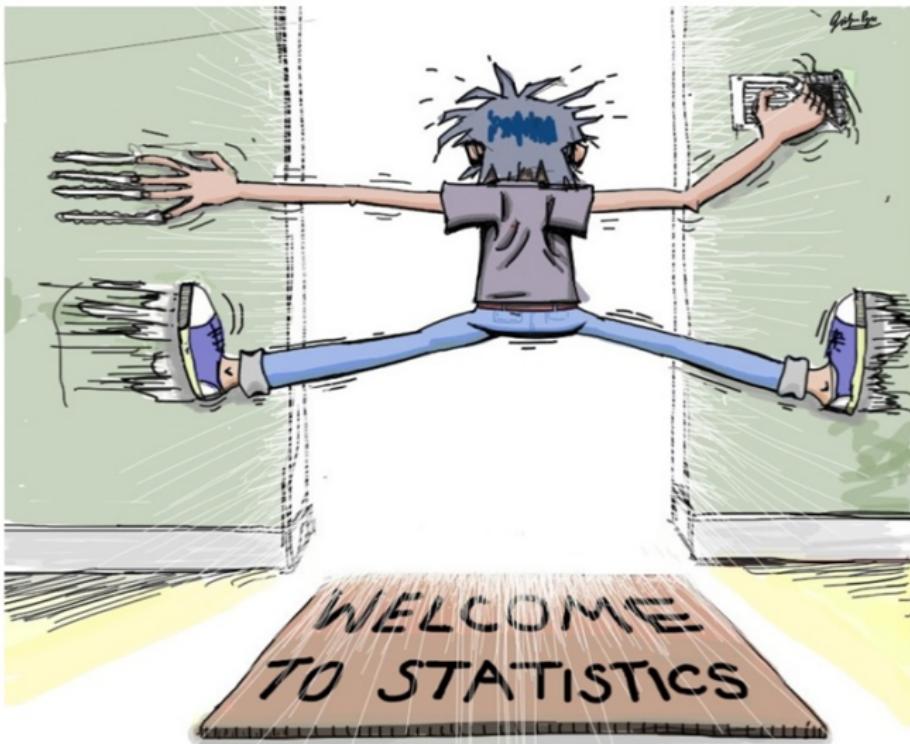


SMI 606: Introduction to Quantitative Data Analysis

Dr. Todd K. Hartman
Sheffield Methods Institute

Welcome!





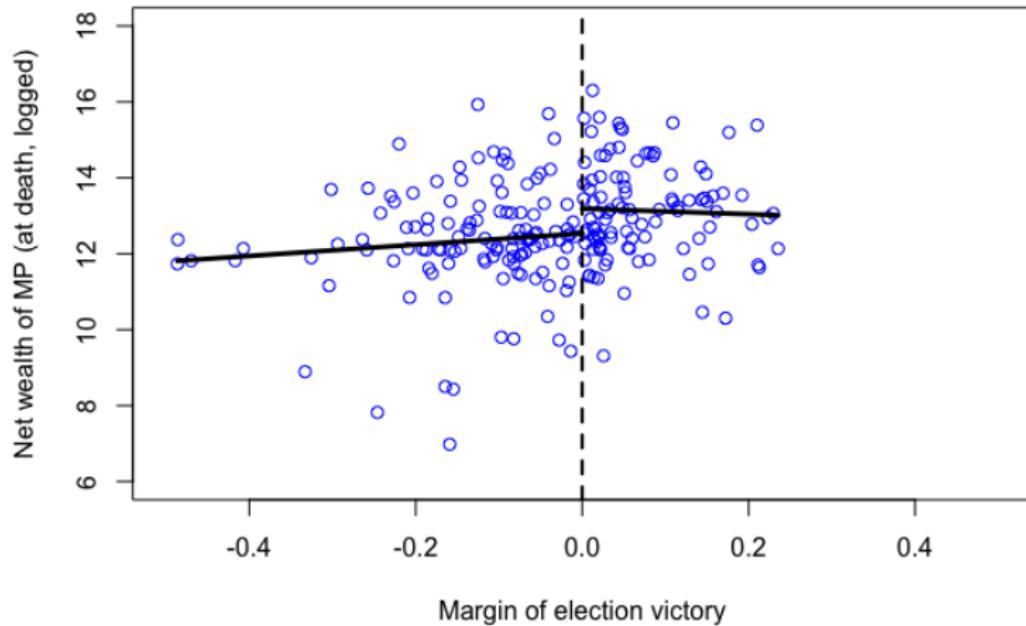
**KEEP
CALM
AND
SHOW US
THE DATA**

Roadmap: Week 1

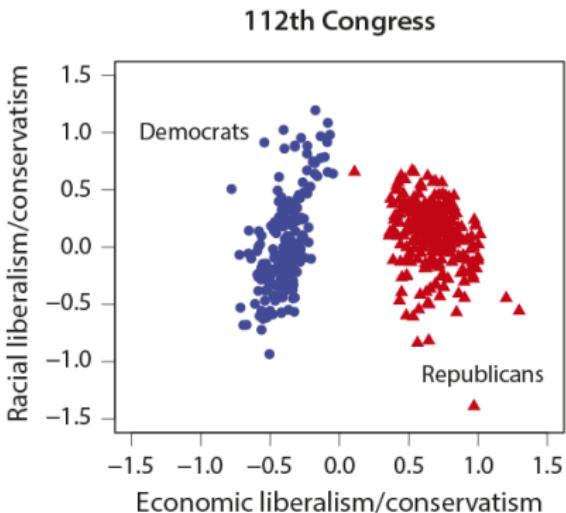
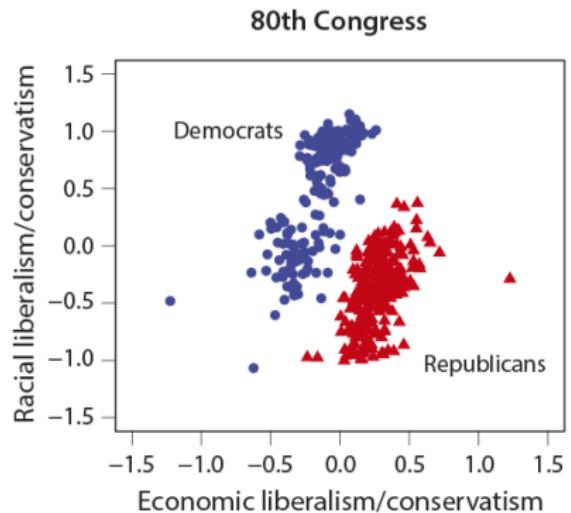
- Some cool things you'll learn!
- Module overview
- Lecture: How do we use data?
- Lab: Getting started in R

