值。



落在左侧区域，于是说 x, y 有线性关系 这一假设没有统计学意义。

$$Pr(>|t|) = 9.50 \times 10^{-8}$$

意义：随机得到此 $|t|$ 的概率是如此之小。

于是有统计规律在里面。

→ “认为回归方程是显著的”

2) F 检验法，构造的统计量不同，当 H_0 成立时，
 $\beta_1 = 0$ ，(即 x, y 无线性关系)

$$F = \frac{\hat{\beta}_1^2 S_{xx}}{\hat{\sigma}^2} \sim F(1, n-2)$$

对于给定的显著性水平 α ，检验拒绝域为 $F > F_{\alpha}(1, n-2)$

3). 相关系数检验法

$R = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}$ ，称为样本相关系数，对于给定的显著性水平 α ，查相关系数临界表可得 $r_{\alpha}(n-2)$

则检验拒绝域为： $|R| > r_{\alpha}(n-2)$

→ regression analysis explores the relationship between

- a quantitative response variable
- and
- one or more explanatory variables.

→ 若只有一个 explanatory: simply linear regression
多个 ——— : multiple ———.

→ 有何用. (linear regression).

→ 比如一个区域有野生的鹿, 也有狼.

鹿易于 measure.

狼更危险些.

→ 于是若知二者有线性关系,
则测鹿即可得狼的数目.
的数目

→ 又比如果种树, 想知其硬度.

但硬度难测.

→ 于是可测 density.

→ density 与硬度成正比.

The first 4 observation (of 36)

X	Density (lb/ft ³)	24.7	39.4	53.4	24.8
Y	Hardness (pound force)	484	1210	1880	427

- X: Explanatory var, independent var.
- Y: Response var, dependent var.

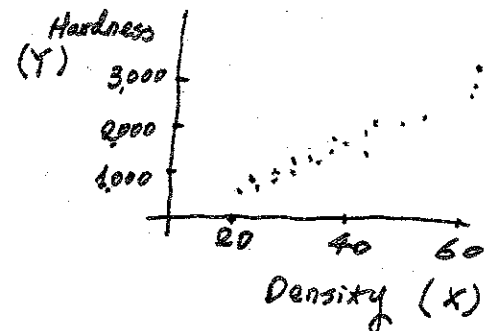
- The data set contains 36, (X, Y) pairs.

$$(X_1, Y_1) = (24.7, 484)$$

⋮

$$(X_4, Y_4) = (24.8, 427)$$

- First to do, plot our data to see what we're working with.

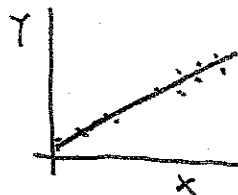


→

Can we use a known value of density (X) to help predict the hardness (Y)?

Yes.

fit a line 拟合直线



→ 你可能会问:

- How to come up with a good fitting line?
- Is a line a good summary between those vars?
- Is the relationship strong enough, such we can use it for prediction?

→ To measure the strength of the ^{linear} relationship, (7) correlation.

We assume a linear relationship between Y & x .

$$\mu_{Y|X} = \beta_0 + \beta_1 x$$

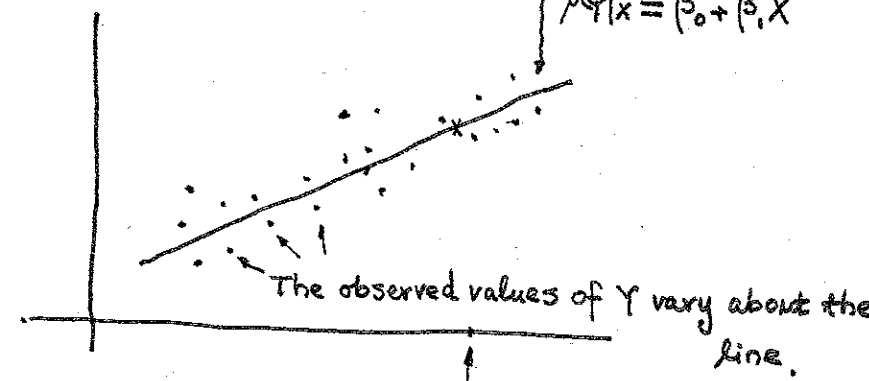
The true mean of Y , given x .

亦可写成:

$$E(Y|X) = \beta_0 + \beta_1 x$$

Expectation of Y given x , $\approx \beta_0 + \beta_1 x$

↓ We assume, this is the true relationship.



β_0 : y intercept

β_1 : the slope

即为 true relationship between Y & x .

$$\mu_{Y|X} = \beta_0 + \beta_1 x$$

→ 给一个 x 值, 则 theoretical mean of Y falls on that line.

