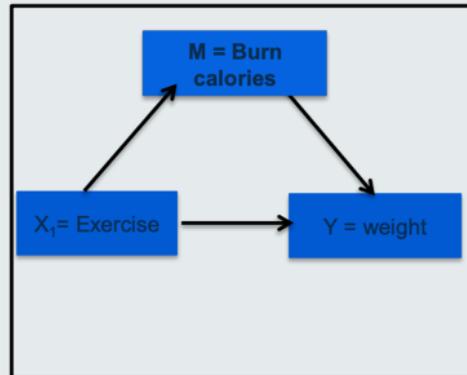
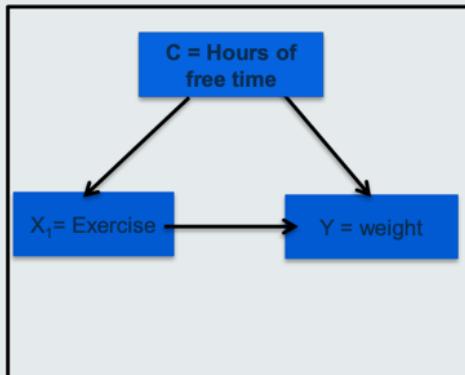


Core Concepts

★ 1. What is Effect Modification (Interaction)?

- **Definition:** When the effect of one variable on an outcome depends on the level of another variable.
- This means the relationship between a predictor (X) and outcome (Y) changes depending on a third variable (Z).
- **Types:**
 - *Categorical* effect modifier (e.g., gender)
 - *Continuous* effect modifier (e.g., age)

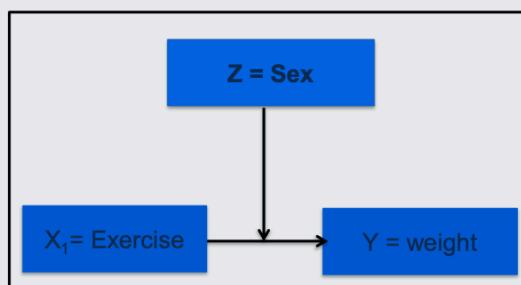


A **confounder** (**C**) has a common effect on the independent and dependent variables. A confounder is **extrinsic to the causal pathway**.

A **mediator** (**M**) is caused by the independent variable which in turn causes the dependent variable. A mediator is **in the causal pathway**

Effect Modification (Interaction)

- The third variable X_2 can have another role.
- X_2 can be a **modifier** (or moderator or have an interaction effect) on the association between Y and X_1 :



We will denote the moderator with letter **Z**

Key point: The association between X_1 and Y is not the same at different values or levels of Z . In other words, a **modifier** is a variable that **alters** the relationship between the independent X_1 and dependent Y variables.

Example: If a man and a woman do the same exercise, the effect on weight is different. **Sex modifies the effect of Exercise on Weight.**

Establishing Effect Modification

To assess the **significance** of an effect modification or interaction:

Step 1:

- A new variable needs to be considered
 - This new term is the **cross-product** between X_1 and the modifier Z. This is called the **Interaction Term**.
 - The new term is noted like $X_1 \times Z$

Step 2:

Consider the original linear regression to test the effect between Y , X_1 and Z
 Add the new variable $X_1 \times Z$ to the regression model: $Y = \beta_0 + \beta_1 X_1 + \beta_2 Z + \beta_3 X_1 \times Z + \epsilon$

In the context of regression analysis, assessing effect modification is the same as assessing interaction effect.

Step 3:

Step 3: Primary focus: Test coefficient β_3 ; $\begin{cases} H_0: \beta_3 = 0 \\ H_1: \beta_3 \neq 0 \end{cases}$ If p value < 0.05 then there is a **significant effect modification**.

If p value < 0.05 we will conclude there is a significant interaction between X_1 and the modifier Z .

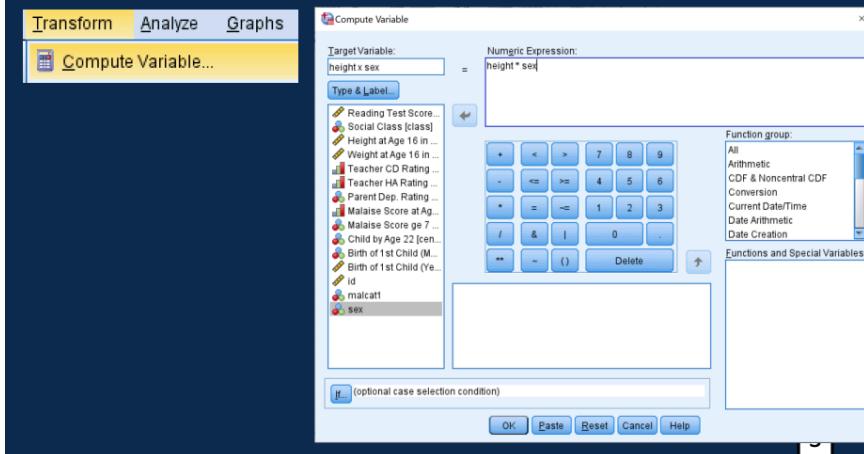
SPSS Slide: 'How to' Steps

Create an interaction term `height_x_sex` from `lecture_9a_data.sav`

1) Use ‘Transform’ -> ‘Compute variable’

2) In “**Target variable**” write the name of your interaction term: “**height_x_sex**”

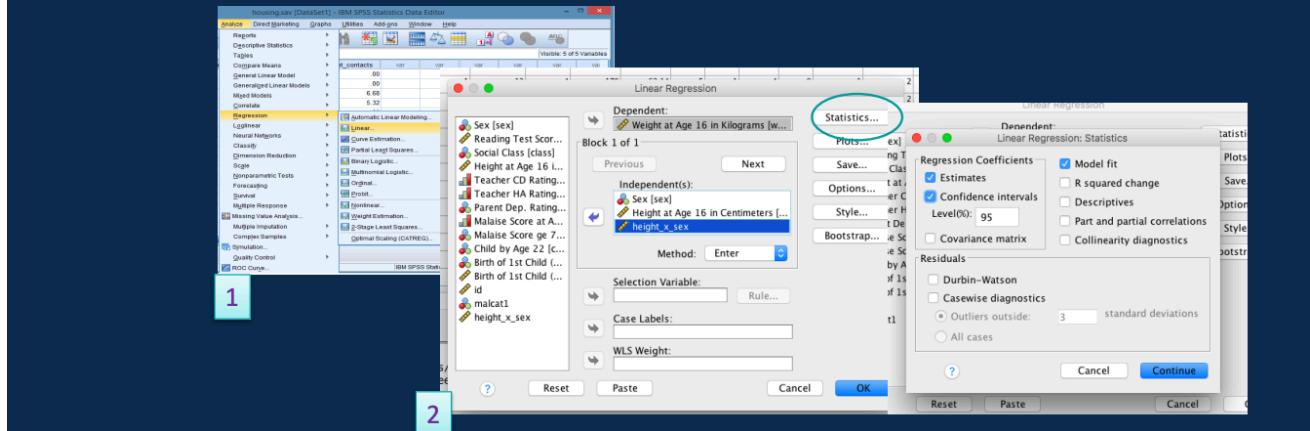
In “Numeric Expression” drag ‘height’ and ‘sex’ separated by a ‘*’



New variable in data set

Estimating the interaction effect height_x_sex in a multiple linear regression model for weight, height and sex from lecture_9a_data.sav data

- 1) Use 'Analyse' -> 'Regression' -> 'Linear'
- 2) In dependent put 'weight' and in independent put 'height', 'sex', 'height_x_sex'



Output and Interpretation

$$weight = b_0 + b_1 height + b_2 sex + b_3 height \cdot sex$$

Coefficients ^a						
Model	Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.	95.0% Confidence Interval for B Lower Bound
	B	Std. Error				
1	(Constant)	-72.014	8.893	-8.098	.000	-89.464
	Height at Age 16 in Centimeters	.771	.052	.640	14.804	.669
	Gender	38.263	12.963	1.982	2.952	.003
	hxs	-.223	.078	-1.870	-2.851	.004

a. Dependent Variable: Weight at Age 16 in Kilograms

- $\beta_1 = 0.77$ is interpreted as the effect of height on weight when sex=0 (boys)
- $\beta_2 = 38.26$ represents the effect of sex on weight when height=0 (not meaningful! because a person's height can not be zero)
- $\beta_3 = -0.22$ represents the difference of the effect of height on weight between girls (sex=1) and boys (sex=0)
- The p-value of the Interaction effect ($\beta_3 = -0.22$) is 0.004, we conclude that height \times sex interaction effect is statistically significant
- The height-weight relationship significantly differs between boys and girls

Interpretation of β 's

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 Z + \beta_3 x_1 \times Z + \epsilon$$

- β_1 is interpreted as the effect of x_1 on Y when $Z = 0$
- Similarly, β_2 represents the effect of Z on Y when $x_1 = 0$
- β_1 and β_2 are no longer useful unless zero values of the respective predictors are of particular interest.
- Both β_1 and β_2 are called **main effects**
- β_3 is interpreted as the difference of the effect of x_1 on Y by levels of Z variable.
- If the hypothesis test for β_3 concludes that it significantly differs from 0, that will imply:
 - Both x_1 and Z are associated with Y
 - The effect of x_1 on Y will depend on Z and vice versa
 - The $x_1 \times Z$ interaction effect is interpreted as the difference of the effect of x_1 between different levels of Z

✍ Interpreting Interaction

- If the interaction term is **statistically significant ($p < 0.05$)**, there is **effect modification**.
- Examine plots and coefficients to understand how the effect changes across groups.

