

FAKECOMPSTAT

An R package for financial data replication

Dylan Phillips and Candace Jens

PURPOSE

- To produce an R package designed to create fully synthetic Compustat and CRSP-style datasets for use in financial paper replication and methodological testing

PACKAGE GOALS

- Generated data must follow accounting and financial principles
 - Ex.Assets = Liabilities + Equity

PACKAGE GOALS

- Generated data must follow accounting and financial principles
 - Ex. $\text{Assets} = \text{Liabilities} + \text{Equity}$
- Generated distributions must be realistic on both the firm and market-levels
 - Must simulate “growth,” cannot use pure randomness

PACKAGE GOALS

- Generated data must follow accounting and financial principles
 - Ex. $\text{Assets} = \text{Liabilities} + \text{Equity}$
- Generated distributions must be realistic on both the firm and market-levels
 - Must simulate “growth,” cannot use pure randomness
- Generated data cannot be matched to any real company
 - Simple data permutation is not sufficient

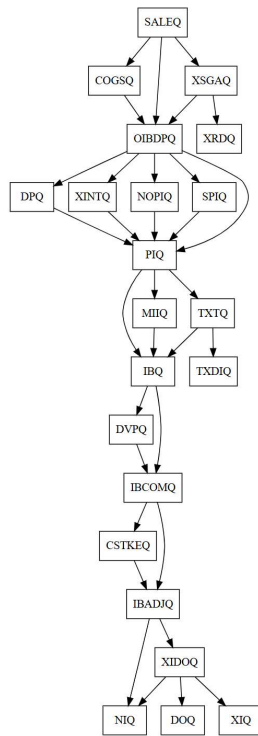
METHODOLOGY

- I. Key variables (assets, sales, stock price) will be simulated along a growth path
 - Starting point is randomly drawn from distribution of values at a sector and year level
 - Subsequent values are created based on distribution of growth

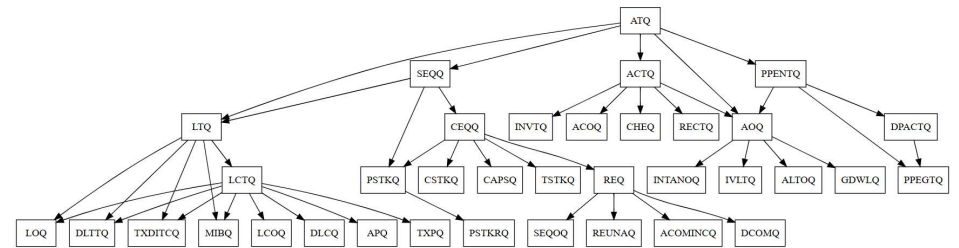
METHODOLOGY

1. Key variables (assets, sales, stock price) will be simulated along a growth path
 - Starting point is randomly drawn from distribution of values at a sector and year level
 - Subsequent values are created based on distribution of growth
2. Remaining variables are simulated using ratios to appropriate key variable

