

[Re] Local alignment statistics - Compare to original table

Nathan Brouwer

10/29/2019

Introduction

Every run of a simulation will present different results. Also, I am coding up my version of this simulation based just on the original article. In the following code I generate plots telling me about how close I got to reproducing the original paper and whether 10000 simulations is enough to get a stable estimate of λ and μ . Also, I used a different statistical approach

Load data

Load original version of the table

```
1 table1 <- read.csv(file = "table1.csv")
```

Load my version of the table

```
1 table1.NLB <- read.csv(file = "table1NLB.csv")
```

Load raw data from simulatoins

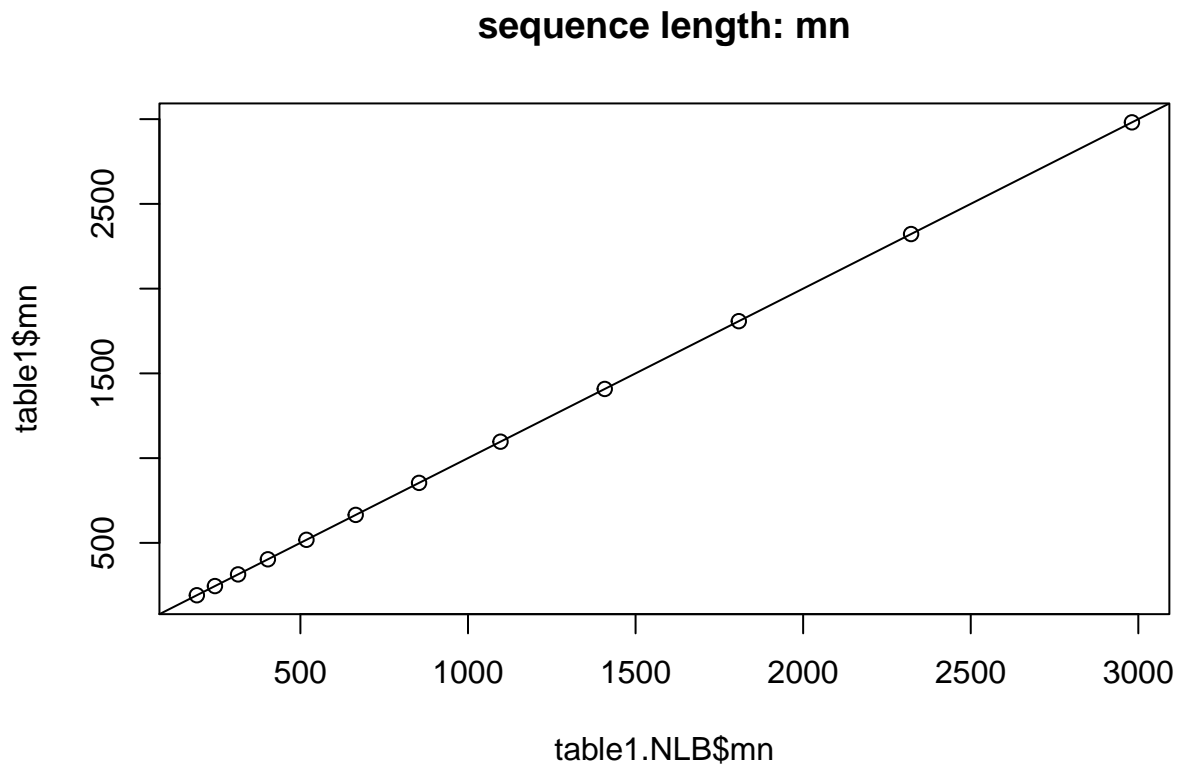
```
1 load("~/1_R/git/blaststats/full_experiment.RData")
```

Make diagnostic plots

Plot the sequence lengths (m,n). THis better work! But always good to make sure a silly typo didn't enter in..

