

STATISTICSII PROJECT REPORT

Submitted by:

Chandrabali Karmakar

SGGW Nr Albumu : 177639

HNEE Matriculation No: 14209583

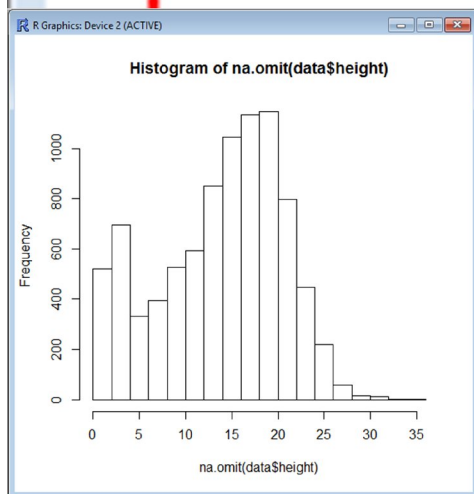
Contents

1) plot a histogram of *height* variable without *NA* records	3
3) check if the average heights of *BRZ* (brzoza = birch) and *SO* (sosna = pine) species are the same.....	4
4) check if the distribution of height for *SO* species follows the normal distribution.....	5
5) using two-way ANOVA - check if there is a significant difference in stand heights between species (species_cd) and age classes (0-20, 21-40, 41-60, 61-80, 81-100, 101-120, above 120 years, based on species_age variable). Use data only for *SO*, *BRZ* and *DB* (dąb = oak) species.....	5
Using R	5
Using SPSS:.....	6
6) using data for *SO* species (without *NA* records) find two variables with a non-linear relationship (e.g. age <> height). Build a nonlinear model for this relationship (choose the best out of at least 3 candidates from example from "CurveCatalog.pdf" publication or from other source). Perform a detailed regression analysis as well as justify the choice of the final model. Please do not limit yourself to the automatic procedure in SPSS, but perform a real process of choosing the proper candidate functions.	9
Using R	9
Using SPSS:.....	13

1) plot a histogram of ***height*** variable without ***NA*** records

Hist(na.omit(height))

```
> hist(na.omit(data$height))
```



2) calculate average height of all stands (not counting ***NAs***) as well as average heights of stands by species

```
> summary(na.omit(height))
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  1.00   10.00   15.00   14.27   19.00   35.00
```

```

> vdata=cbind(height,species_cd)
> dfdata=data.frame(vdata)
> ref_dfdata=na.omit(dfdata)
> names(ref_dfdata)
[1] "height"      "species_cd"
> h=ref_dfdata$height
> s=ref_dfdata$species_cd
> result=tapply(h,s,mean)
> result
      1      2      3      4      5      6      7      8      9     10
12.296875 15.661017  6.942857 14.710037 19.142857  2.750000 12.435897  3.133333 26.333333 18.000000
     11     12     13     14     15     16     17     18     19     20
11.971429 16.000000 20.500000 11.702703 17.425220 13.000000 16.000000 14.161620 27.600000 21.000000

```

3) check if the average heights of *BRZ* (brzoza = birch) and *SO* (sosna = pine) species are the same

levels(species_cd)

1. "IW"
2. "AK"
3. "BK"
4. "BRZ"
5. "DB"
6. "DB.B"
7. "DB.C"
8. "DB.S"
9. "DG"
10. "GB"
11. "JS"
12. "JW"
13. "LP"
14. "MD"
15. "OL"
16. "OL.S"
17. "OS"
18. "SO"
19. "TP"
20. "WZ"

So "BRZ" is the 4th in the list which shows mean as 14.710037 and "SO" is the 18th which gives 14.161620 as the average.

1	2	3	4	5	6	7	8	9
2.296875	15.661017	6.942857	14.710037	19.142857	2.750000	12.435897	3.133333	26.333333
17	18	19	20					
6.000000	14.161620	27.600000	21.000000					

4) check if the distribution of height for *SO* species follows the normal distribution

```
> pine_heights = split(h,s)$'18'
> ks.test(pine_heights,"pnorm",mean(pine_heights),sd(pine_heights))

One-sample Kolmogorov-Smirnov test

data: pine_heights
D = 0.10243, p-value < 2.2e-16
alternative hypothesis: two-sided
```

Is not a normal distribution.

[precise mean and standard variation are used so that the test is not done on a standard gaussian distribution.]

