Website of my Linear Model course
http://www.stat.nthu.edu.tw/~swcheng/Teaching/stat5410/

**Question**

# What is Statistics?

| 哈利波特 | Real Life |
|---|---|
| 占卜學 | Statistics |
| 崔老妮 | Statisticians |
| 水晶球 | Data |
| 未來的資訊 | Information |

*aim of statistics*: provide *insight* by means of *data*

# Basic Procedures of Statistics

• Statistics divides the study of data into *five* steps:

Problem Formulation & Modeling (conceptual )

Data Collection

Statistical Modeling (empirical)

Data Analysis

Decision Making

# When to use regression analysis (linear model)?

• <u>Regression</u>: a statistical tool for investigating the "<u>linearity relationship</u>" between $x$ and $y$.

➤ <u>causal</u> relationship: examine the <u>effects</u> of $x$ on $y$, i.e. <u>how</u> the <u>changes in</u> $x$ *result in* the <u>change in</u> $y$

➤ Even where <u>no sensible causal relationship</u> exists between $x$ and $y$, we may wish to <u>relate them</u> by some sore of <u>mathematical equation</u>

Example



(controllable)
$X_1$ • • • $X_m$

**Exp'tal unit** → **System/ Process** → **Output (response, y)**

$Z_1$ • • • $Z_n$
(uncontrollable)

• <u>data type</u> in regression analysis and some <u>terminologies</u>

<u>A</u> {*response*, *output*, *dependent*} variable $Y$ is modeled or explained by $p$ <u>effects/functions</u> of $m$ {*predictor*, *input*, *independent*, *regressor*} variables <u>$X_1,...,X_m$</u>

➤ <u>$Y$</u>: "<u>approximately</u>" <u>continuous</u>

➤ <u>$X_1,...,X_m$</u>: <u>continuous</u> and <u>discrete</u> (quantitative), <u>categorical</u> (qualitative) (**Q**: <u>example</u>?)

➤ $p$=1, <u>simple</u> regression; $p$>1, <u>multiple</u> regression

➤ $X_1,...,X_m$
    all quantitative         ⇒ multiple regression
    quantitative+qualitative ⇒ analysis of covariance
    all qualitative          ⇒ analysis of variance (ANOVA)

➤ <u>more than one</u> $Y$, <u>multivariate regression</u>

- linear model:

$$Y = \sum_{i=0}^{p} \beta_i \cdot g_i(X_1, \ldots, X_m) \quad + \quad \epsilon$$

$$\boxed{\begin{array}{l} \mathrm{E}(Y|X_1, \ldots, X_m) \\ = \sum_{i=0}^{p} \beta_i \cdot g_i \\ \mathrm{Var}(Y|X_1, \ldots, X_m) \\ = \mathrm{Var}(\epsilon) = \sigma^2 \end{array}}$$

$\underbrace{\text{deterministic component}}_{} \xrightarrow{} \begin{array}{c}\text{mean} \\ \text{function}\end{array} \qquad \begin{array}{c}\text{error:} \\ \underline{\text{random}} \\ \text{component}\end{array} \xrightarrow{} \begin{array}{c}\text{variance} \\ \text{function}\end{array}$

> $X_1, \ldots, X_m$ are regarded as deterministic, i.e., no random phenomenon (when they are random variables, regard the linear model as conditional on $X_1, \ldots, X_m$)

> $g_0(X_1, \ldots, X_m)$, ..., $g_p(X_1, \ldots, X_m)$: *known* functions of $X_1, \ldots, X_m$, called *effects*

> (unknown) *parameters* $\beta_0$, ..., $\beta_p$ enter linearly

> variation due to random error only appears on $y$-axis

- Rationale: a general model for the relationship between $Y$ and $X_1, \ldots, X_m$, is:

$Y = f(X_1, \ldots, X_m) + \varepsilon,$ where $f$ is **unknown and arbitrary**

> **Note:** # of parameters in $f$ is infinite, usually do not have enough data to estimate $f$ directly (globally), we have to assume that it has some more restricted form

> local approximation of $f$ may be achievable by a linear model

- **Note:** Because the predictors can be transformed and combined in any way, linear models are actually very flexible.

## Matrix representation

- Given the data matrix,

| $Y$ | $X_1$ | $X_2$ | ... | $X_m$ |
|-----|-------|-------|-----|-------|
| $y_1$ | $x_{11}$ | $x_{12}$ | ... | $x_{1m}$ |
| $y_2$ | $x_{21}$ | $x_{22}$ | ... | $x_{2m}$ |
| ... | | ... | | |
| $y_n$ | $x_{n1}$ | $x_{n2}$ | ... | $x_{nm}$ |

a row: one group of observations

a column: one variable
(response or predictor)

- We may write a linear model as follows (*functional form*): for $i = 1, 2, \ldots, n$,

$$y_i = \beta_0 + \beta_1 g_1(x_{i1}, \ldots, x_{im}) + \beta_2 g_2(x_{i1}, \ldots, x_{im}) + \ldots + \beta_{p-1} g_{p-1}(x_{i1}, \ldots, x_{im}) + \varepsilon_i,$$

| $Y$ | **1** | $g_1$ | $g_2$ | ... | $g_{p-1}$ |
|-----|-------|-------|-------|-----|-----------|
| $y_1$ | 1 | $g_{11}$ | $g_{12}$ | ... | $g_{1p-1}$ |
| $y_2$ | 1 | $g_{21}$ | $g_{22}$ | ... | $g_{2p-1}$ |
| ... | | | ... | | |
| $y_n$ | 1 | $g_{n1}$ | $g_{n2}$ | ... | $g_{np-1}$ |

a row: one group of observations

a column: response or effect

where $g_{ij} = g_j(x_{i1}, \ldots, x_{im})$

> the expression is (i) ugly notation (ii) conceptually awkward

> matrix/vector notation is more elegant

| | | $\overset{\beta_0}{\overset{\downarrow}{\times}} 1$ | + | $\overset{\beta_1}{\overset{\downarrow}{\times}} g_1$ | + | $\overset{\beta_2}{\overset{\downarrow}{\times}} g_2$ | + | ... | + | $\overset{\beta_{p-1}}{\overset{\downarrow}{\times}} g_{p-1}$ | + | $\varepsilon$ | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $Y$ | = | | | | | | | | | | | | |
| $y_1$ | = | 1 | + | $g_{11}$ | + | $g_{12}$ | + | ... | + | $g_{1p-1}$ | + | $\varepsilon_1$ | a <u>row</u>: one <u>group</u> of |
| $y_2$ | = | 1 | + | $g_{21}$ | + | $g_{22}$ | + | ... | + | $g_{2p-1}$ | + | $\varepsilon_2$ | <u>observations</u> |
| ... | = | | | | | ... | | | | | + | ... | a <u>column</u>: <u>response</u> or |
| $y_n$ | = | 1 | + | $g_{n1}$ | + | $g_{n2}$ | + | ... | + | $g_{np-1}$ | + | $\varepsilon_n$ | <u>effect</u> |

- <u>Matrix form</u> of the linear model:

$$Y = X\beta + \varepsilon,$$

where

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_n \end{bmatrix}, X = \begin{bmatrix} 1 & g_{11} & \cdots & g_{1p-1} \\ 1 & g_{21} & \cdots & g_{2p-1} \\ \cdots & \cdots & \cdots & \cdots \\ 1 & g_{n1} & \cdots & g_{np-1} \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \cdots \\ \beta_{p-1} \end{bmatrix}, \varepsilon = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \cdots \\ \epsilon_n \end{bmatrix}.$$

and $E(\varepsilon) = 0$ and $var(\varepsilon) = \sigma^2 I$ (**Note**: the assumption that <u>errors</u> are <u>normally</u> <u>distributed</u> is <u>not</u> required at the <u>estimation</u> stage)

# Estimating $\beta$

- (<u>ordinary</u>) <u>least square</u> estimator
  - ➤ assume $\varepsilon$ are (i) <u>uncorrelated</u> (ii) <u>equal variance</u> ($Var(\varepsilon) = \sigma^2 I$)
  - ➤ define the <u>best</u> $\hat{\beta}$ as that <u>minimizes</u> <u>sum of squared error</u>: $\varepsilon^T \varepsilon = \sum_{i=1}^n \epsilon_i^2$

$$\varepsilon^T \varepsilon = (Y - X\beta)^T (Y - X\beta) = Y^T Y - 2\beta^T X^T Y + \beta^T X^T X\beta \qquad (*)$$

$$\Rightarrow \text{a } \underline{\text{second-order}} \text{ polynomial of } \beta$$

