

Spurious relationships in Twitter data

Data Science in Techno-Socio-Economic Systems

Spring Semester 2016

Claudio Bruderer

Andrina Nicola

`cbruderer@phys.ethz.ch`

`andrinaa@phys.ethz.ch`

Anna Weigel

`anweigel@phys.ethz.ch`

Abstract

... Claudio ...

include link to website

1 Introduction (Anna)

We have entered the era of Big Data and the wealth of information available to us entails its own new challenges. Relationships between quantities, whether real or without any obvious causal link, will naturally arise if the data sets are large enough (for examples see [1]). Even though we might be inclined to interpret apparent links between measurements, a correlation does not imply causation. One needs to find a way to distinguish spurious relationships between measurements from those with a causal link, and quantify the frequency of chance correlations occurring.

Based on a list of pre-defined keywords and a Twitter data, we aim to show that relationships between quantities without any obvious causal link can be identified in large data sets. We use a Twitter dataset which contains unfiltered tweets between **time range**. For each day, we determine the number of tweets containing one of our pre-defined keywords. We then measure and quantify the time-correlation functions between the popularity of

various keywords. To test our method we recover and quantify expected correlations between synonyms or related keywords. We additionally investigate and discuss the performance of different correlation measures.

2 Data (Claudio)

2.1 Keywords

... List keywords ...

maybe also mention why keywords are not suited, e.g. who or right

2.2 Dataset

... Twitter dataset ...

2.3 Data analysis

... Describe how we reduce Twitter data with pyspark ...

3 Time series analysis (Andrina)

... Correlation measures, k-means clustering, binning/smoothing ...

The main results of the data reduction explained in the previous section, are time series of the occurrence of each of the 100 keywords in the queried twitter dataset. In order to identify relationships between different keywords, we compare their occurrence frequencies using three different methods based on the Pearson correlation coefficient and k means clustering.

We can determine the amount of linear dependence between two time series X_i and X_j using the correlation coefficient ρ . The correlation $\rho(\Delta t)$ between two time series is a function of the time lag Δt between the two and is defined as [2]

$$\rho(\Delta t) = \langle (X_{i,t} - \bar{X}_i)(X_{j,t+\Delta t} - \bar{X}_j) \rangle. \quad (1)$$

For $\Delta t = 0$ we recover the usual correlation coefficient, whereas for $\Delta t \neq 0$, $\rho(\Delta t)$ quantifies the amount of linear dependence between shifted time series. If the two time series are equal i.e. $i = j$, then we recover the autocorrelation function of the time series, which quantifies the amount of linear dependence in the time series itself. In order to compare the correlation between different time series, it is customary to normalise Eq. 1 by the variance of the two time series i.e. [2]

$$\rho(\Delta t) = \frac{\langle (X_{i,t} - \bar{X}_i)(X_{j,t+\Delta t} - \bar{X}_j) \rangle}{\sqrt{\langle (X_{i,t} - \bar{X}_i)(X_{i,t} - \bar{X}_i) \rangle \langle (X_{j,t} - \bar{X}_j)(X_{j,t} - \bar{X}_j) \rangle}}. \quad (2)$$

We can estimate the value of $\rho(\Delta t)$ from the sample covariance and variances.

In order for Eq. 1 to be a valid correlation measure, the values of both time series at all times t need to be identically distributed (IID) random variables. This means that both time series X_i, X_j need to be stationary i.e. that their mean μ , variance σ^2 and autocorrelation function $\rho(\Delta t)$ do not depend explicitly on time t

$$\mu = \langle X_{i,t} \rangle, \quad (3)$$

$$\sigma^2 = \langle (X_{i,t} - \bar{X})^2 \rangle \quad (4)$$

$$\rho(\Delta t) = \frac{\langle (X_{i,t} - \bar{X})(X_{i,t+\Delta t} - \bar{X}) \rangle}{\sqrt{\langle (X_{i,t} - \bar{X})(X_{i,t} - \bar{X}) \rangle \langle (X_{i,t+\Delta t} - \bar{X})(X_{i,t+\Delta t} - \bar{X}) \rangle}}. \quad (5)$$

We do not expect the occurrences of keywords in twitter data to be stationary time series, since the mean occurrence of any keyword will for example depend on world events, its popularity and also on the number of people using twitter. We can assess if any time series is stationary by evaluating its autocorrelation function $\rho(\Delta t)$. For an IID time series, $\rho(\Delta t)$ will be close to zero whereas non stationary time series are characterised by slowly decaying autocorrelation functions. In order to be able to quantify the correlation between non stationary time series we need to make the series stationary prior to evaluating Eq. 1. This process is called pre-whitening. The most basic form of prewhitening of the differencing method, which allows to transform time series with a slow trend to stationary series. A time series with a trend can be written as [2]

$$X_{i,t} = T_{i,t} + R_{i,t}, \quad (6)$$

where $T_{i,t}$ denotes the time series' trend and $R_{i,t}$ are random fluctuations around this trend. If the time series varies slowly, we can remove the trend by considering the differenced time series instead of the initial one. The differenced time series is defined as [2]

$$D_{i,t} = X_{i,t+1} - X_{i,t}. \quad (7)$$

If $T_{i,t}$ is slowly varying, this transformed time series will only capture the random variations around the global trend. The cross correlation between two differenced time series

thus encodes the amount of linear dependence between changes in one time series from its global trend and changes in the second.

In order to find relationships between the keywords we will proceed in two steps: for illustration purposes we will compute the cross correlation between all the time series regardless of stationarity. We will then compute the autocorrelation function of all the series and recompute the correlation coefficient using the differenced time series for non stationary processes and compare these two measures.

4 Results (Anna, Andrina)

...

5 Conclusion (All)

...

6 Acknowledgements

This project was part of the lecture *Data Science in Techno-Socio-Economic Systems* given by Dr. Evangelos Pournaras, Prof. Dr. Dirk Helbing and Dr. Izabela Moise at ETH Zurich during the spring semester 2016. We would like to thank Dr. Izabela Moise and Dr. Rok Roskar for the valuable guidance and helpful discussions. We used APACHE SPARK [3] to query the Twitter data, and PYTHON [4] and BOKEH [5] for the data analysis and visualization. The website was built on a HTML5UP [6] template.

References

- [1] <http://tylervigen.com/spurious-correlations>
- [2] Marcel Dettling, Applied Time Series Analysis, SS 2014, Zurich University of Applied Sciences, https://stat.ethz.ch/education/semesters/ss2014/atsa/Scriptum_v140523.pdf

[3] <http://spark.apache.org>

[4] <https://www.python.org>

[5] <http://bokeh.pydata.org/en/latest/>

[6] <http://html5up.net>