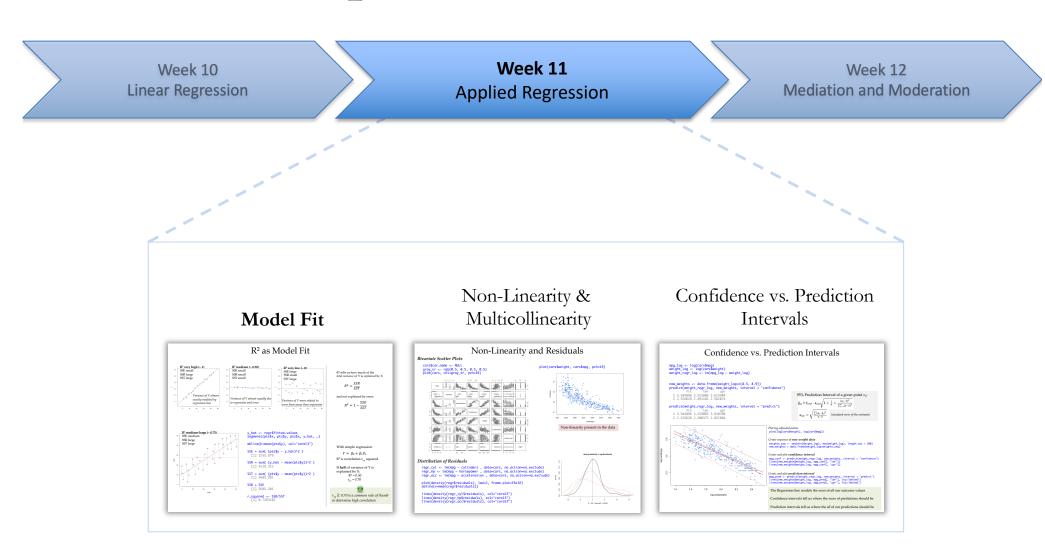
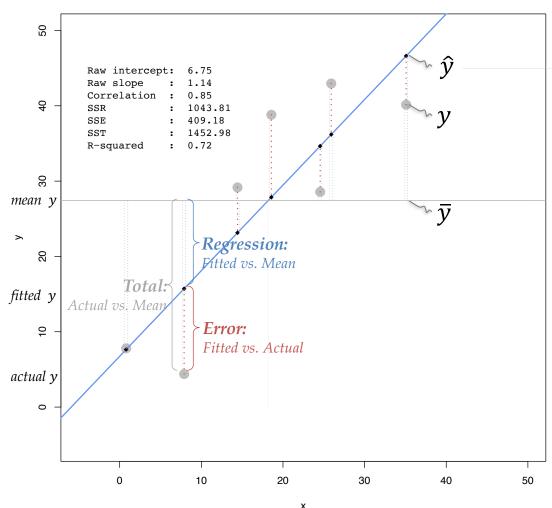
Business Analytics Using Computational Statistics





R² tells us how much of the total variance of Y is explained by the regression model

$$R^2 = \frac{SSR}{SST}$$

or, how much of the *total variance* is *not simply error variance*:

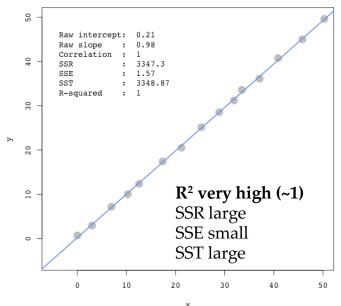
$$R^2 = 1 - \frac{SSE}{SST}$$

```
regr \leftarrow lm(y \sim x, data=pts)
y_hat <- regr$fitted.values</pre>
                                    "fitted" points ŷ
points(pts$x, y_hat, pch=18)
abline(h=mean(pts$y), ...)
segments(pts$x, pts$y, pts$x, y hat, ...)
segments(pts$x-0.2, pts$y, pts$x-0.2, mean(pts$y), ...)
segments(pts$x+0.2, y hat, pts$x+0.2, mean(pts$y), ...)
SSE = sum( (pts\$y - y hat)^2)
[1] 409.1765
SSR = sum( (y_hat - mean(pts$y))^2 )
[1] 1043.807
SST = sum( (pts\$y - mean(pts\$y))^2 )
[1] 1452.983
SSR + SSE
[1] 1452.983
r squared <- SSR/SST
[1] .7183887
```

