```r
#Additional analysis scripts for the manuscript: "Gut microbiome composition is
predictive of incident type 2 diabetes in a population cohort of 5 572 Finnish adults"
by Ruuskanen & Erawijantari et al.
#In this analysis code, participants with type 2 diabetes diagnosis within two years
from baseline have been removed from the data
#Due to sensitive health information, the data in this study are available based on a
written application to the THL Biobank as instructed in:
https://thl.fi/en/web/thl-biobank/for-researchers

if (!requireNamespace("BiocManager")) {
  install.packages("BiocManager")
}

#Use development version of ComplexHeatmap, 2.7.11<
#library(devtools)
#install_github("jokergoo/ComplexHeatmap")
#devtools::install_github("slowkow/ggrepel")

packages <- c("ggplot2", "biomformat", "ggthemes", "phyloseq", "vegan", "uwot",
"patchwork", "microbiome", "tidyverse", "reshape2", "survival", "magrittr", "ggnewscale"
, "propr", "ComplexHeatmap", "maptree", "RColorBrewer", "rms", "viridis", "scales",
"data.table")


is.installed <- function(pkg) {
  new.pkg <- pkg[!(pkg %in% installed.packages()[, "Package"])]
  if (length(new.pkg)) {
    BiocManager::install(new.pkg, ask=F)
  }
  sapply(pkg, require, character.only = TRUE)
}
is.installed(packages)

wideScreen <- function(howWide=Sys.getenv("COLUMNS")) {
  options(width=as.integer(howWide))
}
wideScreen()

theme_set(theme_tufte(base_family = "sans", base_size = 18) + theme(panel.border =
element_rect(colour = "black", fill = NA), axis.text = element_text(colour = "black",
size = 18)))


#All data are included in the THL Biobank release package.
#Phenotype data is loaded from the included R object
load("FR_02_phenotype_data.RData")
#Subset to data which includes the fecal samples
FR02 <- FR02[!is.na(FR02$Barcode),]
row.names(FR02) <- FR02$Barcode
#Construct objects with NCBI data from SHOGUN
if (file.exists("microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data.RDs"))
 {
  ncbi_data <- readRDS(
  "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data.RDs")
} else {
  #Construct the primary phyloseq object and subset to FR02 samples.
  ncbi_data <- biomformat::read_biom(
  "microbiome_predicts_incident_T2D/combined_redist.species.biom") #BIOM table from the
  SHOGUN species-level output
  ncbi_data <- biomformat::biom_data(ncbi_data)
  ncbi_tax_table <- strsplit(row.names(as.matrix(ncbi_data)), ";")
  ncbi_tax_table <- matrix(unlist(ncbi_tax_table), nrow=length(ncbi_tax_table), byrow=T)
  row.names(ncbi_data) <- row.names(ncbi_tax_table)
  ncbi_data <- phyloseq(otu_table(as.matrix(ncbi_data), taxa_are_rows=T), tax_table(
  ncbi_tax_table))
  #Format the tax table.
  colnames(tax_table(ncbi_data)) <- c("Domain", "Phylum", "Class", "Order", "Family",
  "Genus", "Species")
  #Combine the phenotype data with the taxa data
```

```r
54      ncbi_data <- phyloseq(otu_table(ncbi_data), tax_table(ncbi_data), sample_data(FR02))
55      #remove samples with less than 50k reads (total)
56      to_be_pruned <- sample_sums(ncbi_data) > 50000
57      ncbi_data <- prune_samples(to_be_pruned, ncbi_data)
58      #remove pregnant (GRAVID==2) participants
59      ncbi_data <- subset_samples(ncbi_data, GRAVID %in% c(1, NA))
60      #remove participants who have used antibiotics in the last 6 months (BL_USE_RX_J01==1)
61      ncbi_data <- subset_samples(ncbi_data, BL_USE_RX_J01 %in% c(0, NA))
62      #remove participants with prevalent diabetes (PREVAL_DIAB==1)
63      ncbi_data <- subset_samples(ncbi_data, PREVAL_DIAB==0)
64      #remove participants with diabetes indicator values over set guidelines:
        FR02_GLUK_NOLLA >= 7, FR02_GLUK_120 >= 11.1 & HBA1C >= 48 (ignore NA values)
65      ncbi_data <- subset_samples(ncbi_data, FR02_GLUK_NOLLA<7 | is.na(FR02_GLUK_NOLLA))
66      ncbi_data <- subset_samples(ncbi_data, FR02_GLUK_120<11.1 | is.na(FR02_GLUK_120))
67      ncbi_data <- subset_samples(ncbi_data, HBA1C<48 | is.na(HBA1C))
68      #remove participants with type 2 diabetes diagnosis within two years from baseline
69      ncbi_data <- subset_samples(ncbi_data, DIAB_T2_AGEDIFF > 2)
70      #save the final objects
71      saveRDS(ncbi_data,
        "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data.RDs")
72  }
73
74
75      #Functions
76      prediab_cat <- function(pseq){
77          data <- sample_data(pseq)
78          prediab <- ifelse(data$FR02_GLUK_NOLLA >= 5.6 & data$FR02_GLUK_NOLLA < 6.9 | data$
