

THE GUARDIAN

UNIVERSITY OF CALIFORNIA, SAN DIEGO

Extra! Extra! COGS108: DATA SCIENCE IN PRACTICE

GITHUB USERS
HAVE DESIRED
JOB SKILLS!

Inferential
Analysis

DATA SCIENCE
IS THE FUTURE

COGS108 STUDENTS
BLOW PROFESSOR AWAY
WITH INTERESTING FINAL
PROJECTS.



Question

Does prestige increase newspaper readership?

"How would you measure prestige?"

- * wikipedia cites each source
 - ↳ where cited?
- * journalism awards
 - Pulitzer Prize
- * how long in circulation
- * where they get \$ from

readership - physical copies
may not be read

- web user context

- views
- scroll rate

→ ad revenue (proxy)
online

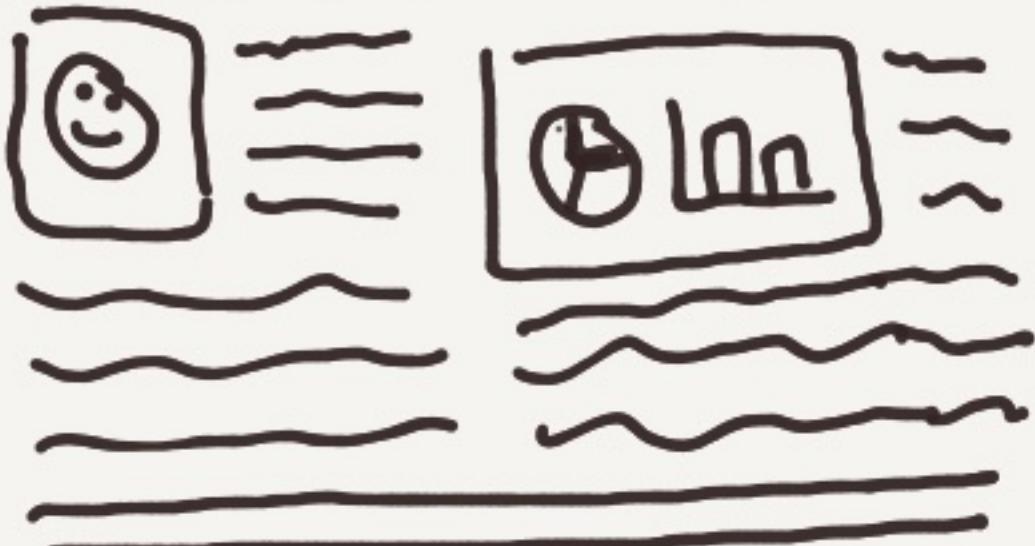
→ # of subscriptions

↳ miss out on "free"
readers

i.e.
FB
questionable

THE GUARDIAN

EXTRA! EXTRA! EXTRA



Question

- what time period?
- what newspapers?
- what specific metrics?

start w/ all news papers
of Politzers) \hookrightarrow mean $\pm 2\sigma$

75
Politzer

newspapers \rightarrow # of politzers

"What is the effect
of x on y ?"

Data Science Question

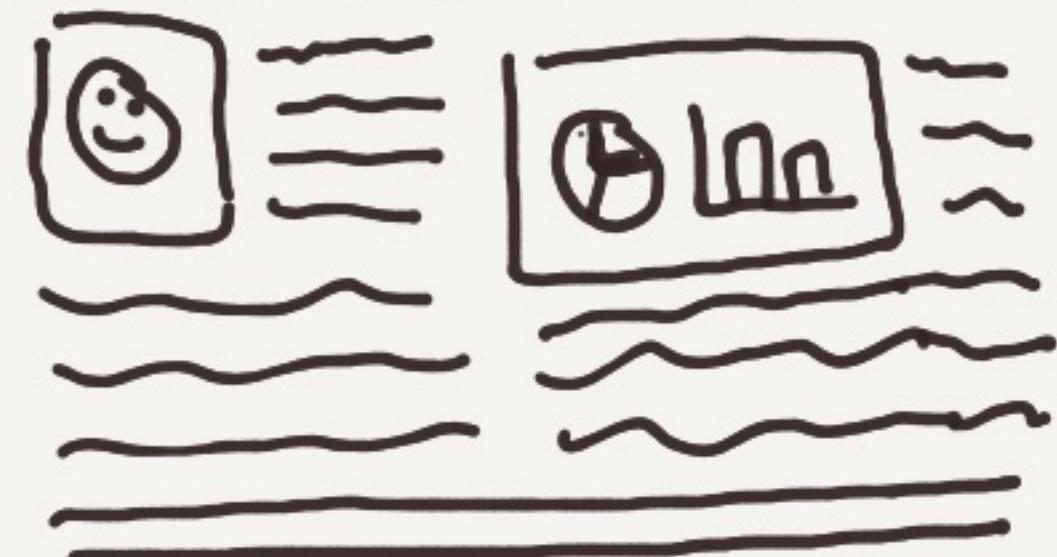
prestige \rightarrow readership
Politzer \rightarrow "subscription + physical copies"

which
newspapers w/ higher
"readership metric" have
more Politzers?

time period: paper
exist?
reach: national vs
local

THE GUARDIAN

EXTRA! EXTRA! EXTRA



Hypothesis

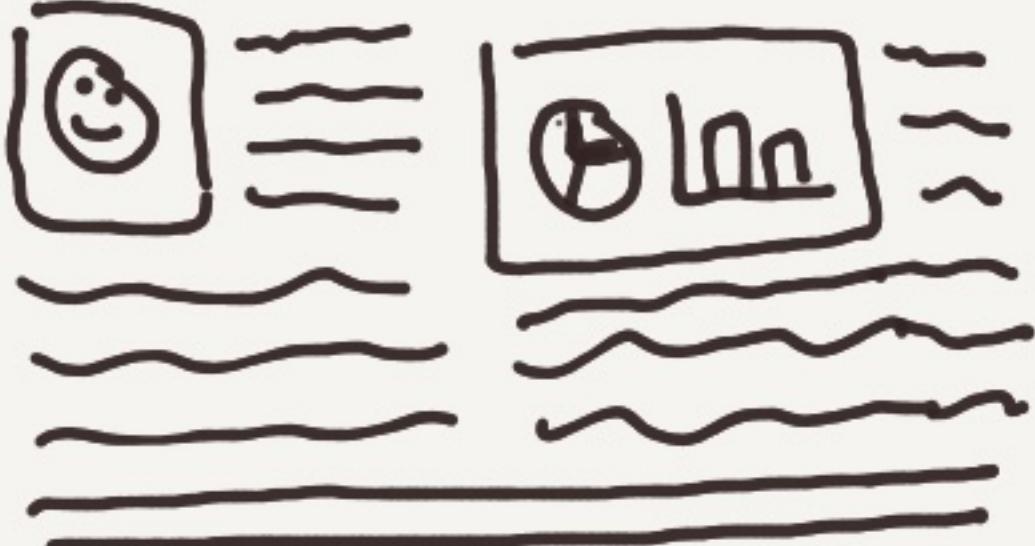
What do you expect the results from this analysis to be?

prestige → readership

- Prestige Pulitzer ≠ readership
 - do readers know how many Pulitzer's read
- effect only seen in Tess papers; New Yorker vs NYT
 - "better reading"
 - "lots of readers"
- will have an effect
opposite, but not super strong
other factors likely playing a larger role

THE GUARDIAN

EXTRA! EXTRA! EXTRA



Data

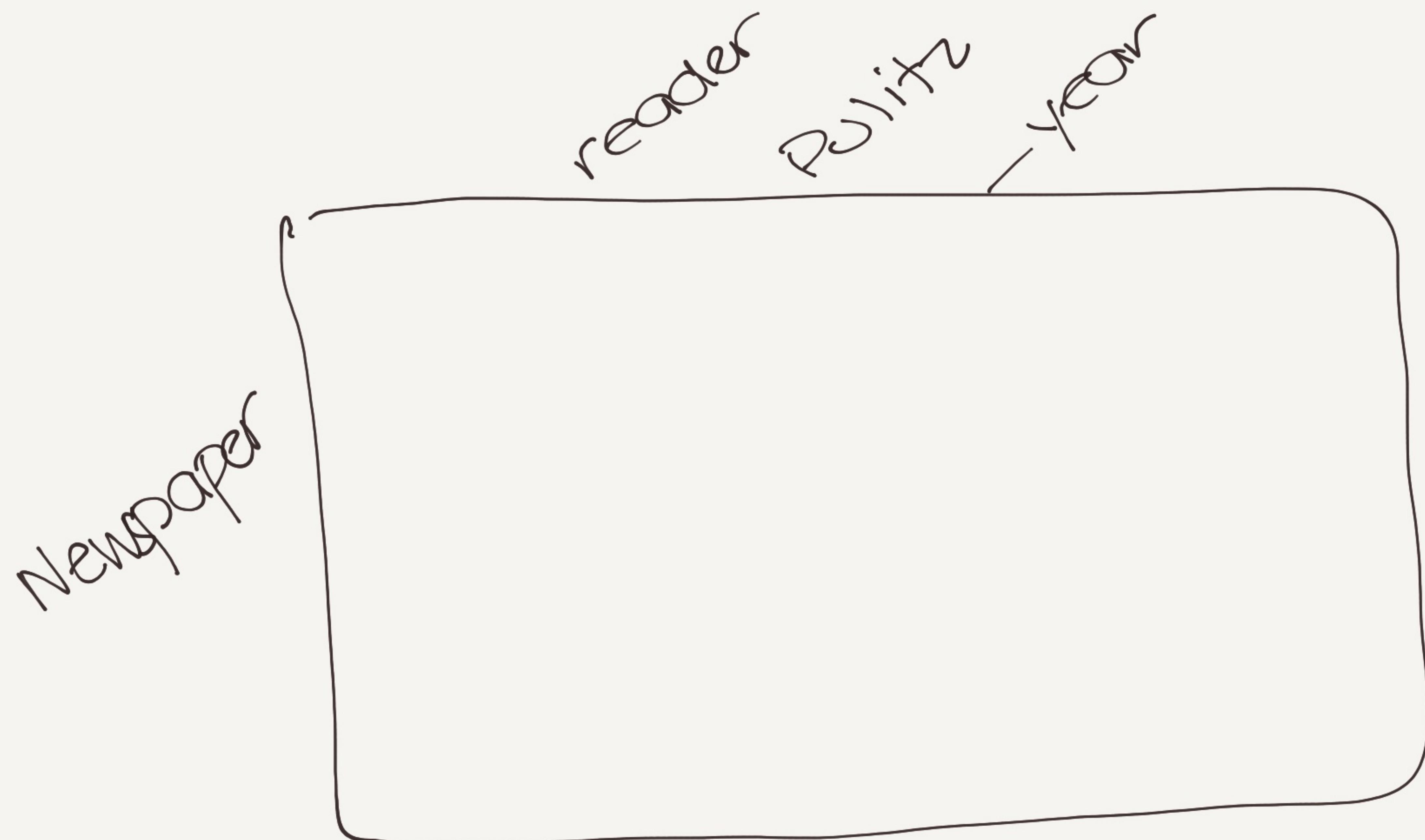
What data/information would you need?