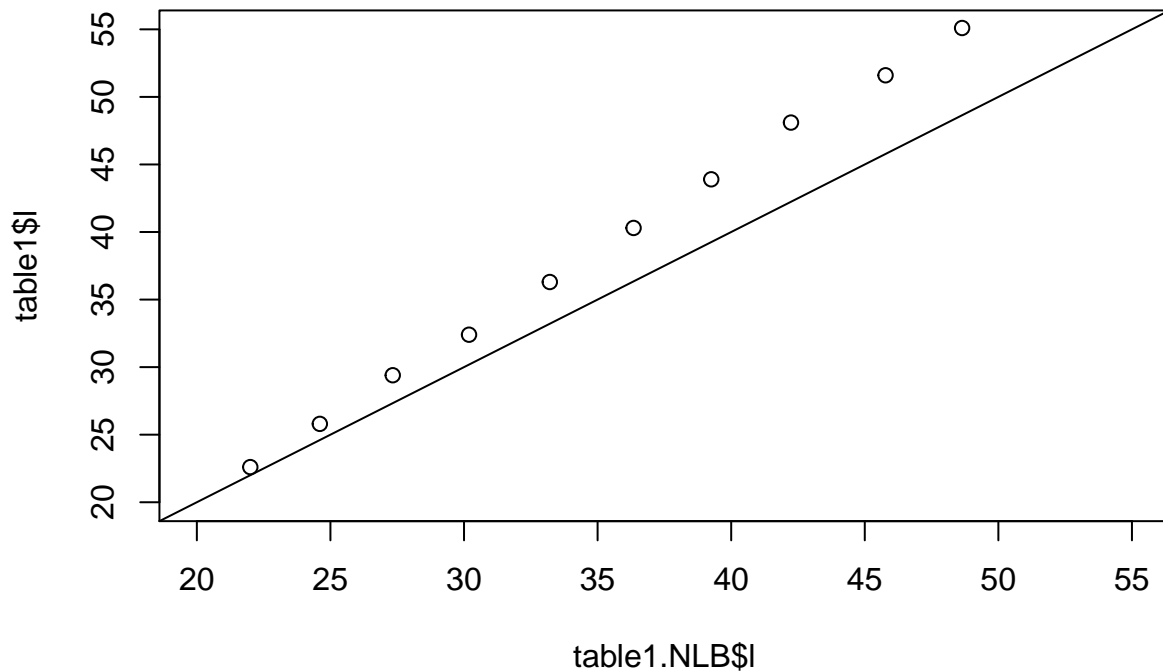
The diagonal line is called a 1:1 and indicates a perfect correlation.

```
1 plot(table1$mn ~ table1.NLB$mn, main = "sequence length: mn")
2 abline(a = 0, b = 1)
```



Plot *alignment* length. This is troubling! Even though highly correlated the points don't fall on the line. Worse, there's **drift**, with the mean lengths of the sequences being close on the short end but the difference grows.

```
1 plot(table1$l ~ table1.NLB$l, xlim = c(20,55),ylim=c(20,55))
2 abline(a = 0, b = 1)
```



```
1 cor(table1$1 , table1.NLB$1)
```

```
1 ## [1] 0.9994748
```

Another way to look at this is to plot sequence length against alignment length. Both of these plots indicate that they were able to build alignments that were somewhat longer than those which I'm building. I'm not sure why this is!