5) using two-way ANOVA - check if there is a significant difference in stand heights between species (species_cd) and age classes (0-20, 21-40, 41-60, 61-80, 81-100, 101-120, above 120 years, based on species_age variable). Use data only for *SO*, *BRZ* and *DB* (dąb = oak) species.

Using R

Step1: Getting the right data without “NA” values

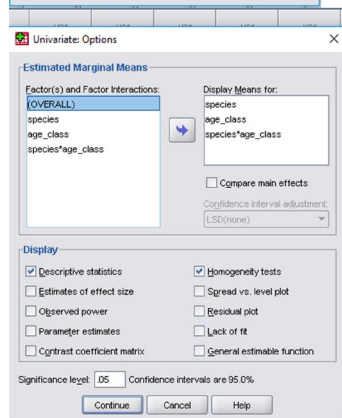
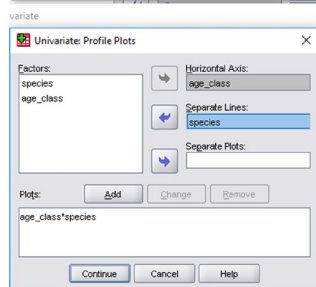
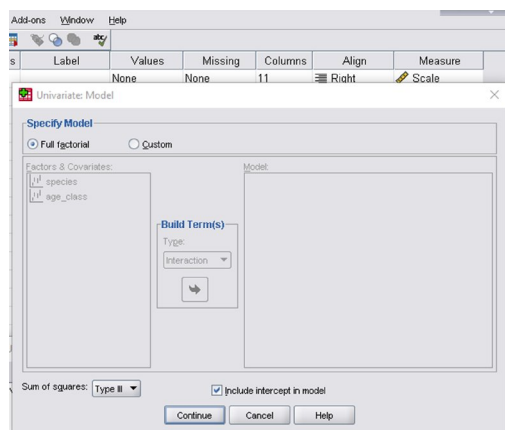
```
> data=read.csv(file.choose())
> names(data)
[1] "adresa_forest"      "stand_struct_cd"    "sub_area"           "site_type_cd"
[5] "storey_cd"          "density_cd"         "standdensity_index" "species_cd"
[9] "species_rank_order" "part_cd"            "species_age"        "bhd"
[13] "height"             "basal_area"         "volume"             "stem_number"
> df = data.frame(data$species_cd,data$species_age,data$height)

> df_subset=subset(df,df$data.species_cd== "SO" | df$data.species_cd == "DB" | df$data.species_cd == "BRZ" )
> df_noNA = na.omit(df_subset)
> names(df_noNA)
[1] "data.species_cd" "data.species_age" "data.height"
> |
```

Note that we had unequal sample sizes, so I put a little more effort to get sample with equal size for each species and age class. But the species “BRZ” does not have the “above 120” age class. However, the data satisfies the assumption of homogeneity of variance through Levene’s test (shown in SPSS result part).

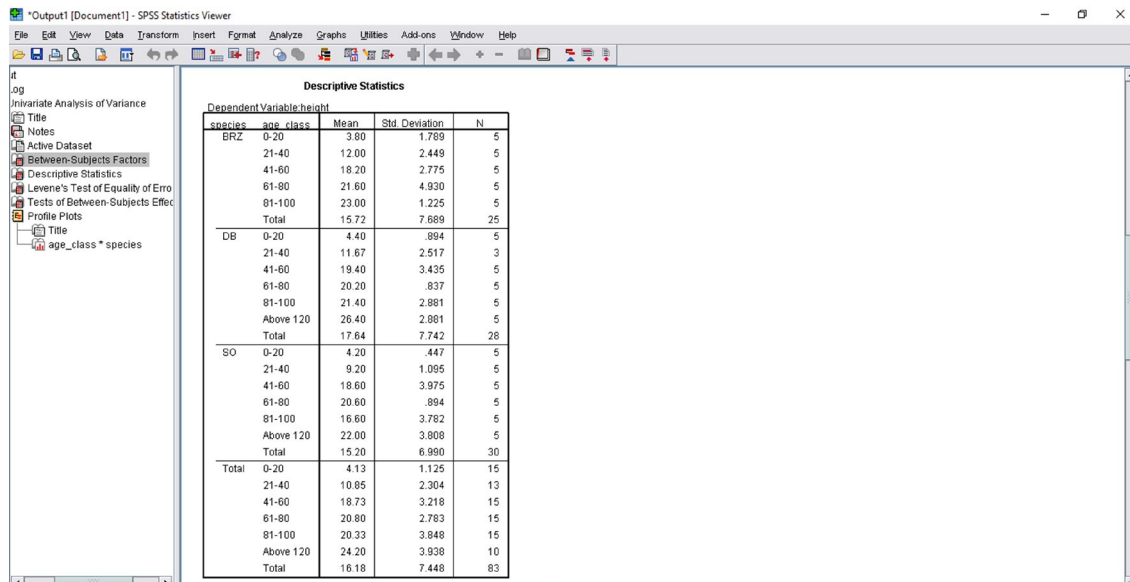
Step 2: Creating categorical variable age_classes from the continuous variable “Species_age”