Regression

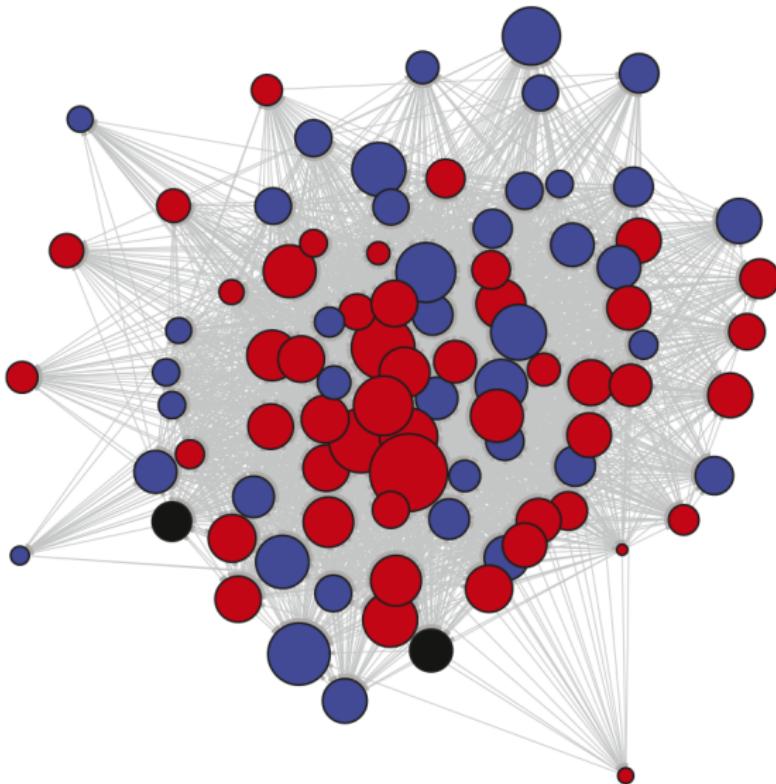
Conservatives



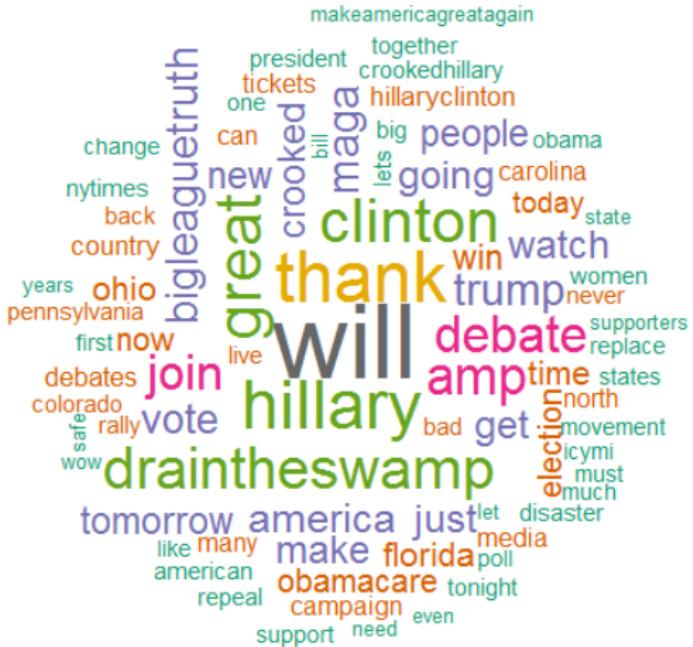
Cluster Analysis



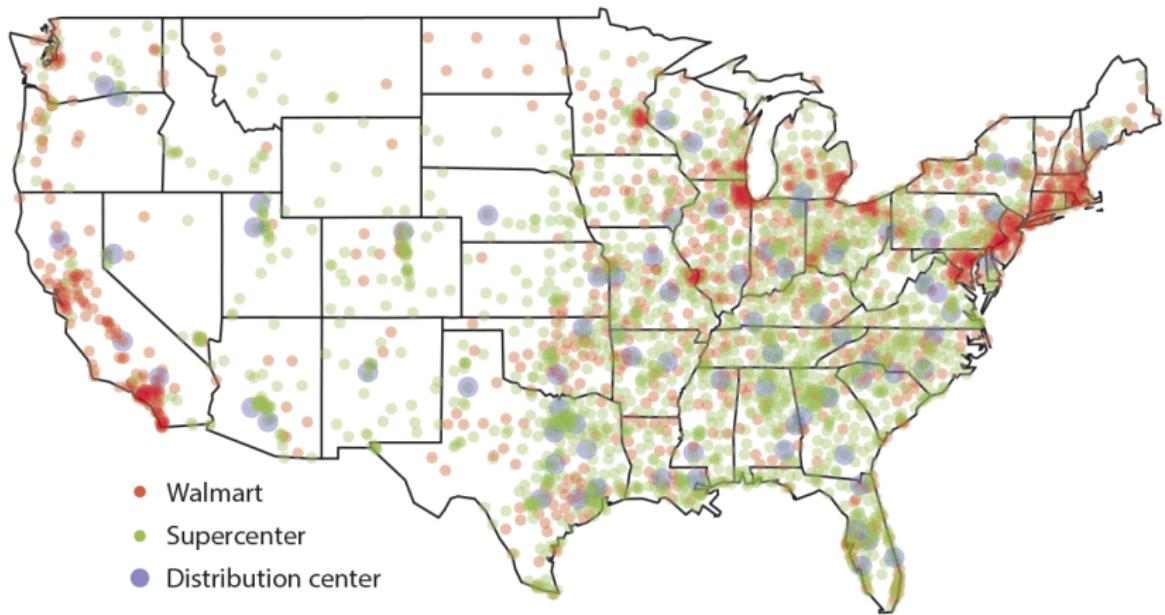
Social Network Analysis



Text Analysis



Spatial Analysis



How Do We Use Data?



① Description

- e.g., census demographics, election results, etc.

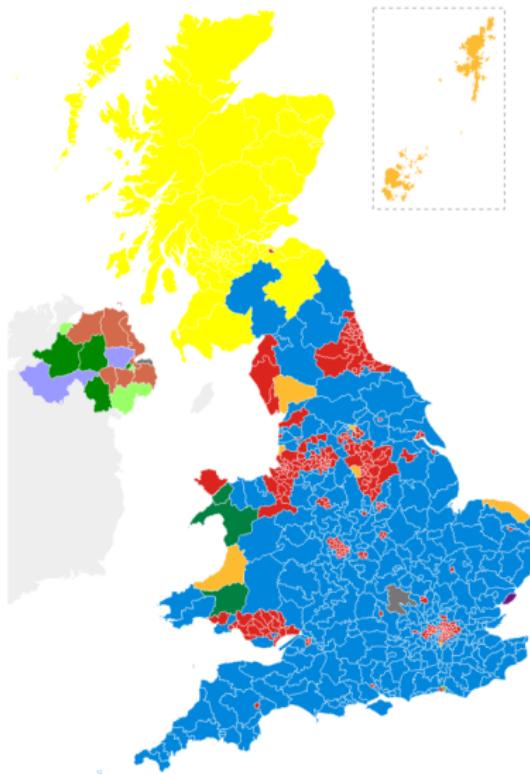
2015 UK General Election Results: Table

Party Standings

Party	Seats	+/-	Votes	%
Conservative	330	+28 ▲	11,300,109	36.8
Labour	232	-25 ▼	9,347,324	30.5
Scottish National Party	56	+50 ▲	1,454,436	4.7
Liberal Democrat	8	-49 ▼	2,415,862	7.9
Democratic Unionist Party	8	0 —	184,260	0.6
Sinn Fein	4	-1 ▼	176,232	0.6
Plaid Cymru	3	0 —	181,704	0.6
Social Democratic and Labour Party	3	0 —	99,809	0.3
UKIP	1	-1 ▼	3,881,064	12.6
Green	1	0 —	1,156,149	3.8
Alliance	0	-1 ▼	61,556	0.2
Respect the Unity Coalition	0	-1 ▼	9,989	0.0
Other	4	0 —	428,012	1.4

Note: The Speaker, who remains neutral, is included in other

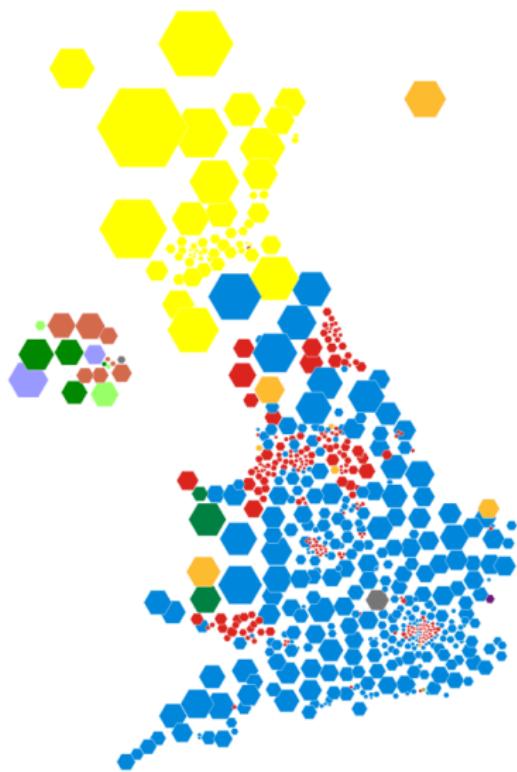
2015 UK General Election Results: Geography Map



