

Conditional Descriptive Statistics

Prottoy A. Akbar

Principles of Empirical Analysis (ECON-A3000)
Lecture 3

- Data, descriptive statistics and causality
 - ① introduction, data
 - ② samples and descriptive statistics
 - ③ **today: conditional descriptive statistics**
 - ④ causality and randomization
 - ⑤ statistical inference
 - ⑥ revealed preferences in observed data
- Quasi-experimental methods

- Data, descriptive statistics and causality
 - ① introduction, data
 - ② samples and descriptive statistics
 - ③ **today: conditional descriptive statistics**
 - ④ causality and randomization
 - ⑤ statistical inference
 - ⑥ revealed preferences in observed data
- Quasi-experimental methods
- **Today's learning objectives.** After this lecture you should understand how to
 - ① characterize conditional distributions
 - ② characterize (linear) relationships between variables
 - ③ apply them to interpret data on income distributions

Conditional descriptive statistics

- Conditional descriptives are statistics of a variables *conditional* on another variables
 - The most important: **conditional expectation**

$$\mathbb{E}[Y|X = x]$$

i.e. expectation of random variable Y when another random variable X takes value x

- Conditional descriptives are statistics of a variables *conditional* on another variables

- The most important: **conditional expectation**

$$\mathbb{E}[Y|X = x]$$

i.e. expectation of random variable Y when another random variable X takes value x

- empirical counterpart: conditional sample average
- All conditional descriptive statistics follow from the **joint distribution** of two or more variables

Summary for variables: earn
by categories of: edul

edul	mean	N
Less/unknown	15527	1807
Secodary	22076	2720
Bachelor	32644	1080
Master	42292	346
Lis./PhD	57950	20
Total	23297	5973

Source: FLEED teaching data
tabstat earn, by(edul) stat(mean N)
alternatively try: tabulate edul, sum(earn)
(see the full code at course website)

Joint distribution (review)

- A simple, yet efficient way to display (small) data of two variables is **cross tabulation**
 - ① the no. rows = no. values that Y can take
 - ② the no. columns = no. values that X can take
 - ③ the cells report no. observations with value (y, x)
- Alternatively, cross tabulation cells may report the share of observations with value (y, x)
- This is the empirical counterpart of the **joint density function**

$$f_{XY}(x, y) = \mathbb{P}(X = x, Y = y)$$

i.e., the probability that random variable X takes the value x *and* that random variable Y takes the value y

edul	woman	
	0	1
Less/unknown	18.07	14.32
Secondary	22.90	21.03
Bachelor	7.03	10.43
Master	2.90	2.96
Lis./PhD	0.27	0.10
		100.00

Source: FLEED teaching data
tabulate edul woman, cell nofreq

Marginal distribution

- The marginal distribution of Y is defined as

$$f_Y(y) = \sum_{x \in X} f_{XY}(x, y)$$

- This is just probability of Y when not taking the value of X into account

edul	woman		Total
	0	1	
Less/unknown	18.07	14.32	32.38
Secodary	22.90	21.03	43.93
Bachelor	7.03	10.43	17.46
Master	2.90	2.96	5.86
Lis./PhD	0.27	0.10	0.37
			100.00

Source: FLEED teaching data
tabulate edul woman, cell nofreq

Marginal distribution

- The marginal distribution of Y is defined as

$$f_Y(y) = \sum_{x \in X} f_{XY}(x, y)$$

- This is just probability of Y when not taking the value of X into account
- Similarly, the marginal distribution of X is

$$f_X(x) = \sum_{y \in Y} f_{XY}(x, y)$$

edul	woman		Total
	0	1	
Less/unknown	18.07	14.32	32.38
Secondary	22.90	21.03	43.93
Bachelor	7.03	10.43	17.46
Master	2.90	2.96	5.86
Lis./PhD	0.27	0.10	0.37
Total	51.17	48.83	100.00

Source: FLEED teaching data
tabulate edul woman, cell nofreq

Conditional distribution

- The conditional distribution of Y is defined as

$$f_{Y|X}(y|x) = \frac{f_{XY}(x, y)}{f_X(x)}$$

i.e., the probability that Y takes value y conditional that X takes value x

- Example: Probability that a working age woman living in Finland in 2010 doesn't have a secondary school degree?

edul	woman		Total
	0	1	
Less/unknown	18.07	14.32	32.38
Secondary	22.90	21.03	43.93
Bachelor	7.03	10.43	17.46
Master	2.90	2.96	5.86
Lis./PhD	0.27	0.10	0.37
Total	51.17	48.83	100.00

Source: FLEED teaching data
tabulate edul woman, cell nofreq

Conditional distribution

- The conditional distribution of Y is defined as

$$f_{Y|X}(y|x) = \frac{f_{XY}(x, y)}{f_X(x)}$$

i.e., the probability that Y takes value y conditional that X takes value x

- Example: Probability that a working age woman living in Finland in 2010 doesn't have a secondary school degree?
 - $\hat{P}(Y = b|X = w) \approx ??$
 - where the "hats" indicate that we are using **estimates** of the population probabilities $\mathbb{P}(\cdot)$

edul	woman		Total
	0	1	
Less/unknown	18.07	14.32	32.38
Secondary	22.90	21.03	43.93
Bachelor	7.03	10.43	17.46
Master	2.90	2.96	5.86
Lis./PhD	0.27	0.10	0.37
Total	51.17	48.83	100.00

Source: FLEED teaching data
tabulate edul woman, cell nofreq

- Let's get back to conditional expectation. When Y is discrete^a, the **conditional expectation function (CEF)** is

$$\mathbb{E}[Y|X = x] = \sum t f_{Y|X}(t|X = x)$$

i.e. **population average of Y holding X fixed**

- in other words: weighted average of Y , where the weight for of each value of Y is the share of sub-population (for whom $X = x$) with this value of Y
- X can also be a vector, i.e., can include many conditioning variables

Summary for variables: earn
by categories of: edul

edul	mean	N
Less/unknown	15527	1807
Secodary	22076	2720
Bachelor	32644	1080
Master	42292	346
Lis./PhD	57950	20
Total	23297	5973

Source: FLEED teaching data
tabstat earn, by(edul) stat(mean N)

^aContinuous version: $\mathbb{E}[Y|X = x] = \int t f_{Y|X}(t|X = x) d(t)$

Conditional expectation

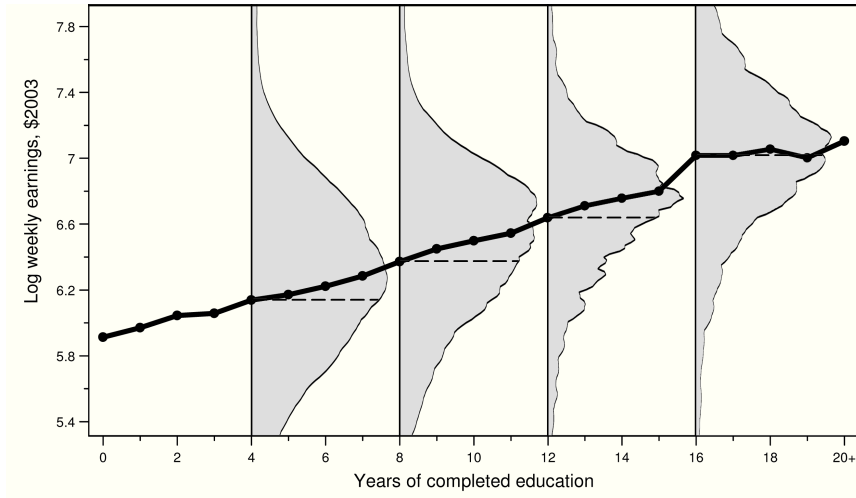


Figure 3.1.1: Raw data and the CEF of average log weekly wages given schooling. The sample includes white men aged 40-49 in the 1980 IPUMS 5 percent file.

