

Analysis of transformations

Project 1
Regression Analysis

Amal Agarwal
Roll. no.: 09D11001

Sagar Chordia
Roll. no.: 09005013

Tuhin Sarkar
Roll. no.: 09007030

Advisor: Prof. Siuli Mukhopadhyay



Department of Mathematics
Indian Institute of Technology, Bombay

2013

Chapter 1

Review of Box-Cox Transformations

1.1 Introduction

The usual techniques for linear regression are justified by the following underlying assumptions:

- simplicity of structure of $E(y)$
- constancy of error variance
- normality of distributions
- independence of observations

However, sometimes these assumptions may not hold. The motivation behind this paper is to derive transformations such that these not only justify these assumptions but rather find, wherever possible, a metric in terms of which the finding may be succinctly expressed.

Tukey (1949, 1950) used transformations to achieve additive in the analysis of variance. To stabilize the variance, the usual method of transformation by Bartlett (1947) is used where the relation between mean and variance is determined. Anscombe (1961) and Tukey and Anscombe (1983) have employed analysis of residuals to detect departures from standard assumptions and have indicated how transformations might be devised from functions of the residuals. There are many problems where both dependent and independent variables are transformed. Box and Tidwell (1962) transformation can be employed without affecting the constancy of variance and normality of error distribution. It is useless to try to linearize a relation which is not monotonic, but a transformation is sometimes useful in such cases.

1.2 General Idea

The general idea is to study transformed variables y^λ indexed by an unknown parameter λ , estimate λ and other parameter of our model. There are two primary lines of analyses. First is where the particular λ is of interest. Second is where we study the factor effects of a choice of λ .

The important family of transformations discussed here are:

$$y^\lambda = \frac{y^\lambda - 1}{\lambda}$$

when $\lambda \neq 0$

$$y^\lambda = \log(y)$$

when $\lambda = 0$

Now this transformation is valid only when $y \geq 0$. There is also a second set of transformations which hold for $y \geq \lambda_2$

$$y^\lambda = \frac{(y + \lambda_2)^{\lambda_1} - 1}{\lambda_1}$$

when $\lambda_1 \neq 0$

$$y^\lambda = \log(y + \lambda_2)$$

when $\lambda_1 = 0$

But the second transformation is nothing but a linear transformation of the first one. Therefore we can analyze the first one itself.

The basic idea is that once we have the transformation we can use our general linear regression model:

$$E[y^\lambda] = a\theta$$

where a is the independent variable matrix and θ are our parameters.

There are two approaches to find λ :

- Use the maximum likelihood method
- Use a Bayesian approach

In the first method we use the likelihood function:

$$\frac{1}{\sqrt{(2\pi)^{n/2}}} \exp \left\{ -\frac{(y^\lambda - a\theta)'(y^\lambda - a\theta)}{2\sigma^2} \right\} J(\lambda; y)$$

$$J(\lambda; y) = \prod_{i=1}^n \left(\frac{dy_i^\lambda}{dy_i} \right)$$

Note that $J(\lambda; y)$ is the Jacobian of transformation from y to y^λ . In this first, we apply "orthodox" large-sample maximum likelihood theory. This approach leads to point estimates of the parameters and to approximate tests and confidence intervals based on chi-squared distribution.

Now, we can use the normal linear regression techniques on the transformed dependent variables, where we get:

$$\hat{\sigma}^2 = y^{\lambda'} a_r y^\lambda / n$$

$$a_r = I - a(a'a)^{-1}a'$$

The log-likelihood function looks like:

$$L_{max}(\lambda) = -\frac{1}{2}n \log \hat{\sigma}^2(\lambda) + \log J(\lambda; y)$$

