

Contents

Simple linear regression and correlation	2
Glossary: Terms and Symbols	2
Simple linear regression	3
Correlations	20
Bivariate Normal Distribution	28
Multiple linear regression	32
Multiple regression overview	32
Assumptions, diagnostics, remedies	44
Transformations	56
Multiplicative models	63
Multicollinearity: definition and effects	71
Sums of squares and F tests	83
Dummy Variables	91

Glossary: Regression Terms and Symbols

Term	ACT	JWHT	Other
Error	ϵ	ϵ	e
Residual or estimated error	e	e	\hat{e}
Intercept parameter	β_0	β_0	α
Intercept estimate	$\hat{\beta}_0$	$\hat{\beta}_0$	b_0 or a
Slope parameter	β_1	β_1	β
Slope estimate	$\hat{\beta}_1$	$\hat{\beta}_1$	b_1 or b
Sum of squared errors	SSE	RSS	
Residual sum of squares			
Total sum of squares	SST	TSS	
Regression sum of squares	SSR		
Error variance (parameter)	σ^2	σ^2	σ_ϵ^2
Variance of the errors			
Mean squared error	MSE		S_ϵ^2
Standard deviation of the errors	σ	σ	σ_ϵ
Residual standard error	s	RSE	$S_\epsilon, \hat{\sigma}$
Root mean squared error			RMSE
Standard error of the estimate			

Introduction to Regression

- *Objective*: to quantify the relationship between an interval-level *response* variable and one or more *predictor* variables from a sample of data.
- **Response variable** (also called **dependent**, **criterion**, or **output** (Y) variable): a variable that we wish to study and is causally dependent on other variables
- Note: we also study categorical dependent variables, but call it the **classification problem**
- **Predictor variables** (also called **independent** (X) variables, **covariates**, or **inputs**): variables that are related to the response variable
 - **Factors** refer to categorical variables
 - * **Binary variable** (aka **dichotomous**): takes two values, generically “success/failure” or “yes/no”
 - * **Nominal variable**: no ordering assumed, e.g., race
 - * **Ordinal variable**: values can be ordered, i.e., $\exists <$ operator, e.g., [Likert scales](#)
 - Numerical variables
 - * **Count**: takes values $0, 1, 2, \dots$
 - * **Amount**: takes non-negative, real values

Other terms

- **Descriptive research:** describe some (sub)population, e.g., with univariate/bivariate descriptive statistics/graphs
- *Why build a predictive model?*
 - **Prediction:** estimate response variable accurately (typically on existing data)
 - **Exploration:** discover which x 's are *associated with y* (**insights** and future hypotheses)
 - **Causal inference (prescription, confirmatory):** how predictors *cause y* (**intervention**)
- **Randomized-controlled experiment:** investigator assigns values of at least some independent variables (**treatments**)
- **Observational study:** independent variables not under the control of data scientist
- **Cross-sectional data:** measures on many sampling units at a fixed point in time
- **Time-series data:** measures recorded at equal-spaced points in time usually on a single sampling unit
- **Panel data:** measures on the same sampling units over time
- **Censored data:** values of outcome (e.g., time of failure, churn, death) only partially known

Simple Linear Regression

Assume that

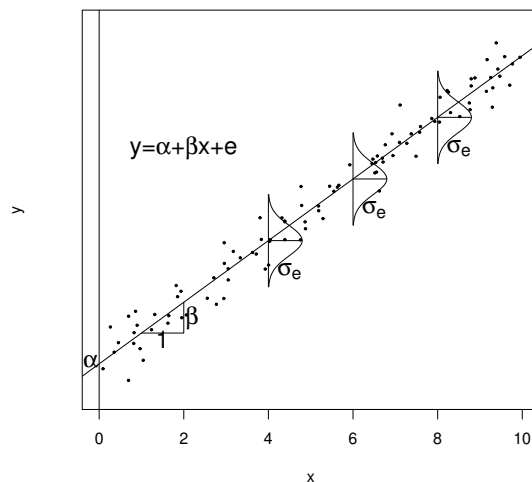
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i,$$

where **error** $\epsilon_i \sim \mathcal{N}(0, \sigma_\epsilon^2)$, and ϵ_i independent of ϵ_j ($\epsilon_i \perp\!\!\!\perp \epsilon_j$) for $i \neq j$ (and $\epsilon_i \perp\!\!\!\perp x_i$ if x random). This model implies that

$$\mathbb{E}(Y|x) = \mu_{Y|x} = \beta_0 + x\beta_1.$$

This is called the **regression of y on x** .

- β_0 : the **intercept**. On average $Y = \beta_0$ when $x = 0$.
- β_1 : the **slope**. Every unit increase in x is associated with an increase in Y of β_1 , *on the average*
- σ_ϵ : **standard deviation of the errors** (σ_ϵ^2 called **error variance**). If errors are normal, empirical rule tells us for any x , 68% of points will fall within one σ_ϵ of the mean ($\beta_0 + \beta_1 x$).



Estimating the Regression Model

- In practice, we don't know values of **parameters** β_0 , β_1 , and σ_ϵ^2 and must estimate them with $\hat{\beta}_0$, $\hat{\beta}_1$, and S_ϵ , respectively
- The estimate of $\mu_{Y_i|x_i}$ is denoted by **fitted, predicted** or **y-hat value** $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$
- The **residual** for observation i is

$$\hat{\epsilon}_i = y_i - \hat{y}_i = y_i - a - bx_i$$

- We choose $\hat{\beta}_0$ and $\hat{\beta}_1$ so that they minimize the **least-squares criterion** (SSE means **sum of squared errors**):

$$\text{SSE} = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2,$$

- The **ordinary least-squares** (OLS) estimates are

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = r \frac{S_y}{S_x} \quad \text{and} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Note, for every SD change in x , there is a change of r SD's in y , on average, where r is the Pearson correlation.

- Estimate σ_ϵ^2 with the *mean squared error*

$$S_\epsilon^2 = \frac{\text{SSE}}{n-2} = S_y^2 (1 - r^2) \frac{n-1}{n-2}$$

σ_ϵ called **residual standard error**

Sampling Distribution of Parameters

- Theorem: **standard errors** are given by

$$S_{\hat{\beta}_1} = \frac{S_\epsilon}{S_x \sqrt{n-1}} \quad \text{and} \quad S_{\hat{\beta}_0} = S_\epsilon \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_x^2(n-1)}}$$

- Theorem:
 - **Unbiasedness**: $\mathbb{E}(\hat{\beta}_1) = \beta_1$ and $\mathbb{E}(\hat{\beta}_0) = \beta_0$
 - $\hat{\beta}_1$ and $\hat{\beta}_0$ have normal **sampling distributions**. The following have t distributions with $n - 2$ degrees of freedom:

$$\frac{\hat{\beta}_1 - \beta_1}{S_{\hat{\beta}_1}} \quad \text{and} \quad \frac{\hat{\beta}_0 - \beta_0}{S_{\hat{\beta}_0}}$$

- **Total sum of squares**: $\text{SST} = \sum (y_i - \bar{y})^2 = (n - 1)S_y^2$
- Definition: The **coefficient of determination** (aka R^2) is the square of the correlation and gives the percentage of the variability of y that is “explained” by x

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

- **Gauss-Markov Theorem**: $\hat{\beta}_1$ has a variance that is “smaller” than that of any other linear, unbiased estimator and is called the *best linear unbiased estimator* (“BLUE”).
- Optional videos deriving OLS theory for simple linear regression [Part 1](#) and [Part 2](#)

Click Ball Point Pens Example

y_i Sales in territory $i = 1, \dots, 40$

x_{i1} Advertising (number of TV spots) in territory i

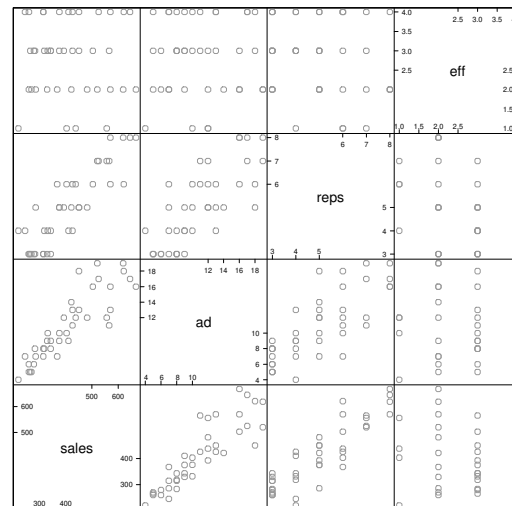
x_{i2} Number of sales reps in territory i

x_{i3} Wholesaler efficiency index in territory i (4=outstanding, 3=good, 2=average, 1=poor)

```
click = data.frame(sales=c(260.3,286.1,279.4,410.8,438.2,315.3,565.1,570.0,426.1,315.0,
  403.6,220.5,343.6,644.6,520.4,329.5,426.0,343.2,450.4,421.8,245.6,503.3,375.7,265.5,
  620.6,450.5,270.1,368.0,556.1,570.0,318.5,260.2,667.0,618.3,525.3,332.2,393.2,283.5,
  376.2,481.8), ad=c(5,7,6,9,12,8,11,16,13,7,10,4,9,17,19,9,11,8,13,14,7,16,9,5,18,18,
  5,7,12,13,8,6,16,19,17,10,12,8,10,12), reps=c(3,5,3,4,6,3,7,8,4,3,6,4,4,8,7,3,6,3,5,
  5,4,6,5,3,6,5,3,6,7,6,4,3,8,8,7,4,5,3,5,5), eff=c(4,2,3,4,1,4,3,2,3,4,1,1,3,4,2,2,4,
  3,4,2,4,3,3,3,4,3,2,2,1,4,3,2,2,2,4,3,3,3,4,2))
```

```
> round(cor(click), 4)
      sales    ad    reps    eff
sales 1.0000 0.8802 0.8818 0.0019
ad    0.8802 1.0000 0.7763 0.0321
reps  0.8818 0.7763 1.0000 -0.1896
eff    0.0019 0.0321 -0.1896 1.0000
```

```
> plot(click)
```



This is a **scatterplot matrix**

Click Ball Point Pens Example

```
import pandas as pd
click = pd.DataFrame({"sales": [260.3, 286.1, 279.4, 410.8, 438.2, 315.3, 565.1, 570.0, 426.1, 315.0,
                                403.6, 220.5, 343.6, 644.6, 520.4, 329.5, 426.0, 343.2, 450.4, 421.8, 245.6, 503.3, 375.7, 265.5,
                                620.6, 450.5, 270.1, 368.0, 556.1, 570.0, 318.5, 260.2, 667.0, 618.3, 525.3, 332.2, 393.2, 283.5,
                                376.2, 481.8],
                      "ad": [5, 7, 6, 9, 12, 8, 11, 16, 13, 7, 10, 4, 9, 17, 19, 9, 11, 8, 13, 14, 7, 16, 9, 5, 18, 18,
                             5, 7, 12, 13, 8, 6, 16, 19, 17, 10, 12, 8, 10, 12],
                      "reps": [3, 5, 3, 4, 6, 3, 7, 8, 4, 3, 6, 4, 4, 8, 7, 3, 6, 3, 5, 5, 4, 6, 5, 3, 6, 5, 3, 6, 7, 6, 4, 3, 8, 8, 7, 4, 5, 3, 5, 5],
                      "eff": [4, 2, 3, 4, 1, 4, 3, 2, 3, 4, 1, 1, 3, 4, 2, 2, 4, 3, 4, 2, 4, 3, 3, 3, 4, 3, 2, 2, 1, 4, 3, 2, 2, 2, 4, 3, 3, 3, 4, 2]})
```

```
# or if you have a csv file
click = pd.read_csv("teach/data/click.csv")
```

```
click.describe()
      sales      ad      reps      eff
count  40.000000  40.000000  40.000000  40.000000
mean   411.287500  10.900000   5.000000   2.825000
std    123.854032   4.307418   1.64862   0.98417
min     220.500000   4.000000   3.000000   1.000000
25%    315.225000   7.750000   3.750000   2.000000
50%    398.400000  10.000000   5.000000   3.000000
75%    507.575000  13.250000   6.000000   4.000000
max     667.000000  19.000000   8.000000   4.000000
```

```
click.corr()
      sales      ad      reps      eff
sales  1.000000  0.880156  0.881778  0.001917
ad      0.880156  1.000000  0.776312  0.032057
reps    0.881778  0.776312  1.000000 -0.189638
eff      0.001917  0.032057 -0.189638  1.000000
```

```
from pandas.tools.plotting import scatter_matrix
scatter_matrix(click, figsize=(10,10))
scatter_matrix(click, figsize=(10,10), diagonal="kde")
```

```
import statsmodels.formula.api as sm
fit = sm.ols(formula="sales ~ ad", data=click).fit()
fit.summary()
```

OLS Regression Results

```
=====
Dep. Variable:          sales    R-squared:                0.775
Model:                  OLS      Adj. R-squared:           0.769
Method:                 Least Squares    F-statistic:        130.6
Date:                  Fri, 01 Jan 2016    Prob (F-statistic):    7.33e-14
Time:                  10:04:17      Log-Likelihood:       -219.21
No. Observations:      40          AIC:                  442.4
```

```

Df Residuals:          38    BIC:          445.8
Df Model:              1
Covariance Type:      nonrobust

```

```

=====
              coef    std err          t      P>|t|      [95.0% Conf. Int.]
-----
Intercept    135.4336     25.907      5.228      0.000      82.989    187.879
ad           25.3077      2.214     11.430      0.000      20.825     29.790
=====
Omnibus:            4.122    Durbin-Watson:           1.721
Prob(Omnibus):      0.127    Jarque-Bera (JB):           2.823
Skew:               0.535    Prob(JB):             0.244
Kurtosis:           3.740    Cond. No.              32.4
=====

```

```

import matplotlib.pyplot as plt
plt.plot(click["ad"], click["sales"], 'o') # o option does not connect dots
plt.plot(click["ad"], fit.fittedvalues)
plt.legend(['Observed', 'Fitted'], loc='upper right')
plt.xlabel('Number of Ads')
plt.ylabel('Sales')
plt.title('Advertising and Sales')
plt.show()

```

```

fit = sm.ols(formula="sales ~ ad+reps+eff", data=click).fit()
fit.summary()

```

OLS Regression Results

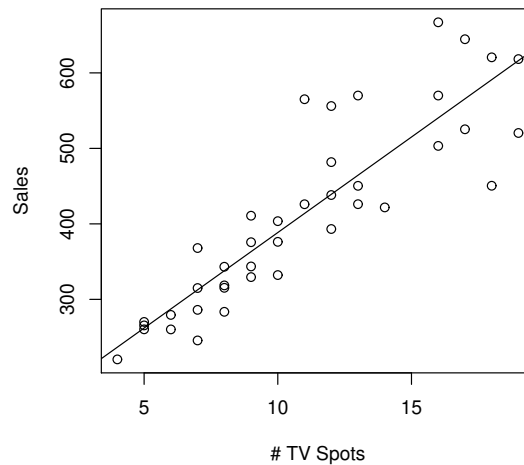
```

=====
Dep. Variable:          sales    R-squared:          0.881
Model:                  OLS      Adj. R-squared:        0.871
Method:                 Least Squares    F-statistic:        89.05
Date:                  Fri, 01 Jan 2016    Prob (F-statistic): 1.02e-16
Time:                  09:48:52    Log-Likelihood:     -206.40
No. Observations:      40    AIC:              420.8
Df Residuals:          36    BIC:              427.6
Df Model:              3
Covariance Type:      nonrobust
=====
              coef    std err          t      P>|t|      [95.0% Conf. Int.]
-----
Intercept     31.1504     34.175      0.911      0.368     -38.160    100.461
ad             12.9682      2.737      4.738      0.000       7.417     18.520
reps          41.2456      7.280      5.666      0.000      26.481     56.010
eff           11.5243      7.691      1.498      0.143      -4.074     27.123
=====
Omnibus:            0.993    Durbin-Watson:           2.104
Prob(Omnibus):      0.609    Jarque-Bera (JB):           0.914
Skew:               0.153    Prob(JB):             0.633
Kurtosis:           2.326    Cond. No.              65.1
=====

```

Estimates From Click Ball Point Pens

```
> fit = lm(sales ~ ad, click)
> plot(click$ad, click$sales,
       xlab="# TV Spots", ylab="Sales")
> abline(fit)
> summary(fit)
> plot(fit) # gives diagnostic plots
```



```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  135.434     25.907    5.228 6.50e-06 ***
ad           25.308       2.214   11.430 7.33e-14 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 59.56 on 38 degrees of freedom
Multiple R-squared:  0.7747, Adjusted R-squared:  0.7687
F-statistic: 130.6 on 1 and 38 DF, p-value: 7.327e-14
```

- Residual standard error: $S_e = 59.56$, standard error of estimate
- R-square = 0.7747: fraction of variation explained by model
- Adj R-sq: 0.7687 adjusted for number of parameters

Interpretation of Output

- The estimated regression model is

$$\hat{y} = 135 + 25.3\text{ad}$$

What does 25.3 tell us? What about 135?

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  135.434      25.907    5.228 6.50e-06 ***
ad           25.308       2.214   11.430 7.33e-14 ***
```

- Standard errors are $S_{\hat{\beta}_0} = 25.91$ and $S_{\hat{\beta}_1} = 2.21$
- A 95% CI for β_1 : $25.3 \pm 2.02 \times 2.21 \approx (20.8, 29.8)$

```
> confint(fit)
              2.5 %    97.5 %
(Intercept) 82.98862 187.87857
ad          20.82538  29.79001
```

- To test the hypotheses (with Type I error rate .05)

$$H_0 : \beta_1 = 0 \quad \text{versus} \quad H_1 : \beta_1 \neq 0,$$

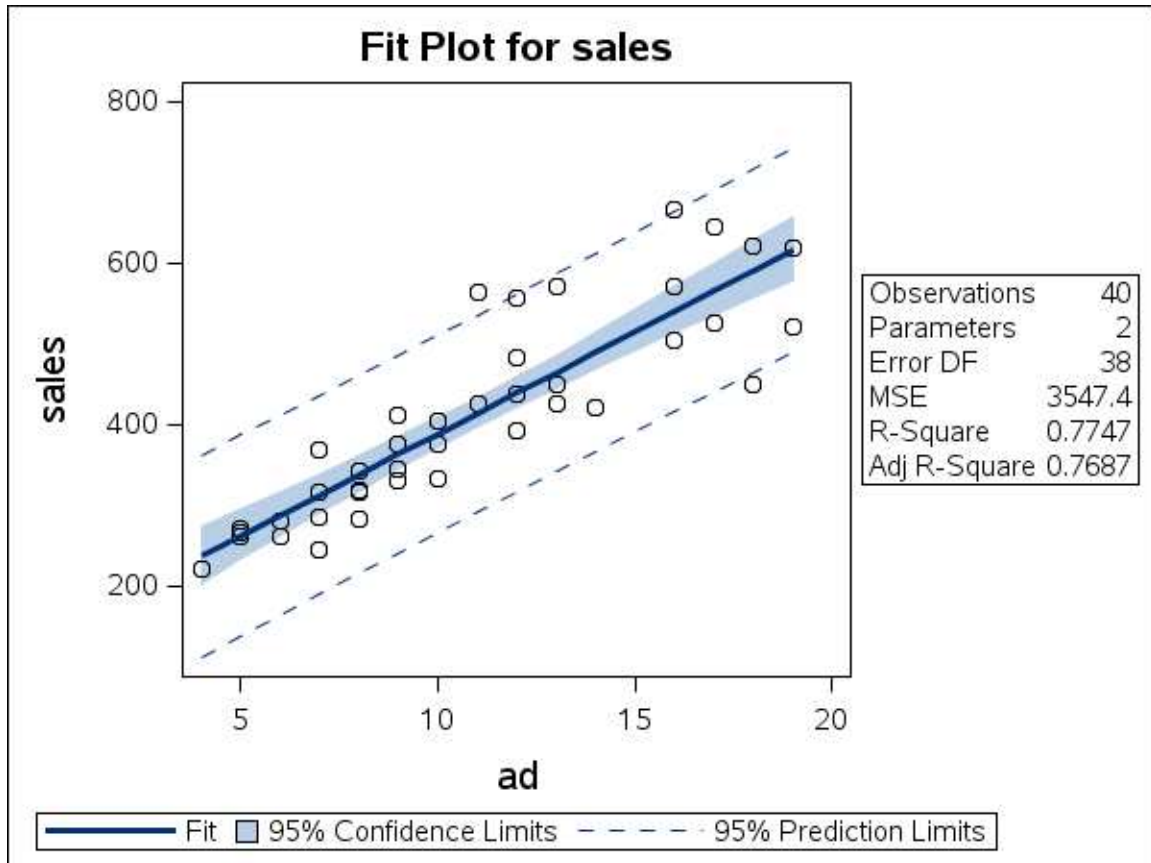
$P\text{-value} = 7.33 \times 10^{-14} < .05$, so reject H_0 and conclude $\beta_1 \neq 0$.

- The expected sales when advertising is 5 (spots) is

$$\hat{y} = 135 + 25.3 \times 5 = 261.97$$

```
> predict(fit, data.frame(ad=5))
261.9721
```

Prediction and Confidence Intervals



- **Confidence interval for mean prediction** (shaded blue area): 95% confidence interval for $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$, i.e., indicates sampling variation of predicted values $\mathbb{E}(Y|x)$.
- **Prediction interval** (dashed lines): indicates where the middle 95% of the distribution of Y for a given x_0 falls. If we knew parameters, it would be $(\beta_0 + \beta_1 x_0) \pm 1.96\sigma_\epsilon$.

Additional Results

- The standard error of a new observation Y given x_0 :

$$S_{Y|x_0} = \sqrt{S_\epsilon^2 \left(1 + \frac{1}{n}\right) + S_{\hat{\beta}_1}^2 (x_0 - \bar{x})^2}$$

Example: find a 95% prediction interval for the mean sales when there are 5 ads. Hint: the 97.5 percentile of a t distribution with $40 - 2 = 38$ degrees of freedom is 2.024.

