

Tutorial on Normalization of Microbiome Data

Contents

1	Introduction	5
1.1	The importance of normalization	5
1.2	The compositional nature of microbiome data	6
2	Importing Data	9
2.1	Global Patterns	9
2.2	Pre-processing Quality Control and Filtering	10
3	Total Sum scaling (TSS)	13
3.1	About TSS	13
3.2	TSS Implementation	14
3.3	TSS on Global Patterns	14
4	Rarefying	17
4.1	About Rarefying	17
4.2	Rarefying implementation	17
4.3	Rarefying on Global Patterns	18
5	DESeq	21
5.1	About DESeq	21
5.2	DESeq Implementation	22
5.3	DESeq on Global Patterns	23
6	TMM (edgeR)	27
6.1	EdgeR TMM implementation	27
6.2	TMM on Global Patterns	28
7	Cumulative sum scaling (CSS)	31
7.1	CSS implementation	31
7.2	CSS on Global Patterns	32
8	GMPR	35
8.1	GMPR Implementation	35
8.2	Global Patterns GMPR	36

9 Wrench	39
9.1 Wrench Implementation	39
9.2 Wrench on Global Patterns	40
10 Comparisions	41

Chapter 1

Introduction

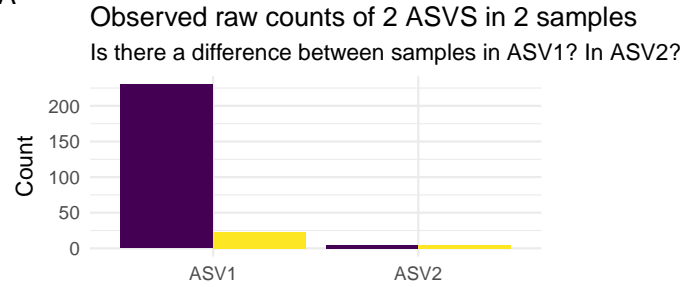
1.1 The importance of normalization

Microbiome data must be normalized before any statistical analysis can be performed. Following the process of sequencing and assigning raw reads into counts per observed and classified identified taxa classes/OTUs/ASVs, the data are in the form of a matrix of read counts. Normalization is the process of transforming raw read count data into data that can be compared between samples. Statistical analysis on this count matrix is then performed depending on the goal of the experiment. Common analysis goals include community-level analysis (alpha/beta diversity), differential abundance testing (the parallel of differential expression testing in gene expression studies), and network analysis.

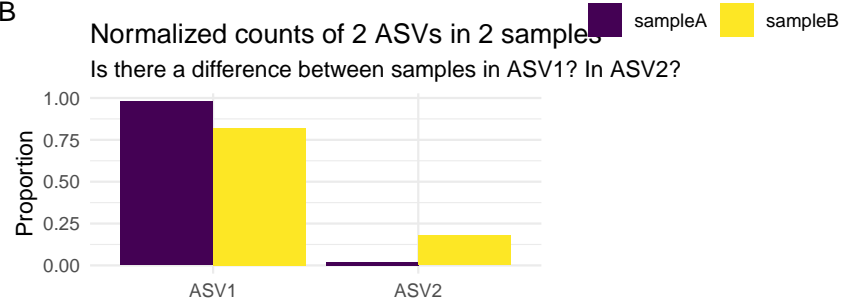
Analysis of composition, differences, connections, etc. should be done based only on true biological aspects. However, technical variation in counts across samples is a given hurdle that must be accounted for. Biases can arise in the sequencing process, sample preparation, contamination, preferential amplification, and can manifest in differences in sparsity and unequal sequencing depths (Salter et al., 2014). An effective normalization strategy should put all samples on equal footing so interpretations are on biological signals, not technical signals such as sequencing depth. Currently, there is no known ‘best’ normalization method that removes all technical artifacts leaving only biological signals.

Due to the sequencing technology, samples will have different sequencing depths, or the sum of all the counts in a sample. Directly comparing raw counts between samples is not possible. To illustrate this, consider the counts of one taxon, labeled ASV1, across two samples shown below. In Sample A, this taxon has a count of 230, and in Sample B, this taxon has a count of 23. Is this taxon differentially abundant between samples?

FigureA



FigureB



Normalization for microbiome data often refers to standardizing sequencing depth across samples. One common approach to this is a scaling-based approach, where a scaling factor is calculated for every sample and the counts for each taxon are divided by the scaling factor for that sample. Figure B shows the same data as figure A, but where each sample has been transformed into proportions by dividing by the total counts for each sample. The difference in ASV1 between samples appears much smaller. However, now there appears to be a difference in ASV2, even though the counts were originally the same. This is because in sample B ASV1 consists of a higher proportion of the total count than in sample A.

This demonstrates the importance of normalization, but also the artifacts that can occur depending on the method.

1.2 The compositional nature of microbiome data

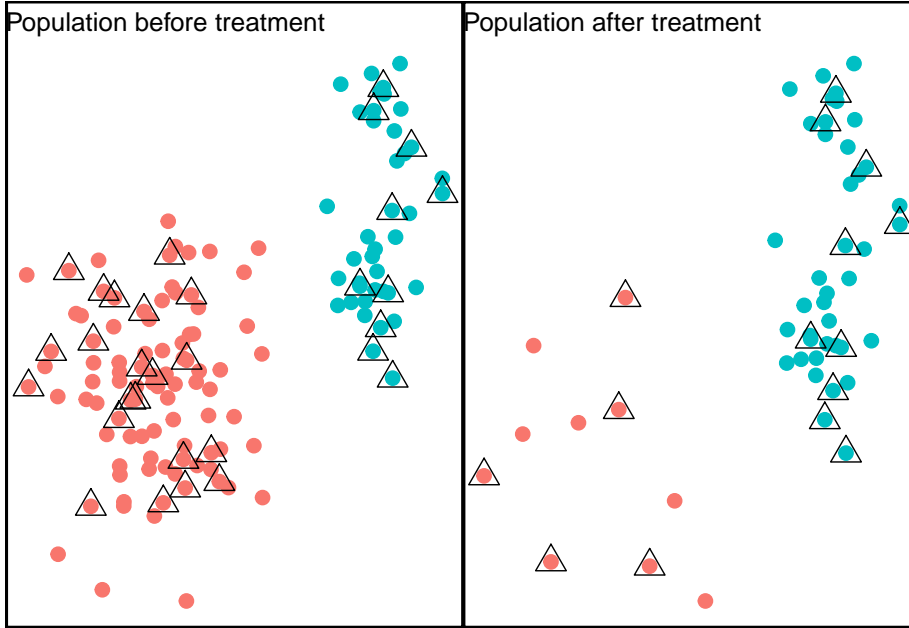
Microbiome data are inherently compositional. The counts of the collection of taxa that make up each sample are constrained by the total sum, or sequencing depth for that sample. This means that the count of each sampled taxon is a portion of a larger whole. Each observed taxon is not independent. As we saw in the above example, before normalization, ASV2 was equal between samples. After converting to proportions, ASV2 no longer appears equal. If there is a difference between two samples it is unclear if that difference is because of a

Table 1.1: Counts of sampled red and blue taxa before and after

Sample	blue	red
Before	10	20
After	10	4

true difference in that taxon or if that taxon is changing because of differences in another taxon. Numerous traditional statistical methods rely on an independence assumption, which is not met with microbiome data. This can lead to spurious correlations that exist only because of the compositional nature and not any true signal.

With library size as the sum constraint for each sample, if we know in a biological system that after an event occurs (treatment), the red taxon decreases, this will change the composition of the sampled blue taxon regardless of its change or lack thereof in the underlying population.



Consider again two samples consisting of red and blue points. We can think of the samples as before and after treatment. In the second plot, the number of red dots in the population and in the observed sample has decreased, but the blue remains the same.

This observed increase in the proportion of blue is due to the compositional nature of the sampled points, and not any true difference in the blue population.

Table 1.2: Proportions of sampled red and blue taxa before and after

Sample	blue	red
Before	0.2857143	0.6666667
After	0.7142857	0.3333333

1.2.1 Log ratio methods

The log-ratio based methodology developed by Aitchison in the 1980s is useful for analyzing compositional data (Aitchison, 1982). Taking the logarithm of ratios can be an appropriate transformation for compositional data, so standard statistical tests can be appropriate again. This transformation removes the issue of standardizing/normalizing different sampling depths. The sampling depth for a given sample will not distort the biological covariance or correlation structure.

This log-ratio method has a drawback, which is the decision of how to define the denominator. One approach to this problem is to use one sample as the reference. This sample should be ‘representative’. The log-ratio transformation is then the ratio of every other taxon to that representative sample. Of course, knowledge of what makes a sample representative is hard to come by and often unknown, and subsequent results can be affected by this choice. This method is frequently called the additive log-ratio approach (alr). The alternative approach is to use the data to create a pseudo-reference sample. This pseudo-reference sample is the geometric mean of the counts of all taxa. This is called the centered log-ratio method.