INCOME STATEMENT



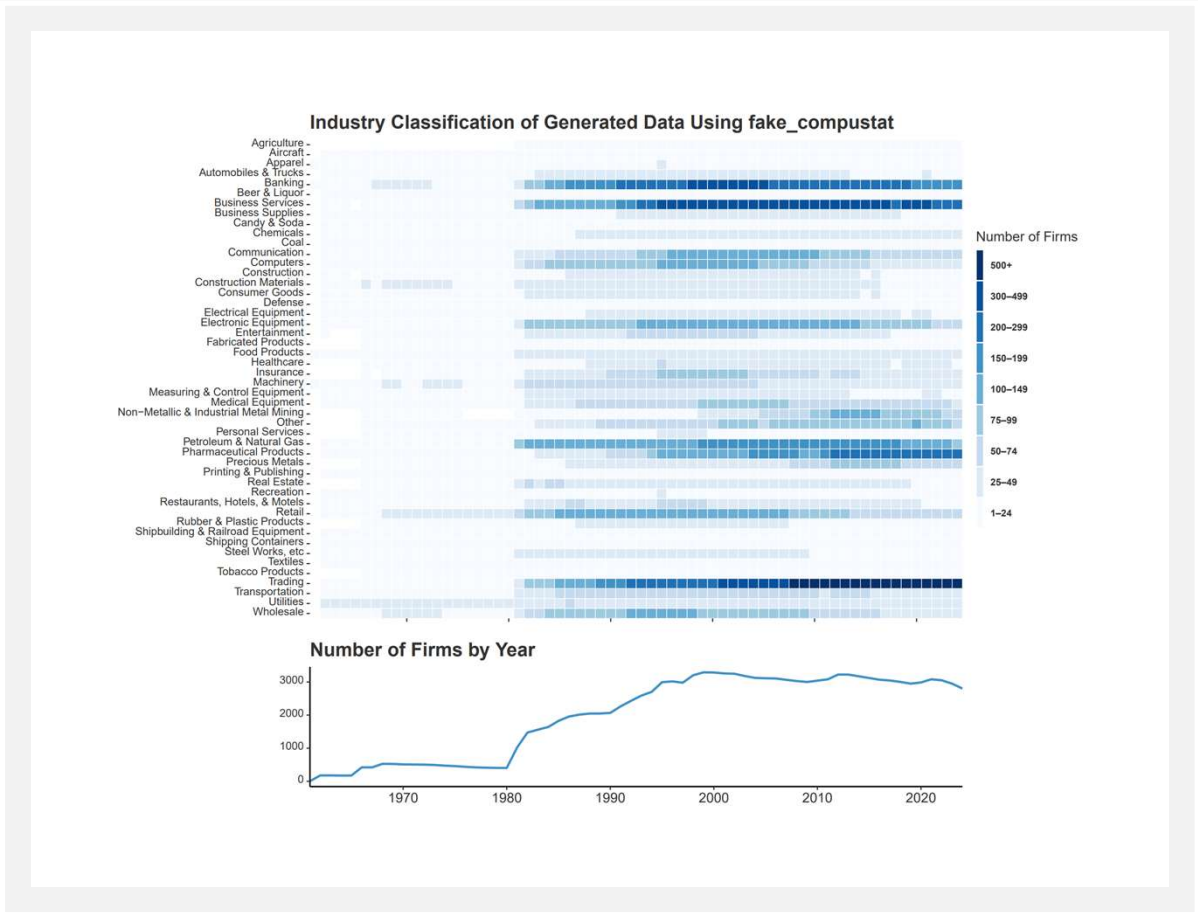
BALANCE SHEET



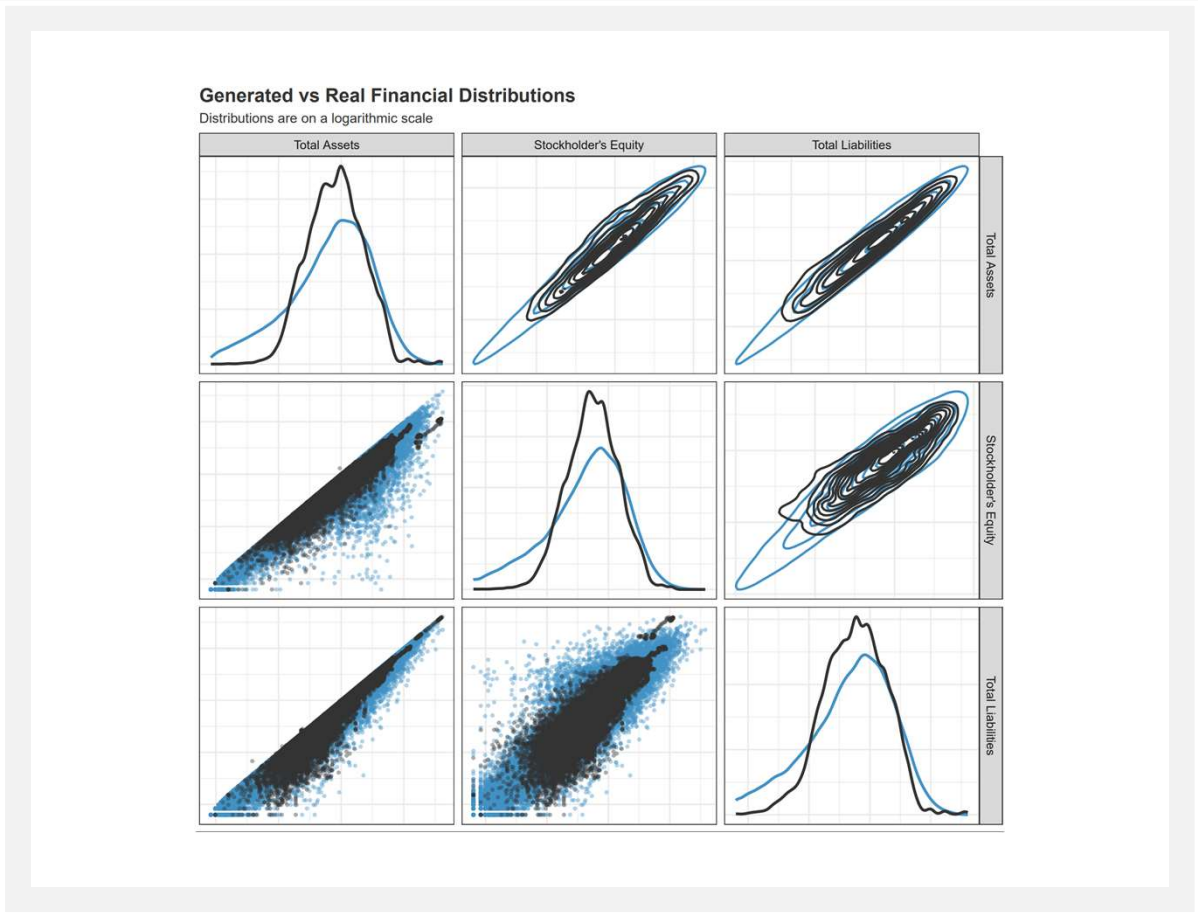
METHODOLOGY

1. Key variables (assets, sales, stock price) will be simulated along a growth path
 - Starting point is randomly drawn from distribution of values at a sector and year level
