

THE GUARDIAN

UNIVERSITY OF CALIFORNIA, SAN DIEGO



Extra! Extra! COGS108: DATA SCIENCE IN PRACTICE

GITHUB USERS
HAVE DESIRED
JOB SKILLS!

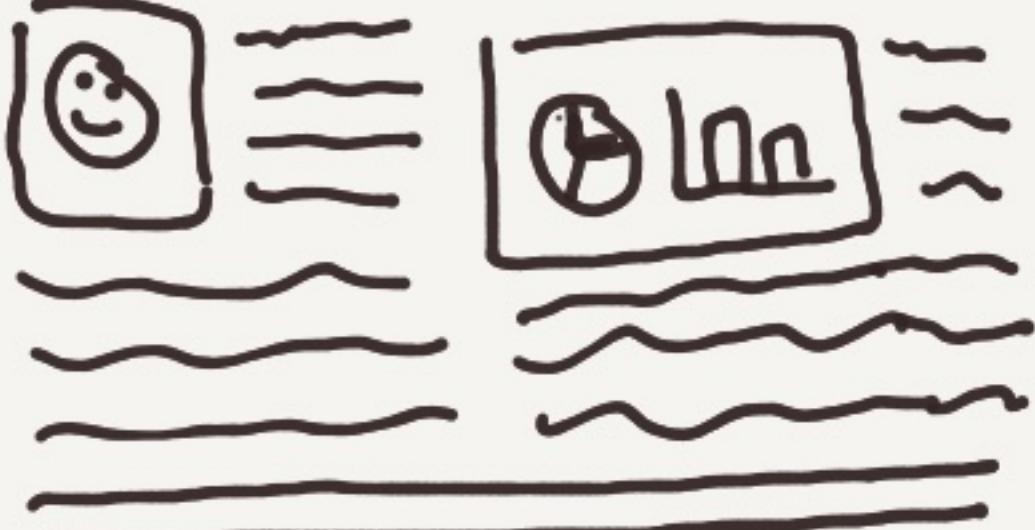
Inferential
Analysis

DATA SCIENCE
IS THE FUTURE

COGS108 STUDENTS
BLOW PROFESSOR AWAY
WITH INTERESTING FINAL
PROJECTS.

THE GUARDIAN

EXTRA! EXTRA! EXTRA



Question

Prestige?

- how long in circulation?
 - older = more prestige
- { blocked in other countries
 - (- hate they receive)

→ Pulitzer prize winners writing for the newspaper

→ circulation?

(→ do we care about readership ability) want to know
o free time?

→ price of subscription ← proxy
↑ prestige $\text{P\$}/\text{md}$

→ revenue

↑ revenue ↑ influence prestige

→ cited by others (in general)
↑ citations ↑ prestigious

Does prestige increase newspaper readership?

readership

→ # of readers

(→ type of readers)

→ online

- how long spent?
- scroll rate

→ circulation?

- (just recycle)

- multiple / paper

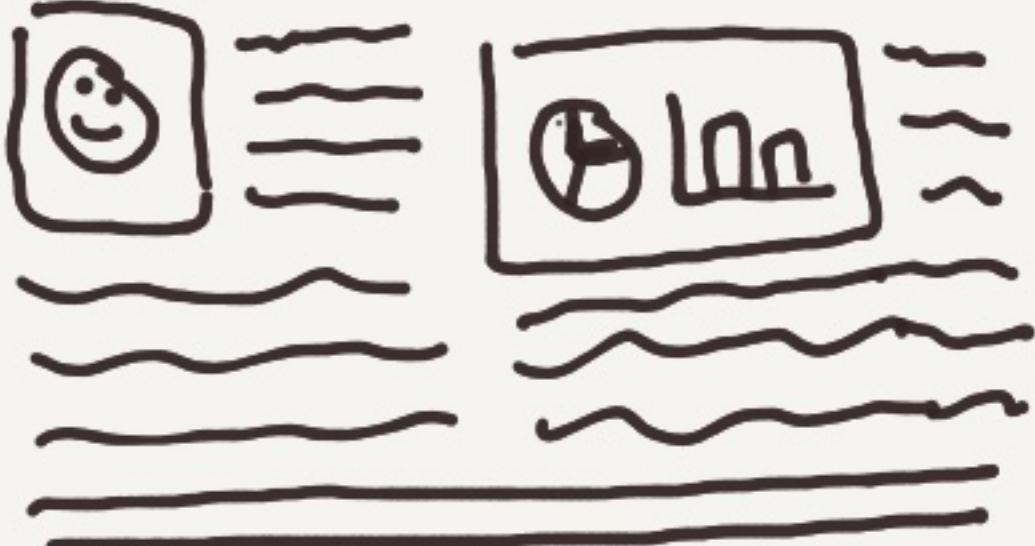
{ est.

→ online?

- subscription
- frequent readers

THE GUARDIAN

EXTRA! EXTRA! EXTRA



Question

Which papers?

What time period?

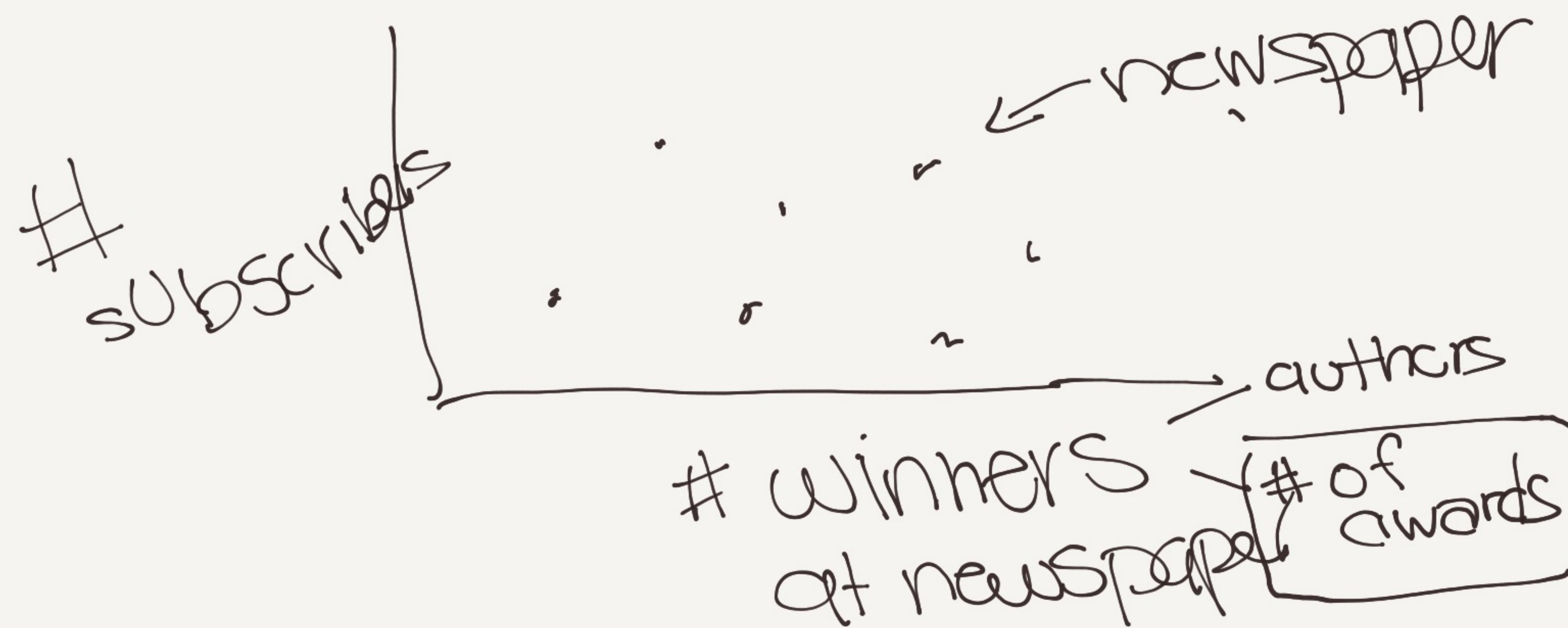
What measured?

What is the effect
of ^{the number of} pulitzers a
newspaper receives
on its # of
readers

Data Science Question

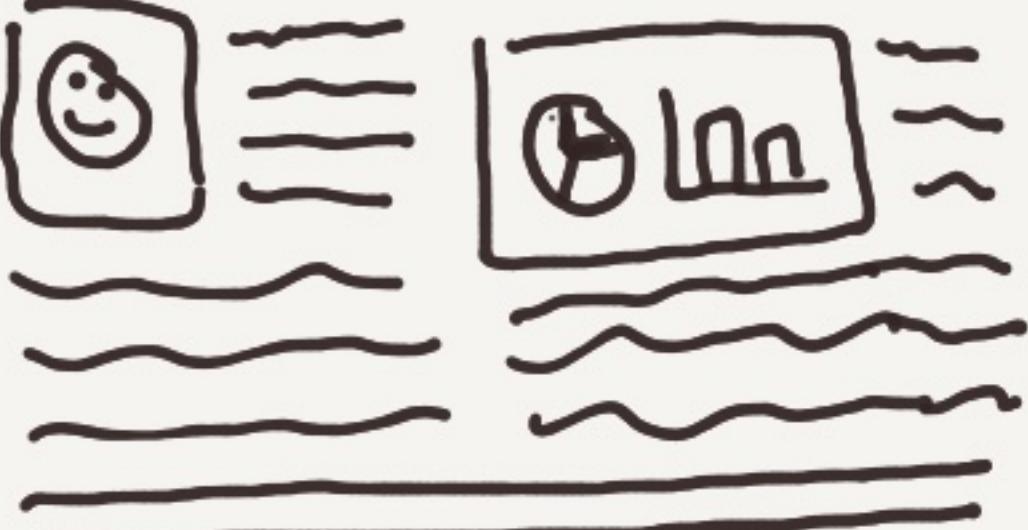
pulitzer → circulation?

How many pulitzer prize
winners at a newspaper
affect # of subscribers?



THE GUARDIAN

EXTRA! EXTRA! EXTRA



Hypothesis

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

pulitzers → circulation
↑ winners ↑ subscribers ↑ quality of writing ↑ confounder

- ↑ winners ↑ subscribers
→ Pulitzer = better / more interesting
→ people would want to read

Popular institutions



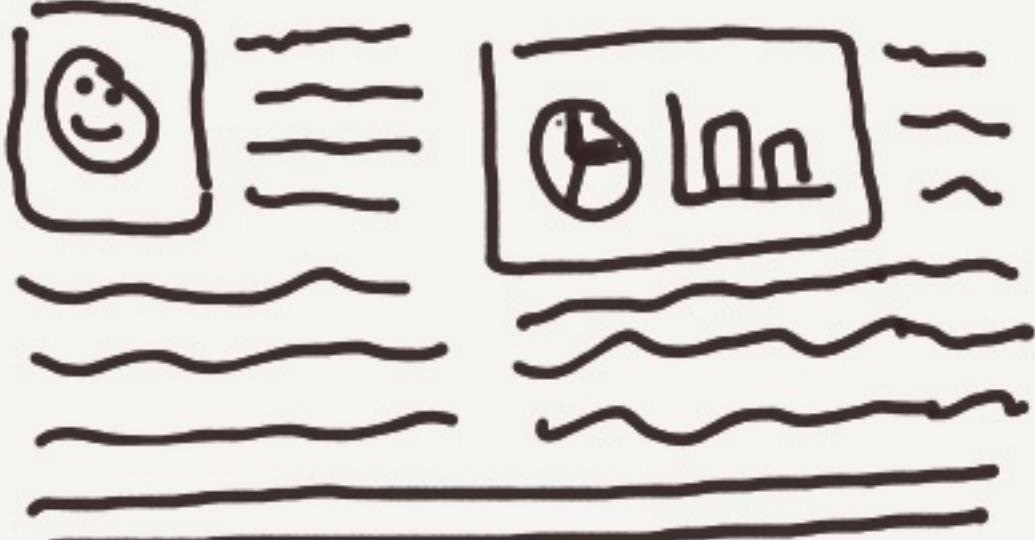
circulation
may NOT
explain
pulitzers

- other things
could be at
play

- No effect
→ other things at play - how extreme clickbait
→ readers don't know about
pulitzers