  - ➤ One method of <u>finding the minimizer</u> is to <u>differentiate</u> $(*)$ w.r.t. $\beta$ and set the derivatives equal to zero

$$\Rightarrow \frac{\partial}{\partial \beta} \epsilon^T \epsilon = -2X^T Y + 2X^T X\beta = 0$$

  - ➤ By <u>calculus</u>, $\hat{\beta}$ is the <u>solution</u> of

$$X^T X\beta = X^T Y \qquad \Leftarrow \text{ called } \textit{normal equation}$$

  - ➤ assume $X^T X$ is <u>non-singular</u>,

$$\hat{\beta} = (X^T X)^{-1} X^T Y \quad \Rightarrow \quad X\hat{\beta} = \underline{X(X^T X)^{-1} X^T Y} \equiv \underline{HY}$$

  - ➤ *predicted values*: $\hat{Y} = X\hat{\beta} = \underline{HY}$
  - ➤ *residuals*: $\hat{\varepsilon} = Y - X\hat{\beta} = Y - \hat{Y} = \underline{(I-H)Y}$
  - ➤ *residual sum of squares* (<u>RSS</u>): $\hat{\varepsilon}^T \hat{\varepsilon} = [Y^T(I-H)^T][(I-H)Y] = \underline{Y^T(I-H)Y}$

➢ OLS $\hat{\beta}$ results from <u>orthogonal projection</u>, makes sense <u>geometrically</u>

➢ (<u>FYI</u>) if $\varepsilon \sim \underline{N}(\mathbf{0}, \sigma^2 I)$, $\hat{\beta}$ is the <u>maximum likelihood</u> estimator

➢ <u>Gauss-Markov thm</u> states $\hat{\beta}$ is <u>BLUE</u> ("Best" Linear Unbiased Estimator)

• <u>mean</u> and <u>covariance matrix</u> of <u>OLS estimator</u> $\hat{\beta}$

$\hat{\beta} = (X^TX)^{-1}X^T\underline{Y}$ is a $p\times 1$ <u>vector</u> of <u>random variables</u>, so

➢<u>mean</u>: $E(\hat{\beta}) = (X^TX)^{-1}X^T\underline{E(Y)} = (X^TX)^{-1}X^T\underline{X}\beta = \beta$ (i.e., *unbiased*)

➢$\underline{Cov(\hat{\beta})} = (X^TX)^{-1}X^T\underline{\sigma^2 I}X(X^TX)^{-1} = \underline{(X^TX)^{-1}\sigma^2}$ ($\Rightarrow$ <u>irrelevant</u> to $\underline{Y}$ and $\beta$ $\Rightarrow$
**Note**: if we can <u>control</u> $\underline{X}$, can <u>decide</u> the <u>var-cov matrix</u> *before* observing $\underline{Y}$ )

Since $\hat{\beta}$ is a random vectors, $\underline{(X^TX)^{-1}\sigma^2}$ is a <u>variance-covariance matrix</u>.

➢$\underline{se(\hat{\beta}_i)} = \sqrt{(X^TX)^{-1}_{ii}}\,\hat{\sigma}$

➢how to calculate the <u>correlation</u> between $\hat{\beta}_i$ and $\hat{\beta}_j$ ?

## Estimating $\sigma^2$

• estimate $\sigma^2$ by $\hat{\sigma}^2 = \underline{\hat{\varepsilon}^T\hat{\varepsilon}}/(n-p) = \underline{RSS}/(n-p)$ $\Rightarrow$ an <u>unbiased</u> estimator

• $\hat{\sigma}^2$ has the <u>minimum</u> <u>variance</u> among <u>all</u> quadratic unbiased estimators of $\underline{\sigma^2}$

• $\hat{\sigma} = \sqrt{RSS/(n-p)}$

• (<u>FYI</u>) if $\varepsilon \sim \underline{N}(\mathbf{0}, \sigma^2 I)$, the <u>maximum likelihood</u> estimator of $\sigma^2$ is $\underline{\hat{\varepsilon}^T\hat{\varepsilon}}/n = \underline{RSS/n}$

## goodness-of-fit: how well does the model fit the data?

• $R^2$, <u>coefficient of determination</u> or **percentage of variance explained**
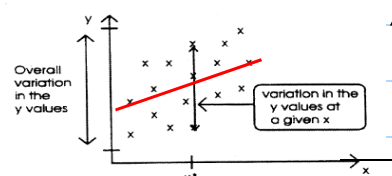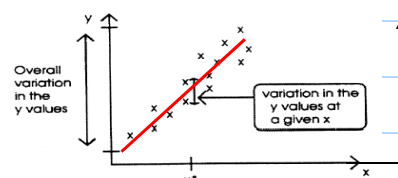
$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum(y_i-\hat{y}_i)^2}{\sum(y_i-\bar{y})^2} = \frac{\sum(\hat{y}_i-\bar{y})^2}{\sum(y_i-\bar{y})^2} = \left(\frac{\sum(y_i-\bar{y})(\hat{y}_i-\bar{y})}{\sqrt{\sum(y_i-\bar{y})^2\sum(\hat{y}_i-\bar{y})^2}}\right)^2$$

$RSS$ is calculated from <u>model</u> with <u>all independent variables</u>,

$TSS$ from model <u>without any independent variables</u>

<u>Interpretation</u> of $R^2$: "<u>proportion</u> of <u>total variation</u> in $y$ that can be <u>explained</u> by the <u>independent variables</u>"

➢ $R$=<u>correlation</u> between $\hat{y}$ and $y$ ; for <u>simple</u> regression, $R$=<u>correlation</u> between $x$ and $y$ (from the <u>geometry</u> viewpoint, ...)

➢ $0 \le R^2 \le 1$, values <u>closer to 1</u> indicate <u>better fits</u>. (what if $\underline{n\approx p}$?)

➢ What is <u>a good</u> value of $R^2$? It <u>depends</u>.

• <u>alternative measure</u> for <u>goodness of fit</u>: $\hat{\sigma}$

➢ it's <u>related</u> to <u>standard error</u> of estimates of $\beta$ and <u>prediction</u>

➢ it's <u>measured</u> in the <u>unit of the response</u> (cf., $R^2$ is <u>free of unit</u>)

# Normality assumption

- **Note**: up till now, **haven't assumed** any distributional form for $\varepsilon$. If we want to perform any hypothesis tests or make any confidence intervals, we will need to do this. The usual assumption is:

$$\varepsilon \sim N(0, \sigma^2 I)$$

➤ model: $Y = X\beta + \varepsilon, \ \varepsilon \sim N(0, \sigma^2 I)$

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

■ **Q**: what does the model describe? e.g.,

$y_x = \beta_0 + \beta_1 x + \varepsilon_x, \ \varepsilon_x$'s $\sim$ i.i.d. $N(0, \sigma^2)$

$\Rightarrow E(y_x) = \beta_0 + \beta_1 x$

$\Rightarrow y_x$'s are independent and $y_x \sim N(\beta_0 + \beta_1 x, \sigma^2)$ at $x = x_i, i = 1, \ldots, n$.

- Some properties of linear models when $\varepsilon \sim N(0, \sigma^2 I)$ :

➤ distribution of $Y$ $[= X\beta + \varepsilon] \sim N(X\beta, \sigma^2 I)$

➤ distribution of $\hat{\beta}$ $[=(X^T X)^{-1} X^T Y] \sim N(\beta, (X^T X)^{-1} \sigma^2)$

➤ distribution of $\hat{\varepsilon}$ $[=(I-H)Y=(I-H)\varepsilon] \sim N(0, (I-H)\sigma^2)$, which has a singular covariance matrix $I-H$ with rank $n-p$ (Note: $\dim(\hat{\varepsilon})=n-p$)

➤ distribution of $RSS$ $[=(n-p)\hat{\sigma}^2 = \hat{\varepsilon}^T \hat{\varepsilon} = \varepsilon^T(I-H)\varepsilon] \sim \sigma^2 \chi^2_{n-p}$

➤ distribution of $\hat{Y}$ $[= X\hat{\beta} = HY] \sim N(X\beta, H\sigma^2)$, which has a singular covariance matrix with rank $p$ (Note: $\dim(\hat{Y})=p$)

➤ $\hat{\beta}$ is independent of $\hat{\sigma}^2$ (Note: $\mathrm{cov}((X^T X)^{-1} X^T Y, (I-H)Y)=0$)

➤ $\hat{Y}$ is independent of $\hat{\varepsilon}$ (Note: $\mathrm{cov}(HY, (I-H)Y)=0$)

➤ distribution of prediction for a new set of predictors, $x_0 = (g_1(x_{10}, \ldots, x_{m0}), \ldots, g_p(x_{10}, \ldots, x_{m0}))^T$

model: $y = \sum_{j=1}^{p} \beta_j \cdot g_j(x_1, \ldots, x_m) + \epsilon$

cf. { fitted model: $\hat{\beta}_j$ $\quad x_0^T \hat{\beta}$

■ mean response v.s. future observation (**Q**: what different?)
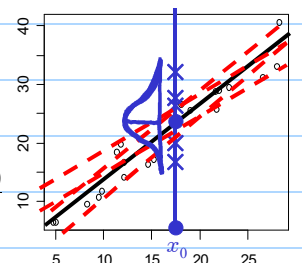
□ Example: average yield when $x=x_0$? and tomorrow's yield when $x=x_0$?