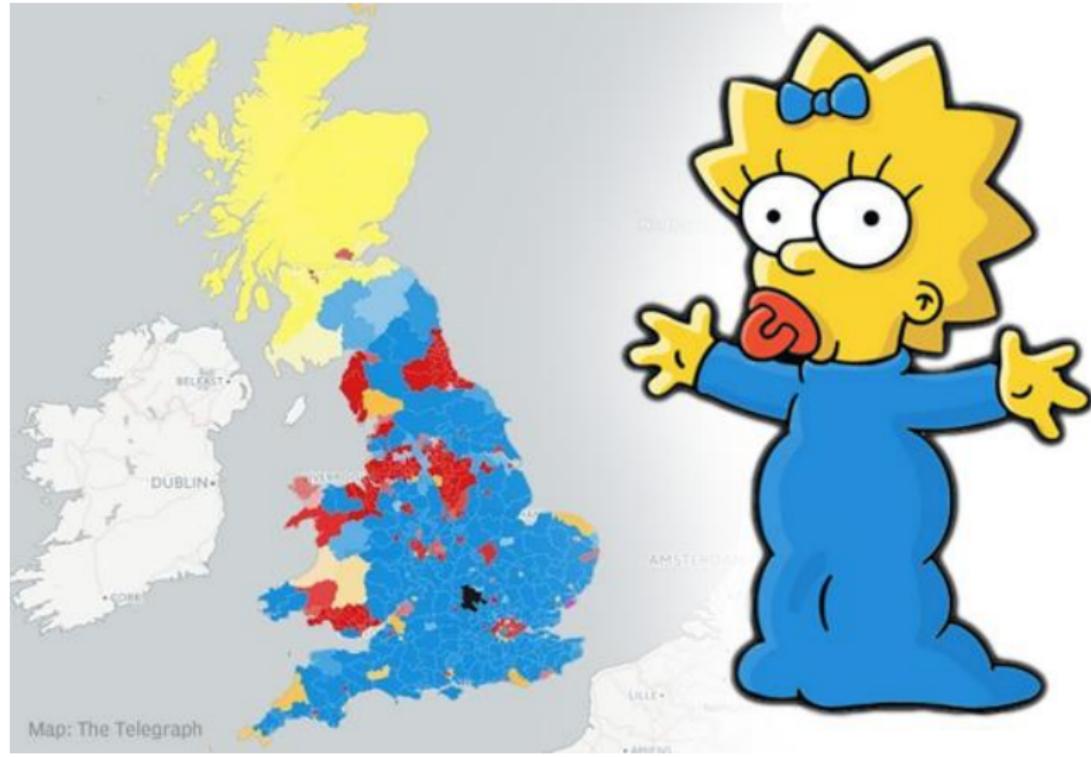
2015 UK General Election Results: Constituency Map



2015 UK General Election Results: Proportional Map



2015 UK General Election Results: Maggie Simpson Map



① Description

- e.g., census demographics, election results, etc.

How Do We Use Data?

- ① Description
 - e.g., census demographics, election results, etc.
- ② Explanation
 - Testing our beliefs about the way the world works

The Theory of Realigning Elections

- Major shift in party support in the electorate
 - Cyclical or generational (i.e., every 30 years)
 - Caused by a major event or issue
 - Witnessed in a 'critical' election

