5. Data Management

Video Link:

https://www.youtube.com/watch?v=9KxSZEWodR0&index=5&list=PL2fQHGEDK7Yyl1W9tgIo8wp YFTDumgc_j

Section 5.1: Code Out Missing Data

Section 5.2: Code in Valid Data

Section 5.3: Logical Response Codes

Section 5.4: Create Secondary Variable from 2 Variables

Section 5.5: Create Secondary Variable from more than 2 Variables

Section 5.6: Collapsing a Variable

Section 5.1: Code Out Missing Data

In this next step, we'll be making decisions about recoding data in ways that will help us answer our question. This process is also known as **data management**. A necessary first step to consider is whether or not you need to code out missing data. Some variables include values that represent uncertain or unknown information. Refer to your personal code book to determine if any of your variables have values that represent unknowns. In order to perform analyses on these variables, we have to set those values to missing. In the following frequency table, the value 9 represents an answer of unknown or not applicable.

-		 							
	1.	E٦	ver	y (day				
	2.	5	to	6	Day	(s)	а	week	
	3.	3	to	4	Day	(s)	a	week	
	4.	1	to	2	Day	(s)	а	week	
	5.	2	to	3	Day	(s)	а	mont	h
	6.	Or	nce	а	mon	th	or	less	
	9.	Ur	nkn	OWI	l				
E	BL.	NZ	A, 1	nev	ver	or	un}	known	j

	Usual Smoking Frequency										
					Cumulative						
		Frequency	Percent	Valid Percent	Percent						
Valid	1	1320	77.4	77.4	77.4						
	2	68	4.0	4.0	81.4						
	3	91	5.3	5.3	86.7						
	4	88	5.2	5.2	91.9						
	5	65	3.8	3.8	95.7						
	6	71	4.2	4.2	99.8						
	9	3	.2	.2	100.0						
	Total	1706	100.0	100.0							

Total1706100.0100.01.To set 9 to missing, click Variable View in the lower left hand corner. You will see the name of
each variable you selected in the first column. Scroll down to the row corresponding to the
variable you need to treat for missing data. Move your cursor to the seventh column titled
Missing. Click in the box with the word None, then click the box with the three dots just to the

right of None .

ta *My	Sorted	Nesarc Data.s	av [DataSet1] - IB	M SPSS Stat	istics Data Edito	r						
<u>File</u>	dit	View Data	<u>T</u> ransform	Analyze	<u>G</u> raphs <u>U</u> ti	lities Add- <u>o</u> ns	Window He	elp				
	H		J 🗠 🦳	¥) * - =	P M		- S			ABG	
		Name	Туре	Width	Decimals	Label	Values	Missing	Columns	Align	Measure	Role
3		IDNUM	Numeric	3	0		None	None	8	🗃 Right	🛷 Scale	🔪 Input
31		AGE	Numeric	2	0		None	None	8	■ Right	🛷 Scale	🔪 Input
235		S2AQ8A	Numeric	2	0		None	None	8	遭 Right	\delta Nominal	🔪 Input
492		SMOKER	Numeric	1	0		None	None	8	🚎 Right	\delta Nominal	🔪 Input
495		CHECK321	Numeric	1	0		None	None	8	疆 Right	🚓 Nominal	🔪 Input
496		S3AQ3B1	Numeric	1	0		None	9	8	I Right	🚓 Nominal	🔪 Input
497		S3AQ3C1	Numeric	2	0		None	99	8	3 Right	I Scale	🔪 Input
		4										
Data Vi	ew 🔪	/ariable View										

2. Click **Discrete** to list one value per box that you want to set to missing for the corresponding variable. To set multiple values that are in ascending order click **Range** plus one optional discrete missing value, then enter the Low: and High: numerical values you want to set as missing for that variable. Click **OK**.

ta Missing Values	ta Missing Values
© <u>N</u> o missing values	◎ <u>N</u> o missing values
Discrete missing values	Discrete missing values
9.000	
© Range plus one optional discrete missing value	Range plus one optional discrete missing value
Low: <u>H</u> igh:	Low: 7 High: 9
Di <u>s</u> crete value:	Di <u>s</u> crete value:
OK Cancel Help	OK Cancel Help

If you run a frequency table on the variables you coded for missing you will see that at the end of the table the value(s) have been set to missing.

					Cumulative
		Frequency	Percent	Valid Percent	Percent
Valid	1	1320	77.4	77.5	77.5
	2	68	4.0	4.0	81.5
	3	91	5.3	5.3	86.8
	4	88	5.2	5.2	92.0
	5	65	3.8	3.8	95.8
	6	71	4.2	4.2	100.0
	Total	1703	99.8	100.0	
Missing	9	3	.2		
Total		1706	100.0		

Usual Smoking Frequency

Return to Video - 05. Data Management

Section 5.2: Code in Valid Data

To avoid asking inappropriate questions, skip patterns are often created in surveys that allow participants to skip questions in which the answer can be logically determined. In this way, missing data on some questions might mean that we can reasonably recover valid information.

For the question, "Did you drink at least one alcoholic drink in the past 12 months?", 16,116 participants said "No". These individuals would not need to be asked the question about how often

they drank alcohol in the past 12 months. Instead, they will be set to blank, or in the case of a SPSS dataset, a dot or period. So, for the variable, how often did you drink alcohol in the past 12 months, it would be reasonable to code this as valid data rather than missing. In the variable below we see that **Missing System** frequency is 180. We can set those to valid responses.