```
predict(fit, data.frame(ad=5), interval="prediction")
```

$$S_{Y|x_0} = \sqrt{3547 \left(1 + \frac{1}{40}\right) + 2.214^2 (5 - 10.90)^2} = 61.70$$

$$262.0 \pm 2.024 \times 61.699 = (137.1, 386.9)$$

- The standard error of a predicted value \hat{Y} given x_0

$$S_{\hat{Y}|x_0} = \sqrt{\frac{S_\epsilon^2}{n} + S_{\hat{\beta}_1}^2 (x_0 - \bar{x})^2}$$

- Example: find a confidence interval for the mean sales when there are 5 ads.

```
predict(fit, data.frame(ad=5), interval="confidence")
```

$$S_{\hat{Y}|x_0} = \sqrt{\frac{3547}{40} + 2.214^2 (5 - 10.90)^2} = 16.104$$

$$262.0 \pm 2.024 \times 16.104 = (229.4, 294.6)$$

Geometric Interpretation of Simple Regression

- Assume \mathbf{x} and \mathbf{y} are mean-centered vectors
- We can show that the regression coefficient is

$$b = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}} = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y}$$

- Recall that the *projection of \mathbf{y} onto \mathbf{x}* is

$$\text{proj}_{\mathbf{x}} \mathbf{y} = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}} \mathbf{x} = b\mathbf{x} = \hat{\mathbf{y}}$$

- This is a *right* triangle with sides \mathbf{y} , and $\hat{\mathbf{y}}$ (which lies on \mathbf{x}) and $\hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}}$. The lengths of the sides are
 - (Hypotenuse) $\|\mathbf{y}\| = \sqrt{\text{SST}}$
 - $\|\hat{\mathbf{y}}\| = \sqrt{\text{SSR}}$
 - $\|\hat{\mathbf{e}}\| = \|\mathbf{y} - \hat{\mathbf{y}}\| = \sqrt{\text{SSE}}$

The sum of squares equality $\text{SST} = \text{SSR} + \text{SSE}$ follows from the Pythagorean theorem

- $\text{SSR} = \|\hat{\mathbf{y}}\|^2 = \|b\mathbf{x}\|^2 = b^2 \mathbf{x}^\top \mathbf{x} = b^2 S_x^2 (n - 1)$

Your Turn

1. You have five machines. The following data gives the age of the machines in years and the annual maintenance cost in thousands of dollars:

```
> machine = data.frame(age=c(2,5,9,3,8), cost=c(6,13,23,5,22))
```

- (a) Draw a scatterplot and describe the relationship.
 - (b) Find the correlation between age and cost. What does it tell you?
 - (c) Find the equation predicting cost from age. **Interpret the slope and intercept.**
 - (d) Superimpose the regression line on your scatterplot.
 - (e) What would you expect the annual maintenance to be for a machine that is 7 years old?
 - (f) What is a typical size for the prediction errors?
 - (g) What fraction of the variation in cost is explained by knowing the age of the machines?
 - (h) Is the relationship between age and cost statistically significant?
 - (i) A colleague suggested to use \$20,000 per year per machine for planning purposes. Perform a test at the 5% level to see if this is reasonable.
2. The amounts of a chemical compound y , which was dissolved in 100 grams of water at various temperatures, x° C, were recorded:

```
dat = data.frame(  
  x=c(0,0,0, 15,15,15, 30,30,30, 45,45,45, 60,60,60, 75,75,75),  
  y=c(8,6,8, 12,10,14, 25,21,24, 31,33,28, 44,39,42, 48,51,44))
```

- (a) Find the equation of the regression line.
 - (b) Graph the line on a scatterplot.
 - (c) Compute and interpret the standard error of the estimate.
 - (d) Compute and interpret the coefficient of determination.
 - (e) Find a 99% CI for β_0
 - (f) Find a 99% CI for β_1
 - (g) Test at the 1% level if the slope differs from 0.
 - (h) Estimate the mean amount that will dissolve in 100 grams of water at 50° C.
 - (i) Find a 99% CI for the mean amount that will dissolve at 50° C.
 - (j) Find a 99% PI for the amount that will dissolve at 50° C.
3. Consider the circulation and the open line rate (price per line for an ad placed just once) for selected large newspapers shown below


```

dat=data.frame(
circ = c(2081995,1374858,1284613,1057536,970051,963069,828236,779259,768288,
691771,663693,657015,645623,533384,528777,514702,492002,486426,
443592,349182),
linerate=c(37.65,18.48,14.50,14.61,16.47,16.07,13.82,13.05,13.78,12.25,10.53,
14.18,12.83,7.81,5.17,11.08,6.58,8.77,6.03,6.77),
row.names=c("WSJ","NY Daily News","USA Today","LA Times","NYT", "NY Post",
"Philadelphia","Chi Tribune","Wash Post","SF Chronicle","Chi Sun Times",
"Detroit News","Detroit Free Press","Long Island Newsday","KC Times",
"Miami Herald","Cleveland","Milwaukee","Houston","Baltimore"))

```

- (a) Create a scatterplot of the open line rate against circulation. Comment.
 - (b) Find and interpret the correlation of open line rate with circulation. Is the correlation reasonable from a business perspective?
 - (c) Find the regression equation to predict open line rate from circulation. Superimpose the regression line on your scatterplot.
 - (d) Test if the association between the open-line rate and circulation is significant ($\alpha = .05$).
 - (e) Find the predicted and residual value for the *New York Times*. Interpret these values. In particular, is the open line rate higher or lower than what you would expect for a newspaper with its circulation?
 - (f) There is an outlier visible in the scatterplot. Let's test to see if it could reasonably be from the same population as the others by treating it as a new observation. Remove *The Wall Street Journal* (WSJ) from the data set and find the regression equation to predict the open line rate from circulation for the other newspapers.
 - (g) Find the two-sided 95% prediction interval (PI) for a new observation, with X_0 being the circulation of *WSJ*.
 - (h) Test whether or not *WSJ* is an outlier by seeing if its open line rate is in the PI.
 - (i) The *milline rate* is defined as the open line rate divided by the circulation, in millions. Thus, it is the cost per line of advertising per million circulation. This adjustment should take care of some of the differences in advertising rates due to circulation. That is, one explanation of the open line rate is that it is proportional to circulation. If it is just proportional, there should be nothing left in the milline rate to be explained by circulation. On the other hand, if there is an additional advantage or penalty to being big, the circulation should help explain the variation in milline rates. Let's use regression analysis to see if there is anything left in the milline rate to be explained by circulation. Draw a scatterplot of milline rate against circulation.
 - (j) Find and interpret the correlation between circulation and milline rate.
 - (k) What percentage of the variation in milline rate is explained by circulation?
 - (l) Test to see if there is a significant relationship between circulation and milline rate.
 - (m) Write a paragraph explaining and interpreting your results.
4. The data below gives mailing-list size (thousands of names) and sales (thousands of dollars) for a group of catalogs.

```
dat = data.frame(
  size=c(168, 21, 94, 39, 249, 43, 589, 41),
  sales = c(5178, 2370, 3591, 2056, 7325, 2449, 15708, 2469))
```

- (a) How strong is the association between these two variables? Find the appropriate summary measure and interpret it.
 - (b) Find the equation to predict sales from the size of the mailing list.
 - (c) What level of sales would you expect for a catalog mailed to 5,000 people?
 - (d) What percent of the variation in the list size can be explained by the fact that some generated more sales than others?
 - (e) Is there a significant relationship between list size and sales? How do you know?
5. JWHT problem 8a,b on pages 121–2 (Hint: see §2.3.4 on page 48–49.) If you use the data from the author’s website you will need to read about the `na.strings` option. Note: omit part c for now. Answer these questions about the output. Type the following to read it in:

```
auto = read.csv("Downloads/auto.csv", na.strings="?")
auto$origin = factor(auto$origin, 1:3, c("US", "Europe", "Japan"))
```

- (a) What is the estimated regression equation?
- (b) What does the slope tell you?
- (c) How much uncertainty is associated with the slope estimate?
- (d) What does the residual standard error tell you?
- (e) Using this model, is there a significant relationship between mpg and horsepower?
- (f) What fraction of the variation in mpg is explained by using this linear function of horsepower?
- (g) What is the predicted mpg associated with a horsepower of 98?
- (h) What is the 95% prediction interval for the predicted mpg associated with a horsepower of 98?
- (i) What is the 99% confidence interval for the mean prediction of mpg when horsepower is 98?
- (j) What is a 90% confidence interval for the slope?
- (k) In looking at the scatterplot and fitted model, note any violations of the model assumptions.

Answers

1. Machine problem. Do in live session?

2. See code below. (a) $\hat{y} = 5.8254 + 0.5676x$. (c) $S_e = 2.57$:

typical size of a residual (different between observed and predicted) (d) $R^2 = 0.973$. (e) $[2.686, 8.965]$. (f) $[0.498, 0.637]$. (g) $H_0 : \beta_1 = 0$ versus $H_1 : \beta_1 \neq 0$, either

$P = 5.7 \times 10^{-14} \ll 0.01$ so reject or $0 \notin [0.498, 0.637]$ so reject H_0 . (h) $\hat{y}(50) = 34.2$. (i) $[32.23569, 36.17701]$. (j) $[26.43824, 41.97446]$.

```
fit = lm(y~x, dat)      # part a
plot(y~x, dat)          # part b
abline(fit)
summary(fit)            # part c, d, g
confint(fit, level=.99) # part f, g
new = data.frame(x=50)
predict(fit, new, interval="conf", level=.99) # part h,i
predict(fit, new, interval="pred", level=.99) # part j
```

3. Newspaper problem. See code below. (a) The relationship is mostly linear but there is an outlier that could have substantial influence. (b) $r = .93$: larger open line rates have a positive association with circulation. (c) $\hat{\text{linerate}} = 0.282 + 0.0000158\text{circ}$; (d) $H_0 : \beta_1 = 0$ versus $H_1 : \beta_1 \neq 0$, $P = 2.9 \times 10^{-9} \ll 0.05$ so reject H_0 . (e) Note that NYT is row 5. See R below to find $y_5 = 16.5$, $\hat{y}_5 = 15.6$ and $\hat{e}_5 = 0.85$; higher. (f) WSJ is paper 1. $\hat{\text{linerate}} = 3.11 + 0.0000117\text{circ}$. (g/h) $37.6 \notin [20.6, 34.2]$, conclude WSJ is an outlier. (i) Not much of a relationship. (j) $r = -0.125$. (k) $R^2 = (-0.125)^2 = 0.0155$. (l) Do either of this: $H_0 : \rho = 0$ versus $H_1 : \rho \neq 0$ or $H_0 : \beta_1 = 0$ versus $H_1 : \beta_1 \neq 0$ $P = 0.6 > 0.05$, so do not reject H_0 .

```
plot(linerate~circ, dat)      # part a
cor(dat)                      # part b
fit = lm(linerate~circ, data=dat) # part c
abline(fit)                   # part c
summary(fit)                  # part c,d
dat[5,]                       # part e
fit$fitted.values[5]
fit$residuals[5]
```

```
fit2 = lm(linerate~circ, data=dat[-1,]) # part f
summary(fit2)
dat[1,]                                # part g
predict(fit2, dat[1,], interval = "pred")
dat$milline=1000000*dat$linerate/dat$circ # part i
plot(milline~circ, dat)
cor(dat)                               # part j
fit3 = lm(milline~circ, dat)           # part k
summary(fit3)
cor.test(dat$milline, dat$circ)
```

4. (a) $r = 0.999$. (b) $\hat{y} = 1393.825 + 24.112\text{size}$. (c) $\hat{y}(5000) = 121954$. (d) $R^2 = 0.997$. (e) $H_0 : \beta_1 = 0$ versus $H_1 : \beta_1 \neq 0$, $P = 5.9 \times 10^{-9} \ll 0.05$ so reject H_0 .
5. JWHT 8. See code below (a) $\text{mpg} = 39.94 - 0.1578\text{horsepower}$. (b) Every unit increase in horsepower is associated with a .1578 decrease in mpg on the average. (c) Standard error is 0.006446. (d) Typical size of residuals. (e) Yes, $P < 2 \times 10^{-16} < .05$. (f) $R^2 = .6059$. (g) 24.47 mpg. (h) 14.8094 to 34.12476. (i) 23.81669 to 25.11747. (j) $[-0.1684719, -0.1472176]$. (k) The scatterplot shows that the relationship is not linear and the error variance is not constant.

```
fit = lm(mpg ~ horsepower, auto) # part a
summary(fit)
plot(mpg ~ horsepower, auto)
abline(fit)
predict(fit, data.frame(horsepower=98), interval="pred")
predict(fit, data.frame(horsepower=98), interval="conf",
       level=.99)
confint(fit, level=.90)
plot(fit) # part c
```

Anscombe's Example

```
dat = data.frame(
  dataset = factor(c(rep(1,11),rep(2,11),rep(3,11),rep(4,11)),
    labels=c("I","II","III","IV")),
  x=c(10,8,13,9,11,14,6,4,12,7,5, 10,8,13,9,11,14,6,4,12,7,5,
    10,8,13,9,11,14,6,4,12,7,5, rep(8,10),19),
  y=c(8.04,6.95,7.58,8.81,8.33,9.96,7.24,4.26,10.84,4.82,5.68,
    9.14,8.14,8.74,8.77,9.26,8.1,6.13,3.1,9.13,7.26,4.74,
    7.46,6.77,12.74,7.11,7.81,8.84,6.08,5.39,8.15,6.42,5.73,
    6.58,5.76,7.71,8.84,8.47,7.04,5.25,5.56,7.91,6.69,12.5)
)
```

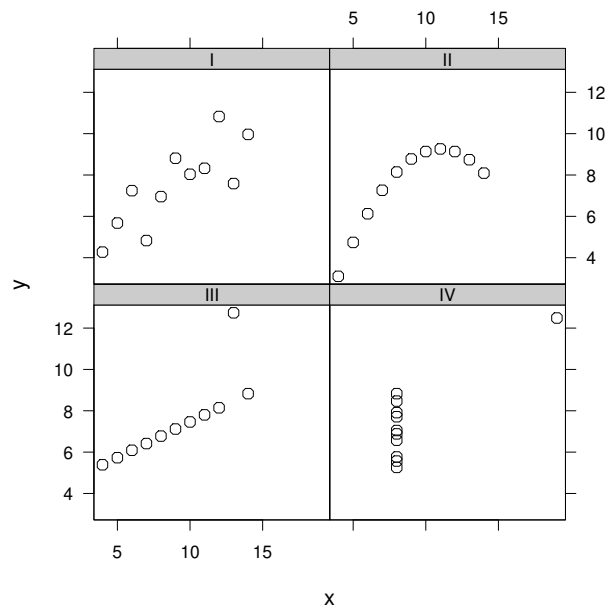
```
> library(lattice)
> xyplot(y~x | dataset, dat)
> tapply(dat$x, dat$dataset, mean)
  I  II III IV
 9   9  9   9

> tapply(dat$y, dat$dataset, mean)
  I      II      III      IV
7.500909 7.500909 7.500000 7.482727

> tapply(dat$x, dat$dataset, var)
  I  II III IV
11 11 11 11

> tapply(dat$y, dat$dataset, var)
  I      II      III      IV
4.127269 4.127629 4.122620 4.151322

> summary(lm(y ~ dataset*x, dat))
```



```

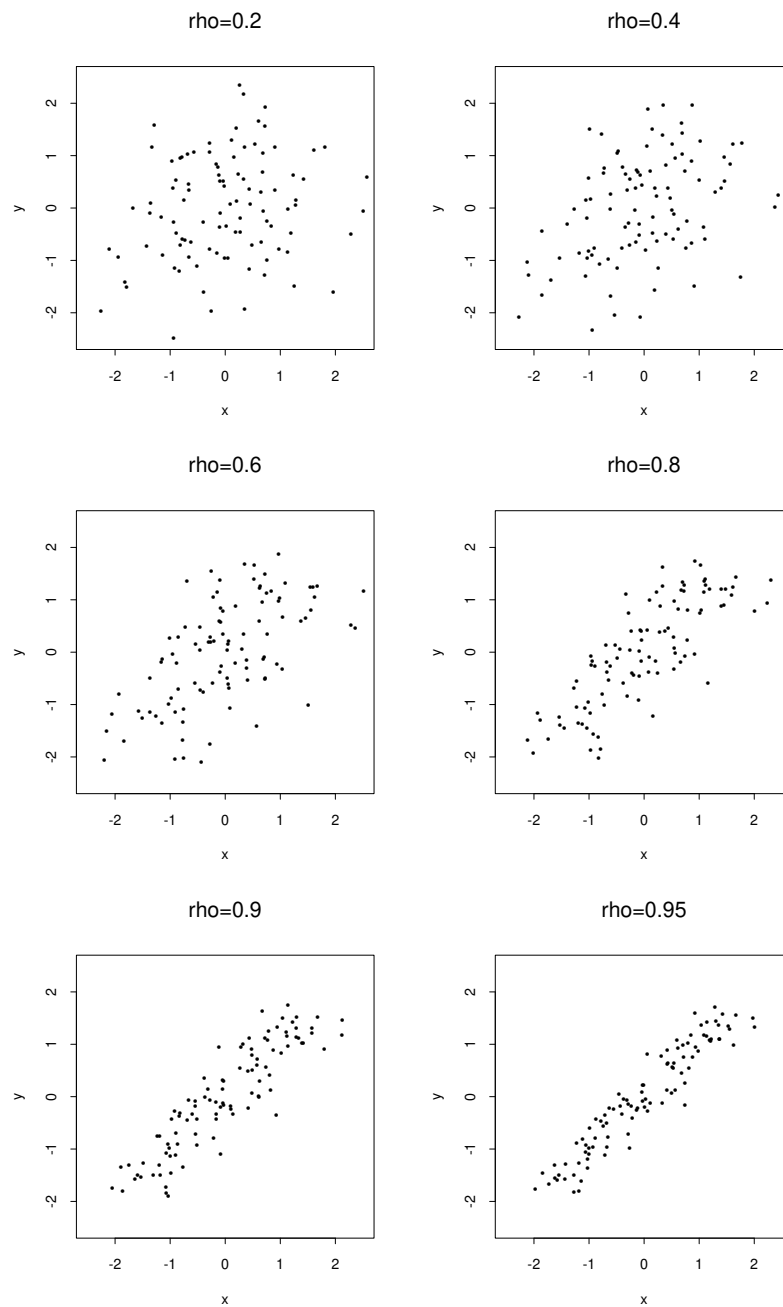
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  3.000e+00  1.125e+00   2.666 0.011431 *
datasetII    8.182e-04  1.592e+00   0.001 0.999593
datasetIII   2.364e-03  1.592e+00   0.001 0.998823
datasetIV   -3.291e-02  1.592e+00  -0.021 0.983618
x             5.001e-01  1.180e-01   4.239 0.000149 ***
datasetII:x  -9.091e-05  1.668e-01  -0.001 0.999568
datasetIII:x -3.636e-04  1.668e-01  -0.002 0.998273
datasetIV:x   1.636e-03  1.668e-01   0.010 0.992229
```

Always look at your data!

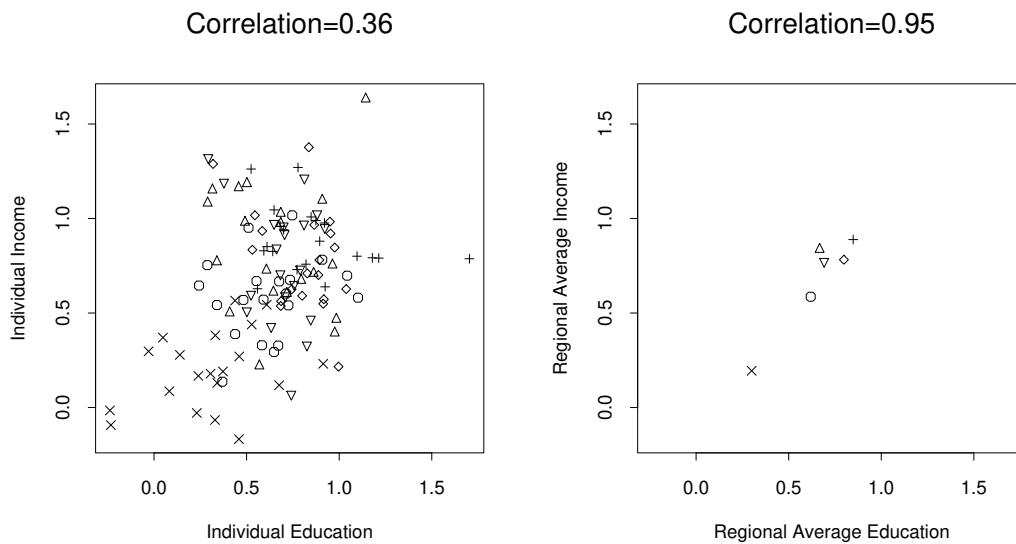
Measures of Association: Correlations

- Commonly used measures:
 - [Pearson](#) product moment correlation: Population value ρ , sample statistic r , $-1 \leq \rho \leq 1$. R, use `cor`, `cor.test`.
- Objective: measure *direction* and *strength* of the *association* between two variables.
 - Positive values ($\rho > 0$) indicate positive association — when one variable increases, so does the other.
 - Negative values ($\rho < 0$) indicate negative association — when one variable increases, the other decreases.
 - Zero values ($\rho \approx 0$) indicate no *linear* association.

Measures of Association: Correlations



Ecological Correlations



- **Ecological correlations** are based on rates or averages and tend to overstate the strength of an association. Beware whenever rates or averages are correlated!
- For example, beware of a correlation involving average store sales

Notation

- Suppose x has mean μ_x and standard deviation σ_x . The **standard units** (also called **Z-scores**) tell how many standard deviations each observation is from the mean:

$$Z_x = \frac{x - \mu_x}{\sigma_x}$$

Use **scale** to compute Z-scores and **cor** to compute correlations in R.

- Let x_1, \dots, x_n be observations from a distribution with mean μ_x and standard deviation σ_x . The *sample mean* and *sample standard deviation* are

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{and} \quad S_x = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

- Let y_1, \dots, y_n be observations from a distribution with mean μ_y and standard deviation σ_y . Likewise let \bar{y} be the sample mean and S_y be the sample standard deviation.

Interpretation and Computation of Correlations

- **Pearson correlation**

$$\begin{aligned}\rho &= \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y} = \frac{\sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)}{\sqrt{\sum_{i=1}^N (x_i - \mu_x)^2 \sum_{i=1}^N (y_i - \mu_y)^2}} \\&= \sum_{i=1}^N \left[\frac{(x_i - \mu_x)}{\sqrt{\sum (x_i - \mu_x)^2}} \times \frac{(y_i - \mu_y)}{\sqrt{\sum (y_i - \mu_y)^2}} \right] \\&= \frac{1}{N} \sum_{i=1}^N \left[\frac{(x_i - \mu_x)}{\sigma_x} \times \frac{(y_i - \mu_y)}{\sigma_y} \right]\end{aligned}$$

Thus, the Pearson correlation coefficient is the *average of the products of the standardized versions of x and y* .

- We estimate ρ with r

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

- Note that Pearson correlations are unaffected by
 - interchanging the two variables
 - replacing a variable x with $ax + b$, where $a > 0$
- There are [formulas](#) for hypothesis tests and CIs, e.g., under $H_0 : \rho = 0$, $t = r\sqrt{(n-2)/(1-r^2)}$ is distributed T_{n-2} .

Geometric Interpretation of Correlation

- Let $\tilde{x}_i = x_i - \mu_x$ and $\tilde{y}_i = y_i - \mu_y$ be *mean-centered* versions of x and y
- Consider $\tilde{x} = (\tilde{x}_1, \dots, \tilde{x}_N)^\top$ and $\tilde{y} = (\tilde{y}_1, \dots, \tilde{y}_N)^\top$ as vectors in N -space. Recall from analytic geometry that the angle θ between the two vectors can be computed using

$$\cos \theta = \frac{\tilde{x}^\top \tilde{y}}{\|\tilde{x}\| \|\tilde{y}\|} = \frac{\sum \tilde{x}_i \tilde{y}_i}{\sqrt{\sum \tilde{x}_i^2} \sqrt{\sum \tilde{y}_i^2}} = \rho$$

The Pearson correlation is thus the *cosine of the angle between the (mean-centered) vectors in N space*.

- Recall that $\cos \theta = 0$ if and only if \tilde{x} is perpendicular (*orthogonal*) to \tilde{y}
- The term *orthogonal* is sometimes used instead of *uncorrelated*
- The coefficient of determination also follows

$$R^2 = 1 - \frac{\text{SSE}}{\text{SST}} = \frac{\text{SSR}}{\text{SST}} = \left(\frac{\text{adj}}{\text{hyp}} \right)^2 = \cos^2 \theta = \rho^2$$

Example: Maintenance Costs

The age in years and maintenance cost in thousands of dollars for 5 machines are provided in the table below.

	Age	Cost	Standard Units		Product
	x	y	Z_x	Z_y	$Z_x \cdot Z_y$
	2	6	-1.247	-1.023	1.275
	5	13	-0.147	-0.105	0.015
	9	23	1.320	1.206	1.592
	3	5	-0.880	-1.154	1.015
	8	22	0.953	1.075	1.025
Mean	5.4	13.8	0	0	0.985
σ	2.728	7.626	1	1	

$$r = \frac{102.40}{\sqrt{37.2 \times 290.8}} = 0.9845$$

Note: I use σ rather than S