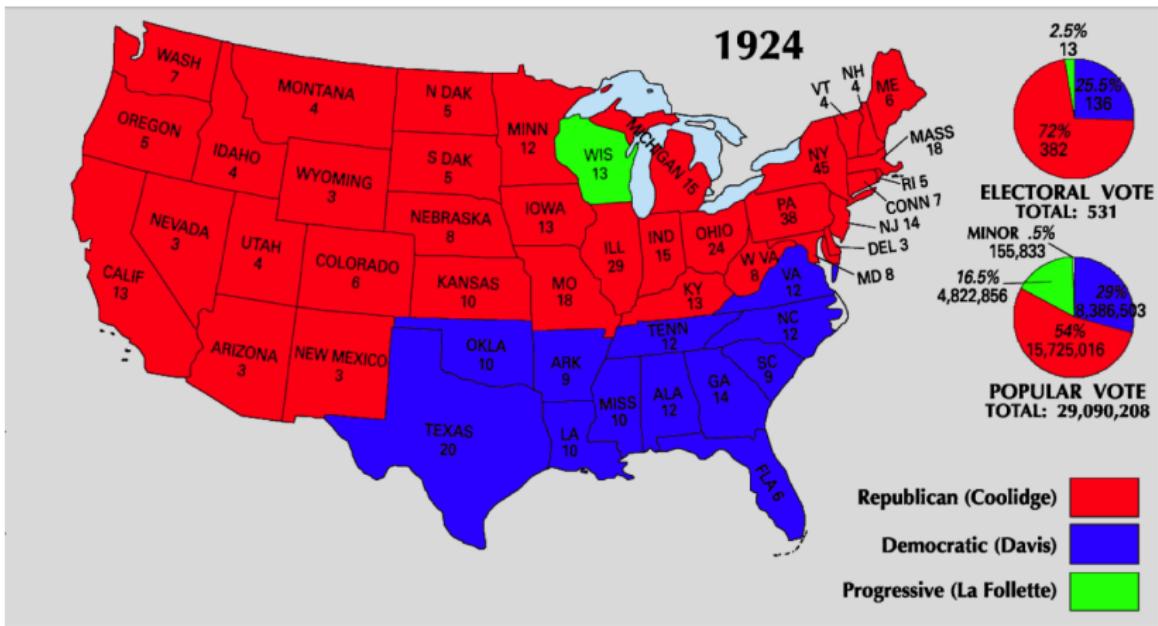


The Theory of Realigning Elections

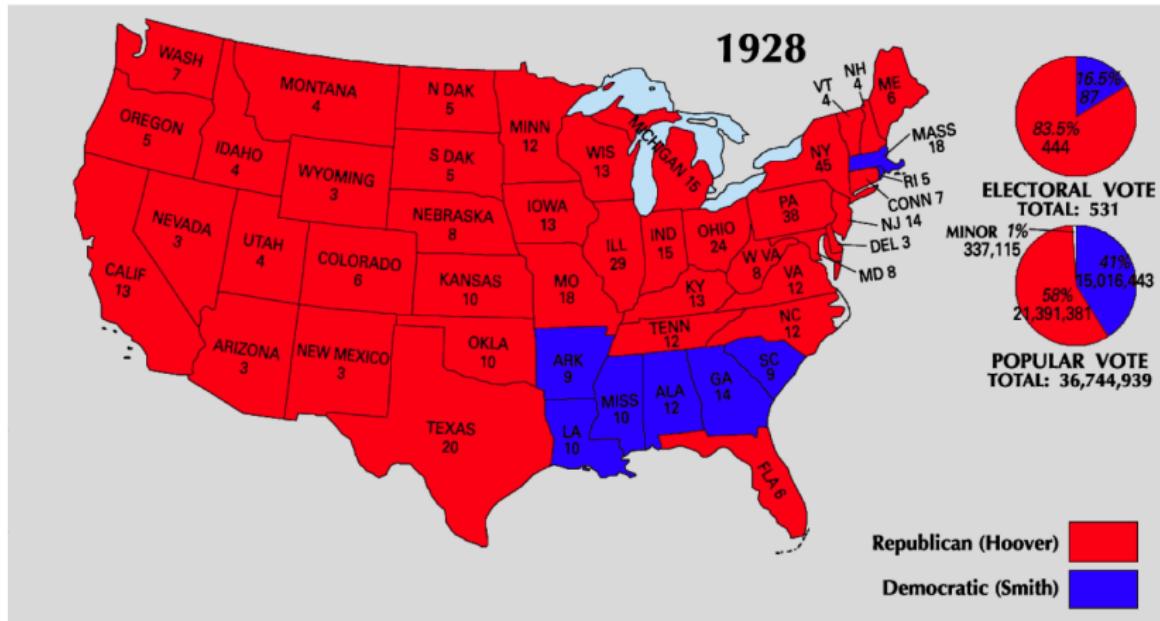
- Major shift in party support in the electorate
 - Cyclical or generational (i.e., every 30 years)
 - Caused by a major event or issue
 - Witnessed in a 'critical' election
