

**Time-Series Analysis of Host Lifestyle and Human Microbiota**

**Final Report**

Timothy Merkel, Max Paszek, Sean Stevens

Drexel University

March 24, 2017

## **Abstract**

We used data provided from the study titled “Host lifestyle affects human microbiota on daily timescales”. The researchers examined the daily human microbiota and linked them to various life circumstances [1]. Our project goal was to analyze the data collected from a Digital Signal Processing perspective. No members of the group have backgrounds in biology or metagenomics, so tried and true analyses for time series data were used. These include, but are not limited to, correlation, arima, and periodicity plots. To achieve this, we utilized the programming language R to acquire the data and process it. The methods to manipulate the data will be discussed in greater detail within the Materials and Methods section.

## **Materials and Methods**

Prior to analyzing the data, we first had to acquire the data found in [1]. In order to accomplish this, we first had to familiarize ourselves with R. Thankfully, the markdowns on the Bioinformatics GitHub page were available to use as reference [2]. The first step was to follow the `r_sra` markdown found on [2]. The first step to download an SRA file to connect to the SRA metadata database, and then download the files from the NCBI. After downloading the SRA files, one must obtain FASTQ files by using the SRA toolkit found on the NCBI's site. A FASTQ file is just a FASTA file that also includes the quality data of the sequence. Recently, FASTQ files have become the standard for storing high-throughput sequencing techniques [3].

After the FASTQ files were obtained, they were then filtered to find higher quality data measures. In [1], the researchers trimmed the first 10 left base pairs off of their sequence data, so we did the same. We also used a truncation length of zero in our program, this yields sequences that are 90 base pairs long. These sequences, of which there are 1,698 were then placed into a sequence table. The sequence table was created by using DADA 2, of which there are very detailed instructions found in the `dada2` markdown on [2]. After completing the DADA 2 workflow, the sequence table below was generated. The code that was used to accomplish all of these steps are published on our GitHub page.

ACAAGGGTGGT CCGAGGGTGGT ACAAGGGTGGT CCG																								
---	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

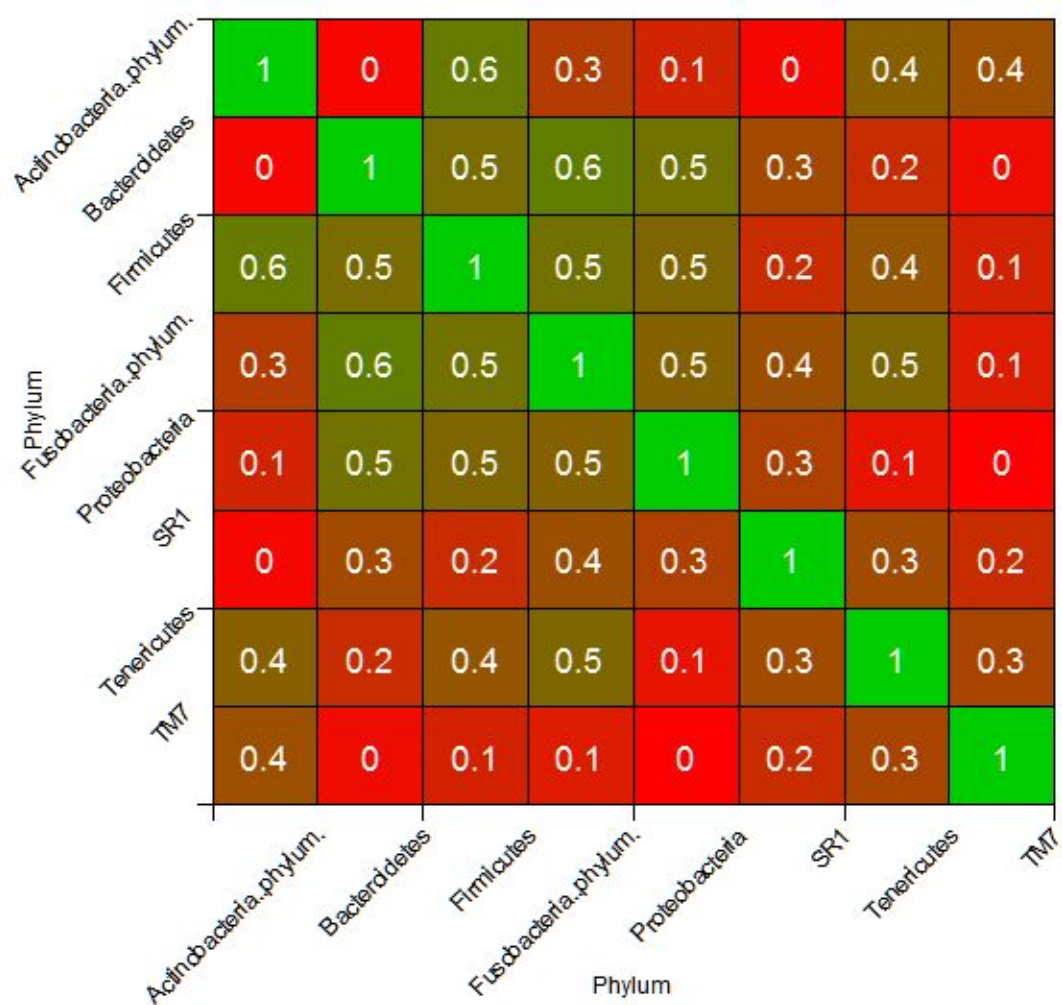
**Figure 1.** The sequence table generated from the FASTQ files by using DADA 2.

