

Bivariate relationships

Week 6

23 February, 2015

Prof. Andrew Eggers



[Home](#) > [About](#) > [Organisation](#) > [Finance and funding](#) > Financial Statements of the Oxford Colleges 2012-13

Financial Statements of the Oxford Colleges (2012-13)

HIS



The financial statements of the 36 colleges of Oxford University for the year ending 31 July 2013 are available as pdfs, together with an aggregated statement of financial activity (SOFA) and aggregated consolidated balance sheet.

The colleges are independent, self-governing and financially autonomous and their accounts are published under the accounting convention developed by the Charity Commissions for use by charities in the UK (the Charity SORP).

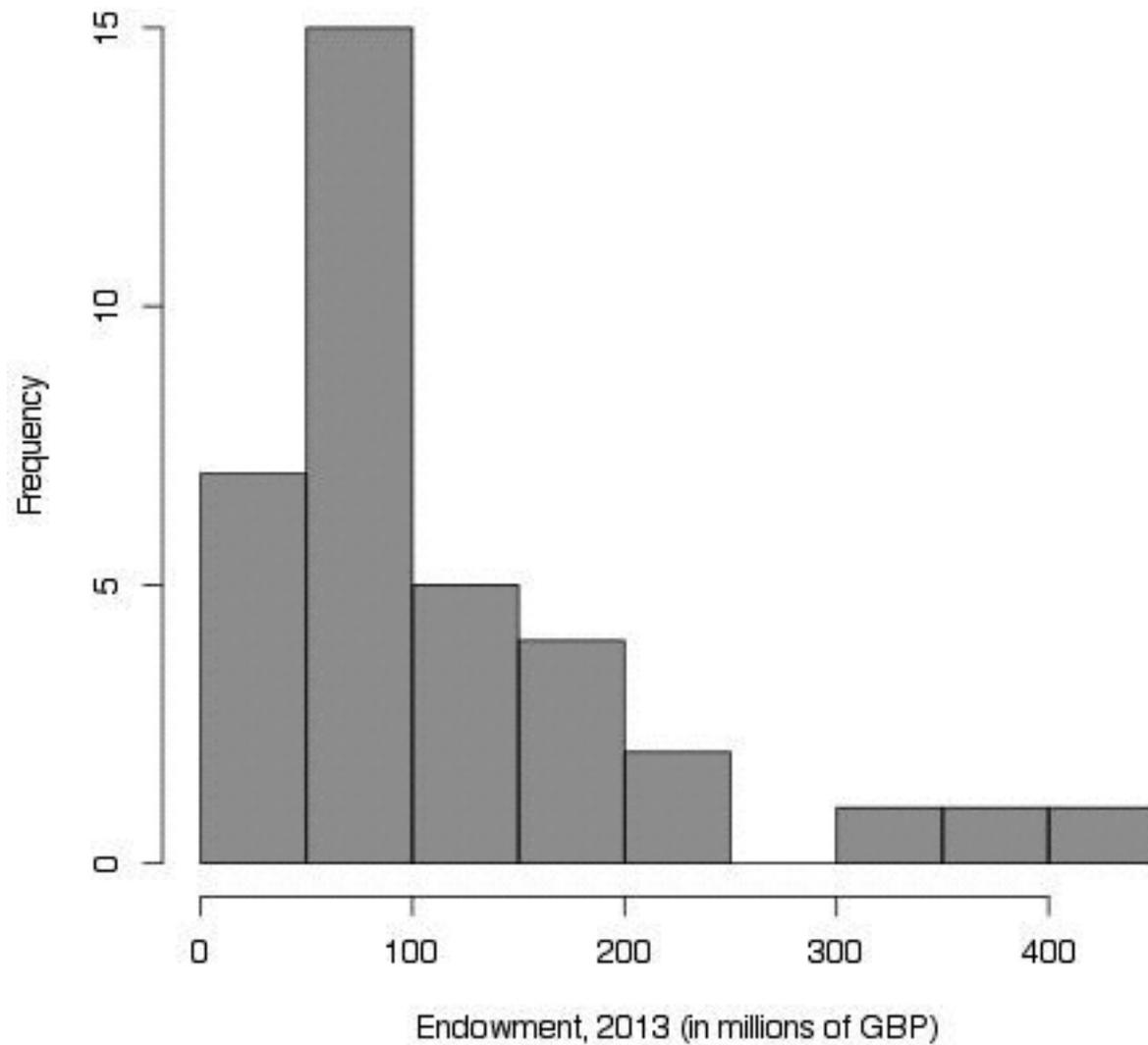
Kellogg College and St Cross College do not have Royal Charters and, for accounting purposes, are departments of the University. As such, their financial results are consolidated into the University's financial statements.

<http://goo.gl/1pJA2r>

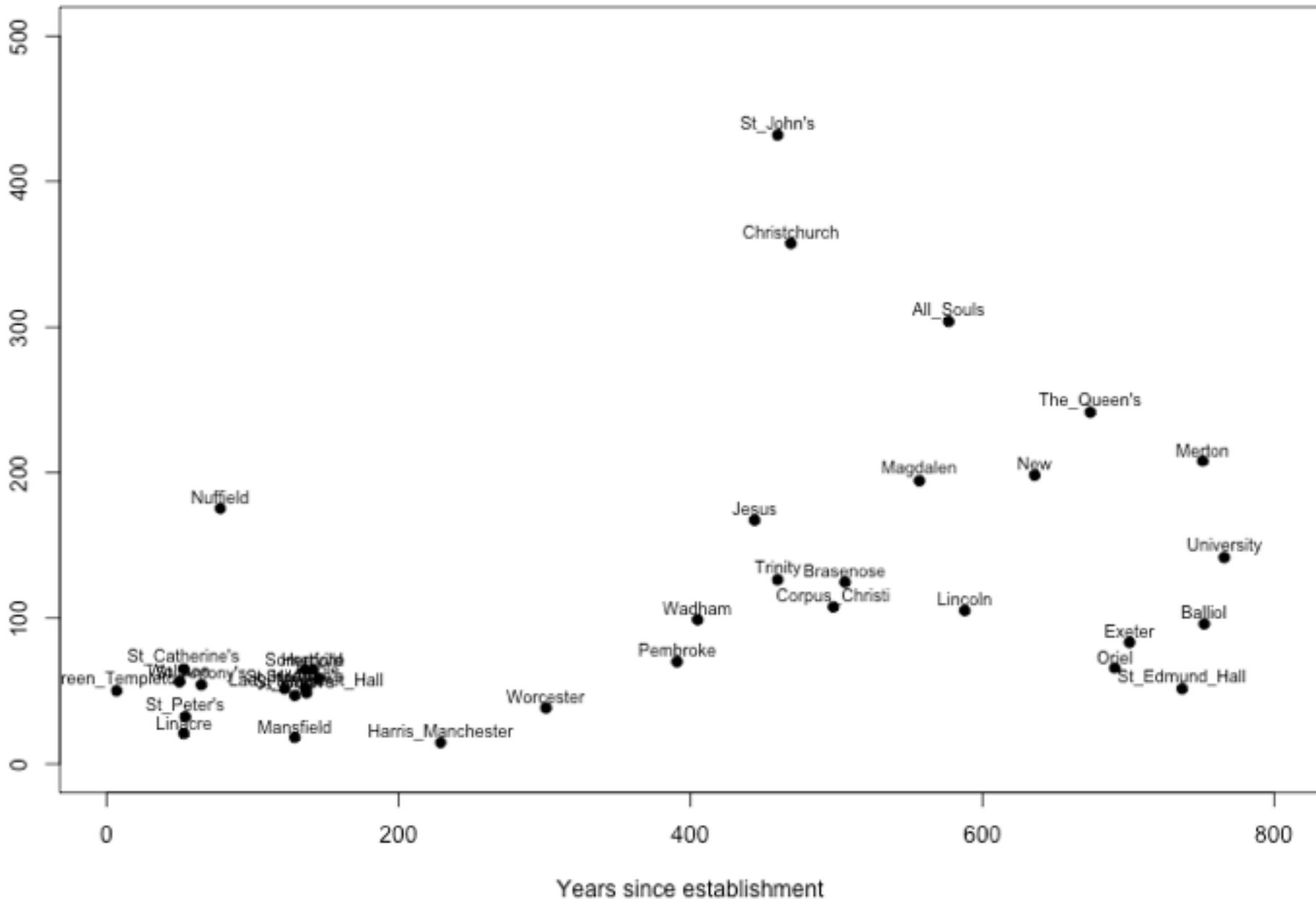
```
> ## data analysis for lecture
> d = read.csv("http://andy.egge.rs/data/college_stats_edited.csv")
> d$rank = 37 - rank(d$endowments, )
>
> # just a table of the data
> d[order(d$endowments, decreasing = T), c("rank", "college", "endowments")]
```

