

## Simple linear regression

We will use a car data set and our dependent variable is mpg (miles per gallon) vs our independent variable car weight.

Data Editor (Edit) - [regression\_auto.dta]

File Edit View Data Tools

make[1] AMC

	make	mpg	weight	weight1	price	foreign	repairs	length
1	AMC	22	2930	2.93	4099	0	3	186
2	AMC	17	3350	3.35	4749	0	3	173
3	AMC	22	2640	2.64	3799	0	3	168
4	Audi	17	2830	2.83	9690	1	5	189
5	Audi	23	2070	2.07	6295	1	3	174
6	BMW	25	2650	2.65	9735	1	4	177
7	Buick	20	3250	3.25	4816	0	3	196
8	Buick	15	4080	4.08	7827	0	4	222
9	Buick	18	3670	3.67	5788	0	3	218
10	Buick	26	2230	2.23	4453	0	3	170
11	Buick	20	3280	3.28	5189	0	3	200

Variables

Filter variables here

<input checked="" type="checkbox"/> Name	Label
<input checked="" type="checkbox"/> make	make of car
<input checked="" type="checkbox"/> mpg	miles per gallon
<input checked="" type="checkbox"/> weight	car weight in pounds
<input checked="" type="checkbox"/> weight1	car weight in 1,000 pounds
<input checked="" type="checkbox"/> price	car price in 1978 dollars
<input checked="" type="checkbox"/> foreign	=1 if car is foreign
<input checked="" type="checkbox"/> repairs	number of car repairs
<input checked="" type="checkbox"/> length	car length

Description of our response variable and predictor variable but lets first create our global variables

```
global ylist mpg
global xlist weight1
```

```
. describe $ylist $xlist
```

variable name	storage type	display format	value label	variable label
mpg	byte	%8.0g		miles per gallon
weight1	float	%9.0g		car weight in 1,000 pounds

Data Summary

```
. summarize $ylist $xlist
```

Variable	Obs	Mean	Std. Dev.	Min	Max
mpg	26	20.92308	4.757504	14	35
weight1	26	3.099231	.6950794	2.02	4.33

```
. summarize $ylist, detail
```

miles per gallon					
Percentiles		Smallest			
1%	14	14			
5%	14	14			
10%	15	15	Obs	26	
25%	17	16	Sum of Wgt.	26	
50%	21		Mean	20.92308	
		Largest	Std. Dev.	4.757504	
75%	23	25			
90%	26	26	Variance	22.63385	
95%	29	29	Skewness	.8806144	
99%	35	35	Kurtosis	4.243808	

## Question

Can the weight of a car statistically significantly predict a car's miles per gallon?

## Hypothesis

H0: Car weight can not statistically significantly predict a car's miles per gallon