The sequence table contains the SRA files names as the rows, then the abundances of the sequences as the columns. This is not the end of the process for preparing the data for R, however. Each SRA filename must be matched up with a sample ID from the study. In order to accomplish this, a metadata table was downloaded from [1]. This table is additional file 18, and contains both sample and nutritional metadata. The accession number for the nucleotides are deposited on the EBI/ENA database under accession number ERP006059. While there may be tools to accomplish the mapping of the taxa to the sequences, nobody in the group was comfortable with using such a tool. Therefore, Excel was utilized to manually manipulate the data and construct a table that contained the abundances of each taxa per each sample ID. The data was broken up into the following categories: saliva samples from patient A, stool samples from patient A, and stool samples from patient B. Tables were also generated

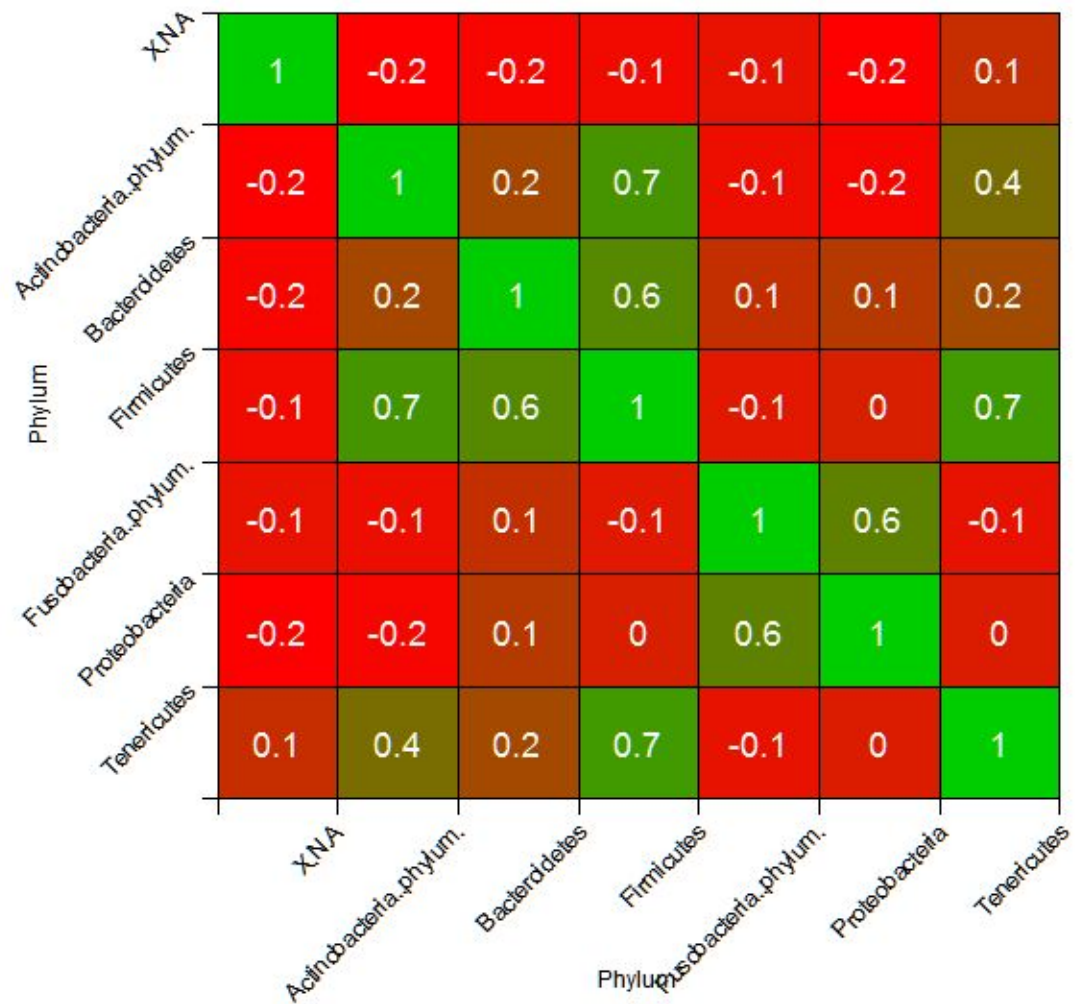
which contained the nutrition information for the tables mentioned previously. The results of this pre-processing then allowed the data to be processed in R.

With the data properly formatted, we were able to visualize and analyze our data by using R. The first plots which were constructed were correlation plots between the various phyla across all samples that were assigned to the sequences in the pre-processing stage.