但真实的 measurements 不会落在此理论预期值线上. 因为 variation.

We account for that variability around the line,
using ε .

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

random error component

represents the fact, Y will vary about the line.

- In the future, to do inference,
have to ~~do~~ make some assumption on
distributions of that random error components.

β_0 & β_1
are parameters,
They are typically
unknown values,

So we want to
estimate them.

We will use sample data to obtain the estimated
regression line

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

↑ ↑
这两个带 hat 的是 sample statistics,

→ We usually use the method of least square to estimate β_0 and β_1

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

→ Assumed relationship between Y and X .

We use data to find the estimated regression line.

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

predicted value of Y

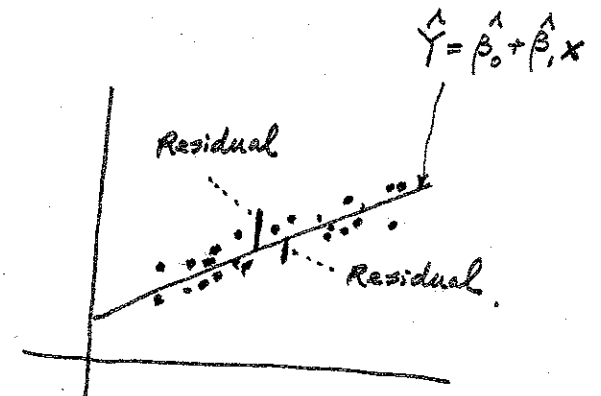
$\hat{\beta}_0$ is a statistic that estimates β_0

$\hat{\beta}_1$ is a _____ β_1

→ Residual = Observed - Predicted

$$e_i = Y_i - \hat{Y}_i \quad \text{--- every observation has a residual associated w. it.}$$

predicted value from the model.



Residual: signed distance

(从点到预测线的)
↓ 负的, 也是有正负号的

What's the best line?

Regression
6.

→ minimizes:

$$\sum |e_i| = \sum |Y_i - \hat{Y}_i|$$

→ Sum of vertical distances
between the points and the line.

→ $\hat{\beta}_0$ & $\hat{\beta}_1$ are chosen to minimize

$$\sum e_i^2 = \sum (Y_i - \hat{Y}_i)^2$$

$$= \sum (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i))^2$$

This is called the method of least squares.

Y_i \hat{Y}_i
→ each observation 的 prediction
的差, 再平方相加.
- 得到的去求最小值.

→ If the error terms are normally distributed,
it's the best possible line to use.

→ minimized summed least square residuals.

$$SS_{XX} = \sum (X_i - \bar{X})^2$$

$$SS_{YY} = \sum (Y_i - \bar{Y})^2$$

$$SP_{XY} = \sum (X_i - \bar{X})(Y_i - \bar{Y})$$

→ Variance of X , $S_X^2 = \frac{SS_{XX}}{n-1}$
samples.

Variance of Y , $S_Y^2 = \frac{SS_{YY}}{n-1}$

$$\text{Cov}(X, Y) = \frac{SP_{XY}}{n-1}$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} \quad (\text{mean of } Y \text{ 减去 } \hat{\beta}_1 \text{ 乘上 mean of } X)$$

$$\hat{\beta}_1 = \frac{\sum P_{xy}}{\sum x^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})(x_i - \bar{x})}$$

$$= \frac{\text{COV}(X, Y)}{\text{Var}(X)}$$

例: R 做 linear regression:

Coefficients:

	Estimate	Std Error	t	Pr(> t)
Intercepts	-1160.50	108.50	-10.69	2.07e-12
Density	57.507	2.279	25.24	< 2e-16

slope,
labeled with explanatory var X.

$$\hat{\beta}_0 = -1160.50$$

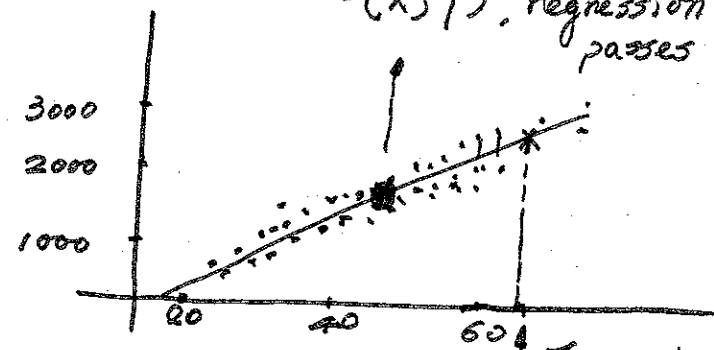
$$\hat{\beta}_1 = 57.507$$

3p least square regression line:

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X$$

符号跟高
往下为负
往上为正

- 所有 sample 点到 regression line 的距离 sum of squares
- (\bar{X}, \bar{Y}) , regression line always passes through that.



