

# การวิเคราะห์องค์ประกอบ

## (Factor analysis)

Suwanna Sayruamyat

Email: [suwanna.s@ku.th](mailto:suwanna.s@ku.th)

Facebook: Suwanna Sayruamyat

Page: [\*\*EatEcon\*\*](#),

Website: [www.eatecon.com](http://www.eatecon.com)

# เป้าหมายของการวิเคราะห์องค์ประกอบ



Data  
summarisations

Data reduction



Variable  
selection

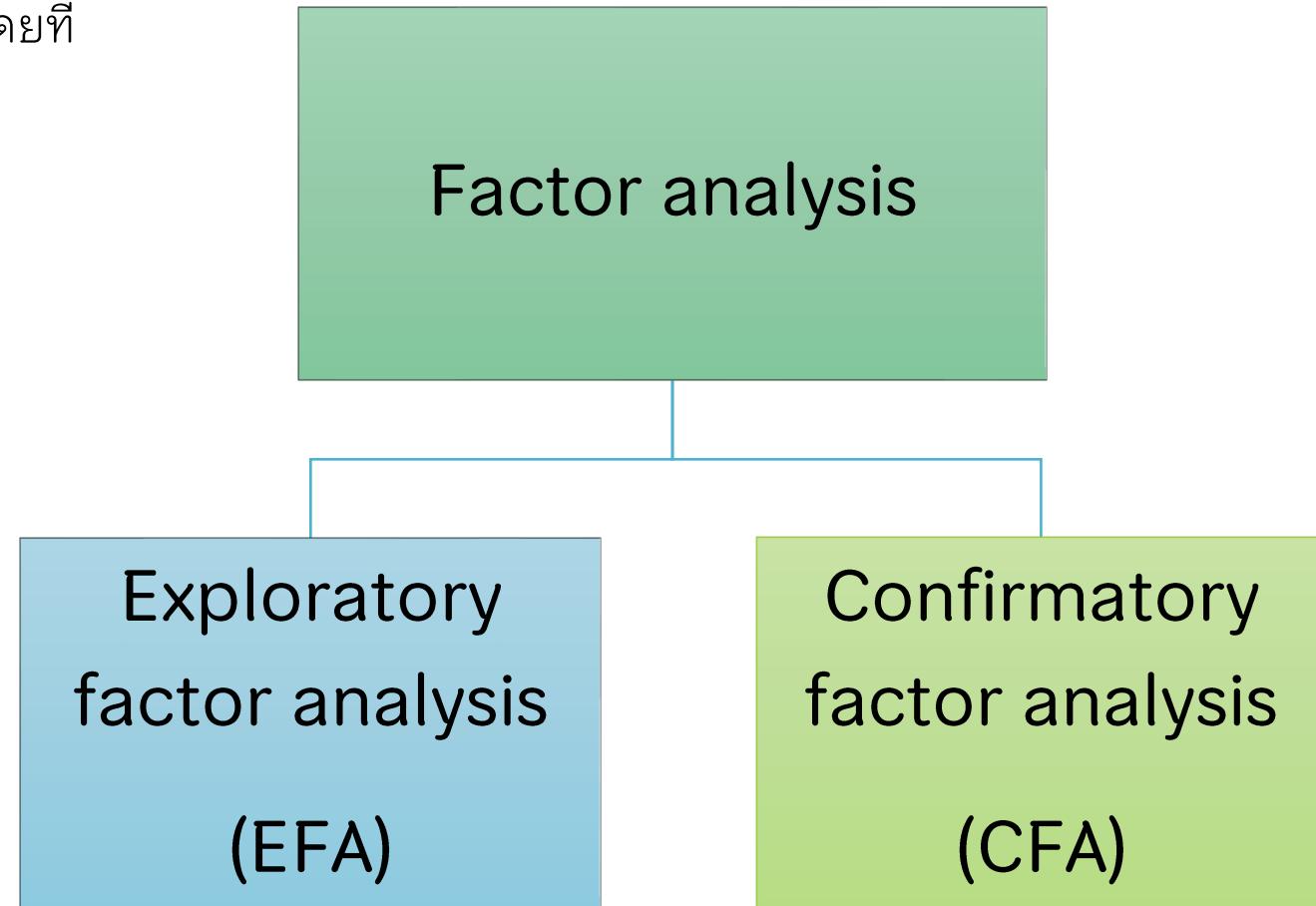
1. การวิเคราะห์องค์ประกอบ (Factor analysis) เป็นเทคนิคในการค้นหาตัวประกอบต่าง ๆ (factors) จากชุดตัวแปรที่มีองค์ประกอบร่วมกัน (สัมพันธ์กัน)
  - เป้าหมายหลักคือจดประสงค์หลักเพื่อกำหนดโครงสร้างพื้นฐานของตัวแปรในการวิเคราะห์
2. ตัวแปรไม่ได้ถูกระบุประเภทว่าเป็นตัวแปรตามหรือตัวแปรอิสระ แต่จะตรวจสอบชุดความสัมพันธ์ที่พึงพาซึ่งกันและกันทั้งหมดระหว่างตัวแปร เพื่อกำหนดชุดของมิติร่วมที่เรียกว่า ปัจจัย (FACTORS)
3. การวิเคราะห์องค์ประกอบได้รับการออกแบบมาเพื่อแสดงคุณลักษณะที่หลากหลายในจำนวนมิติใหม่ที่มีจำนวนน้อยลงโดยมีการสูญเสียข้อมูลน้อยที่สุด

# ประเภทของการวิเคราะห์องค์ประกอบ

EFA: ต้องการจัดกลุ่มตัวแปร โดยที่ยังไม่รู้มาก่อนว่าตัวแปรใดอยู่ภายใต้องค์ประกอบใด

- สำรวจเพื่อให้รู้ว่าตัวแปรใดอยู่ด้วยกันบ้าง

- Summarising data** by grouping correlated variables.
- Investigating sets of measured variables** related to theoretical constructs.
- Preliminary exploration of data (**Data-driven**)



CFA: ต้องการตรวจสอบตัวแปรที่อยู่ในแต่ละกลุ่มว่ามีน้ำหนักเพียงพอที่จะอยู่ในกลุ่มนั้น จริงหรือไม่

- มีกรอบแนวคิดอยู่แล้ว ต้องการยืนยันว่าตัวแปรเหล่านั้นอยู่ในกลุ่มจริง
- ใช้ CFA เพื่อทดสอบว่าโมเดลที่กำหนดไว้เหมาะสมสมกับข้อมูลที่มีมากน้อยเพียงใด
- Testing generalisation of factor structure to new data.
- Making use of only the **measurement model** component of the general SEM.
- It should be based on theory and/or the results of EFA and other psychometric tests.**
- Test of theory against data (**Theory-driven**)

# การวิเคราะห์องค์ประกอบ

- Also called “unrestricted” factor analysis.
- ค้นหาค่าความสัมพันธ์ของปัจจัย (factor loadings) ที่สร้างความสัมพันธ์ระหว่างตัวแปรที่สังเกตได้ให้ดีที่สุด
- จำนวนปัจจัย ( $n$  factors) = จำนวนตัวแปรที่สังเกตได้ ( $n$  of observed variables)
- ตัวแปรทั้งหมดมีความสัมพันธ์กับทุกปัจจัย (factors).
- เก็บรักษาจำนวนปัจจัย  $< n$  ที่ "อธิบาย" ปริมาณความแปรปรวนที่สังเกตได้ในระดับที่น่าพอใจ
- "ความหมาย" ของปัจจัยถูกกำหนดโดยรูปแบบของค่าความสัมพันธ์ของปัจจัย (pattern of loadings).
- No unique solution where  $>1$  factor, rotation used to clarify what each factor measures.

## สมมติฐานการวิเคราะห์

1. ความสัมพันธ์ระหว่าง Factor กับ variable เป็นเชิงเส้นตรง (linear relationship)
2. ข้อมูลที่ใช้ควรเป็น interval scale
3. Factor และ error เป็นอิสระต่อกัน
4. จำนวนข้อมูลที่นำมาวิเคราะห์ต้องมากกว่าจำนวนตัวแปร
5. Multicollinearity in the data is desirable because the aim is to identify interrelated set of variables.
6. Data is not an identity matrix. (ตอบเหมือนๆ กัน)

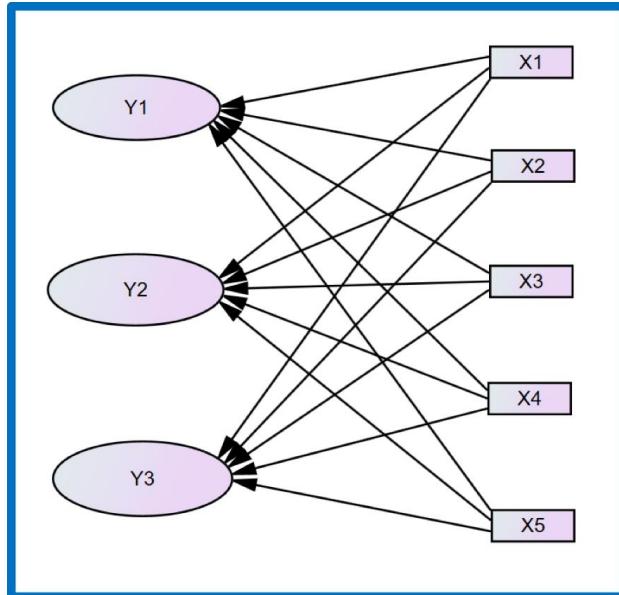
# ขนาดตัวอย่างที่เหมาะสมในการวิเคราะห์

1. อัตราส่วนบันทึกของตัวแปรสำหรับ EFA คือ 1:5 (1 ตัวแปร ต่อจำนวนตัวอย่าง 5 observations)
  - Ex. 20 ตัวแปร ความมีตัวอย่าง  $\geq 100$  observations
2. Ideal condition ratio is 1:20.
3. จำนวนตัวอย่างต้องมากกว่าจำนวนตัวแปรที่วิเคราะห์
4. ขนาดตัวอย่างไม่ควรน้อยกว่า 50 ตัวอย่าง

Sample size	Sufficient factor loading
50	0.75
60	0.70
70	0.65
85	0.60
100	0.55
120	0.50
150	0.45
200	0.40
250	0.35
350	0.30

Source: Hair et al. (2014)

## Principle component analysis

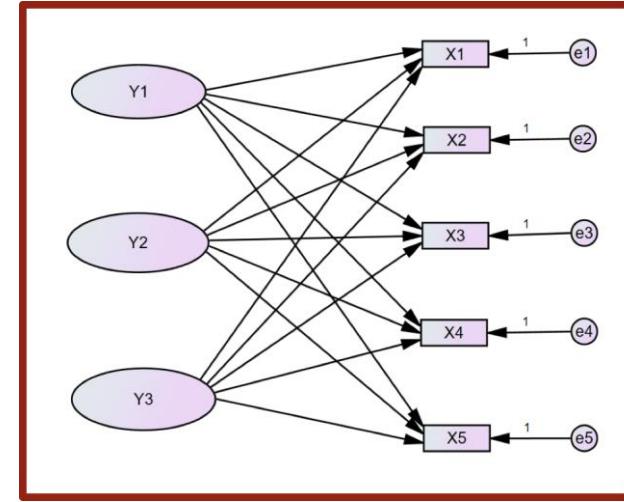


$$Y_i = \sum_{j=1}^p a_{ij} x_j \text{ for } i = 1, 2, \dots, p$$

### Formative

- Direction of causality is from measure to construct
- No reason to expect the measures are correlated
- Indicators are not interchangeable

## Factor analysis



$$X_j = \sum_{i=1}^p b_{ji} Y_i \text{ for } j = 1, 2, \dots, p$$

$$X_j = \sum_{i=1}^p \lambda_{ji} F_i + \lambda F_{j,spec} + e_j \text{ for } j = 1, 2, \dots, p$$

### Reflective

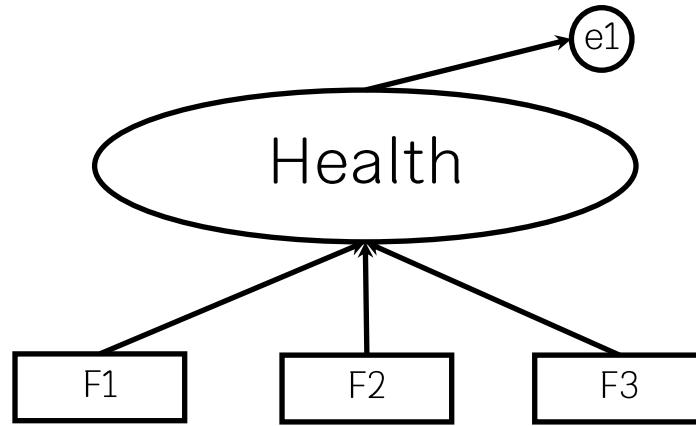
- Direction of causality is from construct to measure
- Measures expected to be correlated
- Indicators are interchangeable

- If you want **summarising** a number of correlating variables in a few new variable with smallest possible loss of information, the **component analysis** is the answer.
- If you want **explaining the correlations** in a data set in form of factors, the **factor analysis** is the answers.
- However, component analysis is less complicated and usually give the same results as exploratory factor analysis. Thus, **most component analysis and EFA both go under the name of factor analysis (Blunch, 2013)**.

# Formative vs. Reflective

## Formative

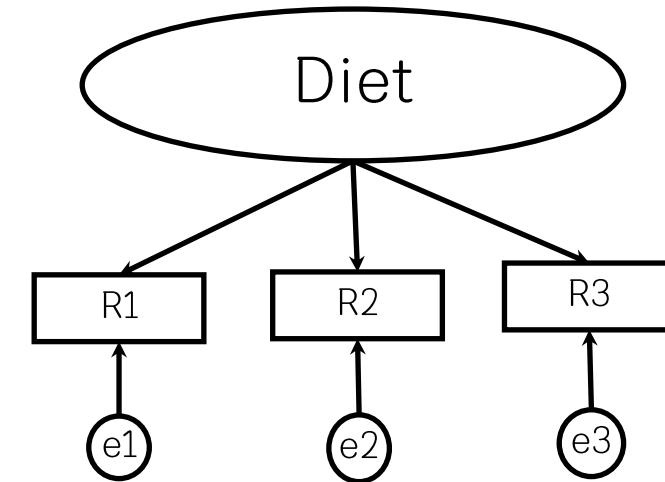
- F1. I have a balanced diet.
- F2. I exercise regularly.
- F3. I get sufficient sleep each night.



Principle component analysis

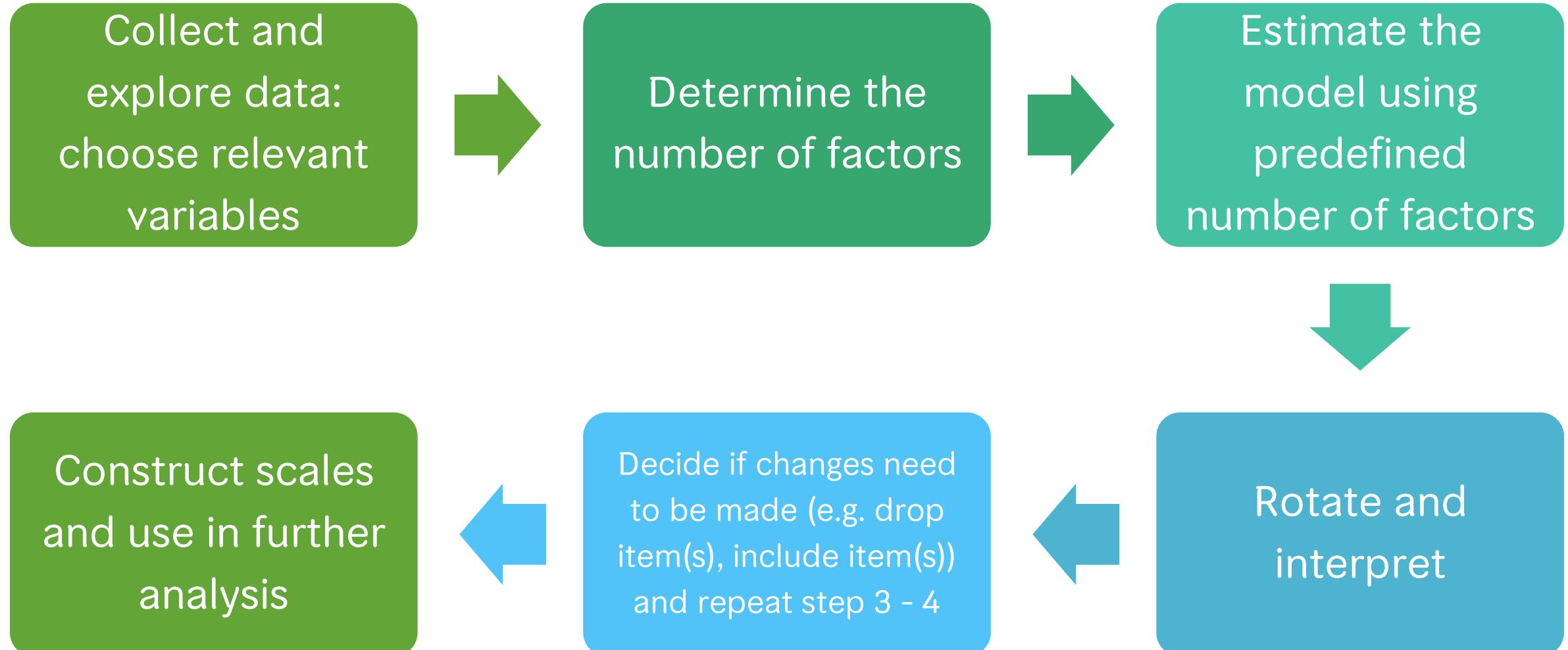
## Reflective

- R1. I eat healthy food.
- R2. I do not eat much junk food.
- R3. I have a balanced diet.



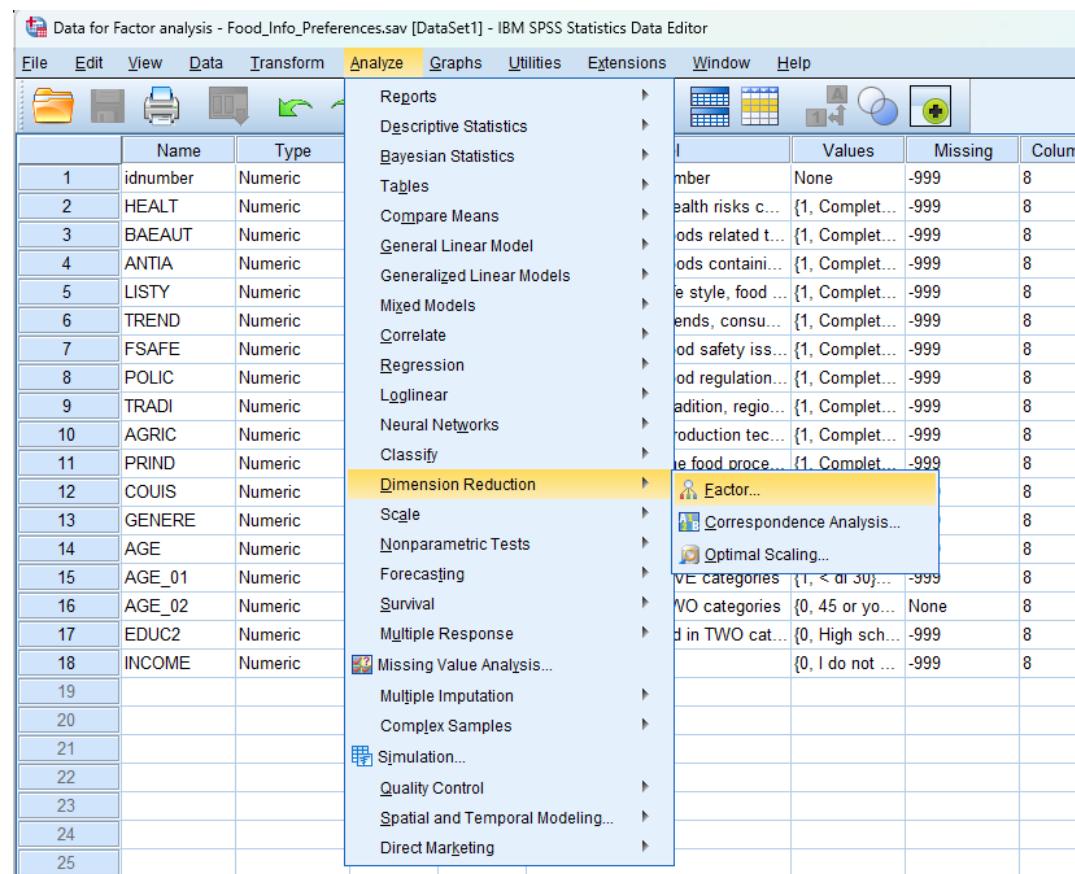
Factor analysis

# Steps in EFA

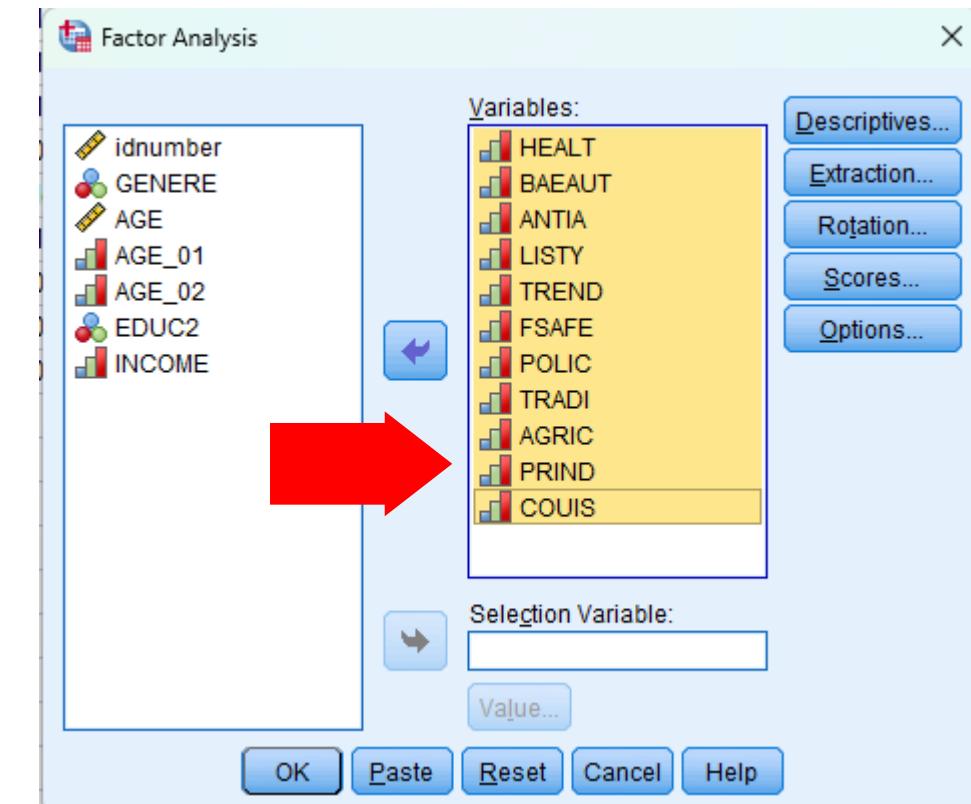


# 1. Download the file Factor analysis – Data for factor analysis

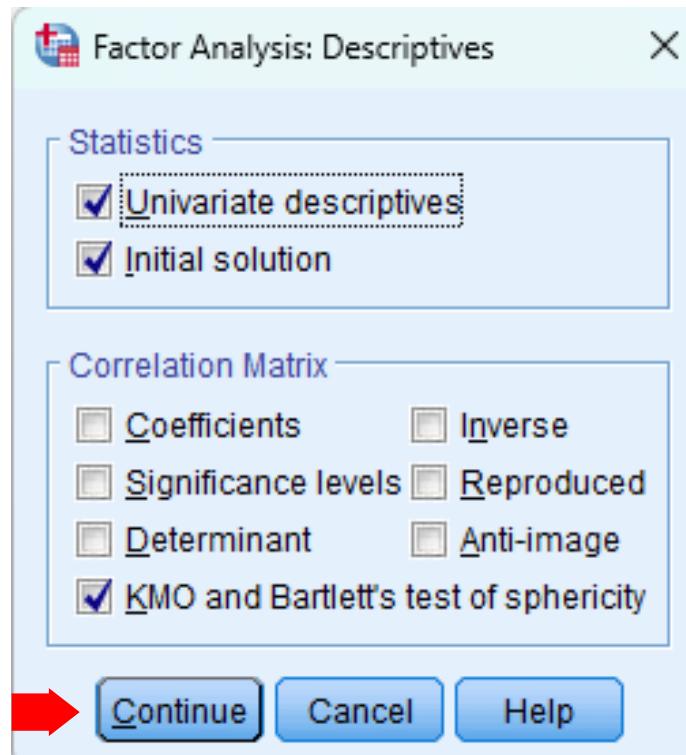
## 2. click the Analyse > Dimension Reduction > Factor



