

Other (non-SL) Regression Models

Hongwei Zhang

<http://www.cs.wayne.edu/~hzhang>



He uses statistics as a drunken man uses lamp-posts — for support rather than for illumination.

--- Andrew Lang

Outline

- Multiple Linear Regression
 - More than one predictor variables
- Categorical Predictors
 - Predictor variables are categories such as CPU type, disk type, and so on
- Curvilinear Regression
 - Relationship is nonlinear
- Transformations
 - Errors are not normally distributed or the variance is not homogeneous
- Outliers
- Common mistakes in regression

Outline

- Multiple Linear Regression
- Categorical Predictors
- Curvilinear Regression
- Transformations
- Outliers
- Common mistakes in regression

Multiple linear regression models

$$y = b_0 + b_1x_1 + b_2x_2 + \cdots + b_kx_k + e$$

- Given a sample of n observations with k predictors

$$\{(x_{11}, x_{21}, \dots, x_{k1}, y_1), \dots, (x_{1n}, x_{2n}, \dots, x_{kn}, y_n)\}$$

$$y_1 = b_0 + b_1x_{11} + b_2x_{21} \cdots b_kx_{k1} + e_1$$

$$y_2 = b_0 + b_1x_{12} + b_2x_{22} \cdots b_kx_{k2} + e_2$$

.

.

.

$$y_n = b_0 + b_1x_{1n} + b_2x_{2n} \cdots b_kx_{kn} + e_n$$

Vector notation

- In vector notation, we have

$$\begin{bmatrix} y_1 \\ y_2 \\ \cdot \\ \cdot \\ \cdot \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{21} & \cdots & x_{k1} \\ 1 & x_{12} & x_{22} & \cdots & x_{k2} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & x_{1n} & x_{2n} & \cdots & x_{kn} \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ \cdot \\ \cdot \\ \cdot \\ b_k \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \cdot \\ \cdot \\ \cdot \\ e_n \end{bmatrix}$$

or

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}$$

- Regression formula:

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})$$

Example 15.1

- Seven programs were monitored to observe their resource demands. In particular, the number of disk I/O's, memory size (in kBytes), and CPU time (in milliseconds) were observed.

CPU Time	Disk I/O's	Memory Size
y_i	x_{1i}	x_{2i}
2	14	70
5	16	75
7	27	144
9	42	190
10	39	210
13	50	235
20	83	400

Example (contd.)

CPU time = $b_0 + b_1(\text{number of disk I/O's}) + b_2(\text{memory size})$

□ In this case:

$$\mathbf{X} = \begin{bmatrix} 1 & 14 & 70 \\ 1 & 16 & 75 \\ 1 & 27 & 144 \\ 1 & 42 & 190 \\ 1 & 39 & 210 \\ 1 & 50 & 235 \\ 1 & 83 & 400 \end{bmatrix}$$

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 7 & 271 & 1324 \\ 271 & 13,855 & 67,188 \\ 1324 & 67,188 & 326,686 \end{bmatrix}$$

Example (contd.)

$$\mathbf{C} = (\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} 0.6297 & 0.0223 & -0.0071 \\ 0.0223 & 0.0280 & -0.0058 \\ -0.0071 & -0.0058 & 0.0012 \end{bmatrix}$$

$$\mathbf{X}^T \mathbf{y} = \begin{bmatrix} 66 \\ 3375 \\ 16,388 \end{bmatrix}$$

- The regression parameters are:

$$\mathbf{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = (-0.1614, 0.1182, 0.0265)^T$$

- The regression equation is:

$$\begin{aligned} \text{CPU time} = & -0.1614 + 0.1182(\text{number of disk I/O's}) + \\ & 0.0265(\text{memory size}) \end{aligned}$$

Example (contd.)

CPU Time	Disk I/O's	Memory Size	Est. CPU time	Error	Error ²
y_i	x_{1i}	x_{2i}	\hat{y}_i	e_i	e_i^2
2	14	70	3.3490	-1.3490	1.8198
5	16	75	3.7180	1.2820	1.6436
7	27	144	6.8472	0.1528	0.0233
9	42	190	9.8400	-0.8400	0.7053
10	39	210	10.0151	-0.0151	0.0002
13	50	235	11.9783	1.0217	1.0439
20	83	400	20.2529	-0.2529	0.0639
Σ	66	271	66.0000	-0.0003	5.3000

□ From the table we see that SSE is:

$$\text{SSE} = \sum e_i^2 = 5.3$$

Example (contd.)

- An alternate method to compute SSE is to use:

$$\text{SSE} = \{\mathbf{y}^T \mathbf{y} - \mathbf{b}^T \mathbf{X}^T \mathbf{y}\}$$

- For this data, SSY and SS0 are:

$$\text{SSY} = \sum y_i^2 = 828$$

$$\text{SS0} = n\bar{y}^2 = 622.29$$

- Therefore, SST and SSR are:

$$\text{SST} = \text{SSY} - \text{SS0} = 828 - 622.29 = 205.71$$

$$\text{SSR} = \text{SST} - \text{SSE} = 205.71 - 5.3 = 200.41$$

Example (contd.)

- The coefficient of determination R^2 is:

$$R^2 = \frac{SSR}{SST} = \frac{200.41}{205.71} = 0.97$$

Thus, the regression explains 97% of the variation of y .

- Standard deviation of errors is:

$$s_e = \sqrt{\frac{SSE}{n-3}} = \sqrt{5.3/4} = 1.2$$

$$C = (X^T X)^{-1}$$

Example (contd.)

