

# Analysis using survey weights

**Ernesto F. L. Amaral**

July 13, 2021

2021 Sociology Quantitative Methods Series



**TEXAS A&M**  
UNIVERSITY.

# Outline

- Inferential statistics
- Sample weights
- Weight options in Stata
- Complex sample cluster design
- Examples of weights in surveys
  - American Community Survey (ACS)
  - General Social Survey (GSS)
- Examples of descriptive statistics



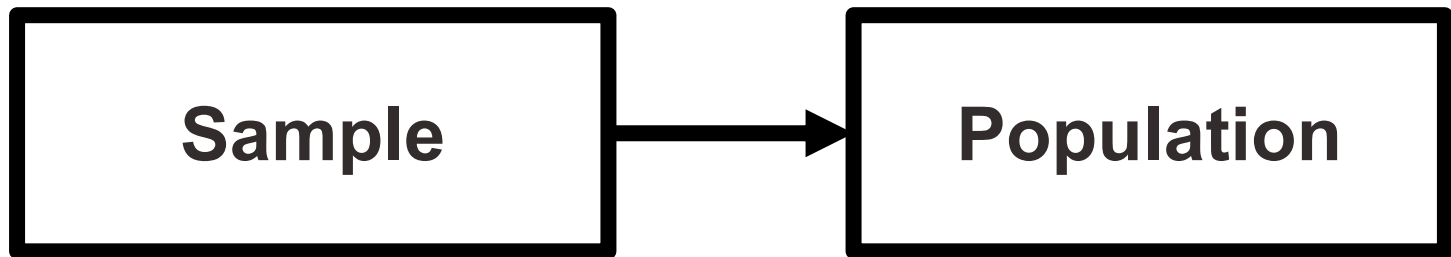
# Inferential statistics

- Social scientists need inferential statistics
  - They almost never have the resources or time to collect data from every case in a population
- Inferential statistics uses data from samples to make generalizations about populations
  - **Population** is the total collection of all cases in which the researcher is interested
  - **Samples** are carefully chosen subsets of the population
- With proper techniques, generalizations based on samples can represent populations

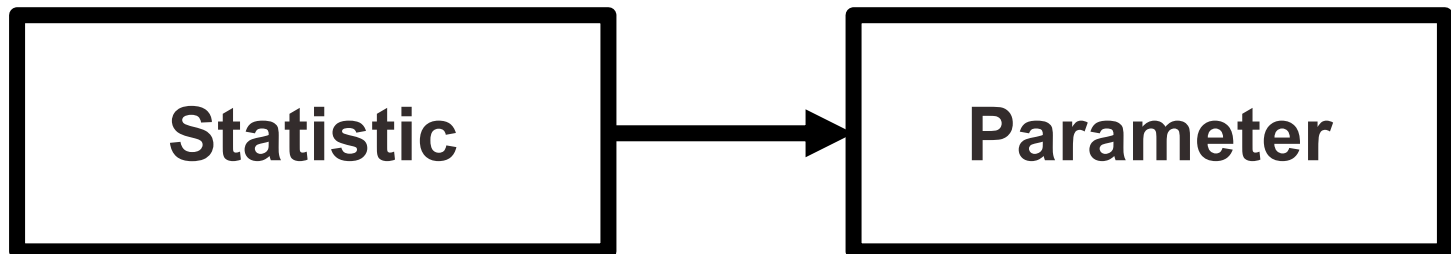


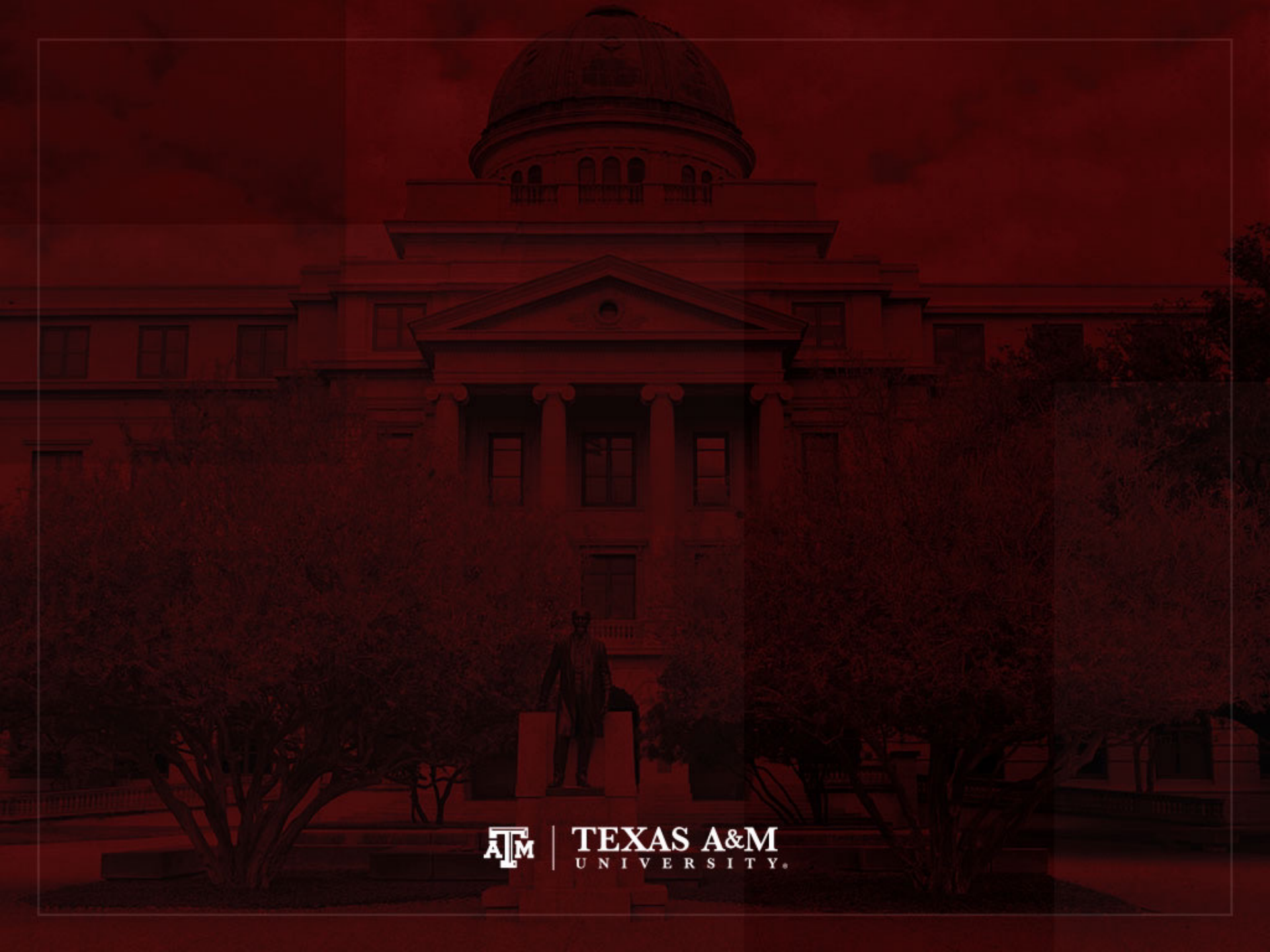
# Basic logic and terminology

- Information from samples is used to estimate information about the population



- Statistics: characteristics of samples
- Parameters: characteristics of populations
- Statistics are used to estimate parameters





TEXAS A&M  
UNIVERSITY.



# Sample weights

Name	Number of observations collected in the survey	Weight to expand to population size	Weight to maintain sample size
José	1	4	0.8
Maria	1	6	1.2
Total	2	10	2

**Sample weight =**  
**Population weight \* (Sum of sample weights / Sum of population weights)**



# Weights for tables

- When we use a sample to estimate the absolute number of people
  - For an area
  - For a specific sub-group
  - We use weights to expand to population size
- If we use a sample to estimate the proportion of people in a specific sub-group
  - And we are not concerned with the absolute value
  - We use weights to maintain the sample size (we focus on percentages)

# Weights for regressions

- In a simple linear regression, the test of statistical significance for a  $\beta$  coefficient ( $t$ -test) is estimated as

$$t = \frac{\hat{\beta}}{SE_{\hat{\beta}}} = \frac{\hat{\beta}}{\sqrt{\frac{MSE}{S_{xx}}}} = \frac{\hat{\beta}}{\sqrt{\frac{RSS}{df * S_{xx}}}} = \frac{\hat{\beta}}{\sqrt{\frac{\sum_i (y_i - \hat{y}_i)^2}{(n - 2) \sum_i (x_i - \bar{x})^2}}}$$

- $SE_{\hat{\beta}}$ : standard error of  $\beta$
- $MSE$ : mean squared error =  $RSS / df$
- $RSS$ : residual sum of squares =  $\sum_i (y_i - \hat{y}_i)^2 = \sum_i \hat{e}_i^2$
- $df$ : degrees of freedom =  $n-2$  for simple linear regression
  - 2 statistics (slope and intercept) are estimated to calculate sum of squares
- $S_{xx}$ : corrected sum of squares for  $x$  (total sum of squares)





# Weights for regressions

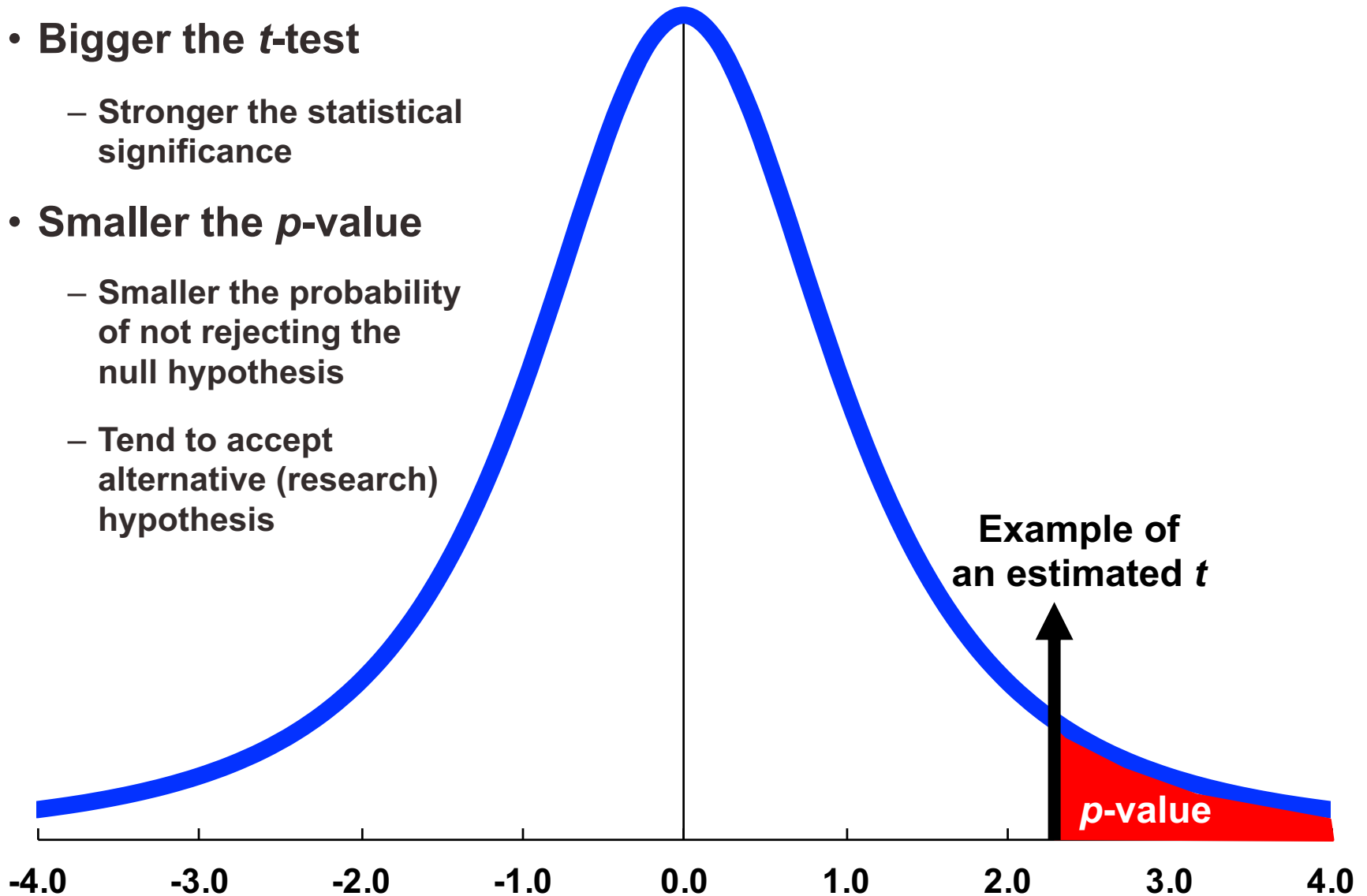
- If we use a weight that expands to the population size ( $N$ ) on regressions
  - We would be incorrectly informing the statistical software that we have a sample with enormous size
  - This would artificially increase the test of statistical significance for the coefficient

$$\uparrow t = \frac{\hat{\beta}}{SE_{\hat{\beta}}} = \frac{\hat{\beta}}{\sqrt{\frac{MSE}{S_{xx}}}} = \frac{\hat{\beta}}{\sqrt{\frac{MSE}{S_{xx}}}} = \frac{\hat{\beta}}{\sqrt{\frac{\sum_i (y_i - \hat{y}_i)^2}{(n-2) \sum_i (x_i - \bar{x})^2}}} \downarrow$$

- We have to inform the weight related to the sample design, but we should maintain the sample size ( $n$ )

# $t$ distribution ( $df = 2$ )

- Bigger the  $t$ -test
  - Stronger the statistical significance
- Smaller the  $p$ -value
  - Smaller the probability of not rejecting the null hypothesis
  - Tend to accept alternative (research) hypothesis