เลือกตัวแปรที่ต้องการใส่ในช่อง variables



# Factor analysis: Descriptives



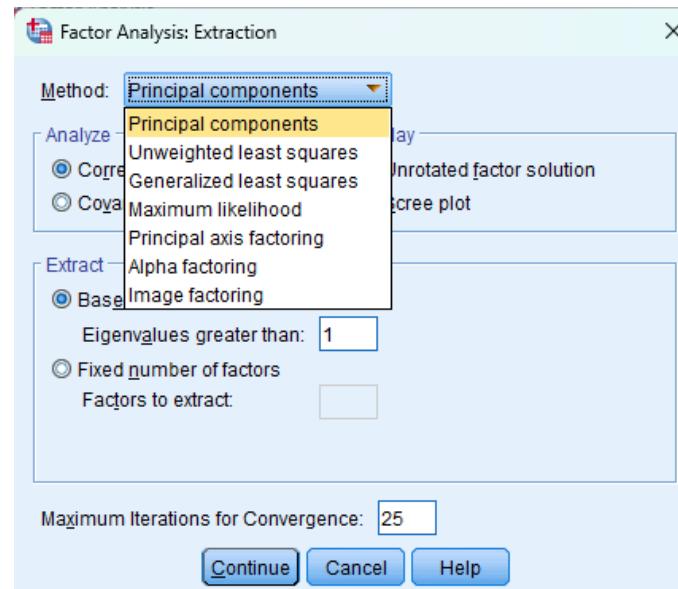
## Statistics

- Univariate descriptive แสดงจำนวนข้อมูล, ค่าเฉลี่ย และค่าเบี่ยงเบนมาตรฐานของตัวแปรแต่ละตัว
- Initial solution แสดงค่า initial communalities, eigenvalue และ percentage of variance explained

## Correlation Matrix

- Coefficients แสดงค่าเมตริกซ์สัมประสิทธิ์สหสัมพันธ์ของตัวแปรทุกคู่
- Significance levels แสดงค่า one-tailed significance level ของการทดสอบค่าสัมประสิทธิ์สหสัมพันธ์ของตัวแปรแต่ละคู่
- Determinant แสดงค่า determinant ของเมตริกซ์สัมประสิทธิ์สหสัมพันธ์
- KMO and Bartlett's test of sphericity แสดงค่า KMO และ Bartlett's test
  - KMO (Kaiser-Meyer-Olkin) เป็นค่าที่ใช้วัดความเหมาะสมของข้อมูลตัวอย่างที่จะนำมาวิเคราะห์โดยเทคนิค Factor Analysis

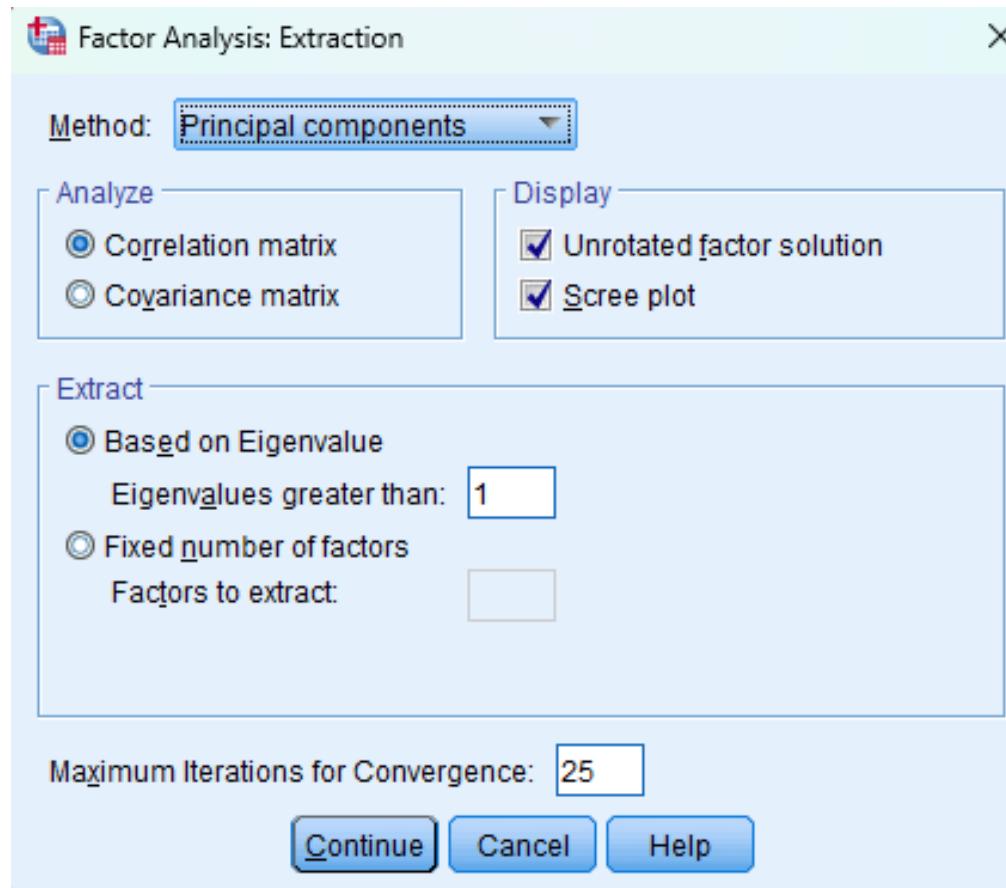
# Factor analysis: Extraction



Method เลือกเทคนิคการวิเคราะห์ปัจจัย แบ่งออกเป็น 2 วิธีหลัก คือ

1. **Principal Component Analysis (PCA)** เป็นวิธีแบบ formative ความนิยมมากที่สุด
2. **Common Factor Analysis (CFA)** เป็นเทคนิคที่มีวัตถุประสงค์เหมือนเทคนิค PCA คือ จะสร้าง Factor เพื่อลดจำนวนตัวแปร แต่หลักเกณฑ์ของ CFA จะพยายามทำให้ค่า แปรปรวนเฉพาะส่วนของ common factor มากที่สุด โดยไม่พิจารณาถึงค่า Unique Factor เทคนิค CFA มีเทคนิคดังนี้
  - 1) Unweighted Least Square เป็นเทคนิคที่ต้องกำหนดจำนวน factor ไว้ແน่อนก่อน แล้วหา Factor pattern matrix ที่ทำให้ผลบวกกำลังสองของรัฐยะห่างระหว่างเมตริกซ์สัมประสิทธิ์ สหสัมพันธ์ที่คำนวณได้จากข้อมูล กับเมตริกซ์สัมประสิทธิ์สหสัมพันธ์ที่สร้างขึ้นใหม่ให้มีค่าน้อยที่สุด
  - 2) Generalized Least Square มีหลักเกณฑ์เหมือนวิธี Unweighted Least Square แต่จะมี การถ่วงน้ำหนักค่าลัมประสิทธิ์สหสัมพันธ์ ด้วยค่าผกผันของ Uniques ของตัวแปรนั้น นั่นคือจะให้น้ำหนักแก่ตัวแปรที่มีค่า Unique สูงน้อยกว่าตัวแปรที่มีค่า unique ต่ำ
  - 3) **Maximum Likelihood Method** วิธีนี้กำหนด factor โดยการประมาณค่าพารามิเตอร์ที่ทำให้ เมตริกซ์สัมประสิทธิ์สหสัมพันธ์ที่คำนวณได้ มีค่าใกล้กับเมตริกซ์ที่ได้จากข้อมูล โดยมีเงื่อนไขว่า ข้อมูลตัวอย่างนั้น (ตัวแปร) ต้องมีการแจกแจงแบบ Multivariate Normal
  - 4) Alpha Method
  - 5) Image Factoring

# Factor analysis: Extraction



## Display

- ✓ Unrotate factor solution เมื่อต้องการให้แสดงผลลัพธ์ของ Factor โดยไม่มีการหมุนแกนปัจจัย โดยผลลัพธ์จะแสดงค่า communality , eigenvalues
- ✓ Scree plot แสดงกราฟค่า eigenvalues โดยเรียงลำดับจากมากไปน้อย โดยใช้ Factor ที่หมุนแกนปัจจัยแล้ว

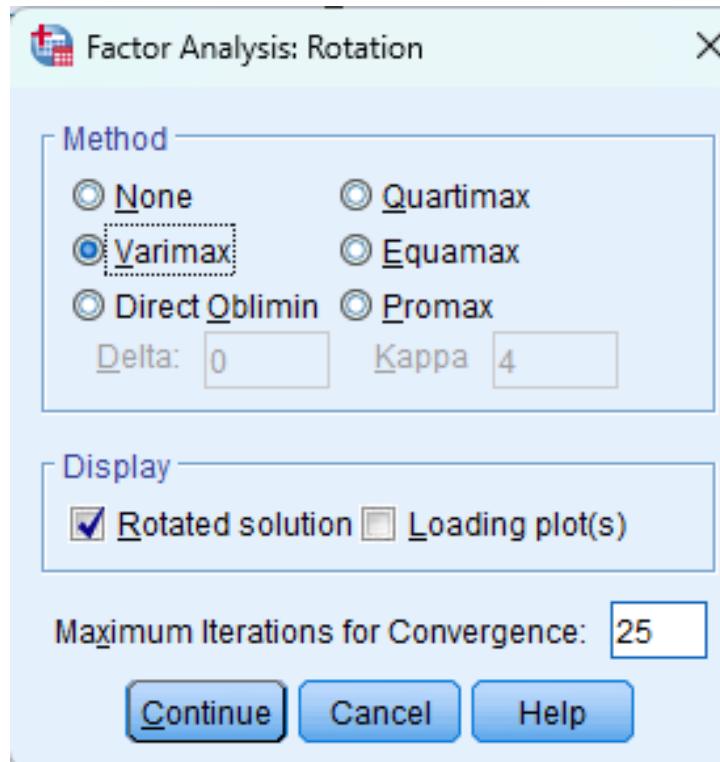
## Extract

- Eigenvalues over: ระบุค่า eigenvalues = 1
- Number of factors: ใส่จำนวน Factor ที่ต้องการ

## Maximum Iterations for Convergence

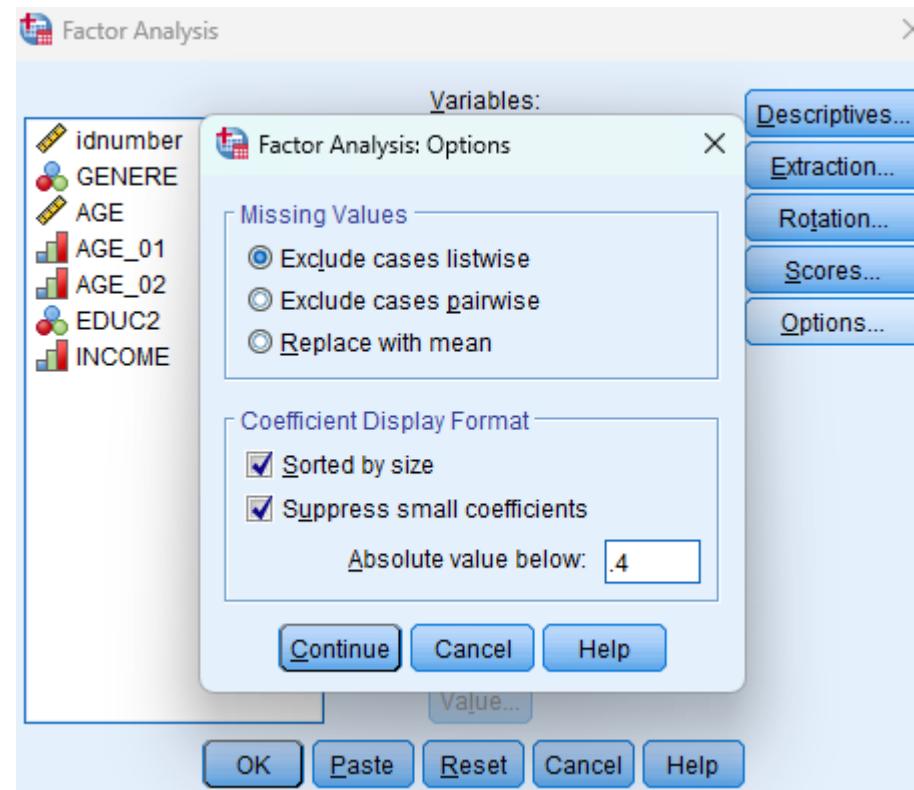
- กำหนดจำนวนรอบสูงสุดของการสกัดปัจจัยโดยโปรแกรม SPSS กำหนดเป็น 25 รอบ หรือเปลี่ยนมากกว่านั้นก็ได้หากข้อมูลขนาดใหญ่

# Factor analysis: Rotation



- 1. Orthogonal Rotation** เหมาะสำหรับ EFA หมุนแกนให้ factor ตั้งฉากกัน ทำให้ factors เป็นอิสระต่อกัน
  - 1) Varimax** เป็นเทคนิคที่ทำให้มีจำนวนตัวแปรที่น้อยที่สุด มีค่า Factor loading มากในแต่ละปัจจัย จึงเป็นวิธีที่นิยมใช้มากที่สุด
  - 2) Quartimax เป็นวิธีที่หมุนแกนปัจจัย โดยจะพยายามทำให้มีจำนวนปัจจัยน้อยที่สุด ใน การอธิบายตัวแปรแต่ละตัว
  - 3) Equamax เป็นเทคนิคที่ใช้เกณฑ์ทั้งของ Varimax และ Quartimax
- 2. Oblique Rotation** เหมาะสำหรับ CFA หมุนแกนเป็นมุ่ง  
แหลม ส่งผลให้ factors สัมพันธ์กัน  
อนุญาตให้ factor ที่กำหนดไม่เป็นอิสระกัน
  - 1) Direct Oblimin
  - 2) Promax

# Factor analysis: Options...



## Missing

- Exclude case listwise จะวิเคราะห์เฉพาะ case ที่มีค่าของทุกตัวแปร
- Exclude case pairwise จะไม่รวม case ที่มี missing ของตัวแปรคู่ได้คู่หนึ่ง
- Replace with mean แทนค่า missing value ด้วยค่าเฉลี่ยของตัวแปรนั้น ๆ และใช้ทุก case ในการวิเคราะห์ปัจจัย

## Coefficient Display Format แสดงค่าสัมประสิทธิ์

- Sorted by size จะแสดงค่า Factor loading เรียงตามลำดับ โดยตัวแปรที่มีค่า Factor loading สูง ๆ ในปัจจัยเดียวกัน จะอยู่ด้วยกัน
- Suppress small coefficients
  - Absolute value below: ..... ระบุค่าที่ต้องการจะไม่แสดงค่าสัมประสิทธิ์สหสัมพันธ์ หรือ Factor loading ที่มีค่าน้อยกว่าที่ระบุ โดยค่าที่จะระบุมีค่า 0 ถึง 1 แนะนำ .3 ขึ้นไปเป็นอย่างน้อย

# Convergent validity

- Convergent validity means that the variables within a single factor are highly correlated. This is evident by the factor loadings.
- Sufficient/significant loadings depend on the sample size of your dataset.
- The table outlines the thresholds for sufficient/significant factor loadings.
- Generally, the smaller the sample size, the higher the required loading.
- **Regardless of sample size, it is best to have loadings greater than 0.500 and averaging out to greater than 0.700 for each factor.**



## The thresholds for sufficient/significant factor loadings

Sample size	Sufficient factor loading
50	0.75
60	0.70
70	0.65
85	0.60
100	0.55
120	0.50
150	0.45
200	0.40
250	0.35
350	0.30

Note: Correlation matrix => highly correlated variables indicate that factor analysis may be an appropriate multivariate statistical technique to explore these variables.

# Output

## Descriptive Statistics

	Mean	Std. Deviation	Analysis N
HEALT	4.57	.671	735
BAEAUT	3.33	1.031	735
ANTIA	4.06	.878	735
LISTY	3.52	.920	735
TREND	3.16	.952	735
FSAFE	4.61	.661	735
POLIC	3.96	.967	735
TRADI	3.90	.858	735
AGRIC	3.71	.912	735
PRIND	3.56	.970	735
COUIS	3.60	.926	735

Correlation Matrix<sup>a</sup>

	HEALT	BAEAUT	ANTIA	LISTY	TREND	FSAFE	POLIC	TRADI	AGRIC	PRIND	COUIS	
Correlation	HEALT	1.000	.260	.398	.015	.017	.408	.248	.101	.164	.195	.055
	BAEAUT	.260	1.000	.435	.225	.291	.155	.136	.212	.121	.139	.280
	ANTIA	.398	.435	1.000	.155	.185	.257	.200	.151	.183	.201	.154
	LISTY	.015	.225	.155	1.000	.411	.120	.119	.428	.256	.205	.429
	TREND	.017	.291	.185	.411	1.000	.123	.237	.346	.276	.211	.323
	FSAFE	.408	.155	.257	.120	.123	1.000	.378	.219	.291	.318	.139
	POLIC	.248	.136	.200	.119	.237	.378	1.000	.303	.421	.412	.087
	TRADI	.101	.212	.151	.428	.346	.219	.303	1.000	.408	.366	.468
	AGRIC	.164	.121	.183	.256	.276	.291	.421	.408	1.000	.772	.269
	PRIND	.195	.139	.201	.205	.211	.318	.412	.366	.772	1.000	.255
	COUIS	.055	.280	.154	.429	.323	.139	.087	.468	.269	.255	1.000

a. Determinant = .043

ค่า determinant > 0 แสดงว่า การวิเคราะห์องค์ประกอบไม่มีปัญหา

# Output: KMO and Barlett's Test

## KMO and Bartlett's Test

Kaiser-Meyer-Olkin Measure of Sampling Adequacy.	.778
Bartlett's Test of Sphericity	
Approx. Chi-Square	2293.830
df	55
Sig.	.000

Note:

- Kaiser-Meyer-Olkin value
  - มีค่าระหว่าง 0-1 ยิ่งเข้าใกล้ 1 ยิ่งดี
  - ค่าที่แนะนำคือ  $>.6$
- Bartlett's Test of Sphericity
  - Bartlett's test should be less (p-value  $<.05$ ).
  - This tests the null hypothesis that the correlation matrix is an identity matrix.
  - You want to reject this null hypothesis (Sig.  $< .05$ ).

## KMO statistics

- Marvelous .90s
- Meritorious .80s
- Middling .70s
- **Mediocre .60s**
- Miserable .50s
- Unacceptable <.50

# Output: Communalities

## Communalities แสดงค่าความร่วมกัน

- ค่า initial เป็นค่าเริ่มต้น/ค่าเดิมก่อนถูกแบ่งใน factor
- Extraction เป็นค่าที่อธิบายว่าตัวแปรนั้นสามารถอธิบายตัวแปรเดิมได้มากเพียงใด
  - ค่า Extraction ยิ่งมากยิ่งดี
  - ค่า Extraction น้อย สามารถเป็นตัวชี้วัดได้ว่า ตัวแปรนั้นควรปรับออกจากการวิเคราะห์ (ค่าระหว่าง 0.0-0.4)
- เช่น ตัวแปร Health หลังจากสกัดปัจจัยแล้ว ความแปรปรวนของตัวแปรถูกอธิบายโดยปัจจัยได้ 64.2%

	Initial	Extraction
HEALT	1.000	.642
BAEAUT	1.000	.617
ANTIA	1.000	.625
LISTY	1.000	.586
TREND	1.000	.475
FSAFE	1.000	.492
POLIC	1.000	.510
TRADI	1.000	.566
AGRIC	1.000	.761
PRIND	1.000	.737
COUIS	1.000	.556

Extraction Method: Principal Component Analysis.