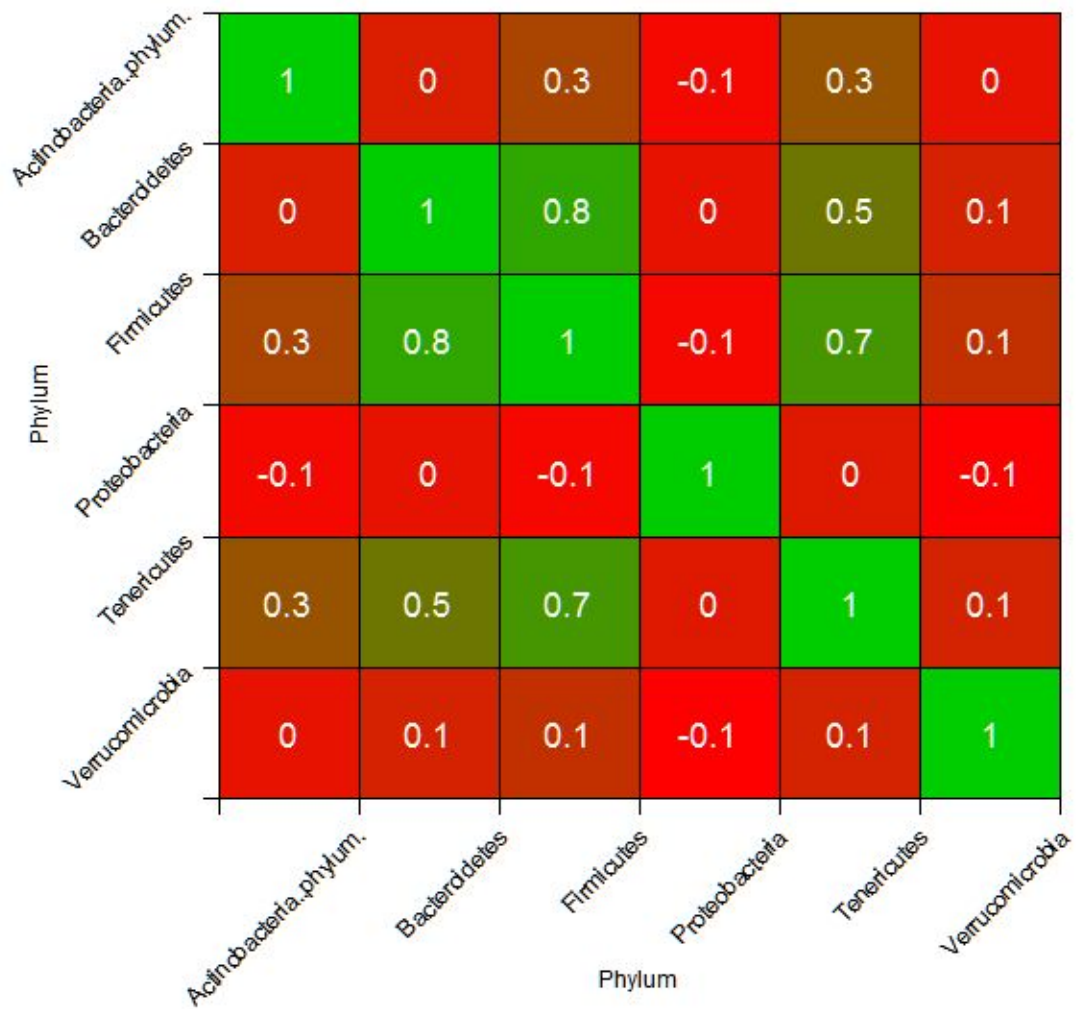
## Results and Discussion



**Figure 2.** The correlation plot of patient A's saliva data across all samples.



**Figure 3.** The correlation plot of patient A's stool data across all samples.



**Figure 4.** The correlation plot of patient B's stool data across all samples.

The `cor()` function in R was utilized to generate the correlation matrices, which were then plotted using the `color2d.matplot()` command.

This data is very interesting because it shows that not only are their correlations between the individual species in the data, but also that taxonomic units even as large as a

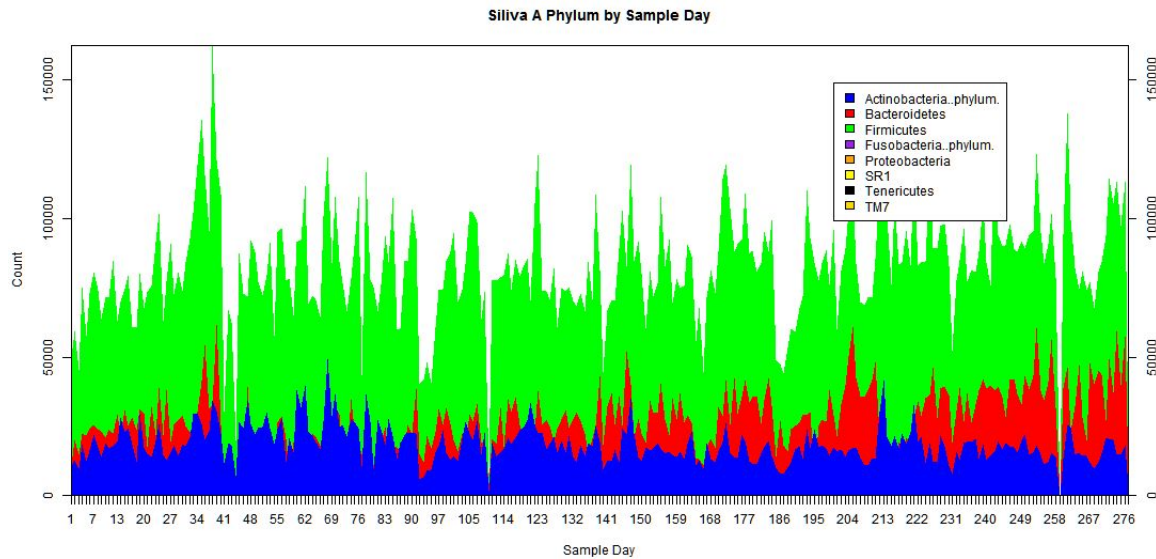


phylum have some pretty significant correlations. For instance, in Figure 2 it is clear that there is a strong positive correlation between Actinobacteria and Firmicutes. This specific example however is not consistent between all 3 samples. In fact, in patient B's stool the same pair has the most negative correlation of the whole sample set. One relationship that stands out is the relationship between Bacteroidetes and Firmicutes. This correlation is very positive in all 3 sample sets with a 0.8, 0.6 and 0.5. This must mean that there is some factor which causes these phylums to change at the same rate. Perhaps they have a symbiotic relationship, or perhaps they both just happen to be sensitive to the same nutritional intake.

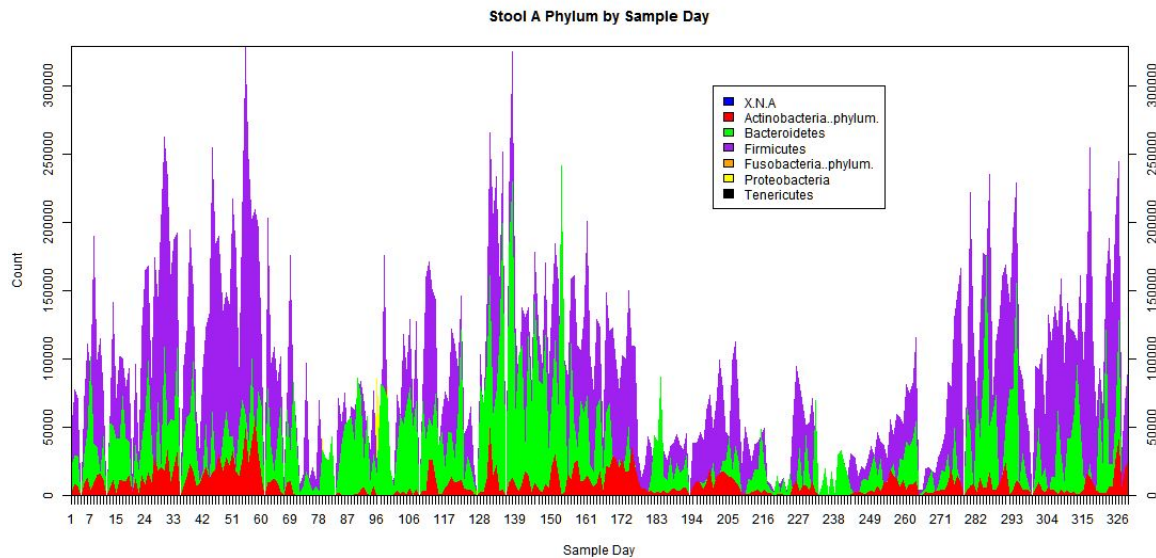
Next the phylum counts per day were visualized. They are shown below in Figure 5, 6, and 7, and were generated using R's `stackpoly()` command. In these stacked area charts one of the most interesting observations is that the saliva of patient A and the stool of patient A contain a very different distribution of phylums. Both sets contain the same three phylums but the quantities change drastically. For example in the saliva of patient A it is clear that the majority of the microorganisms belong to the Firmicutes phylum. In the stool of patient A there was a period between 70 and 150 days where Bacteroidetes dominated.