            FR02_GLUK_120 >= 7.8 & data$FR02_GLUK_120 < 11 | data$HBA1C >= 39 & data$HBA1C < 47
            , 1, 0)
79          prediab <- as.factor(prediab)
80          return(prediab)
81      }
82
83      cox_wrapper <- function(data,
84                              predictors,
85                              covariates,
86                              status,
87                              time_to_event,
88                              alpha_level,
89                              normalize,
90                              test_ph_assumption) {
91      if(normalize) {
92        if(class(data[, predictors]) == "numeric") {
93          x <- data[, predictors]
94          data[, predictors] <- (x - mean(x, na.rm = T))/sd(x, na.rm = T)
95        } else {
96          data[, predictors] <- apply(data[, predictors], 2, FUN = function(x) {(x - mean(x
            , na.rm = T))/sd(x, na.rm = T) })
97        }
98      }
99      ## Formulas ***************************
100     linear_formulas <- lapply(predictors, function(x) {
101       formula_data <- deparse(substitute(data))
102       formula <- paste0("Surv(", formula_data, "$", time_to_event, ", ",formula_data,"$",
            status, ") ~ ",paste(covariates, collapse = "+"), " + ",x)
103       return(formula)
104     }) %>%
105       set_names(predictors)
106     ## Cox regression *********************
107     print("Cox")
108     linear_cox_fit <- lapply(linear_formulas, function(x) {
109       coxph(as.formula(x), data=data, x=TRUE)
110     })
111     ## Check PH assumptions ***************
112     if(test_ph_assumption) {
113       print("PH assumptions")
114       ph_assumption <-  lapply(predictors, function(m) {
115         west <- cox.zph(linear_cox_fit[[m]])
116         p_values <- west$table[, "p"]
```

```r
      # significant cases
      x <- which(p_values < 1)
      if(length(x) == 0) {
        return(NULL)
      }
      df <- data.frame(feature = m, variable_not_ph = names(x), p_value = p_values[x])
    }) %>%
      do.call(rbind, .) %>%
      mutate(p_adj = p.adjust(p_value, "BH")) %>%
      filter(p_value < alpha_level)
  }
  ## Results ****************************
  print("Results")
  results <- lapply(predictors, function(x) {
    df <- summary(linear_cox_fit[[x]])$coefficients %>% as.data.frame()
    df <- df[nrow(df), ] %>%
      select(coef, "se(coef)", "z", "Pr(>|z|)") %>%
      set_colnames(c("coef", "se_coef", "west_stat_value", "p")) %>%
      mutate(west_stat = "Wald")
    df <- df %>%
      mutate(predictor = x)
  }) %>%
    do.call(rbind, .)
  # Multiple westing correction
  results <- results %>%
    mutate(P_adjusted = p.adjust(p, "BH")) %>%
    ungroup() %>%
    group_by(predictor)
  # Results in neat form for presentation
  neat_results <- results %>%
    # filter(p == min(p)) %>%
    ungroup() %>%
    mutate(HR = round(exp(coef),3)) %>%
    mutate(HR_lower_95 = round(exp(coef - 1.96*se_coef), 3),
           HR_upper_95 = round(exp(coef + 1.96*se_coef), 3),
           P = round(p, 5),
           Coefficient = round(coef, 3),
           "Coefficient SE" = round(se_coef, 3)) %>%
    mutate(HR = paste0(HR, " (95% CI, ", HR_lower_95, "-", HR_upper_95, ")")) %>%
    select(Predictor = predictor, Coefficient, "Coefficient SE", HR, "p","P_adjusted",
           "west_stat_value", "west_stat") %>%
    mutate(HR = ifelse(is.na(Coefficient), NA, HR))  %>%
    filter(P_adjusted < alpha_level) %>%
    arrange(P_adjusted) %>%
    set_colnames(c("Predictor", "Coefficient", "Coefficient SE", "HR","P-value" ,"P
(adjusted)", "west Statistic Value", "west Statistic"))
  # Results in a form more convenient for further manipulations
  results <- results %>%
    ungroup %>%
    mutate(PH = exp(coef)) %>%
    mutate(p_adj = P_adjusted) %>%
    mutate(direction = ifelse(coef < 0, "negative", "positive"))
  if(nrow(neat_results) == 0) {
    return(list(results = results))
  }
  if(test_ph_assumption) {
    if(nrow(neat_results) == 0) {
      return(list(results = results, ph_assumption = ph_assumption))
    }
    return(list(neat_results = neat_results,
                results = results,
                ph_assumption = ph_assumption))
  }
  return(list(neat_results = neat_results, results = results))
}

#preprocess data
if (file.exists(
  "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_raw_east.RDs") &&
```

```r