□ same predicted value $x_0^T \hat{\beta}$, but different distributions

■ distribution of prediction error for mean response at $x_0$

$$x_0^T \hat{\beta} - x_0^T \beta \sim N(0, (x_0^T(X^T X)^{-1} x_0)\sigma^2)$$

■ distribution of prediction error for future observations at $x_0$

$$x_0^T \hat{\beta} - (x_0^T \beta + \varepsilon) \sim N(0, (x_0^T(X^T X)^{-1} x_0 + 1)\sigma^2)$$

# hypothesis testings (for $\beta$)

- formulation of hypothesis testing from the view of **comparing models** (**model spaces**)
  - ➢ a model space ≡ the space spanned by columns of some $X$ (model matrix)
  - ➢ consider a large model space, $\Omega$, and a smaller model space, $\omega$, where $\omega \subset \Omega$ (i.e., $\omega$ represents a subset/subspace of $\Omega$). Suppose dimension (# of parameters) of $\Omega$ is $p$ and $\dim(\omega)=q$, where $p>q$.
  - ➢ to answer "which of the model spaces is more adequate" in statistical language $\Rightarrow$ perform the test $H_0$: $\omega$ ($A\beta=c$) v.s. $H_1$: $\Omega\backslash\omega$

- a geometric view of $H_0$: $\omega$ v.s. $H_1$: $\Omega\backslash\omega$

$\hat{\varepsilon}_\omega$ $((n-q)$-dim) $\qquad$ $Y$ $(n$-dim) $\qquad$ $\hat{\varepsilon}_\Omega$ $((n-p)$-dim)

the space that
(1) $\subset \Omega$
(2) $\perp \omega$
(dim $= p-q$)

$\Omega^\perp$ (dim $= n-p$)

$$\left\|\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega\right\|^2$$
$$= \left\|\hat{\varepsilon}_\omega\right\|^2 - \left\|\hat{\varepsilon}_\Omega\right\|^2$$
$$= RSS_\omega - RSS_\Omega$$

$\hat{Y}_\Omega$ $(p$-dim)

$\hat{Y}_\omega$ $(q$-dim)

$\mathbf{0}$

$\hat{\varepsilon}_\omega - \hat{\varepsilon}_\Omega = \hat{Y}_\Omega - \hat{Y}_\omega$ $((p-q)$-dim)

$\Omega$ (dim $= p$) $\qquad$ $\omega$ (dim $= q$)

➢ suppose dimension (# of parameters) of $\Omega$ is $p$ and $\dim(\omega)=q$.

Under the null $H_0$: $\omega$,

$$(RSS_\omega - RSS_\Omega)/\sigma^2 \sim \chi^2_{p-q},$$
$$RSS_\Omega/\sigma^2 \sim \chi^2_{n-p}$$

and they are *independent*.

So, we have $\qquad F = \dfrac{(RSS_\omega - RSS_\Omega)/(p-q)}{RSS_\Omega/(n-p)} \sim F_{p-q,n-p}.$

Therefore, reject if $F > F_{p-q,n-p}^{(\alpha)}$ (usually check if $p$-value $< \alpha$)

➢ **General form**: because the degree of freedom of residuals in a model is the number of observations minus the number of parameters (in $\beta$), this test statistics can be written as:

$$F = \dfrac{(RSS_\omega - RSS_\Omega)/(df_\omega - df_\Omega)}{RSS_\Omega/df_\Omega} \sim F_{df_\omega - df_\Omega, df_\Omega},$$

where $df_\omega = \dim(\omega^\perp) = n-q$ and $df_\Omega = \dim(\Omega^\perp) = n-p$.

➢ The test is widely used in regression and ANOVA. The beauty of this approach is you only need to know the general form.

➢ This test is the likelihood-ratio test.

- Example 1: test of all predictors
  - **Q**: are any of the predictors $g_i$'s useful in predicting the response?
    - $\Omega$: $y = \beta_0 + \beta_1 g_1 + \cdots + \beta_{p-1} g_{p-1} + \epsilon$ , $\dim(\Omega)= p$ , $df_\Omega = n-p$
    - $\omega$: $y = \beta_0 + \epsilon$ , $\dim(\omega)= 1$ , $df_\omega = n-1$
    - $H_0$: $\beta_1 = \beta_2 = \cdots = \beta_{p-1} = 0$   $H_1$: at least one of $\beta_1, \cdots, \beta_{p-1}$ is not zero
    - $RSS_\Omega$: $\hat{\varepsilon}_\Omega^T \hat{\varepsilon}_\Omega = \sum_{i=1}^{n} (y_i - \hat{y}_{i,\Omega})^2$   $RSS_\omega$: $(Y - \bar{Y}\mathbb{1})^T (Y - \bar{Y}\mathbb{1}) = \sum_{i=1}^{n} (y_i - \bar{y})^2$
    - (the overall $F$) $F = \dfrac{(RSS_\omega - RSS_\Omega)/(df_\omega - df_\Omega)}{RSS_\Omega/df_\Omega} = \dfrac{[\Sigma(y_i - \bar{y})^2 - \Sigma(y_i - \hat{y}_i)^2]/(p-1)}{\Sigma(y_i - \hat{y}_i)^2/n-p}$

they are functionally related → cf. $R^2 = 1 - \dfrac{RSS_\Omega}{RSS_\omega} = 1 - 1/(1 + \frac{p-1}{n-p} F)$

- Example 2: testing just one predictor
  - **Q**: Can one particular predictor, say $g_i(\underset{\sim}{x})$, be dropped from the model?
    - $\Omega$: $y = \beta_0 + \cdots + \beta_i g_i + \cdots + \beta_{p-1} g_{p-1} + \epsilon$ , $\dim(\Omega)= p$ , $df_\Omega = n-p$
    - $\omega$: $y = \beta_0 + \cdots + \cancel{\beta_i g_i} + \cdots + \beta_{p-1} g_{p-1} + \epsilon$ , $\dim(\omega)= p-1$ , $df_\omega = n-p+1$
    - $H_0$: $\beta_i = 0$ ($\beta_j \in \mathbb{R}$, for $j \neq i$) $H_1$: $\beta_i \neq 0$ ($\beta_j \in \mathbb{R}$, for $j \neq i$)
    - $F = [(RSS_\omega - RSS_\Omega)/(df_\omega - df_\Omega)]/(RSS_\Omega/df_\Omega) \sim F_{df_\omega - df_\Omega, df_\Omega}$ with $1$, $\hat{\sigma}_\Omega^2$, $1$
  - alternative method $t$-test: $t_i = \hat{\beta}_i/se(\hat{\beta}_i) \sim t_{n-p}$ [Note. $t_i^2 \sim F_{1,n-p}$, and $t_i^2 = F$]