- Standard deviations of the regression parameters are:

$$\text{Estimated std. dev. of } b_0 = s_e \sqrt{c_{00}} = 1.2 \sqrt{0.6297} = 0.9131$$

$$\text{Estimated std. dev. of } b_1 = s_e \sqrt{c_{11}} = 1.2 \sqrt{0.0280} = 0.1925$$

$$\text{Estimated std. dev. of } b_2 = s_e \sqrt{c_{22}} = 1.2 \sqrt{0.0012} = 0.0404$$

- The 90% t-value at 4 degrees of freedom is 2.132.

$$90\% \text{ Conf. interval of } b_0 = -0.1614 \mp (2.132)(0.9131) = (-2.11, 1.79)$$

$$90\% \text{ Conf. interval of } b_1 = 0.1182 \mp (2.132)(0.1925) = (-0.29, 0.53)$$

$$90\% \text{ Conf. interval of } b_2 = 0.0265 \mp (2.132)(0.0404) = (-0.06, 0.11)$$

None of the three parameters is significant at a 90% confidence level.

Example (contd.)

- A single future observation for programs with 100 disk I/O's and a memory size of 550:

$$\begin{aligned}y_{1p} &= b_0 + b_1x_1 + b_2x_2 \\ &= -0.1614 + 0.1182(100) + 0.0265(550) = 26.2375\end{aligned}$$

- Standard deviation of the predicted observation is:

$$s_{y_{1p}} = s_e \sqrt{\{1 + \mathbf{x}^T(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{x}\}} = 1.2\sqrt{1 + 7.4118} = 3.3435$$

- 90% confidence interval using the t value of 2.132 is:

$$26.2375 \mp (2.132)(3.3435) = (19.1096, 33.3363)$$

Example (contd.)

- Standard deviation for a mean of a large number of future observations is:

$$s_{\hat{y}_p} = s_e \sqrt{\{\mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}\}} = 1.2 \sqrt{7.4118} = 3.1385$$

- 90% confidence interval is:

$$26.2375 \mp (2.132)(3.1385) = (19.5467, 32.9292)$$

Analysis of variance (ANOVA)

- Test the hypothesis that SSR is less than or equal to SSE

$$SST = SSY - SS0 = SSR + SSE$$

- Degrees of freedom = Number of independent values required to compute

$$\begin{array}{ccccccc} SST & = & SSY & - & SS0 & = & SSR & + & SSE \\ n - 1 & = & n & - & 1 & = & k & + & (n - k - 1) \end{array}$$

- Assuming “Errors are i.i.d. Normal” & “x's are nonstochastic (i.e., can be measured without errors)” => y's are also normally distributed
- Various sums of squares have a *chi-square distribution* with the degrees of freedom as given above

F-test

- Given SS_i and SS_j with v_i and v_j degrees of freedom, the ratio $(SS_i/v_i)/(SS_j/v_j)$ has an F distribution with v_i numerator degrees of freedom and v_j denominator degrees of freedom.
- Hypothesis that the sum SS_i is less than or equal to SS_j is rejected at α significance level, if the ratio $(SS_i/v_i)/(SS_j/v_j)$ is greater than the $1-\alpha$ quantile of the F-variate.
- This procedure is also known as **F-test**.
- The F-test can be used to check:
Is SSR significantly higher than SSE?
 \Rightarrow Use F-test \Rightarrow Compute $(SSR/v_R)/(SSE/v_e) = MSR/MSE$

F-test (contd.)

$$\text{MSR} = \frac{\text{SSR}}{k} \quad \text{and} \quad \text{MSE} = \frac{\text{SSE}}{n - k - 1}$$

- ❑ MSE = Variance of Error
- ❑ MSR/MSE has $F[k, n-k-1]$ distribution
- ❑ F-test = Null hypothesis that y doesn't depend upon any x_j :
against an alternate hypothesis that y depends upon at least one x_j and therefore, at least one $b_j \neq 0$.
- ❑ If the computed ratio is less than the value read from the table, the null hypothesis cannot be rejected at the stated significance level.
- ❑ In simple regression models,
If the confidence interval of b_1 does not include zero
⇒ Parameter is nonzero
⇒ Regression explains a significant part of the response variation
⇒ F-test is not required.

ANOVA table for multiple linear regression

❑ See Table 15.3 on page 252

Component	Sum of Squares	%Variation	DF	Mean Square	F-Comp.	F-Table
y	$SSY = \sum y^2$		n			
\bar{y}	$SS0 = n\bar{y}^2$		1			
y - \bar{y}	$SST = SSY - SS0$	100	n-1			
Regression	$SSR = SST - SSE$	$100 \left(\frac{SSR}{SST} \right)$	k	$MSR = \frac{SSR}{k}$	$\frac{MSR}{MSE}$	$F_{[1-\alpha; k, n-k-1]}$
Errors	$SSE = \mathbf{y}^T \mathbf{y} - \mathbf{b}^T \mathbf{X}^T \mathbf{y}$	$100 \left(\frac{SSE}{SST} \right)$	n-k-1	$MSE = \frac{SSE}{n-k-1}$		
	$s_e = \sqrt{MSE}$					

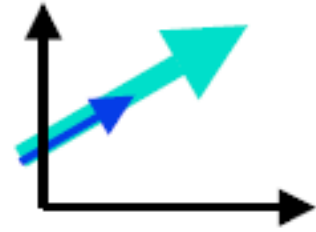
Example 15.2

- For the Disk-Memory-CPU data of Example 15.1
- Computed F ratio > F value from the table
⇒ Regression does explain a significant part of the variation

Component	Sum of Squares	%Variation	DF	Mean Square	F-Comp.	F-Table
y	828.					
\bar{y}	622.					
y- \bar{y}	206.	100.0%	6			
Regression	200.	97.4%	2	100.20	75.40	4.32
Errors	5.32	2.6%	4	1.33		

$$s_e = \sqrt{\text{MSE}} = \sqrt{1.33} = 1.15$$

- Note: Regression passed the F test ⇒ Hypothesis of all parameters being zero cannot be accepted. However, none of the regression parameters are significantly different from zero. This contradiction ⇒ Problem of **multicollinearity**