Interestingly, this lines up very well with the time period where the patient went abroad. In the saliva of patient A this period also sees a change in that the bacteroidetes are drastically reduces. Given the correlation data between bacteroidetes and firmicutes

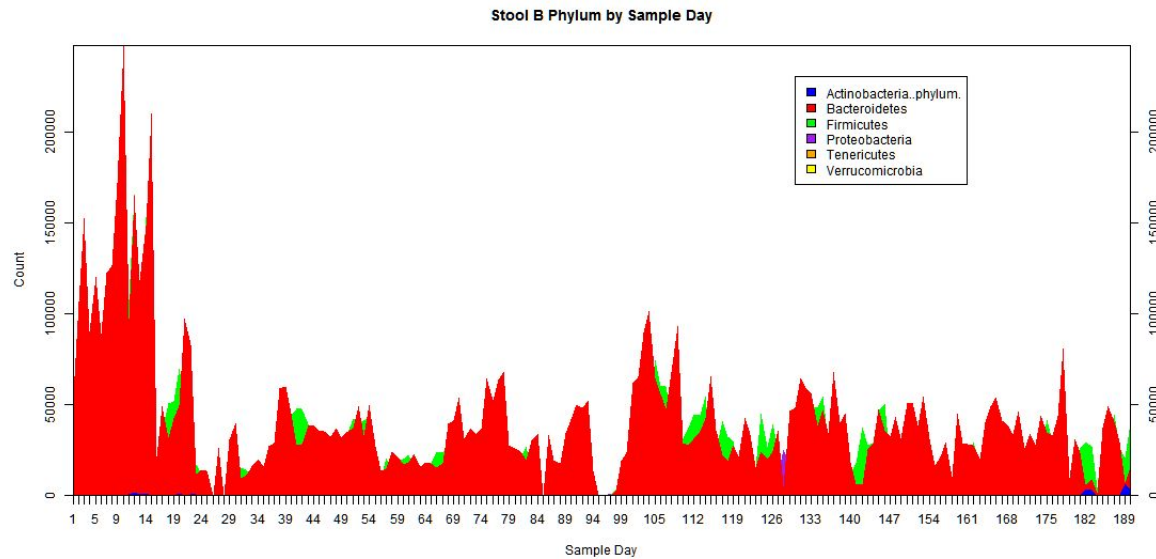
(one of the most correlated in the sample) it is extremely interesting that a decrease in the bacteroidetes in the saliva of patient A corresponds so well with the decrease in Firmicutes in the stool of the same patient over the same period.



**Figure 5.** Patient A's saliva samples grouped by Phylum, per day.



**Figure 6.** Patient A's stool samples grouped by Phylum, per day.



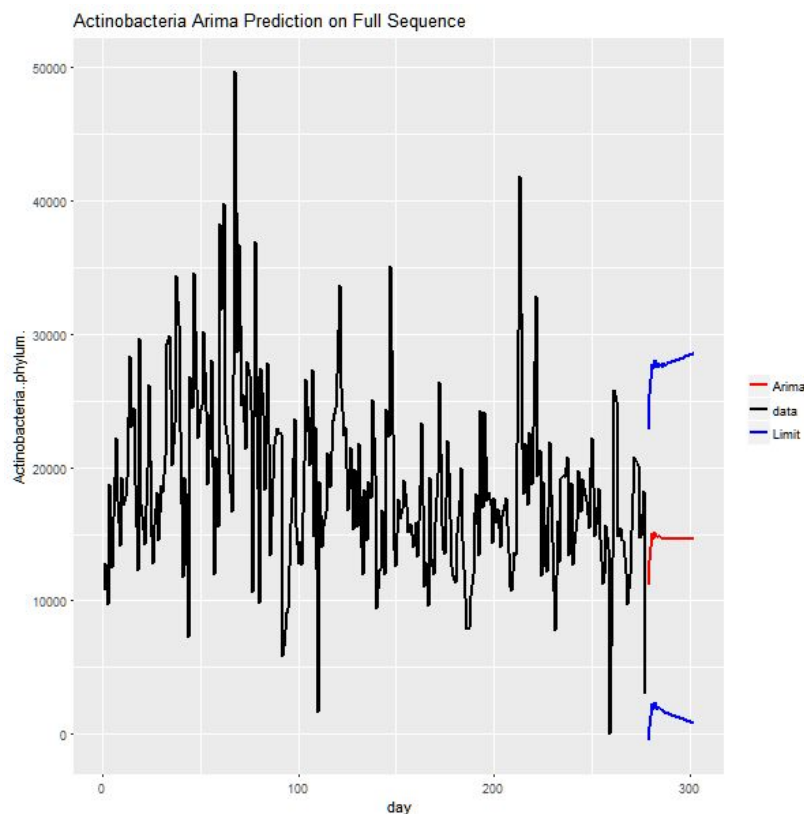
**Figure 7.** Patient B's stool samples grouped by Phylum, per day.