While promising, these log-ratio methods have drawbacks in practice. Microbiome data are often incredibly sparse, with up to 80-90% of count matrices containing zero counts. For ratio transformations, if we have sparse data, the geometric mean can be zero. Then the ratio is undefined, and further, so is the logarithm of zero count taxa.

One solution to this is adding a small pseudo count to every element in the data. This removes problems occurring from having zero counts in the data, but there is not a clear best choice of what pseudo count to use, and it the choice can impact downstream results.

Chapter 2

Importing Data

There are multiple publicly available pre-compiled microbiome data sets. These data sets begin after the bioinformatics pipeline and are matrices of counts of OTUs per sample. These data sets can exist as **phyloseq** objects, a popular R package for microbiome analysis (McMurdie and Holmes, 2013), or as separate tables of counts and metadata.

2.1 Global Patterns

The Global Patterns dataset (Caporaso et al., 2011) is an available dataset in the **phyloseq** package. These data contain samples from 25 different environmental samples and mock communities. The sampling depth of these samples averages 3.1 million total counts. We will use this dataset to work through the different normalization methods.

The following lines load the relevant packages and data.

```
library(tidyverse)
library(phyloseq)

##
## Attaching package: 'phyloseq'

## The following object is masked from 'package:SummarizedExperiment':
##
##     distance

## The following object is masked from 'package:Biobase':
##
##     sampleNames
```

```
## The following object is masked from 'package:GenomicRanges':
##
##      distance

## The following object is masked from 'package:IRanges':
##
##      distance

data("GlobalPatterns")
# examine phyloseq object
GlobalPatterns

## phyloseq-class experiment-level object
## otu_table()   OTU Table:             [ 19216 taxa and 26 samples ]
## sample_data() Sample Data:          [ 26 samples by 7 sample variables ]
## tax_table()   Taxonomy Table:        [ 19216 taxa by 7 taxonomic ranks ]
## phy_tree()    Phylogenetic Tree:     [ 19216 tips and 19215 internal nodes ]
```

2.2 Pre-processing Quality Control and Filtering

In addition to normalization, there are some steps we can perform that ideally remove technical artifacts from the sequencing process that only introduce noise.

These filtering steps commonly consist of filtering out samples with a low total read depth and filtering out taxa that are rarely abundant.

Let's create a filtered version of the Global Patterns dataset. Note that there are only 26 samples, and all have a large library size, so we will not filter out any samples here.

For taxa filtering, we will remove taxa that appear fewer than 5 times in more than half the samples.

```
# Determine which taxa to remove
filter_taxa <- genefilter_sample(GlobalPatterns,
                                filterfun_sample(function(x) x > 5),
                                A=0.5*nsamples(GlobalPatterns))
# Remove those taxa from the GlobalPatterns dataset
# Save as an object with the un-normalized counts
gp_raw <- prune_taxa(filter_taxa, GlobalPatterns)
gp_raw

## phyloseq-class experiment-level object
## otu_table()   OTU Table:             [ 219 taxa and 26 samples ]
## sample_data() Sample Data:          [ 26 samples by 7 sample variables ]
## tax_table()   Taxonomy Table:        [ 219 taxa by 7 taxonomic ranks ]
```

```
## phy_tree()    Phylogenetic Tree: [ 219 tips and 218 internal nodes ]
```

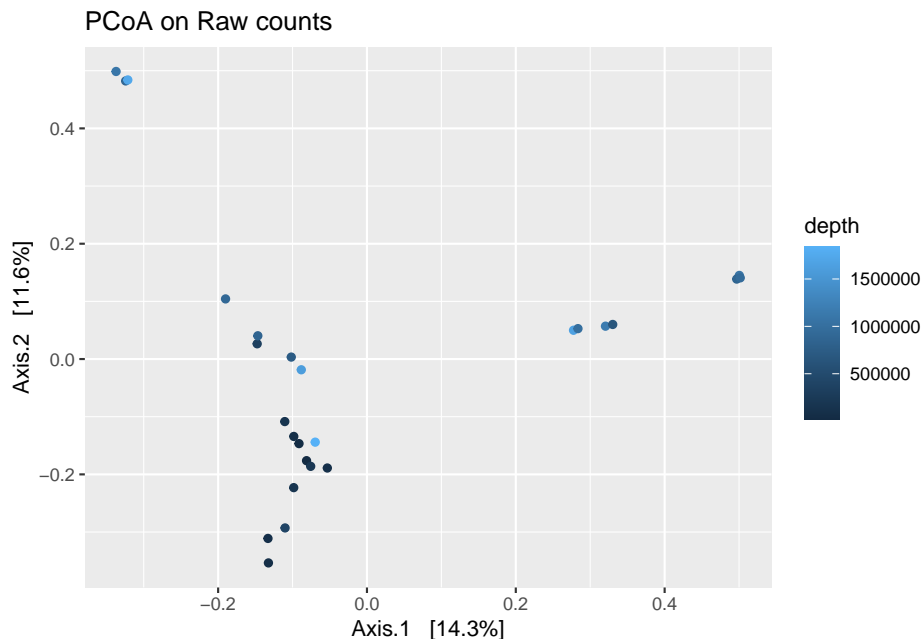
This decreases the number of taxa from 19216 to 219. This is not surprising, because this dataset contains samples from widely different locations (gut, soil, etc), and few taxa are shared among all samples and locations. One potential problem with this approach is the widely different locations, so it is possible that the remaining taxa could be some technical artifact, or could be a general ‘core’ set of taxa shared across the disparate environments.

Additionally, let us save the total sampling depth as the variable `depth` in the metadata for the Global Patterns dataset.

```
gp_raw@sam_data$depth <- sample_sums(gp_raw)
```

We can visualize technical artifacts of sapling depth is by looking at principal coordinates plots using the Bray-Curtis dissimilarity, coloring by sampling depth too see how much variation can be explained by the original sampling depth.

```
gp_raw_dist <- phyloseq::ordinate(gp_raw, "PCoA", "bray")
plot_ordination(gp_raw,
  gp_raw_dist,
  color = "depth",
  title = "PCoA on Raw counts")
```



We don't see any extreme patterns with sampling depth, but this might additionally be due to the differences in different locations might have different sampling depths. This comparison might be more interesting when we only have one location we are sampling from.

Chapter 3

Total Sum scaling (TSS)

3.1 About TSS

The first method described is Total Sum Scaling (TSS). This method is also referred to as Total Count (TC), converting into proportions, or relative abundance. This is a scaling method to normalize the different library sizes across samples. For every entry in the count matrix, we scale by the total read depth of that sample. This converts the counts into the proportion of abundance present in each given sample.

Though a more straightforward method, TSS normalization is not without its drawbacks. In microbiome data, it is common to have numerous low or zero-count observations, and that only a few most common OTUs contribute to the majority of the total sum of the sampling depth. These high-count, frequent, taxa could be an artifact of the sequencing step, where high abundance observations are preferentially sampled. Using these large counts can dominate the scaling factor for each sample. As seen below, we see that the scaling factor for each sample is completely dominated by **ASV1**, if that one taxon were not included in the sample, the scaling factor would be widely different.

Sample	ASV1	ASV2	ASV3	ASV4	TSS Scaling Factor	Scaling factor w/o ASV1
Sample A	10314	34	8	12	10368	54
Sample B	824	23	13	20	880	56

Because this method does not account for the preferential sequencing overabundance of **ASV1** it is possible to see an increase in false positives. However, this is a widely used method, and one of the few normalization methods that completely accounts for differing library sizes, which can be an important consideration depending on the analysis goal. Community level-analysis for example can be library-size dependent (ordination, some dissimilarity measures).

3.2 TSS Implementation

Here, we provide a function that will normalize a `phyloseq` object by Total Sum Scaling. We have the option of keeping the result as proportions (having values 0-1), or transforming to an equal sequencing depth so the results are counts per million.

```
norm_TSS <- function(ps, keep_prop = F){
  # keep as proportions or convert to counts per million?
  scale <- ifelse(keep_prop, 1, 1e6)
  # TSS function
  ps_normed <- phyloseq::transform_sample_counts(ps, function(x) x * scale / sum(x))
  return(ps_normed)
}
```

3.3 TSS on Global Patterns

Using the above function, we apply this normalization to the Global Patterns data.

```
gp_tss <- norm_TSS(gp_raw)
# rename the depth as the scaling factor
sample_data(gp_tss)$scaling_factor <- sample_data(gp_tss)$depth
```

To see the differences between the un-normalized, raw data, and the TSS transformed normalized data, one possible way is to look at ordination plots. Microbiome data are high dimensional, so visualization directly of the data is difficult. Here, let us examine the principal coordinates plot using the Bray-Curtis dissimilarity.