Multicollinearity

- Two lines are said to be collinear if they have the same slope and same intercept. (same line)
 - These two lines can be represented in just one dimension instead of the two dimensions required for lines which are not collinear.
 - Two collinear lines are not independent.
- When two predictor variables are linearly dependent, they are called collinear
- Collinear predictors => Problem of multicollinearity (i.e., contradictory results from various significance tests)
- High Correlation => Eliminate one variable and check if significance improves

Example 15.3

- For the data of Example 15.2, $n=7$, $\sum x_{1i}=271$, $\sum x_{2i}=1324$, $\sum x_{1i}^2=1385$, $\sum x_{2i}^2=326,686$, $\sum x_{1i}x_{2i}=67,188$.

$$\begin{aligned}\text{Correlation}(x_1, x_2) &= R_{x_1 x_2} \\ &= \frac{\sum x_{1i} x_{2i} - \frac{1}{n} (\sum x_{1i}) (\sum x_{2i})}{\left[\sum x_{1i}^2 - \frac{1}{n} (\sum x_{1i}) (\sum x_{1i}) \right]^{1/2} \left[\sum x_{2i}^2 - \frac{1}{n} (\sum x_{2i}) (\sum x_{2i}) \right]^{1/2}} \\ &= \frac{67,188 - \frac{1}{7} (271)(1324)}{\left[1385 - \frac{1}{7} (271)(271) \right]^{1/2} \left[326,686 - \frac{1}{7} (1324)(1324) \right]^{1/2}} = 0.9947\end{aligned}$$

- Correlation is high
 \Rightarrow Programs with large memory sizes have more I/O's
- In Example 14.1, CPU time on number of disk I/O's regression was found significant.

Example (contd.)

- Similarly, as shown in Exercise 14.3, CPU time is regressed on the memory size and the resulting regression parameters are found to be significant.
- Thus, *either* the number of I/O's *or* the memory size can be used to estimate CPU time, but not both.
- **Lesson:**
 - Adding a predictor variable does not always improve regression accuracy.
 - If the variable is correlated to other predictors, it may reduce the statistical accuracy (i.e., more variance) of the regression.
- Try all 2^k possible subsets and choose the one that gives the best results with small number of variables.
- Correlation matrix for the subset chosen should be checked

Outline

- Multiple Linear Regression
- Categorical Predictors
- Curvilinear Regression
- Transformations
- Outliers
- Common mistakes in regression

Regression of categorical (i.e., nonnumerical) predictors

- Note: If all predictor variables are categorical, use one of the experimental design and analysis techniques for statistically more precise (less variant) results
 - Use regression if most predictors are quantitative and only a few predictors are categorical
- Two Categories:
$$x_j = \begin{cases} 0 & \Rightarrow \text{First value} \\ 1 & \Rightarrow \text{Second value} \end{cases}$$

b_j : represents difference in the effect of the two alternatives

 - b_j is insignificant \Rightarrow two alternatives have similar performance
- Alternatively:
$$x_j = \begin{cases} -1 & \Rightarrow \text{First value} \\ +1 & \Rightarrow \text{Second value} \end{cases}$$

b_j : represents the difference from the average response

 - Difference of the effects of the two levels is $2b_j$

Categorical predictors (contd.)

- ❑ Three Categories: Incorrect:

$$x_1 = \begin{cases} 1 & \Rightarrow \text{Type A} \\ 2 & \Rightarrow \text{Type B} \\ 3 & \Rightarrow \text{Type C} \end{cases}$$

This coding implies an order \Rightarrow B is half way between A and C. This may not be true.

- ❑ Recommended: Use two predictor variables

$$x_1 = \begin{cases} 1, & \text{If type A} \\ 0, & \text{Otherwise} \end{cases}$$

$$x_2 = \begin{cases} 1, & \text{If type B} \\ 0, & \text{Otherwise} \end{cases}$$

Categorical predictors (contd.)

Thus, $(x_1, x_2) = (1, 0) \Rightarrow \text{Type A}$

$(x_1, x_2) = (0, 1) \Rightarrow \text{Type B}$

$(x_1, x_2) = (0, 0) \Rightarrow \text{Type C}$

- This coding does not imply any ordering among the types.
Provides an easy way to interpret the regression parameters.

$$y = b_0 + b_1x_1 + b_2x_2 + e$$

Categorical predictors (contd.)

- The average responses for the three types are:

$$\bar{y}_A = b_0 + b_1$$

$$\bar{y}_B = b_0 + b_2$$

$$\bar{y}_C = b_0$$

- Thus, b_1 represents the difference between type A and C.
 b_2 represents the difference between type B and C.
 b_0 represents type C.

Categorical predictors (contd.)

- ❑ Level = Number of values that a categorical variable can take
- ❑ To represent a categorical variable with k levels, define $k-1$ binary variables:

$$x_j = \begin{cases} 1, & \text{If } j\text{th value} \\ 0, & \text{otherwise} \end{cases}$$

- ❑ k th (last) value is defined by $x_1 = x_2 = \dots = x_{k-1} = 0$.
- ❑ b_0 = Average response with the k th alternative.
- ❑ b_j = Difference between alternatives j and k .
- ❑ If one of the alternatives represents the status quo or a standard against which other alternatives have to be measured, that alternative should be coded as the k th alternative.