- e.g., 1896 - 1928 U.S. Presidential Elections
 - Republicans won 7 of 9 elections (except '12 & '16 to Wilson)
 - What happened between 1928 and 1932?



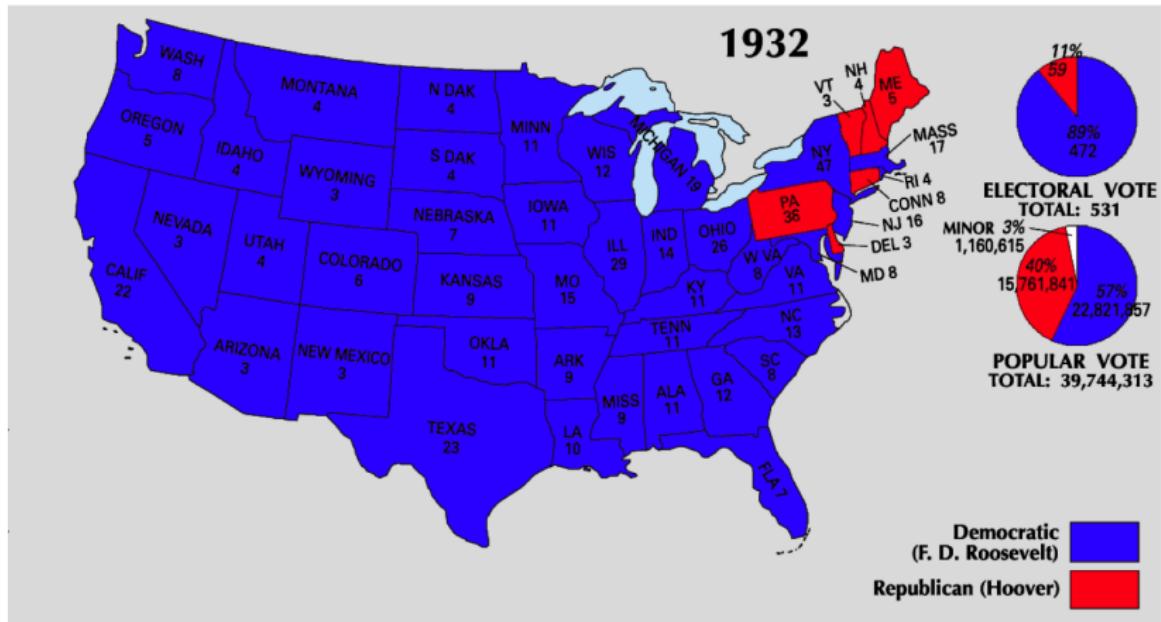
The Theory of Realigning Elections: 1924



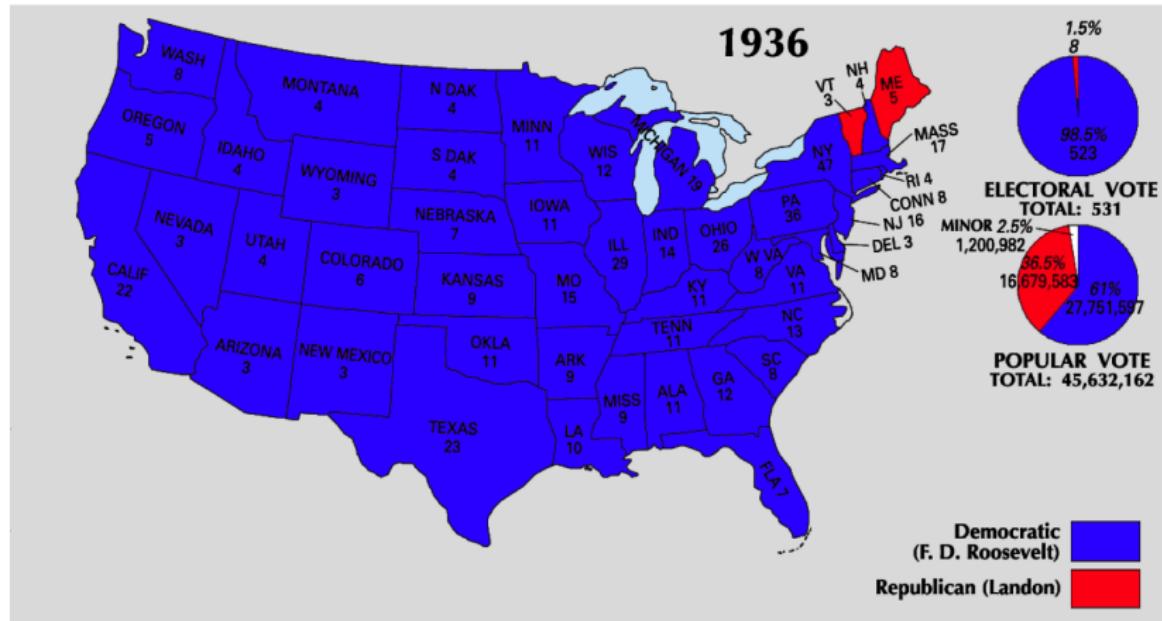
The Theory of Realigning Elections: 1928