Source: Angrist and Pischke (2009).

Example:
Widening U.S. income distribution

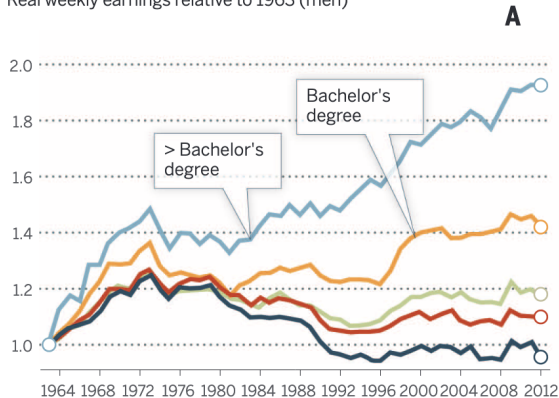
- We now have tools to understand the basic results of the income distribution literature
 - group averages
 - changes over the entire distribution
 - extras: top percent shares, social mobility
- Much of this research is based on tax data
 - available over long time periods and many countries, but earlier periods limited to the top (historically, only the rich paid taxes)
 - tax records never capture all income → ongoing work to deal with the missing parts
- Lot's of work also based on surveys, particularly the Labor Force Survey



Source: [The Economist](#), 28 Nov 2019

Changes in real wage levels of full-time U.S. workers by sex and education, 1963–2012

Real weekly earnings relative to 1963 (men)



Real weekly earnings relative to 1963 (women)

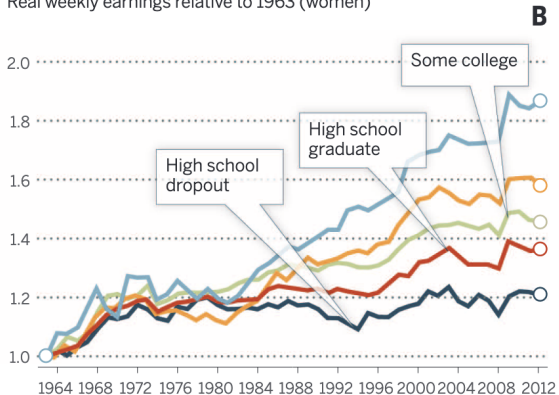
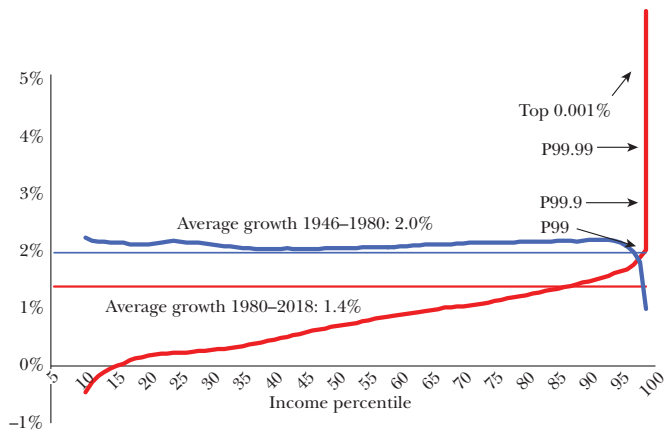


Fig. 6. Change in real wage levels of full-time workers by education, 1963–2012. (A) Male workers, (B) female workers. Data and sample construction are as in Fig. 3.

Source: Autor (2014), Science.

- Estimates over time for $\mathbb{E}[w|E = e, G = G]$, where w is weekly wage, E education level and G is gender. Wages are divided by 1963 group-specific average wages.

Average Annual Income Growth Rates

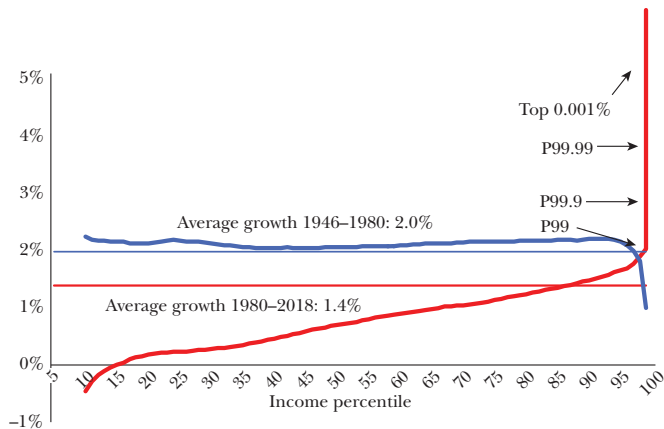


Source: Saez and Zucman (2019b).

Note: This figure depicts the annual real pre-tax income growth per adult for each percentile in the 1946-1980 period (in blue) and 1980-2018 period (in red). From 1946 to 1980, growth was evenly distributed with all income groups growing at the average 2 percent annual rate (except the top 1 percent which grew slower). From 1980 to 2018, growth has been unevenly distributed with low growth for bottom income groups, mediocre growth for the middle class, and explosive growth at the top.

Source: Saez and Zucman (2020), Journal of Economic Perspectives.

Average Annual Income Growth Rates



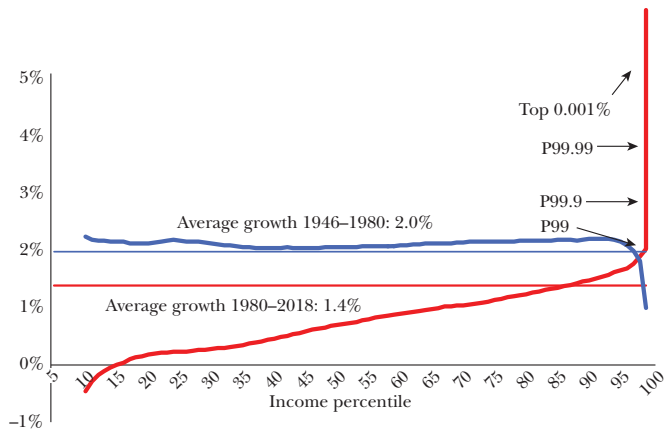
- **1946–1980:** roughly 2% annual income growth across the distribution among "the 99%"

Source: Saez and Zucman (2019b).

Note: This figure depicts the annual real pre-tax income growth per adult for each percentile in the 1946–1980 period (in blue) and 1980–2018 period (in red). From 1946 to 1980, growth was evenly distributed with all income groups growing at the average 2 percent annual rate (except the top 1 percent which grew slower). From 1980 to 2018, growth has been unevenly distributed with low growth for bottom income groups, mediocre growth for the middle class, and explosive growth at the top.

Source: Saez and Zucman (2020), Journal of Economic Perspectives.

Average Annual Income Growth Rates



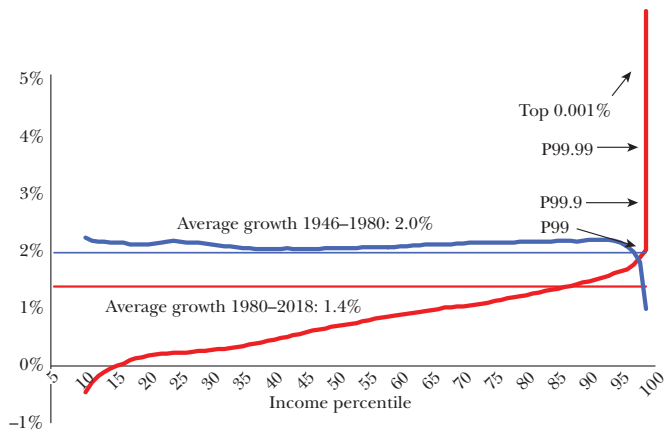
- **1946–1980:** roughly 2% annual income growth across the distribution among "the 99%"
- **1980–2018:** income growth faster among the more wealthy even among "the 99%"; the very top very different than the rest

Source: Saez and Zucman (2019b).