Step2:Specifying Model ,plots and other options

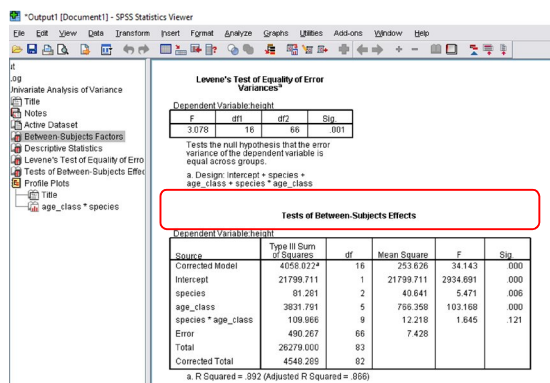
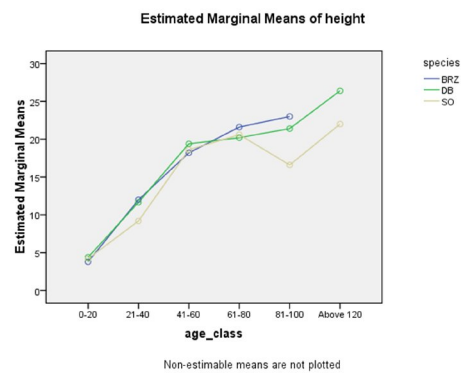


Test result:

Univariate Analysis of Variance		
[DataSet1]		
Between-Subjects Factors		
		N
species	BRZ	25
	DB	28
	SO	30
age_class	0-20	15
	21-40	13
	41-60	15
	61-80	15
	81-100	15
	Above 120	10



Profile Plots



Discussion : Both SPSS and R displays similar results that age and species have significant effect on stand height but not their interaction.

6) using data for *SO* species (without *NA* records) find two variables with a non-linear relationship (e.g. age \leftrightarrow height). Build a nonlinear model for this relationship (choose the best out of at least 3 candidates from example from "CurveCatalog.pdf" publication of from other source). Perform a detailed regression analysis as well as justify the choice of the final model. Please do not limit yourself to the automatic procedure in SPSS, but perform a real process of choosing the proper candidate functions.

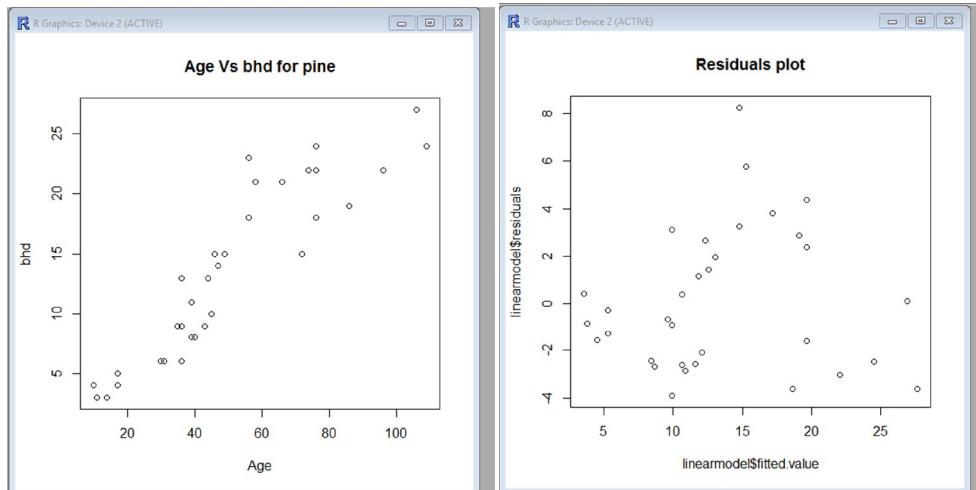
Using R

```
> data=read.csv(file.choose())
> names(data)
 [1] "adrees_forest"      "stand_struct_od"    "sub_area"
 [4] "site_type_od"       "storey_od"          "density_od"
 [7] "standdensity_index" "species_od"         "species_rank_order"
[10] "part_od"            "species_age"        "bhd"
[13] "height"             "basal_area"         "volume"
[16] "stem_number"
> pine_data=subset(data,species_od=="SO")
> names(pine_data)
 [1] "adrees_forest"      "stand_struct_od"    "sub_area"
 [4] "site_type_od"       "storey_od"          "density_od"
 [7] "standdensity_index" "species_od"         "species_rank_order"
[10] "part_od"            "species_age"        "bhd"
[13] "height"             "basal_area"         "volume"
[16] "stem_number"
```