Case study 15.1

- RPC performance on Unix and Argus

$$y = b_0 + b_1x_1 + b_2x_2$$

where, y is the elapsed time, x_1 is the data size and

$$x_2 = \begin{cases} 1 & \Rightarrow \text{UNIX} \\ 0 & \Rightarrow \text{ARGUS} \end{cases}$$

UNIX		ARGUS	
Data Bytes	Time	Data Bytes	Time
64	26.4	92	32.8
64	26.4	92	34.2
64	26.4	92	32.4
64	26.2	92	34.4
234	33.8	348	41.4
590	41.6	604	51.2
846	50.0	860	76.0
1060	48.4	1074	80.8
1082	49.0	1074	79.8
1088	42.0	1088	58.6
1088	41.8	1088	57.6
1088	41.8	1088	59.8
1088	42.0	1088	57.4

Case study (contd.)

Para- meter	Mean	Std. Dev.	Confidence Interval
b_0	36.739	3.251	(31.1676, 42.3104)
b_1	0.025	0.004	(0.0192, 0.0313)
b_2	-14.927	3.165	(-20.3509, -9.5024)

- All three parameters are significant (diff. from 0). The regression explains 76.5% of the variation.
- Per byte processing cost (time) for both operating systems is 0.025 millisecond.
- Set up cost is 36.73 milliseconds on ARGUS, which is 14.927 milliseconds more than that with UNIX.

Differing conclusions

- Case Study 14.1 concluded that there was no significant difference in the set up costs. The per byte costs were different.

Case Study 15.1 concluded that per byte cost is same but the set up costs are different.

- *Which conclusion is correct?*
 - Need system (domain) knowledge. Statistical techniques applied without understanding the system can lead to a misleading result ☹
 - Case Study 14.1 was based on the assumption that the processing as well as set up in the two operating systems are different
 - => four parameters
 - The data showed that the setup costs were numerically indistinguishable.

Differing conclusions (contd.)

- The model used in Case Study 15.1 is based on the assumption that the operating systems have no effect on per byte processing.
 - This will be true if the processing is identical on the two systems and does not involve the operating systems. i.e., only set up requires operating system calls.
 - If this is, in fact, true, then the regression coefficients estimated in the joint model of this case study 15.1 are more realistic estimates of the real world.
 - On the other hand, if system programmers can show that the processing follows a different code path in the two systems, then the model of Case Study 14.1 would be more realistic.

Outline

- Multiple Linear Regression
- Categorical Predictors
- Curvilinear Regression
- Transformations
- Outliers
- Common mistakes in regression

Curvilinear regression

- If the relationship between response and predictors is nonlinear but it can be converted into a linear form
⇒ curvilinear regression.

Example:

$$y = bx^a$$

Taking a logarithm of both sides we get:

$$\ln y = \ln b + a \ln x$$

Thus, $\ln x$ and $\ln y$ are linearly related. The values of $\ln b$ and a can be found by a linear regression of $\ln y$ on $\ln x$.

Other examples

Nonlinear

$$y = a + b/x$$

$$y = a + bx^n$$

$$y = 1/(a + bx)$$

$$y = x/(a + bx)$$

$$y = ab^x$$

Linear

$$y = a + b(1/x)$$

$$y = a + b(x^n)$$

$$(1/y) = a + bx$$

$$(x/y) = a + bx$$

$$\ln(y) = \ln(a) + (\ln(b))x$$

- If a predictor variable appears in more than one transformed predictor variables, the transformed variables are likely to be correlated \Rightarrow multicollinearity.

Try various possible subsets of the predictor variables to find a subset that gives significant parameters and explains a high percentage of the observed variation.

Example 15.4

- Amdahl's law: I/O rate is proportional to the processor speed. For each instruction executed there is one bit of I/O on the average.

System No.	MIPS Used	I/O Rate
1	19.63	288.60
2	5.45	117.30
3	2.63	64.60
4	8.24	356.40
5	14.00	373.20
6	9.87	281.10
7	11.27	149.60
8	10.13	120.60
9	1.01	31.10
10	1.26	23.70

Example (contd.)

- Let us fit the following curvilinear model to this data:

$$\text{I/O Rate} = \alpha(\text{MIPS Rate})^{b_1}$$

- Taking a log of both sides we get:

$$\log(\text{I/O Rate}) = \log(\alpha) + b_1 \log(\text{MIPS Rate})$$

$$b_0 = \log(\alpha)$$

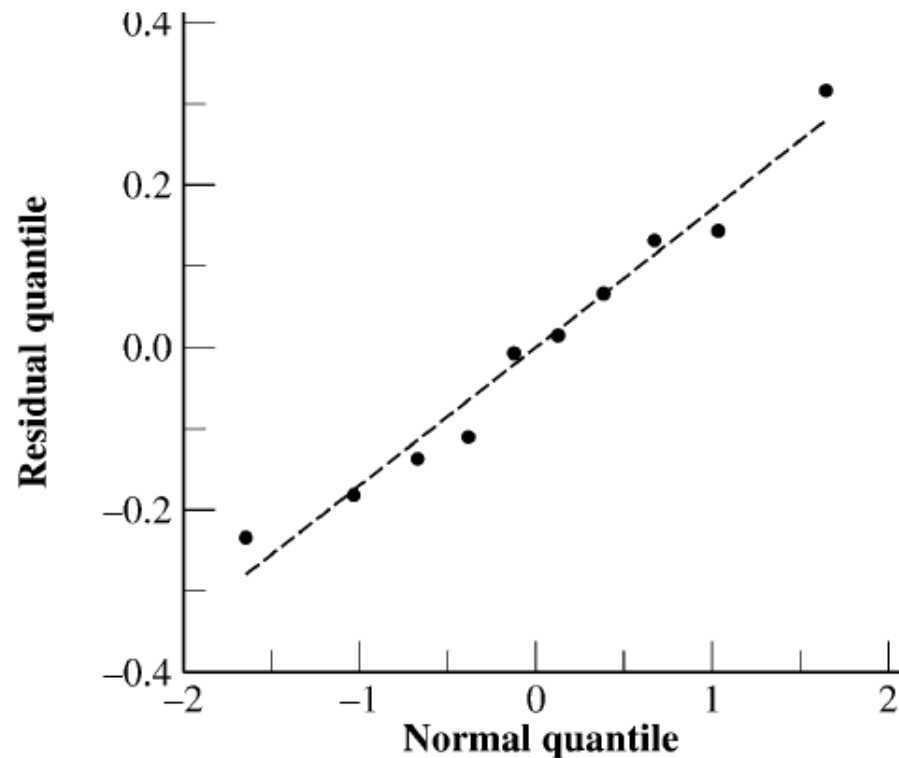
Example (contd.)

