

gd, v. 0.15: Estimating Genetic Diversity and Other Population Genetic Parameters from Aligned DNA Sequence Data

Bernhard Haubold

April 1, 2016

1 Introduction

gd is a computer program that implements routine population genetic analyses. Given a set of aligned sequences, it can compute globally or for sliding windows

1. the number of segregating sites, S ;
2. the number of average pairwise differences, π ;
3. a statistic for testing the neutral model of evolution based on S and π known as Tajima's D , D_T [4].

There are a number of good programs already available to carry out these and many related tasks [3]. gd is intended to combine simplicity with efficiency to take some of the tedium out of genome-scale population genetic analysis.

2 Getting Started

gd was written in C on a computer running Mac OS X; it is intended to run on any UNIX system with a C compiler. However, please contact me at haubold@evolbio.mpg.de if you have problems with the program.

- Unpack the program

```
tar -xvzf gd_XXX.tgz
```

where XXX indicates the version.

- Change into the newly created directory

```
cd Gd_XXX
```

and list its contents

```
ls
```

- Generate gd

```
make
```

- List its options

```
./gd -h
```

- Test the program

```
make test
```

3 Tutorial

This tutorial is intended to demonstrate the usage of `gd` by applying it to simulated data sets. In order to carry out the simulations, you will need to have installed on your computer Richard Hudson's program `ms` [1] and the auxiliary program `ms2dna` by Peter Pfaffelhuber and myself, which is available from

```
http://guanine.evolbio.mpg.de/mlRho/
```

1. Generate one sample of 10 haplotypes with mutation rate $\theta = 4N_e\mu = 100$ and no recombination among 10,000 potentially recombining sites:

```
ms 10 1 -t 100 -r 0 10000 > sample.ms
```

The `-r` option is included here to allow later conversion to DNA sequences 10,000 bp long. `ms` by itself would give the same result without it.

2. Compute standard statistics for this sample using the program `sample_stats`, which is part of the `ms` package:

```
sample_stats < sample.ms
```

3. Convert the haplotypes in `sample.ms` to DNA sequences in FASTA format:

```
ms2dna sample.ms > sample.fasta
```

4. Compute π for the sequences in `sample.fasta` using `gd` and compare the result to that obtained with `sample_stats`:

```
./gd sample.fasta
```

5. Repeat this computation for the number of segregating sites, S

```
./gd -s s sample.fasta
```

and for D_T :

```
./gd -s t sample.fasta
```

6. For each statistic the program can be switched into sliding window mode by using the `-w` option, for example:

```
gd -w 1000 sample.fasta
```

Notice that the distance between windows is 100, that is, one tenth of the window length. This default step length can be changed using the `-S` option. So for printing every window, type

```
gd -w 1000 -S 1 sample.fasta
```

7. The rather verbose output from the last command is usually summarized in a graph and a simple way to draw one is to use the program `graph`, which is part of the GNU `plotutils` package:

```
gd -w 1000 sample.fasta | graph -T X
```

4 Change Log

1. v. 0.4 (April 20, 2010)
 - First version distributed.
2. v. 0.5
 - Fixed polymorphism positions printed using $-P$.
 - Fixed output of sliding window analysis.
3. v. 0.6 (May 3, 2010)
 - Removed printing of error message `no_complete_column_in_window`; now windows without a single complete column are simply not reported.
 - Fixed serious bug in sliding window analysis.
4. v. 0.7 (May 11, 2010)
 - Introduced $-W$ option for setting the minimum number of nucleotides in sliding window.
5. v. 0.8 (October 25, 2010)
 - Fixed comparison of $-W$ option from $>$ to \geq .
 - Fixed positioning of window.
6. v. 0.9 (November 20, 2012)
 - Fixed handling of data with zero polymorphisms.
7. v. 0.10 (December 20, 2012)
 - Removed colon in output with $-s\ s$ (segregating sites).
8. v. 0.11 (April 16, 2013)
 - Fixed $-p$ option.
9. v. 0.12 (October 30, 2013)
 - Fixed treatment of sliding windows in case $-w$ is greater than the length of the alignment.
 - Improved program interface.
10. v. 0.13 (July 10, 2015)
 - Included estimation of the minimum number of recombination events [2].
11. v. 0.14 (July 13, 2015)
 - Sped up computation of minimum number of recombination events.
12. v. 0.15 (July 19, 2015)
 - Fixed computation of the R_M by fixing the function `cmpIntervals`.
 - Made R_m computation optional.