Note: This figure depicts the annual real pre-tax income growth per adult for each percentile in the 1946–1980 period (in blue) and 1980–2018 period (in red). From 1946 to 1980, growth was evenly distributed with all income groups growing at the average 2 percent annual rate (except the top 1 percent which grew slower). From 1980 to 2018, growth has been unevenly distributed with low growth for bottom income groups, mediocre growth for the middle class, and explosive growth at the top.

Source: Saez and Zucman (2020), Journal of Economic Perspectives.

Average Annual Income Growth Rates



- **1946–1980:** roughly 2% annual income growth across the distribution among "the 99%"
- **1980–2018:** income growth faster among the more wealthy even among "the 99%"; the very top very different than the rest
- Next: How is this figure constructed?

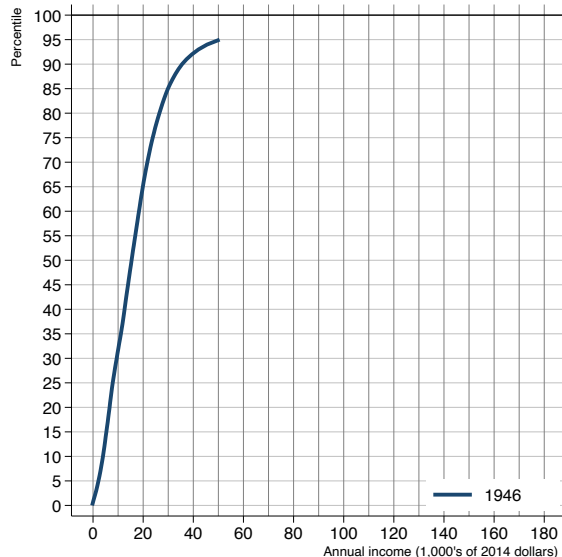
Source: Saez and Zucman (2019b).

Note: This figure depicts the annual real pre-tax income growth per adult for each percentile in the 1946–1980 period (in blue) and 1980–2018 period (in red). From 1946 to 1980, growth was evenly distributed with all income groups growing at the average 2 percent annual rate (except the top 1 percent which grew slower). From 1980 to 2018, growth has been unevenly distributed with low growth for bottom income groups, mediocre growth for the middle class, and explosive growth at the top.

Source: Saez and Zucman (2020), Journal of Economic Perspectives.

The U.S. income distribution, 1962–2014, bottom 95 percentiles

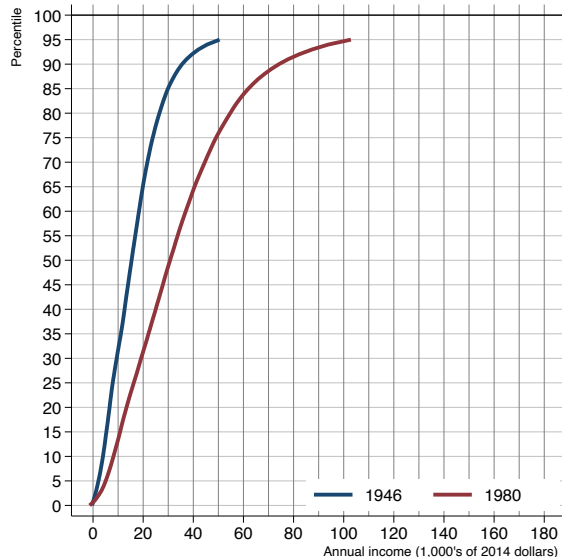
- Let's start with the CDF of income distribution in 1946
 - 90/10 percentile ratio: $\frac{35.5}{3.8} = 9.0$



Source: Piketty, Saez, Zucman (2018) data appendix

The U.S. income distribution, 1962–2014, bottom 95 percentiles

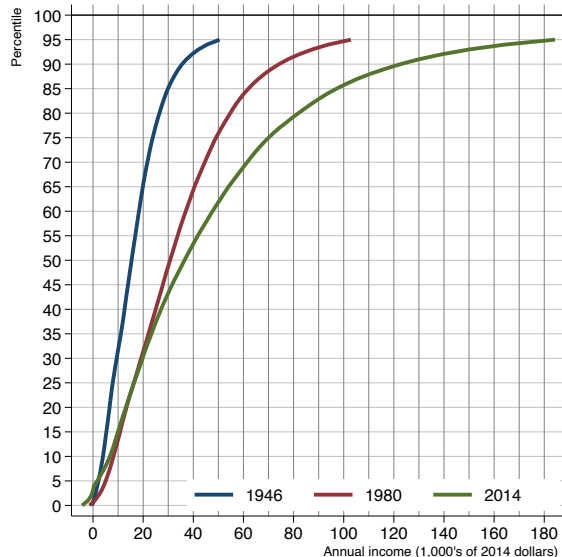
- Let's start with the CDF of income distribution in 1946
 - 90/10 percentile ratio: $\frac{35.5}{3.8} = 9.0$
- Adding the CDF for 1980 income
 - 90/10 percentile ratio: $\frac{74.2}{8.1} = 9.1$



Source: Piketty, Saez, Zucman (2018) data appendix

The U.S. income distribution, 1962–2014, bottom 95 percentiles

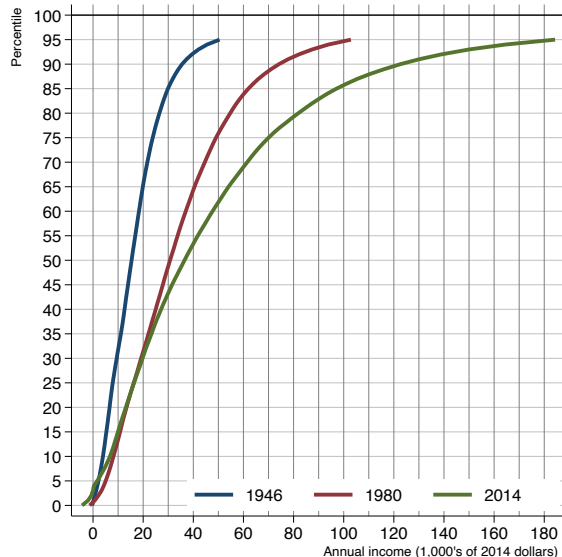
- Let's start with the CDF of income distribution in 1946
 - 90/10 percentile ratio: $\frac{35.5}{3.8} = 9.0$
- Adding the CDF for 1980 income
 - 90/10 percentile ratio: $\frac{74.2}{8.1} = 9.1$
- Adding the CDF for 2014 income
 - 90/10 percentile ratio: $\frac{122.6}{6.7} = 18.2$



Source: Piketty, Saez, Zucman (2018) data appendix

The U.S. income distribution, 1962–2014, bottom 95 percentiles

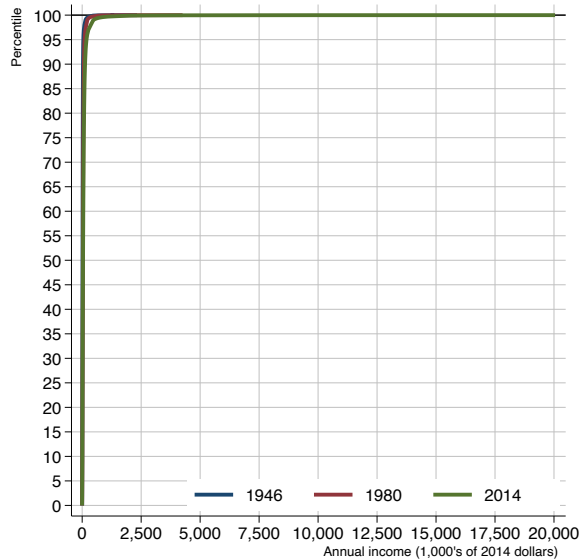
- Let's start with the CDF of income distribution in 1946
 - 90/10 percentile ratio: $\frac{35.5}{3.8} = 9.0$
- Adding the CDF for 1980 income
 - 90/10 percentile ratio: $\frac{74.2}{8.1} = 9.1$
- Adding the CDF for 2014 income
 - 90/10 percentile ratio: $\frac{122.6}{6.7} = 18.2$
- Horizontal distance btw the CDFs = dollar change for each percentile
 - these are not the same *people*; we are comparing percentiles



