

Week 5: Linear Regression

Video 1: Introduction to Linear Regression

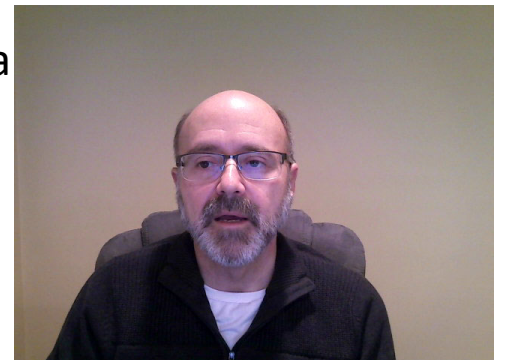
EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

Department of Epidemiology, Harvard TH Chan School of Public Health



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Key Concepts

- Advantages of regression models
- Linear regression
- Log odds scale
- Logistic regression



Advantages of Regression Models

- Generally more efficient than stratification-based methods when data are sparse
- Modelling of a continuous outcomes
- Specify continuous and categorical exposures, confounders and modifiers
- Specify interactions to model effect modification
- Model nonlinear relationships between exposure and outcome and other covariates



Regression Models

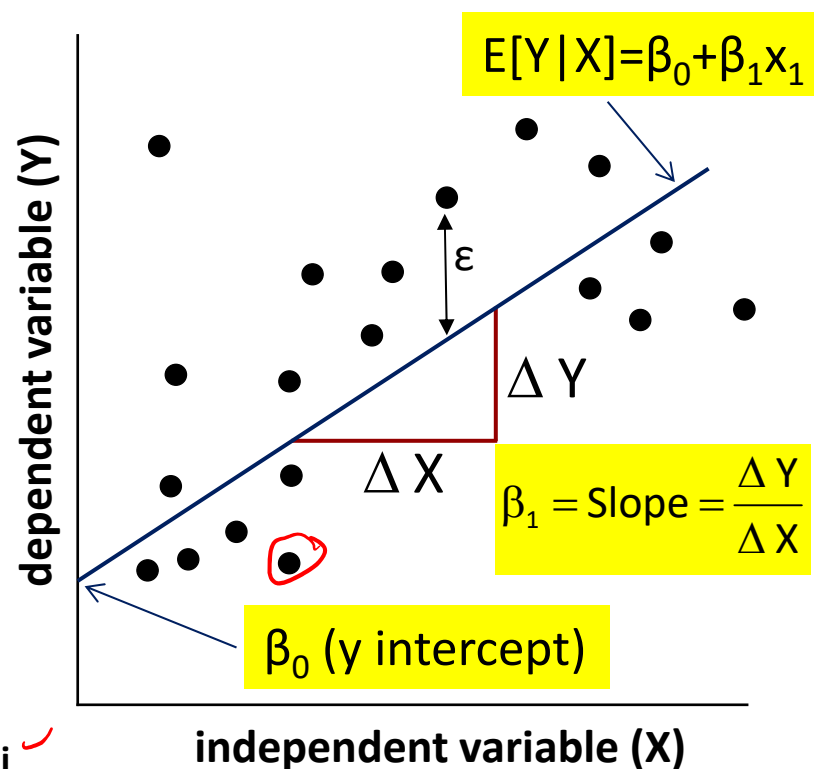
y	=	$\beta_0 + \beta_1 x_1 + \dots \beta_i x_i + \varepsilon$
Dependent		Independent
Predicted		Predictor variables
Response variable		Explanatory variables
Outcome variable		Covariables



Simple Linear Regression

- Predict a continuous dependent (outcome) variable y from a continuous independent (exposure) variable x
- Simple linear regression fits a straight line to the data using the least squares method.

- Regression line: $E[Y|X] = \beta_0 + \beta_1 x_1$
 - Often presented as $y = mx + b$ where
 - $b = y\text{-intercept}$
 - $m = \text{slope} = \Delta y / \Delta x$ (rise/run)



- Individual predicted value: $Y_i = \beta_0 + \beta_1 x_{1i} + \epsilon_i$
 - $\beta_0 = y\text{-intercept}$ (where the line crosses the Y-axis)
 - $\beta_1 = \text{slope} = \Delta y / \Delta x$ = average change in y when x changes by one unit
 - x_1 is a known constant
 - ϵ , the error, is an observation's deviation from the conditional mean, $N(0, \sigma^2)$