5 Listings

The following listings document central parts of the code for `td`.

5.1 The Driver Program: `gd.c`

```
1  /***** gd.c *****/
* Description: Program to quantify genetic
*   diversity.
* Author: Bernhard Haubold, haubold@evolbio.mpg.de
* Date: Wed Feb 17 13:24:02 2010
6  *****/
#include <stdio.h>
#include <stdlib.h>
#include <fcntl.h>
#include <unistd.h>
11 #include "eprintf.h"
#include "interface.h"
#include "genDiv.h"
#include "rmin.h"

16 void runAnalysis(Args *args, int fd) {
    Sequence *seq;
    Alignment *al;
    long i;
    double s, sum, p, winPos;
21    double *arr;

    seq = readFasta(fd);
    al = seq2al(seq);
    if(args->p) {
        printPoly(al);
        return;
    }
    if(args->w) {
        args->w = args->w < al->n ? args->w : al->n;
31        if(args->W == -1)
            args->W = args->w / 2 + 1;
        if(args->s == 's')
            arr = winSs(args, al);
        else if(args->s == 't')
            arr = winTajima(args, al);
36        else
            arr = winPi(args, al);
        winPos = (args->w-1.)/2. + 1;
        for(i=0;i<al->numWin;i++) {
            if(args->s != 't'){
                if(al->winNumNuc[i] >= args->W)
                    printf("%g\t%.6f\n", winPos, arr[i]/al->winNumNuc[i]);
            }else{
                printf("%g\t%.6f\n", winPos, arr[i]);
46            }
            winPos += args->S;
        }
    }else{
        if(args->s == 's') {
            sum = 0;
51            for(i=1;i<al->m;i++)
                }
```

```

        sum += 1./i;
        s = ss(al);
        printf("S:\t%g\tnumsites:\t%ld\tsite:\t%.6f\ttheta_W/site:\t%.6f",
               s,al->numNuc,s/al->numNuc,s/al->numNuc/sum);
56      if(args->r)
        printf("\tR_m:\t%d",rmin(al));
        printf("\n");
    }else if(args->s == 't'){
        printf("D_T:\t%.6f\n",tajima(al));
61    }else{
        p = pi(al);
        printf("pi:\t%.6f\tnumsites:\t%ld\tpi/site:\t%.6f\ttheta_pi/site:\t%.6f",
               p,al->numNuc,p/al->numNuc,p/al->numNuc);
        if(args->r)
            printf("\tR_m:\t%d",rmin(al));
        printf("\n");
66    }
}
}

71 int main(int argc, char *argv[]) {
    Args *args;
    char *version;
    int fd;
    int i;

76     version = "0.15";
     setprogname2("gd");
     args = getArgs(argc, argv);
     if(args->v)
         printSplash(version);
81     if(args->h || args->e)
         printUsage(version);
     if(args->numInputFiles) {
        for(i=0;i<args->numInputFiles;i++) {
            fd = open(args->inputFiles[i],O_RDONLY,0);
            runAnalysis(args,fd);
            close(fd);
        }
    }else{
91        fd = 0;
        runAnalysis(args,fd);
    }
    free(args);
    free(progname());
96    return 0;
}

```

5.2 Diversity Estimation

5.3 genDiv.h

```

***** genDiv.h *****
* Description:
3  * Author: Bernhard Haubold, haubold@evolbio.mpg.de