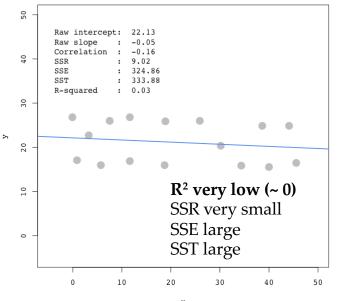
$$R^2 = 0.72$$
: Variance of Y is **mostly** explained by X

Interpreting R²

R² as Model Fit



Variation of Y almost exactly explained by regression Model



Variance of Y more related to error from mean than regression

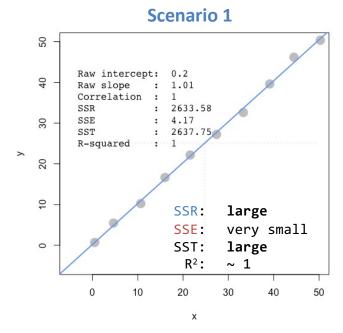
With simple regression: $Y = \beta_0 + \beta_1 X_1$ **R**² is correlation \mathbf{r}_{xy} squared

If half of variance of Y is explained by X:

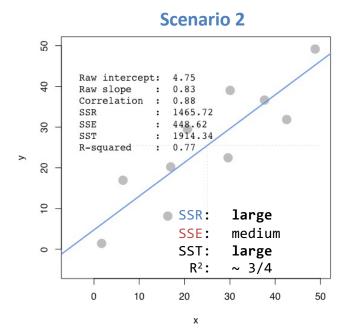
R² = 0.50 implies r_{xy} = 0.70



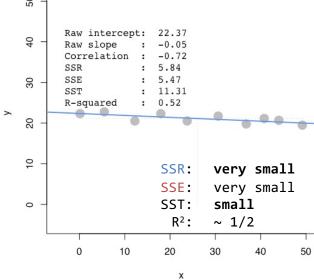
 $r_{xy} \ge 0.70$ is a common rule of thumb to determine high correlation; it suggests that one variable at least half explains the variance of the other



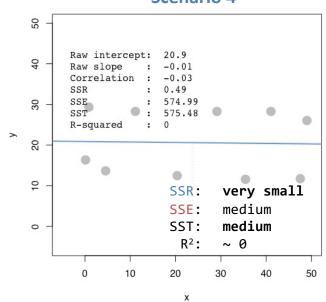
$$R^2 = \frac{SSR}{SST}$$



Scenario 3



Scenario 4



Regression Using Linear Algebraic

Rsq <- SSR/SST # 0.847

Constant values of 1 to capture the intercept

Intercept term
$$\widehat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T y = \begin{bmatrix} \widehat{\beta_0} \\ \widehat{\beta_1} \\ \widehat{\beta_2} \\ \vdots \end{bmatrix}$$

$$\hat{\mathbf{y}} = X\hat{\beta}$$

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

Variance Explained

This is how much of our data's total variance is explained by our regression model

```
prog <- read.csv("programmer salaries.txt", sep="\t")</pre>
x_vars <- c("Experience", "Score", "Degree")</pre>
X <- cbind("(intercept)"=1, as.matrix(prog[, x_vars]))</pre>
        (intercept) Experience Score Degree
  y <- prog[, "Salary"]</pre>
   [1] 24.0 43.0 ...
  30.0
beta hat <- solve(t(X)%*%X) %*% t(X)%*%y
   [,1]
                                 Regression Coefficients
  (intercept) 7.944849
  Experience 1.147582 <
                                 These are the estimated slopes
   Score
             0.196937
                                    of our regression model
             2.280424
   Degree
y_hat <- X %*% beta hat</pre>
            [,1]
                                      Fitted Values
   [1,] 27.89626
                                These are estimated values of y
   [2, ] 37.95204
                                  from our regression model
  [20,] 28.91499
SSR <- sum((y_hat - mean(y))^2) # 507.896
SSE <- sum((y - y_hat)^2) # 91.890
SST \leftarrow sum((y - mean(y))^2) # 599.786
                                                         5
```