Ha: Car weight can statistically significantly predict a car's miles per gallon

## The level of significance

$\alpha = 0.05$

## ASSUMPTIONS

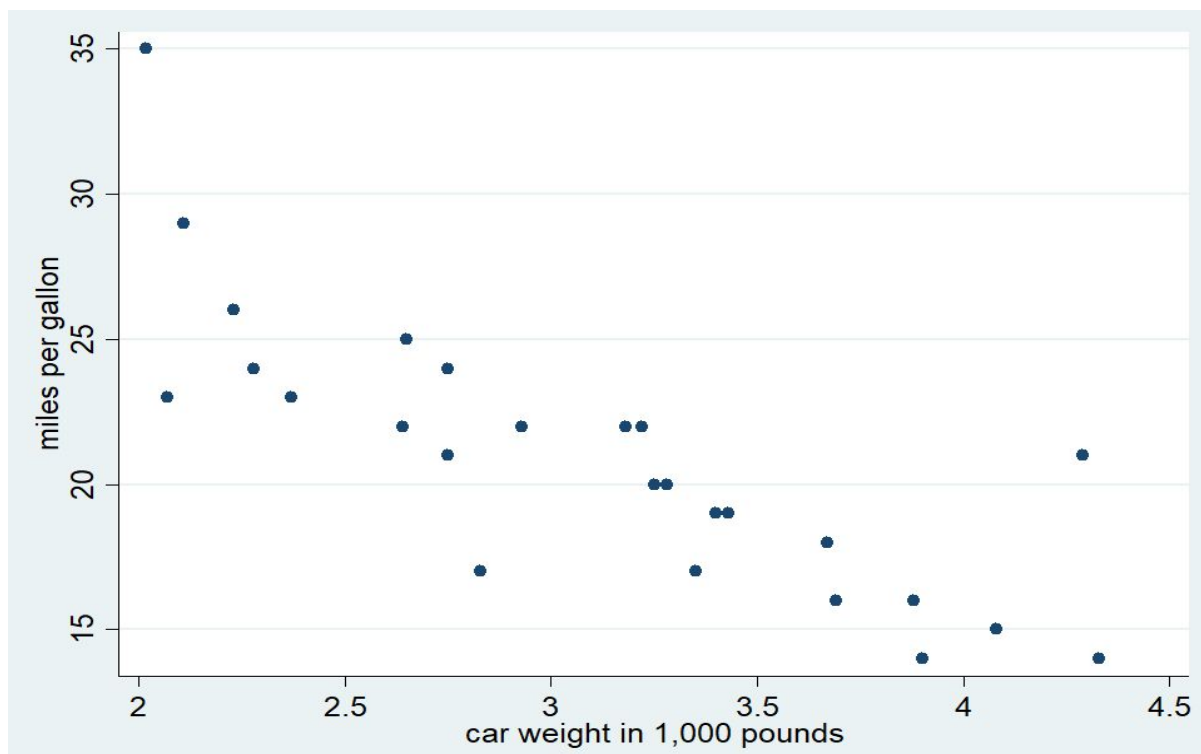
Determine if data meets requirements to perform a linear regression.

**Assumption #1:** Your response variable should be measured on a continuous scale.

**Assumption #2:** Your independent variable should be measured at the continuous or categorical level.

**Assumption #3:** There needs to be a linear relationship between the dependent and independent variables. Let's plot a scatter plot and check

```
. graph twoway (scatter $ylist $xlist)
```



```
. correlate $ylist $xlist  
(obs=26)
```

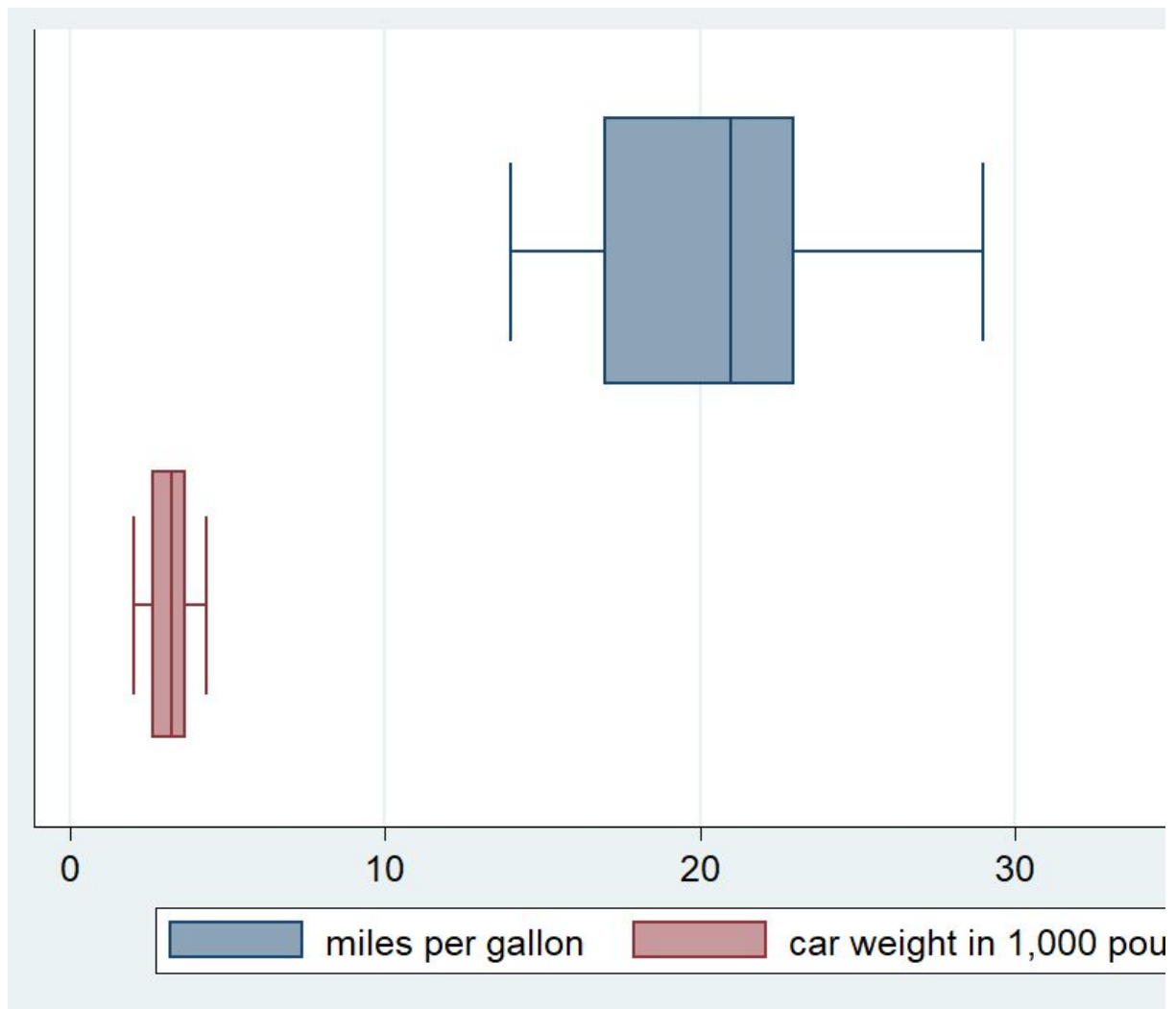
	mpg	weight1
mpg	1.0000	
weight1	-0.8082	1.0000