					Cumulative
		Frequency	Percent	Valid Percent	Percent
Valid	1	76	4.5	5.0	5.0
	2	84	4.9	5.5	10.5
	3	194	11.4	12.7	23.2
	4	229	13.4	15.0	38.2
	5	216	12.7	14.2	52.4
	6	248	14.5	16.3	68.6
	7	134	7.9	8.8	77.4
	8	85	5.0	5.6	83.0
	9	134	7.9	8.8	91.7
	10	118	6.9	7.7	99.5
	99	8	.5	.5	100.0
	Total	1526	89.4	100.0	
Missing	System	180	10.6		
Total		1706	100.0		

Frequency of Drinking past 12 months

1. Click **Transform > Recode into Same Variables**.

ta *I	Nesarc Da	ata_1.sav	[DataSet	2] - IBM SPSS	Statistics Da	ta Editor					
<u>F</u> ile	Edit	View	<u>D</u> ata	<u>T</u> ransform	<u>A</u> nalyze	<u>G</u> raphs	<u>U</u> tilities	Add- <u>o</u>	ns	<u>W</u> indow	<u>H</u> elp
10 :				☐ <u>C</u> omput Program Count V:	e Variable Imability Tr alues withir	ansformatic 1 Cases)n			*	
		IDN	UM	Shift Val	ues				R	CHECK321	S3A0
	1			Recode	into Same	Variables					
	2			Recode	into <u>o</u> ame	valiables					
	3			Recode	Into Differe	nt variables	·		2		
	4			Automat <u>A</u> utomat	ic Recode						
	5			Create E)ummy Var	iables					1
	6			Visual B	inning				2		2
	7		-	Rank Ca	ises						-
	8			🖨 Date and	1 Time Wiz	ard					-
	9		-			aru					-
	10			Create I	i <u>m</u> e Series						
	11			Replace	Missing <u>V</u> a	alues					
	12			🍘 Random	Number G	enerators					
	13			Run Per	iding <u>T</u> rans	forms	Ctrl	+G			
	14			0	1	6.7			2		

2. Select the variable you want to recode from the list on the left, move it to the **Numeric Variables** column on the right using the arrow. Click **Old and New Values...**

ta Recode into Same Variable	25				
AGE CHECK321 CHECK321 SAQ3B1	Numeric <u>V</u> ariables:				
S3AQ3C1	Old and New Values	ion condition)			
OK Paste Reset Cancel Help					

3. In the **Old Value** column on the left hand side click **System-missing** since currently the value is treated as missing indicated by a "." in **Data View**. In the **New Value** column on the right hand side click **Value** then enter the numerical value (i.e., dummy code) you want it to now be. Refer to your personal code book to determine for each of your specific variables what would be an appropriate dummy code to use for each variable you intend to code in valid data. Click **Add** > **Continue.**

🔚 Recode into Same Variables: Old and New Values	
Old Value © <u>V</u> alue:	New Value Value: System-missing
 System-missing System- or <u>u</u>ser-missing Range: through Range, LOWEST through value: Range, value through HIGHEST: 	Ol <u>d</u> > New: SYSMIS> 11 Add Change Remove
◎ All <u>o</u> ther values	
Continue	Cancel Help

If you need to do this for more than one variable, move the variable currently listed under **Numeric Variables** you just 'coded in' as valid back to the left column using the arrow. Click on the name of the next variable and move it to the right column. Click **Old and New Values...** <u>It is extremely important to click in the **Old --> New:** column in the lower right hand side, click on what you previously added, and then click **Remove**. It is unlikely that you would code in valid data on a different variable and use the same dummy code. Failure to do this will result in the variable you are working with now to be coded for valid data how you asked SPSS to code the previous variable. Repeat step #3 above.</u>

Run a frequency table on the variable(s) on which you 'coded in' valid data to confirm that the new valid responses (i.e., dummy code of 11) shows in the table.

-		Frequency	Percent	Valid Percent	Cumulative Percent
Valid	1	76	4.5	4.5	4.5
	2	84	4.9	4.9	9.4
	3	194	11.4	11.4	20.8
l	4	229	13.4	13.5	34.3

Frequency Drank Alcohol Last 12 Months

	5	216	12.7	12.7	47.1
	6	248	14.5	14.6	61.7
	7	134	7.9	7.9	69.6
	8	85	5.0	5.0	74.6
	9	134	7.9	7.9	82.4
	10	118	6.9	6.9	89.4
	11	180	10.6	10.6	100.0
	Total	1698	99.5	100.0	
Missing	99	8	.5		
Total		1706	100.0		

Return to Video - 05. Data Management

Section 5.3: Logical Response Codes

Another very useful step in data management is to give your variables response codes that may be more logical than those they were originally given. For the variable, usual smoking frequency, you can see from the code book entry that lower values mean that the respondents smoked more and higher values mean that the respondents smoked less. Perhaps this seems counterintuitive to you. We could chose to reverse code this variable to be more intuitive so that higher values mean more smoking and lower values mean less smoking.