Obs. No.	x_1	y
1	1.293	2.460
2	0.736	2.069
3	0.420	1.810
4	0.916	2.552
5	1.146	2.572
6	0.994	2.449
7	1.052	2.175
8	1.006	2.081
9	0.004	1.493
10	0.100	1.375

Para-meter	Mean	Std. Dev.	Confidence Interval
b_0	1.423	0.119	(1.20, 1.64)
b_1	0.888	0.135	(0.64, 1.14)

- Both coefficients are significant at 90% confidence level.
- The regression explains 84% of the variation.
- At this confidence level, we can accept the hypothesis that the relationship is linear since the confidence interval for b_1 includes 1.

Example (contd.)



- ❑ Errors in log I/O rate do seem to be normally distributed.

Outline

- Multiple Linear Regression
- Categorical Predictors
- Curvilinear Regression
- Transformations
- Outliers
- Common mistakes in regression

Transformations

- Transformation: Some function of the measured response variable y is used. For example,

$$\sqrt{y} = b_0 + b_1x_1 + b_2x_2 + \cdots + b_kx_k + e$$

Transformation is a subset of the curvilinear regression. *experiment design*
However, the ideas apply to non-regression model as well.

- Physical considerations \Rightarrow Transformation
For example, if response is inter-arrival times y and it is known that the number of requests per unit time ($1/y$) has a linear relationship to a predictor
- If the range of the data covers several orders of magnitude and the sample size is small. That is, if y_{\max}/y_{\min} is large.
- If the homogeneous variance (**homoscedasticity**) assumption of the residuals is violated.

Transformations (contd.)

- scatter plot shows non-homogeneous spread \Rightarrow Residuals are still functions of the predictors
- Plot the standard deviation of residuals at each value of \hat{y} as a function of the mean \hat{y} .
- If s and the mean \bar{y} :

$$s = g(\bar{y})$$

- Then a transformation of the form:

$$w = h(y)$$

$$h(y) = \int \frac{1}{g(y)} dy$$

may help solve the problem

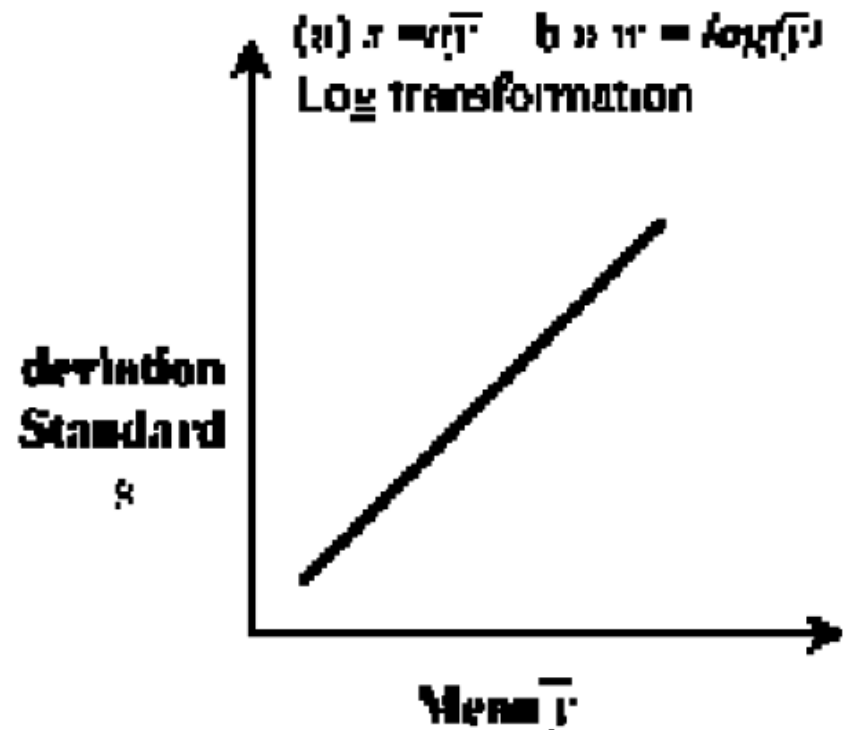
Useful transformations

- Log Transformation: Standard deviation s is a linear function of the mean ($s = a \bar{y}$)

$$w = \ln y$$

and, therefore:

$$h(y) = \int \frac{1}{ay} dy = a \ln y$$



Useful transformations (contd.)

Logarithmic transformation is useful only if the ratio y_{\max}/y_{\min} is large.