Source: Piketty, Saez, Zucman (2018) data appendix

The U.S. income distribution, 1962–2014, full distribution

- CDFs for very skewed distributions are uninformative

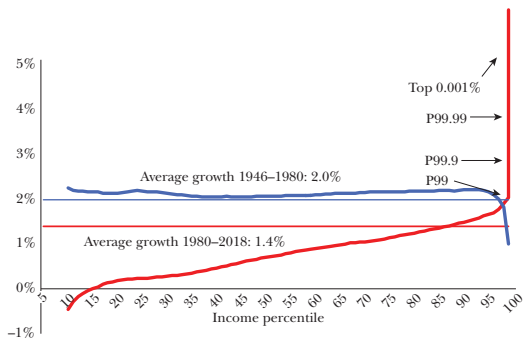


Source: Piketty, Saez, Zucman (2018) data appendix

The U.S. income distribution, 1962–2014, full distribution

- CDFs for very skewed distributions are uninformative ... but changes can nevertheless be made visible

Average Annual Income Growth Rates



Source: Saez and Zucman (2019b).

Note: This figure depicts the annual real pre-tax income growth per adult for each percentile in the 1946–1980 period (in blue) and 1980–2018 period (in red). From 1946 to 1980, growth was evenly distributed with all income groups growing at the average 2 percent annual rate (except for the top 1 percent which grew slower). From 1980 to 2018, growth has been unevenly distributed with low growth for bottom income groups, mediocre growth for the middle class, and explosive growth at the top.

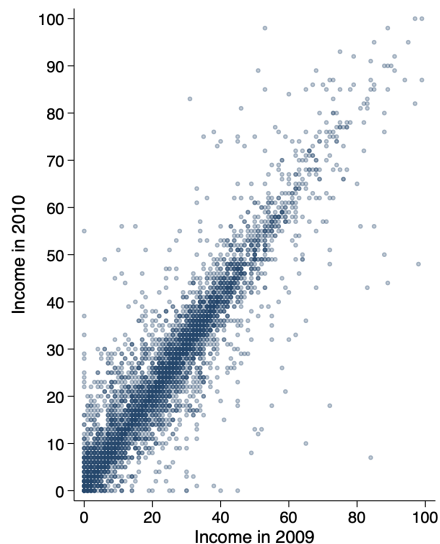
Source: Saez and Zucman (2020), Journal of Economic Perspectives.

Correlation

- Conditional expectation is a powerful way to detect how variables are associated with each other

Scatter plot

- Conditional expectation is a powerful way to detect how variables are associated with each other
- An alternative approach is to show all observations and plot two variables against each other
- Example: persistence of income over time
 - **scatter plot**: each dot in this graph shows each individual's income in 2009 and 2010

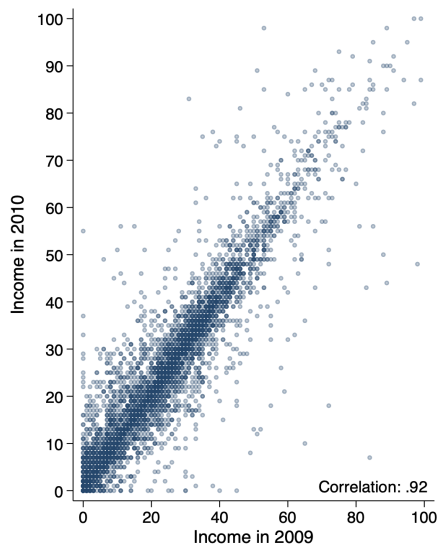


Source: FLEED teaching data

```
scatter earn earn.t1, mcolor(navy%25) msize(vsmall)
```

Scatter plot

- Conditional expectation is a powerful way to detect how variables are associated with each other
- An alternative approach is to show all observations and plot two variables against each other
- Example: persistence of income over time
 - **scatter plot**: each dot in this graph shows each individual's income in 2009 and 2010
- The best known descriptive statistic to characterize how two variables' values are aligned is **correlation**
 - here, the correlation is 0.92
 - next: what does that mean?



Source: FLEED teaching data
`scatter earn_t1, mcolor(navy%25) msize(vsmall)`

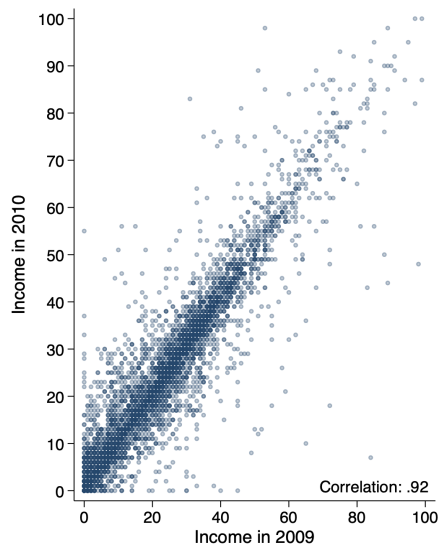
- To get to correlation, we need to first define the **covariance** of Y and X

$$\text{Cov}(X, Y) = \mathbb{E}[X - \mathbb{E}(X)]\mathbb{E}[Y - \mathbb{E}(Y)]$$

... and its empirical counterpart

$$\widehat{\text{Cov}}(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

- Here, the covariance is 256.6
 - a hard number to interpret
 - (unit of measurement is the unit of X times the unit of Y)



Source: FLEED teaching data
`scatter earn_t1, mcolor(navy%25) msize(vsmall)`

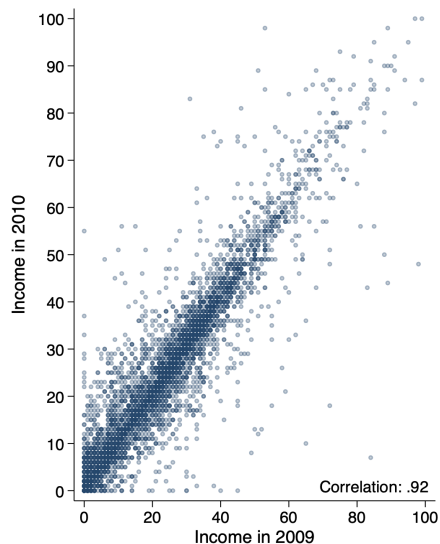
Correlation

- Pearson correlation coefficient is a scaled covariance

$$\text{Cor}(X, Y) = \rho_{X,Y} = \frac{\text{Cov}(X, Y)}{SD(X)SD(Y)}$$

that varies between $-1 \leq \text{Cor}(X, Y) \leq 1$

- just makes the number easier to interpret



Source: FLEED teaching data
`scatter earn_t1, mcolor(navy%25) msize(vsmall)`

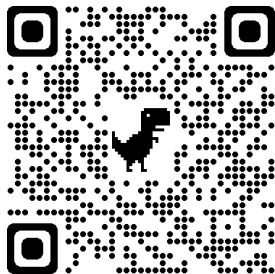
Correlation

- Pearson correlation coefficient is a scaled covariance

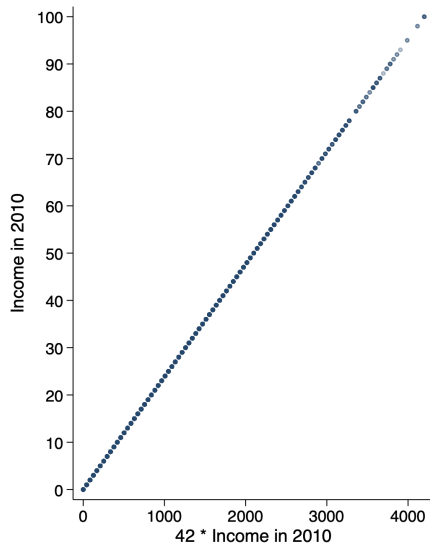
$$\text{Cor}(X, Y) = \rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\text{SD}(X)\text{SD}(Y)}$$

that varies between $-1 \leq \text{Cor}(X, Y) \leq 1$

- More examples
 - correlation = ?



