

Basic biostatistics with SPSS Looking beyond numbers

Jamalludin Ab Rahman MD MPH

Associate Professor

Department of Community Medicine
Kulliyyah of Medicine



Content

1. Describe data
2. Check distribution
3. Compare two means
4. Compare more than two means
5. Compare proportions
6. Non-parametrics
7. Correlation

[2]



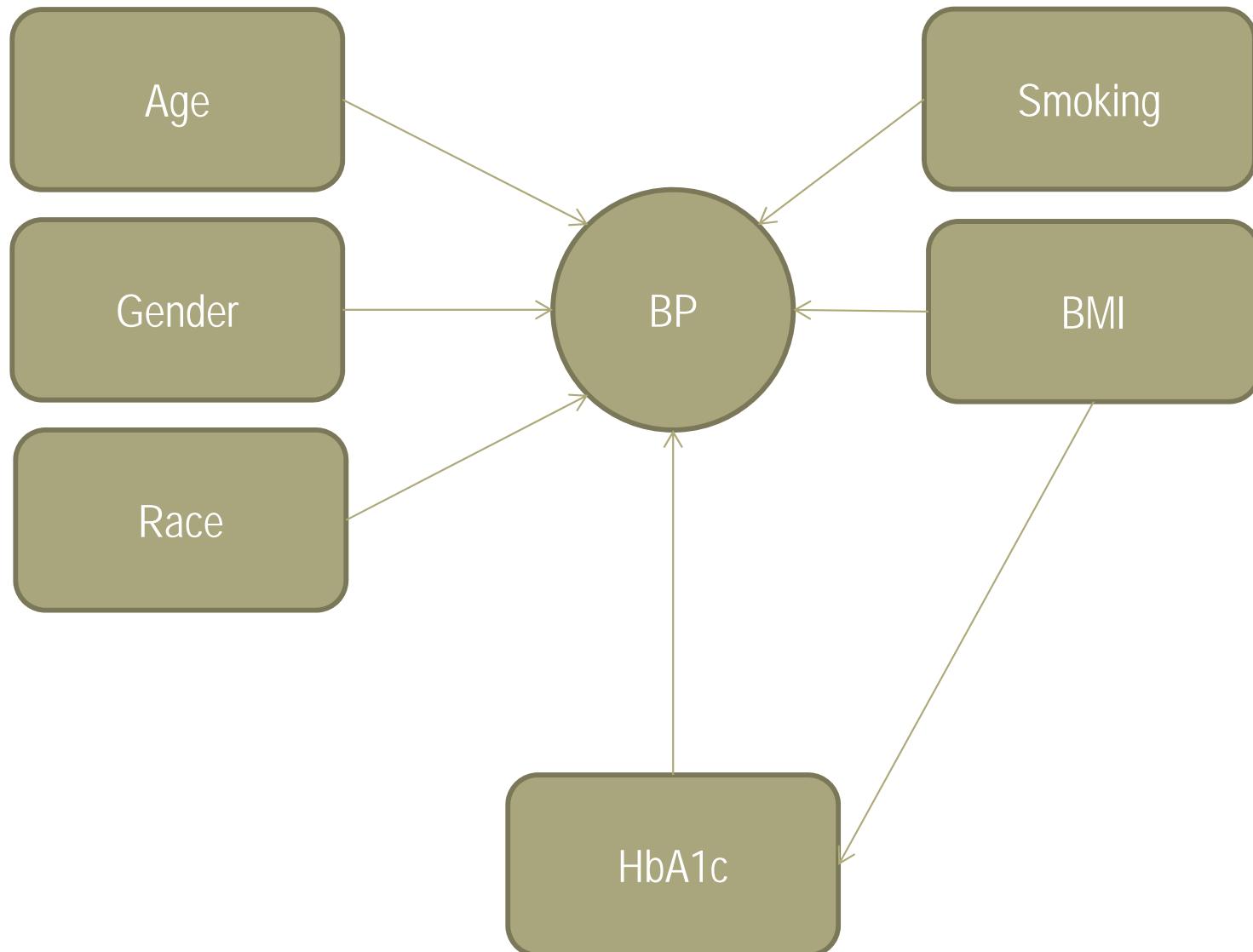
Preliminary

- Use bp2.sav
- Hypothetical data consists of

<i>Variable</i>	<i>Position</i>	<i>Label</i>	<i>Measurement Level</i>
id	1	<none>	Scale
age	2	Age (years)	Scale
gender	3	Gender	Nominal
race	4	Race	Nominal
wt	5	Weight (kg)	Scale
ht	6	Height (cm)	Scale
smoking	7	Smoking status	Nominal
sbp	8	SBP	Scale
dbp	9	DBP	Scale
hba1c	10	HbA1c (%)	Scale



Conceptual framework



Basic biostatistics with SPSS

DESCRIBE DATA & CHECK DISTRIBUTION

[5]



Describe numerical data (*age*)

1

2

3

4

5

6

7

8

www.jamalrahman.net 17/12/2012

Case Processing Summary

	Cases					
	Valid		Missing		Total	
	N	Percent	N	Percent	N	Percent
Age (years)	450	100.0%	0	0.0%	450	100.0%

Check for missing values

Descriptives

		Statistic	Std. Error
Age (years)	Mean	36.57	.534
	95% Confidence Interval for Mean	Lower Bound Upper Bound	35.52 37.62
	5% Trimmed Mean	36.33	
	Median	36.00	
	Variance	128.352	
	Std. Deviation	11.329	
	Minimum	18	
	Maximum	75	
	Range	57	
	Interquartile Range	19	
	Skewness	.281	.115
	Kurtosis	-.636	.230

*Check for Normality
(mean=median, skewness ± 2 , kurtosis ± 2).*

Decide to describe using mean (SD) or median (IQR).

Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
Age (years)	.078	450	.000	.968	450	.000