	rank	college	endowments
28	1	St_John's	432075
4	2	Christchurch	357667
1	3	All_Souls	303896
30	4	The_Queen's	241467
36	5	Merton	208054
17	6	New	198160
15	7	Magdalen	194344
18	8	Nuffield	175415
10	9	Jesus	167333
32	10	University	141872
31	11	Trinity	126350
3	12	Brasenose	124684
5	13	Corpus_Christi	107692
14	14	Lincoln	105245
33	15	Wadham	98935
2	16	Balliol	96044
6	17	Exeter	83383
20	18	Pembroke	70195

College endowments



Size of endowment, 2013 (millions of GBP)



Research questions you might have about Oxford colleges' endowments and age

Descriptive/predictive questions:

- How are age and size of endowment related?
- What is the average endowment of a college that is X years old?



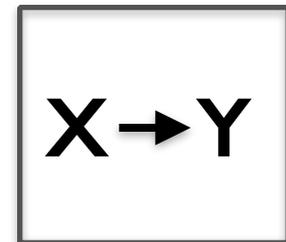
Explanatory questions:

- Why do some colleges have more money than others? (Maybe age is the/an answer.)



Causal questions:

- (What is the effect of greater age on a college's endowment?)



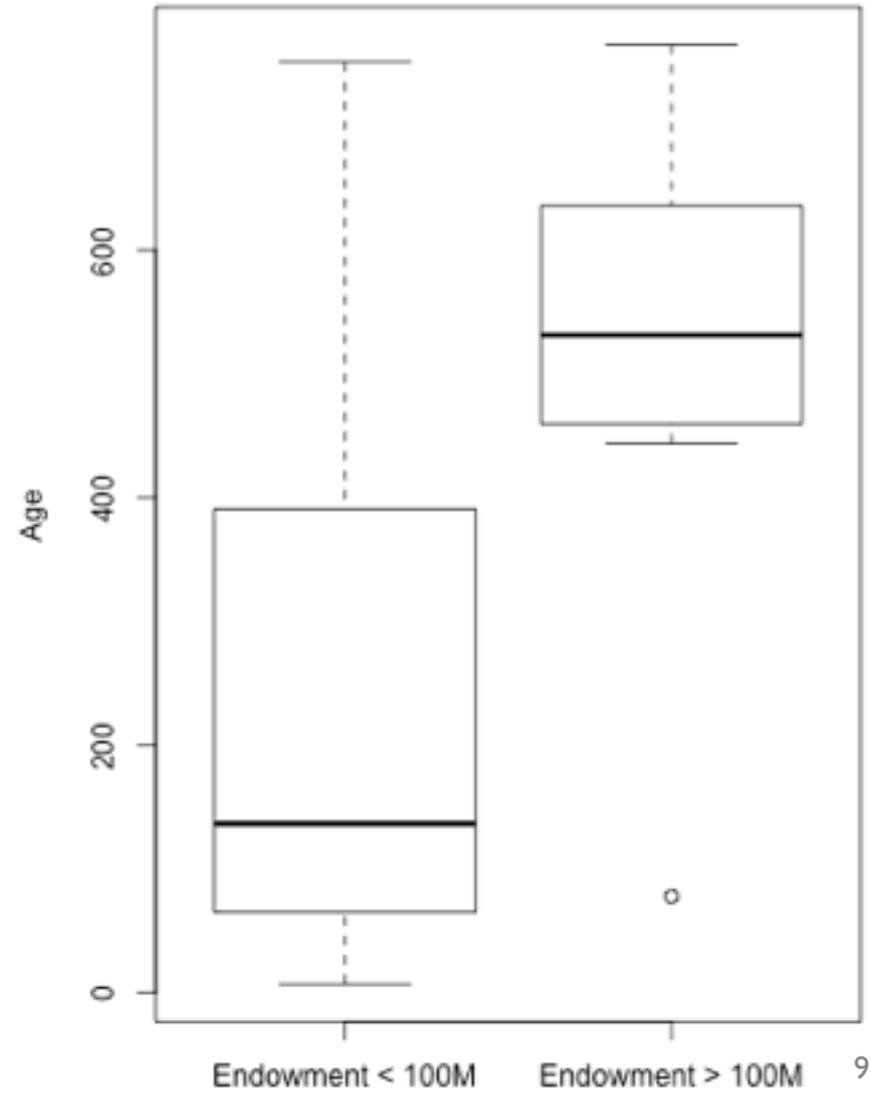
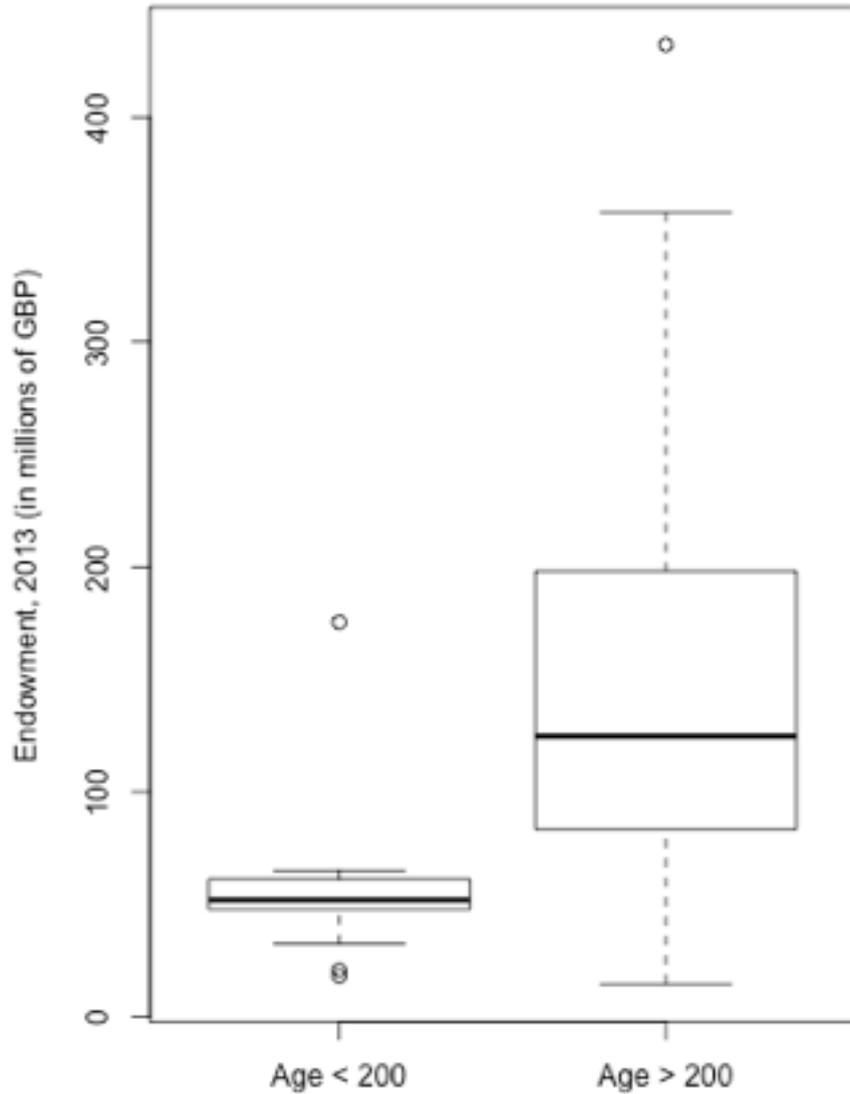
Various ways to summarize a relationship between two variables

