Dr. Azmi Mohd Tamil
Jabatan Kesihatan Masyarakat
Fakulti Perubatan UKM
Based on lecture notes by Dr Azidah Hashim
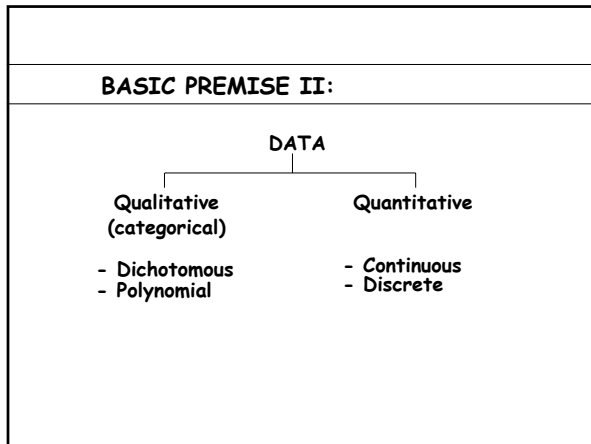
## WHY MULTIVARIATE ANALYSIS?

- Used for analysing complicated data sets
- When there are many Independent Variables (IVs) and/or many Dependent Variables (DVs)
- When IVs and DVs are correlated with one another to varying degrees
- When need to come up with Prediction Model
- Parallels greater complexity of contemporary research

## USING MULTIVARIATE ANALYSIS

- WHICH STATISTICAL PROCEDURE TO USE?
- HOW TO PERFORM CHOSEN PROCEDURE?
- HOW TO INFER FROM RESULTS OBTAINED?
- ANY OTHER ALTERNATIVE APPROACH?

## CHOICE OF APPROPRIATE STATISTICAL METHOD BASED ON:

- Nature of IVs and DVs
- Investigator's Experience
- Personal Preferences
- Ease of Comfort with Methods Used
- Literature Review References
- Consultation with Statistician

## BASIC PREMISE I:

Relationship between

X ⟶ Y

Eg . Smoking ⟶ Lung Cancer
Age ⟶ Hypertension
Maternal ANC ⟶ Birthweight

| INDEPENDENT VARIABLE | DEPENDENT VARIABLE |
|---|---|
| RISK FACTOR | OUTCOME |

## BASIC PREMISE I:

TERM USED:-

| X | Y |
|---|---|
| Independent Variable (IV) | Dependent Variable (DV) |
| Predictor | Outcome |
| Explanatory | Response |
| Risk Factor | Effect |
| Covariates (Continuous) | |
| Factor (Categorical) | |
| Control | |
| Confounders | |
| Nuisance | |

## BASIC PREMISE II:

DATA

Qualitative (categorical)
- Dichotomous
- Polynomial

Quantitative
- Continuous
- Discrete

---

## TERMINOLOGY:

**Univariate Analysis**
- Analysis in which there is a single DV

**Bivariate Analysis**
- Analysis of two variables
- Wish to simply study the relationship between the variables

**Multivariate Analysis**
- Simultaneously analyse multiple DVs and IVs

---

## BASIC PREMISE 111: VARIATIONS OF THE SAME THEME

THE GENERAL LINEAR MODEL:

$Y = a + b_1x_1 + b_2x_2 + b_3x_3 \ldots + b_ix_i$

Used in the following procedures:

ANOVA

ANCOVA

Multiple Linear Regression

Multiple Logistic Regression

Log Linear Regression

Discriminant Function

---

## Rough Guide to Multivariate Methods (1)

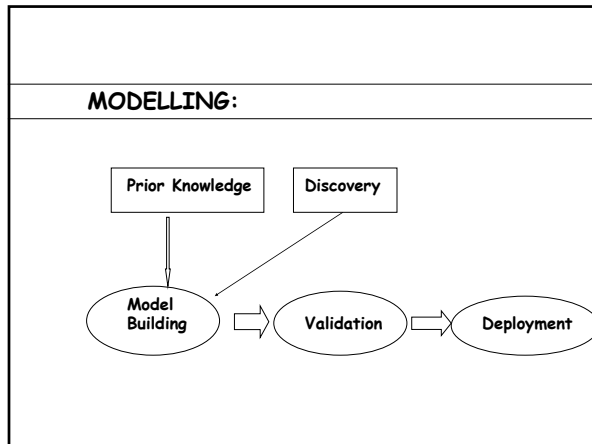| Name | Xs | y |
|---|---|---|
| Regression and Correlation | Continuous (eg. age) | Continuous (eg. BP) |
| Analysis of Variance (ANOVA) | Categorical (eg. SES) | Continuous (eg. BP) |
| Analysis of Covariance (ANCOVA) | Categorical and Continuous (eg age and SES) | Continuous (eg BP) |

