

UCLA Department of Statistics
Statistical Consulting Center

Regression in R
Part I :
Simple Linear Regression

Denise Ferrari & Tiffany Head
denise@stat.ucla.edu tiffany@stat.ucla.edu

Feb 10, 2010





Initial Data Analysis I

Does the data look like as we expect?

Prior to any analysis, the data should always be inspected for:

- Data-entry errors
- Missing values
- Outliers
- Unusual (e.g. asymmetric) distributions
- Changes in variability
- Clustering
- Non-linear bivariate relationships
- Unexpected patterns



Initial Data Analysis II

Does the data look like as we expect?

We can resort to:

- Numerical summaries:
 - 5-number summaries
 - correlations
 - etc.
- Graphical summaries:
 - boxplots
 - histograms
 - scatterplots
 - etc.

Coding Missing Data II

Example: Diabetes in Pima Indian Women

R code for missing data

- Zero should **not** be used to represent missing data
 - it's a valid value for some of the variables
 - can yield misleading results
- Set the missing values coded as zero to NA:

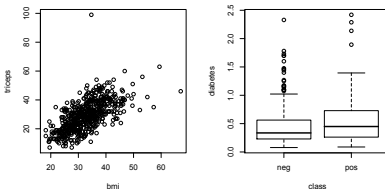
```
1  pima$glucose[pima$glucose==0] <- NA
2  pima$bp[pima$bp==0] <- NA
3  pima$triceps[pima$triceps==0] <- NA
4  pima$insulin[pima$insulin==0] <- NA
5  pima$bmi[pima$bmi==0] <- NA
```


Graphical Summaries

Example: Diabetes in Pima Indian Women

- Bivariate

```
1 # scatterplot
2 plot(triceps~bmi, pima)
3 # boxplot
4 boxplot(diabetes~class, pima)
```



Estimation of unknown parameters I

We want to find the equation of the line that “best” fits the data. It means finding b_0 and b_1 such that the **fitted values** of y_i , given by

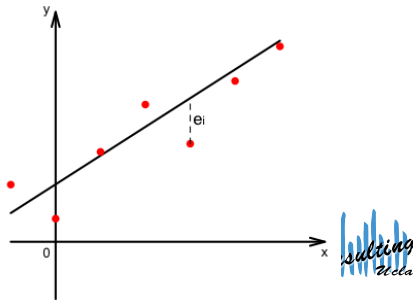
$$\hat{y}_i = b_0 + b_1 x_i,$$

are as “close” as possible to the **observed values** y_i .

Residuals

The difference between the observed value y_i and the fitted value \hat{y}_i is called **residual** and is given by:

$$e_i = y_i - \hat{y}_i$$



Fitted values and residuals

- Fitted values obtained using the function `fitted()`
- Residuals obtained using the function `resid()`

```
1 # Create a table with fitted values and
  residuals
2 data.frame(production, fitted.value=fitted(
  production.lm), residual=resid(production.lm)
  )
```

Case	RunTime	RunSize	fitted.value	residual
1	195	175	195.1152	-0.1152469
2	215	189	198.7447	16.2553496
3	243	344	238.9273	4.0726679
...				
20	172	68	167.3762	4.6237657

$$\hat{y}_1 = 149.75 + 0.26 * 175 = 195.115$$

$$e_1 = 195 - 195.115 = -0.115$$



Fitted values and residuals I

When there are missing data

Missing data need to be handled carefully. Using the `na.exclude` method:

```

1 # Load the package that contains the data
2 library(ISwR)
3 data(thuesen); attach(thuesen)
4 # Option for dealing with missing data
5 options(na.action=na.exclude)
6 # Now fit the regression model as before
7 velocity.lm <- lm(short.velocity~blood.glucose
8 )
9 # Create a table with fitted values and
  residuals
10 data.frame(thuesen, fitted.value=fitted(
  velocity.lm), residual=resid(velocity.lm)

```



Measuring Goodness of Fit I

Coefficient of Determination, R^2

- represents the proportion of the total sample variability explained by the regression model.
- for simple linear regression, the R^2 statistic corresponds to the square of the correlation between Y and X .
- indicates of how well the model fits the data.

From the ANOVA table:

$$R^2 = \frac{12868.4}{(12868.4 + 4754.6)} = 0.7302$$

which we can also find in the regression summary.



