# Week 3
# Contingency Tables & Regression Assumptions

## Applied Statistical Analysis/Quantitative Methods I

Jeffrey Ziegler, PhD

Assistant Professor in Political Science & Data Science
Trinity College Dublin

Fall 2024

# Roadmap for today

**Last class we learned:**

- What is a hypothesis test?
  - ▶ We will test hypotheses about population from sample
- Relationship between CI and null hypothesis testing
- Five steps of hypothesis testing
- Types of errors
- Discussion of 1-sided/2-sided tests

Sometimes, our specific aim is to understand if *X* causes *Y*

- Compare two independent samples
- So, today...

# Transition: From causation to association

**This class:**

- Estimate if two variables are dependent
  - ► Chi-squared test of independence
  - ► Standardized residuals
- Estimating correlation
  - ► Does variation in one explain variation in another
- Bivariate regression
  - ► Assumptions
  - ► Estimation (i.e. drawing the "best" line through data)

We have 2 categorical variables, any relation?

- If we have three samples, our data might look like this

| Variable 1 (Outcome or response) | Variable 2 (Sample or grouping) |
|---|---|
| 1 | 1 |
| 2 | 0 |
| 3 | 1 |
| 5 | 2 |
| 3 | 2 |
| 2 | 0 |
| 4 | 0 |
| ⋮ | ⋮ |

# CROSS-TABS: THE BASICS

Assume we have two variables that are nominal

- Party-ID and racial/ethnicity

- **WARNING:** Calculations are different when variables are ordinal … especially if BOTH are ordinal (see end of chapter 8 in A&F)
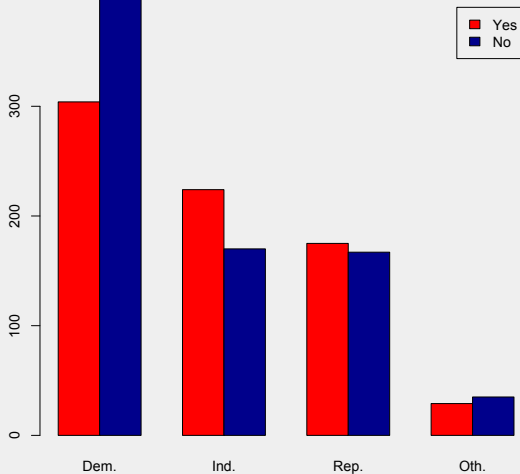
We will use a contingency table, which is usually (at least by me) referred to as a "cross-tabulation"

# Example: Preferences over abortion policy

"Please tell me whether or not you think it should be possible for a pregnant woman to obtain a legal abortion if … the family has a very low income and cannot afford any more children?"

|              | Yes | No  | Total |
|--------------|-----|-----|-------|
| Democrats    | 304 | 398 | 702   |
| Independents | 224 | 170 | 394   |
| Republicans  | 175 | 167 | 342   |
| Other        | 29  | 35  | 64    |
| Total        | 732 | 770 | 1502  |

# PLOT: DATA COMES FROM 1972 GSS



**Not an experiment! What do we do?**

# We want **conditional distribution** of outcome

## How do we get that?

- What is the distribution of the outcome variable conditioned on the independent/grouping variable? Proportions?!

|              | Yes  | No   |
|--------------|------|------|
| Democrats    | 0.43 | 0.57 |
| Independents | 0.57 | 0.43 |
| Republicans  | 0.51 | 0.49 |
| Other        | 0.45 | 0.55 |

Note that the rows total to 100% because the rows indicate the independent variable

- We might have the columns add up to 100% if that was our explanatory variable

# Chi-square test of independence

*Statistical independence*: *Two variables are statistically independent if the conditional distributions of the **population** are identical across categories*[1]

Note: Events can be *logically* or *physically independent* but still *statistically dependent*.