---

## Rough Guide to Multivariate Methods (2)

| Name | Xs | y |
|---|---|---|
| Multiple Linear Regression | Continuous (eg. Age, Ht, Wt) | Continuous (eg. BP) |
| Logistic Regression | Continuous (eg. age) | Categorical (eg. CHD) |
| Logistic Regression | Categorical (eg. sex) | Categorical (eg. CHD) |
| Discriminant Function Analysis | Continuous (eg. Age, Income) | Nominal / Ordinal (eg. Quality of Life) |

---

## BASIC PREMISE IV:

**TWO APPROACHES TO CHOOSING:**

1. Based on Type of Modelling

2. Based on Type of Research Question

## MODELLING:

```
  ┌─────────────────┐   ┌──────────────┐
  │ Prior Knowledge │   │  Discovery   │
  └─────────────────┘   └──────────────┘
          │                    │
          ▼                    │
  ╭──────────╮        ╭────────────╮        ╭────────────╮
  │  Model   │   ⇨   │ Validation │   ⇨   │ Deployment │
  │ Building │        ╰────────────╯        ╰────────────╯
  ╰──────────╯
```

## TWO TYPES OF MODELLING TOOLS:

1. Theory Driven-Hypothesis Testing:
   Attempts to substantiate or disprove preconceived ideas
2. Data Driven:
   Automatically creates model based on patterns found in data

## THEORY-DRIVEN MODELLING TOOLS:

CORRELATIONS
t-TESTS
ANOVA
LINEAR REGRESSION
LOGISTIC REGRESSION
DISCRIMINANT ANALYSIS
FORECASTING METHODS

## DATA-DRIVEN MODELLING TOOLS:

CLUSTER ANALYSIS
FACTOR ANALYSIS
DECISION TREES
DATA VISUALISATION
NEURAL NETWORKS

◁

## BASIC PREMISE IV:

### Types of Research Questions

- Degree of Relationship among Variables
- Significance of Group Differences
- Prediction of Group Membership
- Structure

## RESEARCH QUESTION I:

**Degree of Relationship among Variables**

Statistical technique:

a. Bivariate r (Bivariate correlation and Regression)
b. Multiple R ( Multiple Correlation and Multiple Regression)
c. Sequential R
d. Canonical R
e. Multiway Frequency Analysis

## RESEARCH QUESTION II:

**Significance of Group Differences**

<u>Statistical Techniques:</u>
a. t-test
b. One-way ANOVA
c. Two-way ANOVA
d. Profile Analysis

## RESEARCH QUESTION III

**Prediction of Group Membership**

<u>Statistical technique</u>

a. Discriminant Function
b. Multiway Frequency Analysis (Logit)
c. Logistic Regression

## RESEARCH QUESTION IV

**Structure**

<u>Statistical technique:</u>
a. Principal Component Analysis
b. Factor Analysis
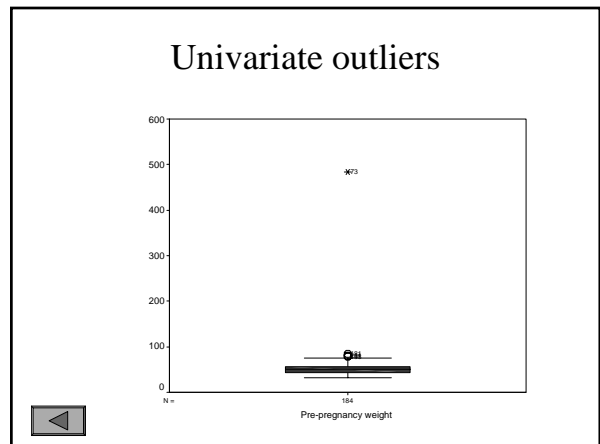c. Structural Equation Modelling

## PRELIMINARY CHECK OF DATA BEFORE MULTIVARIATE ANALYSIS

• <u>Accuracy of Data File</u>
• <u>Honest Correlation</u>
• <u>Missing Data</u>
• <u>Outliers</u>
• <u>Normality, linearity and homoscedasticity</u> ⎱ Regression
• <u>Multicollinearity and Singularity</u> ⎰ Diagnostics
• <u>Common Data Transformations</u>

## ACCURACY OF DATA FILE :

Inspect univariate descriptive statistics for accuracy of input

a. Out-of-range values
b. Plausible means and standard deviation
c. Coefficient of variation
d. Univariate outliers

## Univariate outliers

## HONEST CORRELATIONS

- Inflated Correlation
- Deflated Correlation
- Inaccurately Completed

---

## Inflated Correlation

- If composite variables are to be used and two or more composite variables have the same raw data, <u>correlation</u> can be inflated.
- i.e. correlation between BMI and weight

---

## Deflated Correlation

1) the <u>range</u> of values for one variable is restricted; "relationship between annual average daily traffic count (AADT) and accidents on rural highways and picks a remote region where all AADT's are less than 3000, then if there is a good correlation, he is likely to underestimate it with such a restricted <u>range</u> on one variable.