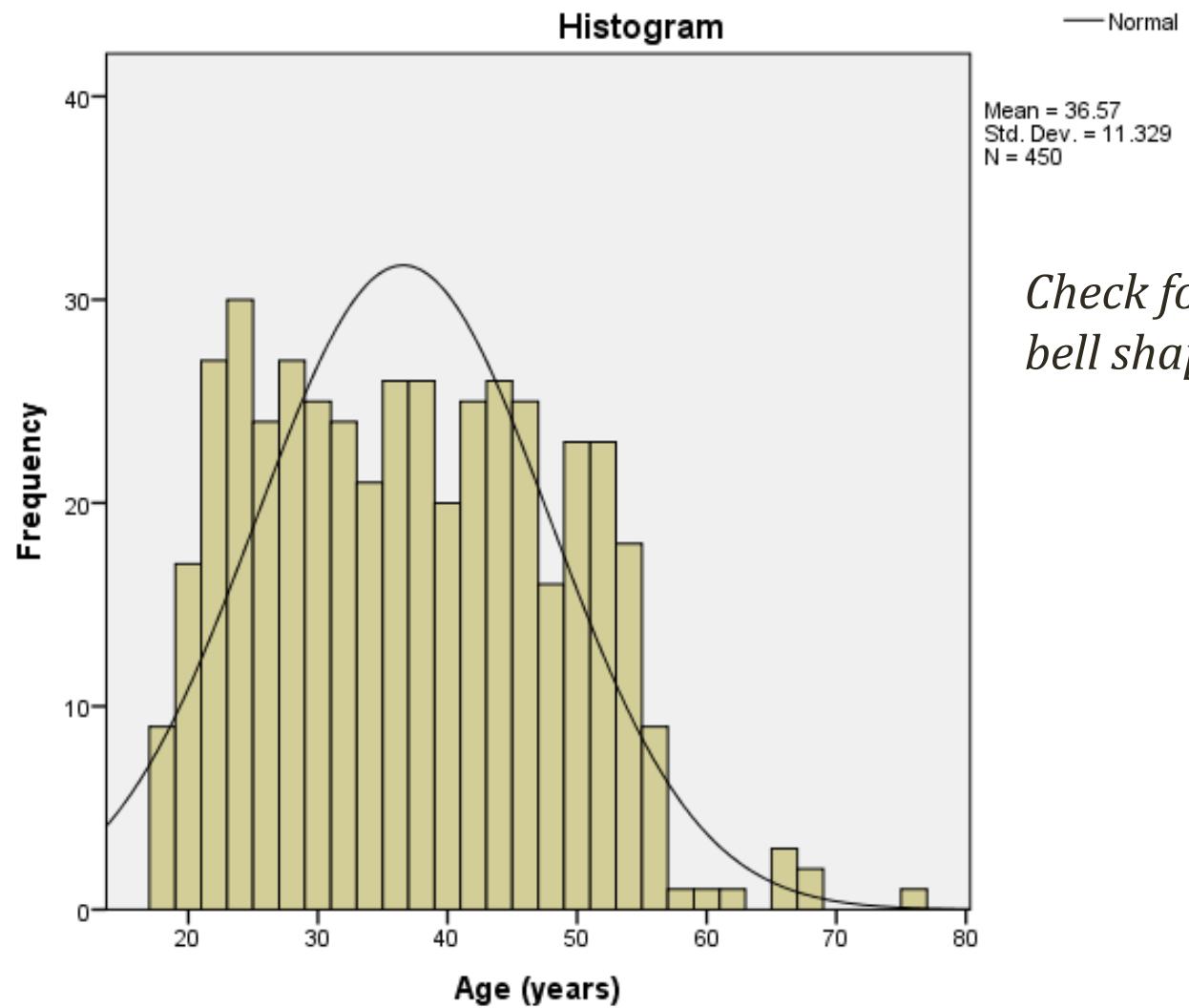
a. Lilliefors Significance Correction

Test Normality. If $P < 0.05$, the distribution is not Normal.*

* Careful when interpreting Normality test

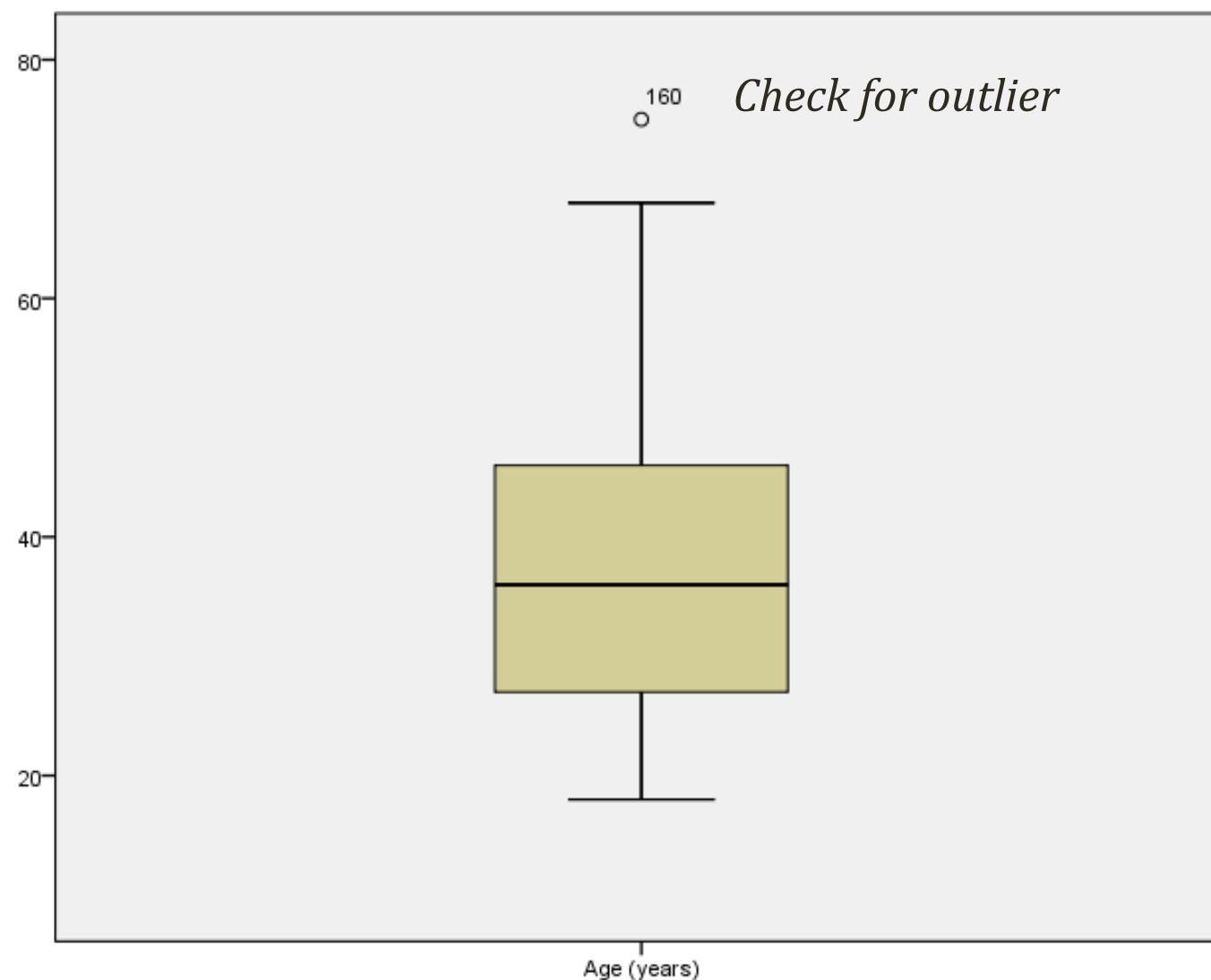


Age (years)