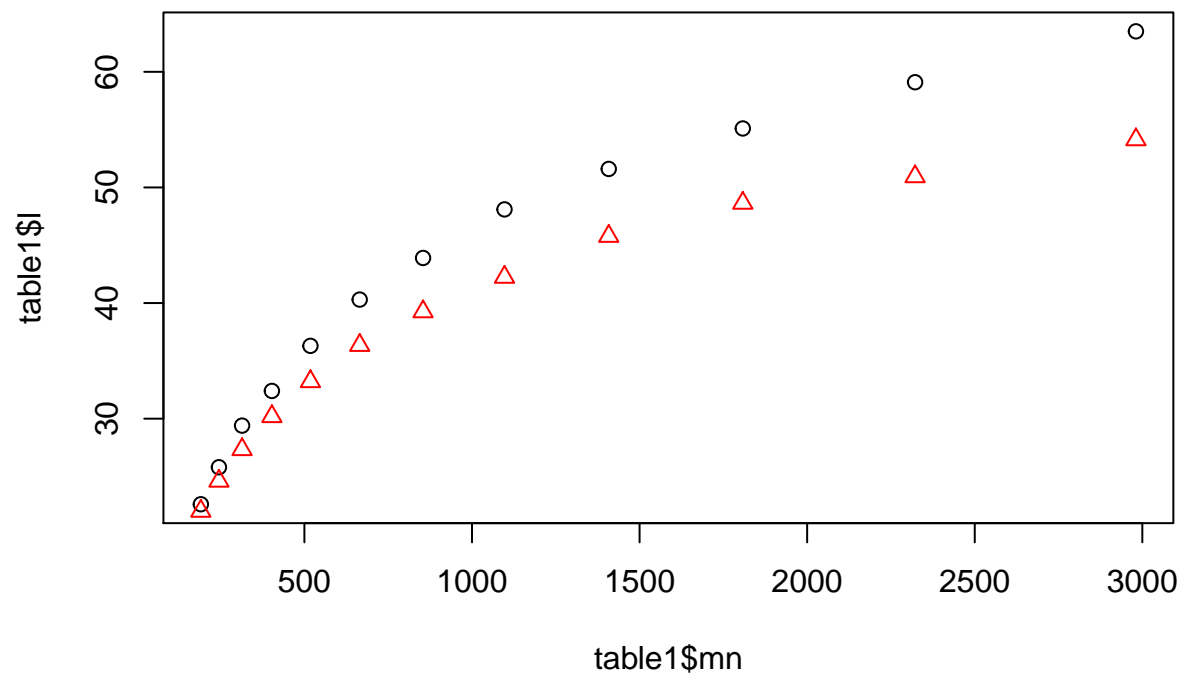
```
1 #plot original data - black circles
```

```
2 plot(table1$1 ~ table1$mn)
```

```
3
```

```
4 #plot my data - read triangels
```

```
5 points(table1.NLB$1 ~ table1.NLB$mn, col = 2, pch =2)
```

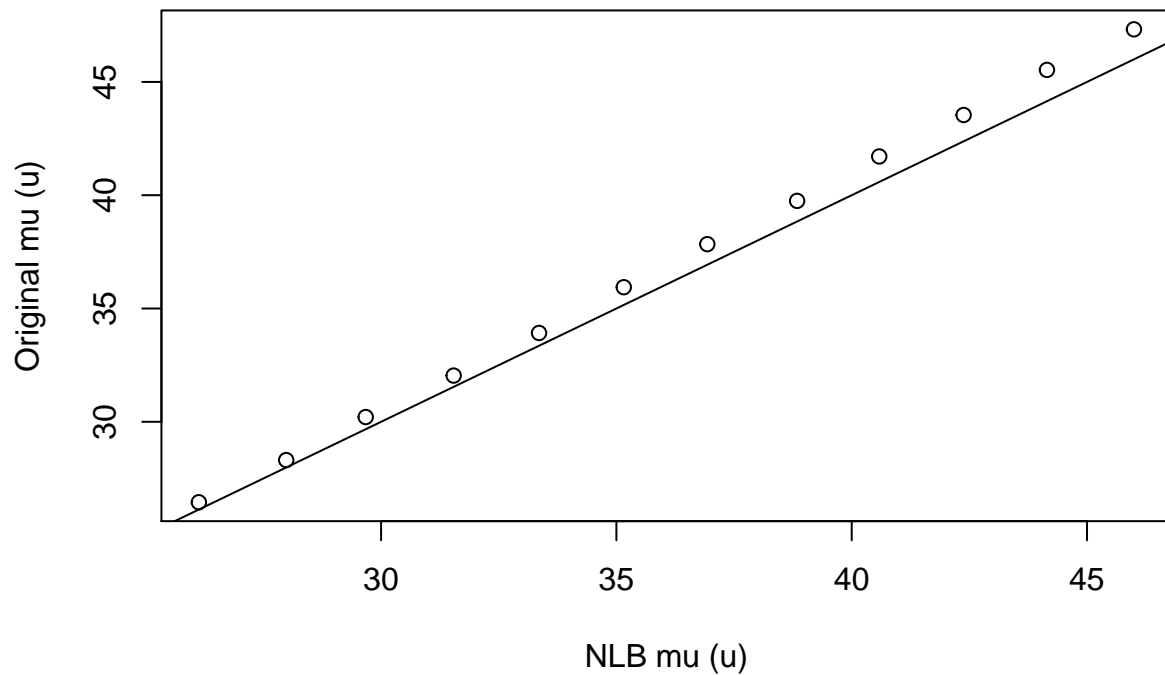


Compare the other paramters. First mu (I call it mu, on their table it u). Same thing, there's drift.

```

1 plot(table1$u ~ table1.NLB$mu,
2       xlab = "NLB mu (u)",
3       ylab = "Original mu (u)")
4 abline(a = 0, b = 1)