- Let $A$ = scored above 90% on entrance exam,

- and $B$ = attends Trinity

These are logically independent (neither one implies the other), but statistically quite dependent, because $Pr(A|B) > Pr(A)$

---

[1] Since $Pr(A \cap B) = Pr(A|B)Pr(B)$, if $A$ is independent of $B$, then $Pr(A \cap B) = Pr(A)Pr(B)$. If this holds, then $B$ is also independent of $A$ since $Pr(B|A) = \frac{Pr(A \cap B)}{Pr(A)} = \frac{Pr(A)Pr(B)}{Pr(A)} = Pr(B)$

# Chi-square test: The intuition

$H_0$: The variables are statistically independent

$H_a$: The variables are statistically dependent

We are going to calculate a test-statistic (the $\chi^2$ statistic) that is distributed according to the $\chi^2$ distribution

$f_{observed} = f_o$ = observed frequency = the raw count (NOT THE %)

$f_{expected} = f_e$ = what we would expect for independent samples

$$= \frac{\text{Row total}}{\text{Grand total}} \times \text{Column total}$$

If $H_0$ is true, then we would expect $f_{observed} = f_{expected}$

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

# Example: Chi-square test

|              | Yes         | No          | Total |
|--------------|-------------|-------------|-------|
| Democrats    | $f_o$=304   | $f_o$=398   | 702   |
| Independents | $f_o$=224   | $f_o$=170   | 394   |
| Republicans  | $f_o$=175   | $f_o$=167   | 342   |
| Other        | $f_o$=29    | $f_o$=35    | 64    |
| Total        | 732         | 770         | 1502  |

# Chi-square test: Example

$$f_{1e} = \frac{\text{row total}}{\text{grand total}} * \text{column total} \tag{1}$$

$$= \frac{702}{1502} * 732 = 342.12 \tag{2}$$

|  | Yes | No | Total |
|---|---|---|---|
| Democrats | $f_o$=304 | $f_o$=398 | 702 |
|  | $\mathbf{f_e = 342.12}$ | $f_e = 359.88$ |  |
| Independents | $f_o$=224 | $f_o$=170 | 394 |
|  | $f_e = 192.12$ | $f_e = 201.98$ |  |
| Republicans | $f_o$=175 | $f_o$=167 | 342 |
|  | $f_e = 166.67$ | $f_e = 175.33$ |  |
| Other | $f_o$=29 | $f_o$=35 | 64 |
|  | $f_e = 31.19$ | $f_e = 32.81$ |  |
| Total | 732 | 770 | 1502 |

$$\chi^2 \;=\; \sum \frac{(f_o - f_e)^2}{f_e} \tag{3}$$

$$=\; \frac{(304 - 342.12)^2}{342.12} + \frac{(398 - 359.88)^2}{359.88} + \ldots + \tag{4}$$

$$\approx\; 19.79 \tag{5}$$

Is the $\chi^2$ statistic "big enough?"

- We are going to need to calculate the degrees of freedom

- This is skewed right and strictly positive

- Always use the upper-tail (no $\times$ 2)

# CALCULATING P-VALUES FOR CHI-SQUARED TESTS

- Frequency $\geq 5$ for all cells

- $df = (rows - 1)(columns - 1)$

- In R: `pchisq(`$\chi^2$`, df = (rows-1)(columns-1), lower.tail=FALSE)`

# Example: Preferences over abortion policy

What is df? $df = 3$

How should we get our p-value?

$$p\text{-value} = \text{pchisq(19.79, df=3, lower.tail=F)} \tag{6}$$
$$= 0.00019 \tag{7}$$

# Standardized residuals

Now, we have evidence that the two variables are not independent

- Where does the deviation from independence take place?

- Why did we reject the null?

- What does it mean?

# Standardized residuals

How far away is each observed value from "expectation"

We need to find the **adjusted residual**:

$$z = \frac{f_{observed} - f_{expected}}{se} = \frac{f_{observe} - f_{expected}}{\sqrt{f_e(1 - \text{row prop.})(1 - \text{column prop.})}}$$

- The denominator is the standard error of the quantity $f_o - f_e$ under the null hypothesis

# Example: Calculating standardized residuals

|                | Yes              | No                | Total |
|----------------|------------------|-------------------|-------|
| Democrats      | $f_o$=304        | $f_o$=398         | 702   |
|                | $f_e = 342.12$   | $f_e = 359.88$    |       |
| Independents   | $f_o$=224        | $f_o$=170         | 394   |
|                | $f_e = 192.12$   | $f_e = 201.98$    |       |
| Republicans    | $f_o$=175        | $f_o$=167         | 342   |
|                | $f_e = 166.67$   | $f_e = 175.33$    |       |
| Other          | $f_o$=29         | $f_o$=35          | 64    |
|                | $f_e = 31.19$    | $f_e = 32.81$     |       |
| Total          | 732              | 770               | 1502  |

$$z_{11} = \frac{304 - 342.12}{\sqrt{342.12(1 - \frac{702}{1502})(1 - \frac{732}{1502})}} \approx -2.395$$