<https://presemo.aalto.fi/empanalysis2024I3>



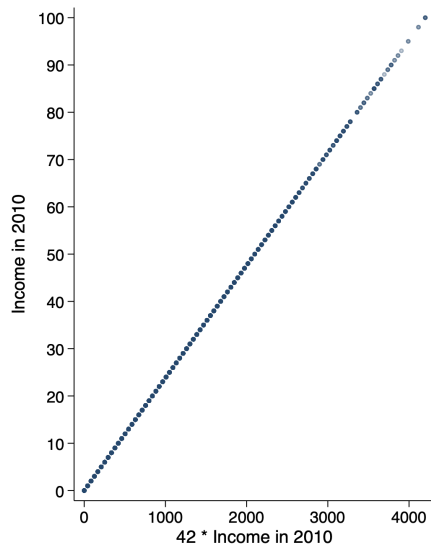
Correlation

- Pearson correlation coefficient is a scaled covariance

$$\text{Cor}(X, Y) = \rho_{X,Y} = \frac{\text{Cov}(X, Y)}{SD(X)SD(Y)}$$

that varies between $-1 \leq \text{Cor}(X, Y) \leq 1$

- More examples
 - correlation 1



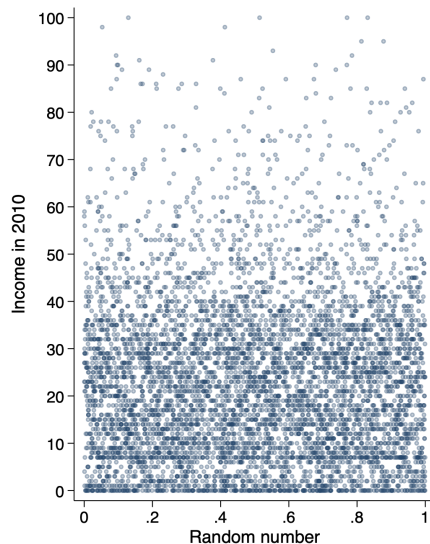
Correlation

- Pearson correlation coefficient is a scaled covariance

$$\text{Cor}(X, Y) = \rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\text{SD}(X)\text{SD}(Y)}$$

that varies between $-1 \leq \text{Cor}(X, Y) \leq 1$

- More examples
 - correlation 1
 - **correlation = ?**



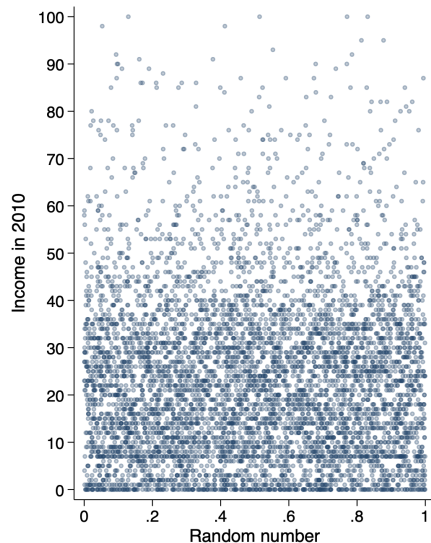
Correlation

- Pearson correlation coefficient is a scaled covariance

$$\text{Cor}(X, Y) = \rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\text{SD}(X)\text{SD}(Y)}$$

that varies between $-1 \leq \text{Cor}(X, Y) \leq 1$

- More examples
 - correlation 1
 - **correlation 0.009**



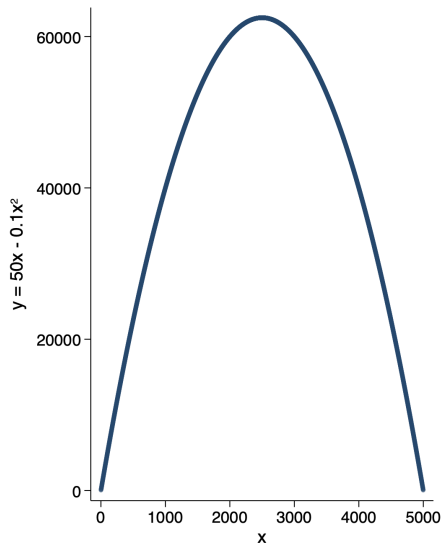
Correlation

- Pearson correlation coefficient is a scaled covariance

$$\text{Cor}(X, Y) = \rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\text{SD}(X)\text{SD}(Y)}$$

that varies between $-1 \leq \text{Cor}(X, Y) \leq 1$

- More examples
 - correlation 1
 - correlation 0.009
 - **correlation = ?**



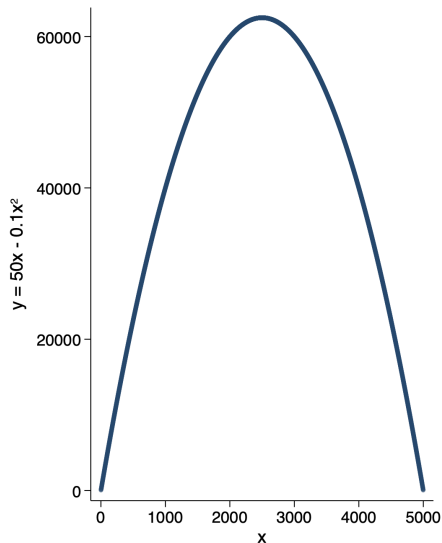
Correlation

- Pearson correlation coefficient is a scaled covariance

$$\text{Cor}(X, Y) = \rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\text{SD}(X)\text{SD}(Y)}$$

that varies between $-1 \leq \text{Cor}(X, Y) \leq 1$

- More examples
 - correlation 1
 - correlation 0.009
 - **correlation 0**



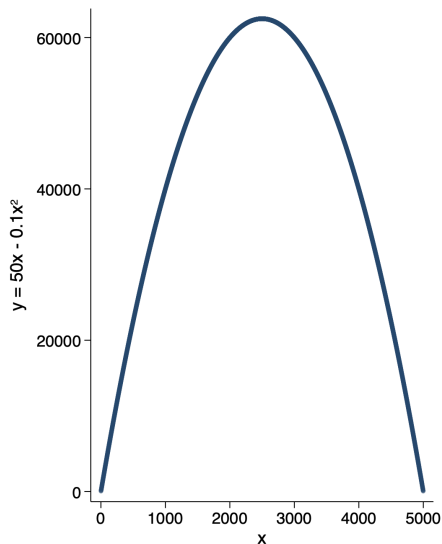
Correlation

- Pearson correlation coefficient is a scaled covariance

$$\text{Cor}(X, Y) = \rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\text{SD}(X)\text{SD}(Y)}$$

that varies between $-1 \leq \text{Cor}(X, Y) \leq 1$

- More examples
 - correlation 1
 - correlation 0.009
 - correlation 0
- Correlation measures a linear dependence
 - the point: possible to have perfect dependence and zero correlation



Regression

- A closely related approach for assessing linear dependence:
bivariate **regression model**

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- A closely related approach for assessing linear dependence: bivariate **regression model**