```



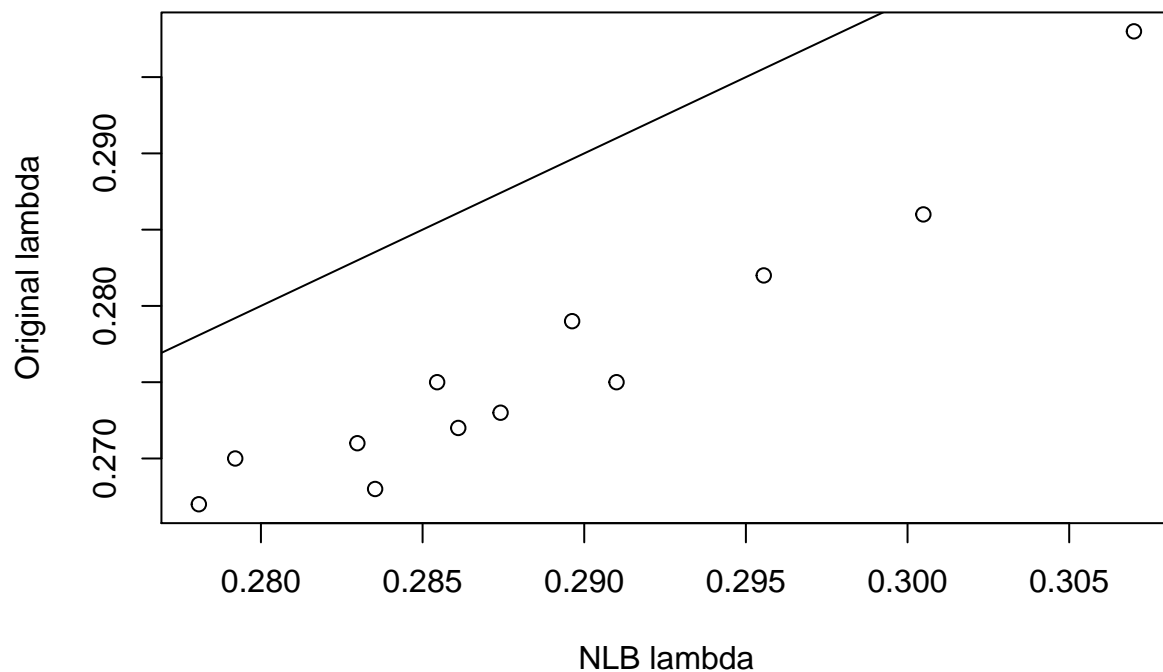
I'll build a quick table to compare the values. There's between a 1 and 3% difference. Not horrible, but not sure why it's happening.

```
1 cbind(table1$mu,
2       table1$u ,
3       table1.NLB$mu,
4       (table1$u - table1.NLB$mu)/table1$u*100)
```

```
1 ##      [,1] [,2]      [,3]      [,4]
2 ## [1,] 191 26.45 26.12708 1.220866
3 ## [2,] 245 28.31 27.98365 1.152786
4 ## [3,] 314 30.21 29.67285 1.778039
5 ## [4,] 403 32.04 31.54146 1.555988
6 ## [5,] 518 33.92 33.35689 1.660106
7 ## [6,] 665 35.94 35.15688 2.178974
8 ## [7,] 854 37.84 36.93168 2.400417
9 ## [8,] 1097 39.75 38.84004 2.289214
10 ## [9,] 1408 41.71 40.58527 2.696556
11 ## [10,] 1808 43.54 42.37834 2.668034
12 ## [11,] 2322 45.53 44.14878 3.033647
13 ## [12,] 2981 47.32 45.99767 2.794443
```