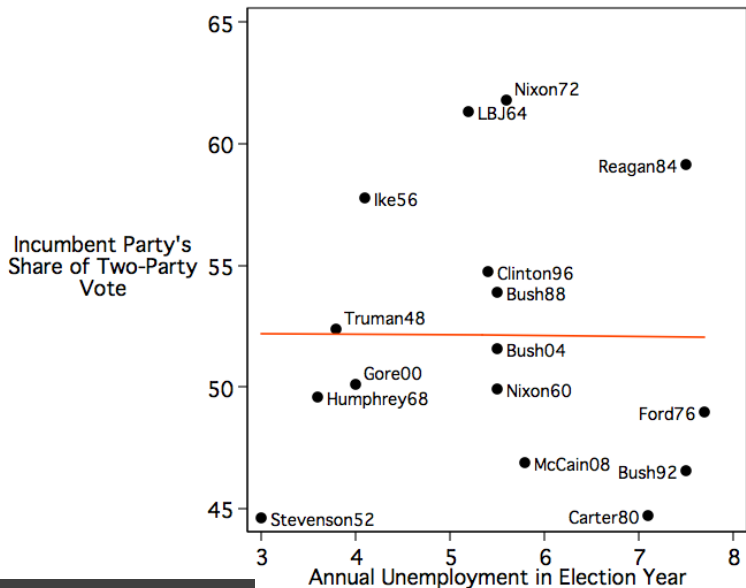We should always visually inspect data using... scatterplots

What are we looking for?

- Form/pattern

- Direction

- Strength

- Outliers

**Correlation!** But, first how do we assess variation in two variables...

# Step 1: Standardizing variation in variables

$$\frac{x - \bar{x}}{s}$$

Example: Populations of European countries

| | $x$ | $\frac{x - \bar{x}}{s}$ |
|---|---|---|
| AT | 8.96mil | ? |
| BE | 11.69mil | ? |
| EL | 10.34mil | ? |
| IE | 5.05mil | ? |
| MT | 0.54mil | ? |
| PL | 41.03mil | ? |

```r
# create vector
x <- c(8.96, 11.69, 10.34, 5.05, 0.54, 41.03)
```

# STANDARDIZING VARIABLES: MEAN AND SD IN R

```r
# create vector
x <- c(8.96, 11.69, 10.34, 5.05, 0.54, 41.03)
# get mean and sd
c(round(mean(x), 2), round(sd(x), 2))
```

```
[1] 12.94 14.35
```

$$\bar{x} = 12.94 \quad s = 14.35$$

# Standardizing variables: Fill in table using R

|     | $x$      | $\frac{x - \bar{x}}{s}$ |
| --- | -------- | --- |
| AT  | 8.96mil  | ?   |
| BE  | 11.69mil | ?   |
| EL  | 10.34mil | ?   |
| IE  | 5.05mil  | ?   |
| MT  | 0.54mil  | ?   |
| PL  | 41.03mil | ?   |

Remember the formula

$$\frac{x - \bar{x}}{s}$$

# Standardizing variables

```
1 # create standardized distance for each observation
2 standardized.x <- (x - mean(x))/sd(x)
3 # view vector of standardized values
4 round(standardized.x, 2)
```

```
[1] -0.28 -0.09 -0.18 -0.55 -0.86 1.96
```

|  | $x$ | $\frac{x - \bar{x}}{s}$ |
|---|---|---|
| AT | 8.96mil | -0.28 |
| BE | 11.69mil | -0.09 |
| EL | 10.34mil | -0.18 |
| IE | 5.05mil | -0.55 |
| MT | 0.54mil | -0.86 |
| PL | 41.03mil | 1.96 |

Computation: Average of the products of the standardized values

$$r = \frac{1}{n-1} \sum \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

What does a positive correlation mean? Negative?

# CORRELATION VISUALIZED

# Facts about correlation

- Linear only

- **Not** causal

- Unit-free

- $-1 \leq r \leq 1$ (Remember that $r$ is the correlation coefficient)

- Sensitive to outliers

Assume we have 2 variables

We need to:

- Ask: Is there an "association" between them?

- Is it statistically significant (we'll discuss this next class)?