Measuring Goodness of Fit II

Adjusted R^2

The adjusted R^2 takes into account the number of degrees of freedom and is preferable to R^2 .

From the ANOVA table:

$$R_{adj}^2 = 1 - \frac{4754.6/18}{(12868.4 + 4754.6)/(18 + 1)} = 0.7152$$

also found in the regression summary.

Attention

Neither R^2 nor R_{adj}^2 give direct indication on how well the model will perform in the prediction of a new observation.

Confidence and prediction bands I

Confidence Bands

Reflect the uncertainty about the regression line (how well the line is determined).

Prediction Bands

Include also the uncertainty about future observations.

Attention

These limits rely strongly on the assumption of normally distributed errors with constant variance and should **not** be used if this assumption is violated for the data being analyzed.

Confidence and prediction bands II

Predicted values are obtained using the function `predict()` .

```
1 # Obtaining the confidence bands:
2 predict(production.lm, interval="confidence")
```

```
          fit      lwr      upr
1  195.1152 187.2000 203.0305
2  198.7447 191.0450 206.4443
3  238.9273 225.4549 252.3998
...
20 167.3762 154.4448 180.3077
```



Confidence and prediction bands III

```

1 # Obtaining the prediction bands:
2 predict(production.lm, interval="prediction")

```

	fit	lwr	upr
1	195.1152	160.0646	230.1659
2	198.7447	163.7421	233.7472
3	238.9273	202.2204	275.6343
...			
20	167.3762	130.8644	203.8881



Confidence and prediction bands IV

For plotting:

```
1 # Create a new data frame containing the  
   values of X at which we want the  
   predictions to be made  
2 pred.frame <- data.frame(RunSize= seq(55, 345,  
   by=10))  
3 # Confidence bands  
4 pc <- predict(production.lm, int="c", newdata=  
   pred.frame)  
5 # Prediction bands  
6 pp <- predict(production.lm, int="p", newdata=  
   pred.frame)  
7  
8
```



Confidence and prediction bands V

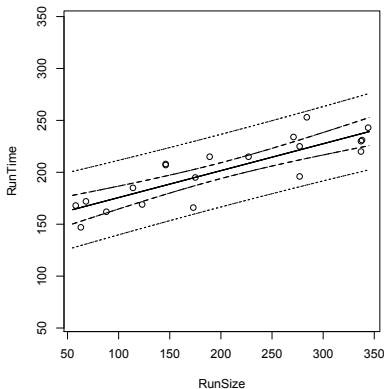
```

9 # Plot
10 require(graphics)
11 # Standard scatterplot with extended limits
12 plot(RunSize, RunTime, ylim=range(RunSize, pp,
    na.rm=T))
13 pred.Size <- pred.frame$RunSize
14 # Add curves
15 matlines(pred.Size, pc, lty=c(1,2,2), lwd=1.5,
    col=1)
16 matlines(pred.Size, pp, lty=c(1,3,3), lwd=1.5,
    col=1)

```



Confidence and prediction bands VI



Dummy Variable Regression

The simple dummy variable regression is used when the **predictor** variable is not quantitative but **categorical** and assumes only two values.



Dummy Variable Regression I

Example: Change over time (Taken from Sheather, 2009)

Loading the Data:

```
1 changeover <- read.table("http://www.stat.tamu
  .edu/~sheather/book/docs/datasets
  /changeover_times.txt", header=T, sep=" ")
```

Variables:

	Method	Changeover	New
1	Existing	19	0
2	Existing	24	0
3	Existing	39	0
...			
118	New	14	1
119	New	40	1
120	New	35	1

Change-over(Y): time (in minutes) required to change the line of food production

New (X): 1 for the new method, 0 for the existing method

We want to be able to test whether the change-over time is different for the two methods.



