

## Introduction

I checked whether the site likelihoods of an inferred tree from the language dataset, have the same distribution as the likelihoods from data that have been indeed produced by the tree itself. I did the same for the parsimony scores of the sites for the maximum likelihood tree.

## The Code

```
1 library(phangorn)
2 rm(list=ls())
3 set.seed(8293) ## change this for further analyses
4
5 cog.tree=read.tree("IE2011_Cognates_rel_ANNOT.nwk")
6 data = read.table("46glossesfull.csv", header=TRUE, row.names
7   =1, colClasses='character')
8 cog.tree$tip.label[cog.tree$tip.label == "Ptg-E"] = "PtgE"
9 datamat = as.matrix(data)
10
11 tab <- apply(datamat, 2, function(x){table(factor(x, levels=c
12   (0,1)))})
13 indsToUse <- apply(tab, 2, function(x){x[1] > 1 && x[2] > 1})
14 datamat <- datamat[,indsToUse]
15 initialdatamat <- datamat
16
17 phy = as.phyDat(datamat, type='USER', levels=c("0","1"),
18   ambiguity="?", names=rownames(data), n=nrow(data), return.
19   index=TRUE)
20
21 dm <- dist.hamming(phy)
22 treeNJ <- NJ(dm)
23 fit <- pml(treeNJ, data=phy)
24 fitJC <- optim.pml(fit, model="GTR", optInv=TRUE, optGamma=
25   TRUE,rearrangement = "NNI")
26
27 pdf("siteMLValues.pdf")
28 plot(density(fitJC$siteLik), col="red", lwd=4, ylim=c(0,
29   0.15), main="ML tree/site likelihoods")
30 siteLikes=c()
31 for(i in 1:300){
32   ## generate some simulated data for the tree
33   d = simSeq(fitJC, l=219)
34   ## calculate the site likelihoods
```

```

29   fit <- pml(tree=fitJC$tree, data=d)
30   ## store the site likelihoods
31   siteLikes = c(siteLikes, fit$siteLik)
32   ## make a plot
33   points(density(fit$siteLik), col="gray", type='l')
34 }
35 points(density(fitJC$siteLik), col="red", lwd=4, type='l')
36 legend("topright", legend=c("real inferred", "simulated true
    tree"), col=c("red" ,"gray"), lwd=2)
37 dev.off()
38
39
40
41
42 dm <- dist.hamming(phy)
43 treeNJ <- NJ(dm)
44 fit <- pml(treeNJ, data=phy)
45 fitJC <- optim.pml(fit, model="GTR", optInv=TRUE, optGamma=
    TRUE, rearrangement = "NNI")
46 fitpars = parsimony(fitJC$tree, data=phy, method = "fitch",
    cost = NULL, site = "site")
47 pdf("siteParsValues.pdf")
48 plot(density(fitpars), col="red", lwd=4, ylim=c(0, 0.4), main
    ="ML trees/Parsimony scores")
49 siteScores=c()
50 i=1
51 for(i in 1:300){
52   ## generate some simulated data for the tree
53   d = simSeq(fitJC, l=219)
54   ## calculate the site likelihoods
55   pars = parsimony(fitJC$tree, data=d, method = "fitch",
    cost = NULL, site = "site")
56   ## store the site likelihoods
57   siteScores = c(siteScores, pars)
58   ## make a plot
59   points(density(pars), col="gray", type='l')
60 }
61 points(density(fitpars), col="red", lwd=4, type='l')
62 legend("topright", legend=c("real inferred", "simulated true
    tree"), col=c("red" ,"gray"), lwd=2)
63 dev.off()