As seen above the two variables appear to be linearly related

**Assumption #4:** You should have independence of observations. which you can easily check using the **Durbin-Watson statistic**, which is a simple test to run using Stata.

```
. dwstat  
  
Durbin-Watson d-statistic( 2, 26) = 1.991873
```

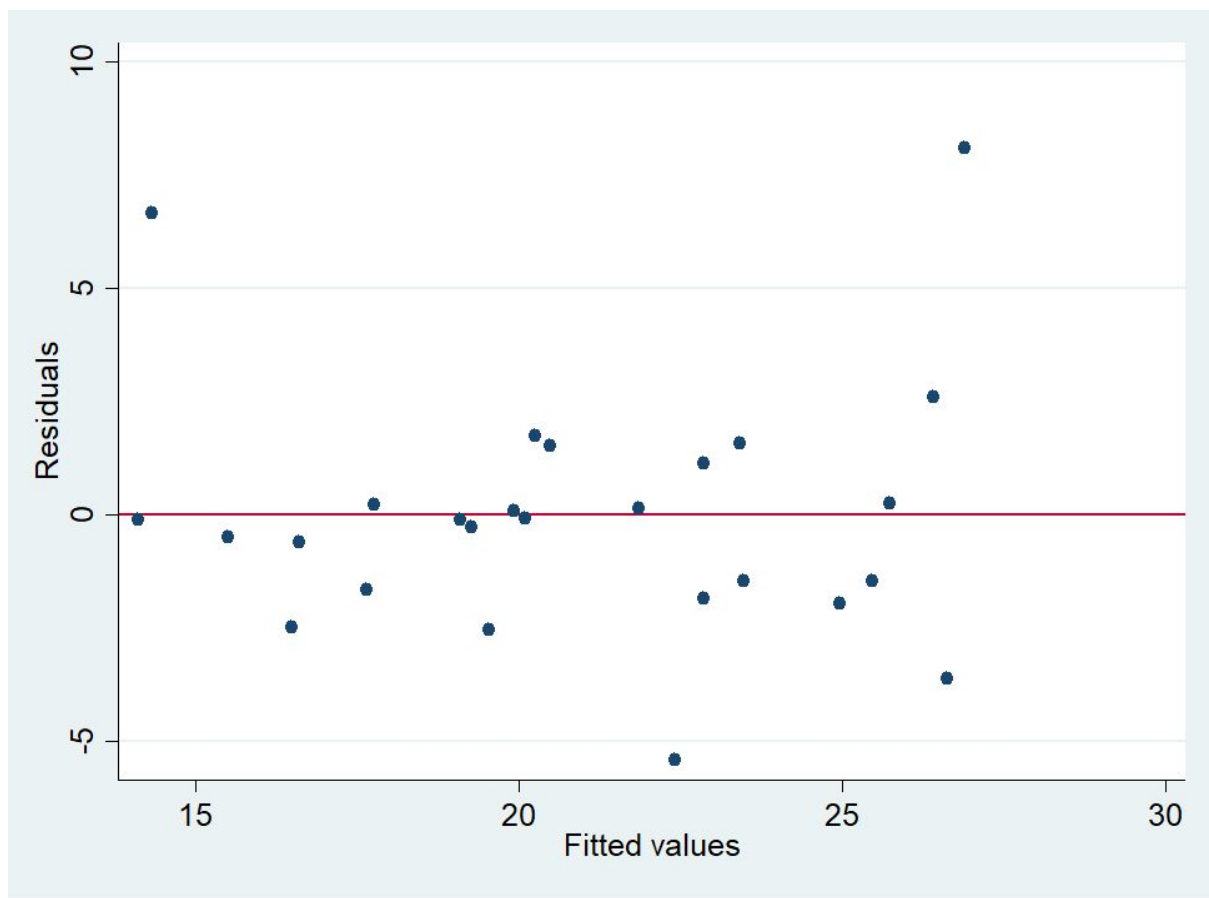
**Assumption #5:** There should be no significant outliers. We use box plot



