

# Simulated and UK tree clustered by UCSD soft. - March 2016

S. Le Vu  
(Dated: March 14, 2016)

## I. INTRO

- Looking for relation between cluster characteristics and heterogenous transmission rates
  - "Time-based" distances from *simulated* coalescent tree are converted to substitutions per site distances with a constant rate from the literature ( $4.3e-3/365$  from *Berry et al. JVI 2007*) (see distributions in Fig. 1)
  - Then matrices of distances are converted in edge lists of pairwise distances with header [ID1, ID2, distance] as needed by UCSD software **hivclustering**
  - Edge lists are inputed in the **hivnetworkcsv** function which returns lists of cluster assignments for thresholds [0.015, 0.02, 0.05, 0.1] for both simulated and UK tree.

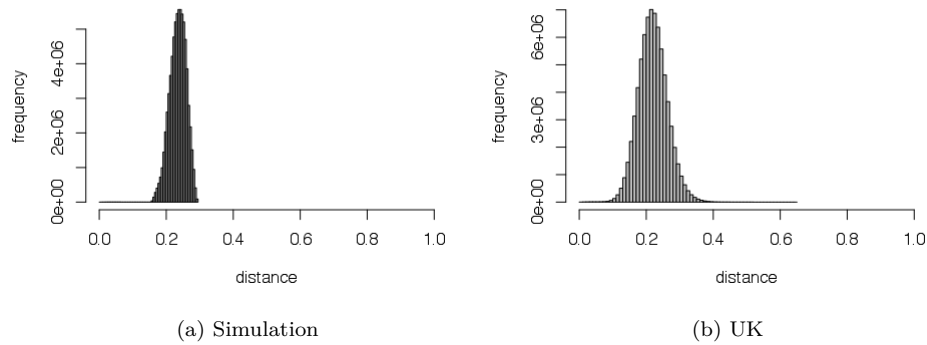


FIG. 1: Distances (subst/site)

```
detail_knitr <- FALSE
source("functions.R")
```

## II. UCSD HIVCLUSTERING

Read saved results from UCSD hivclustering

```
### only cluster members (based on given threshold)
simclus <- readRDS(file = "data/simclus2.rds")[[ "cl" ]]
ukclus <- readRDS(file = "data/ukclus2.rds")[[ "cl" ]]

### only cluster members (based on quantiles)
# simclus <- readRDS(file = "data/simclus.rds")[[3]]
# ukclus <- readRDS(file = "data/ukclus.rds")[[3]]
```

- Now, thresholds for clustering are fixed and not determined by quantiles of distances
- Because of this and because of the applied substitution rate, the number of clusters and mean sizes of clusters are not identical but are of the same magnitude between simulated and real UK data

- Up to threshold = 5%, distributions of cluster size seem to follow a power law for both simulated and UK data