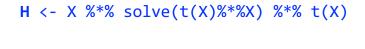
The "Hat" Matrix

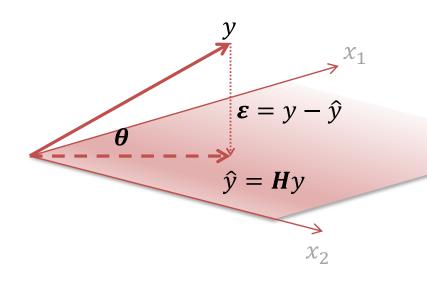
Another way of looking at regression is that it is trying to orthogonally project y to the space spanned by X in "row space"

y_hat <- H %*% y

$$\boldsymbol{H} = X(X^T X)^{-1} X^T$$

$$\hat{y} = Hy$$





H puts the "hat" on y

H orthogonally projects y onto the space of X \hat{y} is the projection of y onto the space of X

$$R = cor(y, \hat{y})$$

$$R^2 = cor(y, \hat{y})^2$$

Regression Coefficients

Our estimated regression coefficients are based on our one sample. They have variance (uncertainty) associated with them...

Simple Regression

Multiple Regression

$$y = \beta_o + \beta_1 x_1$$
$$\beta_1 = r_{yx_1}$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$\beta_1 = \frac{r_{yx_1} - r_{yx_2} r_{x_1 x_2}}{1 - r_{x_1 x_2}}$$

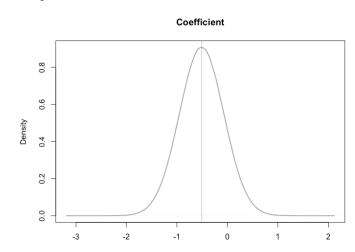
Each beta coefficient b_i could vary with standard error s_{b_i}

$$s_e = \sqrt{MSE} = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n - 2}}$$

Standard Error of coefficient b_i

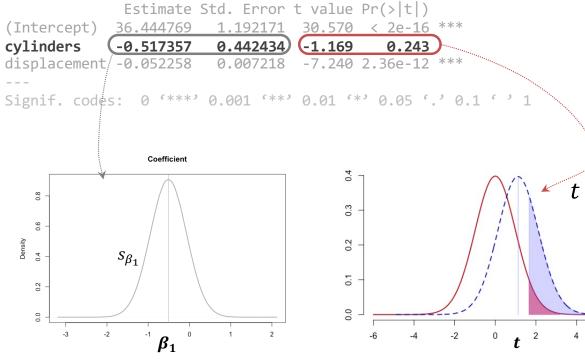
$$s_{\beta_1} = \frac{s_e}{\sum (x_i - \bar{x})^2}$$

(for simple regression)



Coefficients and t-Tests

engine_regr <- lm(mpg ~ cylinders + displacement, data=cars, na.action=na.exclude)
summary(engine regr)</pre>

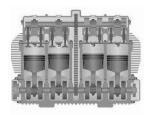


Distribution of cylinders coefficient:

- mean value -0.51
- standard deviation of mean: 0.44 (standard error)

t-test of H_{null}: coefficient is 0

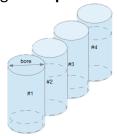
- t-value: -1.169
- p-value: 0.243 (two-tailed)



Engine Cylinders



Engine **Displacement**



Multiple Regression

regr <- lm(mpg ~ cylinders + displacement + horsepower + weight + acceleration + model_year + factor(origin),

data=cars, na.action=na.exclude)

```
Estimate Std. Error t value Pr(>|t|)
(Intercept)
               -1.795e+01 4.677e+00 -3.839 0.000145 ***
cylinders
               -4.897e-01 3.212e-01 -1.524 0.128215
displacement
                2.398e-02 7.653e-03
                                     3.133 0.001863 **
horsepower
               -1.818e-02 1.371e-02 -1.326 0.185488
weight
               -6.710e-03 6.551e-04 -10.243 < 2e-16 ***
acceleration
               7.910e-02 9.822e-02 0.805 0.421101
model_year
               7.770e-01 5.178e-02 15.005 < 2e-16 ***
factor(origin)2 2.630e+00 5.664e-01
                                     4.643 4.72e-06 ***
factor(origin)3 2.853e+00 5.527e-01
                                      5.162 3.93e-07 ***
Multiple R-squared: 0.8242, Adjusted R-squared: 0.8205
```