1. Click Transform > Recode into Different Variables.



2. Choose the variable from the list in the far left column and move it to the middle column by clicking the arrow. In the far right column create a new variable name. This should be unique to the dataset, have no spaces, and be 12 characters or less. Directly below that, label this new variable. It can be the same label as the variable you are coding to be more logical or something new that appropriately describes the variable. Click Change > Old and New Values...

ta Recode into Different Vari	ibles	×
 ✓ AGE ← CHECK321 ✓ IDNUM ← S2AQ3 ← S3AQ3C1 ← SMOKER ← TABP12MDX 	Numeric Variable -> Output Variable: S3AQ3B1> USFREQ Image: Comparison of the second seco	Output Variable <u>Name:</u> USFREQ Label: Usual Frequency When Sr Change
	OK Paste Reset Cancel Help	

3. Refer to your personal code book to verify what the current dummy codes are for the variable you are recoding to be more logical and decide what the new dummy codes will be. Typically you will be reversing the current dummy codes so that lower dummy codes represent categories of lower value and higher dummy codes for higher value. In the **Old Value** column on the left side input the current dummy code under **Value**. In the **New Value** column on the ride side input the new dummy code that is more logical next to **Value**. Click **Add**. Repeat this for each dummy code for the variable you are recoding. Once you are done with this verify that in the lower right column under **Old-->New** you see the original dummy codes going down in ascending order and the new dummy codes that are more logical going down in descending order. Click **Continue**.

Old Value	New Value
© Value:	Value:
	© System-missing
◎ <u>S</u> ystem-missing	© Co <u>p</u> y old value(s)
System- or <u>u</u> ser-missing Range: through Range, LOWEST through value:	Add 1-> 6 2-> 5 3-> 4 3-> 4 4-> 3 5-> 2 6-> 1
© Range, value through HIGHEST:	Output variables are strings Width: 8
O All other values	Convert numeric strings to numbers ('5'->5)

4. If you have more than one variable to recode click the variable name of old and new variable in the middle column and move it back to the left column using the arrow. Repeat step #2 through #3. Make sure to remove or change the Old-->New values in the lower right column otherwise the system will recode the current variable you are data managing the same way it did your previous variable.

ta Recode into Different Vari	ibles	×
 ✓ AGE ✓ CHECK321 ✓ IDNUM ✓ S2AQ3 ✓ S2AQ8A ✓ S3AQ3C1 ✓ SMOKER ✓ TABP12MDX 	Numeric Variable -> Output Variable: S3AQ3B1> USFREQ Image: Contract of the second	Output Variable <u>Name:</u> USFREQ <u>Label:</u> Usual Frequency When Sr Change
	If (optional case selection condition) OK Paste Reset Cancel Help	

5. Once you have recoded to make more logical response categories for all of the variables you decided to recode, Click **OK**. Then at the top menu click **Utilities > Define Variable Sets**.

ta -	My Sorter	d Nesarc Data.sav	[DataSet2] - IB	M SPSS Statist	ics Data b	ditor				
Eile	Edit	<u>View</u> <u>D</u> ata	Transform	Analyze (Braphs	Utilities	Add-ons	Window	He	lp
6				× .	*		Ibres Control Pan	el		
39. 		IDNUM	AGE	52AQ3	S2A0	30	in a Winned			S3AQ3B1
	1-	1	23	2		E Scor	ing wizard			
	2	2	28	1		4 Marg	e Model <u>X</u> ML			
-	3	3	81	2		Calc	ulate with Piv	ct Table		
	4	4	18	2	1	🛃 Data	File Comme	nts		-
-	5-	5	36	2		Marg	e Viewer Tab	les		
	5-	6	34	1		V Defin	e Variable S	ets		
	7	7	19	1		Can	sor Table			4
	8	8	84	1		Colles	Variabia Detr			
	y	9	29	1		€ Use	variable Sets			
فسر	10	10	16	1		Show	v Ali variable	S		+
-	11-	11	68			Proc	ess Data File	IS		1
	12	12	48	1		🐀 Spei	ling			1
-	13-	13	31	1		D Run	Script			
	14	14	55	1		Prod	uction Eacilit			1
	15	15	54	1		Man	Conversion	RUEZ		1
-	16	16	51	2		Map	Conversion (zanty		
	17-	17	38	,		Cus	tom Dialogs		0	5
	18	18	40	2		Exto	nsion Bundle	16	,	1
-	19	19	54	1		10	1		1	1
-	20	20	21	1		5	3		4	(2)
1	21	21	25	1		9	1		-	1

6. Click on **My Variables**, or whatever you previously named your Variable Set. Move the newly created variable(s) from the lower left column to the lower right column using the arrow. Click **Change Set** toward the top left. Click **Close**.

ta Define Variab	le Sets	×
Add Set Change Set Remove Set	Set <u>N</u> ame: My Variables <mark>My Variables</mark>	
TBTYPE2 TBTYPE3 TBTYPE4 TBTYPE5 TRAN12AB TRANP12/ USFREQ WEIGHT WINEECF	ABDEP ABDEP Close Help	

7. Click **Utilities > Use Variable Sets**.

🙀 "My Sorter	d Nesarc Data.sav	DataSet2] - JB	M SPSS Statisti	cs Data E	ditor			
Eile Edit	View Data	Transform	Analyze G	raphs	Utilities	Add-ons	Window	Help
2		-	¥ 🏋		₩ Уана Фомз	ibies Control Pan	el	
1:USFREQ					S ONO	Identifiera		
	IDNUM	AGE	52AQ3	SZAC	- Door	ing Wizard		53AQ38
1	1	23	2		T Mare	Madal VM		
-2-	2	28	1		- marg	te woder Zwit		
_3-	3	81	2		Calc	ulate with Piv	ct Table	
-4	4	18	2		📝 Data	File Comme	nts	
-5-	5	36	2		Merg	e Viewer Tat	les	
-8-	6	34	1		Z Defi	ne Variable S	ets	
7-	7	19	1		Can	e or Table		
-8	8	84	1		Chillen	Mariable Date	2	
	9	29	1		₩ <u>U</u> se	variable Set	3,	
10-	10	16	1		Sho:	w All Variable	s	
_11	11	68	1		Proc	ess Data File	es.	
12	12	48	1		1 Spe	ling		
13-	13	31	1		Run	Script		
_14	14	55	1		22 Drod	uction Enciet		
15-	15	54	1		mo rive	ouron Facilit	F=+	
16-	16	51	2		Мар	Conversion (Junty	
17	17	38	1		Cus	tom <u>D</u> ialogs		1
-18	18	40	2		Extension Bundles			•
_19	19	54	1		10	1		1