2) if an intervening variable mediates between two variables; "relationship between thickness of asphalt and chloride content, then picking only those bridges in the <u>population</u> which have a waterproofing membrane will likely push down the estimate. The waterproofing membrane intervenes, literally and statistically. "

---

## Inaccurately Completed

- Questionnaires inaccurately completed due to lack of time, lack of concern, emotional bias will affect correlation.

◁

---

## MISSING DATA

Seriousness depends on
- Pattern of missing data   <   **Random**   **Non-Random**
- How much is missing
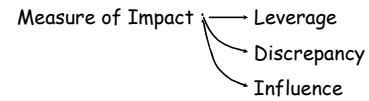- Why is it missing

---

## MISSING DATA

**How to handle missing data**

a. Deleting cases or variables

b. Estimating missing data
   - use of prior knowledge
   - inserting mean values
   - using regression   ◁

c. Using a missing data correlation matrix

d. Treating missing data as data

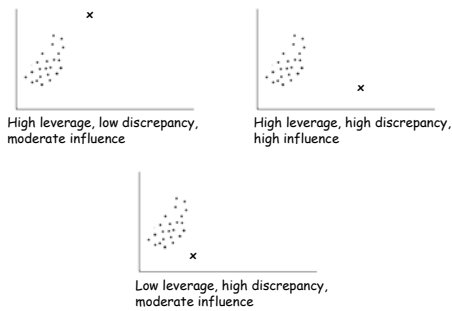e. Repeating analyses with and without missing data

## OUTLIERS

- Cases with such extreme values on one variable or a combination of variables that they distort statistics
- Presence due to
  - incorrect data entry
  - failure to specify missing value codes
  - outlier is not a member of target population
  - distribution for variable in population is more skewed than normal

## OUTLIERS :

Measure of Impact : ── Leverage

Discrepancy

Influence

---

The relationship among leverage, discrepancy and influence



High leverage, low discrepancy, moderate influence

High leverage, high discrepancy, high influence

Low leverage, high discrepancy, moderate influence

## OUTLIERS :

**Dealing with Outliers**

- Find and remedy errors in data entry
- Find and remedy missing values specification
- Deletion
- Retention with alteration

## OUTLIERS :

**Detecting Outliers**

Univariate outliers:

- inspection of z-scores
- graphical methods e.g. histograms, box plots, normal probability plots

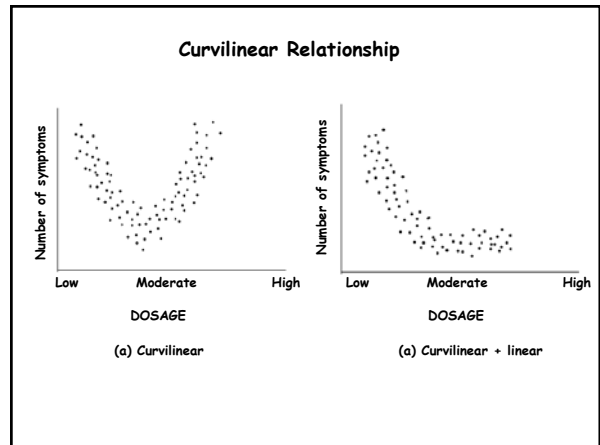Multivariate outliers:

- computation of Mahalanobis distance
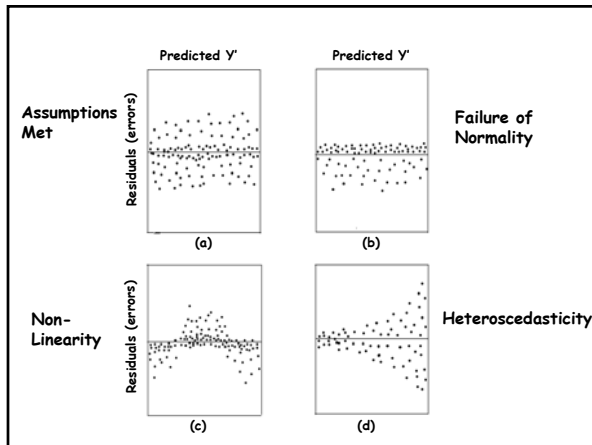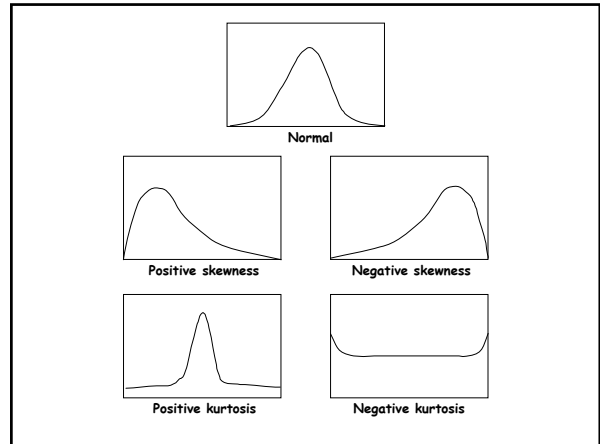
## NORMALITY, LINEARITY AND HOMOSCEDASCITY