Unstandardized Raw Coefficients have units (mpg/dependent units)

Easier to relate coefficient to domain: "-6.71 mpg per pound (lb.)"

But hard to compare coefficients: "0.77 mpg/year" versus "-6.71 mpg/lb"

summary(cars\$weight)

Min. 1st Qu. Median Mean 3rd Qu. Max. 1613 2224 2804 2970 3608 5140

summary(cars\$cylinders)

Min. 1st Qu. Median Mean 3rd Qu. Max. 3.000 4.000 4.000 5.455 8.000 8.000

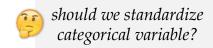
Standardized Multiple Regression

regr_std <- lm(mpg ~ cylinders + displacement + horsepower + weight + acceleration + model_year + **factor(cars\$origin)**, data=as.data.frame(**scale(cars)**), na.action=na.exclude)

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-0.13323	0.03174	-4.198	3.35e-05	***
cylinders	-0.10658	0.06991	-1.524	0.12821	
displacement	0.31989	0.10210	3.133	0.00186	**
horsepower	-0.08955	0.06751	-1.326	0.18549	
weight	-0.72705	0.07098	-10.243	< 2e-16	***
acceleration	0.02791	0.03465	0.805	0.42110	
model_year	0.36760	0.02450	15.005	< 2e-16	***
<pre>factor(cars\$origin)2</pre>	0.33649	0.07247	4.643	4.72e-06	***
<pre>factor(cars\$origin)3</pre>	0.36505	0.07072	5.162	3.93e-07	***

Multiple R-squared: 0.8242, Adjusted R-squared: 0.8205



NOTE: Coefficient significances and model fit (R^2) have not changed

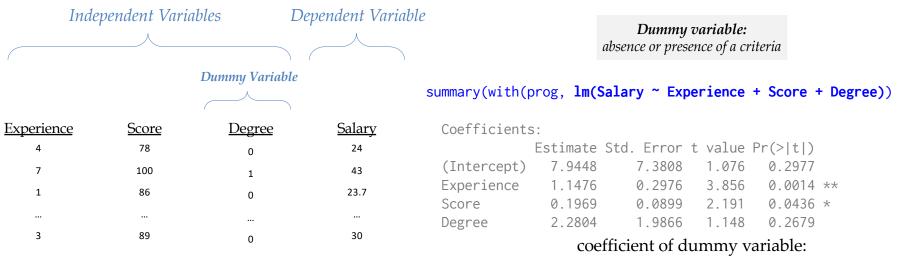
Standardized coefficients are in standard deviations

Relative interpretation of coefficients: "-0.72 sd of mpg per sd weight"

Easier to compare coefficients? "-0.72 sd change from weight" "0.36 sd change from years"

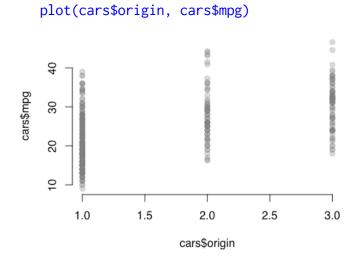
Categorical Variables

Two-level categorical variable: use a single dummy variable



change in Salary (dependent) when Degree (dummy) is present

Multilevel Categorical Variable: use multiple (k-1) dummy variables



summary(lm(mpg ~ factor(origin), data=cars))

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	20.0835	0.4056	49.517	<2e-16 ***
<pre>factor(origin)2</pre>	7.8079	0.8658	9.018	<2e-16 ***
<pre>factor(origin)3</pre>	10.3671	0.8264	12.544	<2e-16 ***