Number of clusters and stats for simulated and UK trees (cluster size 1 does not exist in these (UCSD) outputs)

```
##- number of different clusters by threshold
sapply(simfreqClust, function(x) dim(x)[1])

0.015  0.02  0.05  0.1
1644   2041  2021  1637

sapply(ukfreqClust, function(x) dim(x)[1])

0.015  0.02  0.05  0.1
1199   1374  1460  611

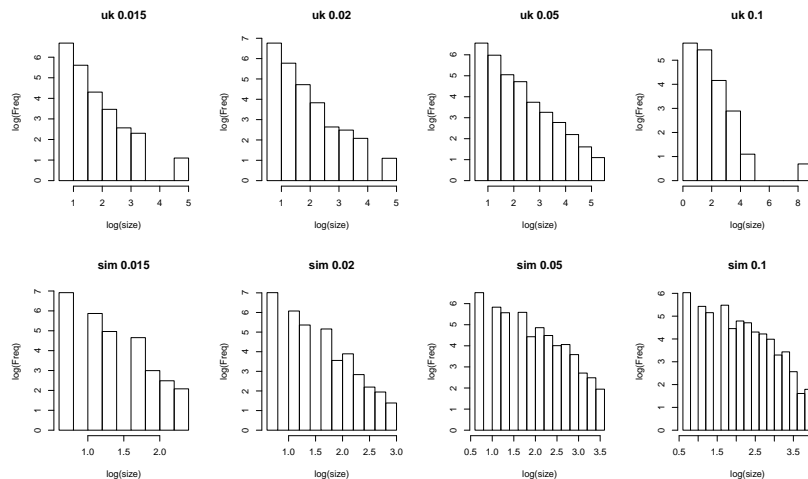
##- cluster size
sapply(simfreqClust, function(x) summary(x$Freq))

      0.015  0.02  0.05  0.1
Min.      2.00  2.000  2.000  2.000
1st Qu.    2.00  2.000  2.000  2.000
Median     2.00  2.000  3.000  5.000
Mean       2.73  3.103  5.147  6.896
3rd Qu.    3.00  3.000  6.000  9.000
Max.      11.00 20.000 35.000 52.000

sapply(ukfreqClust, function(x) summary(x$Freq))

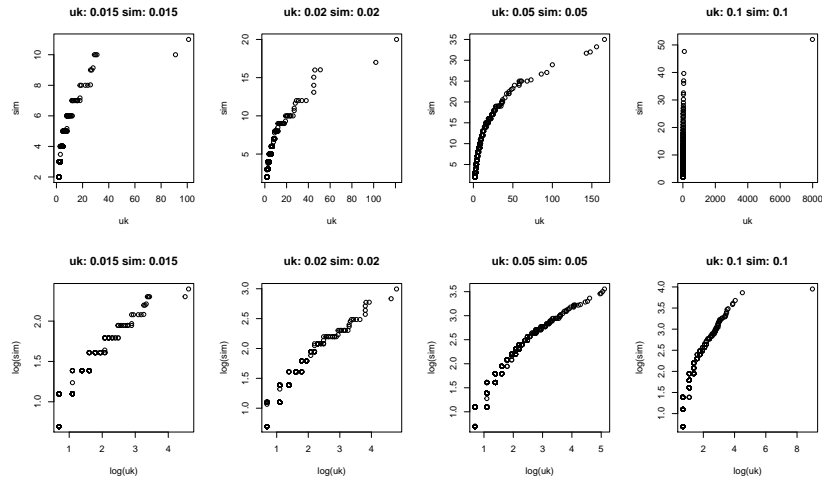
      0.015  0.02  0.05  0.1
Min.      2.000  2.000  2.000  2.00
1st Qu.    2.000  2.000  2.000  2.00
Median     2.000  2.000  3.000  3.00
Mean       3.188  3.531  5.442 17.84
3rd Qu.    3.000  3.000  4.000  4.00
Max.     101.000 121.000 166.000 7986.00
```

Plots of log(size) for UK and simulated clusters



QQ plots UK vs simulated, untransformed and log-log

```
null device
1
```



### III. ASSOCIATIONS

After merging with co-variables allocated from demes states contained in tree, non-clustering individuals are assigned a cluster size of 1.

The proportion of individuals into clusters (0 = out of, 1 = in a cluster) and stats for "size of cluster for each individuals".

```
##-proportion in or out clusters
sapply(l_sim, function(x) round(prop.table(table(x$binclus)),2))
```

```
      0.015 0.02 0.05 0.1
0  0.63 0.48 0.14 0.07
1  0.37 0.52 0.86 0.93
```

```
sapply(l_uk, function(x) round(prop.table(table(x$binclus)),2))
```

```
      0.015 0.02 0.05 0.1
0  0.69 0.6 0.35 0.1
1  0.31 0.4 0.65 0.9
```

```
##- cluster sizes (by individuals having such a size !!)
## probably meaningless
```

```
sapply(l_sim, function(x) summary(x$size))
```

```
      0.015 0.02 0.05 0.1
Min.    1.00 1.000 1.000 1.00
1st Qu. 1.00 1.000 2.000 4.00
Median  1.00 2.000 6.000 10.00
Mean    1.85 2.695 7.914 12.15
3rd Qu. 2.00 3.000 11.000 17.00
Max.   11.00 20.000 35.000 52.00
```

```
sapply(l_uk, function(x) summary(x$size))
```

```
      0.015 0.02 0.05 0.1
Min.    1.000 1.000 1.00 1
1st Qu. 1.000 1.000 1.00 11
Median  1.000 1.000 3.00 7986
Mean    3.941 5.853 18.82 5247
3rd Qu. 2.000 3.000 16.00 7986
Max.   101.000 121.000 166.00 7986
```

