

# DTS206TC Formulas

## Distributions

Exponential Dispersion Family 指数分布族:

$$p(y; \theta, \phi) = \exp\left\{\frac{y\theta - b(\theta)}{\phi} + c(y, \phi)\right\}$$

$\mu = \theta$   $\sigma^2 = \phi$  均值  $\mu$ , 方差  $\sigma^2$

Normal Distribution

$$p(y; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2}(y - \mu)^2\right\}$$

高斯分布的概率密度函数 polf.

Binomial Distribution 伯努利分布

$$p(y; \pi) = \binom{n}{y} \pi^y (1 - \pi)^{n-y}, \text{ where } y = 0, 1, \dots, n, \text{ and } \binom{n}{y} = \frac{n!}{y!(n-y)!}$$

均值  $\pi$  方差  $(1-\pi)$

$ny$  has a Binomial Distribution.

$$p(y; \pi, n) = \binom{n}{ny} \pi^{ny} (1 - \pi)^{n-ny}$$

where,  $y = 0, \frac{1}{n}, \frac{2}{n}, \dots, 1$ .

Poisson Distribution

泊松

$$p(y; \mu) = \frac{\mu^y e^{-\mu}}{y!}$$

均值  $\mu$  方差  $\lambda$ .

where,

- $y = 0, 1, 2, \dots$
- $\mu$  is the parameter of the distribution

Geometric distribution

几何分布的 Binomial Distribution:

$$f(y; \pi) = \pi(1 - \pi)^y, y = 0, 1, 2, \dots$$

Estimation, Inference

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

$$S_{xx} = \sum (x_i - \bar{x})^2$$

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

$$s_R^2 = MSE = \frac{RSS}{n-k}$$

residual variance.

$$y = \beta_0 + \beta_1 x^* \text{ 此时 } k=2$$

$k$  为模型参数数量.

$$\begin{aligned}\sigma^2\{\hat{\beta}_1\} &= \frac{\sigma_{y|x}^2}{\sum(x_i - \bar{x})^2} \\ s^2\{\hat{\beta}_1\} &= \frac{s_R^2}{\sum(x_i - \bar{x})^2}\end{aligned}$$

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\frac{s_R^2}{(n-1)s_x^2}}} \sim t_{n-2}$$

$$\hat{\beta}_1 \pm t_{n-2,\alpha/2} \cdot s\{\hat{\beta}_1\}$$

$$\hat{\beta}_1 \pm t_{n-2,\alpha/2} \sqrt{\frac{s_R^2}{(n-1)s_x^2}}$$

$$\sigma^2\{\hat{\beta}_0\} = \sigma_{y|x}^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right)$$

$$s^2\{\hat{\beta}_0\} = s_R^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right)$$

$$s^2\{\hat{\beta}_0\} = s_R^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2} \right)$$

$$\frac{\hat{\beta}_0 - \beta_0}{s\{\hat{\beta}_0\}} \sim t_{n-2}$$

$$\hat{\beta}_0 \pm t_{n-2,\alpha/2} \cdot s\{\hat{\beta}_0\}$$

$$\frac{\hat{\beta}_0 - \beta_0}{\sqrt{s_R^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2} \right)}} \sim t_{n-2}$$

$$\hat{\beta}_0 \pm t_{n-2,\alpha/2} \sqrt{s_R^2 \left( \frac{1}{n} + \frac{\bar{x}^2}{(n-1)s_x^2} \right)}$$

$$\hat{y}_0 \pm t_{n-2,\alpha/2} \sqrt{s_R^2 \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{(n-1)s_x^2} \right)}$$

$$\hat{y}_0 \pm t_{n-2,\alpha/2} \sqrt{s_R^2 \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{(n-1)s_x^2} \right)}$$

Bayesian topic

$$\begin{aligned}p(\beta|Y) &= \mathcal{N}(\beta|\mu_n, \Sigma_n) \\ \mu_n &= \Sigma_n (\Sigma_0^{-1} \mu_0 + \sigma^{-2} X^T Y) \\ \Sigma_n &= (\Sigma_0^{-1} + \sigma^{-2} X^T X)^{-1}\end{aligned}$$

## Background on Regression Model.

### • Regression Model 回归模型.

x: Independent or Explanatory variable. 独立或解释变量.

y: Dependent or Response variable 依赖或应变量.

从 n 对值的样本  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  研究不同 x 值的合理估计

### • Regression Analysis 回归分析

statistical methodology. 利用两个或者多个定量变量的关系，以便从一个变量  
预测另一个变量或响应变量.  
quantitative variable.

a response or outcome variable can be predicted from other or others.

### • Relations between Variables

Deterministic: Functional Relation 确定性. 功能关系.

↓ expressed by dependent

a mathematical formula

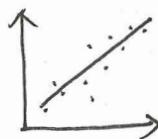
$$y = f(x)$$

↓ independent

Nondeterministic: Statistical Relation 非确定性. 统计关系.

not perfect.

$y = f(x) + \epsilon \leftarrow \epsilon$  is an unknown perturbation  
(random variable)



## Type of Relation

Linear.  $y = a + bx$  /  $y = mx + b$

$\downarrow$   
 $y = \beta_0 + \beta_1 x \rightarrow$  a straight line (linear)

$\beta_1 > 0 \rightarrow$  positive linear relationship

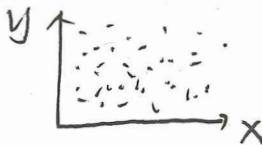
$\beta_1 < 0 \rightarrow$  negative linear relationship

Nonlinear  $f(x) = \log(x)$

$f(x) = x^2 + x + 1$

## Absence of relation

Whenever  $f(x) = c$ , that is, whenever  $f(x)$  doesn't depend on  $x$ .



## Measures of Linear Dependency

Covariance  $\text{协方差}$   $\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$  取决于变量的  
同时量单位

- ① positive linear relation  $\rightarrow \text{cov}$  positive and large.
- ② negative linear relation  $\rightarrow \text{cov}$  negative and large in absolute value
- ③ no relation/significantly ~~linear~~  $\rightarrow \text{cov}$  close to zero

the Correlation Coefficient 相关系数.

$$r(x, y) = \text{cor}(x, y) = \frac{\text{cov}(x, y)}{S_x S_y} \quad S_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \quad S_y^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$$

- ①  $-1 \leq \text{cor}(x, y) \leq 1$  ②  $\text{cor}(x, y) = \text{cor}(y, x)$  ③  $\text{cor}(ax+b, cy+dx) = \text{sign}(a)\text{sign}(c)\text{cor}(x, y)$

- Simple Linear Regression (SLR) Model with Distribution of Error Terms Unspecified.