$t_i^2 = \left( \dfrac{\hat{\beta}_{i,\Omega}}{\sqrt{(X_\Omega^T X_\Omega)_{ii}^{-1}}\ \hat{\sigma}_\Omega} \right)^2 = F \Leftarrow RSS_\omega - RSS_\Omega = \left( \dfrac{\hat{\beta}_{i,\Omega}}{\sqrt{(X_\Omega^T X_\Omega)_{ii}^{-1}}} \right)^2$ (exercise)   $\sqrt{\sigma}$

large $se(\hat{\beta}_i)$ | small $se(\hat{\beta}_i)$ | $0$ $\hat{\beta}_i$

- (sequential) ANOVA ($A$: 3 levels; $B$: 4 levels)
  - anova($y \sim 1 + A + B + A:B$)
    1) test $\omega_1$:model 1 ($y \sim 1$) against $\Omega_1$:model 2 ($y \sim 1+A$) [$df_\omega - df_\Omega = 2$]
    2) test $\omega_2$:model 2 ($y \sim 1+A$) against $\Omega_2$:model 4 ($y \sim 1+A+B$) [$df_\omega - df_\Omega = 3$]
    3) test $\omega_3$:model 4 ($y \sim 1+A+B$) against $\Omega_3$:model 5 ($y \sim 1+A+B+A:B$) [$df_\omega - df_\Omega = 6$]
    - $F = \dfrac{(RSS_\omega - RSS_\Omega)/(df_\omega - df_\Omega)}{RSS_{\text{model 5}}/df_{\text{model 5}}} \sim F_{df_\omega - df_\Omega,\ df_{\text{model 5}}}$
    - invariant to the choice of dummy variables if they generate same $\omega$ and $\Omega$
  - ANOVA could have different results when the order of effect sequence is changed, e.g., anova($y \sim 1 + B + A + A:B$):
    $\alpha$) test $\omega_1$:model 1 ($y \sim 1$) against $\Omega_1$:model 3 ($y \sim 1+B$) [$df_\omega - df_\Omega = 3$]
    $\beta$) test $\omega_2$:model 3 ($y \sim 1+B$) against $\Omega_2$:model 4 ($y \sim 1+B+A$) [$df_\omega - df_\Omega = 2$]
    $\chi$) test $\omega_3$:model 4 ($y \sim 1+B+A$) against $\Omega_3$:model 5 ($y \sim 1+B+A+A:B$) [$df_\omega - df_\Omega = 6$]

$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$
$H_0^3: \beta_3 = 0$   $H_0^2: \beta_2 = 0$
$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$   $H_0: \beta_2 = \beta_3 = 0$   $Y = \beta_0 + \beta_1 X_1 + \beta_3 X_3$
$H_0^{2'}: \beta_2 = 0$   $H_0^{3'}: \beta_3 = 0$
$Y = \beta_0 + \beta_1 X_1$

# Confidence intervals and regions

- Model: $\underline{Y = X\beta + \varepsilon}$, $\underline{\varepsilon \sim N(\mathbf{0}, \sigma^2 I)}$; $\underline{\hat{\beta}}$ : OLS estimator $\Rightarrow \underline{\hat{\beta} \sim N(\beta, (X^TX)^{-1}\sigma^2)}$

  ➢ Confidence region for $\underline{A\beta}$, where $\underline{A}$ is a full rank $d\times p$ matrix and $\underline{d \leq p}$

  $A\hat{\beta} \sim N(A\beta, \underline{A(X^TX)^{-1}A^T\sigma^2}) \Rightarrow [(A\hat{\beta}-A\beta)^T[A(X^TX)^{-1}A^T]^{-1}(A\hat{\beta}-A\beta)]/\sigma^2 \sim \chi^2_{\underline{d}}$,

  $$(n-p)\,\underline{\hat{\sigma}^{2}}/\,\sigma^2 \sim \chi^2_{\underline{n-p}},$$

  and they are **independent**.

  $$[(\underline{A\hat{\beta}-A\beta})^T[A(X^TX)^{-1}A^T]^{-1}(\underline{A\hat{\beta}-A\beta})] / (d\,\underline{\hat{\sigma}^{2}}) \sim F_{\underline{d,n-p}}$$

  ➢ $\underline{100(1-\alpha)\%\text{ confidence region}}$ of $\underline{A\beta}$: collection of $A\beta$'s (or $\beta$) that satisfy

  **general form** $\quad [(A\hat{\beta}-\underline{A\beta})^T[A(X^TX)^{-1}A^T]^{-1}(A\hat{\beta}-\underline{A\beta})] / (d\,\hat{\sigma}^{2}) \leq F_{\underline{d,n-p}}^{(\alpha)}$

  The regions are often $\underline{\text{ellipsoidally shaped}}$ (**Q**: $\underline{\text{why}}$?).

- Examples:

  ➢ confidence region $\underline{\text{for }\beta}$, i.e, $\underline{A = I_{p\times p}}$

  $$(\hat{\beta}-\beta)^T X^T X (\hat{\beta}-\beta) \leq (\underline{p}\,\hat{\sigma}^{2})\, F_{\underline{p,n-p}}^{(\alpha)}$$

  ➢ confidence region of $\underline{\beta_i, \beta_j}$, i.e, $\underline{A} = \begin{pmatrix} 0, \cdots, 0, 1, 0, \cdots, 0, 0, 0, \cdots, 0 \\ 0, \cdots, 0, 0, 0, \cdots, 0, 1, 0, \cdots, 0 \end{pmatrix}$

  $$[(\underline{A\hat{\beta}-A\beta})^T[A(X^TX)^{-1}A^T]^{-1}(\underline{A\hat{\beta}-A\beta})] \leq (\underline{2}\,\hat{\sigma}^{2})F_{\underline{2,n-p}}^{(\alpha)}$$

➢ confidence interval $\underline{\text{for }\beta_i}$, i.e, $\underline{A = (0,...,0,1,0,...,0)}$

$$(\hat{\beta}_i-\beta_i)^2/(X^TX)^{-1}_{ii} \leq \hat{\sigma}^{2} F_{\underline{1,n-p}}^{(\alpha)} \Rightarrow |(\hat{\beta}_i-\beta_i)/(\hat{\sigma}\sqrt{(X^TX)^{-1}_{ii}})| \leq t_{\underline{n-p}}^{(\alpha/2)}$$

$\underline{\text{alternative method}}$:

① $\hat{\beta}_i \sim N(\beta_i, \sigma^2(X^TX)^{-1}_{ii})$, ② $(n-p)\hat{\sigma}^{2}/\sigma^2 \sim \chi^2_{n-p}$, and ③ they are $\underline{\text{independent}}$

$$\Rightarrow \underline{(\hat{\beta}_i-\beta_i)/(\hat{\sigma}\sqrt{(X^TX)^{-1}_{ii}}) \sim t_{n-p}} \qquad \Rightarrow \text{C.I.: } \underline{\hat{\beta}_i} \pm \underline{t_{n-p}^{(\alpha/2)}} \times (\underline{\hat{\sigma}\sqrt{(X^TX)^{-1}_{ii}}}).$$

➢ confidence interval for $\underline{\text{prediction of mean response}}$ at $\underline{x_0}$

$$x_0^T\hat{\beta} - x_0^T\beta \sim N(\mathbf{0}, (x_0^T(X^TX)^{-1}x_0)\sigma^2) \Rightarrow (x_0^T\hat{\beta}-x_0^T\beta)/(\hat{\sigma}\sqrt{x_0^T(X^TX)^{-1}x_0}) \sim t_{n-p}$$

$$\Rightarrow \text{C.I.: } \underline{x_0^T\hat{\beta}} \pm \underline{t_{n-p}^{(\alpha/2)}} \times (\underline{\hat{\sigma}\sqrt{x_0^T(X^TX)^{-1}x_0}})$$

➢ C.I. for $\underline{\text{prediction of future observation}}$ at $\underline{x_0}$

$$x_0^T\hat{\beta} - (x_0^T\beta + \varepsilon) \sim N(\mathbf{0}, (x_0^T(X^TX)^{-1}x_0 + \underline{1})\sigma^2)$$

$$\Rightarrow \text{C.I.: } \underline{x_0^T\hat{\beta}} \pm t_{n-p}^{(\alpha/2)} \times (\underline{\hat{\sigma}\sqrt{1+x_0^T(X^TX)^{-1}x_0}})$$

➢ a $\underline{\text{general form}}$ for $\underline{\text{confidence interval}}$:

**estimate $\pm$ (critical value) $\times$ (standard error of estimate)**

# Interpreting parameter estimates

- **Q**: $Y = X\beta + \varepsilon$, what does $\hat{\beta}$ mean?