8. Confirm that My Variables or whatever you named your Variable Set is still clicked and the others listed are <u>un</u>checked. Click **OK**.

ta Use Variable Sets 🛛 💌
Select variable sets to apply
My variables
Check <u>A</u> ll Uncheck All
Only variables in the selected sets will appear in the Data Editor and in the dialogs.
OK Cancel Help

You will see your newly created variables in Data View and Variable View. **Now, rather than using the original variable, S3AQ3B1 in our analysis, we will use this new variable USFREQ.** It is important to add these details to the code book so there is good documentation of these changes as you move forward.

- There is actually another option that may be even better for giving this variable values that are maximally informative. In a way, dummy codes 1 through 6 are unnecessarily categorical given that how often you smoke is generally a quantitative measurement. And it may be reasonable to recode the variable in order to capture more of its quantitative features.
- To do this, we choose values that reasonably correspond to the number of times each individual smokes in a usual month. So, someone who reports smoking everyday could be said to smoke cigarettes on 30 days in a usual month. Someone who smokes 5 to 6 days a week could be said to smoke 22 days in a usual month. Someone who smokes 3 to 4 days a week could be said to smoke 14 days in a usual month.1 to 2 days a week, 5 days in a usual month. 2 to 3 days a month, 2.5 days in a usual month and once a month or less, 1.
- Although these are estimates, they capture the quantitative nature of the measure and also keep individuals ordered in terms of the frequency with which they smoke. Our new variable is called USFREQMO, which stands for number of days smoked in a usual month.

9. Repeat steps #1 through #8.

ta Recode into Different Vari	bles		— ×
 ✓ AGE ✓ CHECK321 ✓ IDNUM ✓ S2AQ3 ✓ S3AQ3C1 ✓ SMOKER ✓ TABP12MDX ✓ test ✓ USFREQ 	Numeric S3AQ3B	<u>V</u> ariable -> Output Variable: 1> USFREQMO New Values Dtional case selection condition)	Output Variable <u>Name:</u> USFREQMO <u>Label:</u> Number of Days Smoked i Change
	OK Pas	te <u>R</u> eset Cancel Help	

	-New Volue
Old Value	
<u>V</u> alue:	Value:
	© System-missing
◎ <u>S</u> ystem-missing	© Copy old value(s)
◎ System- or <u>u</u> ser-missing	
O Range:	Ol <u>d</u> > New:
	1> 30
through	2->22
inough	3> 14
	Change 4>5
O Range, LOWEST through value:	$\mathbb{R}^{\text{emove}}$ $\begin{array}{c} 5 \longrightarrow 2.5 \\ 6 \longrightarrow 1 \end{array}$
Range, value through HIGHEST:	
	Output variables are strings Width: 8
◎ All <u>o</u> ther values	Convert numeric strings to numbers ('5'->5)
Cont	inue Cancel Help

- So, this even makes more sense, doesn't it? What we've basically done is to spread the measurement of smoking frequency out in a meaningful way. So, that while it is a categorical variable, we are actually getting more information out of it than we had originally been given.
- When we run a frequency table, the table now describes the sample of 1,706 young adult past year smokers, somewhat more clearly, by showing the numbers and percentages of individuals who smoked these approximate number of days in a usual month. You can see that the vast majority of the sample, 1,320 individuals or 77.4% smoked all 30 days.

			-		Cumulative
		Frequency	Percent	Valid Percent	Percent
Valid	1.00	71	4.2	4.2	4.2
	2.50	65	3.8	3.8	8.0
	5.00	88	5.2	5.2	13.2
	14.00	91	5.3	5.3	18.5
	22.00	68	4.0	4.0	22.5
	30.00	1320	77.4	77.5	100.0
	Total	1703	99.8	100.0	
Missing	System	3	.2		
Total		1706	100.0		

Number of Days Smoked in a Usual Month

Return to Video - 05. Data Management

Section 5.4: Create Secondary Variable from 2 Variables

Secondary variables are variables that include information from two or more primary variables. We can create secondary variables by using a mathematical or logical operation on two or more variables.

In this case we want to know the number of cigarettes smoked per month. We know the number of usual number of cigarettes smoked on days that an individual smokes (quantity), S3AQ3C1, and our new variable, USFREQMO, give us an estimate of the number of days smoked in a usual month. If we want to estimate the total number of cigarettes that participants smoked per month, it would make sense to multiply those two variables and get a product that represents the number of cigarettes per day times the number of days smoked per month.