▲ Dummy Variables

- For a categorical variable with k categories, create $k-1$ dummy variables.
- SPSS will do this automatically if the variable is defined as categorical in the regression dialog.

Example:

If you have a 3-category variable (A, B, C), SPSS will generate:

- Dummy1 = 1 if B, 0 otherwise
- Dummy2 = 1 if C, 0 otherwise
(A is reference group)

Dummy Variables

Example: $Y = \text{Income}$; $X_1 = \text{job}$; $Z = \text{born city} ; (\text{London, Manchester, Leicester})$.

Z is converted into two binary dummy variables:

$$\begin{aligned}d_{\text{London}} &= 1, 0, 0 \\d_{\text{Manchester}} &= 0, 1, 0\end{aligned}$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 d_{\text{London}} + \beta_3 d_{\text{Manchester}} + \beta_4 x_1 \times d_{\text{London}} + \beta_5 x_1 \times d_{\text{Manchester}} + \varepsilon$$

Test coefficients β_4 ; $\begin{cases} H_0: \beta_4 = 0 \\ H_1: \beta_4 \neq 0 \end{cases}$ and β_5 ; $\begin{cases} H_0: \beta_5 = 0 \\ H_1: \beta_5 \neq 0 \end{cases}$

SPSS Slide

Download the data that we are going to use during the lecture. The dataset is the **lecture_9b_data.sav**.

	state	vio_crime_rate	murder_rate	metropolitan	white	education	poverty	single	metro_cat
1	AK	761	9.00	41.80	75.20	86.60	9.10	14.30	.00
2	AL	780	11.60	67.40	73.50	66.90	17.40	11.50	1.00
3	AR	593	10.20	44.70	82.90	66.30	20.00	10.70	.00
4	AZ	715	8.60	84.70	88.60	78.70	15.40	12.10	1.00
5	CA	1078	13.10	96.70	79.30	76.20	18.20	12.50	1.00
6	CO	567	5.80	81.80	92.50	84.40	9.90	12.10	1.00
7	CT	456	6.30	95.70	89.00	79.20	8.50	10.10	1.00

The dataset contains data from 51 US states, measuring the crime rates and background measures for each State with respect to their

- **violent crime:** per 100,000 population
- **murder :** per 100,000 population
- **poverty:** percent below the poverty line
- **single:** percentage of lone parents
- **urban:** level of urbanicity

Dummy Variables

Example: Y = crime rate; X_1 = poverty; Z = Urban; (Low, Medium, High)

Only 2 dummy variables (e.g. d_{Low} and d_{Medium}) are needed to represent a variable with 3 levels.

Z is converted into two binary dummy variables:

$$\begin{aligned} d_{\text{Low}} &= 1, 0, 0 \\ d_{\text{Medium}} &= 0, 1, 0 \end{aligned}$$

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 d_{\text{Low}} + \beta_3 d_{\text{Medium}} + \beta_4 x_1 \times d_{\text{Low}} + \beta_5 x_1 \times d_{\text{Medium}} + \varepsilon$$

Test coefficients β_4 ; $\begin{cases} H_0: \beta_4 = 0 \\ H_1: \beta_4 \neq 0 \end{cases}$ and β_5 ; $\begin{cases} H_0: \beta_5 = 0 \\ H_1: \beta_5 \neq 0 \end{cases}$

US crime data. The variable urban is a categorical variable with three levels "Low", "Medium" and "High"

state	urban
AK	Low
AR	Low
IA	Low
ID	Low
KY	Low
ME	Low
AL	Medium
GA	Medium
KS	Medium
MN	Medium
MO	Medium
NC	Medium
AZ	High
CA	High
CO	High
CT	High
DE	High

The variable
urban is a
categorical
variable with
three levels
"Low", "Medium"
and "High"

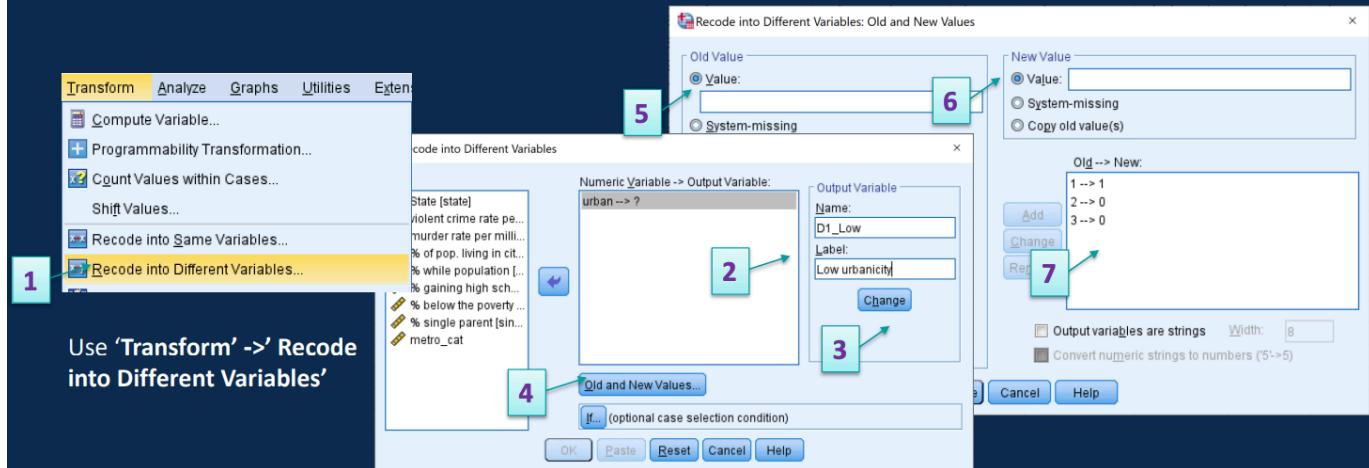
Dummy coding of
urban ($k=3$)

d1	d2	d3
1	0	0
1	0	0
1	0	0
1	0	0
1	0	0
1	0	0
0	1	0
0	1	0
0	1	0
0	1	0
0	1	0
0	1	0
0	0	1
0	0	1
0	0	1
0	0	1
0	0	1

SPSS Slide: 'how to'

Researchers believe there is a relationship between Violent Crime and poverty and the level of urbanicity in an area modifies this effect. The variable `urban` is a categorical variable with three levels "Low", "Medium" and "High" and needs to be converted to dummy variables to include in the regression.

Step 1: Generating dummy variables for 'urban' variable from US crime

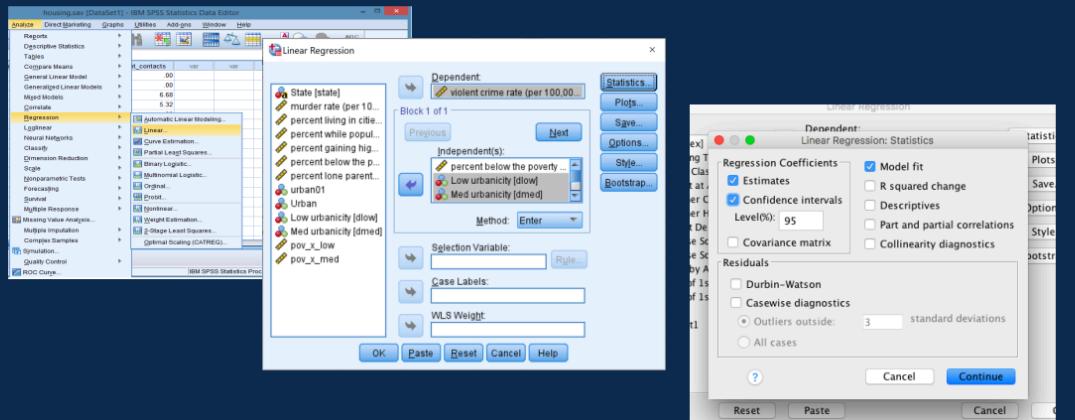


- Create an interaction term `pov_x_low` and `pov_x_med` where high urbanicity is the reference from Lecture_9b_data.sav
- Use 'Transform' -> 'Compute variable'
- In 'Target variable' write the name of your interaction term: "pov_x_low"
- In 'Numeric Expression' drag 'poverty' times (*) 'low' and accept.
- Repeat for "pov_x_med"

	pov_x_low	pov_x_med	var
1	.00	.00	9.10
2	.00	.00	.00
3	.00	.00	20.00
4	.00	.00	.00
5	.00	.00	.00
6	.00	.00	.00
7	.00	.00	.00
8	.00	.00	.00
9	.00	.00	.00
10	.00	.00	.00
11	.00	.00	.00
12	.00	.00	.00
13	.00	.00	.00
14	.00	.00	.00
15	.00	.00	.00
16	.00	.00	.00
17	.00	.00	.00
18	.00	.00	.00
19	.00	.00	.00
20	.00	.00	.00

New variables in data set

- Estimating the interaction effect **pov_x_low** and **pov_x_med** in a multiple linear regression model for crime rate, poverty, dlow and dmed from lecture_9b_data.sav data
- 1) Use 'Analyse' -> 'Regression' -> 'Linear'
- 2) In dependent put 'crime' and in independent put 'poverty', 'dlow', 'dmed', 'pov_x_low' and 'pov_x_med'



