# Can R Notebooks help with reproducibility?

## Introduction

### What is reproducibility?

In science it's by some considered to be the holy grail of statistics. Basically it means doing an experiment or research multiple times by different people or groups, with the same numbers and get the same result.

### Challenges

One of the big challenges of (statistical?/Data?) science has been to reproduce the identical results from an earlier research. In a significant amount of research it has proved hard or flat out impossible to reproduce the results from other peoples and older research for a multitude of reasons. One of the big reason is quite simply that the data used in the research is unobtainable, due to them being "too old" or lost. Another reason is the author doesn't want to share the data or can't find it (again, lost). Other challenges is the usage of different systems, equipment, and so on. Theories around the research and calculations can also change over time.

### Computable Documents

In more modern times a new idea around this have been brought up: computable documents. The idea behind this is when publishing a research paper, one submits the data, equations and calculations together with a document that's computable. That way when other researches wants to run the experiment, they have access to everything, and can see what others did before them. So the question then becomes: can computable documents like R Notebook help with reproducibility?

## Short literature review

According to McNutt (2014) advancing in science is reliant on discoveries that can be trusted, but stated that studies that have been performed, can't be reproduced. In tandum with this, Peng (2011) has brought up the idea that reproducibility should be a requirement for publishing research.

Already in the early 1930's the accessability to the data for a research was aired by Frisch He stated that in statistical research, the raw data should be published (1933). Later in the 1960's it was deemed nearly impossible to reproduce research on economy with big models. Then in 1982 there was a research project done: the Journal of Money, Credit and Banking, where they tried to replicate research article with only the submitted data. This is looking back, very close to reproducibility in modern terms. The results was a staggering 2 of 70 articles could be reproduced. They concluded was because for the most part missing data, documentation and computer systems, as quite few of the authors would provide their data and coding. (Kilde. Arnstein sin e paywall) To counteract this, they came up the solution to store the coding and raw data used in research articles, which in a sense is the foundation of computable documents today.

Kva starta "computable documents" med og kortid? Ka e gjort i forkant ang. R Notebook-programmer og reproduserbarhet? Ka seie forfattar X om Y? Ka seie Forfattar Z om Y? Litt "bakgrunn" for detta med "computable documents".

–> ca. 1 sida +/-

# Discussion

Why is reproducibilty important? What are the issues surrounding it? What steps have been taken to improve the situation? What does scientist say about it?

# Conclusion

Help? Yes, to a degree. Solve/Fix the problem? No, atleast not yet. Too many issues. Different systems, different OS, different software and packages. Theories and calculations can change over time. Packages being updated/outdated. R have "Session Info" to list systems, maybe a library of packages would help? It's seems to be a step in the right direction, but it still remains issues around reproducibility even with the help of computable documents like R Notebook.

# References