```
> machine <- data.frame(age=c(2,5,9,3,8), cost=c(6,13,23,5,22))
> plot(machine)
> cor(machine)      # full correlation matrix
      age      cost
age  1.0000000 0.9845353
cost 0.9845353 1.0000000

> cor.test(machine$age, machine$cost)  # or use cor.test(~age+cost, machine)

Pearson's product-moment correlation

data:  machine$age and machine$cost
t = 9.734, df = 3, p-value = 0.002303
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval: 0.7784347 0.9990256
sample estimates: 0.9845353
```

Bivariate Normal Distribution

- Consider $p = 2$ random variables X_1 and X_2 .
- Let $E(X_j) = \mu_j$, $V(X_j) = \sigma_j^2$, and the correlation be

$$\rho = \frac{\text{Cov}(X_1, X_2)}{\sigma_1 \sigma_2} = \frac{\sigma_{12}}{\sigma_1 \sigma_2},$$

where $\sigma_{12} = \rho \sigma_1 \sigma_2$ is the *covariance* of X_1 and X_2 .

- The *mean vector* of random vector $\mathbf{X} = (X_1, X_2)^\top$ is

$$\boldsymbol{\mu} = E(\mathbf{X}) = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

- The *covariance matrix* is

$$\boldsymbol{\Sigma} = E[(\mathbf{X} - \boldsymbol{\mu})^\top (\mathbf{X} - \boldsymbol{\mu})] = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \rho \sigma_1 \sigma_2 \\ \rho \sigma_1 \sigma_2 & \sigma_2^2 \end{pmatrix}.$$

- The *bivariate normal PDF* is $(-1 < \rho < 1, \sigma_i > 0)$

$$f(\mathbf{x}) = \frac{1}{2\pi \sigma_1 \sigma_2 \sqrt{1 - \rho^2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right],$$

where $\boldsymbol{\Sigma}^{-1}$ is the inverse of $\boldsymbol{\Sigma}$

$$\boldsymbol{\Sigma}^{-1} = \frac{1}{\sigma_1^2 \sigma_2^2 (1 - \rho^2)} \begin{pmatrix} \sigma_2^2 & -\rho \sigma_1 \sigma_2 \\ -\rho \sigma_1 \sigma_2 & \sigma_1^2 \end{pmatrix}.$$

Justification of Elliptical Contours

- Assume uncorrelated variables with mean $(0, 0)^\top$, i.e.,

$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix} \quad \text{and} \quad \Sigma^{-1} = \begin{pmatrix} 1/\sigma_1^2 & 0 \\ 0 & 1/\sigma_2^2 \end{pmatrix}.$$

- A *contour* is the set of points such that $f(\mathbf{x}) = c$ for some fixed c , i.e., $\{\mathbf{x} : f(\mathbf{x}) = c\}$.
- We need to solve

$$c = f(\mathbf{x}) = \frac{1}{2\pi\sigma_1\sigma_2} \exp\left(-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1} \mathbf{x}\right)$$

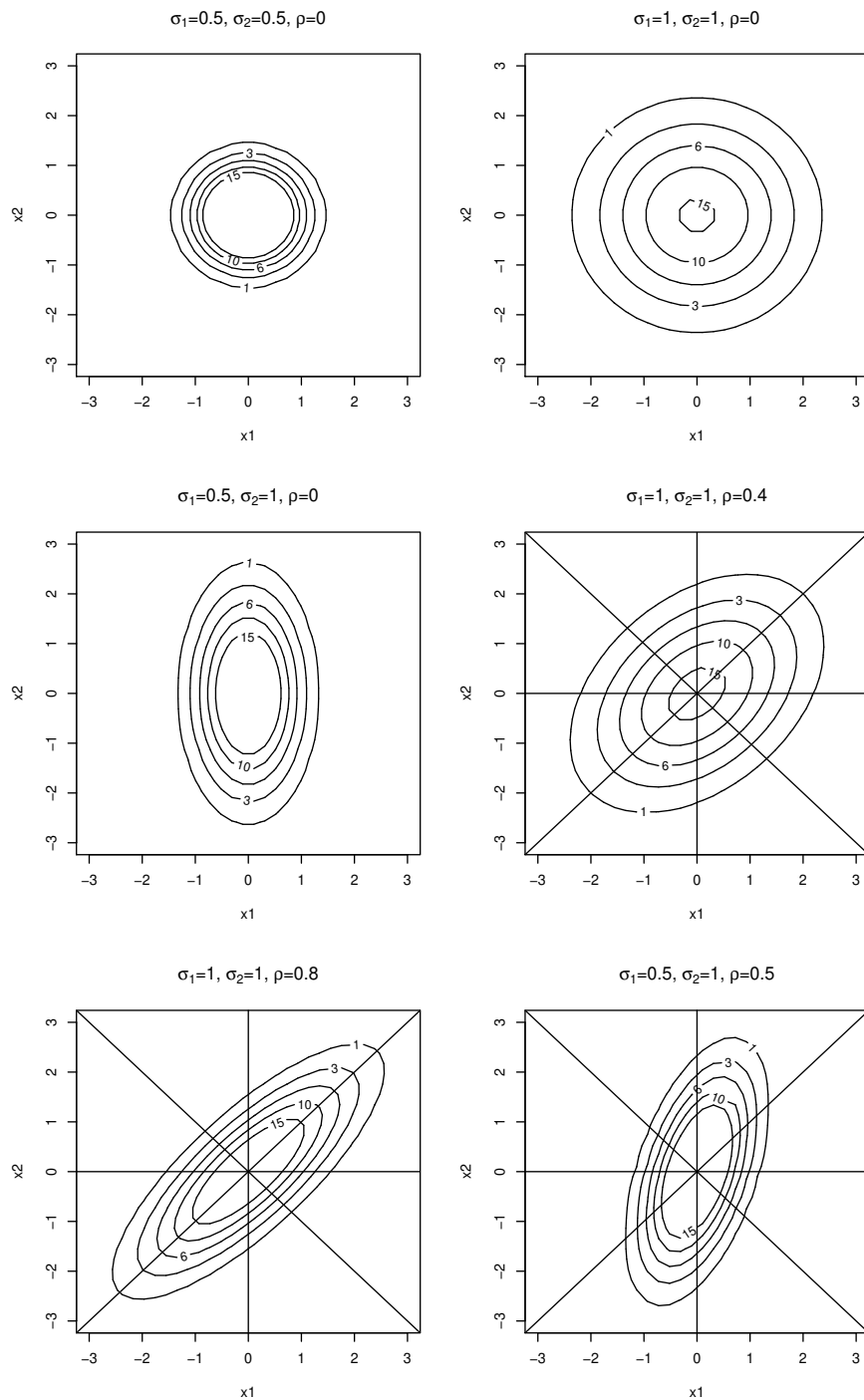
$$2\pi\sigma_1\sigma_2 c = \exp\left(-\frac{1}{2}\mathbf{x}^\top \Sigma^{-1} \mathbf{x}\right)$$

$$-2 \log(2\pi\sigma_1\sigma_2 c) = \mathbf{x}^\top \Sigma^{-1} \mathbf{x} = \frac{x_1^2}{\sigma_1^2} + \frac{x_2^2}{\sigma_2^2}$$

Which we recognize as an ellipse (for sufficiently small c).

- For correlated variables we can show that the ellipse is rotated using the eigenvectors of Σ .

Various Bivariate Normal PDFs



Correlation and regression

- Let X and Y have a bivariate normal distribution with means μ_x and μ_y , and variances σ_x^2 and σ_y^2 , respectively. Their covariance is σ_{xy} and their correlation is $\rho = \sigma_{xy}/(\sigma_x\sigma_y)$.
- Let

$$\beta_1 = \frac{\sigma_y}{\sigma_x}\rho, \quad \beta_0 = \mu_y - \beta_1\mu_x, \quad \text{and} \quad \sigma^2 = \sigma_y^2(1 - \rho^2)$$

Then the following hold:

1. The marginal distributions of X and Y are $\mathcal{N}(\mu_x, \sigma_x^2)$ and $\mathcal{N}(\mu_y, \sigma_y^2)$, respectively
2. The conditional distribution of Y given $X = x$ is

$$Y|x \sim \mathcal{N}(\beta_0 + \beta_1x, \sigma^2)$$

Multiple Linear Regression

Multiple linear regression allows us to study the relationship between a response variable and p predictor variables.

- x_{ij} : value of the j^{th} **predictor variable** on observation i ($i = 1, \dots, n$ and $j = 1, \dots, p$). $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^{\top}$ is an $n \times p$ matrix with row $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ip})^{\top}$ (the extra 1 is for the intercept)
- y_i : value of **dependent variable**, $\mathbf{y} = (y_1, \dots, y_n)^{\top}$.
- β_0 : the **intercept**. When $x = 0$, on average $Y = \beta_0$.
- β_j : **slope coefficient for variable j** . A unit increase in x_j is associated with an increase in y of β_j , *on average* after controlling for other predictors.
- Multiple regression model:

$$y_i = \beta_0 + x_{i1}\beta_1 + \dots + x_{ip}\beta_p + \epsilon_i = \mathbf{x}_i^{\top} \boldsymbol{\beta} + \epsilon_i$$

where $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^{\top}$ and ϵ_i are iid, normal with mean 0 and standard deviation σ_{ϵ} . This implies

$$\mathbb{E}(Y|x_i) = \beta_0 + x_{i1}\beta_1 + \dots + x_{ip}\beta_p$$

- As before, σ_{ϵ}^2 is called the **error variance** and σ_{ϵ} : the **standard deviation of the errors**.
- We write this compactly as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \text{MVN}(\mathbf{0}, \sigma^2 \mathbf{I})$$

Estimating the Regression Model

- We don't know values of **parameters** $\beta_0, \beta_1, \dots, \beta_p$, and σ_ϵ .
- **Estimates** denoted by b_0, b_1, \dots, b_p , and S_ϵ
- $\mathbb{E}(Y_i|x_i)$ called **fitted, predicted**, or “**y-hat**” values:

$$\hat{y}_i = b_0 + b_1x_{i1} + \dots + b_px_{ip}$$

- The **residual** for observation i is

$$\hat{\epsilon}_i = y_i - \hat{y}_i = y_i - \left(b_0 + \sum_{j=1}^p b_jx_{ij} \right)$$

- We choose **b** to minimize the **least-squares criterion** (SSE means **sum of squared errors**):

$$\text{SSE} = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta})$$

- The **ordinary least squares** (OLS) **estimates** of β :

$$\frac{\partial \text{SSE}}{\partial \hat{\beta}} = -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) \implies \hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

- Predicted values are given by

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

where $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is called the **hat matrix**

Some Properties of OLS Estimates

- $\hat{\beta}$ is unbiased ([proof](#))

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}(\mathbf{X}\boldsymbol{\beta} + \mathbf{e}) = \boldsymbol{\beta}$$

- $\mathbb{V}(\hat{\beta}) = \sigma_\epsilon^2 (\mathbf{X}^\top \mathbf{X})^{-1}$

$$\mathbb{V}[(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}] = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{V}(\mathbf{y}) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} = \sigma_\epsilon^2 (\mathbf{X}^\top \mathbf{X})^{-1}$$

- Gauss-Markov Theorem: $\hat{\beta}$ has a covariance matrix that is “smaller” than that of any other linear estimator and is called the *best linear unbiased estimator* (“BLUE”).

- Estimate σ_ϵ^2 with the **mean squared error**

$$S_\epsilon^2 = \text{MSE} = \frac{\text{SSE}}{n - p - 1}$$

and σ_ϵ with the **residual standard error** (a.k.a. **root mean squared error**) $S_\epsilon = \sqrt{S_\epsilon^2}$.

- If $\boldsymbol{\epsilon} \sim N(0, \sigma^2 \mathbf{I})$ and the linear model is correct, then $\hat{\beta}$ has a multivariate normal distribution because it is a linear transformation of normally distributed $\boldsymbol{\epsilon}$:

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) = \boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\epsilon}$$

- $S_{b_j} = S_\epsilon \sqrt{v_j}$ is called the *standard error* of b_j , where v_j is the j^{th} diagonal element of $(\mathbf{X}^\top \mathbf{X})^{-1}$, and $(b_j - \beta_j)/S_{b_j}$ has a t distribution.

- A $(1 - \alpha)\%$ confidence region for $\boldsymbol{\beta}$ is

$$(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \leq \hat{\sigma}^2 \chi_{p+1}^{2(1-\alpha)} \quad (3.15)$$

Geometry of Least Squares

The geometrical interpretation is often useful

- $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ is the orthogonal projection of \mathbf{y} onto the subspace $\mathcal{M} \subset \mathbf{R}^n$ spanned by the columns of \mathbf{X} . This is true even if \mathbf{X} is not of full column rank.
- The columns of \mathbf{X} are basis vectors for \mathcal{M} and $\hat{\boldsymbol{\beta}}$ gives the coordinates of $\hat{\mathbf{y}}$ with respect to this basis
- Hat matrix \mathbf{H} projects n -vectors onto \mathcal{M}
- $\mathbf{P} = \mathbf{I} - \mathbf{H}$ projects onto the orthogonal complement of \mathcal{M}
- Residual vector $\mathbf{Py} = \mathbf{y} - \hat{\mathbf{y}}$ is orthogonal to \mathcal{M}
- The three vectors \mathbf{y} , $\hat{\mathbf{y}}$, and $\mathbf{y} - \hat{\mathbf{y}}$ form a right triangle in n -space, where \mathbf{y} is the hypotenuse
- If \mathbf{y} and columns of \mathbf{X} are mean centered
 - **total sum of squares:** $\mathbf{y}^\top \mathbf{y} = (n - 1)S_y^2$, the squared length of the hypotenuse (S_y^2 is the variance of y)
 - **regression sum of squares:** $\hat{\mathbf{y}}^\top \hat{\mathbf{y}} = \mathbf{y}^\top \mathbf{Hy}$
 - $\text{SSE} = (\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y}^\top \mathbf{Py}$
 - Pythagoras gives us the ANOVA equality

$$\mathbf{y}^\top \mathbf{y} = \mathbf{y}^\top \mathbf{Hy} + \mathbf{y}^\top \mathbf{Py}$$

Basic Model Building Process

Suppose that (1) your interest is confirmatory and (2) the model is correct (i.e., y is really a linear function of the specified x variables plus additive, normal, independent, homoscedastic errors)

1. Inspect your data for outliers, typos, missing values, etc.
 - Generate n , means, mins, and maxs of each variable
 - Generate boxplots or histograms of each variable
 - Generate a scatterplot matrix for small data sets
 - Generate correlation matrix to assess correlations between predictor variables and pairwise correlations with DV
2. Estimate model and check residual and normality plots, VIFs
3. Test **overall significance of the model**
$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_p = 0$$
$$H_1 : \text{at least one } \beta_j \neq 0$$
4. If you can reject H_0 in Step 3, interpret model and test significance of individual coefficients. If you cannot reject H_0 in Step 3, don't try to test individual coefficients.
 - In practice you usually will not know the correct model, which complicates the process substantially.
 - If only goal is prediction, hypothesis tests not used

Estimates from Click Ball Point Pens

```
> fit = lm(sales ~ ad + reps + eff, click)
> summary(fit)
Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   31.150      34.175   0.911    0.368
ad             12.968       2.737   4.738 3.34e-05 ***
reps           41.246       7.280   5.666 1.95e-06 ***
eff            11.524       7.691   1.498    0.143
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 44.42 on 36 degrees of freedom
Multiple R-squared:  0.8812, Adjusted R-squared:  0.8714
F-statistic: 89.05 on 3 and 36 DF,  p-value: < 2.2e-16
```

- Is the regression significant? *Solution:*

$H_0 : \beta_1 = \beta_2 = \beta_3 = 0$ versus H_1 : at least one $\beta_j \neq 0$.

$P = 2.2e - 16 < .05$, so reject H_0 and conclude that at least one of the predictors is predictive.

- State estimated regression equation. *Solution:*

$$\hat{y} = 31.15 + 12.97\text{ad} + 41.25\text{reps} + 11.52\text{eff}$$

- Interpret the coefficient for **reps** (41.25).
- Here $\hat{\beta}_1 = 13.0$ but on page 12, $\hat{\beta}_1 = 25.3$. **Which is right?**

Estimates from Click (Continued)

- Construct at 95% CI for **reps**.

Solution: $41.25 \pm 2.028 \times 7.28 = [26.5, 56.0]$

```
> confint(fit)
              2.5 %    97.5 %
(Intercept) -38.159815 100.46059
ad           7.416798  18.51953
reps        26.480882  56.01037
eff        -4.074175  27.12268
```

- Is **eff** different from zero (use .05 level)?

$H_0 : \beta_3 = 0$ versus $H_1 : \beta_3 \neq 0$.

$$\begin{aligned} P(b_3 > 11.52) &= P\left(T_{36} > \frac{11.52 - 0}{7.6912}\right) \\ &= P(T_{36} > 1.498) = 0.0714 \end{aligned}$$

The P value is thus $2*(1-pt(1.498, 36)) = 0.1429$. We cannot reject H_0 because $0.1428 > 0.05$ and we cannot conclude that wholesaler efficiency affects sales.

- Predict sales for **ad**=4, **reps**=3, **eff**=1. *Solution:*

$$\hat{y} = 31.15 + 12.97 \times 4 + 41.25 \times 3 + 11.52 \times 1 = 218.3$$

