Model on Suicide Analysis and Prediction

"Ming Chu Cheng(Miranda), Wanying Li, Tal Jacobi, Jayant Bishnoi" 5/5/2021

Introduction Nowadays, people are intelligent and continuously seeking a high quality of life. At the same time, people are also faced with an enormous amount

STA 9750 Final Project

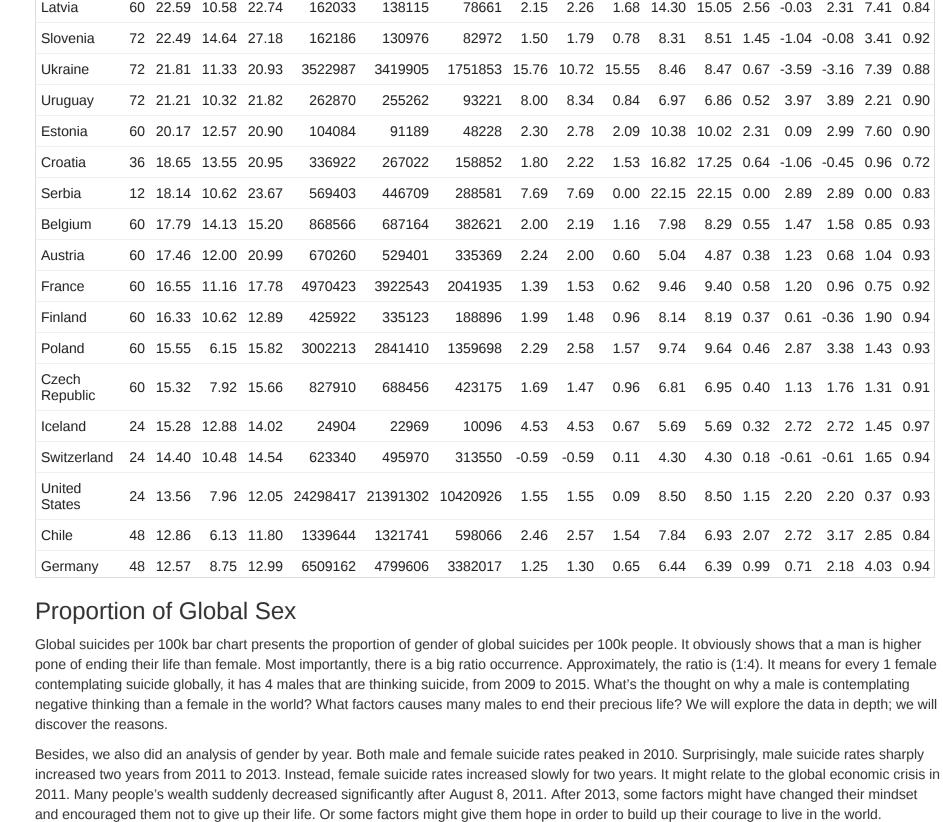
of pressure such as the cost of living, workload, relationships, even freedom.... etc. When people feel helpless, depressed and hopeless, they would contemplate ending their life in order to ease their unhappiness. The topic of suicide is always a complicated subject within society. Therefore, we try to use different features to explore the topic of suicide by using multiple linear regression and modeling analysis. We wish many people would contemplate the consequence of suicide and hope the result of suicide analysis could help some Governors' awareness and decision making through various means of advertising. There are in total of five data sources for our project analysis and prediction, including master.csv, gdp.csv, CPI_total.csv, has 27820 rows and 10 columns which contains historical data from 1985 to 2016. As we hope our data analysis has more diversification and

Unemployment_total_data.csv and Happiness.csv. All the resource comes from World Bank and Kaggle. The largest data sets of Suicides (master) accuracy, we find more component data sets and join them together to become our new data sets, encompassing historical data from 1985 to 2020. As the data in earlier years and some countries are relatively scarce and drop all NA, we select an analysis time frame from 2009 to 2015 which has 2940 observations and 28 columns which is called complete_join for our main data set. The complete_join data set comprises many global market index and rates. Consequently, we select different types of components to perform deep analysis and prediction, such as Suicides per 100k, our dependent variable. Population, GDP, CPI, Unemployment Rate, Social support, Freedom to make life choice, Perceptions of corruption, Confidence in national government, Democratic quality are our predictors. Order data cleaning procedures includes conversion of column types from the csv file, removal of unnecessary symbol, rename column name and pivot_longer the column and value and then join them together. Summary of the suicides rate of the countries from 2009 to 2015