To get λ from this, plot the maximized log-likelihood function for a series of trial λ values. From the plot the maximizing value of λ can be taken as:

$$L_{max}(\hat{\lambda}) - L_{max}(\lambda) < \frac{1}{2}\chi_\nu^2(\alpha)$$

Here ν is the number of degrees of freedom of λ

The above results can be simplified by using the further transformation given by:

$$z^\lambda = y^\lambda / J^{1/n}$$

We now discuss how to approach the problem of finding λ using the Bayesian analysis. Consider

$$p(y|\theta, \sigma^2) = \frac{1}{(2\pi)^{n/2}\sigma^n} \exp\left\{-\frac{\nu s^2(\lambda) + (\theta - \hat{\theta}_\lambda)'a'a(\theta - \hat{\theta}_\lambda)}{2\sigma^2}\right\} J(\lambda; y)$$

Where $\hat{\theta}_\lambda$ is the least square estimate of θ for given λ .

In this method we use the Bayes theorem to get the posterior distribution of λ by using the following:

$$I(\lambda|y) = \int_{-\infty}^{\infty} d(\log \sigma) \int_{-\infty}^{\infty} d\theta p(y|\theta, \sigma^2)$$

$$p(\lambda|y) = K'_y \frac{I(\lambda|y)p_0(\lambda)}{\{J(\lambda; y)\}^{(n-\nu)/n}}$$

Note K'_y is a normalizing constant. Here is also assumed that θ and $\log \sigma$ are effectively uniform over the range where likelihood is appreciable. By taking $p_0(\lambda)$ as uniform (prior distribution), we maximize the posterior distribution. That is required λ .

1.3 Further Analysis of Transformations

This analysis indicates the following things:

- how simple a model we are justified in using
- what weight is given to considerations (simple expectation structure, constant variance and normal distributions) in choosing λ
- whether different transformations are really needed

The basic ideas behind this analysis is to use the two approaches as before. We first discuss the log-likelihood method. Consider the constraint C . Now:

$$L_{max}(\lambda|C) = L_{max}(\lambda) + \{L_{max}(\lambda|C) - L_{max}(\lambda)\}$$

The two terms are examined separately. Similarly when more number of constraints are imposed.

Now in the Bayesian case:

$$p(\lambda|C) = p(\lambda) \frac{p(C|\lambda)}{p(C)}$$

Chapter 2

Simulation

2.1 Aim

For non-trivial application of Box-cox transformation, data should not satisfy all properties of Linear normal regression model. So we consider various datasets as cases where one or more assumptions are violated.

2.2 Box-Cox transformation

We apply box-cox transformation on given dataset. To estimate λ we use maximum likelihood method. Note X is generated from uniform distribution and hence rank of X is most of times equal to N number of data points. Our choice of dataset ensures that we can use maximum likelihood method.

On evaluation of concepts in section 2 it can be shown that $\hat{\lambda}$ can be found by minimizing following function

$$RSS(V(\lambda)) = [V(\lambda) - \hat{V}(\lambda)]^t [V(\lambda) - \hat{V}(\lambda)]$$