**Need for Multivariate Normality**

- assumption that each variable and all linear combinations of the variables are normally distributed
- robustness to violation of assumption still inconclusive

## NORMALITY

- assessed via statistical or graphical methods

- 2 components: skewness and kurtosis

- if non-normal, consider transformation

---

Normal

Positive skewness     Negative skewness

Positive kurtosis     Negative kurtosis

---

Predicted Y'     Predicted Y'

Assumptions Met (a)     Failure of Normality (b)

Residuals (errors)

Non-Linearity (c)     Heteroscedasticity (d)

Residuals (errors)

---

## Curvilinear Relationship

Number of symptoms

Low     Moderate     High
DOSAGE
(a) Curvilinear

Number of symptoms

Low     Moderate     High
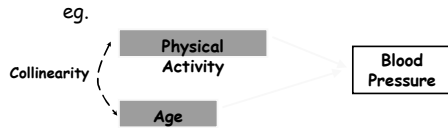DOSAGE
(a) Curvilinear + linear

---

## LINEARITY

- assumption of a straight line relationship between 2 variables

- Nonlinearity is diagnosed from

  • residuals plots or

  • bivariate scatterplots

---

## HOMOSCEDASTICITY

- assumption that the variability in scores for one continuous variable is roughly the same at all values of another continuous variable

- failure due to

  • non normality of one of the variables or

  • one variable is related to some transformation of the other

---

## COLLINEARITY

- concerns the relationship of the IVs to one another and does not directly involve the response variable

eg.

Collinearity → Physical Activity → Blood Pressure

Age

## PRESENCE OF COLLINEARITY CAUSES

- Unstable regression coefficient estimates
- Large estimates of coefficient variances
- Wide 95% Confidence Limits
- Large p-values
- Large Standard errors

## MULTICOLLINEARITY AND SINGULARITY
### (related to a Correlation Matrix)

Multicollinearity : Variables are too highly correlated (> 0.90)

Singularity : Variables are redundant; Matrix cannot be inversed the variables are perfectly correlated.

expose the redundancy of variables and the need to remove variables from the analysis.

## PROBLEMS WITH MULTICOLLINEARITY AND SINGULARITY

- inflate size of error

- weaken analysis

## COMMON DATA TRANSFORMATION
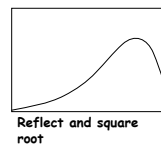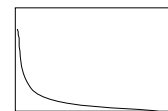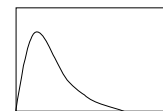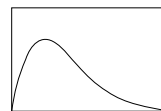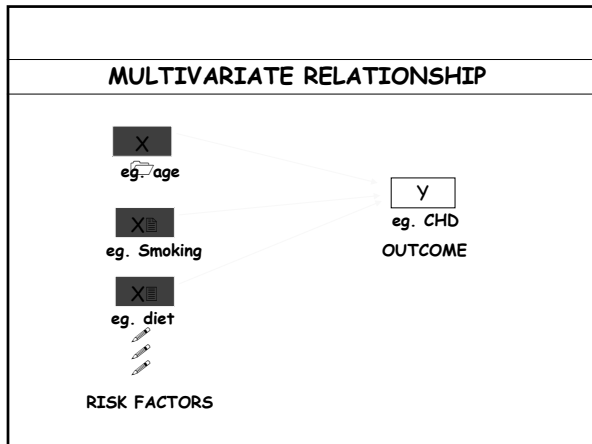
· recommended as a remedy for outliers and for failures of normality, linearity and homoscedasticity

## TYPES OF TRANSFORMATIONS

Square root

Logarithm

Inverse

Reflect and square root

Reflect and logarithm

Reflect and inverse

8

## MULTIVARIATE RELATIONSHIP



| X |
| :---: |
| eg. age |

| X_2 |
| :---: |
| eg. Smoking |

| X_3 |
| :---: |
| eg. diet |

**RISK FACTORS**

| Y |
| :---: |
| eg. CHD |

**OUTCOME**

---

## APPLICATIONS OF MULTIVARIATE ANALYSIS

a) The primary purpose is to study the effect on variable Y of changes in a particular single variable $X_1$, but it is recognised that Y may be affected by several other variables $X_2, X_3, ...$ The effect on Y of simultaneous changes in $X_1, X_2, X_3, ...$ must therefore, be studied.

b) Which of a set of variables $X_1, X_2, X_3, ...$ has apparently most influence on the outcome variable Y.

c) To predict the outcome (i.e. variable Y) in future individuals.