Now lambda - this is kinda ugly!

```
1 plot(table1$lambda ~ table1.NLB$lambda,
2       xlab = "NLB lambda",
3       ylab = "Original lambda")
4 abline(a = 0, b = 1)
```



```

1 cor(table1$lambda , table1.NLB$lambda)
2 cbind(table1$mn,table1$lambda , table1.NLB$lambda)
3 table1$lambda - table1.NLB$lambda

```

```

1 ## [1] 0.9610492
2 ##      [,1] [,2]      [,3]
3 ## [1,] 191 0.298 0.3070028
4 ## [2,] 245 0.286 0.3004887
5 ## [3,] 314 0.282 0.2955524
6 ## [4,] 403 0.275 0.2909991
7 ## [5,] 518 0.279 0.2896259
8 ## [6,] 665 0.273 0.2874119
9 ## [7,] 854 0.272 0.2861039
10 ## [8,] 1097 0.275 0.2854506
11 ## [9,] 1408 0.268 0.2835316
12 ## [10,] 1808 0.271 0.2829821
13 ## [11,] 2322 0.267 0.2780796
14 ## [12,] 2981 0.270 0.2792049
15 ## [1] -0.009002775 -0.014488661 -0.013552359 -0.015999118 -0.010625867
16 ## [6] -0.014411887 -0.014103860 -0.010450583 -0.015531567 -0.011982117
17 ## [11] -0.011079607 -0.009204876

```

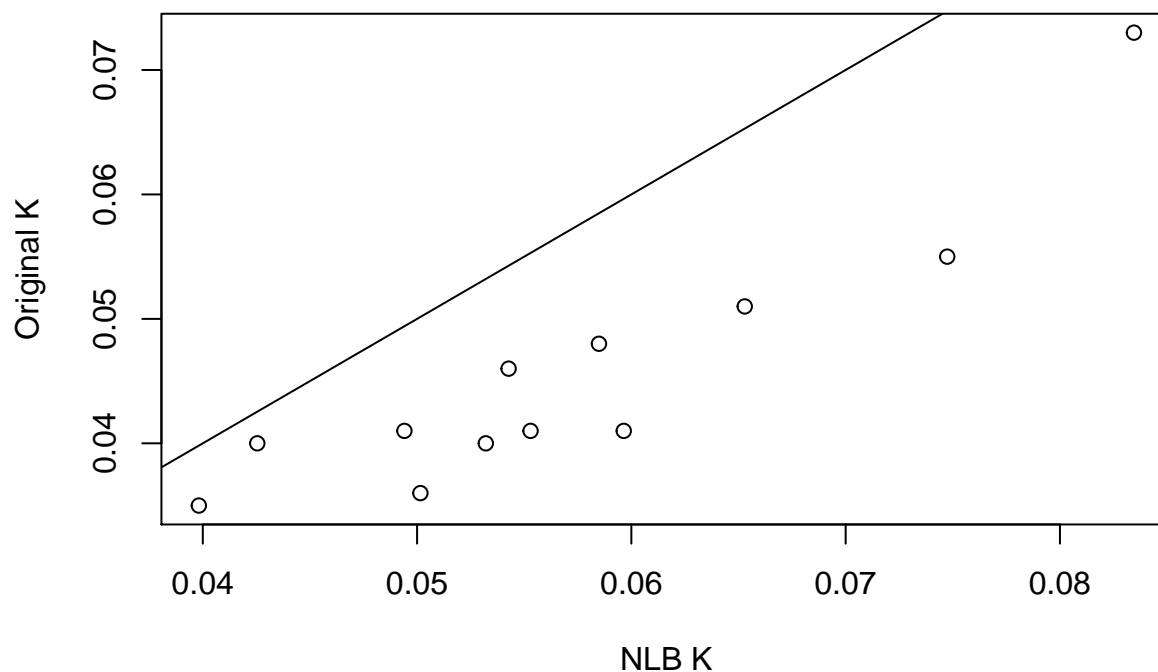
K is pretty bad too.

```

1 plot(table1$K ~ table1.NLB$K,
2       xlab = "NLB K",
3       ylab = "Original K")