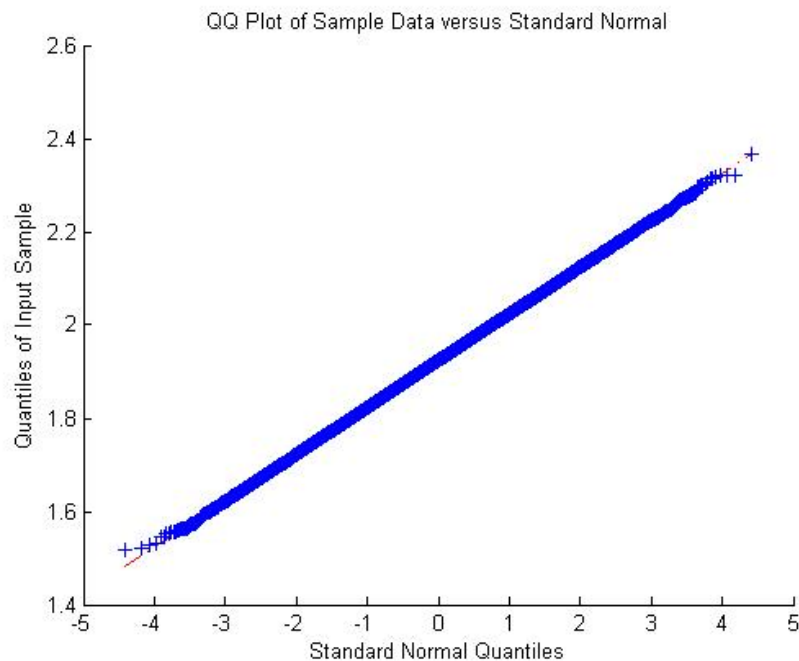
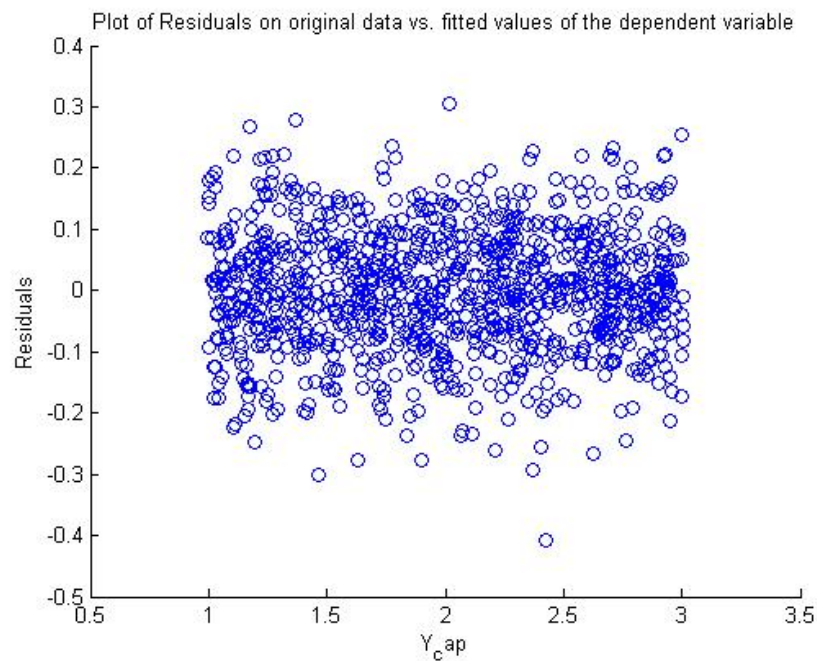
where

$$\begin{aligned}\hat{V}(\lambda) &= X(X^t X)^{-1} X^t V(\lambda) \\ V(\lambda) &= \begin{bmatrix} V_1(\lambda) \\ V_2(\lambda) \\ \vdots \\ V_n(\lambda) \end{bmatrix} \\ V_i(\lambda) &= (Y_i^\lambda - 1) / \lambda \widehat{Y}^{\lambda-1}, \lambda \neq 0 \\ V_i(\lambda) &= \widehat{Y} \log(Y_i), \lambda = 0 \\ \widehat{Y} &= (Y_1 Y_2 \cdots Y_n)^{1/N}\end{aligned}$$

2.3 Analysis of the Simulation Results

2.3.1 Normal linear data with constant variance over X

This dataset satisfies all conditions of linear regression. For this data we plot graph of semi studentized residuals vs X value and QQplot for data. We also perform lack of fit test of both data sets.