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- Y is the dependent variable (or outcome)

- A closely related approach for assessing linear dependence: bivariate **regression model**

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- Y is the dependent variable (or outcome)
- X is the independent variable (or regressor)
 - **observed** in data

- A closely related approach for assessing linear dependence: bivariate **regression model**

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- Y is the dependent variable (or outcome)
- X is the independent variable (or regressor)
 - **observed** in data
- ϵ is the residual (or "error term")
 - represents the relevant **unobserved factors**
 - defined to have $\mathbb{E}[\epsilon] = 0$

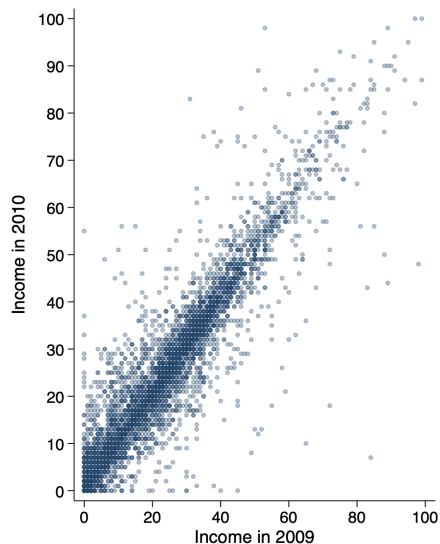
- A closely related approach for assessing linear dependence: bivariate **regression model**

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- Y is the dependent variable (or outcome)
- X is the independent variable (or regressor)
 - **observed** in data
- ϵ is the residual (or "error term")
 - represents the relevant **unobserved factors**
 - defined to have $\mathbb{E}[\epsilon] = 0$
- parameters: β_0 (constant), β_1 (regression coefficient)

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- *Question:* How should we set β_0 and β_1 to best describe the data?

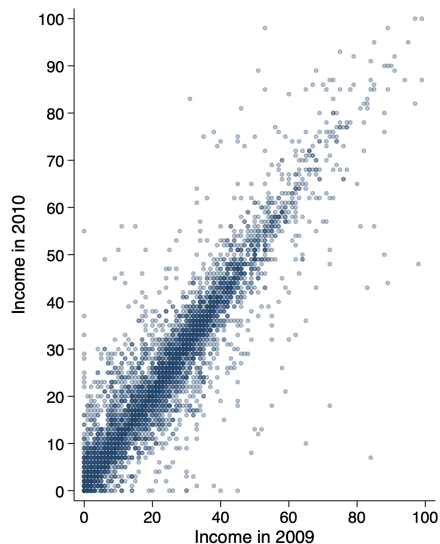


Source: FLEED teaching data
scatter earn earn.t1

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- *Question:* How should we set β_0 and β_1 to best describe the data?
- *One answer:* Ordinary Least Squares (OLS)

$$\arg \min_{\beta_0, \beta_1} \sum_{i=1}^n [Y_i - (\beta_0 + \beta_1 X_i)]^2$$



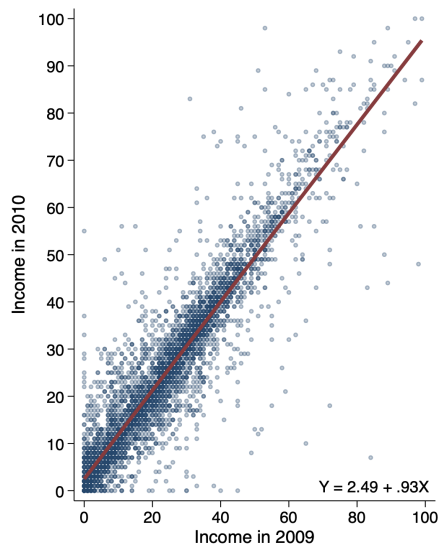
Source: FLEED teaching data
scatter earn earn.t1

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- *Question:* How should we set β_0 and β_1 to best describe the data?
- *One answer:* Ordinary Least Squares (OLS)

$$\arg \min_{\beta_0, \beta_1} \sum_{i=1}^n [Y_i - (\beta_0 + \beta_1 X_i)]^2$$

- In words: let's find the values of β_0 and β_1 that minimize (the square of) the difference between observed data and regression model's prediction



Source: FLEED teaching data
the code is available at the course's website

$$Y = \beta_0 + \beta_1 X + \epsilon$$

- *Question:* How should we set β_0 and β_1 to best describe the data?
- *One answer:* Ordinary Least Squares (OLS)

$$\arg \min_{\beta_0, \beta_1} \sum_{i=1}^n [Y_i - (\beta_0 + \beta_1 X_i)]^2$$

- In words: let's find the values of β_0 and β_1 that minimize (the square of) the difference between observed data and regression model's prediction
 - here, the answer is: $\hat{\beta}_0 = 2.49$, $\hat{\beta}_1 = 0.93$

Source	SS	df	MS	Number of obs	=	5,777
Model	1390738.85	1	1390738.85	F(1, 5775)	=	33626.24
Residual	238846.737	5,775	41.3587423	Prob > F	=	0.0000
				R-squared	=	0.8534
				Adj R-squared	=	0.8534
Total	1629585.58	5,776	282.130468	Root MSE	=	6.4311

earn	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
earn_t1	.9383461	.0051171	183.37	0.000	.9283147 .9483776
_cons	2.487598	.1438088	17.30	0.000	2.205679 2.769518

Source: FLEED teaching data
regress earn earn_t1

- Turns out that correlation and bivariate regression are closely related, namely:

$$\beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

- Turns out that correlation and bivariate regression are closely related, namely:

$$\beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

- Compare to Pearson correlation coefficient:

$$\rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

- Turns out that correlation and bivariate regression are closely related, namely:

$$\beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

- Compare to Pearson correlation coefficient:

$$\rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

- In our example
 - $\hat{\beta}_0 = 2.49, \hat{\beta}_1 = 0.93$
 - $\hat{\rho}_{X,Y} = 0.92$

- Turns out that correlation and bivariate regression are closely related, namely:

$$\beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

- Compare to Pearson correlation coefficient:

$$\rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

- In our example
 - $\hat{\beta}_0 = 2.49$, $\hat{\beta}_1 = 0.93$
 - $\hat{\rho}_{X,Y} = 0.92$
- Here, $\hat{\rho}_{X,Y} \approx \hat{\beta}_1$ because $\text{Var}(X) \approx \text{Var}(Y)$

- Turns out that correlation and bivariate regression are closely related, namely:

$$\beta_1 = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

- Compare to Pearson correlation coefficient:

$$\rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$$

- In our example
 - $\hat{\beta}_0 = 2.49$, $\hat{\beta}_1 = 0.93$
 - $\hat{\rho}_{X,Y} = 0.92$
- Here, $\hat{\rho}_{X,Y} \approx \hat{\beta}_1$ because $\text{Var}(X) \approx \text{Var}(Y)$
- In other applications numerical values may differ ... but this is just a matter of different scaling
 - i.e., both measure essentially the same thing

- A complementary way to think about inequality is based on the idea of equality of opportunities
 - the extent to which people compete on a “level playing field” vs. inherit their position

- A complementary way to think about inequality is based on the idea of equality of opportunities
 - the extent to which people compete on a “level playing field” vs. inherit their position
- An incomplete, but powerful measure

$$\mathbb{E}[p_c | P_p = p_p]$$

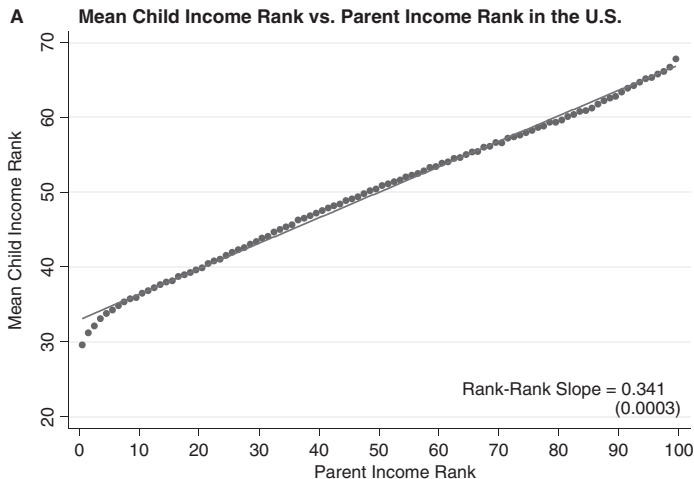
where p_c is the child's position in (lifetime) income distribution and p_p is her parent's position