```
> predict(fit, data.frame(ad=4, reps=3, eff=1))
1
218.2842
```

Comparing Regression Coefficients

- *Question:* Which of the variables is more “important” in explaining sales?
- *Answer:* The coefficients are not directly comparable because of differences in units of measurement.
- *Ideal solution:* convert to commensurate units, e.g., dollars.
- *Possible solution:* **Standardized regression coefficients** (all variables standardized $z = (x - \bar{x})/s_x$ before the analysis to have mean 0 and variance 1). The “unit” of measurement is now the standard deviation (units cross out)

```
> Zclick = as.data.frame(scale(click[,1:4]))
> fit = lm(sales ~ ad + reps + eff -1, Zclick) # -1 drops intercept
> summary(fit)
Coefficients:
      Estimate Std. Error t value Pr(>|t|)
ad      0.45101    0.09390   4.803 2.59e-05 ***
reps    0.54902    0.09559   5.744 1.40e-06 ***
eff     0.09157    0.06028   1.519  0.137
```

- *Possible solution:* compare t scores ($t = b_j/S_{b_j}$, reported in the fourth column of output, units also cross out)
- See [Bring \(1994\)](#) for discussion

Standardized Regression Coefficients

- Theorem: the standardized regression coefficient is $b_j S_{x_j} / S_y$, where S_{x_j} and S_y are the standard deviations of x_j and y , respectively, and b_j is the (unstandardized) regression estimate for variable j .
- Standardized regression coefficients ...
 - Also called “*beta*” coefficients
 - Interpreted as a standard deviation increase in x_j is associated with “beta” standard deviations in y
 - Equals correlation r between x and y when $p = 1$ predictor
- Do not use “beta” coefficients blindly:
 - If x_j is a 0-1 variable, then the standard deviation is $\sqrt{\bar{x}_j(1 - \bar{x}_j)}$
 - If you are analyzing a designed experiment, you select the x_j values and thus the standard deviations
 - “Beta” values are still a function of other variables in the model (multicollinearity)
 - “Beta” values do not consider costs

Your Turn

1. A commercial real estate company evaluates vacancy rates, square footage, rental rates, and operating expenses for commercial properties in a large metropolitan area in order to provide clients with quantitative information upon which to make rental decisions. The data in `commercial.txt` are taken from 81 suburban commercial properties that are the newest, best located, most attractive, and expensive for five specific geographic areas. It includes their age (`x1`), operating expenses and taxes (`x2`), vacancy rates (`x3`), total square footage (`x4`), and rental rates (`y`).¹
 - (a) Read the commercial data set into R and run basic descriptive statistics (counts, mins, maxs, means). Do the descriptives make sense? Hint:

```
> comm = read.table("commercial.txt", header=T)
> summary(comm)
```
 - (b) Produce a scatterplot matrix and correlation matrix. Discuss the relationships between the variables. Hint: use `plot(comm)` and `cor(comm)`.
 - (c) Regress rental rates on the four predictor variables. State the estimated regression equation. Hint:

```
> fit = lm(y ~ x1 + x2 + x3 + x4, comm)
> summary(fit)
```
 - (d) Test whether the overall model is significant. State the null and alternative, P -value and decision. Hint: `summary(fit)`
 - (e) What fraction of variation in rental rates is explained by these predictor variables?
 - (f) If the overall regression model is significant in the previous part then test whether each of the individual regression coefficients equals 0 at the 5% level against a 2-sided alternative, i.e., $H_0 : \beta_j = 0$ for $j = 1, 2, 3, 4$. For each predictor, state the P -value and your decision.
 - (g) Assume that the regression model you have estimated is appropriate. Three properties with the following characteristics did not have any rental information available.

¹Solution: (a) The summary statistics make sense. Age ranges from 0 to 20, which is reasonable for real estate properties. Operating expenses and taxes are positive. The vacancy rate is between 0 and 1. Square footage looks reasonable, as do rental rates. (b) The first thing to note is the correlations with y . The older the property, the lower the rent. The higher the expenses, the higher the rent. Vacancy rate has a positive correlation, but it is very weak. Square footage has a positive correlation with rent. There are also correlations among the predictor variables, especially between age and expenses (positive), size and expenses (positive), and expenses and vacancy rate (negative). (c) $\hat{y} = 12.2 - 0.142 \text{ x1} + 0.282 \text{ x2} + 0.619 \text{ x3} + 0.00000792 \text{ x4}$. (d) $H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$ versus H_1 : at least one $\beta_j \neq 0$. $P = 7.27 \times 10^{-14} < .05$ so reject H_0 . At least one predictor is related to rental rate. (e) $R^2 = .5847$. (f) $H_0 : \beta_j = 0$ versus $H_1 : \beta_j \neq 0$. `x1`, `x2`, and `x4` all have $P < .05$. `x3` has $P = .57$ and we cannot reject the null or conclude that vacancy rate has an effect on the rental rate. (g) Property 1: 15.15, 12.85, 17.44; Property 2: 15.54, 13.25, 17.84; Property 3: 16.91, 14.53, 19.29.

	Property 1	Property 2	Property 3
x1	4	6	12
x2	10	11.5	12.5
x3	0.1	0	0.32
x4	80,000	120,000	340,000

Predict the rental rate and compute separate prediction intervals for the rental rates using 95% confidence. **Briefly tell what the prediction interval tells you.** Hint:

```
> newx = data.frame(x1=c(4,6,12), x2=c(10,11.5,12.5), x3=c(.1,0,.32),
  x4=10000*c(8,12,34))
> predict(fit, newx, interval="prediction")
```

2. In a small-scale experimental study of the relation between degree of brand liking (**y**) and the moisture content (**moisture**) and sweetness (**sweetness**) of the product, the results in the data below were obtained from the experiment based on a completely randomized design.²

```
brand = data.frame(
  liking=c(64,73,61,76,72,80,71,83,83,89,86,93,88,95,94,100),
  moisture=c(4,4,4,4,6,6,6,6,8,8,8,8,10,10,10,10),
  sweetness=c(2,4,2,4,2,4,2,4,2,4,2,4,2,4,2,4) )
```

- (a) Read the brand data set into R and run basic descriptive statistics (counts, mins, maxs, means). Do the descriptives make sense?
- (b) Produce a scatterplot matrix and correlation matrix. Discuss the relationships between the variables.
- (c) Regress liking on the two predictor variables. State the estimated regression equation.
- (d) What fraction of variation in liking is explained by these predictor variables?
- (e) Test whether the overall model is significant. State the null and alternative, P -value and decision.
- (f) If the overall regression model is significant in the previous part then test whether each of the individual regression coefficients equals 0 at the 5% level against a 2-sided alternative, i.e., $H_0 : \beta_j = 0$ for $j = 1, 2$. For each predictor, state the P -value and your decision.
- (g) Assume that the regression model you have estimated is appropriate. Predict the liking when **moisture**=5 and **sweetness**=4. Find a prediction interval and separately a confidence interval for the estimated mean value using 99% confidence.

²(a) They make sense. All variables are positive. (b) Moisture has a stronger positive correlation with liking than sweetness. There is no correlation between sweetness and moisture because the data are from an experiment with an orthogonal design. (c) $\hat{y} = 37.650 + 4.425\text{moisture} + 4.375\text{sweetness}$. (d) $R^2 = .9521$. (e) $H_0 : \beta_1 = \beta_2 = 0$ versus H_1 : at least one $\beta_j \neq 0$. $P = 2.658 \times 10^{-9} < .05$ so reject H_0 . At least one predictor is related to liking. (f) $H_0 : \beta_j = 0$ versus $H_1 : \beta_j \neq 0$. **moisture** and **sweetness** have $P < .05$. (g) Use `predict(fit, data.frame(moisture=5, sweetness=4), interval="conf", level=.99)` to find 77.275, 73.88, 80.67.

3. JWHT problem 9(a)–(c) on page 122. Find the correlation and scatterplot matrices and regress mpg on all other variables except for name. Hint: when finding correlations see the `use="pair"` option. Answer these questions.³

- (a) Based on the scatterplots, comment on the relationships between the predictors and mpg.
- (b) What is the correlation between mpg and displacement and what does it tell you?
- (c) Is there a statistically significant relationship between the predictors and the response?
- (d) Which predictors appear to have a statistically significant relationship to the response?
- (e) What does the slope coefficient for the year variable suggest?
- (f) What does the slope coefficient for the displacement variable suggest?

³Solution: (a) There are nonlinear relationships between mpg and displacement, horsepower, weight and acceleration. There may be other nonlinear relationships. (b) $r = -.7763$, so larger displacement is associated with smaller mpg. (c) Yes, $P < 2.2 \times 10^{-16}$. (d) Displacement, weight, year and origin. (e) $b = .75$ suggests that high gas mileage improves over time. (f) $b = 0.0199$ suggests that larger displacement is associated with higher gas mileage. This contradicts part b because of multicollinearity. It is an example of a sign flip.

Assumptions, Diagnostics, Remedies

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \epsilon_i, \quad \mathbb{E}(\epsilon_i) = 0, \mathbb{V}(\epsilon_i) = \sigma^2, \epsilon_i \text{ normal, uncorrelated}$$

Assumption	Diagnostic	Problem	Remedies
Linear function	Residual plot	Biased \hat{y}	Transform x
$\mathbb{V}(\epsilon_i)$ constant	Residual plot	OLS not BLUE	Transform y or GLM
Correlated errors	Autocorrelations Plot $\hat{\epsilon}_t$ vs. t Runs and DW tests	Estimates unbiased Variances wrong	Transform: Cochrane-Orcutt
Errors normal	QQ plot	t/F assume it	Transform y or GLM
Outliers [†]	Leverage plot	Unstable estimates	
Collinearity [†]	Correlations/VIFs	Inflated variances	Discuss next week

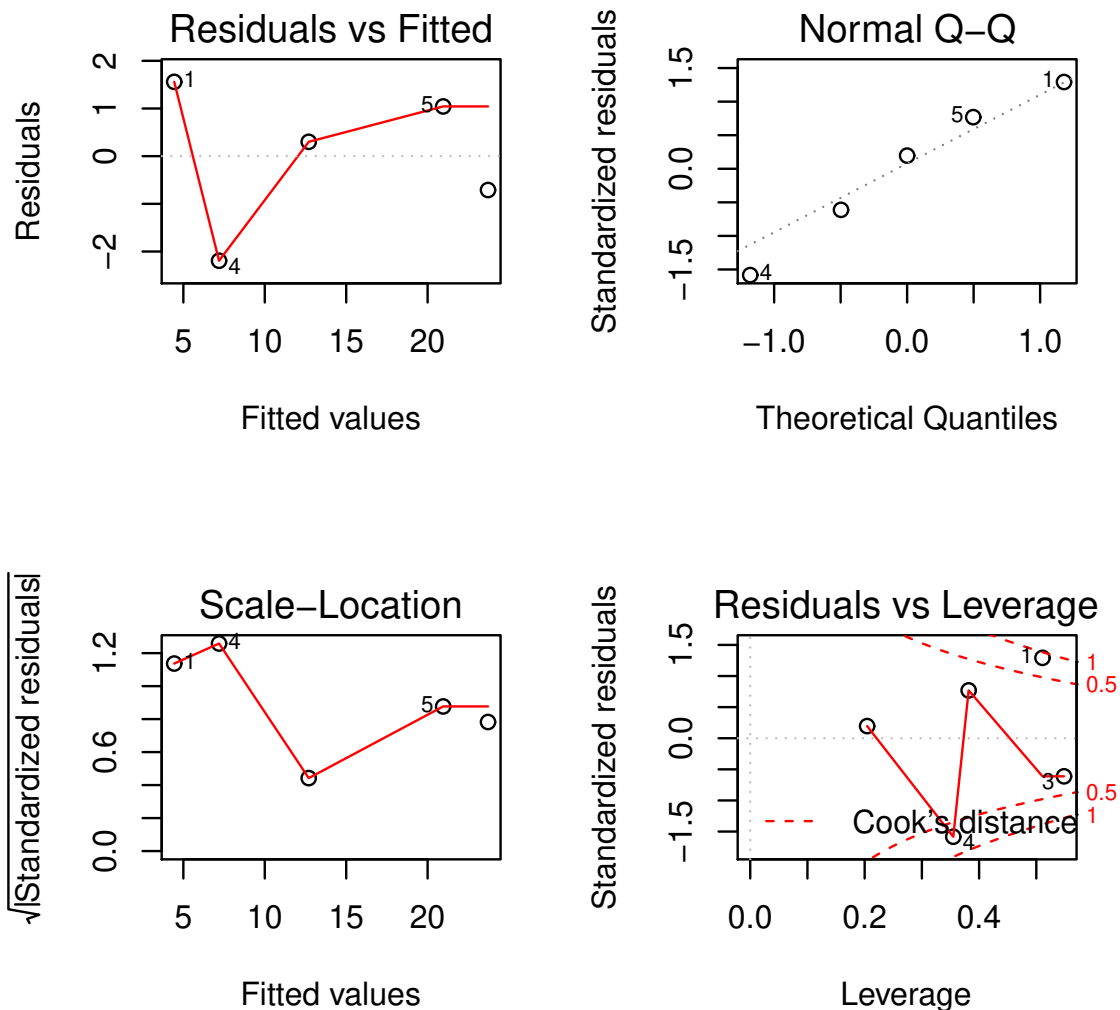
[†] = Not an assumption but it can cause problems

- To avoid omitted variables, start with conceptual framework
- Fixing functional form is usually my next priority
- Constant $\mathbb{V}(\epsilon_i)$: having efficient (BLUE) estimates is not critical when sample sizes are large, but I still worry about this
- Whenever observations come from different time periods (time series) you should suspect correlated errors
- Normal errors: with even modest sample sizes this is a lower priority because of the [central limit theorem](#)
- Outliers: Determine why the value is an outlier
 - Erroneous value: fix it or drop it
 - Correct but extreme value: many considerations. Transform amounts/counts (log or root) will reduce influence, or there are also [robust](#) versions of regression

Residual and QQ Plots in R

- `plot(fit)` gives diagnostic plots of `lm` objects.
- Alternatively, `plot(fit, which=1)` gives residuals `plot(fit, which=2)` gives QQ plots, etc. See `?plot.lm`

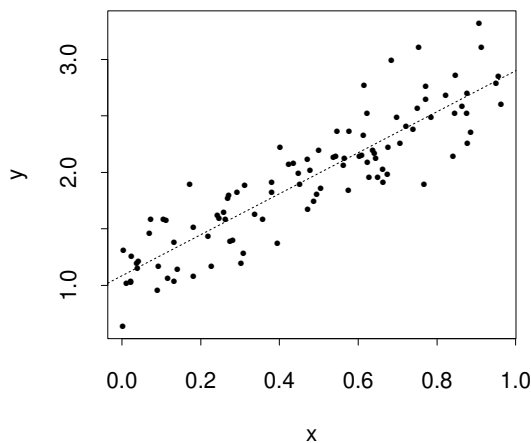
```
> machine = data.frame(age=c(2,5,9,3,8), cost=c(6,13,23,5,22))
> fit = lm(cost ~ age, machine)
> par(mfrow=c(2,2)) # show 2*2 grid of plots. Use c(1,1) for one plot per page
> plot(fit)
```



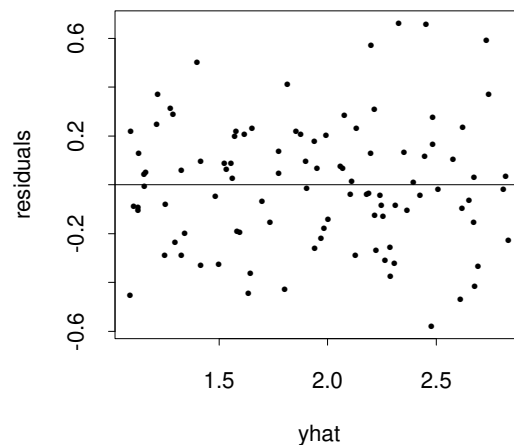
Ideal Regression Model

True model: $y = 1 + 2x + \epsilon$

Raw data



Residual Plot



Least-squares fit: $\hat{y} = 1.09 + 1.81x$

- Definition of residuals

$$\hat{\epsilon} = y - \hat{y}$$

- Residual plots ($\hat{\epsilon}$ against \hat{y}) help us to understand how well the model fits the data. Use `plot(fit, which=1)` in R.
- Here, residuals do not follow any pattern
- Variance of residuals does not depend on \hat{y} (homoscedasticity)
- This is an ideal residual plot

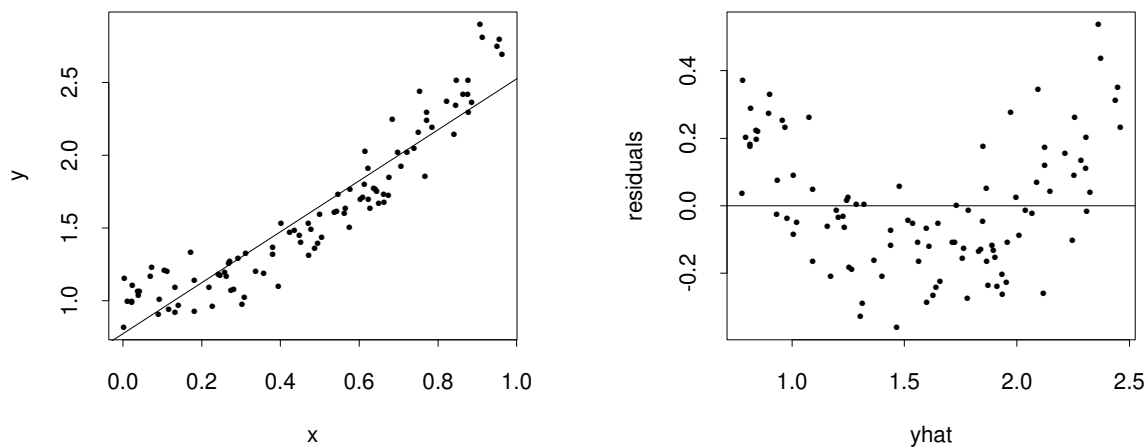
Model misspecification I

- Suppose the true model is

$$y = 1 + 2x^2 + \epsilon$$

- We estimate the *incorrect* (misspecified) model:

$$y = \beta_0 + \beta_1 x + \epsilon$$



- True relationship is nonlinear (curved)
- Fitted line does not describe the data well
- Pattern in residual plot indicates model misspecification
- But, Variance of residuals does *not* depend on \hat{y} (error variance still homoscedastic)

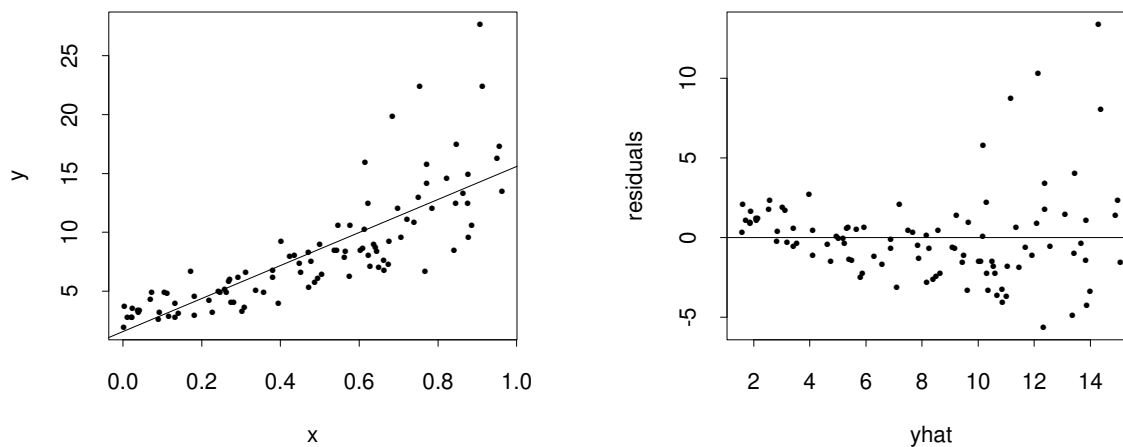
Model misspecification II

- Suppose the true model is

$$y = \exp(1 + 2x + \epsilon)$$

- We estimate the *incorrect* (misspecified) model:

$$y = \beta_0 + \beta_1 x + \epsilon$$



- Residuals show a pattern: first they are mostly positive, then mostly negative, then roughly centered at 0.
- Variance of the residuals increases with \hat{y} indicating heteroscedasticity
- Note that

$$\log(y) = 1 + 2x + \epsilon$$

Review of Key Results

- Multiple linear regression model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \mathbb{E}(\boldsymbol{\epsilon}) = \mathbf{0}, \quad \text{and} \quad \mathbb{V}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I} = \mathbb{V}(\mathbf{y})$$

- We estimate $\boldsymbol{\beta}$ with OLS estimates \mathbf{b} , and $\boldsymbol{\epsilon}$ with $\hat{\boldsymbol{\epsilon}}$

$$\mathbf{b} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad \text{and} \quad S_\epsilon^2 = \sum \hat{\epsilon}_i^2 / (n - p - 1)$$

- Estimates of predicted values and residuals are given by

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{b} = \mathbf{H}\mathbf{y} \quad \text{and} \quad \hat{\boldsymbol{\epsilon}} = \mathbf{y} - \mathbf{X}\mathbf{b}$$

where $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is the **hat matrix**

- Theorem: $\mathbb{V}(\hat{\mathbf{y}}) = \sigma^2 \mathbf{H}$ and $\mathbb{V}(\hat{\boldsymbol{\epsilon}}) = \sigma^2 (\mathbf{I} - \mathbf{H})$
- Definition: the **leverage** of observation i is h_{ii} , where

$$0 \leq h_{ii} \leq 1 \quad \text{and} \quad \sum_{i=1}^n h_{ii} = p$$

As a rule of thumb, leverages greater than twice their average, i.e., $h_{ii} > 2p/n$, are considered large

- For simple linear regression

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{j=1}^n (x_j - \bar{x})^2}$$

Standardized Residuals

- Problem 1: estimated residuals $\hat{\epsilon}$ don't have constant variance (even if model, which assumes homoscedastic errors, is true).
- The **standardized residual**, sometimes called the **studentized residual** is obtained by dividing \hat{e} by its estimated standard deviation

$$r_i = \frac{\hat{\epsilon}_i}{S_\epsilon \sqrt{1 - h_{ii}}}$$

Also sometimes called **internal standardized residuals**.

- Problem 2: if there are outliers, S is inflated, which deflates all r_i . One solution is to omit observation i and reestimate the model giving prediction $\hat{y}_{(i)}$ and MSE $S_{(i)}^2$.
- The **deleted** (or **external**) residual and **studentized deleted/external residual** are

$$d_i = y_i - \hat{y}_{(i)} \quad \text{and} \quad \frac{d_i}{S_{(i)} \sqrt{1 - h_{ii}}}$$

- **Cook's distance**

$$D_i = \frac{\sum_j (\hat{y}_j - \hat{y}_{j(i)})^2}{p S_e^2}$$

As a rule of thumb, **values greater than 1 are considered large**.

Leverage in R: Machine Example

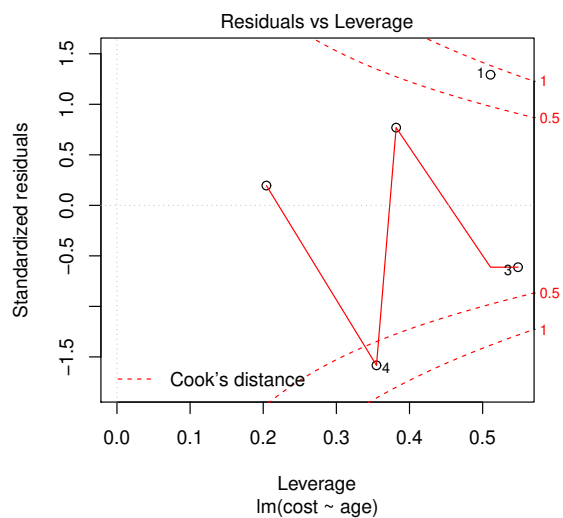
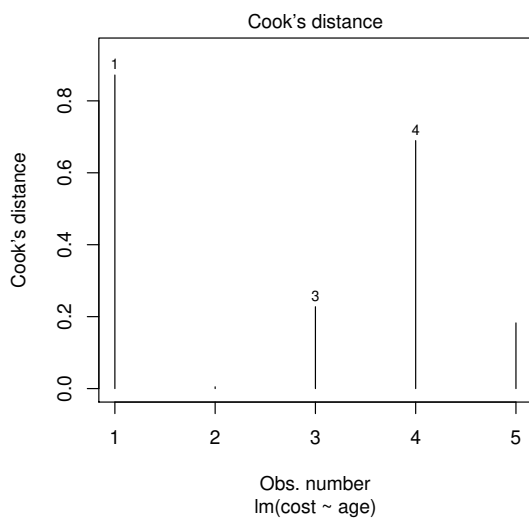
```
> fit = lm(cost~age, machine)
> plot(fit, which=c(4,5))

