

Regression (I)

- 1 A simple linear model
- 2 Ordinary Least Squares (OLS)
- 3 Regression output
- 4 Draft No. 2

FiveThirtyEight

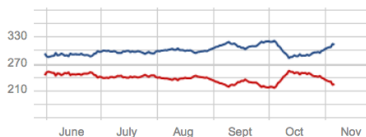
Nate Silver's Political Calculus

313.0
+14.0 since Oct. 30

**Electoral
vote**

225.0
-14.0 since Oct. 30

270 to win

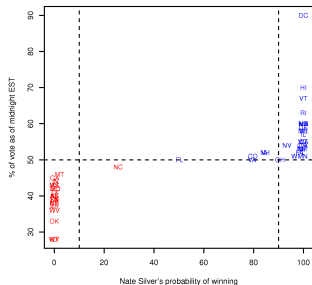


90.9%
+13.5 since Oct. 30

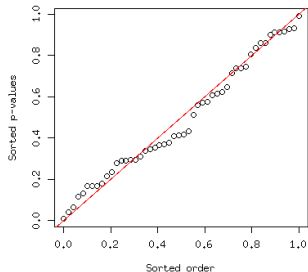
**Chance of
Winning**

9.1%
-13.5 since Oct. 30

50%

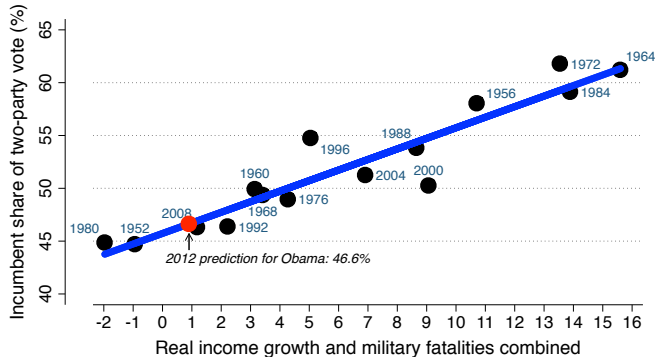


538 p-values



Obama's re-election prospect under bread and peace voting

October 26 2012 update based on projections of Oct-Nov 2012 conditions



Combination of real growth and fatalities weights each variable by its estimated coefficient.
 Estimated effects of fatalities on vote shares: -0.7% in 2008 (Iraq), -7.4% in 1968 (Vietnam),
 -9.7% in 1952 (Korea); negligible in 1964, 1976, 2004, 2012, and null in other years.
 Source: www.douglas-hibbs.com October 26 2012

To what extent can trust in government be predicted from variations in economic growth?

DV: Trust in government

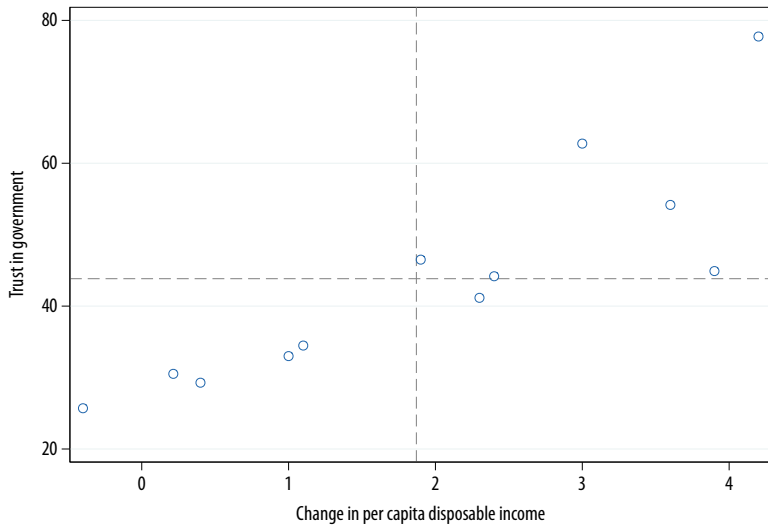
“Just about always/Most of the time”
(American National Election Studies)