```

```

* Date: Tue Feb 16 21:58:07 2010
***** */
#ifndef GENDIV
#define GENDIV
8
#include "sequenceData.h"
#include "interface.h"

typedef struct alignment{
13    char **al;           /* nucleotide alignment */
    long m;              /* number of sequences in alignment */
    long n;              /* number of nucleotides per sequence in alignment */
    char **headers;      /* headers of sequences */
    long *poly;          /* polymorphic positions in alignment */
18    long *nuc;          /* positions consisting solely of canonical nucleotides
                           */
    char *polInd;        /* is position polymorphic? */
    char *nucInd;        /* does position consist of nucleotides only? */
    int *winNumNuc;      /* number of canonical nucleotides per window */
    int numWin;          /* the number of windows */
23    float *polyFreq;   /* frequency of minor allele */
    long numPoly;        /* number of polymorphic positions in alignment */
    long numNuc;         /* number of positions consisting solely of canonical
                           nucleotides */
}Alignment;

28 Alignment *seq2al(Sequence *seq);
    double pi(Alignment *al);
    double ss(Alignment *al);
    double *winPi(Args *args, Alignment *al);
    double *winSs(Args *args, Alignment *al);
33    double *winTajima(Args *args, Alignment *al);
    void printPoly(Alignment *al);
    void findPoly(Alignment *al);
    double tajima(Alignment *al);

38 #endif

```

5.4 genDiv.c

```

***** genDiv.c *****
2 * Description: Routines for measuring genetic
* diversity.
* Author: Bernhard Haubold, haubold@evolbio.mpg.de
* Date: Tue Feb 16 21:33:24 2010
***** */

7 #include <stdio.h>
#include <stdlib.h>
#include <assert.h>
#include <math.h>
#include <limits.h>
12 #include "interface.h"
#include "genDiv.h"
#include "queue.h"
#include "eprintf.h"

```

```

17 void prepareWin(Alignment *al, int winLen, int stepLen);

double td(double a1, double a2, double s, double p, int n);

void printPoly(Alignment *al) {
22   int i, j, c;

   if(al->poly == NULL)
     findPoly(al);
   /* Print positions of polymorphisms */
   printf(">Positions_frequencies:%ld\n",al->numPoly);
27   for(i=0;i<al->numPoly;i++)
     printf("%ld\t%.3f\n",al->poly[i]+1,al->polyFreq[i]);
   /* Print alleles at polymorphic positions */
   for(i=0;i<al->m;i++) {
32     printf("%s\n",al->headers[i]);
     c=0;
     for(j=0;j<al->numPoly;j++) {
       printf("%c",al->al[i][al->poly[j]]);
       c++;
       if(c==60) {
         printf("\n");
         c = 0;
       }
     }
     if(c)
       printf("\n");
   }
}

47 /* pi: number of average pairwise differences per site */
double pi(Alignment *al) {
  int i, j, k;
  double s1, s2;

  if(al->poly == NULL) {
    findPoly(al);
  }
  s2 = 0;
  for(i=0;i<al->m-1;i++) {
57    for(j=i+1;j<al->m;j++) {
      s1 = 0;
      for(k=0;k<al->numPoly;k++) {
        if(al->al[i][al->poly[k]] != al->al[j][al->poly[k]])
          s1++;
      }
      s2 += s1;
    }
  }
62 }

67 s2 /= al->m * (al->m - 1) / 2;

return s2;

```

```

    }

/* tajima: Tajima's D */
72 double tajima(Alignment *al) {
    double s, p, a1, a2, n;
    int i;

    s = ss(al);
77    p = pi(al);
    n = al->m;

    a1 = 0;
    a2 = 0;
82    for(i=1;i<n;i++) {
        a1 += 1.0/i;
        a2 += (1.0/i/i);
    }

87    return td(a1, a2, s, p, n);
}

/* ss: number of segregating sites per site */
double ss(Alignment *al) {
92
    if(al->poly == NULL)
        findPoly(al);

    return (double)al->numPoly;
97 }

double *winPi(Args *args, Alignment *al) {
    int i, j, k, numWin;
    int lb, rb;
102   double *arr, *pol, factor;
    double sum;

    arr = (double *)emalloc(al->n*sizeof(double));
    pol = (double *)emalloc(al->n*sizeof(double));
107   if(al->poly == NULL)
        findPoly(al);
    if(al->polInd == NULL)
        prepareWin(al,args->w,args->S);
    /* compute average number of mismatches for all positions */
112   factor = al->m*(al->m-1)/2;
    for(i=0;i<al->n;i++) {
        pol[i] = 0;
        if(al->polInd[i]){
            for(j=0;j<al->m-1;j++)
                for(k=j+1;k<al->m;k++)
                    if(al->al[j][i] != al->al[k][i])
                        pol[i]++;
                    pol[i] /= factor;
            }
        }
122   numWin = 0;
}