suicide rate countries is in Europe or near Europe. The highest average rate of suicides per 100K people is the Lithuania country. Comparing with the United States, Lithuania has 2.5 times of average of per 100k people. But its population is less than the United States 6 times. Meanwhile,

comparing with others top nine countries, they also have 1.3 - 1.85 times higher comparing with the United States' average of per 100k people. It seems that most people in these counties, are feeling helpless and hopeless every day. Lithuania 72 34.17 19.06 32.22 239889 197766 1.64 13.80 13.57 2.34 0.63 3.54 7.08 0.91 113892 2.36 2.20 Belarus 36 25.05 12.29 25.75 749644 732032 372172 16.46 18.12 2.52 2.37 0.50 2.68 0.98 1.00 0.63 0.90 Hungary 60 23.96 13.74 27.27 788480 629552 386399 3.19 3.93 2.17 10.22 11.00 1.30 1.55 1.86 1.81 0.88

The summary chart provides some summary statistics of our analysis data set. It has contained 9 predictors and the dependent variables (suicide 100k pop) average, median and standard deviation. The table starts with descending by the suicide 100k pop. We can see that most of the top ten



14 -13 -

16 -

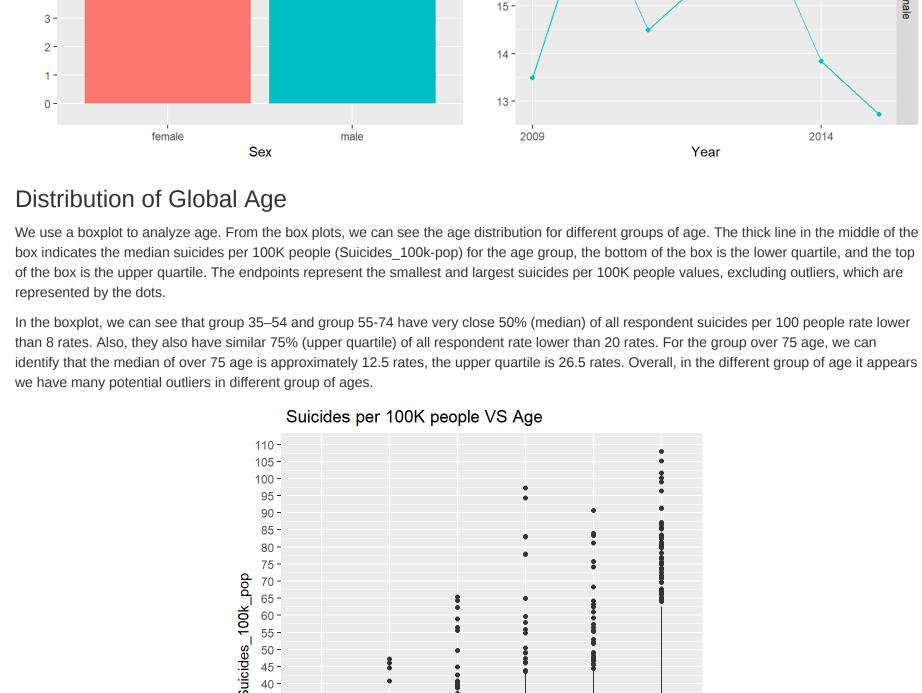
15 -

Global suicides per 100k by Sex

4.0 12 -11 -10 -100k 9 per

Trends Over Time, by Sex

Suicides per 100k 16 -



• yi = dependent variable • xi = explanatory predictors

Where, for i=n observations:

• β0 = y-intercept

Model Comparision

linear suicides:

log1p_linear_suicides:

0.4

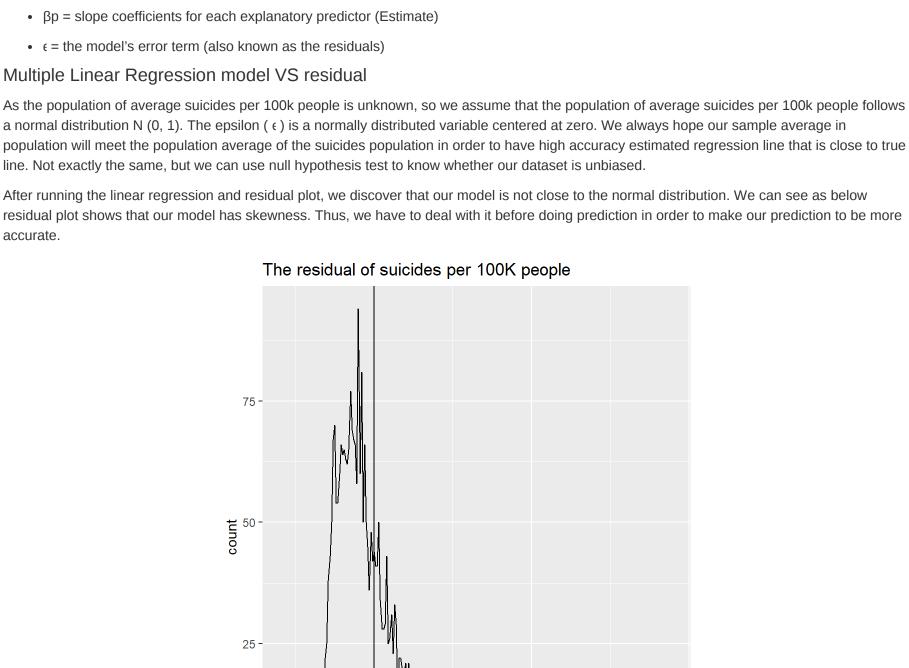
0.3 -

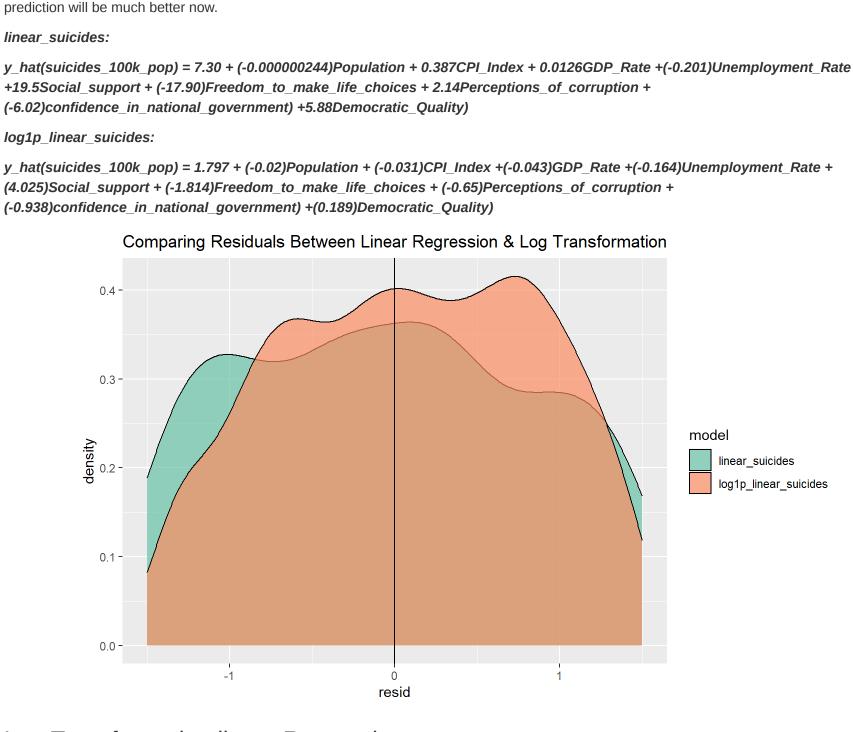
Term

35 -30 -25 -20 -15 -10 -

residual plot shows that our model has skewness. Thus, we have to deal with it before doing prediction in order to make our prediction to be more accurate. The residual of suicides per 100K people

75 -





suicide_resid

We try to use base log transformation, polynomial regression and log1p transformation to adjust our model. Please see the following plot. We can

see that the residual decreases frequently and the distribution of my fitted model moves much closer to the normal distribution N(0,1). Our

100



0.56 0.50 0.15 0.072 0.57 0.49 0.49 0.68 0.42 0.071 our 2.7. It is

log1p(Population)

0.5

 $log1p(Perceptions_of_corruption)$

Decision Tree

Component+Residual(log1p(Suicides_100k_pop) Component+Residual(log1p(Suicides_100k 0.5 1.0 1.5 2.0 2.5 3.0 $0.40 \quad 0.45 \quad 0.50 \quad 0.55 \quad 0.60 \quad 0.65$ 0.40 0.45 0.50 0.55 0.60 0.65 $log1p(Unemployment_Rate)$ log1p(Social_support) $log1p(Freedom_to_make_life_choices)$ _100k_pop) ponent+Residual(log1p(Suicides_100k_pop) lent+Residual(log1p(Suicides_100k_