IV: Economic performance

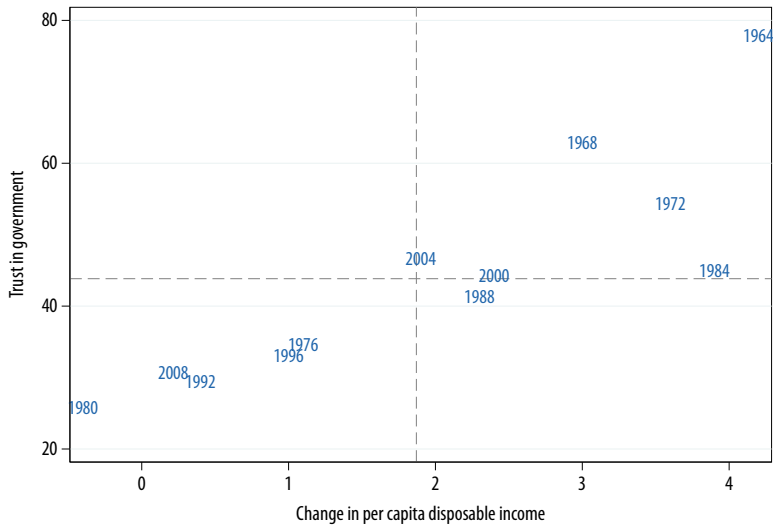
Change in disposable income per
capita (Bureau of Economic Analysis)

*Example kindly provided by John
Sides.*

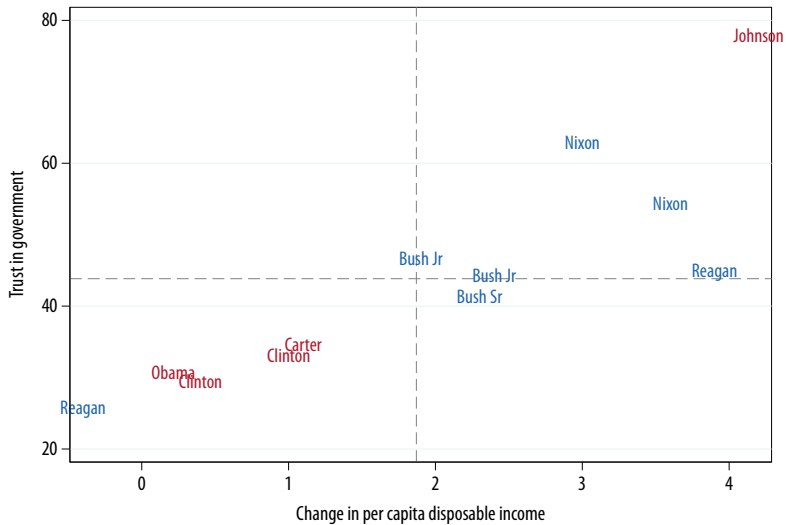




Dashed lines at averages. Pearson correlation $\rho = .86$ significant at $p < .01$.



Dashed lines at averages. Pearson correlation $p = .86$ significant at $p < .01$.



Dashed lines at averages. Pearson correlation $\rho = .86$ significant at $p < .01$.

Simple linear regression

Equations

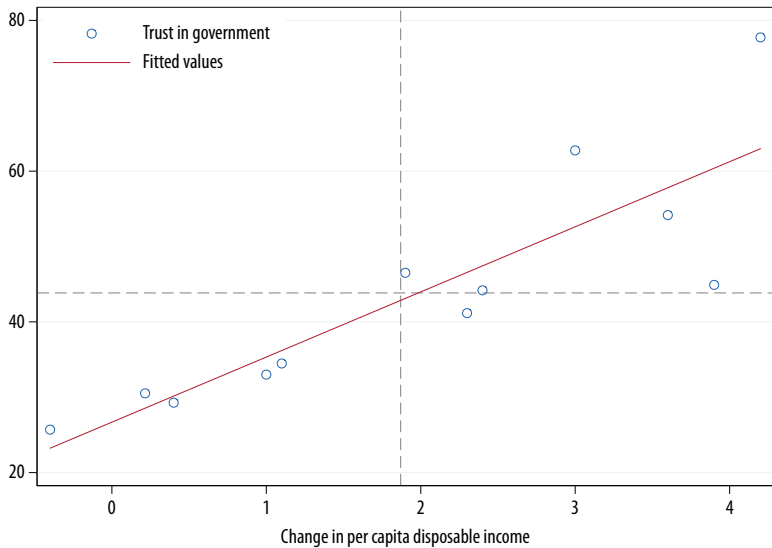
$$Y = \alpha + \beta X + \epsilon \quad \hat{Y} = \alpha + \beta X \quad \epsilon = Y - \hat{Y}$$