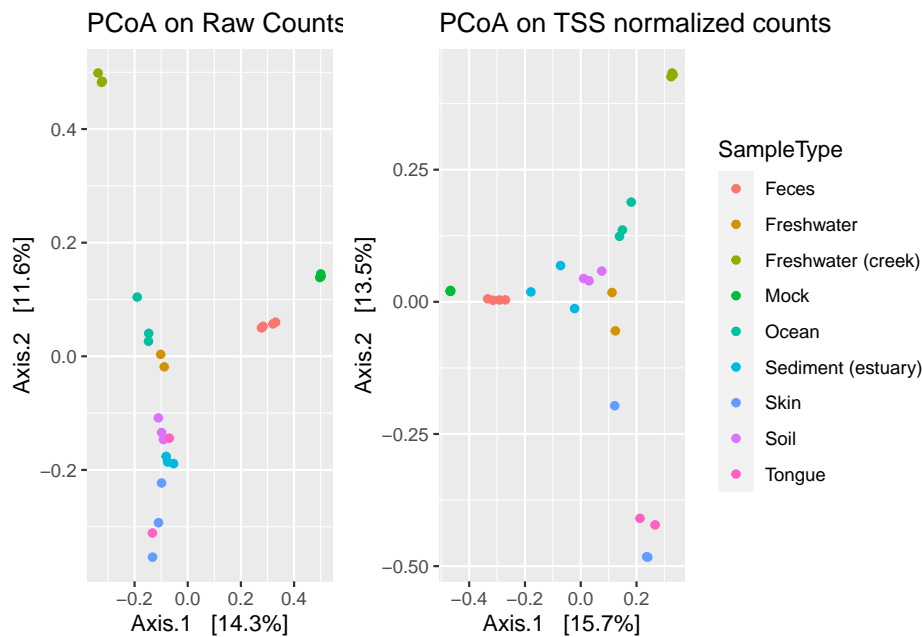
First calculate the distance matrices, using the `phyloseq` function `ordinate()`

```
gp_raw_dist <- phyloseq::ordinate(gp_raw, "PCoA", "bray")
gp_tss_dist <- phyloseq::ordinate(gp_tss, "PCoA", "bray")
```

Now plot the two ordinations. Even before normalization, the different communities are clearly clustered.

(Note to self: perhaps choose a different dataset to use for walk-through, currently using Global Patterns since it is small and quick for computations, but harder to see differences.)

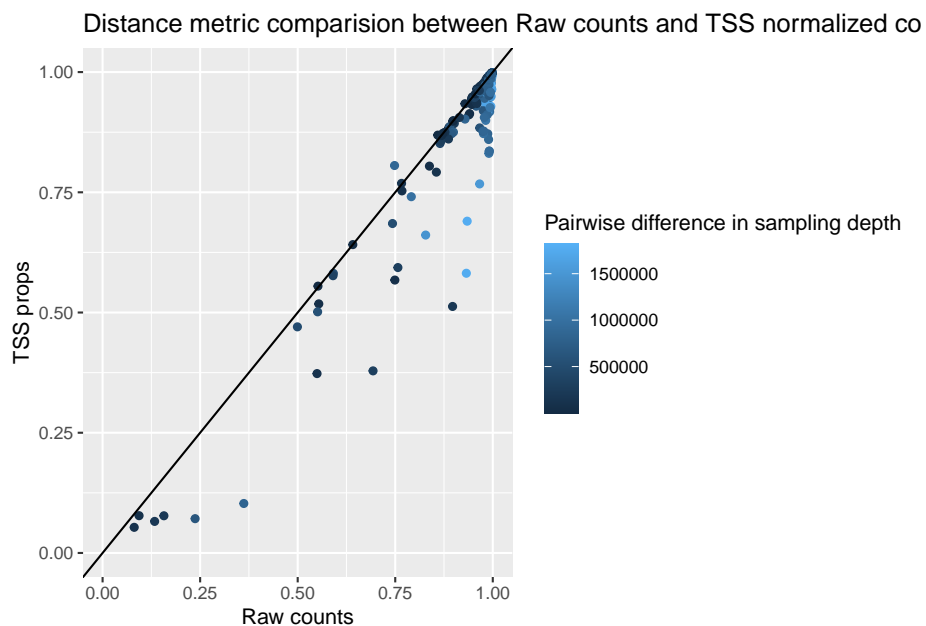
```
plot_ordination(gp_raw, gp_raw_dist, color = "SampleType",
  title = "PCoA on Raw Counts") +
plot_ordination(gp_tss, gp_tss_dist, color = "SampleType",
  title = "PCoA on TSS normalized counts") +
plot_layout(guides = 'collect')
```



We can also compare the values of the distance matrices before and after normalization to see how the normalization method is impacting different types of points.

```
# Function to visualize potential differences and changes after normalization methods
plot_norm_changes <- function(data_norm, data_raw, dist_method = "bray", x_lab = "raw", y_lab = "norm") {
  # calculate the Bray-Cutris distance matrix for the raw data, the normalized data,
  # and calculate the pairwise difference between the original library sizes between samples
  plot <- data.frame(raw = as.numeric(phyloseq::distance(data_raw, dist_method)),
                    norm = as.numeric(phyloseq::distance(data_norm, dist_method)),
                    diff = as.numeric(dist(get_variable(data_raw, "depth")))) %>%
    ggplot(aes(x = raw, y = norm, color = diff)) +
      geom_point() +
      geom_abline() +
      ggtitle(title) +
      xlab(x_lab) + ylab(y_lab) +
      labs(color = "Pairwise difference in sampling depth") +
      xlim(c(0,1)) + ylim(c(0,1))
  return(plot)
}

plot_norm_changes(gp_tss, gp_raw,
  x_lab = "Raw counts", y_lab = "TSS props",
  title = "Distance metric comparison between Raw counts and TSS normalized counts")
```



Points below the line are pairs of samples that are marked as more similar after normalization. Points above the line are marked as more different after normalization. Values closer to 1 are 'more different'. Unsurprisingly, pairs that had larger original differences in sampling depth were marked as more different on the raw, un-normalized data, and became marked as more similar after TSS normalization.

Chapter 4

Rarefying

4.1 About Rarefying

Rarefying is another common normalization technique that standardizes the library size across samples that was originally used in ecology. This method standardizes the read depth across all samples. To perform this method we first choose a minimum library size. Looking at rarefaction/collectors curves, or using a certain percentile can guide choosing this cutoff. Then all samples that have a read depth below this cutoff are discarded. Thus this method has a built-in filtering step. Next, we sample without replacement of the size of the chosen cutoff. It can be a standalone method or combined with other methods and transformations.

This is a very commonly used method, but it has also been criticized (McMurdie and Holmes, 2014). First of all, it throws away valid data, and this results in a loss of power and an increase in false positives. Rare taxa can be removed in this approach too. It is however encouraged when we have widely different library sizes as it can lower the false discovery rate (Weiss et al., 2017), and has also been shown to perform well in community-level analysis (McKnight et al., 2019), as it completely standardizes the read count depth, and some methods are sensitive to differences in read count. Rarefying has been shown to separate by biological signal in ordination methods based on presence/absence.

4.2 Rarefying implementation

The following function returns a rarefied `phyloseq` object. We can either pass in the minimum sampling depth as a second argument, or use the default minimum depth of the samples.

```
norm_rarefy <- function(phyloseq, depth = min(sample_sums(phyloseq))) {
  return(phyloseq::rarefy_even_depth(phyloseq, sample.size = depth))
}
```

4.3 Rarefying on Global Patterns

We use the above function to rarefy the Global Patterns data. The first difficulty is choosing a minimum sampling depth. The Global Patterns dataset already has a very high sampling depth for all samples, so we will chose the lowest as the minimum depth to rarefy to. Since we chose the minimum sampling depth, no samples have been dropped. In data sets where we have low sampling depth there is a balance between how many samples to drop and how low to set the minimum depth to.

```
gp_rare <- norm_rarefy(gp_raw)
```

```
## You set `rngseed` to FALSE. Make sure you've set & recorded
## the random seed of your session for reproducibility.
## See `?set.seed`

## ...
```

We can check that indeed all samples now have the same sampling depth, which is 15905. Note that the highest sampling depth in this dataset was almost 2 million, so we have discarded a lot of data to reduce to 15905.

```
max(sample_sums(gp_raw))
```

```
## [1] 1842380
```

```
sample_sums(gp_rare)
```

```
##      CL3      CC1      SV1  M31Fcsw  M11Fcsw  M31Plmr  M11Plmr  F21Plmr
##  15905  15905  15905   15905   15905   15905   15905   15905
## M31Tong M11Tong LMEpi24M SLEpi20M  AQC1cm  AQC4cm  AQC7cm    NP2
##  15905  15905  15905   15905   15905   15905   15905   15905
##      NP3      NP5  TRRsed1  TRRsed2  TRRsed3    TS28    TS29  Even1
##  15905  15905  15905   15905   15905   15905   15905   15905
##  Even2  Even3
##  15905  15905
```