### A. Naive regressions on simulation

Linear regression with co-variates as ordinal or categorical

```
reg.sum(ls = listclus, reg = lm, model = lm_model_ordinal)

$model
[1] "scale(size) ~ scale(age) + scale(stage) + scale(time) + scale(risk)"

$parameter
              0.015      0.02      0.05      0.1
(Intercept)  3.8e-14 -2.8e-14 -9.4e-16  2.0e-15
scale(age)   -1.3e-02 -8.5e-03 -1.4e-02 -6.5e-03
scale(stage) -9.7e-02 -8.2e-02 -3.3e-02 -2.3e-02
scale(time)   2.8e-01  3.0e-01  3.0e-01  2.3e-01
scale(risk)  -8.2e-03 -9.0e-03 -3.2e-02 -2.3e-02

$pvalue
              0.015 0.02 0.05 0.1
(Intercept)
scale(age)
scale(stage) *** *** *** **
scale(time)  *** *** *** ***
scale(risk)          *** **

$r.squared
      0.015  0.02  0.05  0.1
0.0948 0.1010 0.0926 0.0573

reg.sum(ls = listclus, reg = lm, model = lm_model_factor)

$model
[1] "size ~ factor(stage) + factor(risk) + factor(age)"

$parameter
              0.015  0.02  0.05  0.1
(Intercept)   2.400  3.500  9.60 14.00
factor(stage)2 -0.330 -0.470 -0.75 -0.96
factor(stage)3 -0.350 -0.470 -0.49 -0.66
factor(stage)4 -0.550 -0.770 -1.10 -1.40
factor(stage)5 -0.620 -0.980 -1.40 -1.50
factor(risk)2  -0.034 -0.065 -0.56 -0.59
factor(age)2   -0.069 -0.120 -0.45 -1.00
factor(age)3   -0.200 -0.310 -1.10 -1.70
factor(age)4   -0.170 -0.260 -0.86 -1.20

$pvalue
              0.015 0.02 0.05 0.1
(Intercept) *** *** *** ***
factor(stage)2 *** *** *** **
factor(stage)3 *** *** *  *
factor(stage)4 *** *** *** ***
factor(stage)5 *** *** *** ***
factor(risk)2          *** *
factor(age)2           *
factor(age)3 ** ** *** ***
factor(age)4 ** * ** **
```

```

$r.squared
  0.015    0.02    0.05    0.1
0.01740 0.01320 0.00689 0.00429

##- logistic
##- model: clus ~ age + stage + time + risk
##- care = 1 for all at diagnosis

# logfit <- lapply(simli , function(x){
#   summary(glm(formula = logit_model_std,
#               data = x, family = binomial(link = "logit")))
# })

```

- The size of cluster of each individuals is always explained by overall stage and time of sampling
- For every threshold of clustering, stage 1 is more likely to belong to large clusters than any other stages
- Age 1 is more likely to belong to large clusters than age 3 or age 4
- Low (!) risk level is associated with larger clusters, only for higher thresholds (fewer and larger clusters on overall)
- Small part of the variance is explained by the variables

Logistic regression with ordinal and categorical variables

```

reg.sum(ls = listclus, reg = glm, model = logit_model_ord, family = binomial(link = "logit"))

$model
[1] "binclus ~ scale(age) + scale(stage) + scale(time) + scale(risk)"

$parameter
      0.015    0.02    0.05    0.1
(Intercept) -0.600  0.086  2.100  3.00
scale(age)   -0.059 -0.085 -0.068 -0.10
scale(stage) -0.270 -0.240 -0.170 -0.14
scale(time)   0.600  0.620  0.880  0.99
scale(risk)   -0.013 -0.068 -0.130 -0.13

$pvalue
      0.015 0.02 0.05 0.1
(Intercept) *** *** *** ***
scale(age)   ** ***  *  *
scale(stage) *** *** *** ***
scale(time)  *** *** *** ***
scale(risk)      *** *** ***

reg.sum(ls = listclus, reg = glm, model = logit_model_fact, family = binomial(link = "logit"))

$model
[1] "binclus ~ factor(stage) + factor(risk) + factor(age)"

$parameter
      0.015    0.02    0.05    0.1
(Intercept)  0.260  1.00  2.50  3.50
factor(stage)2 -0.420 -0.44 -0.14 -0.22
factor(stage)3 -0.500 -0.48 -0.18 -0.16
factor(stage)4 -0.860 -0.76 -0.41 -0.45
factor(stage)5 -0.940 -0.91 -0.63 -0.58