Arima plots for the three most common phyla in patient A's saliva were generated by using R's `arima()` command and plotted using the `ggplot()` command. Two different window lengths were used. In the first set the arima model was trained on the full data set and some future data was guessed. The arima parameters used on this set were (3,1,3) with seasonality parameters (0,1,1). These results can be found in Figure 8, 9, and 10. In the second set the arima model was trained on the first 40 days and then was used to attempt to model the time period when the patient was abroad. The arima parameters used on this set were also (3,1,3) with seasonality parameters (0,1,1). These results can be found in Figure 11, 12, and 13.

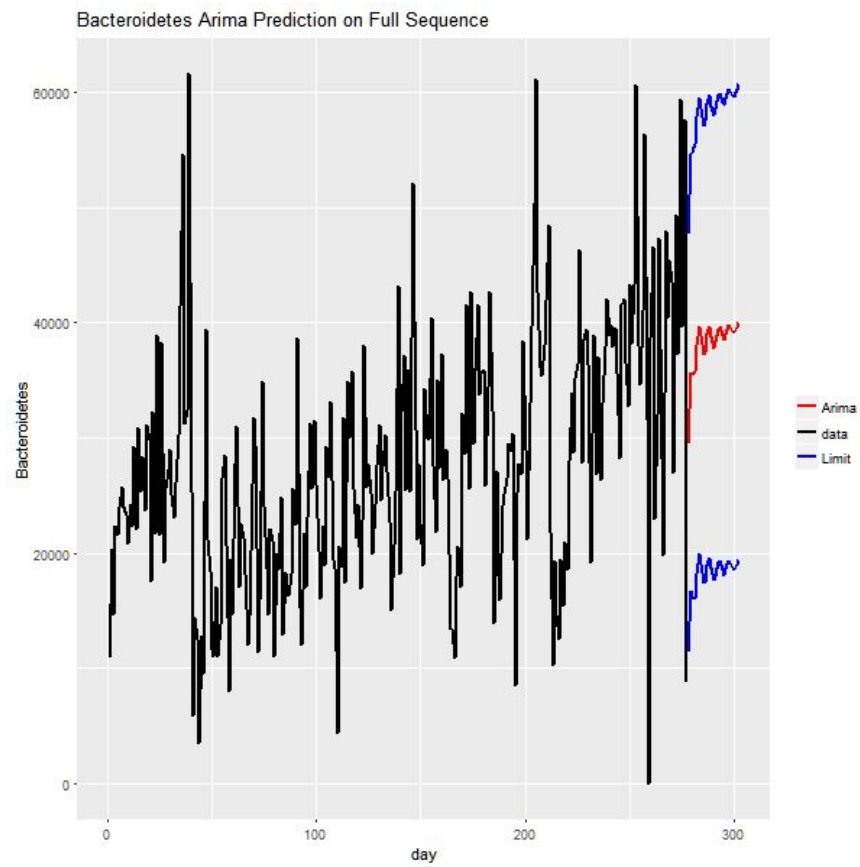
For the set that was trained on the full data the arima model seems reasonable. Each one has found a mean value that seems consistent with the current trend

of the data. The upper and lower limit do seem to line up reasonably well with the possible spikes of the data. This data is not incredibly interesting. The only conclusion there is to draw from this is that the arima model is working relatively well.

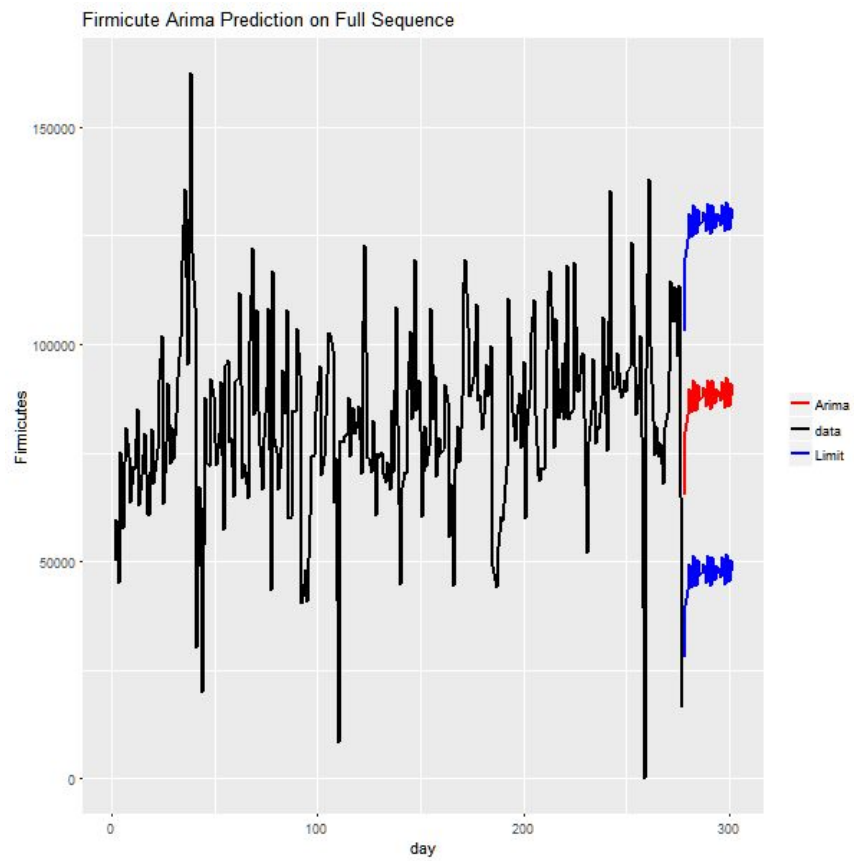
The second set is much more interesting. In each of these figures it is clear that the arima model is not correct. The reason for this is that the patient has just gone abroad causing a change in habits, nutrition consumption, and living conditions. As a result the microorganisms inside their saliva became very different than the training set. Upon returning home their habits returned to normal and the arima model trained on the first 40 days would likely work once again.