decisions

- # of Pulitzer
 - winners + nominees
- measure of read
 - a single year? over time?
- time period
 - internet popular?
- Which newspapers
 - TOP 50?
 - SD from the mean



How would you want this
information to be stored?
(STRUCTURED)





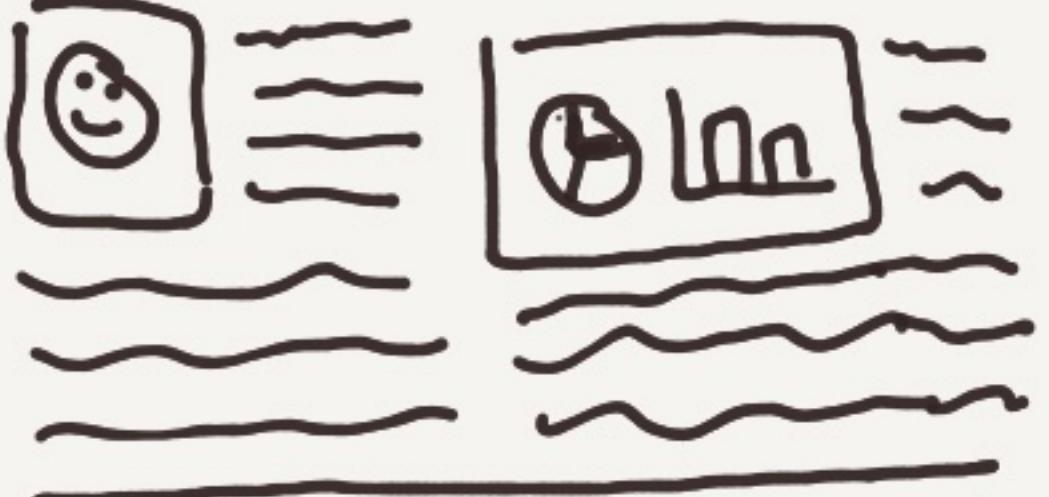
Python

What Python tools would you need to use?

- pandas - dataframe ←
- stats package?
- numpy? - helpful for numbers
- matplotlib - plots
 - o pandas
 - o seaborn

THE GUARDIAN

EXTRA! EXTRA! EXTRA



Dataset

— 50 rows, 7 columns

Newspaper	2004 Daily circulation	2013 Daily circulation	Change in Daily Circulation	Pulitz. 1990- 2003	Pulitz 2004-14	Pulitz 1990- 2014
USA Today	2,192,098	1,674,306	-24%	1	1	2
WSJ	2,101,017	2,378,827	+13%	30	20	50
NYT	:	:	:	:	:	:
LAT	:	:	:	:	:	:
:	:	:	:	:	:	:
:	:	:	:	:	:	:
!	:	:	:	:	:	:
.	:	:	:	:	:	:
:	:	:	:	:	:	:
:	:	:	:	:	:	:
Investor's Business Daily	215,735	157,161	-27%	0	1	1



Q: → What is the effect of prestige (as measured by Pulitzers won) on newspaper readership?

H: - no effect

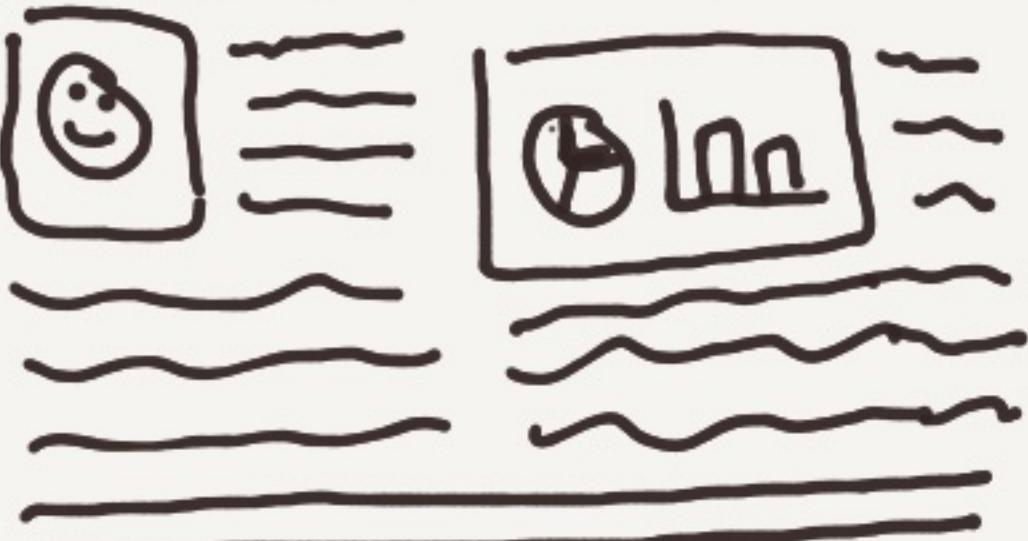
- readers don't care/know
- other things more important

- pos. effect

Data:

THE GUARDIAN

EXTRA! EXTRA! EXTRA

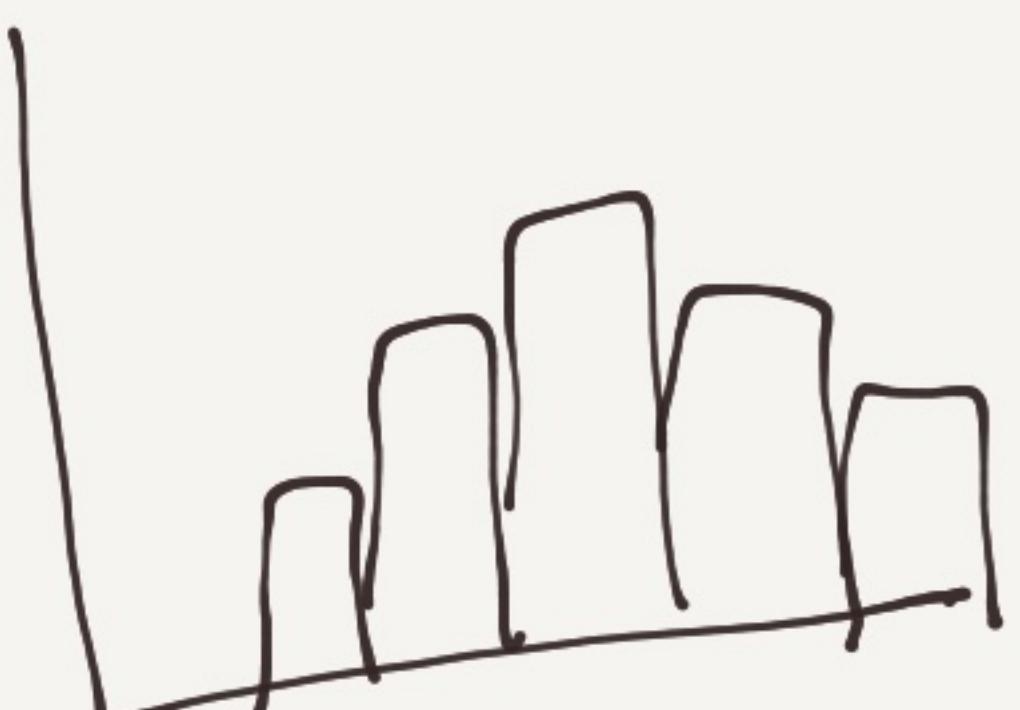


EDA

- overall trend
 - daily circulation
 - o Pulitzer
 - o pick year + look
 - o scatterplot

- before relationship
 - o check linear regression assumptions
 - normality dist
 - log transform?