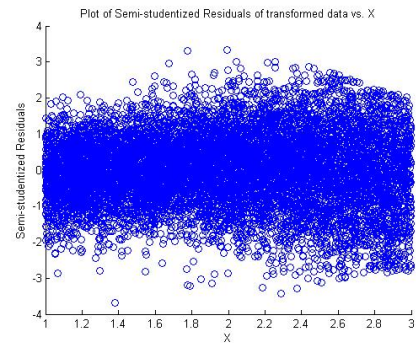
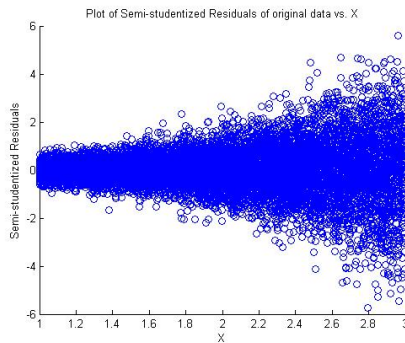
As we can see residual plot is spread uniformly across X over range of Y. QQplot is linear and lack of fit test confirms that model is linear. Thus when we perform box-cox transformation on data which doesn't satisfy some of the linear regression conditions should return graphs similar to above one in ideal case. The estimated parameters of model are:

$$\beta_0=1$$

$$\beta_1=2$$

2.3.2 Normal linear data with heteroscedasticity

As we can see in residual plot the variance changes as X changes. We plot residual graph (semi-studentized residual vs X) of original and transformed data.



From the uniform span of residual plot in transformed graph, it is clear that box-cox transformation is helpful in removing the heteroscedasticity from the data. The estimated parameters of model are:

$$\beta_0=-2.13$$

$$\beta_1=7.6$$

$$\lambda_{max}=0.34$$

2.3.3 Normal non-linear data with constant variance over X

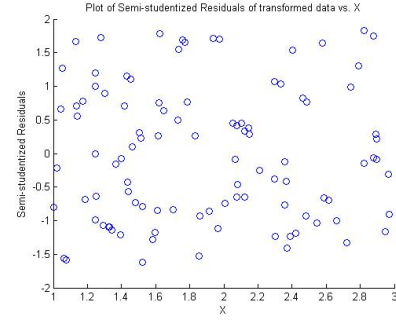
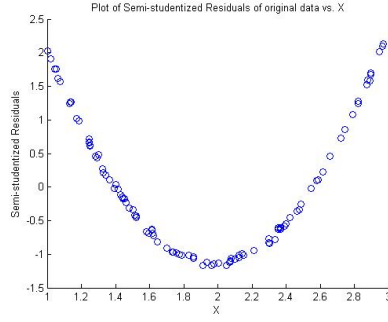
Here we have used some non-linear function to generate Y with respect to X. Here we have used quadratic function. Again we plot residual graph (semi-studentized residual vs X) for original and transformed data.

From the uniform span of residual plot in transformed graph, it is clear that box-cox transformation is helpful in removing the non-linearity from the data. The estimated parameters of model are:

$$\beta_0=50.4$$

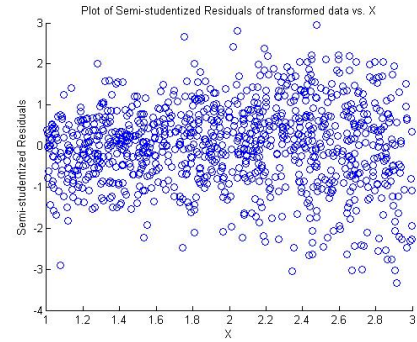
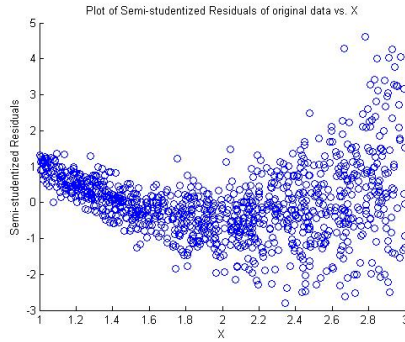
$$\beta_1=4.5$$

$$\lambda_{max}=0.48$$



2.3.4 Normal non-linear data with heteroscedasticity

Original data follows quadratic equation with heteroscedasticity of variance. We plot residual graph (semi-studentized residual vs X) for original and transformed data.