- Show the scatterplot!
- Compare boxplots of one variable across categories of the other
- Show/report mean of one variable across categories of the other (binned means or kernel average smoother)
- Report covariance/correlation
- Report regression coefficient(s)
- Report predicted values of one variable based on the other from regression

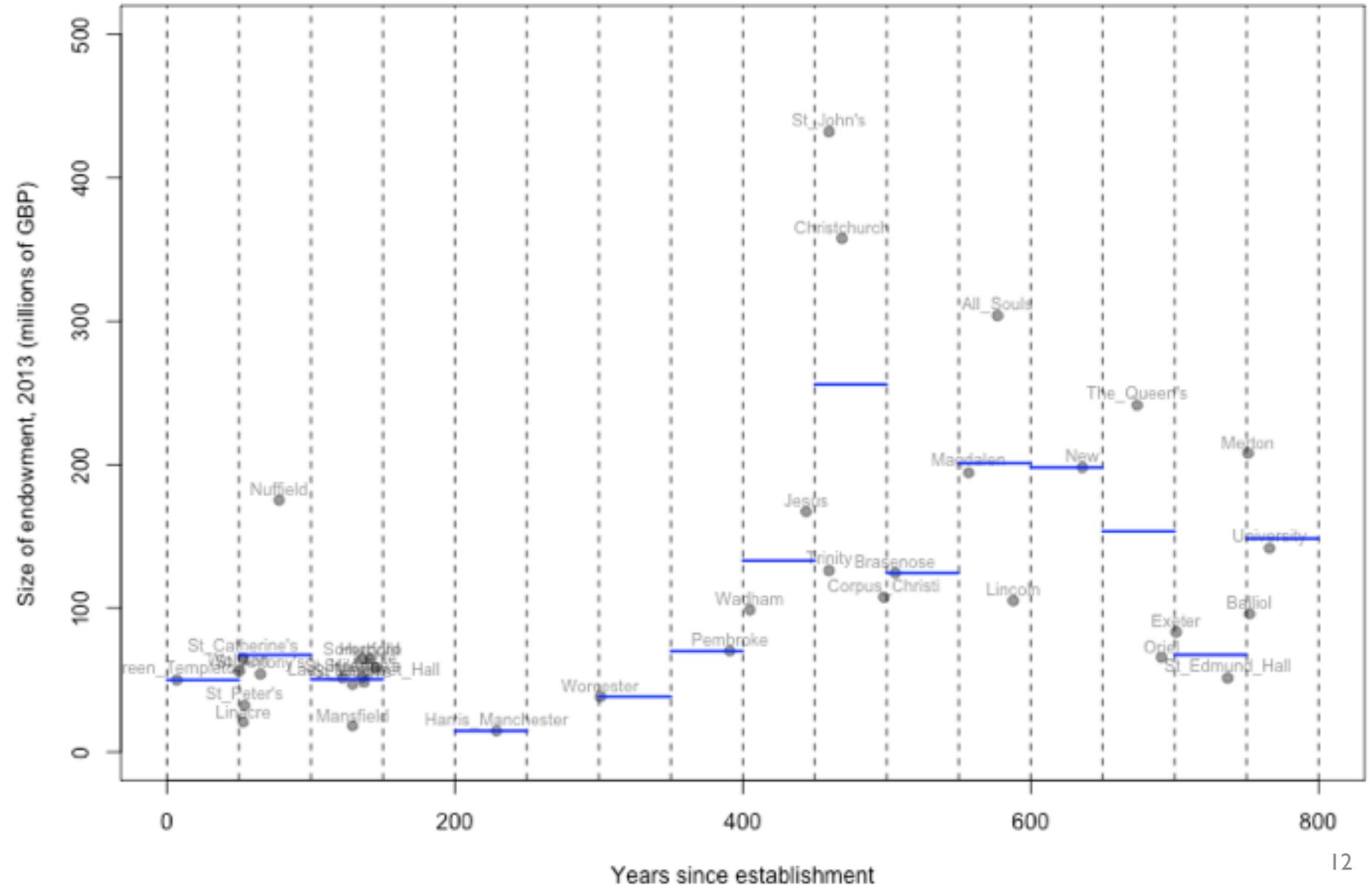
But first: one minute to provide us with some data!