Basic regression model  $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$

$\uparrow$        $\uparrow$   
 系數 / 參數  
 coefficients / parameters

$y_i$  = the value of the dependent variable for the  $i$ -th observation

$x_i$  = the value of the independent variable for the  $i$ -th observation

$\epsilon_i$  : the error for the  $i$ -th observation,

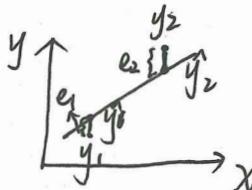
mean  $E\{\epsilon_i\} = 0$  and variance  $\sigma^2 \epsilon_i = \sigma^2$  方差 (各樣本數指與平均數之差的平方和)  
 $\beta_0, \beta_1$  are parameters or coefficients.  $\beta_0$ : intercept 截距 平方和  
 $\beta_1$ : slope 斜率.

The parameter's to estimate are:  $\beta_0, \beta_1, \sigma$ .

### SLR Model

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_i \leftarrow \text{the regression line}$$

residual 残差  $e_i = y_i - \hat{y}_i$   
 有  $\epsilon$  號



• Hypothesis of the SLR Model 假设.

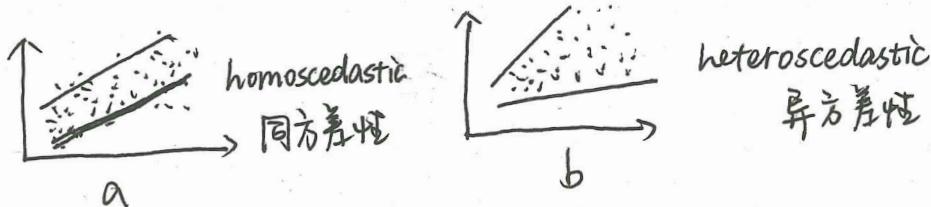
线性 Linearity: The existing relation between  $x$  and  $y$  is linear.  
 data look reasonably straight  $f(x) = \beta_0 + \beta_1 x$

齐性 Homogeneity: The mean value of the error is zero.  
 $E[\epsilon_i] = 0$

方差齐性 Homoscedasticity: The variance of the errors is constant

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

数据分散, 疏密恒定. dispersion of the data must be constant.



独立性 Independence: The observations are independent.

数据必须独立.  $E[\epsilon_i \epsilon_j] = 0$

Any  $y_i$  and  $y_j$  are uncorrelated  
 two response

①  $y_i$  a random variable : (i) constant term  $\beta_0 + \beta_1 x$  (ii) random term  $\epsilon$

$$② E\{y\} = E\{\beta_0 + \beta_1 x + \epsilon\} = \beta_0 + \beta_1 x + E\{\epsilon\} = \beta_0 + \beta_1 x$$

③  $y_i$  或者大于或者小于  $(y = \beta_0 + \beta_1 x + \epsilon_i)$

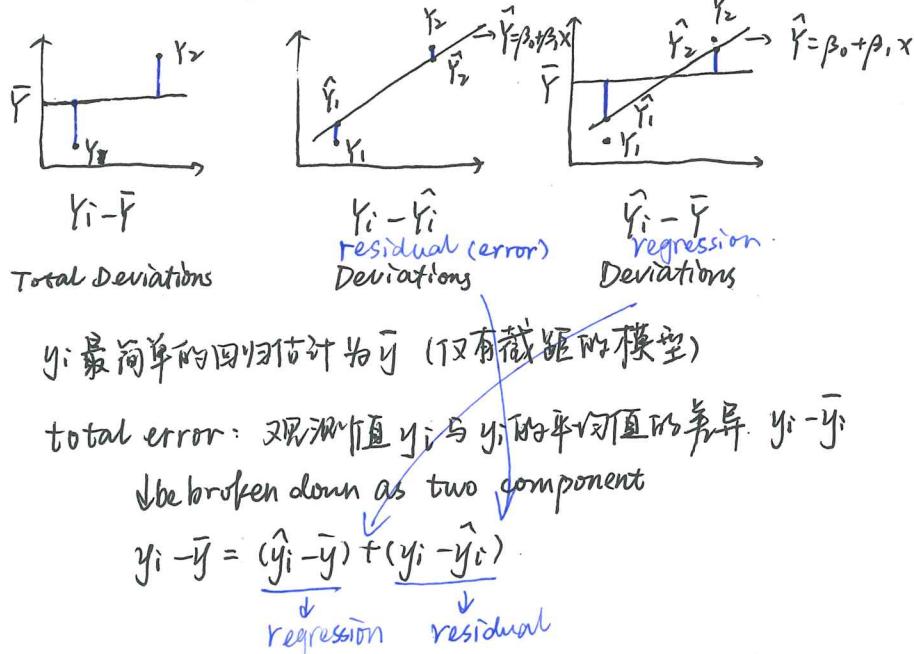
④ The error term  $\epsilon_i$  are assumed to have constant variance  $\sigma^2$ .

$\epsilon_i$  and  $\epsilon_j$  are uncorrelated. so are the responses.

⑤ cause  $\epsilon_i$  and  $\epsilon_j$  not correlated  $\rightarrow y_i$  and  $y_j$ .

## Total Deviation Partition.

↓ 检查回归线的拟合度.



$y_i$  最简单的回归估计为  $\bar{y}$  (仅有截距的模型)

total error: 观测值  $y_i$  与  $\bar{y}$  的平均值的差异  $y_i - \bar{y}$

↓ be broken down as two component

$$y_i - \bar{y} = (\hat{y}_i - \bar{y}) + (y_i - \hat{y}_i)$$

↓ regression      ↓ residual

↓ 这种变化导致了回归分析的基本方程.

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$

$$SSTO = SSR + SSE$$

↓ Total sum of squares 表示样本点与样本均值之间的垂直距离的平方和

所有  $y$  相同,  $SSTO = 0$ .  $y$  的变化越大,  $SSTO$  越大  
(离散程度)

SSE Measure of after predictor effect  $\downarrow$

↓ error sum of squares  $\downarrow$   $\hat{y}_i$  是  $y$  的变化度量.

$$SSE = \sum (y_i - \hat{y}_i)^2$$

SSR Regression sum of square (回归平方和)  $SSTO - SSE = SSR$

↓ regression sum of squares  $SSTO$  和  $SSE$  之间的差为  $SSR$ .

$$SSR = \sum (\hat{y}_i - \bar{y})^2$$

ANOVA Table 方差分析表. 用于汇总回归结果.

↓

the analysis of variance

$n$  为样本大小.

$p$  为模型中包含的预测变量的数量.

自由度.

Source of Variation

Sum of Squares

df

MS

Variance Ratio (F)

p-value: F 统计量的显著性水平.

Regression

SSR

p

MSR

$$F = \frac{MSR}{MSE}$$

$$Pr(F_{p, n-p-1} > F)$$

Error

SSE

n-p- $\beta_0$

MSE

Total

SSTO

n-1

MS 平均平方和

每个来源的平方和除以  
其自由度.

$$MSR = \frac{SSR}{p}$$

Regression Mean Square

$$MSE = \frac{SSE}{n-p-1}$$

Error Mean Square

F-statistic.