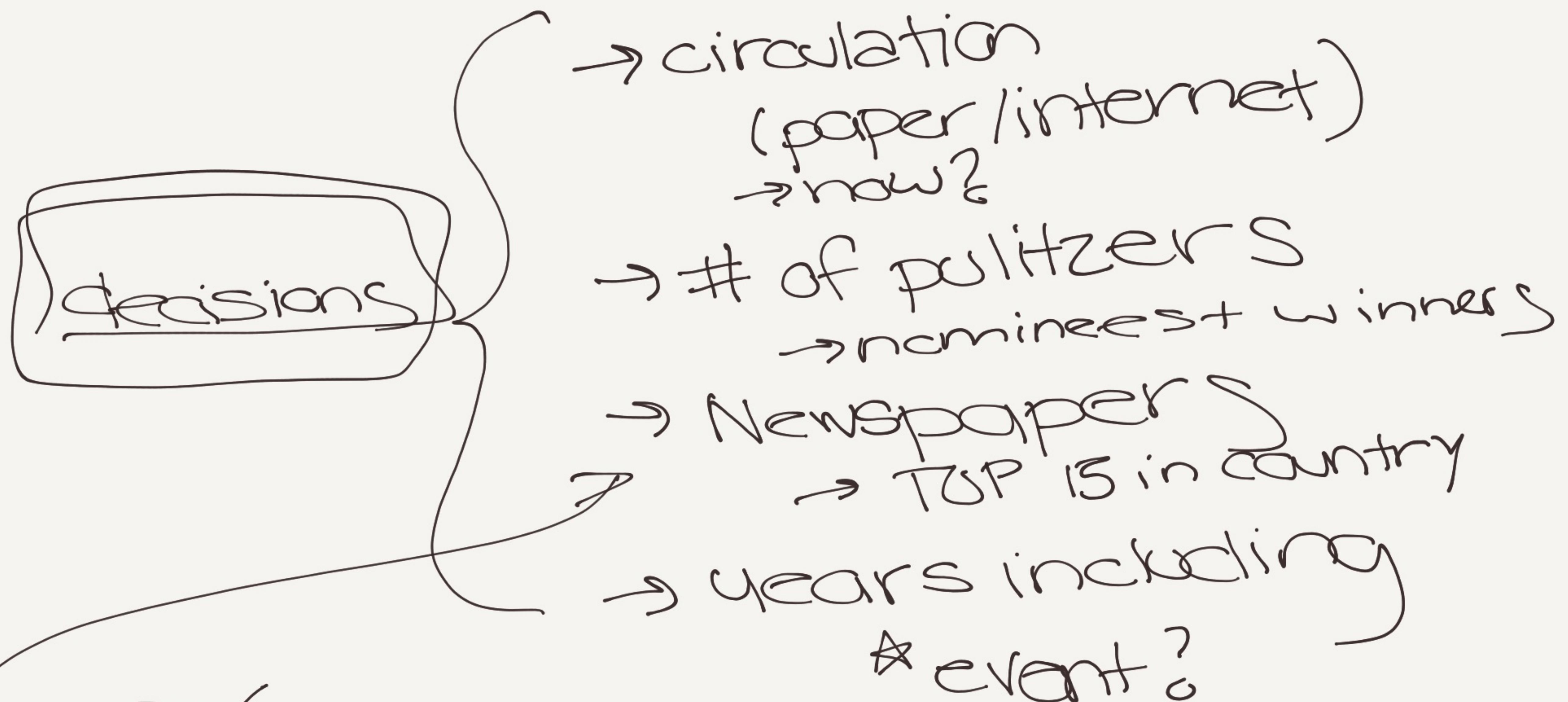
THE GUARDIAN

EXTRA! EXTRA! EXTRA



What data/information would you need?

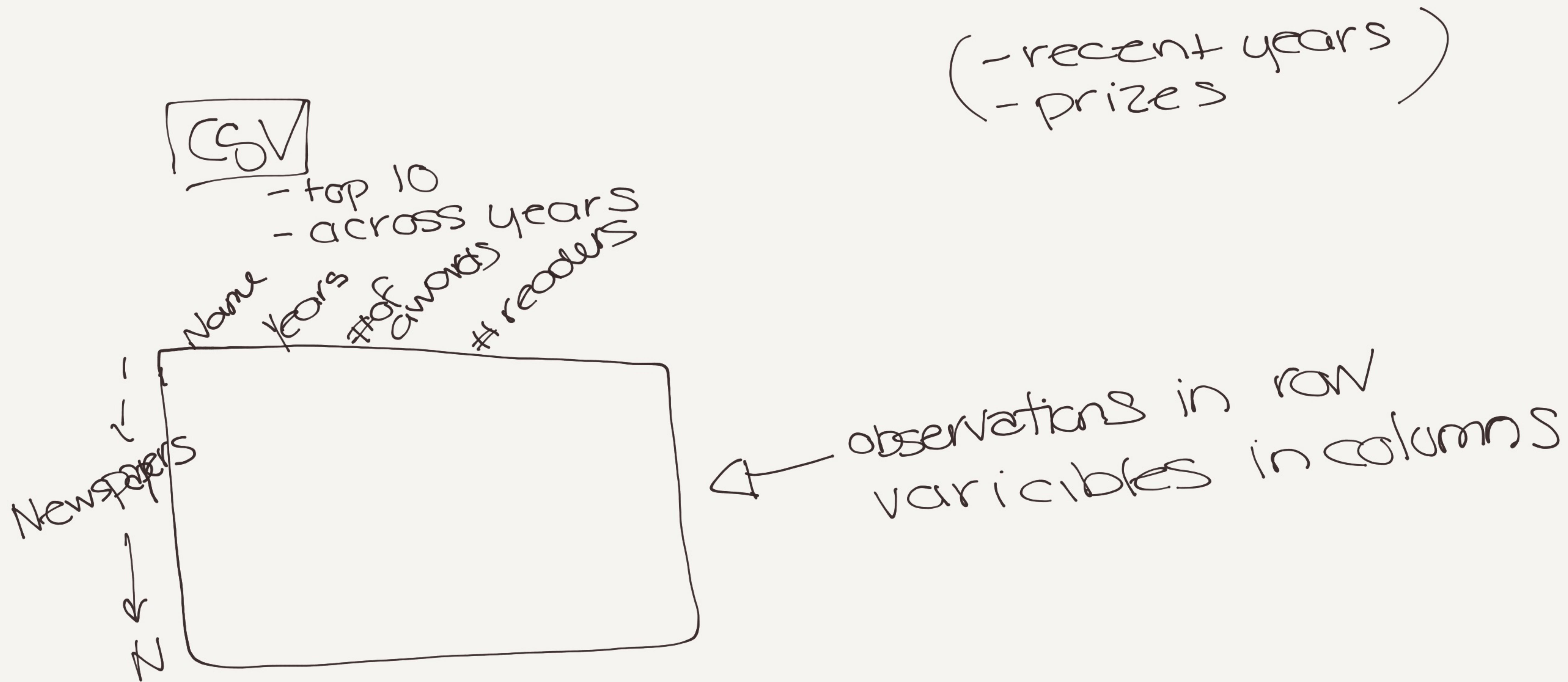
Data



- any newspaper w/ pulitzer?
- at least 5
- \bar{x} pulitzers
 $\bar{x} - 2 * SD$?

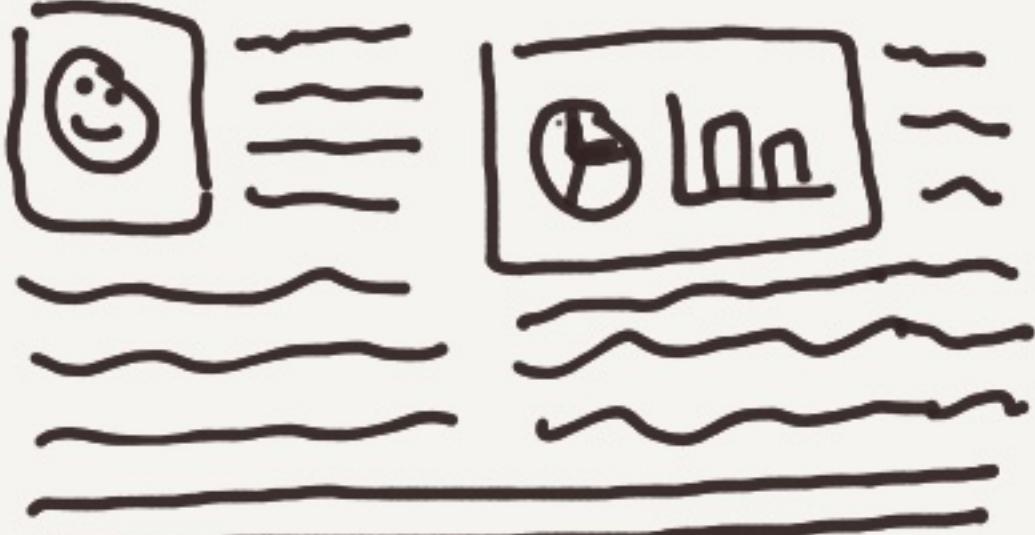


How would you want this
information to be stored?
(STRUCTURED)



THE GUARDIAN

EXTRA! EXTRA! EXTRA



Python

What Python tools would you need to use?

- pandas - DataFrame

- BeautifulSoup
(web scraping)

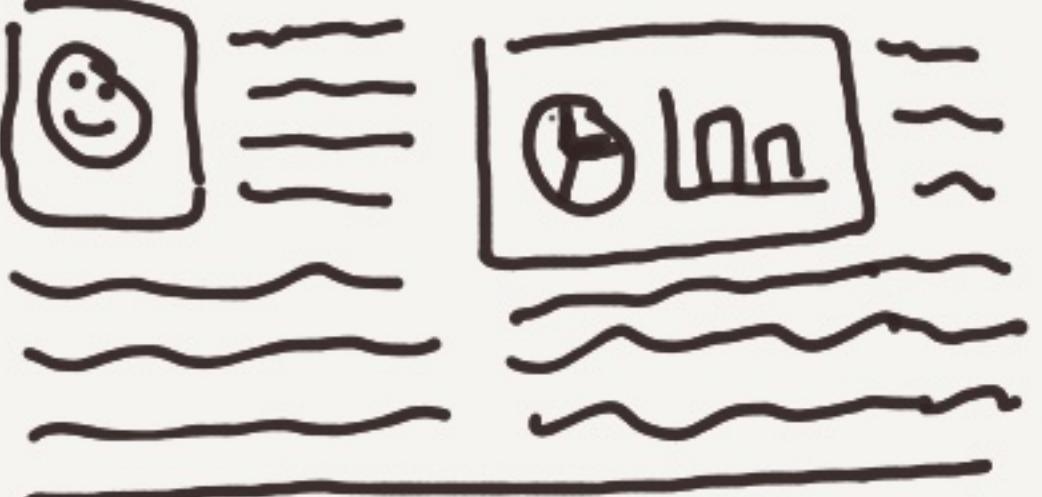
- seaborn (graphs)
→ matplotlib
→ pandas

→ Do the analysis?
- correlation → pandas
- numpy(np)

→ stats (statistics)
package

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 Pulitzer's won) on newspaper readership?

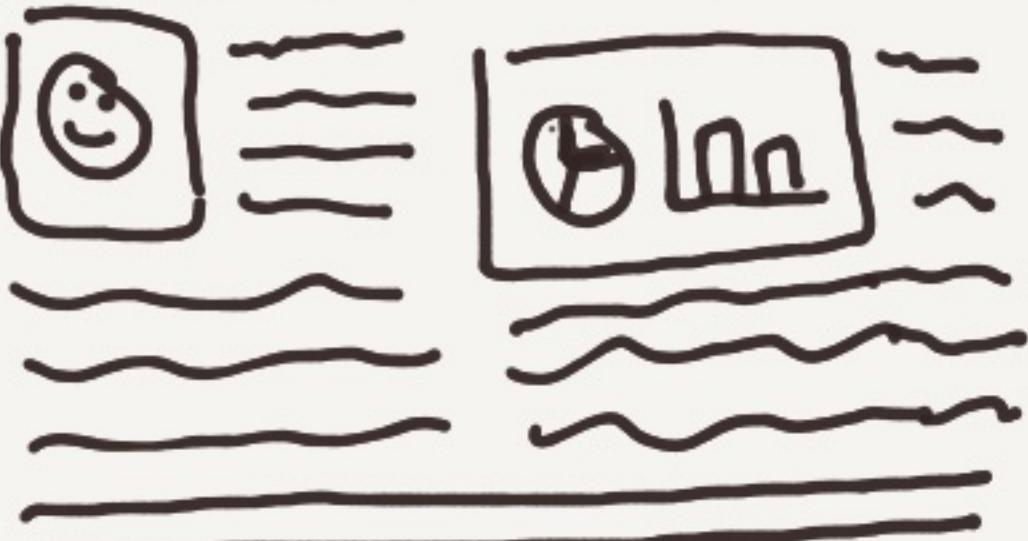
H: - no effect
o readers don't care/know
o other things more important

- pos. effect

Data:

THE GUARDIAN

EXTRA! EXTRA! EXTRA



EDA

Boxplot?