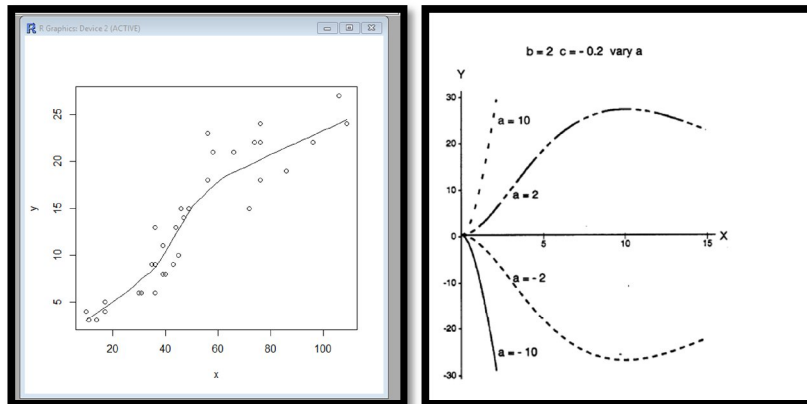
```
> raw_data=cbind(pine_data$species_age,pine_data$bhd)
> df=data.frame(raw_data)
> df_noNA = na.omit(df)
> mydata=df_noNA
Error: object 'df_noNA' not found
> raw_data=cbind(pine_data$species_age,pine_data$bhd)
> df=data.frame(raw_data)
> df_noNA = na.omit(df)
> mydata=df_noNA
```

```
> number=(1:200)
> mysample=sample(number,40,replace=T)
> x=X1[mysample]
Error: object 'X1' not found
> x=mydata$X1[mysample]
> y=mydata$X2[mysample]
```

```
> plot(y~x,main="Age Vs bhd for pine" ,xlab="Age",ylab="bhd")
> linearmodel = lm(y~x)
> plot(linearmodel$residuals~linearmodel$fitted.value , main="Residuals plot")
> 
```



The residuals are not evenly dispersed which says that I need a nonlinear model for this relationship.



The scatter plot seems to be an exponential function graph . But to ensure the correct model , I try a power function and a polynomial function and finally choose one among the three.

I. the exponential model:

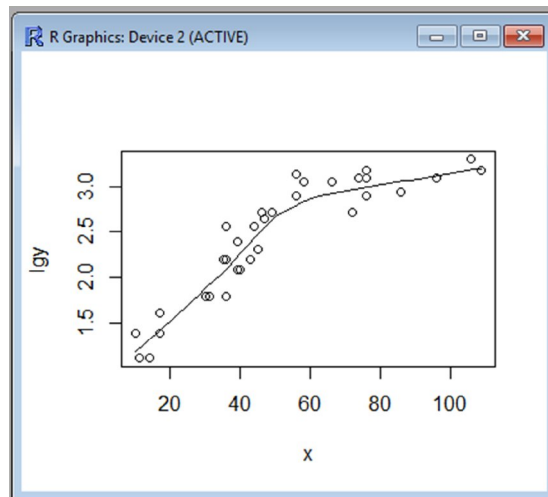
$$Y = ae^{bx}$$

A log transform is needed.

Linearized model and parameters: $\ln(Y) = b_0 + b_1X$

$a = eb_0$ $b = b_1$

Description: The parameter **a** is the Y-intercept; the parameter **b** is the shape parameter of the curve.