> # use lm.influence function to get leverages in R
> lm.influence(fit)$hat
      1      2      3      4      5
0.5107527 0.2043011 0.5483871 0.3548387 0.3817204
> sum(lm.influence(fit)$hat)
[1] 2

> # check our work with matrix inversion
> X=cbind(1,machine$age)
> diag(X %*% solve(t(X) %*% X) %*% t(X))
[1] 0.5107527 0.2043011 0.5483871 0.3548387 0.3817204

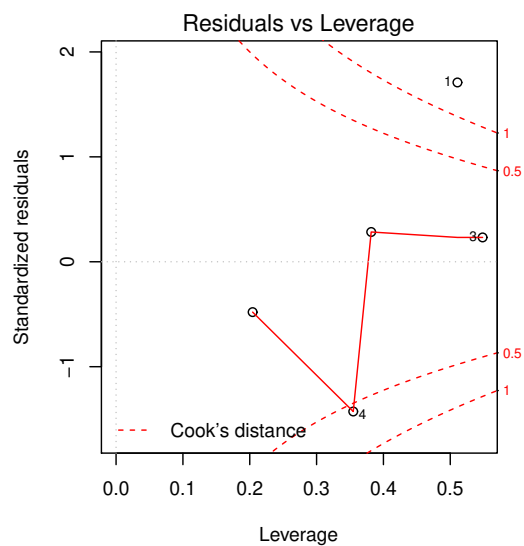
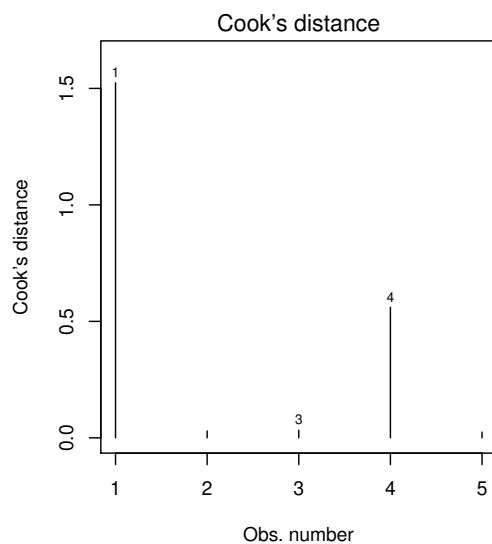
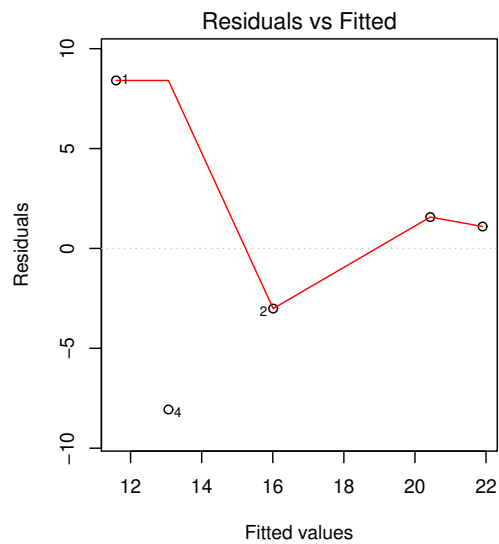
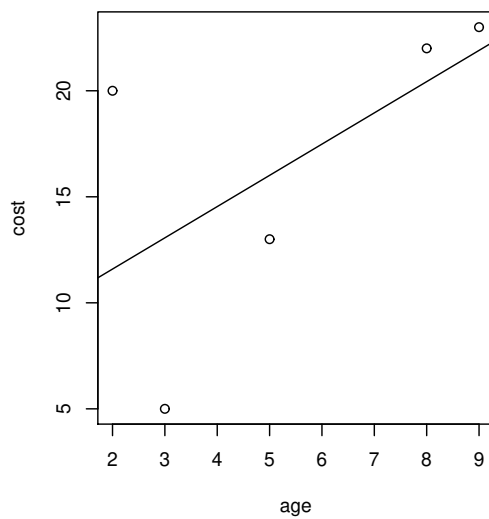
> # check our work with simple formula
> xbar = mean(machine$age)
> 1/5+(machine$age-xbar)^2/sum((machine$age-xbar)^2)
[1] 0.5107527 0.2043011 0.5483871 0.3548387 0.3817204
```

None are considered large since all are less than $2(2/5) = 0.8$



Modified Machine Example

```
> machine2 = data.frame(age=c(2,5,9,3,8), cost=c(20,13,23,5,22)) # change y1=20 from 6
> fit = lm(cost~age, machine2)
> plot(machine2); abline(fit)
> plot(fit)
```



Your Turn

For each of the problems below, examine the residual, Q-Q and Cook's distance plots. Are there problems? What are you looking for and why?

1. Problem 2 on page 16.
2. Problem 3 on page 16.
3. Problem 4 on page 16.
4. Commercial property problem on page 41.
5. Brand problem on page 41.

Answers

- | | |
|--|--|
| 1. No problems. | is not what one usually sees with amount variables (variance usually increases with the mean). There's some non-normality and case 1 has Cook's distance greater than 1. |
| 2. WSJ has a Cook's distance greater than 1 suggesting it is influential. WSJ and USA deviate from normality. | |
| 3. You don't have many cases, but the residual plot has an unusual shape. Perhaps there is heteroscedasticity, but it looks like the residuals are larger for small fitted values, which | 4. The residual plot shows no patterns. |
| | 5. There may be an inverted-U shaped relationship in the residual plot, but the sample is tiny and the fit is not bad. |

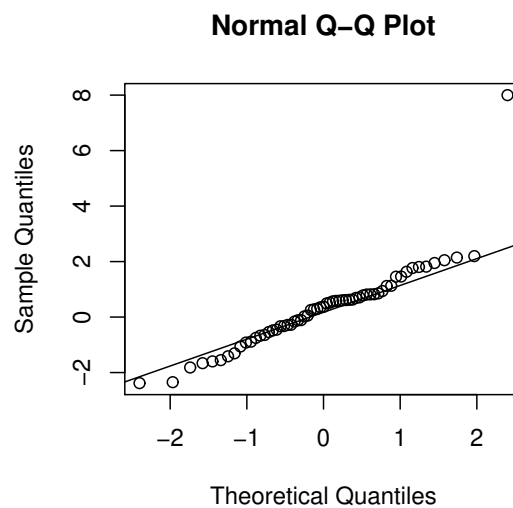
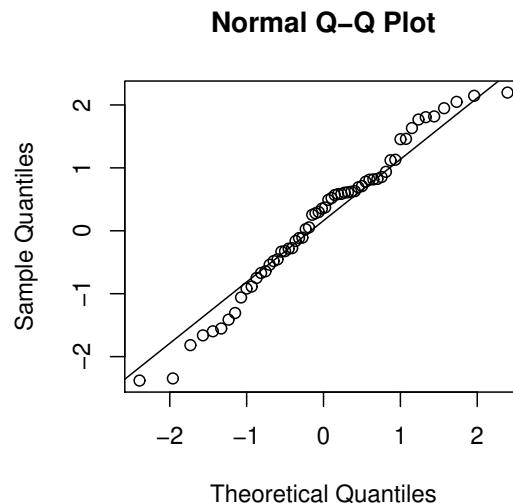
Evaluating Normality

- When the CLT has not “converged” then the population distribution of residuals must be normal for you to use the t distribution.
- Evaluate normality using a *normal probability plot*, which plots the observed quantile against normal quantiles (“Q-Q plot”).
- Points falling on a line indicate normality.

```
> set.seed(12345)
> z = rnorm(60)

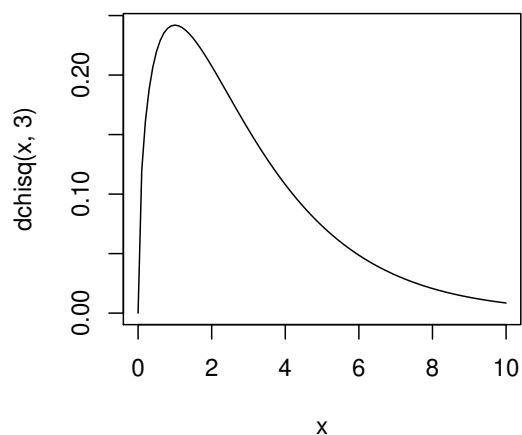
> # data from normal distribution
> qqnorm(z); qqline(z)

> # now we add an outlier
> qqnorm(c(z, 8)); qqline(c(z,8))
```



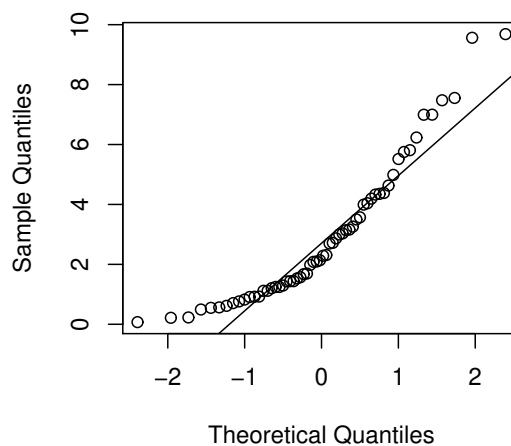
Evaluating Normality

```
> x = seq(0, 10, .1)
> plot(x, dchisq(x, 3), type="l")
```

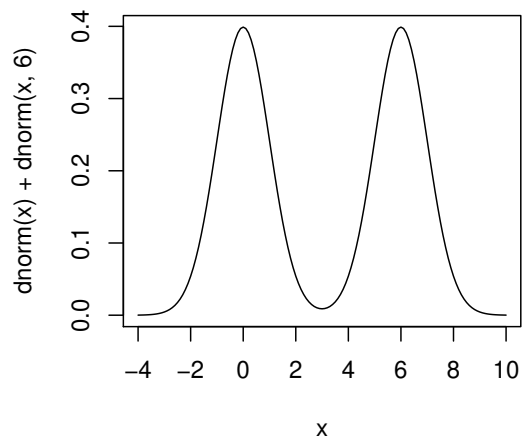


```
> rx = rchisq(60, 3)
> qqnorm(rx); qqline(rx)
```

Normal Q-Q Plot

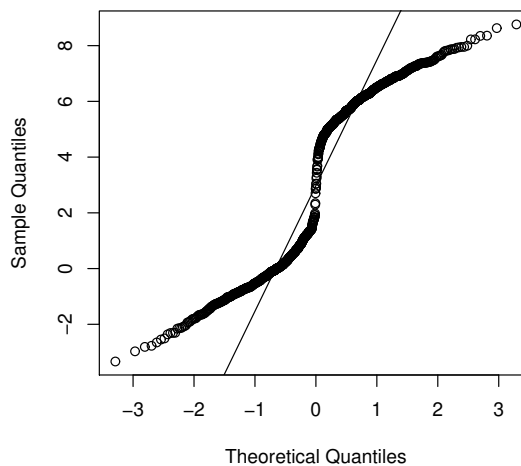


```
> x = seq(-4, 10, .1)
> plot(x, dnorm(x)+dnorm(x,6), type="l")
```



```
> x = c(rnorm(100), rnorm(100, 6))
> qqnorm(x); qqline(x)
```

Normal Q-Q Plot



Transformations

- Transformations change the functional relationship between dependent and predictor variables
- Two reasons for transformations:
 - Heteroscedasticity / non-symmetric error distributions ($\mathbb{E}(\epsilon_i) = 0$ but $\text{Skewness}(\epsilon_i) \neq 0$): transform the *dependent* variable
 - Underlying relationship nonlinear: Transform the *predictor* variable(s)
- Outline of transformation lecture
 1. Transformations of dependent variable
 2. Identifying nonlinear relationships
 - (a) Scatterplots
 - (b) Compare R^2 values for models using various transformations (e.g., Tukey's ladder of re-expressions)
 3. Other transformations based on combinations of variables

Heteroscedasticity

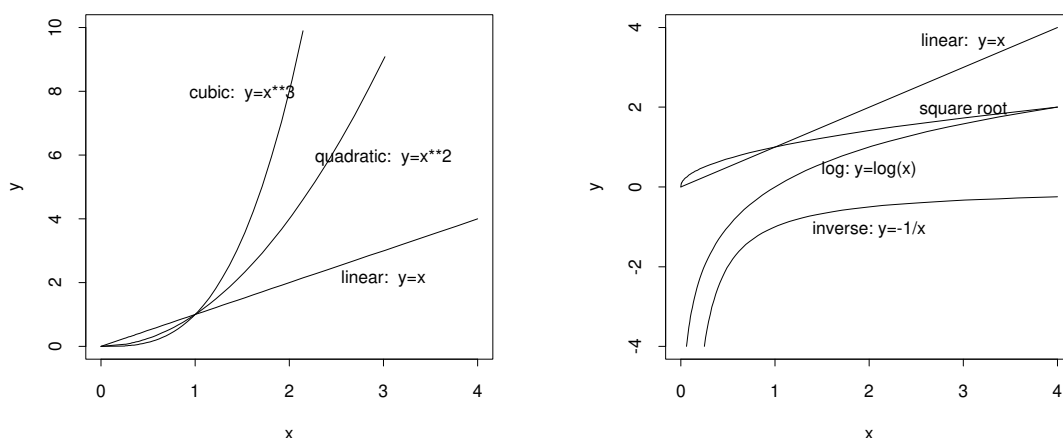
- Slide 32 gave model assumptions, including homoscedasticity
- If a model has heteroscedastic error variance, the least-squares estimates will still be unbiased, but will not be BLUE.
- Modeling heteroscedastic data: Use ...
 - a variance-stabilizing transformation. For count and amount dependent variables, use the logarithm or square root as a variance-stabilizing transformation. Take **logs** when the **standard deviation** of the errors is proportional to the mean, and **square roots** when the **variance** is proportional to the mean.
 - a different model, e.g., Poisson or logistic regression
 - **weighted least squares** (WLS): let $w_i \propto 1/\mathbb{V}(y_i)$ and $\mathbf{W} = \text{diag}(w_1, \dots, w_n)$. Then the WLS estimate is BLUE:

$$\hat{\boldsymbol{\beta}}_{\text{WLS}} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{y}.$$

When \mathbf{W} is unknown (usually), we use *iteratively reweighted least squares* (IRLS)

Tukey's Ladder of Transformations

Consider models of the form $y = \beta x^k$ over $[0, \infty)$



Note: “Returns” refer to the first derivative (slope)

k	Function	Slope	Nature of “Returns”
3	$y = \beta x^3$	$dy/dx = 3\beta x^2$	Increasing
2	$y = \beta x^2$	$dy/dx = 2\beta x^1$	Increasing
1	$y = \beta x^1$	$dy/dx = \beta x^0$	Constant
1/2	$y = \beta \sqrt{x}$	$dy/dx = \beta x^{-1/2}/2$	Decreasing
0	$y = \beta \log x$	$dy/dx = \beta x^{-1}$	Decreasing
-1	$y = \beta x^{-1}$	$dy/dx = -\beta x^{-2}$	Decreasing

- For $k < 1$ the slope approaches 0, but never changes sign
- You may need to shift variables, e.g. $\log(x + 1)$ or $1/(x + 1)$

Tukey’s First Aid Re-Expressions⁴

“Choosing exactly the right re-expression for a particular quantity may not be easy. To try to do a good job, we may have to (1) sense rather weak indications from the data in hand, (2) draw on experience with other bodies of data, or (3) lean on subject-matter knowledge. Even all three may not suffice. Both because we may not be prepared to try hard to choose our re-expression, or because we have too little information for anyone to choose reliably, we need rules of thumb that can provide “first aid,” that can lead us to re-expressions that are almost always not bad — and usually pretty good.

Four rules will deal quite effectively with most of our needs, namely:

1. Take logs of an amount or count (if there are zeros or infinities, we may need to deal with them; see the next section)
2. Take logits or folded logs of fractions or percents; use some multiple of

$$\log \left(\frac{p}{1-p} \right)$$

...

These rules are not supposed to be a final answer—just as first aid for the injured is no substitute for a physician—but they offer a safe beginning.”

Also see discussion of “Tukey’s ladder of re-expressions” in Tukey (1977), *EDA*, pp. 90–1.

⁴Mosteller and Tukey (1977), *Data Analysis and Regression*, p. 109

Cereal Problem

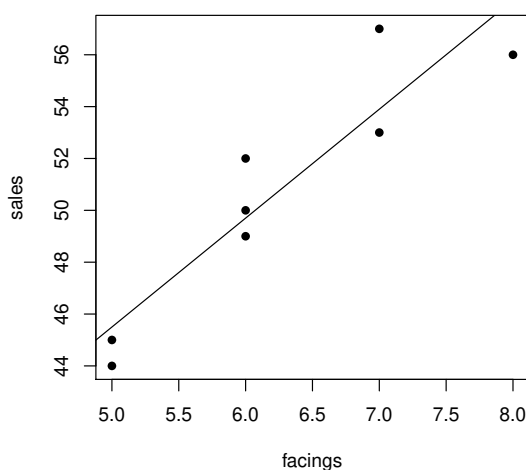
A cereal manufacturer believes that there is an association between cereal sales and the number of facings the cereal has on each stores' shelves. Eight stores were surveyed to test this hypothesis.

