

RegressHaplo Example

Sivan Leviyang

February 12, 2017

The folder **data/** contains a BAM and associated index file for a paired-end NGS dataset. The reads are synthetic, constructed using the ART simulation package.

Running the Pipeline

Make sure you load RegressHaplo

```
library(RegressHaplo)
```

then run the RegressHaplo pipeline on the BAM file.

```
bam_file <- "data/example.bam"
out_dir <- "output/"
dir.create(out_dir)
full_pipeline(bam_file, out_dir, start_pos=500, end_pos=1500, num_trials=700)
```

We have chosen to restrict reconstruction to reference positions 500 through 1500 and to optimize the penalized regression 700 times. The haplotypes returned will cover the full consensus, composed of 2500 positions, but positions outside the 500 – 1500 region will be set to consensus values. Here we restrict to such a short region for the sake of an example with relatively short run time. On a i7-gen5 machine, it took about 3 minutes to run the code directly above.

Parsing RegressHaplo Results

In this case, we directed RegressHaplo to place output files in the **output/** directory. The final output file is **output/final_haplo.fasta**. The fasta file can be accessed directly; here we use Biostrings to read in and see the sequences.

```
haps <- readDNASTringSet("output/final_haplo.fasta")
haps
```

```
## DNASTringSet object of length 4:
##      width seq                                     names
## [1]  2500 ATGGGATGTCTTGGGAATCAGCT...GTGGAGCTATTTCCATGAGGCGG haplotype1_0.5498
## [2]  2500 ATGGGATGTCTTGGGAATCAGCT...GTGGAGCTATTTCCATGAGGCGG haplotype2_0.2607
## [3]  2500 ATGGGATGTCTTGGGAATCAGCT...GTGGAGCTATTTCCATGAGGCGG haplotype3_0.1246
## [4]  2500 ATGGGATGTCTTGGGAATCAGCT...GTGGAGCTATTTCCATGAGGCGG haplotype4_0.0649
```

But RegressHaplo also provides several functions to view and analyze the reconstruction.

```
# determine the positions on the reference that were considered variable
get_variable_positions.pipeline(out_dir)

## [1]  558  588  634  648  677  705  796  860  973  986 1003 1031 1043 1085 1153
## [16] 1279 1300 1307 1337 1491
```

```

# get haplotype information
info <- get_fasta.pipeline(out_dir)
# info is a list containing the elements haplotypes and freq
# freq gives the frequency of the reconstructed haplotypes
info$freq

## [1] 0.5498 0.2607 0.1246 0.0649

# info$haplotype is a character vector containing the haplotypes.
class(info$haplotypes)

## [1] "character"

length(info$haplotypes)

## [1] 4

# we can use Biostrings to see the haplotypes since the sequences are rather long
DNASTringSet(info$haplotypes)

## DNASTringSet object of length 4:
##      width seq                                     names
## [1]  2500 ATGGGATGTCTTGGGAATCAGCT...GTGGAGCTATTTCCATGAGGCGG haplotype1_0.5498
## [2]  2500 ATGGGATGTCTTGGGAATCAGCT...GTGGAGCTATTTCCATGAGGCGG haplotype2_0.2607
## [3]  2500 ATGGGATGTCTTGGGAATCAGCT...GTGGAGCTATTTCCATGAGGCGG haplotype3_0.1246
## [4]  2500 ATGGGATGTCTTGGGAATCAGCT...GTGGAGCTATTTCCATGAGGCGG haplotype4_0.0649

```

Evaluating Regression Performance

After running RegressHaplo, check if the solutions have captured a range of haplotype reconstructions. The solutions summary data.frame describes the results of the 700 trials. Each row in the data.frame corresponds to a single trial and provides K (the number of reconstructed haplotypes), fit (a fit measure, lower is better), and rho (the penalty parameter). Typically many trials generate the same result - meaning they find the same local minimum.

```

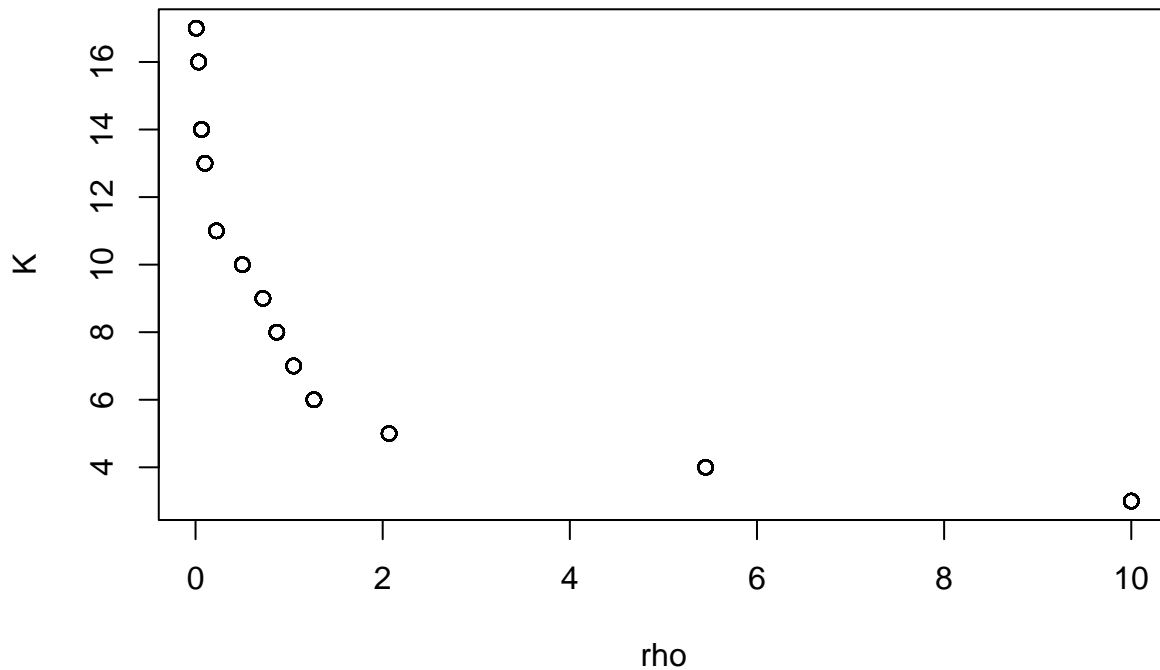
df <- get_solutions_summary.pipeline(out_dir)
# let's look at every 100th solution
df[seq(from=100,to=700,by=100),]

##      rho K      fit solution_number
## 100 0.0330938 16 0.2717402          100
## 200 0.1004420 13 0.2720341          200
## 300 0.5000084 10 0.2742871          300
## 400 0.8685114  8 0.2770060          400
## 500 1.2648552  6 0.2825098          500
## 600 5.4503410  4 0.3195976          600
## 700 10.0000000  3 0.5373616          700