- Estimate "expected values" for an outcome variable given a set of covariates

# Starting off: Some preliminaries

*Y* = Response variable/Dependent variable/
Outcome variable/Explained variable/Left-hand side

*X* = Explanatory variable/Independent variable/
Treatment Variable/Right-hand side

We want to know: How might *Y* and *X* be related?

- We can visualize the relationship with a line!

## Linear Model

$$Y = \alpha + \beta X$$

$\alpha =$ Y-intercept and $\beta =$ slope of the line

$Y =$ outcome vector and $X =$ vector of predictor

# How do we write the regression line?



- $Y = 0.8 + \frac{2}{5}X$
- Interpret $\alpha$: = 0.8 = value of $Y$ when $X = 0$
- Interpret $\beta$: = 0.4 = a 1 unit $\uparrow$ in $X$ is associated with a 0.4 unit $\uparrow$ in $Y$
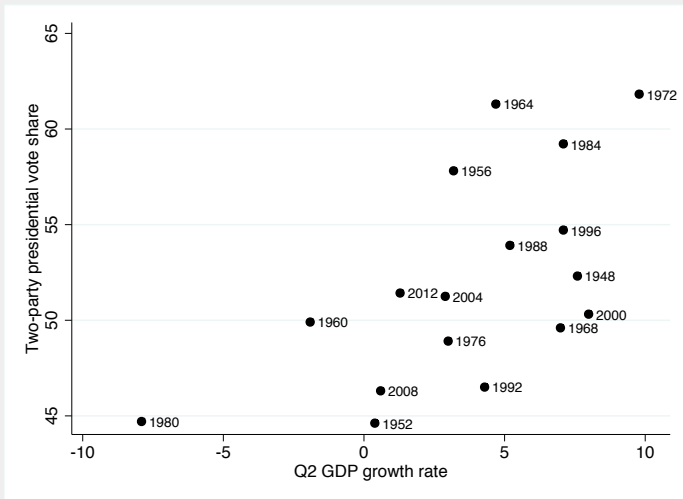
```
1 # create and view
     components of
     regression lines
2 # we have two variables
     X and Y
3 X <- runif(100, -10,
     10)
4 Y <- 0.8 + (2/5)*X
5 # plot relationship (
     with a line)
6 plot(X, Y, type="l",
     xlim = c(-10, 10),
     ylim = c(-5,5))
```
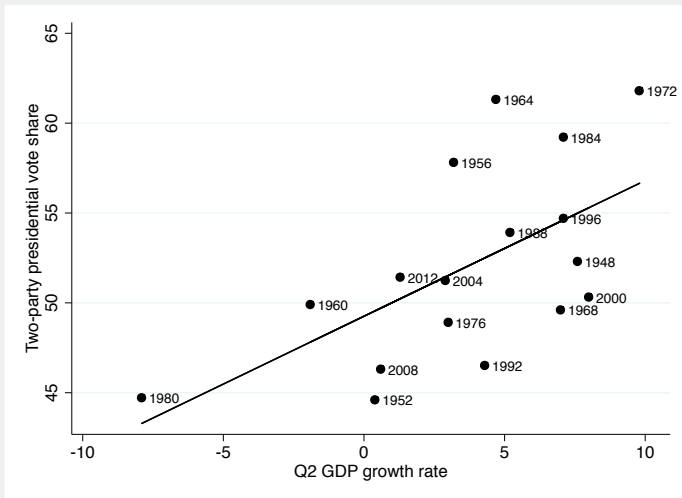
# Interpreting a regression line



- $\alpha$: E[$Y$] when $X$ is zero
- $\beta > 0$: Positive relationship between $X$ and $Y$
- $\beta < 0$: Negative relationship between $X$ and $Y$
- $\beta = 0$: Null relationship between $X$ and $Y$

Is there an association?