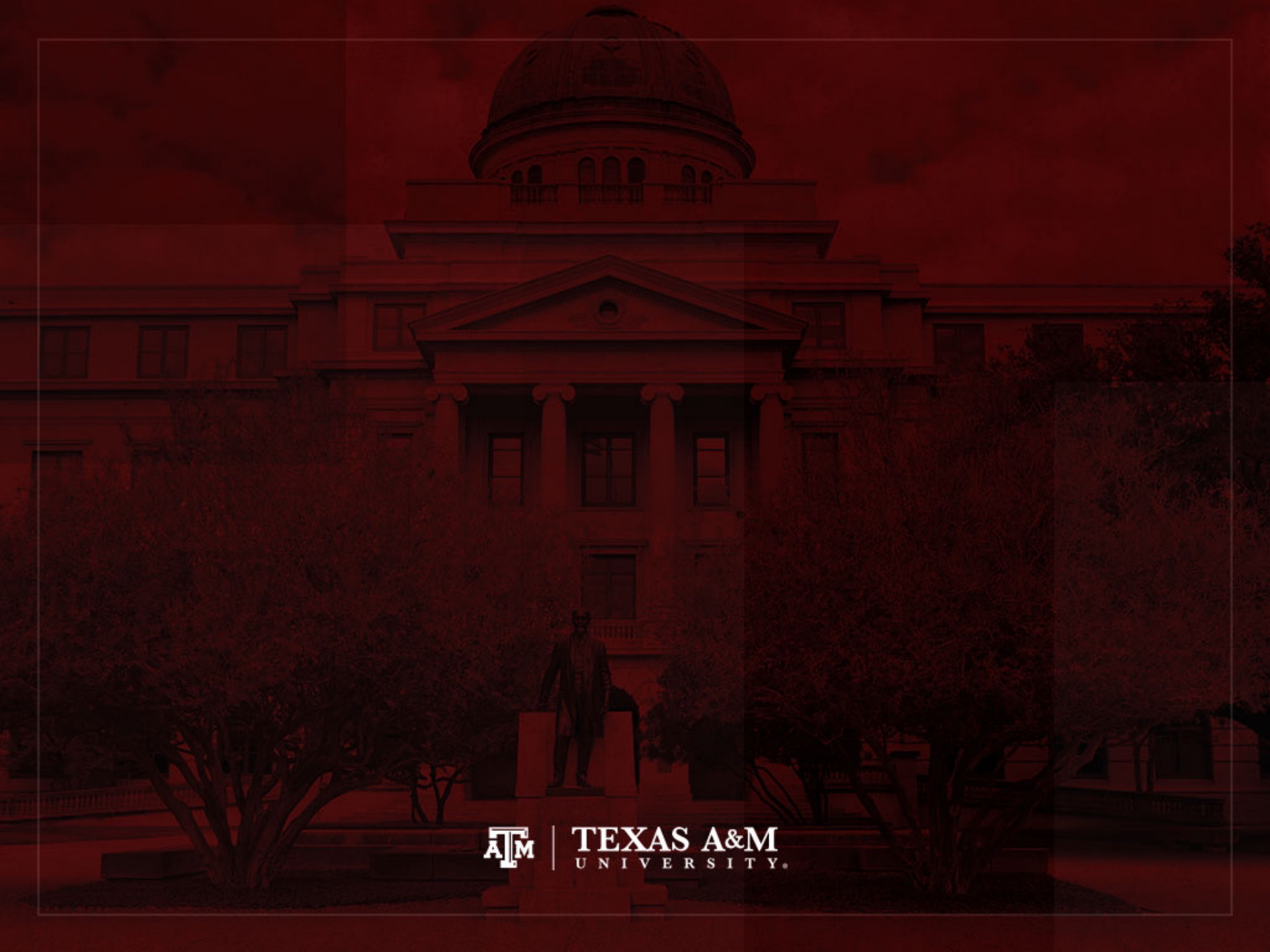
# Decisions about hypotheses

Hypotheses	$p < \alpha$	$p > \alpha$
Null hypothesis ( $H_0$ )	Reject	Do not reject
Alternative hypothesis ( $H_1$ )	Accept	Do not accept

- **$p$ -value** is the probability of not rejecting the null hypothesis
- If a statistical software gives only the two-tailed  $p$ -value, divide it by 2 to obtain the one-tailed  $p$ -value

Significance level ( $\alpha$ )	Confidence level (success rate)
0.10 (10%)	90%
0.05 (5%)	95%
0.01 (1%)	99%
0.001 (0.1%)	99.9%





TEXAS A&M  
UNIVERSITY.

# Extract of 2018 ACS microdata

	year	strata	cluster	perwt	hhwt	sex	age	income
1	2018	360248	2.018012e+12	56.00	56.00	Male	46	28000
2	2018	360248	2.018012e+12	51.00	51.00	Male	20	5000
3	2018	360248	2.018012e+12	76.00	76.00	Female	84	0
4	2018	360248	2.018012e+12	55.00	55.00	Female	18	1200
5	2018	360248	2.018012e+12	143.00	143.00	Female	56	1500
6	2018	360248	2.018012e+12	198.00	198.00	Male	31	10000
7	2018	360248	2.018012e+12	48.00	48.00	Female	19	2000
8	2018	360248	2.018012e+12	48.00	48.00	Male	25	7000
9	2018	360248	2.018012e+12	65.00	65.00	Female	18	0
10	2018	360248	2.018012e+12	53.00	53.00	Female	18	15000
11	2018	360248	2.018012e+12	17.00	17.00	Male	63	0
12	2018	360248	2.018012e+12	39.00	39.00	Female	18	4000
13	2018	360248	2.018012e+12	104.00	104.00	Male	21	1000
14	2018	360248	2.018012e+12	200.00	200.00	Male	40	80000
15	2018	360248	2.018012e+12	20.00	20.00	Male	33	0
16	2018	360248	2.018012e+12	59.00	59.00	Male	19	2900
17	2018	360248	2.018012e+12	56.00	56.00	Male	55	0
18	2018	360248	2.018012e+12	77.00	77.00	Male	18	9000
19	2018	360248	2.018012e+12	16.00	16.00	Female	41	1100
20	2018	360248	2.018012e+12	46.00	46.00	Male	33	0



# Frequency weight in Stata

- **FWEIGHT**

- Expands survey size to the population size
- Indicates the number of duplicated observations
- Used on tables to generate frequencies
- Can be used in frequency distributions only when weight variable is discrete (no fractional numbers)

```
tab x [fweight = weight]
```





# "Importance" weight in Stata

- **IWEIGHT**

- Indicates the "importance" of the observation in some vague sense
- Has no formal statistical definition
- Any command that supports iweights will define exactly how they are treated
- Intended for use by programmers who want to produce a certain computation
- Can be used in frequency distributions even when weight variable is continuous (fractional numbers)

```
tab x [iweight = weight]
```



# Analytic weight in Stata

- **AWEIGHT**

- Inversely proportional to the variance of an observation
- Variance of the  $j$ th observation is assumed to be  $\sigma^2/w_j$ , where  $w_j$  are the weights
- For most Stata commands, the recorded scale of aweights is irrelevant
- Stata internally rescales frequencies, so sum of weights equals sample size

```
tab x [aweight = weight]
```

```
regress y x1 x2 [aweight = weight]
```



# More about analytic weight

- Observations represent averages and weights are the number of elements that gave rise to the average

group	x	y	n
1	3.5	26.0	2
2	5.0	20.0	3

- Instead of

group	x	y
1	3	22
1	4	30
2	8	25
2	2	19
2	5	16

- Usually, survey data is collected from individuals and households (not as averages)
  - Thus, aweights are not appropriate for most cases



# Sampling weight in Stata

- **PWEIGHT**

- Denote the inverse of the probability that the observation is included due to the sampling design
- Variances, standard errors, and confidence intervals are estimated with a more precise procedure
- Indicated for statistical regressions to estimate robust standard errors
  - Obtain unbiased standard errors of OLS coefficients under heteroscedasticity (i.e., residuals not randomly distributed)
  - Robust standard errors are usually larger than conventional ones

**`regress y x1 x2 [pweight = weight]`**



# Summary of Stata weights

WEIGHTS IN FREQUENCY DISTRIBUTIONS		
Weight unit of measurement	Expand to population size	Maintain sample size
Discrete	fweight	aweight
Continuous	iweight	

WEIGHTS IN STATISTICAL REGRESSIONS should maintain sample size	
Robust standard error	Adjusted R <sup>2</sup> , TSS, ESS, RSS
pweight	aweight
reg y x, robust	outreg2



# Example of 2018 ACS weights

. tab sex

Sex	Freq.	Percent	Cum.
Male	1,574,618	48.98	48.98
Female	1,639,921	51.02	100.00
Total	3,214,539	100.00	

. tab sex [iweight=perwt]

Sex	Freq.	Percent	Cum.
Male	161,072,404	49.23	49.23
Female	166,095,035	50.77	100.00
Total	327,167,439	100.00	

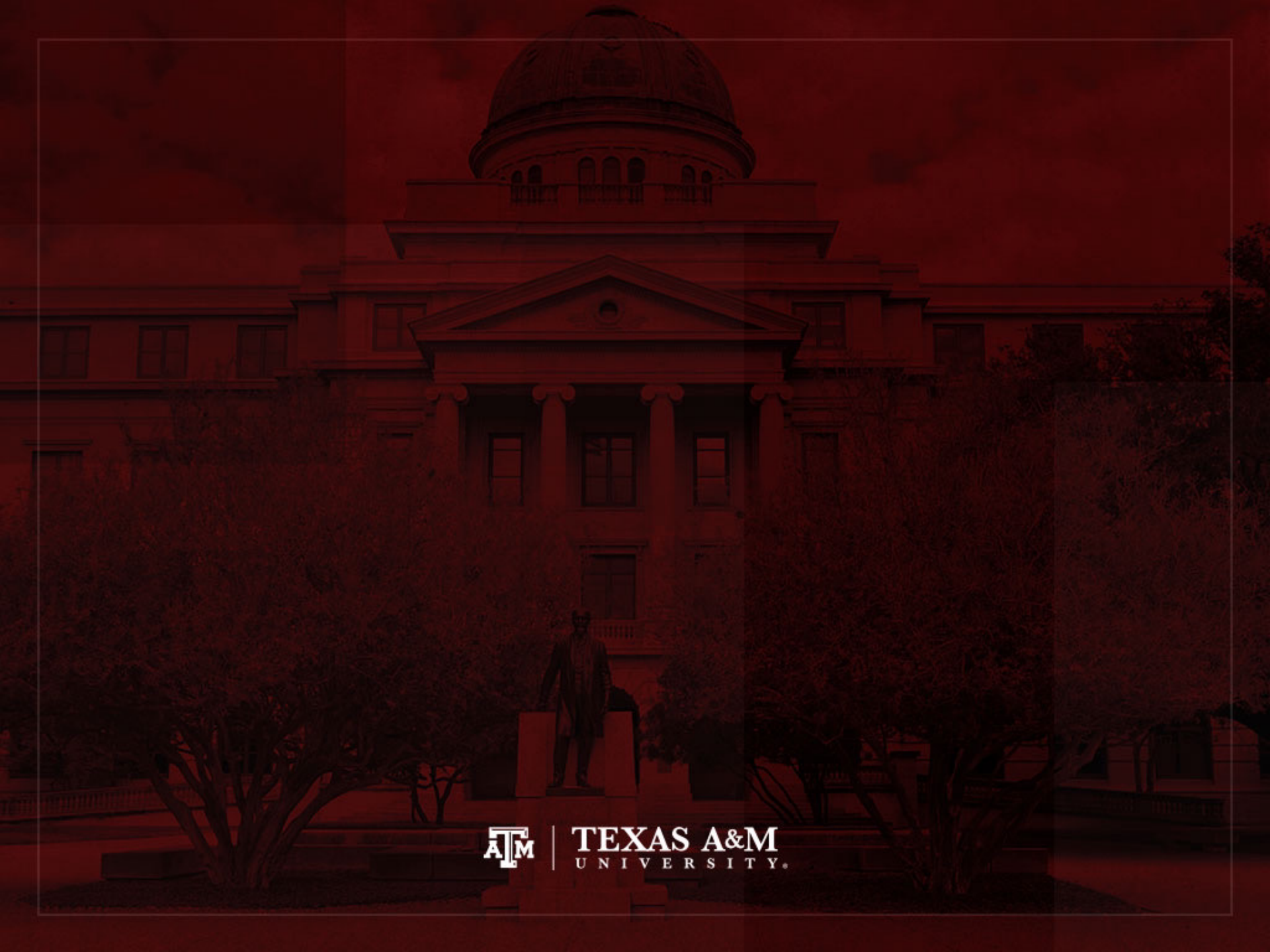
. tab sex [aweight=perwt]

Sex	Freq.	Percent	Cum.
Male	1,582,595	49.23	49.23
Female	1,631,944	50.77	100.00
Total	3,214,539	100.00	

. tab sex [fweight=perwt]

Sex	Freq.	Percent	Cum.
Male	161,072,404	49.23	49.23
Female	166,095,035	50.77	100.00
Total	327,167,439	100.00	





TEXAS A&M  
UNIVERSITY.

# Complex sample cluster design

- To calculate standard errors correctly, variables for sample cluster design must be used
  - Without design variables, Stata will assume a simple random sample and underestimate standard errors
- Strata are created based on the lowest level of geography available in each sample
  - We use additional statistical techniques that account for the complex sample design to produce correct standard errors and statistical tests



# Cluster design for tables

- If we want to estimate a confidence interval for a sample statistic (mean or proportion), we need to inform the complex survey design
- **Confidence interval** is a range of values used to estimate the true population parameter
- **Confidence level** is the success rate of the procedure to estimate the confidence interval
- Larger confidence levels generate larger confidence intervals

# Confidence level, $\alpha$ , and Z

Confidence level (1 - $\alpha$ ) * 100	Significance level alpha ( $\alpha$ )	$\alpha / 2$	Z score
90%	0.10	0.05	$\pm 1.65$
<b>95%</b>	<b>0.05</b>	<b>0.025</b>	<b><math>\pm 1.96</math></b>
99%	0.01	0.005	$\pm 2.58$
99.9%	0.001	0.0005	$\pm 3.32$
99.99%	0.0001	0.00005	$\pm 3.90$