**Component Matrix<sup>a</sup>**

	Component		
	1	2	3
AGRIC	.718		-.493
PRIND	.700		-.480
TRADI	.669		
POLIC	.577		
COUIS	.558	-.467	
TREND	.546		
LISTY	.539	-.520	
FSAFE	.528	.461	
HEALT	.403	.626	
BAEAUT	.478		.623
ANTIA	.487		.521

Extraction Method: Principal Component Analysis.

a. 3 components extracted.

# ควรเลือกจำนวน factor เท่าไร

Total Variance Explained										
Component	Initial Eigenvalues			Extraction Sums of Squared Loadings			Rotation Sums of Squared Loadings			
	Total	% of Variance	Cumulative %	Total	% of Variance	Cumulative %	Total	% of Variance	Cumulative %	
1	3.594	32.674	32.674	3.594	32.674	32.674	2.364	21.487	21.487	
2	1.586	14.416	47.090	1.586	14.416	47.090	2.352	21.386	42.873	
3	1.387	12.610	59.701	1.387	12.610	59.701	1.851	16.828	59.701	
4	.803	7.301	67.001							
5	.778	7.072	74.073							
6	.602	5.475	79.548							
7	.561	5.104	84.652							
8	.515	4.679	89.331							
9	.494	4.494	93.824							
10	.458	4.161	97.985							
11	.222	2.015	100.000							

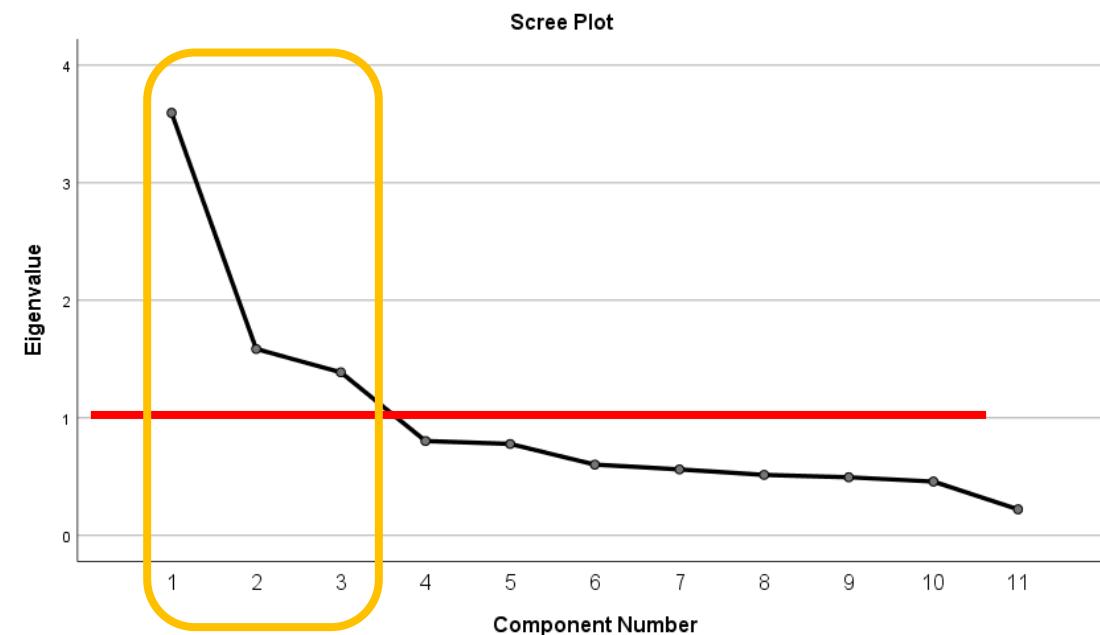
Extraction Method: Principal Component Analysis.

$$\text{Initial Eigenvalue} = \frac{\text{total variance}}{\text{total no. of component}}$$

Factor1 accounts for 32.67% of the total variance (3.594/11)

Factor2 accounts for 14.41% of the total variance (1.586/11)

Factor3 accounts for 12.61% of the total variance (1.383/11)



- Total หมายถึงค่า Eigenvalues ซึ่งคือความแปรปรวนทั้งหมดของตัวแปรเดิมที่อธิบายได้โดยปัจจัยนั้น ๆ เช่นปัจจัย 1 มีค่า Eigenvalues เท่ากับ 3.594 แสดงว่าปัจจัย 1 สามารถนำมาแทนตัวแปรเดิมได้ 3.594 ตัว (ดังนั้นจึงควรพิจารณาปัจจัยที่ Eigenvalues มากกว่า 1)
- 3 factors แรก สามารถอธิบายได้มากถึง 59.7% ของความแปรปรวนทั้งหมด
  - หมายความว่า ตัวแปรใหม่ที่ได้สามารถครอบคลุมข้อมูลเดิมของตัวแปรได้ 59.7%

# Output: Component and Rotated Component Matrix

**Component Matrix<sup>a</sup>**

	Component		
	1	2	3
AGRIC	.718		-.493
PRIND	.700		-.480
TRADI	.669		
POLIC	.577		
COUIS	.558	-.467	
TREND	.546		
LISTY	.539	-.520	
FSAFE	.528	.461	
HEALT	.403	.626	
BAEAUT	.478		.623
ANTIA	.487		.521

Extraction Method: Principal Component Analysis.

a. 3 components extracted.

**Rotated Component Matrix<sup>a</sup>**

	Component		
	1	2	3
PRIND	.828		
AGRIC	.822		
POLIC	.676		
FSAFE	.521		.469
LISTY		.761	
COUIS		.735	
TREND		.668	
TRADI		.643	
ANTIA		.767	
HEALT		.740	
BAEAUT		.429	.651

Extraction Method: Principal Component Analysis.

Rotation Method: Varimax with Kaiser Normalization.

a. Rotation converged in 6 iterations.

- ค่าที่แสดงเป็นค่า Factor Loading หรือค่าที่แสดงความสัมพันธ์ระหว่างตัวแปรแต่ละตัวกับ Factor
- เมื่อ ตัวแปร FSAFE มีค่า Factor Loading ของปัจจัยแรกเท่ากับ .521 มากกว่าค่า Factor Loading ของปัจจัยที่ 3 ที่ .469
- Factor loading พิจารณาเปรียบเทียบเฉพาะค่า absolute ไม่พิจารณาเครื่องหมาย

**Component Transformation Matrix**

Component	1	2	3
1	.658	.623	.424
2	.329	-.744	.582
3	-.678	.243	.694

Extraction Method: Principal Component Analysis.

Rotation Method: Varimax with Kaiser Normalization.

# การตั้งชื่อตัวแปร

**Rotated Component Matrix<sup>a</sup>**

	Component		
	1	2	3
PRIND	.828		
AGRIC	.822		
POLIC	.676		
FSAFE	.521		.469
LISTY		.761	
COUIS		.735	
TREND		.668	
TRADI		.643	
ANTIA		.767	
HEALT		.740	
BAEAUT		.429	.651

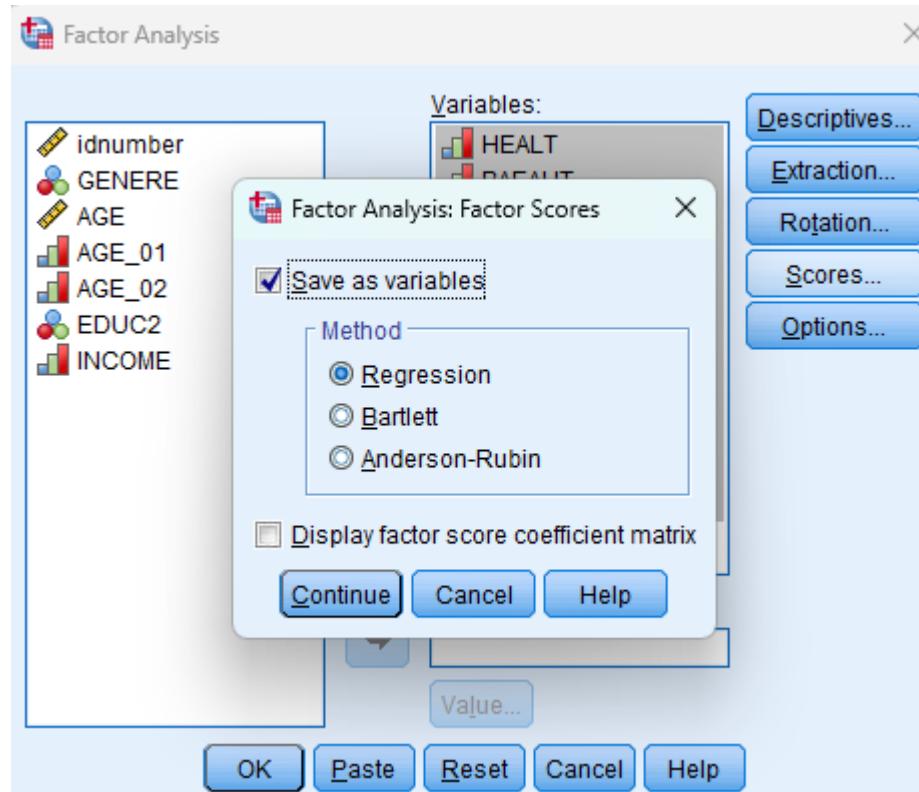
Extraction Method: Principal Component Analysis.

Rotation Method: Varimax with Kaiser Normalization.

a. Rotation converged in 6 iterations.

Information about the food processing industry and innovations in terms of products and processes	PRIND
Information about the production techniques used in the primary sector	AGRIC
Information about food regulations, affecting consumer choices and the food industry	POLIC
Information about food safety issues caused by bacteria and other substances	FSAFE
Information about lifestyle, food tourism and eating out	LISTY
Information about Italian and international cuisine, food culture and good living	COUIS
information about trends, consumption evolution, food fads, and underscoring ethnicity, cultural, social diversity of Italian polulation	TREND
Information about tradition, regional typical products and quality foods that are disappearing from the litalian market	TRADI
Information about food containing anti ageing properties	ANTIA
Information about health risks caused by obesity, anorexia nervosa bulimia and other illesses liked to food	HEALT
Information about foods related to health and beauty	BEAUT

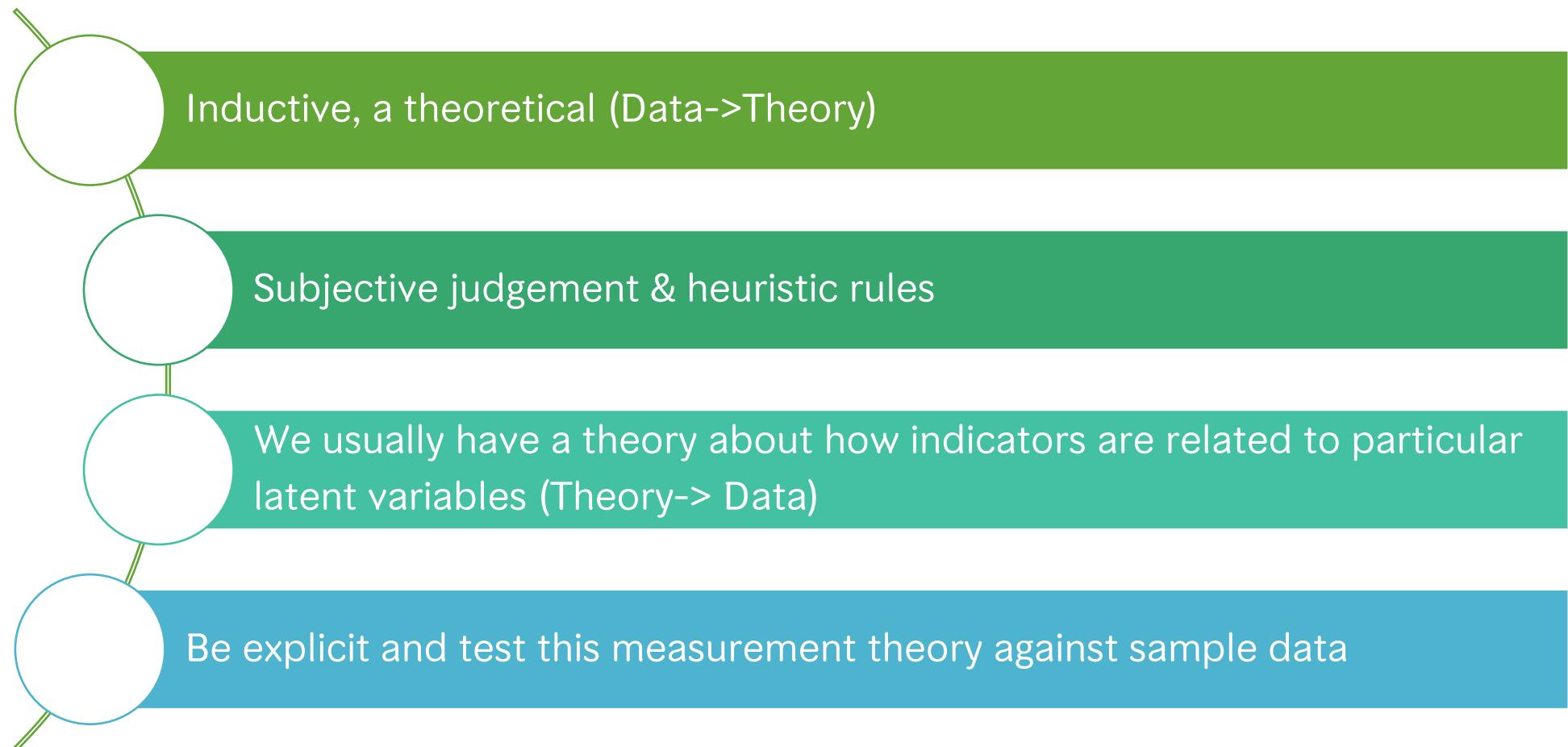
# Factor analysis: Scores...



## Save as variables

- เมื่อเลือกทางเลือกนี้จะเป็นการ save Factor score ในรูปของตัวแปร (1 factor = 1 new variable)
- Factor score มีวิธีการคำนวณให้เลือก 3 วิธี
  - Regression** โดยวิธีนี้ให้ค่าแปรปรวนเท่ากับ (สัมประสิทธิ์สหสัมพันธ์ระหว่างค่า Factor score ที่ประมาณได้กับค่า Factor score จริง)
  - Bartlett
  - Anderson-Rubin

# Limitations of EFA



# Factor analysis via Stata

Suwanna Sayruamyat

# Syntax

- Factor analysis of data

**factor** varlist [if] [in] [weight] [, method options]

# ตัวอย่างการวิเคราะห์

## วิเคราะห์ตัวแปร v1 to v5

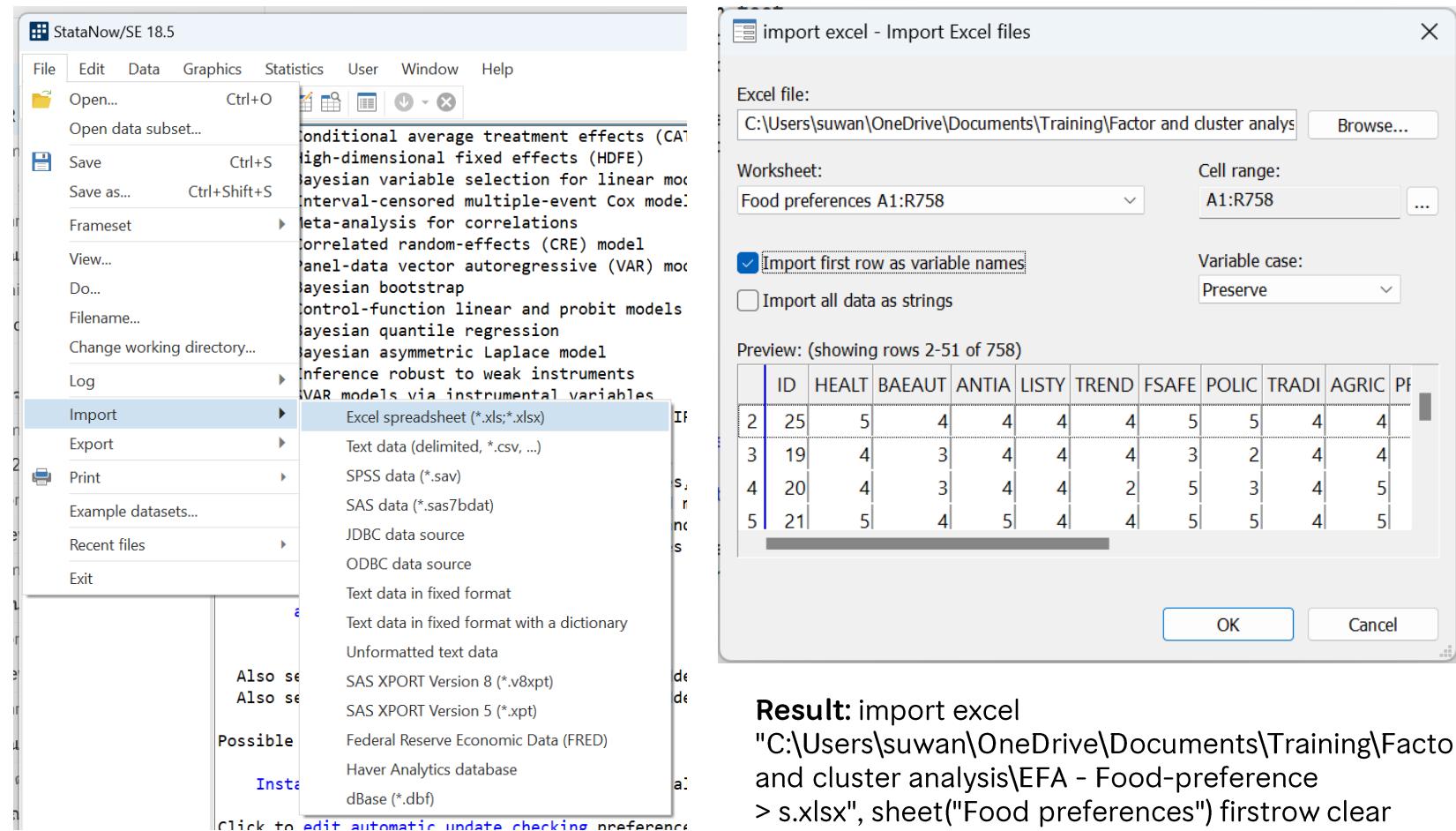
### Syntax

Principal-factor analysis ด้วย	<ul style="list-style-type: none"> <li>factor v1 v2 v3 v4 v5</li> <li>factor v1-v5</li> </ul>
Principal-factor analysis และกำหนดให้จัดกลุ่ม 3 กลุ่ม	factor v1-v5, factors(3)
Principal-component factor analysis	factor v1-v5, pcf
Maximum-likelihood factor analysis	factor v1-v5, ml

- After running factor you need to rotate the factor loads to get a clearer pattern, just type **rotate** to get a final solution.
- To create the new variables, after **factor, rotate** you type **predict**.
  - predict factor1 factor2 /\*or whatever name you prefer to identify the factors\*/**

# Practice: Factor analysis in stata

- Import excel file to Stata
  - File > Import > Excel spreadsheet (\*.xls; \*.xlsx)
  - หน้าต่าง import excel
    - เลือก Browse... เลือกไฟล์ที่ต้องการ
    - เลือก Import first row as variable names
    - OK



**Result:** import excel

"C:\Users\suwan\OneDrive\Documents\Training\Factor and cluster analysis\EFA - Food-preference > s.xlsx", sheet("Food preferences") firstrow clear (18 vars, 757 obs)

# Principal-factor analysis

## Command

```
. factor HEALT BAEAUT ANTIA LISTY TREND FSAFE POLIC TRADI AGRIC PRIND COUIS  
(obs=757)
```

วิธีที่วิเคราะห์

Factor analysis/correlation  
Method: principal factors  
Rotation: (unrotated)

Number of obs = 757  
Retained factors = 7  
Number of params = 55

จำนวนตัวอย่างในการวิเคราะห์

ผลการวิเคราะห์เบื้องต้น ค่า Eigenvalue > 1 มีเพียง factor1 ที่มีค่ามากกว่า 1

- Eigenvalue คือ ผลรวมความแปรปรวนของแต่ละ Factor โดยที่ผลรวมของ eigenvalues = ผลรวมความแปรปรวนของ Total factors.
- Kaiser criterion suggests to retain those factors with eigenvalues  $\geq 1$ .