Incumbent party vote = 39.3 + 0.75 $Q_2$ GDP

# THINK FORMALLY ABOUT HOW TO DRAW "BEST" LINE

Let our data be the dyads $(Y_i, X_i)$, $i = 1, \ldots, n$

**Assumption #3:** There's a linear relationship between the variables:

$$E(Y_i) = \alpha + \beta X_i$$

However, we know there is error, so

$$Y_i = \alpha + \beta X_i + \epsilon_i$$

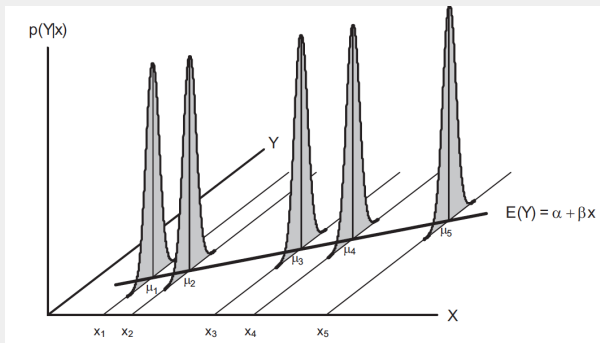**Assumption #4:** $\epsilon_i \sim N(0, \sigma^2)$

This is equivalent to:

$$Y_i \sim N(\alpha + \beta X_i, \sigma^2)$$

# Reminder: Assumptions for linear regression

(1) Randomized data generation

(2) Independent observations

# NEW: ASSUMPTIONS FOR LINEAR REGRESSION

(3) Linearity: Population means of *y* at different values of *x* have a straight-line relationship with *x*, i.e. $\mu_{y|x} = \beta_0 + \beta_1 x$

(4) Normality and Constant Variance: Population values of *y* at each value of *x* follow a **normal** distribution, with the **same** standard deviation $\sigma$ at each *x* value (constant variance in *y* for all *x*)

First, create data:

```r
1  # create data
2  # (1) draw N=100 from uniform distribution
3  # w/ min=0 and max=1
4  X <- runif(100, 0, 1)
5  # (2) draw corresponding outcome for N=100
6  # w/ mean=0 and sd=1
7  # embed correlation w/ X so we know what
8  # the relationship should be
9  Y <- 2 + X*1.5 + rnorm(100, 0, 1)
10 # run bivariate regression
11 temp.model <- lm(Y~X)
```

# Visualizing perfectly normal errors

Notice we have an imperfect model of "true" relationship between *Y* and *X*:

**Table:** Estimated relationship between *Y* and *X*

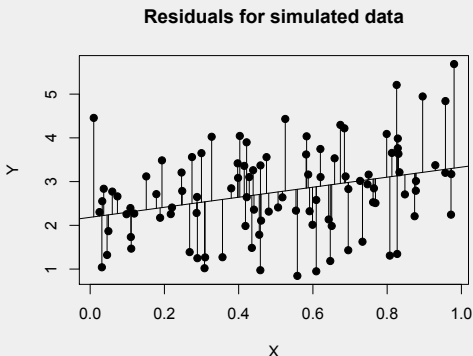|           | Outcome = *Y* |
|-----------|:-------------:|
| X         | 1.72$^{***}$  |
|           | (0.34)        |
|           |               |
| Intercept | 2.02$^{***}$  |
|           | (0.19)        |

*Notes*: OLS regression coefficients shown with standard errors in parentheses. $^{***}p < 0.001$, $^{**}p < 0.01$, $^{*}p < 0.05$. In all models, N=100.

# VISUALIZING PERFECTLY NORMAL ERRORS
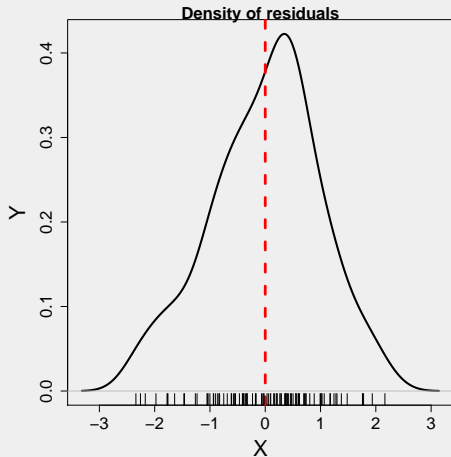
Plot residuals:

```
1  # plot residuals
2  plot(X, Y, pch=19, main
     ="Residuals for
     simulated data")
3  # display estimated
     regression line
4  abline(temp.model)
5  # show how far each
     prediction is from
6  # the estimated
     regression line
7  preds <- predict(temp.
     model)
8  segments(X,Y,X,preds)
```



**Residuals for simulated data**

# VISUALIZING PERFECTLY NORMAL ERRORS

Check that residuals are zero in expectation ($\epsilon_i \sim N(0, \sigma^2 = 1)$):