0.3

We split the data 80/20, and use 80% of data to create a train set, and 20% to make predictions. According to the test data set, we can see the overall probability of suicide is 4%. If the sex is male, 50% of males will have a suicide probability of 4%. If Perception of Corruptions rate is smaller than 0.93, the chance of suicide is 2%. If Perception of Corruptions rate is larger than 0.93, the change of suicide is also 2%. If Democratic Quality rate is smaller than 0.75, the chance of suicide is 1%. If it is larger than 0.75, the chance of suicide is 1%. If Freedom to make life choice rate is

The decision tree model we build helps us identify the most important factors that affect suicide rate. Through the model, we can easily find out that

larger than or equal to 0.68, the chance of suicide is 0. Instead, if it is smaller than 0.68, the chance of suicide is 1%.

 $log1p(Confidence_in_national_government)$

 $log1p(CPl_Index)$

people are highly valued "perception of corruption", "democratic quality" and freedom to make life choices' The ultimate goal of our predicting model is to detect and visualize the relationship/patterns among all the variables, as well as determine whether an event (Suicide Y/N) will occur or ## Rows: 5,880 ## Columns: 7 ## \$ Confidence_in_national_government <dbl> 0.37, 0.37, 0.37, 0.37, 0.37, 0.37, 0.3... ## \$ Perceptions_of_corruption <dbl> 0.88, 0.88, 0.88, 0.88, 0.88, 0.8... ## \$ Social_support <dbl> 0.68, 0.68, 0.68, 0.68, 0.68, 0.6... ## \$ Democratic_Quality <dbl> -0.32, -0.32, -0.32, -0.32, -0.32... ## \$ Freedom_to_make_life_choices <dbl> 0.44, 0.44, 0.44, 0.44, 0.44, 0.4... <chr> "male", "male", "male", "male", "... ## \$ Sex <chr> "No", "No", "No", "No", "No", "No... ## \$ Suicide ## [1] 4704 33 100% yes -Sex = female - no 0.04 50% Perceptions_of_corruption < 0.93 0.02 8% **Democratic_Quality < 0.75** 0.01 3% Freedom_to_make_life_choices >= 0.68 0.00 0.02 0.01 0.00 0.01 50% 42% 5% 1% 2% Step-Wise Model Step-wise model which is a combination of forward and backward selection, consists of iteratively adding and removing predictors in our fitted model in order to find the subset of predictors in our data set resulting in the best performing model with lower prediction error. After running the step-wise model, our best performing model is Population, CPI Index, Unemployment Rate, Social Support, Freedom to make life choices, Confidence in national government and Democratic Quality. R squared is 0.084, Adj.R squared is 0.0829. all partial coefficient and p-value as below table. The most influence is "Freedom to make life choice" to Suicides_100K_pop. When "Freedom to make life choice" increases one unit, it controls Suicide_100k_pop decreases 18.714 units, holding other predictors constant. This outcome is 2.46 times of "confidence in national government" decreasing. Also, it is 88.27 times of "Unemployment Rate". The table also reveals the p-value of all predictors are smaller than 0.05 which means all rejects the null hypothesis = 0, all slope coefficient are statistically significant. It has strong evidence that these predictors can control the suicides 100k pop for a unit change. $y_hat(suicides_100k_pop) = 9.558 + (0)Population + (0.374)CPI_Index + (-0.212)Unemployment_Rate + (20.46)Social_suport + (0.374)CPI_Index + (-0.212)Unemployment_Rate + (20.46)Social_suport + (0.374)CPI_Index + (-0.212)Unemployment_Rate + (20.46)Social_suport + (-0.212)Unemployment_Rate + (-0.46)Social_suport + (-0.46)$ (-18.714)Freedom_to_make_life_choices + (-7.592)confidence_in_national_government) +(5.597)Democratic_Quality)

Estimate

9.558

0.374

-0.212

20.46

-18.714

-7.592

5.597

In the distribution of a standardized Residual plot with Suicides as below, shows that the distribution is not exactly a normal distribution and it has a little skewness. But comparing with our model which didn't do the log1p transformation, it is much closer to the population average N(0, 1). This

Distribution of Standardized Residuals

0

SE

2.821

0.053

0.048

3.095

2.123

1.536

0.399

0

T-statistic

3.389

-4.755

7.033

-4.431

6.611

-8.813

-4.944

14.041

6

P-value

< 0.001

< 0.001

< 0.001

< 0.001

< 0.001

< 0.001

< 0.001

< 0.001

Term

(Intercept)

Population

CPI_Index

Unemployment_Rate

Democratic_Quality

Freedom_to_make_life_choices

Confidence_in_national_government

Step-wise model density

prediction model would be more accurate.