```

```

factor(risk)2 -0.034 -0.16 -0.29 -0.30
factor(age)2 -0.092 -0.21 -0.25 -0.45
factor(age)3 -0.290 -0.39 -0.41 -0.64
factor(age)4 -0.300 -0.43 -0.42 -0.65

```

```

$pvalue
      0.015 0.02 0.05 0.1
(Intercept)    **   ***   ***   ***
factor(stage)2  ***   ***
factor(stage)3  ***   ***   .
factor(stage)4  ***   ***   ***   ***
factor(stage)5  ***   ***   ***   ***
factor(risk)2   ***   ***   ***
factor(age)2    *     .     *
factor(age)3    ***   ***   **   **
factor(age)4    ***   ***   **   ***

```

All variables are significantly associated with being into a cluster for all thresholds

## B. Regressions on down-sampled simulation

- For each clustering threshold,
- sample one individual by cluster
- re-sample 100 times
- sample size equals number of clusters with each non-clustered individuals counting for one

### 1. First type of analysis

- apply 100 linear regressions on each sample
- calculate proportion of p-values < 0.05 over the 100 iterations

```

### over thresholds

## ordinal variables
dd_ord <- sapply( listclus, function(x) {
  downsample(df = x,
             lm_model = lm_model_ordinal,
             var = c("time", "age", "care", "stage", "risk"),
             iter = 100)
})

## percent of signif paramater
t(do.call(rbind, dd_ord["percent.signif", ] ))

      0.015 0.02 0.05 0.1
(Intercept)  0.00 0.00 0.00 0.00
scale(age)   0.28 0.51 0.13 0.12
scale(risk)  0.02 0.30 0.63 0.25
scale(stage) 1.00 1.00 0.84 0.61
scale(time)  1.00 1.00 1.00 1.00

## categorical variables
dd_cat <- sapply( listclus, function(x) {
  downsample(df = x,
             lm_model = lm_model_factor,
             var = c("time", "age", "care", "stage", "risk"),
             iter = 100)
})

```

```

}
)
## percent of signif paramater
t(do.call(rbind, dd_cat["percent.signif", ] ))
      0.015 0.02 0.05 0.1
(Intercept)      1.00 1.00 1.00 1.00
factor(age)2      0.09 0.21 0.15 0.11
factor(age)3      0.49 0.66 0.49 0.39
factor(age)4      0.44 0.66 0.40 0.43
factor(risk)2      0.02 0.35 0.66 0.48
factor(stage)2     1.00 1.00 0.16 0.23
factor(stage)3     1.00 1.00 0.22 0.22
factor(stage)4     1.00 1.00 0.69 0.69
factor(stage)5     1.00 1.00 0.91 0.77