USFREQMO

Usual smoking days per month

1 = Once a month or less 2.5 = 2 - 3 Day(s) a month 6 = 1 - 2 Day(s) a week (1.5 x 4 wk) 14 = 3 - 4 Day(s) a week (3.5 x 4 wk) 22 = 5 - 6 Day(s) a week (5.5 x 4 wk) 30 = Every day

S3AQ3C1

Usual quantity when smoked cigarettes

1-98 usual cigarettes/day

1. To do this, click **Transform > Compute Variable**.

🖬 м	🔒 My Sorted Nesarc Data.sav [DataSet2] - IBM SPSS Statistics Data Editor										
File	Edit	View	<u>D</u> ata	Transform	<u>A</u> nalyze	Graphs	<u>U</u> tilities	Add-o	ns	Window	<u>H</u> elp
2				Compute Program	Variable	ansformati	on		ł	*5	
14 :				Count Vo	luoc within	Casas					
		IDNU	JM	Chift Value	indes within	04565			R	CHECK321	S3AQ3B1
	1		1	Shiit Valu	ies				3		
	2		2	🔤 Recode i	nto <u>S</u> ame \	/ariables			3		
	3		3	🔤 <u>R</u> ecode i	nto Differer	nt Variable	S		3		
	4		4	🛐 <u>A</u> utomati	c Recode				3		
	5	1	5	Create D	ummy Vari	ables			3		
	6		6	Visual <u>B</u> i	nning				3		
·	7		7	Rank Ca	ses				3		
	8		8	Date and	Time Wiza	rd			3	-	
	99		9		i ma Oariaa	iru			3	-	
	0		10		i <u>m</u> e Series.				3		
	4		11	Replace	Missing <u>V</u> a	lues			1	1	1
	2		12	🍘 Random	Number <u>G</u>	enerators.			2	2	1
	3	ļ	13	Run Pen	ding <u>T</u> rans	forms	Ctrl	+G	3		
1	4		14	55		1	4		1	1	1
1	5		15	54		1	2		1	1	1
	6		16	51		2	11		3		.

2. Enter the name of the new variable under **Target Variable** on the far upper left hand corner. For our example, let's call this new variable NUMCIGMO_EST, which stands for the estimate of the number of cigarettes smoked per month. Click **Type & Label...** directly below the variable name you just typed in.

ta Compute Variable		
Compute Variable	-	Numeric Expression: + < > 7 8 9 - <= >= 4 5 6 * = ~= 1 2 3 / & 1 0 ** ~ () Delete
(optional case selection	on condi	ion)
		OK Paste Reset Cancel Help

3. Type in an appropriate label for the secondary variable you are creating. Click **Continue**.

🔄 Compute Variable: Type and 📧						
Label						
Label: s smoked per month						
© <u>U</u> se expression as label						
Туре						
<u> </u>						
© <u>S</u> tring <u>W</u> idth: 8						
Continue Cancel Help						

In the lower left hand column find the name of the first variable you are using to create this secondary variable. Use the arrow to move it to the upper right column titled Numeric Expression. Since we are multiplying the two variables together we will insert an asterisk "*" then the second variable name. Then click OK.

Target Variable: NUMCIGMO_EST Type & Label V AGE CHECK321 DNUM S2A03 S2A03 S3A03C1 S3A03C1 S3A03C1 SA0CKER TABP12MDX USFREQ USFREQ USFREQ () Delete Functions and Special Variables: (ptional case selection condition)	tariable Compute Variable		X
(optional case selection condition)	Compute Variable Target Variable NUMCIGMO_EST Type & Label AGE CHECK321 DINUM SAQA3 SAQ381 SAQ381 SAQ361 SMOKER TABP12MDX USFREQ VISFREQ VISFREQMO	= USFREQMO * S3AQ3C1	Function group: All Arithmetic CDF & Noncentral CDF Conversion Current Date/Time Date Arithmetic Date Creation
OK Paste Reset Cancel Help	(optional case select	OK Paste Reset Cancel Help	

- To create your secondary variable you may need to add "+", subtract "-", multiple "*", divide "/" or a series of those steps.
 - 5. Repeat steps 5.3 #5 through #8 to make sure that your secondary variable you just created to show up in Data View and Variable View.
 - 6. If you need to create another secondary variable repeat steps #1 through #5 making sure to change the **Target Variable** name, **Label**, and **Numeric Expression**.

If we generate a frequency distribution for this new variable, NUMCIGMO_EST, you see it is a quantitative variable that ranges from 1 to 2,940 with 9 missing observations. So, how can we check to make sure that this new secondary variable was created as we intended?

7. Click **Analyze > Reports > Case Summaries**.

🐚 *M	😭 "My Sorted Nesarc Data.sav (DataSei2) – IBM SPSS Statistics Uata Editor													
File	Edit	View	Data	Transform	Analyze	Graphs	Utilities	Add-on	15	Window	Help			
6					Repo Desc	ts iptive Stat	istics	•		Codebook			A	0
3 : NUN	CIGMO	_ESI			Ta <u>b</u> le	s		•	5					
		IDNU	JM	AGE	Comp	are Mean	s	,		Case Sun	intanes		1 TA	BP12MDX
1			1	23	Gene	al Linear	Model	•		Eeport Su	mmaries	in Rows	1.1	0
2			2	28	Gene	ralized Lin	ear Models	•		Report Su	mmaries	in Columns	1.1	0
3			3	81	Mixed	Models		•	3					0
4			4	18	Corre	late		,	3					0
&			5	36	Regie	ssion		,	3					0
6			6	34	Logir	car		•	3					0
7			7	19	Class	ity		,	3					0
8			8	84	Dime	nsion Red	luction	,	3		_			0
9			9	29	Scale			,	3					0
16	1		10	18	Nonp	arametric	Tests	,	3					0
11	-		11	68	Ecrec	asting		,	1		1	1	20	1
			12	48	Suniv	al		,	2		2	1	5	0
-12			13	31	Multin		100	,	β					0
14	-		14	55	- Dimut	ne neopor	196		1		1	1	20	0
-16			15	54	Ginus Ouell				1		1	1	20	1
-10	÷		16	51	Qualit	y Control		,	3		1			0
17			17	38	ROCO	Curve			2		2	5	2	0
18	-		18	40		2	11		1		1	1	20	1
16			19	54		1	10		1		1	1	10	0
20			20	21		1	5		3					0