- get out a device (smartphone, laptop, tablet)
- go to <http://andy.egge.rs/form.html>
- click on the link
- fill out the survey (5 questions)
- put away your device and look up

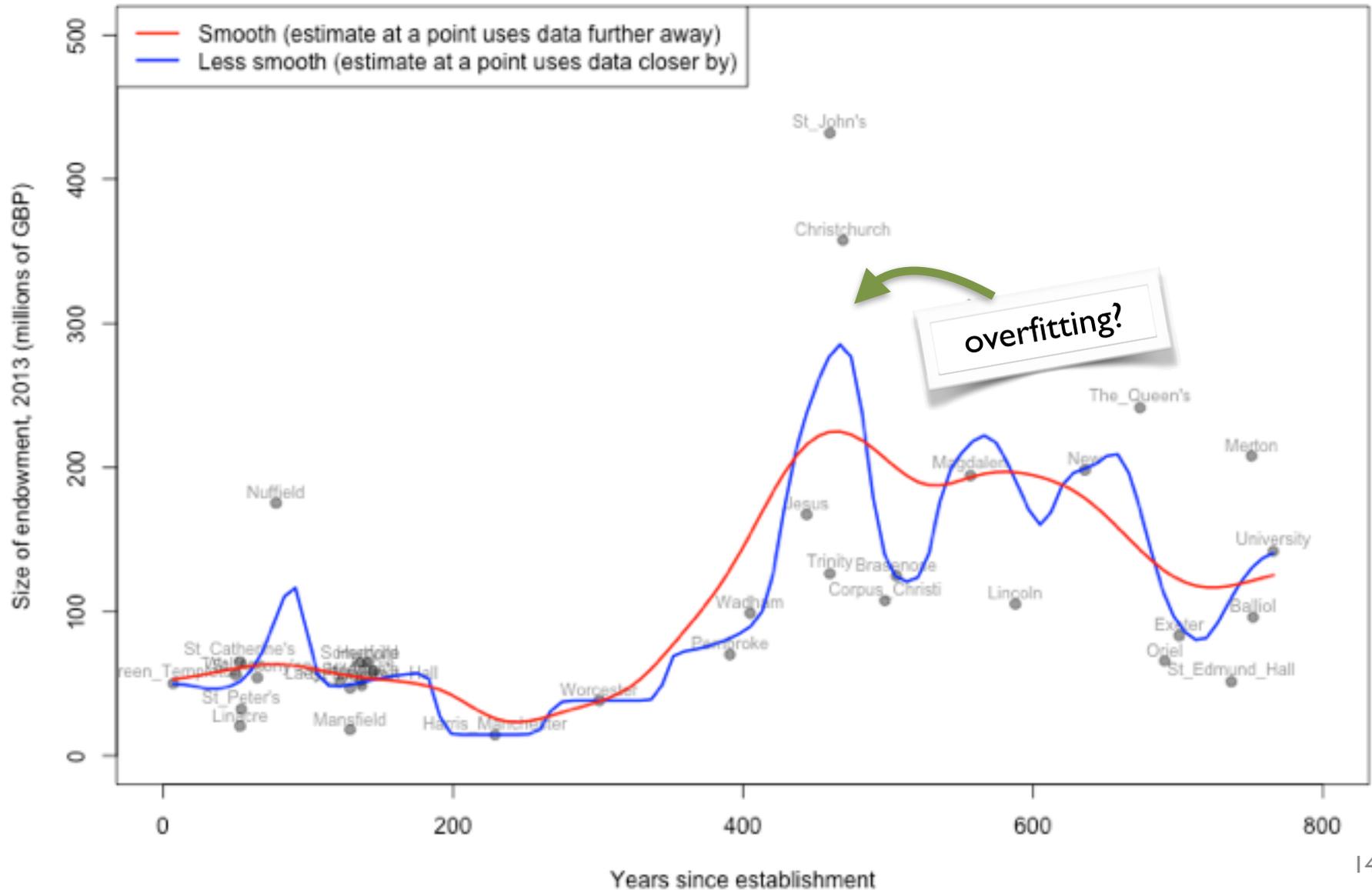
Boxplots: which do you prefer?



Mean endowment within 50-year intervals



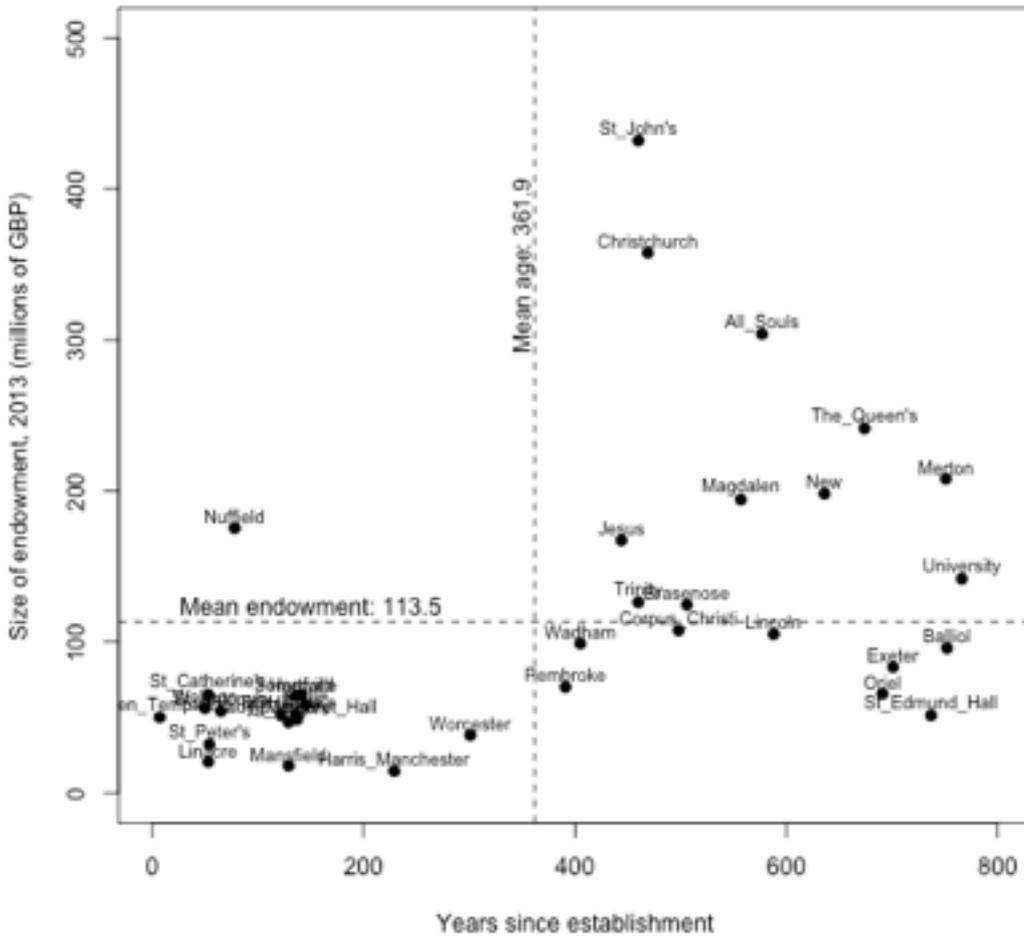
Kernel smoother: estimate at a point is (weighted) average of nearby points