```

```
4 abline(a = 0, b = 1)
```



```
1 cor(table1$K , table1.NLB$K)
2 cbind(table1$mn,table1$K, table1.NLB$K)
3 table1$K - table1.NLB$K
```

```
1 ## [1] 0.9151261
2 ##      [,1] [,2]      [,3]
3 ## [1,] 191 0.073 0.08345393
4 ## [2,] 245 0.055 0.07473939
5 ## [3,] 314 0.051 0.06529082
6 ## [4,] 403 0.041 0.05964555
7 ## [5,] 518 0.048 0.05848812
8 ## [6,] 665 0.041 0.05529516
9 ## [7,] 854 0.040 0.05320699
10 ## [8,] 1097 0.046 0.05427142
11 ## [9,] 1408 0.036 0.05015451
12 ## [10,] 1808 0.041 0.04940808
13 ## [11,] 2322 0.035 0.03981575
14 ## [12,] 2981 0.040 0.04254238
15 ## [1] -0.010453931 -0.019739391 -0.014290823 -0.018645554 -0.010488116
16 ## [6] -0.014295165 -0.013206989 -0.008271416 -0.014154507 -0.008408083
17 ## [11] -0.004815749 -0.002542378
```

Upshot - somethings wrong. I need to check that errors in older versions that I identified actually got fixed!.

Method of moments versus MLE

The original paper used the “method of moments” statistical approach to estimate the parameters. The `extRemes` package do this method, but `EnvStats` does not. Both `EnvStats` and `extRemes` do maximum likelihood estimation (MLE), so I can confirm that my results are dependent on the package I use, which is a good idea since estimating parameters of extreme value distributions is not an everyday statistical task and - most importantly - one I am not very familiar with.

The following code fits a bunch of variations in the estimation method, both MLE from two

```
1 #install.packages("EnvStats")
2 library(EnvStats)

##
## Attaching package: 'EnvStats'

## The following objects are masked from 'package:stats':
##
##   predict, predict.lm

## The following object is masked from 'package:base':
##
##   print.default

1 library(extRemes)

## Loading required package: Lmoments
## Loading required package: distillery

##
## Attaching package: 'extRemes'

## The following objects are masked from 'package:EnvStats':
##
##   devd, pevd, qevd, revd

## The following objects are masked from 'package:stats':
##
##   qqnorm, qqplot

1 # methods of moments - EnvStats
2 mme.191 <- eevd(random.scores.191$score.i, method = "mme") # methods of moments
3 mmue.191 <- eevd(random.scores.191$score.i, method = "mmue") # method of moments, unbiased estimator of
4 pwme.191.a <- eevd(random.scores.191$score.i, method = "pwme") # method of moments, probability-weighted
5 pwme.191.b <- eevd(random.scores.191$score.i, method = "pwme",
6                     pwme.method = "plotting.position") # method of moments, probability-weighted vs 2
7
8 # max like - EnvStats
9 mle.191 <- eevd(random.scores.191$score.i, method = "mle") # mle
10
11 # max like - extRemes
12 fevd.191 <- fit.gumbel.191 <- fevd(random.scores.191$score.i,
13                                   type = "Gumbel",
14                                   method = "MLE")
```

I'll eyeball the comparison between statistical methods. They are very close for μ but scale ($1/\lambda$) varies a lot more. There definitely appears to be a difference in between method of moments and max like.


```

1 rbind(eevd.mme = mme.191$parameters,
2       eevd.mmue = mmue.191$parameters,
3       eevd.pwme.a = pwme.191.a$parameters,
4       eevd.pwme.b = pwme.191.b$parameters,
5       eevd.mle = mle.191$parameters,
6       fevd.mle = fevd.191$results$par)

```

```

1 ##           location    scale
2 ## eevd.mme   26.11266 3.303690
3 ## eevd.mmue  26.11256 3.303855
4 ## eevd.pwme.a 26.12827 3.276643
5 ## eevd.pwme.b 26.12776 3.277528
6 ## eevd.mle   26.12708 3.257298
7 ## fevd.mle   26.12708 3.257300

```

How many replications needed?

While not relevant to understanding why my results are a bit off, its important to think about how many iterations of the experiment are necessary to get good estimates of μ , λ and k . As is, it take about 5 hours to run all these simulations. Maybe its possible to use fewer.

