

Chapter 3: Multiple Linear Regression

3.1 Description of the Data and Model $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$, contains no systematic information for Y that is not already captured. 一系列定義: Least Square: $SSE = S(\beta_0, \beta_1, \dots, \beta_p) = \sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 = ||y - X\beta||^2$, Least Square Solution: $\hat{\beta} = (X'X)^{-1} X'y$. Fitted equation: $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_p X_p$, ordinary least squares residuals: $e_i = y_i - \hat{y}_i$. unbiased estimate of σ^2 is $\hat{\sigma}^2 = \frac{SSE}{n-p-1} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-p-1}$ (for no intercept model, only divide by n-p), called sum of squared residuals (RSS)

從 simple lin reg to multiple: 例如 $\hat{Y} = 15.3276 + 0.7803X_1 - 0.0502X_2$, 想 翻 -0.0502 出黎嘅步驟係: 1. Fit Y vs X_1 , 2. Fit X_2 vs X_1 , 3. Fit e_{Y, X_1} vs e_{X_2, X_1} 解釋: 第一步計到嘅 residual 係 part of Y not linearly related to X_1 , 第二步計到嘅係 part of X_2 not linearly related to X_1 , 第三步係 take out effect of X_1 to give coeff

3.3 **Scaling and Centering**: **Scale**: for w or w/o intercept model, 兩種 scaling 方法:

1. Unit-length Scaling: $\tilde{Z}_j = (X_j - \bar{x}_j)/L_j$, $L_j = \sqrt{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}$, same as y (just replace j by y), \tilde{Z}_j has 0 mean and 1 length. $Cor(X_j, X_k) = \sum_{i=1}^n \tilde{z}_{ij} \tilde{z}_{ik}$
2. Standardizing: $\tilde{X}_j = \frac{X_j - \bar{x}_j}{s_j}$, $s_j = \sqrt{\frac{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}{n-1}}$, Center: only for w/ intercept, $X_j - \bar{x}_j$ to make mean as 0
Original to standardized by: $\beta_j = (s_y/s_j)\hat{\beta}_j$

3.4 Assumptions for Multi Linear Regression
1. $\epsilon_1, \dots, \epsilon_n \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$ 2. $X'X$ is invertible
If assumptions hold, we have 2 results: 1. $\hat{\beta}_j \sim N(\beta_j, \sigma^2 c_{jj})$ 2. $W = SSE/\sigma^2 \sim \chi^2(n-p-1)$, $\hat{\beta}_j$ independent $\hat{\sigma}^2 = \sum (y_i - \hat{y}_i)^2 / (n-p-1)$
Multiple Corr. Coeff: $Cor(Y, Y) = \frac{\sum (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum (y_i - \bar{y})^2 \sum (\hat{y}_i - \bar{\hat{y}})^2}}$, $R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$ % of the total variation in Y can be accounted for by the p predictors
Drawback of R^2 : more variables will yield a larger R^2 . so we have adjusted R^2 : $R_a^2 = 1 - \frac{SSE/(n-p-1)}{SST/(n-1)}$, meaning $R_a^2 = 1 - \frac{n-1}{n-p-1} (1 - R^2)$, 而 R_a^2 不可解作這個

3.5 Inference for Individual Regression Coefficients
 $H_0: \beta_j = \beta_j^0$ vs $H_1: \beta_j \neq \beta_j^0$, reject if $|t_j| = \frac{\hat{\beta}_j - \beta_j^0}{s.e.(\hat{\beta}_j)} \geq t_{(n-p-1, \alpha/2)}$ or $p(|t_j|) \leq \alpha$
Conf. Interv.: $\hat{\beta}_j \pm t_{(n-p-1, \alpha/2)} \times s.e.(\hat{\beta}_j)$
Even stat. not suff., a constant should always included