Check for symmetrical bell shape curve

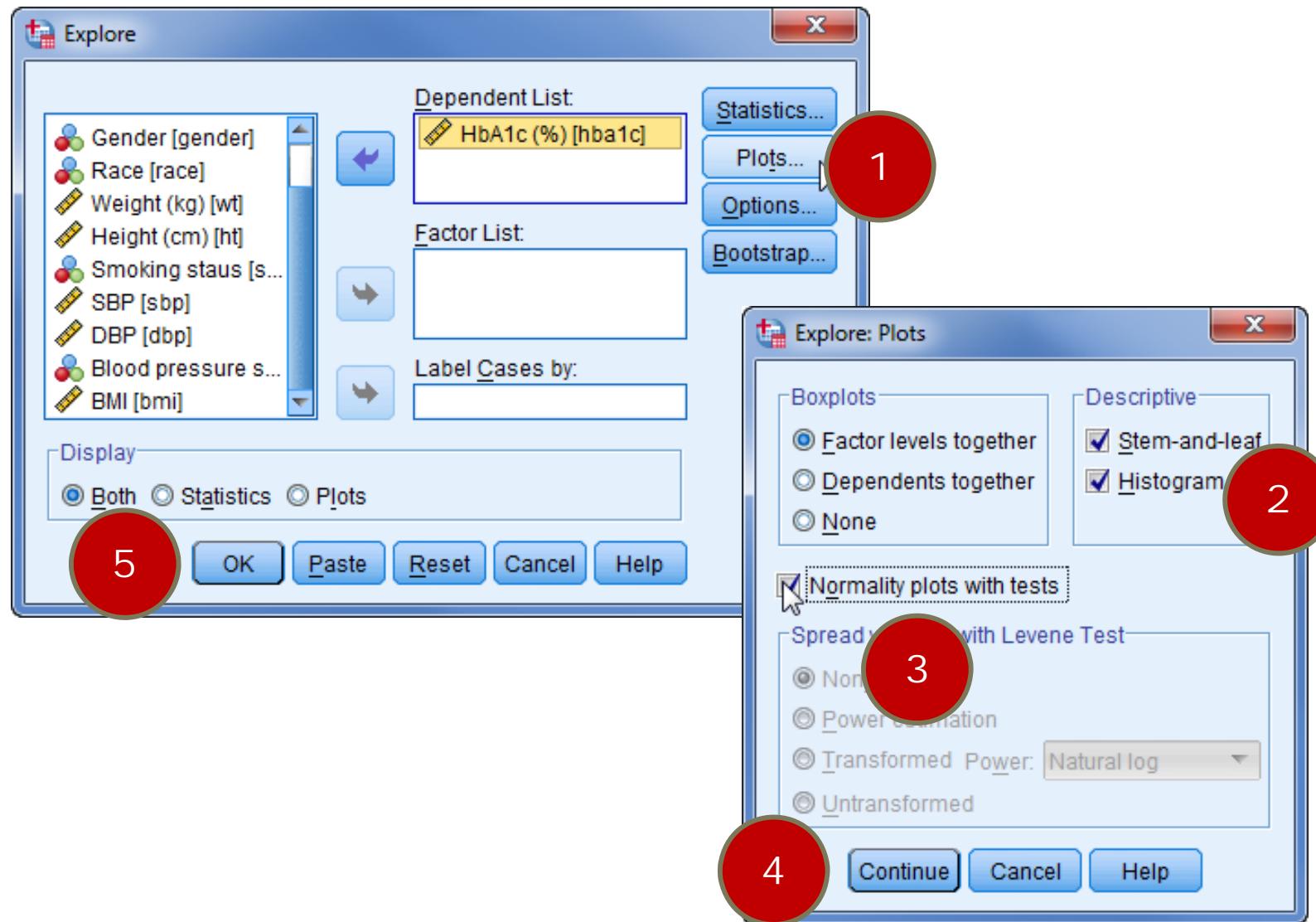




[9]



Describe numerical data (*HbA1c*)



(10)



Case Processing Summary

	Cases					
	Valid		Missing		Total	
	N	Percent	N	Percent	N	Percent
HbA1c (%)	450	100.0%	0	0.0%	450	100.0%

Descriptives

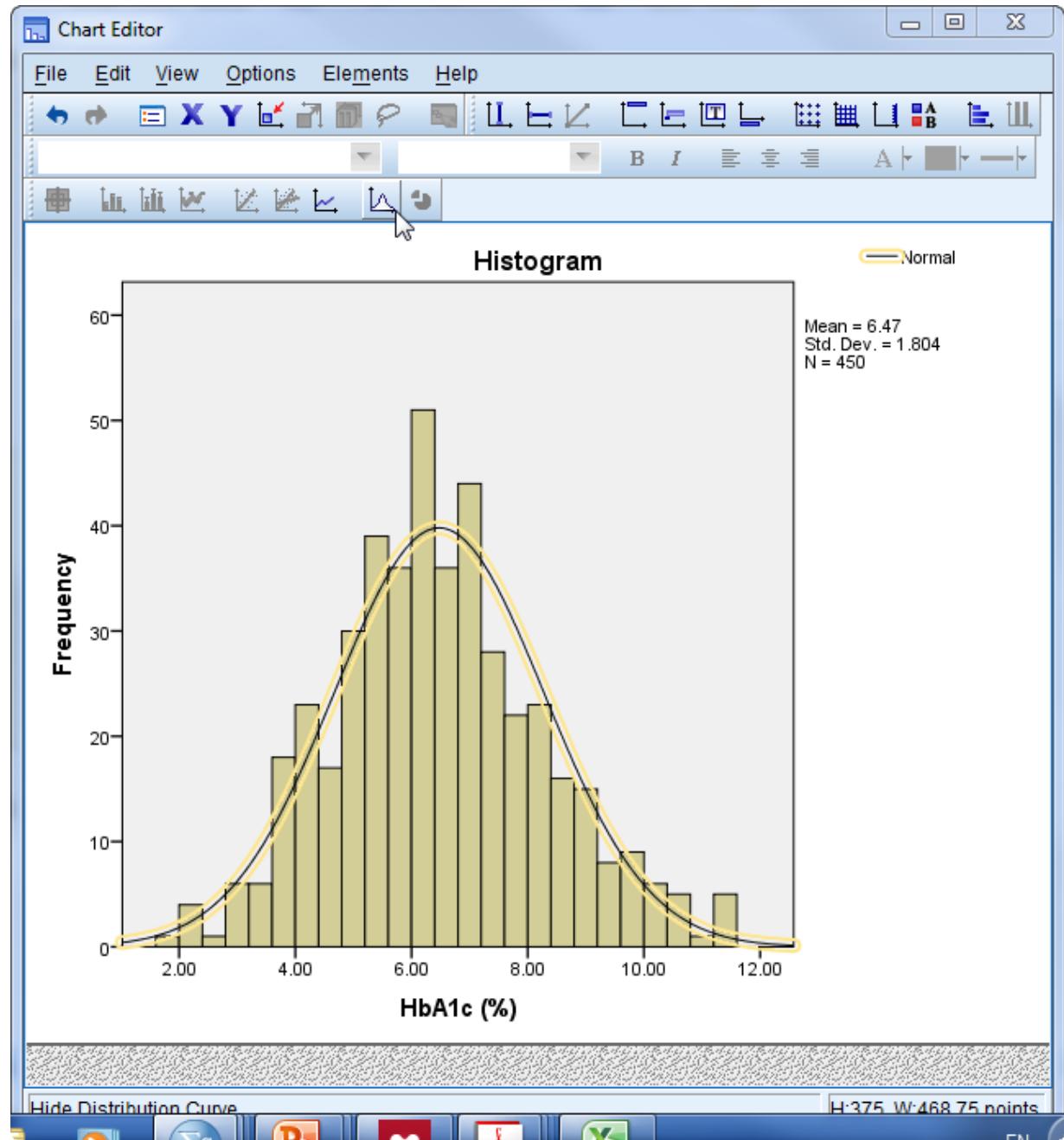
		Statistic	Std. Error
	Mean	6.4745	.08505
HbA1c (%)	95% Confidence Interval for Mean	Lower Bound Upper Bound	6.3073 6.6416
	5% Trimmed Mean	6.4395	
	Median	6.3511	
	Variance	3.255	
	Std. Deviation	1.80425	
	Minimum	1.91	
	Maximum	11.58	
	Range	9.66	
	Interquartile Range	2.26	
	Skewness	.280	.115
	Kurtosis	.009	.230

Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
HbA1c (%)	.049	450	.011	.992	450	.021

a. Lilliefors Significance Correction





[12]



Describe categorical data (*gender*)

The screenshot shows the SPSS interface with several numbered callouts:

- Callout 1: A red circle points to the "Analyze" menu item.
- Callout 2: A red circle points to the "Descriptive Statistics" option under the Analyze menu.
- Callout 3: A red circle points to the "Frequencies..." option in the Descriptive Statistics submenu.
- Callout 4: A red circle points to the "Variable(s)" list in the "Frequencies" dialog box, where "Gender [gender]" is selected.
- Callout 5: A red circle points to the "Display frequency tables" checkbox at the bottom of the dialog box.