# Confidence intervals from samples

*c.i. = sample estimate  $\pm$  margin of error*

*c.i. = sample estimate  $\pm$  score of confidence level \* standard error*

- Sample mean ( $\bar{x}$ ), standard deviation (s),  $n < 30$

$$c.i. = \bar{x} \pm t \left( \frac{s}{\sqrt{n}} \right) \quad df = n - 1$$

- Sample mean ( $\bar{x}$ ), standard deviation (s),  $n \geq 30$

$$c.i. = \bar{x} \pm Z \left( \frac{s}{\sqrt{n - 1}} \right)$$

- Sam. proportion ( $P_s$ ), pop. proportion ( $P_u$ ),  $n \geq 30$

$$c.i. = P_s \pm Z \sqrt{\frac{P_u(1 - P_u)}{n}}$$



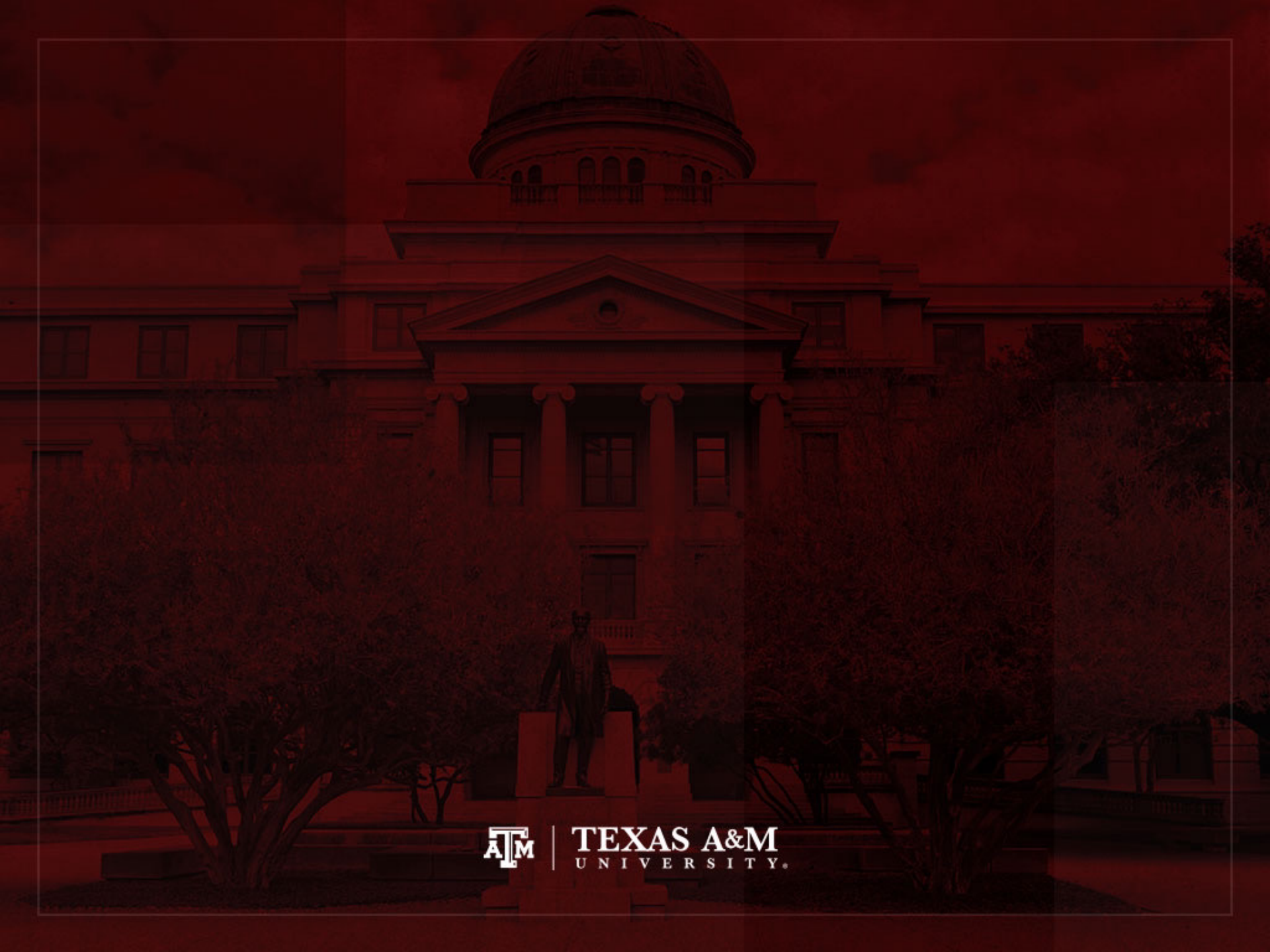
# Cluster design for regressions

- We also need to inform cluster design for regressions, because the  $t$ -test utilizes standard errors

$$t = \frac{\hat{\beta}}{SE_{\hat{\beta}}} = \frac{\hat{\beta}}{\sqrt{\frac{MSE}{S_{xx}}}} = \frac{\hat{\beta}}{\sqrt{\frac{RSS}{df * S_{xx}}}} = \frac{\hat{\beta}}{\sqrt{\frac{\sum_i (y_i - \hat{y}_i)^2}{(n - 2) \sum_i (x_i - \bar{x})^2}}}$$

- $SE_{\hat{\beta}}$ : standard error of  $\beta$
- $MSE$ : mean squared error =  $RSS / df$
- $RSS$ : residual sum of squares =  $\sum_i (y_i - \hat{y}_i)^2 = \sum_i \hat{e}_i^2$
- $df$ : degrees of freedom =  $n-2$  for simple linear regression
- $S_{xx}$ : corrected sum of squares for  $x$  (total sum of squares)





TEXAS A&M  
UNIVERSITY.

# Weights in ACS

- In the American Community Survey (ACS) PERWT indicates how many persons in the U.S. population are represented by a given person in an IPUMS sample

[https://usa.ipums.org/usa-action/variables/PERWT#description\\_section](https://usa.ipums.org/usa-action/variables/PERWT#description_section)

- HHWT indicates how many households in the U.S. population are represented by a given household in an IPUMS sample
  - Users should also be sure to select one person (e.g., PERNUM = 1) to represent the entire household

[https://usa.ipums.org/usa-action/variables/HHWT#description\\_section](https://usa.ipums.org/usa-action/variables/HHWT#description_section)



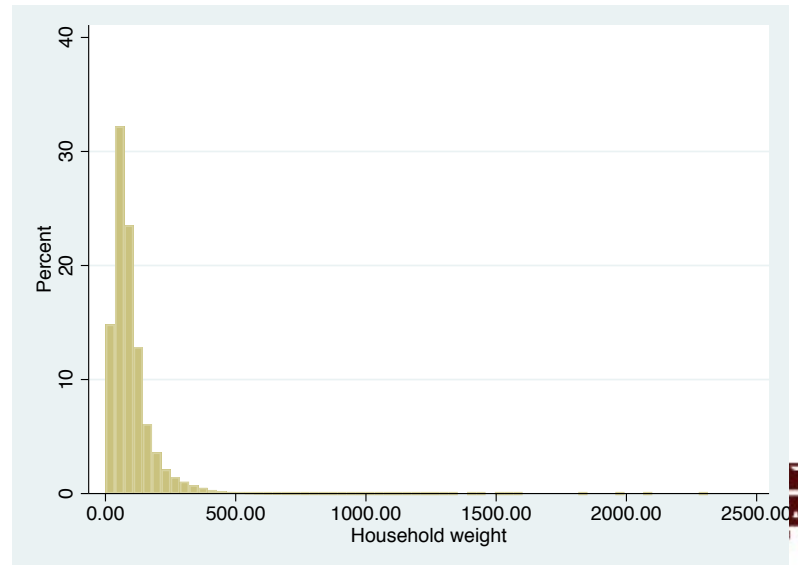
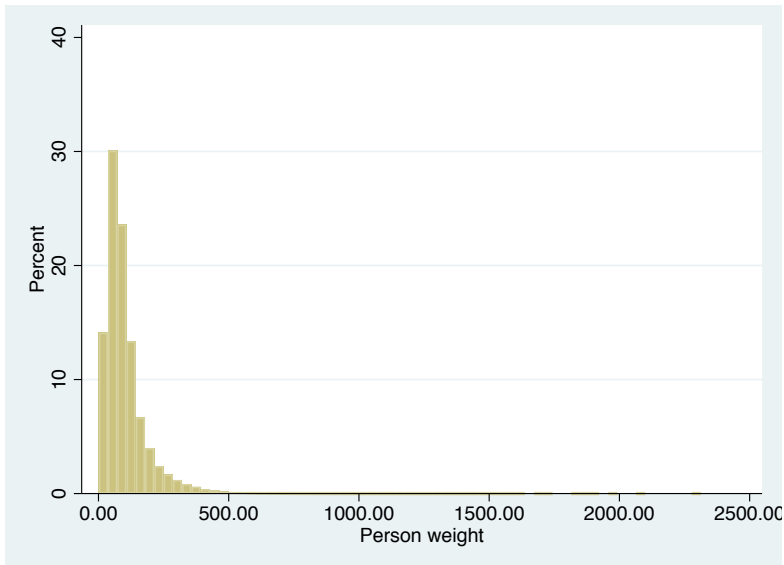
# Summary of 2018 ACS weights

. sum perwt, d

Person weight				
Percentiles		Smallest		
1%	10	1		
5%	19	1		
10%	29	1	Obs	3,214,539
25%	52	1	Sum of Wgt.	3,214,539
50%	80		Mean	101.7774
		Largest	Std. Dev.	83.93534
75%	124	1916		
90%	195	1990	Variance	7045.14
95%	263	2097	Skewness	2.845116
99%	427	2313	Kurtosis	17.99265

. sum hhwt if pernum==1, d

Household weight				
Percentiles		Smallest		
1%	8	1		
5%	16	1		
10%	25	1	Obs	1,410,976
25%	48	1	Sum of Wgt.	1,410,976
50%	73		Mean	91.85967
		Largest	Std. Dev.	75.18581
75%	112	1837		
90%	173	1990	Variance	5652.906
95%	234	2097	Skewness	2.88203
99%	386	2313	Kurtosis	19.09996



# ACS has a cluster sample

- All IPUMS samples are cluster samples
  - Samples are not individual-level samples
  - They are samples of households or dwellings
  - Individuals are sampled as parts of households
    - Information about all individuals within the same household
- Samples are also stratified to some degree
  - U.S. Census Bureau divides population into strata based on key characteristics
  - Sample separately from each stratum
  - Each stratum is proportionately represented in the final sample



# ACS variables for cluster design

- Sampling weight (PERWT or HHWT)
  - It is chosen based on type of research question
- Household strata (STRATA)
  - Integrated variable that represents the impact of the sample design stratification on the estimates of variance and standard errors
  - In the 2005 onward ACS samples, strata are defined as unique Public Use Micro-data Areas (PUMA)
- Household cluster (CLUSTER)
  - Integrated variable which uniquely identifies each household record in a given sample





# ACS complex sample design

- Account for ACS sample design in Stata  
`svyset cluster [pweight=perwt], strata(strata)`

```
. svyset cluster [pweight=perwt], strata(strata)
```

```
      pweight: perwt  
      VCE: linearized  
Single unit: missing  
Strata 1: strata  
SU 1: cluster  
FPC 1: <zero>
```

- After "svyset," you should indicate survey design with the option "svy" for commands that estimate standard errors

```
      svy: mean y  
      svy: reg y x1 x2
```





# Mean income

```
. mean income [pweight=perwt]
```

Mean estimation

Number of obs = 2,642,681

	Mean	Std. Err.	[95% Conf. Interval]	
income	31175.11	39.98542	31096.74	31253.48

```
. svy: mean income
```

(running mean on estimation sample)

Survey: Mean estimation

Number of strata = 2,351

Number of obs = 2,642,681

Number of PSUs = 1408111

Population size = 262,216,823

Design df = 1,405,760

	Mean	Linearized Std. Err.	[95% Conf. Interval]	
income	31175.11	40.99966	31094.75	31255.47



# For subpopulations

- We use the following approach to conduct subpopulation analysis without compromising the data design structure
  - We produce estimates for the population of interest, while incorporating the full sample design information for variance estimation
- Example: only people with 15–64 years of age

```
svyset cluster [pweight=perwt], strata(strata)
svy, subpop(if age>=15 & age<=64): mean var1
```



# Mean income

```
. svy: mean income
(running mean on estimation sample)
```

Survey: Mean estimation

```
Number of strata =    2,351      Number of obs   =    2,642,681
Number of PSUs   =  1408111      Population size =  262,216,823
Design df        =    1,405,760
```

	Mean	Linearized Std. Err.	[95% Conf. Interval]	
income	31175.11	40.99966	31094.75	31255.47

```
. svy, subpop(if income!=.): mean income
(running mean on estimation sample)
```

Survey: Mean estimation

```
Number of strata =    2,351      Number of obs   =    3,214,539
Number of PSUs   =  1410976      Population size =  327,167,439
Subpop. no. obs  =    2,642,681
Subpop. size     =  262,216,823
Design df        =    1,408,625
```

	Mean	Linearized Std. Err.	[95% Conf. Interval]	
income	31175.11	41.00232	31094.74	31255.47

If we consider that missing cases are part of the population, we need to inform that subpopulation is only non-missing cases



# Mean income (15–64)

```
. svy, subpop(if age>=15 & age<=64): mean income
(running mean on estimation sample)
```

Survey: Mean estimation

Number of strata =	2,351	Number of obs =	3,175,157
Number of PSUs =	1410150	Population size =	323,036,047
		Subpop. no. obs =	2,004,091
		Subpop. size =	209,809,274
		Design df =	1,407,799

	Linearized			
	Mean	Std. Err.	[95% Conf. Interval]	
income	36736.34	48.39971	36641.48	36831.2

If we consider that missing cases are part of the population, we need to inform that subpopulation is only non-missing cases

```
. svy, subpop(if age>=15 & age<=64 & income!=.): mean income
(running mean on estimation sample)
```

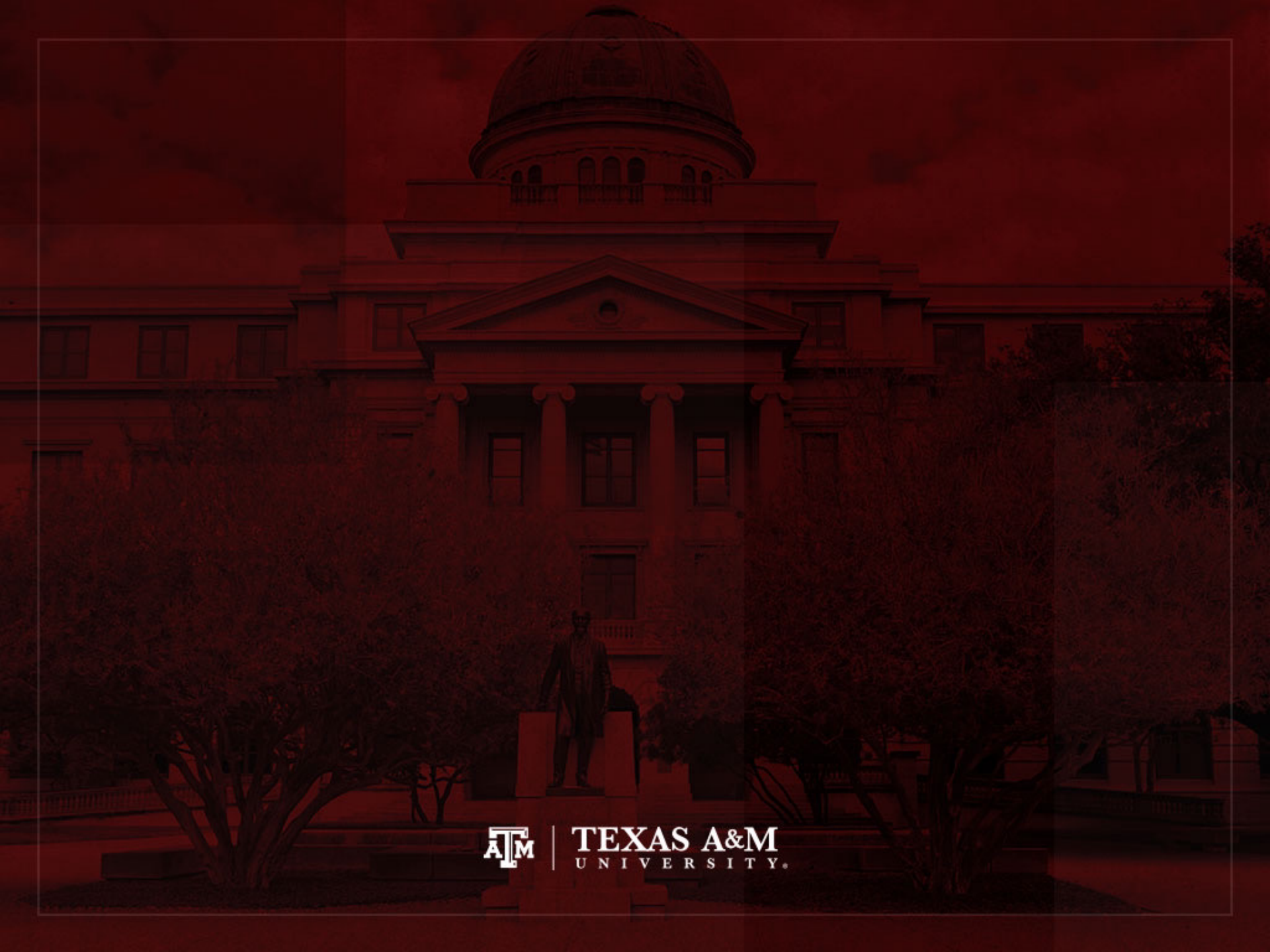
Survey: Mean estimation

Number of strata =	2,351	Number of obs =	3,214,539
Number of PSUs =	1410976	Population size =	327,167,439
		Subpop. no. obs =	2,004,091
		Subpop. size =	209,809,274
		Design df =	1,408,625

	Linearized			
	Mean	Std. Err.	[95% Conf. Interval]	
income	36736.34	48.40061	36641.48	36831.21







TEXAS A&M  
UNIVERSITY.