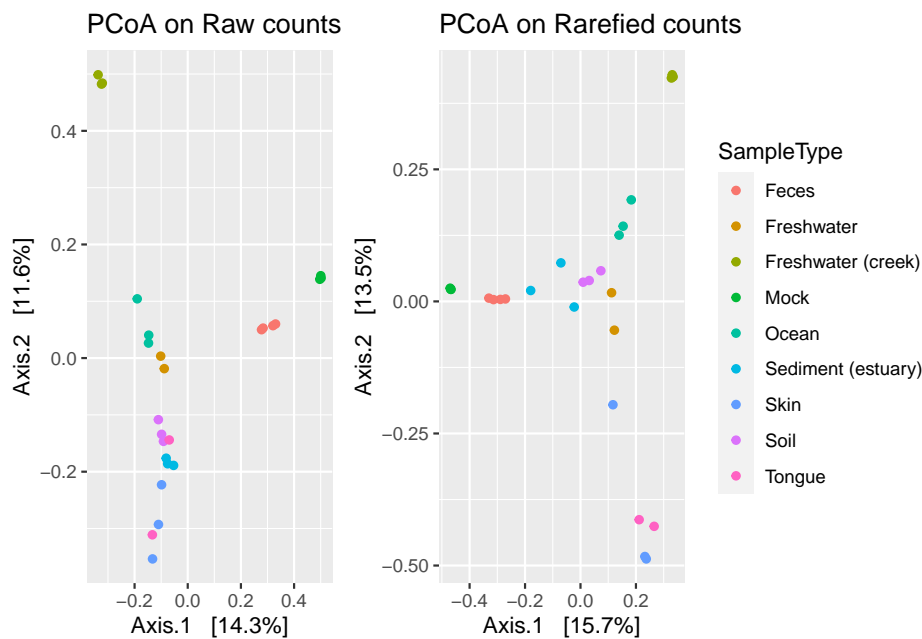
We can again compare the PCoA plots between rarefied and raw counts, coloring by sample type to view clusters.

```
plot_ordination(gp_raw,
  phyloseq::ordinate(gp_raw, "PCoA", "bray"),
  color = "SampleType",
```

```

    title = "PCoA on Raw counts") +
plot_ordination(gp_rare,
  phyloseq::ordinate(gp_rare, "PCoA", "bray"),
  color = "SampleType",
  title = "PCoA on Rarefied counts")+
plot_layout(guides = 'collect')

```

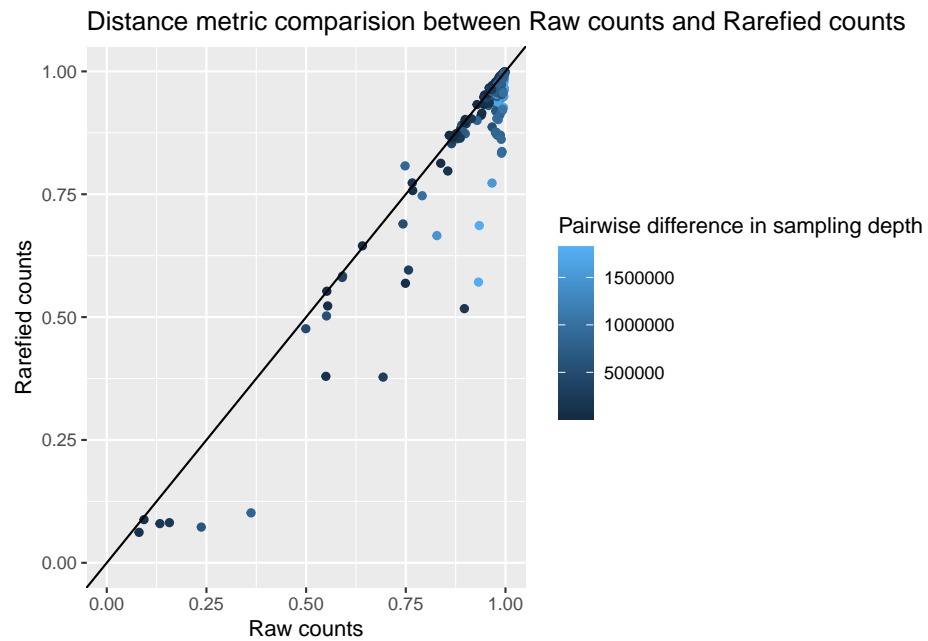


Now examine how the distance matrices change before/after normalization. We see a similar pattern to TSS when distance matrices calculated from rarefied counts are compared to those from the raw counts.

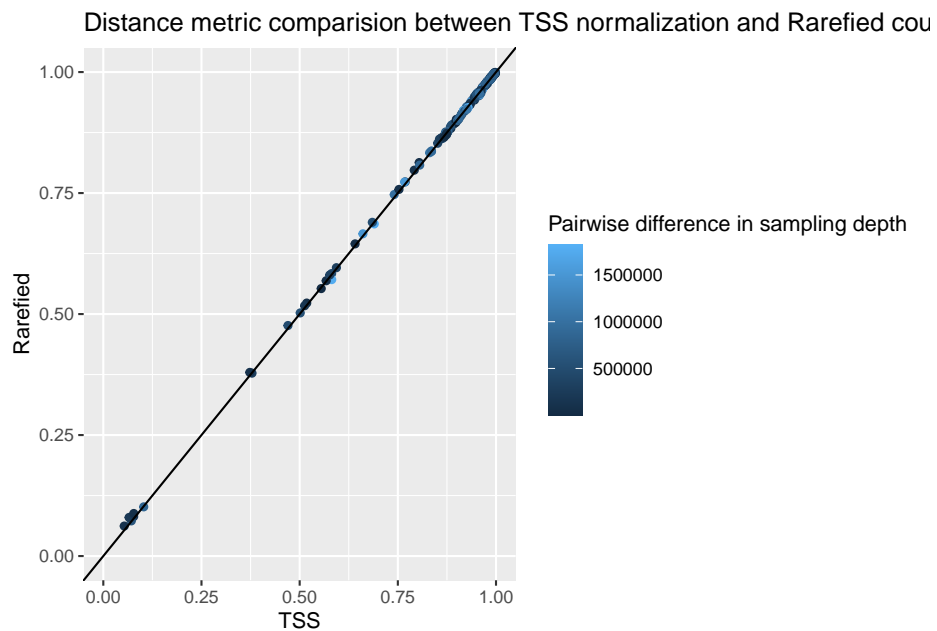
```

# Identify any samples filtered in rarefying process
rare_samples <- sample_names(gp_rare)
gp_raw_match <- prune_samples(rare_samples, gp_raw)
plot_norm_changes(gp_rare, gp_raw_match,
  x_lab = "Raw counts", y_lab = "Rarefied counts",
  title = "Distance metric comparison between Raw counts and Rarefied counts ")

```



```
## Compare to tss
gp_tss_match <- norm_TSS(gp_raw_match)
plot_norm_changes(gp_rare, gp_tss,
  x_lab = "TSS", y_lab = "Rarefied",
  title = "Distance metric comparison between TSS normalization and Rarefied counts")
```



Chapter 5

DESeq

5.1 About DESeq

DESeq2 includes a normalization method for adjusting for differing library sizes across samples (Anders and Huber, 2010). This method also can account for differences in library composition. This method has also been called MED, RLE, or DESeq in the literature.

DESeq2 first takes the natural logarithm of every entry in the count matrix. Due to this, all entries with zero will be set to negative infinity. Next, the row average is calculated (geometric average), so we have a vector of average counts for each taxon. Taking the log first should avoid undue influence by extreme outliers. All taxa with an average of infinity are removed. This step will remove all taxa with zero read count in one or more samples. This can be a problem in microbiome data. Next, we subtract the average log value from the $\log(\text{counts})$, this gives a log ratio. This is equivalent to the ratio of the reads in each sample to the average across all samples. Next, we calculate the median of the log-ratios for each sample. These medians are converted to scaling factors for each sample by exponentiation. An extension of this method, denoted ‘poscounts’, has been suggested, which instead of taking the geometric mean of the logged counts for each taxon, we take the n -th root of the product of the non-zero counts.

This method assumes that the taxon of median absolute abundance is not differentially abundant, which is more likely true for the RNA-Seq it was developed for, but may not be true for microbiome studies, especially when there are more study groups, or we are analyzing higher taxonomic levels.

An additional option can be used to perform a variance stabilizing transformation on the count matrix before normalizing with the above size factors. This method calculates a dispersion-mean relationship and then transforms the data. The result ideally is an abundance matrix that is approximately homoskedastic

or constant variance across the range of mean values. The package also includes an option for a ‘rlog’ transform, which they recommend over the variance stabilizing method in the case when there is a large difference in library sizes.

If differential abundance is of interest to calculate, DESeq uses a negative binomial distribution to model differential abundances. It is possible to provide the size factors calculated by another method to DESeq to perform differential analysis.