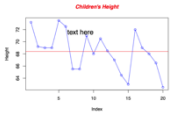
3.6 Test of Hypothesis in Linear Model and Prediction
Framework: Reject H_0 if $F = \frac{[SSE(RM) - SSE(FM)] / (p+1-k)}{SSE(FM) / (n-p-1)} \geq F_{(p+1-k, n-p-1, \alpha)}$, meaning RM does not give as good a fit as FM.
FM has p+1 params (including β_0), $SSE(FM)$'s DF = n-p-1, RM has k params (including β_0), $SSE(RM)$'s DF = n-k
Note that $SSE(RM) \geq SSE(FM)$ since additional params cant increase SSR.
1. $H_0: Y = \beta_0 + \epsilon$ vs $H_1: FM$
 $F = \frac{SSR/p}{SSE/(n-p-1)} = \frac{MSR}{MSE} = \frac{R_p^2/p}{(1-R_p^2)/(n-p-1)}$, since $SSE(RM) = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - \bar{y})^2 - SST(FM)$, $SSE(RM) - SSE(FM) = SST(FM) - SSE(FM) = SSR(FM)$
Table 3.4 Analysis of Variance (ANOVA) Table in Multiple Regression

2. (下列 model 假設有 6 個 params) $H_0: Y = \beta_0 + \beta_1 X_1 + \beta_3 X_3 + \epsilon$, $H_1: FM$, $F = \frac{(R_p^2 - R_q^2)/(p-q)}{(1-R_p^2)/(n-p-1)}$, FM has p params, RM has q params 注意呢到同上面唔同

特殊情况: 如果 $q=p-1$ ($RM=FM-1$), 或者 simple lin re, $F = t_j^2 = \frac{\hat{\beta}_j^2}{Var(\hat{\beta}_j)} < F_{(1, n-p-1, \alpha)}$
3. $H_0: \beta_1 = \beta_3$ (other $\beta = 0$), $H_1: Y: \beta_0 + \beta_1 X_1 + \beta_3 X_3 + \epsilon$
呢題做法: $RM: Y = 0.5X_1 - 0.5X_3 \sim 1$, $FM: Y \sim X_1 + X_3$, and use $p=2, k=1$ do F test.
4. $H_0: \beta_1 = \beta_3$, other $\beta = 0$, $H_1: Y: \beta_0 + \beta_1 X_1 + \dots + \beta_6 X_6 + \epsilon$, RM 同上面一樣, FM 係全部 6 個 var, use $p=6, k=1$
=== Prediction ===
Prediction interval: $\hat{y}_0 \pm t_{(n-p-1, \alpha/2)} \text{ s.e. } (\hat{y}_0)$, s.e. $(\hat{y}_0) = \hat{\sigma} \sqrt{1 + x_0' (X'X)^{-1} x_0}$
Confidence interval: $\hat{\mu}_0 \pm t_{(n-p-1, \alpha/2)} \text{ s.e. } (\hat{\mu}_0)$, s.e. $(\hat{\mu}_0) = \hat{\sigma} \sqrt{x_0' (X'X)^{-1} x_0}$

R corner:

```
plot(c = 1:20,
      y = GaltonFamilies$childHeight[1:20],
      type="n",
      col="blue",
      xlab = "Index",
      ylab = "Height")
abline(c=mean(GaltonFamilies$childHeight[1:20]),
       col="red")
text(c = 10, y = 72,
      label = "text here",
      adj=c(1,1),
      size = 1.5)
```


Basic: rev() reverse vector, rbind, cbind=combine vectors, t(X) transpose, matrix(c)/1:9/0, nrow=3, ncol=3, subset(df, df1='remove'), all.equal(a,b, tolerance)
Stat: confint(lm, level=0.99), coef(), resid(), fitted(), rnorm(), runique(), lgen R.V of dist
Qqplot: qqPlot(lm), DTITS: ols_plot_dffits(lm), Cook: ols_plot_cooks_bar(lm), PR: ols_plot_resid_pot(lm), leverage_cal: hatvales(lm), abline(lm, lwd=width)