Dummy Variable Regression II

Example: Change over time (Taken from Sheather, 2009)

```
1 attach(changeover)
2 # Summary:
3 summary(changeover)
```

Method	Changeover	New
Existing:72	Min. : 5.00	Min. :0.0
New :48	1st Qu.:11.00	1st Qu.:0.0
	Median :15.00	Median :0.0
	Mean :16.59	Mean :0.4
	3rd Qu.:21.00	3rd Qu.:1.0
	Max. :40.00	Max. :1.0

We need to recode the X variable (New) to factor :



Dummy Variable Regression III

Example: Change over time (Taken from Sheather, 2009)

```
1 changeover$New <- factor(changeover$New)
2 summary(changeover)
```

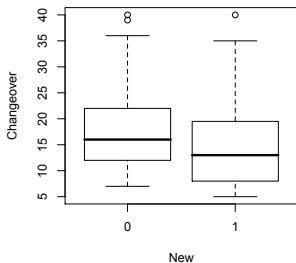
Method	Changeover	New
Existing:72	Min. : 5.00	0:72
New :48	1st Qu.:11.00	1:48
	Median :15.00	
	Mean :16.59	
	3rd Qu.:21.00	
	Max. :40.00	

Dummy Variable Regression IV

Example: Change over time (Taken from Sheather, 2009)

Plotting the data:

```
1 plot(Changeover~New)
```



Dummy Variable Regression V

Example: Change over time (Taken from Sheather, 2009)

Fitting the linear regression:

```
1 # Fit the linear regression model
2 changeover.lm <- lm(Changeover ~ New, data =
  changeover)
3 # Extract the regression results
4 summary(changeover.lm)
```

Dummy Variable Regression VI

Example: Change over time (Taken from Sheather, 2009)

The output looks like this:

Call:

```
lm(formula = Changeover ~ New, data = changeover)
```

Residuals:

Min	1Q	Median	3Q	Max
-10.861	-5.861	-1.861	4.312	25.312

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	17.8611	0.8905	20.058	<2e-16 ***
New1	-3.1736	1.4080	-2.254	0.0260 *

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 7.556 on 118 degrees of freedom

Multiple R-squared: 0.04128, Adjusted R-squared: 0.03315

F-statistic: 5.081 on 1 and 118 DF, p-value: 0.02604



Dummy Variable Regression VII

Example: Change over time (Taken from Sheather, 2009)

Analysis of the results:

- There's significant evidence of a reduction in the mean change-over time for the new method.
- The estimated mean change-over time for the new method ($X = 1$) is:

$$\hat{y}_1 = 17.8611 + (-3.1736) * 1 = 14.7 \text{ minutes}$$

- The estimated mean change-over time for the existing method ($X = 0$) is:

$$\hat{y}_0 = 17.8611 + (-3.1736) * 0 = 17.9 \text{ minutes}$$



Diagnostics

Assumptions

The assumptions for simple linear regression are:

- Y relates to X by a linear regression model:

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

- the errors are independent and identically normally distributed with mean zero and common variance



Diagnostics

What can go wrong?

Violations:

- In the linear regression model:
 - linearity (e.g. quadratic relationship or higher order terms)
- In the residual assumptions:
 - non-normal distribution
 - non-constant variances
 - dependence
 - outliers

Checks:

- ⇒ look at plot of residuals vs. X
- ⇒ look at plot of residuals vs. fitted values
- ⇒ look at residuals Q-Q norm plot



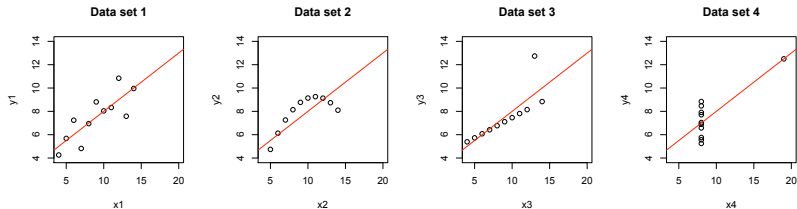
Validity of the regression model II

Example: The Anscombe's data sets (Taken from Sheather, 2009)

```
1  # Fitting the regressions
2  a1.lm <- lm(y1~x1, data=anscombe)
3  a2.lm <- lm(y2~x2, data=anscombe)
4  a3.lm <- lm(y3~x3, data=anscombe)
5  a4.lm <- lm(y4~x4, data=anscombe)
6
7  #Plotting
8  # For the first data set
9  plot(y1~x1, data=anscombe)
10 abline(a1.lm, col=2)
```

Validity of the regression model III

Example: The Anscombe's data sets (Taken from Sheather, 2009)



For all data sets, the fitted regression is the same:

$$\hat{y} = 3.0 + 0.5x$$

All models have $R^2 = 0.67$, $\hat{\sigma} = 1.24$ and the slope coefficients are significant at $< 1\%$ level. To check that, use the `summary()` function on the regression models.