↑

$$F = \frac{MSR}{MSE}$$

$$Pr(F_{p, n-p-1} > F)$$

F Test for SLR ( $H_0: \beta_1 = 0$  for SLR)

对于  $MSR$  的期望是  $E(MSR) = \sigma_y^2 + \beta_1^2 \sum (x_i - \bar{x})^2$

对于  $MSE$  的期望是  $E(MSE) = E(\sigma_y^2 | x) = \sigma_y^2 | x$

两个方差相等 ( $\sigma_1^2 = \sigma_2^2$ ) 的情况下, 为了获得的比值服从  $F$  分布  $\frac{MSR / \sigma_1^2}{MSE / \sigma_2^2} \sim F_{p, n-p-1}$

在回归线的真实斜率非零的假设下 ( $H_0: \beta_1 = 0$ )  $MSR$  和  $MSE$  是  $\sigma_y^2 | x$  的独立估计.

因此  $F^* = \frac{MSR}{MSE} \sim F_{p, n-p-1}$  回归平方与残差均方比不服从具有  $p$  和  $n-p-1$  个自由度的  $F$  分布

F检验用于检验包含协变量的模型是否与仅包含截距项的模型相比, 导致残差平方和显著减少.

(covariates)

带有自变量  $x$  的回归模型

$$Y = \beta_0 + \beta_1 x + \epsilon$$

仅有常数项的模型

$$Y = \beta_0 + \epsilon$$

如果原假设为真 ( $H_0: \beta_1 = 0$ ) F比的期望值应该为1. 如果  $H_0$  错误, F比的期望值应该大于1.

决定系数 Coefficient of Determination  $R^2$

$$R^2 = \frac{SSR}{SSTO} = 1 - \frac{SSE}{SSTO} \quad 0 \leq SSE \leq SSTO \Rightarrow 0 \leq R^2 \leq 1$$

①  $R^2 = 0$  时: regression line explains none of the variability in  $y$  and the regression line is no better than using  $\bar{y}$  as our predictor of  $y$ .

回归线不如  $\bar{y}$  用单预测  $y$ .

②  $R^2 = 1$ : a perfect fit and regression line explains all ~~the~~ of the variability. 这种情况下所有的数据点都落在回归线上且没有残差

Residual variation

不衡量斜率的大小或衡量直线模型的适当性.

not measure magnitude of the slope

or measure the appropriateness of the straight-line model.

$R^2$  很大时并不意味着适当的模型.

③ high  $R^2$  doesn't mean that there is a good linear fit between predictor and output. 不一定有很好的线性拟合=实际曲线关系的近似(不高的相关性)

④ low  $R^2$  doesn't mean that there is ~~a~~ no relationship between input and output

$y$  的方差 variance of  $y$ .  $\hat{y}^2$  可能有关系, 但线性函数无法解释这些关系.

$$\hat{y}^2 = \frac{\sum (y_i - \bar{y})^2}{n-1} = \frac{SSTO}{n-1}$$

由于计算和包含均值而产生  $\bar{y}$  线性约束

linear constraint

$\hat{y}^2$  (Residual standard error)<sup>2</sup>

$$\hat{y}^2 = \frac{\sum (y_i - \hat{y}_i)^2}{n-2} = \frac{SSE}{n-2}$$

n-2 degrees of freedom

$\beta_0$  和  $\beta_1$  估计的线性约束

$$(n-1) - (n-2) = 1 + 1 = 2$$

1 degree of freedom:  $\beta_1$  依赖一个参数, 即斜率  $\beta_1$ .

因为减少了平均响应, mean response. 截距参数就消失了. 只保留 slope 参数.

自由度

F-statistic.

↑

$$Pr(F_{p, n-p-1} > F)$$

- A T-test will tell if a single variable is statistically significant.

- An F-test will tell if a group of variables are jointly significant.

$H_0: \beta_1 = 0$  vs  $\beta_1 \neq 0$ .

Two-sided test  $H_0: \beta_1 = 0$ .  $H_1: \beta_1 \neq 0$ .

T-test statistic from before

$$t^* = \frac{\beta_1 - \bar{\beta}_1}{\text{SE}(\beta_1)}$$

ANOVA statistic

进行 F 检验再进行 t 检验.

判断  $\alpha$  是否显著

判断  $\beta_1$  是否显著

t-Test:

$$t = \frac{\beta_1}{\text{SE}(\beta_1)} \text{ where } t \sim t_{n-p-1} \text{ under } H_0$$

F-Test.

$$\text{If } x \sim t_n, \text{ then } x^2 \sim F_{1, n}$$

## Motivation.

The partial F-test is used to in model building and variable selection to help decide if a variable or term can be removed from a model without making the model significantly worse.

在使得模型不明显变差的情况下，使用F检验用于模型构建和变量选择。

帮助确定是否可以用于从模型中删除变量或项。

Partial F-test is used to compare Nested model. 用于比较嵌套模型。

Full (or complete) model 包括所有预测因子, predictors of interest.

Reduced model. 完整模型的某子集 subset

在 MLR Model 的背景下. 总体 F 检验: Overall F-test

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

整组自变量是否对  $y$  的预测有显著贡献

$$H_1: \text{at least one of the } \beta_k \neq 0$$

Full model  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$

Reduced model  $y = \beta_0$

偏 F 检验: Partial F test 比较完整模型和简化模型的 SSE.

$$H_0: \beta_1^* = \beta_2^* = \dots = \beta_k^* = 0$$

测试一组变量，添加一些自变量是否显著增加了  $y$  的预测。

$$H_1: \text{at least one of the } \beta_k^* \neq 0$$

超过了模型中已存在的其他自变量所实现的预测。

Full model  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + (\beta_{p+1}^* x_{p+1} + \dots + \beta_k^* x_k)$

Reduced model  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$

对于 Reduced Model

简化模型是否会显著增加，从而降低模型的预测能力。误差平方和 SSE

对于 Full Model

模型中包含 (新的)自变量是否会降低误差并显著提高模型的预测能力  
或 是否显著改善了模型。

## 更简单的模型

某种程度上，SSE 可以用来总结模型的拟合程度。

$H_0$ : No significant difference in SSE of Full and Reduced Models

$H_1$ : Full Model has a significantly lower SSE than the Reduced Model.

## F Test Statistic

change in the number of parameters.

$$F^* = \frac{(SSE_{\text{Reduced}} - SSE_{\text{Full}})/k}{MSE_{\text{Full}}} \sim F_{k, n-p-k-1}$$

越大  $\Rightarrow$  两个模型之间的 SSE 越大

k: 增加了 k 个自变量 x

p: 原来有 p 个自变量 x

Covariance 协方差 对于线性相关性的度量. Linear dependency.

$$\text{cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

正线性相关 cov 非负大. 反线性相关 cov 非正大.

负线性相关 cov 负且绝对值大.

没有关系变量之间 cov 为0 或趋近于0.

The Correlation Coefficient 相关系数.  $\rightarrow$  协方差的标准化.

$$r(x, y) = \text{cor}(x, y) = \frac{\text{cov}(x, y)}{s_x s_y} \quad s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \quad s_y^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$$

$$-1 \leq \text{cor}(x, y) \leq 1 \quad \text{cor}(x, y) = \text{cor}(y, x)$$

$$\text{cor}(ax+b, cy+d) = \text{sign}(a)\text{sign}(c) \text{cor}(x, y)$$

取 a 和 c 同符号 (+)

$$\text{Residual } e_i = y_i - \hat{y}_i$$

Hypothesis  $\rightarrow$  SLR Model

$$\{ f(x) = \beta_0 + \beta_1 x$$

$E[e_i] = 0$  mean value of the error is 0

$\text{Var}(e_i) = \sigma^2$  variance error is constant

$E[e_i e_j] = 0$  observations are independent.