Output and Interpretation

Coefficients ^a						
Model	Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.	95.0% Confidence Interval for B Lower Bound
	B	Std. Error				Upper Bound
1	(Constant)	-296.662	179.306		-1.655	.105
	percent below the poverty line	74.694	12.195	.776	6.125	.000
	Low urbanicity	263.137	468.308	.194	.562	.577
	Med urbanicity	481.824	337.218	.467	1.429	.160
	pov_x_low	-55.812	30.105	-.650	-1.854	.070
	pov_x_med	-58.218	22.288	-.876	-2.612	.012

a. Dependent Variable: violent crime rate (per 100,000 people)

crime

$$= -296.662 + 74.694\text{poverty} + 263.137\text{low} + 481.824\text{med} - 55.812\text{pov * low} - 58.218\text{pov * med}$$

The Coefficient of **pov × low** interaction is -55.812 , $p=0.070$

The Coefficient of **pov × med** interaction is -58.218 , $p=0.012$

Effect of poverty on crime decreases in low and medium urbanised areas compared to high urbanised areas

The mean crime rate at average poverty level (mean = 14.2588) for low urbanised states = $-296.662 + 74.694 \times 14.2588 + 263.137 - 55.812 \times 14.2588 = 235.71$ per 100,000 people.

The mean crime rate at average poverty level (mean = 14.2588) for med urbanised states = $-296.662 + 74.694 \times 14.2588 + 481.824 - 58.218 \times 14.2588 = 420.09$ per 100,000 people, only the interaction between poverty and med urbanised areas showed a significant effect.

Interaction & Type of Variables

Interaction between variables where both independent variables (x_1 and Z) are either categorical or continuous is handled in the same way, i.e., by creating cross-product terms:

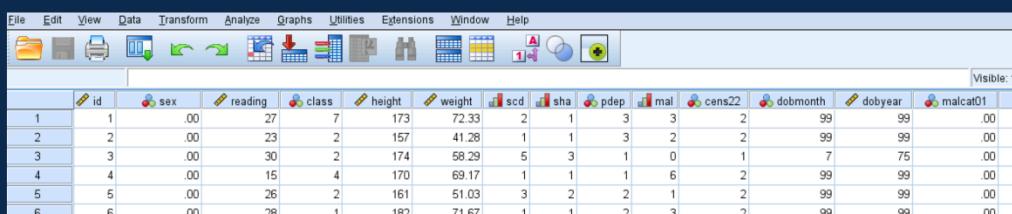
- **continuous × continuous**
- **categorical × categorical**

Example: Continuous × Continuous Interaction

- In Lecture_9a_data.sav, The dataset contains data from 1000 individuals, from the National Child Development Study (NCDS), both height and reading scores are **continuous** variables
- There is **no reason** to believe that reading score will affect weight, but let's see an example involving reading score to demonstrate how we can investigate interactions when the two independent variables are continuous.
- We are interested in testing if reading score modifies the effect of height on weight
- This will require **computing a new variable** – the cross-product of height and reading score (as we did before for height x sex) **height × reading**
- And then **including the product term** as an additional predictor in a regression model

SPSS Slide

Download the data that we are going to use during the lecture. The dataset is the **lecture_9a_data.sav**.



The screenshot shows the SPSS Data View window. 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 data view itself has a header row with column names: id, sex, reading, class, height, weight, scd, sha, pdep, mal, cens22, dobmonth, dobyear, and malcat01. Six data rows are visible, each containing numerical values for these variables. Row 1: id=1, sex=.00, reading=27, class=7, height=173, weight=72.33, scd=2, sha=1, pdep=3, mal=3, cens22=2, dobmonth=99, dobyear=99, malcat01=.00. Row 2: id=2, sex=.00, reading=23, class=2, height=157, weight=41.28, scd=1, sha=1, pdep=3, mal=2, cens22=2, dobmonth=99, dobyear=99, malcat01=.00. Row 3: id=3, sex=.00, reading=30, class=2, height=174, weight=58.29, scd=5, sha=3, pdep=1, mal=0, cens22=1, dobmonth=7, dobyear=75, malcat01=.00. Row 4: id=4, sex=.00, reading=15, class=4, height=170, weight=69.17, scd=1, sha=1, pdep=1, mal=6, cens22=2, dobmonth=99, dobyear=99, malcat01=.00. Row 5: id=5, sex=.00, reading=26, class=2, height=161, weight=51.03, scd=3, sha=2, pdep=2, mal=1, cens22=2, dobmonth=99, dobyear=99, malcat01=.00. Row 6: id=6, sex=.00, reading=28, class=1, height=182, weight=71.67, scd=1, sha=1, pdep=2, mal=3, cens22=2, dobmonth=99, dobyear=99, malcat01=.00.

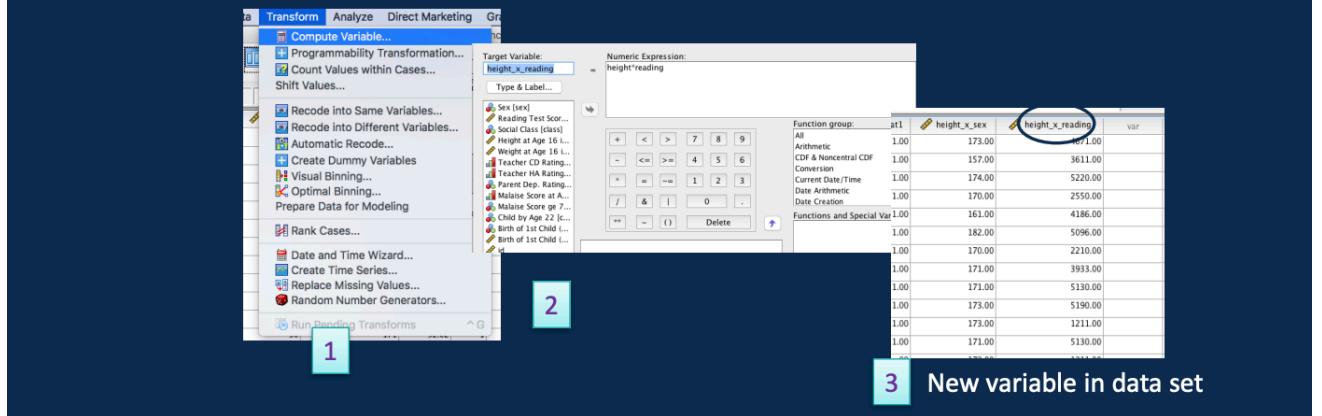
	id	sex	reading	class	height	weight	scd	sha	pdep	mal	cens22	dobmonth	dobyear	malcat01
1	1	.00	27	7	173	72.33	2	1	3	3	2	99	99	.00
2	2	.00	23	2	157	41.28	1	1	3	2	2	99	99	.00
3	3	.00	30	2	174	58.29	5	3	1	0	1	7	75	.00
4	4	.00	15	4	170	69.17	1	1	1	6	2	99	99	.00
5	5	.00	26	2	161	51.03	3	2	2	1	2	99	99	.00
6	6	.00	28	1	182	71.67	1	1	2	3	2	99	99	.00

The dataset contains data from 1000 individuals, from the National Child Development Study (NCDS) with respect to their

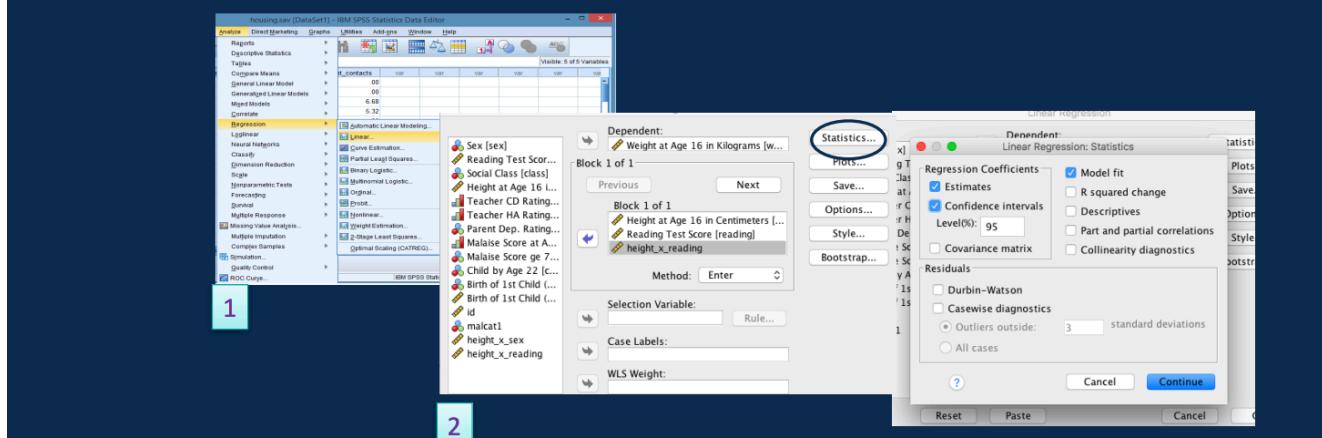
- **sex**: gender of child (0=male, 1=female)
- **height**: height in cm at age 16
- **weight**: weight in kg at age 16
- **reading**: reading score
- **mal**: malaise (a feeling of general discomfort/uneasiness) score
- **class**: general classification of social class (7 Categories)

SPSS Slide: 'How to' Steps

- Create an interaction term `height_x_reading` from `ncds.sav` data
- Use 'Transform' -> 'Compute variable'
- In 'Target variable' write the name of your interaction term: "`height_x_reading`"
- In 'Numeric Expression' drag 'height' times (*) 'reading' and accept.