The "Frequencies" dialog box also includes buttons for "OK", "Paste", "Reset", "Cancel", and "Help".



Frequency Table

Gender

		Frequency	Percent	Valid Percent	Cumulative Percent
Valid	Female	240	53.3	53.3	53.3
	Male	210	46.7	46.7	100.0
	Total	450	100.0	100.0	



Basic biostatistics with SPSS

COMPARE TWO MEANS

[15]



Age vs. BP

1. Analyze menu open.

2. Compare Means option selected.

3. Independent-Samples T Test... option selected.

4. First Independent-Samples T Test dialog box.

5. Second Independent-Samples T Test dialog box.

6. Define Groups dialog box open.

7. Group 1 set to 1.

8. Continue button clicked.

9. OK button clicked.



Group Statistics

Blood pressure status	N	Mean	Std. Deviation	Std. Error Mean
Age (years) High	179	36.76	11.491	.859
Normal	271	36.45	11.241	.683

Independent Samples Test

	Levene's Test for Equality of Variances		t-test for Equality of Means						
	F	Sig.	t	df	Sig. (2-tailed)	Mean Difference	Std. Error Difference	95% Confidence Interval of the Difference	
								Lower	Upper
Age (years)	Equal variances assumed	.004	.950	.283	448	.777	.310	1.092	-1.837 2.456
				.282	375.276	.778	.310	1.097	-1.848 2.467

Is there age difference between hypertensive and normal subjects?

36.8 (11.5) years vs. 36.4 (11.2) years ($t=0.283$, $df=448$, $P=0.777$).

Choose $P=0.777$ instead of $P=0.778$ if P for Levene's test > 0.05 (meaning equal variances assumed)



Basic biostatistics with SPSS

COMPARE MORE THAN TWO MEANS

[18]



Race vs. BMI

The screenshot illustrates the steps to run a One-Way ANOVA analysis in SPSS:

1. Click on the Analyze menu.
2. Click on Compare Means.
3. Click on One-Way ANOVA... (highlighted with a yellow box).
4. In the One-Way ANOVA dialog box:
 - Dependent List: BMI [bmi] (highlighted with a red circle).
 - Factor: Race [race] (highlighted with a red circle).
 - OK button (highlighted with a red circle).
 - Paste, Reset, Cancel, Help buttons.
5. In the One-Way ANOVA dialog box:
 - Dependent List: BMI [bmi] (highlighted with a red circle).
 - Factor: Race [race] (highlighted with a red circle).
 - OK button (highlighted with a red circle).
 - Paste, Reset, Cancel, Help buttons.
6. Click on Options... (highlighted with a red circle).
7. In the One-Way ANOVA: Options dialog box:
 - Statistics:
 - Descriptive (checked)
 - Fixed and random effects
 - Homogeneity of variance test (checked)
 - Brown-Forsythe
 - Welch
 - Means plot
 - Missing Values:
 - Exclude cases analysis by analysis (radio button selected)
 - Exclude cases listwise

Continue, Cancel, Help buttons.
8. Click on Continue.
9. Click on OK.

Below the dialog boxes, a portion of the SPSS data editor is visible, showing variables Age, Gender, Weight, Height, Smoking status, SBP, DBP, HbA1c, and Blood pressure, along with the Factor variable Race and its levels Malay and Male.



Descriptives

BMI

	N	Mean	Std. Deviation	Std. Error	95% Confidence Interval for Mean		Minimum	Maximum
					Lower Bound	Upper Bound		
Malay	234	26.2103	5.91398	.38661	25.4486	26.9720	18.90	66.10
Chinese	182	27.3701	6.69782	.49648	26.3904	28.3497	18.49	62.53
Indian	32	27.1894	7.19586	1.27206	24.5950	29.7838	19.56	48.90
Total	448	26.7514	6.34901	.29996	26.1619	27.3409	18.49	66.10

Test of Homogeneity of Variances

BMI

Levene Statistic	df1	df2	Sig.
2.333	2	445	.098

Is there differences of mean BMI between Malay, Chinese or Indian? 26.2 vs. 27.3 vs. 27.2?

ANOVA

BMI

	Sum of Squares	df	Mean Square	F	Sig.
Between Groups	144.315	2	72.158	1.796	.167
Within Groups	17874.221	445	40.167		
Total	18018.536	447			

$$F(2; 445) = 1.796, P = 0.167 \text{ i.e. } P > 0.05.$$



Basic biostatistics with SPSS

COMPARE PROPORTIONS

[21]



Gender vs. BP

The image shows the SPSS software interface with several windows open. The main window is titled "Crosstabs: Cell Display". The "Row(s)" field contains "Gender [gender]" and the "Column(s)" field contains "Blood pressure status [...]".

The "Cells..." button in the main dialog is highlighted with a red circle labeled 6. In the "Percentages" section of the sub-dialog, the "Row" option is checked, highlighted with a red circle labeled 7. The "Statistics..." button in the main dialog is also highlighted with a red circle labeled 8.

Red numbered circles 1 through 5 are placed on the menu path: Analyze > Descriptive Statistics > Crosstabs... .



Crosstabs

Row(s): Gender [gender]

Column(s): Blood pressure status [...]

Layer 1 of 1

Previous Next

Display layer variables in table layers

Display clustered bar charts
 Suppress tables

OK Paste Reset Cancel Help

Exact...
Statistics...
Cells...
Format...
Bootstrap...

Crosstabs: Statistics

Chi-square Correlations

Nominal

Contingency coefficient
 Phi and Cramer's V
 Lambda
 Uncertainty coefficient

Ordinal

Gamma
 Somers' d
 Kendall's tau-b
 Kendall's tau-c

Nominal by Interval

Eta

Kappa
 Risk
 McNemar

Cochran's and Mantel-Haenszel statistics

Test common odds ratio equals:

Continue Cancel Help

9

10

11

12

(23)



Gender * Blood pressure status

			Crosstab		
			Blood pressure status		Total
Gender	Female	Count	Normal	High	
		% within Gender	145 60.4%	95 39.6%	240 100.0%
Gender	Male	Count	126 60.0%	84 40.0%	210 100.0%
		% within Gender			
Total		Count	271	179	450
		% within Gender	60.2%	39.8%	100.0%

Who has higher prevalence of high blood pressure? Male or female? 40% vs. 39.6%.