The Theory of Realigning Elections: 1932



The Theory of Realigning Elections: 1936

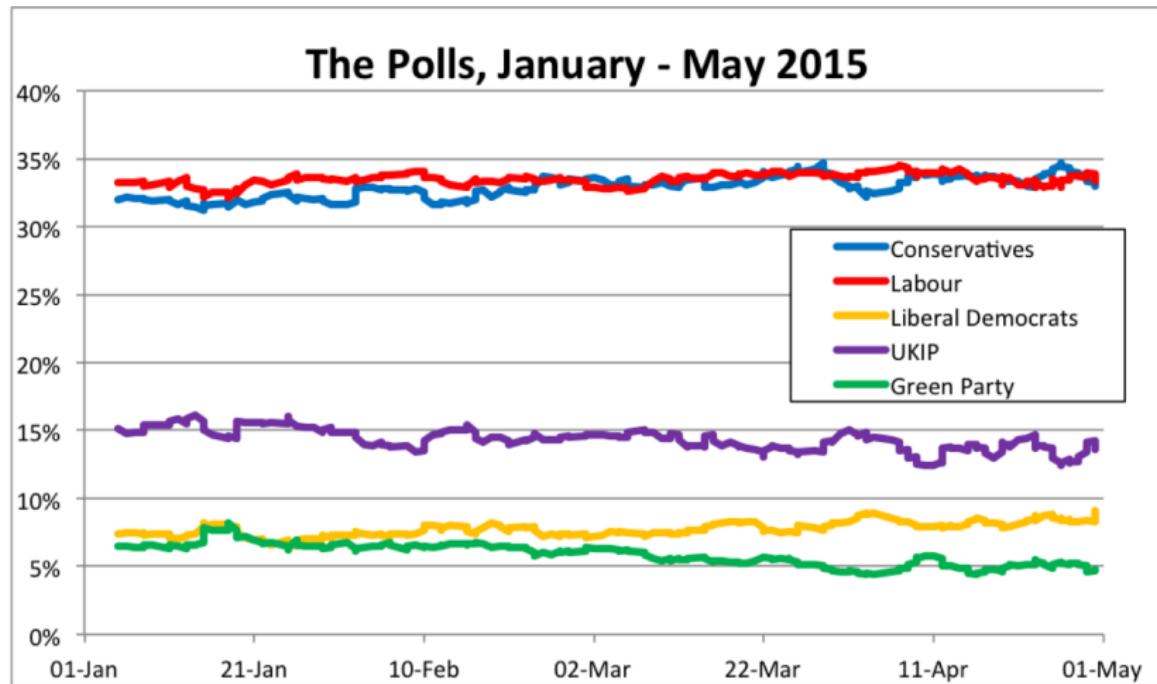


How Do We Use Data?

- ① Description
 - e.g., census demographics, election results, etc.
- ② Explanation
 - Testing our beliefs about the way the world works

- ① Description
 - e.g., census demographics, election results, etc.
- ② Explanation
 - Testing our beliefs about the way the world works
- ③ Prediction
 - e.g., who will win the upcoming election?

Forecasting the 2015 UK General Election: Polls

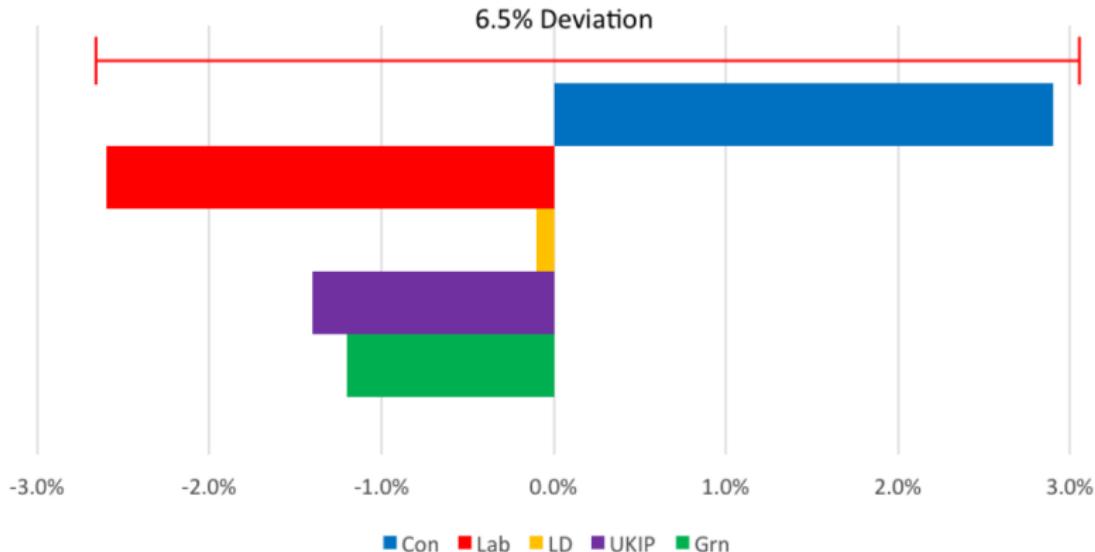


Forecasting the 2015 UK General Election: Models

	Five Thirty Eight	May 2015	Elections Etc	The Guardian	ACTUAL
Conservatives	278	273	285	273	331
Labour	267	268	262	273	232
Liberal Democrats	27	28	25	27	8
SNP	53	56	53	52	56
UKIP	1	2	3	3	1
Other	24	23	23	22	22
RESULT	HUNG	HUNG	HUNG	HUNG	CONSERVATIVE

Forecasting the 2015 UK General Election: Oops!