From the uniform span of residual plot in transformed graph, it is clear that box-cox transformation is helpful in removing the non-linearity and heteroscedasticity simultaneously. The estimated parameters of model are:

$$\beta_0 = -1.27$$

$$\beta_1 = 7.4$$

$$\lambda_{max} = 0.29$$

The matlab code for this case is given in next chapter.

2.4 Outlier detection and removal

To remove outliers we use semi-studentized graph. All data points whose semi-studentized variance is greater than +3 or less than -3 we term them as outliers and remove them from original dataset. Then we perform the box-cox transformation for transformed modified data and display its semi-studentized residual plot again.

2.5 Lack of Fit test

We did lack of fit test for Case 3 and Case 4 to check linearity of transformed data. The test results were negative for the initial model but after transformation it gave positive results thus justifying the aim of the transformation. We did hypothesis test for $\beta_0 = 0$ and $\beta_1 = 0$ vs.

corresponding alternatives for 95 percent confidence coefficient that is $\alpha = 0.05$ Calculated t-statistic in each case was much greater leading to rejection of null hypothesis. Thus we conclude $\beta_0 \neq 0$ and $\beta_1 \neq 0$ after transformation.

2.6 Unresolved case of Non-normality

We performed transformations on non-normal data with different distributions like chi-square, t, F and beta. However in each case the qq plot for the transformed data was devoid of any improvement compared to the original data thus refuting the assumption given in the paper that Box-Cox transformation can be used to make non-normal data again normal.

2.7 Conclusions

This paper gives an idea of a transformation that can be employed to modify dependent variables so that they are workable (standard linear regression techniques can be applied). This can be used for removing non-linearity and heteroscedasticity from the data as seen from the Simulation results. Unfortunately this fails to make non-normal data again normal. At the end of the paper there is also some discussion as to how we can also use a transformation where even the independent variables are modified.