      file.exists(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_raw_west.RDs") &&
183   file.exists("microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_main.RDs")
       && file.exists("microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca.RDs")
      &&
184   file.exists("microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_east.RDs")
       && file.exists(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_east.RDs") &&
185   file.exists("microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_west.RDs")
       && file.exists(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_west.RDs") &&
186   file.exists(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_data_east.RDs") &&
      file.exists(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_data_west.RDs")) {
187     ncbi_data_raw_east <- readRDS(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_raw_east.RDs")
188     ncbi_data_raw_west <- readRDS(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_raw_west.RDs")
189     ncbi_data_main <- readRDS(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_main.RDs")
190     ncbi_data_east <- readRDS(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_east.RDs")
191     ncbi_data_west <- readRDS(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_west.RDs")
192     ncbi_pca <- readRDS(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca.RDs")
193     ncbi_pca_east <- readRDS(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_east.RDs")
194     ncbi_pca_west <- readRDS(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_west.RDs")
195     ncbi_pca_data_east <- readRDS(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_data_east.RDs")
196     ncbi_pca_data_west <- readRDS(
      "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_data_west.RDs")
197   } else {
198     #Limit taxa to core in the east (EAST) set
199     core_ncbi_taxa <- core(prune_samples(meta(ncbi_data)$EAST == 1, ncbi_data) %>%
      transform("compositional"), detection = .1/100, prevalence = 1/10) %>% taxa_names()
200     ncbi_data_main <- prune_taxa(core_ncbi_taxa, ncbi_data)
201     #divide non-transformed data to east/west (EAST/WEST) sets
202     ncbi_data_raw_east <- prune_samples(meta(ncbi_data_main)$EAST == 1, ncbi_data_main)
203     ncbi_data_raw_west <- prune_samples(meta(ncbi_data_main)$EAST == 0, ncbi_data_main)
204     #CLR-transform raw counts
205     ncbi_data_main <- transform(ncbi_data_main, "clr")
206     #calculate additional variables
207     PREDIAB <- prediab_cat(ncbi_data_main)
208     NON_HDL <- sample_data(ncbi_data)$KOL - sample_data(ncbi_data)$HDL
209     #calculate diversity
210     ncbi_diversity <- estimate_richness(ncbi_data, measures = c("Observed", "Shannon"))
211     #reduce metadata to useful columns
212     useful_variables <- c("BL_AGE", "BMI", "MEN", "SYSTM", "CURR_SMOKE", "TRIG",
      "INCIDENT_DIAB_T2", "DIAB_T2_AGEDIFF", "EAST")
213     sample_data(ncbi_data_main) <- sample_data(ncbi_data_main)[,sample_variables(
      ncbi_data_main) %in% useful_variables]
214     #separate transformed and curated data to east/west (EAST/WEST) sets
215     ncbi_data_east <- prune_samples(meta(ncbi_data_main)$EAST == 1, ncbi_data_main)
216     ncbi_data_west <- prune_samples(meta(ncbi_data_main)$EAST == 0, ncbi_data_main)
217     #calculate 10 first PCAs with full community
218     ncbi_data_raw_clr <- transform(ncbi_data, "clr")
219     ncbi_pca <- ordinate(ncbi_data_raw_clr, "RDA")
220     ncbi_pca_data <- as.data.frame(scores(ncbi_pca, choices = c(1:10))$sites)
221     ncbi_pca_east <- ordinate(prune_samples(meta(ncbi_data_raw_clr)$EAST == 1,
      ncbi_data_raw_clr), "RDA")
222     ncbi_pca_data_east <- as.data.frame(scores(ncbi_pca_east, choices = c(1:10))$sites)
223     ncbi_pca_west <- ordinate(prune_samples(meta(ncbi_data_raw_clr)$EAST == 0,
      ncbi_data_raw_clr), "RDA")
224     ncbi_pca_data_west <- as.data.frame(scores(ncbi_pca_west, choices = c(1:10))$sites)
225     #combine with additional data
```

```r
226    sample_data(ncbi_data_main) <- cbind(sample_data(ncbi_data_main), PREDIAB, NON_HDL,
       ncbi_diversity, ncbi_pca_data)
227    sample_data(ncbi_data_east) <- cbind(sample_data(ncbi_data_east), PREDIAB = PREDIAB[
       which(meta(ncbi_data_main)$EAST == 1)], NON_HDL = NON_HDL[which(meta(ncbi_data_main)$
       EAST == 1)], ncbi_diversity[which(meta(ncbi_data_main)$EAST == 1),],