```

- For first two threshold, first stage of infection and recent time of sampling are always (100%) associated with cluster size
  - As sizes of clusters increase, young age and low-risk tend to be associated with cluster size (in maximum 47% and 63% of samples respectively)
2. Second type of analysis
- calculate the mean of co-variates over the 100 iterations
  - apply one linear regression: size ~ mean(covariates)

```

## mean of co-variates by cluster
# str(dd["mean.sample",])
mean.down <- dd_ord["mean.sample",]

##- linear regression ordinal
# lapply(mean.down, function(x) {
#   summary(lm(lm_model_ordinal, data = x))})
reg.sum(ls = mean.down, reg = lm, model = lm_model_ordinal)

$model
[1] "scale(size) ~ scale(age) + scale(stage) + scale(time) + scale(risk)"

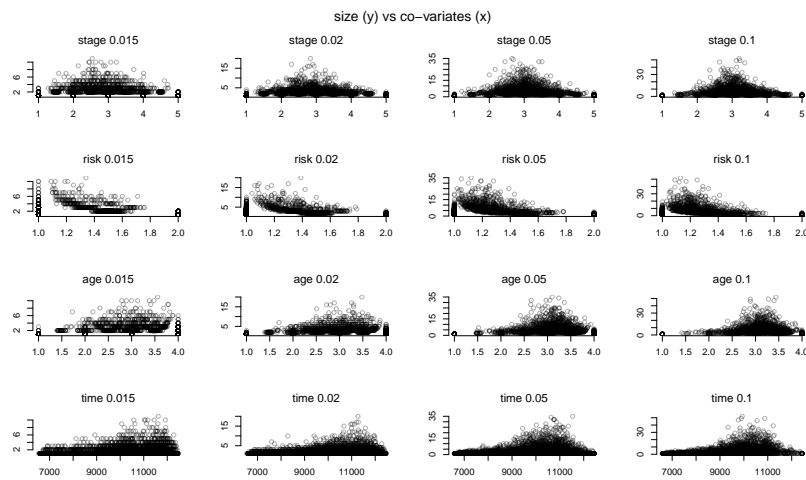
$parameter
      0.015      0.02      0.05      0.1
(Intercept) 3.9e-15 1.4e-14 -8.7e-15 -3.1e-15
scale(age)   -1.6e-02 -2.4e-02 -1.6e-02 -1.9e-02
scale(stage) -8.4e-02 -8.1e-02 -4.8e-02 -5.1e-02
scale(time)  1.8e-01 2.0e-01 3.1e-01 3.4e-01
scale(risk)  -6.4e-03 -1.9e-02 -3.8e-02 -3.1e-02

$pvalue
      0.015 0.02 0.05 0.1
(Intercept)
scale(age)      *
scale(stage)   *** ***  **  **
scale(time)    *** ***  *** ***
scale(risk)      .  *   .

$r.squared
      0.015      0.02      0.05      0.1
0.0435 0.0523 0.1060 0.1280

## plot
size.vs.covar(mean.down)
# dev.off()

```



- First stages of infection and recent time of sampling are always (100%) associated with larger cluster size

### C. On real UK data

Same process ...

#### 1. Naive regressions for real UK data

Linear regression

```
#### just on low and high threshold (but not too high !)
# li <- listUKclus[ 1:(length(listUKclus)-1) ]
lm_model_uk = "scale(size) ~ scale(ageddiag) + scale(sqrt(cd4)) + scale(ydiag)"
# lapply(li, function(x) summary(lm(lm_model_std, data = x)))
reg.sum(ls = listUKclus, reg = lm, model = lm_model_uk)

$model
[1] "scale(size) ~ scale(ageddiag) + scale(sqrt(cd4)) + scale(ydiag)"

$parameter
              0.015    0.02    0.05    0.1
(Intercept)   0.0022  0.0021  0.0024  0.0013
scale(ageddiag) -0.0390 -0.0530 -0.0072  0.0970
scale(sqrt(cd4)) 0.0380  0.0490  0.0480  0.0130
scale(ydiag)    0.1200  0.1500  0.1100 -0.2500

$pvalue
              0.015 0.02 0.05 0.1
(Intercept)
scale(ageddiag) ***  ***      ***
scale(sqrt(cd4)) ***  ***  ***
scale(ydiag)    ***  ***  ***  ***

$r.squared
              0.015 0.02 0.05 0.1
0.0146 0.0234 0.0134 0.0583
```