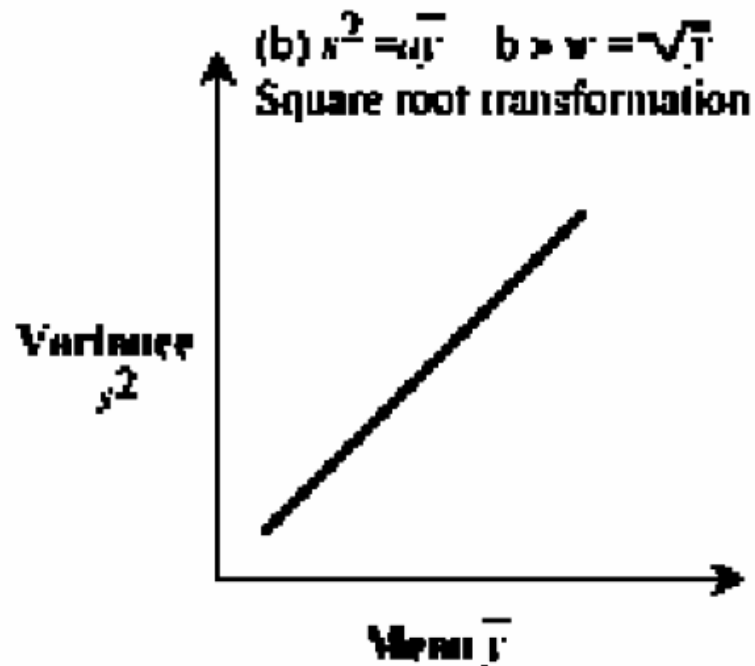
For a small range the log function is almost linear.

s.t. transformation
does not help much

- Square Root Transformation: For a Poisson distributed variable: $s = \sqrt{y}$

Variance versus mean
will be a straight line

$w = \sqrt{y}$ helps stabilize
the variance.



Useful transformations (contd.)

- Arc Sine Transformation: If y is a proportion or percentage, $\sin^{-1} \sqrt{y}$ may be helpful.
- Omega Transformation: This transformation is popularly used when the response y is a proportion.

$$w = 10 \log_{10} \left(\frac{y}{1-y} \right)$$

- The transformed values w 's are said to be in units of *deci-Bells*. The term comes from signaling theory where the ratio of output power to input power is measured in dBs.
- Omega transformation converts fractions between 0 and 1 to values between $-\infty$ to $+\infty$.
- This transformation is particularly helpful if the fractions are very small or very large.
- If the fractions are close to 0.5, a transformation may not be required.

Useful transformations (contd.)

- ❑ **Power Transformation:** y^a is regressed on the predictor variables.
 - Standard deviation of residuals s_e is proportional to \hat{y}^{1-a}

Useful transformations (contd.)

■ A short summary

Relationship between s and \bar{y}	Transformation
$s \propto \bar{y}$	$w = \ln(y)$ or $w = \ln(y + c)$
$s \propto \bar{y}^{1/2}$	$w = y^{1/2}$
$s \propto \bar{y}^a$	$w = y^{1-a}$ or $w = (y + c)^{1-a}$
$s \propto \bar{y}^2$	$w = \frac{1}{y}$
$s \propto 1 - \bar{y}^2$	$w = \ln \left(\frac{1+y}{1-y} \right)$
$s \propto \bar{y}(1 - \bar{y})$	$w = \ln \left(\frac{y}{1-y} \right)$
$s \propto (1 + \bar{y})\sqrt{\bar{y}}$	$w = \sin^{-1} \sqrt{y}$

Shifting: $y+c$ (with some suitable c) may be used in place of y .

- Useful if there are negative or zero values and if the transformation function is not defined for these values.

Box-cox family of transformations

- If the value of the exponent a in a power transformation is not known, Box-Cox family of transformations can be used:

$$w = \begin{cases} \frac{y^a - 1}{ag^{a-1}}, & a \neq 0 \\ (\ln y)g, & a = 0 \end{cases}$$

Where g is the geometric mean of the responses:

$$g = (y_1 y_2 \cdots y_n)^{1/n}$$

- The Box-Cox transformation has the property that w has the same units as the response y for all values of the exponent a .
- All real values of a , positive or negative can be tried.

The transformation is continuous even at zero, since:

$$\lim_{a \rightarrow 0} \frac{y^a - 1}{ag^{a-1}} = (\ln y)g$$

Box-cox transformations

- Use a that gives the smallest SSE.
- Use simple values for a. If $a=0.52$ is found to give the minimum SSE and the SSE at $a=0.5$ is not significantly higher, the latter value may be preferable.
- $100(1-\alpha)$ confidence interval for a: all the a for which the SSE is <

$$\text{SSE}_{\min} \left(1 + \frac{t_{[1-\alpha/2; \nu]}^2}{\nu} \right)$$

Where, SSE_{\min} is the minimum SSE, and ν is the number of degrees of freedom for the errors.

If the confidence interval for a includes $a = 1$, then the hypothesis that the relationship is linear cannot be rejected
 \Rightarrow No need for the transformation.

Case study 15.2: garbage collection

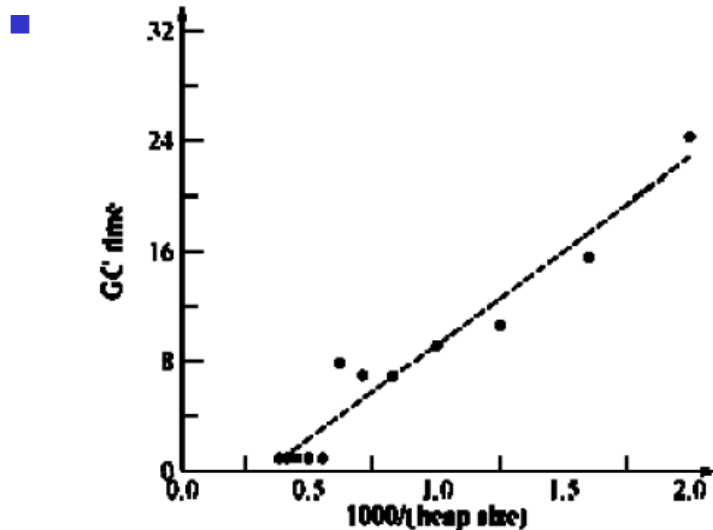
- The garbage collection time for various values of heap sizes.