       ncbi_pca_data_east)
228    sample_data(ncbi_data_west) <- cbind(sample_data(ncbi_data_west), PREDIAB = PREDIAB[
       which(meta(ncbi_data_main)$EAST == 0)], NON_HDL = NON_HDL[which(meta(ncbi_data_main)$
       EAST == 0)], ncbi_diversity[which(meta(ncbi_data_main)$EAST == 0),],
       ncbi_pca_data_west)
229    saveRDS(ncbi_data_raw_east,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_raw_east.RDs")
230    saveRDS(ncbi_data_raw_west,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_raw_west.RDs")
231    saveRDS(ncbi_data_main,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_main.RDs")
232    saveRDS(ncbi_data_east,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_east.RDs")
233    saveRDS(ncbi_data_west,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_data_west.RDs")
234    saveRDS(ncbi_pca,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca.RDs")
235    saveRDS(ncbi_pca_east,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_east.RDs")
236    saveRDS(ncbi_pca_west,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_west.RDs")
237    saveRDS(ncbi_pca_data_east,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_data_east.RDs")
238    saveRDS(ncbi_pca_data_west,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/ncbi_pca_data_west.RDs")
239  }
240
241  #Filter features based on corrected p-values in the east dataset
242  #set variables
243  alpha_level <- 0.05 #to filter
244  status <- "INCIDENT_DIAB_T2"
245  time_to_event <- "DIAB_T2_AGEDIFF"
246  ncbi_cox_data_east <- cbind(meta(ncbi_data_east), as.matrix(t(otu_table(ncbi_data_east
       ))))
247  predictors <- c("Shannon", "Observed", colnames(ncbi_pca_data_east), taxa_names(
       ncbi_data_east))
248  covariates <- c("BL_AGE", "BMI", "MEN", "SYSTM", "NON_HDL", "CURR_SMOKE", "TRIG")
249  splines <- TRUE
250  normalize <- TRUE
251  test_ph_assumption <- FALSE
252  #Cox regression with previously defined function
253  set.seed(11235)
254  ncbi_cox_east <- cox_wrapper(data = ncbi_cox_data_east,
255                      predictors = predictors,
256                      covariates = covariates,
257                      alpha_level = alpha_level,
258                      status = status,
259                      time_to_event = time_to_event,
260                      normalize = normalize,
261                      test_ph_assumption = test_ph_assumption)
262
263  ncbi_cox_results_east <- merge(ncbi_cox_east$neat_results, as.data.frame(ncbi_data_east@
       tax_table@.Data), by.x="Predictor", by.y="row.names")
264  ncbi_cox_results_east <- ncbi_cox_results_east[order(-ncbi_cox_results_east$Coefficient)
       ,]
265  ncbi_cox_results_east$Species <- gsub("s__", "", ncbi_cox_results_east$Species)
266  ncbi_cox_results_east$Species <- gsub("_", " ", ncbi_cox_results_east$Species)
267  ncbi_cox_results_east$Family <- gsub("f__", "", ncbi_cox_results_east$Family)
268
269  #Correlations and clustering between the associated taxa in east data
270  otu_table_assoc_taxa <- as.data.frame(otu_table(prune_taxa(ncbi_cox_east$neat_results$
       Predictor, ncbi_data_raw_east)))
271  rownames(otu_table_assoc_taxa) <- ncbi_cox_results_east$Species[match(rownames(
       otu_table_assoc_taxa), ncbi_cox_results_east$Predictor)]
```

```r
272    set.seed(11235)
273    proprmatrix <- propr(t(otu_table_assoc_taxa), metric = "rho", p = 100)
274    clusters_assoc <- hclust(dist(proprmatrix@matrix), method = "ward.D2")
275    #Compute the Kelley-Gardner-Sutcliffe penalty function for a hierarchical cluster tree,
       to determine optimal number of clusters
276    op_k <- kgs(clusters_assoc, dist(proprmatrix@matrix), maxclus = 20)
277    op_k <- as.numeric(names(op_k[which(op_k == min(op_k))]))
278    cluster_ids <- cutree(tree = clusters_assoc, k = op_k)
279    svg("microbiome_predicts_incident_T2D/2y_exclusion_analysis/clusters.svg", width=10,
       height=10)
280    plot(clusters_assoc)
281    rect.hclust(clusters_assoc, k = op_k, border = 2:7)
282    dev.off()
283
284    heatmap_annotation <- data.frame(Species = rownames(proprmatrix@matrix), Cluster =
       cluster_ids)
285    heatmap_annotation$Predictor <- ncbi_cox_results_east$Predictor[match(heatmap_annotation
       $Species, ncbi_cox_results_east$Species)]
286
287    #Clustering correlating significant taxa for east and west data
288    #Combine read counts of clusters and calculate their CLR values
