

CSS Pitch

“Correlation is not causation” –
Exploring & Visualizing Correlations

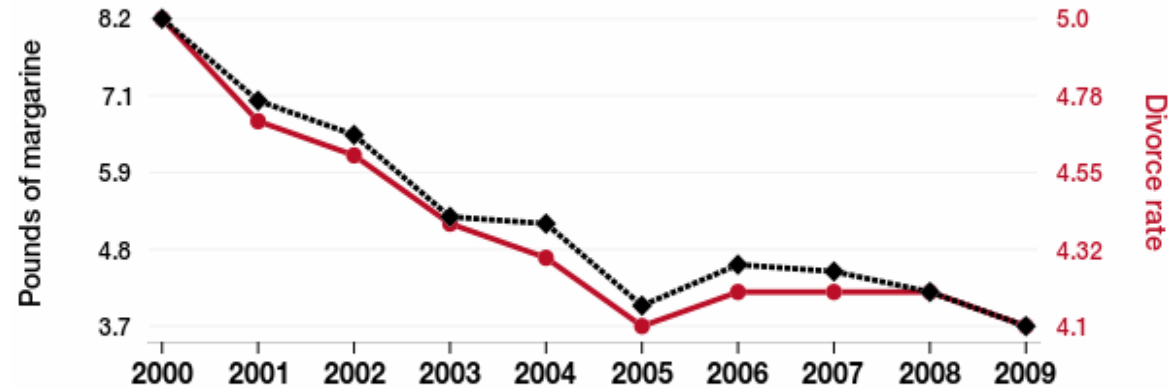
Nils Düpont
CRC 1342, Project INF

✉ duepont@uni-bremen.de

Per capita consumption of margarine

correlates with

The divorce rate in Maine

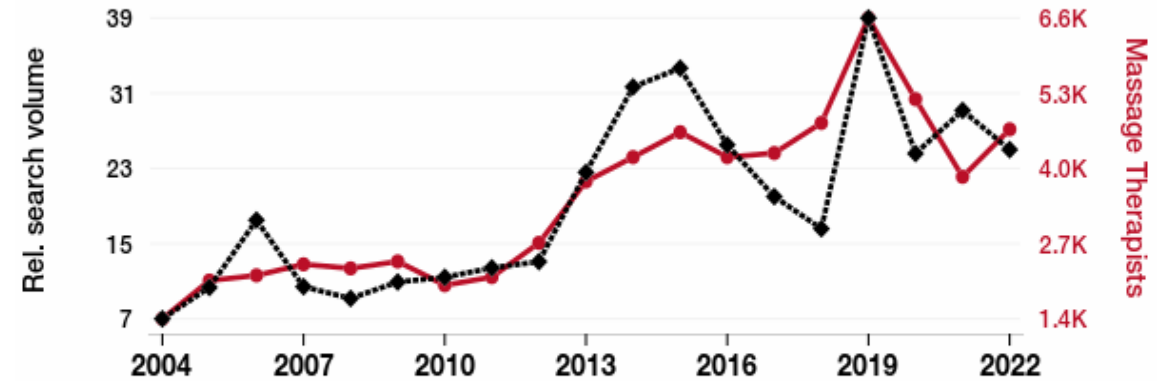


◆ Per capita consumption of margarine in the United States · Source: US Department of Agriculture
 ● The divorce rate in Maine · Source: CDC National Vital Statistics
 2000-2009, $r=0.993$, $r^2=0.985$, $p<0.01$ · tylervigen.com/spurious/correlation/5920

Google searches for 'funny cat videos'

correlates with

The number of massage therapists in New York



◆ Relative volume of Google searches for 'funny cat videos' (United States, without quotes) · Source: Google Trends
 ● BLS estimate of massage therapists in New York · Source: Bureau of Labor Statistics
 2004-2022, $r=0.867$, $r^2=0.751$, $p<0.01$ · tylervigen.com/spurious/correlation/9000

For more spurious correlations check: <https://tylervigen.com/spurious-correlations>

- Correlation is not causation, but...
- ...an essential part in the social sciences
 - in inductive data exploration
 - in proposing new hypotheses
 - in theory building
- Correlation analysis gains importance as the volume of data increases
- Correlation analysis not new to the social sciences
 - usually conducted in statistical software like R or Stata
 - becomes “cumbersome” as the number of variables increases

[illegible]

- “Correlation is not causation” project pitch
 - Up to ten variables
 - Interactive → dashboard?
 - Different ways of visualizing correlations → new ones?
 - Verbalize correlations → integrate generative AI models?
- The data set: “Varieties of Democracy” (<https://www.v-dem.net/>)
 - V-Dem captures the multidimensionality of *democracy*
 - high-level: electoral, liberal, participatory, deliberative, and egalitarian democracy
 - 4000+ variables, 27.000+ country-year observations
 - almost complete data (← quite rare in social sciences)