- Estimating the interaction effect `height_x_reading` in a multiple linear regression model for `weight`, `height` and `reading` from `lecture_9_a_data.sav` data
- 1) Use 'Analyse' -> 'Regression' -> 'Linear'
- 2) In dependent put '`weight`' and in independent put '`height`', '`reading`', '`height_x_reading`



Output and Interpretation

- The Coefficient of **height × reading** interaction is - 0.005
- **Negative interaction** effect means that:
 - Effect of height decreases as reading scores increases, and
 - Effect of reading scores decreases as height increases
- However, the **height × reading** interaction is **not significant** ($p=0.286$) The height-weight relationship does not significantly differ by reading score

Model	Coefficients ^a						
	Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.	95.0% Confidence Interval for B	
	B	Std. Error				Lower Bound	Upper Bound
1	(Constant)	-67.302	19.900	-3.382	.001	-106.352	-28.251
	Height at Age 16 in Centimeters	.748	.119	.622	.265	.514	.983
	Reading Test Score	.858	.800	.582	1.072	.284	-.712
	hxr	-.005	.005	-.588	-1.068	.286	-.015

a. Dependent Variable: Weight at Age 16 in Kilograms

$$\text{weight} = -67.302 + 0.748\text{height} + 0.858\text{reading} - 0.005\text{height} * \text{reading}$$

- **Negative interaction coefficient** (-0.005) = *moderating effect*.
- When one variable increases, the effect of the other variable **gets smaller**.
- The variables **suppress** each other's effects slightly.

Presenting Continuous × Continuous Interactions: Tabular Format

- The effect for height on weight is: $\beta_1 + \beta_3 \times \text{reading}$
- The effect for reading on weight is: $\beta_2 + \beta_3 \times \text{height}$
- For example, in the model for NCDS data, effect of height can be calculated at different values (e.g., quartiles) of reading scores, and vice-versa:

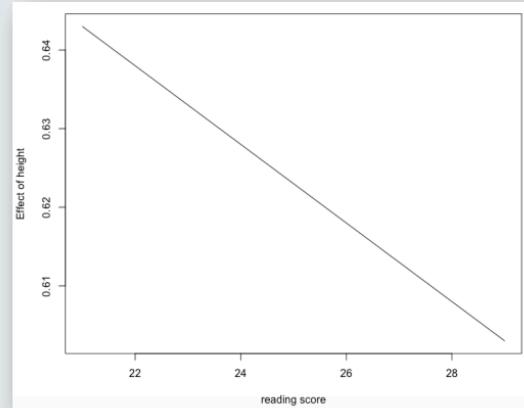
$$\text{weight} = -67.302 + 0.748\text{height} + 0.858\text{reading} - 0.005\text{height} * \text{reading}$$

Reading score (quartiles)	Effect for height: $0.748 - 0.005 \times \text{reading}$
reading= 21 (first quartile)	$0.748 - 0.005 \times 21 = 0.643 \text{ kg/cm}$
reading= 27 (median)	$0.748 - 0.005 \times 27 = 0.613 \text{ kg/cm}$
reading= 29 (3 rd quartile)	$0.748 - 0.005 \times 29 = 0.603 \text{ kg/cm}$

*Similar table can be created for the effect of reading scores at varying values of height

Presenting Continuous × Continuous Interactions: Graphical Format

Reading score score (quartiles)	Effect for height:
21	0.643
27	0.613
29	0.603



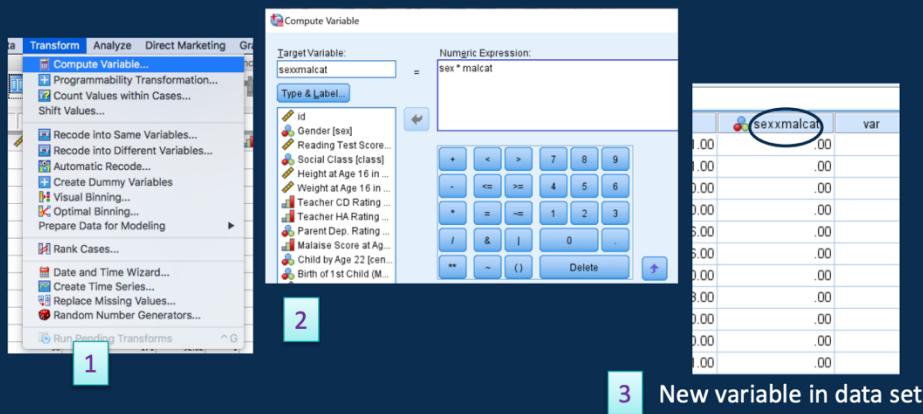
- The plot shows the effect of **height** as a function of reading scores
- Effect of height **decreases** as reading scores **increases**
- Similar plot can be created for the effect of reading scores as a function of height

Example: Categorical × Categorical Interaction

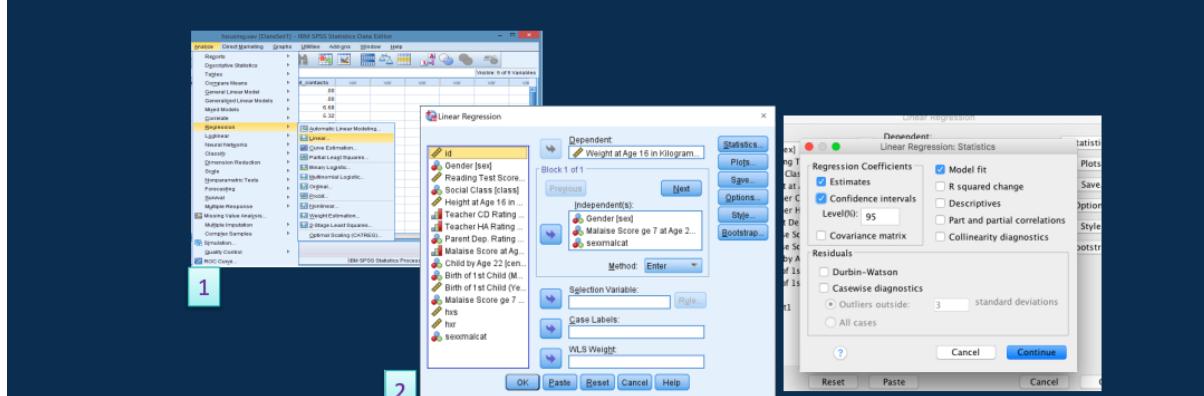
- In NCDS data, **sex** and **malcat** are two **categorical** (binary) variables
- The variable **malcat (0=low, 1=high)** represents a categorised version (median split) of the continuous variable malaise scores (**mal**) (a feeling of general discomfort/uneasiness)
- Suppose we are interested in testing the **sex × malcat** interaction
- As before, this will require computing a new variable – the **cross-product of sex and malcat**, and including it as an additional predictor in the regression model.

- Create an interaction term **sex_x_malcat** from **lecture_9a_data.sav**.

- Use ‘Transform’ -> ‘Compute variable’
- In ‘Target variable’ write the name of your interaction term: “**sex_x_malcat**”
- In ‘Numeric Expression’ drag ‘Gender’ times (*) ‘malcat’ and ‘Ok’.



- Estimating the interaction effect **sex_x_malcat** in a multiple linear regression model for weight, sex and malcat from **lecture_9_a_data.sav** data
- 1) Use ‘Analyse’ -> ‘Regression’ -> ‘Linear’
- 2) In dependent put ‘weight’ and in independent put ‘sex’, ‘malcat’, ‘sex_x_malcat’



Output and Interpretation

Model	Coefficients ^a							
	Unstandardized Coefficients		Standardized Coefficients	t	Sig.	95.0% Confidence Interval for B		
	B	Std. Error	Beta			Lower Bound	Upper Bound	
1	(Constant)	59.534	.435	136.706	.000	58.680	60.389	
	Gender	-4.676	.626	-.242	.7474	.000	-5.904	-3.448
	Malaise Score ≥ 7 at Age 22 0 = No, 1 = Yes	-.324	1.892	-.010	-.171	.864	-4.038	3.390
	sexmalcat	.690	2.238	.018	.309	.758	-3.701	5.082

a. Dependent Variable: Weight at Age 16 in Kilograms

$$\text{weight} = 59.534 - 4.676 \text{ sex} - 0.324 \text{ malcat} + 0.690 \text{ sex * malcat}$$

- Coefficient of **sex × malcat** interaction = 0.690
- Positive interaction effect means that:
 - Effect of gender is **higher** for high (=1) category of malaise score, and
 - Effect of malaise score is **higher** for girls (sex=1) than for boys (sex=0)
- The **sex × malcat** interaction is **not significant** ($p=0.758$)