Ŋ

0.3

0.2

0.1

0.0

-2

Density

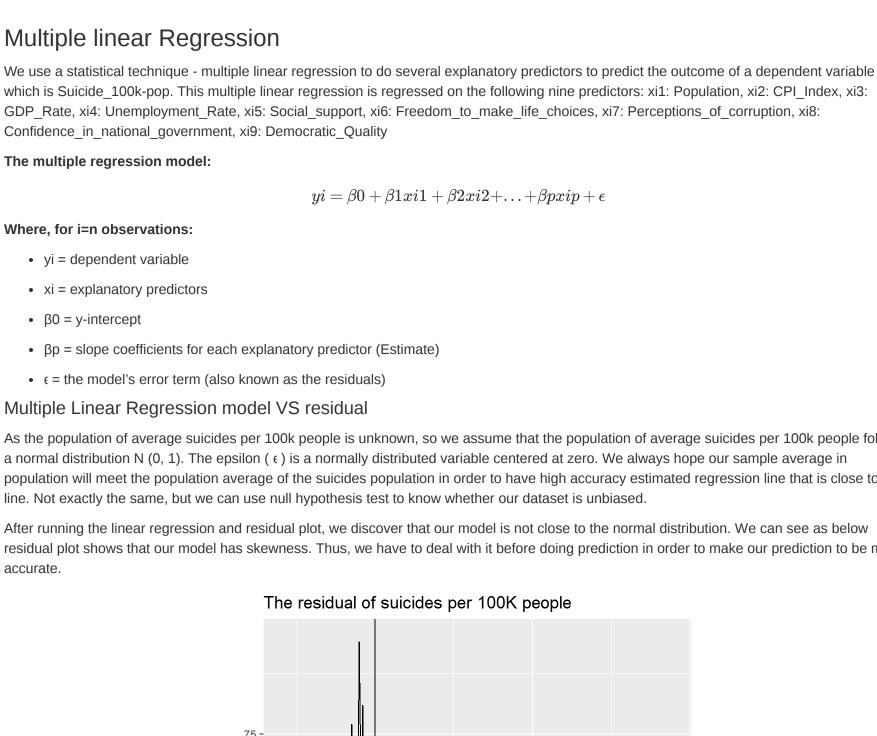
Social support

Conclusion In conclusion, our log1p regression model with low R square values can also be a good model. It is because some field of study have an inherently greater amount of unexplainable variation such as confidence in national government, Freedom to make life choices, etc. To explain human behavior generally is more difficult to predict than things that are physical. At least in our model analysis and prediction, we have tried our best to analyze and predict using diversification to make the outcome accurate.

2

0

sresid We find that global suicides are prone to males instead of females. The age between 35-75 is the in most negative thinking. They maybe faced with an enormous amount of pressure. Our prediction reveals that freedom and confidence in national government are the most influenced subjects on people to have an contemplate suicide, instead of say related with economy. Thus, we hope this result of suicide analysis and prediction could help some Governors' awareness that the country's wealth are their people. If they live in the world feeling hopeless and helpless, no one will be willing to contribute their ability and knowledge to their country to build up a prosperous country.



Age

model linear_suicides log1p_linear_suicides 0.1 Log Transformation linear Regression

After comparing with residual, we decide to use log1p transformation to adjust our model. The log1p model as below:

(4.025)Social_suport + (-1.814)Freedom_to_make_life_choices + (-0.65)Perceptions_of_corruption +

(-0.938)confidence_in_national_government) +(0.189)Democratic_Quality)

0.077

CPI_Index

Population

0.014

0.04

GDP_Rate

0.082

0.12

0.26

0.11

0.081

0.039

0.37

0.035

0.25

0.21

0.48

0.036

0.13

0.15

0.22

0.012

0.038

0.34

0.41

log1p(GDP_Rate)

log1p(Democratic_Quality)