Multilevel categorical (k levels): *k-1 dummies*

coefficient of dummy variables:

factor(origin)2 (JP): change in mpg relative to origin 1 (US) factor(origin)3 (EU): change in mpg relative to origin 1 (US)

Regression Assumptions

Requirements for Regression:

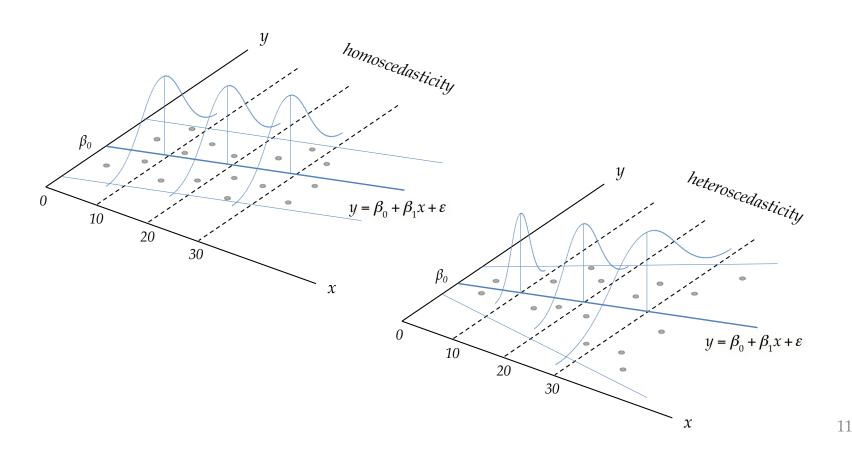
- 1. The error terms (*E*) are: random, normally distributed, and have a mean of zero
- 2. The variance of \mathcal{E} is the same for all values of x
- 3. The values of \mathcal{E} are independent across values of x

Implications:

The mean values of y are on the regression line Standard errors of $\hat{\beta}$ are symmetrically distributed

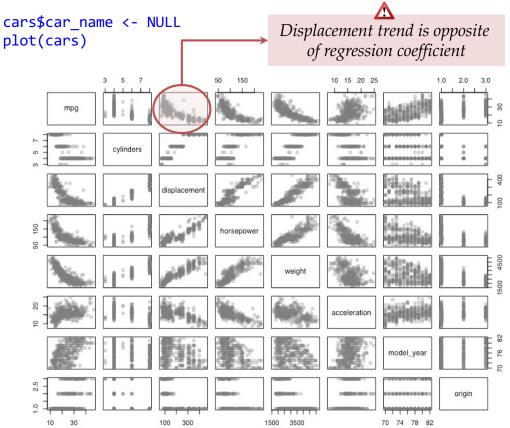
The distribution of y is the same across values of x

The values of y are independent of each other

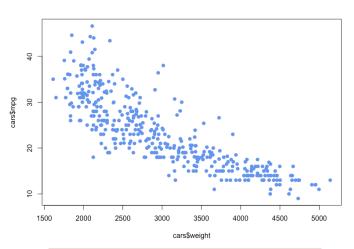


Diagnosing our Regression

Bivariate Scatter Plots



plot(cars\$weight, cars\$mpg, pch=19)



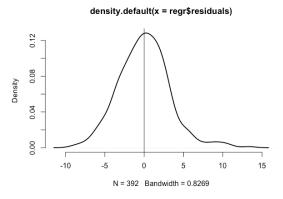
⚠ Non-linearity present in the data

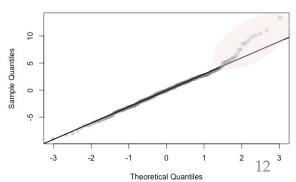
Normality of Residuals (assumption 1)

density plot of residuals
plot(density(regr\$residuals), ...)
abline(v=mean(regr\$residuals))

normal quantiles vs. empirical quantiles
qqnorm(regr\$residuals, ...)
qqline(regr\$residuals, ...)