BREAK

Week 5: Linear Regression

Video 2: Linear Regression Example

EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

Department of Epidemiology, Harvard TH Chan School of Public Health

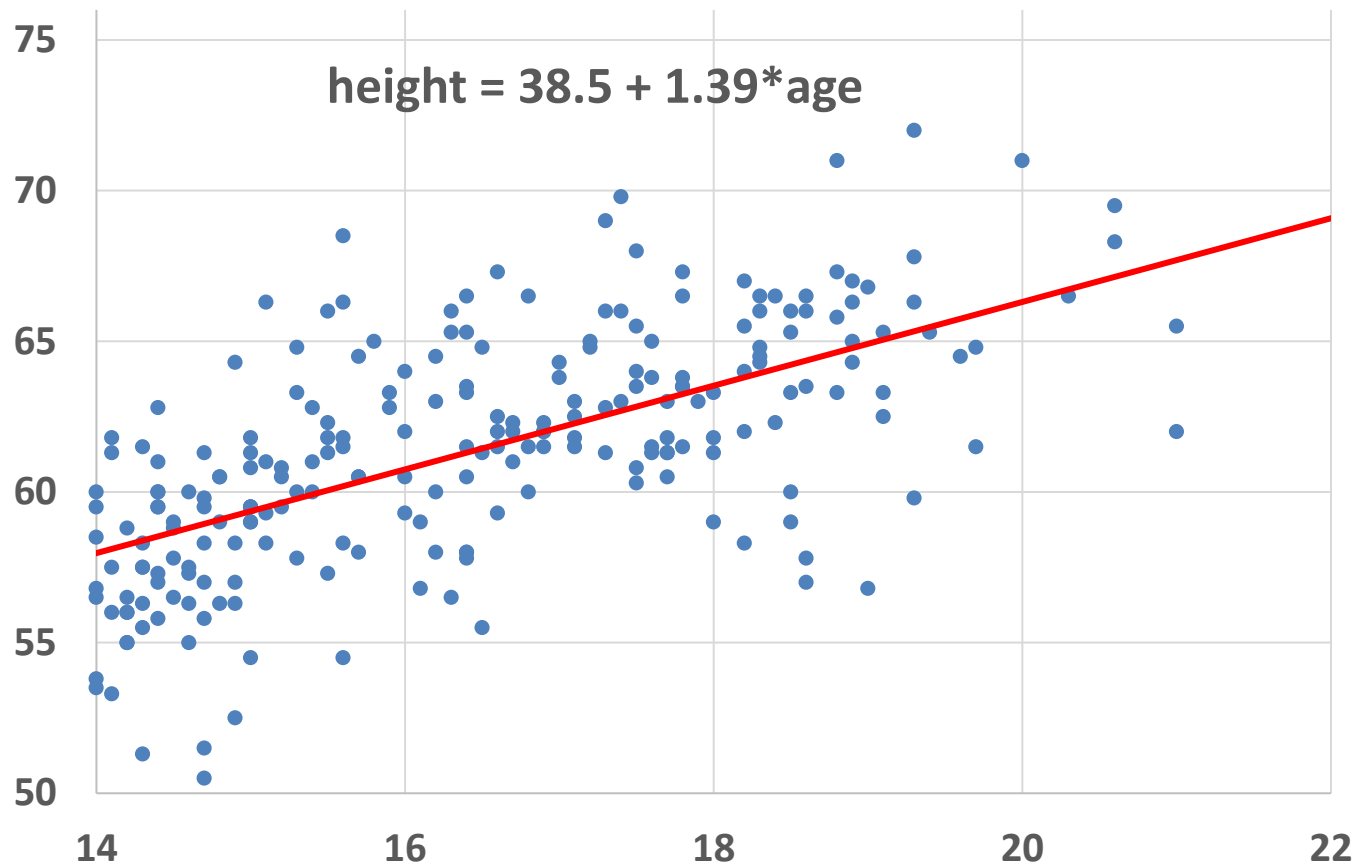


HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



Linear Regression

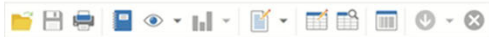
Age versus Height



Linear Regression Example

- Download the HeightWeight dataset from the Regression I module on Canvas
 - The dataset is available in multiple formats
 - CSV
 - Excel
 - R
 - SAS
 - Stata





History

Filter commands here

Command

There are no items to show.

Variables

Filter variables here

Name	Label
------	-------

There are no items to show.

Properties

Variables	
Name	
Label	
Type	
Format	
Value label	
Notes	
Data	
Frame	default
Filename	
Label	
Notes	
Variables	0
Observations	0
Size	0
Memory	64M
Sorted by	



History



Filter commands here

Command

1 use "C:\Users\Mu...