289    taxa_clusters <- merge(heatmap_annotation[c("Cluster")], ncbi_cox_results_east[c(
       "Species", "Predictor")], by.x = "row.names", by.y = "Species")
290    taxa_clusters$Cluster <- factor(taxa_clusters$Cluster, levels = 1:length(unique(
       taxa_clusters$Cluster)))
291
292    cluster_phylo_east <- ncbi_data_raw_east
293    cluster_phylo_west <- ncbi_data_raw_west
294    index_taxa <- c()
295    for (cluster in levels(taxa_clusters$Cluster)) {
296    taxa_to_merge <- taxa_clusters$Predictor[which(taxa_clusters$Cluster == cluster)]
297    cluster_phylo_east <- merge_taxa(cluster_phylo_east, taxa_to_merge, archetype=1)
298    cluster_phylo_west <- merge_taxa(cluster_phylo_west, taxa_to_merge, archetype=1)
299    index_taxa[cluster] <- taxa_to_merge[1]
300    }
301    cluster_phylo_east <- transform(cluster_phylo_east, "clr")
302    cluster_phylo_west <- transform(cluster_phylo_west, "clr")
303    #Retain only clusters
304    cluster_phylo_east <- prune_taxa(index_taxa, cluster_phylo_east)
305    cluster_phylo_west <- prune_taxa(index_taxa, cluster_phylo_west)
306    taxa_names(cluster_phylo_east) <- paste0("Cluster_", taxa_clusters$Cluster[match(
       taxa_names(cluster_phylo_east), taxa_clusters$Predictor)])
307    taxa_names(cluster_phylo_west) <- paste0("Cluster_", taxa_clusters$Cluster[match(
       taxa_names(cluster_phylo_west), taxa_clusters$Predictor)])
308
309    #test the individual taxa and clusters in the east data
310    #set variables
311    alpha_level <- 1 #to include everything in the results
312    status <- "INCIDENT_DIAB_T2"
313    time_to_event <- "DIAB_T2_AGEDIFF"
314    ncbi_cox_data_east_2 <- cbind(meta(ncbi_data_east), as.matrix(t(otu_table(ncbi_data_east
       ))), as.matrix(t(otu_table(cluster_phylo_east))))
315    predictors <- c(ncbi_cox_results_east$Predictor, taxa_names(cluster_phylo_east), "PC1")
316    covariates <- c("BL_AGE", "BMI", "MEN", "SYSTM", "NON_HDL", "CURR_SMOKE", "TRIG")
317    splines <- TRUE
318    normalize <- TRUE
319    test_ph_assumption <- FALSE
320    #Cox regression with previously defined function
321    set.seed(11235)
322    ncbi_cox_east_2 <- cox_wrapper(data = ncbi_cox_data_east_2,
323                        predictors = predictors,
324                        covariates = covariates,
325                        alpha_level = alpha_level,
326                        status = status,
327                        time_to_event = time_to_event,
328                        normalize = normalize,
329                        test_ph_assumption = test_ph_assumption)
330
331    ncbi_cox_results_east_2 <- data.frame(ncbi_cox_east_2$neat_results)
```

```
332    ncbi_cox_results_east_2 <- merge(ncbi_cox_results_east_2[c("Predictor", "Coefficient",
       "HR", "P.value")], ncbi_cox_results_east[c("Predictor", "Family", "Species")], by =
       "Predictor", all = TRUE)
333    ncbi_cox_results_east_2 <- ncbi_cox_results_east_2[order(-ncbi_cox_results_east_2$
       Coefficient),]
334    ncbi_cox_results_east_2$Set <- "East"
335
336    #test the individual taxa and clusters in the west data
337    #use same variables as for previous model run (thus not repeated here)
338    ncbi_cox_data_west <- cbind(meta(ncbi_data_west), as.matrix(t(otu_table(ncbi_data_west
       ))), as.matrix(t(otu_table(cluster_phylo_west))))
339    #Cox regression with previously defined function
340    set.seed(11235)
341    ncbi_cox_west <- cox_wrapper(data = ncbi_cox_data_west,
342                                 predictors = predictors,
343                                 covariates = covariates,
344                                 alpha_level = alpha_level,
345                                 status = status,
346                                 time_to_event = time_to_event,
347                                 normalize = normalize,
348                                 test_ph_assumption = test_ph_assumption)
349
350    ncbi_cox_results_west <- data.frame(ncbi_cox_west$neat_results)
351    ncbi_cox_results_west <- merge(ncbi_cox_results_west[c("Predictor", "Coefficient", "HR"
       , "P.value")], ncbi_cox_results_east[c("Predictor", "Family", "Species")], by =
       "Predictor", all = TRUE)
352    ncbi_cox_results_west <- ncbi_cox_results_west[order(-ncbi_cox_results_west$Coefficient)
       ,]
353    ncbi_cox_results_west$Set <- "West"
354
355    #save results