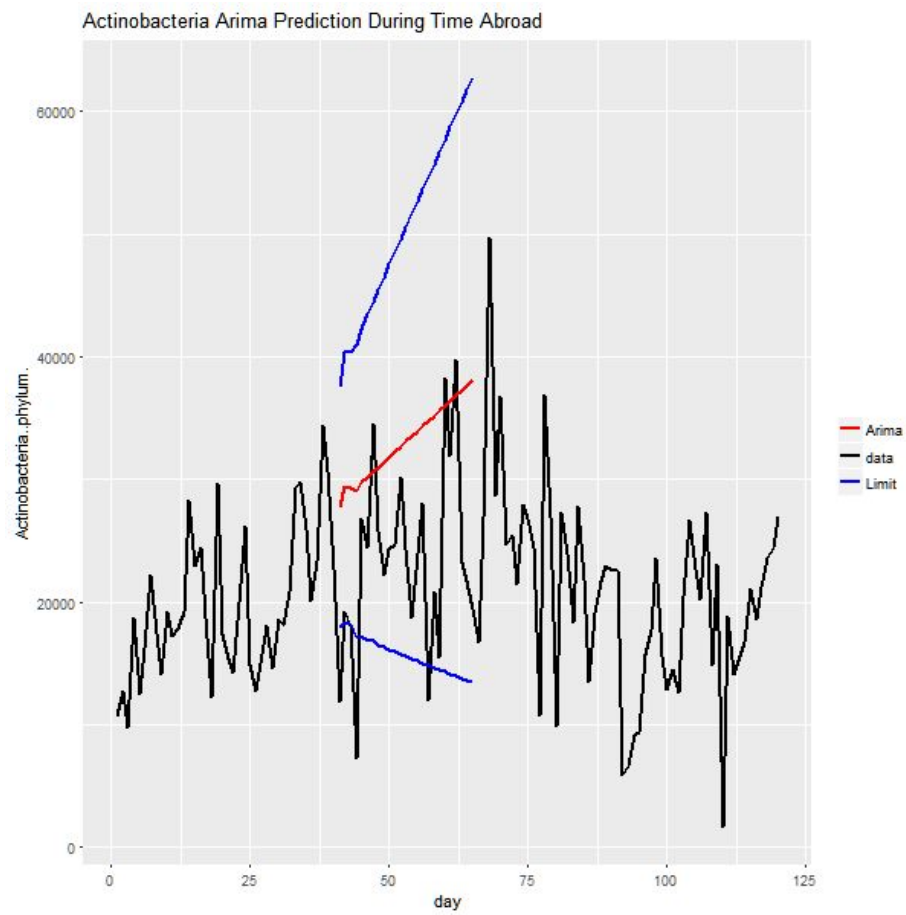
**Figure 8.** Patient A's saliva actinobacteria count arima forecast.



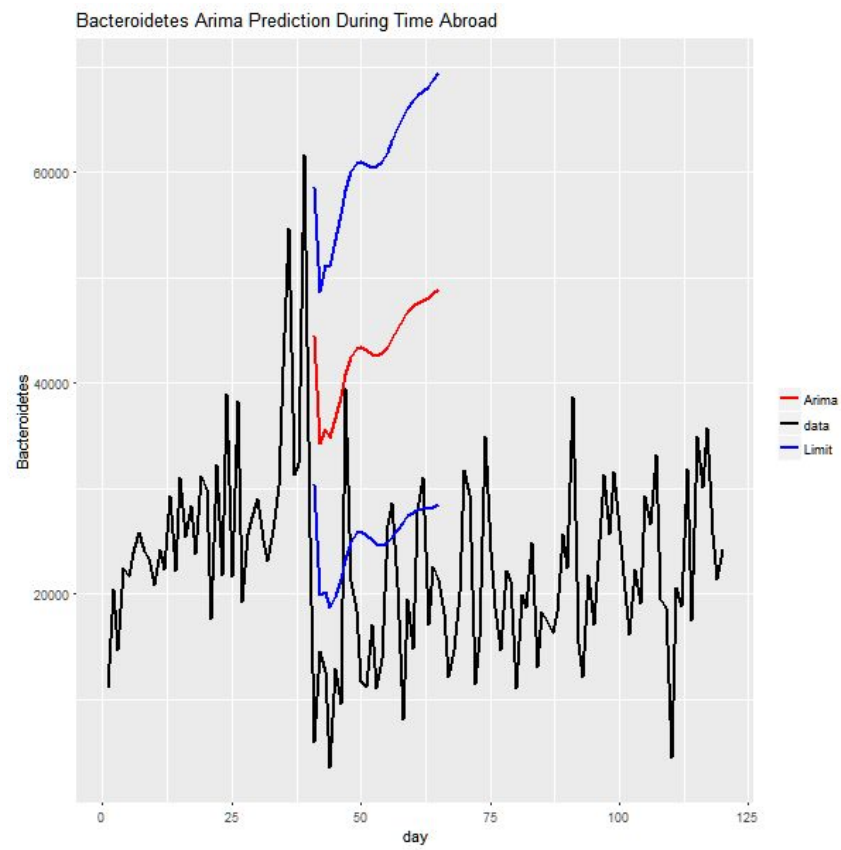
**Figure 9.** Patient A's saliva bacteroidites count arima forecast.