Density of residuals

```
1  plot(density(Y-preds),
      main="Density of
      residuals",
2    ylab="Y", xlab="X",
      cex.axis=1.5, cex.lab
      =2, cex.main=1.5, lwd
      =3)
```

# Implications of what we've done so far

- We have reduced all of data to a simplified model

- We have three parameter $(\alpha, \beta, \sigma^2)$ we now need to estimate using our data

- Once we have our parameter estimates, we want to then make inferences

  - ▶ Just like before, we will set up hypotheses

  - ▶ Just like before, we will summarize how well the data supports these hypotheses

But before we can do anything else, we need to choose estimates:

$$\hat{\alpha}, \hat{\beta}, \hat{\sigma}^2$$
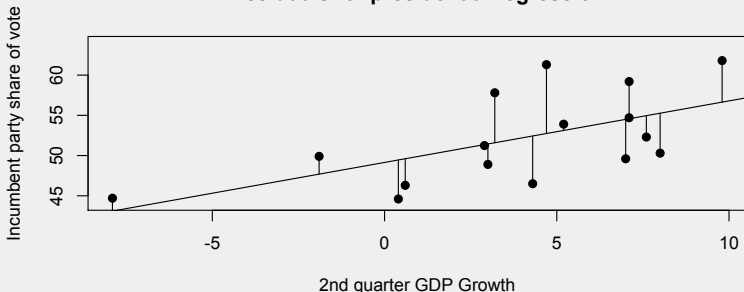
# PICK PARAMETERS TO MINIMIZE ERROR: DEFINE $\epsilon$

Let's define observed residual for observation $i$ as $e_i$ (our "error")

- Difference between our "best guess" for value of $Y_i$ given $X_i$ & what was actually observed (similar to $\chi^2$ test)
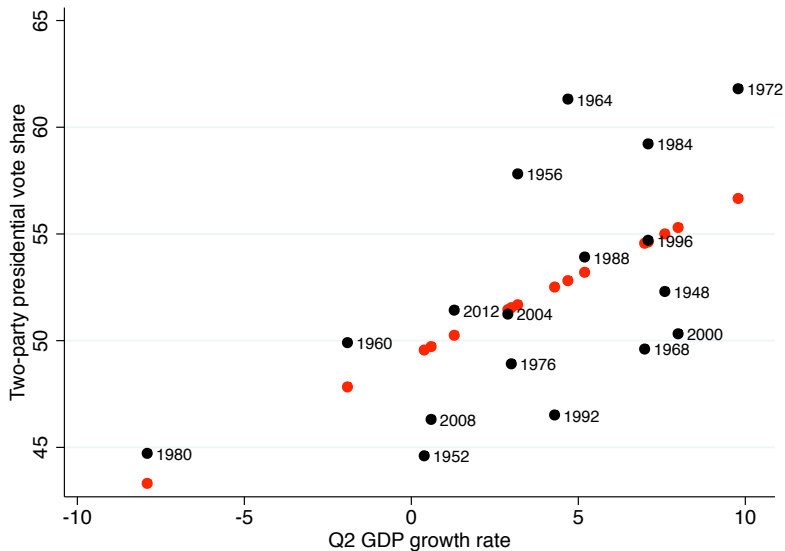
## Residuals

$$e_i = (Y_i - \hat{Y}_i) = (Y_i - \hat{\alpha} - \hat{\beta}X_i)$$



**Residuals for presidential regression**

Incumbent party share of vote

2nd quarter GDP Growth

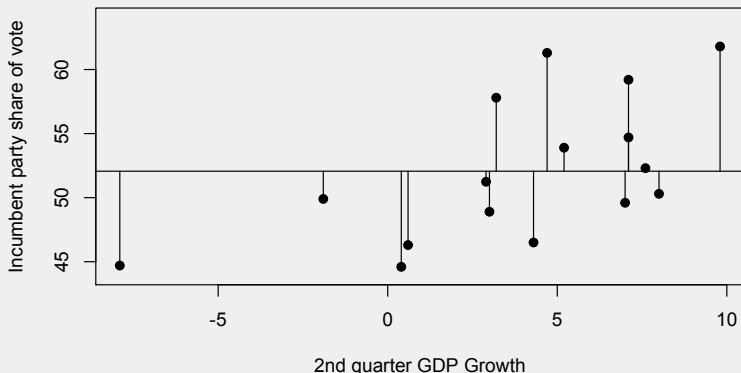# ANOTHER LOOK AT RESIDUALS FROM OUR EXAMPLE

# Again, think back to drawing the "best" line

Intuition: "Best" line is the one that reduces the greatest error

- We might think to first check sum of the errors $\sum_{i=1}^{n}(Y_i - \hat{Y}_i)$, but that's not actually a good criterion, why? (errors equal out)