Deviation of the Pre-election Polls and Actual Vote Share



Measuring the World: Variables

- A *variable* is a quantitative measurement

- A *variable* is a quantitative measurement
- Two types of variables
 - ① Dependent (a.k.a. outcome) variable
 - In politics: e.g., vote choice, political attitudes, etc.
 - In criminology: e.g., crime rates, incarceration, etc.
 - In management: e.g., organizational performance, consumer behaviour, etc.

- A *variable* is a quantitative measurement
- Two types of variables
 - ① Dependent (a.k.a. outcome) variable
 - In politics: e.g., vote choice, political attitudes, etc.
 - In criminology: e.g., crime rates, incarceration, etc.
 - In management: e.g., organizational performance, consumer behaviour, etc.
 - ② Independent (a.k.a. predictor) variable(s)
 - e.g., gender, age, race/ethnicity, party affiliation, etc.

- A *dependent variable* is what we want to explain
i.e., its variation *depends* on one or more variables

- A *dependent variable* is what we want to explain
i.e., its variation *depends* on one or more variables
 - In equations, generally referred to as y ;

- A *dependent variable* is what we want to explain
i.e., its variation *depends* on one or more variables
 - In equations, generally referred to as y ;
 - *Independent variables* are used to explain a dependent variable
i.e., it will be used to predict the dependent variable
 - In equations, generally referred to as x ;
- Thus, $x_i \rightarrow y_i$ (x_i predicts or causes y_i)

① Nominal Variable (a.k.a. Categorical Variable)

- Attributes are defined by NAME only (i.e., no inherent ordering to the categories)
- e.g., country, gender, religion, etc.

① Nominal Variable (a.k.a. Categorical Variable)

- Attributes are defined by NAME only (i.e., no inherent ordering to the categories)
- e.g., country, gender, religion, etc.

② Ordinal Variable

- Categories are properly ORDERED
- But distance between values is not consistent
- e.g., Likert scales, polling questions, education, etc.

① Nominal Variable (a.k.a. Categorical Variable)

- Attributes are defined by NAME only (i.e., no inherent ordering to the categories)
- e.g., country, gender, religion, etc.

② Ordinal Variable

- Categories are properly ORDERED
- But distance between values is not consistent
- e.g., Likert scales, polling questions, education, etc.

③ Continuous or Pseudo-Continuous Variable

- Values are ordered and distance is consistent
- e.g., time, fiscal spending, etc.
- Interval scale (e.g., IQ test has no real 0)
- Ratio scale (e.g., Kelvin, where zero is meaningful)

Descriptive Statistics: The Mean

- Sum of values divided by total # of values

$$\bar{y} = \frac{y_1 + y_2 + y_3 + \dots + y_n}{n}$$

Mean

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

The Mean: An Example

- 2012 London Olympics: Men's 100m Final

Rank	Athlete	Country	Result
1	Usain BOLT	JAM	9.63
2	Yohan BLAKE	JAM	9.75
3	Justin GATLIN	USA	9.79
DQ	Tyson GAY	USA	9.80
5	Ryan BAILEY	USA	9.88
6	Churandy MARTINA	NED	9.94
7	Richard THOMPSON	TTO	9.98
8	Asafa POWELL	JAM	11.99

The Mean: An Example

- What is the mean 100m time?
- Sum of values divided by total # of values

$$\begin{aligned}\bar{y} &= \frac{9.63 + 9.75 + 9.79 + 9.80 + 9.88 + 9.94 + 9.98 + 11.99}{8} \\ &= \frac{80.76}{8} \\ &= 10.095\end{aligned}$$

An Outlier: Asafa Powell



Descriptive Statistics: The Median

- The middle value in a distribution
- Less sensitive to outliers (i.e., extreme values)
- To determine:
 - Order the values from low to high
 - For odd # of values, median is middlemost value
 - For even # of values, median is mean of two middle values

The Median: An Example

- What is the median 100m time?
- Recall, order the values and select middlemost value

Order the race times: 9.63, 9.75, 9.79, 9.80, 9.88, 9.94, 9.98, 11.99

$$\text{Median} = \bar{y} = \frac{9.80 + 9.88}{2} = 9.84$$

Comparing the Mean & Median: YouGov Poll

YouGov

YouGov®

Sample 1000 Adult Interviews
Conducted March 27 - 29, 2015
Margin of Error ±3.9%

1. Do you believe in love at first sight?

Yes	51%
No	35%
Not sure	14%

Comparing the Mean & Median: YouGov Poll

2. Do you believe there is such a thing as a soulmate?

Yes	69%
No	18%
Not sure	13%

Comparing the Mean & Median: YouGov Poll

3. Do you think a person has only one soulmate or do you think a person can have more than one?

Asked of people who believe there is such a thing as a soulmate

Only one soulmate	25%
More than one soulmate	33%
Not sure	10%
Does not believe in or isn't sure about soulmates	31%

Comparing the Mean & Median: YouGov Poll

4. How many soulmates do you think a person can have?

Asked of people who believe there is such a thing as a soulmate

mean	3,741,057
median	2

Descriptive Statistics: The Variance

- Sum of squared deviations from the mean divided by total # of values

$$s^2 = \frac{(y_1 - \bar{y})^2 + (y_2 - \bar{y})^2 + \dots + (y_n - \bar{y})^2}{n}$$