Chi-Square Tests

	Value	df	Asymp. Sig. (2-sided)	Exact Sig. (2-sided)	Exact Sig. (1-sided)
Pearson Chi-Square	.008 ^a	1	.928		
Continuity Correction ^b	.000	1	1.000		
Likelihood Ratio	.008	1	.928		
Fisher's Exact Test				1.000	.502
Linear-by-Linear Association	.008	1	.928		
N of Valid Cases	450				

a. 0 cells (0.0%) have expected count less than 5. The minimum expected count is 83.53.

b. Computed only for a 2x2 table

$$\chi^2 = 0.928, df = 1, P = 0.928 \text{ i.e. } P > 0.05.$$



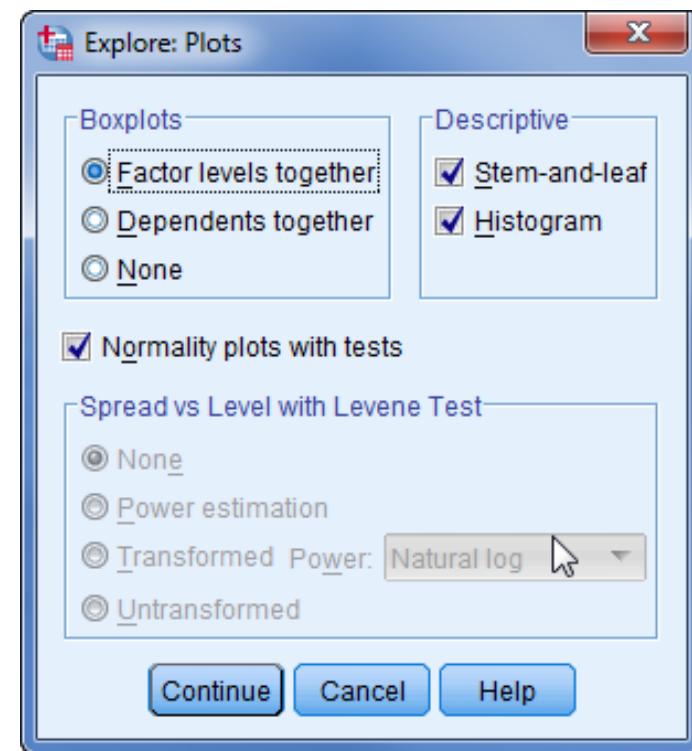
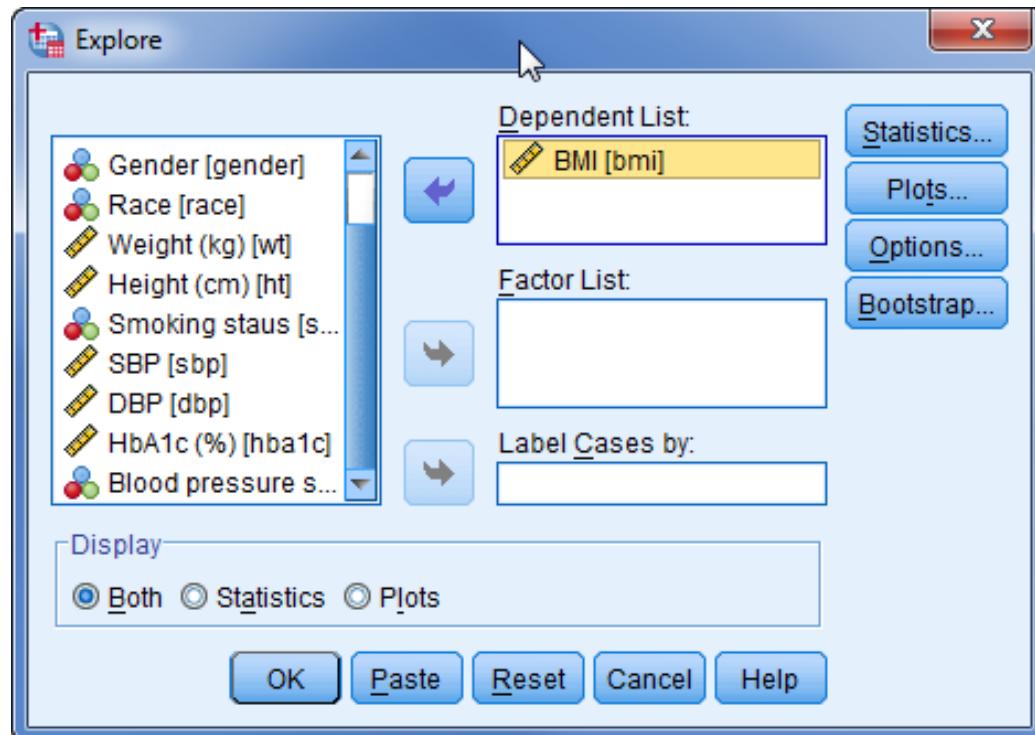
Basic biostatistics with SPSS

NON-PARAMETRICS

[25]



Describe BMI



Case Processing Summary

	Cases					
	Valid		Missing		Total	
	N	Percent	N	Percent	N	Percent
BMI	448	99.6%	2	0.4%	450	100.0%

Descriptives

		Statistic	Std. Error
BMI	Mean	26.7514	.29996
	95% Confidence Interval		
	for Mean	Lower Bound	26.1619
		Upper Bound	27.3409
	5% Trimmed Mean	26.1638	
	Median	25.0194	
	Variance	40.310	
	Std. Deviation	6.34901	
	Minimum	18.49	
	Maximum	66.10	
	Range	47.60	
	Interquartile Range	6.20	
	Skewness	1.896	.115
	Kurtosis	5.746	.230

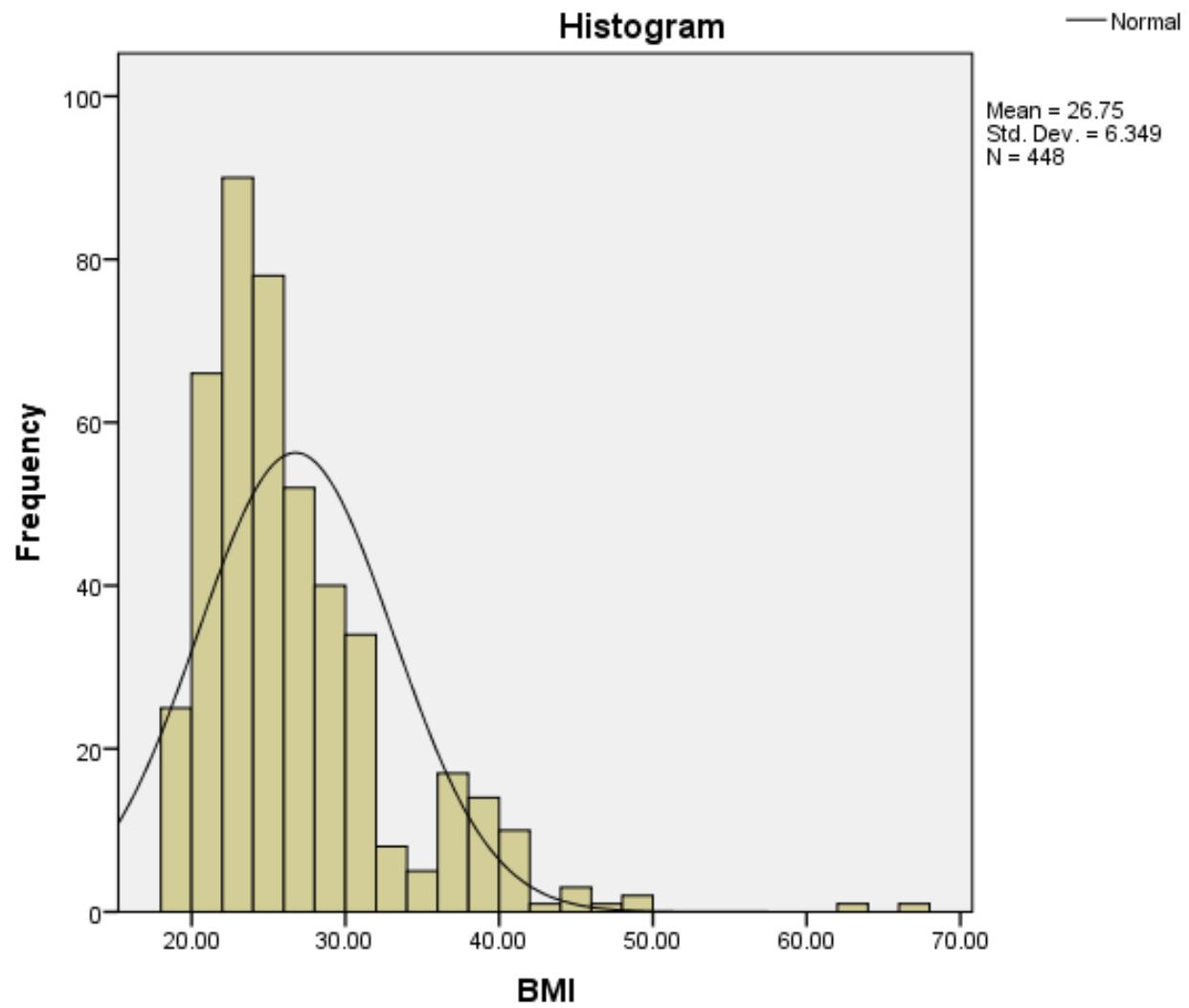
BMI is not Normally distributed.