```
> dat = data.frame(facings=c(5,6,6,7,5,7,6,8), sales = c(45,50,52,53,44,57,49,56))
> plot(dat, pch=16)
> fit = lm(sales ~ facings, data = dat)
> abline(fit)
> summary(fit)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	24.500	4.541	5.395	0.00167 **
facings	4.200	0.718	5.849	0.00110 **

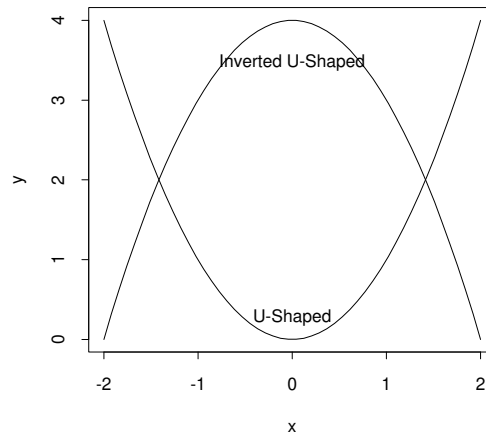
Residual standard error: 1.966 on 6 degrees of freedom
Multiple R-squared: 0.8508, Adjusted R-squared: 0.8259
F-statistic: 34.22 on 1 and 6 DF, p-value: 0.001102



How do you like this model? ([Video solution](#))

Polynomial Transformations

Consider models the form $y = \beta_0 + \beta_1x + \beta_2x^2$



- The min/max value occurs at $x_{\text{opt}} = -\beta_1/2\beta_2$
- U-shaped when $\beta_2 > 0$
- Higher-order polynomials can also be used, e.g., cubic $\beta_0 + \beta_1x + \beta_2x^2 + \beta_3x^3$, but they are problematic
 - Curvature can change several times — few theories postulate this.
 - x, x^2, x^3 highly correlated — mean center or standardize x 's before fitting model.
 - Interpretation complicated — don't interpret estimates of individual terms. Use F -tests to gauge significance and plots to interpret effects.

Purification case

The “techies” (scientists) in the laboratory have been lobbying you, and management in general, to include just one more laboratory step. They think it’s a good idea, although you have some doubt because one of them is known to be good friends with the founder of the start-up biotechnology company that makes the reagent used in the reaction. But if adding this step works as expected, it could help immensely in reducing production costs. The trouble is, the test results just came back and they don’t look so good. Discussion at the upcoming meeting between the technical staff and management will be spirited, so you’ve decided to take a look at the data.

Your firm is anticipating government approval from the Food and Drug Administration (FDA) to market a new medical diagnostic test made possible by monoclonal antibody technology, and you are part of the team in charge of production. Naturally, the team has been investigating ways to increase production yields or lower costs.

The proposed improvement is to insert yet another reaction as an intermediate purifying procedure. This is good because it focuses resources down the line on the particular product you want to produce. But it shares the problem of any additional step in the laboratory: one more manipulation, one more intervention, one more way for something to go wrong. In this particular case, it has been suggested that, while small amounts of the reagent may be helpful, trying to purify too well will actually decrease the yield and increase costs.

The design of the test was to have a series of test production runs, each with a different amount of purifier, including one test run with the purification step omitted entirely (i.e., 0 purifier). The order of the tests was randomized so that any time trends would not be mistakenly interpreted as being due to purification.

```
purify = data.frame(x = 0:10,
  y=c(13.39,11.86,27.93,35.83,28.52,41.21,37.07,51.07,51.69,31.37,21.26))
```

1. Regress `yield` on `amount`. Is the regression significant? Based on this test alone, do you recommend including a purifying step in the process?
2. Generate a scatterplot of `yield` against `amount`. Comment.
3. Modify your regression model as appropriate. Based on the revised analysis, do you recommend including a purifying step in the process?

[Video solution](#)

Multiplicative Models

- Multiplicative models have the form

$$y = \beta_0 x_1^{\beta_1} x_2^{\beta_2} \epsilon$$

- Taking natural logs of both sides we get

$$\log y = \log \beta_0 + \beta_1 \log x_1 + \beta_2 \log x_2 + \log \epsilon,$$

which has the form of a multiple linear regression model, regressing $\log y$ on $\log x_1$ and $\log x_2$. We estimate this model under the assumption that $\log \epsilon$ is normal with $\mathbb{V}(\log \epsilon) = \sigma^2$ constant.

- Predictions: we estimate (untransformed) y with

$$\exp(b_0 + b_1 \log x_1 + b_2 \log x_2 + S_\epsilon^2/2)$$

- The coefficient of variation of e (rather than the variance) is

$$\frac{\sqrt{\mathbb{V}(y)}}{\mathbb{E}(y)} = \sqrt{\exp(\sigma^2) - 1},$$

which doesn't depend on i (and is thus constant).

- Even without the assumption that ϵ is log normal, we will see that if $\sigma_y \propto \mu_y$ then logging y “stabilizes” the error variance, making it constant. Logging y is called a **variance stabilizing transformation**.

Interpreting β_1

$$y = \beta_0 x^{\beta_1} \quad \frac{dy}{dx} = \beta_0 \beta_1 x^{\beta_1-1} \quad \frac{d^2y}{dx^2} = \beta_0 \beta_1 (\beta_1 - 1) x^{\beta_1-2}$$

- Interpretation (ignoring error term, and assuming $\beta_0 > 0$ and $x > 0$)
 - $\beta_1 = 1 \implies y \propto x$ (linear, proportional returns)
 - $0 < \beta_1 < 1 \implies$ changes in y *decrease* as x increases (concave downward)
 - $\beta_1 > 1 \implies$ changes in y *increase* as x increases (concave upward)
- In economic applications, β_j is the **elasticity** of y with respect to x_j — the expected percentage change in y of a 1% change in x_j , all else being equal. Let dx be an infinitesimal change in x and dy be the corresponding change in y . Then the elasticity is

$$\frac{dy/y}{dx/x} = \frac{dy}{dx} \cdot \frac{x}{y} = \beta_0 \beta_1 x^{\beta_1-1} \cdot \frac{x}{\beta_0 x^{\beta_1}} = \beta_1$$

- We estimate β_1 by logging both sides. The base of the logarithm does not matter, but natural logs are usually used.

Background: Lognormal Distribution

- If y has a normal distribution with mean μ and variance σ^2 , then $y^* = \exp(y)$ has a **log-normal distribution**, i.e., $\log y^* = y$ is normal.
- Theorem: the mean and variance of y^* are

$$\mathbb{E}(y^*) = \exp(\mu + \sigma^2/2)$$

$$\mathbb{V}(y^*) = \exp(2\mu + \sigma^2)(\exp(\sigma^2) - 1)$$

- The coefficient of variation of y^* is

$$\begin{aligned} \frac{\sqrt{\mathbb{V}(y^*)}}{\mathbb{E}(y^*)} &= \frac{\sqrt{\exp(2\mu + \sigma^2)(\exp(\sigma^2) - 1)}}{\exp(\mu + \sigma^2/2)} \\ &= \sqrt{\exp(\sigma^2) - 1}, \end{aligned}$$

which does not depend on μ

Business Failure Case

Consider the slightly scary topic of business failures. This problem analyzes data from each state on the number of failed businesses and the population in thousands for each of the 50 states and the District of Columbia (51 observations).

```
busfail = data.frame(  
  row.names=c("AL", "AK", "AZ", "AR", "CA", "CO", "CT", "DE", "DC", "FL", "GA", "HI", "ID",  
    "IL", "IN", "IA", "KS", "KY", "LA", "ME", "MD", "MA", "MI", "MN", "MS", "MO", "MT", "NE",  
    "NV", "NH", "NJ", "NM", "NY", "NC", "ND", "OH", "OK", "OR", "PA", "RI", "SC", "SD", "TN",  
    "TX", "UT", "VT", "VA", "WA", "WV", "WI", "WY"),  
  pop=c(4187, 599, 3936, 2424, 31211, 3566, 3277, 700, 578, 13679, 6917, 1172, 1099, 11697,  
    5713, 2814, 2531, 3789, 4295, 1239, 4965, 6012, 9478, 4517, 2643, 5234, 839, 1607, 1389,  
    1125, 7879, 1616, 18197, 6945, 635, 11091, 3231, 3032, 12048, 1000, 3643, 715, 5099,  
    18031, 1860, 576, 6491, 5255, 1820, 5038, 470),  
  fail=c(841, 108, 2064, 186, 19695, 1542, 1093, 137, 200, 5088, 2350, 305, 350, 2094, 1091, 507,  
    1069, 841, 664, 383, 1540, 2720, 2546, 921, 322, 1230, 173, 399, 568, 617, 2843, 448, 6916, 1194,  
    145, 2127, 1440, 969, 3124, 344, 392, 175, 1209, 7096, 351, 173, 1738, 2025, 315, 1224, 90)  
)
```

1. Make a scatterplot of business failures against population and superimpose a regression line. Describe the relationship. Comment on whether the linear model appears to hold.
2. Make a scatterplot of the log of business failures against the log of population. Superimpose a regression line. Does the linear model hold better with the logged data?
3. Regress the log of failures on the log of population. State estimated regression equation.
4. Test at the 5% level to see whether there is a significant relationship between the logs of failure and population. Explain.
5. Test whether the population slope for the logs is significantly different from 1 or not. What does this tell you?
6. Illinois had 2,094 failures with a population of 11,697 (thousand). Find the predicted log failures for Illinois.
7. Estimate the unlogged failures and identify Illinois on your scatterplot. (Hint: use `identify` function in R.)

Your Turn

1. The table below shows the level of investment and the results obtained by the important players in fiber-optics cable for long-distance communications.

```
fiber = data.frame(invest=c(1300,500,130,2000,1200,110,40,60,57,500,90,90),  
  miles=c(1700,650,110,1200,2400,165,72,45,85,650,50,87))
```

- (a) Find the regression equation predicting circuit miles from investment.
 - (b) Draw the scatterplot and residuals. Discuss whether the linear model holds.
 - (c) Examine Cook's distance. Are there influential observations, as indicated by having Cook's distance greater than 1?
 - (d) Regress the log of circuit miles on the log of investment. State the estimated equation.
 - (e) Draw the scatterplot and residuals for the log model. Discuss whether the linear model holds.
 - (f) Do firms that spend more achieve significantly more circuit miles? State the null and alternative, P -value and decision using the 5% level of significance.
 - (g) Test at the 5% level whether the coefficient for $\log(\text{investment})$ is different from 1, which indicates that investment is proportional to miles. What do values not equal to 1 indicate in terms of economies of scale?
 - (h) Predict $\log(\text{miles})$ from an investment of \$1000.
 - (i) Predict unlogged miles from an investment of \$1000.
2. A research analyst for an oil company wants to develop a model to predict miles per gallon based on highway speed. An experiment is designed in which a test car is driven at speeds ranging from 10 miles per hour to 75 miles per hour in increments of 5 MPH. Two replicates were observed for each speed:

```
speed = data.frame(mph=rep(seq(10,75,by=5), 2),  
  mpg=c(4.8,8.6,9.8,13.7,18.2,19.9,22.4,21.3,20.5,18.6,14.4,12.1,10.1,8.4,  
    5.7,7.3,11.2,12.4,16.8,19,23.5,22,19.7,19.3,13.7,13,9.4,7.6))
```

- (a) Make a scatterplot of MPG against MPH. Based on the plot, do you suggest transforming the data? If, so which transformation do you suggest?
- (b) Regress MPG on MPH (do not include transformations yet). Report (i) the estimated regression equation, (ii) the P -value testing the overall significance of the model, and (iii) a residual plot of residuals versus predicted values.
- (c) Add appropriate transformations to your model. Report (i) the estimated regression equation, (ii) the P -value testing the overall significance of the model, and (iii) a residual plot with comments about the fit of your model.

- (d) Using the model you developed in the previous part, estimate gas milage for a car traveling at 62 miles per hour.
- (e) At what speed is gas milage maximized?
- (f) Produce a scatterplot with the two fits superimposed (linear, quadratic).
3. Use the auto data set from JWHT problem 3.9 on page 136.
- (a) The `origin` variable is categorical, where 1=US, 2=Europe and 3=Japan. We'll cover dummies next week, but for now type the following command to make it a factor variable and assign meaningful labels:
- ```
auto$origin = factor(auto$origin, 1:3, c("US", "Europe", "Japan"))
```
- (b) Regress `mpg` on `origin`, `weight` and `year`. Examine the diagnostic plots and comment on which assumptions of the linear model, if any, are violated.
- (c) Regress `log(mpg)` on `origin`, `log(weight)`, `year` and `year squared`. Examine the diagnostic plots and the summary. For you to think about but not turn in: why would year have this effect for year?
- (d) Describe the effect of year on `log(mpg)`, i.e., is it U-shaped, inverted-U shaped, or linear? If it is nonlinear, where is the minimum or maximum. Draw a graph showing the effect.
- (e) What does the coefficient for `log(weight)` tell you?

## Answers

1. Fiber problem. (a)  $\text{miles}(\text{hat}) = 101.813 + 0.986 \text{ invest}$ ; (b) The linear model does not hold because the variance of the residuals increases with the mean of miles. (c) Yes, observation 4 is influential, and 5 is nearly "influential." (d)  $\log(\text{miles}) = 0.06803 + 1.00735 \log(\text{invest})$ . (e) The residuals have more constant variance and no obvious pattern. (f)  $H_0 : \beta = 0$  versus  $H_1 : \beta \neq 0$ ,  $P = 1.02 \times 10^{-6} / 2 < .05$ , so we reject  $H_0$ . (g)  $H_0 : \beta = 1$  versus  $H_1 : \beta \neq 1$ . A 95% confidence interval for the slope is  $1 \in [0.79, 1.22]$ . Since 1 is in this interval, we cannot reject the null hypothesis that the slope is 1. It is plausible that miles are proportional to investment. Recall that  $\text{investment} = e^a \times \text{invest}^b$ . (h) `predict(fit, data.frame(invest=1000)) = 7.026552` (i)  $\exp(7.026552 + 0.1938/2)$ , where 0.1938 is the MSE from the ANOVA table.
2. Speed problem. (a) The scatterplot shows an inverted-U shaped relationship, but no heteroscedasticity. We should add a quadratic term for speed. (b)  $\text{mpg} = 12.75 + 0.039 \text{ speed}$ .  $P = 0.468 > 5\%$  so we cannot reject  $H_0 : \beta = 0$ . The residuals show a strong pattern indicating that the model is misspecified. (c) Add a quadratic term:  
 $\text{mpg} = -7.56 + 1.27\text{speed} - 0.0145\text{speed}^2$ .  
 $P = 2.338 \times 10^{-14} < 5\%$ , so we reject  $H_0 : \beta_1 = \beta_2 = 0$ . The residuals are not perfect, but the magnitude of the pattern is reduced. It looks like MPG increases linearly until about 40MPH. The quadratic model provides a substantially better fit ( $R^2 = .9188$ ) than the linear model ( $R^2 = .02039$ ), but the quadratic model could be improved further. (e) `predict(fit, data.frame(speed=62)) = 15.54618` (g)  $-1.27/(2 \times -0.0145) = 43.8$  MPH
3. JWHT 3.9 (b) The residual plot shows a pattern, with consistently positive values, then negative, then positive again. This indicates that the fit can be improved. (d) The quadratic coefficient  $0.0019 > 0$  indicating a U shape, with a min value at  $0.2559684/(2 \times 0.0019051) = 67.17978$ , i.e., around 1967. To make a graph, note that year ranges from 70 to 82 and see code below to see the effect. (e) The coefficient is negative, which indicates that as the weight of the car increases the milage decreases.
- ```
fit = lm(mpg~mph, speed)
fit2 = lm(mpg~mph+I(mph^2), speed)
# not part of exercise, but interesting
fit3 = lm(mpg~as.factor(mph), speed)
plot(speed); abline(fit)
lines(10:75, predict(fit2, data.frame(mph=10:75)), col=2)
lines(speed$mph[1:14], fit3$fit[1:14], type="b", col=4, pch=16)
```
- ```
fit = lm(mpg ~ weight + year, auto) # part b
plot(fit, which=1) # part b
fit = lm(log(mpg) ~ log(weight) + year + I(year^2)
+ origin, data=auto) # part c
summary(fit) # part c
drop1(fit, test="F") # part c
x = 70:82 # part d
plot(x, -0.255968375*x + 0.001905147*x^2, type="l")
```

## Newfood case

Mr. Conrad Ulcer, newly appointed New Products Marketing Director for Concorn Kitchens, was considering the possibility of marketing a new highly nutritional food product with widely varied uses. This product could be used as a snack, a camping food, or as a diet food. The product was to be generically labeled Newfood.

Because of this wide range of possible uses, the company had great difficulty in defining the market. The product was viewed as having no direct competitors. Early product and concept tests were very encouraging. These tests led Mr. Ulcer to believe that the product could easily sell 2 million cases (24 packages in a case) under the proposed marketing proposal involving a 24-cent package price and an advertising program involving \$3 million in expenditures per year. There were no capital expenditures required to go national, since manufacturing was to be done on a contract-pack basis.

Because there was considerable uncertainty among Concorn Management as to either probable first-year and subsequent-year sales, or the best introductory campaign, Ulcer decided that a six-month market test would be conducted. The objectives of the test were to:

- Better estimate first-year sales.
- Study certain marketing variables to determine an optimal — or at least better — introductory plan.
- Estimate the long-run potential of the product

These objectives were accomplished through the controlled introduction of the product into four markets. Conditions were experimentally varied within the grocery stores in each of the four markets. Sales were measured with a store audit of a panel of stores. Preliminary results had been obtained. Now it was up to Mr. Ulcer to understand their implications on the introduction of Newfood.

## Design of Experimental Study

The three variables included in the experimental design were price, advertising expenditures, and location of the product within the store. Three prices were tested (24 cents, 29 cents, and 34 cents), two levels of advertising (a simulation of a \$3 million introduction and a \$6 million plan), and two locations (placing the product in the bread section versus the instant breakfast section). Prices and location were varied across stores within cities while advertising was varied across cities. The advertising was all in the form of TV spots. The levels were selected so that they would stimulate on a local basis the impact that could be achieved from national introduction programs at the \$3 million and \$6 million expenditure levels. Due to differential costs between markets and differential costs between spot and network (to be used in national introduction), an attempt was made to equate (and measure) advertising inputs of gross advertising impressions generated, normalized for

market size. Unfortunately, it was not possible to achieve exactly the desired levels. This was due to the problem of non-availabilities of spots in some markets and discrepancies between estimates of TV audiences made at the time the test was being planned and the actual audiences reached at the time the commercials were actually run.

In the selection of cities and stores for the tests, attempts were made to match stores on such variables as store size, number of checkout counters, and characteristics of the trading area. Because it was not certain that adequate matches had been achieved, Ulcer decided to obtain measurements on some of these variables for possible use in adjusting for differences in cell characteristics. He also felt that it might be possible to learn something about the relationships between these variables and sales, and that this information would be of assistance in planning the product introduction into other markets.

```
newfood = data.frame(
 sales=c(225,323,424,268,224,331,254,492,167,226,210,289,204,288,245,161,161,
 246,128,154,163,151,180,150),
 price=c(24,24,24,24,24,24,24,24,29,29,29,29,29,29,34,34,34,34,34,34,34),
 ad=c(0,0,1,1,0,0,1,1,0,0,1,1,0,0,1,1,0,0,1,1,0,0,1,1),
 loc=c(0,0,0,0,1,1,1,1,0,0,0,0,1,1,1,1,0,0,0,0,1,1,1,1),
 income=c(7.3,8.3,6.9,6.5,7.3,8.3,6.9,6.5,6.5,8.4,6.5,6.2,6.5,8.4,6.5,6.2,
 7.2,8.1,6.6,6.1,7.2,8.1,6.6,6.1),
 volume=c(34,41,32,28,34,41,23,37,33,39,30,27,37,43,30,19,32,42,29,24,32,36,29,24),
 city=c(3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2,3,4,1,2))
```

### Questions for class discussion

1. Compute the correlation matrix. How do you explain the 0 correlations (e.g., between location and advertising)?
2. Run a regression of **sales** (the first two months sale) on **price** alone. Next, on **price** and **ad**. Finally on **price**, **ad**, and **loc**. Thus, you will have three regressions. What happens to the coefficients of **price** in the three regressions? What happens to the coefficients of **ad** in the two regressions? Explain.
3. Run a regression of **sales** against **price**, **ad**, **loc**, and **volume**. What happens to the coefficients of **price**, **ad**, and **loc** which you found in the third regression in question 2 above? Which coefficient changes the most with the introduction of store size? Why does this happen?
4. Finally, run a regression of **sales** against **price**, **ad**, **loc**, **volume**, and **income**. What changes do you observe between these results and that of the fourth regression? Explain.
5. What additional regression runs, if any, should be made to complete the analysis of these data?

## Effects of Model Misspecification: Omitting Relevant Predictor

---

- Suppose we fit the model

$$y = \beta_0 + x\beta_1 + \epsilon$$

- But the true model is

$$y = \beta_0 + x\beta_1 + z\beta_2 + \epsilon$$

- Theorem: (**Omitted variable bias**) If  $x$  and  $z$  are correlated,  $b_1$  will be biased, with the direction of the bias depending on the sign of the correlation between  $x$  and  $z$  and the sign of  $\beta_2$ .

$$\mathbb{E}(b_1) = \beta_1 + \beta_2 r_{xz} \frac{S_z}{S_x}$$

where  $S_x$  and  $S_z$  are the sample standard deviations of  $x$  and  $z$   
 Direction of Bias in  $b_1$

| Sign of correlation<br>between $x$ and $z$ | if $\beta_2 > 0$ | if $\beta_2 < 0$ |
|--------------------------------------------|------------------|------------------|
| $\text{corr}(x, z) = r_{xz} > 0$           | Upward bias      | Downward bias    |
| $\text{corr}(x, z) = r_{xz} = 0$           | No change        | No change        |
| $\text{corr}(x, z) = r_{xz} < 0$           | Downward bias    | Upward bias      |

- Note that if  $x$  and  $z$  are uncorrelated (orthogonal design), we can add or drop variables without changing the other coefficients.