Chapter 2 - Simple Linear Regression

Chapter 4: Regression Diagnostics: Detection on Model Violations

4.1 - The Standard Regression Assumptions
1. **Linearity**: The model that relates Y to X_1, X_2, \dots, X_p is assumed to be linear in the regression parameters $\beta_0, \beta_1, \dots, \beta_p$, namely: $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \epsilon$
Check linearity assumption: Scatter plot of Y versus X, they should be linear plot
2. **Errors**: The errors $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ are assumed to be independently and identically distributed (iid) normal random variables each with mean 0 and a common variance σ^2
可以拆成四個 assumptions: 2.1: **Normality**: error is normal distributed, 2.2: **errors have mean 0**, 2.3: **Constant Variance**/Homogeneity: errors have same variance σ^2 , 2.4: **Independent-errors** assumptions (如果唔符合就會有 auto-correlation problem): errors are independent of each other
3. Predictors: 3.1: predictors X_1, X_2, \dots, X_p are **nonrandom**, the values $x_{1j}, x_{2j}, \dots, x_{nj}$ are assumed fixed or selected in advanced.
3.2. The values $x_{1j}, x_{2j}, \dots, x_{nj}$ are **measured without error**, 這個很難被滿足
3.3 Predictors are **linear independent** of each other, 為咗 $(X^T X)^{-1}$ exists
4. **Observations**: all observations are equally reliable and have an approximately equal role in determining regression results

4.2 Various Types of **Residuals**: $e_i = y_i - \hat{y}_i$, fitted value $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip}$
or written as: $\hat{y}_i = p_{i1} y_1 + \dots + p_{in} y_n$, $\hat{y} = X(X'X)^{-1} X'y$. and P = highlighted.
Leverage value: P_{ii} 即是 P 的對角線, 如果係 simple lin. reg., $p_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{\sum (x_i - \bar{x})^2}$. 簡單計下得到: $Cov(e) = \sigma^2(I - P)$, $Var(e_i) = \sigma^2(1 - p_{ii})$, 而我地想呢個 $Var(e_i)$ 一樣, 所以需要使用 studentized residual
Internal Studentized Residual (Default, rstandard in R): $r_i = \frac{e_i}{\hat{\sigma} \sqrt{1 - p_{ii}}}$, 呢個不符合 t-distribution
External Studentized Residual (rstudent in R): $r_i^* = \frac{e_i}{\hat{\sigma}_{(i)} \sqrt{1 - p_{ii}}}$, 符合 t 分佈, df = n-p-2, 因為與 $\hat{\sigma}_{(i)}$ 與 e_i 獨立
 $\hat{\sigma}^2 = \frac{SSE}{n-p-1}$, $\hat{\sigma}_{(i)}^2 = \frac{SSE_{(i)}}{n-p-2}$, which is omitting the i-th observation, they are both unbiased estimates of σ^2 .
 r_i 與 r_i^* 的關係 is one-to-one related: $r_i^* = r_i \sqrt{\frac{n-p-2}{n-p-1-r_i^2}}$