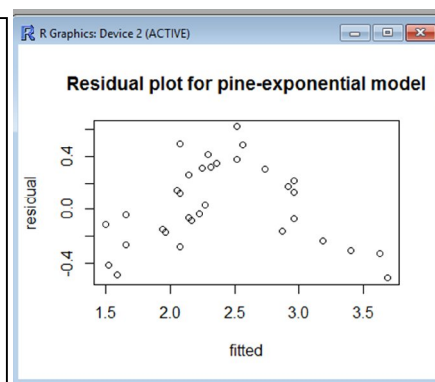


```
Call:
lm(formula = lgy ~ x)

Residuals:
    Min       1Q   Median       3Q      Max
-0.51307 -0.24399 -0.05115  0.27050  0.61923

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.274907   0.105315   12.11 1.31e-14 ***
x             0.022167   0.001853   11.96 1.86e-14 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.3254 on 38 degrees of freedom
Multiple R-squared:  0.7902,    Adjusted R-squared:  0.7847
F-statistic: 143.2 on 1 and 38 DF,  p-value: 1.864e-14
```



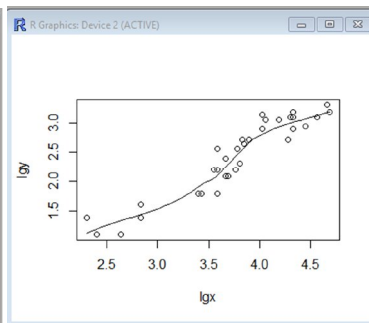
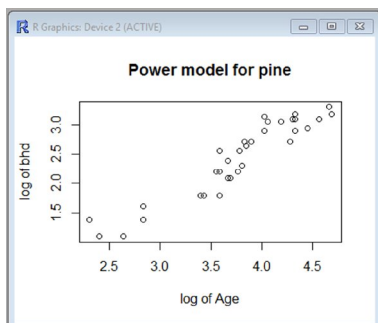
The p-values for intercept and slope are low, indicating that both of them are significantly different from zero.