8. Remembering the order you put your variables into the **Numeric Expression** when you created this secondary variable, move the variables along with the newly created secondary variable from the left column to the upper right column. In the lower left corner <u>un</u>click **Limit cases to first** and **Show only valid cases** then click **OK**.

ta Summarize Cases		×				
 ✓ AGE CHECK321 ✓ IDNUM S2AQ3 S2AQ8A S3AQ3B1 SMOKER TABP12MDX USFREQ 	Variables: ✓ USFREQMO ✓ S3AQ3C1 ✓ NUMCIGMO_EST Grouping Variable(s):	Options				
 ✓ Display cases ☐ Limit cases to first 100 ☐ Show only valid cases ☐ Show <u>c</u>ase numbers OK <u>Paste</u> <u>Reset</u> Cancel Help 						

This will create a table in your output that will list every single participant id number in the first column, the values for each of the participant's responses for the variables you used to create the secondary variable and the secondary variable itself. (Note the middle of the table was

omitted to save space and we show only showing the beginning and the end.) The rows represent individual observations similar to the dataset itself. And the columns show the values for the specific variables. Looking over these values for a handful of individuals, we can see that, in fact, our new secondary variable does indeed show the value of number of days smoked per month times number of cigarettes smoked per day. Remember that whenever you are conducting data management, it is important to find a way to check for errors at each step of the process.

	Case Sum	maries					
	Number of Days Smoked in a Usual Month	Usual Quantity When Smoked Cigarettes	number of cigarettes smoked per month	1689 1690	5.00 30.00	4 20	20.00 600.00
1	30.00	3	90.00	1691	30.00	40	1200.00
2	22.00	3	66.00	1692	30.00	10	300.00
3	30.00	10	300.00	1693	30.00	10	300.00
4	30.00	10	300.00	1694	30.00	20	600.00
5	30.00	20	600.00	1695	30.00	10	300.00
6	30.00	5	150.00	1696	30.00	20	600.00
7	30.00	8	240.00	1697	30.00	20	600.00
8	30.00	1	30.00	1698	20.00	20	00.000
9	30.00	10	300.00	1600	30.00	20	000.00
10	30.00	20	600.00	1699	30.00	10	300.00
11	5.00	2	10.00	1700	5.00	2	10.00
12	30.00	3	90.00	1701	30.00	40	1200.00
13	14.00	5	70.00	1702	30.00	5	150.00
14	30.00	1	30.00	1703	1.00	3	3.00
15	30.00	98	2940.00	Total N	1703	1697	1697
16	30.00	20	600.00		 1705	1007	1007

Return to Video - 05. Data Management

Section 5.5: Create Secondary Variable from more than 2 Variables

So what if you want to combine more than 2 variables? A good example of this would be creating a single secondary variable to characterize race/ethnicity from a number of separate primary variables available in the NESARC dataset. Race or ethnicity are measured by a series of questions, coded 1 if "Yes" and 2 if "No". Participants in this sample could have indicated more than one race or ethnicity, we could decide to characterize those participants who indicate multiple racial or ethnic groups separately from those who could be characterized with a single ethnicity. Below you can see the 6 variables we will need to use to create one secondary variable to represent race/ethnicity.

S1Q1C	8308 34785	HISPANIC OR LATINO ORIGIN 1. Yes 2. No
SIQIDI	1304 41789	"AMERICAN INDIAN OR ALASKA NATIVE" CHECKED IN MULTIRACE CODE 1. Yes 2. No
S1Q1D2		"ASIAN" CHECKED IN MULTIRACE CODE
	1334 41759	1. Yes 2. No
S1Q1D3		"BLACK OR AFRICAN AMERICAN" CHECKED IN MULTIRACE CODE
	8600 34493	1. Yes 2. No
S1Q1D4		"NATIVE HAWAIIAN OR OTHER PACIFIC ISLANDER" CHECKED IN MULTIRACE CODE
	363 42730	1. Yes 2. No
		·
S1Q1D5		"WHITE" CHECKED IN MULTIRACE CODE
	32789 10304	1. Yes 2. No

- 1. Add the 6 variables listed above to your Variable Set by repeating steps 5.3 #5 through #8.
- 2. For each of these 6 variables the response of "No" is dummy coded as a 2, but we need it to be dummy coded as a 0 so that when we sum the dummy codes for the 6 variables for each participant a sum of 2 or more would represent the participant marking "Yes" on 2 or more of the ethnicity questions. To recode these variables repeat steps 5.2 #1 through #3. Since you are doing the same recode to each of the 6 variables you can add all the variable names under Numeric Variables at one time then just once put in the Old Value for each of the 6 variables put a 2 and in the New Value put a 0. This will make the response of "No" dummy coded as a 0.
- 3. Next we need to sum the variables to get a new variable, which we'll call NUMETHNIC that indicates the number of race or ethnicity variables that were endorsed. Repeat steps 5.4 #1 through #4, but instead of multiplying you will add the 6 variables listed above using a "+" sign between each variable name.