# Weights in GSS

- The General Social Survey (GSS) targets the adult population (18+) living in U.S. households
- Due to the adoption of the sub-sampling design of non-respondents, a weight must be employed when using the GSS 2004 and after
- There are three continuous weight variables
  - WTSS
  - WTSSNR
  - WTSSALL
- They all maintain the original sample size, even in frequency distributions with "iweight"





# WTSS

- WTSS variable takes into consideration
  - Sub-sampling of non-respondents
  - Number of adults in the household
- In years prior to 2004, a value of one is assigned to all cases, so they are effectively unweighted
  - Number of adults can be utilized to make this adjustment for years prior to 2004

# WTSSNR

- WTSSNR variable takes into consideration
  - Sub-sampling of non-respondents
  - Number of adults in the household
  - Differential non-response across areas
- In years prior to 2004, a value of one is assigned to all cases, so they are effectively unweighted
  - Number of adults can be utilized to make this adjustment for years prior to 2004
  - Area non-response adjustment is not possible

# WTSSALL

- WTSSALL takes WTSS and applies an adult weight to years before 2004
- The weight value of WTSSALL is the same as WTSS for 2004 and after
- Researchers who use the GSS data before or after 2004 may consider using the WTSSALL weight variable

```
tab x [aweight = wtssall]
```

```
sum x [aweight = wtssall]
```



# GSS has a cluster sample

([https://gssdataexplorer.norc.umd.edu/pages/show?page=gss%2Fstandard\\_error](https://gssdataexplorer.norc.umd.edu/pages/show?page=gss%2Fstandard_error))

- First- and second-stage units are selected with probabilities proportional to size
  - Size is defined by number of housing units
- Third-stage units (housing units) are selected to be an equal-probability sample
  - This results in roughly the same number of housing units selected per second-stage sampling unit



# GSS variables for cluster design

([https://gssdataexplorer.norc.umd.edu/pages/show?page=gss%2Fstandard\\_error](https://gssdataexplorer.norc.umd.edu/pages/show?page=gss%2Fstandard_error))

- There are two design variables
  - VSTRAT
  - VPSU
- First-stage unit
  - VSTRAT: Variance Stratum
  - National Frame Areas (NFAs): one or more counties
- Second-stage unit
  - VPSU: Variance Primary Sampling Unit
  - Segments: block, group of blocks, or census tract



# GSS complex sample design

([https://gssdataexplorer.norc.umd.edu/pages/show?page=gss%2Fstandard\\_error](https://gssdataexplorer.norc.umd.edu/pages/show?page=gss%2Fstandard_error))

- Account for GSS sample design in Stata

```
svyset [weight=wtssall], strata(vstrat) psu(vpsu) singleunit(scaled)
```

- After "svyset," you should indicate survey design with the option "svy" for commands that estimate standard errors

```
svy: mean y
```

```
svy: reg y x1 x2
```





# Strata with single sampling unit

([https://gssdataexplorer.norc.umd.edu/pages/show?page=gss%2Fstandard\\_error](https://gssdataexplorer.norc.umd.edu/pages/show?page=gss%2Fstandard_error))

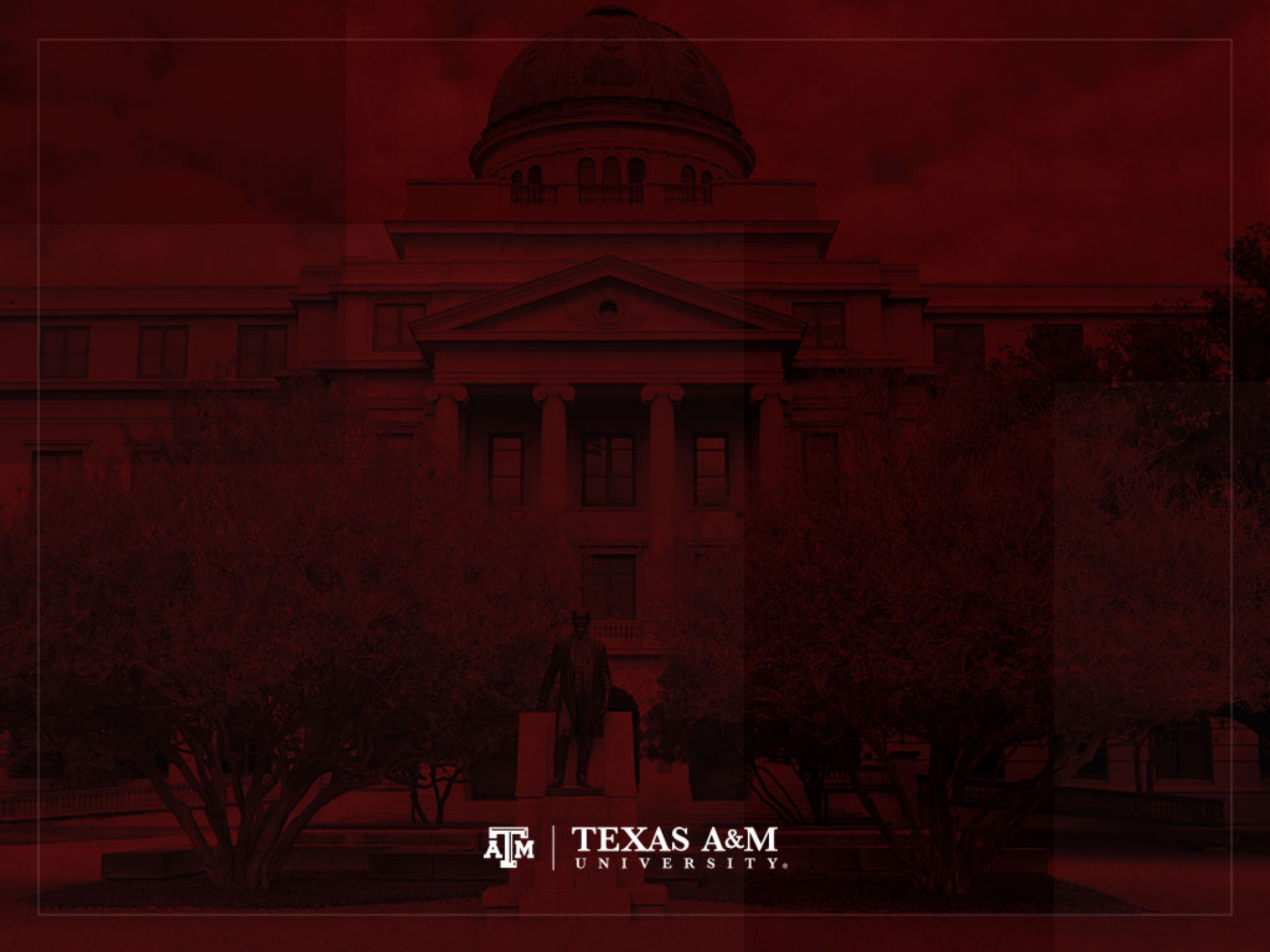
- VSTRAT and VPSU were created with a minimum of three respondents within a cell
  - If all cases are missing on a variable, you get an error message in Stata
  - "Missing standard error because of stratum with single sampling unit"
- It is recommended to utilize the "subpop" option for any subdomain analyses (e.g., for males)

```
svy, subpop(if sex==1): tab x
```

- You can also specify that strata with one sampling unit are "centered" at grand mean instead of stratum mean

```
svyset [weight=wtssall], strata(vstrat) psu(vpsu) singleunit(centered)
```





TEXAS A&M  
UNIVERSITY.

# Descriptive statistics

- Nominal-level variable
- Ordinal-level variable
- Interval-ratio-level variable



# Nominal-level variable (Example: 2018 ACS in Stata)

```
. tab raceth [fweight=perwt]
```

raceth	Freq.	Percent	Cum.
White	197,034,851	60.22	60.22
African American	40,373,281	12.34	72.56
Hispanic	59,740,273	18.26	90.82
Asian	18,662,293	5.70	96.53
Native American	2,170,486	0.66	97.19
Other races	9,186,255	2.81	100.00
Total	327,167,439	100.00	

```
. count if raceth!=.  
3,214,539
```



# Edited table

**Table 1. Distribution of U.S. population by race/ethnicity, 2018**

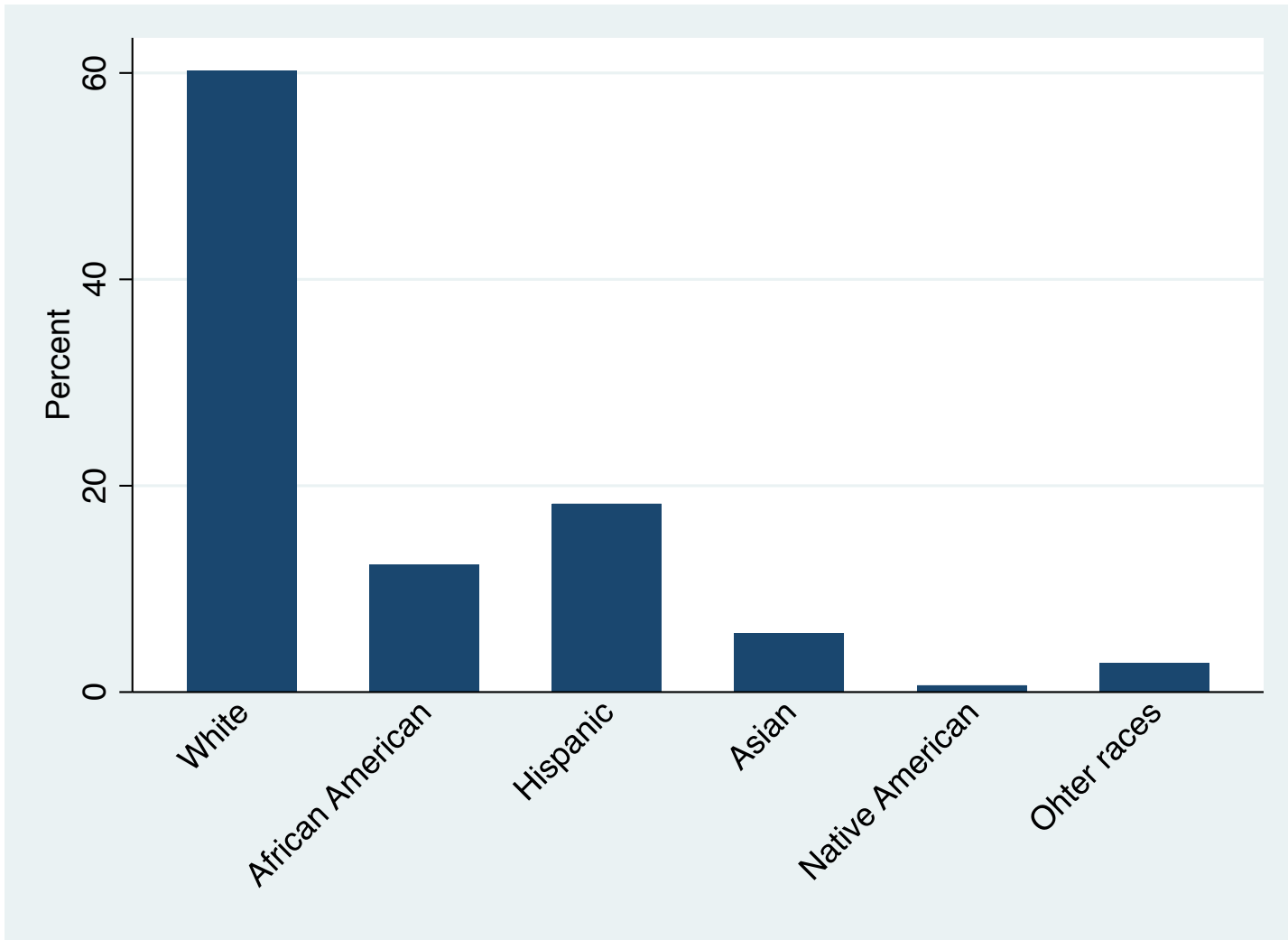
<b>Race/ethnicity</b>	<b>Percentage</b>
Non-Hispanic White	60.22
Non-Hispanic African American	12.34
Hispanic	18.26
Non-Hispanic Asian	5.70
Non-Hispanic Native American	0.66
Other races	2.81
<b>Total</b>	<b>99.99</b>
Population size (N)	327,167,439
Sample size (n)	3,214,539

Source: 2018 American Community Survey.



# Column graph for race/ethnicity, 2018

```
graph bar [fweight=perwt], over(raceth,  
label(angle(45))) ytitle("Percent")
```



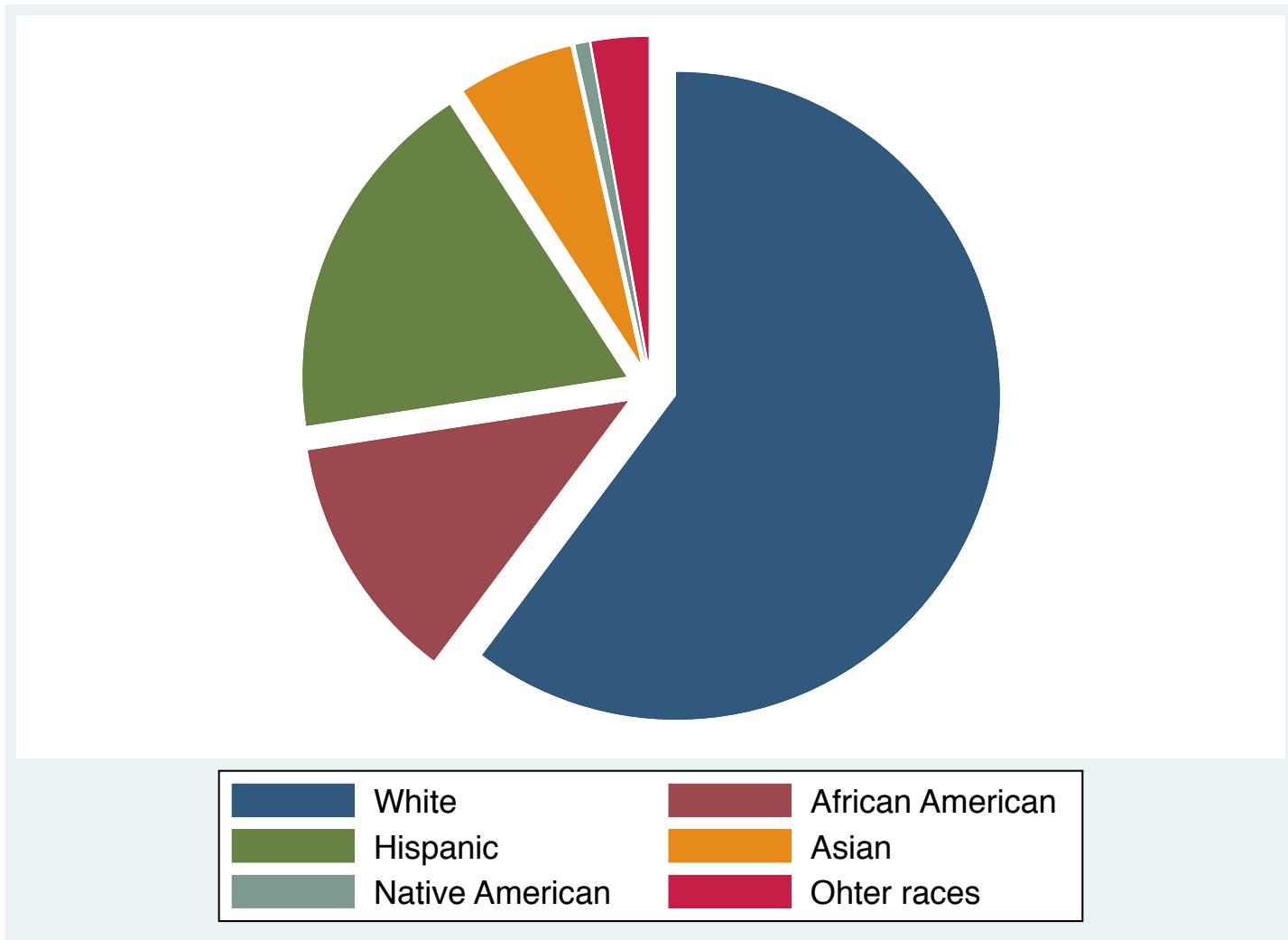
Source: 2018 American Community Survey.