Presenting Categorical × Categorical Interactions

- Effect of each variable can be estimated at each level of the other variable
- For example, effect of gender can be calculated at low and high levels of malaise scores
- Effect of sex on weight = $\beta_1 + \beta_3 \times \text{malcat} = -4.676 + 0.690 \times \text{malcat}$

$$\text{weight} = 59.534 - 4.676 \text{ sex} - 0.324 \text{ malcat} + 0.690 \text{ sex * malcat}$$

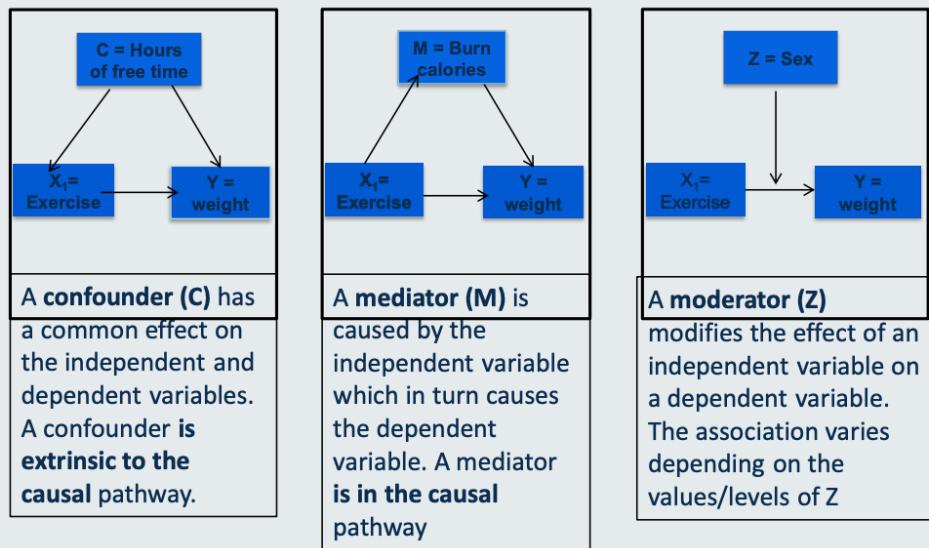
malcat = -4.676 + 0.690 × malcat	
Low (=0)	-4.676 kg ←
High (=1)	-3.986 kg ←

Difference in mean weight between girls (sex=1) and boys (sex=0) at low malaise scores

Difference in mean weight between girls (sex=1) and boys (sex=0) at high malaise scores

Confounding vs Mediation vs Interaction

- Both confounder, mediator and moderator, are third variables that explain a part (or most) of the association between an independent and dependent variable.



🔍 Outliers and Influential Points

⚠ 1. Outliers

- Standardized residuals $> |2|$ are potential outliers
- Use: `Plots > Save standardized residuals` in SPSS

⚠ 2. Influential Points

- Measured with **Cook's Distance** and **Leverage**
- SPSS: Save them via `Save > Influence Statistics`
- Thresholds:
 - Cook's D** > 1 may indicate high influence
 - Leverage** $> (2k+2)/n$

Outliers and Influential Points

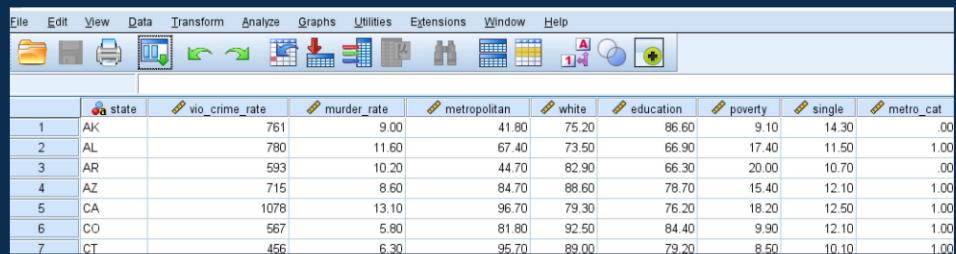
- An **outlier** is an observation that lies an abnormal distance from other values in a random sample from a population.
- Outliers can be problematic for many statistical analyses because they can cause tests to either miss significant findings or distort real results.
- Finding outliers depends on subject-area knowledge and an understanding of the data collection process.

Outliers and Influential Points in Regression

- All outliers are not harmful. Some outliers influence the regression model more than the others
- Outliers with large influence on the fitted regression model are called **influential observations**
- Influential observations need special attention as they may distort the actual relationship

SPSS Slide

Download the data that we are going to use during the lecture. The dataset is the [lecture_9b_data.sav](#).



	state	vio_crime_rate	murder_rate	metropolitan	white	education	poverty	single	metro_cat
1	AK	761	9.00	41.80	75.20	86.60	9.10	14.30	.00
2	AL	780	11.60	67.40	73.50	66.90	17.40	11.50	1.00
3	AR	593	10.20	44.70	82.90	66.30	20.00	10.70	.00
4	AZ	715	8.60	84.70	88.60	78.70	15.40	12.10	1.00
5	CA	1078	13.10	96.70	79.30	76.20	18.20	12.50	1.00
6	CO	567	5.80	81.80	92.50	84.40	9.90	12.10	1.00
7	CT	456	6.30	95.70	89.00	79.20	8.50	10.10	1.00

The dataset contains data from 51 US states, measuring the crime rates and background measures for each State with respect to their

- **violent crime**: per 100,000 population
- **murder** : per 100,000 population
- **poverty**: percent below the poverty line
- **single**: percentage of lone parents
- **urban**: level of urbanicity

Outliers

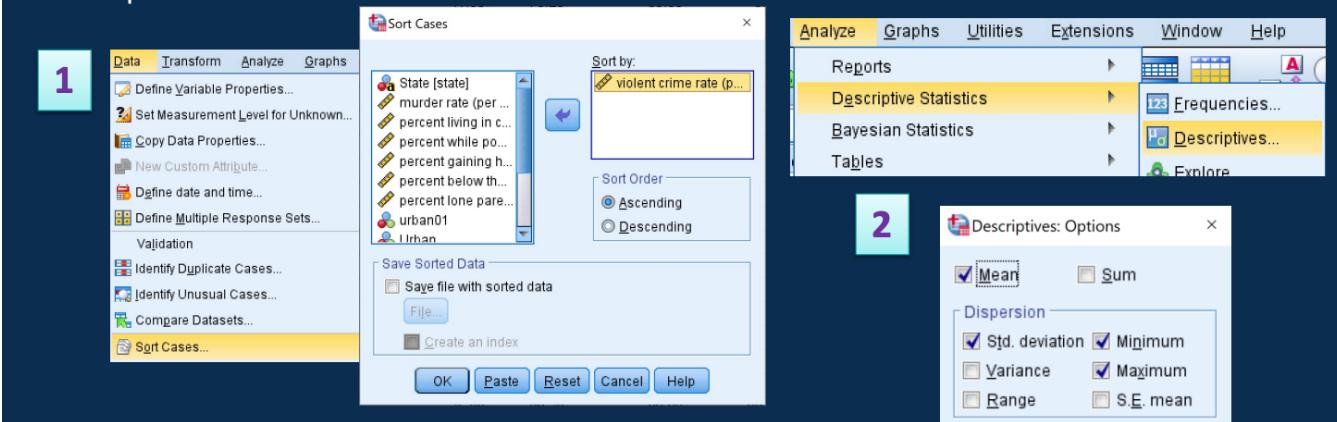
Variable	With Outlier		Without Outlier		Mean Difference	Std. Difference
	Mean	Std. Deviation	Mean	Std. Deviation		
violent crime rate (per 100,000 people)	612.84	441.1	566.66	295.9	46.18	145.2
murder rate (per 100,000 people)	8.33	11.0	6.92	4.6	1.40	6.4

- To demonstrate how much a single outlier can affect the results, let's examine the effect of a potential outlier in the lecture_9b_data.sav.
- The table above shows the mean and standard deviation for violent crime and murder rate with and without the potential outlier.
- From the table, it's easy to see how a single outlier can distort the data summaries. A single value changes the mean crime rate by 46.18 (per 100 000) and the standard deviation by a large amount 145.2.

SPSS Slide: Finding Outliers and Influential Points

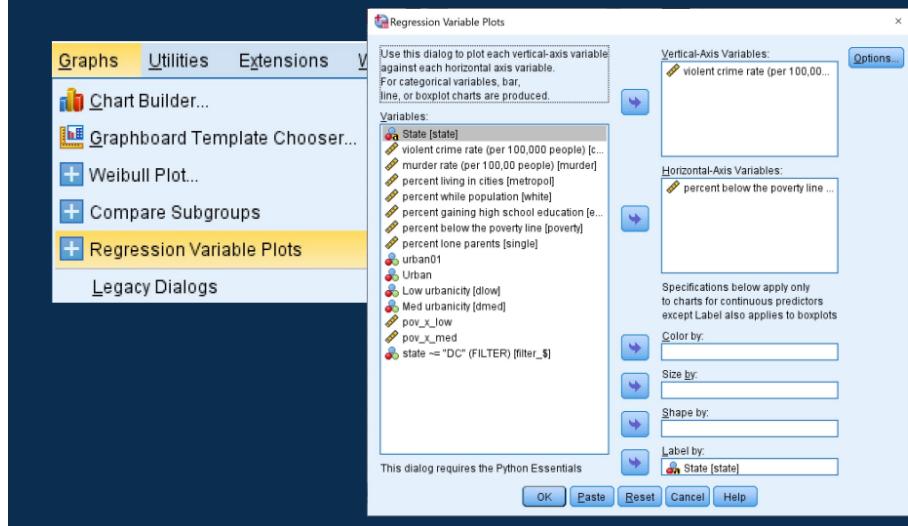
Sorting Your Datasheet to Find Outliers

- Sorting your datasheet is a simple but effective way to highlight unusual values. Simply sort your data sheet for each variable and then look for unusually high or low values.
- Alternatively, when asking for “Descriptives” ask for the minimum and maximum to be included in the output



Graphing Your Data to Identify Outliers

- Boxplots, histograms, and scatterplots can highlight outliers



In SPSS you are able to now create a Regression variable plot which shows a scattergraph of two variables and a box plot of their data.