In the code take the data for $m = n = 191$ and fit an extreme value distribution model at a range of sample sizes, from 100 to 10,000

```

1 # how many rows in output?
2 n.rows <- nrow(random.scores.191)
3
4
5 # sequence of values - I don't want to do this for all 10,000 rows
6 ## I'll go by multiples of 100
7 i.s <- seq(from = 100, to = nrow(random.scores.191), by = 100)
8
9 #data frame for storage
10 cumulative.params <- data.frame(i = i.s,
11                                loc = rep(NA, length(i.s)),
12                                scale = rep(NA, length(i.s)))
13
14 for(i in 1:length(i.s)){
15   i.s.working <- i.s[i]
16   fit.gumbel.191 <- fevd(random.scores.191$score.i[1:i.s.working],
17                         type = "Gumbel",
18                         method = "MLE")
19
20   loc.param.191.i <- fit.gumbel.191$results$par[1]
21   scale.param.191.i <- 1/fit.gumbel.191$results$par[2]
22
23   cumulative.params$loc[i] <- loc.param.191.i
24   cumulative.params$scale[i] <- scale.param.191.i
25 }

```

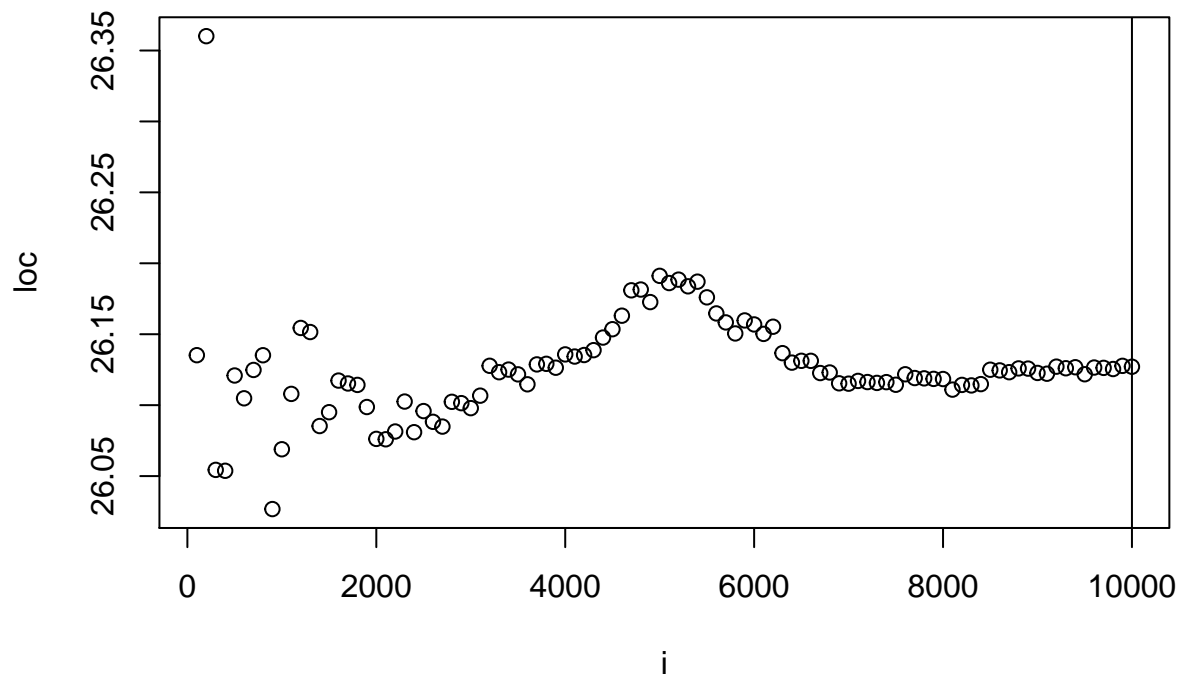
Plot location parameter againsts the number of samples. It gets pretty level by 10, though there does appear to be some drift upward.

```

1 plot(loc ~ i, data = cumulative.params)
2 abline(h = 26.45)

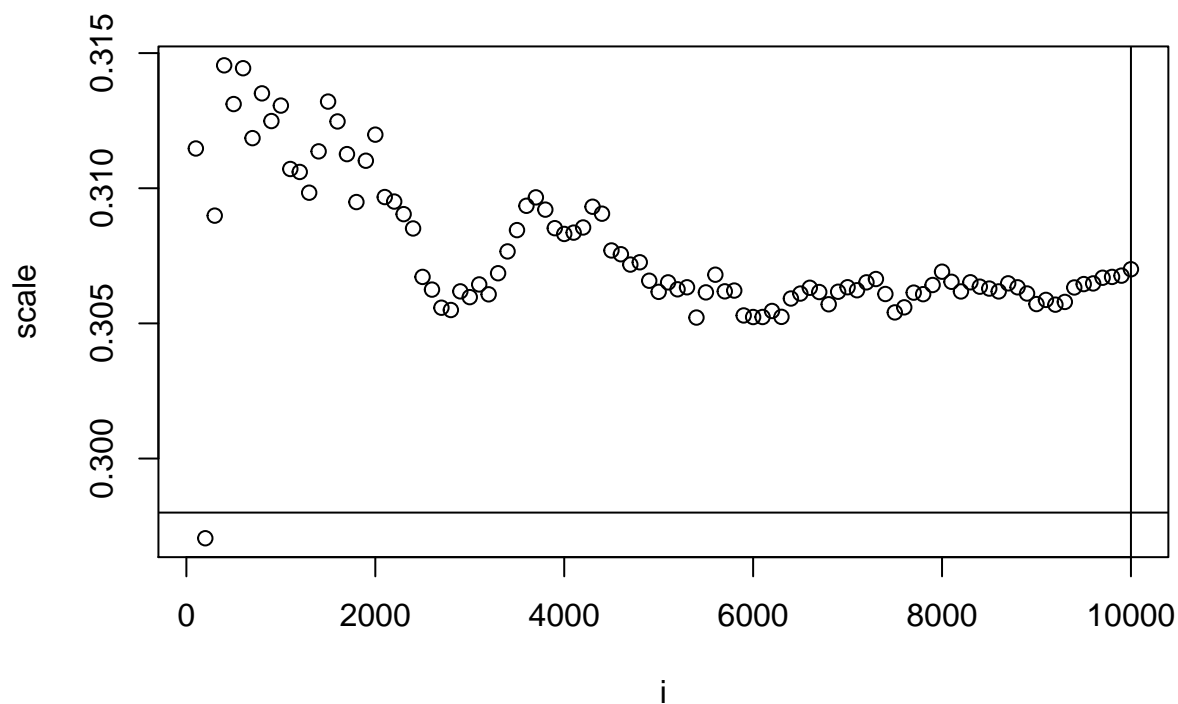
```

```
3 abline(v = 10000)
```