Intergenerational mobility

- A complementary way to think about inequality is based on the idea of equality of opportunities
 - the extent to which people compete on a “level playing field” vs. inherit their position
- An incomplete, but powerful measure

$$\mathbb{E}[p_c | P_p = p_p]$$

where p_c is the child's position in (lifetime) income distribution and p_p is her parent's position



Children born in 1980–82. Their income is the mean of 2011–2012 family income (when the child is approximately 30 years old). Parent income is mean family income from 1996 to 2000. Children are ranked relative to other children in their birth cohort, and parents are ranked relative to all other parents. *Source: Chetty, Hendren, Kline and Saez (2014), Quarterly Journal of Economics.*

- If the conditional expectation function (CEF) of Y is linear in X , then:

$$\mathbb{E}[Y|X = x] = \beta_0 + \beta_1 x$$

- If the conditional expectation function (CEF) of Y is linear in X , then:

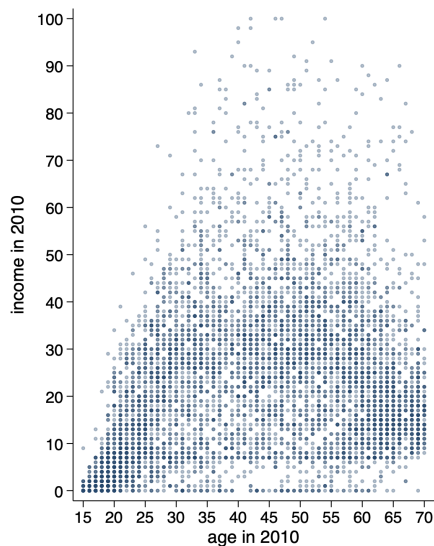
$$\mathbb{E}[Y|X = x] = \beta_0 + \beta_1 x$$

- Even if CEF is not linear, regression still provides an approximation
 - specifically, regression is the best minimum mean squared error linear approximation of CEF (more about this in later courses)
 - for many (not all) applications, this is good enough ... particularly when using multivariate regression to make it more flexible (next example)

Example: Age and income

Association between age and income

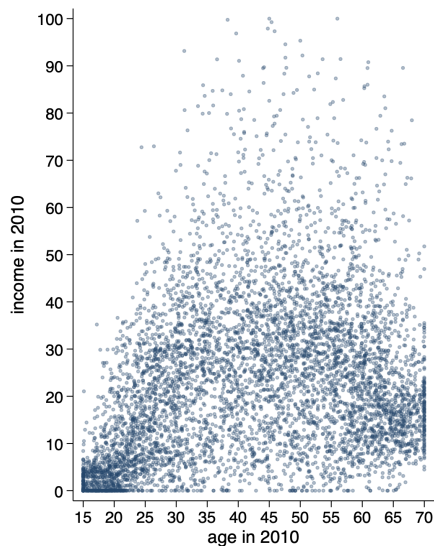
- *Question:* How does income vary with age?
 - scatter plot of the full data



Source: FLEED teaching data
scatter earn age

Association between age and income

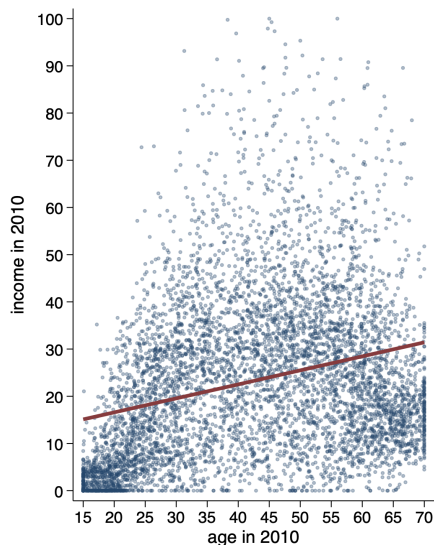
- *Question:* How does income vary with age?
 - scatter plot of the full data
 - adding a little bit of noise sometimes makes the pattern more visible



Source: FLEED teaching data
scatter earn age, jitter(10)

Association between age and income

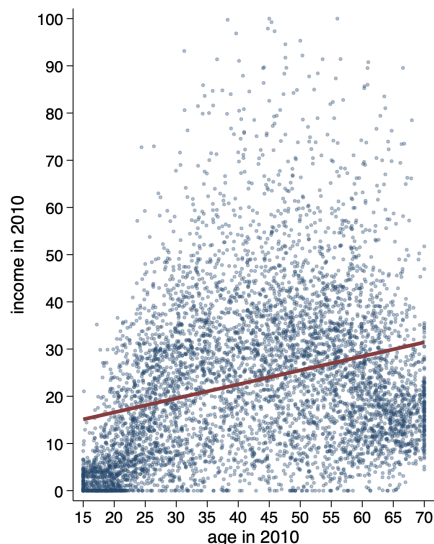
- *Question:* How does income vary with age?
 - scatter plot of the full data
 - adding a little bit of noise sometimes makes the pattern more visible
- Let's use the measures of linear dependence
 - $\hat{\rho}_{X,Y} = 0.28$
 - estimating regression $Y = \beta_0 + \beta_1 X + \epsilon$ yields parameter estimates of $\hat{\beta}_0 = 10,654$, $\hat{\beta}_1 = 297$
 - ▶ note that these estimates are in euros, while the figure's y-axis is in thousands of euros



Source: FLEED teaching data
the code is available at the course's website

Association between age and income

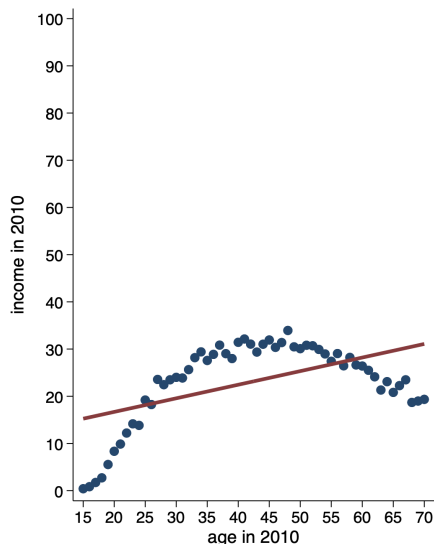
- *Question:* How does income vary with age?
 - scatter plot of the full data
 - adding a little bit of noise sometimes makes the pattern more visible
- Let's use the measures of linear dependence
 - $\hat{\rho}_{X,Y} = 0.28$
 - estimating regression $Y = \beta_0 + \beta_1 X + \epsilon$ yields parameter estimates of $\hat{\beta}_0 = 10,654$, $\hat{\beta}_1 = 297$
 - ▶ note that these estimates are in euros, while the figure's y-axis is in thousands of euros
- Are these helpful summary statistics?
 - what do they imply for $\mathbb{E}[Y|X = x]$?