5.2 DESeq Implementation

Here we provide two normalization functions using DESeq methods. The first calculates the RLE normalization using the `poscounts` option for microbiome data. The second calculates the variance stabilizing transformation.

```
norm_DESeq_RLE_poscounts <- function(ps, group = 1){
  require(DESeq2, quietly = T)
  # keep arbitrary design for normalization
  # Convert to DESeq object
  ps_dds <- phyloseq_to_deseq2(ps, ~1)
  # Calculate the size factors (scaling)
  ps_dds <- estimateSizeFactors(ps_dds, type = "poscounts")
  # Extract counts
  counts <- DESeq2::counts(ps_dds, normalized = T)
  # Convert back to phyloseq
  otu <- otu_table(counts, taxa_are_rows = T)
  sam <- access(ps, "sam_data")
  sam$scaling_factor <- sizeFactors(ps_dds)
  tax <- access(ps, "tax_table")
  phy <- access(ps, "phy_tree")
  ps_DESeq <- phyloseq(otu, sam, tax, phy)
  return(ps_DESeq)
}
```

```
norm_DESeq_vs <- function(ps, group = 1){
  require(DESeq2, quietly = T)
  ps_dds <- phyloseq_to_deseq2(ps, ~ 1)
  ps_dds <- estimateSizeFactors(ps_dds, type = "poscounts")
  # Variance transformation
  ps_dds <- estimateDispersions(ps_dds)
  abund <- getVarianceStabilizedData(ps_dds)
  # don't allow deseq to return negative counts
  # add the minimum count to all values
```

```

# another option is to replace negative counts with 0
abund <- abund + abs(min(abund))
otu <- otu_table(abund, taxa_are_rows = T)
sam <- access(ps, "sam_data")
tax <- access(ps, "tax_table")
phy <- access(ps, "phy_tree")
ps_DESeq <- phyloseq(otu,sam,tax,phy)
return(ps_DESeq)
}

```

5.3 DESeq on Global Patterns

Perform DESeq RLE normalization as well as DESeq variance stabilized transformation on Global Patterns:

```
gp_deseq_rle <- norm_DESeq_RLE_poscounts(gp_raw)
```

```
## converting counts to integer mode
```

```
gp_deseq_vs <- norm_DESeq_vs(gp_raw)
```

```
## converting counts to integer mode
```

```
## gene-wise dispersion estimates
```

```
## mean-dispersion relationship
```

```
## -- note: fitType='parametric', but the dispersion trend was not well captured by the
##      function: y = a/x + b, and a local regression fit was automatically substituted.
##      specify fitType='local' or 'mean' to avoid this message next time.
```

```
## final dispersion estimates
```

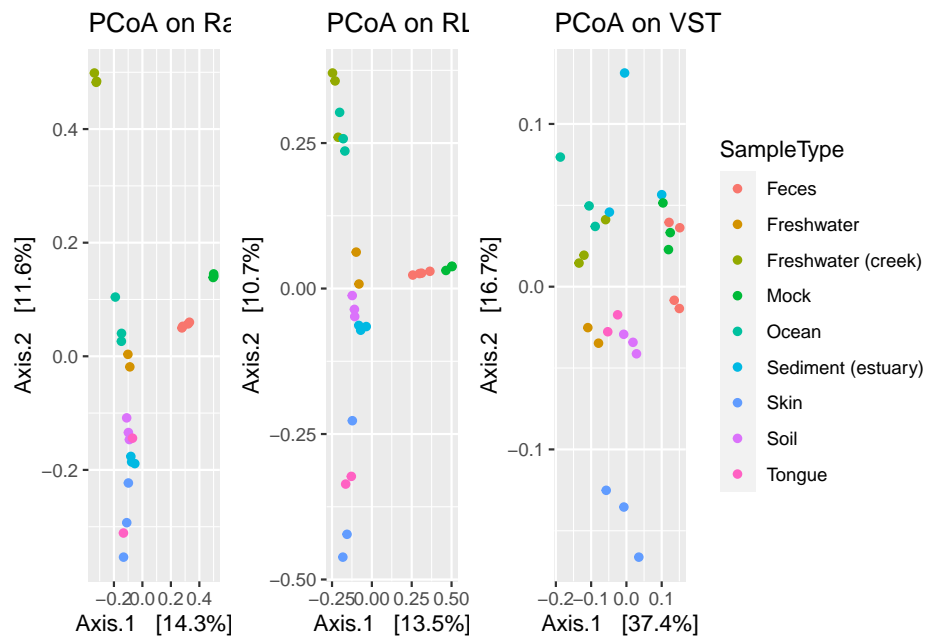
Examine principal coordinate plots between raw data and both DESeq normalized data.

```

# First calculate distance matrices
gp_rle_dist <- phyloseq::ordinate(gp_deseq_rle, "PCoA", "bray")
gp_vs_dist <- phyloseq::ordinate(gp_deseq_vs, "PCoA", "bray")

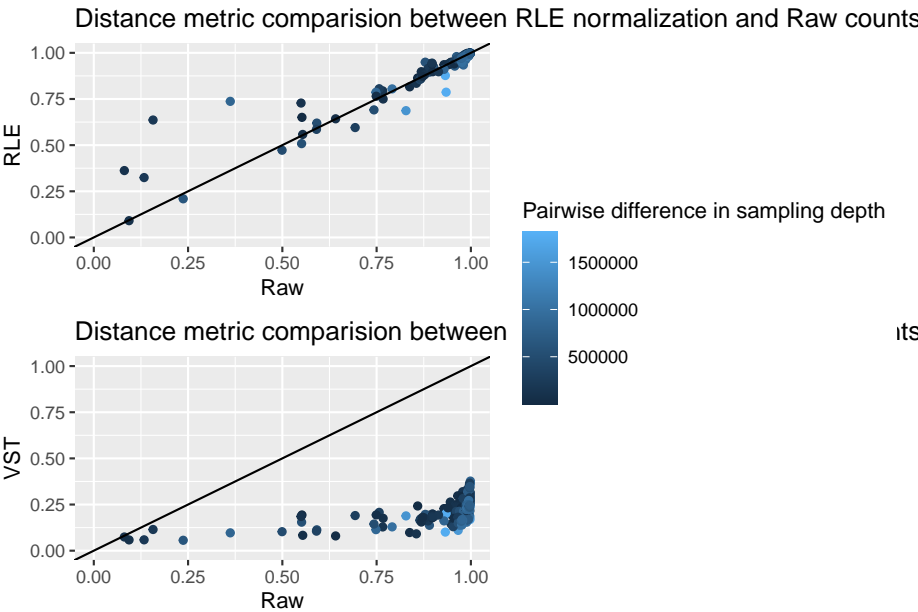
# Plot ordinations
plot_ordination(gp_raw, gp_raw_dist, color = "SampleType", title = "PCoA on Raw data") +
plot_ordination(gp_deseq_rle, gp_rle_dist, color = "SampleType", title = "PCoA on RLE") +
plot_ordination(gp_deseq_vs, gp_vs_dist, color = "SampleType", title = "PCoA on VST") +
plot_layout(guides = 'collect')

```



See how dissimilarity matrices differ from raw counts and each DESeq transformation.

```
plot_norm_changes(gp_deseq_rle, gp_raw,
  x_lab = "Raw", y_lab = "RLE",
  title = "Distance metric comparison between RLE normalization and Raw",
plot_norm_changes(gp_deseq_vs, gp_raw,
  x_lab = "Raw", y_lab = "VST",
  title = "Distance metric comparison between VST normalization and Raw",
plot_layout(guides = 'collect')
```

Chapter 6

TMM (edgeR)

TMM (Trimmed median of m-values) is another method borrowed from RNA-Seq analysis, and implemented in **edgeR** (Robinson and Oshlack, 2010). This method uses, or calculates a reference sample, and compares all other samples to the reference sample. The size factor is the mean of the log-ratios after excluding the highest count taxa and taxa with the largest fold change. As taxa with zero counts are excluded, a pseudo count is needed. Additionally, there is the `TMMwsp` option which is encouraged as it is more robust to zero counts. Positive counts are reused to increase the number of features when we compared. The singleton positive counts are paired up in decreasing order of size and then a modified TMM method is applied to the re-ordered libraries.

6.1 EdgeR TMM implementation

```
norm_TMM <- function(physeq, group = 1, method="TMM", pseudocount = 1, ...){
  require("edgeR", quietly = T)
  require("phyloseq", quietly = T)
  # Enforce orientation.
  if( !taxa_are_rows(physeq) ){ physeq <- t(physeq) }
  x = as(otu_table(physeq), "matrix")
  # Add one to protect against overflow, log(0) issues.
  x = x + pseudocount
  # Check `group` argument
  if( identical(all.equal(length(group), 1), TRUE) & nsamples(physeq) > 1 ){
    # Assume that group was a sample variable name (must be categorical)
    group = get_variable(physeq, group)
  }
  # Define gene annotations (`genes`) as tax_table
```

```

taxonomy = tax_table(physeq, errorIfNULL=FALSE)
if( !is.null(taxonomy) ){
  taxonomy = data.frame(as(taxonomy, "matrix"))
}
# Now turn into a DGEList
y = DGEList(counts=x, group=group, genes=taxonomy, remove.zeros = TRUE)
# Calculate the normalization factors
d = edgeR::calcNormFactors(y, method=method)
# Check for division by zero inside `calcNormFactors`
if( !all(is.finite(d$samples$norm.factors)) ){
  stop("Something wrong with edgeR::calcNormFactors on this data,
    non-finite $norm.factors, consider changing `method` argument")
}
scalingFactor <- d$samples$norm.factors * d$samples$lib.size / 1e6
dataNormalized <- t(t(otu_table(physeq)) / scalingFactor)
#dataNormalized <- cpm(d)
otu <- otu_table(dataNormalized, taxa_are_rows = T)
sam <- access(physeq, "sam_data")
sam$scaling_factor <- scalingFactor
tax <- access(physeq, "tax_table")
phy <- access(physeq, "phy_tree")
ps_edgeR <- phyloseq(otu,sam,tax,phy)

return(ps_edgeR)
}