  Some matters needing attention about $\hat{\beta}$ :
  - $\hat{\beta}$ have units [e.g., fuel consumption data, fitted model:
    fuel $= 154.19 + (-4.23)$Tax $+ (0.47)$Dlic $+ (-6.14)$Income $+ (18.54)\log_2$(Miles)]

  - sign of $\hat{\beta}$ : direction of the relationship between the term and the response
  - interpretation of estimated value (see next two slides)
  - better to also consider values
    contained in its confidence interval
  - causality or association
  - the parameters $\beta$
    - some $\beta_i$'s have physical interpretation, especially those from a conceptual model [e.g., attach weights $x$ to a spring and measure the extension $y$]
      $\Rightarrow$ unfortunately, such cases are rare
    - usually, $\beta_i$'s do not have such physical interpretation
      $\Rightarrow$ in the case, the model $Y = X\beta + \varepsilon$ is only an *empirical model*, i.e., a convenience for representing a complex reality within the range of $X \Rightarrow$
      the real meaning of a particular $\beta_i$ is not obvious, interpretation is difficult

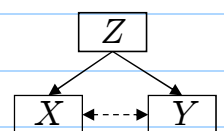- Some interpretations of parameter estimates
  - a naive interpretation:
    **"A unit increase in $X_i$ will *cause* an average change of $\hat{\beta}_i$ in $Y$"** $\Leftarrow$ causality statement
    [e.g., $Y$: annual income, and $X$: years of education]
    - **Q**: what if there exist lurking variables?
      [e.g., $X$: shoe size, $Y$: reading abilities, $Z$: age of child]
      $\Rightarrow$ causal conclusion is doubtful
    - **Q**: what if the roles of predictor and response are mistakenly switched?
      [e.g., $Y$: fire damage, and $X$: numbers of firefighters called out]
    - **Q**: what if some important effects are not included in model?
      - $X$ fixed. $E(\hat{\beta}_1) = \beta_1 + (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_2 \beta_2$
      - $X$ random. true model: $E(Y \mid \mathbf{X}_1, \mathbf{X}_2) = \mathbf{X}_1 \beta_1 + \mathbf{X}_2 \beta_2$ ,
        fitted model: $E(Y \mid \mathbf{X}_1) = \mathbf{X}_1 \beta_1$
        $E(Y \mid \mathbf{X}_1) = \mathbf{X}_1 \beta_1 + E(\mathbf{X}_2 \mid \mathbf{X}_1) \beta_2$
        $Var(Y \mid \mathbf{X}_1) = \sigma^2 + \beta_2^T Var(\mathbf{X}_2 \mid \mathbf{X}_1) \beta_2$
    - even though we have all important variables in the model
      and no lurking variables, there still are problems, e.g.:
      $y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon = \beta_0 + (\beta_1 - \beta_2) X_1 + \beta_2(X_1 + X_2) + \varepsilon$
    - in a properly designed experiment, the naive interpretation is
      more reasonable (because of its use of orthogonal designs and
      randomization); but for observational data, it's often questionable.

➢ an alternative interpretation

" **A unit increase in $X_i$ with all the other (specified) terms _held constant_ will be associated with an average change of $\hat{\beta}_i$ in $Y$** "

- ▪ **Q**: can other terms be held constant? e.g.
  - ▫ $X_1$ and $X_2$ are highly correlated
  - ▫ consider the model $E(Y)=\beta_0+\beta_1 X_1+\beta_2 X_2+\beta_3 X_1 X_2=\beta_0+(\beta_1+\beta_3 X_2)X_1+\beta_2 X_2$
- ▪ it requires the specification of the other terms/effects.

  **Q**: what will happen in the analysis when strong collinearity exists between effects?

  $\Rightarrow$ estimates and tests of $\beta_i$'s may significantly change according to _what other effects are included_. It makes the interpretation almost impossible. In some cases, the problem can be removed by redefining the terms into new linear combinations that may be easier to interpret.

➢ an interpretation from prediction viewpoint

regarding the parameters and their estimates as fictional quantities, and concentrating on prediction enable a rather cautious interpretation of $\hat{\beta}$:

given $(g_{1,0},...,g_{i,0},...,g_{p-1,0}) \rightarrow \hat{y}_0$ , observe $(g_{1,0},...,g_{i,0}\underline{+1},...,g_{p-1,0}) \rightarrow \hat{y}_0 + \hat{\beta}_i$

- ▪ prediction is more stable than parameter estimation
- ▪ directly interpretable and success may be measured in future
- ▪ dangers of extrapolation, be cautious when $x_0$ is outside the range of $X$

## Mean structure

- idea: data are generated from an underlying system, which is assumed to have the form: $y = f(x_1, ..., x_m) + \varepsilon$, where $f$ is **_unknown_**.
- regression _approximates_ the mean structure $f$ by a linear combination of (known) _base functions_ $g_i(x_1,..., x_m)$'s, $i=1, ..., p$, i.e.,

$$f \longleftarrow \sum_{i=1}^{p} \beta_i \cdot g_i(x_1,\ldots,x_m)$$

  ➢ when the structure of $f$ is simple and almost linear, it can be approximated by a simple structure with fewer terms, e.g.,

$$E(y) = f \approx \beta_0 + \beta_1 x_1 + \cdots + \beta_m x_m$$

  - ▪ **Q**: nature is simple?
  - ▪ **Q**: are there sufficient data to support/fit a complex model?

  ➢ when $f$ is complex and non-linear $\Rightarrow$ need more terms to get a good approximation

  - ▪ more parameters, need more degrees of freedom, i.e., more data
  - ▪ e.g., 2 levels, only linear effects; 3 levels, linear and quadratic effects
  - ▪ **Q**: what other complex models?

- base functions for quantitative and qualitative predictors $x_i$'s are defined in different ways

# Polynomial regression

- <u>one predictor</u> case:　　$y = \beta_0 + \beta_1 \underline{x} + \beta_2 \underline{x^2} + ... + \beta_d \underline{x^d} + \varepsilon$

- <u>two predictors</u> $x_1, x_2$ case:

$$y = \beta_0 + \beta_1 \underline{x_1} + \beta_2 \underline{x_2} + \beta_{11} \underline{x_1^2} + \beta_{22} \underline{x_2^2} + \beta_{12} \underline{x_1 x_2} + \varepsilon \quad (d=2, \underline{2^{nd}\text{-order}} \text{ model})$$

➢ the <u>cross-product</u> term $\underline{x_1 x_2}$ can be interpreted as an "<u>interaction</u>" effect, e.g.,

$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2$, where $x_1, x_2 \in \{-1, 1\}$

$\underline{x_1 = +1} \Longrightarrow \overline{E(y)} = \overline{(\beta_0 + \beta_1)} + \overline{(\beta_2 + \beta_3)} \underline{x_2}$

$\underline{x_1 = -1} \Longrightarrow \overline{E(y)} = \overline{(\beta_0 - \beta_1)} + \overline{(\beta_2 - \beta_3)} \underline{x_2}$

➢ models for <u>more predictors</u> can be similarly extended

$$y = \beta_0 + \sum_{i=1}^{m}\beta_{1,i} \underline{x_i} + \sum_{i=1}^{m}\beta_{2,i} \underline{x_i^2} + \sum_{1 \le i < j \le m}\beta_{3,ij} \underline{x_i x_j} + \epsilon$$

- <u>orthogonal polynomials</u>

➢ <u>polynomial terms</u> can cause <u>numerical instability</u> (especially when $d$ large) and <u>collinearity</u>

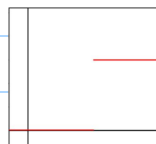➢ example: <u>$2^{nd}$-order model</u>



NTHU STAT 5230, 2025, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

# broken stick (line) regression (segmented regression)

- <u>Recall</u>. polynomial regression: suitable for <u>smooth</u> mean structure, but <u>cannot</u> capture <u>local</u> <u>abrupt change</u> (example?)



- suppose the <u>break</u> occurs at the <u>known</u> value $c$, define the <u>base function</u> (where $c$ is called a <u>knot</u>):
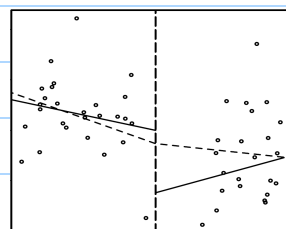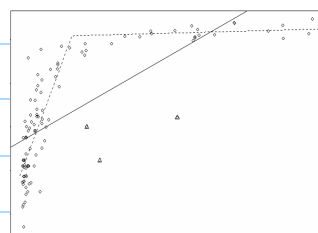
$$d_c(x) = \begin{cases} 1, & \text{if } x > c, \\ 0, & \text{if } x \le c. \end{cases}$$



- model:　$y = \beta_0 + \beta_1 \underline{x} + \beta_2 \underline{(x-c)d_c(x)} + \varepsilon$



$$E(y) = \begin{cases} \beta_0 + \beta_1 x, & \text{if } x \le c, \\ (\beta_0 - \beta_2 c) + (\beta_1 + \beta_2) x, & \text{if } x > c, \end{cases}$$



➢ the <u>two lines</u> <u>meet</u> at <u>$c$</u> $\Rightarrow$ <u>continuous fit</u>

➢ notice only <u>3 parameters</u> in the <u>model</u> $\Rightarrow$ <u>one degree of freedom</u> is <u>saved</u> because of the <u>continuity restriction</u>