```

```

/* take care of first window */
rb = 0;
sum = 0;
127 while(rb < args->w && rb < al->n)
    sum += pol[rb++];
arr[numWin++] = sum;
/* scan remaining windows */
lb = 0;
132 while(rb < al->n-args->S+1) {
    for(i=0;i<args->S;i++) {
        sum += pol[rb+i];
        sum -= pol[lb+i];
    }
    rb += args->S;
    lb += args->S;
    arr[numWin++] = sum;
}
free(pol);
142 return arr;
}

/* winSs: sliding window analysis of segregating sites */
double *winSs(Args *args, Alignment *al) {
147 int i, numWin, numPol, rb, lb;
double *arr;

arr = (double *)emalloc(al->n*sizeof(double));
if(al->poly == NULL)
    findPoly(al);
152 if(al->polInd == NULL)
    prepareWin(al,args->w,args->S);
numWin = 0;
/* take care of first window */
rb = 0;
157 numPol = 0;
while(rb < args->w && rb < al->n)
    numPol += al->polInd[rb++];
arr[numWin++] = numPol;
/* scan remaining windows */
162 lb = 0;
while(rb < al->n-args->S+1) {
    for(i=0;i<args->S;i++) {
        if(al->polInd[rb+i])
            numPol++;
        if(al->polInd[lb+i])
            numPol--;
    }
    rb += args->S;
167 lb += args->S;
    arr[numWin++] = numPol;
}
172 return arr;
}

```

177

```

/* td: compute Tajima's D */
double td(double a1, double a2, double s, double p, int n) {
    double b1, b2, c1, c2, e1, e2;
    double d;
182
    b1=(n+1)/3./(n-1);
    b2=2.* (n*n+n+3.)/(9.*n*(n-1));
    c1=b1-1./a1;
    c2=b2-(n+2)/(a1*n)+a2/a1/a1;
187
    e1=c1/a1;
    e2=c2/(a1*a1+a2);

    if (s>0)
        d=(p-s/a1)/sqrt(e1*s+e2*s*(s-1));
192
    else
        d=0.;

    return d;
}
197
double *winTajima(Args *args, Alignment *al) {
    long i;
    double *pArr, *sArr, *dArr;
    double a1, a2;
202
    a1 = 0;
    a2 = 0;
    for (i=1;i<al->m;i++) {
        a1 += 1.0/i;
        a2 += (1.0/i/i);
207
    }

    pArr = winPi(args, a1);
    sArr = winSs(args, a1);
212
    dArr = (double *)emalloc((al->n-args->w+1)*sizeof(double));
217
    for (i=0;i<al->n-args->w+1;i++)
        dArr[i] = td(a1, a2, sArr[i], pArr[i], al->m);
    }
    return dArr;
}

Alignment *seq2al(Sequence *seq) {
222
    int i;
    Alignment *al;

    al = (Alignment *)emalloc(sizeof(Alignment));
    al->m = seq->numSeq;
227
    al->headers = seq->headers;
    al->n = seq->borders[0];
    al->al = (char **)emalloc(al->m*sizeof(char *));
    al->al[0] = seq->seq;
    for (i=1;i<seq->numSeq;i++)
        al->al[i] = seq->seq[i];
}
```

```

232     al->al[i] = seq->seq + seq->borders[i-1] + 1;

al->numPoly = 0;
al->poly = NULL;
al->polInd = NULL;
237 al->nucInd = NULL;
return al;
}

/* findPoly: fills an array of positions that are polymorphic
242 *      and consist solely of the four canonical nucleotides.
*/
void findPoly(Alignment *al) {
    int i, j, c, p;
    int *dic;
247
    dic = getRestrictedDnaDictionary(NULL);
    al->poly = (long *)emalloc(al->n*sizeof(long));
    al->nuc = (long*)emalloc(al->n*sizeof(long));
    al->polyFreq = (float *)emalloc(al->n*sizeof(float));

252
    al->numPoly = 0;
    al->numNuc = 0;
    for(i=0;i<al->n;i++) {
        p = 0;
        if(dic[(int)al->al[0][i]])
            c = 1;
        else{
            c = 0;
            continue;
        }
        for(j=1;j<al->m;j++) {
            if(!dic[(int)al->al[j][i]]){
                c = 0;
                break;
            }
            if(al->al[0][i] != al->al[j][i])
                p = 1;
        }
        if(c) {
            al->nuc[al->numNuc++] = i;
            if(p) {
                al->poly[al->numPoly] = i;
                al->polyFreq[al->numPoly] = 1.;
                for(j=1;j<al->m;j++)
                    if(al->al[j][i] == al->al[0][i])
                        al->polyFreq[al->numPoly]++;
                al->polyFreq[al->numPoly] /= (float)al->m;
                if(al->polyFreq[al->numPoly] > 0.5)
                    al->polyFreq[al->numPoly] = 1. - al->polyFreq[al->numPoly];
                al->numPoly++;
            }
        }
    }
}