histograms
- individual dist



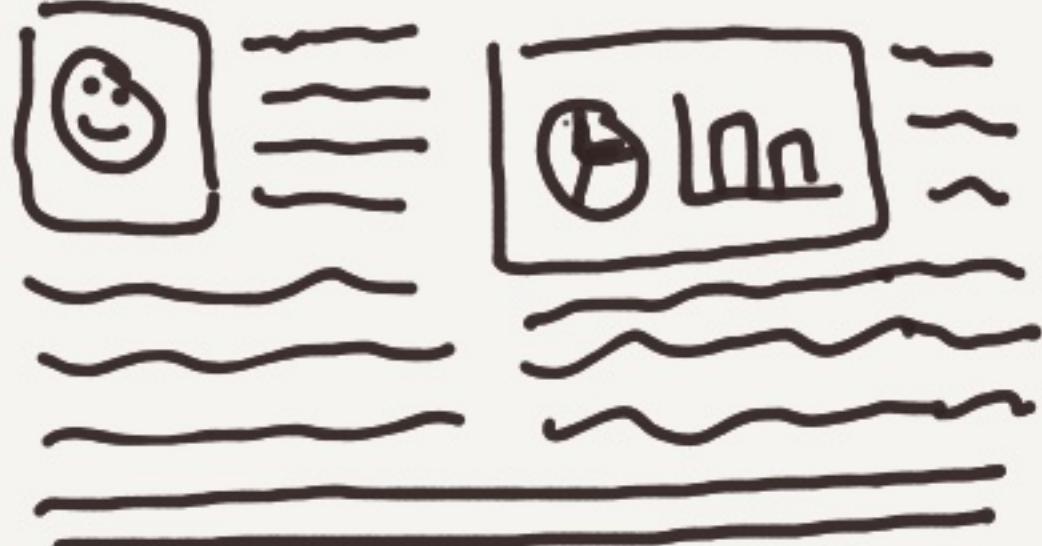
Newspaper	2004 Daily circulation	2013 Daily circulation	Change in Daily Circulation	Pulitz. 1990-2003	Pulitz 2004-14	Pulitz 1990-2014
USA Today	2,192,098	1,674,306	-24%	1	1	2
WSJ	2,101,017	2,378,827	+13%	30	20	50

- How would you explore the clatci?
- change in Pulitzer + change in circulation
 -
 -
- scatterplot

 - missingness?
 - * no missing data

THE GUARDIAN

EXTRA! EXTRA! EXTRA



EDA

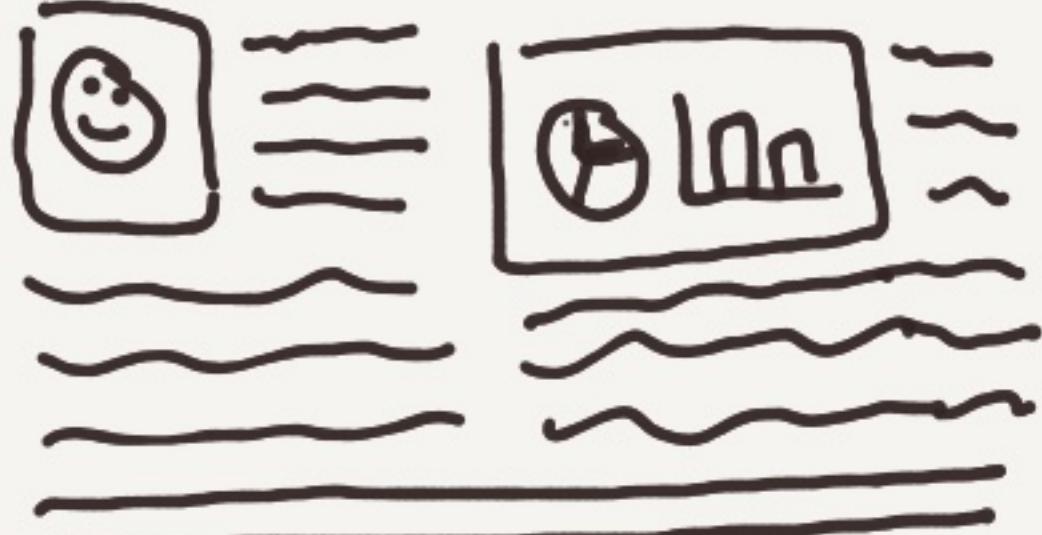
what exploratory visualizations would you generate?

- check extremes
 - really high/low?

Newspaper	2004 Daily circulation	2013 Daily circulation	Change in Daily Circulation	Pulitz. 1990-2003	Pulitz 2004-14	Pulitz 1990-2014
USA Today	2,192,098	1,674,306	-24%	30	20	50
WSJ	2,101,017	2,378,827	+13%			

THE GUARDIAN

EXTRA! EXTRA! EXTRA



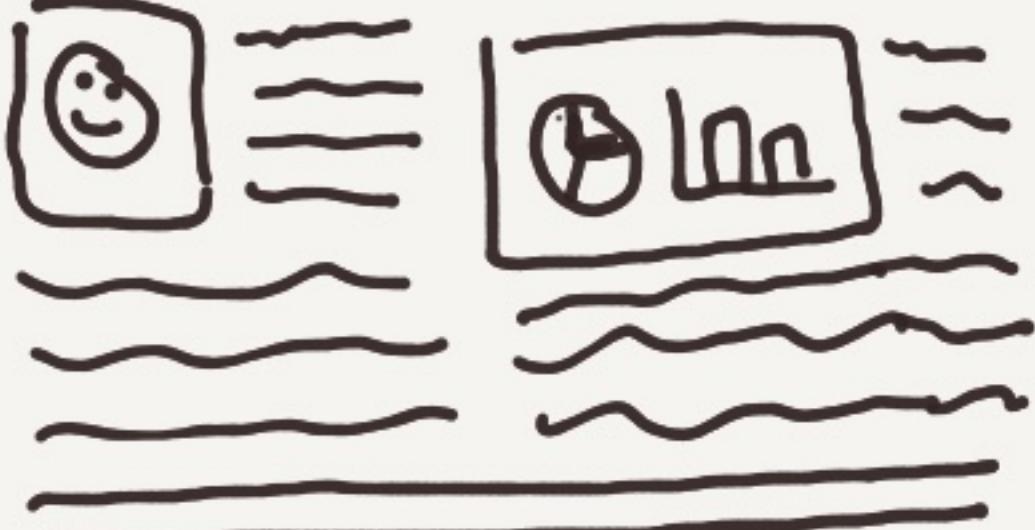
EDA

what exploratory
visualizations would
you generate?

Newspaper	2004 Daily circulation	2013 Daily circulation	Change in Daily circulation	Pulitz. 1990-2003	Pulitz 2004-14	Pulitz 1990-2014
USA Today	2,192,098	1,674,306	-24%	30	20	50
WSJ	2,101,017	2,378,827	+13%			

THE GUARDIAN

EXTRA! EXTRA! EXTRA



Data

What if our data didn't look as we anticipated?

- outliers
- non-normal distributions

- transformations

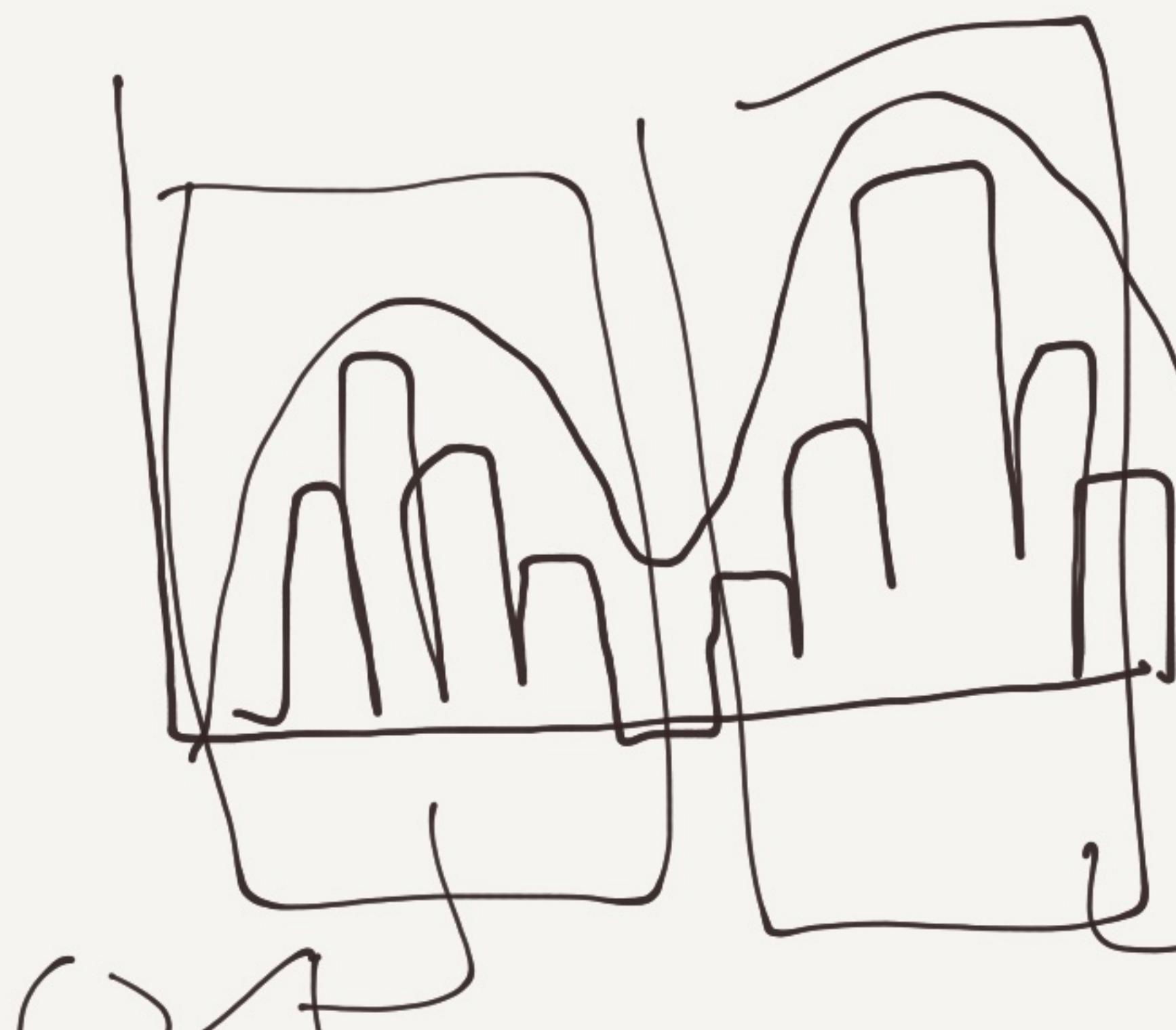
- log
- χ^2
-

↳ check normality test for normality

- bimodal



↳ stratification



"recode"