```

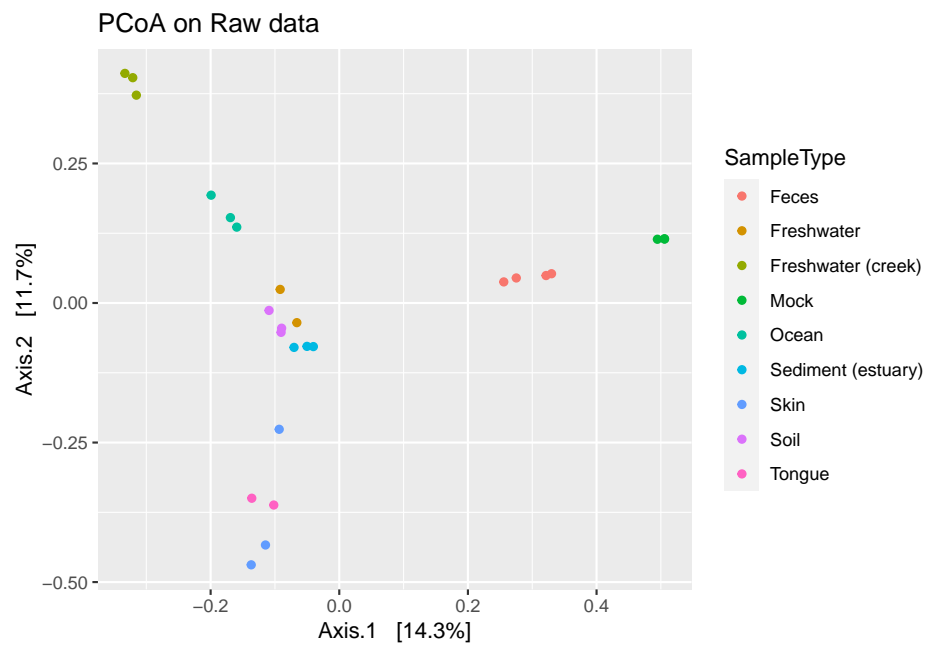
6.2 TMM on Global Patterns

Perform normalization:

```
gp_tmm <- norm_TMM(gp_raw)
```

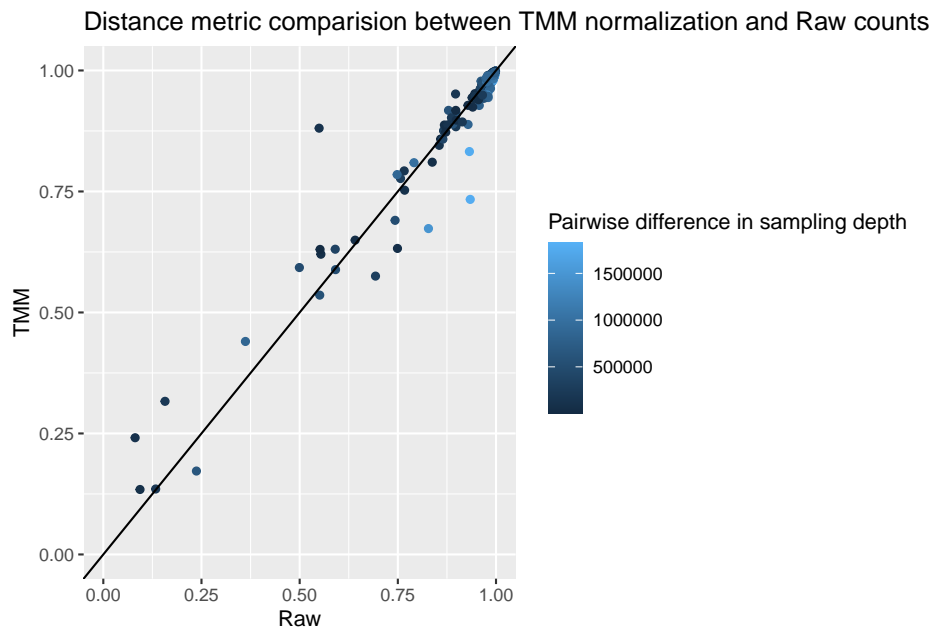
View PCoA plots

```
plot_ordination(gp_tmm, phyloseq::ordinate(gp_tmm, "PCoA", "bray") , color = "SampleType")
```



View how TMM normalization changes distance metrics differently than raw counts.

```
plot_norm_changes(gp_tmm, gp_raw,  
                  x_lab = "Raw", y_lab = "TMM",  
                  title = "Distance metric comparision between TMM normalization and Raw counts ")
```



Chapter 7

Cumulative sum scaling (CSS)

Cumulative sum scaling is a scaling normalization method, developed for marker gene sequencing. It is intended to account for undersampling and correct biases from preferentially amplified features in a sample-specific manner (Paulson et al., 2013). This method assumes that count distributions should be roughly equivalent and independent to each other up to the given quantile which is chosen as the smallest value at which instability is found. This method is an extension to UQ scaling where a quantile is specified. If there is high count variability the assumption may not be met. This is not a method that accounts for compositionality. CSS normalization initially showed improvements in separating samples biologically in ordination, it was shown to be an artifact of unequal application of log transformation across methods (Costea et al., 2014).

7.1 CSS implementation

```
norm_CSS <- function(ps){  
  require(metagenomeSeq)  
  # Convert to metagenomeSeq data type  
  ps.metaG<-phyloseq_to_metagenomeSeq(ps)  
  p_stat = cumNormStatFast(ps.metaG)  
  ps.metaG = cumNorm(ps.metaG, p = p_stat)  
  ps.metaG.norm <- MRcounts(ps.metaG, norm = T)  
  # Convert back to phyloseq with normalized counts  
  otu <- otu_table(ps.metaG.norm, taxa_are_rows = T)  
  sam <- access(ps, "sam_data")  
}
```

```

sam$scaling_factor <- normFactors(ps.metaG)/1e6
tax <- access(ps, "tax_table")
phy <- access(ps, "phy_tree")
ps_CSS <- phyloseq(otu,sam,tax,phy)
return(ps_CSS)
}

```

7.2 CSS on Global Patterns

Perform CSS normalization:

```
gp_css <- norm_CSS(gp_raw)
```

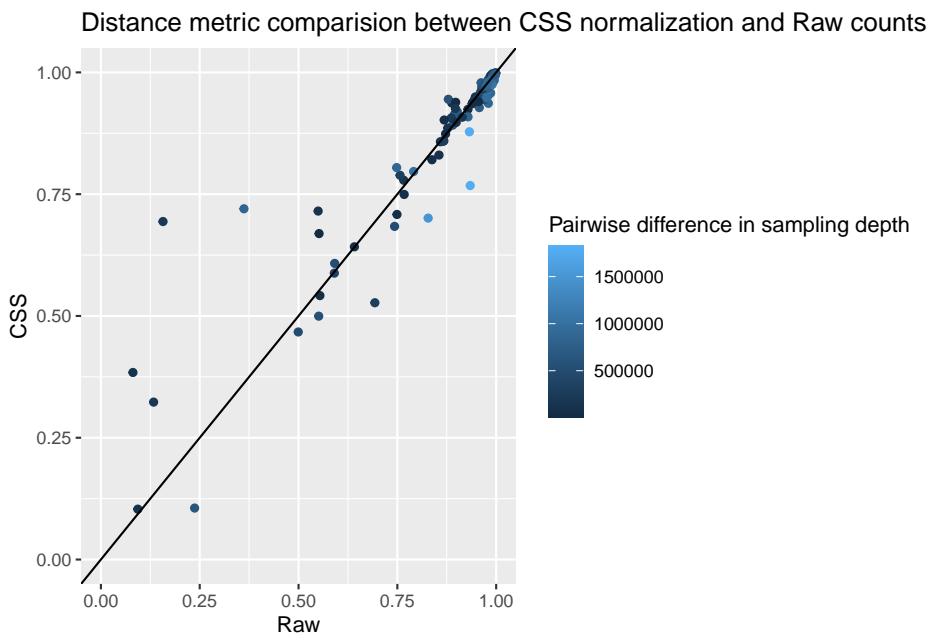
Default value being used.

