# Measuring User Behavior

I'm always looking for different ways to look at data. Things that will give me insights into how users act as they go about their daily activities. Having done this for a while now I can say that people typically like routines as they go about their work, but how they go about their work is largely dictated by their job role. Examples may be the developer that often creates archive files or the sales person that authenticates from multiple locations per day.

Operating under the assumption that people will generally look the same from day to day it would make sense to understand why someone may deviate from their normal routine, especially from a security context. One way we can do this is to measure their own, unique behavior and apply some anomaly detection to highlight days of interest.

For the example I'm going to show, I wanted to focus on insider from a data theft standpoint, but this method can be applied to many different scenarios as long as you can define the points of interest to observe. The hypothesis that I'm using for insider data theft is:

1. Users who knowingly steal data will often use deceptive actions.
2. They will perform actions that are new or rare for them.
3. They will use uncommon exfiltration paths.
4. Rare actions across multiple phases can be identified.

The phases (categories) I'm using here are Deception, data staging and exfiltration.

Below I'm pasting screenshots of the relevant portion of the Jupyter notebook I'm using. The cool thing about using Jupyter is it's independent of the log source. I can use this as long as the data can be exported or accessed via api. Code blocks below are commented to describe the function.



Below is an example rule from an external rule file.

```
# Attempted file copy over bluetooth
f4 = df[df['Application'].str.contains('fsquirt.exe')]
if f4.shape[0] > 0:
    shapesize4 = 5
else:
    shapesize4 = 0
if shapesize4 == 0:
    print("Rule 4 not found")
else:
    print("Rule 4 detected.  Attempted file copy over bluetooth.")
```





Below would be the output from the anomaly detection.



The fields that are being measured are:

Rule_1 - Rule_6: Detection rules from rules file
Archive_Count: The number of archives created by day
Mail_Count: The number by day of emails sent to a personal email provider with an attachment
Rename_Count: The number of files renamed in a day
NTU_Count: The number of files uploaded to an external source in a day
Weighted: The sum by day of additional values given to rule detections
Day_Count: The total count of all events by day

This method won't tell you 100% that anomaly is malicious.  But you may be able to infer what the user was doing and the likely hood of needing further investigation.  An example may be 2 archive events + 2 email events may result in 2 zip files being emailed to a personal email address.  That may or may not need investigating depending on the users role.