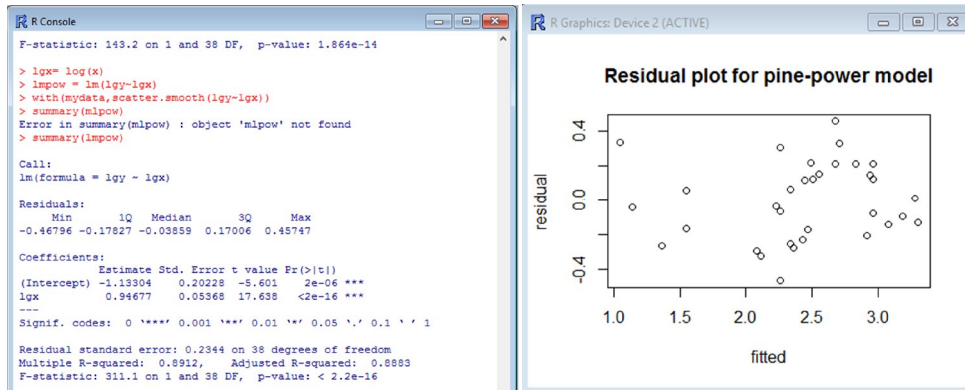
II. The power model

Functional form: $Y = aX^b$

Linearized model and parameters: $\ln(Y) = b_0 + b_1 \ln(X)$

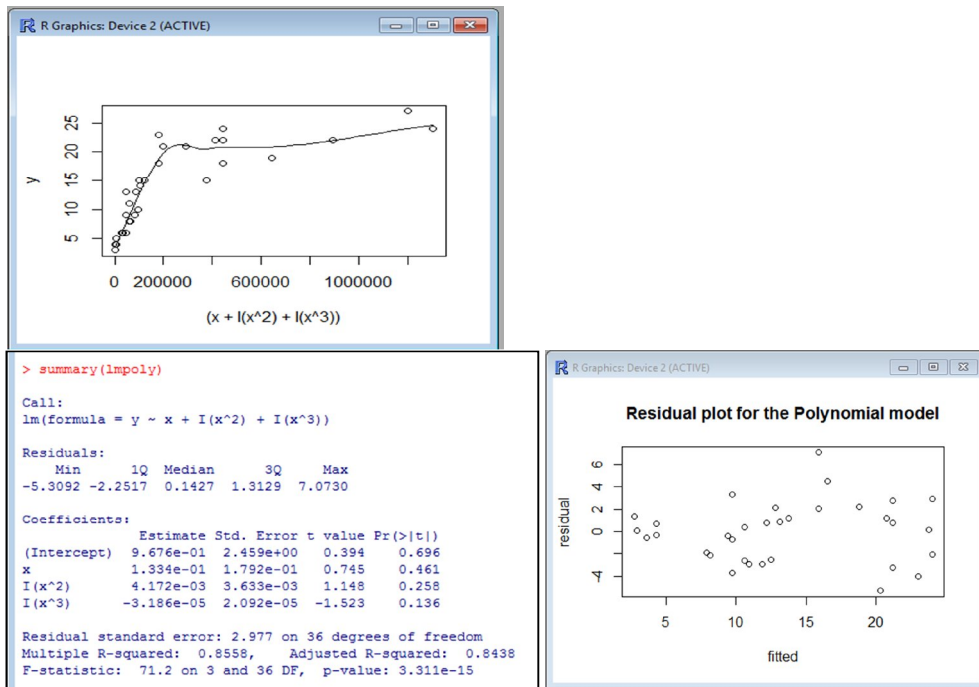
$a = eb_0$ $b = b_1$





III The polynomial model

The graph may have a cubic term.
So I use a polynomial of degree 3.



This model fails as the p-values are high meaning that the estimates are not significantly different from zero.

Diagnostics:

Model comparison

Model	Fitted plot	Adjusted R Square	Residual plot
Exponential model		0.7843	The best
Power model		0.8883	Better than the linear model
Polynomial model	Fails		

So I can select the power model. The equation for power model would be:

Functional form: $Y = aX^b$

Linearized model and parameters: $\ln(Y) = b_0 + b_1 \ln(X)$

$$a = e^{b_0} \quad b = b_1$$

From R, we have got

$$b_0 = -1.3304$$

$$b_1 = 0.94677$$

$$\ln(Y) = -1.3304 + 0.94677 * \ln(X)$$

⇒ (Raised to the power with base e)

$$\Rightarrow Y = 0.26 * X^{(0.94677)}$$

$$\Rightarrow BHD = 0.26 * (AGE)^{(0.94677)}$$

Example 1:

Age = 44

$$BHD = 0.26 * (44)^{(0.94677)} = 9.35$$

Example 2:

Age = 39

$$BHD = 0.26 * (39)^{(0.94677)} = 8.34$$

```

R Console
F-statistic: 143.2 on 1 and 38 DF, p-value: 1.864e-14

> lgx= log(x)
> lmpow = lm(lgy~lgx)
> with(mydata, scatter.smooth(lgy~lgx))
> summary(lmpow)
Error in summary(lmpow) : object 'lmpow' not found
> summary(lmpow)

Call:
lm(formula = lgy ~ lgx)

Residuals:
    Min       1Q   Median       3Q      Max
-0.46796 -0.17827 -0.03859  0.17006  0.45747

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) -1.13304    0.20228  -5.601   2e-06 ***
lgx           0.94677    0.05368  17.638  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2344 on 38 degrees of freedom
Multiple R-squared:  0.8912,    Adjusted R-squared:  0.8883
F-statistic: 311.1 on 1 and 38 DF, p-value: < 2.2e-16

```

Using SPSS:

Now I check my result from R with SPSS:

Iteration History ^a			
Iteration Number	Residual Sum of Squares	Parameter	
		A	B
1.0	253487.400	.260	.947
1.1	95494.116	.339	.960
2.0	95494.116	.339	.960
2.1	95048.580	.340	.956
3.0	95048.580	.340	.956
3.1	95048.557	.340	.956
4.0	95048.557	.340	.956
4.1	95048.557	.340	.956

Derivatives are calculated numerically.

Parameter Estimates				
Parameter	Estimate	Std. Error	95% Confidence Interval	
			Lower Bound	Upper Bound
A	.340	.008	.325	.355
B	.956	.005	.946	.967

Correlations of Parameter Estimates		
	A	B
A	1.000	-.995
B	-.995	1.000

ANOVA ^a			
Source	Sum of Squares	df	Mean Squares
Regression	2251273.443	2	1125636.721
Residual	95048.557	8785	10.819
Uncorrected Total	2346322.000	8787	
Corrected Total	612739.425	8786	

Dependent variable: bhd

a. $R^2 = 1 - (\text{Residual Sum of Squares}) / (\text{Corrected Sum of Squares}) = .845$.

SPSS performed 8 iterations and found the equation as :

$$\text{BHD} = 0.34 * (\text{AGE})^{(0.956)}$$

Discussion: Values I have retrieved with R are used as initial values in SPSS which seems to further improve the equation.