**A very bad line with residuals that sum to zero**

# Defining "best" as minimizing SSE

For many good statistical reasons, we'll say any line that reduces the **"Sum of Squared Error"** is equivalent to having the "best" line (defined as *"most efficient unbiased estimator"*)

- For right now, trust me!

Sum of Squared Error

$$SSE = \sum_{i=1}^{n}(Y_i - \hat{Y}_i)^2 = \sum_{i=1}^{n}(Y_i - \hat{\alpha} - \hat{\beta}X_i)^2$$

We are going to minimize SSE with respect to $\hat{\alpha}$ and $\hat{\beta}$ (using calculus in the background!)

- With these parameters, we can draw "best" lines

# ESTIMATORS FOR $\alpha$ AND $\beta$

$$\hat{\beta} = \frac{\sum_{i=1}^{n} \left( (X_i - \bar{X})(Y_i - \bar{Y}) \right)}{\sum_{i=1}^{n}(X_i - \bar{X})^2}$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

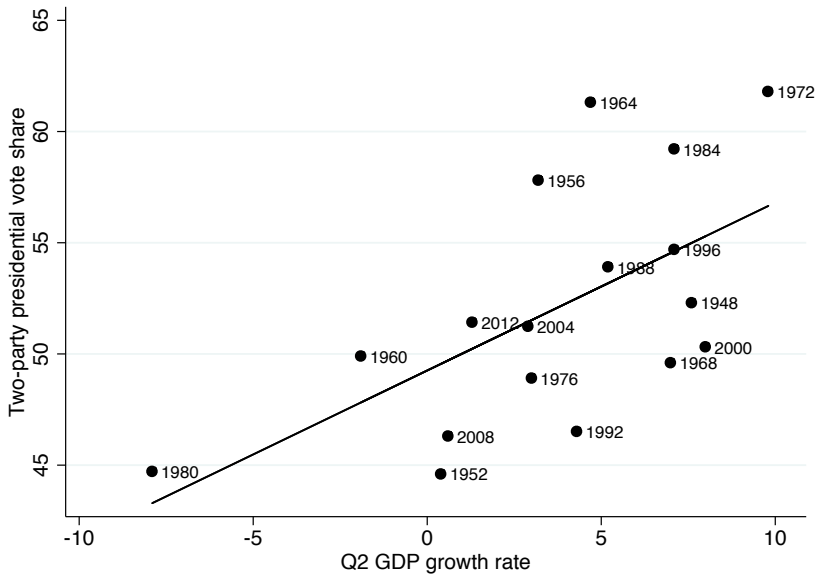Both of these are functions of the data

For our presidential election data:
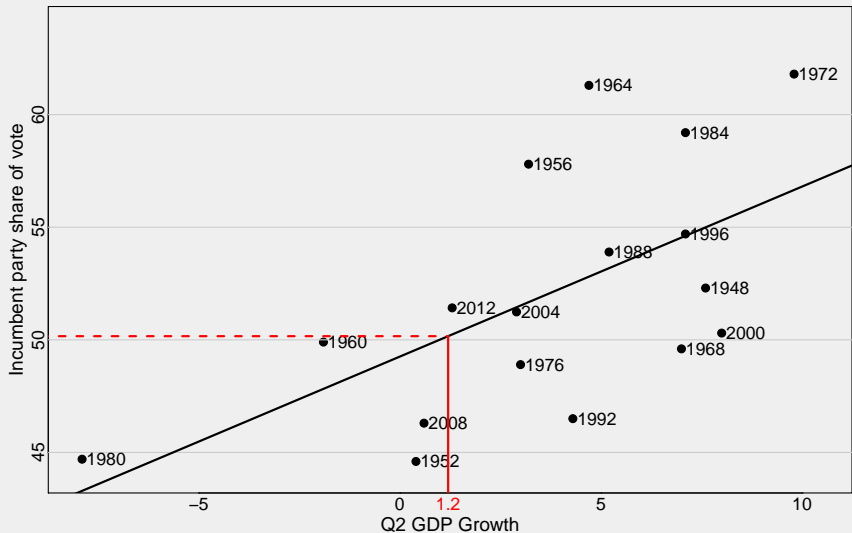
- $\hat{\alpha} = 39.3$

- $\hat{\beta} = 0.75$

How do we interpret these estimates?