① $\frac{SD}{\text{from the mean}}$

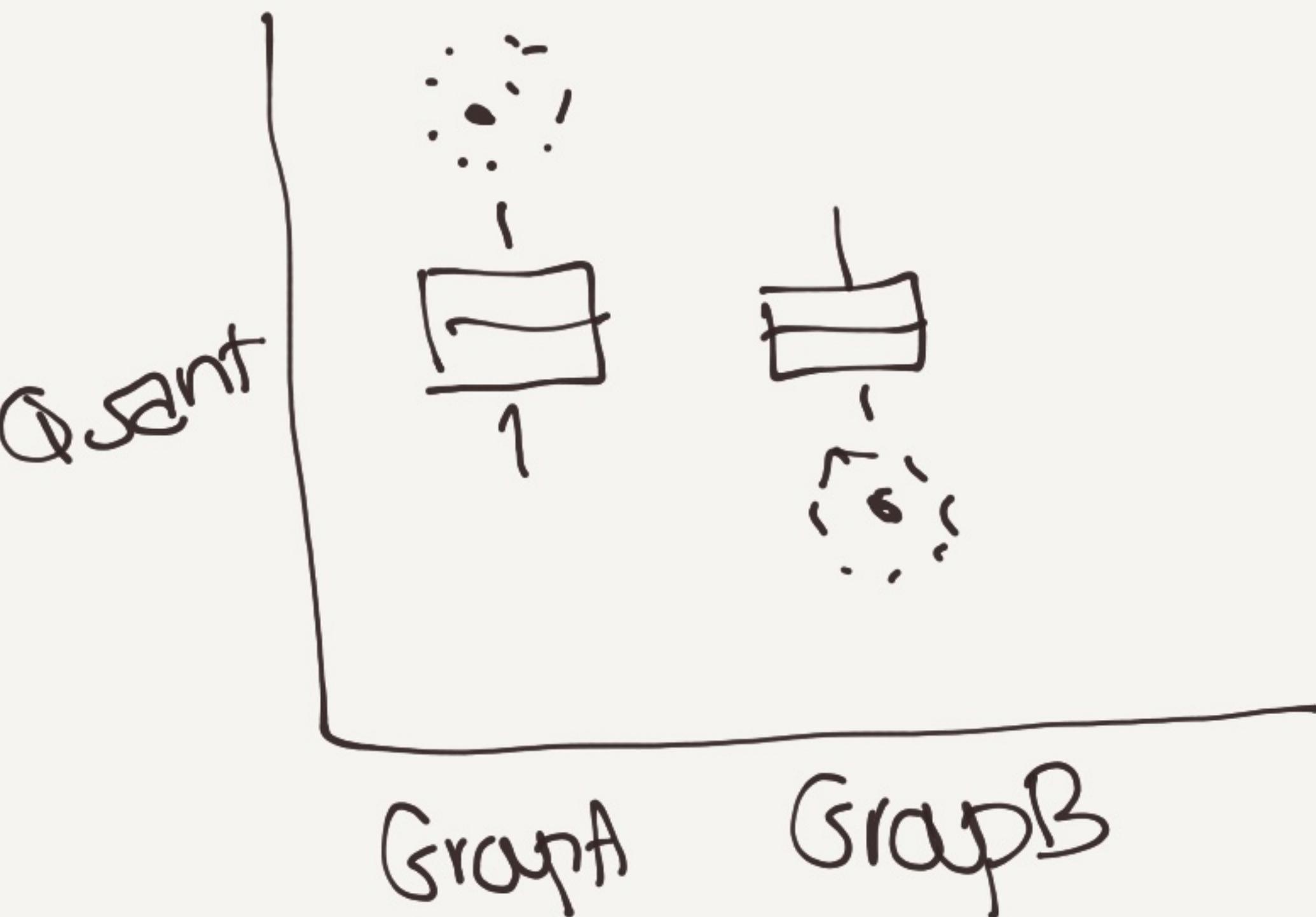
? magazines

✓ check for
outliers?

→ get rid of
them?

②

✓ - check for
missing/null
data
(No missing
data)



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?

③ ↑ circulation ↑ pulitzer prizes

↑ 30 newspr

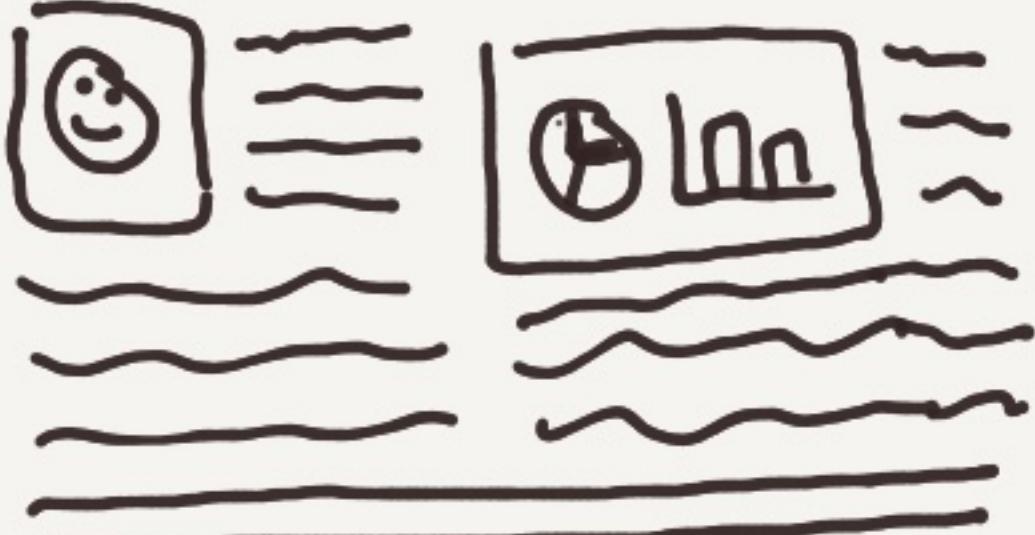
- correlation
- regression
- relationship
b/w 2 var.

③

outlier

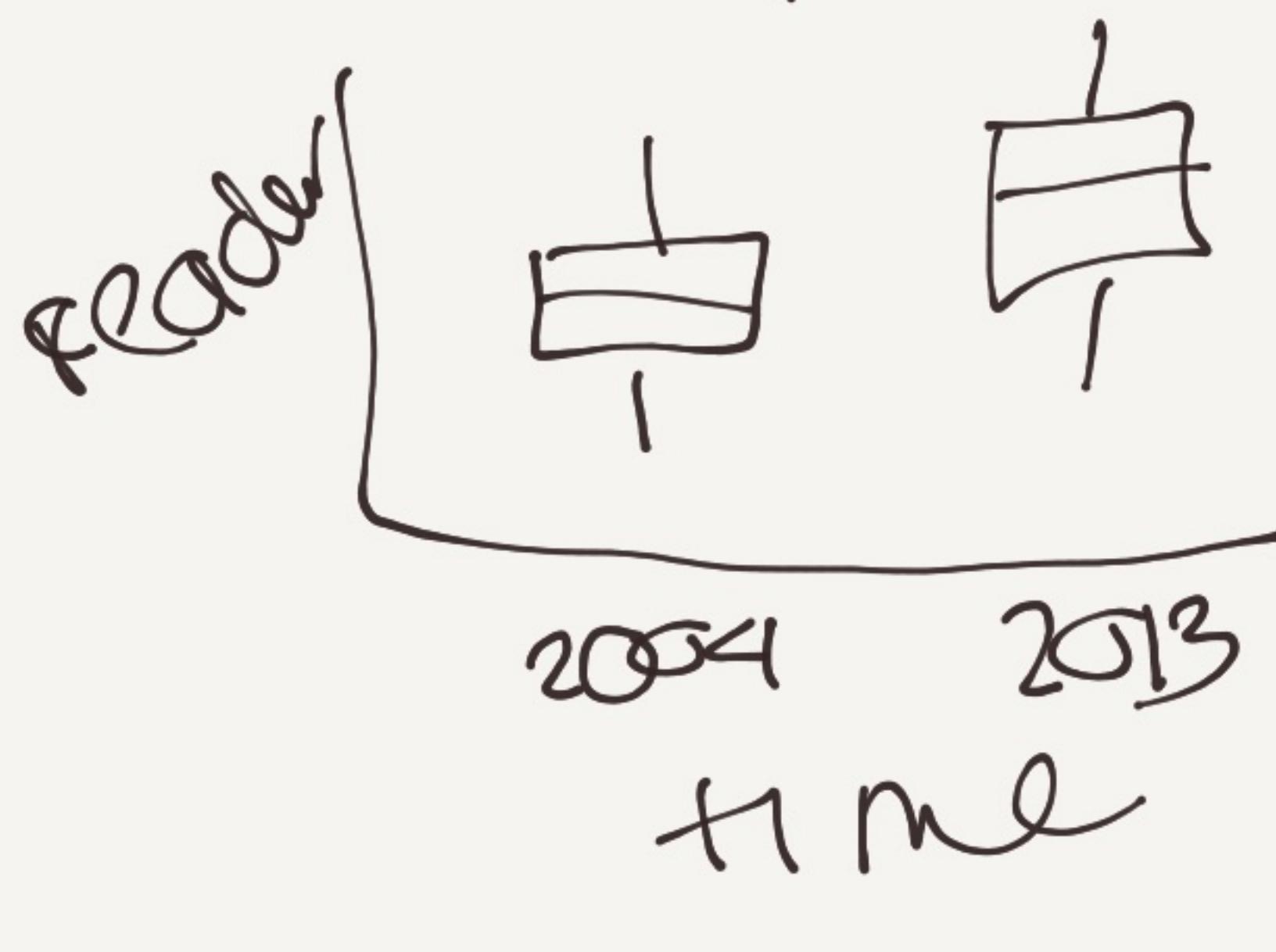
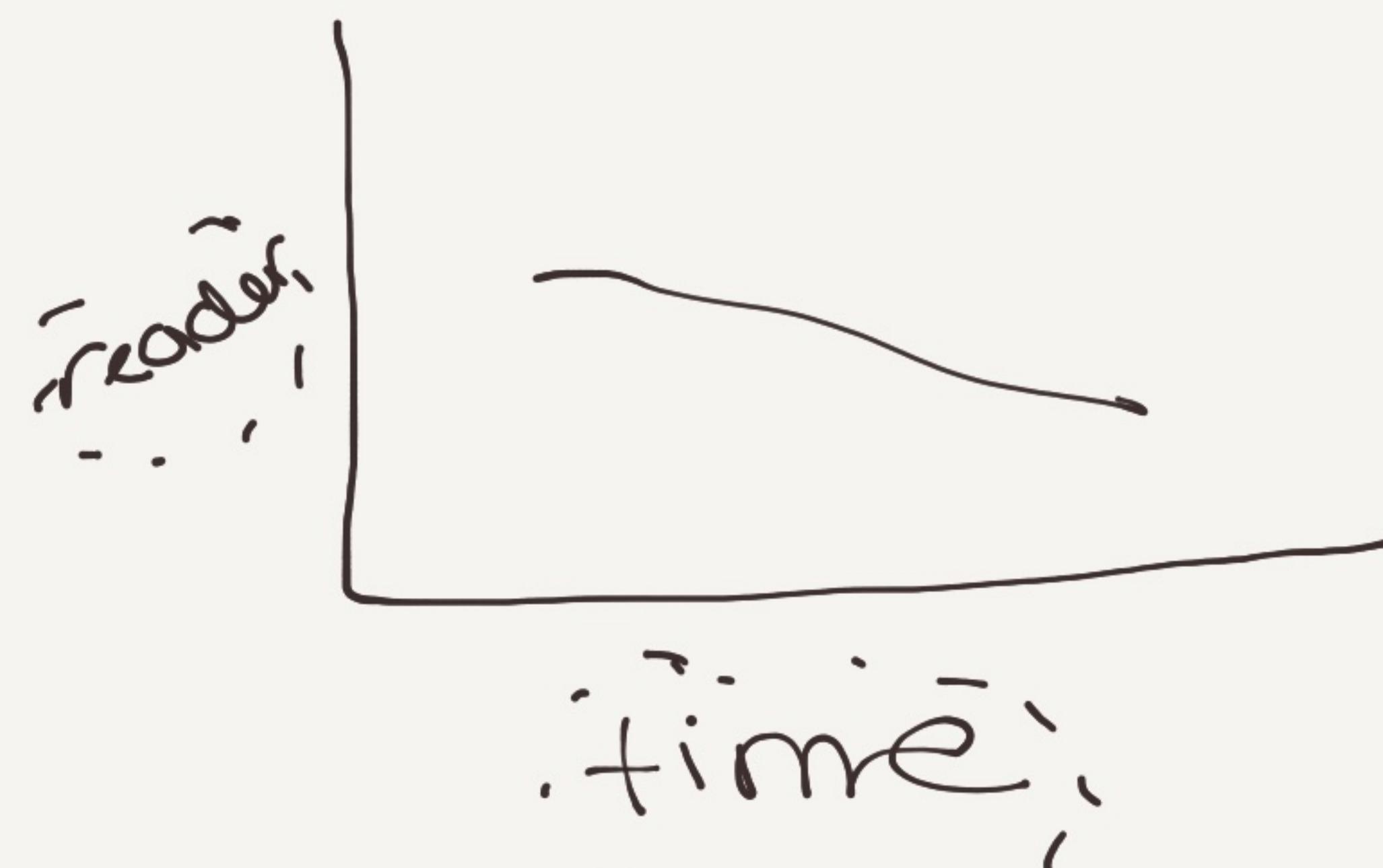
THE GUARDIAN