0.16

0.40

0.16

0.02

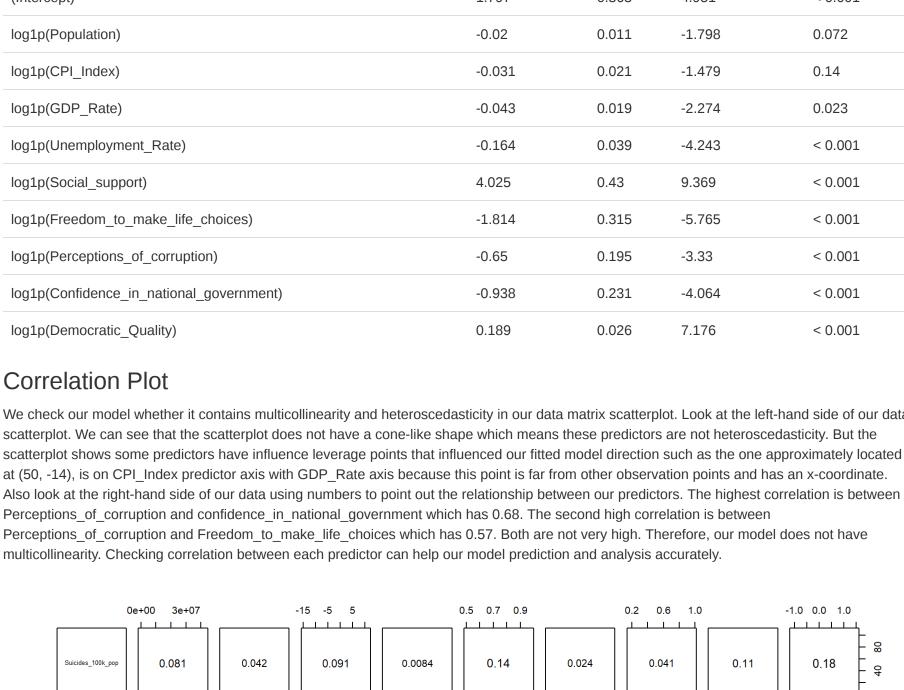
- 20

 $y_hat(suicides_100k_pop) = 1.797 + (-0.02)Population + (-0.031)CPI_Index + (-0.043)GDP_Rate + (-0.164)Unemployment_Rate + (-0.164$

life choices, holding other predictors constant. This indicates negative relationship between Freedom to make life choices and

The outcome of our model as following table shows that β6 = -1.814 will influence the suicides 100k pop for a unit change of Freedom to make

Suicide 100k pop which means one unit increase of Freedom to make life choices, Suicides 100k pop will decrease 1.814 unit holding other



0 40 80	0 20 40	0 10 20	0.4 0.7	0.1 0.4 0.7	
Variance Inflation	Factor (VIF)				
			·		
As per VIF for a regression mod	•				
model does not have multicollin	-	=	_		
a very low outcome. The highes	·	•		•	
the same as 2.5. This outcome	is matching the outcome of	f the previous correlation plot	t. We can be sure	that our fitted model does not	have
Multicollinearity.					
	Variance Inflation Factor(VIF) table				
			X		
	log1n(F	log1p(Population)			
	0 1 1	log1p(CPI_Index)			
		log1p(GDP_Rate)			
		log1p(Unemployment_Rate)			
		log1p(Social_support)			
		log1p(Freedom_to_make_life_choices)			
		log1p(Perceptions_of_corruption)			
		log1p(Confidence_in_national_government)			
		log1p(Democratic_Quality)			
	8-F(-		1.9		
Pocidual					
Residual					
In the following Component res	idual plots, we can see wh	ether each predictor has a lin	ear relationship to	Suicides per 100k people. Ti	ne blue dash
line is the best fit line. The pink	•	•	•		
with Suicides per 100k people 6		-			-
prediction the most.	3	01	,	,	
•			D. (
(c		mponent + Residual f			
مُّ مِنْ مِنْ مِنْ مِنْ مِنْ مِنْ مِنْ مِنْ				0 0 00 00	
	00 ok		6 6 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0		, (S) ,(S)
			S S	0 00 00 00	######################################
	OS SOS SOS SOS SOS SOS SOS SOS SOS SOS		8 o licid ←		•
)code (%)	No.		ns)d	8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8	8
000		9 8 8 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9	log1		
	nal()		nal(8
onent+Residual(log1p(Suicides_100k_pop))	onent+Residual(log1p(Suicides_100k_pop)		o o o o o o o o o o o o o o o o o o o		
£	A THE STATE OF THE		nt H		1
	onel -2		onel -		.