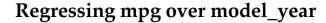
⚠ Residuals are not normally distributed



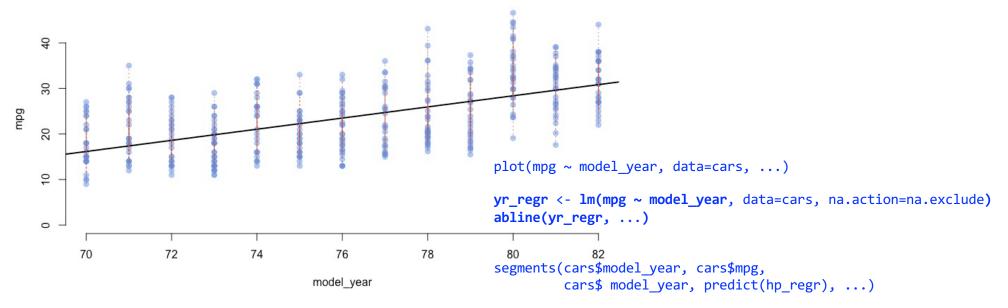


Normal Q-Q Plot

Linear Relationships

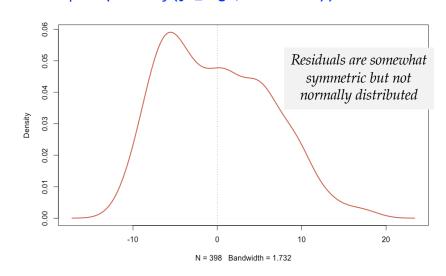


lm(mpg ~ model_year)



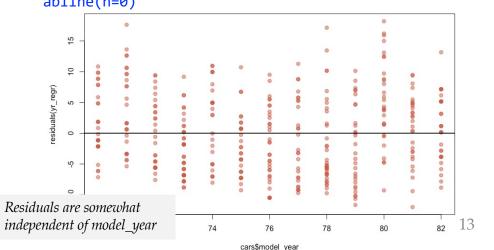
Distribution of residuals

plot(density(yr_regr\$residuals))



Residuals as a function of independent variable

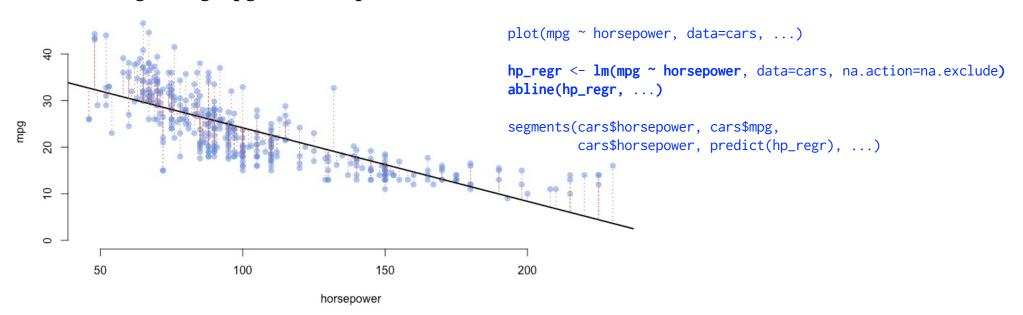
plot(cars\$model_year, resid(yr_regr), col="red", lwd=2)
abline(h=0)



Non-linear Relationships

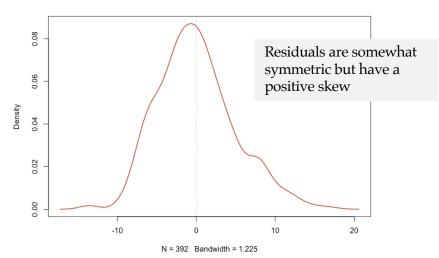
Regressing mpg over horsepower

lm(mpg ~ horsepower)



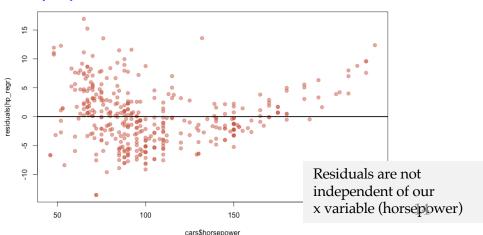
Distribution of residuals

plot(density(hp_regr\$residuals))