# 1. Multicollinearity Definition and Effects

---

- Definition: **Multicollinearity** is when the predictor variables are highly correlated with one another (Note:  $Y$  is not mentioned here). When more than one  $X$ 's move together, it is difficult to sort out their separate effects.
- What are some effects? You cannot sort out what is doing what.
  - *Unstable coefficients*. High estimated standard errors for one or more slope coefficients.
    - \* Implication: Low  $t$ -ratio, so sometimes we cannot reject the null hypothesis that  $\beta = 0$ . At the extreme, the model as a whole may be significant, while none of the individual slope parameters are!
    - \* Coefficients can change wildly when variables are added or dropped from the model.
  - *Incorrect signs*. Slope estimates can have signs that are not consistent with intuition.
- Note: predicted values not affected directly



## Predicted Values Unaffected

---

- Suppose the true model is

$$y = 2x_1 + x_2$$

- We have  $n = 3$  observations:

| $x_1$ | $x_2$ | $y$ | e |
|-------|-------|-----|---|
| 1     | 2     | 4   | 0 |
| 2     | 4     | 8   | 0 |
| 3     | 6     | 12  | 0 |

Note that  $x_2 = 2x_1$ , so that the two columns are perfectly correlated

- The true model fits perfectly, but so do many others, e.g.,

$$y = 4x_1 + 0x_2 \implies x_2 \text{ no effect}$$

$$y = 0x_1 + 2x_2 \implies x_1 \text{ no effect}$$

$$y = 6x_1 - x_2 \implies x_2 \text{ negative effect}$$

**Predictions correct for all these choices of parameter estimates, but substantive interpretation completely different!**

- Which  $(b_1, b_2)$  is correct? We can't tell from these data.
- What happens if we use our fitted model to *extrapolate*? e.g., estimate  $y$  for  $x_1 = 1$  and  $x_2 = 1$ . The correct answer is 3, but the other models gives estimates 4, 2, and 5, respectively. **Extrapolation is especially questionable when using a model estimated from multicollinear data.**

## 2. Detecting Multicollinearity

---

- Compute a correlation matrix of the predictor variables and possibly a scatterplot matrix. Large correlations indicate multicollinearity could be a problem.
- Unstable coefficients or incorrect signs
- **Tolerance** is  $1 - R^2$  from this regression (fraction of variance *unexplained* by the model)
- **Variance inflation factor** (VIF) is  $1/\text{tolerance}$ .

```
> install.packages("car") # do this only once
> library(car) # you must first download it from CRAN
> fit = lm(sales~price+ad+loc+volume+income, newfood)
> vif(fit)
 price ad loc volume income
1.079882 2.697664 1.005447 3.447143 3.367158

> # where does VIF for ad come from?
> summary(lm(ad ~ price + loc+volume+income, newfood))

...
Multiple R-squared: 0.6293, Adjusted R-squared: 0.5513

> 1/(1-.6293)
[1] 2.697599
```

- Interpretation
  - VIF=1: no multicollinearity
  - $2 < \text{VIF} < 5$  beware
  - $5 < \text{VIF} < 10$  substantial multicollinearity
  - $\text{VIF} \geq 10$  severe multicollinearity

## Living with Multicollinearity

---

1. If the objective is primarily to make good *predictions*, then you could do nothing, although you may be better off using stepwise, ridge/lasso, or principal components regression (PCR). Data mining applications often fall into this category.
2. If the objective is to *interpret* regression coefficients, then
  - (a) If possible, avoid multicollinearity with an *orthogonal* design, where predictor variables are uncorrelated
  - (b) Understand why you have multicollinearity
    - **Pipe**, e.g.,  $x_1 \rightarrow x_2 \rightarrow y$ . If you are studying  $x_1 \rightarrow y$  then do not control for  $x_2$  because it blocks path.
    - Predictors manifestations of common, underlying **latent construct**, e.g.,  $w \rightarrow x_1$  and  $w \rightarrow x_2$ . Often, estimate  $w$  and use it instead of  $x_1$  and  $x_2$ .
    - “Back-door” confound (**fork**). Include control to block back-door path, e.g., if  $w \rightarrow y$  and  $w \rightarrow x \rightarrow y$  then control for  $w$  to study  $x \rightarrow y$ .
    - We usually<sup>5</sup> do not control for **colliders**, e.g., if  $x \rightarrow w$  and  $y \rightarrow w$ , then do not control for collider  $w$  when studying  $x \rightarrow y$ .

---

<sup>5</sup>As we will see in a few weeks, there can be complicated situations where we must have a collider as a control, but then we will have to add other control(s) to fix the problems the collider creates.

## Model Specification Issues

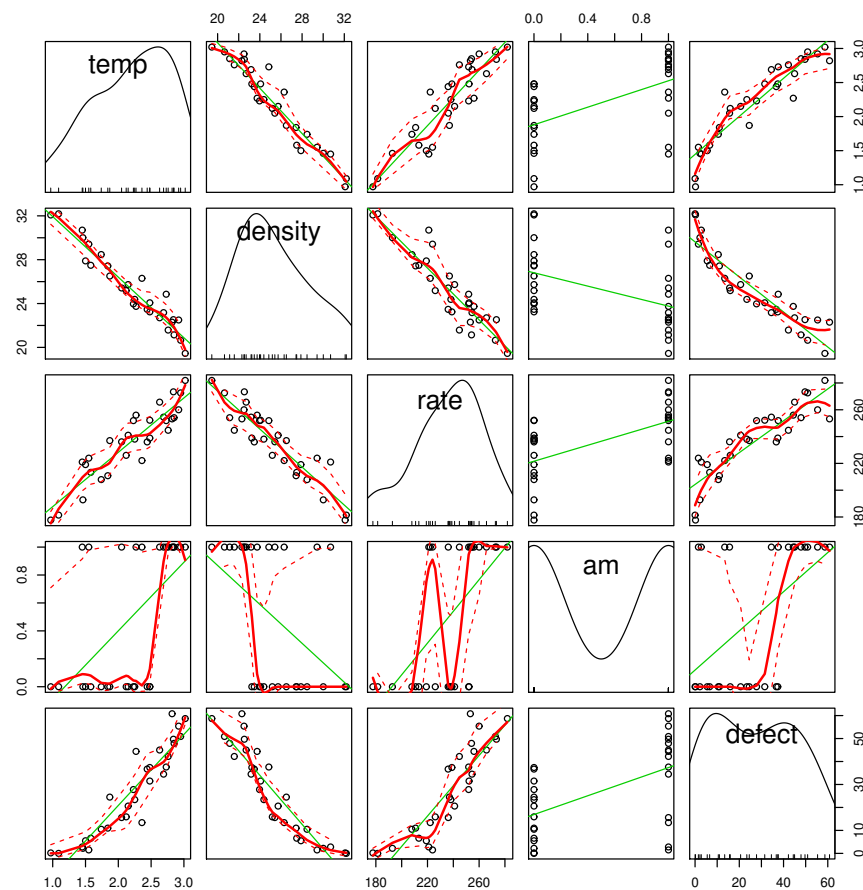
---

- For successful use of multiple regression, sufficient knowledge *about the subject domain* is required to identify relevant predictor variables and their functional relationship with the dependent variable.
- If the only goal is predictive, it should make sense for a variable to be in the model, e.g., does it make sense to have customer ID number as a predictor variable (it's probably a proxy for something else, e.g., tenure)?
- If the goal is confirmatory, **begin with a conceptual framework** (directed acyclic graph or DAG). Include a variable if
  - It is a decision variable
  - The variable helps to control for important causal factors (forks), e.g., seasonality or competitive actions. Remember the omitted variable bias theorem.
  - Unless you are really sure of your framework, do robustness checks by adding/dropping variables.

# Scatterplots with CAR

The car package offers improved scatterplot matrices:

```
> library(car) # do this if you haven't ready done so
> scatterplotMatrix(~temp+density+rate+am+defect, quality)
```



- The red lines are *smoothers*, which trace the middle of the distribution of the vertical variable conditional on the horizontal variable. Green lines are robust regression lines.
- The ticks are called a *rug* and show individual observations.
- The graphs on the diagonal are density estimates (like a histogram).

## Quality Control Case

Everybody seems to disagree about just why so many parts have to be fixed or thrown away after they are produced. Some say that it's the temperature of the production process, which needs to be held constant (within a reasonable range). Others claim that it's clearly the density of the product, and that if we could only produce a heavier material, the problems would disappear. Then there is Ole, who has been warning everyone forever to take care not to push the equipment beyond its limits. This problem would be the easiest to fix, simply by slowing down the production rate; however, this would increase costs. Interestingly, many of the workers on the morning shift think that the problem is "those inexperienced workers in the afternoon," who, curiously, feel the same way about the morning workers.

Ever since the factory was automated, with computer network communication and bar code readers at each station, data have been piling up. You've finally decided to have a look. After your assistant aggregated the data by 4-hour blocks and then typed in the AM/PM variable, you found the following description of the variables:

- **temperature**: measures the temperature variability as a standard deviation during the time of measurement
- **density**: indicates the density of the final product
- **rate**: rate of production
- **am**: 1 indicates morning and 0 afternoon
- **defect**: average number of defects per 1000 produced

### Discussion Questions

1. Generate a scatterplot matrix and a correlation matrix. Interpret the correlations. What "obvious conclusions" can you draw?
2. Run a multiple regression predicting defect rate from the other four variables. Is the overall model significant? Which predictors, if any, are significant? Compute and interpret variance inflation factors. What "obvious conclusions" can you draw?
3. Predict defect from each of the predictor variables separately, e.g., **defect** from **temp**, **defect** from **density**, **defect** from **rate**, etc. Which of the predictors are significant in the simple linear regressions?
4. Perform further analysis as needed. What action do you recommend? Why? Hint: think about the causal relationships between the variables.
5. To compare the two shifts, would it be appropriate to perform an independent-sample  $t$ -test (i.e.,  $H_0$  : AM defect rate = PM defect rate)? How is this different from the multiple regression approach? Which is preferred? Discuss.

6. How would you present your findings to a client?

```
quality = data.frame(
 temp=c(.97,2.85,2.95,2.84,1.84,2.05,1.5,2.48,2.23,3.02,2.69,2.63,1.58,2.48,2.25,
 2.76,2.36,1.09,2.15,2.12,2.27,2.73,1.46,1.55,2.92,2.44,1.87,1.45,2.82,1.74),
 density=c(32.08,21.14,20.65,22.53,27.43,25.42,27.89,23.24,23.97,19.45,23.17,
 22.7,27.49,24.07,24.38,21.58,26.3,32.19,25.73,25.18,23.74,24.85,30.01,
 29.42,22.5,23.47,26.51,30.7,22.3,28.47),
 rate=c(177.7,254.1,272.6,273.4,210.8,236.1,219.1,238.9,251.9,281.9,254.5,265.7,
 213.3,252.2,238.1,244.7,222.1,181.4,241,226,256,251.9,192.8,223.9,260.0,236,
 237.0,221,253.2,207.9),
 am=c(0,1,1,1,0,1,0,0,0,1,1,1,0,0,0,1,1,0,0,0,1,1,0,1,1,0,0,1,1,0),
 defect=c(.2,47.9,50.9,49.7,11,15.6,5.5,37.4,27.8,58.7,34.5,45,6.6,31.5,23.4,
 42.2,13.4,0,20.6,15.9,44.4,37.6,2.2,1.5,55.4,36.7,24.5,2.8,60.8,10.5)
)
```

# Your Turn

---

1. Use the auto data set from JWHT problem 3.9 on page 122. Type the following:

```
auto = read.table("auto.txt", header=T) # auto.txt on Canvas
auto$origin = factor(auto$origin, 1:3, c("US", "Europe", "Japan"))
```

The data set is also in the ISLR library. Since you are changing the origin variable it might be best not to touch the ISLR file.

- (a) Regress `mpg` on `cylinders`, `displacement`, `weight`, and `year`. Comment on the signs of the estimated coefficients and note which are significantly different from 0. What is value of  $R^2$ ?
  - (b) Compute the variance inflation factors. What do they tell you?
  - (c) Drop `weight` from the model. What happens to the parameter estimates and  $R^2$ ?
  - (d) Drop `weight` and `displacement` from the model. What happens to the parameter estimates and  $R^2$ ?
2. Use the data set `part.csv`, available from canvas. This question investigates how participation in a social media contest about a brand affect future spending on the brand. A brand sponsored a social media contest. Customers in the company's database were invited to write about their relationship with the company on a social media forum. Those who participated by writing at least one word on the forum received a reward worth approximately \$1, and the dummy variable `tx` indicates whether or not a customer participated. In total, 7089 customers participated, and there is a matched control group of 7089 consistent of customers who did not participate, but had similar purchase activities prior to the contest. The total sample size is thus  $2 \times 7089 = 14,178$ . The variable `y` records the amount spent by each customer in the week following the contest. The variable `x` gives the amount spent per week prior to the contest and will be used as a control variable to account for differences in customer loyalty. Finally, the `wc` variable gives the word count of the entries, where `wc` = 0 for all who did not participate. Word count measures *cognitive elaboration*. Note that `tx` = (`wc` > 0).
- (a) **Model 1:** regress  $\log(y + 1)$  on  $\log(x + 1)$  and `tx`. Give the output.
  - (b) **Model 2:** regress  $\log(y + 1)$  on  $\log(x + 1)$ , `tx` and  $\log(\text{wc} + 1)$ . Give the output.
  - (c) Use the following notation in answering the questions:

$$\log(y + 1) = \beta_0 + \beta_1 \log(x + 1) + \beta_2 \text{tx} + \beta_3 \log(\text{wc} + 1) + e,$$

where  $\beta_3$  is constrained to be 0 in Model 1 and  $\log$  is the natural log. Based on Model 1, does participation have a significant effect on future spending? Explain. Note: to receive full credit you should state null and alternative hypotheses and do something to determine whether  $H_0$  can be rejected at the 5% level.



- (d) Using Model 1, post-period *spending* is how many *times* greater for those who participate than for those who do not? Note that this question asks about *spending* and not log spending. Another way to ask this question is, suppose there are two people with identical pre-period spending, but one participates and the other does not. If  $y_1$  is the post-contest spending of a participant, and  $y_0$  is the post-contest spending of a non-participant, how many times greater is  $(y_1 + 1)$  than  $(y_0 + 1)$ ?
- (e) Is  $(y + 1)$  proportional to  $(x + 1)$ , i.e., is spending in the week after the contest proportional to pre-contest spending? How do you know? Note: to receive full credit, state a null and alternative hypothesis and do something to determine whether or not  $H_0$  can be rejected.
- (f) Why is the magnitude of the  $tx$  variable so different in Model 2 (0.050) than in Model 1 (0.244)?
- (g) Now consider Model 2. How do the results from Model 2 change your conclusions about how participation affects future spending. I am looking for you to summarize the key learnings from Model 2 succinctly.
- (h) Generate the normal probability plot for Model 2. What specifically does the plot tell you, and how does it (i.e., what the plot tells you) affect the conclusions you have drawn from the previous parts.
- (i) What do the results of this analysis suggest the company should do in the future when designing social media contests?
3. JWHT problem 3.14a–f on page 125. For part (c)–(e), are the parameters “covered” by the 95% confidence intervals?

## Answers

1. JWHT 3.9. (a) It is odd that displacement has a positive sign, although it is not significant. Only weight and year are significant.  $R^2 = .8091$ . (b) The VIF values for all variables but year are large, indicating multicollinearity: cars that weigh more tend to have larger engine displacement and more cylinders. (c) The coefficients change drastically. Cylinders goes from  $-.29$  to  $-.62$ . Displacement goes from  $0.00497$  to  $-0.0415$  and now has a positive, significant slope. Year changes less.  $R^2$  decreases a little bit to  $0.7423$ . The other variables do the much of the explaining previously done by weight. (d) Cylinders is now significant, and  $R^2$  reduces only slightly to  $.7135$ .

```
part a
Call: lm(mpg ~ cylinders + displacement + weight + year)
```

```
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) -14.076941 4.055159 -3.471 0.000575 ***
cylinders -0.289589 0.329225 -0.880 0.379611
displacement 0.004973 0.006701 0.742 0.458425
weight -0.006702 0.000572 -11.717 < 2e-16 ***
year 0.764751 0.050684 15.089 < 2e-16 ***
```

```
Residual standard error: 3.436 on 392 degrees of freedom
Multiple R-squared: 0.8091, Adjusted R-squared: 0.8072
F-statistic: 415.5 on 4 and 392 DF, p-value: < 2.2e-16
```

```
part b
 cylinders displacement weight year
10.524432 16.406259 7.888061 1.173000
```

```
part c
```

```
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) -18.199719 4.688296 -3.882 0.000122 ***
cylinders -0.620910 0.380657 -1.631 0.103658
displacement -0.041545 0.006265 -6.632 1.1e-10 ***
year 0.699324 0.058461 11.962 < 2e-16 ***
```

```
Residual standard error: 3.988 on 393 degrees of freedom
Multiple R-squared: 0.7423
```

```
part d
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) -17.30285 4.93534 -3.506 0.000507 ***
cylinders -3.00405 0.13223 -22.718 < 2e-16 ***
year 0.75289 0.06098 12.347 < 2e-16 ***
```

```
Multiple R-squared: 0.7135
```

2. Participation problem. (c)  $H_0 : \beta_2 = 0$  versus  $H_1 : \beta_2 \neq 0$ . The  $P$ -value is less than 5%, so reject  $H_0$ . (d)  $e^{0.24438} = 1.2768$ . (e)  $H_0 : \beta_1 = 1$  versus  $H_1 : \beta_1 \neq 1$ . Students should either find a 95% confidence interval:  $0.80318 \pm 1.96 \times 0.01205 = [0.7796, 0.8268]$  and note that it does not cover 1, or find the  $P$ -value to reject  $H_1$ . Note:  $(0.80318 - 1)/0.01205 = -16.33$ , so  $P < .05$ . (f) Multicollinearity. We expect there to be a positive correlation between  $tx$  and  $\log(wc+1)$  since  $tx=0$  implies  $wc=0$ . The effect of  $tx$  in the Model 1 is overstated. (g) The  $tx$  variable is not significant in Model 2, but the word count is. This indicates that participation by itself is not important, and it is the amount of cognitive elaboration that affects future spending.

(h) (1) The plot indicates that the residuals are not normal.  
 (2) The analysis is unaffected because OLS does not assume normality and the CLT implies the sampling distributions of the slopes will be normal. (i) The positive association between word count and future spending suggest that they should do things to encourage more cognitive elaboration in future contests. Students do not need to say this, but regression establishes that there is a relationship, but not causation; there could be other factors that cause someone to elaborate more and increase their spending.

```
part a
Call: lm(log(y + 1) ~ log(x + 1) + tx)
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.31705 0.04155 -7.631 2.47e-14 ***
log(x + 1) 0.80318 0.01205 66.657 < 2e-16 ***
tx 0.24438 0.02845 8.591 < 2e-16 ***

Residual standard error: 1.693 on 14175 df
Multiple R-squared: 0.2406, Adjusted R-squared: 0.2405
F-statistic: 2246 on 2 and 14175 DF, p-value: < 2.2e-16

part b
Call: lm(log(y + 1) ~ log(x + 1) + tx + log(wc + 1))
```

```
Coefficients:
 Estimate Std. Error t value Pr(>|t|)
(Intercept) -0.30823 0.04166 -7.398 1.46e-13 ***
log(x + 1) 0.80026 0.01209 66.168 < 2e-16 ***
tx 0.05039 0.07657 0.658 0.51053
log(wc + 1) 0.07382 0.02706 2.729 0.00637 **

Residual standard error: 1.693 on 14174 df
Multiple R-squared: 0.241, Adjusted R-squared: 0.2409
F-statistic: 1500 on 3 and 14174 DF, p-value: < 2.2e-16
```

3. JWHT problem 3.14. (a)  $\beta_0 = \beta_1 = 2$ ,  $\beta_2 = 0.3$ ,  $\sigma = 1$ ,  $y = 2 + 2x_1 + 0.3x_2 + e$ . (b)  $\text{cor}(x_1, x_2) = 0.83512$ . (c) Barely reject  $H_0 : \beta_1 = 0$  ( $P = .0487$ ), but not  $H_0 : \beta_2 = 0$  ( $P = .3754$ ). (d) Reject. (e) Reject. (f) This illustrates the omitted variable bias. The predictors are highly correlated and can serve as proxies for one another. The  $\beta_2$  coefficient is not significant in part c because of high standard errors, but is significant in part e because it explains some of the variation due to  $x_1$ . (g)  $\beta_1 = 2$  is in both intervals, but  $\beta_2 = .3$  is only the interval for the first one. This is because of the omitted variable bias. The estimate of  $\beta_1$  is also biased, but not enough to cause the interval not to cover the parameter.

## Measuring model performance and variable importance

---

- The *full model* is  $y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \epsilon$ .
- The variation left unexplained by the full model is given by the *residual sum of squares* or *deviance*:

$$\text{SSE} = \sum_{i=1}^n \hat{\epsilon}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- When  $\beta_1 = \cdots = \beta_p = 0$  we call it the *intercept* or *null* model, and it turns out  $\hat{y}_i = \bar{y}$ , the mean of  $y$ . The variation left unexplained by the null model is the *total sum of squares*:

$$\text{SST} = \sum_{i=1}^n (y_i - \bar{y})^2 = (n - 1)S_y^2$$

- Usually the full model explains more than the null model

```
> deviance(lm(sales~1, newfood)) # null model leaves 184K unexplained
[1] 184066

> deviance(lm(sales~ad, newfood)) # model with ads leaves 182K unexplained
[1] 181544.5

> deviance(lm(sales~ad+volume, newfood)) # model with both leaves 87K unexplained
[1] 87246.89
```

- We can say that **ad** explains  $184066 - 181544.5 = 2,521.5$ .
- **ad** and **volume** together explain  $184066 - 87247 = 96,819$ .