Tests of Normality

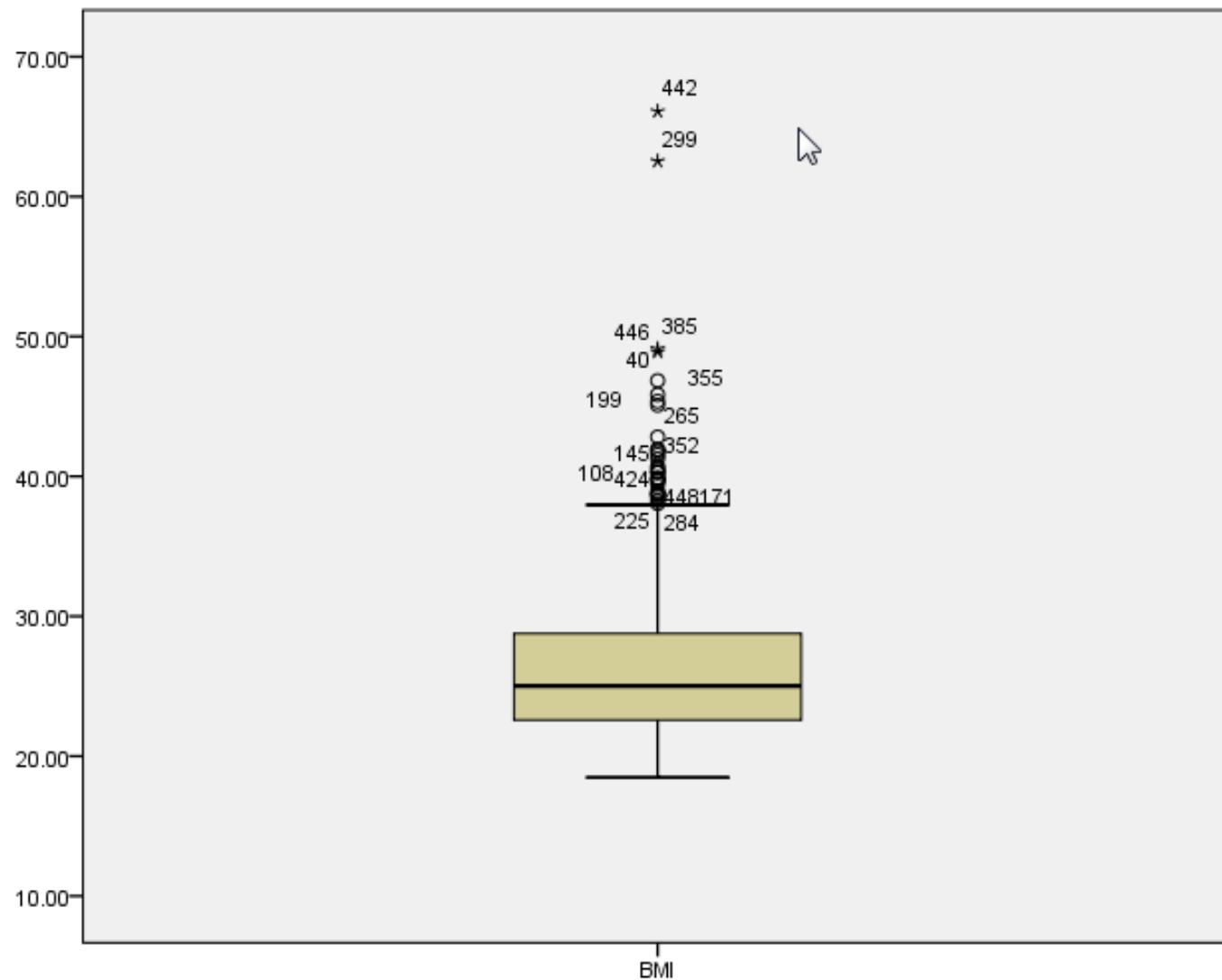
	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
BMI	.137	448	.000	.844	448	.000

a. Lilliefors Significance Correction





Observed Value



[29]



Gender vs. BMI

1. Analyze menu open.

2. Nonparametric Tests selected.

3. Legacy Dialogs selected.

4. 2 Independent Samples... selected.

5. Test Variable List: BMI [bmi] selected.

6. Grouping Variable: gender(??) selected.

7. Define Groups... button.

8. Group 1: 1 selected.

9. Group 2: 0 selected.

10. OK button.

	race	wt
140	Malay	56
146	Malay	55
147	Malay	56
149	Malay	50
150	Malay	86
151	Malay	65
153	Malay	55
157	Malay	69
162	Malay	147
163	Malay	51
164	Malay	149
165	Malay	52
166	Malay	147
167	Malay	150
171	Malay	149
172	Malay	168
173	Malay	159
174	Malay	61
178	Malay	73
180	Malay	70
181	Malay	70
184	Malay	50
186	Malay	50
187	Female	Malay
188	Female	Malay
189	Female	Malay
191	Female	Malay