EXTRA! EXTRA! EXTRA



EDA

(change)
increase in
readership
over time



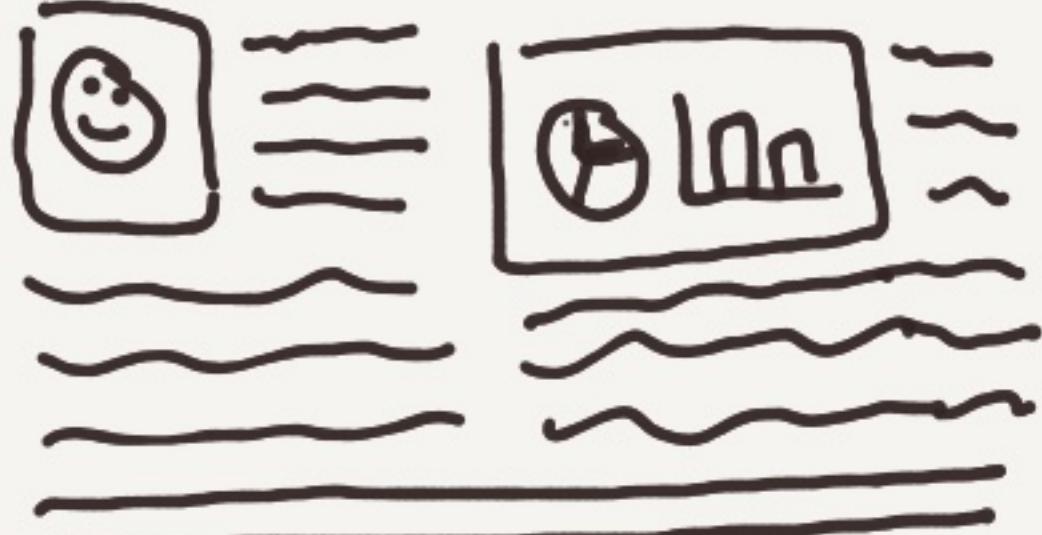
what exploratory visualizations would you generate?

- single variable
- histograms for each column
 - new papers?
- density plot
- boxplots {two variables}
- scatterplots {variables}