## The anova and drop1 commands

```
> attach(newfood)
> var(sales)*(nrow(newfood)-1) # TSS
[1] 184066
> sum((sales-mean(sales))^2) # TSS
[1] 184066

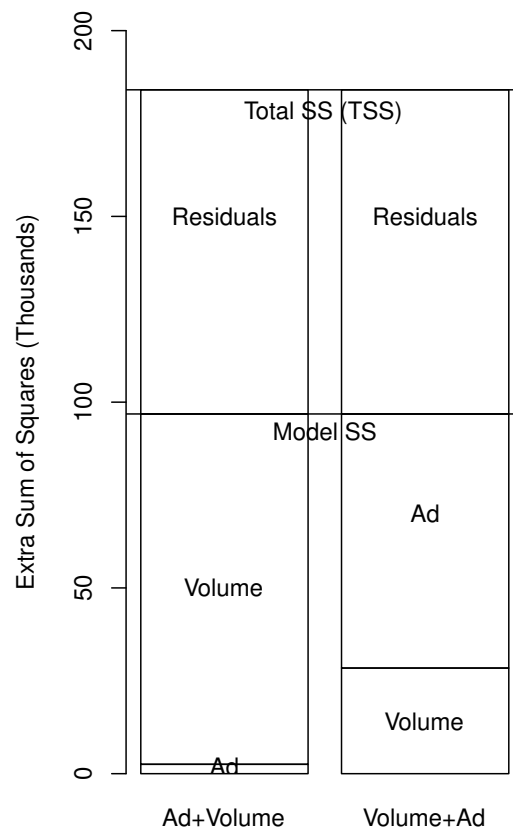
> fit = lm(sales ~ ad + volume, newfood)
> anova(fit) # extra SS
Analysis of Variance Table
Response: sales
 Df Sum Sq Mean Sq F value Pr(>F)
ad 1 2521 2521 0.6069 0.4446
volume 1 94298 94298 22.6971 0.0001
Residuals 21 87247 4155

> drop1(fit) # partial SS
Single term deletions
Model: sales ~ ad + volume
 Df Sum of Sq RSS
<none> 87247
ad 1 68391 155638
volume 1 94298 181544

> fit2 = lm(sales ~ volume+ad, newfood)

> anova(fit2) # extra SS
Analysis of Variance Table
Response: sales
 Df Sum Sq Mean Sq F value Pr(>F)
volume 1 28428 28428 6.8426 0.0161
ad 1 68391 68391 16.4614 0.0006
Residuals 21 87247 4155

> drop1(fit2) # partial SS
Single term deletions
Model: sales ~ volume + ad
 Df Sum of Sq RSS
<none> 87247
volume 1 94298 181544
ad 1 68391 155638
```



- TSS = RSS for intercept model = sum of extra SS
- Extra SS (**anova**): change in SS if term added
- Partial SS (**drop1**): change in SS if term dropped
- For last term in, Extra SS = Partial SS

## The $F$ test of “overall significance”

---

- Recall how to test the following hypothesis:

$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_p = 0$$

$$H_1 : \text{at least one } \beta_j \neq 0$$

```
> summary(fit)
Call: lm(formula = sales ~ ad + volume, data = newfood)

 Estimate Std. Error t value Pr(>|t|)
(Intercept) -324.31 116.95 -2.773 0.011396 *
ad 159.25 39.25 4.057 0.000567 ***
volume 14.87 3.12 4.764 0.000105 ***

Residual standard error: 64.46 on 21 degrees of freedom
Multiple R-squared: 0.526, Adjusted R-squared: 0.4809
F-statistic: 11.65 on 2 and 21 DF, p-value: 0.0003941
```

- This test compares the full model ( $H_1$ ) with the null model

$$F = \frac{\frac{SST - SSE}{p}}{\frac{SSE}{n - p - 1}} = \frac{\frac{184066 - 87246.89}{2}}{\frac{87246.89}{21}} = 11.652 = \left( \frac{\frac{\Delta SSE}{\Delta df}}{S_e^2} \right)$$

- The  $P$ -value can be found in R:

```
> 1 - pf(11.652, 2, 21)
[1] 0.0003941411
```

- This foreshadows an important application of the  $F$  test

## The $F$ test for a single predictor

---

- Now consider testing  $H_0 : \beta_2 = 0$  versus  $H_1 : \beta_2 \neq 0$
- Assuming the null is true ( $\beta_2 = 0$ ) we get  $y = \beta_0 + \beta_1 x_1 + e$ , which will be called the *reduced* model
- The **summary** output give a  $t$  test and shows  $P = .000105$ .
- We can equivalently perform an  $F$  test, which will have more general uses later in the course:

```
> anova(fit)
Analysis of Variance Table
Response: sales
 Df Sum Sq Mean Sq F value Pr(>F)
ad 1 2521 2521 0.6069 0.4446424
volume 1 94298 94298 22.6971 0.0001048 ***
Residuals 21 87247 4155

> drop1(fit, test="F")
Single term deletions
Model: sales ~ ad + volume
 Df Sum of Sq RSS F value Pr(>F)
<none> 87247
ad 1 68391 155638 16.461 0.0005666 ***
volume 1 94298 181544 22.697 0.0001048 ***
```

- The following  $F$  has 1, 21 df

$$F = \frac{\frac{\Delta \text{SSE}}{\Delta df}}{S_e^2} = \frac{\frac{94298}{1}}{\frac{87247}{21}} = \frac{94298}{4155} = 22.6971$$

```
> 1-pf(22.6971, 1, 21)
[1] 0.0001047984
```

- Why is the **ad** not significant in the **anova** output?

## Summary of key points

---

- Model selection involves picking a model between the null and full model
- We use SSE to measure what is unexplained by a model
- The **anova** command generates *extra sum of squares*, telling how much SSE is reduced as we add terms to a model one at a time. They depend on the order of the terms.
- The **drop1** command generates *partial sum of squares*, telling how much SSE is increased when we drop each term. They do not depend on order.
- The  $F$  test allows for hypothesis testing between the full and reduced models

# Your Turn

---

1. Consider the click ballpoint pens data given on page 8.
  - (a) How much variation is left unexplained by the intercept model? (this will be called the *null deviance*)
  - (b) How much variation is explained by adding **ad** to the intercept model?
  - (c) How much additional variation is explained by adding **reps** to a model that already has **ad** in it?
  - (d) How much additional variation is explained by adding **eff** to a model that already has **ad** and **reps** in it?
  - (e) How much variation is unexplained by a model having all three predictors?
  - (f) How much less variation is explained if we drop **ad** from a model with all three predictors in it?
  - (g) Compute  $R^2$  for the three-predictor model “by hand” using only the numbers you have found above. Confirm your answer by having R compute it.
  - (h) Compute adjusted  $R^2$  by hand and confirm it.
  - (i) Compute the  $F$  statistic for the overall test of significance by hand.
  - (j) Compute the  $F$  statistic to test  $H_0 : \beta_1 = 0$  by hand.
2. Consider the commercial properties data.
  - (a) Obtain the analysis of variance table that decomposes the regression sum of squares into extra sums of squares associated with  $X_4$ ; with  $X_1$  given  $X_4$ ; with  $X_2$  given  $X_1$  and  $X_4$ ; and with  $X_3$  given  $X_1, X_2$  and  $X_4$ . Hint use the **lm** and **anova** functions.
  - (b) Test whether  $X_3$  can be dropped from the regression model given that  $X_1, X_2$  and  $X_4$  are retained. Use the  $F$  test statistic and level of significance of .01. State the null and alternative hypotheses, test statistic,  $P$ -value and decision.
  - (c) Test whether both  $X_2$  and  $X_3$  can be dropped from the regression model given that  $X_1$  and  $X_4$  are retained; use  $\alpha = 0.01$ . State the null and alternative hypotheses, test statistic,  $P$ -value and decision. Hint: use the **pf** function to find  $P$  values.
  - (d) Find the variance inflation factors for the full model with all four predictors in the model. What do they tell you?
3. Consider the brand problem.
  - (a) Find the variance inflation factors for the full model with both predictors in the model. What do they tell you?
  - (b) Obtain the analysis of variance table that decomposes the regression sum of squares into extra sums of squares associated with **moisture**; and with **sweetness** given **moisture**?



- (c) Obtain the analysis of variance table that decomposes the regression sum of squares into extra sums of squares associated with **sweetness**; and **moisture** given **sweetness**. What do you notice?
  - (d) Regress liking on moisture content only. How does the estimate of  $\beta_1$  in the previous part compare with the estimate in the model with both predictors?
4. Consider the quality control data set discussed in class.
- (a) How much variation is left unexplained by the intercept model? (this will be called the *null deviance*)
  - (b) How much variation is explained by adding **rate** to the intercept model?
  - (c) How much additional variation is explained by adding **am** to a model that already has **rate** in it?
  - (d) How much variation is unexplained by a model having both predictors?
  - (e) How much less variation is explained if we drop **rate** from a model with both predictors in it?
  - (f) Compute  $R^2$  for the two-predictor model “by hand” using only the numbers you have found above. Confirm your answer by having R compute it.
  - (g) Compute the  $F$  statistic for the overall test of significance by hand.
  - (h) Using the two-variable model, compute the  $F$  statistic to test  $H_0 : \beta_1 = 0$  by hand (where  $\beta_1$  is for rate) (hint: it is in the **drop1** output).
5. Return to the model you estimated for the Auto problem of the Your Turn on page 67. Fit this model,

```
auto = read.table("auto.txt", header=T) # auto.txt on Canvas
auto$origin = factor(auto$origin, 1:3, c("US", "Europe", "Japan"))
fit = lm(log(mpg)~ log(weight)+year+I(year^2)+origin, auto)
```

- (a) Interpret the effect of origin on  $\log(\text{mpg})$ . Which origin has the best gas mileage? Worst? Rank them in order of gas mileage from least to greatest.
- (b) After controlling for the other variables, what is the difference in  $\log(\text{mpg})$  between Japan and Europe, on the average?
- (c) Is there a significant difference between the log gas mileage for US and Japan after controlling for the other variables?
- (d) You should see that there are two dummy variables for the origin variable. If origin were dropped from the model (i.e., the two dummies were set equal to 0), by how much would RSS increase?
- (e) Can you reject the null hypothesis that both origin dummies are 0, so that none of the origin levels have different effects?
- (f) As an extra challenge, perform an  $F$  test for whether **year** and its quadratic effect can both be dropped from the model, i.e., test  $H_0 : \beta_2 = \beta_3 = 0$ . You can do this by fitting the reduced model (call it **fit2**) then use the **anova(fit2, fit)** command or include **poly(year, 2, raw=T)** and use **drop1**.

## Answers

1. (a) SST = 598253; (b) 463451, Hint: `fit = lm(sales ~ ad + reps + eff, click)` and then `anova(fit0)`; (c) 59327; (d) 4431; (e) 71044; (f) 44295, Hint: `drop1(fit, test="F")`; (g)  $1 - 71044/598253 = .8812$ , Hint: `summary(fit)`; (h)  $1 - (71044/36)/(598253/39) = .8714$ ; (i)  $((598253 - 71044)/3)/(71044/36) = 89.05$ ; (j)  $(44295/1)/(71044/36) = 22.45$ , Hint: see `drop1` output.

2. (a) Hint: `fit = lm(y ~ x4+x1+x2+x3, comm)` then `anova(fit)`; (b)  $H_0 : \beta_3 = 0$  versus  $H_1 : \beta_3 \neq 0$ ,  $P = 0.5704$ , we cannot reject  $H_0$ . (c)  $H_0 : \beta_2 = \beta_3 = 0$  versus  $H_1 : \beta_2 \neq 0$  or  $\beta_3 \neq 0$ . From the output below the  $P$ -value is less than 0 and so we reject  $H_0$ .

```
> fit2 = lm(y ~ x1+x4, comm)
> 1-pf(((deviance(fit2)-deviance(fit))/2) / (deviance(fit)/76), 2, 76)
[1] 6.682136e-05
```

(d) Hint: `vif(fit)`. The VIF values are 1.41, 1.24, 1.65, and 1.32. All are fairly close to 1 indicating that multicollinearity is not a serious problem.

3. Brand problem. (a) The VIFs are both 1 indicating uncorrelated predictors; (b) The extra SS for moisture are 1566 and for sweetness 306; (c) They are the same as in the previous part; (d) The coefficient is the same.
4. Quality control problem. (a) 10929.29. (b) 8566.9. (c) 7.1. (d) 2355.3. (e) 5440.5. (f)  $1 - 2355.3/10929.29 = 0.7844965$ . (g)  $((10929.29 - 2355.3)/2)/(2355.3/27) = 49.14$ . (h)  $(5440.5/1)/(2355.3/27) = 62.37$ .
5. Auto problem. (a) The base is US. Europe had 0.0668 higher log(mpg) than the US on average, and Japan had 0.032 higher log(mpg) than the US, on average. Europe had the best gas milage, followed by Japan, then the US. (b)  $0.06683 - 0.03197 = 0.03486$ . (c)  $H_0 : \beta_5 = 0$  versus  $H_1 : \beta_5 \neq 0$ ,  $P = 0.075 > .05$  so we cannot reject  $H_0$ . (d) See `drop1`: 0.1857. (e)  $P = .00085 < .05$ , so we can reject  $H_0 : \beta_4 = \beta_5 = 0$ . (f)

```
fit = lm(log(mpg)~ log(weight)+year+I(year^2)+origin, auto)
summary(fit); drop1(fit)
fit2 = lm(log(mpg)~ log(weight)+origin, auto)
anova(fit2, fit)
fit3 = lm(log(mpg)~ log(weight)+poly(year,2, raw=T) +origin, auto)
summary(fit3); drop1(fit3, test="F")
```

# Dummy Variables

---

- Question: How do we include nominal variables in a regression?
- Answer: Use dummy (also called indicator) variables
- Example: Quality Case. Let AM be a *dummy* variable that takes the value 1 if the observation comes from the morning shift and 0 otherwise (afternoon shift). As a first step we could regress defects on AM.

```
> fit = lm(defect~ am, quality)
> summary(fit)
...
 Estimate Std. Error t value Pr(>|t|)
(Intercept) 16.920 4.308 3.927 0.00051 ***
am 20.440 6.093 3.355 0.00229 **
```

Interpretation of 20.44: every unit increase in AM (i.e., going from PM to AM) is associated with a 20.44 change in defect rate. The AM shift has 20.44 more defects per thousand than the PM shift.

- This is equivalent to an independent-sample  $t$ -test with equal variances assumed:

```
> t.test(defect~am, quality, var.equal=T)
data: defect by am
t = -3.3547, df = 28, p-value = 0.002295
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -32.920676 -7.959324
sample estimates:
mean in group 0 mean in group 1
 16.92 37.36
```

The mean difference is  $37.36 - 16.92 = 20.44$

- We reject  $H_0 : \beta = 0$ , or equivalently,  $H_0 : \mu_{AM} = \mu_{PM}$  because  $0.0023 < 0.05$ .

## Dummy Variables: Wholesaler Efficiency

---

| If the Wholesaler is | Dummy Variable Coding |      |             |
|----------------------|-----------------------|------|-------------|
|                      | Fair                  | Good | Outstanding |
| Fair                 | 1                     | 0    | 0           |
| Good                 | 0                     | 1    | 0           |
| Outstanding          | 0                     | 0    | 1           |
| Poor                 | 0                     | 0    | 0           |

```
> fit = lm(sales~ad+reps+as.factor(eff), data=click)
> drop1(fit, test="F")
 Df Sum of Sq RSS AIC F value Pr(>F)
<none> 71018 311.27
ad 1 41227 112245 327.58 19.7376 8.955e-05 ***
reps 1 54607 125625 332.09 26.1433 1.226e-05 ***
as.factor(eff) 3 4457 75475 307.71 0.7112 0.552

> summary(fit)
 Estimate Std. Error t value Pr(>|t|)
(Intercept) 45.051 36.631 1.230 0.227
ad 13.063 2.940 4.443 8.96e-05 ***
reps 40.948 8.009 5.113 1.23e-05 ***
as.factor(eff)2 9.239 27.916 0.331 0.743
as.factor(eff)3 20.283 29.344 0.691 0.494
as.factor(eff)4 33.260 28.440 1.169 0.250
```

$$\hat{y} = 45 + 13\text{ad} + 41\text{reps} + 9\text{fair} + 20\text{good} + 33\text{out}$$

- First test  $H_0$  : all  $\beta_j = 0$  using **drop1** in R.
- **as.factor(eff)** indicates that **eff** should be treated as categorical (i.e., create dummies).
- What null hypothesis does the  $P$ -value for OUT test ( $P = .2503$ )? What does this test mean in English?

## SPSS Estimates

Note: In practice first look at the ANOVA table (next page).

Parameter Estimates

| Parameter | B              | Std Error | t      | Sig   |
|-----------|----------------|-----------|--------|-------|
| Intercept | 78.311         | 26.993    | 2.901  | 0.006 |
| AD        | 13.063         | 2.940     | 4.443  | 0.000 |
| REPS      | 40.948         | 8.009     | 5.113  | 0.000 |
| [EFF=1]   | -33.260        | 28.440    | -1.169 | 0.250 |
| [EFF=2]   | -24.020        | 19.339    | -1.242 | 0.223 |
| [EFF=3]   | -12.976        | 18.643    | -0.696 | 0.491 |
| [EFF=4]   | 0 <sup>a</sup> | .         | .      | .     |

a. This parameter is set to zero because it is redundant.

- Estimated regression equation:  $\hat{y} = 78.3 + 13\text{ad} + 41\text{reps} - 33\text{Poor} - 24\text{Fair} - 13\text{Good} + 0\text{Out}$
- Why are estimates different than on page 92? Note that **ad** and **reps** are identical.
- What null hypothesis does the  $P$ -value for **Poor** test ( $P = .250$ )? What does this mean in English?
- Note that SPSS/Minitab/SAS do not give standardized regression coefficients for dummies. Why?

## General Linear Test for Multiple Betas

---

```
> drop1(fit, test="F")
sales ~ ad + reps + as.factor(eff)
 Df Sum of Sq RSS AIC F value Pr(>F)
<none> 71018 311
ad 1 41227 112245 328 19.7376 8.955e-05 ***
reps 1 54607 125625 332 26.1433 1.226e-05 ***
as.factor(eff) 3 4457 75475 308 0.7112 0.552

> deviance(fit)
[1] 71017.78

> anova(fit)
 Df Sum Sq Mean Sq F value Pr(>F)
ad 1 463451 463451 221.8787 < 2.2e-16 ***
reps 1 59327 59327 28.4032 6.414e-06 ***
as.factor(eff) 3 4457 1486 0.7112 0.552
Residuals 34 71018 2089
```

- The “ad” line tests  $H_0 : \beta_1 = 0$  and is equivalent to the  $t$  test on the previous page. Likewise for the “reps” line.
- The “eff” line tests  $H_0 : \beta_3 = \beta_4 = \beta_5 = 0$  (i.e., the three dummies for **eff** are all 0 meaning all levels of wholesaler efficiency are the same) versus  $H_1$  : at least one of  $\beta_3$ ,  $\beta_4$ , or  $\beta_5$  is different from 0.
- The “< none >” line gives RSS for the full model

# How to Handle Missing Values

---

```
agemiss = data.frame(
 age = c(NA,NA,35,NA,81,39,20,25,62,NA,45,57,36,39,NA,48,36,NA,NA,30,
 78,35,NA,20,26,28,44,30,31,32,72,33,33,NA,55,37,36,43,40,NA),
 y = c(2.9,2.8,8.4,2.8,4.5,8.3,9.4,9.1,5.6,2.9,7.4,6.3,7.7,8.1,3.2,6.5,
 7.9,3.0,3.0,9.0,5.1,8.8,3.4,9.5,8.9,8.3,7.4,8.3,8.6,8.7,5.3,8.3,
 7.8,3.2,6.6,8.4,8.6,7.8,7.6,3.7))
```

Solution: treat missing as a separate category and include a dummy:

1. Create dummy `xmiss` that equals 1 when `x` is missing and 0 otherwise.
2. When `x` is missing, set `x=0` (or impute with a regression)
3. Regress `y` on both `x` and `xmiss`

```
> agemiss$xmiss=is.na(agemiss$age) # 1. create dummy
> agemiss$age[is.na(agemiss$age)] = 0 # 2. set missings to 0
> plot(agemiss$age, agemiss$y)
> fit = lm(y ~ age + xmiss, agemiss) # 3. regression
> summary(fit)
```

Coefficients:

```
 Estimate Std. Error t value Pr(>|t|)
(Intercept) 11.040189 0.163265 67.62 <2e-16 ***
age -0.080755 0.003737 -21.61 <2e-16 ***
xmissTRUE -7.950189 0.191438 -41.53 <2e-16 ***
Residual standard error: 0.3161 on 37 degrees of freedom
```

- When age is present,  $y = 11.04 - 0.08\text{age} - 7.95(0)$ .
- When age is missing,  $y = 11.04 - 0.08(0) - 7.95(1)$