# EXAMPLE: INTERPRETING ESTIMATES

# WHO WILL WIN THE ELECTION?

If $Q_2$ GDP $\Delta$ = 1.2 ...

## Example: Bivariate regression "by hand"

| $X_i$ | $Y_i$ |
|------|------|
| 3.8  | 3.5  |
| 3.0  | 3.3  |
| 3.5  | 4.0  |
| 2.8  | 2.3  |
| 2.4  | 1.8  |
| 2.7  | 2.7  |

Find $\hat{\alpha}$ and $\hat{\beta}$

$$\hat{\beta} = \frac{\sum_{i=1}^{n}\left((X_i - \bar{X})(Y_i - \bar{Y})\right)}{\sum_{i=1}^{n}(X_i - \bar{X})^2}; \quad \hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$$

You need to compute each component:

$$
\begin{array}{ccccc}
Y_i & X_i & (Y_i - \bar{Y}) & (X_i - \bar{X}) & (Y_i - \bar{Y})(X_i - \bar{X}) \\
\hline
\vdots & \vdots & \vdots & & \vdots \\
\hline
\end{array}
$$

(1) Create data in R:

```
1 # bivariate regression by hand
2 regressMat <- as.data.frame(matrix(c(3.8, 3.5, 3.0, 3.3,
    3.5, 4.0, 2.8, 2.3, 2.4, 1.8, 2.7, 2.7), nrow=6, byrow
    = T))
3 colnames(regressMat) <- c("X", "Y")
```

$\bar{Y} = 2.93$

```
1 mean(regressMat$Y)
```

$\bar{X} = 3.03$

```
1 mean(regressMat$X)
```

$\sum Y_i = 17.6$

```
1 sum(regressMat$Y)
```

$\sum X_i = 18.2$

```
1 sum(regressMat$X)
```

$\sum (Y_i - \bar{Y})(X_i - \bar{X}) = 1.863$

```
1 sum((regressMat$Y -
2     mean(regressMat$Y))
3   *(regressMat$X-
4     mean(regressMat$X)))
```

$\sum (X_i - \bar{X})^2 = 1.373$

```
1 sum((regressMat$X-
2     mean(regressMat$X))^2)
```

# CALCULATE EACH COMPONENT "BY HAND" IN R

| $Y_i$ | $X_i$ | $(Y_i - \bar{Y})$ | $(X_i - \bar{X})$ | $(Y_i - \bar{Y})(X_i - \bar{X})$ |
|------|------|--------|--------|--------|
| 3.8 | 3.5 | 0.767 | 0.567 | 0.435 |
| 3.0 | 3.3 | -0.033 | 0.367 | -0.012 |
| 3.5 | 4.0 | 0.467 | 1.067 | 0.498 |
| 2.8 | 2.3 | -0.233 | -0.633 | 0.147 |
| 2.4 | 1.8 | -0.633 | -1.133 | 0.717 |
| 2.7 | 2.7 | -0.333 | -0.233 | 0.078 |

$$\hat{\beta} = \frac{1.863}{1.373} = 1.357$$

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X} = 2.933 - 1.357(3.033) = -1.183$$

# CHECK WITH lm() IN R

```
1 lm(Y~X, data=regressMat)
```

**Table:** Estimated relationship between *Y* and *X*

|              | Outcome = *Y* |
| ------------ | ------------- |
| X            | $1.36^{*}$    |
|              | (0.38)        |
| Intercept    | $-1.18$       |
|              | (1.18)        |

*Notes*: OLS regression coefficients shown with standard errors in parentheses.
$^{***}p < 0.001$, $^{**}p < 0.01$, $^{*}p < 0.05$.

We'll talk about standard errors next week!

# WRAP-UP

## Today we...

1. Estimate if two variables are dependent

   ▶ Chi-squared test of independence

   ▶ Standardized residuals

2. Correlations

3. Simple linear regression:

   ▶ Assumptions

   ▶ Estimation

# Class business

- Read the required and suggested online materials

- Work on Problem set # 1!

- These slides are available on the course website