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



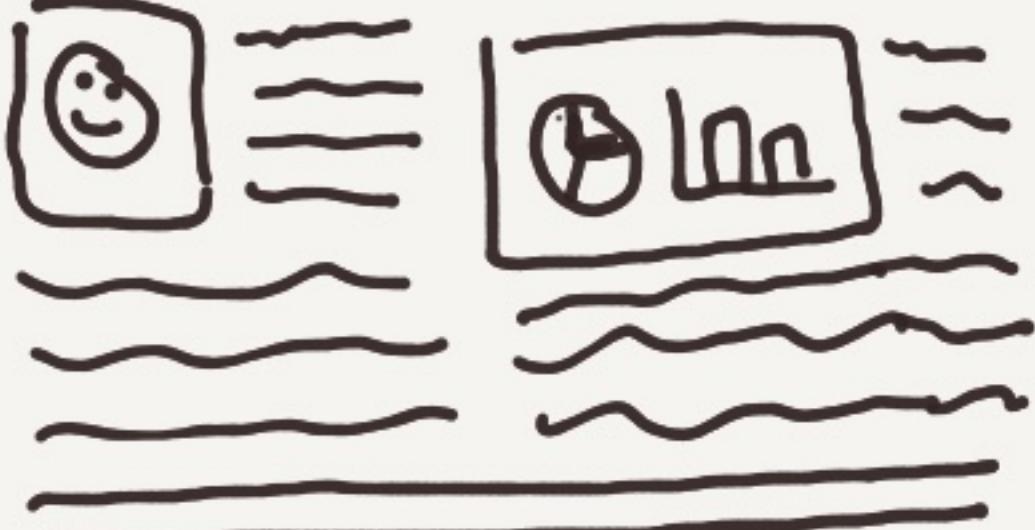
EDA

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

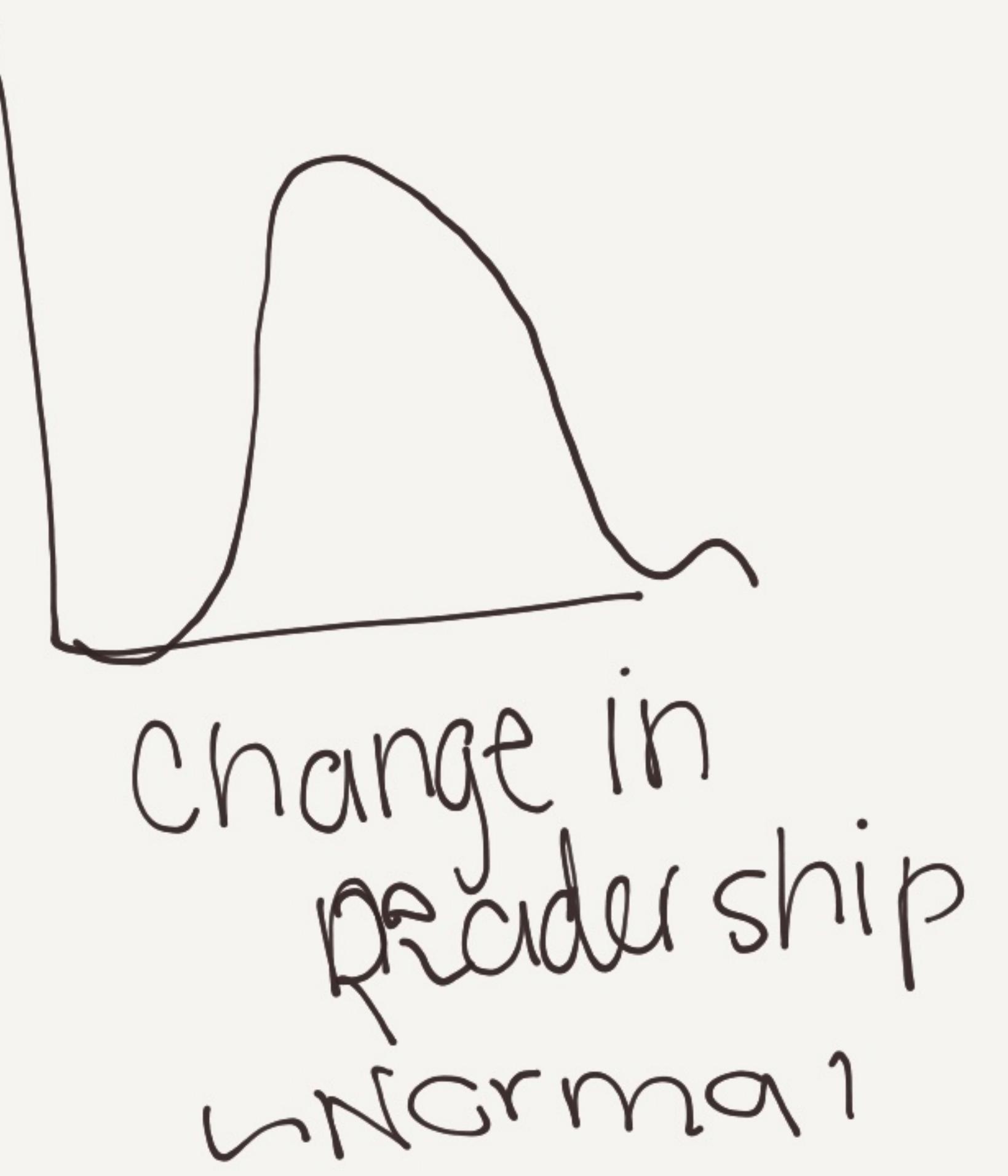
EXTRA! EXTRA! EXTRA



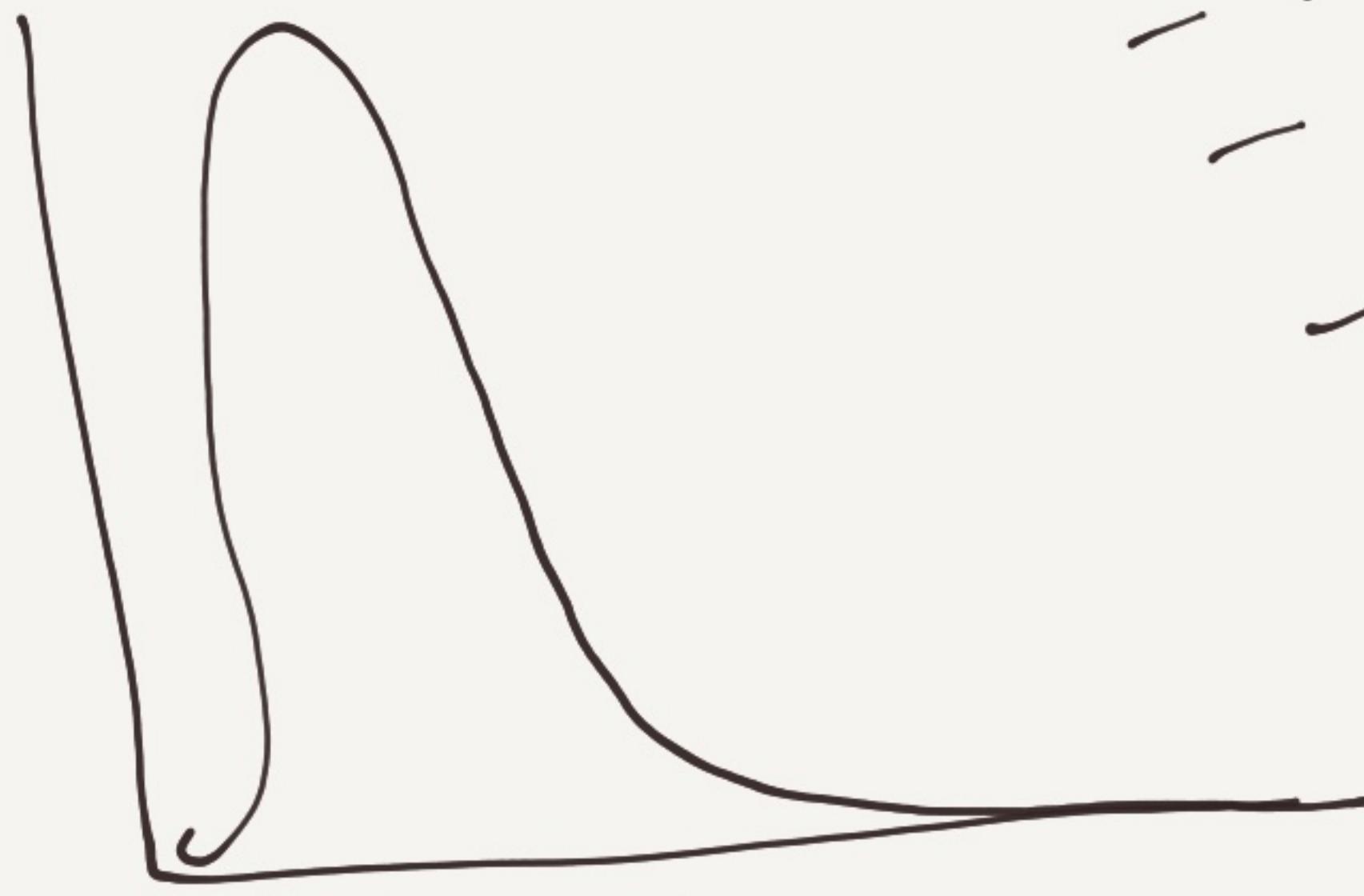
Data



stratify



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



Pulitzers
skewed

↳ log-transform

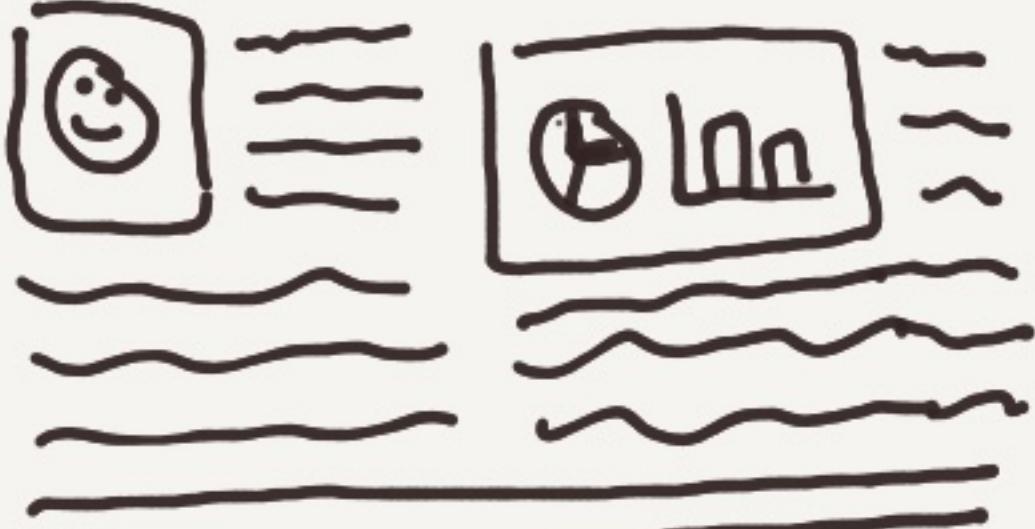


assumptions

- X. - multicolinearity
- homoskedasticity
- "Normal"
-

THE GUARDIAN

EXTRA! EXTRA! EXTRA

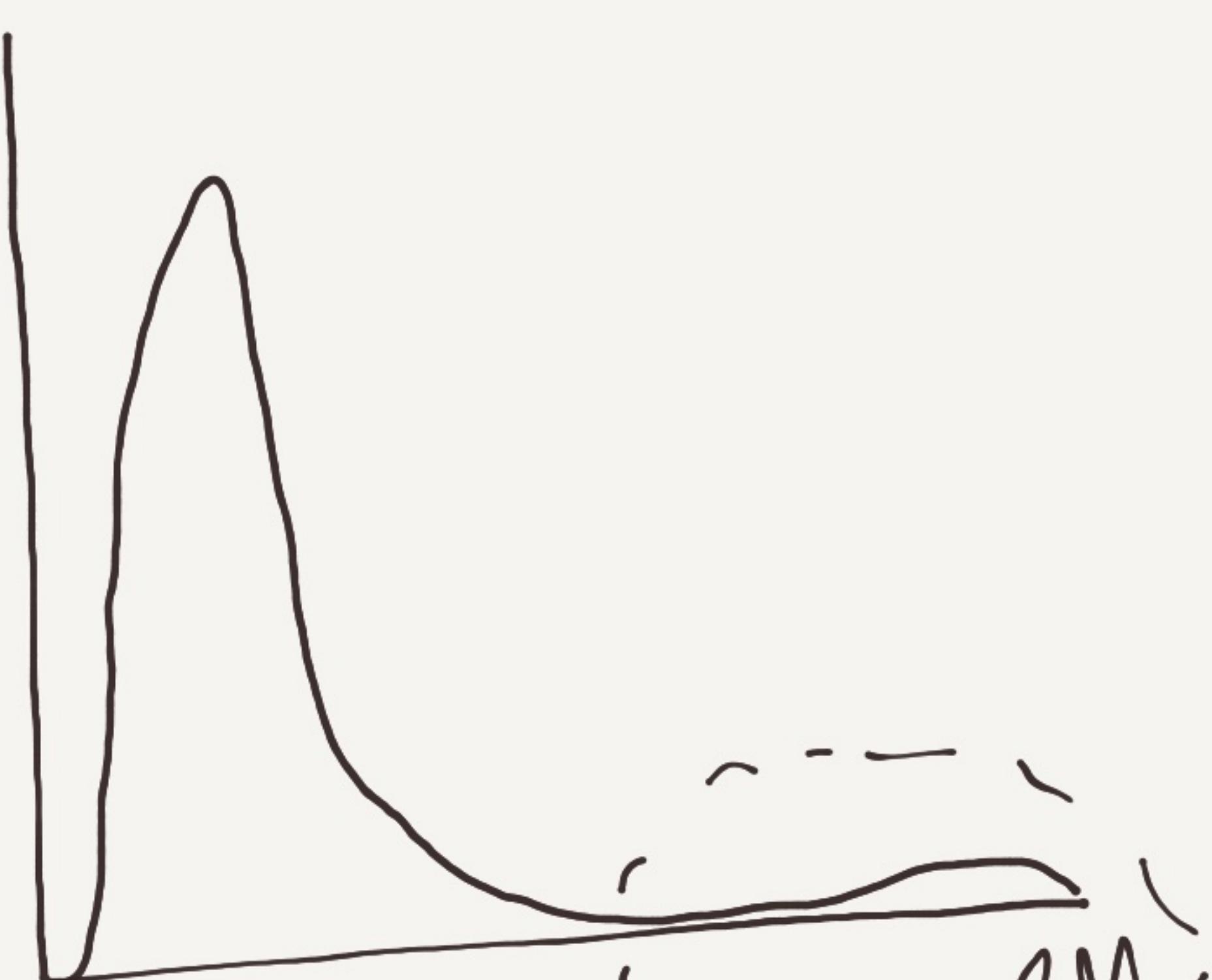


EDA

EDA: circulation (Readership)

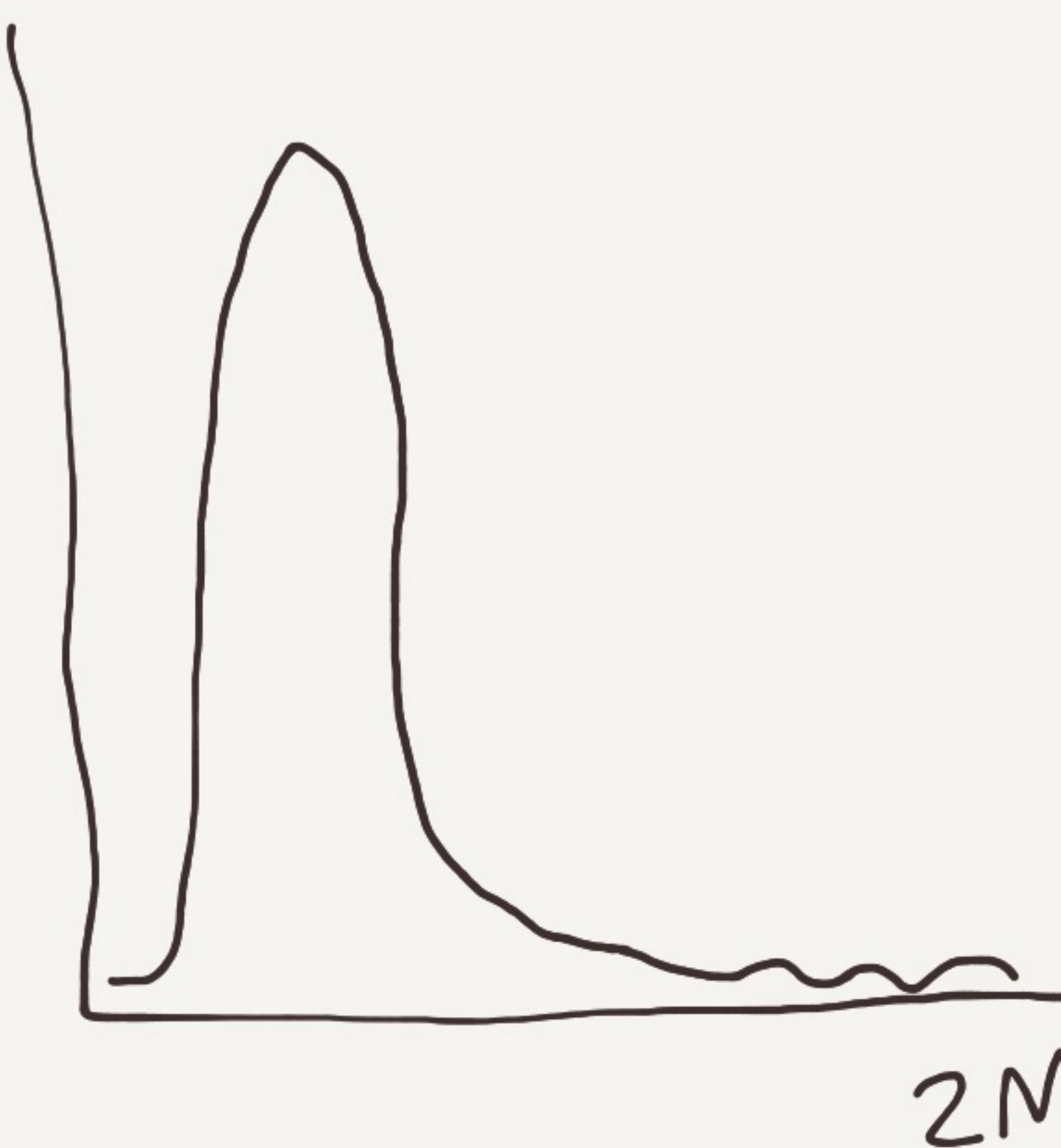
trend?
interpretation?
looks weird?

Daily Circulation
2004



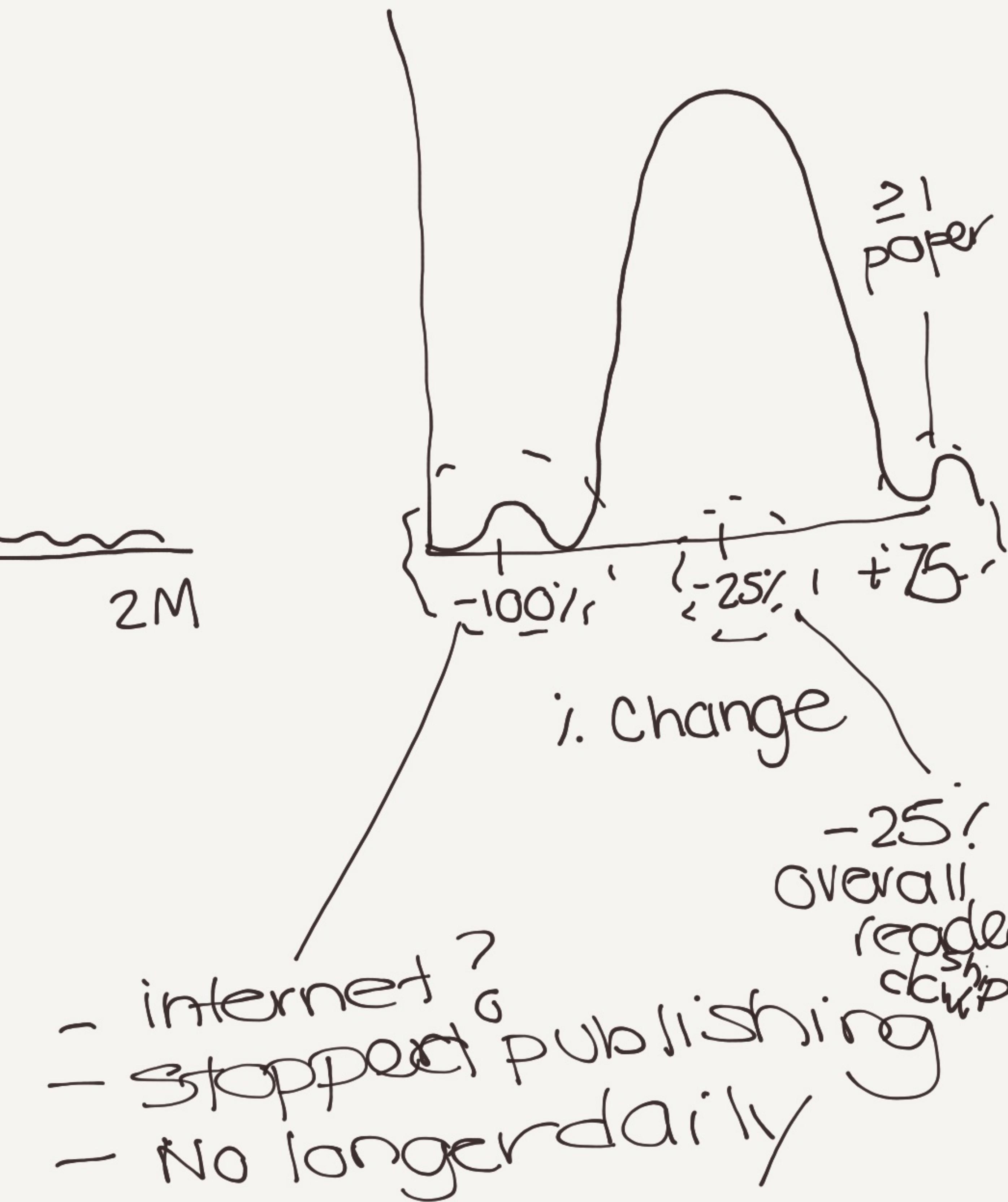
most
data
for newspapers
close to
zero

Daily Circulation
2013



outliers?