356    results_out_east <- rbind(data.frame(ncbi_cox_east$neat_results), data.frame(
       ncbi_cox_east_2$neat_results[which(grepl("Cluster", ncbi_cox_east_2$neat_results$
       Predictor)),]))
357    results_out_west <- data.frame(ncbi_cox_west$neat_results)
358    results_out_east <- merge(results_out_east[c("Predictor", "Coefficient", "HR", "P.value"
       , "P..adjusted.")], as.data.frame(ncbi_data_east@tax_table@.Data)["Species"], by.x=
       "Predictor", by.y="row.names", all.x = TRUE)
359    results_out_west <- merge(results_out_west[c("Predictor", "Coefficient", "HR", "P.value"
       , "P..adjusted.")], as.data.frame(ncbi_data_west@tax_table@.Data)["Species"], by.x=
       "Predictor", by.y="row.names", all.x = TRUE)
360    results_out_west$P..adjusted. <- NA
361    results_out_east[which(grepl("Cluster", results_out_east$Predictor)),]$P..adjusted. <- NA
362    results_out <- merge(results_out_east, results_out_west, by="Predictor", suffixes=c(
       ".east",".west"))
363    result_order <- results_out[rev(order(results_out$Coefficient.east)),]$Predictor
364    result_order <- c("PC1", paste0("Cluster_", 1:6),result_order[which(grepl("sp",
       result_order))])
365    results_out <- results_out[match(result_order, results_out$Predictor),]
366    results_out[-which(is.na(results_out$Species.east)),"Predictor"] <- as.character(
       results_out[-which(is.na(results_out$Species.east)),"Species.east"])
367    results_out$Predictor <- gsub("s__", "", results_out$Predictor)
368    results_out$Predictor <- gsub("_", " ", results_out$Predictor)
369    names(results_out) <- gsub("\\.\\.", ".", names(results_out))
370    results_out <- results_out[,!names(results_out) %in% c("Species.east", "Species.west",
       "P.adjusted.west")]
371    results_out[c("P.value.east", "P.adjusted.east", "P.value.west")] <- lapply(results_out[
       c("P.value.east", "P.adjusted.east", "P.value.west")], function (x) round(x, 4))
372    write.csv(results_out,
       "microbiome_predicts_incident_T2D/2y_exclusion_analysis/Table_S2.csv", row.names=F)
373
374    #plot heatmap of taxa associations, clustering, and hazard ratios in the east data
375    newnames <- lapply(rownames(proprmatrix@matrix),function(x) bquote(italic(.(x))))
376    heatmap_annotation$HR <- gsub("([0-9]\\.[0-9]*)[[:space:]].*", "\\1",
       ncbi_cox_results_east$HR[match(heatmap_annotation$Predictor, ncbi_cox_results_east$
       Predictor)])
377    heatmap_annotation$HR <- round(as.numeric(as.character(heatmap_annotation$HR)), 1)
378    heatmap_annotation$HR <- factor(heatmap_annotation$HR, levels = rev(seq(0.8, 1.2, 0.1)))
379
```

```
380    ann_colors <- list(HR = brewer.pal(n = 5, name = "BrBG"), Cluster = brewer.pal(n = 12,
       name = "Paired")[-seq(0,12,2)])
381    names(ann_colors$HR) <- levels(heatmap_annotation$HR)
382    names(ann_colors$Cluster) <- c("1", "2", "3", "4", "6", "5")
383    ann_colors$Cluster <- factor(ann_colors$Cluster, levels = ann_colors$Cluster[c(4,2,6,5,3
       ,1)])
384    heatmap_colors <- rev(brewer.pal(n = 10, name = "RdBu"))
385    heatmap_colors[c(5,6)] <- "#FFFFFF"
386    svg("microbiome_predicts_incident_T2D/2y_exclusion_analysis/correlations.svg", width=15
       , height=15)
387    pheatmap(proprmatrix@matrix, labels_row = as.expression(newnames), labels_col =
       as.expression(newnames), annotation_row = heatmap_annotation[4], treeheight_row = 0,
       annotation_col = heatmap_annotation[2], annotation_colors = ann_colors, cutree_rows =
       op_k, cutree_cols = op_k, clustering_method = "ward.D2", color = heatmap_colors, breaks
       = seq(-1, 1, length.out = 11), legend_breaks = seq(-1, 1, length.out = 11), cellwidth=10
       , cellheight=10)
388    dev.off()
389
390    #plot HR of both west and east data
391    ncbi_cox_results <- rbind(ncbi_cox_results_east_2, ncbi_cox_results_west)
392
393    Species <- c()
394    Family <- c()
395    Set <- c()
396    Facet <- c()
397    HR <- c()
398    HR1 <- c()
399    HR2 <- c()
400
401    for (i in 1:length(ncbi_cox_results$Predictor)){
402      Species[[i]] <- ifelse(is.na(ncbi_cox_results$Species[i]), sub("_", " ",
         ncbi_cox_results$Predictor[i]), as.character(ncbi_cox_results$Species[i]))
403      Family[[i]] <- ifelse(is.na(ncbi_cox_results$Family[i]), NA, as.character(
         ncbi_cox_results$Family[i]))
404      HR[[i]] <- str_split(ncbi_cox_results$HR[i]," ")[[1]][1]