View how TMM normalization changes distance metrics differently than raw counts.

```

plot_norm_changes(gp_css, gp_raw,
  x_lab = "Raw", y_lab = "CSS",
  title = "Distance metric comparision between CSS normalization and Raw

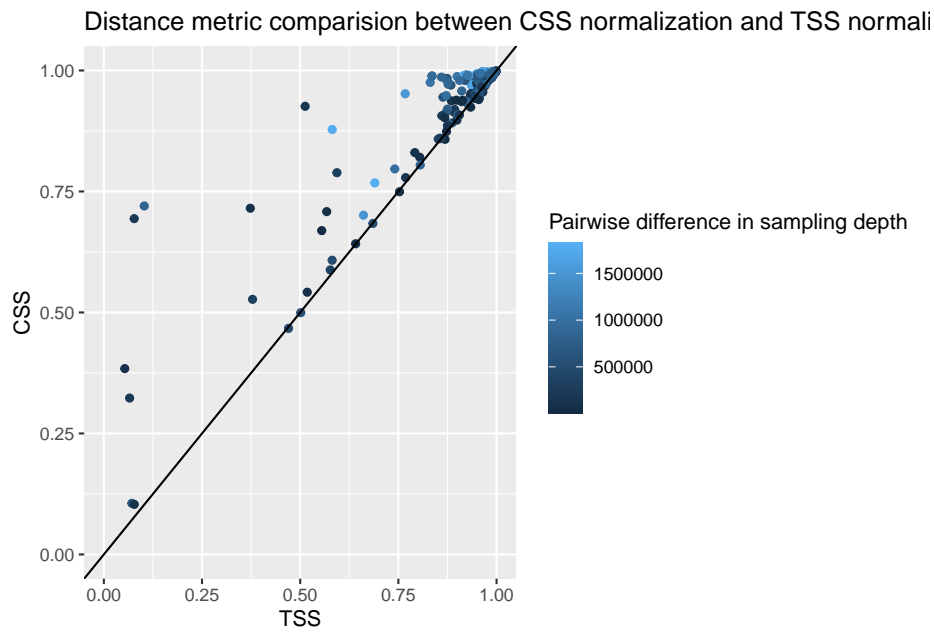
```



```

plot_norm_changes(gp_css, gp_tss,
  x_lab = "TSS", y_lab = "CSS",
  title = "Distance metric comparision between CSS normalization and T

```

CSS normalization appears to consider pairs as more different than TSS normalization, and pairs with high sequencing depth differences even more so.

Chapter 8

GMPR

A recent extension of the RLE DESeq method is the Geometric mean of Pairwise ratios (GMPR) approach (Chen et al., 2018). This method reverses the steps of RLE, and instead calculates the median count ratio of the non-zero counts between pairs of samples as although only a small number of taxa are likely to be shared for every sample, it is more likely that there are many shared taxa between pairs. It then uses the pairwise results to calculate the size factor for each sample. This method has slow computation, but is robust to differential and outlier taxa. It addresses sparsity, but not composition. The size factors can be inputted to DESeq and a VST transformation applied additionally. It is a newer method, and has unfortunately not been included in many benchmarking studies, although initial results show it to be more powerful than DESeq, not surprisingly, as it uses more data, as zero counts do not need to be discarded. It assumes there is a large invariant portion of the count data, similar to other methods.

8.1 GMPR Implementation

```
norm_GMPR <- function(ps, scale = 1e6){
  require(GMPR, quietly = T)
  # Calculate GMPR size factor
  # Row - features, column - samples
  otu <- as(otu_table(ps), "matrix")
  if(taxa_are_rows(ps)){otu <- t(otu)}
  otu_df = as.data.frame(otu)
  otu.tab <- matrix(otu, ncol = ncol(otu))
  gmpR.size.factor <- GMPR::GMPR(otu_df, intersect_no = 4, min_ct = 2)
```

```

# normalize
otu.tab.norm <- t(t(otu) / (gmpr.size.factor/scale))

# convert back to PS
sam <- access(ps, "sam_data")
sam$scaling_factor <- gmpr.size.factor
tax <- access(ps, "tax_table")
phy <- access(ps, "phy_tree")
ps_GMPR <- phyloseq(otu_table(otu.tab.norm, taxa_are_rows = F), sam, tax, phy)
return(ps_GMPR)
}

```

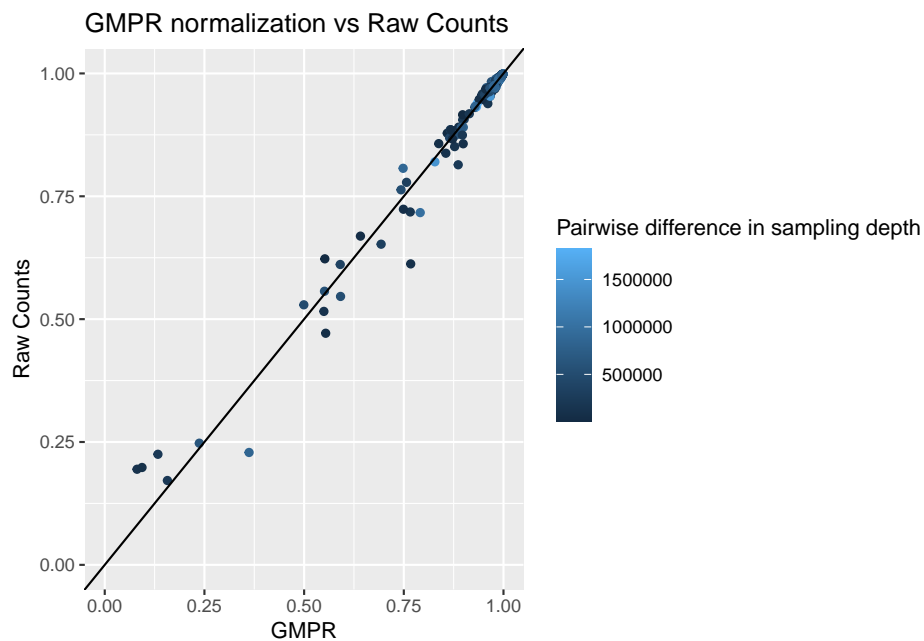
8.2 Global Patterns GMPR

```

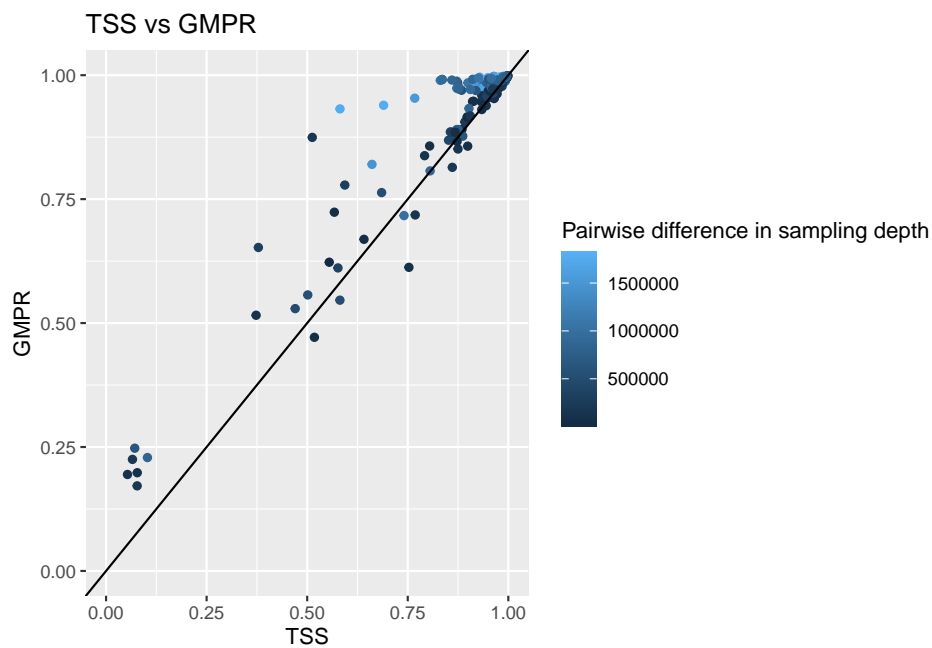
# Normalize
gp_gmpr <- norm_GMPR(gp_raw)

# Compare to other methods like before
plot_norm_changes(gp_gmpr, gp_raw,
  x_lab = "GMPR", y_lab = "Raw Counts",
  title = "GMPR normalization vs Raw Counts")

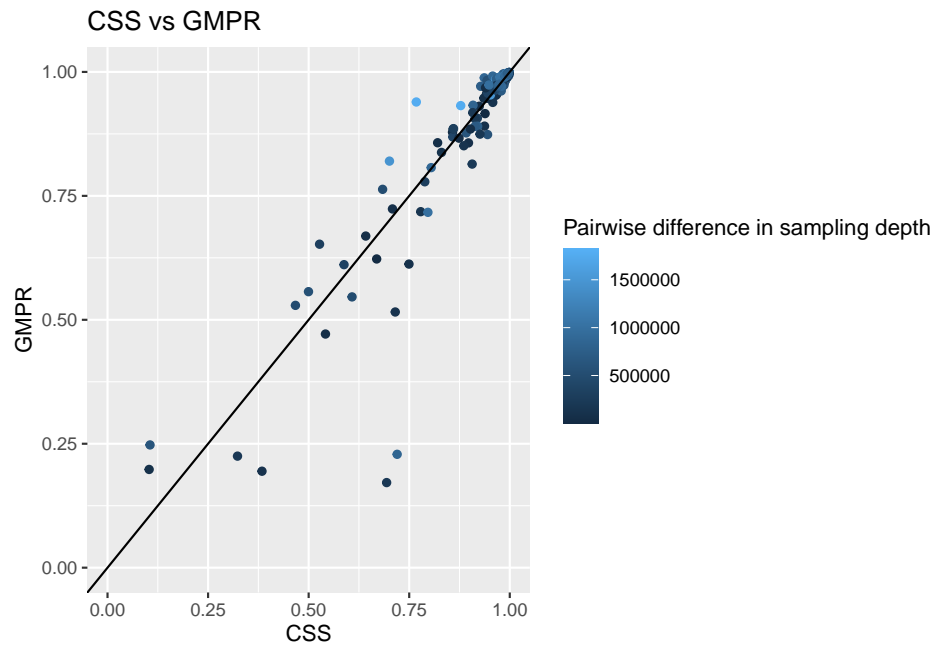
```



```
plot_norm_changes(gp_gmpr, gp_tss,  
                  x_lab = "TSS", y_lab = "GMPR",  
                  title = "TSS vs GMPR")
```



```
plot_norm_changes(gp_gmpr, gp_css,  
                  x_lab = "CSS", y_lab = "GMPR",  
                  title = "CSS vs GMPR" )
```