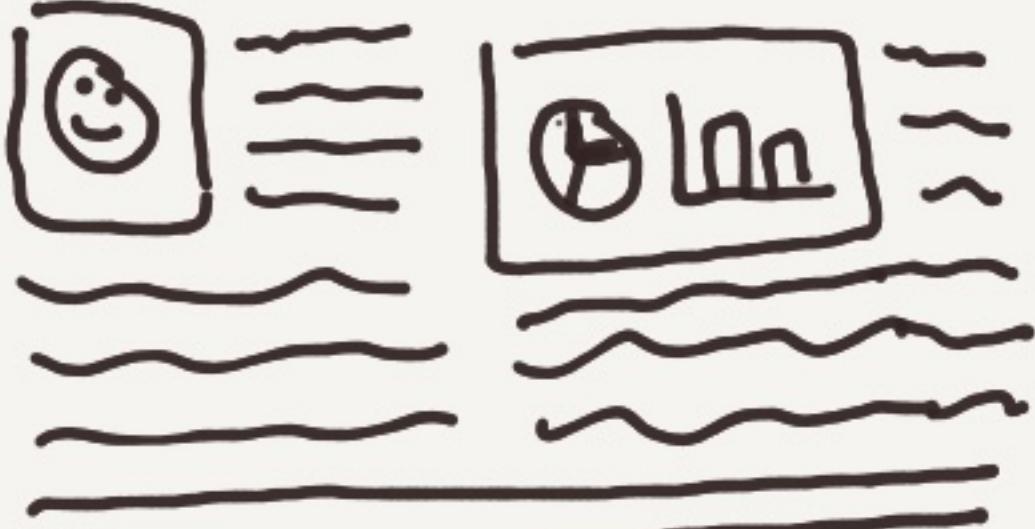
change in Daily
circulation
2004-2013



- internet?
- stopped publishing
- No longer daily

THE GUARDIAN

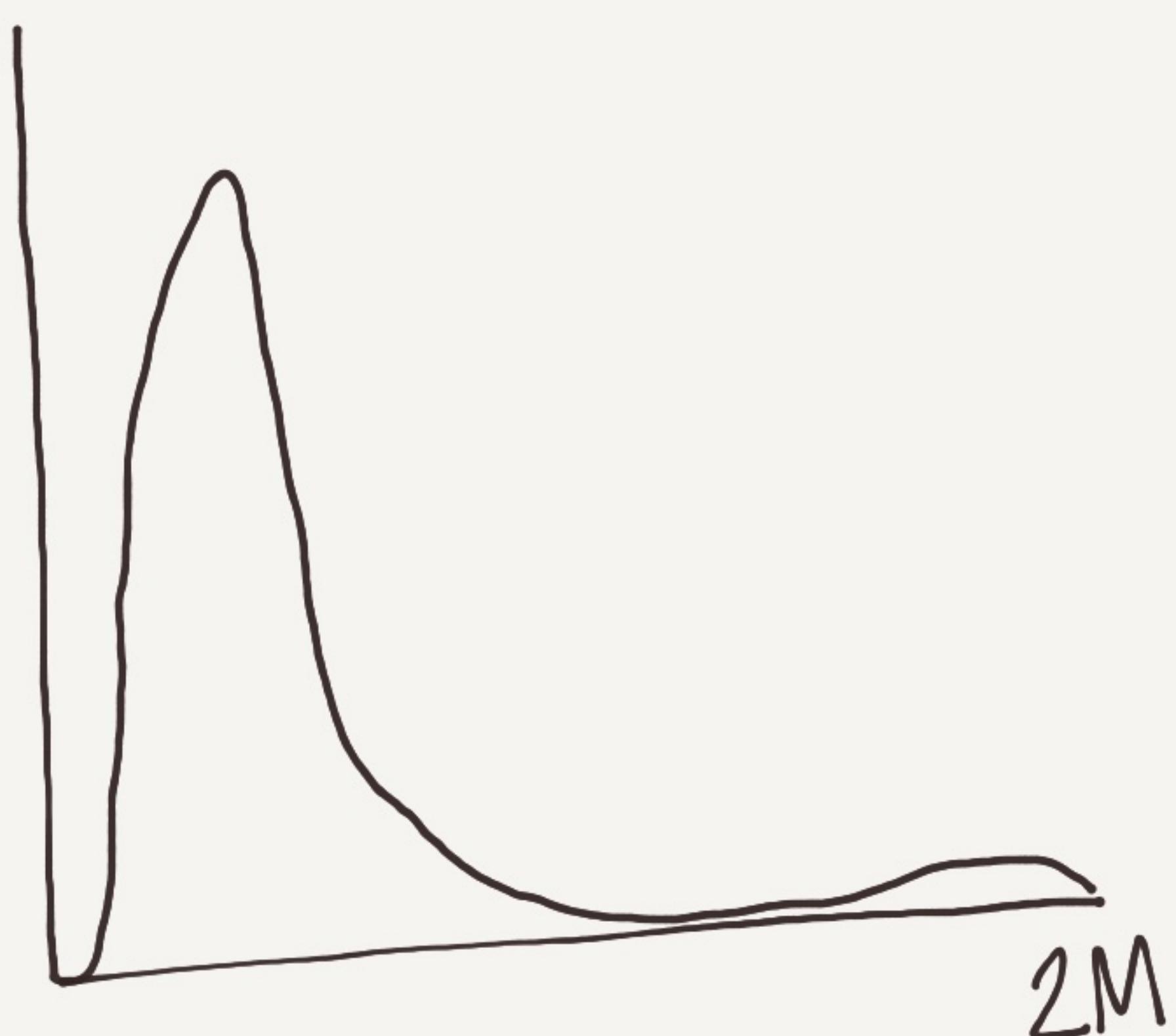
EXTRA! EXTRA! EXTRA



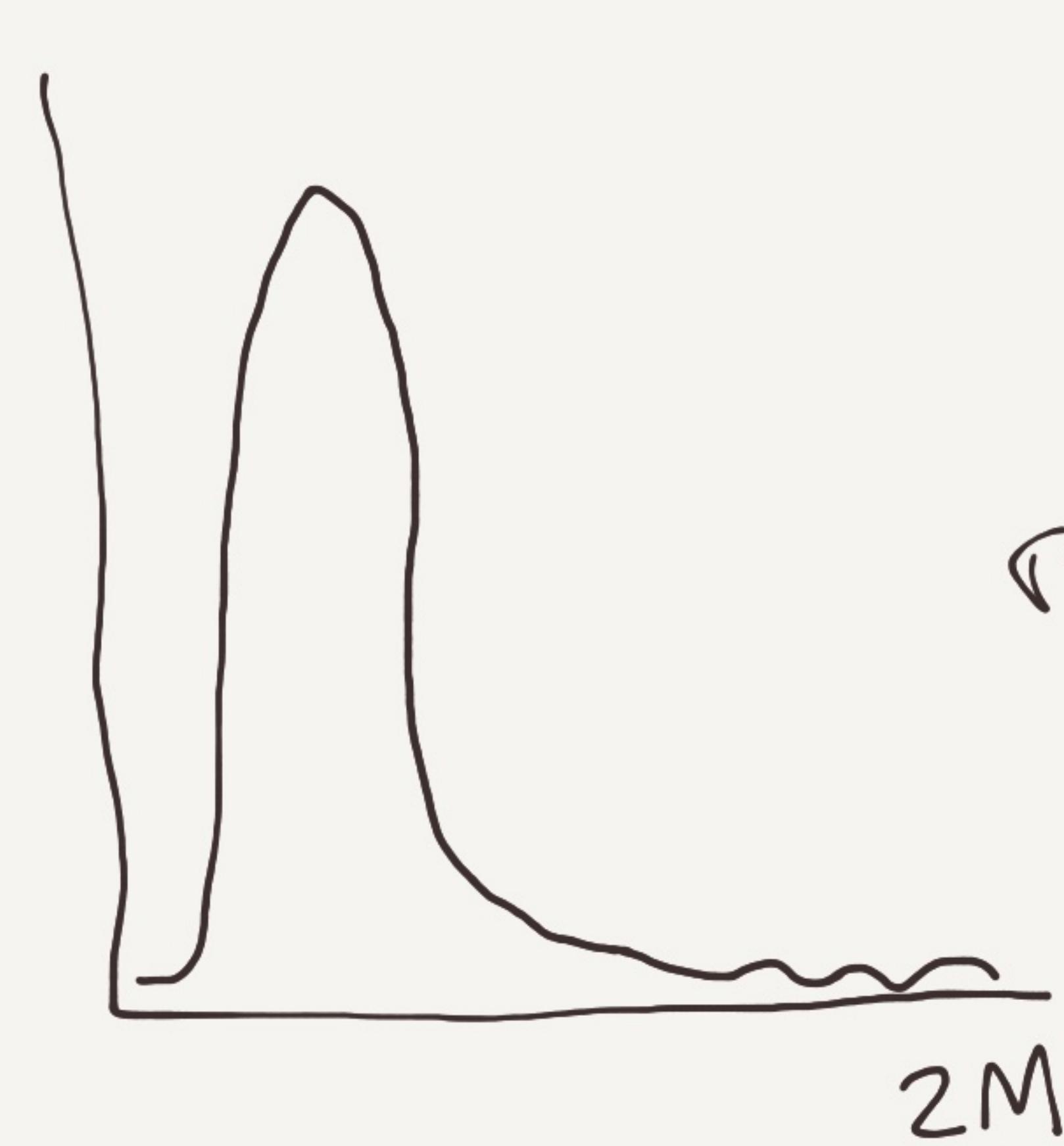
EDA

EDA: circulation (Readership)
prestige → readership

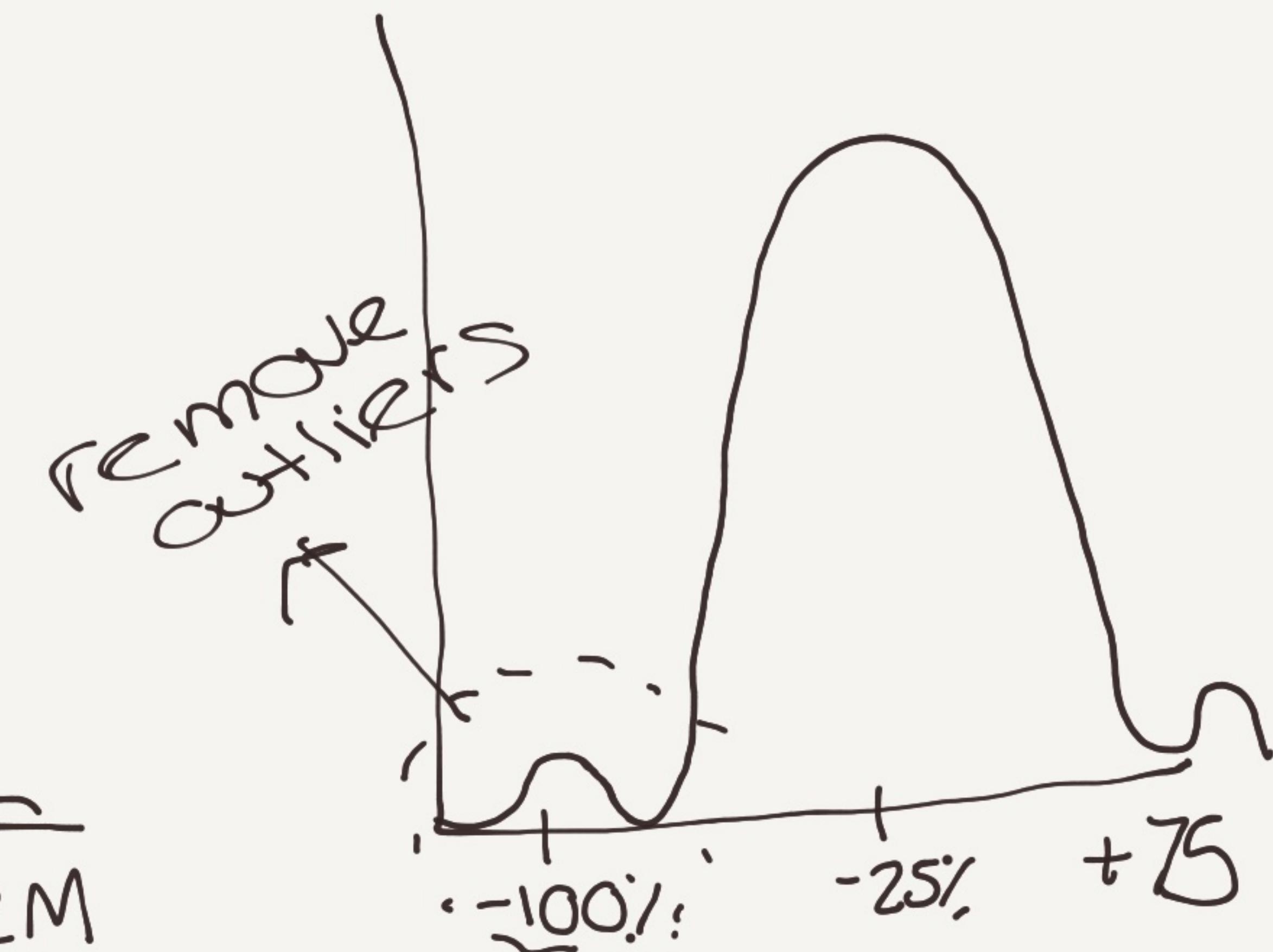
Daily Circulation
2004



Daily Circulation
2013



Change in Daily
Circulation
2004-2013



Clicker
Question :

Which variable would you use
for this analysis?