Now the scale parameter. This remains noisier

```
1 plot(scale ~ i, data = cummulative.params)
2 abline(h = 0.298)
3 abline(v = 10000)
```



How many replications needed? $m = n = 2981$

```

1  # how many rows in output?
2  n.rows <- nrow(random.scores.2981)
3
4
5  # sequence of values - I don't want to do this for all 10,000 rows
6  ## I'll go by multiples of 100
7  i.s <- seq(from = 100, to = nrow(random.scores.2981), by = 100)
8
9  #data frame for storage
10 cumulative.params <- data.frame(i = i.s,
11                                loc = rep(NA, length(i.s)),
12                                scale = rep(NA, length(i.s)))
13
14 for(i in 1:length(i.s)){
15   i.s.working <- i.s[i]
16   fit.gumbel.191 <- fevd(random.scores.2981$score.i[1:i.s.working],
17                         type = "Gumbel",
18                         method = "MLE")
19
20   loc.param.191.i <- fit.gumbel.191$results$par[1]
21   scale.param.191.i <- 1/fit.gumbel.191$results$par[2]
22

```

```

23 cumulative.params$loc[i] <- loc.param.191.i
24 cumulative.params$scale[i] <- scale.param.191.i
25 }
26
27
28 par(mfrow = c(2,1))
29 plot(loc ~ i, data = cumulative.params)
30 abline(h = 26.45)
31
32 plot(scale ~ i, data = cumulative.params)
33 abline(h = 0.298)

```