Young age, high level of CD4 and recent time of diagnosis are associated with larger cluster size



```
##- model: clus ~ age + stage + time + risk
##- care = 1 for all at diagnosis
## ex.
logit_model_uk = "binclus ~ scale(agediag) + scale(sqrt(cd4)) + scale(ydiag)"

reg.sum(ls = listUKclus, reg = glm, model = logit_model_uk, family = binomial(link = "logit"))

$model
[1] "binclus ~ scale(agediag) + scale(sqrt(cd4)) + scale(ydiag)"

$parameter
              0.015  0.02  0.05    0.1
(Intercept)  -0.88 -0.47  0.71  2.2000
scale(agediag) -0.17 -0.19 -0.14  0.0980
scale(sqrt(cd4)) 0.22  0.22  0.25  0.2000
scale(ydiag)    0.73  0.75  0.72 -0.0032

$pvalue
              0.015  0.02  0.05  0.1
(Intercept)    ***   ***   ***   ***
scale(agediag)  ***   ***   ***   **
scale(sqrt(cd4)) ***   ***   ***   ***
scale(ydiag)    ***   ***   ***
```

Same associations with binary cluster membership

## 2. Down-sampled regressions for real UK data

```
### over thresholds

## ordinal variables
dd_ord_uk <- sapply( listUKclus, function(x) {
  downsample(df = x,
             lm_model = lm_model_uk,
             var = c("agediag", "cd4", "ydiag"),
             iter = 100)
})

## percent of signif paramater
t(do.call(rbind, dd_ord_uk["percent.signif", ] ))

              0.015  0.02  0.05  0.1
(Intercept)    0.00  0.00  0.00  0.03
scale(agediag)  0.59  0.70  0.25  0.04
scale(sqrt(cd4)) 0.77  0.89  0.73  0.10
scale(ydiag)    1.00  1.00  1.00  0.09
```

```
## mean of co-variables by cluster
# str(dd["mean.sample",])
mean.down_uk <- dd_ord_uk["mean.sample",]
# head(mean.down_uk[[2]])
##- linear regression ordinal
# lapply(mean.down, function(x) {
#   summary(lm(lm_model_ordinal, data = x))})
reg.sum(ls = mean.down_uk, reg = lm, model = lm_model_uk)
```

```

$model
[1] "scale(size) ~ scale(agediag) + scale(sqrt(cd4)) + scale(ydiag)"

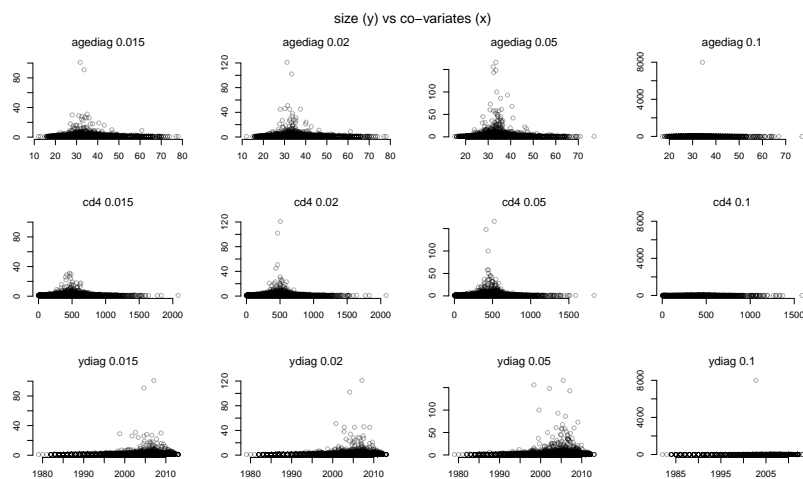
$parameter
              0.015  0.02  0.05  0.1
(Intercept) -0.018 -0.017 -0.043 -2.4e-02
scale(agediag) -0.017 -0.021 -0.011 -9.7e-06
scale(sqrt(cd4)) 0.034  0.039  0.042  1.4e-03
scale(ydiag)    0.073  0.082  0.075  2.0e-03

$pvalue
              0.015  0.02  0.05  0.1
(Intercept)    **      .   ***   ***
scale(agediag)  *      *
scale(sqrt(cd4)) ***   ***   ***   **
scale(ydiag)    ***   ***   ***   ***

$r.squared
      0.015  0.02  0.05  0.1
0.01610 0.00973 0.01330 0.01390

## plot
size.vs.covar(l = mean.down_uk, depvar = "size",
indepvar = c("agediag", "cd4", "ydiag"))
# dev.off()

```



CD4 and year of diagnosis systematically associated with cluster size / membership.