ความแตกต่างระหว่างค่า Eigenvalue ของ factor นั้น กับค่า Eigenvalue ของ factor ถัดไป

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	1.53937	0.74450	0.5778	0.5778
Factor2	0.79488	0.19475	0.2984	0.8762
Factor3	0.60012	0.10878	0.2253	1.1014
Factor4	0.49134	0.26457	0.1844	1.2859
Factor5	0.22677	0.18208	0.0851	1.3710
Factor6	0.04469	0.03726	0.0168	1.3878
Factor7	0.00742	0.21031	0.0028	1.3905
Factor8	-0.20289	0.02636	-0.0762	1.3144
Factor9	-0.22925	0.04391	-0.0861	1.2283
Factor10	-0.27316	0.06200	-0.1025	1.1258
Factor11	-0.33516	.	-0.1258	1.0000

LR test: independent vs. saturated:  $\chi^2(55) = 1235.60$  Prob>chi2 = 0.0000

- ค่า Factor loading แสดงน้ำหนักและ correlation ระหว่างตัวแปรกับแต่ละ factor model
- ค่า factor loading ยิ่งมากยิ่งมีบทบาทใน factor model
- ค่า factor loading  $< 0$  แสดงถึงผิดเบี่ยงต่อ factor model นั้น

Factor loadings (pattern matrix) and unique variances

Variable	Factor1	Factor2	Factor3	Factor4	Factor5	Factor6	Factor7	Uniqueness
HEALT	0.0253	0.0247	0.5117	0.0791	0.0752	-0.0545	0.0104	0.7219
BAEAUT	0.0090	0.6528	-0.0361	0.0599	0.0265	-0.0001	0.0117	0.5680
ANTIA	-0.0103	0.4329	-0.1270	0.0598	0.2768	0.0003	-0.0244	0.7156
LISTY	0.0656	-0.0070	0.1049	0.0536	-0.0173	0.1932	-0.0147	0.9439
TREND	0.0738	0.4032	0.0280	-0.0387	-0.3315	0.0106	0.0176	0.7194
FSAFE	0.0297	0.0818	0.5459	0.0331	0.0275	0.0177	-0.0045	0.6922
POLIC	0.5966	-0.0634	-0.0738	0.3527	-0.0066	0.0021	-0.0009	0.5101
TRADI	0.4354	-0.0269	-0.0398	-0.1089	0.1129	0.0198	0.0660	0.7788
AGRIC	0.7324	-0.0186	0.0016	-0.1885	0.0785	0.0202	-0.0042	0.4211
PRIND	0.6445	0.0597	0.0448	-0.1842	-0.1056	-0.0522	-0.0407	0.5296
COUIS	0.1746	-0.0519	-0.0447	0.5165	-0.0522	-0.0206	0.0015	0.6949

Cumulative แสดงถึงค่าความแปรปรวนสะสมจากค่าความแปรปรวนรวมทั้งหมด (total variance)

- Uniqueness คือความแปรปรวนของตัวแปรนั้น ๆ ที่ไม่ได้แชร์กับตัวแปรตัวอื่น เทียบได้กับค่า Communalities ในผลของ SPSS
  - เช่น HEALT มีค่า Uniqueness = 72.19% หมายความว่า ตัวแปรนี้มีความเป็นตัวของตัวเองถึง 72.19% ใน Overall factor model.
  - Note: ยิ่ง Uniqueness มีค่ามาก ยิ่งเกี่ยวข้องกับตัวแปรอื่นใน factor model น้อย

# Rotate

MAEAE & MAB Training

Rotation method: ด้วยวิธี

orthogonal varimax (by default)

- วิธีนี้อยู่บนข้อสมมติว่า factor ที่สร้างขึ้นจะเป็นอิสระต่อกัน (not correlated to each other)
- หมาย味着รับสร้างตัวชี้วัด หรือตัวแปรใหม่ที่ต้องความเป็นอิสระต่อกัน

Factor analysis/correlation		Number of obs = 757		
Method: principal factors		Retained factors = 7		
Rotation: orthogonal varimax (Kaiser off)		Number of params = 55		
Factor	Variance	Difference	Proportion	Cumulative
Factor1	1.41688	0.67155	0.5318	0.5318
Factor2	0.74533	0.13271	0.2798	0.8116
Factor3	0.61262	0.01726	0.2300	1.0416
Factor4	0.59536	0.31569	0.2235	1.2650
Factor5	0.27967	0.23401	0.1050	1.3700
Factor6	0.04566	0.03661	0.0171	1.3871
Factor7	0.00906	.	0.0034	1.3905

LR test: independent vs. saturated: chi2(55) = 1235.60 Prob>chi2 = 0.0000

Variable	Factor1	Factor2	Factor3	Factor4	Factor5	Factor6	Factor7	Uniqueness
HEALT	0.0065	-0.0116	0.0211	0.5251	-0.0246	-0.0334	0.0052	0.7219
BAEAUT	-0.0035	0.6365	-0.0034	0.0306	0.1609	-0.0026	0.0059	0.5680
ANTIA	-0.0113	0.5086	-0.0009	-0.0512	-0.1512	0.0021	-0.0132	0.7156
LISTY	0.0423	-0.0179	0.0589	0.1009	0.0248	0.1993	0.0028	0.9439
TREND	0.0697	0.2837	-0.0247	0.0231	0.4406	0.0038	-0.0031	0.7194
FSAFE	0.0243	0.0203	-0.0266	0.5501	0.0435	0.0392	-0.0057	0.6922
POLIC	0.4408	-0.0075	0.5424	-0.0217	-0.0222	0.0163	0.0066	0.5101
TRADI	0.4495	0.0042	0.0434	-0.0314	-0.0958	0.0174	0.0828	0.7788
AGRIC	0.7557	-0.0070	0.0671	0.0037	-0.0481	0.0280	0.0166	0.4211
PRIND	0.6648	0.0106	0.0463	0.0339	0.1471	-0.0430	-0.0397	0.5296
COUIS	-0.0127	-0.0008	0.5519	0.0126	0.0052	-0.0085	-0.0048	0.6949

	Factor1	Factor2	Factor3	Factor4	Factor5	Factor6	Factor7
Factor1	0.9402	0.0126	0.3376	0.0253	0.0281	0.0144	0.0168
Factor2	0.0114	0.9470	-0.0972	0.0854	0.2936	-0.0030	-0.0107
Factor3	0.0095	-0.1257	-0.1037	0.9822	0.0853	0.0362	-0.0069
Factor4	-0.3356	0.1062	0.9263	0.1188	-0.0577	0.0230	-0.0100
Factor5	0.0513	0.2755	-0.0857	0.1079	-0.9482	0.0151	0.0563
Factor6	-0.0091	0.0001	-0.0207	-0.0397	0.0184	0.9942	0.0956
Factor7	-0.0211	-0.0054	0.0087	0.0062	0.0546	-0.0964	0.9935

ตาราง factor rotation matrix นี้แสดง correlation matrix ระหว่าง factor 1 to factor 7

ตาราง Rotate factor loadings หรือ pattern matrix แสดงค่า factor loadings ที่ผ่านการหมุนแล้วที่มีความชัดเจนยิ่ง โดยแต่ละตัวแปรจะมีความชัดเจนว่าหมาย含義 สมที่จะอธิบาย factor model ได้

- เช่น HEALT มีค่า factor loading สูงที่สุดใน Factor4

NOTE: หากต้องการการหมุนด้วย Oblique rotation ให้พิมพ์คำสั่ง

# Oblique rotation ให้พิมพ์คำสั่ง Rotate, promax

Rotation method: ด้วยวิธี Oblique rotation

- วิธีนี้อยู่บนข้อสมมติว่า factor ที่สร้างขึ้น **ไม่เป็นอิสระต่อกัน**
- เหมาะสมสำหรับการวิเคราะห์สมการเชิงโครงสร้าง (Structural equation modelling: SEM)

Factor analysis/correlation			Number of obs = 757
Method: principal factors			Retained factors = 7
Rotation: oblique promax (Kaiser off)			Number of params = 55
Factor	Variance	Proportion	Rotated factors are correlated
Factor1	1.46331	0.5493	
Factor2	1.00100	0.3757	
Factor3	0.78654	0.2952	
Factor4	0.73739	0.2768	
Factor5	0.59978	0.2251	
Factor6	0.40796	0.1531	
Factor7	0.18399	0.0691	

LR test: independent vs. saturated: chi2(55) = 1235.60 Prob>chi2 = 0.0000

Rotated factor loadings (pattern matrix) and unique variances

Variable	Factor1	Factor2	Factor3	Factor4	Factor5	Factor6	Factor7	Uniqueness
HEALT	-0.0139	0.0170	0.0295	0.0033	0.5446	-0.0391	-0.0396	0.7219
BAEAUT	-0.0432	0.0252	0.0097	0.6032	0.0322	0.1948	-0.0074	0.5680
ANTIA	0.0490	-0.0288	-0.0066	0.5352	-0.0252	-0.1580	0.0021	0.7156
LISTY	0.0173	-0.0112	0.0166	-0.0179	0.0272	0.0185	0.2190	0.9439
TREND	0.0020	-0.0009	-0.0064	0.1932	-0.0224	0.4635	0.0087	0.7194
FSAFE	0.0255	-0.0169	-0.0335	0.0207	0.5331	0.0144	0.0451	0.6922
POLIC	0.3160	0.0213	0.5247	0.0015	-0.0239	-0.0147	0.0152	0.5101
TRADI	0.2735	0.2312	0.0009	0.0277	-0.0026	0.0008	-0.0125	0.7788
AGRIC	0.7134	0.0631	0.0003	-0.0017	0.0014	-0.0447	0.0242	0.4211
PRIND	0.7165	-0.0813	0.0115	-0.0289	0.0182	0.0890	-0.0312	0.5296
COUIS	-0.1197	-0.0205	0.5781	0.0065	0.0148	0.0084	-0.0076	0.6949

Factor rotation matrix

	Factor1	Factor2	Factor3	Factor4	Factor5	Factor6	Factor7
Factor1	0.9632	0.7542	0.5274	-0.0099	0.0456	0.1663	0.2017
Factor2	-0.0153	0.1513	-0.0206	0.0051	-0.0387	0.0022	0.8979
Factor3	0.0112	-0.1381	-0.1014	-0.1534	0.9864	0.1642	0.3300
Factor4	0.0301	-0.2654	-0.1129	0.9369	0.0883	0.4816	-0.0143
Factor5	-0.2642	-0.1622	0.8333	0.1077	0.1040	-0.1149	0.1995
Factor6	0.0185	0.4369	-0.0635	0.2947	0.0688	-0.8356	0.0331
Factor7	-0.0294	0.3188	0.0063	-0.0118	0.0073	0.0450	-0.0549

เนื่องจาก factor model มีความสัมพันธ์กัน ตรงนี้จะไม่สะท้อนค่าความแปรป่วนได้ ออกมานะ

factor HEALT BAEAUT ANTIA LISTY TREND FSAFE POLIC TRADI AGRIC PRIND COUIS  
 rotate, promax  
**predict factor1 factor2**

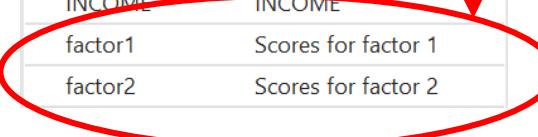
Variables	
Name	Label
ID	ID
HEALT	HEALT
BAEAUT	BAEAUT
ANTIA	ANTIA
LISTY	LISTY
TREND	TREND
FSAFE	FSAFE
POLIC	POLIC
TRADI	TRADI
AGRIC	AGRIC
PRIND	PRIND
COUIS	COUIS
GENDER	GENDER
AGE	AGE
AGE_01	AGE_01
AGE_02	AGE_02
EDUC2	EDUC2
INCOME	INCOME
factor1	Scores for factor 1
factor2	Scores for factor 2

ประมาณการค่า factor1 และ factor3 ด้วย  
 วิธี Regression จากข้อมูลการหมุนด้วยวิธี  
 promax

. predict factor1 factor2  
 (option regression assumed; regression scoring)

Scoring coefficients (method = regression; based on promax(3) rotated factors)

Variable	Factor1	Factor2	Factor3	Factor4	Factor5	Factor6	Factor7
HEALT	0.00391	-0.02579	0.01180	-0.01421	0.36361	0.00046	0.07525
BAEAUT	0.00516	-0.11028	-0.00757	0.45865	0.02083	0.19641	-0.00934
ANTIA	-0.00394	0.04446	-0.00090	0.32701	-0.04414	-0.12803	-0.01277
LISTY	0.00639	0.01485	0.02565	-0.01005	0.05782	0.02396	0.19857
TREND	0.03742	-0.16137	-0.01392	0.11308	0.01528	0.38369	0.00994
FSAFE	0.00576	-0.05101	-0.01241	-0.00448	0.39330	0.06148	0.16075
POLIC	0.16451	0.16784	0.46225	-0.00508	-0.02123	-0.02424	0.11326
TRADI	0.15238	0.19556	0.01123	0.00984	-0.02190	-0.05735	0.02062
AGRIC	0.44454	0.40856	0.03014	-0.01069	0.00197	-0.02864	0.08748
PRIND	0.31692	0.11674	0.00761	-0.01532	0.03591	0.22450	-0.04568
COUIS	-0.03264	-0.02254	0.34921	0.00179	0.01078	-0.00665	0.05343



# Principal Component analysis

MAEAE & MAB Training

Command

```
. factor HEALT BAEAUT ANTIA LISTY TREND FSAFE POLIC TRADI AGRIC PRIND COUIS, pcf  
(obs=757)
```

วิธีที่วิเคราะห์

Factor analysis/correlation		Number of obs = 757
Method: principal-component factors		Retained factors = 5
Rotation: (unrotated)		Number of params = 45

- ผลการวิเคราะห์เบื้องต้น ค่า Eigenvalue  $> 1$
- บ่งบอกจำนวน Factor model ที่เหมาะสมเท่ากับ 5

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	2.16996	0.62586	0.1973	0.1973
Factor2	1.54410	0.11877	0.1404	0.3376
Factor3	1.42533	0.19954	0.1296	0.4672
Factor4	1.22579	0.19904	0.1114	0.5787
Factor5	1.02675	0.03568	0.0933	0.6720
Factor6	0.55107	0.21039	0.0501	0.7221
Factor7	0.77468	0.19162	0.0704	0.8325
Factor8	0.58306	0.09591	0.0530	0.8855
Factor9	0.48715	0.06443	0.0443	0.9298
Factor10	0.42273	0.07334	0.0384	0.9682
Factor11	0.34939	.	0.0318	1.0000

LR test: independent vs. saturated: chi2(55) = 1235.60 Prob>chi2 = 0.0000

Factor loadings (pattern matrix) and unique variances

Variable	Factor1	Factor2	Factor3	Factor4	Factor5	Uniqueness
HEALT	0.0389	0.0536	0.7966	0.0547	-0.2395	0.3006
BAEAUT	0.0144	0.8685	-0.0638	0.0839	-0.0466	0.2322
ANTIA	-0.0136	0.6322	-0.1955	0.1271	-0.5500	0.2433
LISTY	0.1007	-0.0052	0.2210	0.1199	0.3836	0.7795
TREND	0.0996	0.5898	0.0183	-0.0826	0.6665	0.1909
FSAFE	0.0428	0.1383	0.8295	-0.0166	-0.0447	0.2887
POLIC	0.7164	-0.0820	-0.0644	0.4703	-0.0001	0.2548
TRADI	0.5834	-0.0456	-0.0764	-0.2348	-0.2203	0.5481
AGRIC	0.8228	-0.0267	-0.0202	-0.2584	-0.0840	0.2480
PRIND	0.7471	0.0732	0.0241	-0.2695	0.1110	0.3510
COUIS	0.2397	-0.0731	-0.0125	0.8732	0.0593	0.1711

# PCA, rotated by varimax

Factor analysis/correlation  
 Method: principal-component factors  
 Rotation: orthogonal varimax (Kaiser off)

Number of obs = 757  
 Retained factors = 5  
 Number of params = 45

Factor	Variance	Difference	Proportion	Cumulative
Factor1	2.04013	0.62275	0.1855	0.1855
Factor2	1.41738	0.04479	0.1289	0.3143
Factor3	1.37259	0.02437	0.1248	0.4391
Factor4	1.34822	0.13461	0.1226	0.5617
Factor5	1.21361	.	0.1103	0.6720

LR test: independent vs. saturated: chi2(55) = 1235.60 Prob>chi2 = 0.0000

Rotated factor loadings (pattern matrix) and unique variances

Variable	Factor1	Factor2	Factor3	Factor4	Factor5	Uniqueness
HEALT	0.0072	0.8279	0.0309	0.0240	-0.1112	0.3006
BAEAUT	-0.0192	0.0494	0.7402	0.0061	0.4659	0.2322
ANTIA	-0.0038	-0.0249	0.8648	0.0166	-0.0889	0.2433
LISTRY	0.0053	0.1611	-0.2474	0.1752	0.3203	0.7795
TREND	0.0500	-0.0203	0.0817	-0.0356	0.8935	0.1909
FSAFE	0.0149	0.8362	-0.0260	-0.0322	0.1005	0.2887
POLIC	0.5033	-0.0298	-0.0016	0.7005	-0.0199	0.2548
TRADI	0.6501	-0.0289	0.0671	-0.0211	-0.1533	0.5481
AGRIC	0.8654	0.0147	-0.0102	0.0515	-0.0101	0.2480
PRIND	0.7771	0.0348	-0.0490	0.0211	0.2025	0.3510
COUIS	-0.0901	0.0065	0.0168	0.9056	-0.0185	0.1711

Factor rotation matrix

	Factor1	Factor2	Factor3	Factor4	Factor5
Factor1	0.9295	0.0391	-0.0096	0.3586	0.0756
Factor2	-0.0115	0.1164	0.7964	-0.0854	0.5872
Factor3	-0.0404	0.9786	-0.1913	-0.0200	0.0616
Factor4	-0.3534	0.0315	0.1237	0.9256	-0.0464
Factor5	-0.0965	-0.1623	-0.5602	0.0838	0.8022

Cumulative แสดงถึงค่าความแปรปรวนสะสม  
 จากค่าความแปรปรวนรวมทั้งหมด (total variance)

ณ ที่นี่ ค่า Cumulative สะสมของ factor1 – factor 5 เท่ากับ 0.6720

- หมายความว่า factor ทั้ง 5 สามารถอธิบาย  
 ความแปรปรวนของข้อมูลตั้งต้นได้ 67.20%

# Maximum-likelihood factor analysis

```
. factor HEALT BAEAUT ANTIA LISTY TREND FSAFE POLIC TRADI AGRIC PRIND COUIS, ml
(obs=757)
number of factors adjusted to 6
Iteration 0: Log likelihood = -140.95172
Iteration 1: Log likelihood = -120.99582
Iteration 2: Log likelihood = -81.019606
Iteration 3: Log likelihood = -57.163546
Iteration 4: Log likelihood = -54.597485
Iteration 5: Log likelihood = -27.462889
Iteration 6: Log likelihood = -20.96795
Iteration 7: Log likelihood = -17.697479
Iteration 8: Log likelihood = -14.698728
Iteration 9: Log likelihood = -11.558032
Iteration 10: Log likelihood = -9.7977926
Iteration 11: Log likelihood = -6.0614349
Iteration 12: Log likelihood = -2.2917898
Iteration 13: Log likelihood = -2.2916691
Iteration 14: Log likelihood = -2.2916675
Iteration 15: Log likelihood = -2.2916675
```

ผลการวิเคราะห์เบื้องต้น  
ค่า Eigenvalue > 1 in factor 1 to factor 4

## MAEAE & MAB Training

Factor analysis/correlation	Number of obs = 757
Method: maximum likelihood	Retained factors = 6
Rotation: (unrotated)	Number of params = 51
	Schwarz's BIC = 342.681
Log likelihood = -2.291667	(Akaike's) AIC = 106.583

Warning: Solution is a **Heywood case**; that is, invalid or boundary values of uniqueness.

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	1.66066	0.53104	0.2523	0.2523
Factor2	1.12962	-0.00451	0.1716	0.4239
Factor3	1.13413	-0.13426	0.1723	0.5963
Factor4	1.26839	0.41977	0.1927	0.7890
Factor5	0.84862	0.30826	0.1289	0.9179
Factor6	0.54036	.	0.0821	1.0000

LR test: independent vs. saturated:  $\chi^2(55) = 1235.60$  Prob> $\chi^2 = 0.0000$

LR test: 6 factors vs. saturated:  $\chi^2(4) = 4.53$  Prob> $\chi^2 = 0.3388$

(tests formally not valid because a **Heywood case** was encountered)

### Factor loadings (pattern matrix) and unique variances

Variable	Factor1	Factor2	Factor3	Factor4	Factor5	Factor6	Uniqueness
HEALT	0.0122	0.0982	0.3855	-0.0638	0.0087	-0.0422	0.8356
BAEAUT	-0.0054	0.1033	0.0927	0.6564	0.7416	-0.0000	0.0000
ANTIA	-0.0078	0.0023	-0.0157	0.1565	0.4103	-0.0203	0.8064
LISTY	0.0665	-0.0095	0.1216	0.0102	-0.0527	0.0706	1.0000
TREND	0.0240	0.3551	0.0097	0.8689	-0.3440	0.0000	0.0000
FSAFE	-0.0343	0.2896	0.9525	-0.0864	-0.0151	-0.0000	0.0000
POLIC	0.9867	-0.1544	0.0448	0.0245	-0.0056	-0.0000	0.0000
TRADI	0.2690	0.0974	-0.0513	-0.0404	0.0267	0.3880	0.7627
AGRIC	0.4742	0.3648	-0.0728	-0.1361	0.0770	0.5970	0.2559
PRIND	0.4943	0.8321	-0.1971	-0.1512	0.0386	-0.0000	0.0000
COUIS	0.3730	-0.2010	0.0755	0.0541	-0.0286	-0.1619	0.7848

## factor HEALT BAEAUT ANTIA LISTY TREND FSAFE POLIC TRADI AGRIC PRIND COUIS, ml factor(4)

factor HEALT BAEAUT ANTIA LISTY TREND FSAFE POLIC TRADI AGRIC PRIND COUIS, ml

```
Factor analysis/correlation
Number of obs      =      757
Method: maximum likelihood
Retained factors   =        6
Rotation: (unrotated)
Number of params   =      51
Schwarz's BIC      =    342.681
Log likelihood     = -2.291667
(Akaike's) AIC     =    106.583
```

Warning: Solution is a [Heywood case](#); that is, invalid or boundary values of uniqueness.