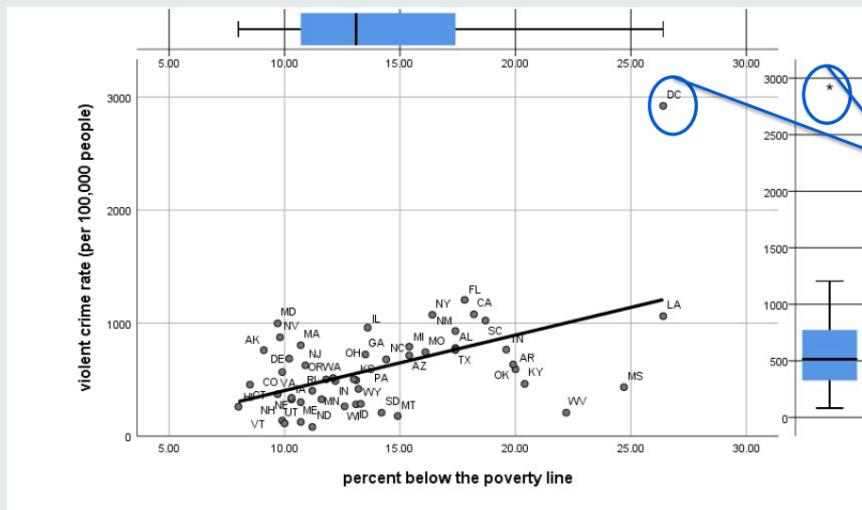
Graphs->Regression Variable Plot -> put dependent variable in the vertical axis, and the independent variable in the horizontal axis

Label by "state" so you can identify any outliers.

Output: Finding Outliers and Influential Points

Graphing Your Data to Identify Outliers

- Boxplots, histograms, and scatterplots can highlight outliers



Finding Outliers and Influential Points

Tukey's Method: Using the Interquartile Range

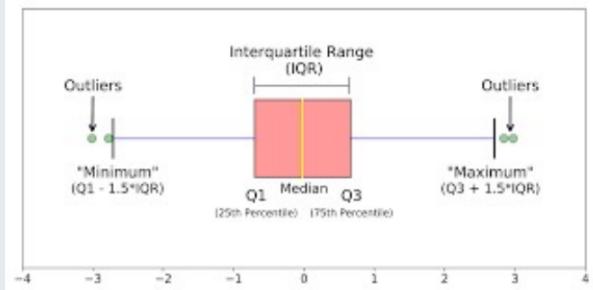
The **IQR** is the middle 50% of the dataset. It's the range of values between the third quartile and the first quartile ($Q_3 - Q_1$).

We can take the IQR, Q_1 , and Q_3 values to calculate the following outlier fences for our dataset: lower outer, lower inner, upper inner, and upper outer.

These fences determine whether data points are outliers and whether they are **mild** or **extreme**.

Extreme outliers tend to lie more than **3** times the interquartile range (below the first quartile or above the third quartile), and

Mild outliers lie between **1.5** and three times the interquartile range (below the first quartile or above the third quartile).



Finding Outliers and Influential Points

Example:

Statistics		
murder rate (per 100,00 people)		
N	Valid	51
	Missing	0
Mean		8.3275
Median		6.6000
Minimum		-9.00
Maximum		78.50
Percentiles	25	3.8000
	50	6.6000
	75	10.3000

$$Q_1 = 3.80$$

$$Q_3 = 10.30$$

$$IQR = Q_3 - Q_1$$

$$IQR = 6.5$$

$$\text{Lower Outer} = Q_1 - 3 \times IQR = -15.7$$

$$\text{Lower Inner} = Q_1 - 1.5 \times IQR = -5.95$$

$$\text{Upper Inner} = Q_3 + 1.5 \times IQR = 20.5$$

$$\text{Upper Outer} = Q_3 + 3 \times IQR = 29.8$$

Order your data:

-9 murder rate for 'IL' is a **mild outlier** as it lies between the lower inner and outer limits

78.50 murder rate for "DC" is a **extreme outlier** as it lies outside of the upper outer limit

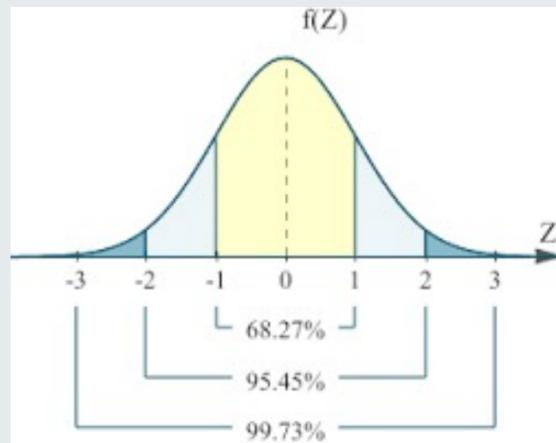
Finding Outliers and Influential Points

Using the Standard Deviation

The **standard deviation (SD)** is a reasonable method to detect outliers when the data distribution is symmetric such as the normal distribution.

68%, 95%, and 99.7% of the data from a normal distribution are within 1, 2, and 3 standard deviations of the mean, respectively.

If data follows a normal distribution, this helps to estimate the likelihood of having extreme values in the data, so that the observation **two or three standard deviations** away from the mean may be considered as an outlier in the data.



Outliers and Influential Observations

Using Standardised Residuals

The good thing about standardized residuals is that they quantify how large the residuals are in standard deviation units, and therefore can be easily used to identify outliers: An observation with an **Absolute standardised residual** that is **larger than 3** (in absolute value) is deemed by some to be an outlier.

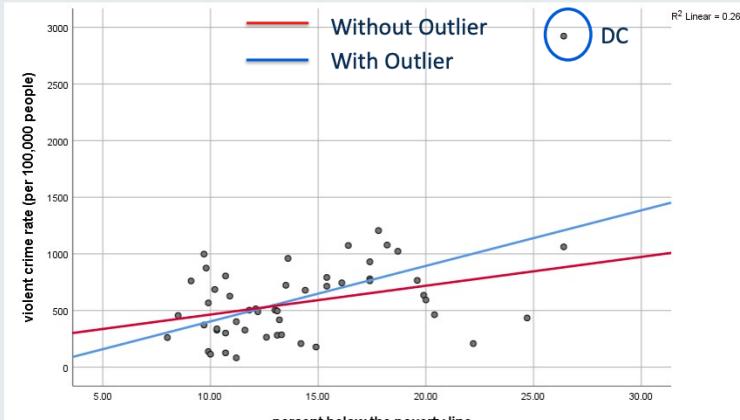
Outliers and Influential Observations

DFBETA and DFFIT

- **DFBETA** and **DFFIT** are two diagnostic measures for flagging influential observations
- For a given observation, **DFBETA** measures the **change in the estimated coefficient β_j** due to deleting that observation
 - Standardised **DFBETA** is defined as **DFBETA divided by the SE (est β_j)** for the adjusted dataset
- For a given observation, **DFFIT** measures the **change in the predicted value (\hat{y})** due to deleting that observation
 - Standardised **DFFIT** is defined as **DFFIT divided by SE(\hat{y})** for the adjusted data
- A general guideline:
 - Absolute **standardised DFBETA > 1** suggests **influential observations**
 - Absolute **standardised DFFIT > 1** suggests **influential observations**

Influential Observations

Consider the following Scatterplot from Lecture_9b_data.sav showing the US crime rate. The figure shows that the state DC is an outlier. Crime rate is very high in DC compared to the other states

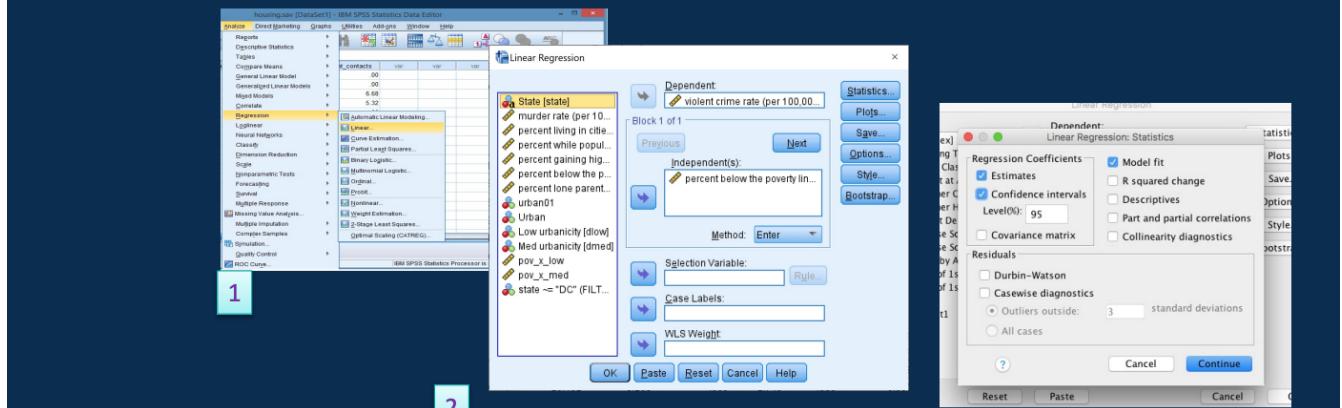


The slope estimated by including the outlier is much higher than the slope estimated with the outlier removed.