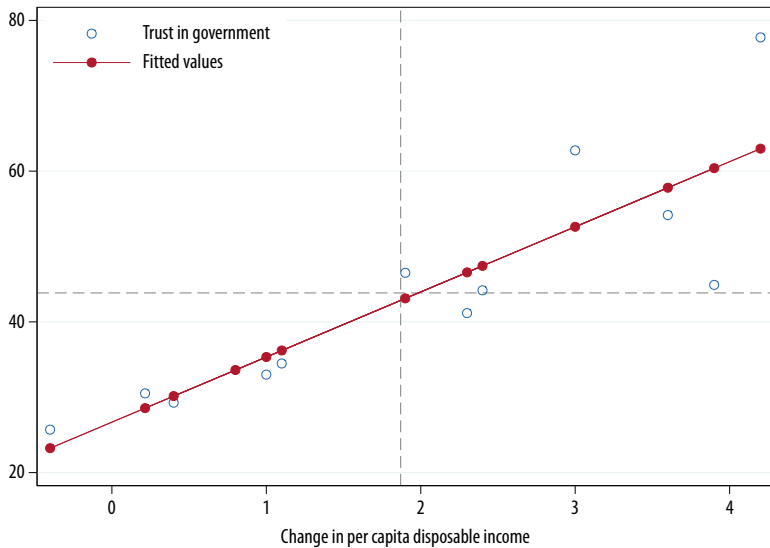
Parameters

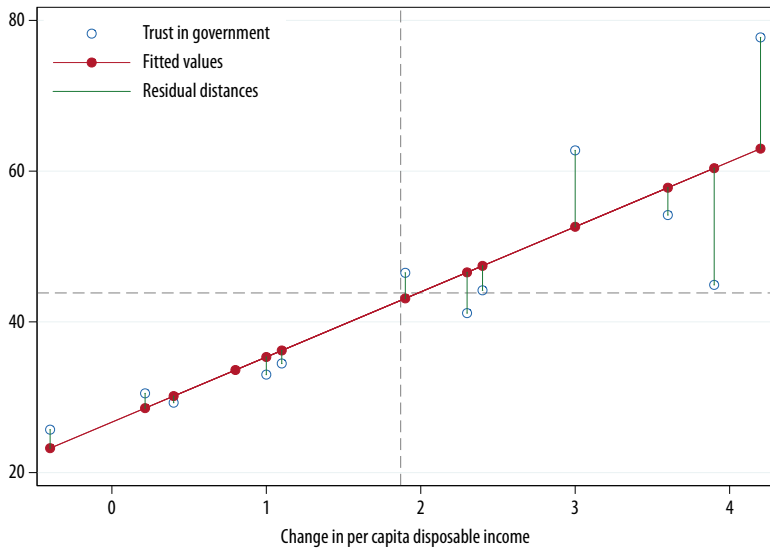
- Y is the dependent variable and \hat{Y} its predicted value
- X is the independent variable used as a predictor of Y
- α is the **constant** (intercept)
- β is the **regression coefficient** (slope)
- ϵ is the **error term** (residuals)

Warning

The model assumes a *linear, additive* relationship.