**Figure 10.** Patient A's saliva firmicutes count arima forecast.

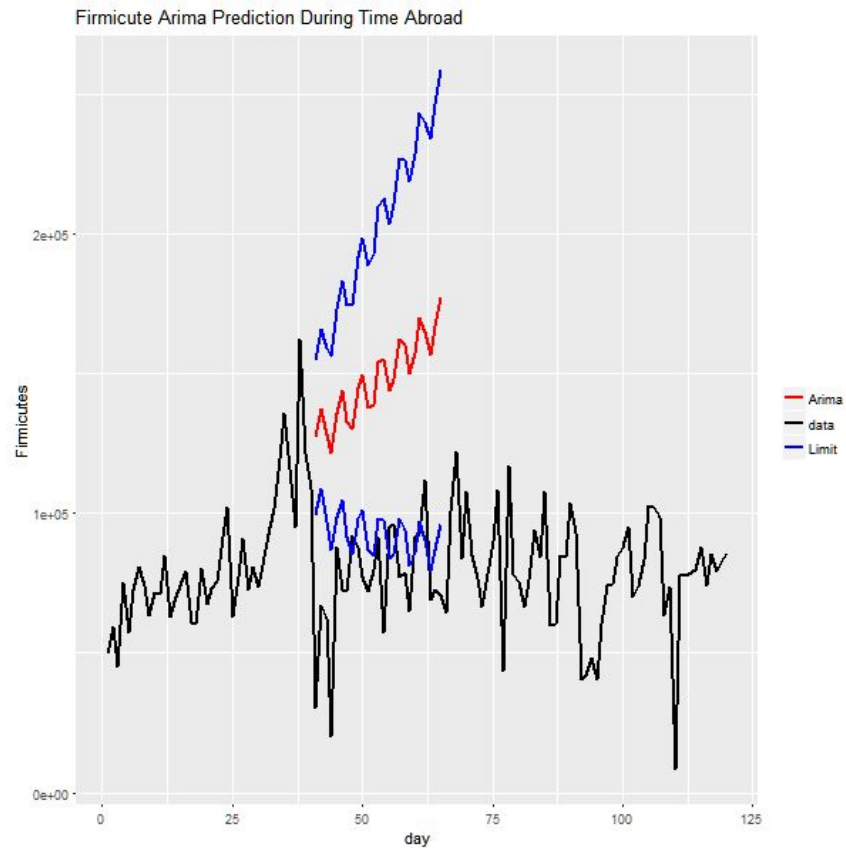


**Figure 11.** Patient A's saliva actinobacteria count arima forecast against truth while abroad.



**Figure 12.** Patient A's saliva bacteriodite count arima forecast against truth while abroad.





**Figure 13.** Patient A’s saliva firmicutes count arima forecast against truth while abroad.

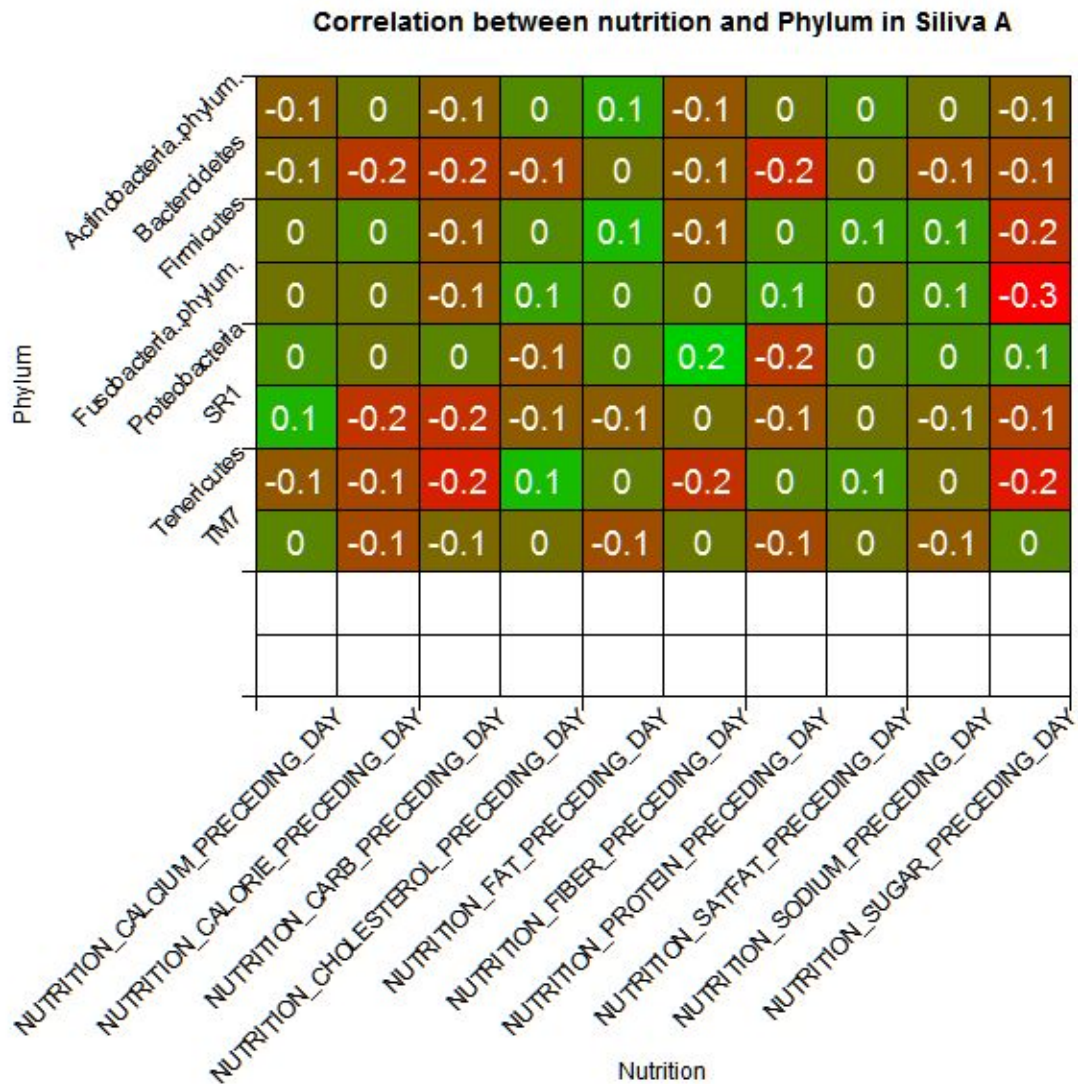
While all of the above factors are just merely representations of bacterial concentration against time, the data was also correlated with the nutritional intake that was recorded by the participants. Only patient A recorded their nutritional intake, so the correlations could only be found in patient A’s saliva and stool. Figure 14 and 15 below show the results of these correlations.