Important Features  $\rightarrow$  SLR Model.

$$E\{y\} = E\{\beta_0 + \beta_1 x + \epsilon\} = E\{\epsilon\} + \beta_0 + \beta_1 x = \beta_0 + \beta_1 x$$

Discrete Random Variable 离散随机变量.

Cumulative probability distribution function  $\rightarrow$  c.d.f.

累加概率分布函数.



取值范围 [0, 1]

单调递增

离散

随机变量小于等于特定值的概率.

$$\text{期望 } E(Y) = \sum_{i=1}^n y_i p_i \quad \mu_Y \leftarrow \text{可用此表示 加权平均值.}$$

Expected value.

总体方差 Variance.

$$\text{Var}(Y) = \sigma_Y^2 = E[(Y - \mu_Y)^2] = \sum_{i=1}^n (y_i - \mu_Y)^2 p_i$$

样本方差 Sample Variance

$$s_Y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 \quad \text{sample mean}$$

标准差 Standard deviation  $\sigma_Y = \sqrt{s_Y^2}$  square root of variance.

$$\text{Var}(Y) = \frac{1}{N} \sum_{i=1}^N (y_i - \mu_Y)^2 \text{ 此时 } p_i \text{ 都一样.}$$

Bessel's correction

除以  $n-1$  而不是  $n$  提供对总方差的无偏估计.

Continuous Random Variables 连续随机变量.

Probability density function (pdf) 概率密度函数.

概率密度函数随机变量. 各个取值发生的相对概率.

$$P(a \leq Y \leq b) = \int_a^b f_Y(y) dy$$

$$P(-\infty \leq Y \leq \infty) = 1 \Rightarrow \int_{-\infty}^{\infty} f_Y(y) dy = 1$$

Standardization VS Normalization?

标准化

$$X \sim N(\mu, \sigma^2) \quad \begin{array}{l} \text{公式} \\ \downarrow \quad \downarrow \\ \frac{|x - \mu|}{\sigma} \end{array}$$

均值 方差

保留了原始数据 极端值

形状和分布.

归一化

将特征的值映射到一个固定区间. 例如  $[0, 1]$  或  $[-1, 1]$

常用方法:

$$\text{Min-Max} = \frac{(x - \min)}{(\max - \min)}$$

$$Z\text{-score} = \frac{(x - \mu)}{\sigma}$$

Residuals VS Error?

观测值与模型中  $y_i$  的差值. 观测值与真实值的差值.

最小二乘法. Least Squares

$$\textcircled{1} \text{ 设 } \hat{y}_i = \beta_0 + \beta_1 x_i;$$

$$\textcircled{2} \text{ 计算 } \bar{x}, \bar{y}$$

$$\textcircled{3} \text{ 计算 } \beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} S_{xx}$$

$$\textcircled{4} \text{ 计算 } \beta_0 = \bar{y} - \beta_1 \bar{x}$$

\textcircled{5} 得出 SLRM

Arithmetic Mean 算术平均值.

$$\text{Mean} = \frac{\sum x_i}{n} \quad \begin{array}{l} \leftarrow \text{raw data value} \\ \leftarrow \text{sample size} \end{array}$$

连续.

$$E(Y) = \mu_Y = \int y f_Y(y) dy \text{ 积分(求面积.)}$$

$$\text{Var}(Y) = \sigma_Y^2 = \int (y - \mu_Y)^2 f_Y(y) dy$$

$$\text{Var}(X) = E(X^2) - E(X)^2$$

Distribution 分布.

Normal Distribution 正态分布  $N(\mu, \sigma^2)$

$$\text{pdf: } f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}.$$

For standard normal distribution 标准正态分布  $N(0, 1)$

the standard normal PDF 用  $\phi$  表示.  $\phi(c) = \Phi'(c)$

the standard normal CDF 用  $\Phi$  表示.  $\Phi(c) = P(Z \leq c), Z \sim N(0, 1)$

The Chi-Squared Distribution 卡方分布 (遇到连续类型问题设时需要)

M 平方独立标准正态分布随机变量的总和服从以具有 M 自由度的卡方分布.

The sum of M squared independent standard normal distributed random variables follows a chi-squared distribution with M degrees of freedom.

$$Z_1^2 + \dots + Z_M^2 = \sum_{m=1}^M Z_m^2 \sim \chi_M^2 \text{ with } Z_m \stackrel{i.i.d.}{\sim} N(0, 1)$$

自由度为 M. 众数为  $M-2$ . 期望为 M 方差  $\text{Var} \approx 2M$

for  $M \geq 2$

$$\text{pdf: } f(x) = \begin{cases} \frac{1}{2^{\frac{M}{2}} \Gamma(\frac{M}{2})} \cdot x^{\frac{M}{2}-1} \cdot e^{-\frac{x}{2}}, & x > 0 \\ 0, & \text{其他.} \end{cases}$$

概率分布

## The F Distribution F 分布

两个独立的  $\chi^2$  分布随机变量以自由度 M 和 n 的比率为 F 分布

$$\frac{W/M}{V/n} \sim F_{M,n} \text{ with } W \sim \chi_M^2, V \sim \chi_n^2$$

分子自由度 M. 简称为  $F_{M,n}$   
分母自由度 n

如分子自由度很大. 以  $F_{M,n}$  分布近似于  $F_{M,\infty}$  分布

对于正态总体来说,  $W/M \sim F_{M,\infty}$ ,  $W \sim \chi_M^2$

两个总体的方差比较可以用 F 分布来检验

$$F = \frac{\frac{\sum S_1^2}{n_1 - 1}}{\frac{\sum S_2^2}{n_2 - 1}} = \frac{\sum (X_i - \bar{X})^2}{n_1 - 1} / \frac{\sum (X_j - \bar{X})^2}{n_2 - 1}$$

Ex 2 / Exercise  
Bernoulli Distributed random variable 伯努利分布随机变量.

binomial distribution. 伯努利试验中的成功次数服从二项分布

二项分布.  $k \sim B(n, p)$

↑

the number of success

$$f(k) = P(k) = \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k}$$

## The Student t Distribution t 分布

$Z$  为 standard normal variate 标准正态变量.

$Z$  与  $W$  独立.  $M$  为自由度.

$$\frac{Z}{\sqrt{W/M}} =: X \sim t_M \quad X \text{ follows a Student t distribution}$$

$t_M$  分布取决于  $M$ .  $M \uparrow$  t 分布越像正态分布.

当  $M > 30$ . t 分布近似于正态分布. 有近似.

t<sub>∞</sub> 分布为标准正态分布.