# Pie graph for race/ethnicity, 2018

```
graph pie [fweight=perwt], over(raceth) pie(_all, explode)
```



Source: 2018 American Community Survey.



# Ordinal-level variable

## (Example: 2018 ACS in Stata)

```
. tab educgr [fweight=perwt]
```

educgr	Freq.	Percent	Cum.
Less than high school	97,758,814	29.88	29.88
High school	92,183,547	28.18	58.06
Some college	60,822,461	18.59	76.65
College	47,865,798	14.63	91.28
Graduate school	28,536,819	8.72	100.00
Total	327,167,439	100.00	

```
. count if educgr!=.  
3,214,539
```



# Edited table

**Table 1. Distribution of U.S. population by educational attainment, 2018**

<b>Educational attainment</b>	<b>Percentage</b>
Less than high school	29.88
High school	28.18
Some college	18.59
College	14.63
Graduate school	8.72
<b>Total</b>	<b>100.00</b>
Population size (N)	327,167,439
Sample size (n)	3,214,539

Source: 2018 American Community Survey.



# Interval-ratio-level variable (Example: 2018 ACS in Stata)

```
. table year [fweight=perwt] if income!=0, c(min income p25 income p50 income p75 income max income)
```

Census year	min(income)	p25(income)	med(income)	p75(income)	max(income)
2018	4	16400	35000	61000	718000

```
. table year [fweight=perwt] if income!=0, c(iqr income sd income mean income)
```

Census year	iqr(income)	sd(income)	mean(income)
2018	44600	62143.93	50043.98

```
. count if income==. | income==0  
1,640,226
```



# Survey design for income

```
. ***Complex survey design
. svy, subpop(if income!=. & income!=0): mean income
(running mean on estimation sample)
```

Survey: Mean estimation

Number of strata =	2,351	Number of obs =	3,214,539
Number of PSUs =	1410976	Population size =	327,167,439
		Subpop. no. obs =	1,574,313
		Subpop. size =	163,349,075
		Design df =	1,408,625

	Linearized			
	Mean	Std. Err.	[95% Conf. Interval]	
income	50043.98	59.74195	49926.89	50161.07

```
.
. ***Estimate standard deviation
. estat sd
```

	Mean	Std. Dev.
income	50043.98	61547.67



# Edited table

**Table 1. Descriptive statistics of respondents' wage and salary income, U.S. population, 2018**

<b>Statistics</b>	<b>Income</b>
Mean	50,043.98
Minimum	4.00
25th percentile	16,400.00
Median	35,000.00
75th percentile	61,000.00
Maximum	718,000.00
Range	717,996.00
Interquartile range	44,600.00
Standard deviation	61,547.67
<b>Population size (N)</b>	<b>163,349,075</b>
<b>Sample size (n)</b>	<b>1,574,313</b>
Missing cases	1,640,226

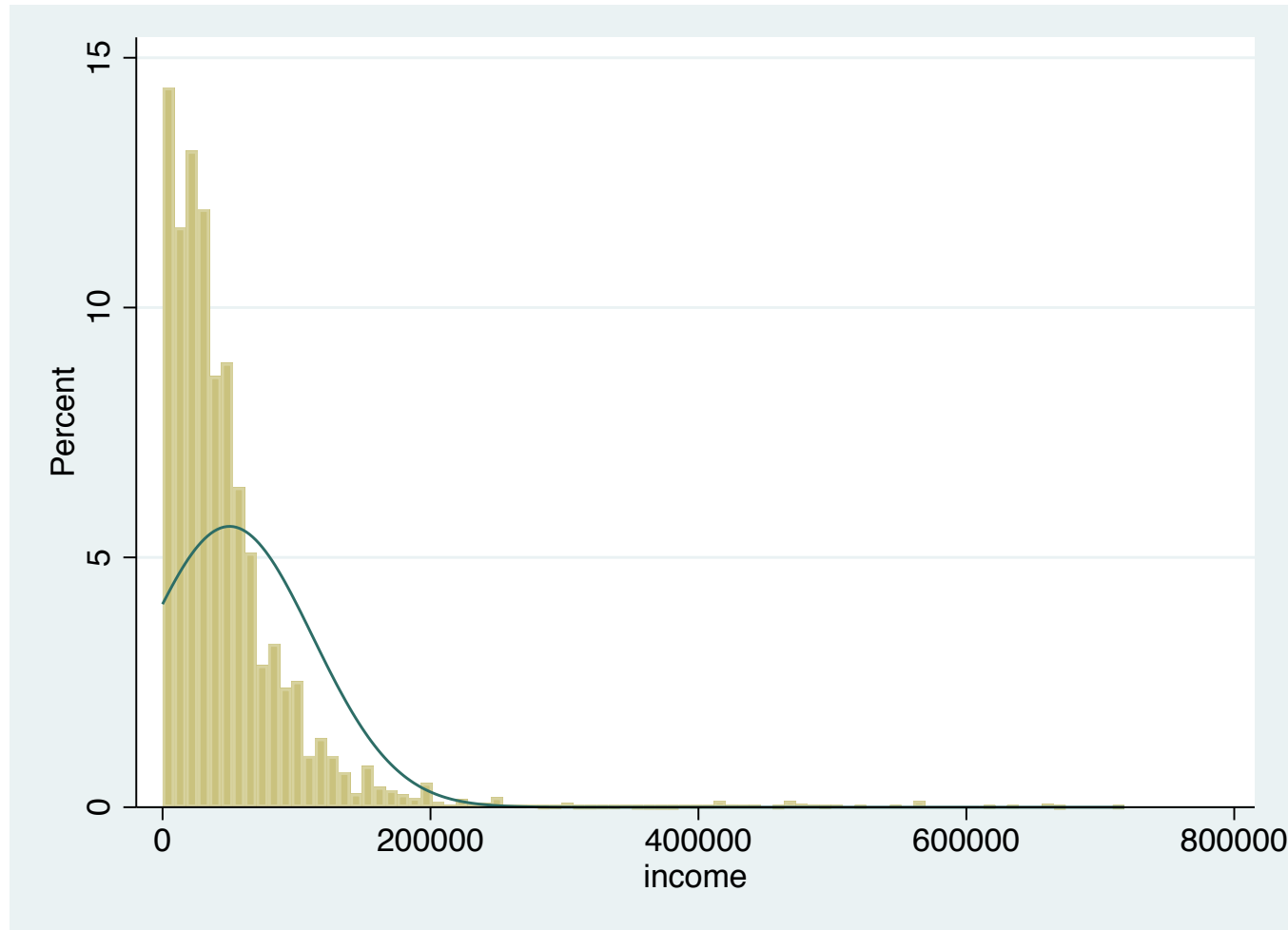
Source: 2018 American Community Survey.





# Histogram of wage and salary income, U.S. population, 2018

```
hist income [fweight=perwt] if income!=0, percent normal
```



Source: 2018 American Community Survey.

Obs.: Only people with some wage and salary income are included (different than zero).



# Wage and salary income by sex, 2018 ACS

```
. ***Income  
. table year [fweight=perwt] if income!=0, c(mean income p50 income)
```

Census year	mean(income)	med(income)
2018	50043.98	35000

```
. ***Income by sex  
. table female [fweight=perwt] if income!=0, c(mean income p50 income)
```

female	mean(income)	med(income)
Male	59014.14	40000
Female	40294.34	30000

# Wage and salary income by race/ethnicity, 2018 ACS

```
. ***Income by race/ethnicity  
. table raceth [fweight=perwt] if income!=0, c(mean income p50 income)
```

raceth	mean(income)	med(income)
White	55289.18	40000
African American	37183.63	29000
Hispanic	36236.16	27500
Asian	64154.23	43000
Native American	34851.55	27000
Other races	44162.79	30000



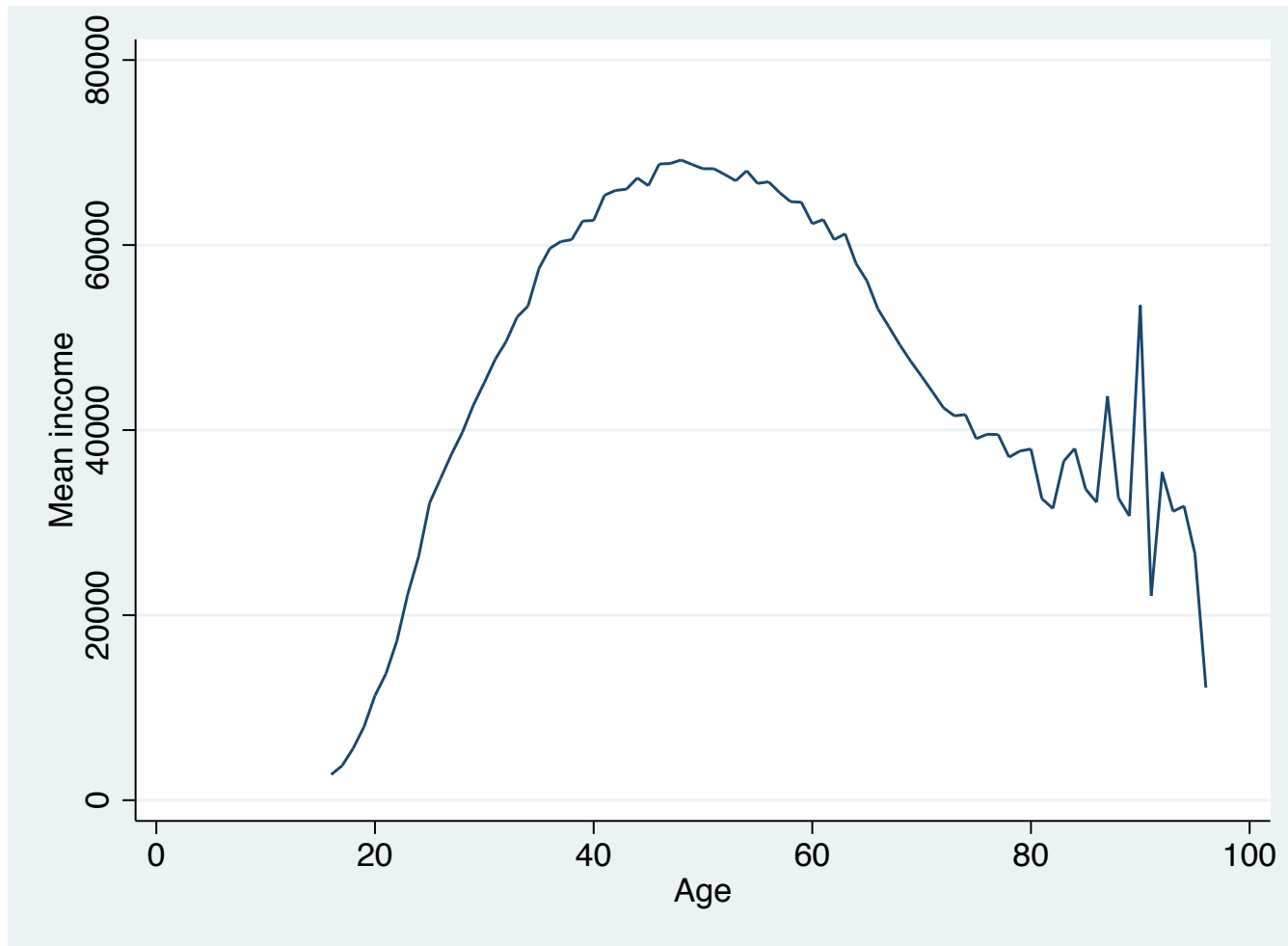
# Wage and salary income by education, 2018 ACS

```
. ***Income by educational attainment  
. table educgr [fweight=perwt] if income!=0, c(mean income p50 income)
```

educgr	mean(income)	med(income)
Less than high school	22750.89	18000
High school	34055.76	27000
Some college	39607.05	30300
College	67654.84	50000
Graduate school	98541.49	72000