 - Subsequent values are created based on distribution of growth
2. Remaining variables are simulated using ratios to appropriate key variable
3. Data is mutated to fulfill principle and distribution goals

RESULTS

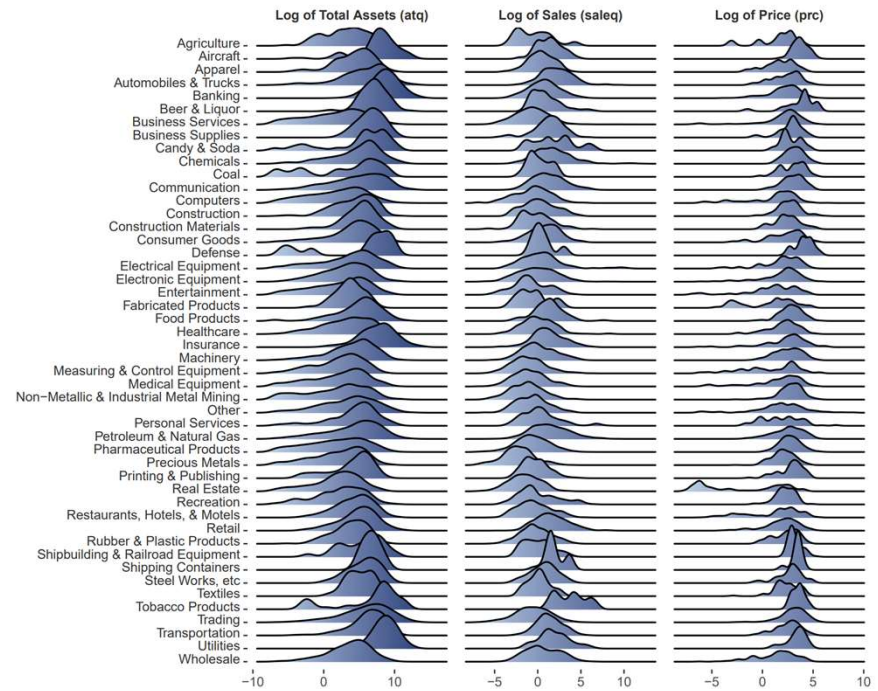


RESULTS



RESULTS

Select Variable Distributions by Industry



THANK YOU