Ordinary Least Squares (OLS)

Error term

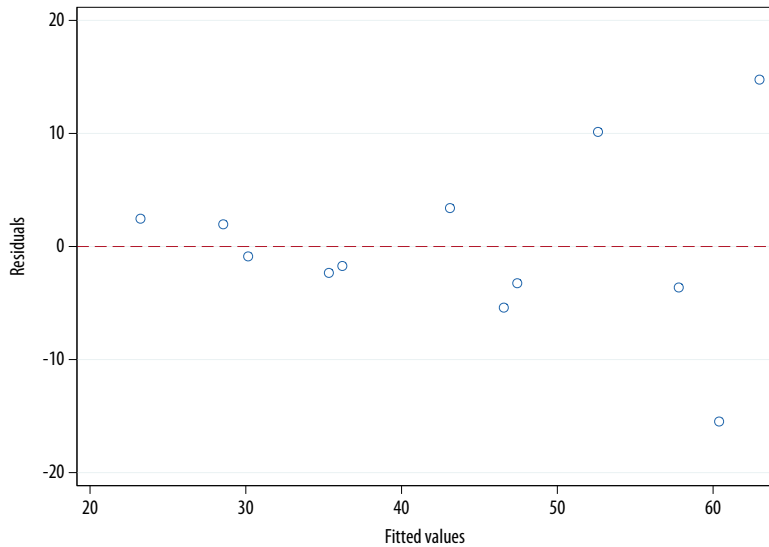
In a simple linear model $Y = \alpha + \beta X + \epsilon$, the regression coefficient β is calculated as to minimize the **residual sum of squares**

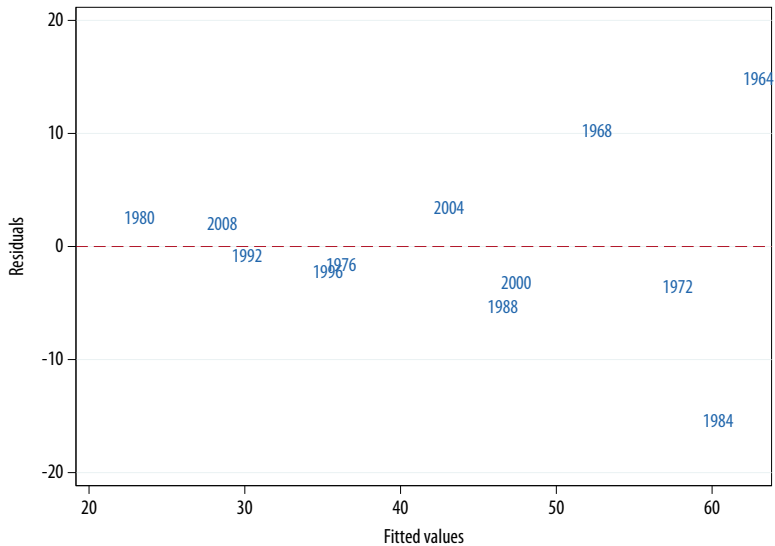
$$RSS = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^n \epsilon^2$$

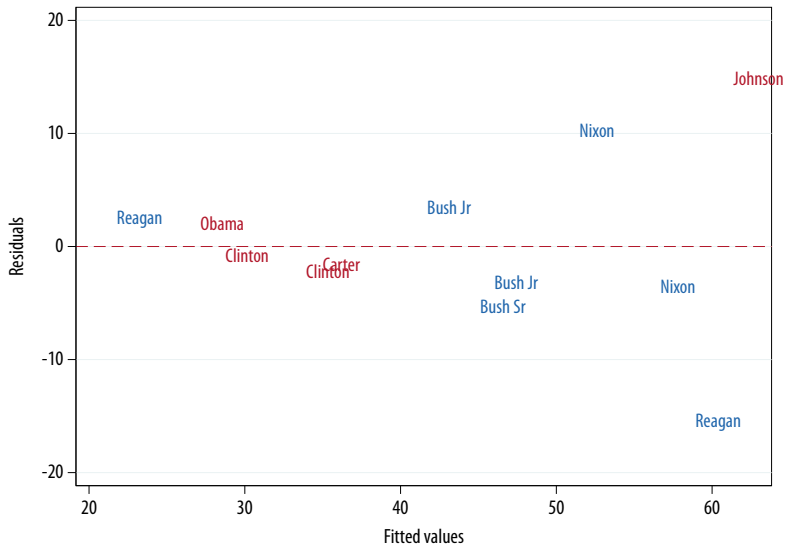
where $Y_i - \hat{Y}_i$ is the residual (or error term) of each observation.