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	1.66066	0.53104	0.2523	0.2523
Factor2	1.12962	-0.00451	0.1716	0.4239
Factor3	1.13413	-0.13426	0.1723	0.5963
Factor4	1.26839	0.41977	0.1927	0.7890
Factor5	0.84862	0.30826	0.1289	0.9179
Factor6	0.54036	.	0.0821	1.0000

LR test: independent vs. saturated:  $\chi^2(55) = 1235.60$  Prob> $\chi^2 = 0.0000$

LR test: 6 factors vs. saturated:  $\chi^2(4) = 4.53$  Prob> $\chi^2 = 0.3388$

(tests formally not valid because a [Heywood case](#) was encountered)

Factor loadings (pattern matrix) and unique variances

Variable	Factor1	Factor2	Factor3	Factor4	Factor5	Factor6	Uniqueness
HEALT	0.0122	0.0982	0.3855	-0.0638	0.0087	-0.0422	0.8356
BAEAUT	-0.0054	0.1033	0.0927	0.6564	0.7416	-0.0000	0.0000
ANTIA	-0.0078	0.0023	-0.0157	0.1565	0.4103	-0.0203	0.8064
LISTY	0.0665	-0.0095	0.1216	0.0102	-0.0527	0.0766	1.0000
TREND	0.0240	0.3551	0.0097	0.8689	-0.3440	0.0000	0.0000
FSAFE	-0.0343	0.2896	0.9525	-0.0864	-0.0151	-0.0000	0.0000
POLIC	0.9867	-0.1544	0.0448	0.0245	-0.0056	-0.0000	0.0000
TRADI	0.2690	0.0974	-0.0513	-0.0404	0.0267	0.3880	0.7627
AGRIC	0.4742	0.3648	-0.0728	-0.1361	0.0770	0.5970	0.2559
PRIND	0.4943	0.8321	-0.1971	-0.1512	0.0386	-0.0000	0.0000
COUIS	0.3730	-0.2010	0.0755	0.0541	-0.0286	-0.1619	0.7848

```
Factor analysis/correlation
Number of obs      =      757
Method: maximum likelihood
Retained factors   =        4
Rotation: (unrotated)
Number of params   =      38
Schwarz's BIC      =    328.398
(Akaike's) AIC     =    152.482
```

Warning: Solution is a [Heywood case](#); that is, invalid or boundary values of uniqueness.

Factor	Eigenvalue	Difference	Proportion	Cumulative
Factor1	1.18638	-0.32786	0.2309	0.2309
Factor2	1.51425	0.22757	0.2947	0.5255
Factor3	1.28673	0.13559	0.2504	0.7759
Factor4	1.15134	.	0.2241	1.0000

LR test: independent vs. saturated:  $\chi^2(55) = 1235.60$  Prob> $\chi^2 = 0.0000$

LR test: 4 factors vs. saturated:  $\chi^2(17) = 75.76$  Prob> $\chi^2 = 0.0000$

(tests formally not valid because a [Heywood case](#) was encountered)

Factor loadings (pattern matrix) and unique variances

Variable	Factor1	Factor2	Factor3	Factor4	Uniqueness
HEALT	0.3972	0.0581	-0.0156	-0.0336	0.8375
BAEAUT	0.0650	0.0009	0.9979	0.0000	0.0000
ANTIA	-0.0271	-0.0109	0.4085	0.0014	0.8323
LISTY	0.1041	0.0855	-0.0293	0.0246	0.9804
TREND	0.0468	-0.0027	0.3504	0.0183	0.8747
FSAFE	0.9970	0.0760	-0.0144	-0.0000	0.0000
POLIC	-0.1137	0.9935	0.0013	-0.0000	0.0000
TRADI	-0.0455	0.2433	-0.0001	0.3586	0.8101
AGRIC	-0.0004	0.4071	-0.0042	0.8469	0.1169
PRIND	0.0220	0.3513	-0.0074	0.5105	0.6154
COUIS	-0.0340	0.4027	0.0004	-0.2069	0.7939

# การวิเคราะห์จัดกลุ่ม (Cluster analysis)

Suwanna Sayruamyat

Email: [suwanna.s@ku.th](mailto:suwanna.s@ku.th)

Facebook: Suwanna Sayruamyat

Page: [\*\*EatEcon\*\*](#),

Website: [www.eatecon.com](http://www.eatecon.com)

# ความหมายของการวิเคราะห์จัดกลุ่ม

- การจัด case: คน สัตว์ สิ่งของ หรือองค์กร ฯลฯ เป็นการจัดตัวเปรียบเป็นกลุ่มย่อย ตั้งแต่ 2 กลุ่มขึ้นไป
  - กลุ่ม (cluster) เดียวกันจะมี case คล้าย ๆ กัน
  - case ที่ต่างกันจะอยู่คนละกลุ่มกัน
- นิยมใช้จัดกลุ่มและสร้าง profile of responses

# ข้อตกลงเบื้องต้นสำหรับการวิเคราะห์จัดกลุ่ม

ไม่ทราบจำนวนมาก่อนว่ามีกี่กลุ่ม

ไม่ทราบมาก่อนว่าใครอยู่กลุ่มใด

คนหนึ่งคนอยู่ได้เพียงกลุ่มเดียว

ใช้หลายตัวแปรเป็นปัจจัยในการแบ่งกลุ่ม

# Methods

ເປັນຂຶ້ນຕອນ

## Hierarchical procedure

### Hierarchical cluster

- is the most common method. It generates a series of models with cluster solutions from 1 (all cases in one cluster) to n (each case is an individual cluster).
- Hierarchical cluster also works with variables as opposed to cases; it can cluster variables together in a manner somewhat similar to factor analysis. In addition, hierarchical cluster analysis can **handle nominal, ordinal, and scale data**; however, it is not recommended to mix different levels of measurement.

ເປັນຂຶ້ນຕອນ

## Non - hierarchical procedure

### K-means cluster

- is a method to quickly cluster large data set. **The researcher define the number of clusters in advance.** This is useful to test different models with a different assumed number of clusters

### Two-step cluster

- analysis identifies groupings by running pre-clustering first and then by running hierarchical methods. Because it uses a quick cluster algorithm upfront, it **can handle large data sets** that would take a long time to compute with hierarchical cluster methods. In this respect, it is a combination of the previous two approaches. Two-step clustering can handle scale and ordinal data in the same model, and it automatically selects the number of clusters.

# Hierarchical vs Non hierarchical methods

## Hierarchical clustering

- No decision about the number of clusters
- Problems when data contain a high level of error
- Can be very slow
- Initial decision are more influential (one-step only)

## Non hierarchical clustering

- Faster, more reliable
- Need to specify the number of clusters (arbitrary)
- Need to set the initial seeds (arbitrary)

## Suggested approach

1. First perform a hierarchical method to define the number of clusters
2. Then use the k-means procedure to form the clusters

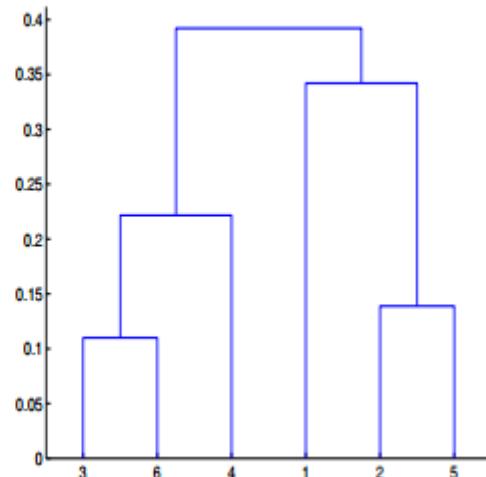
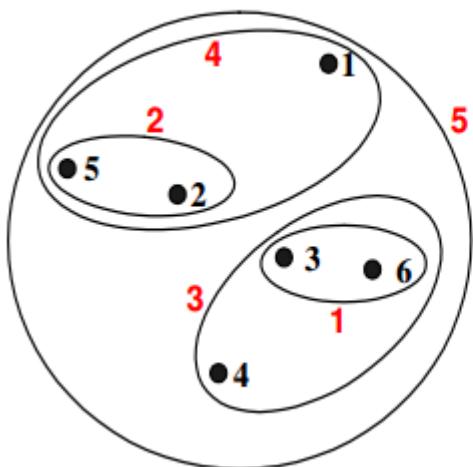
# ขั้นตอนในการวิเคราะห์การจัดกลุ่ม

1. เลือกวิธีวัดระยะห่าง (Distance measure)
2. เลือกวิธีการจัดกลุ่ม (Clustering algorithm)
3. ระบุจำนวนกลุ่ม (Determine the number of clusters)
4. ตรวจสอบความเหมาะสมของผลการวิเคราะห์ (Validate the analysis)

## Cluster analysis: basic steps

1. Apply Ward's methods on the principal components score
2. Check the agglomeration schedule
3. Decide the number of clusters
4. Apply the k-means method

# Hierarchical cluster analysis: HCA



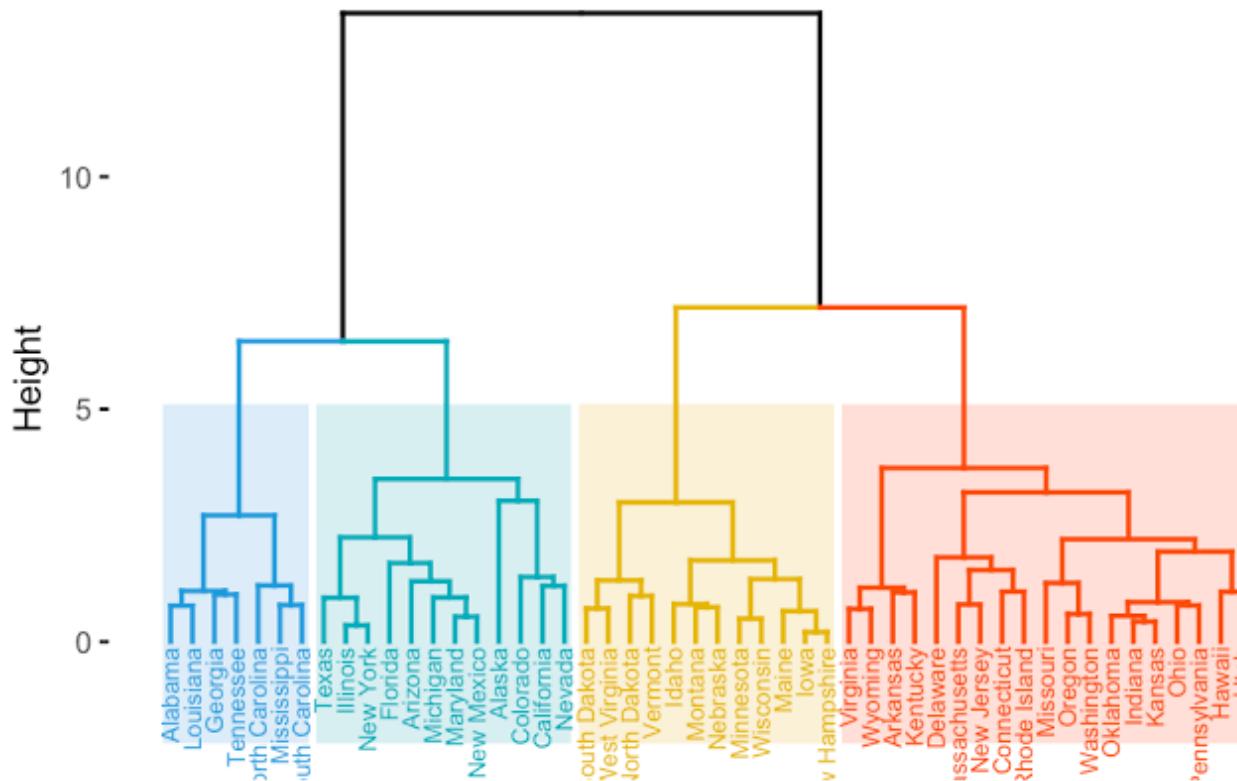
การวิเคราะห์จัดกลุ่มตามลำดับชั้น เป็นการจัดกลุ่มที่นิยมใช้แบ่งกลุ่ม Case หรือกลุ่มตัวแปร โดยมีเงื่อนไข ดังนี้

1. เหมาะกับข้อมูลขนาดเล็ก (จำนวนเคส  $< 200$  เคส/ตัวแปร)
  - ตัวอย่าง  $\Rightarrow$  แบ่งกลุ่ม (Classify cases)
  - ตัวแปร  $\Rightarrow$  ทดสอบความสัมพันธ์ระหว่างตัวแปร
2. ไม่จำเป็นต้องทราบจำนวนกลุ่มมาก่อน
3. ไม่จำเป็นต้องทราบว่าตัวแปรใดหรือเคสใดอยู่กลุ่มใดก่อน
4. ชนิดตัวแปรที่เหมาะสมคือ nominal, ordinal, and scale data และไม่ควรผสมชนิดของตัวแปร

# HCA: เกณฑ์ในการจัดกลุ่ม

แยกแต่ละตัวอย่างไปยังกลุ่ม ด้วยระยะทาง (Distance) เริ่มต้นโดยแบ่งแยก 2 กลุ่มที่เหมือนกันมากที่สุดไปเรื่อย ๆ จนครบทุกตัวอย่าง

Cluster Dendrogram

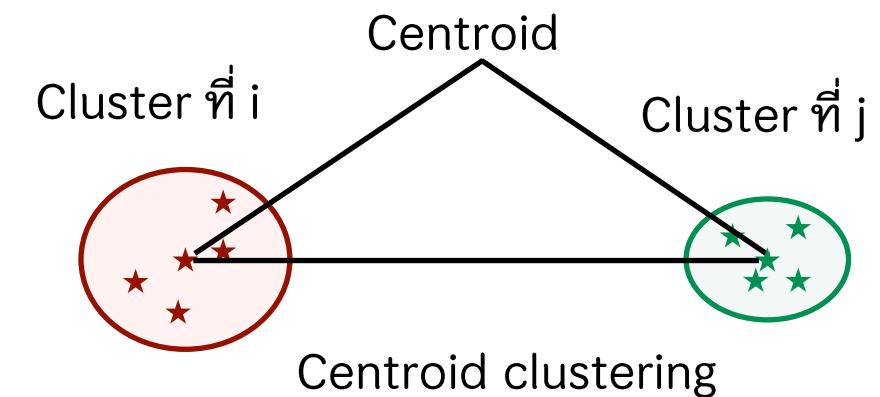
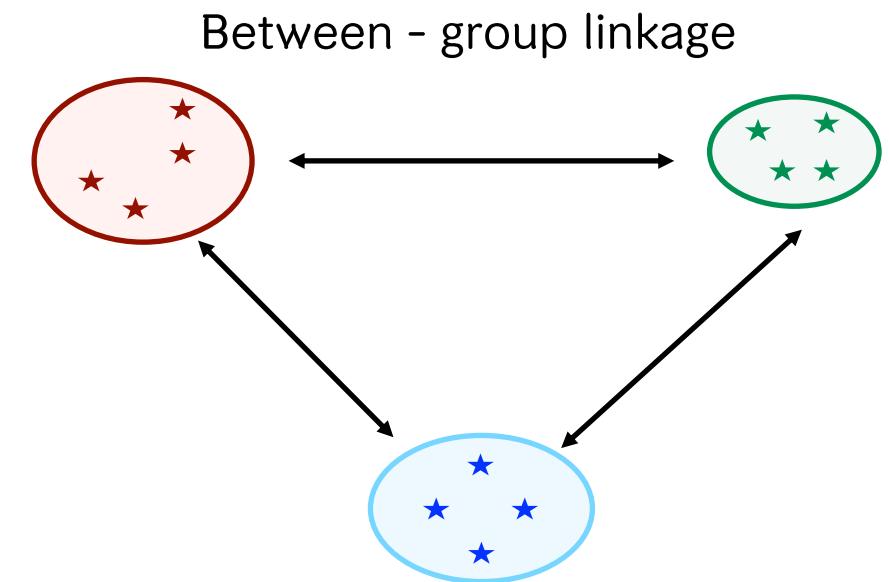


ขั้นตอนของเทคนิค Hierarchical cluster analysis สำหรับการแบ่งกลุ่มเคลส

- เลือกตัวแปรหรือปัจจัยที่คาดว่ามีอิทธิพลที่ทำให้เคลสต่างกัน ตัวแปรจะทำให้สามารถแบ่งกลุ่มเคลสได้ชัดเจน
- เลือกวิธีการวัดระยะห่างระหว่างเคลสแต่ละคู่ หรือเลือกวิธีการคำนวณเพื่อวัดค่า ความคล้ายของเคลสแต่ละคู่
- เลือกเกณฑ์ในการรวมกลุ่มหรือรวม Cluster

# HCA: เกณฑ์ในการจัดกลุ่ม

1. Between - group linkage (or average linkage between group)
2. Within-group linkage (or average linkage within groups method) - วิธีนี้จะรวม cluster เป็นลำดับกันถ้าระยะห่างเฉลี่ยระหว่างทุกค่าใน cluster นั้นๆ มีค่าน้อยที่สุด
3. Centroid clustering - รวม 2 cluster เป็นลำดับกันโดยพิจารณาจากระยะห่างของจุดกลางของ cluster 2 cluster
4. Ward's method - พิจารณาค่า sum of the squared within-cluster distance โดยจะรวม cluster ที่ทำให้ค่า sum of the squared within-cluster distance เพิ่มขึ้นน้อยที่สุด โดยค่า square within-cluster distance คือค่า square Euclidean distance ของแต่ละเดสกับ cluster mean



**Table 1: Items used to measure consumers attitudes towards shopping.**

Variables name	Attitude items	Type of data
1. NOR	Number of respondent	Scale
2. FUN	Shopping is fun	Scale
3. BUDGET	Shopping is bud for your budget	Scale
4. EATINGOUT	I combine shopping with eating out	Scale
5. BESTBUYS	I try to get the best buys when shopping	Scale
6. NO_CARE	I don't care about shopping	Scale
7. PRICE	You can save a lot of money by comparing prices	Scale
8. GENDER		Nominal
9. EDUCATION		Ordinal
10. INCOME		Scale

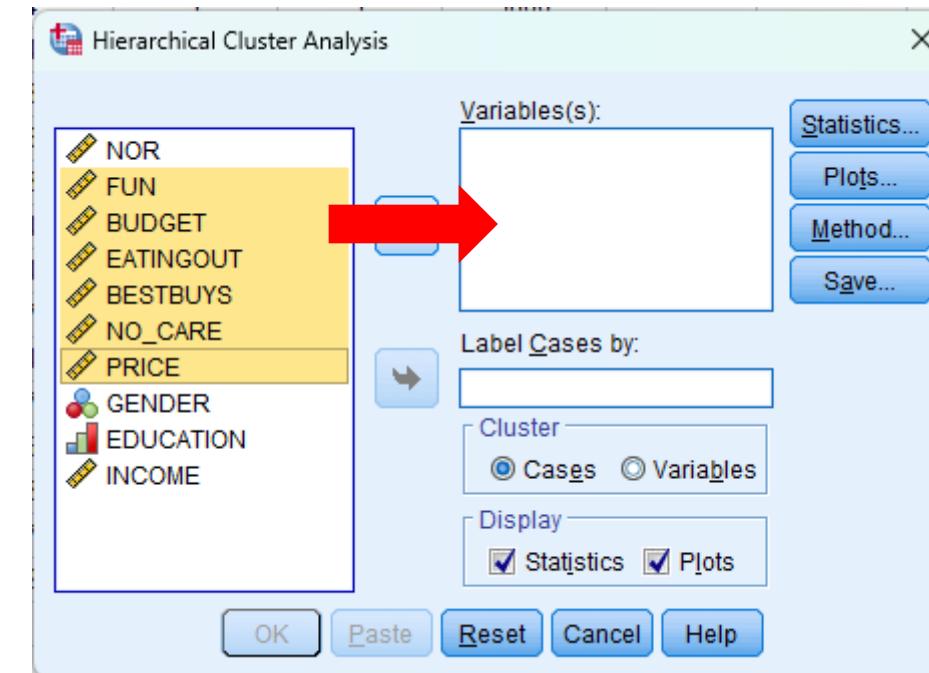
# Analyze > Classify > Hierarchical cluster...

Data for cluster analysis - Attitudes\_Shopping.sav [DataSet2] - IBM SPSS Statistics Data Editor

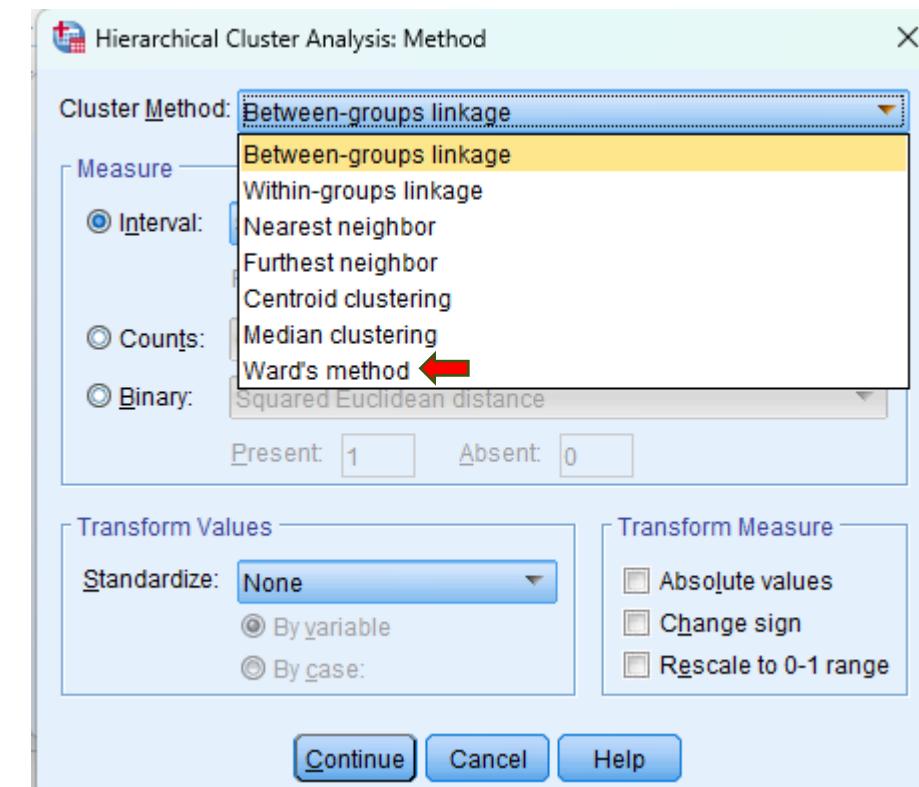
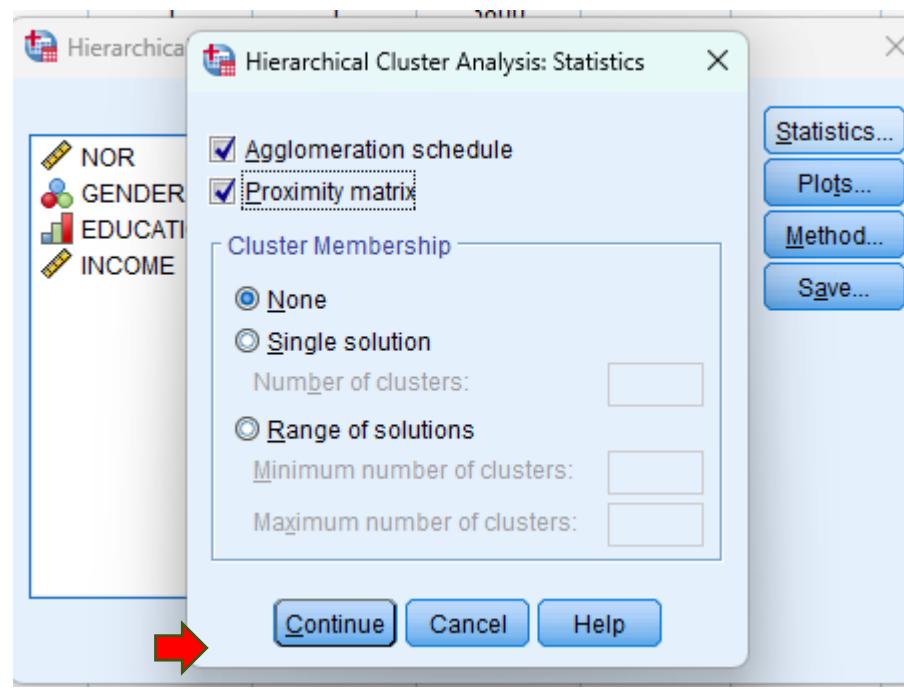
File Edit View Data Transform Analyze Graphs Utilities Extensions Window Help

	NOR	FUN
1	1	6
2	2	2
3	3	7
4	4	4
5	5	1
6	6	6
7	7	5
8	8	7
9	9	2
10	10	3
11	11	1
12	12	5
13	13	2
14	14	4
15	15	6
16	16	3
17	17	4
18	18	3
19	19	4
20	20	2
21		
22		
23		

Reports  
Descriptive Statistics  
Bayesian Statistics  
Tables  
Compare Means  
General Linear Model  
Generalized Linear Models  
Mixed Models  
Correlate  
Regression  
Loglinear  
Neural Networks  
Classify  
TwoStep Cluster...  
K-Means Cluster...  
Hierarchical Cluster...  
Cluster Silhouettes  
Tree...  
Discriminant...  
Nearest Neighbor...  
ROC Curve...  
ROC Analysis...