Chapter 9

Wrench

Wrench is a recent normalization method developed for microbiome data (Kumar et al., 2018). This method includes compositional bias correction for sparse datasets. This method uses a hurdle log-normal distribution to estimate the normalization factors (the location estimate for the group). For this method, we assume abundances are drawn from a hurdle Log-Gaussian distribution, and the scaling factor used is essentially the location estimate for the group.

9.1 Wrench Implementation

```
norm_wrench <- function(ps, group_col){
  require(Wrench, quietly = T)
  if( identical(all.equal(length(group_col), 1), TRUE) & nsamples(ps) > 1 ){
    # Assume that group was a sample variable name (must be categorical)
    group = get_variable(ps, group_col)
  }
  otu_tab <- otu_table(ps)
  W <- wrench(otu_tab, group)

  compositionalFactors <- W$ccf
  normalizationFactors <- W$nf

  normed_otu <- otu_tab/(normalizationFactors/1e6)
  otu <- otu_table(normed_otu, taxa_are_rows = T)
  sam <- access(ps, "sam_data")
  sam$scaling_factor <- normalizationFactors
  tax <- access(ps, "tax_table")
  phy <- access(ps, "phy_tree")
}
```

```
ps_wrench <- phyloseq(otu,sam,tax,phy)

  return(ps_wrench)
}
```

9.2 Wrench on Global Patterns

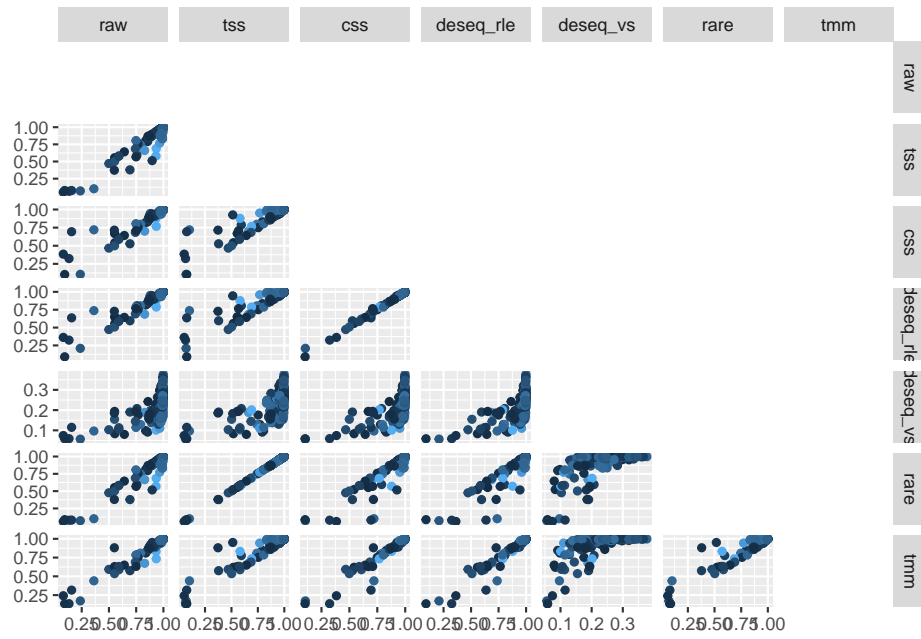
```
gp_wrench <- norm_wrench(gp_raw, group_col = "SampleType")
```


Chapter 10

Comparisons

```
library(GGally)
```

```
## Registered S3 method overwritten by 'GGally':  
##   method from  
##   +.gg      ggplot2  
data.frame(raw = as.numeric(phyloseq::distance(gp_raw, "bray")),  
            tss = as.numeric(phyloseq::distance(gp_tss, "bray")),  
            css = as.numeric(phyloseq::distance(gp_css, "bray")),  
            deseq_rle = as.numeric(phyloseq::distance(gp_deseq_rle, "bray")),  
            deseq_vs = as.numeric(phyloseq::distance(gp_deseq_vs, "bray")),  
            rare = as.numeric(phyloseq::distance(gp_rare, "bray")),  
            tmm = as.numeric(phyloseq::distance(gp_tmm, "bray")),  
            diff = as.numeric(dist(get_variable(gp_raw, "depth")))) %>%  
  ggpairs(columns = 1:7, upper = "blank",  
          diag = "blank", ggplot2::aes(colour=diff))
```



Bibliography

- Aitchison, J. (1982). The statistical analysis of compositional data. *Journal of the Royal Statistical Society: Series B (Methodological)*, 44(2):139–160. [_eprint: https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2517-6161.1982.tb01195.x](https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2517-6161.1982.tb01195.x).
- Anders, S. and Huber, W. (2010). Differential expression analysis for sequence count data. *Genome Biology*, 11(10):R106.
- Caporaso, J., Lauber, C., Walters, W., Berg-Lyons, D., Lozupone, C., Turnbaugh, P., Fierer, N., and Knight, R. (2011). Global patterns of 16s rna diversity at a depth of millions of sequences per sample. *Proceedings of the National Academy of Sciences*, 108(Supplement 1):4516–4522. Publisher: National Academy of Sciences Section: Colloquium Paper PMID: 20534432.
- Chen, L., Reeve, J., Zhang, L., Huang, S., Wang, X., and Chen, J. (2018). Gmpr: A robust normalization method for zero-inflated count data with application to microbiome sequencing data. *PeerJ*, 6:e4600. Publisher: PeerJ Inc.
- Costea, P., Zeller, G., Sunagawa, S., and Bork, P. (2014). A fair comparison. *Nature Methods*, 11(4):359–359. Bandiera_abtest: a Cg_type: Nature Research Journals Number: 4 Primary_atype: Correspondence Publisher: Nature Publishing Group Subject_term: Data mining;Metagenomics;Statistical methods Subject_term_id: data-mining;metagenomics;statistical-methods.
- Kumar, M., Slud, E., Okrah, K., Hicks, S., Hannenhalli, S., and Corrada Bravo, H. (2018). Analysis and correction of compositional bias in sparse sequencing count data. *BMC Genomics*, 19(1):799.
- McKnight, D., Huerlimann, R., Bower, D., Schwarzkopf, L., Alford, R., and Zenger, K. (2019). Methods for normalizing microbiome data: An ecological perspective. *Methods in Ecology and Evolution*, 10(3):389–400. [_eprint: https://besjournals.onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.13115](https://besjournals.onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.13115).
- McMurdie, P. and Holmes, S. (2013). phyloseq: An r package for reproducible interactive analysis and graphics of microbiome census data. *PLOS ONE*, 8(4):e61217. Publisher: Public Library of Science.

- McMurdie, P. and Holmes, S. (2014). Waste not, want not: Why rarefying microbiome data is inadmissible. *PLOS Computational Biology*, 10(4):e1003531. Publisher: Public Library of Science.
- Paulson, J., Stine, O., Bravo, H., and Pop, M. (2013). Differential abundance analysis for microbial marker-gene surveys. *Nature Methods*, 10(12):1200–1202. Bandiera_abtest: a Cg_type: Nature Research Journals Number: 12 Primary_atype: Research Publisher: Nature Publishing Group Subject_term: Data mining;Metagenomics;Statistical methods Subject_term_id: data-mining;metagenomics;statistical-methods.
- Robinson, M. and Oshlack, A. (2010). A scaling normalization method for differential expression analysis of rna-seq data. *Genome Biology*, 11(3):R25.
- Salter, S., Cox, M., Turek, E., Calus, S., Cookson, W., Moffatt, M., Turner, P., Parkhill, J., Loman, N., and Walker, A. (2014). Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. *BMC Biology*, 12(1):87.
- Weiss, S., Xu, Z., Peddada, S., Amir, A., Bittinger, K., Gonzalez, A., Lozupone, C., Zaneveld, J., Vázquez-Baeza, Y., Birmingham, A., Hyde, E., and Knight, R. (2017). Normalization and microbial differential abundance strategies depend upon data characteristics. *Microbiome*, 5(1):27.