Covariance: a measure of linear association

$$\text{Cov}(x, y) = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

“mean of y”



College	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$
Mansfield	-232.9	-95.4	22,218.1
Nuffield	-283.9	61.9	-17,579.2
St. John's	98.1	318.6	31,256.5
Wadham	43.1	-14.6	-627.6

etc., for all the colleges, and take the average of the final column. Or, more simply:

```
> cov(d$age, d$endow.mill)
[1] 11828.5
```

Two facts about covariance

$$\text{Cov}(x, y) = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

FACT 1:

$$\text{Cov}(x, y) = \text{Cov}(y, x)$$

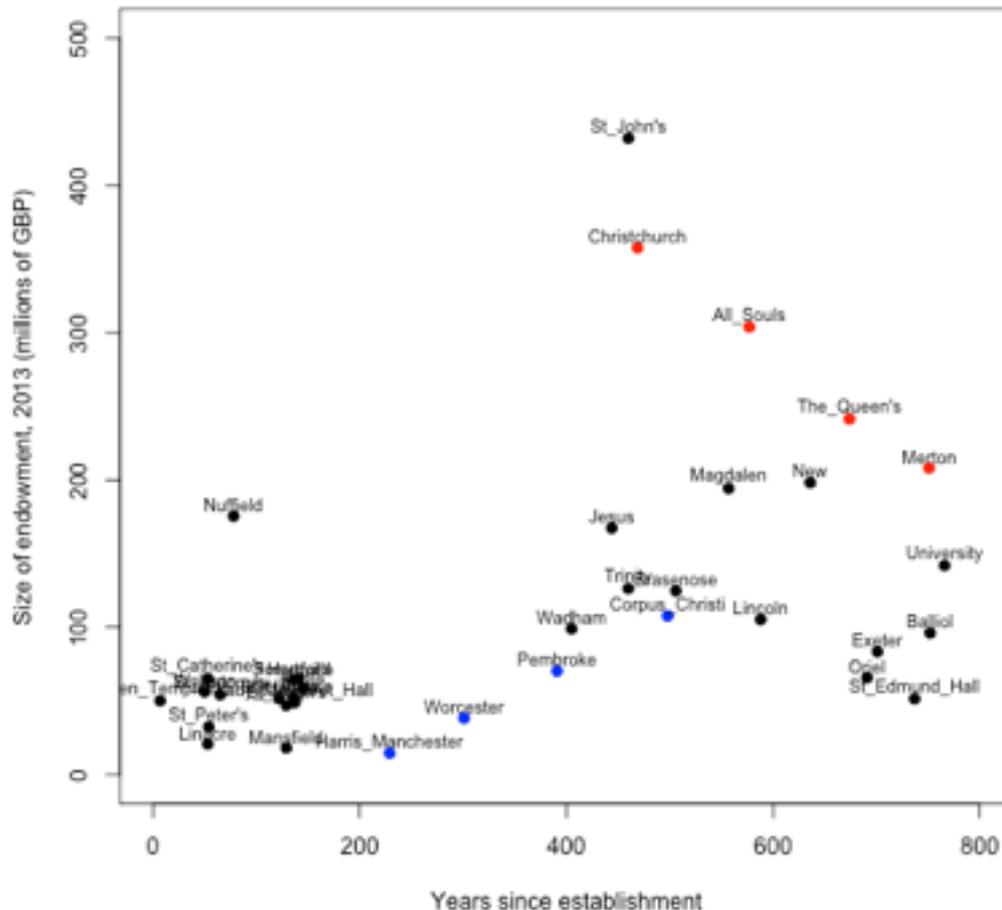
FACT 2:

$$\text{Cov}(x, x) = \text{Var}(x)$$

Correlation: a scale-invariant measure of linear association

$$\text{Cor}(x, y) = \frac{\text{Cov}(x, y)}{\text{sd}(x)\text{sd}(y)}$$

“standard deviation of y”



<i>Sample</i>	<i>Cov(x,y)</i>	<i>Cor(x,y)</i>
All	11,838.5	0.471
All, with y in 1000s of GBP	11,828,503.0	0.471
Red only	-8093.7	-0.99785
Blue only	4697.3	0.99992
Red and blue	15083.7	0.658

Five facts about correlation

$$\text{Cor}(x, y) = \frac{\text{Cov}(x, y)}{\text{sd}(x)\text{sd}(y)}$$

FACT 1:

$$\text{Cor}(x, y) \in [-1, 1] \quad \forall \quad x, y$$

FACT 2:

$$\text{Cor}(x, y) = \text{Cov}\left(\frac{x}{\text{sd}(x)}, \frac{y}{\text{sd}(y)}\right)$$

FACT 3:

$$\text{Cor}(x, y) = \text{Cor}(y, x)$$

FACT 4:

The x,y correlation for a line of points will be 1 or -1 unless it is vertical or horizontal, in which case it is undefined.

FACT 5:



Finally: regression!

ANTHROPOLOGICAL MISCELLANEA.

REGRESSION *towards* MEDIOCRITY *in* HEREDITARY STATURE.

By FRANCIS GALTON, F.R.S., &c.

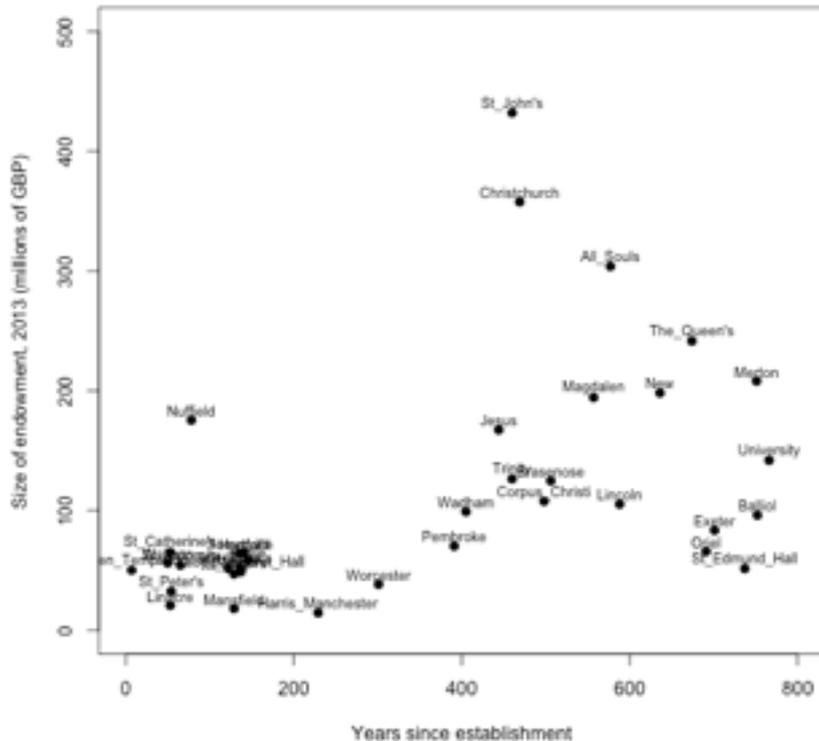
[WITH PLATES IX AND X.]