# Mean income by age, U.S. population, 2018

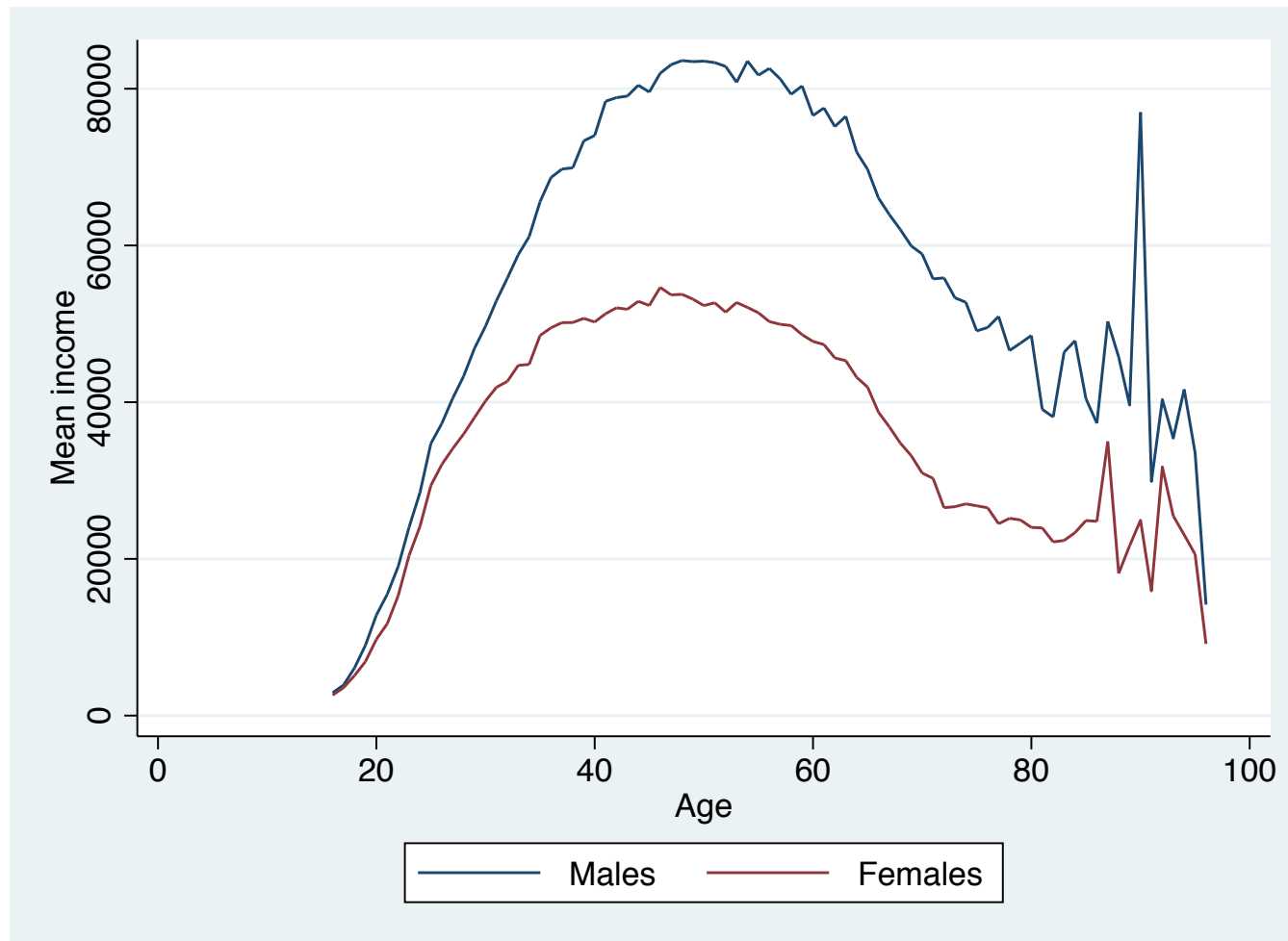


Source: 2018 American Community Survey.

Obs.: Only people with some wage and salary income are included (different than zero).



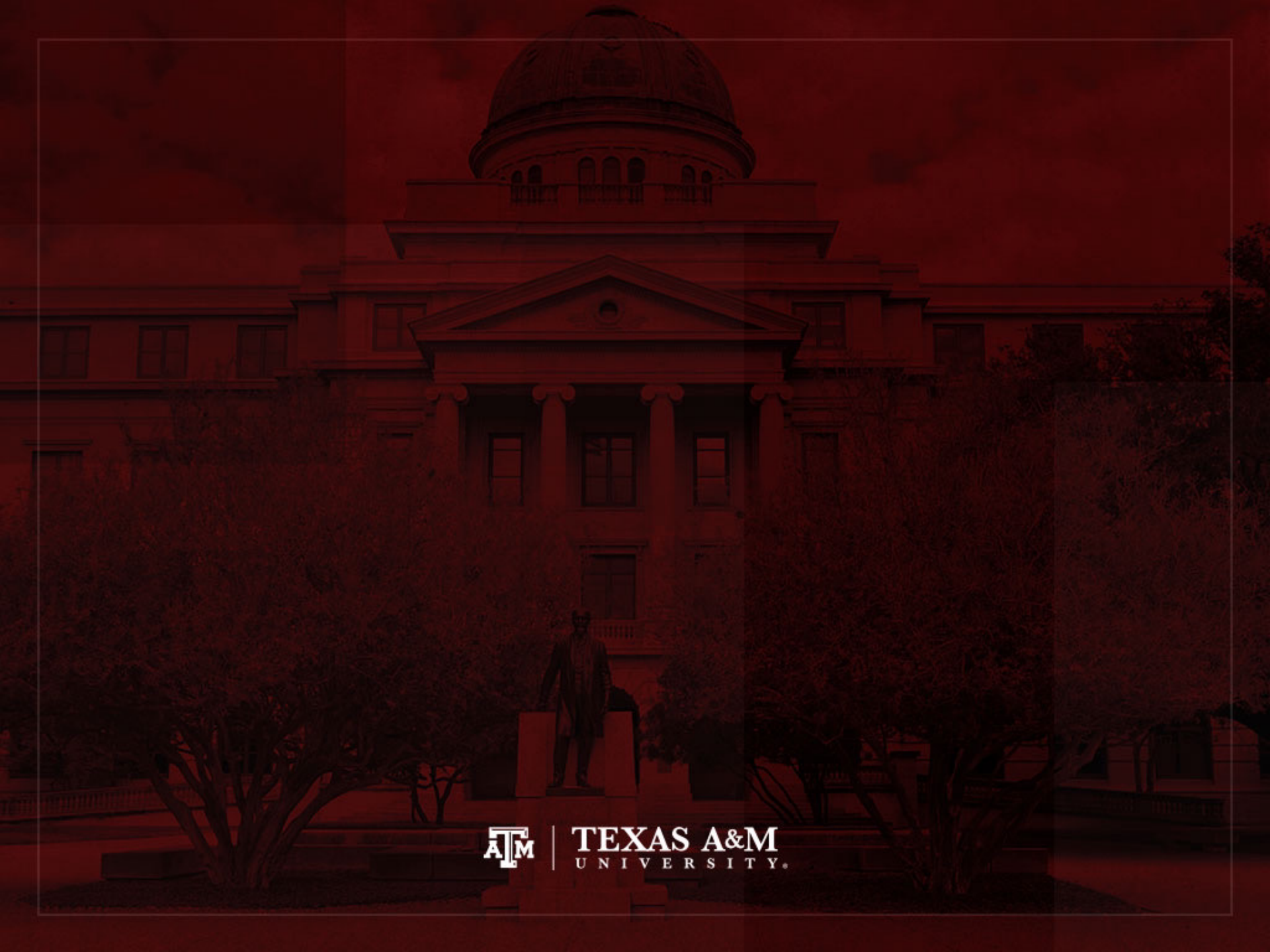
# Mean income by age and sex, U.S. population, 2018



Source: 2018 American Community Survey.

Obs.: Only people with some wage and salary income are included (different than zero).

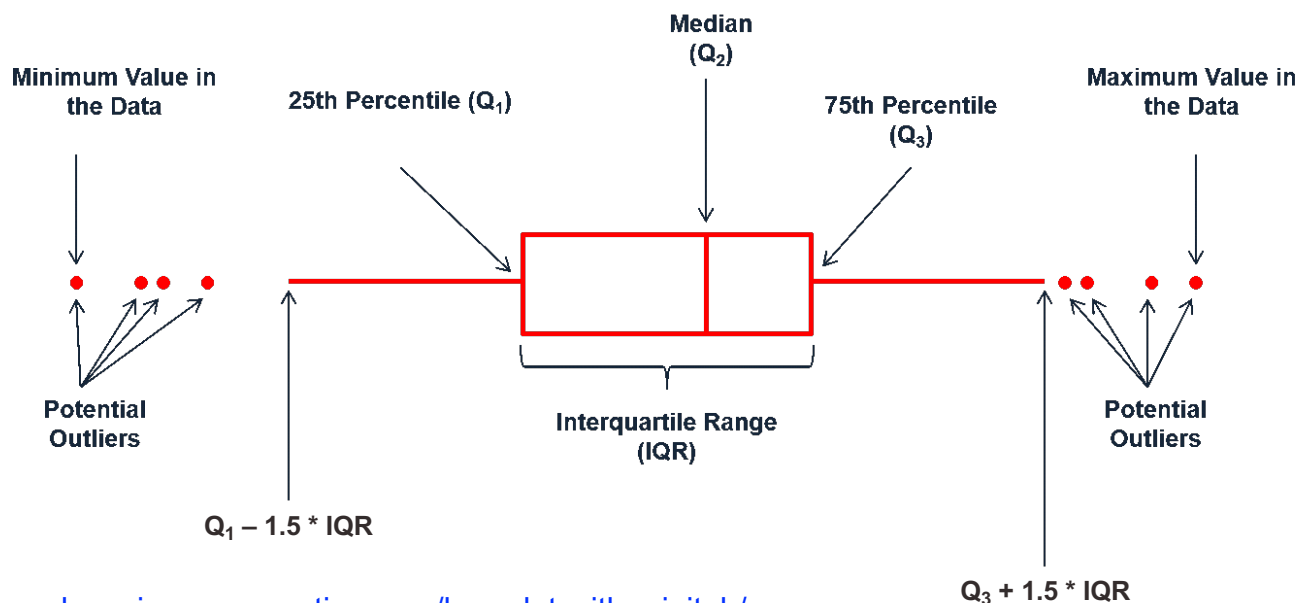




TEXAS A&M  
UNIVERSITY.

# Boxplots

- Boxplot is also known as "box and whiskers plot"
  - It provides a way to visualize and analyze dispersion
  - Useful when comparing distributions
  - It uses median, range, interquartile range, outliers
  - Easier to read all this information than in tables



Source: <https://www.leansigmacorporation.com/box-plot-with-minitab/>

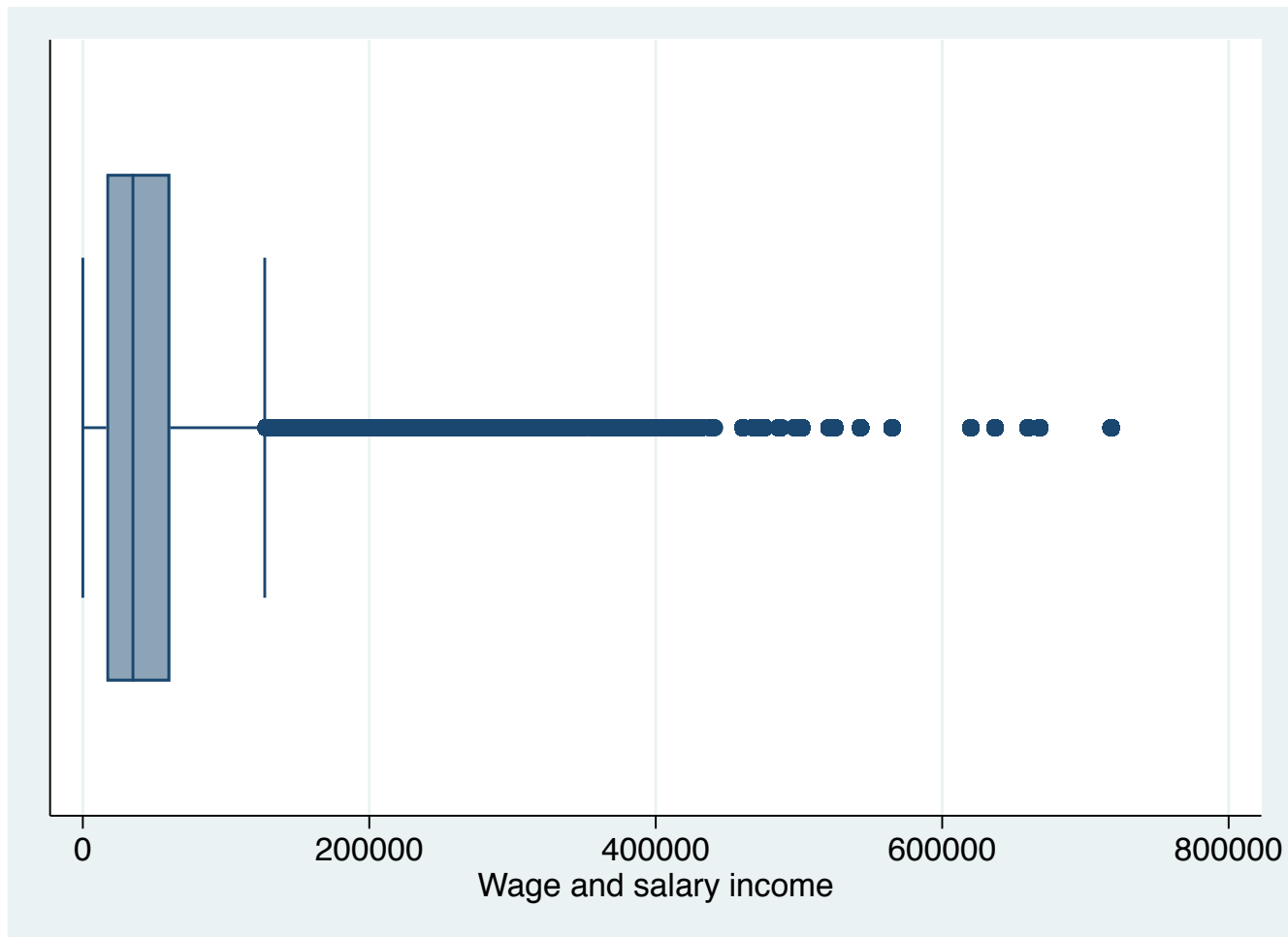
# Example: 2018 ACS in Stata

- Generate box plot for respondents' wage and salary income

```
graph hbox income if income!=0 [fweight=perwt],  
        ytitle(Wage and salary income)
```

# Edited figure

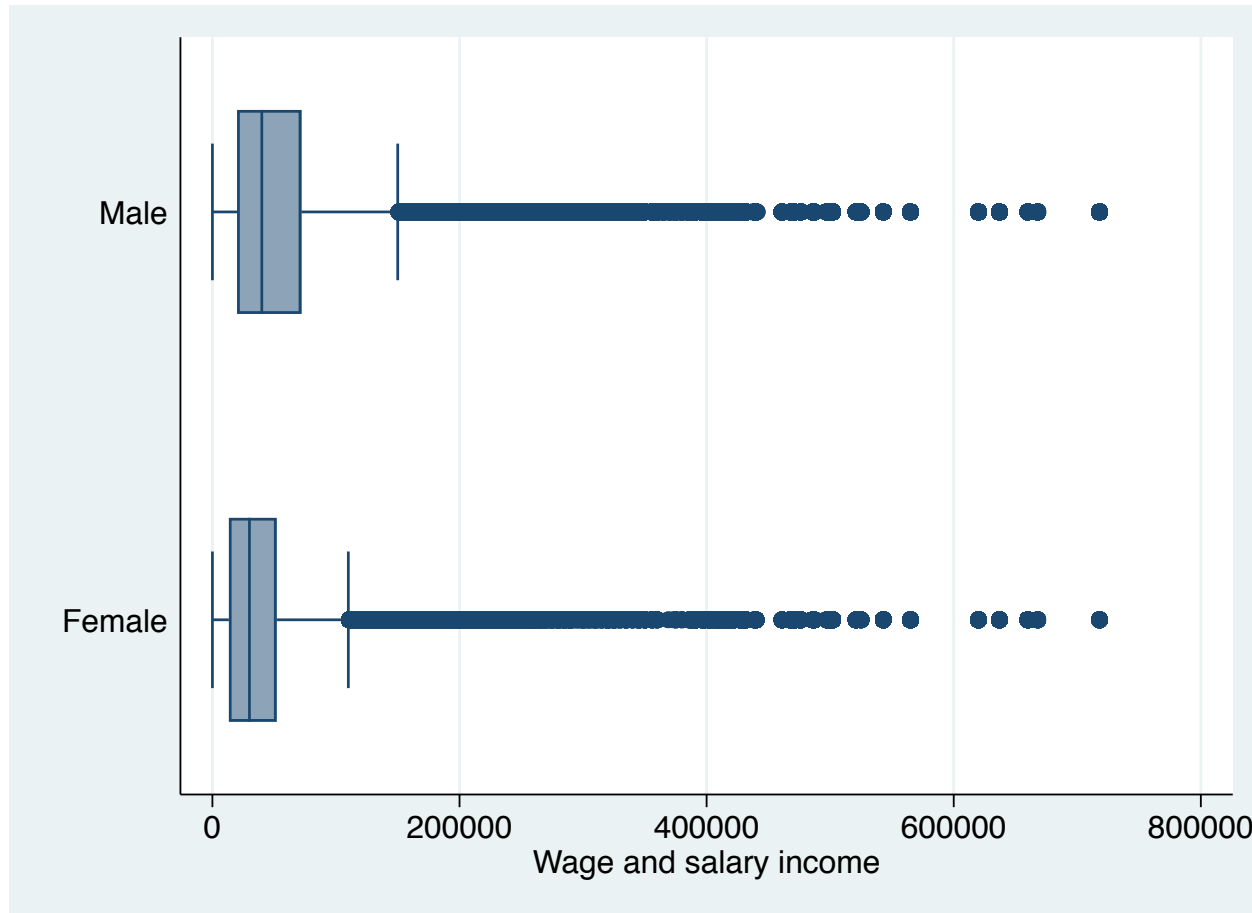
**Figure 1. Distribution of respondents' wage and salary income, U.S. population, 2018**



Source: 2018 American Community Survey.