The slope difference is the **influence** of the outlier state DC

SPSS Slide: 'How to' Steps

- Researchers believe that the state of DC is giving a distorted understanding of the Crime – poverty relationship. They have decided to run an analysis including this potential outlier and without it to check the level of influence
- 1) Use 'Analyse' -> 'Regression' -> 'Linear'
- 2) In dependent put 'crime' and in independent put 'poverty',



Step 2: Create a filter to exclude DC

Go to Data → Select Cases.

Choose: "If condition is satisfied", then click "If..."

In the condition window, enter something like:

```
perl
```

Copy

Edit

```
state ≠ "DC"
```

or, if DC is coded as a number, use:

```
rust
```

Copy

Edit

```
state_code ≠ [DC's code]
```

Click Continue, then OK.

⌚ This step tells SPSS to temporarily exclude DC from any analysis.

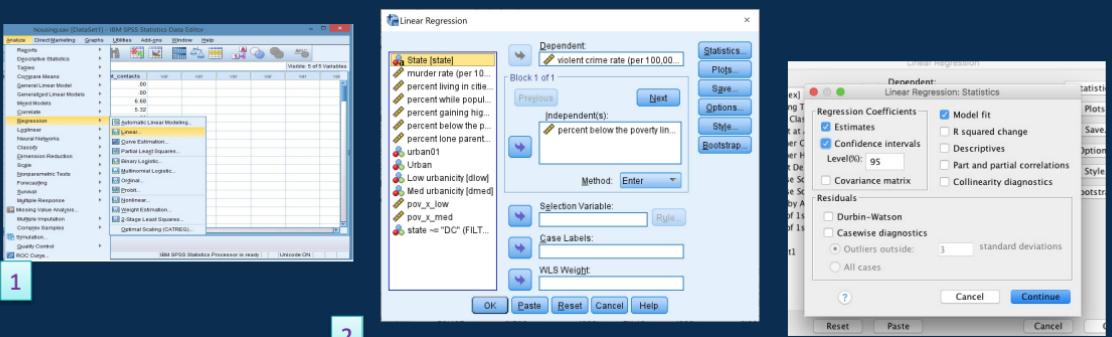
- Try something like:

```
nginx
```

```
State ≈ "DC"
```

SPSS Slide: 'How to' Steps

- Researchers believe that the state of DC is giving a distorted understanding of the Crime – poverty relationship. They have decided to run an analysis including this potential outlier and without it to check the level of influence. Use 'Select Cases' option under the 'Data' Menu to remove the outlier from the analysis. Re-run the regression
- 1) Use 'Analyse' -> 'Regression' -> 'Linear'
- 2) In dependent put 'crime' and in independent put 'poverty',



Output and Interpretation

The first table were generated for data in all US states, whilst the second table was generated excluding DC state.

All data		Coefficients^a <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th rowspan="2">Model</th> <th colspan="2">Unstandardized Coefficients</th> <th rowspan="2">Standardized Coefficients Beta</th> <th rowspan="2">t</th> <th rowspan="2">Sig.</th> <th colspan="2">95.0% Confidence Interval for B</th> </tr> <tr> <th>B</th> <th>Std. Error</th> <th>Lower Bound</th> <th>Upper Bound</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>(Constant)</td> <td>-86.201</td> <td>176.990</td> <td>-487</td> <td>.628</td> <td>-441.876</td> <td>269.474</td> </tr> <tr> <td></td> <td>percent below the poverty line</td> <td>49.025</td> <td>11.828</td> <td>.510</td> <td>4.145</td> <td>.000</td> <td>25.256</td> <td>72.794</td> </tr> </tbody> </table> <p>a. Dependent Variable: violent crime rate (per 100,000 people)</p>	Model	Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.	95.0% Confidence Interval for B		B	Std. Error	Lower Bound	Upper Bound	1	(Constant)	-86.201	176.990	-487	.628	-441.876	269.474		percent below the poverty line	49.025	11.828	.510	4.145	.000	25.256	72.794
Model	Unstandardized Coefficients			Standardized Coefficients Beta	t				Sig.	95.0% Confidence Interval for B																					
	B	Std. Error	Lower Bound			Upper Bound																									
1	(Constant)	-86.201	176.990	-487	.628	-441.876	269.474																								
	percent below the poverty line	49.025	11.828	.510	4.145	.000	25.256	72.794																							
Excluding state DC		Coefficients^a <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th rowspan="2">Model</th> <th colspan="2">Unstandardized Coefficients</th> <th rowspan="2">Standardized Coefficients Beta</th> <th rowspan="2">t</th> <th rowspan="2">Sig.</th> <th colspan="2">95.0% Confidence Interval for B</th> </tr> <tr> <th>B</th> <th>Std. Error</th> <th>Lower Bound</th> <th>Upper Bound</th> </tr> </thead> <tbody> <tr> <td>1</td> <td>(Constant)</td> <td>209.920</td> <td>135.613</td> <td>1.548</td> <td>.128</td> <td>-62.748</td> <td>482.588</td> </tr> <tr> <td></td> <td>percent below the poverty line</td> <td>25.452</td> <td>9.260</td> <td>.369</td> <td>2.749</td> <td>.008</td> <td>6.833</td> <td>44.072</td> </tr> </tbody> </table> <p>a. Dependent Variable: violent crime rate (per 100,000 people)</p>	Model	Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.	95.0% Confidence Interval for B		B	Std. Error	Lower Bound	Upper Bound	1	(Constant)	209.920	135.613	1.548	.128	-62.748	482.588		percent below the poverty line	25.452	9.260	.369	2.749	.008	6.833	44.072
Model	Unstandardized Coefficients			Standardized Coefficients Beta	t				Sig.	95.0% Confidence Interval for B																					
	B	Std. Error	Lower Bound			Upper Bound																									
1	(Constant)	209.920	135.613	1.548	.128	-62.748	482.588																								
	percent below the poverty line	25.452	9.260	.369	2.749	.008	6.833	44.072																							

- ➡ DFBETA for the poverty variable can be calculated as the difference of the coefficient from the full model to the adjusted models ($49.025 - 25.452 = 23.573$) standardized by accounting for the full model and adjusted model error and covariance (not covered in this course).
- ➡ To estimate the standardized DFBETA and standardized DFFIT, we select this option from SPSS as it is shown in the next slide
- ➡ As per SPSS, standardized DFBETA for the coefficient for poverty is = 2.75

1

SPSS Slide: 'How to' Steps

- Dfbeta and Dffit in SPSS
- 1) Use 'Analyse' -> 'Regression' -> 'Linear'
- 2) In dependent put 'crime' and in independent put 'poverty',
- 3) Click on 'Save'

The image shows a sequence of four screenshots illustrating the SPSS steps for calculating DFBETA and DFFIT:

- Screenshot 1:** The SPSS menu bar is shown with "Analyze" selected. The "Regression" option is expanded, and "Automatic Linear Modeling..." is highlighted.
- Screenshot 2:** The "Linear Regression" dialog box is open. The "Dependent" field contains "violent crime rate (per 100,000)". The "Independent(s)" field contains "percent below the poverty line". The "Method" dropdown is set to "Enter". The "Save..." button is circled in green.
- Screenshot 3:** The "Linear Regression: Save" sub-dialog box is open. Under "Predicted Values", "Unstandardized" is checked. Under "Influence Statistics", "DfBeta(s)" and "DFFIT" are checked. The "Save..." button is circled in green.
- Screenshot 4:** The output window displays a table of results. The last row shows the value 2.74990, which is circled in red.

Quiz:

Given the model: $\text{years_married} = 6 - 2.4 \text{ problems_inlaw_family} - 3.1 \text{ cheating_with_others} - 3.7 \text{ problems_inlaw_family} * \text{cheating_with_others}$

Select one:

- 1. The effect of problems_inlaw_family on years_married is $\beta_1 = -2.4$ ×
- 2. The effect of problems_inlaw_family on years_married is $= \beta_1 + (\beta_3 * \text{cheating_with_others}) = -2.4 + (-3.7 * \text{cheating_with_others})$
- 3. The effect of problems_inlaw_family on years_married is $= 6 - 2.4 - 3.1 + (-3.7 * \text{cheating_with_others})$

Your answer is incorrect.

The correct answer is: The effect of problems_inlaw_family on years_married is $= \beta_1 + (\beta_3 * \text{cheating_with_others}) = -2.4 + (-3.7 * \text{cheating_with_others})$

❓ What the question is asking:

What is **the effect of problems_inlaw_family on years_married ?**

➡ The key point is: **there is an interaction term between problems_inlaw_family and cheating_with_others .**

This means the effect of **problems_inlaw_family depends on the level of cheating_with_others .**

Question 5

Complete

Not graded

Flag

The effect of motivation depended on whether there was a clear standard for excellence.