Leverage (or influential) points and outliers I

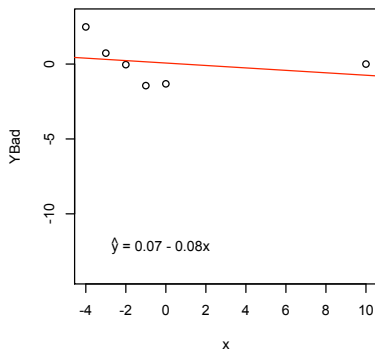
Leverage points

Leverage points are those which have great influence on the fitted model, that is, those whose x -value is distant from the other x -values.

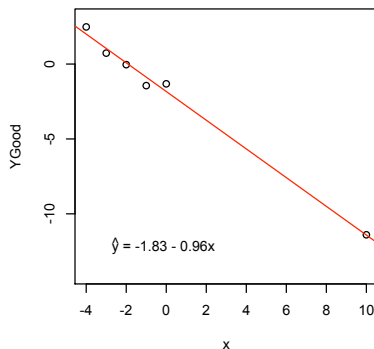
- Bad leverage point: if it is also an **outlier**, that is, the y -value does not follow the pattern set by the other data points.
- Good leverage point: if it is **not** an outlier.

Leverage (or influential) points and outliers II

Bad Leverage Point



Good Leverage Point





Standardized residuals I

Standardized residuals are obtained by dividing each residual by an estimate of its standard deviation:

$$r_i = \frac{e_i}{\hat{\sigma}(e_i)}$$

To obtain the standardized residuals in R, use the command `rstandard()` on the regression model.

Leverage Points

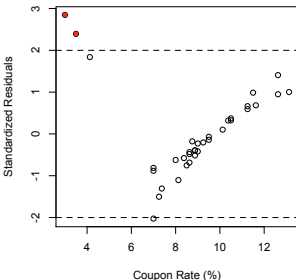
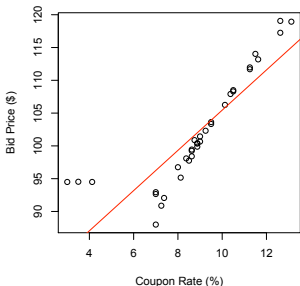
- Good leverage points have their standardized residuals within the interval $[-2, 2]$
- **Outliers** are leverage points whose standardized residuals fall outside the interval $[-2, 2]$



Leverage (or influential) points and outliers II

How to deal with them

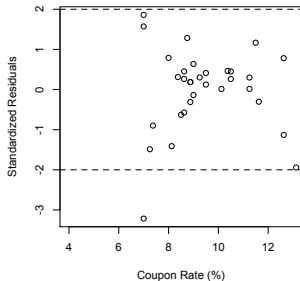
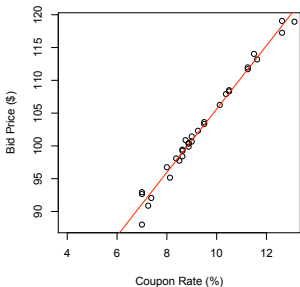
Data set containing outliers:



Leverage (or influential) points and outliers III

How to deal with them

After their removal:



Normality and constant variance of errors

Normality and Constant Variance Assumptions

These assumptions are necessary for inference:

- hypothesis testing
- confidence intervals
- prediction intervals

⇒ Check the Normal Q-Q plot of the standardized residuals.

⇒ Check the Standardized Residuals vs. X plot.

Note

When these assumptions do not hold, we can try to correct the problem using data transformations.

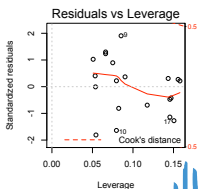
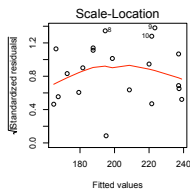
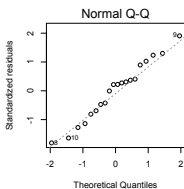
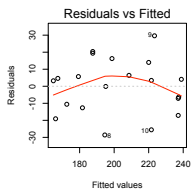
Normality and constant variance checks

Example: Production Runs

```

1 # Regression model
2 production.lm <- lm(RunTime~RunSize, data=
  production)
3 # Residual plots
4 plot(production.lm)

```



Example of correction: non-constant variance I

Example: Cleaning Data (Taken from Sheather, 2009)

Variables:

Rooms (Y): number of rooms cleaned

Crews (X): number crews

We want to be able to model the relationship between the number of rooms cleaned and the number of crews.