One thing that is interesting about the saliva data (figure 14) is that there are a lot of zeros present in this data set. While there are no strong positive or negative

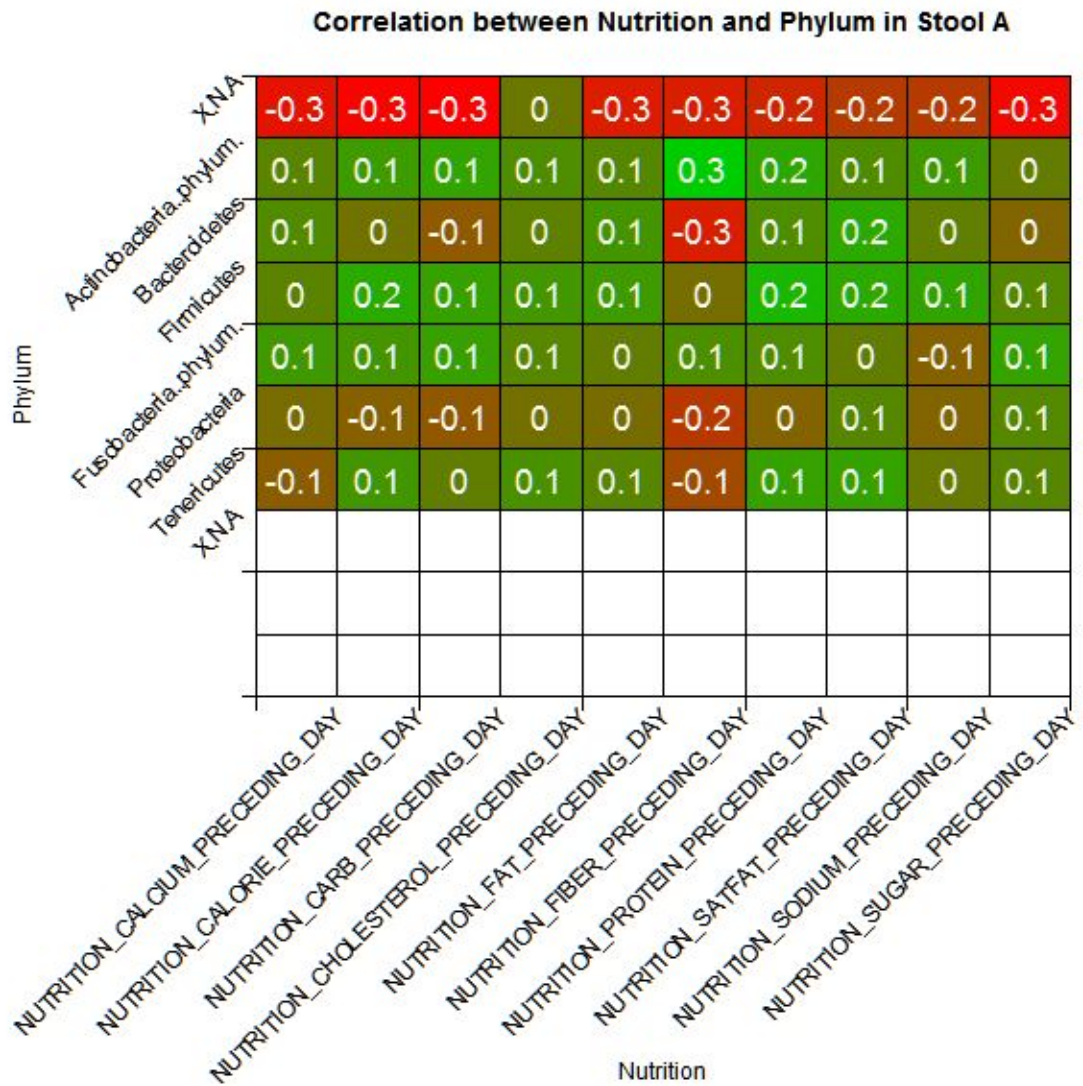
correlations a zero means that the phylum does not respond at all to that stimulus. One example is that the only phylums which respond to saturated fat are Tenericutes and Firmicutes. All others present in this data set are not affected at all.

The nutritional intake seems to affect the stool of patient A significantly more than the saliva. This makes sense because the stool is a better representation of the stomach than the saliva. Nutritional intake collects in the stomach so those microorganisms will respond to in more vibrantly. Still there are no outstanding correlations. However, as seen in figure 15, it is clear that each of these phylums does have a response to most the nutrition types with the maximum correlation being between Actinobacteria and the previous day's fiber intake.

The phylum level seems to be too high level for this analysis. At the phylum level these correlations are all relatively close to zero. In order to increase the value of this analysis one option would be to look at the data on a more granular scale such as genus or family. This could allow the truly positive and negative correlations to show themselves without having to be held down by the other members of their phylum. Another way to potentially find more interesting results could be to introduce a lag in the system. Perhaps it takes a few days before the effect of the nutrition can really be seen. This would cause the correlations to appear as though nothing is happening. By introducing a lag one could attempt to combat this effect.



**Figure 14.** Correlation between nutrition and phylum in patient A's saliva.



**Figure 15.** Correlation between nutrition and phylum in patient A's stool.

## **Conclusion**

In conclusion, we were able to visualize the data from the study in various methods. While our biology knowledge limited, we were still able to draw meaningful conclusions from the data. It is easy to see that certain life events such as travel can cause drastic changes in the microbiota of any person. It is clear that there is a correlation between different phyla in both the saliva and stool of the patient tested in this study. It is also seen that there is a correlation between nutritional intake and certain organisms. The analysis of nutritional correlations in this study was likely to high level to see the deeper trends, but there are a few indicators that they are present. For future work, the fft analysis that was starting to be explored in the R code would be fleshed out to be able to draw conclusions. Overall this project helped the team to understand some biological topics as well as develop a deeper understanding of R as a scripting language.

## References

- [1] David, Lawrence A., Arne C. Materna, Jonathan Friedman, Maria I. Campos-Baptista, Matthew C. Blackburn, Allison Perrotta, Susan E. Erdman, and Eric J. Alm. "Host Lifestyle Affects Human Microbiota on Daily Timescales." *Genome Biology* 15.7 (2014). Print
- [2] Sw1. "Sw1/Bioinformatics." GitHub. 02 Mar. 2017. Web. 22 Mar. 2017.
- [3] Cock, P. J. A., C. J. Fields, N. Goto, M. L. Heuer, and P. M. Rice. "The Sanger FASTQ File Format for Sequences with Quality Scores, and the Solexa/Illumina FASTQ Variants." *Nucleic Acids Research* 38, no. 6 (12, 2009): 1767-771. doi:10.1093/nar/gkp1137.