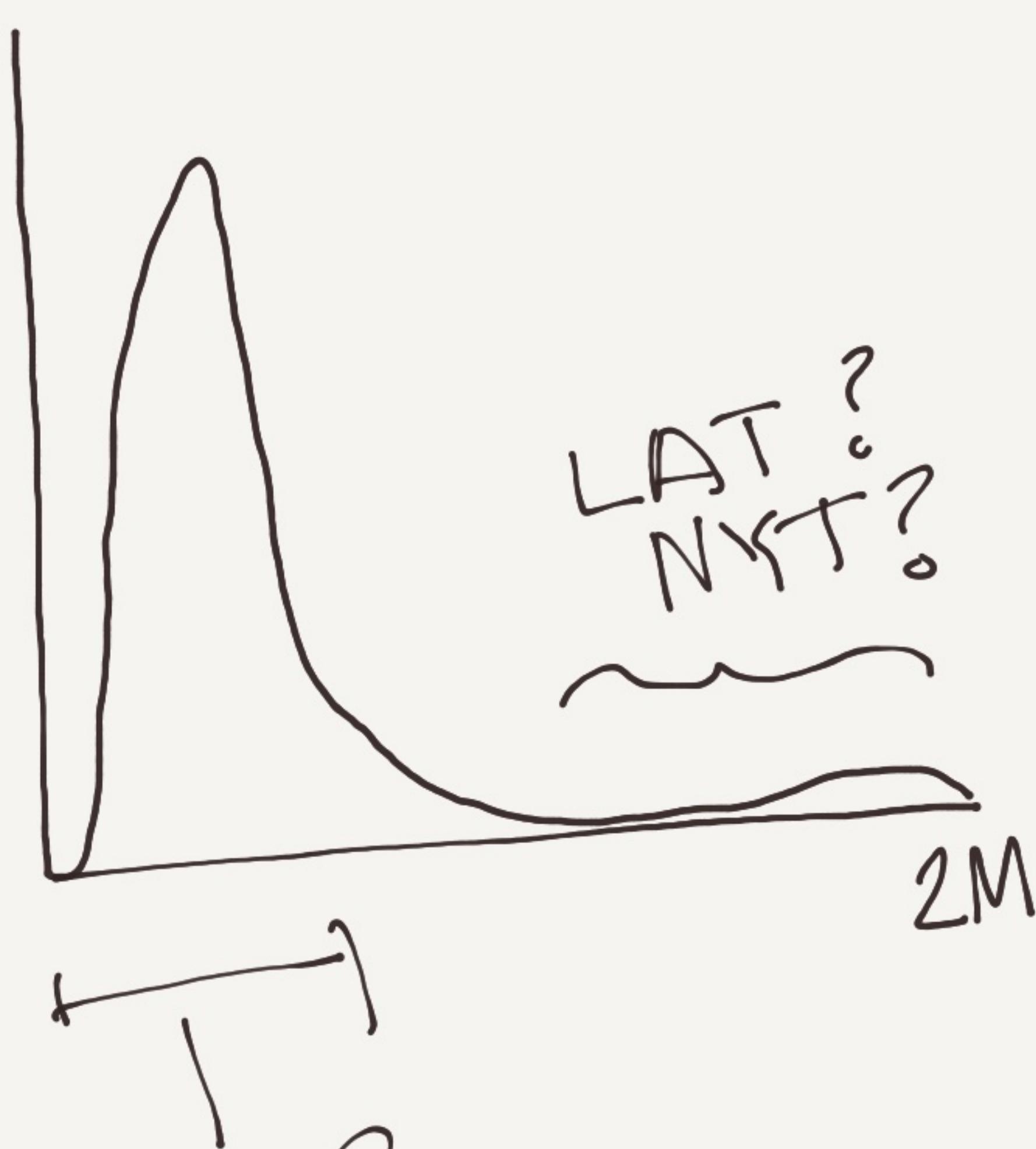


↳ interpretation



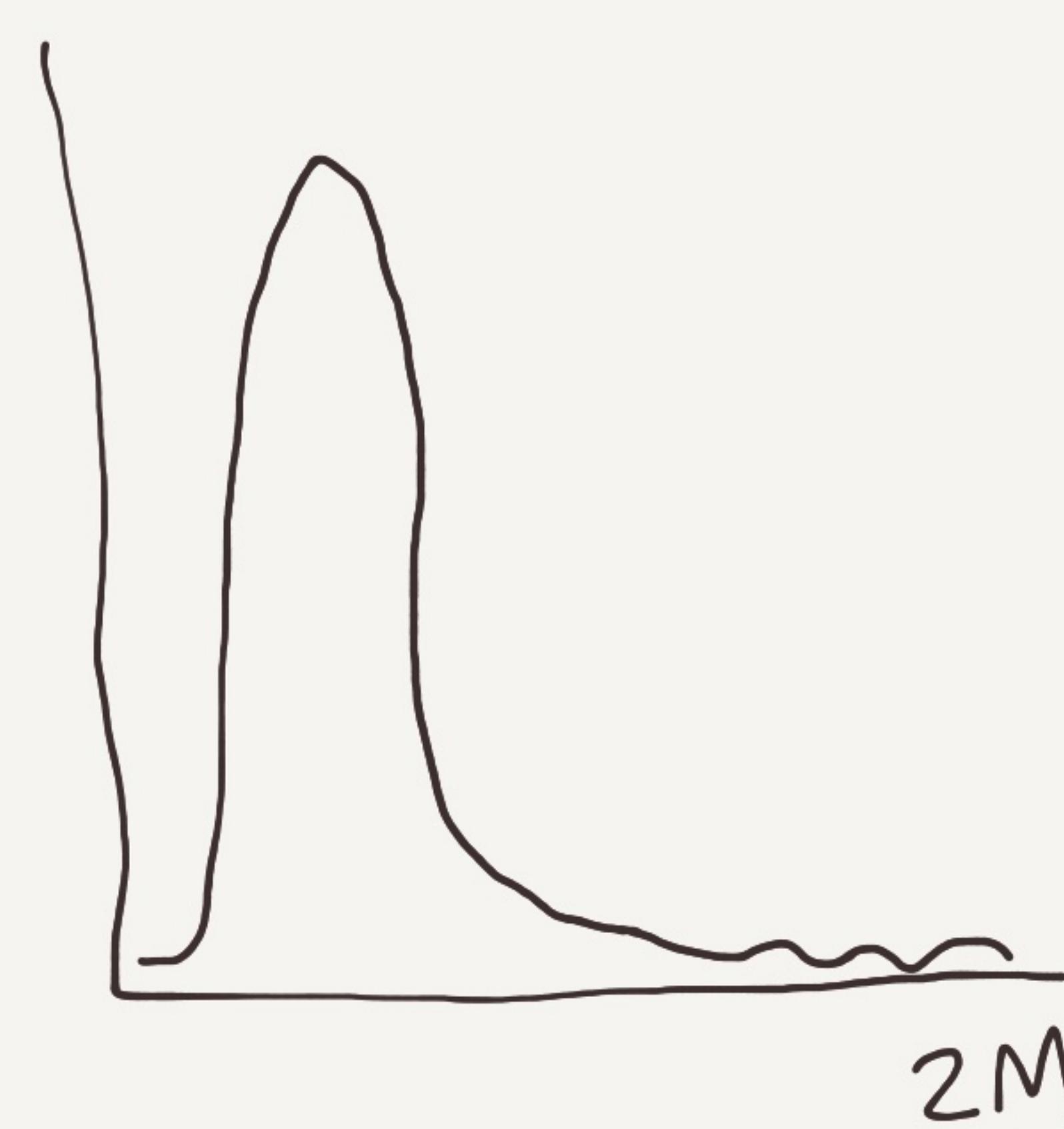
EDA: circulation (Readership)