---

## MATHEMATICAL MODEL – the General Linear Model

$\{Outcome\} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + ...$ ⇓

Y                in Linear Regression

Logit P          in Logistic Regression

$\ln \left\{ \dfrac{h(t)}{h_0(t)} \right\}$   in Cox Regression for survival analysis

**$\beta_k$ measures the relationship between Y and that particular $X_k$, adjusting (i.e. controlling) for all the other X's**

---

## GENERAL GOALS IN ANY MULTIVARIATE ANALYSIS

1. Best Fit

2. Most Parsimonious (Occams Razor)
   -use as simple model as possible

3. Biologically reasonable

---

## CENTRAL QUESTION IN ANY MULTIVARIATE ANALYSIS

Does the model that includes the variable(s) in question tell us more about the outcome variable than does a model that does not include those variable(s)?

$Y \sim \boxed{X_1, X_2, X_3} \ \ \vdots X_4, X_5 \vdots$

**Does addition of these 2 add significantly to prediction of outcome variable?**

---

## MULTIPLE REGRESSION

**Relationship**



| X |
| :---: |
| eg. age |

| X_2 |
| :---: |
| eg. weight |

| X_3 |
| :---: |
| eg. physical activity |

| Y |
| :---: |
| eg. SBP |

**Y , $X_1$ , $X_2$ , $X_3$ are continuous variables**

## ANOVA Table for Linear Regression

$\hat{Y}_i - \bar{Y}$ = amount at $X_i$ unexplained by regression

$\hat{Y}_i - \bar{Y}$ = amount at $X_i$ explained by regression

$Y_i - \bar{Y}$ = total amount unexplained at $X_i$

$\hat{y} = \hat{\beta}_0 - \hat{\beta}_1 X_i$

$\hat{y} = \hat{\beta}_0 - \hat{\beta}_1 X$

For one observation, it is clear that :

$$Y_i - \bar{Y} = (\hat{Y}_i - \hat{Y}_i) + (Y_i - \bar{Y})$$

---

**It can be shown that:**

$$\sum_{i=1}^{\blacksquare} (y_i - \boxtimes)^2 = \sum_{i=1}^{\blacksquare} (\hat{y}_1 - y_1)^2 + \sum_{i=1}^{\blacksquare} (\hat{y}_1 - \boxtimes)^2$$

| SSY | = | SSE | + | SSR |
|---|---|---|---|---|
| ‖ | | ‖ | | ‖ |
| Total Sum of Squares about the mean | | Sum of Squares due to error | | Sum of Squares explained by Regression |
| ‖ | | ‖ | | |
| Total Unexplained variation | | Residual Sum of Squares | | |
| | | ‖ | | |
| | | Sum of Squares about the Regression Line | | |

---

## Visual Interpretation of SSY and SSE



SSE represents the variation of the data points seen from A

SSY represents the variation of the data points seen from B

When $\beta_1 = 0$, then SSE = SSY

The steeper the regression line, the greater will be SSY compared to SSE

---

## Visual Interpretation of SSR



SSR may be interpreted as the variation of hypothetical data points sitting exactly on the regression line, when seen from C

---

## General Approach for Assessing the Significance of Predictor Variable(s) in the Model:

i) Compare Observed vs Expected (with)

   = SSE (with)

ii) Compare Observed vs Expected (without)

   = SSE (without)

iii) ✄ = SSE (without) - SSE (with)

---

## Application of the General Approach to Multiple Regression Model

i) **Compare O vs E(with):**

$$\sum_{i=1}^{\blacksquare} \{Y_i - \hat{Y}_i(with)\}^2 = SSE(with)$$

Observed value of response variable y

Predicted( or expected) value of y, according to the model which includes the predictor variable(s) in question

ii) **Compare O vs E(without):**

$$\sum_{i=1}^{\blacksquare} \{Y_i - \hat{Y}_i(without)\}^2 = SSE(without)$$

Observed value of response variable y

Predicted( or expected) value of y, according to the model which excludes the predictor variable(s) in question

iii) **Difference between i) and ii):**

   ⬕ = SSE(without) - SSE(with)

$r^2 = SSR/SSY$

So $r^2$ is the proportion of the total variation which can be explained by the linear regression model.

e.g. for r = 0.5, only 25% of the total observed variation can be explained by the linear regression model. It takes r > 0.7 to make $r^2$ > 0.5, i.e. more than 50% of the total variation explained by the linear regression model.

When r = 0, none of the observed variation can be explained by the linear regression model.

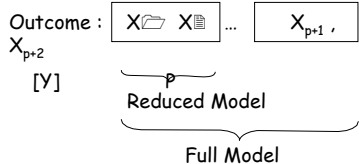When r = 1, all of the observed variation can be explained by the linear regression model.