. clear

. use "C:\Users\Murray\Dropbox\EPI202 Fall\EPI202 Fall 2019\Regression data set\HeightWeight.dta"

.

Command

I

Variables



Filter variables here

Name Label

gender

age

height

weight

female

femage

Properties



Lock < >

Variables

Name

Label

Type

Format

Value label

Notes

Data

Frame

default

Filename

HeightWeight.dta

Label

Notes

Variables

6

Observations

237

Size

6.71K

Memory

64M

Sorted by

height	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
age	1.457297	.1107962	13.15	0.000	1.239011	1.675583
_cons	37.42901	1.828191	20.47	0.000	33.82719	41.03083

BREAK

Week 5: Linear Regression

Video 3: Linear Regression Adjusted for Covariates

EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

Department of Epidemiology, Harvard TH Chan School of Public Health



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH

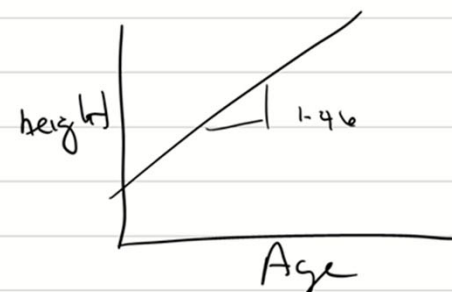


height	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
age	1.457297	.1107962	13.15	0.000	1.239011	1.675583
_cons	37.42901	1.828191	20.47	0.000	33.82719	41.03083

$$Y = \beta_0 + \beta_1 X_1 + \varepsilon$$

$$E(\text{height} | \text{Age}) = \beta_0 + \beta_1 (\text{age})$$

$$E(\text{height} | \text{Age}) = 37.43 + 1.46 (\text{age})$$



$$E(\text{height} | 15) = 37.43 + 1.46(15)$$

Unnamed.pdf* - PDF Annotator

File Edit Tool View Extras Window Help

Page Width

Unnamed.pdf*

Pen

female Age height

Modified List of loaded documents

Type here to search

1 of 1

1:54 PM 11/29/2019

Unnamed.pdf* - PDF Annotator

File Edit Tool View Extras Window Help

Page Width

Pause 00:00:00 Select Area Audio Record Pointer

Unnamed.pdf*

Pen

female Age height

height	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
female	-1.533874	.5042084	-3.04	0.003	-2.527242	-.5405058
_cons	62.06	.3457925	179.47	0.000	61.37874	62.74126

Modified

Type here to search

1 of 1

2:05 PM 11/29/2019



History

Filter commands here

Command

1 regress height age

2 regress height fe...

. regress height age

Source	SS	df	MS	Number of obs	=	236
Model	1545.44634	1	1545.44634	F(1, 234)	=	173.00
Residual	2090.37277	234	8.93321697	Prob > F	=	0.0000
				R-squared	=	0.4251
				Adj R-squared	=	0.4226
Total	3635.81911	235	15.4715707	Root MSE	=	2.9888

height	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
age	1.457297	.1107962	13.15	0.000	1.239011 1.675583
_cons	37.42901	1.828191	20.47	0.000	33.82719 41.03083

. regress height female

Source	SS	df	MS	Number of obs	=	236
Model	138.324876	1	138.324876	F(1, 234)	=	9.25
Residual	3497.49423	234	14.9465566	Prob > F	=	0.0026
				R-squared	=	0.0380
				Adj R-squared	=	0.0339
Total	3635.81911	235	15.4715707	Root MSE	=	3.8661

height	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
female	-1.533874	.5042084	-3.04	0.003	-2.527242 -.5405058
_cons	62.06	.3457925	179.47	0.000	61.37874 62.74126

Command

Variables

Filter variables here

Name Label

gender

age

height

weight

female

femage

Properties

Variables

Name

Label

Type

Format

Value label

Notes

Data

Frame default

Filename HeightWeight.dta

Label

Notes

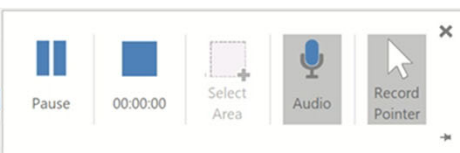
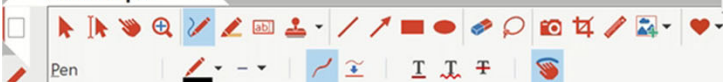
Variables 6