THIS memoir contains the data upon which the remarks on the Law of Regression were founded, that I made in my Presidential Address to Section H, at Aberdeen. That address, which will appear in due course in the Journal of the British Association, has already been published in "Nature," September 24th. I reproduce here the portion of it which bears upon regression, together with some amplification where brevity had rendered it obscure, and I have added copies of the diagrams suspended at the meeting, without which the letterpress is necessarily difficult to follow. My object is to place beyond doubt the existence of a simple and far-reaching law that governs the hereditary transmission of, I believe, every one of those simple qualities which all possess, though in unequal degrees. I once before ventured to draw attention to this law on far more slender evidence than I now possess.



Bivariate regression: drawing a line near some points

$$y_i = \alpha + \beta x_i$$



No way to draw a line **through** all of these points!

How can we produce a line “close” to the points?

- “Eyeball it”
- Draw a line that minimizes distance from the points to the line
 - Minimize the **total** distance from the points to the line (or equivalently, **average** distance)
 - Minimize the **median** distance from the points to the line
- Draw a line that minimizes total/average or median **vertical** distance from the points to the line (i.e. “prediction error”, if we think of x predicting y)

Linear regression (ordinary least squares) minimizes the total/average squared vertical distance from the points to the line.

Bivariate regression: more formally

$$y_i = \hat{\alpha} + \hat{\beta}x_i + \hat{u}_i$$

Linear regression (OLS) solves:

$$\min_{\hat{\alpha}, \hat{\beta}} \sum_i \hat{u}_i^2$$

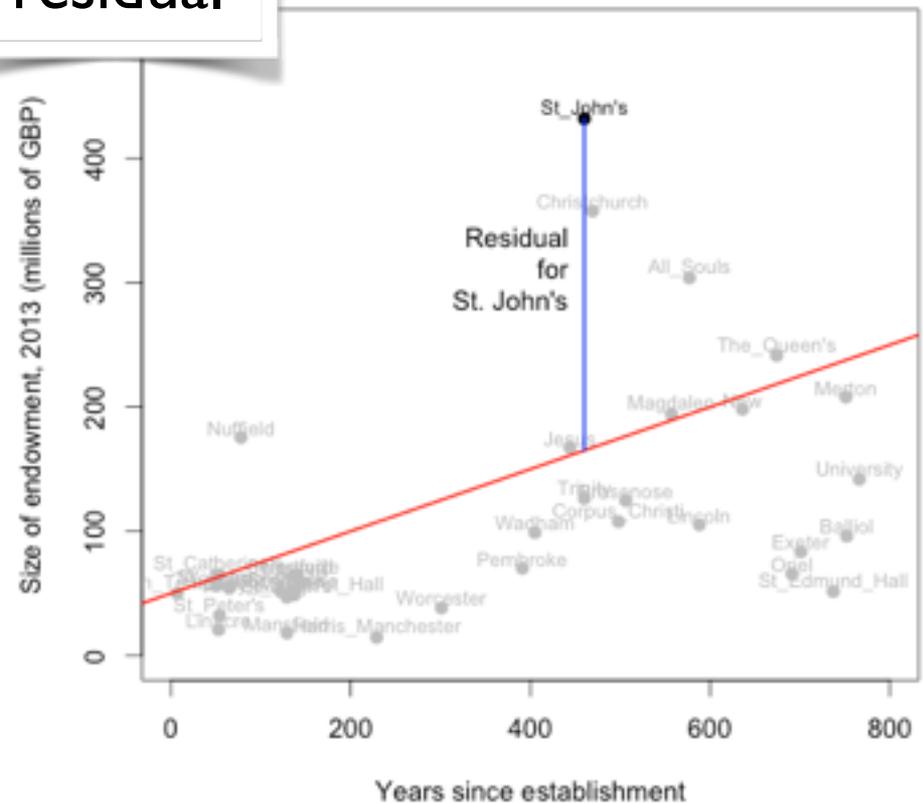
Minimization (remember calculus?) gives:

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$$

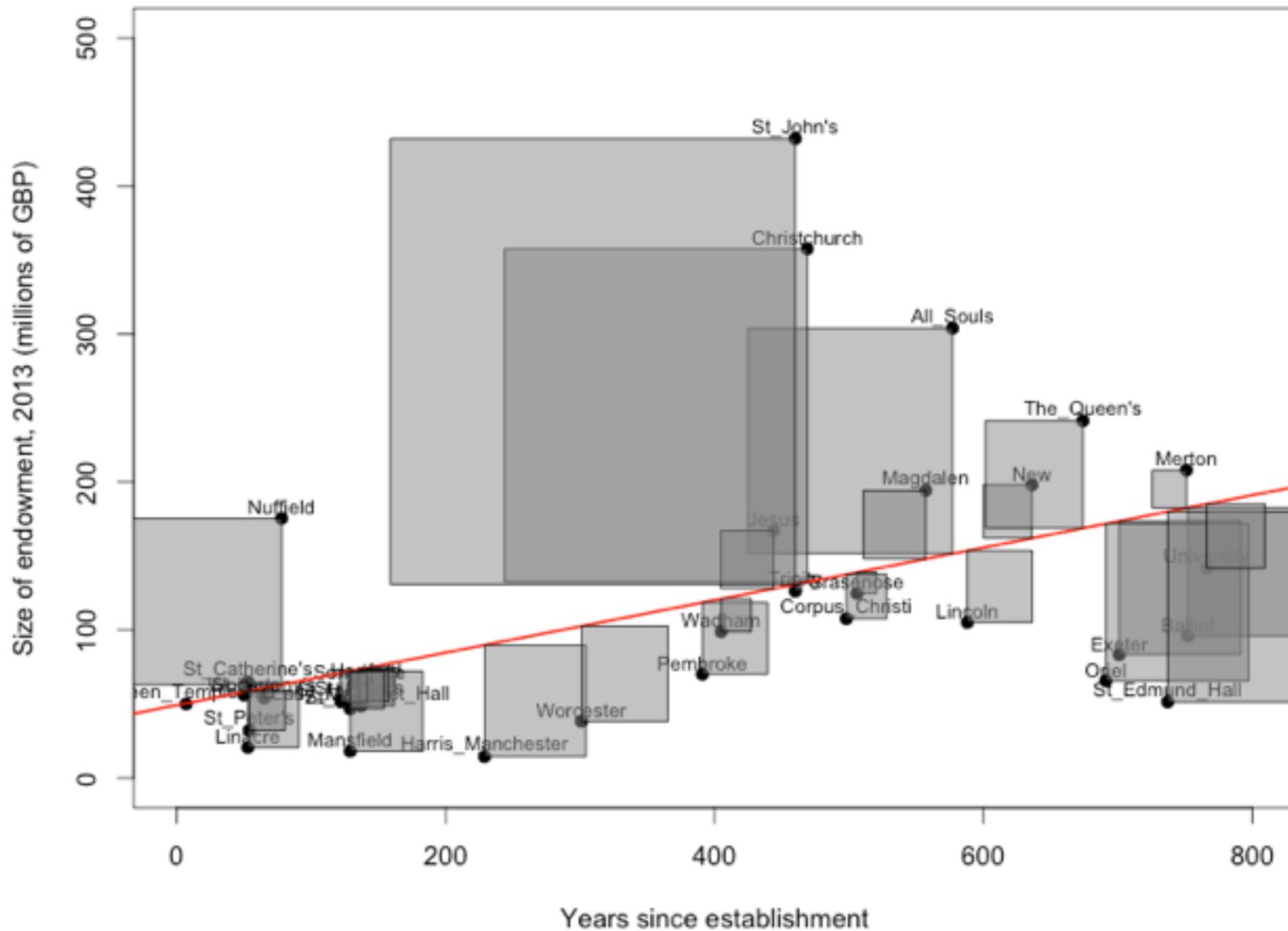
$$\hat{\beta} = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