---

ANOVA Table For Straight-Line Regression

| Source | Sum of Squares (SS) | Degrees of Freedom | Mean square (MS) | Variance Ratio |
|--------|------|------|------|------|
| Regression | SSR | 1 | MSR = (SSR/1) | F = MSR/MSE |
| Residual | SSE | n – 2 | MSE = (SSE/n-2) | |
| Total | SSY | n – 1 | | |

$$\hat{\sigma}^2 = MSE$$

$$r^2 = SSR/SSY$$

---

## BASIC PREMISE

Outcome : $\boxed{X_1 \quad X_2 \quad ... \quad X_{p+1} ,}$ $X_{p+2}$

[Y]  $\underbrace{\qquad}_{p}$  Reduced Model

$\underbrace{\qquad\qquad\qquad}$  Full Model

For Multiple Regression : Use F Statistic

For Logistic Regression : Use LR Statistic (likelihood ratio)

---

## STATISTICAL ASSUMPTIONS

**Assumptions:**

i) $E(Y | x_1, x_2, ... x_k) \equiv \mu_{Y | x_1, x_2, ... x_k}$

$= \beta_0 + \beta_1 x_1 + ... + \beta_k x_k$ {Assumption of Linearity}

ii) $VAR(Y | x_1, x_2, ... x_k) \equiv \sigma^2_{Y | x_1, x_2, ... x_k} = \sigma^2$ {Assumption of Homoscedasticity}

iii) distribution

---

## TESTS available

i) Test for Significant Overall Regression

Y : $X_1 \quad X_2 \quad X_3$

ii) Test for Addition of Single Variable
   - Partial F Test

Y : $\underbrace{X_1 , X_2}_{\text{p=2; predictor variable already in the model}} \quad \underbrace{X_3}_{\text{k=1; additional variable in question}}$

iii) Test for additional of a group of variables
   - Multiple-Partial F test

Y : $\underbrace{X_1 ,}_{\text{p=1; predictor variable already in the model}} \quad \underbrace{X_2 , X_3}_{\text{k=2; additional variables in question}}$

---

## Important Applications of :

A) The Partial F-test

Y : $\boxed{\underbrace{C_1 \quad C_2 \quad .. \quad C_p}_{\text{p controlling variables (confounders)}}} \quad \boxed{\underbrace{S}_{\text{main study variables}}}$

b) The Multiple-Partial F-test

Y : $\underbrace{X_1 , X_2 , X_3}_{\substack{\text{variables already} \\ \text{in the model}}} \quad \underbrace{X_1^2 , X_2^2 , X_3^3}_{\substack{\text{additional "higher order"} \\ \text{variables in question}}}$

Y : $\underbrace{X_1 , X_2 , X_3}_{\substack{\text{variables already} \\ \text{in the model}}} \quad \underbrace{X_1 X_2 , X_1 X_3 , X_2 X_3}_{\substack{\text{additional interaction} \\ \text{variables in question}}}$

## CONFOUNDING AND INTERACTION

Confounding



Confounding is the distortion of a risk factor-disease relationship brought about by the association of other factors with both risk factor and disease, the latter associations with the disease being causal. These factors are called Confounding Factors or "Confounders"

---

## Examples of Confounding



---

## INTERACTION



**Interacting Factor ( = Effect Modifier)**

Interaction exists when the primary relationship of interest between a risk factor and a disease is different at different levels of the interacting factor (also known as effect modifier)

---

## Examples of Interaction



---

## How to Detect Existence of Confounding & Interaction:



$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 (X_1 X_2)$

**Product term**

Statistical testing with the F tests could be used to evaluate the existence of interaction for any given model. In the above example, a partial F test of $H_0 : \beta_3 = 0$ could be used.

---

## CONFOUNDING

Have 2 models

$Y = \beta_0 + \beta_1 x_1 + \varepsilon$

$Y = \beta^*_0 + \beta^*_1 x_1 + \beta^*_2 x_2 + \varepsilon$

with respect to $x_2$,

$\beta_1$ = β crude (ignore $x_2$)

$\beta^*_1$ = β adjusted (adjust for $x_2$)

β crude ≠ β adjusted if confounding is present

12

## CONFOUNDING AND INTERACTION

· Confounding and interaction are different phenomena

· A variable may be both a confounder and an interactor, or only one of the two or neither

· Interaction should be assessed before confounding

· The use of adjusted estimate that controls for confounding is recommended only when there is no meaningful interaction

· If strong interaction is found, an adjustment for confounding is inappropriate. Instead there should be separate results for separate categories of the effect modifier.
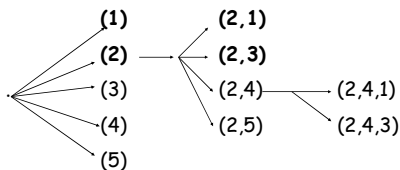
---

## AUTOMATIC ELIMINATION STRATEGIES

With a large number of inter-related predictor variables, it often becomes quite difficult to sort out the meaning of the individual regression coefficients