405      HR_range <- str_split(ncbi_cox_results$HR[i]," ")[[1]][4]
406      HR1[[i]] <- str_split(HR_range,"-")[[1]][1]
407      HR2_bef <- str_split(HR_range,"-")[[1]][2]
408      HR2[[i]] <- substr(HR2_bef,1,nchar(HR2_bef)-1)
409      Set[[i]] <- ncbi_cox_results$Set[i]
410      Facet[[i]] <- ifelse(is.na(ncbi_cox_results$Family[i]), "Grouping", "Taxa")
411      HRdf <- data.frame(Species = Species,
412                         Family = Family,
413                         Set = Set,
414                         Facet = Facet,
415                         HR = HR,
416                         HR1 = HR1,
417                         HR2 = HR2)
418    }
419
420    family_color_map <- data.frame(Color = c("chartreuse2", "#7b562e", "#9bb940", "#c5bb9a"
       , "darkred", "#ff4ae3", "#339a00", "#d78343", "darkblue", "#5f96d6", "black"),
421     Family = c("Akkermansiaceae", "Bacteroidaceae", "Clostridiaceae", "Eggerthellaceae",
        "Eubacteriaceae", "Lachnospiraceae", "Oscillospiraceae", "Rickenellaceae",
        "Ruminococcaceae", "Sutterellaceae", NA))
422
423    HRdf$Species <- factor(HRdf$Species, levels = c(paste0("Cluster ", 6:1), "PC1",
       as.character(HRdf[which(HRdf$Set %in% "East" & HRdf$Facet %in% "Taxa"),][order(HRdf[
       which(HRdf$Set %in% "East" & HRdf$Facet %in% "Taxa"),]$HR),]$Species)))#order features
       by effect size in the east data
424    p <-  ggplot(data = HRdf, aes(y = Species, x = as.numeric(as.character(HR)), color =
       Family)) +
425          geom_pointrange(aes(xmin=as.numeric(as.character(HR1)), xmax=as.numeric(
          as.character(HR2))), lwd = 1) +
426          scale_x_continuous(limits = c(0.6, 1.5)) +
427          scale_color_manual(name = "Family", values = as.character(family_color_map$Color))
            +
428          guides(color = guide_legend(override.aes = list(size = 1.4))) +
429          xlab("HR") + ylab("Species") +
```

```r
430          geom_vline(xintercept=c(1.0), linetype="dotted") +
431          theme(axis.text.y = element_text(face = "italic"), legend.text = element_text(face
                  = "italic"), axis.title.y = element_blank()) +
432          facet_grid(Facet~Set, scales = "free")
433
434     ggsave("microbiome_predicts_incident_T2D/2y_exclusion_analysis/HR_comparison.svg", plot=
        p, units="in", width=15, height=10)
435
436     #Plot Kaplan-Meier curves
437     kp_predictors <- ncbi_cox_results_west$Predictor[which(ncbi_cox_results_west$P.value <
        0.05)]
438     kp_covariates <- covariates
439     kp_time_to_event <- time_to_event
440     kp_status <- status
441     kp_data <- ncbi_cox_data_west[,which(colnames(ncbi_cox_data_west) %in% c(kp_status,
        kp_time_to_event, kp_predictors, kp_covariates))]
442     kp_time <- seq(0, max(kp_data$DIAB_T2_AGEDIFF), by = .01)
443     kp_list <- list(NULL)
444     for (time in 1:length(kp_time)) {
445     kp_table <- lapply(kp_predictors, function(x) {
446         return_table <- data.frame(groupkm(kp_data[x], Surv(kp_data$DIAB_T2_AGEDIFF, kp_data
                $INCIDENT_DIAB_T2), g=4, u=kp_time[time], pl=FALSE))
447         return_table$Predictor <- x
448         return_table$quantile <- c(1:4)
449         return(return_table)
450      })
451     kp_table <- do.call(rbind, kp_table)
452     kp_table$time <- kp_time[time]
453     kp_list[[time]] <- kp_table
454     }
455
456     kp_list <- do.call(rbind, kp_list)
457     kp_predictors <- recode(kp_predictors, 'sp2673' = "[Clostridium] citroniae", 'sp2671' =
        "[Clostridium] bolteae", 'sp2697' = "Tyzzerella nexilis", 'sp2638' = "[Ruminococcus]
        gnavus")
458     kp_predictors <- gsub("_", " ", kp_predictors)
459     kp_list$Predictor <- recode(kp_list$Predictor, 'sp2673' = "[Clostridium] citroniae",
        'sp2671' = "[Clostridium] bolteae", 'sp2697' = "Tyzzerella nexilis", 'sp2638' =
        "[Ruminococcus] gnavus")
460     kp_list$Predictor <- gsub("_", " ", kp_list$Predictor)
461     kp_list$Predictor <- factor(kp_list$Predictor, levels = kp_predictors)
462
463     p <-  ggplot(data = kp_list, aes(y = KM, x = time, group = quantile)) +
464         geom_line(aes(color = quantile)) +
465         geom_vline(aes(xintercept = 2), linetype = "dashed", color ="gray") +
466         scale_color_viridis(labels = c("Min to Q1", "Q1 to Q2", "Q2 to Q3", "Q3 to max")) +
467         scale_y_continuous(breaks = pretty_breaks()) +