Daily Circulation
2004

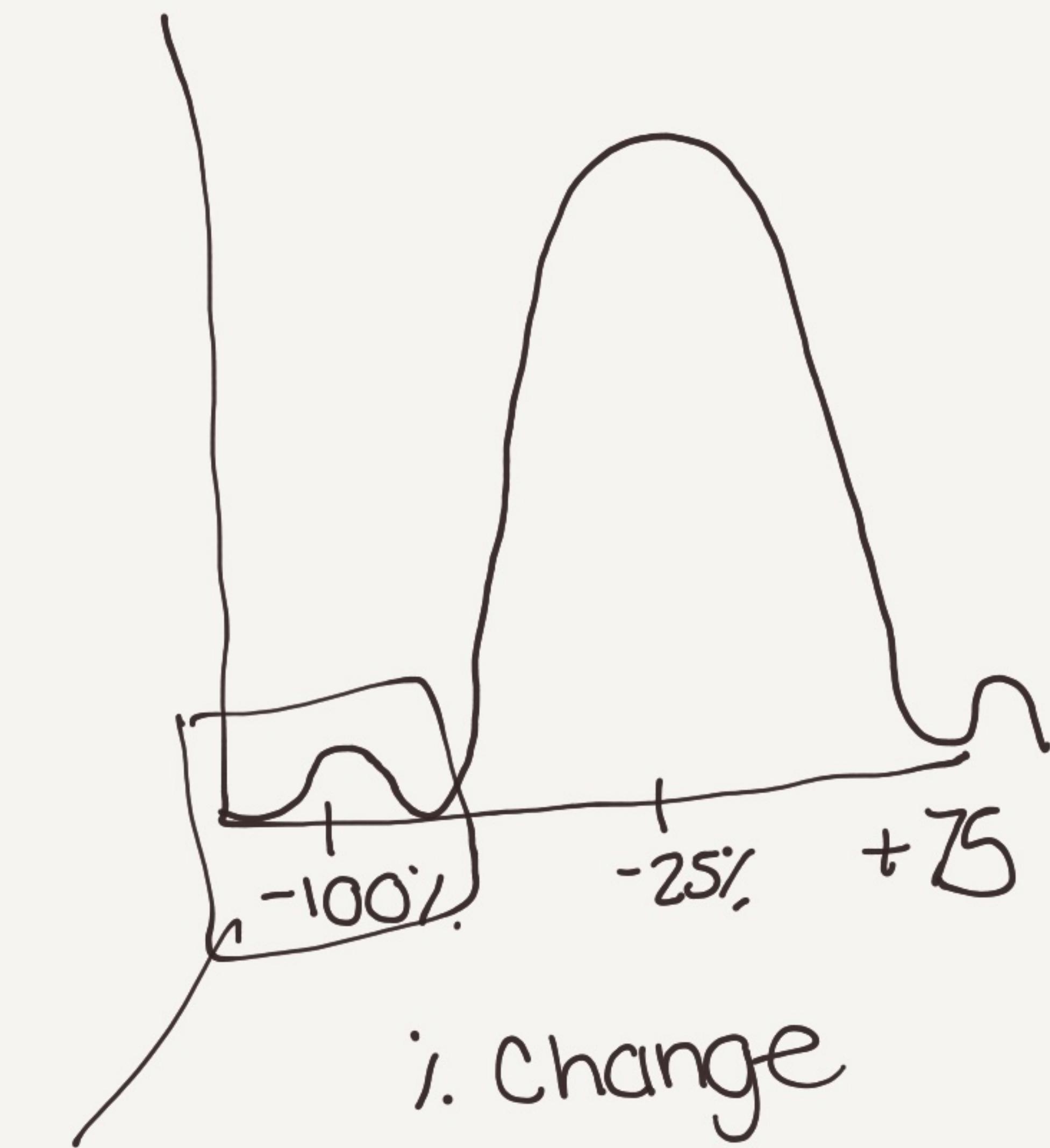


lots of publications w/ lower circulation

Daily Circulation
2013



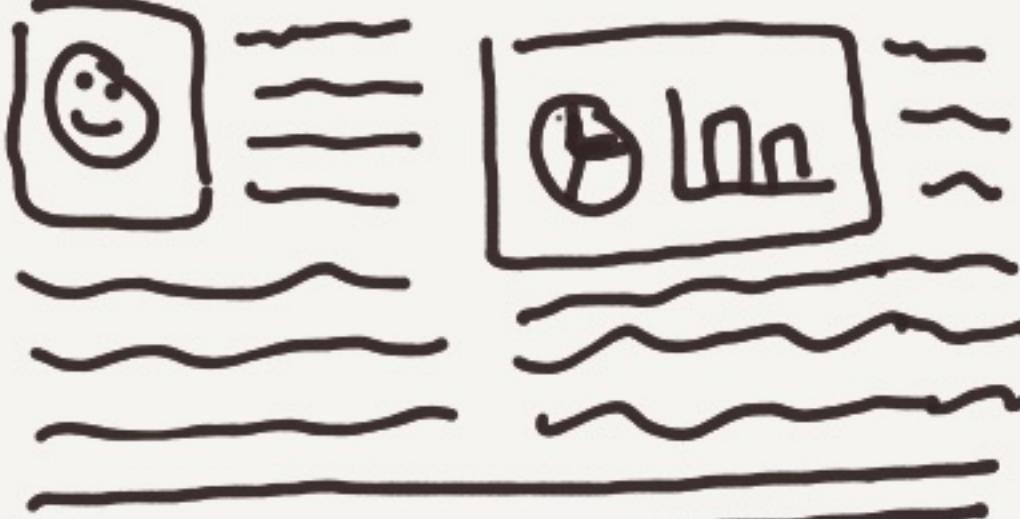
Change in Daily Circulation 2004-2013



Outliers
- Out of business
- out of daily circulation

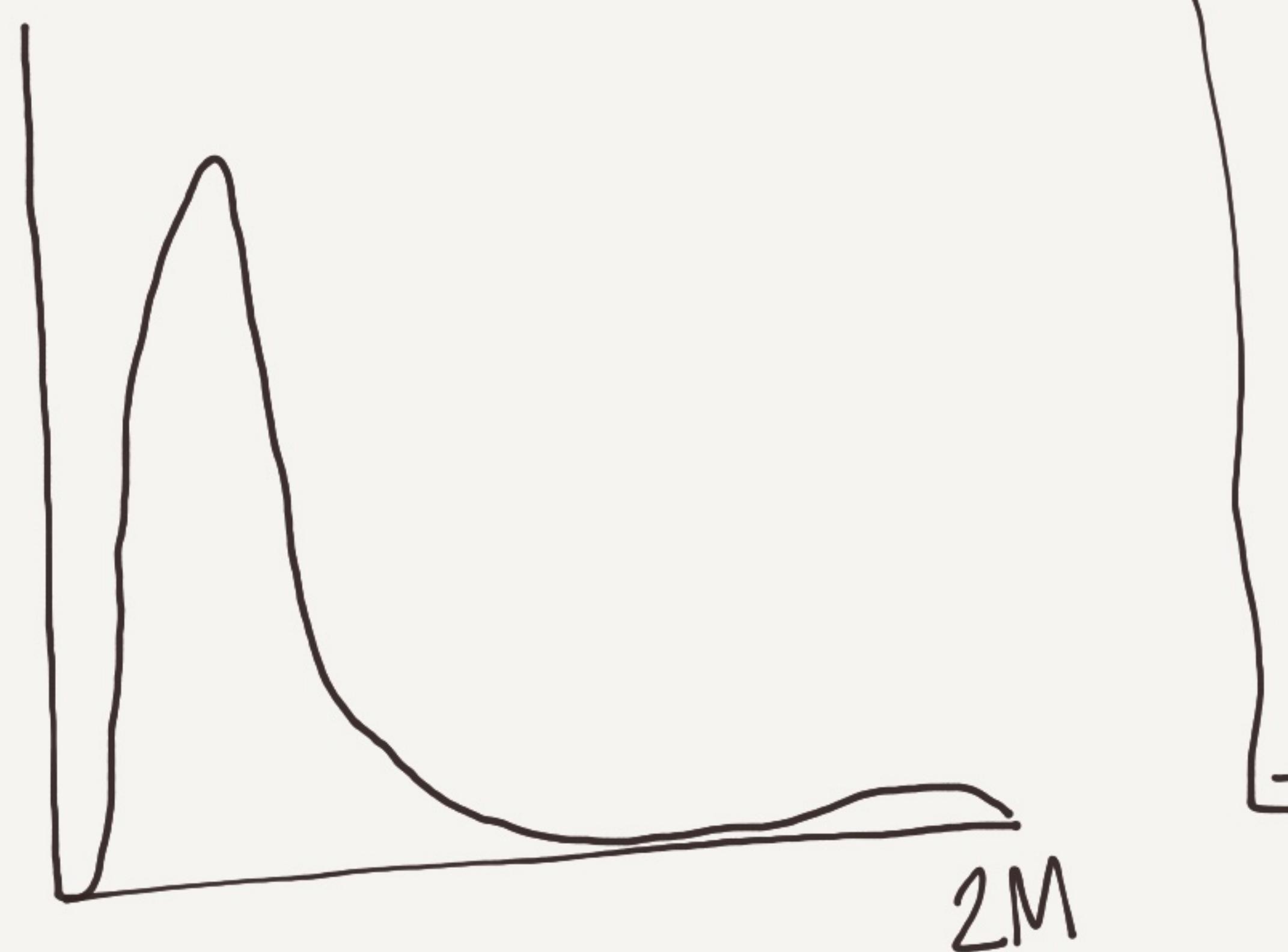
THE GUARDIAN

EXTRA! EXTRA! EXTRA

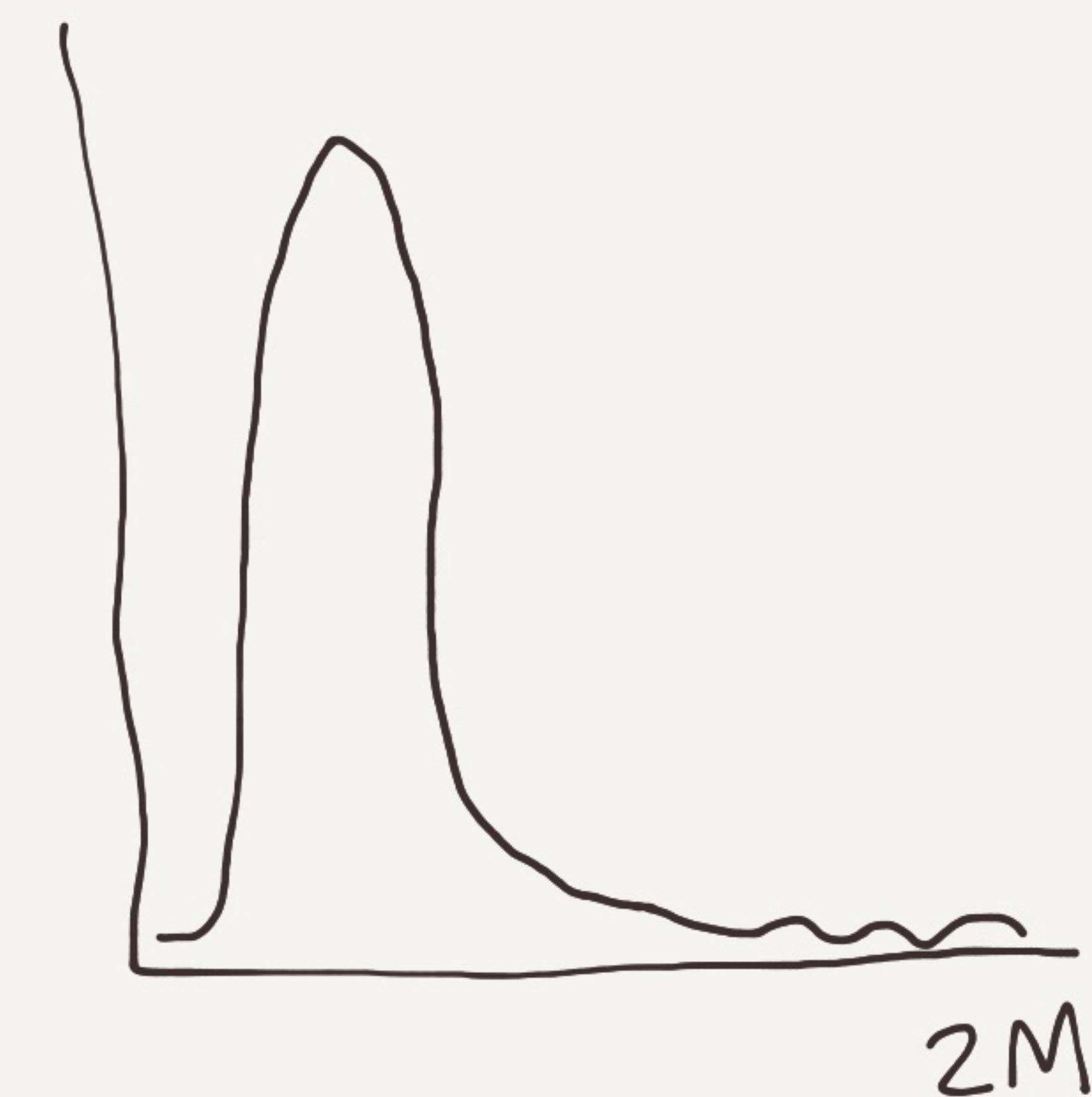