Heap Size	Garbage Collection Time	Heap Size	Garbage Collection Time
500	594.34	1600	63.64
600	247.42	1800	1.00
800	114.24	2000	1.00
1000	85.64	2200	1.00
1200	49.60	2400	1.00
1400	50.30	2600	1.00

Case study (contd.)

- The analyst hypothesizes that

$$(\text{Time})^{1/2} = b_0 + \frac{b_1}{\text{Heap Size}}$$

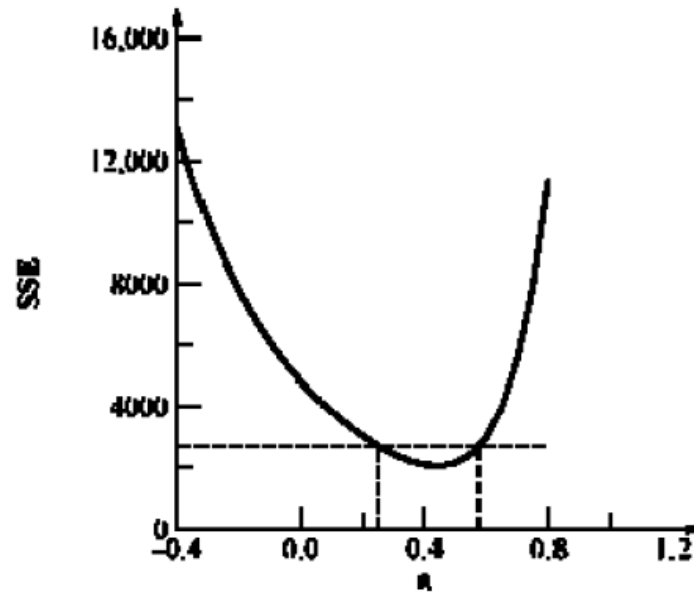


The points do not appear to be close to the straight line.

Case study (contd.)

- Is exponent on time different than a half?

=> Use Box-Cox transformations with “a” ranging from -0.4 to 0.8



- The minimum SSE of 2049 occurs at $a = 0.45$.

Case study (contd.)

- Since 0.95-quantile of a t variate with 10 degrees of freedom is 1.812

$$\begin{aligned}SSE &= 2049 \left(1 + \frac{(1.812)^2}{10} \right) \\ &= 2721.8\end{aligned}$$

- The $SSE = 2271$ line intersects the curve at $a = 0.2465$ and $a = 0.5726$.
- 90% confidence interval for a is $(0.2465, 0.5726)$. Since the interval includes 0.5, we cannot reject the hypothesis that the exponent is 0.5.

Outline

- Multiple Linear Regression
- Categorical Predictors
- Curvilinear Regression
- Transformations
- Outliers
- Common mistakes in regression

Outliers

- Any observation that is *atypical* of the remaining observations *may* be considered an outlier.
 - Including the outlier in the analysis may change the conclusions significantly.
 - Excluding the outlier from the analysis may lead to a misleading conclusion, if the outlier in fact represents a correct observation of the system behavior.
- A number of statistical tests have been proposed to test if a particular value is an outlier.
 - Most of these tests assume a certain distribution for the observations.
 - If the observations do not satisfy the assumed distribution, the results of the statistical test would be misleading.
- Easiest way to identify outliers is to look at the scatter plot of the data.

Outliers (contd.)

- Any value significantly away from the remaining observations should be investigated for possible experimental errors.
 - Other experiments in the *neighborhood of the outlying observation* may be conducted to verify that the response is typical of the system behavior in that operating region.

Once the possibility of errors in the experiment has been eliminated, the analyst may decide to include or exclude the suspected outlier based on the intuition.

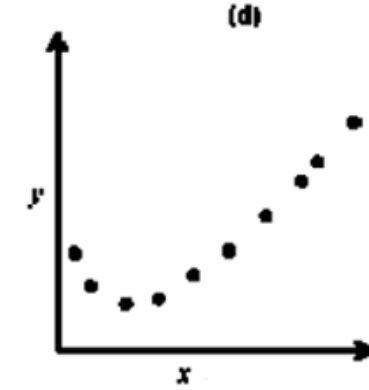
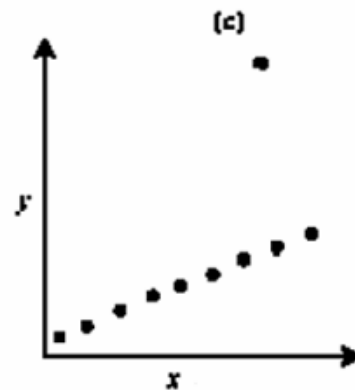
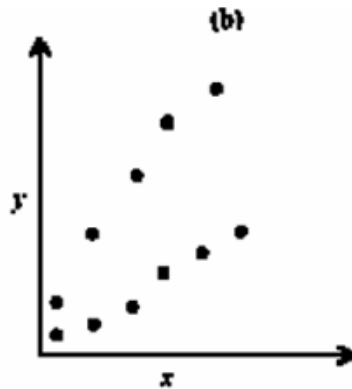
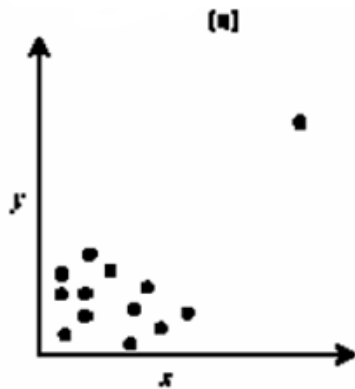
- One alternative is to repeat the analysis *with and without* the outlier and state the results separately.
- Another alternative is to *divide the operating region into two (or more) sub-regions* and obtain a separate model for each sub-region.