Mann-Whitney Test

Ranks

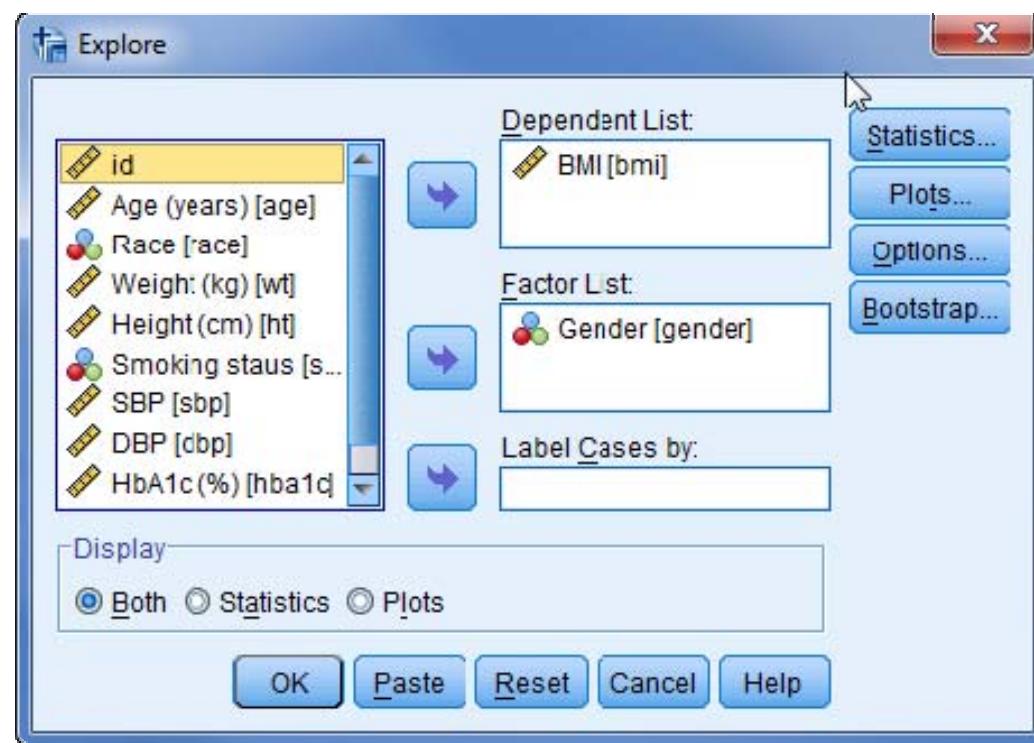
Gender	N	Mean Rank	Sum of Ranks
BMI Female	238	168.59	40124.00
Male	210	287.87	60452.00
Total	448		

Test Statistics^a

	BMI
Mann-Whitney U	11683.000
Wilcoxon W	40124.000
Z	-9.731
Asymp. Sig. (2-tailed)	.000

a. Grouping Variable: Gender

*Mann-Whitney tests between 'rank'.
Don't have to cite the mean 'rank'. They
have no meaning in the actual context.
Instead describe their medians.*





$23.3 (IQR=4.8)$ vs. $27.5 (IQR=11.0) \text{ kg/m}^2$, Z (Mann-Whitney) = 9.731 , $P < 0.001$

Gender		Statistic	Std. Error
Female	Mean	24.0304	.22031
	95% Confidence Interval for Mean	Lower Bound Upper Bound	23.5963 24.4644
	5% Trimmed Mean	23.8684	
	Median	23.3391	
	Variance	11.552	
	Std. Deviation	3.39877	
	Minimum	18.49	
	Maximum	32.82	
	Range	14.33	
	Interquartile Range	4.75	
	Skewness	.684	.158
	Kurtosis	-.314	.314
	Mean	29.8352	.51236
Male	95% Confidence Interval for Mean	Lower Bound Upper Bound	28.8252 30.8453
	5% Trimmed Mean	29.2533	
	Median	27.5397	
	Variance	55.127	
	Std. Deviation	7.42476	
	Minimum	19.92	
	Maximum	66.10	
	Range	46.18	
	Interquartile Range	11.02	
	Skewness	1.463	.168
	Kurtosis	3.319	.334

Race vs. BMI

1 Analyze Direct Marketing Graphs Utilities Add-ons Window Help

- Reports
- Descriptive Statistics
- Tables
- Compare Means
- General Linear Model
- Generalized Linear Models
- Mixed Models
- Correlate
- Regression
- Loglinear
- Neural Networks
- Classify
- Dimension Reduction
- Scale
- Nonparametric Tests
- Forecasting
- Survival
- Multiple Response
- Missing Value Analysis
- Multiple Imputation
- Complex Samples
- Quality Control
- ROC Curve...

	race	wt
22	Female	Malay
48	Female	Malay
47	Female	Malay
21	Female	Malay
		50
		50

2 One Sample...
3 Independent Samples...
4 Related Samples...
5 Legacy Dialogs

6 Test Variable List: BMI [bmi]
7 Grouping Variable: race(??)
8 Define Range
9 Continue
10 OK

Tests for Several Independent Samples

Kruskal-Wallis H Median
Jonckheere-Terpstra

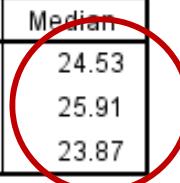
Several Independent Samp...

Range for Grouping Variable

Minimum: 0 Maximum: 2



		BMI	
		Count	Median
Race			
	Malay	234	24.53
	Chinese	184	25.91
	Indian	32	23.87



Kruskal-Wallis Test

Ranks

	Race	N	Mean Rank
BMI	Malay	234	214.29
	Chinese	182	238.32
	Indian	32	220.58
	Total	448	



Test Statistics^{a,b}

	BMI
Chi-Square	3.557
df	2
Asymp. Sig.	.169

a. Kruskal Wallis
Test

b. Grouping
Variable: Race

(34)