EDA

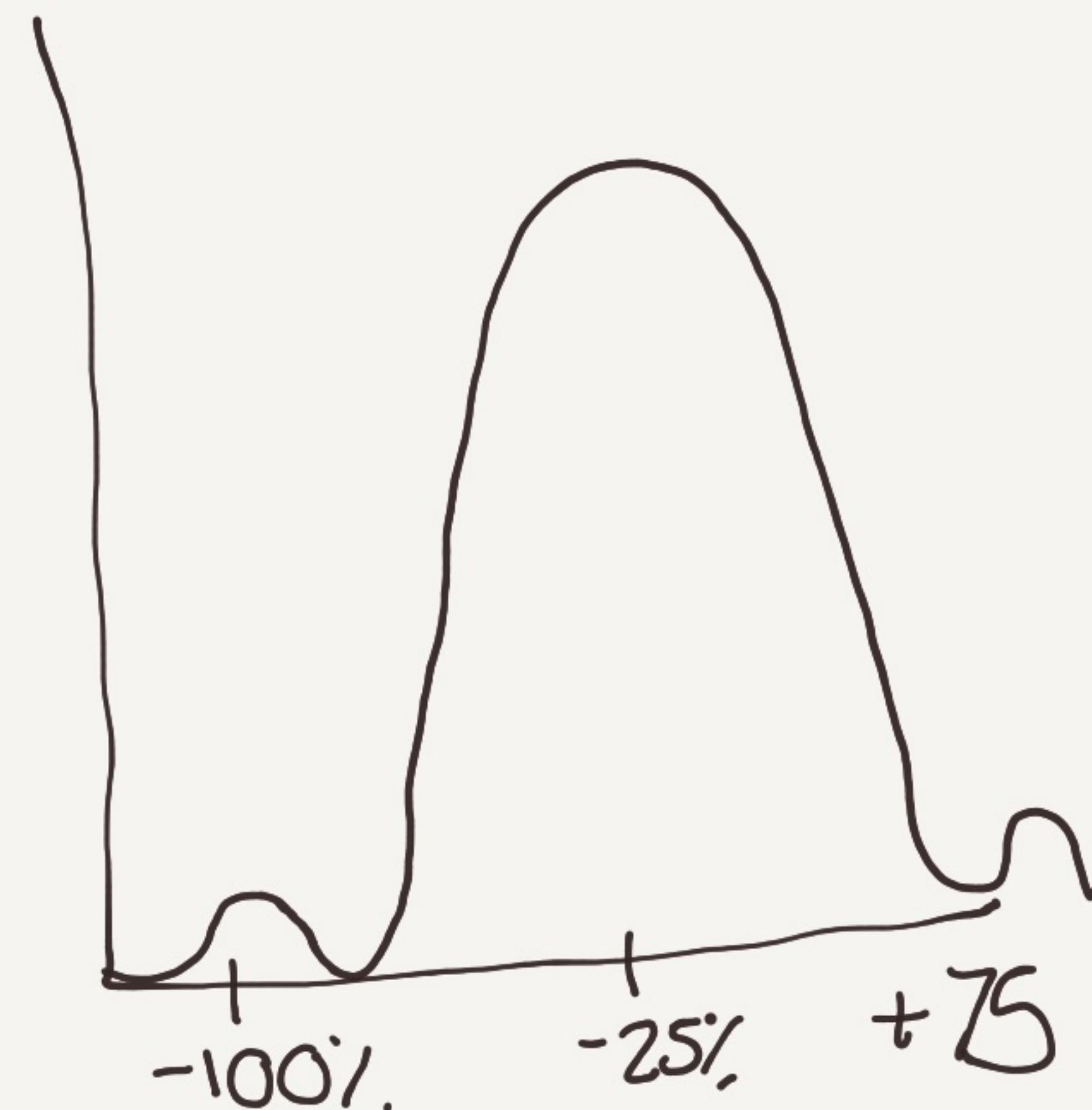
Daily Circulation
2004



Daily Circulation
2013



Change in Daily
Circulation
2004-2013



Clicker
Question :

Which variable would you use
for this analysis?

i. change



A Daily circulation, 2004 (18%)

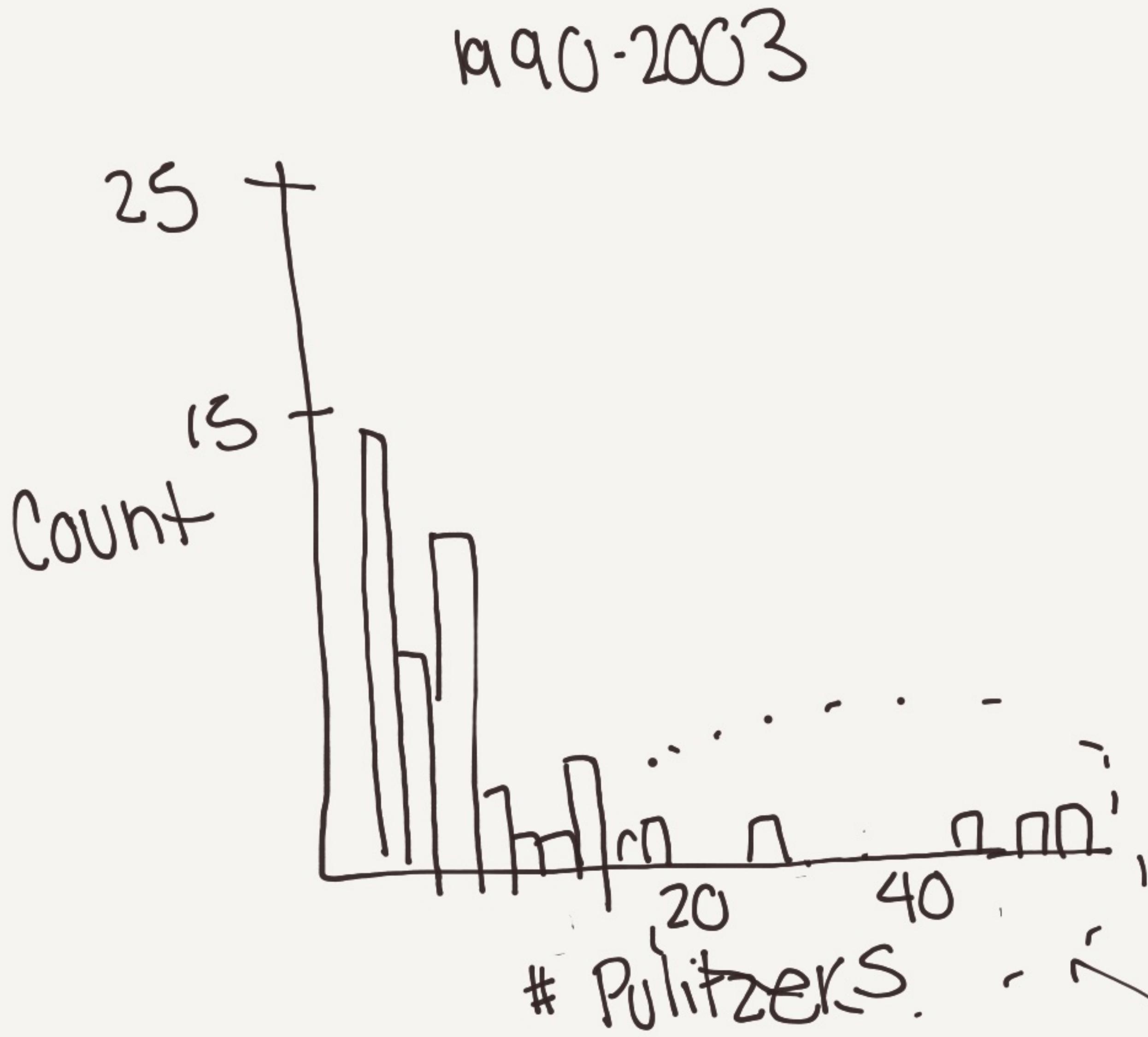
B Daily circulation, 2013 (17%)

C Change in Daily circulation (64%)