Outline

- Multiple Linear Regression
- Categorical Predictors
- Curvilinear Regression
- Transformations
- Outliers
- Common mistakes in regression

Common mistakes in regression

1. *Not verifying that the relationship is linear.*
2. *Relying on automated results without visual verification*



- In all these cases, $R^2 = \text{High}$
- High R^2 is necessary but not sufficient for a good model.

Common mistakes (contd.)

3. *Attaching importance to numerical values of regression parameters.*

$$\text{CPU time in seconds} = 0.01 (\text{Number of disk I/O's}) + 0.001 (\text{Memory size in kilobytes})$$

0.001 is too small \Rightarrow memory size can be ignored

$$\text{CPU time in milliseconds} = 10 (\text{Number of disk I/O's}) + 1 (\text{Memory size in kilobytes})$$

$$\text{CPU time in seconds} = 0.01 (\text{Number of disk I/O's}) + 1 (\text{Memory size in bytes})$$

4. *Not specifying confidence intervals for the regression parameters.*
5. *Not specifying the coefficient of determination.*

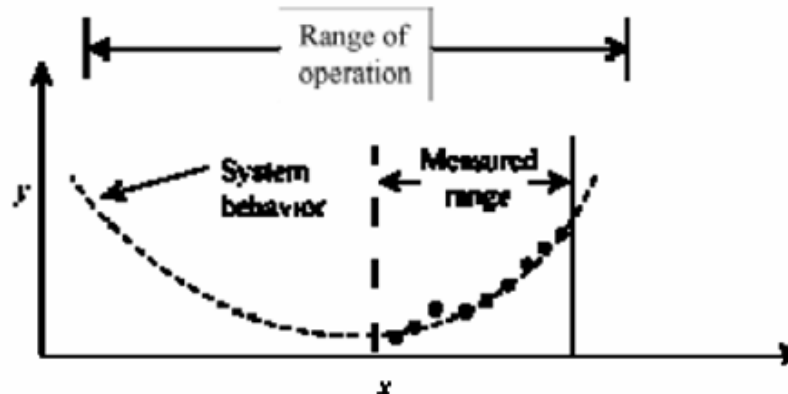
Common mistakes (contd.)

6. *Confusing the coefficient of determination and the coefficient of correlation*
R=Coefficient of correlation, R^2 = Coefficient of determination
 $R=0.8$, $R^2=0.64$
 \Rightarrow Regression explains only 64% of variation and not 80%.
7. *Using highly correlated variables as predictor variables.*
Analysts often start a multi-linear regression with as many predictor variables as possible
 \Rightarrow severe multicollinearity problems.
8. *Using regression to predict far beyond the measured range.*
Predictions should be specified along with their confidence intervals
9. *Using too many predictor variables.*
 k predictors $\Rightarrow 2^k-1$ subsets

Common mistakes (contd.)

Subset giving the minimum **SSE** is the *best*. But, other subsets that are close may be used instead for practical or engineering reasons. For example, if the second best has only one variable compared to five in the best, the second best may be the preferred model.

10. *Measuring only a small subset of the complete range of operation, e.g., 10 or 20 users on a 100 user system.*



Common mistakes (contd.)

11. Assuming that a good predictor variable is also a good control variable.

- Correlation \Rightarrow Can predict with a high precision
 \nRightarrow Can control response with predictor
- For example, the disk I/O versus CPU time regression model can be used to predict the number of disk I/O's for a program given its CPU time.
 However, reducing the CPU time by installing a faster CPU will not reduce the number of disk I/O's.
- w and y both controlled by x
 $\Rightarrow w$ and y highly correlated and would be good predictors for each other.

Common mistakes (contd.)

- The prediction works both ways:
w can be used to predict y and vice versa.
- The control often works only one way:
x controls y but y may not control x.

Summary

- Multiple Linear Regression
- Categorical Predictors
- Curvilinear Regression
- Transformations
- Outliers
- Common mistakes in regression

Summary of “performance evaluation”

- Common mistakes and how to avoid them
- Selection of techniques and metric
- Workload characterization techniques
- Introduction to experiment design
- 2^k factorial design
- One-factor experiments
- General full factorial design with k factors
- Introduction to simulation
- Summarizing measured data
- Comparing systems using sample data
- Regress models: simple linear regression, non-SL regress

Further reading

- *Richard Draves, Jitendra Padhye, Brian Zill, "Comparison of routing metrics for static multi-hop wireless networks", ACM SIGCOMM'04*
- *S. Corson, J. Macker, "Mobile Ad hoc Networking (MANET): Routing Protocol Performance Issues and Evaluation Considerations", IETF RFC 2501*
- *Hongwei Zhang, Anish Arora, Prasun Sinha, "Link Estimation and Routing in Sensor Network Backbones: Beacon-based or Data-driven?", IEEE Transactions on Mobile Computing, 2009*
- *Dongjin Son, Bhaskar Krishnamachari, John Heidemann, "Experimental Analysis of Concurrent Packet Transmissions in Low-Power Wireless Networks", ACM SenSys'06*
- *Lili Qiu, Yin Zhang, Feng Wang, Mi Kyung Han, Ratul Mahajan, "A General Model of Wireless Interference", ACM MOBICOM'07*

Exercise

Time to encrypt or decrypt a k-bit record was measured on a uniprocessor as well as on a multi-processor. The times in milliseconds are shown below. Using a log transformation and the method for categorical predictors fit a regression model and interpret the results.

k	Uniprocessor	Multiprocessor
128	93	67
256	478	355
512	3408	2351
1024	25,410	17,022