Source: FLEED teaching data
the code is available at the course's website

Association between age and income

- *Question:* How does income vary with age?
 - scatter plot of the full data
 - adding a little bit of noise sometimes makes the pattern more visible
- Let's use the measures of linear dependence
 - $\hat{\rho}_{X,Y} = 0.28$
 - estimating regression $Y = \beta_0 + \beta_1 X + \epsilon$ yields parameter estimates of $\hat{\beta}_0 = 10,654$, $\hat{\beta}_1 = 297$
 - ▶ note that these estimates are in euros, while the figure's y-axis is in thousands of euros
- Are these helpful summary statistics?
 - what do they imply for $\mathbb{E}[Y|X = x]$?
- Compare to sample average by age
 - these are **nonparametric** estimates for $\mathbb{E}[Y|X = x]$
 - any ideas about how to improve the fit?



Source: FLEED teaching data
the code is available at the course's website

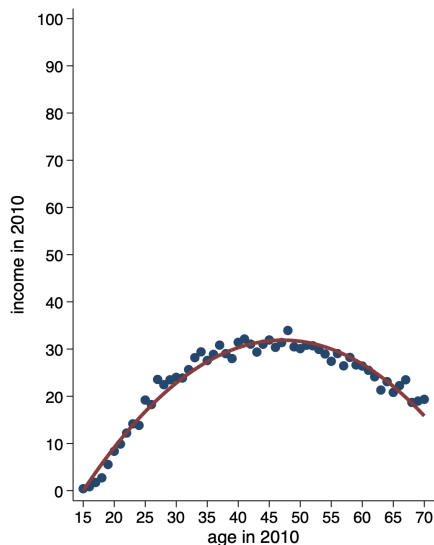
Association between age and income

- Let's use a multivariate regression model:

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \epsilon$$

- Now, the estimates that best fit the data best are:

$$\hat{\beta}_0 = -37,549, \hat{\beta}_1 = 2.857, \hat{\beta}_2 = -31$$



Source: FLEED teaching data
the code is available at the course's website

Association between age and income

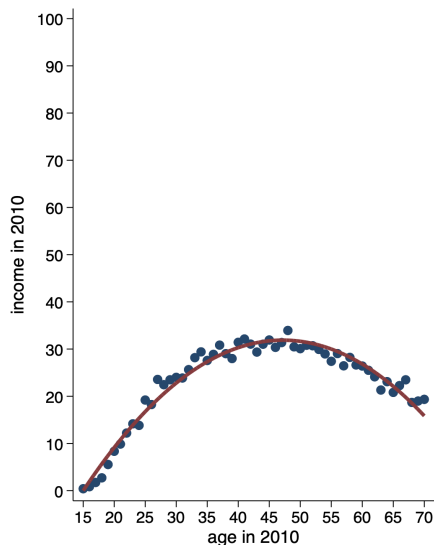
- Let's use a multivariate regression model:

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \epsilon$$

- Now, the estimates that best fit the data best are:

$$\hat{\beta}_0 = -37,549, \hat{\beta}_1 = 2.857, \hat{\beta}_2 = -31$$

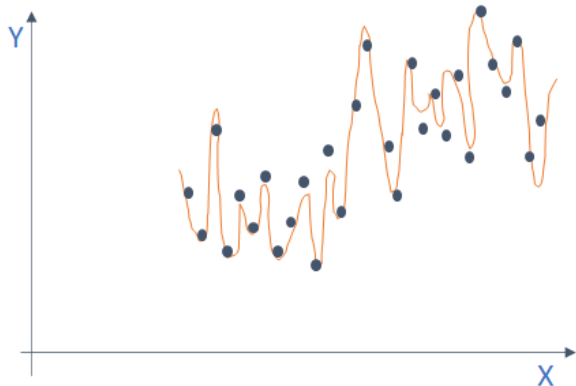
- Are these helpful summary statistics?
 - seems pretty good for approximating $\mathbb{E}[Y|X = x]$ within the 15–70 age range (the figure)
 - less so outside this age range, e.g., suggest that expected income of a new-born would be -37,549€
- General lesson: looking at the data in several ways almost always a good idea



Source: FLEED teaching data
the code is available at the course's website

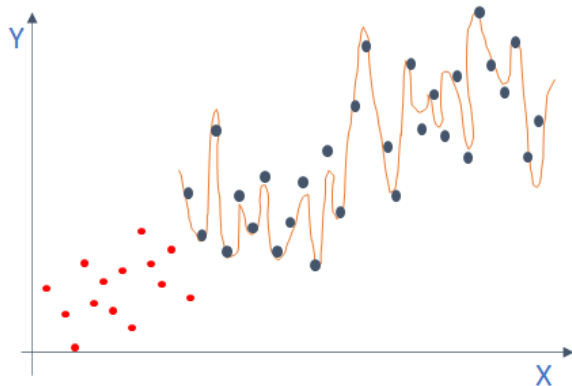
We can fit more complicated regression models

- E.g. we can fit the sample data even better using regressions with higher order polynomials.



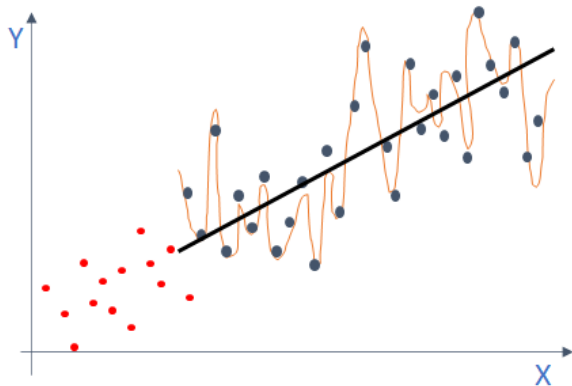
But we don't want to over-fit

- E.g. we can fit the sample data even better using regressions with higher order polynomials.
- But an over-fitted relationship might not generalize
 - e.g. with a different sample of the population.



But we don't want to over-fit

- E.g. we can fit the sample data even better using regressions with higher order polynomials.
- But an over-fitted relationship might not generalize
 - e.g. with a different sample of the population
 - likely to be worse than a linear fit



- Today we learned the basics tools for characterizing joint distributions and associations between variables.
- You should now know well the following concepts:
 - marginal and conditional distribution
 - conditional expectation function
 - scatter plots
 - covariance and correlation
 - regression, ordinary least square (OLS)
- Stata code for today's examples are on [MyCourses/More Materials/](#).
- **In-class worksheet 2** due on MyCourses before next lecture.

Up next: Causal questions

- Thus far, we have focused on **descriptive** questions
 - aim: measure the actual state of the world
 - "what is joint distribution of X and Y ?"

Up next: Causal questions

- Thus far, we have focused on **descriptive** questions
 - aim: measure the actual state of the world
 - "what is joint distribution of X and Y ?"
- We would often need to evaluate the **impact** of X on Y , e.g.
 - education on earnings
 - marketing campaign on sales
 - carbon tax on emissions
 - R&D subsidy on innovation
 - fiscal stimulus on unemployment

Up next: Causal questions

- Thus far, we have focused on **descriptive** questions
 - aim: measure the actual state of the world
 - "what is joint distribution of X and Y?"
- We would often need to evaluate the **impact** of X on Y, e.g.
 - education on earnings
 - marketing campaign on sales
 - carbon tax on emissions
 - R&D subsidy on innovation
 - fiscal stimulus on unemployment
- These are **causal** questions
 - aim: compare *counterfactual* states of the world
 - "how would Y change if we changed X?"
 - ▶ we typically refer to Y as "outcome" and to X as "treatment"

- **Pre-class assignment 3**
 - to prepare for next lecture (causality and research design)
 - case of high-occupancy vehicle restriction in Jakarta
- **Homework 1** deadline: Jan 17 at 13:00
 - Don't wait till the last minute!
 - An important skill when working with data is to learn to troubleshoot efficiently.
 - This learning often involves being "stuck".
- Use the course **Slack** channel to seek help and help others in the class
 - Quicker than waiting for private responses from the TA or me
 - Recall extra incentive: bonus points for active participation