- Greek letters without hats: parameters in the **population**.
 - Greek letters with “hats”*: estimates of those parameters produced from the sample.
- * In Andrea’s lecture, Roman letters.

“residual”

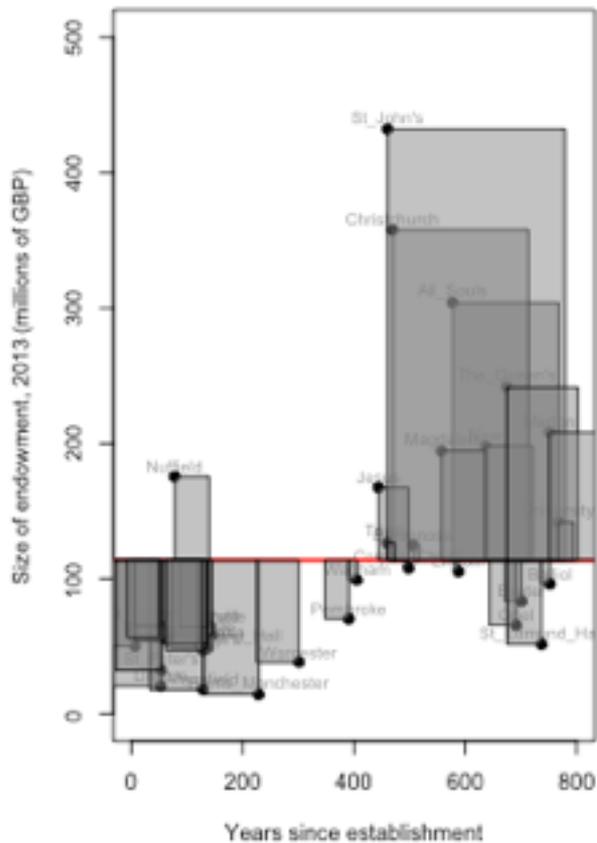


Visualizing the sum of squared residuals

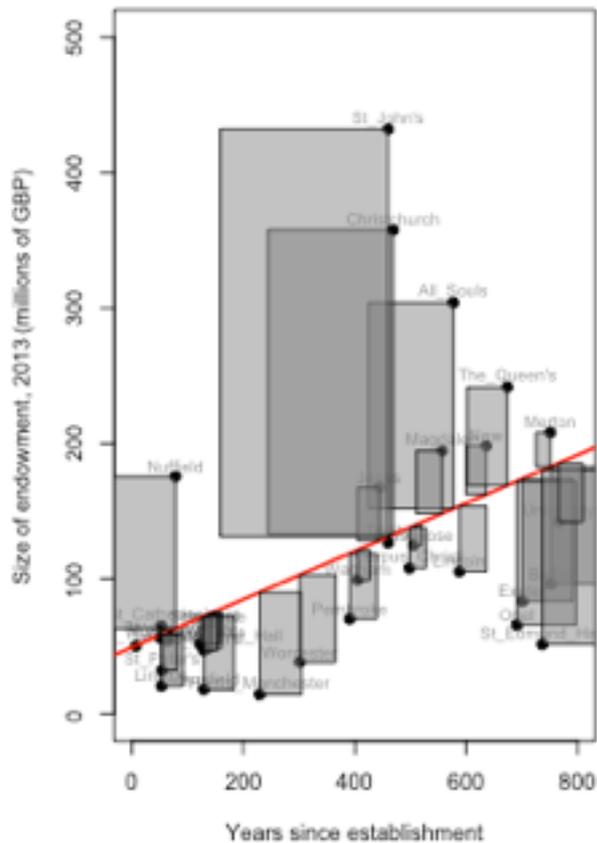


Visualizing the minimization of the sum of squared residuals

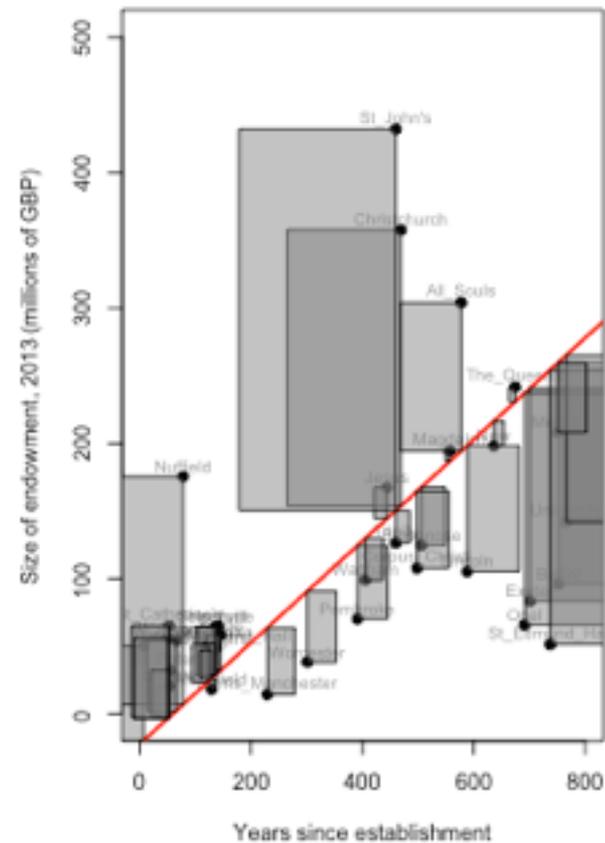
Intercept: 113.5; slope: 0
Sum of squared residuals: 330797.9