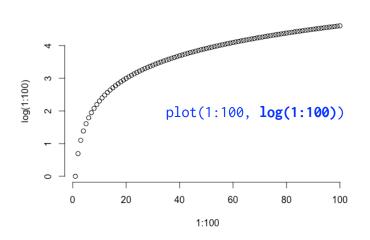


Residuals as a function of independent variable

plot(cars\$horsepower, resid(hp_regr), col="red", lwd=2)
abline(h=0)



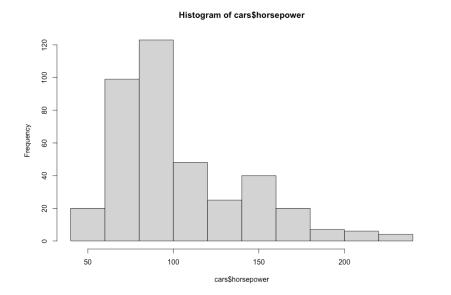
Log Transforming Variables



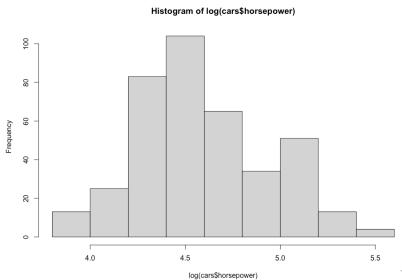
Logarithm is the inverse of exponentiation

Taking the log of variables can give them more symmetric distributions

hist(cars\$horsepower)



hist(log(cars\$horsepower))

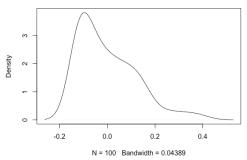


Non-linear model: $y = b^{-x}$

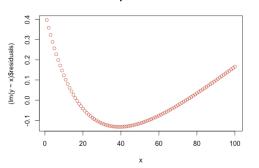
Non-linear relationships pose problems for linear regression

residuals are skewed

 $density.default(x = Im(y \sim x) \\ residuals)$



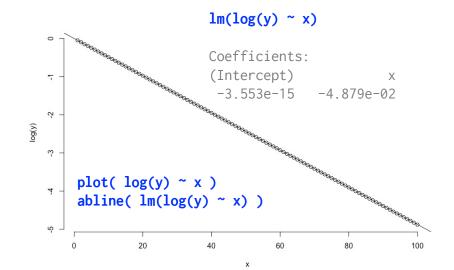
residuals have relationship with independent variable



$$plot(x, (lm(y \sim x)sesiduals))$$

Linear model: log(y) = bx

Log transformation: $log(b^x) = log(b) \cdot x$



We can *log transform* one or both sides of our regression to get linear relationships

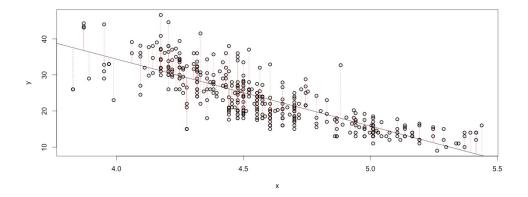
$lm(y \sim log(x))$

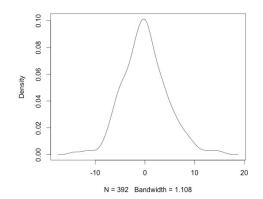
1% change in X leads to $\beta/100$ units change in Y

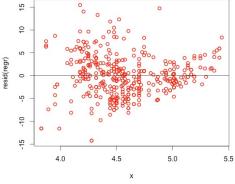
lm(mpg ~ log(horsepower))

Coefficients:
(Intercept) log(horsepower)
108.70 -18.58

1% more horsepower leads to 0.18 decrease in mpg







$lm(log(y) \sim x)$

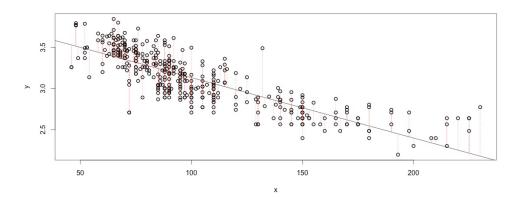
one unit change in X leads to β*100% change in Y