当  $M > 1$ . t<sub>M</sub> 分布随机变量  $X$  有期望  $E(X) = 0$ ,  $M > 1$

当  $M > 2$  有方差  $Var(X) = \frac{M}{M-2}$ ,  $M > 2$

## Lecture 6:

### 线性回归.

C.I.

置信区间

对真值总体均值具有一定置信水平的范围估计.

平均响应; 通过回归方程来估计.

$$\hat{y}_0 = \beta_0 + \beta_1 x_0$$

在估计平均响应( $\hat{y}_0$ )的情况下, 置信区间可以计算为:

$$\hat{y}_0 \pm t_{n-2, 1-\frac{\alpha}{2}} * SE(\hat{y}_0) \Rightarrow \beta_0 + \beta_1 x_0 \pm t_{n-2, 1-\frac{\alpha}{2}} * SE(\hat{y}_0)$$

- $\hat{y}_0$ : 表示估计的平均响应
- $t_{n-2, 1-\frac{\alpha}{2}}$ : 表示分布的临界值.
- $SE(\hat{y}_0)$ : 表示标准误差估计的平均响应.

$$\hat{y}_0 \pm t_{n-2, 1-\frac{\alpha}{2}} \sqrt{S^2 R \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{(n-1)S_x^2} \right)}$$

$$S^2 R = \frac{\sum e_i^2}{n-2}. \text{ (无偏估计).}$$

$$\frac{1}{2\pi\sigma} \int_{-\infty}^{\infty} e^{-\frac{(y-\mu)^2}{2\sigma^2}} dx$$

$\ln b = \text{blne.}$



口↑, 曲线更窄  
口↓, 曲线更宽.

of a Mean Response;

平均响应.

均响应?

极大似然.

$$Pr(x|M, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

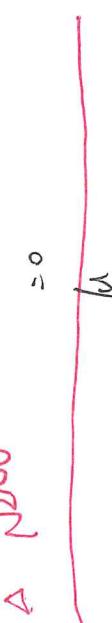
$\mu$ : 均值.

$\sigma$ : 方差.

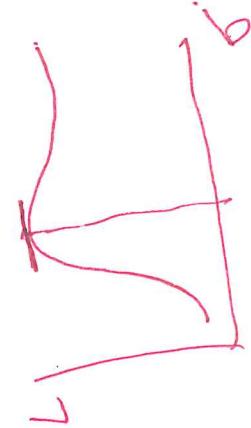
$\mu, \sigma^2$  (参数).



$$\frac{1}{\Delta} = \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2} \right)^N$$



$\mu$  (参数).



区别 C.I. 和 P.I. (PPT 10 原话):

We need to distinguish between P.I. and C.I.;  
P.I. is prediction intervals, where we predict outcome  
of one single experiment, and confidence intervals.  
(预测单个实验结果, 以及置信区间).  
C.I. is where we predict the mean value of future  
outcomes. (预测未来结果的平均值).

P.I.

预测区间

提供了一个范围的估计值, 在该范围内, 总体的单个观察值预计会以一定的置信度水平下降.

$$\hat{y}_0 \pm t_{n-2, 1-\frac{\alpha}{2}} * SE(\hat{y}_0)$$

•  $\hat{y}_0$ : 表示给定  $x_0$  的  $y$  预测值.

•  $t_{n-2, 1-\frac{\alpha}{2}}$ : 表示分布的临界值.

•  $SE(\hat{y}_0)$ : 表示预测值的标准误差.

## Hypothesis Test

线性!!!

① on the slope / linearly / SLM ( $y = \beta_0 + \beta_1 x$ ) etc.

没有  $\beta_1$ , 就是常数的了. 所以: Assume

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

$H_0$  的 reject 区域

The rejection region for the null hypothesis is:

要背 \*

$$\left| \frac{\hat{\beta}_1}{\sqrt{S_R^2/(n-1)S_x^2}} \right| > t_{n-2, \alpha/2}$$

two-tails

题目中应用 like: 先算  $\left| \frac{\hat{\beta}_1}{\sqrt{S_R^2/(n-1)S_x^2}} \right| = \dots = \dots > ? < t_{n-2, \alpha/2}$

② 3点二:

同样, 如果 the value zero is outside of the CI for  $\beta_1$  at  $(1-\alpha)$  <sup>one-tail</sup> level, we reject the null hypothesis at this level.

The p-value of the test is:

③  $P\text{-value} = 2 \Pr(t_{n-2} > \left| \frac{\hat{\beta}_1}{\sqrt{S_R^2/(n-1)S_x^2}} \right|)$  越来越远越好!

if p-value  $< \alpha$ , we can reject  $H_0$ .

∴ Hence,

Probability value

R 语言解读:

57.22% of the variability in  $y$  is explained by the linear relationship with  $x$ .

Multiple R-squared: 多重决定系数, 因变量的比值 e.g. 0.5722.

Adjusted R-squared: 调整后决定系数 考虑了自变量, 我们可以用 Adjusted R 检正.

F-statistic: F 统计, 如果 F-statistic 的值  $< 0.05$ .

我们认为模型显著的, 至少有一个自变量对因变量有影响.

df = n - (1+p) 只要它对应的 p-value  $< 0.05$ ,

模型显著

F-stat.: xx on P and n-(1+p) DF.

\* 表示小于多少, 只要 p 值  $< \alpha$  level

(P 小于 0.001 ...)

(" Significance test ")

④ 假设检验 F\*

$F^*$  分布在  $F(1, n-2)$  when  $H_0$  holds,

将该错误风险控制在 I error 时, F statistic 符合

$$F^* \sim F(1, n-2) \text{ 分布} \Rightarrow$$

① if  $F^* \leq F(1-\alpha; 1, n-2)$  conclude  $H_0$

② if  $F^* > F(1-\alpha; 1, n-2)$  conclude  $H_1$

## Lecture 13 和 14

GLM 一般线性 (Generalized!)

→ 不止解释一个变量

general linear model

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i \quad \text{误差.} \Rightarrow 正态分布}$$

响应变量 response variable  $y_i$

解释 explanatory  $x_j, j=1 \dots p.$

我们假设  $\epsilon_i$  是独立且同分布的

$$E(\epsilon_i) = 0$$

$$\text{var}(\epsilon_i) = \sigma^2.$$

$$\epsilon_i \sim N(0, \sigma^2)$$

指教族有正态、泊松、二项等。



为进行模型拟合和推断，

应该使用最大似然 (MLE) (Maximum likelihood Estimation)

GLM 的三个组成部分:  $(y, \phi)$

- ① Random Component  $\Rightarrow$  response variable  $y$  and assume a probability distribution for it  
这里  $y$  是假设来自指数族分布的
- ② Systematic  $\Rightarrow$  specify 解释  $\phi$  what the predictor  $(x_1, x_2, \dots)$   $\Rightarrow$  (线性) manner  $\Rightarrow$  Predictor.  
或  $\beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$
- ③ Link  $\Rightarrow$  the relation between the mean and expected value of random component and systematic component  
 $\eta = \mu = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$   
linear / non-linear

Link function  $g(\cdot)$ :