**Assumption #6:** Your data needs to show homoscedasticity, which is where the variances along the line of best fit remain similar as you move along the line.

One of the major assumptions given for type ordinary least squares regression is the homogeneity in the case of variance of the residuals. In the case of a well-fitted model, if you plot residual values versus fitted values, you should not see any particular pattern. Now, what if the variance given by the residuals is not a constant? In this case, the **residual variance** is called **heteroscedastic**. The most commonly used way to detect heteroscedasticity is by plotting residuals versus predicted values. We should not have any side narrower than the other.

```
. rvfplot, yline(0)
```



We can also use non-graphical commands as shown below. The first test on heteroscedasticity given by `imtest` is the White's test and the second one given by `hettest` is the Breusch-pagan test.

Both test the null hypothesis that the variance of the residuals are homogenous.

```
. estat imtest
```

Cameron & Trivedi's decomposition of IM-test

Source	chi2	df	p
Heteroskedasticity	6.39	2	0.0410
Skewness	2.57	1	0.1090
Kurtosis	3.36	1	0.0667
Total	12.32	4	0.0151

```
. estat hettest
```

Breusch-Pagan / Cook-Weisberg test for heteroskedasticity

Ho: Constant variance

Variables: fitted values of mpg

```
chi2(1)      =      1.48
Prob > chi2   =      0.2239
```

**Assumption #7:** Finally, you need to check that the residuals (errors) of the regression line are approximately normally distributed.

First we run the regression as shown below

```
. reg $ylist $xlist
```

Source	SS	df	MS	Number of obs	=	26
Model	369.567767	1	369.567767	F(1, 24)	=	45.19
Residual	196.278387	24	8.17826611	Prob > F	=	0.0000
				R-squared	=	0.6531
				Adj R-squared	=	0.6387
Total	565.846154	25	22.6338462	Root MSE	=	2.8598