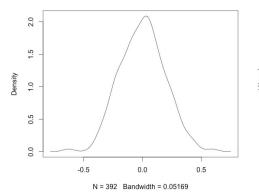
lm(log(mpg) ~ horsepower)

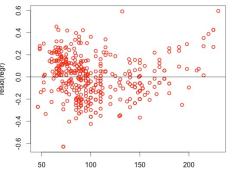
Coefficients:

(Intercept) horsepower 3.864467 -0.007334

One more horsepower leads to 0.73% decrease in mpg







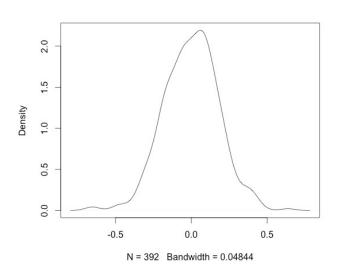
$lm(log(y) \sim log(x))$

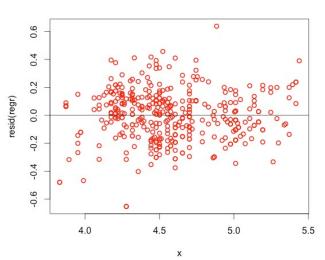
1% change in X leads to β % change in Y

$lm(log(mpg) \sim log(horsepower))$

Coefficients:

(Intercept) log(horsepower)
6.9606 -0.8418





Multicollinearity

```
round( cor(cars[, -8], use="pairwise.complete.obs"), 2 )
                 mpg cylinders displacement horsepower weight acceleration model_year
                1.00
 mpg
 cylinders
               -0.78
                          1.00
                          0.95
 displacement -0.80
                                       1.00
 horsepower
               -0.78
                          0.84
                                       0.90
                                                  1.00
 weight
               -0.83
                          0.90
                                       0.93
                                                  0.86
                                                         1.00
 acceleration 0.42
                         -0.51
                                      -0.54
                                                  -0.69
                                                        -0.42
                                                                       1.00
 model_year
                0.58
                         -0.35
                                      -0.37
                                                  -0.42 -0.31
                                                                       0.29
                                                                                  1.00
```

High multicollinearity: cylinders, displacement, horsepower, weight

weight is the only significant variable among collinear set

```
mpg = \beta_1 displacement
```

$mpg = \beta_1 displacement + \beta_2 weight$

Remove 8th column ("country of origin")

Symptoms of Multicollinearity

When new independent terms added:

- Change in size of coefficient
- Change in standard error of coefficient (change in std error, t-value, p-value)

Variance Inflation Factor (VIF)

Diagnosing Multicollinearity

Given a regression with k independent variables:

$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k$	(T.)	Estimate Std. Error t value Pr(> t)			la alla	Coefficient (slope)	-0.4897
	(Intercept) cylinders			-3.839 0.000145 ** -1.524 0.128215		Std. Error of slope	0.3212
	displacement	2.398e-02	7.653e-03	3.133 0.001863 **	k	Variance of coeff.	0.1032
regr <- lm(mpg ~ cylinders + displacement + horsepower + weight + acceleration + model_year + factor(origin),)							

Multicollinearity inflates the variance of coefficients by:

$$VIF_j = \frac{1}{1 - R_i^2}$$

 R_j^2 is the R^2 of regressing the j^{th} independent variable over all other k-1 independent variables:

$$X_j = b_0 + b_2 X_1 + \cdots + b_k X_k$$

inflated variance reduces significance of coefficients

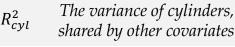
cylinders

e.g., VIF of cylinders in our regression:

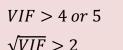
```
cylinders = b_0 + b_1 displacement + \cdots + b_k origin
  cylinder_regr <- lm(cylinders ~ displacement + horsepower + weight + acceleration + model_year + factor(origin),
                        data=cars, na.action=na.exclude)
                                                                              The variance of cylinders,
  r2_cylinders <- summary(cylinder_regr)$r.squared</pre>
```



What can we do about highly collinear independent variables?



High multicollinearity:





 X_k shares more than half its variance with other independent variables