# Income by sex, 2018

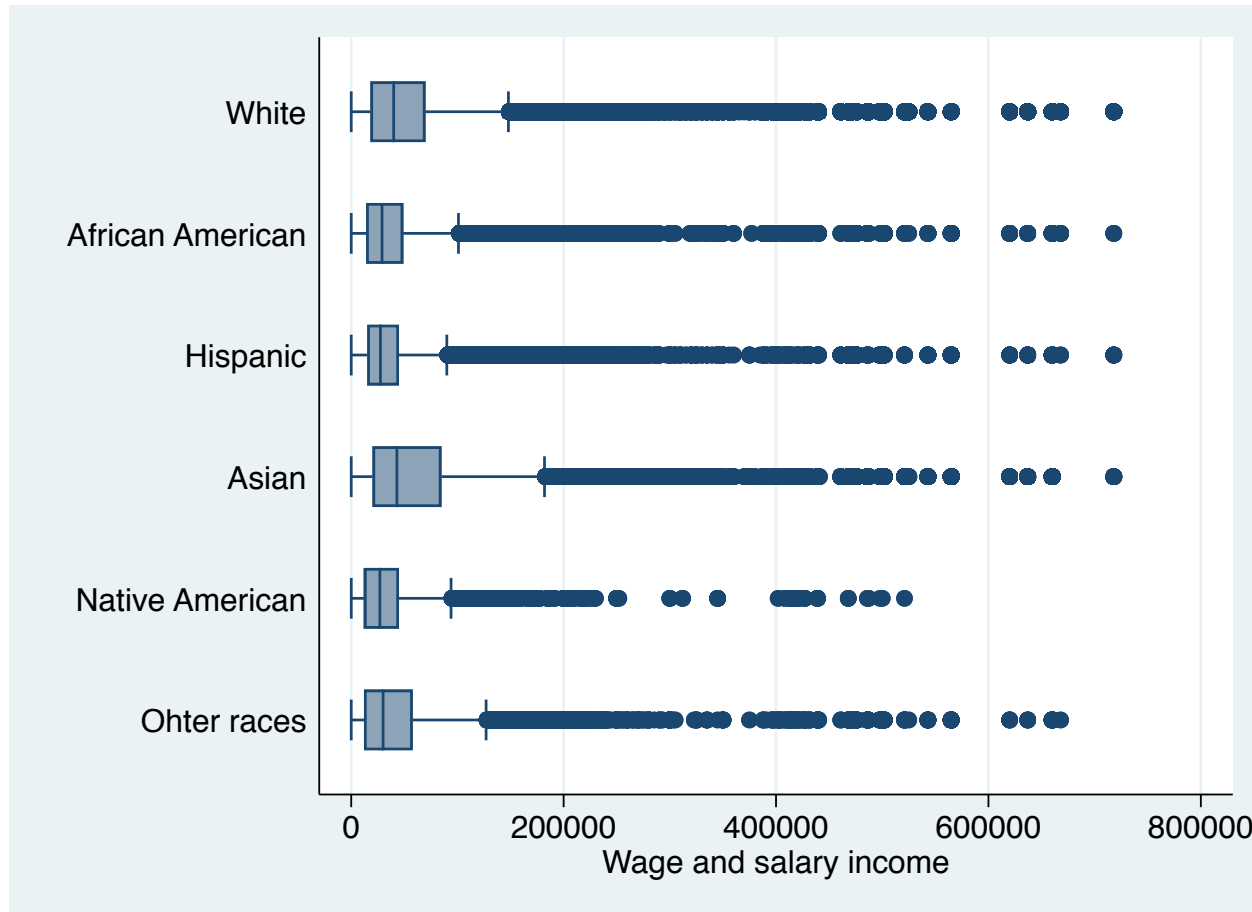
```
graph box income if income!=0 [fweight=perwt],  
    over(female) ytitle(Wage and salary income)
```



Source: 2018 American Community Survey.

# Income by race/ethnicity, 2018

```
graph box income if income!=0 [fweight=perwt],  
    over(raceth) ytitle(Wage and salary income)
```



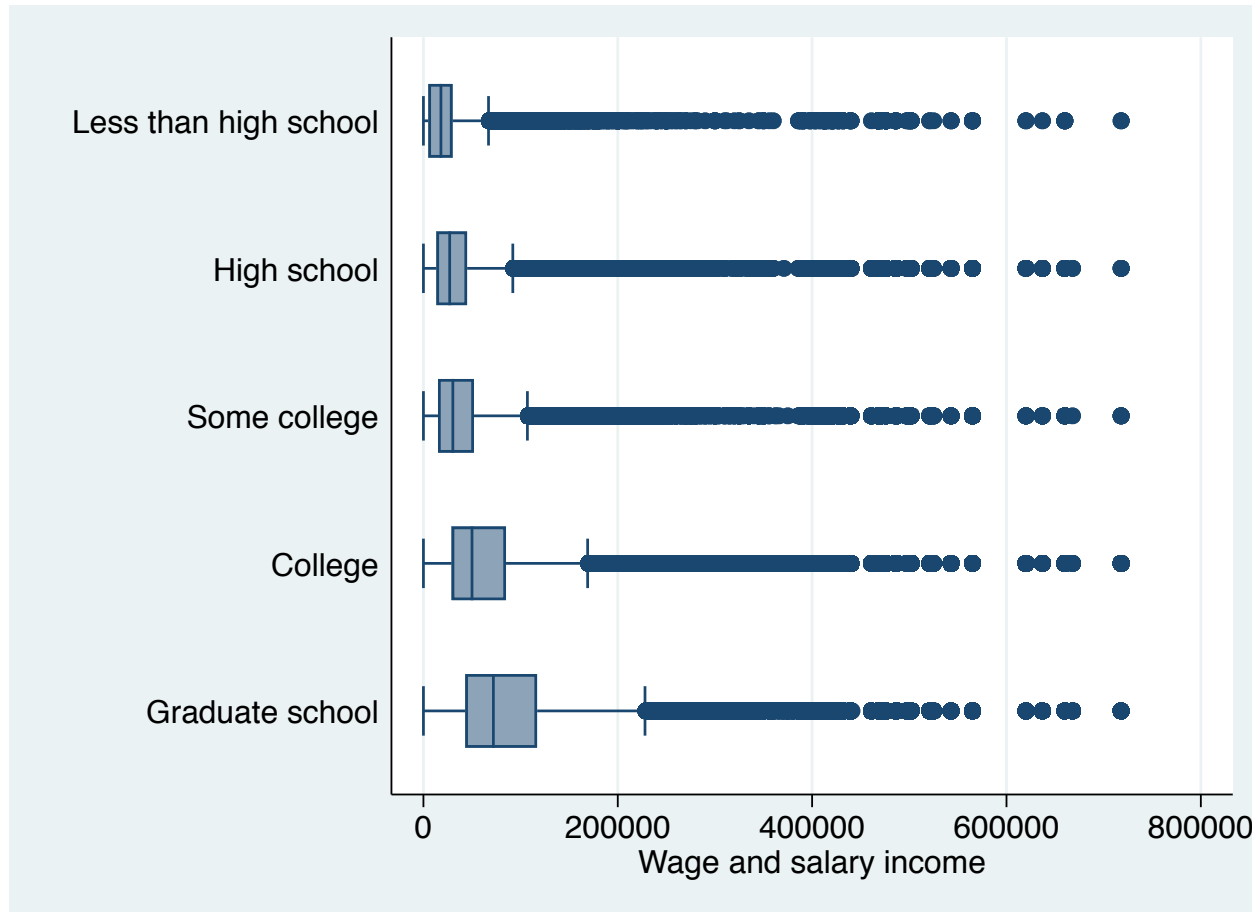
Source: 2018 American Community Survey.





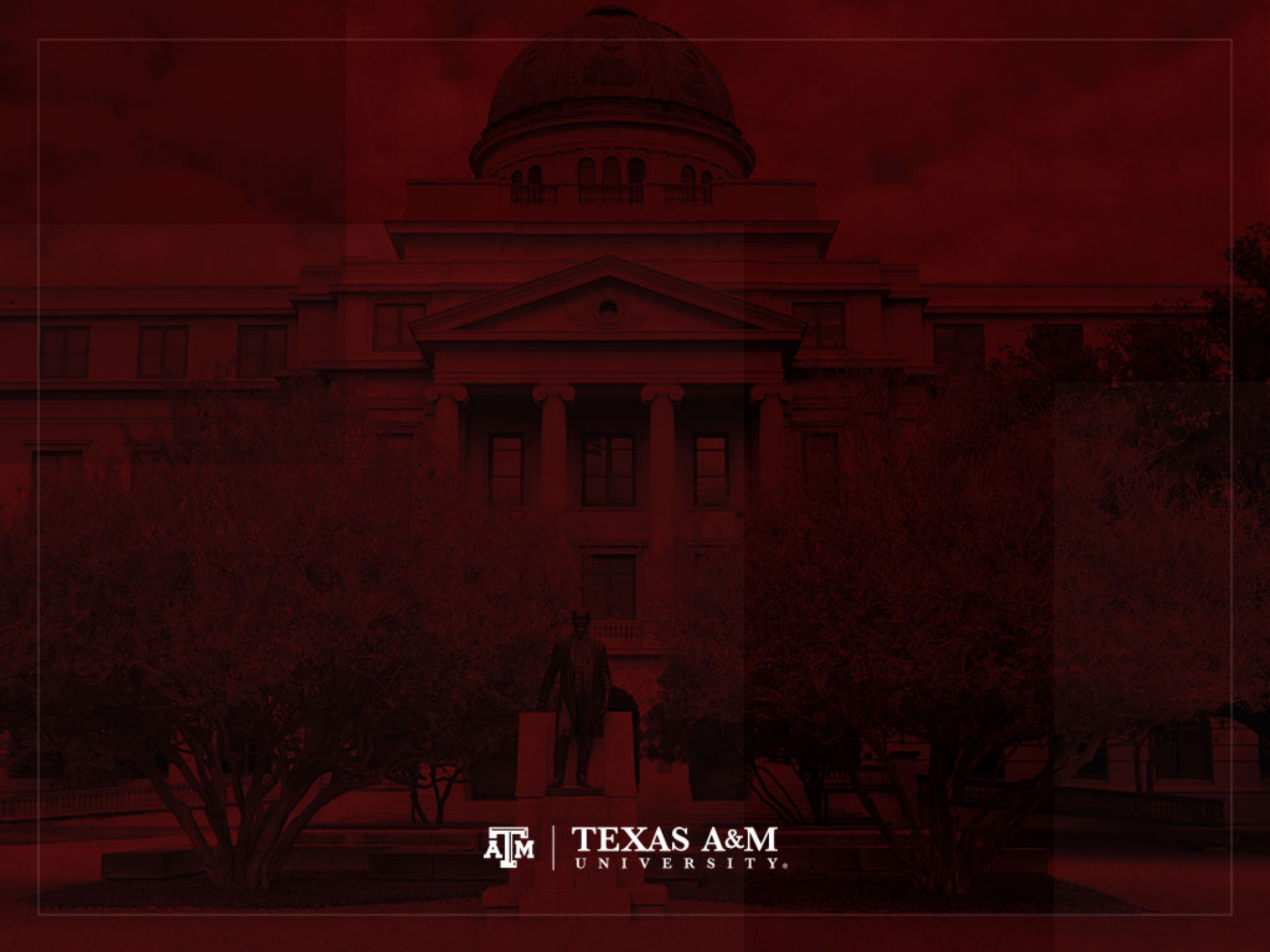
# Income by education, 2018

```
graph box income if income!=0 [fweight=perwt],  
over(educgr) ytitle(Wage and salary income)
```



Source: 2018 American Community Survey.





TEXAS A&M  
UNIVERSITY.