Compute Variable Target Variable: NUMETHNIC Type & Label	Numgric Expression: = \$101C + \$101D1 + \$101D2 + \$101D3 + \$101D4 +	S101D5	
 ✓ AGE ✓ CHECK321 ✓ ETHRACE2A ✓ IDNUM ✓ NUMCIGMO_EST ✓ \$101D1 ✓ \$101D2 ✓ \$101D3 ✓ \$101D3 ✓ \$101D5 ✓ \$2AQ3 ✓ \$2AQ3A ✓ \$3AQ3C1 ✓ \$MOKER ✓ TABP12MDX ✓ USFREQ ✓ USFREQ ✓ USFREQ ✓ USFREQ ✓ (optional case selector) 	• < > 7 8 9 • < > 7 8 9 • < = >= 4 5 6 • = ~= 1 2 3 7 8 1 0 . • < ~ () Delete •	Function group:	Compute Variable: Type and Label Label: s that were endorsed Label: s that were endorsed Use expression as label Type Numeric String Width: 8 Continue Cancel Help

- 4. Add the NUMETHNIC variable to your Variable Set by repeating steps 5.3 #5 through #8. This could seem tedious, but throughout the rest of the semester this will save <u>you</u> a lot of time once we are analyzing the variables by eliminating searching through thousands of variables looking for your specific variable names.
- 5. Repeat steps 5.4 #1 through #3 naming the new variable ETHNICITY and properly labeling it.
- 6. Remove the variables currently listed under **Numeric Expression**. Move your NUMETHNIC variable from the left hand column to the **Numeric Expression** column in the upper right corner. Start with dummy code 3 by typing it into the **Numeric Expression**. Then click **If...** in the lower left hand corner.

tai Compute Variable		
Target Variable:		Num <u>e</u> ric Expression:
Type & Label	=	J
AGE	•	

7. Click **Include if case satisfies condition:** Move the variable SC1Q1C to the upper right column from the left column using the arrow then type "= 1" click **Continue**.



8. Repeat steps #6 through #7 five more times, once for each of the 5 remaining variables used to create NUMETHNIC. Note it is important to start with dummy code 3, then go to 4, 5, 6, 7 then do step 9 for dummy codes 1 and 2 to work around SPSS order of operation. You will increase the dummy code value you enter in **Numeric Expression** by one with each repeat of step 5. If you refer back to the code book you will see that for each of the 6 variables we used to create NUMETHNIC a dummy code of 1 represents "Yes". We will need to tell SPSS for each of those variables dummy code of 1 to compute a different dummy code going in ascending order for the secondary variable ETHNICITY we are creating. We will start with 2 and end with 7 as 1 is already representing multiracial.

The second time you complete step #7 and #8 the below pop up window will appear each time you compute a new dummy code level for the Ethnicity variable. Click **OK** each time.

Target Variable Numeric Expression: ETHNICITY = Type & Label 4	Include all cases Include all cases Include if case satisfies condition: Include if case satisfies condition: Include if case satisfies condition:
CHECK321	VUMCIGMO_ESI
Target Variable: Numeric Expression: ETHNICITY = Type & Label 5 V AGE Image: CHECK321	Image: Compute Variable: If Cases Image: Check321 Image: Check321

ta Compute Variable	
Target Variable: Numeric Expression: ETHNICITY = Type & Label 6 ✓ AGE ✓ ✓ CHECK321 ✓	Compute Variable: If Cases AGE CHECK321 CHECK32 CHECK321 CHECK32 CHECK
Target Variable: Numeric Expression: ETHNICITY = Type & Label 7 Image: AGE Image: AGE Image: CHECK321 Image: AGE	Image: Compute Variable: If Cases Image: Check321

9. Repeat steps #6 & #7 using the dummy code of 1 for Numeric Expression. Click Include if case satisfies condition: type "NUMETHNIC >1" which tells SPSS that if the sum of the 6 variables we used to create NUMETHNIC is 2 or more than the new secondary variable I am creating called ETHNICITY should be dummy coded as a 1. This will represent a category of multiple racial or ethnic groups endorsed. Click Continue > OK>OK.

Compute Variable	ta Compute Variable: If Cases
Target Variable: Num <u>eric Expression:</u> ETHNICITY = Type & Label 1 Ø AGE Image: CHECK321 G CHECK321 Image: CHECK321 Function Function	AGE CHECK321 C

10. Repeat step 8 for the last variable used to create NUMETHNIC.

tai Compute Variable	ta Compute Variable: If Cases
Target Variable: Num <u>e</u> ric Expression: ETHNICITY = Type & Label 2 ✓ AGE ✓ ✓ CHECK321 ✓	 ✓ AGE ◇ CHECK321 ◇ ETHRACE2A ✓ IDNUM ✓ NUMCIGMO_EST ◇ NUMETHNIC ○ Include all cases ③ Include if case satisfies condition: S1Q1C = 1

11. Add the ETHNICITY variable you just created to your Variable Set by repeating steps 5.3 #5 through #8.

12. Repeat steps 5.3 #7 through #8 using the 6 variables you used to create NUMETHNIC, as well as NUMETHNIC and ETHNICITY to check that how you intended the variable to be created is what actually happened. A portion of the table is below.

Case Summaries								
	S1Q1C	S1Q1D1	S1Q1D2	S1Q1D3	S1Q1D4	S1Q1D5	Sum of Endorsed Ethnicities	Ethnicity of Participant
1	0	0	0	1	0	0	1.00	5.00
2	1	0	0	0	0	1	2.00	1.00
3	0	0	0	0	0	1	1.00	7.00
4	0	0	0	0	0	1	1.00	7.00
5	0	0	0	0	0	1	1.00	7.00
6	0	0	0	0	0	1	1.00	7.00
7	0	0	0	1	0	0	1.00	5.00
8	0	0	0	0	0	1	1.00	7.00
9	1	0	0	0	0	1	2.00	1.00
10	0	0	0	0	0	1	1.00	7.00
11	0	0	1	0	0	0	1.00	4.00
12	0	0	0	1	0	0	1.00	5.00
13	0	0	0	0	0	1	1.00	7.00
14	1	0	0	0	0	1	2.00	1.00
15	0	0	0	1	0	0	1.00	5.00
16	1	0	0	0	0	1	2.00	1.00
17	0	0	0	0	0	1	1.00	7.00
18	1	0	0	0	0	1	2.00	1.00
19	1	0	0	0	0	1	2.00	1.00
20	0	0	0	0	0	1	1.00	7.00
21	0	0	0	0	0	1	1.00	7.00
22	1	0	0	0	0	1	2.00	1.00

You can see in the table below we now have one variable that represents all ethnicities.