```

Check that the different rho have produced a range of K values

```
plot(df$rho, df$K, xlab="rho", ylab="K")
```



We have captured solutions with K values (the number of haplotypes reconstructed) ranging from 3 to 14, although $K = 6, 10, 11$ are missing. This is a good spread of K values; often a particular K value will be missed.

We can check the ρ values used

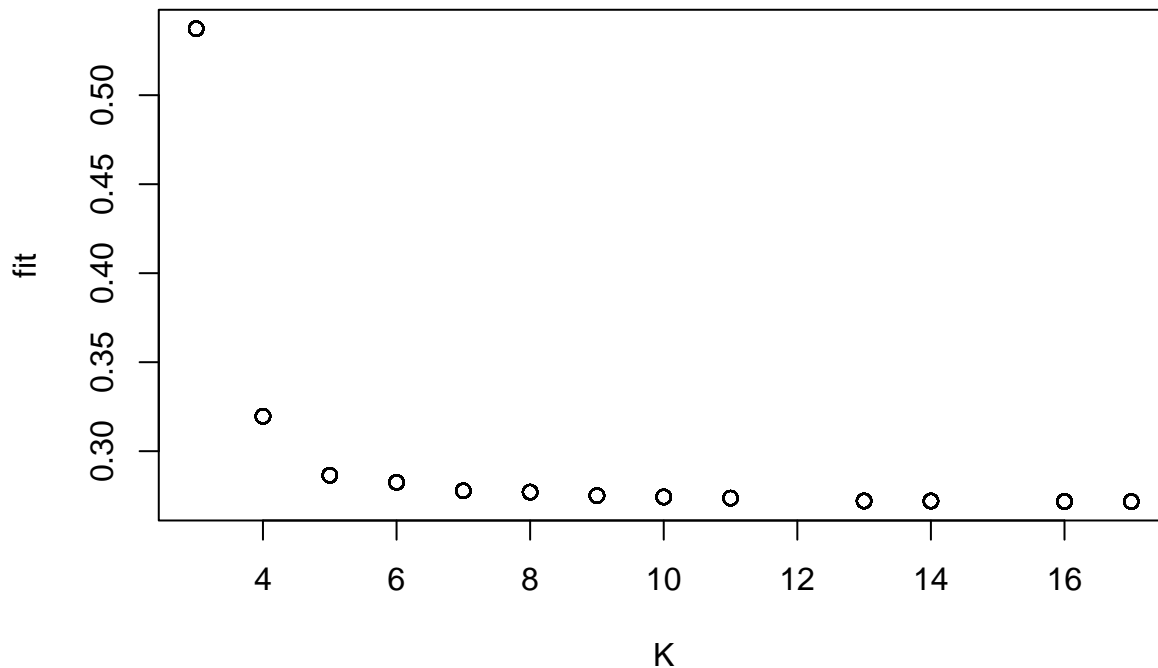
```
unique(df$rho)
```

```
## [1] 0.006662124 0.033093798 0.063243822 0.100441966 0.223186835
## [6] 0.500008391 0.719685673 0.868511374 1.048113134 1.264855217
## [11] 2.068545054 5.450341008 10.000000000
```

If further K values are desired, additional ρ values can be specified through the `full_pipeline` ρ parameter, see help for specifics. Increasing ρ will lower K .

We can also check whether the final solution chosen, $K = 4$, reflects a good tradeoff between over and under fitting.

```
# K values vs fit
plot(df$K, df$fit, xlab="K", ylab="fit")
```



Here $K = 5$ would have produced a better fit, but $K = 4$ seems reasonable.

Finally, we can visualize the accuracy of the haplotype reconstruction in terms of the predicted vs sampled frequencies at each position. To do this, we pick a solution. For example, for $K = 4$, we can pick the solution with the lowest fit.

```
df4 <- dplyr::filter(df, K==4) %>%
  dplyr::arrange(fit)
# the best solution is now the in the first row of df4
solution <- df4$solution_number[1]
solution
```

```
## [1] 598
```

And given a particular solution, we can visualize the frequency of particular nucleotides at each variable position. Here we show A and C, but any combination of A,C,G,T,d,i are possible. The variable positions shown can also be varied, see help.

```
out <- solution_accuracy.pipeline(out_dir,
                                solution,
                                nucs=c("A", "C"),
                                plot=T)
```