4.3 Graphical Methods: Before and After Fitting a Model

咩圖	咩 assumption	點先叫好
Histogram	Normality	Bell shape
Dot/Box Plot	No outlier in observations	No point far away
Pairwise Scatter Plot Matrix, Y against X_1, \dots, X_n (plot before fit, w/corr coeff)	Linearity	呢個圖無用, 因為有可能幾個 predictor 加埋就有 linearity
Q-Q plot (Residual vs normal scores)	Normality assumption, $r_i \sim N(0,1)$	大 sample size 下, d point 貼住條線 of y=x
Scatter $(r_{(j)}, z_{(j-0.5)/n})$	Linearity/Constant Var/Unrelated with X	無 nonlinear/越大越發散
Residual vs Predictors	Error Mean = 0	在某 fitted value 的 mean(residual) ≈ 0
Residual vs Fitted	Constant Variance	Same Spread of residuals
Fitted vs Predictor + Fitted Line	No outlier when fitting line	呢個圖無用 只可以肉眼睇分 outlier
Index plot of internal studentized residuals	No outlier in Y-space (Response Var)	$r_i < 2/3$
Index plot of leverage	No outlier in X-space (p 個 predictors)	$p_{ii} < 2(p+1)/n$
下列方法原理 = Deleting i-th observation and see the change in fitted value*:		
Index plot of Cook's Distance $C_i = \frac{r_i^2}{p+1} \times \frac{p_{ii}}{1-p_{ii}}$	No influential point	$C_i < F(0.5, p+1, n-p-1)$ $C_i < 4/n$ (from online)
Index plot of DFITS _i = $r_i^* \sqrt{\frac{p_{ii}}{1-p_{ii}}}$	No influential point	$ DFITS_i < 2\sqrt{(p+1)/(n-p-1)} \sim 2\text{sqr}t(p/n)$ (online)
Index plot of Hadi, $H_i = \frac{p_{ii}}{1-p_{ii}} + \frac{p+1}{1-p_{ii}} \frac{d_i^2}{1-d_i^2}$	No influential point	No threshold
Potential-Residual Plot, $\frac{p_{ii}}{1-p_{ii}} \text{ vs } \frac{p+1}{1-p_{ii}} \frac{d_i^2}{1-d_i^2}$	No outlier in X/Y space as well as influential point	\nwarrow = outlier in X-space \searrow = outlier in Y-space \nearrow = influential points

*Masking: detect 唔到係 outlier, Swamping: detect 錯咗 as outlier
Math corner for above graphs:
Q-Q plot: $r_{(j)}$ 係 order statistics, $\#\{i: r_i \leq r_{(j)}\} = j$, $\Phi(z_{(j-0.5)/n}) = \frac{j-0.5}{n}$ where $\Phi(\cdot)$ is the c.d.f. of $N(0,1)$, Leverage: $\frac{\sum_{i=1}^n p_{ii}}{n} = tr(P) = tr(X(X^T X)^{-1} X^T) = tr((X^T X)^{-1} X^T X) = tr(I_{p+1}) = p+1$
Ordinary Residual and Leverage 關係: $p_{ii} + \frac{e_i^2}{SSE} \leq 1$, Hadi 個 $d_i = e_i / \sqrt{SSE}$
Cook's Distance Big Formula: $C_i = \frac{\sum_{j=1}^n (\hat{y}_j - \hat{y}_{j(i)})^2}{\hat{\sigma}^2(p+1)}$, $\hat{\sigma}^2$ obtained from LSE using all obser.
Potential Function: $p_{ii} / (1 - p_{ii})$
DFITS Big Formula: $DFITS_i = \frac{\hat{y}_i - \hat{y}_{i(0)}}{\hat{\sigma}_{(0)} \sqrt{p_{ii}}}$
Cal 機 program: 1. 3x3 matrix multi & inverse, 先入 0 代表乘法, 入 1 代表 inverse. 然後由左至右, 上至下輸入第一個 matrix, (inverse 嘅話 X? 會出 determinant) 然後逐欄輸入第二個 matrix, 會逐欄出答案. 如果就咁求 inverse, 就將第二個 matrix 入做 I 就得 | 2. SSReg, SSError, R^2. 先係 Reg lin 入咗然後過黎 prog 2 | 3. stat increase. $\Sigma x^2, \Sigma x, n, \Sigma y^2, \Sigma y$ 保持不變 | 4. (or given 兩點) $\Sigma x^2 = 14, \Sigma x = 6, n = 3, \Sigma y^2 = 38, \Sigma y = 10, \Sigma xy = 23$ 求直線. 先去 reg lin clear stat 然後 (輸入已知兩點), prog4 順住入, 然後 Svar 睇翻 a, b 係乜, 條線就係 $y = ax + b$

Ch1 Intro