# dummy variable (indicator variable, coding)

- categorical (qualitative) predictors
  - ➢ nominal v.s. ordinal
  - ➢ examples: male/female, treatment/control, eye colors, blocks, ...
  - ➢ qualitative in nature:  ▪ values are symbols, no quantitative meaning
    - ▪ no value exist between categories
  - ➢ **Q**: what properties can we explore for qualitative predictor?
    - category $i \to y_{ij}$ , $\mu_i = E(y_{ij}) \Rightarrow$ can only study *difference* between $\mu_i$'s
    - (cf., quantitative predictor)
  - ➢ **Q**: how to fit these predictors into the format of linear regression model
    $Y = X\beta + \varepsilon$? $\Rightarrow$ Ans: dummy variables
- one dichotomous predictor: two categories
  - ➢ for a dichotomous predictor $C$ with two categories $c_1$ and $c_2$, define a dummy
    variable $d$:  $d(C) = \begin{cases} 0, & \text{if } C = c_1 , \\ 1, & \text{if } C = c_2 . \end{cases}$
  - ➢ for a data set with response $y$, one quantitative predictor $x$, and one qualitative
    predictor $C$ (dummy variable $d$), possible models are:
    - model 1: $y = \beta_0 + \beta_1 d + \varepsilon$,  model 2: $y = \beta_0 + \beta_1 x + \varepsilon$,
    - model 3: $y = \beta_0 + \beta_1 d + \beta_2 x + \varepsilon$,  model 4: $y = \beta_0 + \beta_1 x + \beta_2 xd + \varepsilon$,
    - model 5: $y = \beta_0 + \beta_1 d + \beta_2 x + \beta_3 xd + \varepsilon$

- ➢ **Q**: how to interpret $\beta_i$'s in models 1~5?
  - ▪ model 1:  $y = \beta_0 + \beta_1 d + \epsilon$

    $C = c_1$ :  $\mu_1 = E(y|d = 0) = \beta_0$  $\beta_0 = \mu_1$
    $C = c_2$ :  $\mu_2 = E(y|d = 1) = \beta_0 + \beta_1$  $\Rightarrow$  $\beta_1 = \mu_2 - \mu_1$

  - ▪ model 2:  $y = \beta_0 + \beta_1 x + \epsilon$

  - ▪ model 3:  $y = \beta_0 + \beta_1 d + \beta_2 x + \epsilon$

    $C = c_1$ :  $\mu_{1,x} = E(y|d = 0, x) = \beta_0 + \beta_2 x$
    $C = c_2$ :  $\mu_{2,x} = E(y|d = 1, x) = (\beta_0 + \beta_1) + \beta_2 x$

    $\beta_0 = \mu_{1,0}$ (intercept in $c_1$ group)
    $\Rightarrow$ $\beta_1 = \mu_{2,x} - \mu_{1,x}$ (difference of intercepts)
    $\beta_2 = $ slope (same slope in two categories)

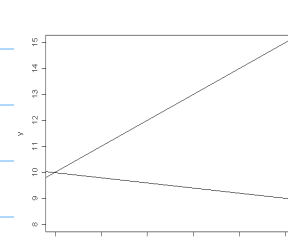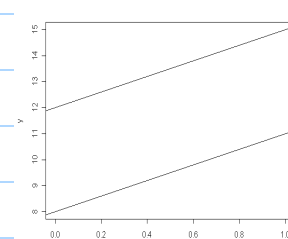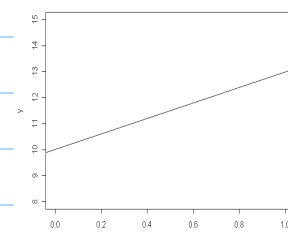  - ▪ model 4:  $y = \beta_0 + \beta_1 x + \beta_2(d \cdot x) + \epsilon$
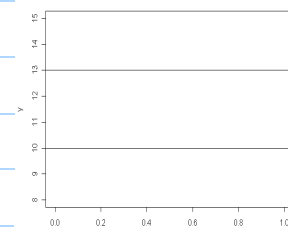    $C = c_1$ :  $\mu_{1,x} = E(y|d = 0, x) = \beta_0 + \beta_1 x$
    $C = c_2$ :  $\mu_{2,x} = E(y|d = 1, x) = \beta_0 + (\beta_1 + \beta_2)x$

    $\beta_0 = \mu_{1,0} = \mu_{2,0}$ (same intercept in two categories)
    $\Rightarrow$ $\beta_1 = $ slope of category $c_1$
    $\beta_2 = $ difference in slopes

- model 5: $y = \beta_0 + \beta_1 \underline{d} + \beta_2 \underline{x} + \beta_3 (\underline{d \cdot x}) + \epsilon$

$\underline{C = c_1}: \quad \underline{\mu_{1,x}} = E(y|\underline{d = 0}, \underline{x}) = \underline{\beta_0 + \beta_2 x}$

$\underline{C = c_2}: \quad \underline{\mu_{2,x}} = E(y|\underline{d = 1}, \underline{x}) = \underline{(\beta_0 + \beta_1)} + \underline{(\beta_2 + \beta_3)x}$

$\Rightarrow \quad$
$\begin{aligned}
\underline{\beta_0} &= \underline{\mu_{1,0}} \text{ (\underline{intercept} of category } \underline{c_1}) \leftarrow \underline{\text{reference}} \\
\underline{\beta_2} &= \underline{\text{slope}} \text{ of category } \underline{c_1} \leftarrow \underline{\text{reference}} \\
\underline{\beta_1} &= \underline{\text{difference}} \text{ in } \underline{\text{intercepts}} \\
\underline{\beta_3} &= \underline{\text{difference}} \text{ in } \underline{\text{slopes}}
\end{aligned}$

➤ alternative coding of dummy variable (better orthogonality)

$$d(\underline{C}) = \begin{cases} \underline{-1}, & \text{if } \underline{C = c_1}, \\ \underline{1}, & \text{if } \underline{C = c_2}. \end{cases}$$

**Q**: how to interpret $\underline{\beta_i}$'s in models 1~5 under this coding?

- model 1: $y = \beta_0 + \beta_1 \underline{d} + \epsilon$

$\underline{C = c_1}: \quad \underline{\mu_1} = E(y|\underline{d = -1}) = \underline{\beta_0 - \beta_1}$
$\underline{C = c_2}: \quad \underline{\mu_2} = E(y|\underline{d = 1}) = \underline{\beta_0 + \beta_1}$
$\Rightarrow \quad \begin{aligned} \underline{\beta_0} &= \underline{(\mu_1 + \mu_2)/2} \\ \underline{\beta_1} &= \underline{(\mu_2 - \mu_1)/2} \end{aligned}$

➤ analysis strategy: start from the full model (model 5) if there are enough degrees of freedom, and then test if some terms can be eliminated

➤ identical methodology applies for more than 2 categories and more quantitative predictors

➤ *ANalysis of COVAriance*: testing model 3 ($\Omega$) against model 2 ($\omega$) (more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called *covariate* and is expected to have the same effect in all categories. The difference between categories is assumed to be an additive effect.

- one polytomous predictor: more than two categories

➤ for $k$ categories, $k-1$ dummy variables are needed to depict the difference between categories (one parameter is used to represent constant term)

➤ various coding of dummy variables: 4 categories $c_1, c_2, c_3, c_4$ example

| treatment coding | | | | Helmert coding | | | | sum coding | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $d_1$ | $d_2$ | $d_3$ | | $d_1$ | $d_2$ | $d_3$ | | $d_1$ | $d_2$ | $d_3$ |
| $c_1$ | 0 | 0 | 0 | $c_1$ | −1 | −1 | −1 | $c_1$ | −1 | −1 | −1 |
| $c_2$ | 1 | 0 | 0 | $c_2$ | 1 | −1 | −1 | $c_2$ | 1 | 0 | 0 |
| $c_3$ | 0 | 1 | 0 | $c_3$ | 0 | 2 | −1 | $c_3$ | 0 | 1 | 0 |
| $c_4$ | 0 | 0 | 1 | $c_4$ | 0 | 0 | 3 | $c_4$ | 0 | 0 | 1 |

➤ consider the model: $y = \beta_0 + \beta_1 \underline{d_1} + \beta_2 \underline{d_2} + \beta_3 \underline{d_3} + \epsilon$

- properties of treatment coding:

$\underline{C = c_1}: \quad \underline{\mu_1} = E(y|\underline{d_1 = 0}, \underline{d_2 = 0}, \underline{d_3 = 0}) = \underline{\beta_0}$
$\underline{C = c_2}: \quad \underline{\mu_2} = E(y|\underline{d_1 = 1}, \underline{d_2 = 0}, \underline{d_3 = 0}) = \underline{\beta_0 + \beta_1}$
$\underline{C = c_3}: \quad \underline{\mu_3} = E(y|\underline{d_1 = 0}, \underline{d_2 = 1}, \underline{d_3 = 0}) = \underline{\beta_0 + \beta_2}$
$\underline{C = c_4}: \quad \underline{\mu_4} = E(y|\underline{d_1 = 0}, \underline{d_2 = 0}, \underline{d_3 = 1}) = \underline{\beta_0 + \beta_3}$

$\Rightarrow \quad \begin{aligned} \underline{\beta_0} &= \underline{\mu_1} \\ \underline{\beta_1} &= \underline{\mu_2 - \mu_1} \\ \underline{\beta_2} &= \underline{\mu_3 - \mu_1} \\ \underline{\beta_3} &= \underline{\mu_4 - \mu_1} \end{aligned}$

- ◻ treats $c_1$ as a reference
- ◻ it is convenient if a "standard" categories exists
- ◻ $d_1$, $d_2$, and $d_3$ are mutually orthogonal, but not orthogonal to constant term

- ■ properties of Helmert coding: $y = \beta_0 + \beta_1 \underline{d_1} + \beta_2 \underline{d_2} + \beta_3 \underline{d_3} + \epsilon$

$C = c_1 :$ $\quad \mu_1 = E(y|d_1 = -1, d_2 = -1, d_3 = -1) = \beta_0 - \beta_1 - \beta_2 - \beta_3$

$C = c_2 :$ $\quad \mu_2 = E(y|d_1 = 1, d_2 = -1, d_3 = -1) = \beta_0 + \beta_1 - \beta_2 - \beta_3$

$C = c_3 :$ $\quad \mu_3 = E(y|d_1 = 0, d_2 = 2, d_3 = -1) = \beta_0 + 2\beta_2 - \beta_3$

$C = c_4 :$ $\quad \mu_4 = E(y|d_1 = 0, d_2 = 0, d_3 = 3) = \beta_0 + 3\beta_3$

$$\Rightarrow \quad \begin{aligned} \beta_0 &= \frac{\mu_1 + \mu_2 + \mu_3 + \mu_4}{4} \equiv \bar{\mu} \\ \beta_1 &= \frac{\mu_2 - \mu_1}{2} \\ \beta_2 &= \frac{\mu_3 - ((\mu_1 + \mu_2)/2)}{3} \\ \beta_3 &= \frac{\mu_4 - ((\mu_1 + \mu_2 + \mu_3)/3)}{4} \end{aligned}$$

- ◻ constant term, $d_1$, $d_2$, and $d_3$ are orthogonal when there are equal # of observations in each categories
- ◻ hard to interpret parameters
- ◻ may suitable for ordinal qualitative predictor

NTHU STAT 5230, 2025, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

- ■ properties of sum coding: $y = \beta_0 + \beta_1 \underline{d_1} + \beta_2 \underline{d_2} + \beta_3 \underline{d_3} + \epsilon$

$C = c_1 :$ $\quad \mu_1 = E(y|d_1 = -1, d_2 = -1, d_3 = -1) = \beta_0 - \beta_1 - \beta_2 - \beta_3$

$C = c_2 :$ $\quad \mu_2 = E(y|d_1 = 1, d_2 = 0, d_3 = 0) = \beta_0 + \beta_1$

$C = c_3 :$ $\quad \mu_3 = E(y|d_1 = 0, d_2 = 1, d_3 = 0) = \beta_0 + \beta_2$

$C = c_4 :$ $\quad \mu_4 = E(y|d_1 = 0, d_2 = 0, d_3 = 1) = \beta_0 + \beta_3$

$$\Rightarrow \quad \begin{aligned} \beta_0 &= \frac{\mu_1 + \mu_2 + \mu_3 + \mu_4}{4} \equiv \bar{\mu} \\ \beta_1 &= \mu_2 - \bar{\mu} \\ \beta_2 &= \mu_3 - \bar{\mu} \\ \beta_3 &= \mu_4 - \bar{\mu} \end{aligned}$$

- ◻ $\beta_0$ represent overall mean
- ◻ compare each category with the overall mean
- ◻ lesser orthogonal

➤ Note: the choice of coding does not affect the $R^2$, $\hat{\sigma}$ and overall $F$-test (to test $H_0$: $\beta_1 = \beta_2 = \beta_3 = 0$, the three codings have same $\omega$ and $\Omega$)

➤ the overall $F$-test is one-way ANOVA (ANalysis Of VAriance)

➤ **Q**: how to work with quantitative predictors? $\Rightarrow$ identical methodology as in 2 categories case. **Q**: how to interpret parameters in the case?

- two qualitative predictors (say, $A$: 3 categories $a_1, a_2, a_3$; $B$: 4 categories, $b_1, b_2, b_3, b_4$)
  - ➤ number of different category combinations = $3 \times 4 = 12$,
    denote their means as $\mu_{ij}$, $i=1, 2, 3$ and $j=1, 2, 3, 4$, i.e.,
    $$y_{ijk} = \mu_{ij} + \epsilon_{ijk}, \quad k = 1, 2, \ldots, n_{ij},$$
    $n_{ij}$ = number of observations in category $A=a_i$ and $B=b_j$
  - ➤ **Q**: how to depict the difference between $\mu_{ij}$'s?
    consider the following linear models:
    - model 1: $E(y_{ijk}) = \beta_0$
    - model 2: $E(y_{ijk}) = \beta_0 + \beta_1 d_1^A + \beta_2 d_2^A$
    - model 3: $E(y_{ijk}) = \beta_0 + \beta_1 d_1^B + \beta_2 d_2^B + \beta_3 d_3^B$
    - model 4: $E(y_{ijk}) = \beta_0 + \beta_1 d_1^A + \beta_2 d_2^A + \beta_3 d_1^B + \beta_4 d_2^B + \beta_5 d_3^B$
      **Q**: how to perform interaction coding? what is interaction?

      *interaction plot*: replace $\mu_{ij}$'s
      by cell means
      $$\bar{y}_{ij} = \sum_{k=1}^{n_{ij}} y_{ijk} / n_{ij}$$

      2-factor interaction
    - model 5:
    $$E(y_{ijk}) = \beta_0 + \beta_1 d_1^A + \beta_2 d_2^A + \beta_3 d_1^B + \beta_4 d_2^B + \beta_5 d_3^B + \sum_{i=1}^{2} \sum_{j=1}^{3} \beta_{ij} d_{ij}$$
    # of parameters: $1 + 2 + 3 + 6 = 12$

- identical methodology applies for more qualitative (3-factor interaction, 4-factor interaction, …) and quantitative predictors (similar modeling to what in LNp.1-26~27)

# Transformation

- transformation of response
  - ➤ Box-Cox transformation family: $t_\lambda(y) = \begin{cases} (y^\lambda - 1)/\lambda, & \text{if } \lambda \neq 0, \\ log(y), & \text{if } \lambda = 0. \end{cases}$
    - $t_\lambda(y)$ is continuous in $\lambda$: for fixed $y > 0$,
      $$\lim_{\lambda \to 0} t_\lambda(y) = \lim_{\lambda \to 0} (y^\lambda - 1)/\lambda = \lim_{\lambda \to 0} (y^\lambda \log(y))/1 = \log(y)$$
    - $\lambda = 1 \Rightarrow$ no transformation, $\lambda = 0 \Rightarrow$ log, $\lambda \neq 0$ or $1 \Rightarrow$ power transformation
    - model: $t_\lambda(y) = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I)$
      - parameters: $\lambda, \beta, \sigma$
      - can write down likelihood for estimation and testing of $\lambda$
      - choice of transformation becomes a estimation/test problem
    - the log-likelihood of $\lambda$ is
      $$L(\lambda) = (-n/2) \, log(RSS_\lambda / n) + (\lambda - 1) \, \Sigma \, log(y_i)$$
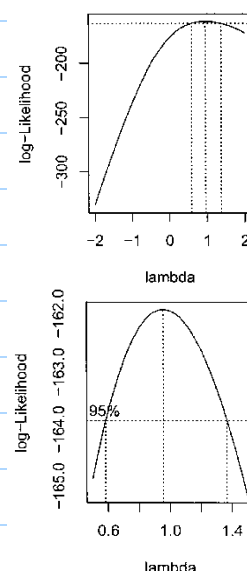      goodness of fit          adjustment
      where $RSS_\lambda$ = residual sum of square when using $t_\lambda(y)$ as response, i.e.,
      $$RSS_\lambda = [t_\lambda(y)]^T (I - H) t_\lambda(y)$$