Parameter estimation

$$\beta = \frac{\text{Cov}(X, Y)}{\text{Var}_X} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \quad \alpha = \bar{Y} - \beta \bar{X}$$







reg y x

. regress trust income

Source	SS	df	MS
Model	1908.80221	1	1908.80221
Residual	643.906248	10	64.3906248
Total	2552.70846	11	232.064405

Number of obs = 12
 F(1, 10) = 29.64
 Prob > F = 0.0003
 R-squared = 0.7478
 Adj R-squared = 0.7225
 Root MSE = 8.0244

trust	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
income	8.639373	1.586767	5.44	0.000	5.103836	12.17491
_cons	26.69501	3.888016	6.87	0.000	18.03197	35.35805

Top left: ANOVA table. Top right: model fit.
 Bottom: regression coefficients.

Interpretation of reg output

Number of observations N , significance test $H_0 : \beta = 0$, coefficient of determination R^2 .

regress brack income

Source	SS	df	MS
Model	1089.88225	1	1089.88225
Residual	945.588338	10	94.5588338
Total	2035.47059	11	203.224599

Sum of Squares	1089.88225
df	1
Prob > F	0.0003
R-squared	0.7478
Adj R-squared	0.7225
Root MSE	9.7244

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
income	8.439575	1.388787	6.08	0.0003	5.58388 11.2953
_cons	26.49581	5.888826	4.50	0.0003	14.67297 38.31865

Goodness of fit

$$R^2 = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \frac{\text{residual sum of squares}}{\text{total sum of squares}}$$

As predicted variance increases, $RSS \rightarrow 0$ and $R^2 \rightarrow 1$, indicating a more efficient fit.

Number of obs =	12
F(1, 10) =	29.64
Prob > F =	0.0003
R-squared =	0.7478
Adj R-squared =	0.7225
Root MSE =	9.7244

Sanity check

The most important statistic here is the actual number of observations in the model.

Interpretation of regression coefficients

A regression coefficient estimates the variation in Y predicted by a change in one unit of X (recall that $Y = \alpha + \beta X + \epsilon$)

regress trust income

Source	SS	df	MS	Number of obs =
Model	1780.88223	1	1780.88223	12
Residual	832.99028	10	83.299028	
Total	2613.87251	11	237.624773	

R-squared	0.6775	Prob > F	0.0000
Adj R-squared	0.6578	Prob > F	0.0000
F-statistic	21.39	Prob > F	0.0000
Root MSE	9.127		

	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
income	8.639373	1.586767	5.44	0.000	5.103836 12.17491
_cons	26.69501	3.888016	6.87	0.000	18.03197 35.35805

trust	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
income	8.639373	1.586767	5.44	0.000	5.103836	12.17491
_cons	26.69501	3.888016	6.87	0.000	18.03197	35.35805

- The **coefficient** is the slope β of the regression line and the **constant** is its intercept, the coordinate of origin $\alpha = \hat{Y}_{X=0}$.
- The **standard error**, t -value and p -value test whether the coefficient is significantly different from 0.

Draft No. 2

Univariate statistics

- Introduction
- Datasets
- Distributions
- Estimation

Assignment No. 1

corrected }
revised }
appended }

Bivariate statistics

- Significance
- Crosstabs
- Correlation
- Simple OLS

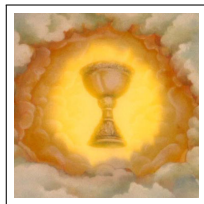
Assignment No. 2



Statistical modelling

- Regression
- Extensions
- Diagnostics
- Conclusion

Final paper



How to proceed

Revise Draft No. 1

- go through corrections
- remove technical content
- **rewrite until concision**

First steps

- between DV and IVs
- between two IVs

Write up **substantive results** as sentences; cite significance tests and other statistics in brackets, e.g. ($\rho = .7, p < .05$).

Thanks for your attention

Project

- Correct and improve first draft
- Finalize association tests and interpretations

Readings

- *Stata Guide*, Sec. 11
- *Making History Count*, ch. 4

Practice

- Replicate do-file
- Include simple OLS results in your second draft