			_		Cumulative
		Frequency	Percent	Valid Percent	Percent
Valid	1.00	49	2.9	2.9	2.9
	2.00	335	19.6	19.6	22.5
	3.00	16	.9	.9	23.4
	4.00	37	2.2	2.2	25.6
	5.00	199	11.7	11.7	37.3
	6.00	10	.6	.6	37.9
	7.00	1060	62.1	62.1	100.0
	Total	1706	100.0	100.0	

Ethnicity of Participant

Return to Video - 05. Data Management

Section 5.6: Collapsing a Variable

Once you've created the secondary variables that you feel you will need, you can also consider whether any of your quantitative variables or categorical variables need to be further grouped or binned.

As an example, we're going to show you how we might choose to group age in the NESARC sample. Rather than looking at age as a quantitative variable, we want to compare age groups categorically. We could examine the number of observations at each age and make some decisions about how we could best make groups.

Here's the frequency table for age in our young adult past year smoking sample. The cumulative percent column displays percentiles that can help you divide the sample into a specific number of equal parts. If we wanted 4 age groups, we would look for cumulative percentages that are closest to 25%, 50%, and 75%. If we wanted three age groups, we would look for cumulative percentages that are closest to 33% of the sample, and 66% of the sample.

Age of Participant							
		_			Cumulative		
		Frequency	Percent	Valid Percent	Percent		
Valid	18	161	9.4	9.4	9.4		
	19	200	11.7	11.7	21.2		
	20	221	13.0	13.0	34.1		
	21	239	14.0	14.0	48.1		
	22	228	13.4	13.4	61.5		
	23	231	13.5	13.5	75.0		
	24	241	14.1	14.1	89.2		
	25	185	10.8	10.8	100.0		
	Total	1706	100.0	100.0			

Let's create a categorical variable that divides the sample into roughly three equal size age groups. We chose the age cut points of 20 and 22 because these correspond roughly to the 33rd and 66th percentile.

1. Repeat steps 5.3 #1 and #2. We will name this new variable AGEGROUP. Remember to remove what is currently in the middle column, **Input Variable --> Output Variable**, using the arrow to move the variables back to the left column.

ta Recode into Different Varia	ibles	X			
 CHECK321 ETHRACE2A ✓ IDNUM ✓ NUMCIGMO_EST S101D1 S101D2 S101D3 S101D4 S101D5 S2A03 S2A03 S3A03B1 ✓ S3A03C1 	Numeric <u>V</u> ariable -> Output Variable: AGE> ?	Output Variable <u>Name:</u> <u>AGEGROUP</u> <u>Label:</u> e of participants in groups <u>Change</u>			
OK Paste Reset Cancel Help					

2. Repeat steps 5.3 #3 and #4. Remember to remove the **Old --> New** values in the lower right column before beginning. Click **Continue > OK**.

Old Value	-New Value
© <u>V</u> alue:	Value:
	◎ System-missing
© System-missing	© Copy old value(s)
System- or user-missing	Ol <u>d</u> > New:
	18 thru 20> 1
through	21 thru 22> 2
	23 thru 25> 3
© Range, LOWEST through value:	<u>Change</u> <u>Remove</u>
© Range, value through HIGHEST:	
	Output variables are strings Width: 8
All other values	Convert numeric strings to numbers ('5'->5)

- 3. Repeat steps 5.3 #5 through #8 to add the AGEGROUP variable you just created to your Variable Set.
- 4. It is important to always verify that how you intended to collapse a variable worked properly. Run a frequency distribution on AGE and AGEGROUP. Confirm that the when you add the frequency of each level (i.e., Age) you collapsed into agroup (dummy codes 1, 2, and 3) match. When we add 161 + 200 + 221 = 582 which is category 1 for our new categorical variable AGEGROUP. You should verify each level this way as it is possible that one level collapsed correctly and another did not.

Age of Participant							
		Frequency	Percent	Valid Percent	Percent		
Valid	18	161	9.4	9.4	9.4		
	19	200	11.7	11.7	21.2		
	20	221	13.0	13.0	34.1		
	21	239	14.0	14.0	48.1		
	22	228	13.4	13.4	61.5		
	23	231	13.5	13.5	75.0		
	24	241	14.1	14.1	89.2		
	25	185	10.8	10.8	100.0		
	Total	1706	100.0	100.0			

					Cumulative
		Frequency	Percent	Valid Percent	Percent
Valid	1.00	582	34.1	34.1	34.1
	2.00	467	27.4	27.4	61.5
	3.00	657	38.5	38.5	100.0
	Total	1706	100.0	100.0	

Now that you have seen a number of different data management decisions that could be made, I hope that you're ready to take on your own which you will do in the next assignment. How much data management you will need for your own variables will depend on the variables that you've selected and the decisions that you make about them. Some may choose to only code out missing data, while others may need to or want to do more. Data management is a part of the research process that you can and will return to again and again as you learn more and are able to make better decisions.