```

## Results

The main results are presented in the two figures below

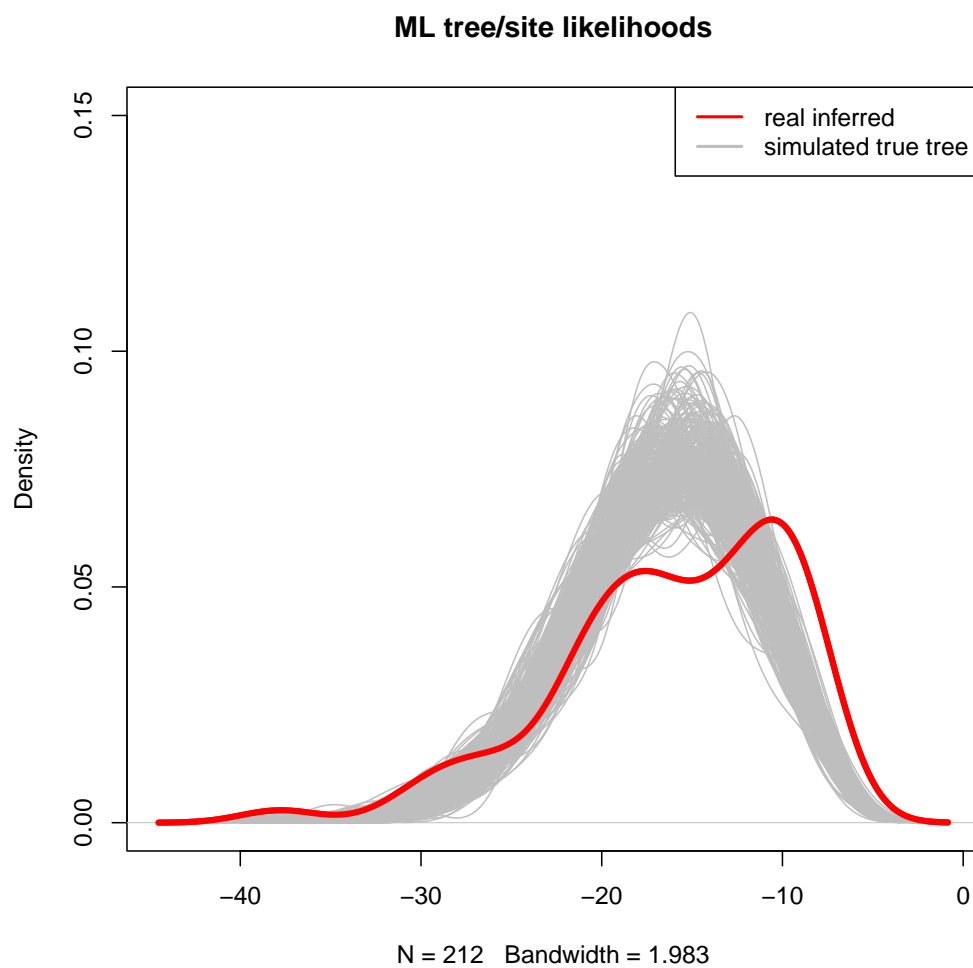


Figure 1: gray color: simulated data on the same ML tree as the inferred tree. Red color: real data on the inferred tree. Each line represents the density plot of the site likelihoods.

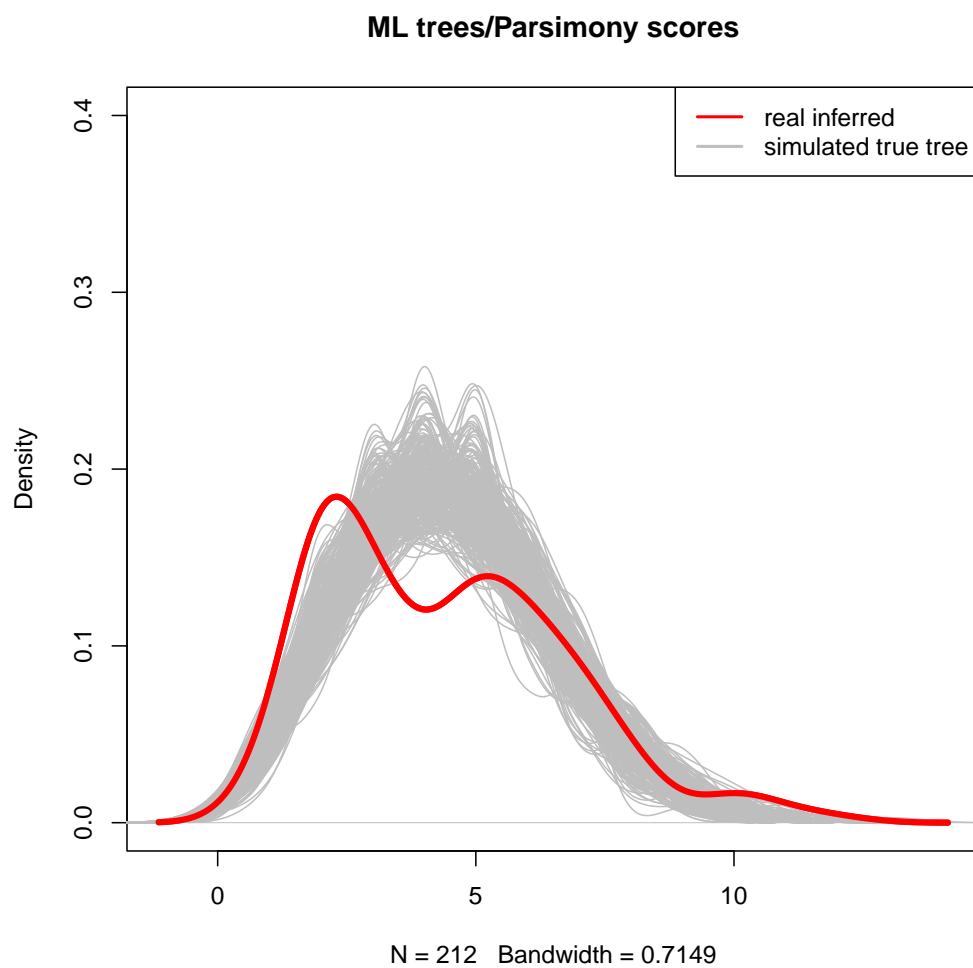


Figure 2: gray color: simulated data on the same ML tree as the inferred tree. Red color: real data on the inferred tree. Each line represents the density plot of the site parsimony scores.

## Conclusion

It seems that the red distribution is qualitatively different than the gray ones. This may mean that the tree inferred from the language data, in fact, might be a non-correct tree. However, I understand that this is not a right way to show this. It's interesting that the distribution shows two peaks. Thus, it might be that in fact I have a mixture of two distributions, i.e., some sites produce a good score and some other a much worse score.