Intercept: 49.2; slope: 0.18
Sum of squared residuals: 257295.9



Intercept: -22.2; slope: 0.38
Sum of squared residuals: 348290.6



Relationships among covariance, correlation, and slope of regression

Recall that the OLS solution for the **slope** is:
$$\hat{\beta} = \frac{\text{Cov}(x, y)}{\text{Var}(x)}$$

Variances are always positive, so $\text{Cov}(x, y)$, $\text{Cor}(x, y)$ and $\hat{\beta}$ always have the same **sign**. For example:

```
> cov(d$age, d$endow.mill)
[1] 11828.5
> cor(d$age, d$endow.mill)
[1] 0.4713768
> coef(lm(d$endow.mill ~ d$age))[2]
      d$age
0.1775421
```

In one (important) special case, $\text{Cov}(x, y)$, $\text{Cor}(x, y)$ and $\hat{\beta}$ are **exactly the same**:

```
> d$std.age = d$age/sd(d$age)
> d$std.endow.mill = d$endow.mill/sd(d$endow.mill)
> cov(d$std.age, d$std.endow.mill)
[1] 0.4713768
> cor(d$std.age, d$std.endow.mill)
[1] 0.4713768
> coef(lm(d$std.endow.mill ~ d$std.age))[2]
      d$std.age
0.4713768
```

Correlation tells you how close to linear a relationship is. Correlation of 0.99 means “very linear!” no matter what units we are talking about.

Slope coefficient in a regression tells you what a slope usually tells you: how much does the line go up (or down) when you move to right by one unit? It depends on units.

Interpreting our regression coefficients

Suppose we had randomly left out a college.

We can use our regression estimate to predict its endowment!

- If it was created this year: 49.2 million
- If it was created 100 years ago: $49.2 + 100 \times 0.178 = 67$ million

(Note correlation does not yield a prediction.)

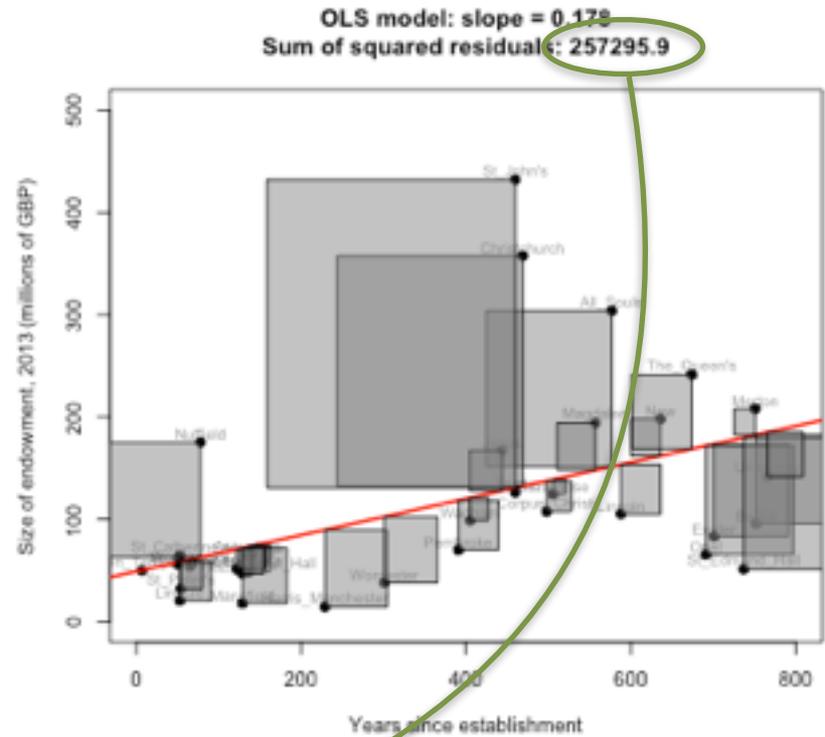
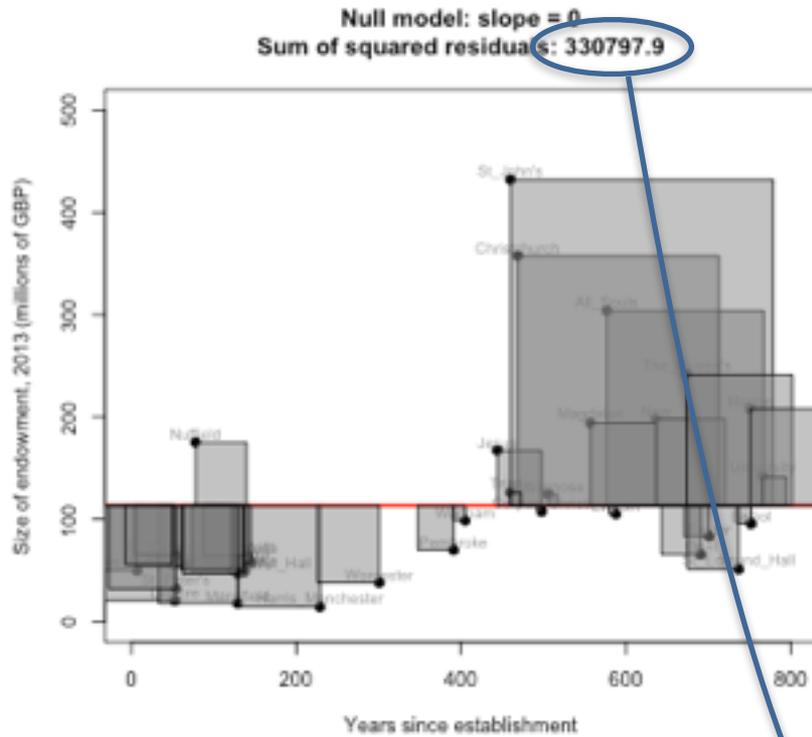
What about statements like these:

- “The effect of age on endowment is 0.178.”
- “A 100-year increase in age leads to an increase of 17.8 million in endowment.”

	<i>Dependent Variable: Endowment in millions</i>
Intercept	49.2 (25.2)
Age of college	0.178** (0.057)
Number of observations	36
R ²	0.222

Interpreting R²

R² measures how much of the variation in y is explained by x.



$$R^2 = 1 - \frac{RSS}{TSS}$$

TSS is closely related to the variance of y:

```
> var(d$endow.mill)*(nrow(d) - 1)  
[1] 330797.9
```

Looking ahead

- Rest of this week: data labs, part 3!
- Next two weeks: multivariate regression and statistical inference