- ■ estimation of $\lambda$: choose $\lambda$ to fit data well using *maximum likelihood*.
  - ❑ can compute $L(\lambda)$ for various values of $\lambda$ and compute $\hat{\lambda}$ exactly to maximize $L(\lambda)$
  - ❑ but usually $\hat{\lambda}$ is not a nice round number, e.g., $\hat{\lambda} = -0.17$. It would be hard to explain what this new response means.
  - ❑ to avoid this, maximize $L(\lambda)$ over a grid of values, such as $\{-2, -1, -1/2, 0, 1/2, 1, 2\}$. This helps with interpretation.
  - ❑ for $\hat{\lambda}$ outside $[-2, 2]$, pay more attention on whether such transformation is required
  - ❑ **Q**: why not just minimize $RSS_\lambda$ to estimate $\lambda$?
- ■ test of $\lambda$: is the transformation really necessary?
  - ❑ we can answer the question form a C.I. for $\lambda$
  - ❑ *likelihood ratio test* ($H_0$: $\lambda = \lambda_0$ vs. $H_A$: $\lambda \neq \lambda_0$):
    $$-2[L(\lambda_0) - L(\hat{\lambda})] \sim \chi_1^2 \ \text{under } H_0$$
  - ❑ a $100(1-\alpha)\%$ C.I. for $\lambda$ can be formed by:
    $$\{\lambda \mid L(\lambda) > L(\hat{\lambda}) - (1/2)\chi_1^2(1-\alpha)\}$$
  - ❑ is $\lambda=1$ in the C.I.? if so, may as well stay with no transformation.
  - ❑ if rounding $\hat{\lambda}$, check that rounded value is in the C.I.

# Generalized Least Square (GLS)

- model: $Y=X\beta+\varepsilon$, $E(\varepsilon)=0$ and $\text{var}(\varepsilon)=\sigma^2 I \Rightarrow \varepsilon$ uncorrelated and constant variance
- Consider the case $\text{var}(\varepsilon)=\sigma^2\Sigma$, where $\Sigma(\neq I)$ is *known* but $\sigma^2$ is unknown, i.e., we know the *correlation* and *relative variance* between the errors but we don't know the absolute scale
- Because $\Sigma_{n\times n}$ is symmetric and positive definite, we can write $\Sigma=SS^T$, where $S$ is an $n\times n$ nonsigular matrix (by Cholesky or spectral decompositions)

$$Y=X\beta+\varepsilon \Rightarrow S^{-1}Y=S^{-1}X\beta+S^{-1}\varepsilon \Rightarrow Y'=X'\beta+\varepsilon', \text{ where}$$

$$Y'=S^{-1}Y, X'=S^{-1}X, \varepsilon'=S^{-1}\varepsilon, \text{ and}$$

$$E(\varepsilon')=0 \text{ and } \text{var}(\varepsilon')=\text{var}(S^{-1}\varepsilon)=S^{-1}\text{var}(\varepsilon)S^{-T}=S^{-1}\sigma^2 SS^T S^{-T}=\sigma^2 I$$

  $\Rightarrow$ For $Y'$ and $X'$, the assumption in ordinary least square is satisfied

- GLS: find $\beta$ that minimize

$$\varepsilon'^T\varepsilon'=(Y'-X'\beta)^T(Y'-X'\beta)=(Y-X\beta)^T S^{-T}S^{-1}(Y-X\beta)=(Y-X\beta)^T\Sigma^{-1}(Y-X\beta)$$

$$\Rightarrow \hat{\beta}=(X'^T X')^{-1}X'^T Y'=(X^T\Sigma^{-1}X)^{-1}X^T\Sigma^{-1}Y$$

$$\Rightarrow \text{var}(\hat{\beta})=\sigma^2(X'^T X')^{-1}=\sigma^2(X^T\Sigma^{-1}X)^{-1}$$

  GLS is like OLS regressing $Y'=S^{-1}Y$ on $X'=S^{-1}X$

- The practical problem is that $\Sigma$ may not be known. It's usually necessary to make some assumptions and examine the residuals to estimate $\Sigma$ (check IRWLS)

# Weighted Least Square (WLS)

- Sometimes, the errors are uncorrelated, but have unequal variance where the form of the inequality is known ($\Rightarrow \Sigma$ is diagonal, it's a special case of GLS), example:
  - ➤ error variance proportional to a function of predictors [e.g., $\text{var}(\varepsilon_i) = x_i^2 \sigma^2$]
  - ➤ data with replicates, which show a pattern of unequal variance [e.g., $\text{var}(\varepsilon_i) \approx$ sample variance of observations with same $x_i$ ]
  - ➤ the observed $y_i$'s are actually averages of several observations. [e.g., suppose $y_i$ is the average of $n_i$ observations, $\text{var}(\varepsilon_i) = \sigma^2/n_i$ ]

- $\varepsilon$: uncorrelated, but not constant variance $\Rightarrow \Sigma$ is diagonal. Write

$$\Sigma = \begin{pmatrix} 1/w_1 & 0 & \cdots & 0 \\ 0 & 1/w_2 & \cdots & 0 \\ . & . & . & . \\ 0 & 0 & \cdots & 1/w_n \end{pmatrix} \Rightarrow \Sigma^{-1} = \begin{pmatrix} w_1 & 0 & \cdots & 0 \\ 0 & w_2 & \cdots & 0 \\ . & . & . & . \\ 0 & 0 & \cdots & w_n \end{pmatrix}$$
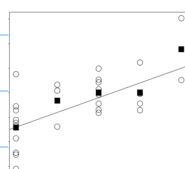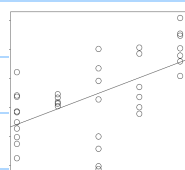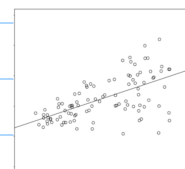
where $w_i$'s ($\propto 1/\text{var}(\varepsilon_i)$) are called *weights*.

low weight $\Leftrightarrow$ high variance; high weight $\Leftrightarrow$ low variance

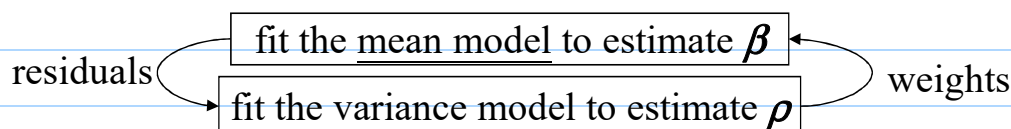- $S = \text{diag}(1/\sqrt{w_1}, ..., 1/\sqrt{w_n})$, then $\Sigma = SS^T$
  - $\Rightarrow$ OLS regress $S^{-1}Y$ (i.e., $\sqrt{w_i}\, y_i$ ) on $S^{-1}X$ ( $\sqrt{w_i}\, x_i$ ) (Note. the column of ones, i.e., intercept needs to be replaced with $\sqrt{w_i}$)
  - $\Rightarrow$ convenient for regression package without a weighted options

- **Q**: Why observations with smaller variance should be multiplied by heavier weight? intuitive interpretation?

- iteratively re-weighted least squares (IRWLS): In all the previous examples, weights (or $\Sigma$ in GLS) are assumed known. **Q**: what if $\text{var}(\varepsilon_i)$ is not completely known, what weights should we use? **Q**: where can you find the information of weights?
  - ➤ model the mean response for $Y$, $\text{E}(Y) = X\beta$
  - ➤ model the variance in $Y$, $\text{var}(Y) = f(X, \rho)$, where $\rho$ are parameters for the variance model

residuals — fit the mean model to estimate $\beta$ — weights — fit the variance model to estimate $\rho$

  - ➤ Example: $\text{var}(\varepsilon_i) = \rho_0 + \rho_1 x_{i1}$
    1. start with $w_i = 1$
    2. use weighted least square to estimate $\beta$
    3. use the residuals to estimate $\rho_0$ and $\rho_1$, perhaps by regressing residuals$^2$ on $x_1$
    4. re-compute the weights and go to 2. Continue until convergence

    Problems: converge? how is the inference about $\beta$ affected? d.f.=? ...etc

  - ➤ alternative approach: jointly estimate the mean and variance parameters using likelihood based method (in R, use `gls()` function in the `nlme` library)