连接  $x$  和  $\mu$

$$g(\mu_i) = \eta_i$$

$$g^{-1}(\eta_i) = \mu_i$$

$$\text{满足 } g(E[y_i|x_i]) = \sum_{j=1}^p \beta_j x_{ij} \text{ 或 } (or X\beta) \quad \text{linear predictor}$$

$$g(E[y|x]) = X\beta.$$

$$E[y_i|x_i] = g^{-1}(X\beta)$$

$n_i$ : 试验次数  $p_i$ : 成功概率

$y_i/n_i \Rightarrow$  成功的相对比例

$$E[y_i] = p_i \Rightarrow \text{期望}$$

$$\text{Var}\left(\frac{y_i}{n_i}\right) = \frac{1}{n_i} p_i (1-p_i)$$

$$V(\mu_i) = \mu_i (1-\mu_i)$$

$$g(\mu_i) = \log\left(\frac{\mu_i}{1-\mu_i}\right)$$

↓  
be.

$$\Rightarrow \begin{cases} \text{对于正态分布的 general} \\ \eta_i = \beta_0 + \beta_1 x_{i1} + \dots \\ \text{link: } g(\mu_i) = \eta_i \\ \text{variance } V(\mu_i) = 1 \end{cases}$$

$$\phi V(\mu)$$

离散系数  
discrete  
dispersion parameter

泊松  
Poisson

$$E[y_i] = \lambda_i$$

$$\text{Var}[y_i] = \lambda_i$$

$$g(\mu_i) = \log(\lambda_i)$$

$$\|U\|_p^p = \sqrt[p]{\|\beta_1\|^p + \dots + \|\beta_p\|^p}$$

一般来说 对方差的处理是

$$\log(y_i) = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$\Rightarrow E(\log y_i) = \beta_0 + \beta_1 x_i$$

$$\log E(y_i) = \beta_0 + \beta_1 x_i, \quad V(y_i) = \mu$$

Transformed - data

$$\log(y_i) = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$E(\log y_i) = \beta_0 + \beta_1 x_i$$

GLMs

$$\log E(y_i) = \beta_0 + \beta_1 x_i$$

$$E(y_i) = e^{\beta_0 + \beta_1 x_i}$$

指数家族

$$f(y|\theta, \phi) = \exp \left\{ \frac{y\theta - b(\theta)}{\phi} + c(y, \phi) \right\} \Rightarrow \text{和 variance 相关}$$

$$\begin{cases} E[y] = b'(\theta) \\ \text{var}(y) = b''(\theta) \phi \end{cases}$$

	g(.) canonical if $\theta_i = \eta_i$
Gaussian	PDF CLT MLE NT (正态检验)
Binary	$\Rightarrow$ linear regression odds ratio interpretation
Poisson	$\Rightarrow$ rate ratio interpretation

Normal Distribution:

$$\theta = \mu, \quad b(\theta) = \frac{1}{2}\mu^2, \quad C(y, \phi) = -\frac{y^2}{2\phi^2} - \frac{1}{2}\ln(2\pi\phi^2)$$

$$\begin{aligned} E(y) &= b'(\theta) = \frac{\partial}{\partial \mu} \frac{1}{2}\mu^2 = \mu \\ \text{Var}(y) &= b''(\theta) \phi = \left( \frac{\partial}{\partial \mu} b'(\theta) \right) \phi = \phi^2 = \sigma^2 \end{aligned}$$

这就可以用链式法则求导 因为是嵌套

均值

$$\text{以步对 } -r \log(1-e^\theta)$$

$$h(\theta) = (-e^\theta) g(\mu) = -r \log(\mu)$$

$$\textcircled{1} \text{ 先求 } g(\mu) = g'(\mu) = -r \left( \frac{1}{\mu} \right) \times \mu' = -\frac{r}{\mu}$$

$$\textcircled{2} \text{ 求 } h(\theta) = -e^\theta$$

$$\textcircled{3} \text{ 链式法则 } f(\theta)' = g'(h(\theta)) \times h'(\theta)$$

$$\begin{aligned} &= \frac{+r}{1-e^\theta} \times (e^\theta) \\ &\text{Binomial Distribution} \end{aligned}$$

$$P(y; x) = \binom{n}{y} \pi^y (1-\pi)^{n-y}$$

$$\binom{n}{y} = \frac{n!}{y!(n-y)!}$$

$n \Rightarrow \text{trial}$   $y \Rightarrow \text{success}$

$$\text{推导: } P(y; \pi, n) = \binom{n}{y} \pi^y (1-\pi)^{n-y}$$

$$= \exp \left( \frac{y \log(\frac{\pi}{1-\pi}) + \log(1-\pi)}{n} + \log(n!) \right)$$

$$\boxed{\theta = \log(\frac{\pi}{1-\pi})} \quad \boxed{b(\theta) = \log(1+\exp(\theta))} \quad \boxed{a(\phi) = 1/n}$$

$$C(y, \phi) = \log(n!)$$

$$\text{过程: } \exp \left\{ \log \binom{n}{y} + ny \log \pi + (n-y) \log(1-\pi) \right\}$$

$$= \exp \left\{ \log \binom{n}{y} + ny \log \pi + n \log(1-\pi) - ny \log(1-\pi) \right\}$$

$$= \exp \left\{ \log \binom{n}{y} + n \left[ \log \pi - \log(1-\pi) \right] y + \log(1-\pi) \right\}$$

$$= \exp \left\{ \log \binom{n}{y} + \frac{n \log \frac{\pi}{1-\pi}}{n} y + \log(1-\pi) \right\}$$

$$\begin{aligned} &\because \theta = \frac{\log \pi}{1-\pi} \\ &e^\theta = \frac{\pi}{1-\pi} \end{aligned}$$

$$\begin{aligned} &e^\theta - e^\theta \pi = \pi \\ &\pi = \frac{e^\theta}{1+e^\theta} \end{aligned}$$

$$\begin{aligned} &\log(1+e^\theta) \\ &= \log(1+e^\theta) \end{aligned}$$

Link function:

geometric distribution

几何分布

为什么推导?