有了 prediction line,
知道个 X, 便可以预测
其 predicted Y

For least square regression:

→ The residual sum to 0 $\sum e_i = 0$

↳ mathematically, every time, least square sum to zero.

→ The line always passes through (\bar{X}, \bar{Y})

The sample linear regression model:

$$\mu_{Y|X} = \beta_0 + \beta_1 X$$

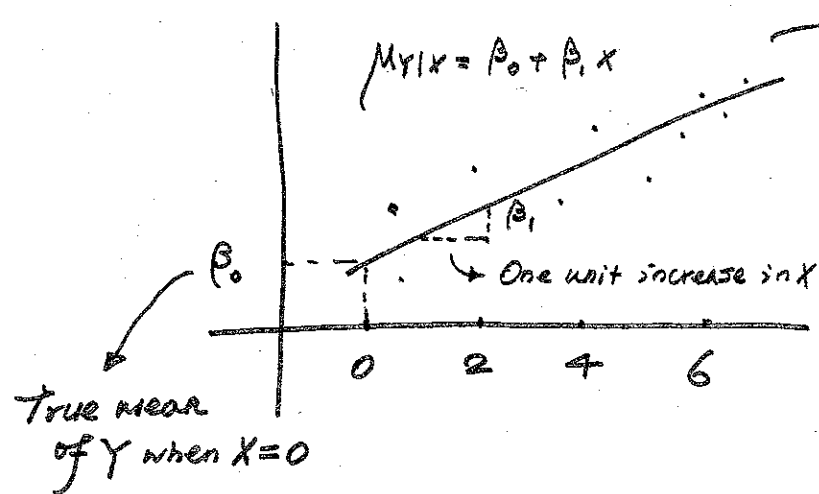
→ ~~The~~ ^{true} mean of Y given X , 等于 $\beta_0 + \beta_1 X$

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

↳ Y value itself does not fall precisely on this line,
↳ a random variable components.

β_0, β_1 are parameters

当 $X=0$ 时
 $\mu_{Y|X} = \beta_0$
即 β_0 是 true mean of Y ,
when $X=0$.
→ 这 practically 可能没用,
因为 $X=0$ 可能比较 rare.
 β_1 : 当 X 变化 1 时,
true mean of Y 变化 β_1



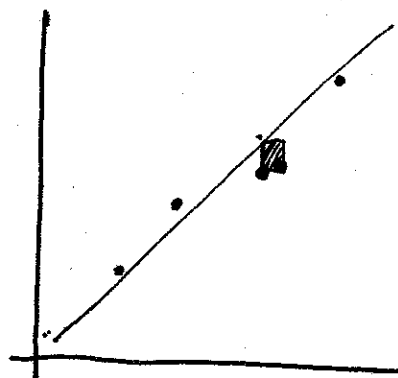
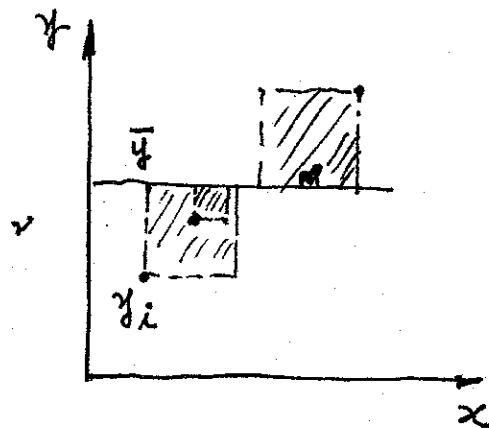
- 通常我们不可能知道
- these true parameters in any given situation.
 - We're going to get sample data and estimate β_0 & β_1 .

线性回归基本假设:

1. 随机误差项是一只期望值或均值为0的R.V.
2. 对于 explanatory var 之所有 observation, 随机误差有相同方差.
3. 随机误差彼此不相关.
4. Explanatory var 是确定性变量, 不是R.V., 与随机误差之间彼此相互独立.
5. 随机误差服从正态分布.

0 同方差性

"灵童独存显"



total sum of squares.

正比于 variance of the data

$$SS_{tot} = \sum_i (y_i - \bar{y})^2$$

最no-brain的,是用
均值这条平行于x轴的线来预测。
例。

$$SS_{reg} = \sum_i (y_i - \hat{y})^2$$

经过 linear regression 后,
误差变小了。

$$R^2 \equiv \frac{SS_{tot} - SS_{reg}}{SS_{tot}}$$

↑
未做线性回归的误差平方和

↑
regression 后误差平方和

$$= 1 - \frac{SS_{reg}}{SS_{tot}}$$

→ 是做差后。
便是线性回归模型
所做贡献。