# HCA: Statistics



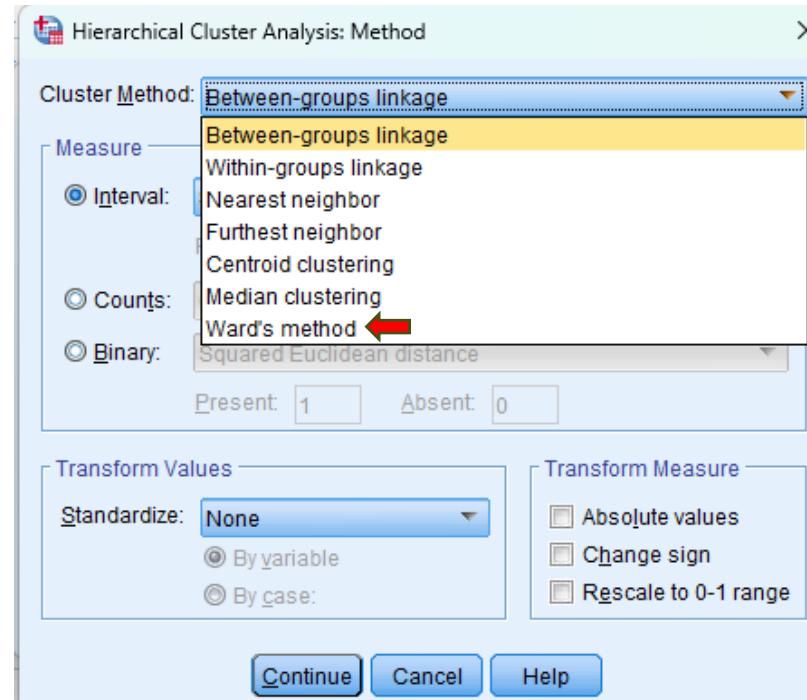
# Agglomerative clustering method

## Linkage methods

- Single linkage (minimum distance)
- Complete linkage (maximum distance)
- Average linkage

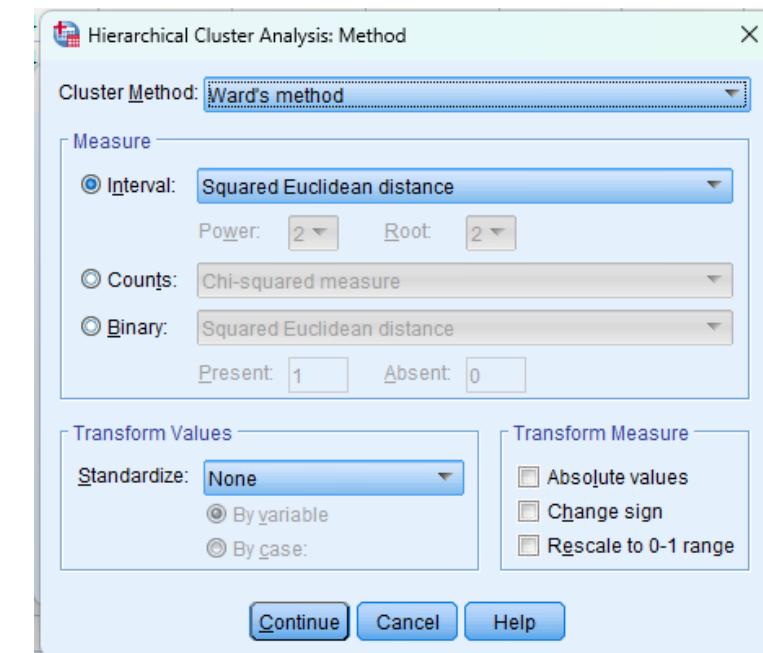
## Centroid method

- The distance between two clusters is defined as the difference between the centroids (cluster averages)



## Ward's method

1. Compute sum of squared distances within clusters
2. Aggregate clusters with the minimum increase in the overall sum of squares



# Proximity Matrix: Square Euclidean Distance

**Proximity Matrix**

Squared Euclidean Distance

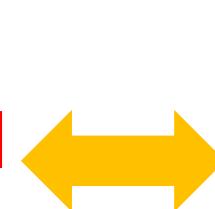
Case	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
1	.000	64.000	8.000	31.000	69.000	3.000	5.000	5.000	48.000	48.000	60.000	7.000	65.000	46.000	13.000	48.000	9.000	56.000	56.000	69.000
2	64.000	.000	68.000	31.000	7.000	47.000	39.000	77.000	8.000	18.000	4.000	35.000	3.000	36.000	49.000	28.000	55.000	24.000	44.000	9.000
3	8.000	68.000	.000	43.000	83.000	11.000	11.000	3.000	64.000	56.000	70.000	11.000	61.000	58.000	19.000	58.000	23.000	70.000	60.000	79.000
4	31.000	31.000	43.000	.000	44.000	20.000	22.000	36.000	31.000	5.000	39.000	12.000	34.000	3.000	22.000	5.000	24.000	17.000	7.000	50.000
5	69.000	7.000	83.000	44.000	.000	52.000	42.000	90.000	5.000	33.000	3.000	46.000	16.000	51.000	58.000	41.000	52.000	41.000	69.000	10.000
6	3.000	47.000	11.000	20.000	52.000	.000	2.000	8.000	35.000	33.000	47.000	4.000	50.000	31.000	10.000	33.000	8.000	45.000	43.000	54.000
7	5.000	39.000	11.000	22.000	42.000	2.000	.000	10.000	29.000	31.000	37.000	4.000	40.000	33.000	14.000	31.000	6.000	47.000	45.000	46.000
8	5.000	77.000	3.000	36.000	90.000	8.000	10.000	.000	69.000	53.000	79.000	10.000	72.000	49.000	18.000	51.000	16.000	71.000	53.000	90.000
9	48.000	8.000	64.000	31.000	5.000	35.000	29.000	69.000	.000	24.000	4.000	31.000	17.000	38.000	45.000	32.000	41.000	24.000	56.000	5.000
10	48.000	18.000	56.000	5.000	33.000	33.000	31.000	53.000	24.000	.000	28.000	21.000	19.000	4.000	39.000	2.000	39.000	14.000	8.000	35.000
11	60.000	4.000	70.000	39.000	3.000	47.000	37.000	79.000	4.000	28.000	.000	37.000	9.000	48.000	51.000	38.000	49.000	30.000	60.000	7.000
12	7.000	35.000	11.000	12.000	46.000	4.000	4.000	10.000	31.000	21.000	37.000	.000	34.000	23.000	8.000	23.000	10.000	31.000	27.000	48.000
13	65.000	3.000	61.000	34.000	16.000	50.000	40.000	72.000	17.000	19.000	9.000	34.000	.000	39.000	52.000	29.000	58.000	29.000	41.000	16.000
14	46.000	36.000	58.000	3.000	51.000	31.000	33.000	49.000	38.000	4.000	48.000	23.000	39.000	.000	39.000	2.000	37.000	22.000	6.000	55.000
15	13.000	49.000	19.000	22.000	58.000	10.000	14.000	18.000	45.000	39.000	51.000	8.000	52.000	39.000	.000	43.000	16.000	43.000	41.000	68.000
16	48.000	28.000	58.000	5.000	41.000	33.000	31.000	51.000	32.000	2.000	38.000	23.000	29.000	2.000	43.000	.000	35.000	24.000	8.000	47.000
17	9.000	55.000	23.000	24.000	52.000	8.000	6.000	16.000	41.000	39.000	49.000	10.000	58.000	37.000	16.000	35.000	.000	59.000	49.000	68.000
18	56.000	24.000	70.000	17.000	41.000	45.000	47.000	71.000	24.000	14.000	30.000	31.000	29.000	22.000	43.000	24.000	59.000	.000	24.000	31.000
19	56.000	44.000	60.000	7.000	69.000	43.000	45.000	53.000	56.000	8.000	60.000	27.000	41.000	6.000	41.000	8.000	49.000	24.000	.000	73.000
20	69.000	9.000	79.000	50.000	10.000	54.000	46.000	90.000	5.000	35.000	7.000	48.000	16.000	55.000	68.000	47.000	68.000	31.000	73.000	.000

This is a dissimilarity matrix

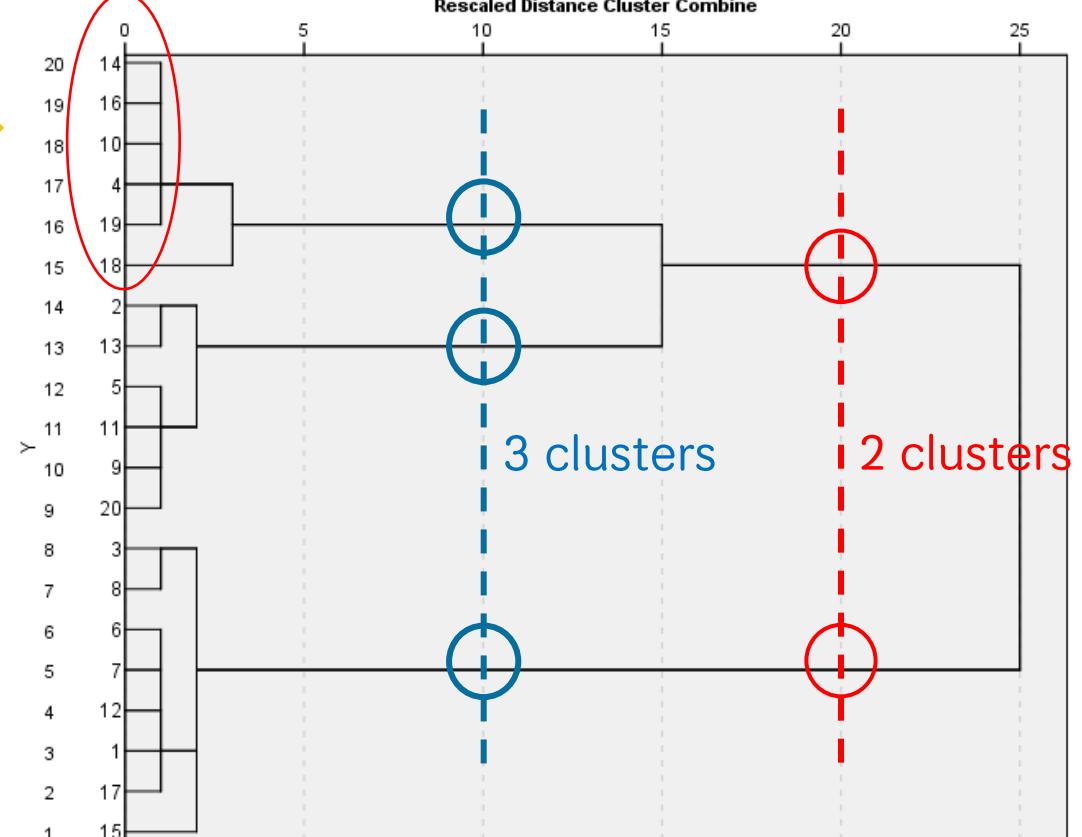
# Ward's Linkage

**Agglomeration Schedule**

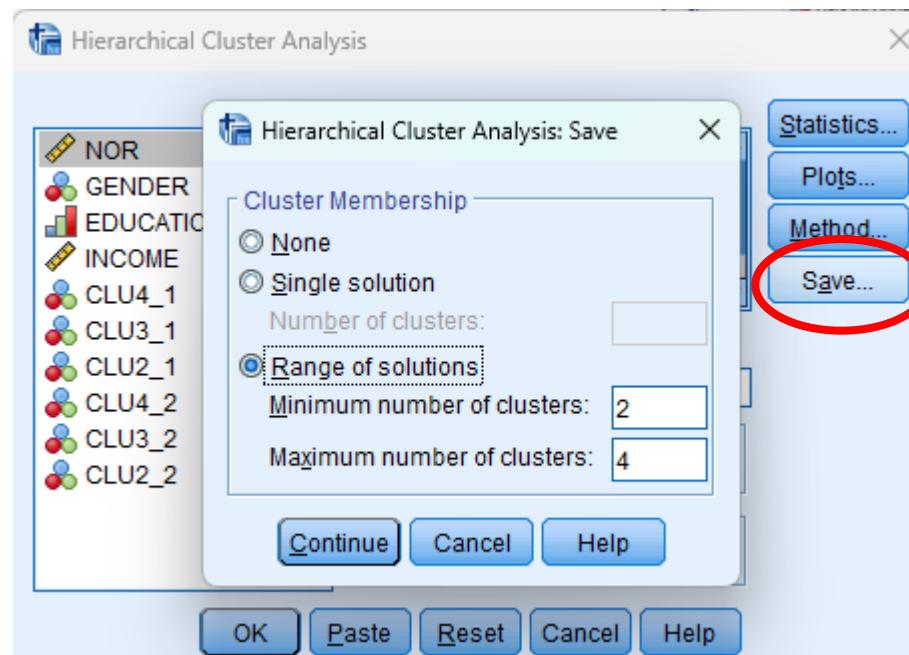
Stage	Cluster Combined		Coefficients	Stage Cluster First Appears		Next Stage
	Cluster 1	Cluster 2		Cluster 1	Cluster 2	
1	14	16	1.000	0	0	6
2	6	7	2.000	0	0	7
3	2	13	3.500	0	0	15
4	5	11	5.000	0	0	11
5	3	8	6.500	0	0	16
6	10	14	8.167	0	1	9
7	6	12	10.500	2	0	10
8	9	20	13.000	0	0	11
9	4	10	15.583	0	6	12
10	1	6	18.500	0	7	13
11	5	9	23.000	4	8	15
12	4	19	27.750	9	0	17
13	1	17	33.100	10	0	14
14	1	15	41.333	13	0	16
15	2	5	51.833	3	11	18
16	1	3	64.500	14	5	19
17	4	18	79.667	12	0	18
18	2	4	172.667	15	17	19
19	1	2	328.600	16	18	0



**Dendrogram using Ward Linkage**



# Save membership



The screenshot shows the SPSS Data Editor window with the title bar '\*Data for cluster analysis - Attitudes\_Shopping.sav [DataSet2] - IBM SPSS Statistics Data Editor'. The menu bar includes File, Edit, View, Data, Transform, Analyze, Graphs, Utilities, Extensions, Window, and Help. Below the menu is a toolbar with various icons. The main area is a data grid with 14 rows and 8 columns. The columns are labeled: Name, Type, Width, Decimals, Label, and Values. The data rows are numbered 1 to 14. The variables listed are: NOR, FUN, BUDGET, EATINGOUT, BESTBUYS, NO\_CARE, PRICE, GENDER, EDUCATION, INCOME, CLU4\_1, CLU3\_1, CLU2\_1, and an unnamed variable at the bottom. The variables CLU4\_1, CLU3\_1, and CLU2\_1 are circled in red.

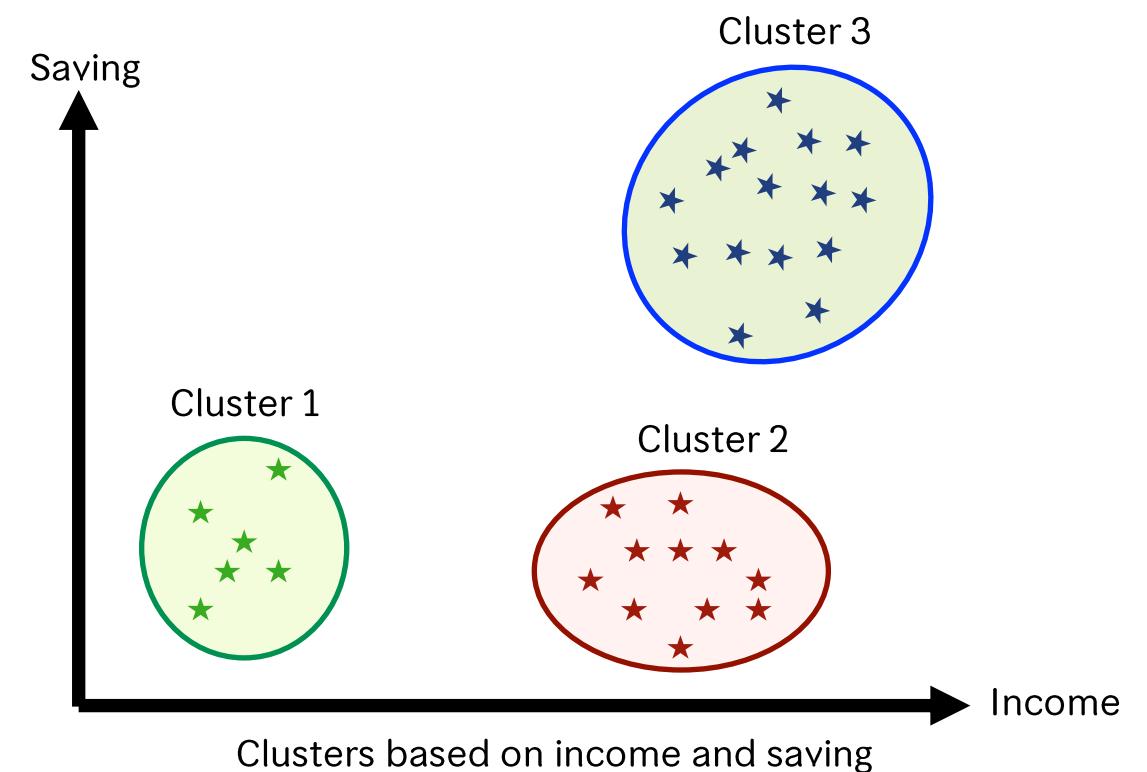
	Name	Type	Width	Decimals	Label	Values
1	NOR	Numeric	8	0	Number of resp...	None
2	FUN	Numeric	8	0	Shopping is fun	{1, Extreme...
3	BUDGET	Numeric	8	0	Shopping is bu...	{1, Extreme...
4	EATINGOUT	Numeric	8	0	I combine shop...	{1, Extreme...
5	BESTBUYS	Numeric	8	0	I try to get the ...	{1, Extreme...
6	NO_CARE	Numeric	8	0	I don't care abo...	{1, Extreme...
7	PRICE	Numeric	8	0	You can save a...	{1, Extreme...
8	GENDER	Numeric	8	0	Gender	None
9	EDUCATION	Numeric	8	0	Education	None
10	INCOME	Numeric	8	0	Net monthly inc...	None
11	CLU4_1	Numeric	8	0	Ward Method	None
12	CLU3_1	Numeric	8	0	Ward Method	None
13	CLU2_1	Numeric	8	0	Ward Method	None
14						

ตัวแปรที่จัดกลุ่มๆ กัน  
2-4 กลุ่มจะปรากฏขึ้น

# K-means cluster analysis: KCA

KCA is a method to quickly cluster large data sets. **The researcher define the number of clusters in advance.** This is useful to test different models with a different assumed number of clusters.