# Age-sex structure in Stata

```
***Generate five-year age groups variable - automatically
egen age5y = cut(age), at(0,5,10,15,20,25,30,35,40,45,50,55,60,65,70,75,80,85,100)
table age5y, contents(min age max age count age)

***Generate male variable (opposite of female variable)
gen male=!female
tab male female, m nolabel

***Generate variables with male and female totals by five-year age groups
sort age5y
by age5y: egen maletotal=total(male)
by age5y: egen femaletotal=total(female)

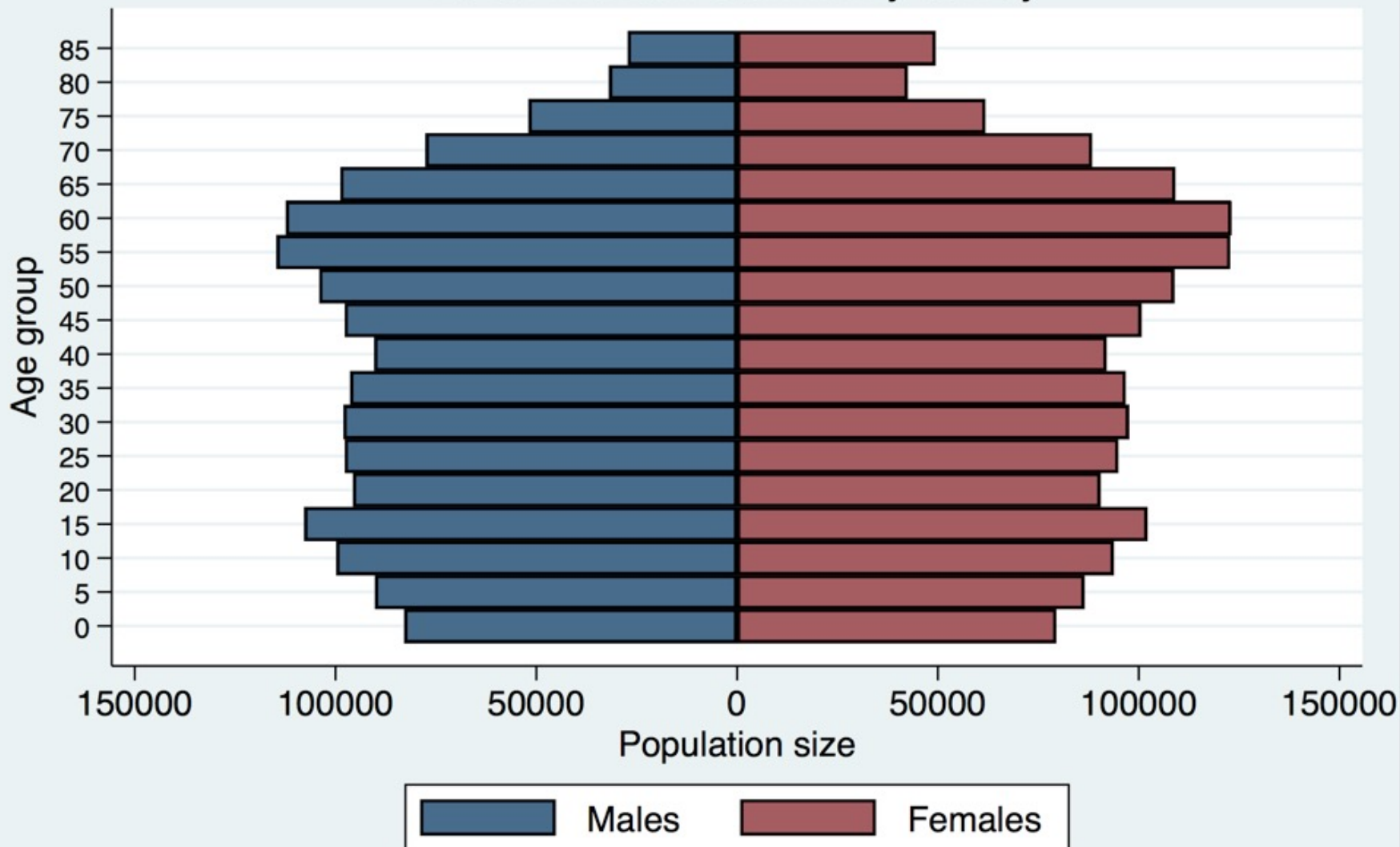
***Replace male total by negative value
replace maletotal=-maletotal

***Age-sex structure
twoway bar maletotal age5y [fweight=perwt], horizontal barwidth(5) fcolor(navy) lcolor(black) lwidth(medium) || ///
bar femaletotal age5y [fweight=perwt], horizontal barwidth(5) fcolor(maroon) lcolor(black) lwidth(medium) ///
legend(label(1 Males) label(2 Females)) ///
ylabel(0(5)85, angle(horizontal) valuelabel labsize(*.8)) ///
yttitle("Age group") ///
xlabel(-150000 "150000" -100000 "100000" -50000 "50000" 0 50000 100000 150000) ///
xttitle("Population size") ///
title("Age-sex structure, United States") ///
subtitle("2018 American Community Survey")
```



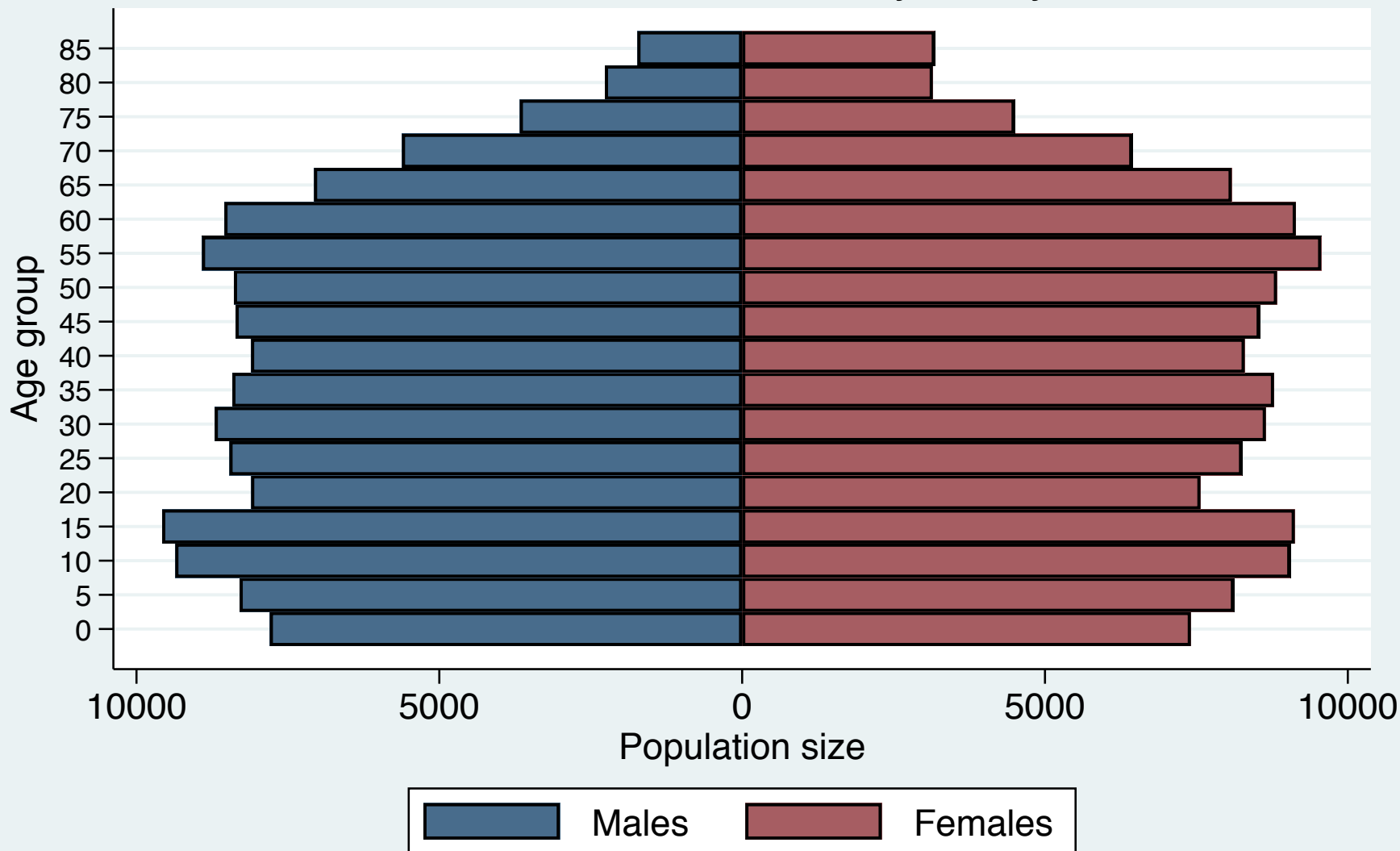
# Age-sex structure, United States

## 2018 American Community Survey



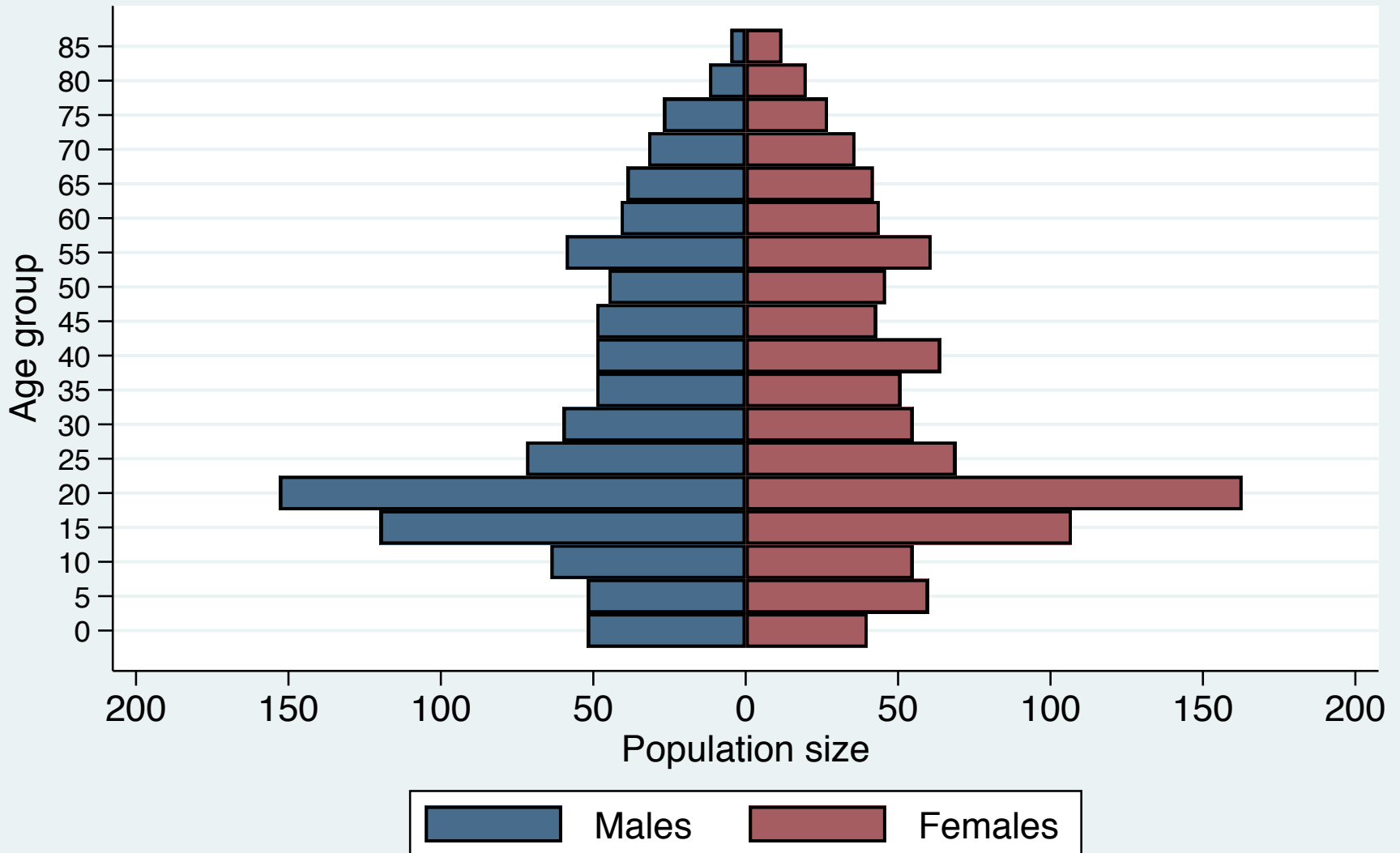
# Age-sex structure, Texas

## 2018 American Community Survey



# Age-sex structure, Brazos county

## 2018 American Community Survey





# Stata practice time

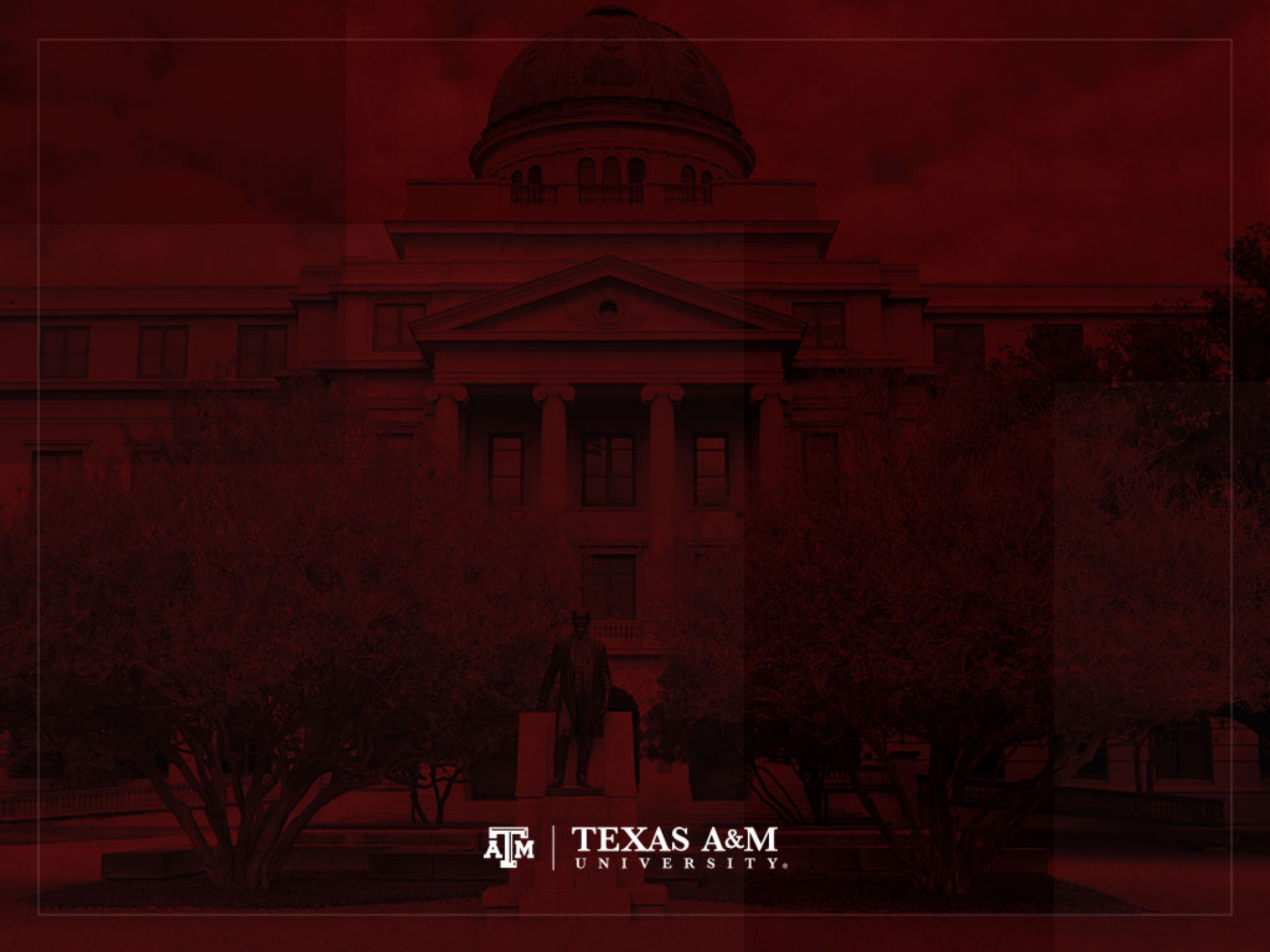
- Additional material on introduction to social statistics using Stata

<http://www.ernestoamaral.com/stata2020a.html>

- You can run all Stata commands that were used in this lecture using this DO-file

<http://www.ernestoamaral.com/docs/Stata2020a/Stata02.txt>





TEXAS A&M  
UNIVERSITY.