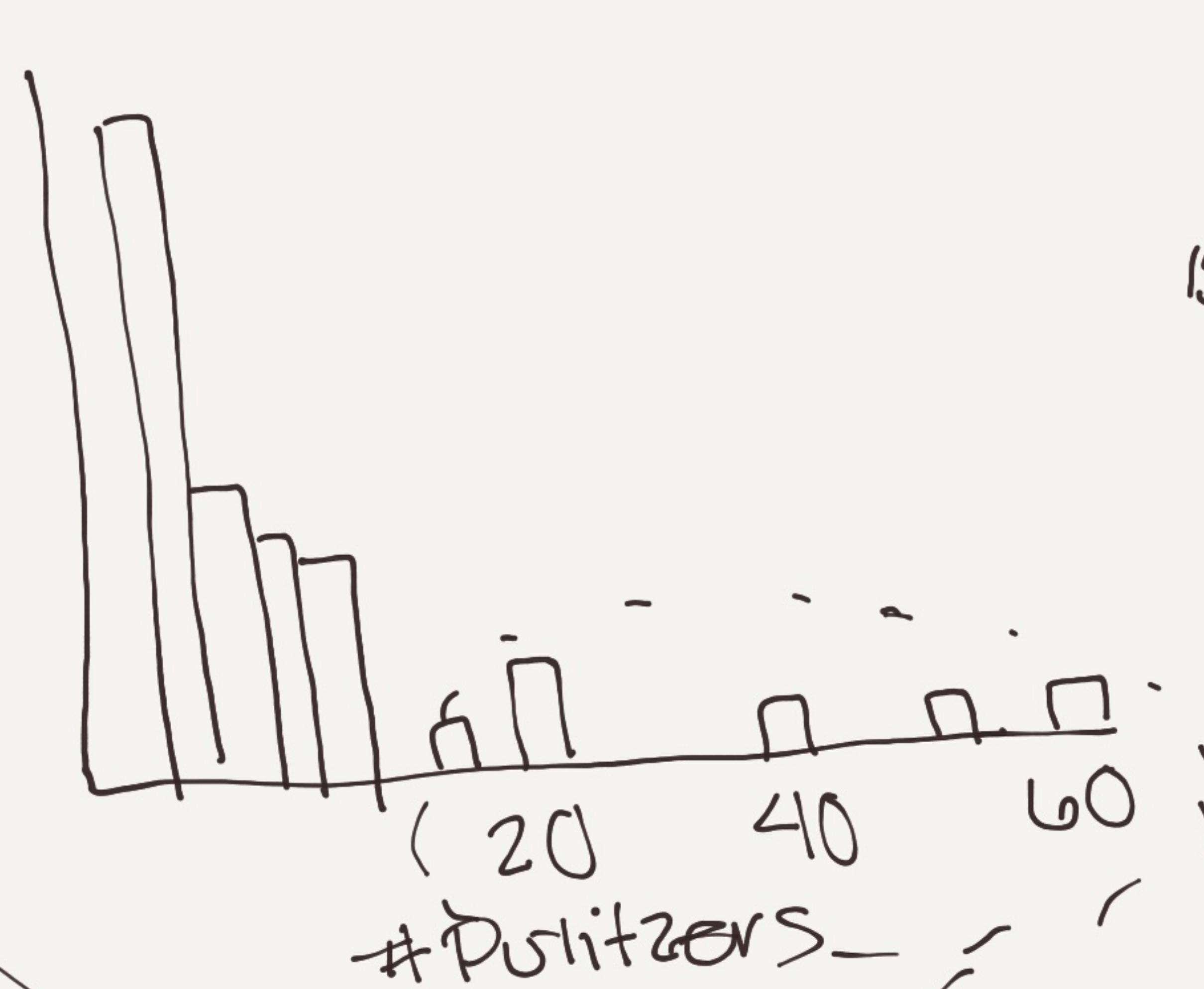


EDA : Pulitzer

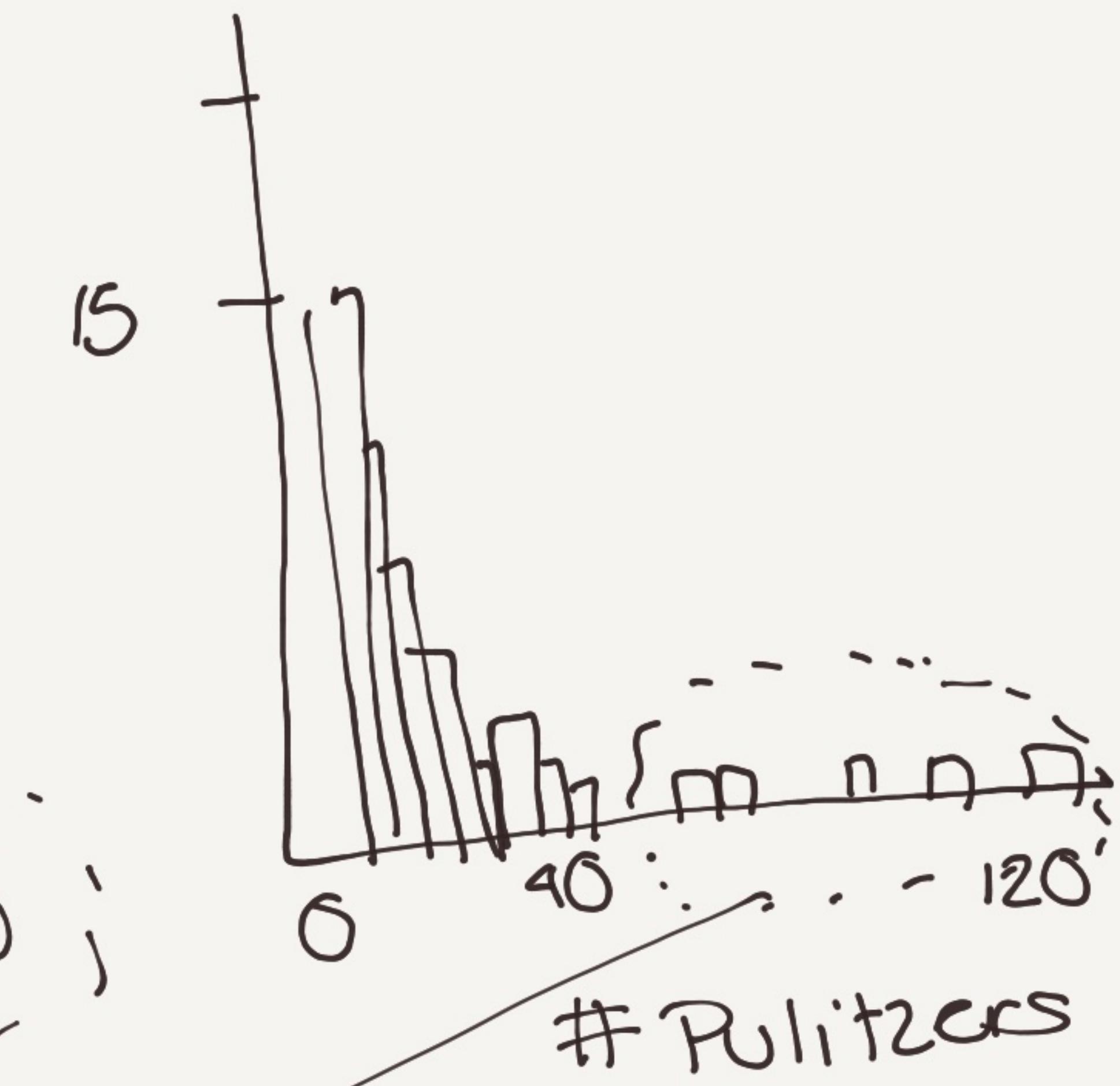
↳ # Winners + Finalists



2004-2014



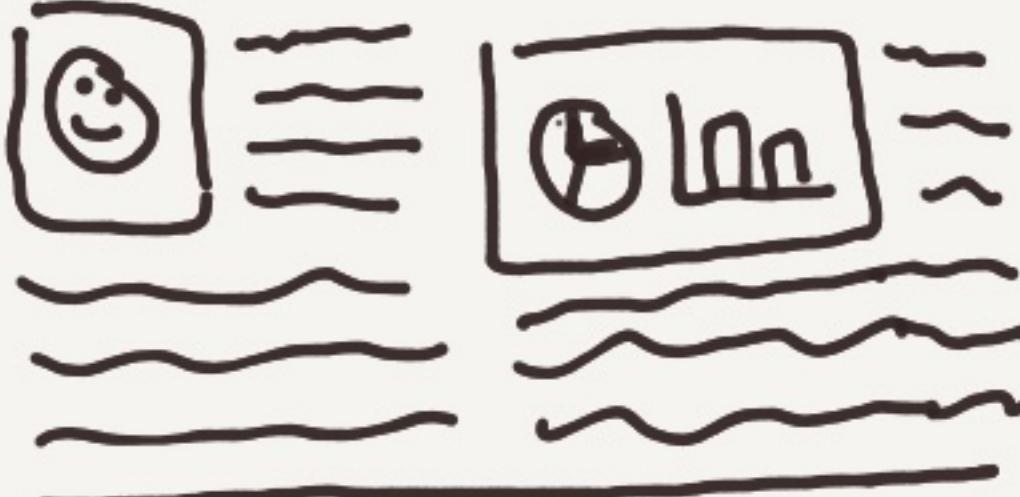
1990-2014



a few publications
get way more
Pulitzers

THE GUARDIAN

EXTRA! EXTRA! EXTRA



Analysis

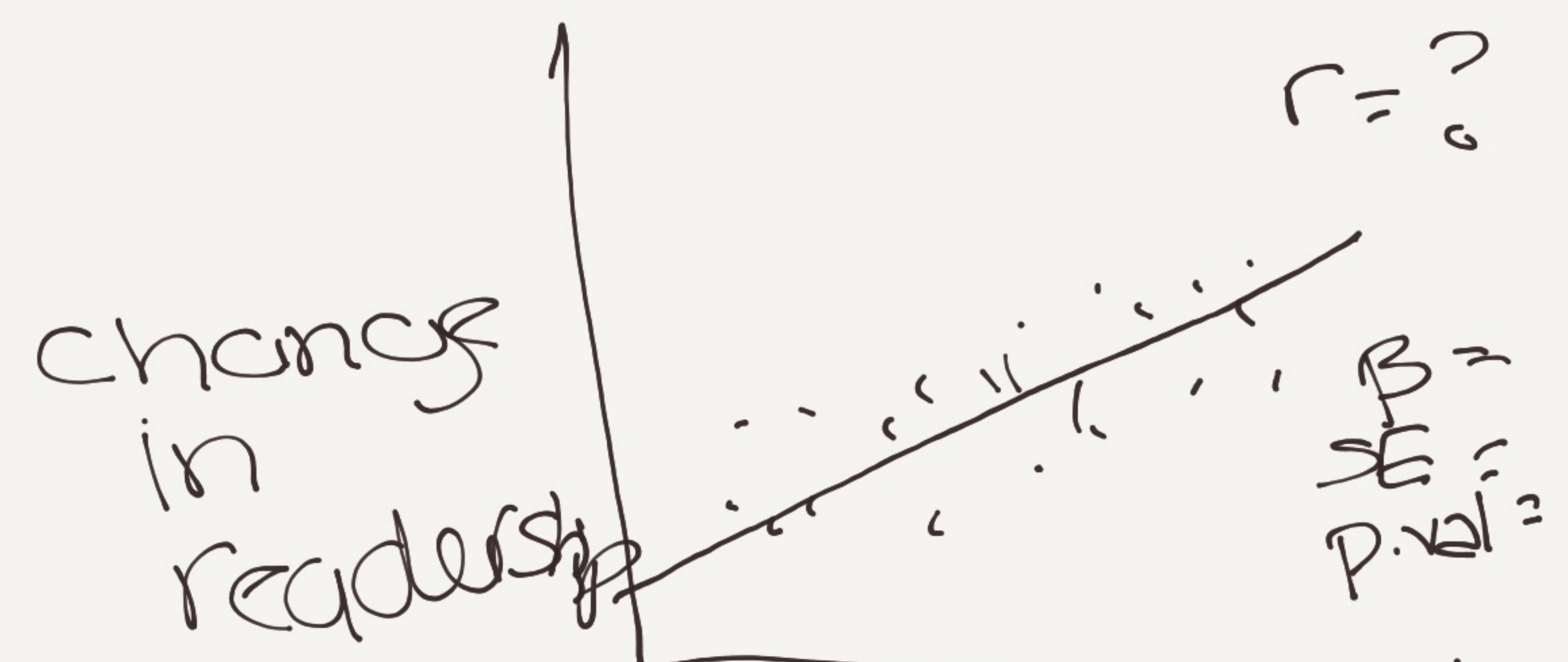
How would you answer the question?

what
the
effect?

prestige → readership
(Poltziers)

H_0 : There is no effect (correlation)

H_a : There is an effect (correlation)

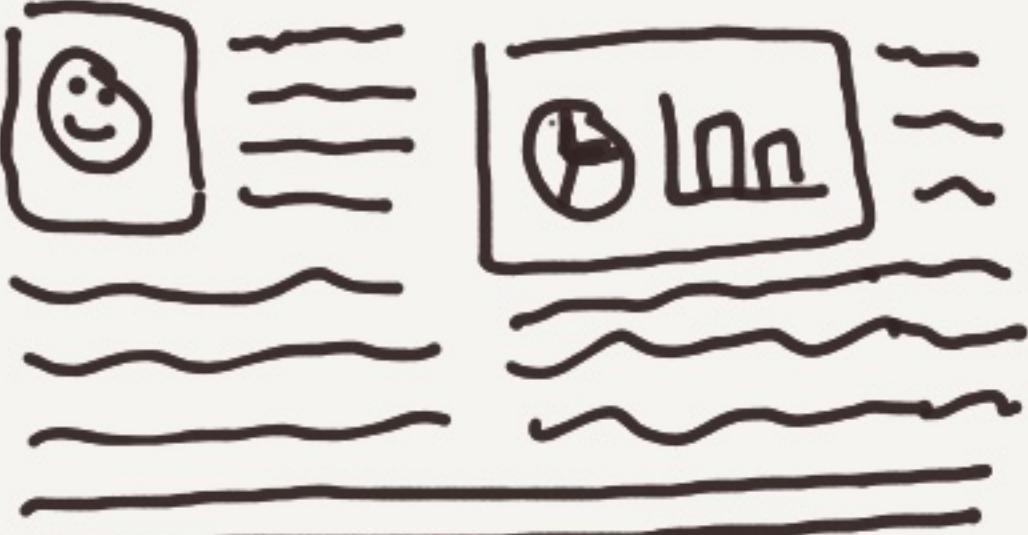


#Poltziers

- year
- across all years?

THE GUARDIAN