- กำหนดจำนวนกลุ่ม (K) ที่ต้องการก่อน และจัดตัวอย่างเข้ากลุ่ม
- ใช้กับข้อมูลขนาดใหญ่ ( $> 200$ )
- เหมาๆ กับตัวแปรค่าต่อเนื่อง (Scale)
- ใช้ค่าเฉลี่ยเพื่อจัด Case เข้ากลุ่ม K
- กลุ่มต่างกันน้อย กลุ่มต่างกันมาก
- หาระยะห่างด้วยวิธี Euclidean distance



Note: หากมีตัวแปรที่มีการแจกแจงไม่ปกติ ให้ทำการแปลงเป็นค่ามาตรฐานก่อน (Z-score)

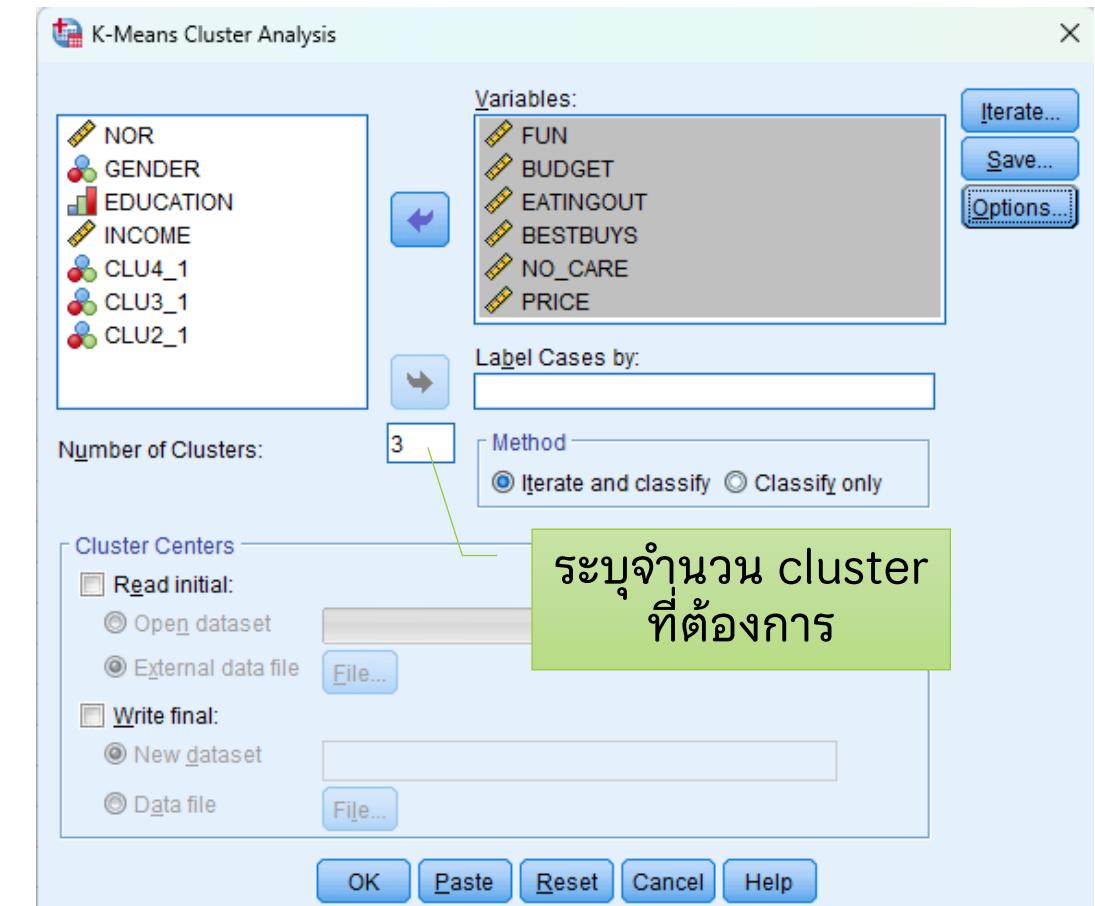
# KCA: Analyze > Classify > K-means Cluster...

\*Data for cluster analysis - Attitudes\_Shopping.sav [DataSet2] - IBM SPSS Statistics Data Editor

File Edit View Data Transform Analyze Graphs Utilities Extensions Window Help

Reports Descriptive Statistics Bayesian 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 Simulation... Quality Control Spatial and Temporal Modeling... Direct Marketing

	NOR	FUN
1	1	6
2	2	2
3	3	7
4	4	4
5	5	1
6	6	6
7	7	5
8	8	7
9	9	2
10	10	3
11	11	1
12	12	5
13	13	2
14	14	4
15	15	6
16	16	3
17	17	4
18	18	3
19	19	4
20	20	2
21		
22		
23		
24		



# KCA: Output

**Initial Cluster Centers**

	Cluster		
	1	2	3
Shopping is fun	4	2	7
Shopping is bud for your budget	6	3	2
I combine shopping with eating out	3	2	6
I try to get the best buys when shopping	7	4	4
I don't care about shopping	2	7	1
You can save a lot of money by comparing prices	7	2	3

**Distances between Final Cluster Centers**

Cluster	1	2	3
1		5.568	5.698
2	5.568		6.928
3	5.698	6.928	

**Cluster Membership**

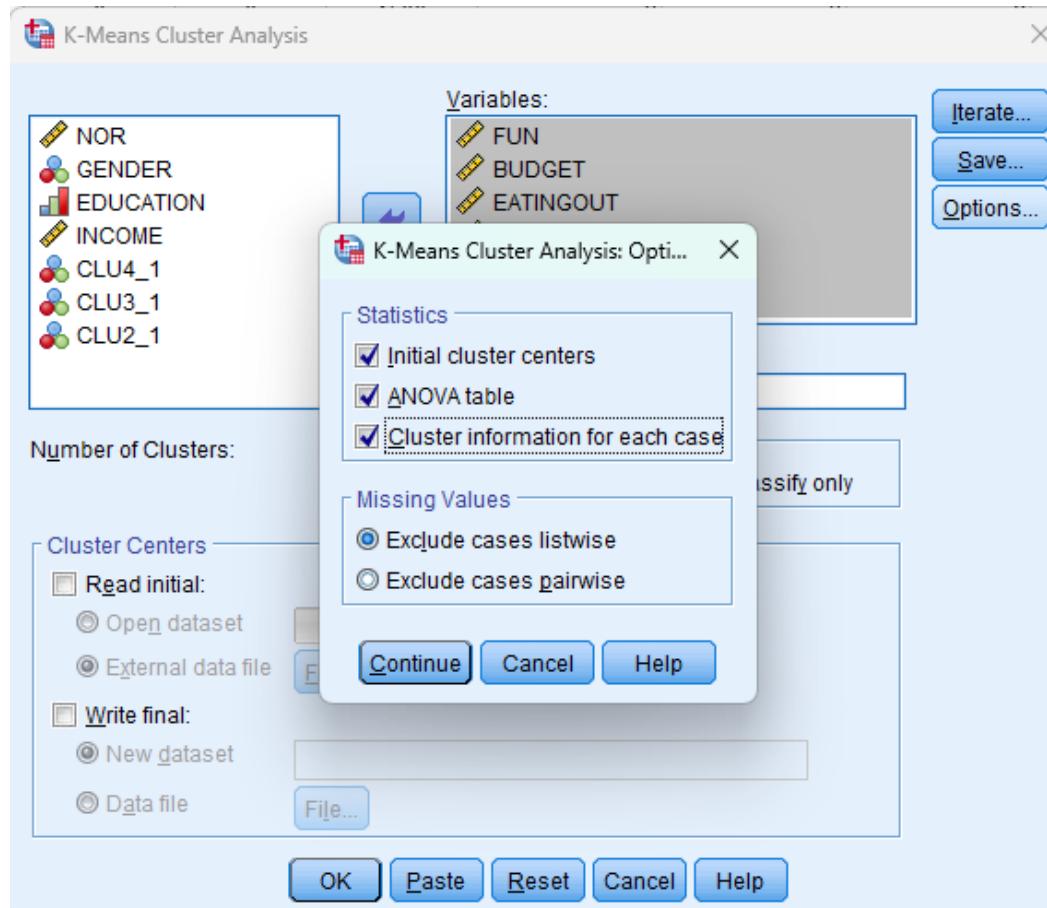
Case Number	Cluster	Distance
1	3	1.414
2	2	1.323
3	3	2.550
4	1	1.404
5	2	1.848
6	3	1.225
7	3	1.500
8	3	2.121
9	2	1.756
10	1	1.143
11	2	1.041
12	3	1.581
13	2	2.598
14	1	1.404
15	3	2.828
16	1	1.624
17	3	2.598
18	1	3.555
19	1	2.154
20	2	2.102

**Number of Cases in each Cluster**

Cluster	1	6.000
	2	6.000
	3	8.000
<b>Valid</b>		<b>20.000</b>
<b>Missing</b>		<b>.000</b>

แสดงจำนวนตัวอย่างของแต่ละกลุ่ม

# KCA: Options



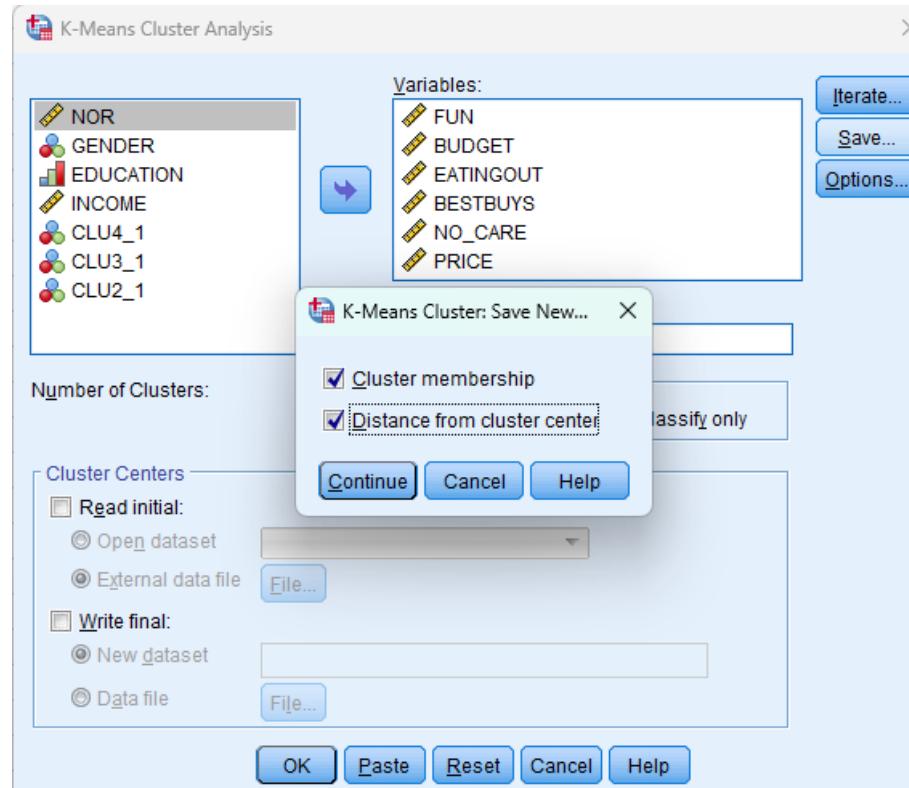
## ANOVA table

- ใช้ชี้วัดว่าตัวแปรใดมีผลมากที่สุดต่อการกำหนดกลุ่มที่กำหนดไว้
- ตัวแปรที่มีค่า F-stat สูงสุด แสดงว่าเป็นตัวแปรที่ใช้แบ่งกลุ่มได้ดีที่สุด

	Cluster		Error		F	Sig.
	Mean Square	df	Mean Square	df		
Shopping is fun	29.108	2	.608	17	47.888	.000
Shopping is bud for your budget	13.546	2	.630	17	21.505	.000
I combine shopping with eating out	31.392	2	.833	17	37.670	.000
I try to get the best buys when shopping	15.713	2	.728	17	21.585	.000
I don't care about shopping	22.537	2	.816	17	27.614	.000
You can save a lot of money by comparing prices	12.171	2	1.071	17	11.363	.001

The F tests should be used only for descriptive purposes because the clusters have been chosen to maximize the differences among cases in different clusters. The observed significance levels are not corrected for this and thus cannot be interpreted as tests of the hypothesis that the cluster means are equal.

# KCA: Save



ได้ตัวแปร 2 ตัว

- QCL1 ตัวแปรที่กำหนดว่าใครอยู่กลุ่มใคร
- QCL2 ตัวแปรที่บอกระยะห่างตัวแต่ละคนจากคลัสเตอร์

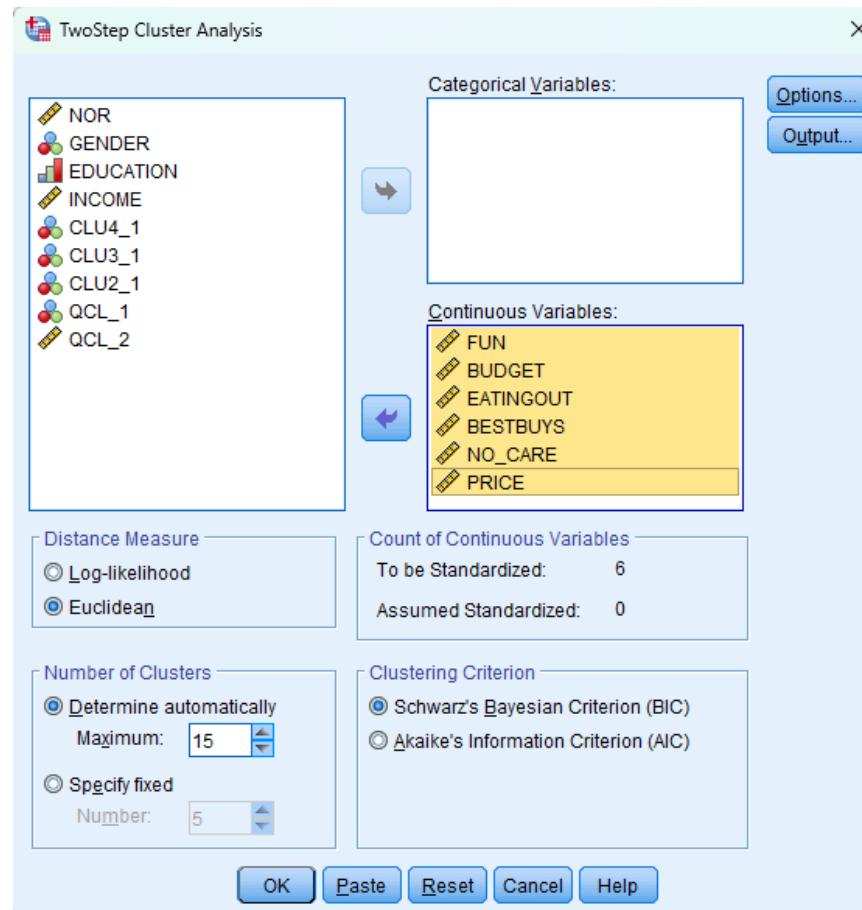
shopping.sav [DataSet2] - IBM SPSS Statistics Data Editor

	Width	Decimals	Label	Values	Missing	Columns	Align	Measure	Role	
	8	0	Number of resp...	None	None	8	Center	Scale	Input	
	8	0	Shopping is fun	{1, Extreme...}	None	8	Center	Scale	Input	
	8	0	Shopping is bu...	{1, Extreme...}	None	8	Center	Scale	Input	
	8	0	I combine shop...	{1, Extreme...}	None	8	Center	Scale	Input	
	8	0	I try to get the ...	{1, Extreme...}	None	8	Center	Scale	Input	
	8	0	I don't care abo...	{1, Extreme...}	None	8	Center	Scale	Input	
	8	0	You can save a...	{1, Extreme...}	None	8	Center	Scale	Input	
	8	0	Gender	None	None	8	Center	Nominal	Input	
	8	0	Education	None	None	8	Center	Ordinal	Input	
	8	0	Net monthly inc...	None	None	8	Center	Scale	Input	
	8	0	Ward Method	None	None	10	Right	Nominal	Input	
	8	0	Ward Method	None	None	10	Right	Nominal	Input	
	8	0	Ward Method	None	None	10	Right	Nominal	Input	
	14	QCL_1	Numeric	8	0	Cluster Number...	None	None	10	Right
	15	QCL_2	Numeric	20	5	Distance of Ca...	None	None	22	Right

# Two-step cluster analysis

- analysis identifies groupings by running pre-clustering first and then by running hierarchical methods. Because it uses a quick cluster algorithm upfront, **it can handle large data sets** that would take a long time to compute with hierarchical cluster methods. In this respect, it is a combination of the previous two approaches. **Two-step clustering can handle scale and ordinal data in the same model**, and it automatically selects the number of clusters.

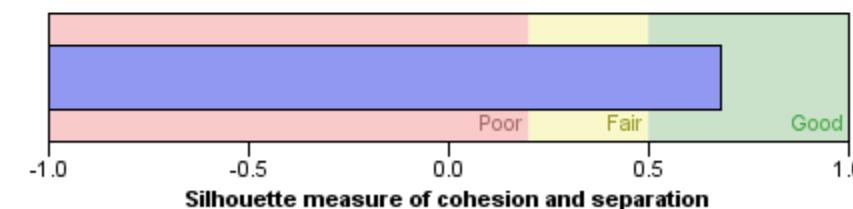
# Two-step output



## Model Summary

<b>Algorithm</b>	TwoStep
<b>Inputs</b>	6
<b>Clusters</b>	3

## Cluster Quality



# Create cluster membership variable

Data view  
↓

The screenshot shows the SPSS TwoStep Cluster Analysis dialog box and its associated Data View and Variable View.

**TwoStep Cluster Analysis Dialog:**

- Output:**  Pivot tables,  Charts and tables in Model Viewer.
- Variables:** NOR, GENDER, EDUCATION, INCOME, CLU4\_1, CLU3\_1.
- Evaluation Fields:** (Empty)
- Working Data File:**  Create cluster membership variable.
- XML Files:** Export final model (Name: [empty]), Export CF tree (Name: [empty]).
- Buttons:** Continue, Cancel, Help, OK, Paste, Reset, Cancel, Help.

A red arrow points to the "Create cluster membership variable" checkbox in the Working Data File section.

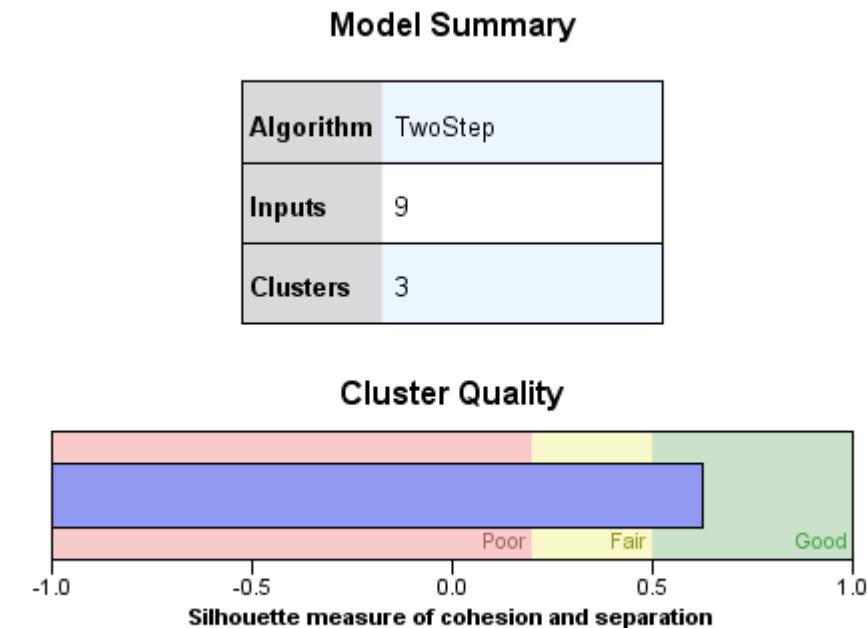
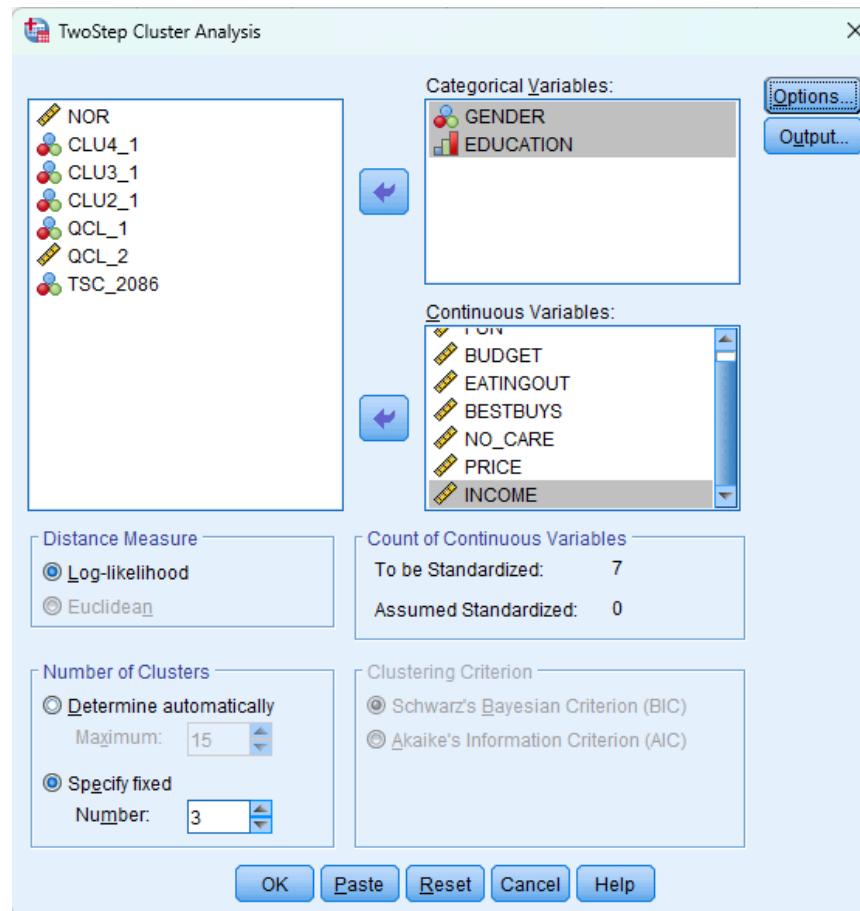
**Data View:** A table showing cluster assignments (CLU4\_1, CLU3\_1, CLU2\_1, QCL\_1, QCL\_2) and a cluster number (TSC\_2086) for 16 cases. A red arrow points to the TSC\_2086 column.