A Daily circulation, 2004 (27)

B Daily circulation, 2013 (13)

C Change in Daily circulation (66)

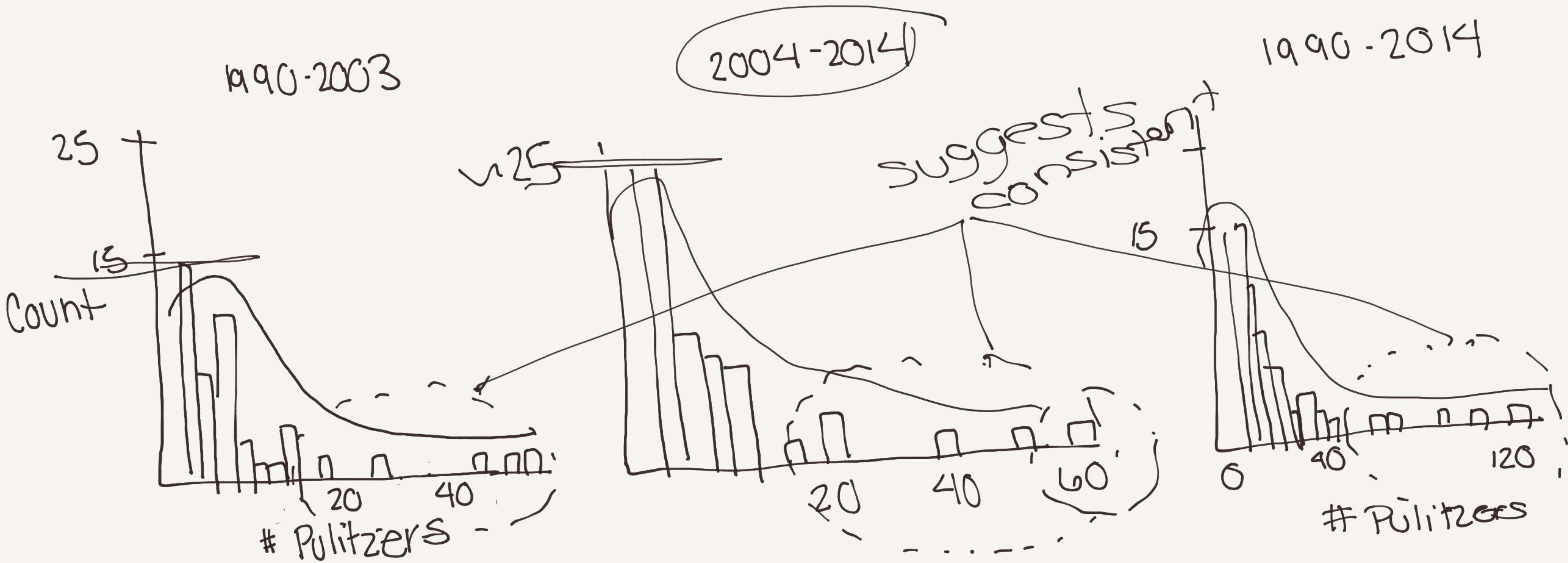
i. change

↓
Normal
distribution



EDA : Pulitzer

↳ # Winners + Finalists



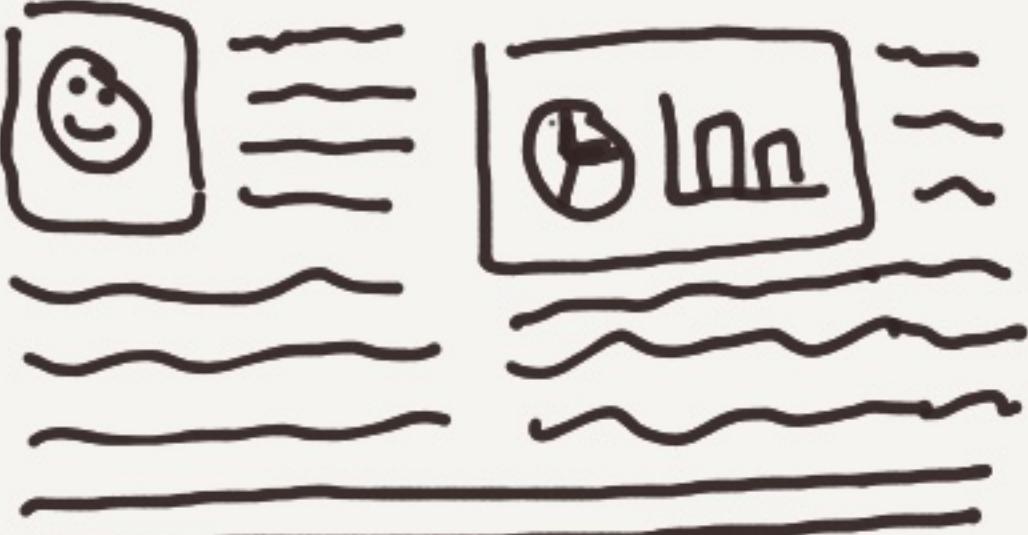
values
close
to lower
limit
(skew)

↳ more outlets?
news

new award
categories?
→ Pulitzer

THE GUARDIAN

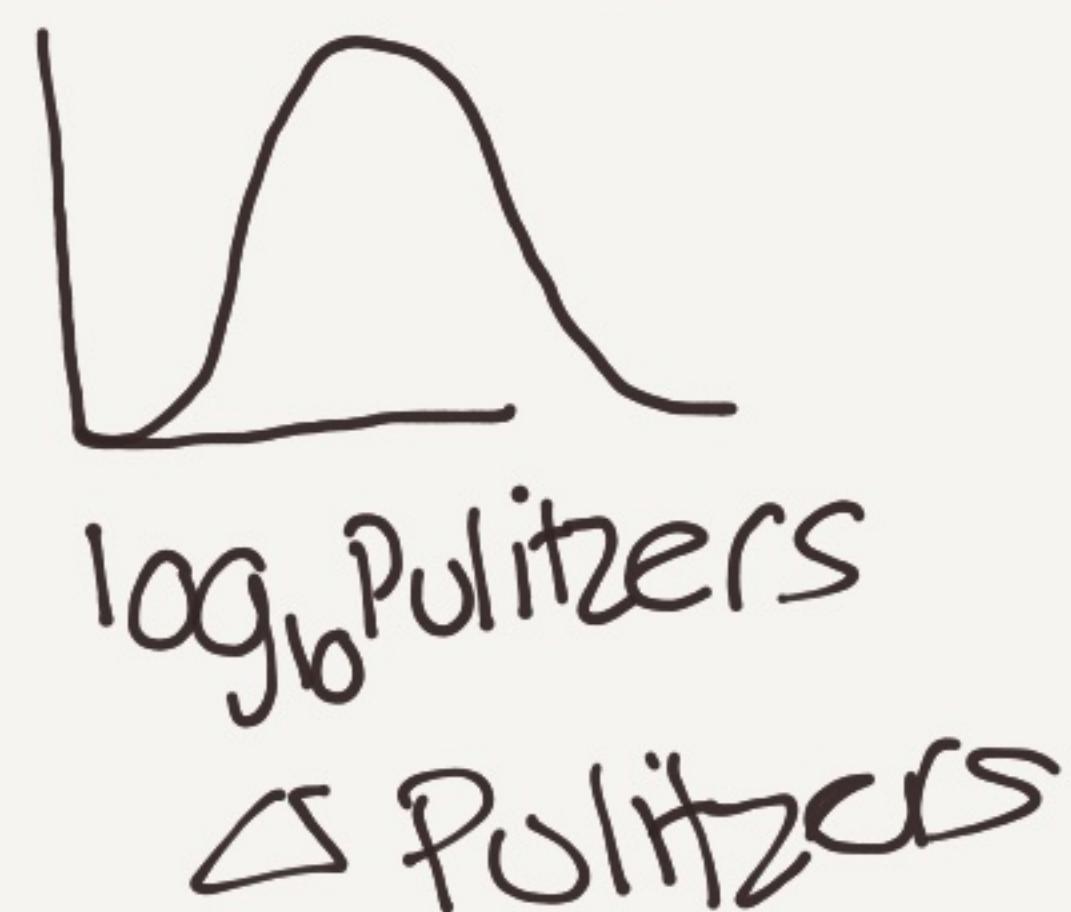
EXTRA! EXTRA! EXTRA



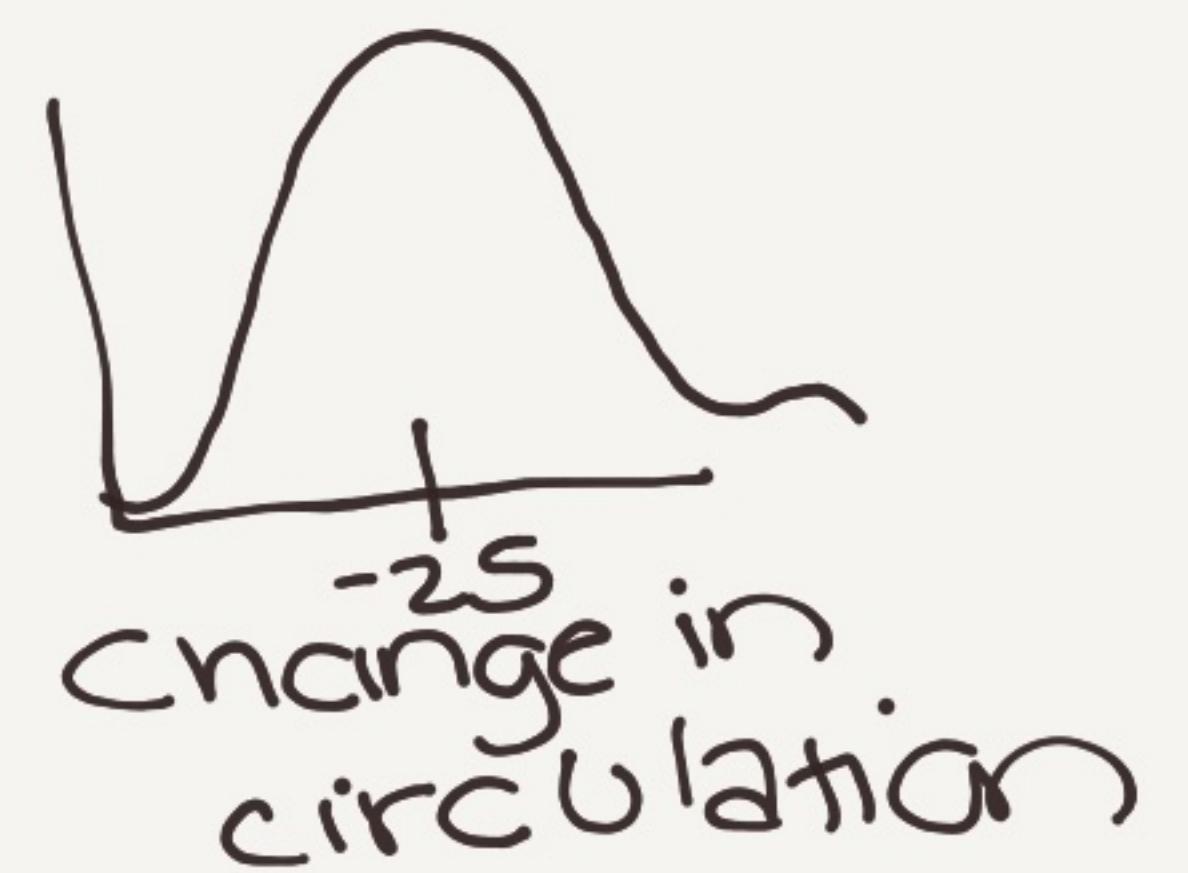
Analysis

How would you answer
the question?