```

```

157     al->poly = (long *)erealloc(al->poly,al->numPoly*sizeof(long));
158     al->nuc = (long *)erealloc(al->nuc,al->numNuc*sizeof(long));
159     if(al->numPoly)
160         al->polyFreq = (float *)erealloc(al->polyFreq,al->numPoly*sizeof(float))
161             );
162     free(dic);
163 }

/* prepareWin: prepare sliding window analysis
 *      should be preceded by findPoly, though this
167 *      is checked for
 */
void prepareWin(Alignment *al, int winLen, int stepLen){
    int i, j, rb, lb, numNuc;

170     if(al->poly == NULL)
        findPoly(al);
    al->polInd = (char *)emalloc(al->n*sizeof(char));
    al->nucInd = (char *)emalloc(al->n*sizeof(char));

177     /* mark canonical nucleotides */
    j = 0;
    for(i=0;i<al->numNuc;i++) {
        while(j<al->nuc[i])
            al->nucInd[j++] = 0;
182        al->nucInd[j++] = 1;
    }
    for(i=j;i<al->n;i++)
        al->nucInd[i] = 0;

187     /* mark polymorphisms */
    j = 0;
    for(i=0;i<al->numPoly;i++) {
        while(j<al->poly[i])
            al->polInd[j++] = 0;
192        al->polInd[j++] = 1;
    }
    for(i=j;i<al->n;i++)
        al->polInd[i] = 0;

200     /* count canonical nucleotides per window */
    al->winNumNuc = (int *)emalloc(al->n*sizeof(int));
    /* take care of first window */
    rb = 0;
    numNuc = 0;
    al->numWin = 0;
    while(rb < winLen && rb < al->n)
        numNuc += al->nucInd[rb++];
    al->winNumNuc[al->numWin++] = numNuc;
    /* scan remaining windows */
    lb = 0;
    while(rb < al->n-stepLen+1) {
        for(i=0;i<stepLen;i++) {

```

```

    if(al->nucInd[rb+i])
        numNuc++;
    if(al->nucInd[lb+i])
        numNuc--;
342    }
    rb += stepLen;
    lb += stepLen;
    al->winNumNuc[al->numWin++] = numNuc;
347}
}