Observations 236

Size 6.91K

Memory 64M

Sorted by



height	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
age	1.465677	.1068379	13.72	0.000	1.255185	1.676169
female	-1.627297	.3758752	-4.33	0.000	-2.367845	-.8867485
_cons	38.0569	1.768544	21.52	0.000	34.57252	41.54128

BREAK

Week 5: Linear Regression

Video 4: Linear Regression with Interaction Terms to Account for Effect Measure Modification

EPI202 – Epidemiologic Methods II

Murray A. Mittleman, MD, DrPH

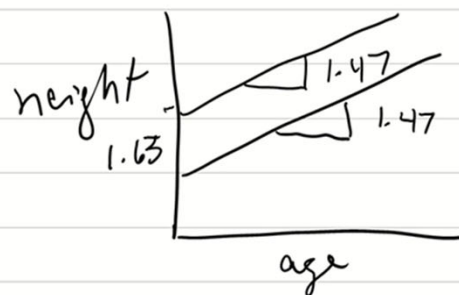
Department of Epidemiology, Harvard TH Chan School of Public Health



HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH



$$\text{height} = 38.06 + 1.47(\text{age}) - 1.63(\text{female})$$





History

Filter commands here

Command

1 regress height age

2 regress height fe...

3 regress height ag...

. regress height female

Source	SS	df	MS	Number of obs	=	236
Model	138.324876	1	138.324876	F(1, 234)	=	9.25
Residual	3497.49423	234	14.9465566	Prob > F	=	0.0026
				R-squared	=	0.0380
				Adj R-squared	=	0.0339
				Root MSE	=	3.8661
Total	3635.81911	235	15.4715707			

height	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
female	-1.533874	.5042084	-3.04	0.003	-2.527242 - .5405058
_cons	62.06	.3457925	179.47	0.000	61.37874 62.74126

. regress height age female

Source	SS	df	MS	Number of obs	=	236
Model	1701.08304	2	850.541519	F(2, 233)	=	102.43
Residual	1934.73607	233	8.30358829	Prob > F	=	0.0000
				R-squared	=	0.4679
				Adj R-squared	=	0.4633
				Root MSE	=	2.8816
Total	3635.81911	235	15.4715707			

height	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
age	1.465677	.1068379	13.72	0.000	1.255185 1.676169
female	-1.627297	.3758752	-4.33	0.000	-2.367845 -.8867485
_cons	38.0569	1.768544	21.52	0.000	34.57252 41.54128

.