概念

12.15

## Partial F-test

Lecture 12:

感兴趣的预测因子

Full (or complete) model: include all predictors of interest

Reduced model: some subset of the full model

Partial F-test is used to compare Nested models

F statistic 是

Lecture 15

simplest test  
3种似然理论的推导

Wald Tests  $\Rightarrow$

inferential approach 极限分配系数

depend on the estimated coefficients  $\phi$  is known!!! and standard errors

Score  $\Rightarrow$

likelihood ratio  $\Rightarrow$

Compare two nested GLMs, are nonlinear functions of  $\beta$  that must be solved iteratively

迭代优化方法

Fisher-scoring Method  $\Rightarrow$

solve Score equations (J<sup>-1</sup>Hessian matrix)

Newton-Raphson Method  $\Rightarrow$

For models with the canonical link  
iterate  
① start with initial estimate  
② calculate working responses and working weights  
③ repeat 2 and 3 till convergence  
Iteratively Re-weighted Least Squares (IRLS)

极大似然：

概念：在所有候选模型中，选择使观察到的数据出现概率最大的模型作为最佳模型。

公式： $\hat{\theta} = \arg \max_{\theta} P(D|\theta)$ ；

所有可能的 $\theta$ 值中选择使似然函数 $P(D|\theta)$ 最大的 $\theta$ 值作为参数 $\theta$ 的估计值。

如果我们想要知道高斯分布的两个参数；  
估计

首先在给定的 $x_1, x_2, \dots, x_n$ 服从 $N(\mu, \sigma^2)$ 的高斯分布；

那么在 $\mu$ 和 $\sigma^2$ 的所有可能值中，极大似然法会选择使 $x_1, x_2, \dots, x_n$ 在高斯分布 $N(\mu, \sigma^2)$ 下出现概率最大 $\mu$ 和 $\sigma^2$ 作为参数估计值。

$$f(x_1, x_2, \dots, x_n | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left\{-\frac{\sum(x_i - \mu)^2}{2\sigma^2}\right\}$$

$$\Pr(x_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_i - \mu)^2}{2\sigma^2}\right\} \Rightarrow$$

two parameters

①.  $\mu$ : 确定正态分布均值的位置。



$\mu_2 > \mu_1$ ; 较大的 $\mu$ 值将分布的均值向右移动。



$\mu_3 > \mu_2 > \mu_1$ ; 较大的 $\mu$ 值将分布的均值向右移动。

②  $\sigma$ : 标准差。

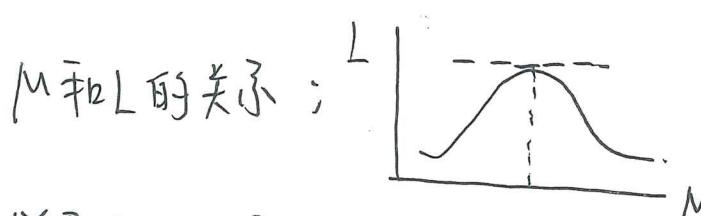


$\sigma$ 较小：正态曲线更高更窄。

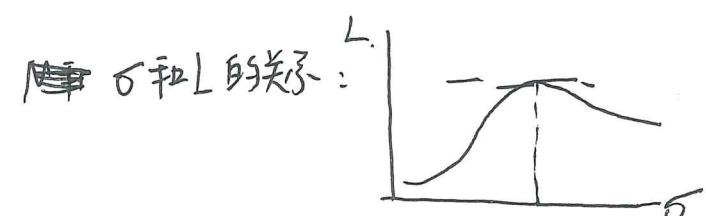
$\sigma$ 较大：正态曲线更低更宽。

现在来计算正态分布的似然：

$$L(N, \sigma | x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \rightarrow \text{找到最佳参数, } \mu \text{ 和 } \sigma.$$



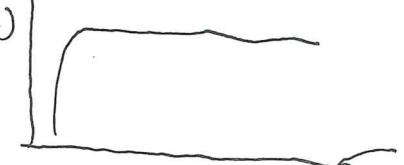
当求解 $\mu$ 的时候，把 $\sigma$ 看作常数；



$$L(\mu, \sigma | x_1, x_2, \dots, x_n) = L(\mu | x_1) \times \dots \times L(\mu | x_n).$$

$$= \underbrace{\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_1 - \mu)^2}{2\sigma^2}\right\}}_{\text{常数}} \times \dots \times \underbrace{\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_n - \mu)^2}{2\sigma^2}\right\}}_{\text{常数}}$$

$\sigma$  和  $\log(L)$  的关系：  $\log(L)$



$L' = 0$  时，找到  $L$  的最大值时 $\mu$ 是多少。

求导前 ①  $\ln\left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_i - \mu)^2}{2\sigma^2}\right\}\right)$

$$1 \text{ (Step 1)} \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} \times e^{-(x_1-\mu)^2/2\sigma^2} \right)$$

$$2) \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \ln \left( e^{-(x_1-\mu)^2/2\sigma^2} \right).$$

$$3) \ln \left[ (2\pi\sigma^2)^{-\frac{1}{2}} \right] - \frac{(x_1-\mu)^2}{2\sigma^2} + \ln(e).$$

$$4) \frac{1}{2} \ln(2\pi\sigma^2) - \frac{(x_1-\mu)^2}{2\sigma^2}$$

$$5) = -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln\sigma^2 - \frac{(x_1-\mu)^2}{2\sigma^2}$$

$$6) = -\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln\sigma^2 - \frac{(x_1-\mu)^2}{2\sigma^2}$$

求导 ( $\sigma$  为常数).

这是化简后我们要求导的表达式:

$$\frac{\partial}{\partial \mu} \ln[L(\mu, \sigma | x_1, \dots, x_n)]$$

$$= 0 - 0 + \frac{2(x_1-\mu)}{2\sigma^2} + \dots + \frac{x_n-\mu}{\sigma^2}.$$

$$= \frac{(x_1+\dots+x_n)-\mu n}{\sigma^2}.$$

$$= \frac{1}{\sigma^2} [(x_1+\dots+x_n)-\mu n].$$

求导 ( $\mu$  为常数).

$$\frac{\partial}{\partial \sigma} \ln[L(\mu, \sigma | x_1, \dots, x_n)].$$

$$= 0 - \frac{n}{\sigma} - \frac{(x_1-\mu)^2}{2\sigma^3} - \dots - \frac{(x_n-\mu)^2}{2\sigma^3}$$

常数.

$$= -\frac{n}{\sigma} + \frac{(x_1-\mu)^2}{2\sigma^3} + \dots + \frac{(x_n-\mu)^2}{2\sigma^3}.$$

$$= -\frac{n}{\sigma} + (x_1-\mu)^2 \sigma^{-3} + \dots + (x_n-\mu)^2 \sigma^{-3}.$$

$$= -\frac{n}{\sigma} + \sigma^{-3} [(x_1-\mu)^2 + \dots + (x_n-\mu)^2].$$

求导公式:

$$(\log a^x)' = \frac{1}{x \ln a}$$

$$(a^x)' = a^x \ln a.$$

$$(f \cdot g)'(x) = f(x)g'(x) + g(x)f'(x).$$

$$\frac{f(x)'}{g(x)'} = \frac{f(x)g'(x) + f'(x)g(x)}{g(x)^2}.$$

$$\ln[L(\mu, \sigma | x_1, \dots, x_n)]$$

有  $n \geq 2 \Rightarrow (-\frac{1}{2})$ .