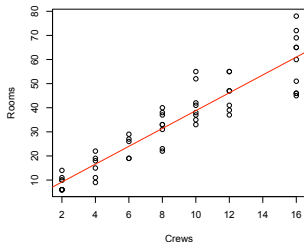
```
1 # Load the data
2 cleaning <- read.table("http://www.stat.tamu.edu/~sheather/book/docs/datasets/cleaning.txt", h=T, sep="")
3 attach(cleaning)
```



Example of correction: non-constant variance II

Example: Cleaning Data (Taken from Sheather, 2009)

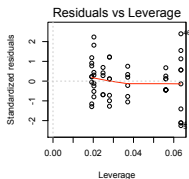
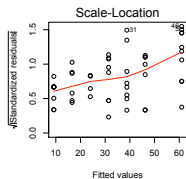
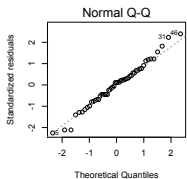
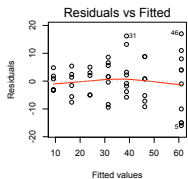
```
1 # Regression model
2 cleaning.lm <- lm(Rooms
  ~ Crews, data=cleaning)
3 # Plotting data and
  regression line
4 plot(Rooms~Crews)
5 abline(cleaning.lm, col=2)
```



Example of correction: non-constant variance III

Example: Cleaning Data (Taken from Sheather, 2009)

- 1 # Diagnostic plots
- 2 `plot(cleaning.lm)`



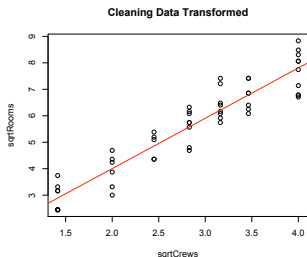
Example of correction: non-constant variance IV

Example: Cleaning Data (Taken from Sheather, 2009)

```

1 # Applying square root
  transformation (counts)
2 sqrtRooms <- sqrt(Rooms)
3 sqrtCrews <- sqrt(Crews)
4 # Regression model on the
  transformed data
5 sqrt.lm <- lm(sqrtRooms
  ~sqrtCrews)
6 # Plotting data and
  regression line
7 plot(sqrtRooms~sqrtCrews)
8 abline(sqrt.lm, col=2)

```



Online Resources for R

Download R: <http://cran.stat.ucla.edu>

Search Engine for R: rseek.org

R Reference Card: <http://cran.r-project.org/doc/contrib/Short-refcard.pdf>

UCLA Statistics Information Portal:

<http://info.stat.ucla.edu/grad/>

UCLA Statistical Consulting Center <http://scc.stat.ucla.edu>



Online Resources for R

- Next Week:
 - Nonlinear Regression (Feb 15, Monday)
- Week After Next:
 - Survival Analysis in R (February 22, Monday)
 - Spatial Statistics in R (February 24, Wednesday)
- For a schedule of all mini-courses offered please visit:
<http://scc.stat.ucla.edu/mini-courses>

Exercise in R I

Airfares Data (Taken from Sheather, 2009)

The data set for this exercise can be found at:

`http://www.stat.tamu.edu/~sheather/book/docs/datasets/airfares.txt`

It contains information on one-way airfare (in US\$) and distance (in miles) from city A to 17 other cities in the US.



Exercise in R II

Airfares Data (Taken from Sheather, 2009)

- 1 Fit the regression model given by:

$$\text{Fare} = \beta_0 + \beta_1 \text{Distance} + \epsilon$$

- 2 Critique the following statement:

The regression coefficient of the predictor variable (Distance) is highly statistically significant and the model explains 99.4% of the variability in the Y-variable (Fare).

Thus this model is highly effective for both understanding the effects of Distance on Fare and for predicting future values of Fare given the value of the predictor variable.

- 3 Does the regression model above seem to fit the data well? If not, describe how the model can be improved.