Command

Variables

Filter variables here

Name Label

gender

age

height

weight

female

femage

Properties

< >

Variables

Name

Label

Type

Format

Value label

Notes

Data

Frame default

Filename HeightWeight.dta

Label

Notes

Variables 6

Observations 236

Size 6.91K

Memory 64M

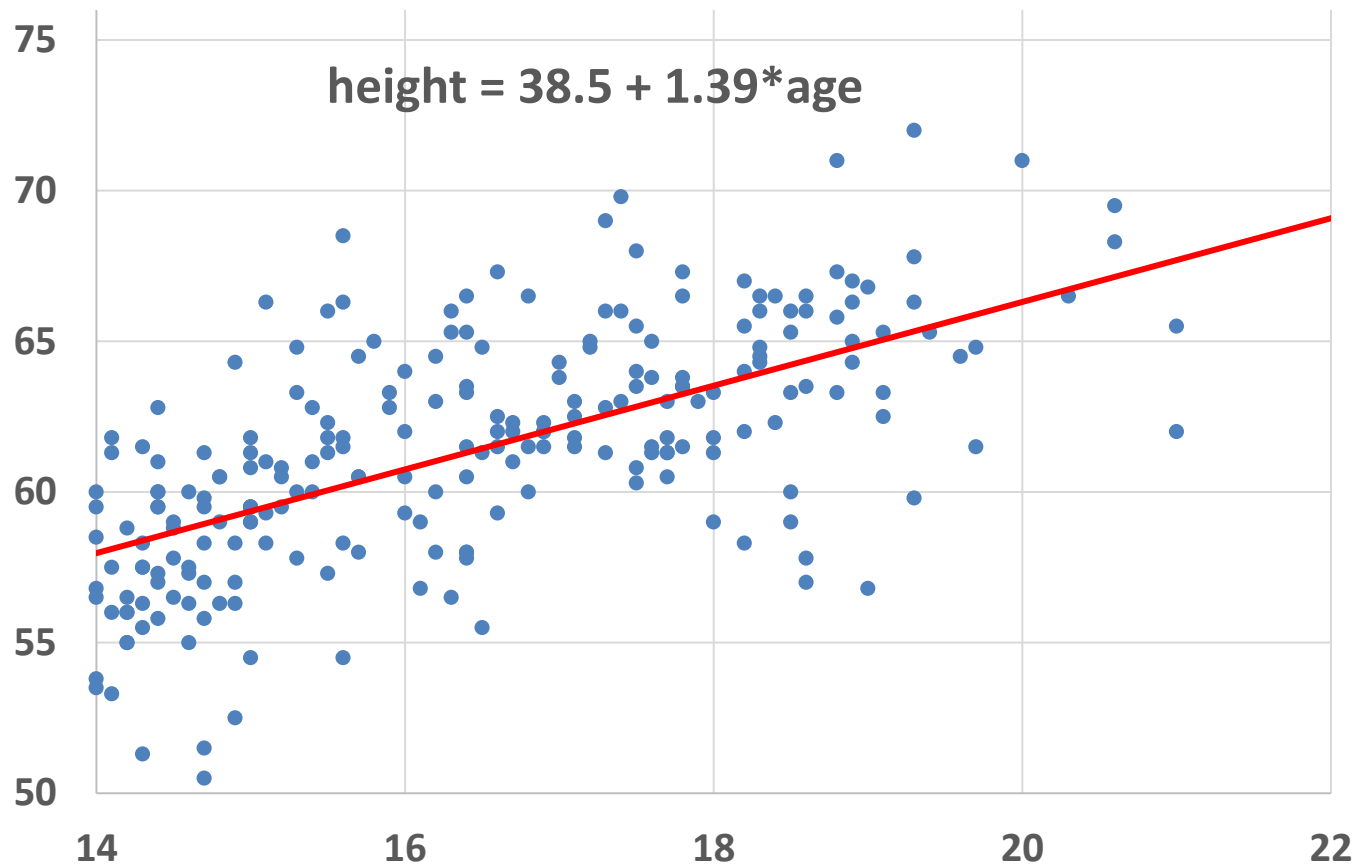
Sorted by

height	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
age	1.918139	.1449452	13.23	0.000	1.632562	2.203716
female	13.31463	3.393413	3.92	0.000	6.628783	20.00047
femage	-.9106005	.2056254	-4.43	0.000	-1.315732	-.5054687
_cons	30.64702	2.386661	12.84	0.000	25.94473	35.34932

$$\text{height} = \beta_0 + \beta_1(\text{age}) + \beta_2(\text{female}) + \beta_3(\text{age} * \text{female})$$

Linear Regression

Age versus Height

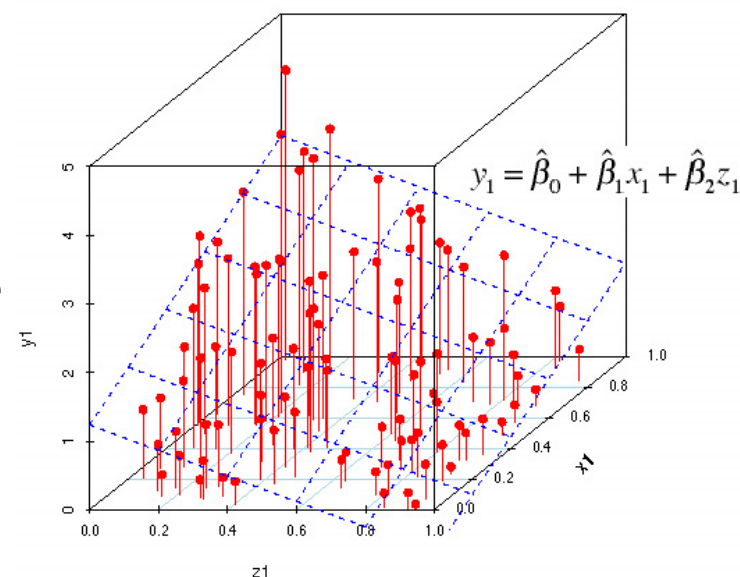


Multiple Linear Regression

- Multiple linear regression models predict a continuous dependent (outcome) variable y from continuous and categorical independent variables x_i :

$$E[Y|X] = \beta_0 + \beta_1 x_1 + \dots + \beta_i x_i$$

- The regression line is the best-fit line through the points in the data
 - $\beta_0, \beta_1, \dots, \beta_k$ are parameters
 - X_1, X_2, \dots, X_k are known constants
 - β_k = change in average outcome (difference in mean outcome) per unit change in X_i holding all other X 's constant



BREAK