	CLU4_1	CLU3_1	CLU2_1	QCL_1	QCL_2	TSC_2086
1	1	1	1	3	1.41421	3
2	2	2	2	2	1.32288	1
3	1	1	1	3	2.54951	3
4	3	3	2	1	1.40436	2
5	2	2	2	2	1.84842	1
6	1	1	1	3	1.22474	3
7	1	1	1	3	1.50000	3
8	1	1	1	3	2.12132	3
9	2	2	2	2	1.75594	1
10	3	3	2	1	1.14261	2
11	2	2	2	2	1.04083	1

**Variable View:** A table showing variable details for 16 variables. A red arrow points to the last row (TSC\_2086).

Variable	Name	Type	Width	Decimals	Label	Values	Missing	Columns	Align	Measure	Role
1	NOR	Numeric	8	0	Number of respondent	None	None	8	Center	Scale	Input
2	FUN	Numeric	8	0	Shopping is fun	{1, Extremel...	None	8	Center	Scale	Input
3	BUDGET	Numeric	8	0	Shopping is bud for your budget	{1, Extremel...	None	8	Center	Scale	Input
4	EATINGOUT	Numeric	8	0	I combine shopping with eating out	{1, Extremel...	None	8	Center	Scale	Input
5	BESTBUYS	Numeric	8	0	I try to get the best buys when shopping	{1, Extremel...	None	8	Center	Scale	Input
6	NO_CARE	Numeric	8	0	I don't care about shopping	{1, Extremel...	None	8	Center	Scale	Input
7	PRICE	Numeric	8	0	You can save a lot of money by comparing prices	{1, Extremel...	None	8	Center	Scale	Input
8	GENDER	Numeric	8	0	Gender	None	None	8	Center	Nominal	Input
9	EDUCATION	Numeric	8	0	Education	None	None	8	Center	Ordinal	Input
10	INCOME	Numeric	8	0	Net monthly income	None	None	8	Center	Scale	Input
11	CLU4_1	Numeric	8	0	Ward Method	None	None	10	Right	Nominal	Input
12	CLU3_1	Numeric	8	0	Ward Method	None	None	10	Right	Nominal	Input
13	CLU2_1	Numeric	8	0	Ward Method	None	None	10	Right	Nominal	Input
14	QCL_1	Numeric	8	0	Cluster Number of Case	None	None	10	Right	Nominal	Input
15	QCL_2	Numeric	20	5	Distance of Case from its Classification Cluster Center	None	None	10	Right	Scale	Input
16	TSC_2086	Numeric	10	0	TwoStep Cluster Number	{-1, Outlier ...	None	8	Right	Nominal	Input

# Two-step output เมื่อใส่ตัวแปรรวมกัน



# การจัดกลุ่มจะมีความแตกต่างกันเมื่อใส่ตัวแปรต่างกัน

	PRICE	GENDER	EDUCATION	INCOME	CLU4_1	CLU3_1	CLU2_1	QCL_1	QCL_2	TSC_2086	TSC_2705
1	3	1	1	3000	1	1	1	3	1.41421	3	1
2	4	0	0	2000	2	2	2	2	1.32288	1	2
3	3	1	1	3500	1	1	1	3	2.54951	3	1
4	6	0	0	1500	3	3	2	1	1.40436	2	3
5	4	1	0	2300	2	2	2	2	1.84842	1	2
6	4	1	1	4000	1	1	1	3	1.22474	3	1
7	4	1	1	3800	1	1	1	3	1.50000	3	1
8	4	1	0	4500	1	1	1	3	2.12132	3	1
9	3	1	0	2600	2	2	2	2	1.75594	1	2
10	6	0	0	1600	3	3	2	1	1.14261	2	3
11	3	0	0	2200	2	2	2	2	1.04083	1	2
12	4	1	1	3600	1	1	1	3	1.58114	3	1
13	4	0	0	2400	2	2	2	2	2.59808	1	2
14	7	0	1	1800	3	3	2	1	1.40436	2	3
15	4	1	1	5000	1	1	1	3	2.82843	3	1
16	7	0	0	1650	3	3	2	1	1.62447	2	3
17	5	0	0	2100	1	1	1	3	2.59808	3	3
18	3	0	0	1400	4	3	2	1	3.55512	2	3
19	7	0	0	1600	3	3	2	1	2.15381	2	3
20	2	1	0	2500	2	2	2	2	2.10159	1	2

# Cluster analysis via STATA

Suwanna Sayruamyat

# การเตรียมข้อมูลและการติดตั้งคำสั่ง

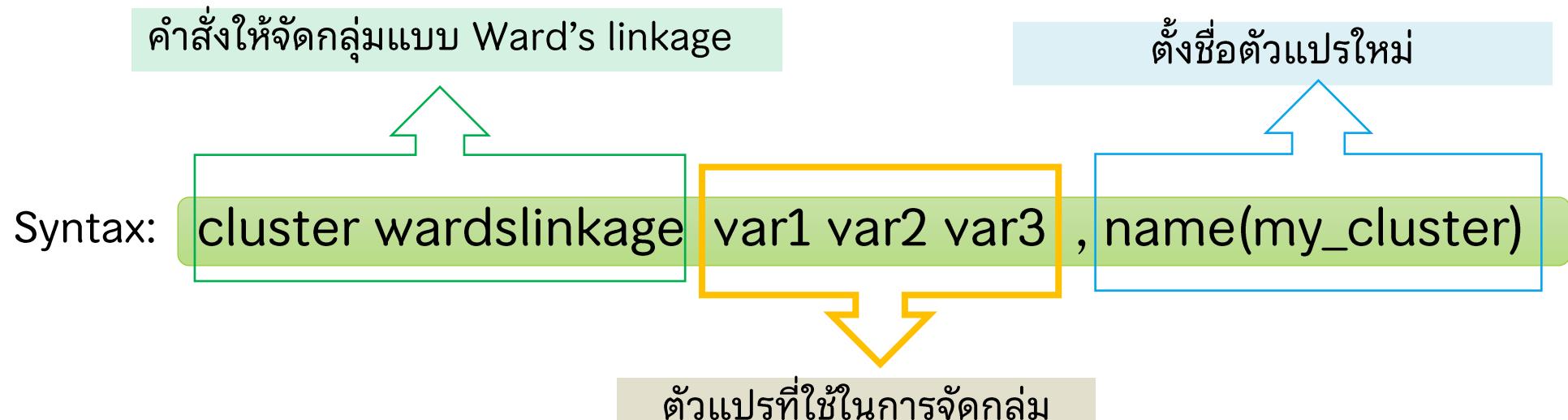
## 1. นำเข้าข้อมูล:

- การนำเข้าข้อมูลตัวอย่าง (\*.dta, \*.csv)

## 2. ทำความสะอาดข้อมูล:

- การจัดการข้อมูลสูญหาย (Missing values) และค่าผิดปกติ (Outliers)

# การวิเคราะห์แบบ Hierarchical Clustering



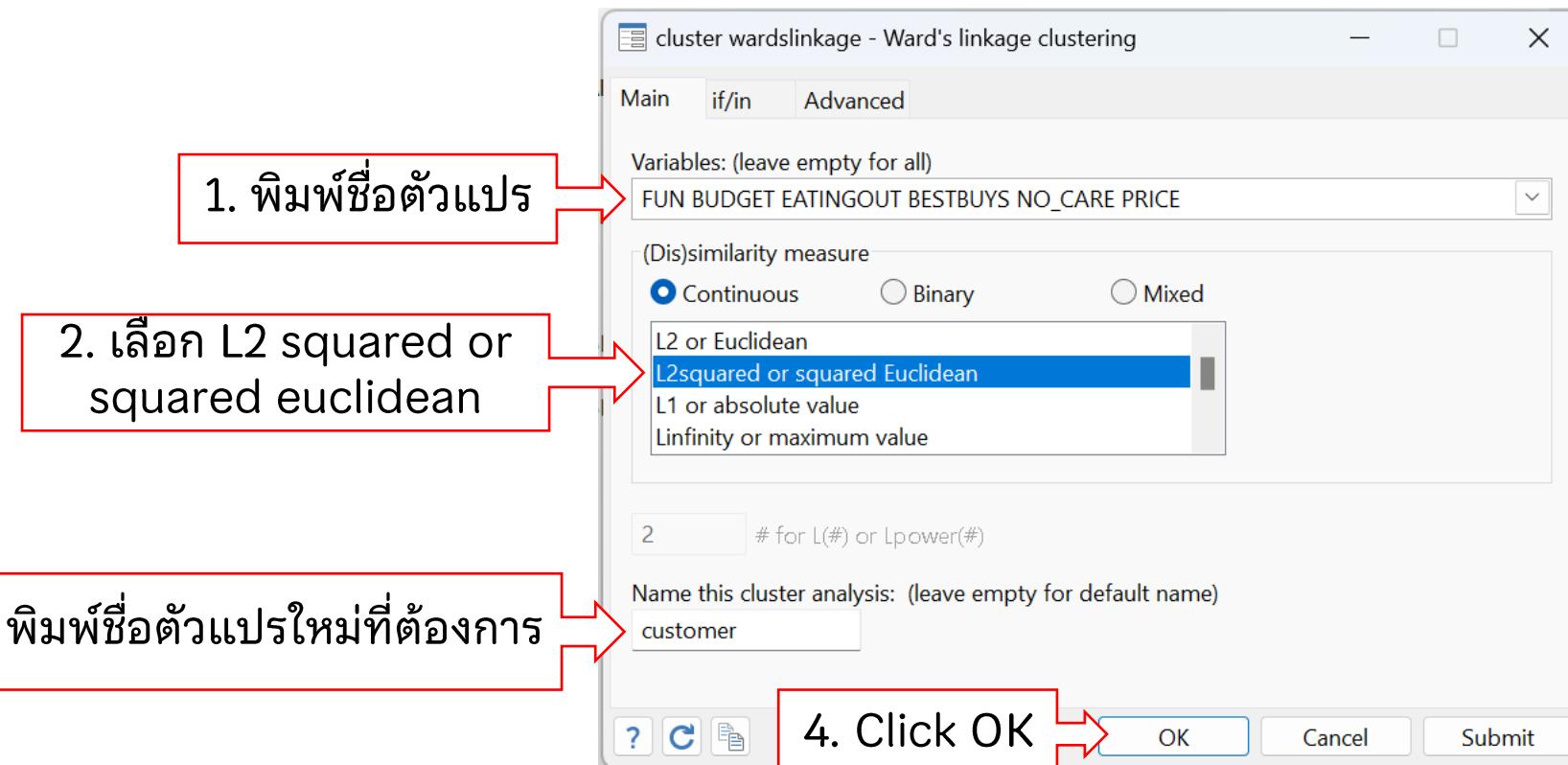
## Linkage options in Stata

<i>linkage</i>	Description
<u>singlalinkage</u>	single-linkage cluster analysis
<u>averagelinkage</u>	average-linkage cluster analysis
<u>completelinkage</u>	complete-linkage cluster analysis
<u>waveragelinkage</u>	weighted-average linkage cluster analysis
<u>medianlinkage</u>	median-linkage cluster analysis
<u>centroidlinkage</u>	centroid-linkage cluster analysis
<u>wardslinkage</u>	Ward's linkage cluster analysis

# การวิเคราะห์แบบ Hierarchical Clustering using Menu

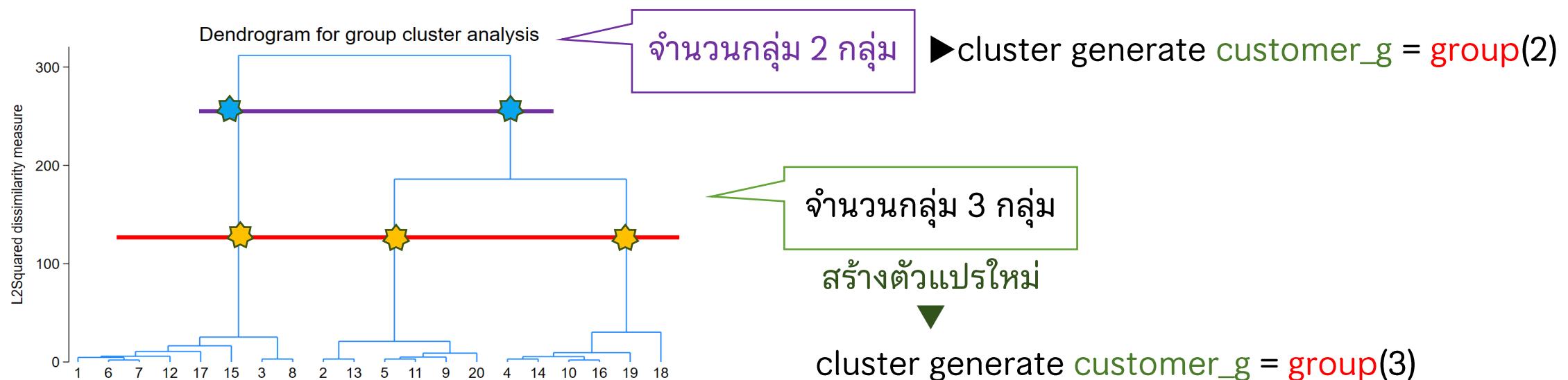
## Using Menu for cluster wardslinkage

Statistics > Multivariate analysis > Cluster analysis > Cluster data > Ward's linkage

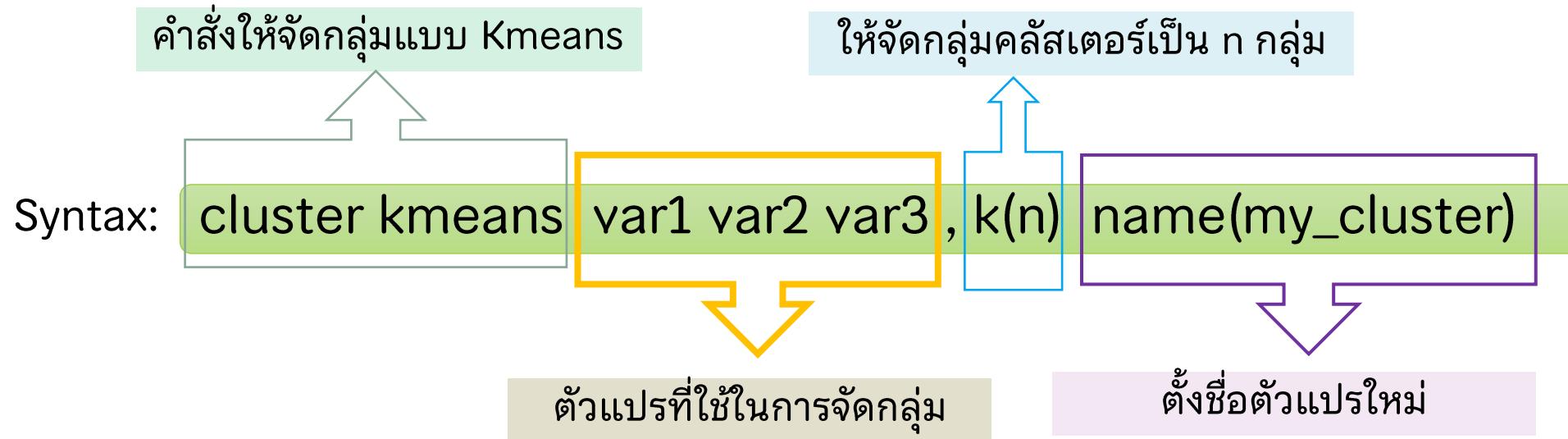


# การวิเคราะห์แบบ Hierarchical Clustering

- คำสั่ง:
- ▶ cluster wardslinkage FUN BUDGET EATINGOUT BESTBUYS NO\_CARE PRICE, name(**group**)
  - ▶ cluster dendrogram
  - ▶ cluster generate **customer\_g** = group(3) → สร้างตัวแปรใหม่ ชื่อ customer\_g จากตัวแปร group จำนวน 3 กลุ่ม



# การวิเคราะห์แบบ K-means Clustering

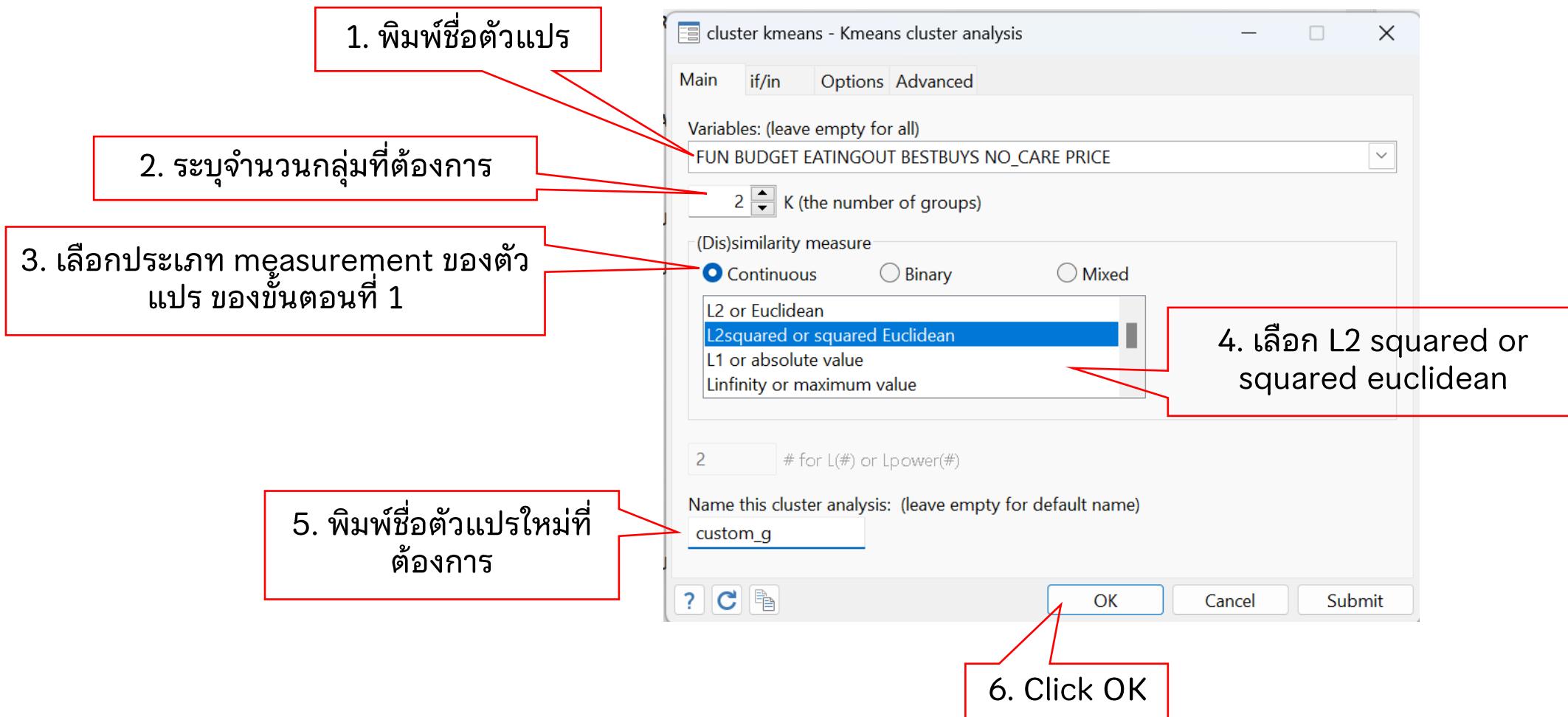


## Using Menu

**cluster kmeans**

Statistics > Multivariate analysis > Cluster analysis > Cluster data > Kmeans

# การวิเคราะห์แบบ K-means Clustering using Menu



# การเปรียบเทียบกลุ่ม: ใช้คำสั่ง tabstat

- ▶ ใช้คำสั่ง tabstat หรือ oneway เพื่อเปรียบเทียบลักษณะของแต่ละกลุ่ม

tab FUN custom\_g

FUN	Cluster ID		Total
	1	2	
1	0	2	2
2	2	2	4
3	3	0	3
4	3	1	4
5	0	2	2
6	0	3	3
7	0	2	2
Total	8	12	20

# การเปรียบเทียบกลุ่ม: One-way ANOVA

`oneway FUN customer_group, tabulate`

customer_group	Summary of FUN		
	Mean	Std. dev.	Freq.
1	5.75	1.0350983	8
2	1.6666667	.51639778	6
3	3.5	.54772256	6
Total	3.85	1.8994459	20

Source	Analysis of variance				
	SS	df	MS	F	Prob > F
Between groups	58.2166667	2	29.1083333	47.89	0.0000
Within groups	10.3333333	17	.607843137		
Total	68.55	19	3.60789474		

Bartlett's equal-variances test:  $\chi^2(2) = 3.4075$  Prob> $\chi^2 = 0.182$

# Post hoc testing of One-way ANOVA

ตัวแปรตาม  
(dependent variable)

ตัวแปรต้น  
(independent variable)

Pairwise comparisons of means with equal variances

`pwmean FUN, over(customer_group) mcompare(tukey) effects`

FUN	Contrast	Std. err.	Tukey		Tukey	
			t	P> t	[95% conf. interval]	
customer_group	2 vs 1	-4.083333	.4210553	-9.70	0.000	-5.163491 -3.003176
	3 vs 1	-2.25	.4210553	-5.34	0.000	-3.330157 -1.169843
	3 vs 2	1.833333	.4501271	4.07	0.002	.6785967 2.98807

เปรียบเทียบตัวแปร FUN ระหว่าง  
กลุ่ม 2 vs กลุ่ม 1  
กลุ่ม 3 vs กลุ่ม 1  
กลุ่ม 3 vs กลุ่ม 1

ค่า P-value แสดง ว่าแต่ละคู่มีความแตกต่าง  
กันอย่างมีนัยสำคัญทางสถิติที่ 0.01

# Assignment

- Using file: Data for practice (Download: <https://www.eatecon.com/courses/business-plan-workshop/>)
  - ใช้ข้อมูลที่กำหนดให้วิเคราะห์ Factor analysis โดยเลือกตัวแปรในการวิเคราะห์ไม่น้อยกว่า 12 ตัวแปร
  - ใช้ข้อมูลที่กำหนดให้วิเคราะห์ Cluster analysis โดยเลือกตัวแปรในการวิเคราะห์จำนวนอย่างน้อย 6 ตัวแปร
  - อธิบายผลการวิเคราะห์ให้เข้าใจโดยจะวิเคราะห์ด้วยโปรแกรม SPSS หรือ Stata ก็ได้
  - เขียนด้วยลายมือ โดยสามารถเขียนใน iPad ได้ โดยรวมเป็นไฟล์เดียว และ save in pdf file ขนาดไฟล์ไม่เกิน 100MB
  - Submission Link: <https://forms.gle/ndh1fNjiZKu4tLH66>
  - กำหนดส่ง 30 สิงหาคม 2568 เวลา 20.00 น.