Variance

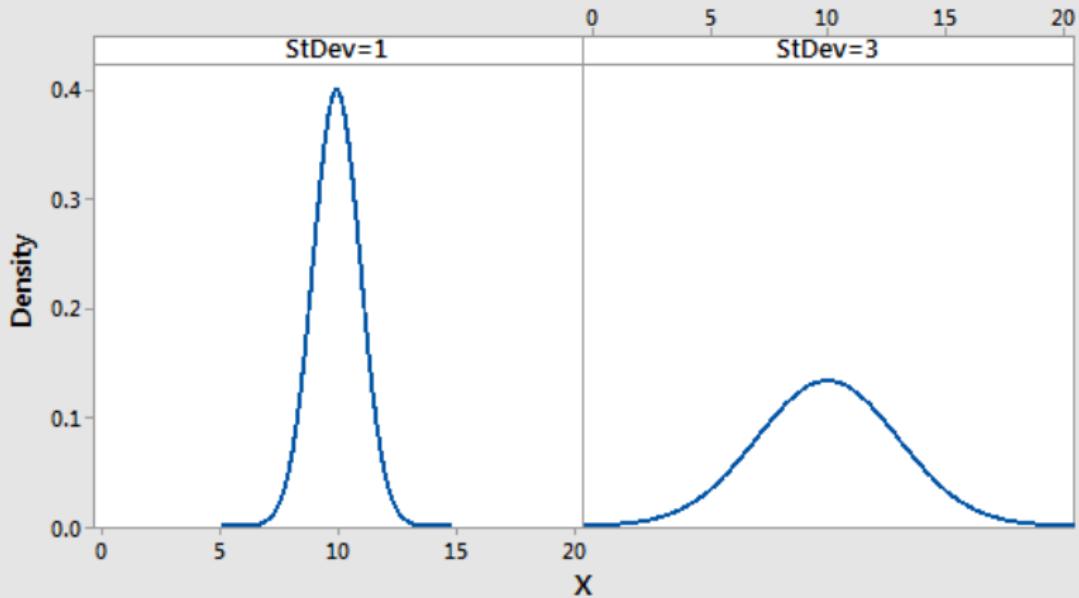
$$s^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

Standard Deviation

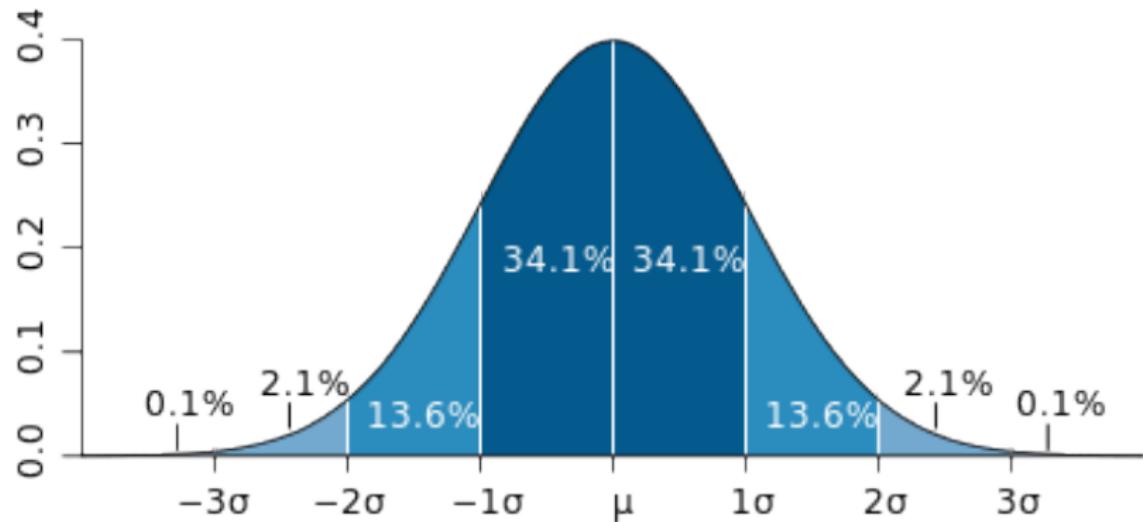
$$s = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}$$

Low and High Variation

Variation Within Samples Low and High Variability



Standard Deviation Distribution



It's All Greek to Me

GREEK ALPHABET

By Ben Crowder • bencrowder.net • Last modified 2 May 2012

Aα

ALPHA [a]
ἀλφα

Bβ

BETA [b]
βῆτα

Γγ

GAMMA [g]
γάμμα

Δδ

DELTA [d]
δέλτα

Εε

EPSILON [e]
εψιλόν

Ζζ

ZETA [dz]
ζήτα

Ηη

ETA [ɛ]
ἦτα

Θθ

THETA [tʰ]
θῆτα

Ιι

IOTA [i]
ἰῶτα

Κκ

KAPPA [k]
καππα

Λλ

LAMBDA [l]
λάμβδα

Μμ

MU [m]
μῦ

Νν

NU [n]
νῦ

Ξξ

XI [ks]
ξεῖ

Οο

OMICRON [o]
ὸ μικρὸν

Ππ

PI [p]
πεῖ

Ρρ

RHO [r]
ῥῶ

Σσς

SIGMA [s]
σιγμα

Ττ

TAU [t]
τᾶτ

Υυ

UPSILON [u]
ὐ ψιλόν

Φφ

PHI [pʰ]
φεῖ

Χχ

CHI [kʰ]
χεῖ

Ψψ

PSI [ps]
ψεῖ

Ωω

OMEGA [ɔ:
ῳ μέγα