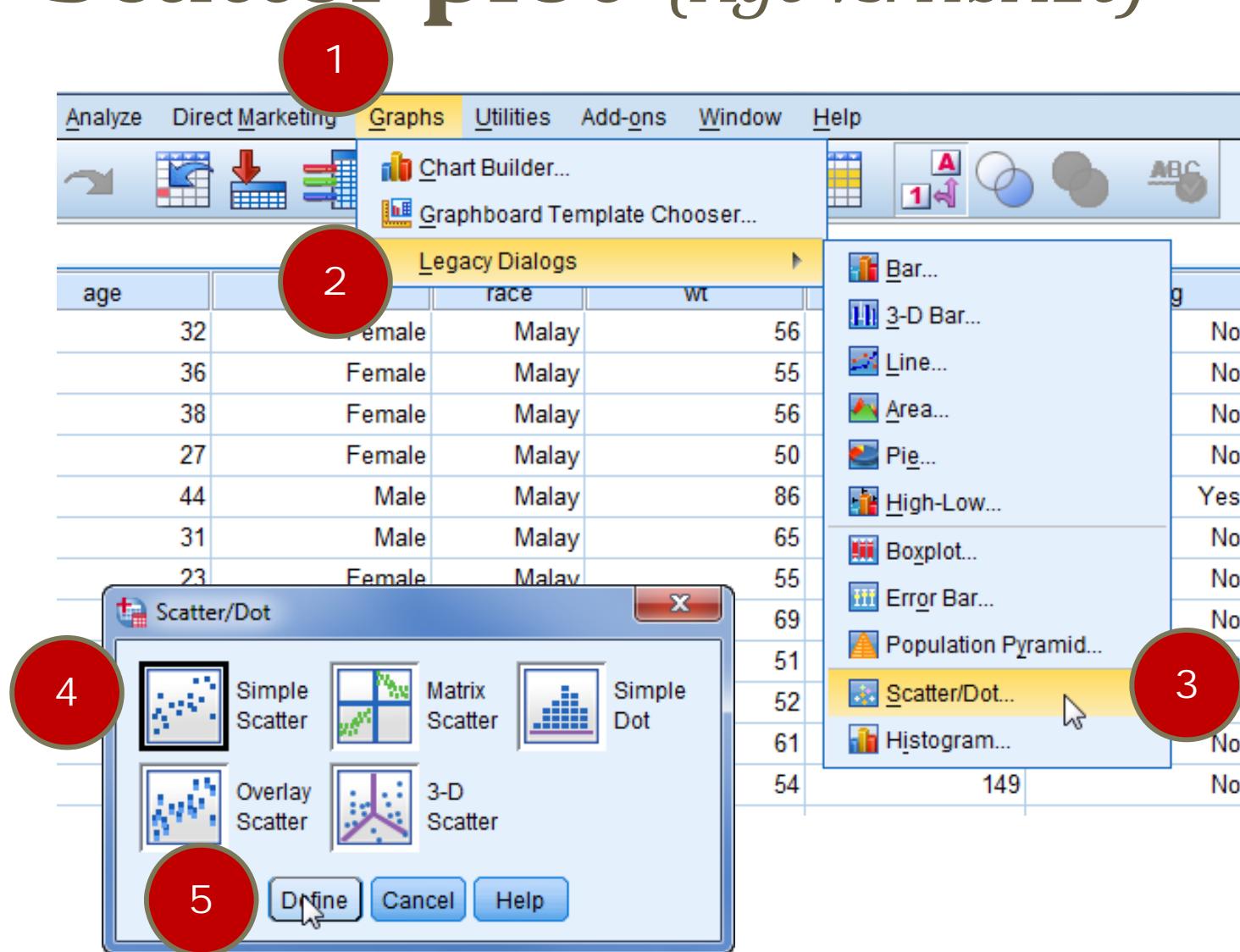
Basic biostatistics with SPSS

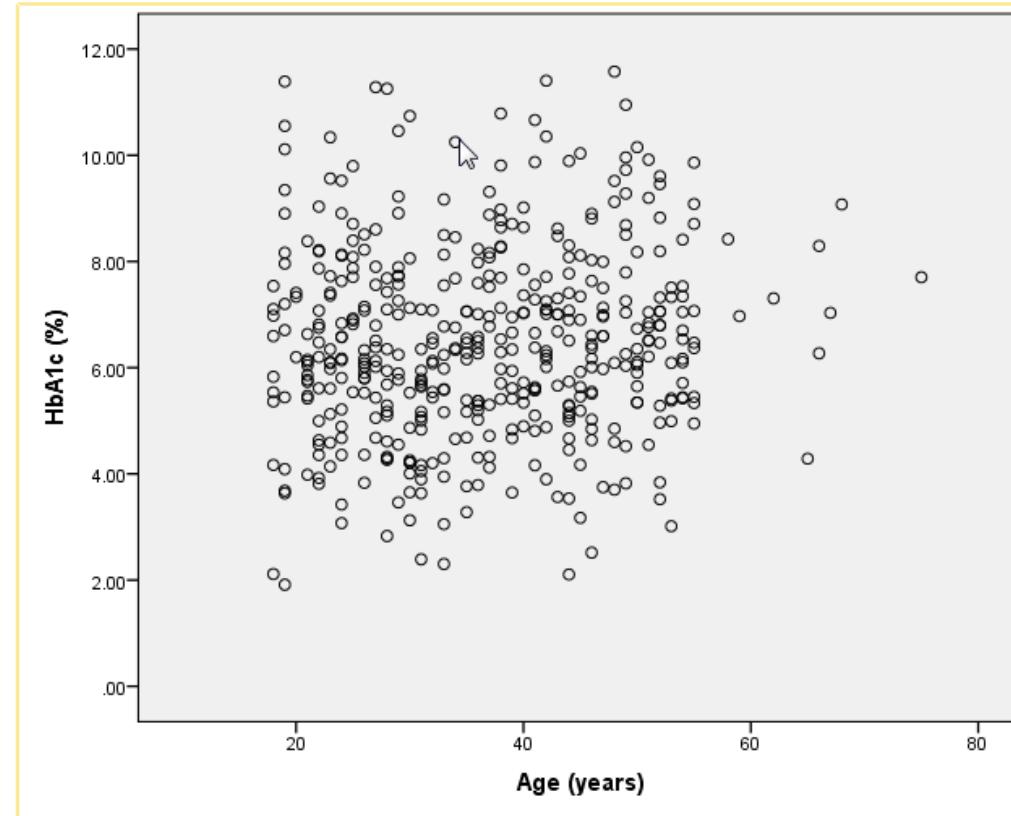
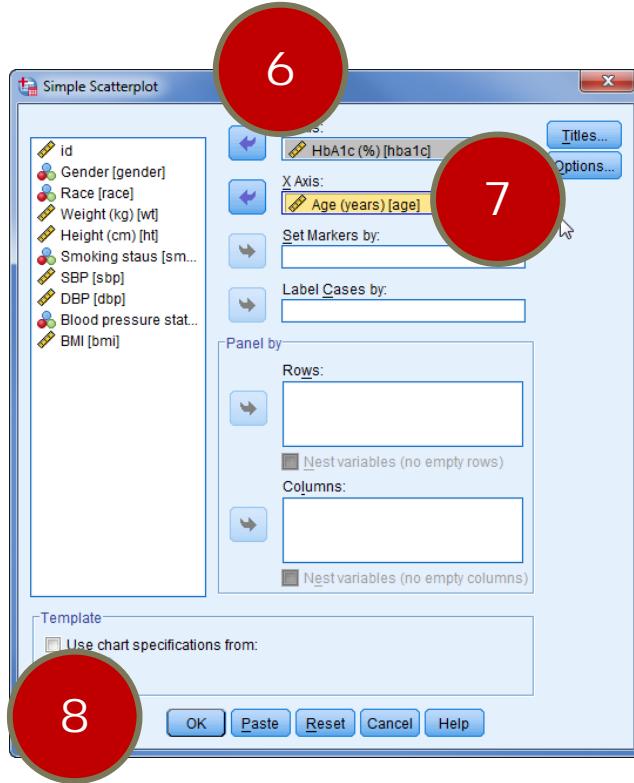
CORRELATION

[35]



Scatter plot (Age vs. HbA1c)





Describe the relationship of two numerical variables.

Dependent variable in Y-axis, explanatory (exposure) in x-axis.



Age vs. HbA1c

The screenshot shows the SPSS interface with several numbered callouts:

- 1: The Analyze menu is open, highlighting the "Correlate" option.
- 2: The "Bivariate..." option under the "Correlate" menu is selected.
- 3: The "Bivariate Correlations" dialog box is open, showing the "Variables" list.
- 4: In the "Variables" list, "Age (years) [age]" and "HbA1c (% [hba1c])" are selected.
- 5: The "OK" button at the bottom of the dialog box is highlighted.

The background of the SPSS interface shows a data view with columns "race" and "wt". The "race" column contains values like "Malay" and "Malay", while the "wt" column contains numerical values such as 56, 55, 56, 50, and 86.



Correlations				
	Age (years)	HbA1c (%)	BMI	
Age (years)	Pearson Correlation	1	.096*	-.016
	Sig. (2-tailed)		.043	.734
	N	450	450	448
HbA1c (%)	Pearson Correlation	.096*	1	.314**
	Sig. (2-tailed)	.043		.000
	N	450	450	448
BMI	Pearson Correlation	-.016	.314**	1
	Sig. (2-tailed)	.734	.000	
	N	448	448	448

*. Correlation is significant at the 0.05 level (2-tailed).
 **. Correlation is significant at the 0.01 level (2-tailed).

There is a weak correlation between age & HbA1c ($r=0.96, P=0.43$).

$< 0.3 = \text{Weak}$
 $0.3-0.7 = \text{Moderate}$
 $> 0.7 = \text{Strong}$