468         guides(color = guide_legend(override.aes = list(size = 1.4)), fill = "none") +
469         xlab("Time (years)") + ylab("Survival without type 2 diabetes") +
470         labs(color = "Relative\nabundance\nrange") +
471         facet_wrap(~ Predictor)
472     ggsave("microbiome_predicts_incident_T2D/2y_exclusion_analysis/KP_plot.svg", plot=p,
        units="cm", width=30, height=20)
473
474     #Plot distributions of the quartiles (for inlays in the KP-plot)
475     quartile_data <- lapply(kp_data[ncbi_cox_results_west$Predictor[which(
        ncbi_cox_results_west$P.value < 0.05)]],
476         function(x) data.frame(x_value = density(x)$x,
477                                y_value = density(x)$y,
478                                quartile = factor(paste0("Q",findInterval(density(x)$x,
                                    quantile(x, prob=c(0, 0.25, 0.5, 0.75, 1)), all.inside=T)))))
479     quartile_data <- data.frame(rbindlist(quartile_data, idcol="Predictor"))
480     quartile_data$Predictor <- recode(quartile_data$Predictor, 'sp2673' = "[Clostridium]
        citroniae", 'sp2671' = "[Clostridium] bolteae", 'sp2697' = "Tyzzerella nexilis",
        'sp2638' = "[Ruminococcus] gnavus")
481     quartile_data$Predictor <- gsub("_", " ", quartile_data$Predictor)
482     quartile_data$Predictor <- factor(quartile_data$Predictor, levels = kp_predictors)
483
484     p <- ggplot(quartile_data, aes(x_value, y_value)) +
```

```r
        geom_line() +
        geom_ribbon(aes(ymin=0, ymax=y_value, fill=quartile)) +
        scale_fill_viridis(labels = c("Q1", "Q2", "Q3", "Q4"), discrete=T) +
        guides(fill = "none") +
        theme(axis.title = element_blank()) +
        facet_wrap(~ Predictor)
ggsave("microbiome_predicts_incident_T2D/2y_exclusion_analysis/KP_plot_quartiles.svg",
plot=p, units="cm", width=30, height=20)

#Correlations and clustering between the associated taxa in west data
otu_table_assoc_taxa_west <- as.data.frame(otu_table(prune_taxa(ncbi_cox_west$
neat_results$Predictor, ncbi_data_raw_west)))
rownames(otu_table_assoc_taxa_west) <- ncbi_cox_results_west$Species[match(rownames(
otu_table_assoc_taxa_west), ncbi_cox_results_west$Predictor)]
set.seed(11235)
proprmatrix_west <- propr(t(otu_table_assoc_taxa_west), metric = "rho", p = 100)
clusters_assoc_west <- hclust(dist(proprmatrix_west@matrix), method = "ward.D2")
#Compute the Kelley-Gardner-Sutcliffe penalty function for a hierarchical cluster tree,
to determine optimal number of clusters
op_k_west <- kgs(clusters_assoc_west, dist(proprmatrix_west@matrix), maxclus = 20)
op_k_west <- as.numeric(names(op_k_west[which(op_k_west == min(op_k_west))]))
cluster_ids_west <- cutree(tree = clusters_assoc_west, k = op_k_west)
svg("microbiome_predicts_incident_T2D/2y_exclusion_analysis/clusters_west.svg", width=10
, height=10)
plot(clusters_assoc_west)
rect.hclust(clusters_assoc_west, k = op_k_west, border = 2:7)
dev.off()

#plot heatmap of taxa associations, clustering, and hazard ratios in the west data
newnames_west <- lapply(rownames(proprmatrix_west@matrix),function(x) bquote(italic(.(x
))))

#clusters are identical in membership of taxa in the same cluster, so we can just copy
the cluster annotation from east data to west data to keep cluster colors and naming
consistent
heatmap_annotation_west <- heatmap_annotation
#get correct hazard ratios for west data
heatmap_annotation_west$HR <- gsub("([0-9]\\.[0-9]*)[[:space:]].*", "\\1",
ncbi_cox_results_west$HR[match(heatmap_annotation_west $Predictor, ncbi_cox_results_west
$Predictor)])
heatmap_annotation_west$HR <- round(as.numeric(as.character(heatmap_annotation_west $HR
)), 1)
heatmap_annotation_west$HR <- factor(heatmap_annotation_west $HR, levels = rev(seq(0.8,
1.2, 0.1)))

svg("microbiome_predicts_incident_T2D/2y_exclusion_analysis/correlations_west.svg",
width=15, height=15)
pheatmap(proprmatrix_west@matrix, labels_row = as.expression(newnames_west), labels_col
= as.expression(newnames_west), annotation_row = heatmap_annotation_west[4],
treeheight_row = 0, annotation_col = heatmap_annotation_west[2], annotation_colors =
ann_colors, cutree_rows = op_k_west, cutree_cols = op_k_west, clustering_method =
"ward.D2", color = heatmap_colors, breaks = seq(-1, 1, length.out = 11), legend_breaks =
 seq(-1, 1, length.out = 11), cellwidth=10, cellheight=10)
dev.off()

save.image("microbiome_predicts_incident_T2D/2y_exclusion_analysis/Analysis.RData")
```