Chapter 3

Matlab Code

```
clear all; clc;
n = 1000; alpha = 0.05; X1 = ones(n, 1); X2 = 1 + 2 * rand(n, 1); X = [X1, X2];
beta0 = 1; beta1 = 2;
var = 0.1*exp(X2); er = normrnd(0, var, n, 1);
Y = beta0 + beta1*(X2.^2) + er;
[b, bint, r, rint, stats] = regress(Y, X); Ycap = b(1) + b(2) * (X2);
rs = r/(sqrt(stats(4))); figure, scatter(X2, rs); title('Plot of Semi-studentized Residuals of original data');
studentizedResiduals');
m1 = 10; y = n+1; erlf = zeros(m1, (y - 1)); Ylf = zeros(m1, (y - 1)); Ymeanlf = zeros(1, (y - 1)); YSSPElf = zeros(m1, (y - 1)); YSSElf = zeros(m1, (y - 1)); for j = 1 : (y - 1) for i = 1 : m1 erlf(i, j) = normrnd(0, var(j)); Ylf(i, j) = beta0 + beta1 * (X2(j)) + erlf(i, j); end Ymeanlf(j) = mean(Ylf(:, j)); for i = 1 : m1 YSSPElf(i, j) = (Ylf(i, j) - Ymeanlf(j))^2; YSSElf(i, j) = (Ylf(i, j) - (Ycap(j)))^2; end end SSPE = sum(sum(YSSPElf)); SSE = sum(sum(YSSElf)); SSLF = SSE - SSPE; MSLF = SSLF/((y - 1) - 2); MSPE = SSPE/((m1 * (y - 1)) - (y - 1)); Fstar = MSLF/MSPE; F = finv((1 - alpha), ((y - 1) - 2), ((m1 * (y - 1)) - (y - 1))); if Fstar > F disp('Regression Function is not linear for original model'); else
k = 0; for i = 1:n if (abs(rs(i))) > 3 k = [i, k]; end end
s = length(k); km = zeros(1, s-1); for i = 1:s-1 km(i) = k(i); end
X2temp1 = X2; vartemp1 = var; ertemp1 = er; Ylftemp1 = Ylf; for i = 1 : n for j = 1 : s - 1 if (km(j) == i) X2temp1(i) = 0; vartemp1(i) = 0; ertemp1(i) = 0; Ylftemp1(:, i) = zeros(m1, 1); end end
X2temp2 = 0; vartemp2 = 0; ertemp2 = 0; Ylftemp2 = zeros(m1, 1); for p = 1 : n if (X2temp1(p) == 0) X2temp2 = [X2temp1(p); X2temp2]; end if (vartemp1(p) == 0) vartemp2 = [vartemp1(p); vartemp2]; end if (ertemp1(p) == 0) ertemp2 = [ertemp1(p); ertemp2]; end if (Ylftemp1(:, p) == zeros(m1, 1)) Ylftemp2 = [Ylftemp1(:, p); Ylftemp2]; end
y = length(X2temp2); X2m = zeros((n - s + 1), 1); varm = zeros((n - s + 1), 1); erm = zeros((n - s + 1), 1); Ylfm = zeros(m1, (n - s + 1)); for i = 1 : y - 1 X2m(i) = X2temp2(i); varm(i) = vartemp2(i); erm(i) = ertemp2(i); Ylfm(:, i) = Ylftemp2(:, i); end
X1m = ones(y - 1, 1); Xm = [X1m, X2m];
Ym = beta0 + beta1 * (X2m.^2) + erm;
m2 = 100000; mid = round((y-1)/2); erd = normrnd(0, var(mid), m2, 1); Yd = beta0 + beta1 * (X2(mid)) + erd; figure, qqplot(Yd);
GM = geomean(Ym); V = @(lambda)((Ym.^lambda) - 1)./(lambda * (GM^lambda - 1)); Vcap = @(lambda)(Xm*(inv(Xm'*Xm))*(Xm'*V(lambda))); RSS = @(lambda)((V(lambda) - Vcap(lambda))'*(V(lambda) - Vcap(lambda))); lambdamax = fminbnd(RSS, -100, 100); W =
```

```

V(lambda_max); W_d = (((Y_d.^lambda_max) - 1)./(lambda_max * (GM(lambda_max - 1)))); W_lf =
(((Y_lf.^lambda_max) - 1)./(lambda_max * (GM(lambda_max - 1))));
[b_t, bint_t, r_t, rint_t, stats_t] = regress(W, X_m); W_cap = b_t(1) + b_t(2) * (X_2m); r_s_t = r_t / (sqrt(stats_t(4)));
figure, scatter(X_2m, r_s_t); title('Plot of Semi-studentized Residuals of transformed data vs. X'); xlabel('X'); ylabel('Semi-studentized Residuals');
figure, qqplot(W_d);
W_mean_lf = zeros(1, (y-1)); W_SSPE_lf = zeros(m1, (y-1)); W_SSE_lf = zeros(m1, (y-1));
for j = 1 : (y-1) W_mean_lf(j) = mean(W_lf(:, j)); for i = 1 : m1 W_SSPE_lf(i, j) =
(W_lf(i, j) - W_mean_lf(j))^2; W_SSE_lf(i, j) = (W_lf(i, j) - W_cap(j))^2; end
end SSPE_t = sum(sum(W_SSPE_lf)); SSLF_t = SSE_t - SSPE_t; MSLF_t = SSLF_t / ((y-1) - 2); MSPE_t =
SSPE_t / ((m1 * (y-1)) - (y-1)); F_star_t = MSLF_t / MSPE_t; F_t = finv((1-alpha), ((y-1) - 2), ((m1 * (y-1)) - (y-1))); if F_star_t > F_t disp('Regression Function is not linear after transformation');

```

Bibliography

- [1] G.E.P. Box and D.R. Cox, An Analysis of Transformations, Journal of the Royal Statistical Society. Series B (Methodological), Vol.26, No.2 (1964), pp.211-252
- [2] MATLAB, version 7.10.0.499 (R2010a)