- Need for Automatic Elimination Strategies

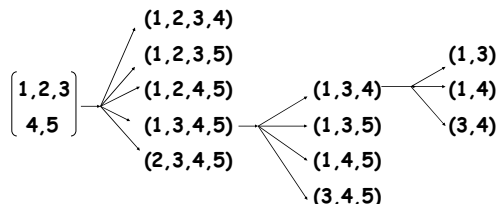---

## AUTOMATIC ELIMINATION STRATEGIES

a) Step-up (=Forward selection) Strategy



The equation starts out empty and IVs are added one at a time provided they meet the statistical criteria for entry

---

## AUTOMATIC ELIMINATION STRATEGIES
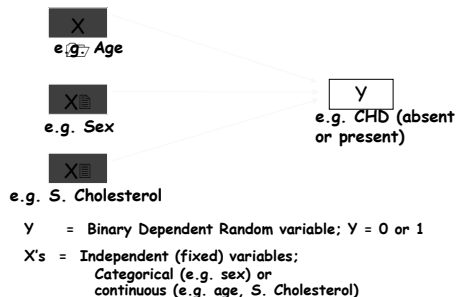
b) Step-down ( = Backward Elimination) Strategy



The equation starts out with all IVs entered and they are deleted one at a time if they do not contribute significantly to regression

---

## AUTOMATIC ELIMINATION STRATEGIES

c) Stepwise Regression

- a compromise between a) and b)

- equation starts out empty and IVs are added one at a time if they meet statistical criteria but they may also be deleted at any step where they no longer contribute significantly to regression.

- considered the surest path to the best prediction equation

---

## Logistic Regression



e.g. Age

e.g. Sex

e.g. CHD (absent or present)

e.g. S. Cholesterol

Y = Binary Dependent Random variable; Y = 0 or 1

X's = Independent (fixed) variables;
Categorical (e.g. sex) or
continuous (e.g. age, S. Cholesterol)

## Mathematical Model for Logistic Regression

a) Explicit form:

$$\Pr\{Y=1\,|\,x\} = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

b) Logit Form

$$\text{Logit } \Pr\{Y=1|x\} = \beta_0 + \beta_1 x$$

$$\equiv \quad Y = \beta_0 + \beta_1 x_1 + \ldots$$

**Remember!   ( $Y = a + bx$ )**

---

## LOGISTIC REGRESSION

**Relationship between Y and X**



(Disease) 1

Pr {Y=1}

(Well) 0 ___ x

This S shaped curve is that of a logistic distribution
- thought to agree with real world situation

---

## Variation of the logistic relationship:

One Independent Dichotomous Variable

X $\longrightarrow$ Y

Exporure (0,1)          Outcome (0,1)

[Independent Variable]     [Dependent variable]

0 = Lower Risk → Nonsmoker     0 = Normal → well
1 = Higher Risk → Smoker       1 = Not normal → sick

---

## Variations of the Logistic Relationship:

One Independent Polytomous Variable

X $\longrightarrow$ Y

Smoking Status              Disease/ Non disease
1 = Nonsmoker
2 = Light Smoker
3 = Heavy smoker

Have to convert to <u>DUMMY VARIABLES</u>

---

## CONSTRUCTION OF DUMMY VARIABLES

|                    | Dummy Variables | |
|--------------------|-----------|-----------|
| **Smoking Status** | Smoke (1) | Smoke (2) |
| 1 = Non smoker     | 0         | 0         |
| 2 = Light smoker   | 1         | 0         |
| 3 = Heavy smoker   | 0         | 1         |

---

## MULTIPLE LOGISTIC REGRESSION MODEL

a) Explicit Form:

$$\Pr\{Y=1|x_1, x_2, x_3, \ldots x_p\} = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \ldots + \beta_p x_p)}}$$

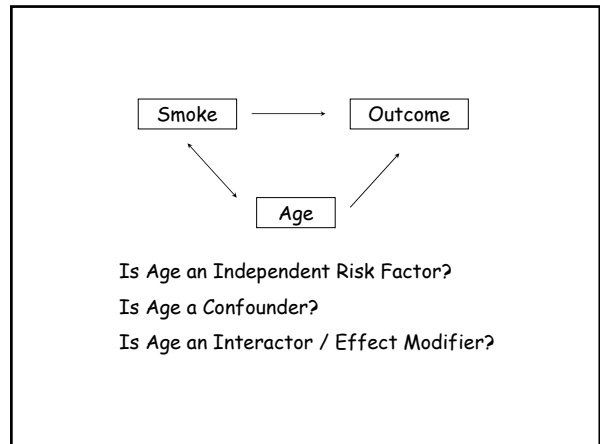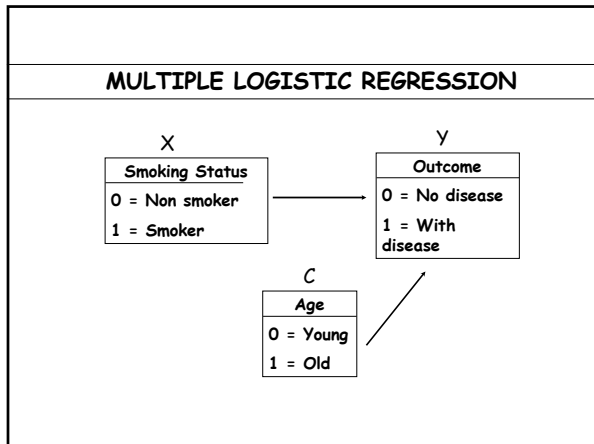B) Logit Form:

$$\text{Logit} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \ldots + \beta_p x_p$$

where Logit P = $\log_e$ ( p / 1- p)

compare with

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \ldots + \beta_p x_p$$

## MULTIPLE LOGISTIC REGRESSION

X

| Smoking Status |
| --- |
| 0 = Non smoker |
| 1 = Smoker |

Y

| Outcome |
| --- |
| 0 = No disease |
| 1 = With disease |

C

| Age |
| --- |
| 0 = Young |
| 1 = Old |

---

Smoke ⟶ Outcome

Age

Is Age an Independent Risk Factor?

Is Age a Confounder?

Is Age an Interactor / Effect Modifier?

---

### Is Age an Interactor / Effect Modifier?

Use Product-Term i.e. Age Smoke

e.g.

$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \ldots$ Reduced model

$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 \ldots$ Full model

Compare Full Model with Reduced Model

---

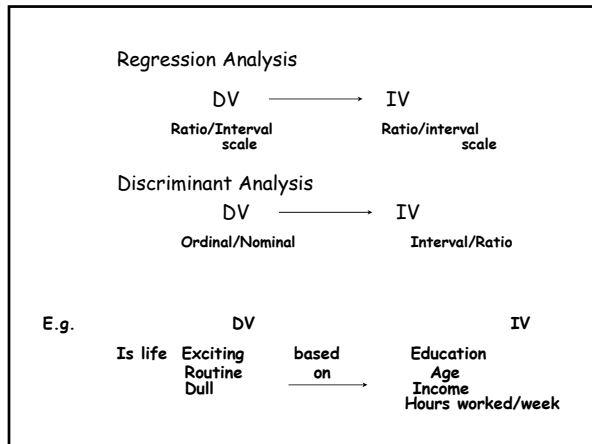### COMPUTER SELECTION OF PREDICTORS

- Forward LR
- Backward LR
- Stepwise

---

### OTHER MULTIVARIATE ANALYSIS

- **Discriminant Function Analysis**
- **Factor Analysis including Principal Component Technique**
- **Cluster Analysis**

---

### DISCRIMINANT FUNCTION ANALYSIS

- A technique for deciding into which category of a variable a case is most likely to fall i.e. to predict group membership

- Compute " discriminant scores" for each case to predict what group it is in

- Normally, have only two discriminant groups

## Slide 1

Regression Analysis

DV $\longrightarrow$ IV

Ratio/Interval scale        Ratio/interval scale

Discriminant Analysis

DV $\longrightarrow$ IV

Ordinal/Nominal        Interval/Ratio

E.g.        DV        IV

Is life   Exciting    based    Education
       Routine     on      Age
       Dull $\longrightarrow$ Income
                  Hours worked/week

## Slide 2

### FACTOR ANALYSIS

· **A technique for condensing many variables into a few underlying constructs**

·**Identify Unifying Concepts**

e.g.     From a 100 item test, we can allocate to 4 distinct abilities:

            Verbal skills

            Mathematical aptitude

            Reasoning ability

            Perceptual speed

## Slide 3

### FACTOR ANALYSIS

· **In SPSS, by default, uses PRINCIPAL COMPONENT Technique to extract factors**

· **Other extraction methods:**

       Principal-axis factoring

       Unweighted least squares

       Maximum likelihood

       Alpha method

       Image factoring

## Slide 4

### STEPS IN FACTOR ANALYSIS

1. Compute correlation matrix
2. Factor extraction
3. Rotation
4. Compute scores for each factor

## Slide 5

### FACTOR ANALYSIS : 2 TYPES

· Exploratory
· Confirmatory

## Slide 6

### CLUSTER ANALYSIS

· Used to find natural groupings within data

· Identify similarities and differences among them

· Use "distances" to reflect similarity and/or dissimilarity

· Multidimensional scaling also uses this concept

## LOGLINEAR REGRESSION

- Extension of CrossTabulation and Chi Square Statistic for Independence
- Difficulty in crosstabulating for > 2 variables
- Use a LOGLINEAR Model

## OTHER SPSS STATISTICS FUNCTIONS

- General Linear Model
- Reliability Analysis
- Multidimensional Scaling
- Probit Analysis
- Survival Analysis
- ETC

## GENERAL LINEAR MODEL

- GLM Factorial Analysis
- GLM Univariate Analysis - combination of Regression and ANOVA
- GLM Multivariate Analysis
- GLM Repeated Measures
- Variance Components

## THANK YOU