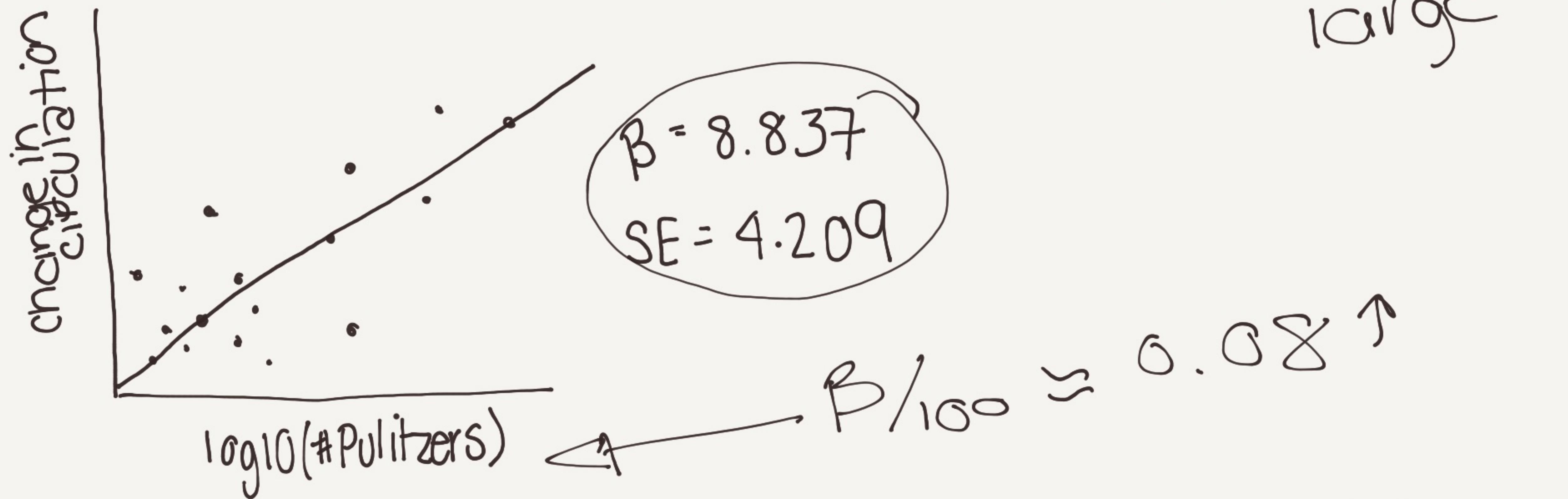
EXTRA! EXTRA! EXTRA



INTERPRETATION

CLICKER
★

GIVEN THE FOLLOWING RESULTS, WHAT
WOULD YOU CONCLUDE?

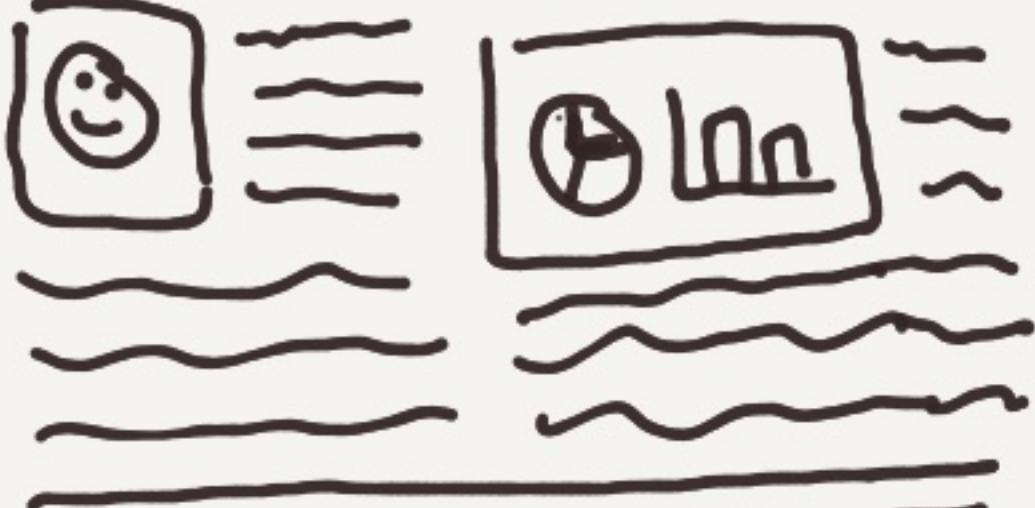


-A; SE looks large

- (70%) A. Prestige increases circulation
 - (12%) B. Prestige decreases circulation
 - (7%) C. Prestige does not affect circulation
 - (11%) D. Something else
- ↳ correlation

THE GUARDIAN

EXTRA! EXTRA! EXTRA



Limitations

What limitations are there
in this analysis?

- small database +
- omitted variables
 - quality of newspapers
 - confounders?
- other prestige?
 - citations