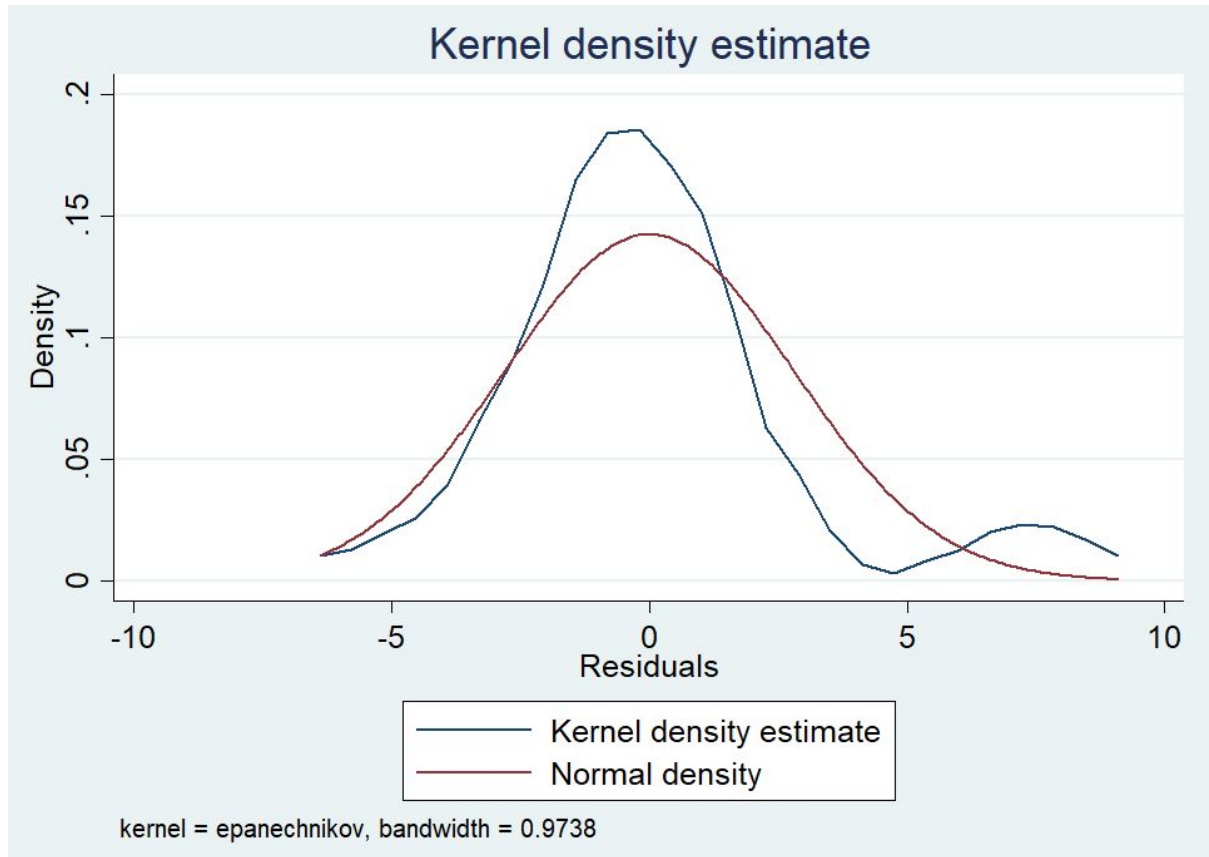
mpg	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
weight1	-5.531496	.8228604	-6.72	0.000	-7.229797 -3.833196
_cons	38.06646	2.611177	14.58	0.000	32.67726 43.45566

We the use the predict command to generate residuals

```
. predict elhat, resid
```

We finally use `kdensity` command to plot a kernel density plot with a normal option requesting that a normal density be overlaid on the plot.

```
. kdensity elhat, normal  
(n() set to 26)
```



We can also use Shapiro-Wilk W test for normal data as shown below. The p value is based on the assumption that the distribution is normal.

```
. swilk elhat
```

Shapiro-Wilk W test for normal data

Variable	Obs	W	V	z	Prob>z
elhat	26	0.89653	2.959	2.223	0.01311

As seen above this assumption fails as the test is significant.

## Simple linear interpretation

```
. reg $ylist $xlist
```

Source	SS	df	MS	Number of obs	=	26
Model	369.567767	1	369.567767	F(1, 24)	=	45.19
Residual	196.278387	24	8.17826611	Prob > F	=	0.0000
				R-squared	=	0.6531
				Adj R-squared	=	0.6387
Total	565.846154	25	22.6338462	Root MSE	=	2.8598

mpg	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
weight1	-5.531496	.8228604	-6.72	0.000	-7.229797 -3.833196
_cons	38.06646	2.611177	14.58	0.000	32.67726 43.45566

A linear regression established that weight of a car could statistically significantly predict mpg (miles covered by the car per gallon),  $F(1, 24) = 45.19$ ,  $p < .05$  and the weight of a car accounted for 65.3% of the explained variability in miles covered by the car per gallon. The regression equation was: predicted miles covered by the car per gallon =  $38.066 - 5.531 \times$  (weight of the car).

## Plotting a regression line