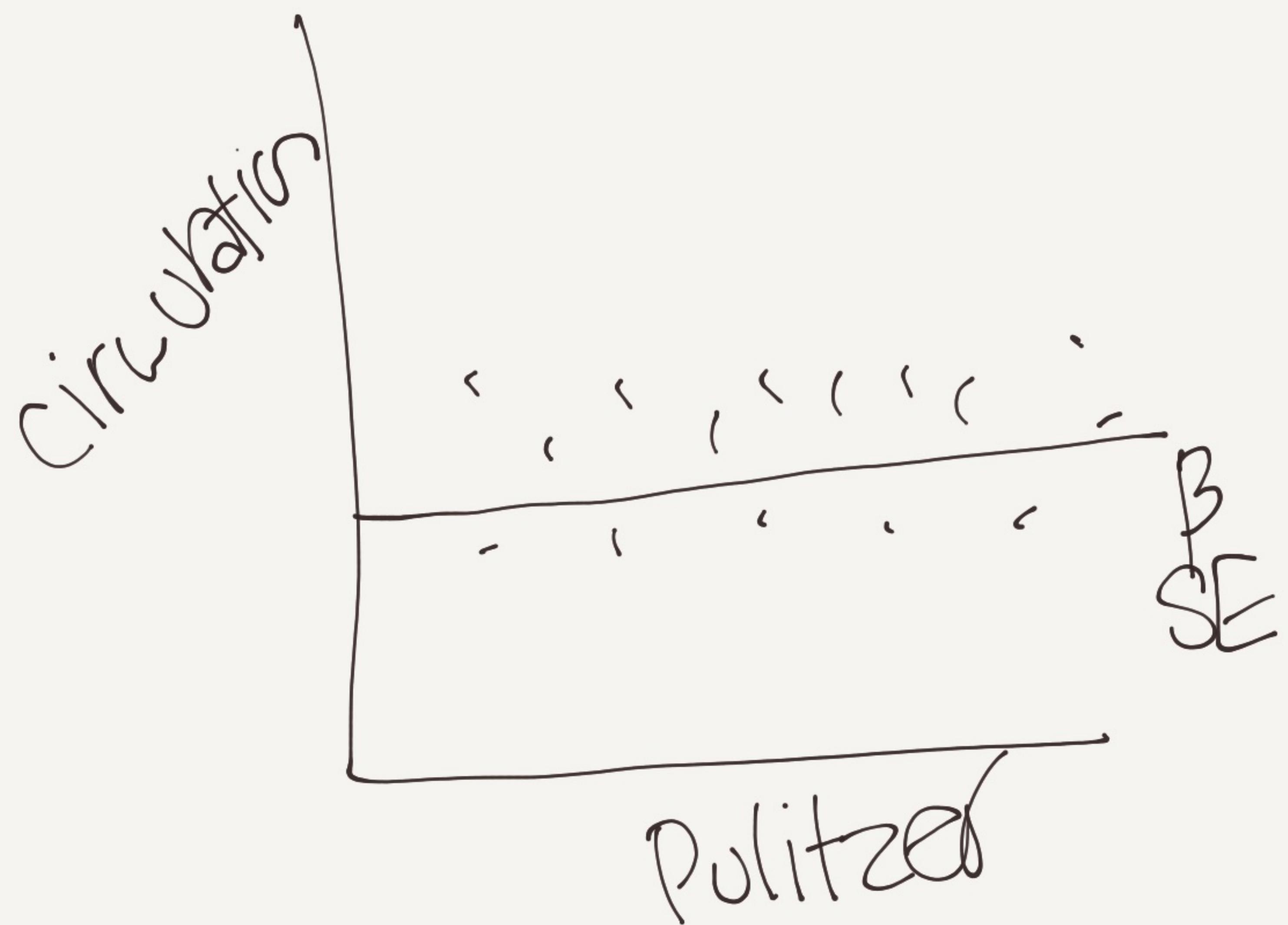
prestige \rightarrow readership



\log_{10} Pulitzer
 Δ Pulitzer



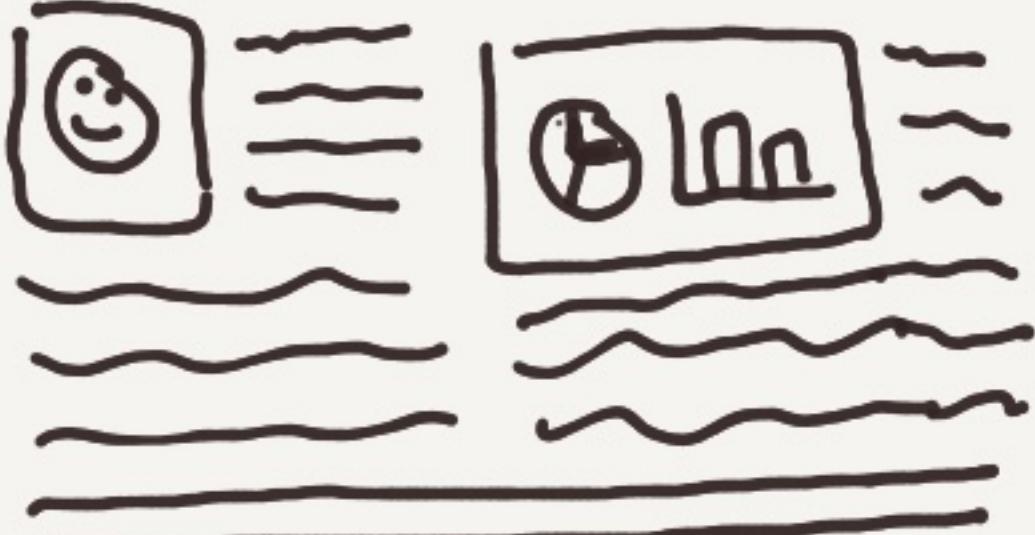
change in
circulation



outcome: readership
predictor: pulitzer prestige

THE GUARDIAN

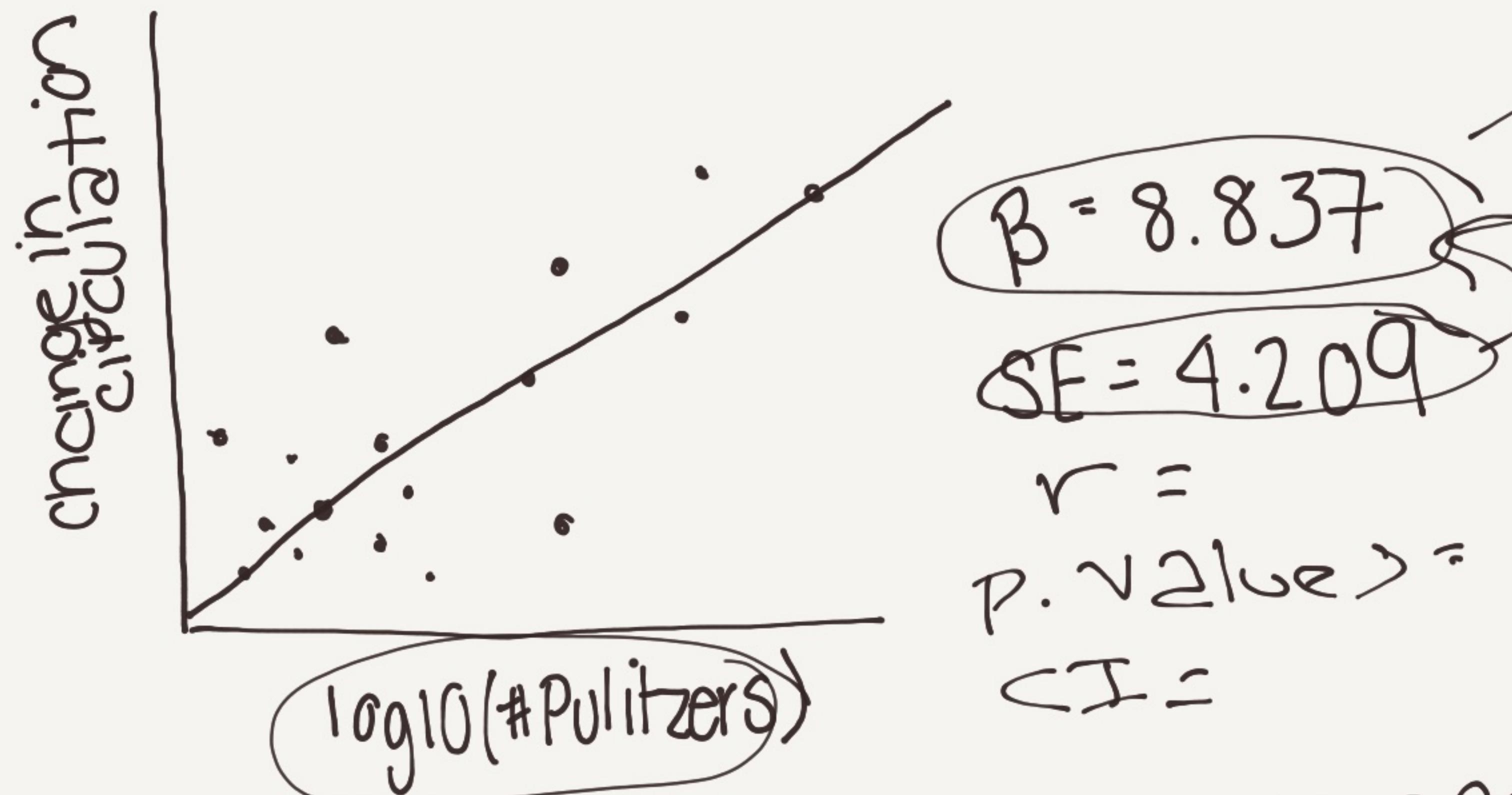
EXTRA! EXTRA! EXTRA



INTERPRETATION

CLICKER
★

GIVEN THE FOLLOWING RESULTS, WHAT
WOULD YOU CONCLUDE?



$$\beta/100$$

forever
~ 0.08!
change
incirc.

$$r =$$

P. value >=

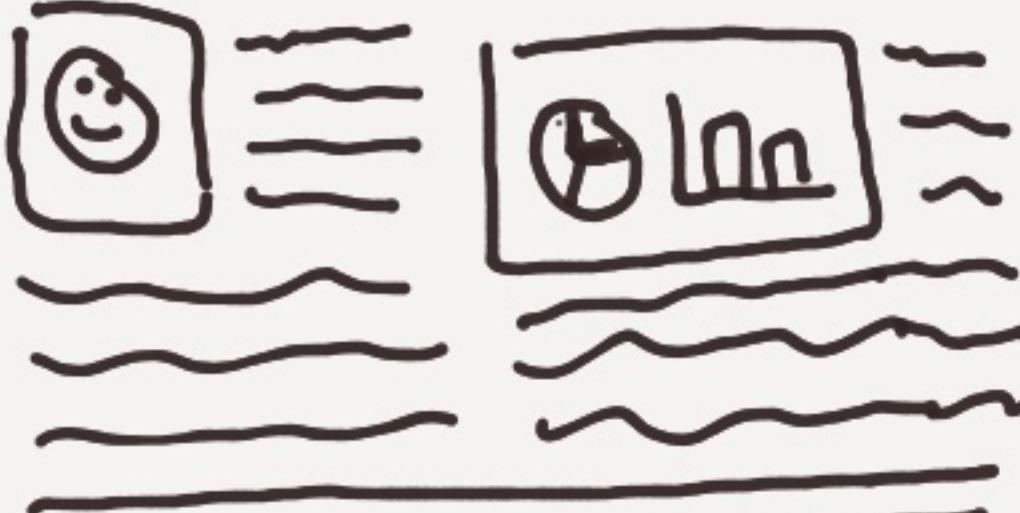
CI =

An increase in prestige tends to increase circulation

- 82 A. Prestige increases circulation
- 5 B. Prestige decreases circulation
- 9 C. Prestige does not affect circulation
- 7 D. Something else
- ↳ correlation

THE GUARDIAN

EXTRA! EXTRA! EXTRA



Limitations

What limitations are there in this analysis?

pre . → change in
#Pulitzers circulation

- account for other variables (confounders)
 - o ↓ in newspaper (internet measurements offed)
- prestige?
 - other measure →
 - citations
- measures of readership
- line graph
→ year... missing