Indicate whether this sentence implies that there was a main effect of motivation, a main effect of standard, or an interaction effect, and give a reason why.

Concept	Meaning
Main Effect	Independent influence of one variable
Interaction Effect	When one variable's effect depends on another

Practical Quiz:

1)

Please estimate Standardised DFBETA and DFFIT measures for the model with no interaction and fills the gaps.

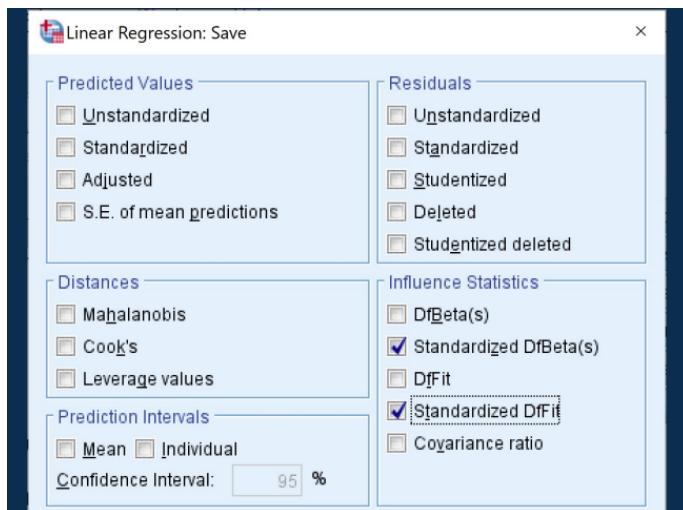
'stateanx presents ✓ influential value(s) because the absolute standardised DFBETA and DFFIT are ✓ . epiNeur presents ✓ influential value(s) because the absolute standardised DFBETA and DFFIT are ✓ '

◆ 2. Run a Linear Regression (no interaction model)

- Go to Analyze → Regression → Linear
- Dependent variable:
- Independent variables: and (! exclude interaction terms for this question)

◆ 3. Request DFBETA and DFFIT

- In the Linear Regression dialog box, click on Save...
- Tick the boxes for:
 - Standardized DFBETAs (for each predictor)
 - Standardized DFFIT
- Click Continue, then OK



- New columns will appear in your SPSS Data View:
- If any value > 1, it's considered **influential**.

- Absolute **standardised DFBETA > 1** suggests **influential** observations
- Absolute **standardised DFFIT > 1** suggests **influential** observations

🧠 Example for Interpretation:

Say you found:

- SDFB_stateanx : has 2 values > 1 → 2 influential values, DFBETA is **higher than 1**
- SDFB_epiNeur : all values < 1 → 0 influential values, DFBETA is **lower than 1**

Option A: Use Descriptive Stats

1. Go to **Analyze → Descriptive Statistics → Frequencies**
2. Select **SDB1_1** and **SDB2_1**
3. Click **Statistics** → tick **Minimum** and **Maximum** → OK
4. Look for any values where the **absolute value > 1**

2)

Use a multiple linear regression model to assess if stateanx is an effect modifier of the epiNeur-bdi association.

What does the regression coefficient of stateanx_X_epiNeur tells you?

Regression coefficient of stateanx_X_epiNeur represents **the interaction effect** ✓
between stateanx and epiNeur. The p-value **0.001** ✓ suggests that
the interaction effect ✓ is statistically **significant** ✓ . This implies
that both variables **independent** ✗ , but their effects are
not independent ✓ of each other. Effect of stateanx
depends ✓ on epiNeur and vice-versa.

statistically [significant]. This implies that both variables [jointly affect depression], but
their effects are [not independent] of each other. Effect of stateanx [depends] on
epiNeur and vice-versa.

3)

Based on the interpretation of the estimated coefficients for stateanx and epiNeur from the fitted interaction model, fill the gaps:

The coefficients of the variables stateanx and epiNeur do not ✓
 carry their usual interpretations because of the presence of an interaction (cross-product) term involving these variables. For the interaction model, the coefficient of stateanx can be interpreted as the effect of stateanx on bdi when
stateanx different from 0 ✗ . The estimated coefficient do ✓
✗ implies in people with 0 ✓ neuroticism symptoms, one unit increase in stateanx leads to 0 ✗ units increase in bdi.
 Similarly, the coefficient of epiNeur represents the effect of epiNeur on bdi when
stateanx=0 ✓ . Both coefficients are not significantly equal to 0 ✗ as the pvalue for the test for the stateanx's beta coefficient is 0.038 ✗ and for epiNeur's coefficient is 0.435 ✓ .

The coefficients of the variables stateanx and epiNeur [do not] carry their usual interpretations because of the presence of an interaction (cross-product) term involving these variables. For the interaction model, the coefficient of stateanx can be interpreted as the effect of stateanx on bdi when [epiNeur=0]. The estimated coefficient [0.038] implies in people with [0] neuroticism symptoms, one unit increase in stateanx leads to [0.038] units increase in bdi. Similarly, the coefficient of epiNeur represents the effect of epiNeur on bdi when [stateanx=0]. Both coefficients are not significantly [different from 0] as the pvalue for the test for the stateanx's beta coefficient is [0.537] and for epiNeur's coefficient is [0.435].

Model	Coefficients ^a						
	Unstandardized Coefficients		Standardized Coefficients Beta	t	Sig.	95.0% Confidence Interval for B	
	B	Std. Error				Lower Bound	Upper Bound
1	(Constant)	.064	2.186	.029	.977	-4.243	4.370
	epiNeur	-.148	.189	-.125	-.782	.435	-.519
	stateanx	.038	.061	.075	.619	.537	-.082
	neur_x_anx	.015	.005	.751	3.279	.001	.006
							.157
							.024

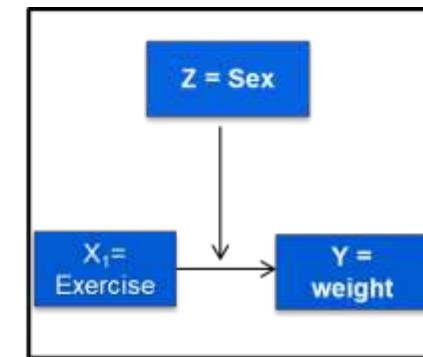
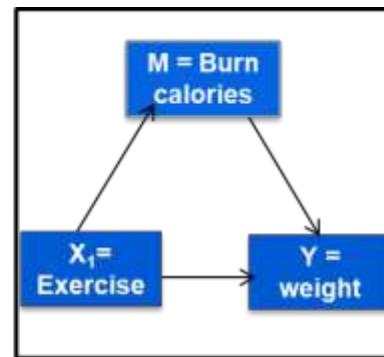
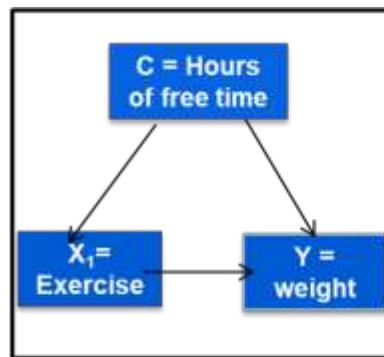
a. Dependent Variable: bdi

- The **p-value** tells us whether the observed effect is statistically significant — i.e., unlikely due to chance.
- A **p-value > 0.05** means we cannot confidently say that the effect is different from 0.
- So:
 - For **epiNeur**, with a p-value of 0.435 → we can't say it significantly affects BDI when stateanx = 0.
 - For **stateanx**, with a p-value of 0.537 → we can't say it significantly affects BDI when epiNeur = 0.

Topic 9 Knowledge Check Quiz Results

Question 1

- Let's focus on the 3 variables (y , x_1 and x_2) case. Which three different roles for x_2 have been discussed?
- confounder, mediator, moderator



A **confounder (C)** has a common effect on the independent and dependent variables. A confounder is **extrinsic to the causal pathway**.

A **mediator (M)** is caused by the independent variable which in turn causes the dependent variable. A mediator is **in the causal pathway**.

A **moderator (Z)** modifies the effect of an independent variable on a dependent variable. The association varies depending on the values/levels of Z .

Question 2

- Which of the next examples is a case of modification?
- If a man and a woman have the same amount of water per week, the effect on their weight is different.

Question 3

- To test moderation...
- A new variable needs to be considered. This new term is the cross-product between X1 and the modifier Z . It is denoted as $X1 \times Z$

Question 4

- Given the following model: $Y=B_0+B_1 X_1+B_2 Z +B_3 X_1 \times Z + E$
- B_1 and B_2 are no longer useful unless zero values of the respective predictors are of particular interest.

Question 5

Given the model: $\text{years_married} = 6 - 2.4 \text{problems_inlaw_family} - 3.1 \text{cheating_with_others} - 3.7 \text{problems_inlaw_family} * \text{cheating_with_others}$

The effect of $\text{problems_inlaw_family}$ on years_married is $= \beta_1 + (\beta_3 * \text{cheating_with_others}) = -2.4 + (-3.7 * \text{cheating_with_others})$

Question 6

The effect of motivation depended on whether there was a clear standard for excellence.

Indicate whether this sentence implies that there was a main effect of motivation, a main effect of standard, or an interaction effect, and give a reason why.