$$\begin{aligned} &= -\frac{1}{2} \ln(2\pi) - \frac{\ln(\sigma)}{2} - \frac{(x_1-\mu)^2}{2\sigma^2} - \dots - \frac{(x_n-\mu)^2}{2\sigma^2} \\ &= -\frac{1}{2} \ln(2\pi) - n \ln(\sigma) - \frac{(x_1-\mu)^2}{2\sigma^2} - \dots - \frac{(x_n-\mu)^2}{2\sigma^2} \end{aligned}$$

$$\frac{1}{2\sigma^2} (x_1-\mu)^2.$$

~~$\frac{\partial}{\partial \mu} (x_1^2 - 2x_1\mu + \mu^2)'$~~

$$0 - 2x_1 + 2\mu$$

$$- \frac{\mu - x_1}{\sigma^2} = + \frac{x_1 - \mu}{\sigma^2}$$

$$\frac{\partial}{\partial \mu} \ln[L(\mu, \sigma | x_1, \dots, x_n)] = \frac{1}{\sigma^2} [(x_1+\dots+x_n) - \mu n]$$

$$\frac{\partial}{\partial \sigma} \ln[L(\mu, \sigma | x_1, \dots, x_n)] = -\frac{n}{\sigma} + \frac{1}{\sigma^3} [(x_1-\mu)^2 + \dots + (x_n-\mu)^2]$$



为了找最大  $L$ ; 令  $\frac{\partial}{\partial \mu} \ln = 0$ .

$$\frac{1}{\sigma^2} [(x_1+\dots+x_n) - \mu n] = 0.$$

$$\mu = \frac{x_1+\dots+x_n}{n}$$

因此  $\mu$  的最大似然 ~~估计~~ 是测量值的平均值.

$$\sum \frac{\partial}{\partial \sigma} \ln = 0$$

$$\frac{n}{\sigma} + \frac{1}{\sigma^3} [(x_1-\mu)^2 + \dots + (x_n-\mu)^2] = 0.$$

$$-n + \frac{1}{\sigma^2} [(x_1-\mu)^2 + \dots + (x_n-\mu)^2] = 0.$$

$$\frac{1}{\sigma^2} [(x_1-\mu)^2 + \dots + (x_n-\mu)^2] = n.$$

$$[(x_1-\mu)^2 + \dots + (x_n-\mu)^2] = n\sigma^2$$

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{n}$$

$\sigma$  的最大似然估计是测量值的标准差.

19

## Bayesian Linear regression model

Data:  $D = (x_1, y_1) \dots (x_n, y_n)$

Predict:  $y^*$  for a new unseen  $x^*$

① Likelihood: Define the likelihood

$$\underline{P(Y|\theta, X)}$$

② Prior: Define the prior  $P(\theta)$

③ Learning: Do inference by applying Bayes's theorem

$$P(\theta|Y, X) \propto P(Y|\theta, X)P(\theta)$$

④ Prediction: Compute  $\overset{\text{predictive}}{\text{distribution}}$  by marginalizing 边缘计算.

$$P(y^*|x^*, Y, X) = \int p(y^*|\theta, x^*)p(\theta|Y, X)d\theta$$

For posterior

$$y_i = x_i^T \beta + \epsilon_i \rightarrow \text{BLR}$$

$$\epsilon_i \sim N(0, r^{-1}) \quad r = \sigma^2 \quad \beta \sim P(\beta)$$

$$\left\{ \begin{array}{l} P(Y|\beta) = N(Y; X\beta, \sigma^2 I) \\ P(\beta) = N(\beta; \mu_0, \Sigma_0) \end{array} \right. \begin{array}{l} \text{likelihood} \\ \text{prior} \end{array}$$

VS

predictive 预测

$$P(y^*|x^*)$$

$$P(y^*|\beta) = N(y^*; X\beta, \sigma^2 I)$$

$$P(\beta|Y) = N(\beta; \mu_n, \Sigma_n) \quad \text{是一样的, 除了}$$

$$P(y^*|Y) = N(y^*; M^*, S^*)$$

$$M^* = X^* \mu_n$$

$$S^* = \sigma^2 + X^* \Sigma_n X^*$$

Regularization least squares fitting procedure

$$RSS = \sum_{i=1}^n (y_i - \sum_{j=0}^P \beta_j x_{ij})^2$$

过拟合  $\Rightarrow$  在训练模型上很好  
在其它测试集上表现很差

Regularization: allows complex models to be trained by data sets of limited size without severe overfitting, essentially by limiting the effective

~~Ridge regularization~~

Lasso:  $\Rightarrow$  线性回归方法

$\hookrightarrow$  正则化

在有限的数据集上训练复杂模型 without overfitting

model complexity  
 $\downarrow$  by 限制 effective  
 mode complexity

# Bayesian Linear Regression

学习参数模型就是找到最适合模型训练的值

OLS

MLE

$\Rightarrow$  unknown and fixed

The data were generated randomly

Frequentist approach

$$\alpha > 0 \quad p \sim \text{Beta}(\alpha, \alpha) \Rightarrow \text{prior}$$

$$P(\theta | D) = \frac{P(D | \theta) P(\theta)}{P(D)}$$

↓  
观测数据

$$P(y) = \int P(y, \theta) d\theta = \int p(y|\theta) p(\theta) d\theta$$

$P(\theta)$  prior  $\Rightarrow$  before collected data

$P(\theta|D)$  posterior  $\Rightarrow$  after inferring data

$P(D|\theta)$  likelihood of data

$P(D)$  the marginal likelihood

$P(y_*|y)$  posterior predictive

$P(\theta|y)$  posterior  $P(y)$  marginal likelihood

$P(\theta)$  prior

example 得病

90% 的人吃汉堡

100,000 之一的人有病

假设 - 半的人吃汉堡

$\Rightarrow$  求一个吃汉堡的人患病可能性

① 设  $k_j = \text{有 } k_j \text{ 病}$   $H = \text{吃汉堡}$

$$P(H=\text{yes} | k_j=\text{yes}) = 90\%$$

$$P(k_j=\text{yes}) = 0.001\%$$

$$P(H=\text{yes}) = 50\%$$

$$P(k_j=\text{yes} | H=\text{yes})$$

		H=yes	H=no	=	$P(k_j=\text{yes}, H=\text{yes})$
$k_j=\text{Yes}$	9	1		<del>10000</del> $P(H=\text{yes})$	
	499991	499999		$P(k_j=\text{yes}, H=\text{yes})$	

$$= P(k_j=\text{yes}, H=\text{yes})$$

$\rightarrow$  吃汉堡人的总数  $H=\text{yes}$

$$P(k_j=\text{yes}, H=\text{yes}) + P(k_j=\text{no}, H=\text{yes})$$

MAP maximum a posteriori (MAP)

Bayesian approach + regularisation 正则化

MAP is  $\theta$  后验 posterior reaches its max.

$$\begin{aligned}\hat{\theta}_{\text{MAP}} &= \arg \theta \max p(\theta|y) \\ &= \arg \theta \max p(y|\theta) p(\theta) \\ &= \arg \theta \max [\ln p(y|\theta) + \ln p(\theta)]\end{aligned}$$
$$p(\theta|y) = \frac{p(y|\theta) p(\theta)}{p(y)}$$

$L^2 \Rightarrow$  正则化 regularised linear regression

$$\hat{\theta}_{L^2} = \arg \theta \max [\ln p(y|\theta) + \ln p(\theta)]$$

当  $p(\theta) \propto \|\theta\|_2^2$  时,  $\theta$  的 MAP 和  $L^2$  的正则化 regularised 估计相同

$p(\theta) \sim N(\beta; M_\theta, \Sigma_\theta)$