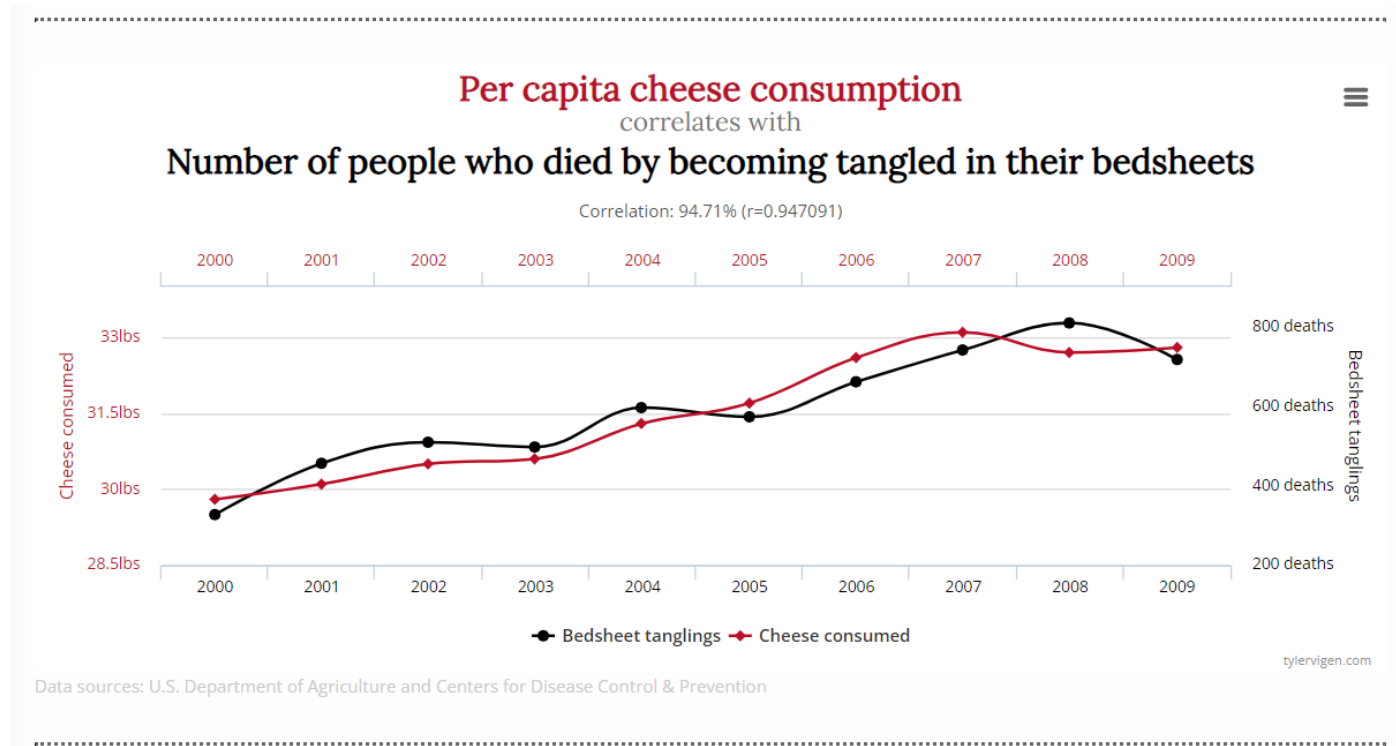*Exercise 7*

# Homework exercise

*To be solved at home before the exercise session.*

---

1.  a. Go to the website which lists pairs of variables that have no causal relationship but still exhibit a large correlation. Pick one of the datasets and figure out how the data is presented, i.e., how are the plots constructed from the $(x_i, y_i)$-data (the plots are *not* scatter plots of the two variables in question), how are individual pairs $(x_i, y_i)$ represented in the plots and what are the lines going through the points?



How the plots are constructed: There are two datasets: (xi, yi)1 represents (cheese consumption per capita, year) and (xi, y1)2 represents (Death by bedsheets, year). Then these two plots are merged together, creating a sense of correlation. The lines going through the points is the change of yi (death counts and chess consumption) with respect to the years

b. Let $x, y, \varepsilon$ be random variables such that,

$$y = x + \varepsilon,$$

where $\text{Var}(x) = 1$, $\text{Var}(\varepsilon) = \sigma^2 > 0$ and $x$ and $\varepsilon$ are independent (interpretation: $x$ and $y$ have a perfect linear relationship but the observed value of $y$ is contaminated with the noise/measurement error $\varepsilon$ having variance $\sigma^2$). Compute the Pearson correlation $\rho$ between $x$ and $y$ and investigate how it behaves when $\sigma^2$ is increased. Interpret this behavior.

We have: $\text{var}(x) = E[x^2] - E[x]^2 = 1 \Rightarrow \delta_x = \sqrt{\text{var}(x)} = 1$

$\text{var}(\varepsilon) = E[\varepsilon^2] - E[\varepsilon]^2 = \sigma^2 \Rightarrow \delta_\varepsilon = \sigma$

The Pearson correlation coefficient: $\rho(x,y) = \dfrac{\text{cov}(xy)}{\delta_x \delta_y} = \dfrac{\delta_{xy}}{\delta_x \delta_y}$

○ We have: $\text{var}(y) = \text{var}(x + \varepsilon) = \text{var}(x) + \text{var}(\varepsilon) + 2\text{cov}(x, \varepsilon)$

$= 1 + \sigma^2 \quad (\text{cov}(x, \varepsilon) = 0 \text{ because independent})$

$\Rightarrow \delta_y = \sqrt{\text{var}(y)} = \sqrt{1 + \sigma^2}$

○ For $\text{cov}(x,y)$: $\text{cov}(x,y) = \text{cov}(x, x + \varepsilon)$

$\Rightarrow \text{cov}(x,y) = E[x(x + \varepsilon)] - E[x]E[x + \varepsilon]$

$= E[x^2] + E[x\varepsilon] - E[x]^2 - E[x]E[\varepsilon]$

$= E[x^2] - E[x]^2 + E[x]E[\varepsilon] - E[x]E[\varepsilon]$

$= \text{var}(x) = 1 \qquad (\text{Since } x \& \varepsilon \text{ are independent})$

$\Rightarrow \rho(x,y) = \dfrac{\text{cov}(x,y)}{\delta_x \delta_y} = \dfrac{1}{1(\sqrt{1 + \sigma^2})} = \dfrac{1}{\sqrt{1 + \sigma^2}}$

$\Rightarrow$ If $\sigma^2$ is increased, $\rho(x,y)$ will decrease. This behavior means that the more varied the noise error $\varepsilon$, the less $y$ is correlated with $x$.

Indeed: $\lim\limits_{\sigma^2 \to \infty} \dfrac{1}{\sqrt{1 + \sigma^2}} = 0 \Rightarrow$ if $\sigma^2$ is very big, $x$ and $y$ are not correlated at all