```

5.5 rmin.c

```

/****** rmin.c *****/
2   * Description: Calculating the minimum number of
*   recombination events using the algorithm
*   in Appendix 2 of Hudson & Kaplan (1985).
*   Reference: Hudson, R. R. and Kaplan, N. L.
*   (1985). Statistical properties of the number
7   * of recombination events in the history of a
*   sample of DNA sequences. Genetics, 111:147-
*   164.
*   Author: Bernhard Haubold, haubold@evolbio.mpg.de
*   Date: Fri Jul 10 10:00:43 2015
12  *****/
#include <stdio.h>
#include <stdlib.h>
#include "eprintf.h"
#include "rmin.h"
17 #include "genDiv.h"

/* numAlleles: number of alleles in sample when considering
*   positions i and j
*/
22 short numAlleles(Alignment *a, int i, int j){
    int k, allele[4], x;

    allele[0] = 1;
    for(k=1;k<4;k++)
27        allele[k] = 0;
    for(k=1;k<a->m;k++) {
        x = 0;
        if(a->al[k][i] != a->al[0][i])
            x += 2;
32        if(a->al[k][j] != a->al[0][j])
            x += 1;
        allele[x] = 1;
    }
    x = 1;
37        for(k=1;k<4;k++)
            x += allele[x];

return x;
}

```

42

```

Interval *newInterval(int i, int j){
    Interval *in;
    in = (Interval *)emalloc(sizeof(Interval));
    in->start = i;
    in->end = j;

    return in;
}

int cmpIntervals(const void *a, const void *b){

    Interval *ia = *(Interval * const *)a;
    Interval *ib = *(Interval * const *)b;
    return (int)(ia->start - ib->start);
}

int sorted(Interval **interval, int n){
    int i;

    for(i=1;i<n;i++)
        if(interval[i]->start < interval[i-1]->start)
            return 0;

    return 1;
}

int rmin(Alignment *al){
    int i, j, n, end, rm, numInt, c;
    short enclosing;
    Interval **intervals1, **intervals2;

    if(al->poly == NULL)
        findPoly(al);

    n = al->numPoly;
    intervals1 = (Interval **)emalloc(n*(n-1)/2*sizeof(Interval *));
    c = 0;
    for(i=0;i<n-1;i++)
        for(j=i+1;j<n;j++)
            if(numAlleles(al,al->poly[i],al->poly[j]) == 4) {
                if(al->poly[i] < al->poly[j])
                    intervals1[c++] = newInterval(al->poly[i],al->poly[j]);
                else
                    intervals1[c++] = newInterval(al->poly[j],al->poly[i]);
            }
    if(c)
        intervals1 = (Interval **)erealloc(intervals1,c*sizeof(Interval *));
    else
        intervals1 = (Interval **)erealloc(intervals1,sizeof(Interval *));
    intervals2 = (Interval **)emalloc(c*sizeof(Interval *));
    qsort(intervals1, c, sizeof(Interval *), cmpIntervals);
    numInt = 0;
}

```

```

97     for(i=0;i<c;i++) {
  enclosing = 0;
  for(j=i+1;j<c;j++) {
    if(intervals1[i]->start <= intervals1[j]->start && intervals1[i]->end
        >= intervals1[j]->end) {
      enclosing = 1;
      break;
    }
  }
  if(!enclosing) {
    for(j=0;j<numInt;j++) {
      if(intervals1[i]->start <= intervals2[j]->start && intervals1[i]->
          end >= intervals2[j]->end) {
        enclosing = 1;
        break;
      }
    }
    if(!enclosing)
      intervals2[numInt++] = intervals1[i];
    }
  }
  if(numInt) {
    rm = 1;
    end = intervals2[0]->end;
  }else
    rm = 0;
  for(i=1;i<numInt;i++) {
    if(intervals2[i]->start == intervals2[i-1]->start)
      end = intervals2[i]->end;
    if(intervals2[i]->start >= end) {
      rm++;
      end = intervals2[i]->end;
    }
  }
  for(i=0;i<c;i++)
    free(intervals1[i]);
  free(intervals1);
  free(intervals2);
  return rm;
}

```

References

- [1] R. R. Hudson. Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioinformatics*, 18:337–338, 2002.
- [2] R. R. Hudson and N. L. Kaplan. Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics*, 111:147–164, 1985.
- [3] J. Rozas, J. C. Sánchez-DelBarrio, X. Messeguer, and R. Rozas. DnaSP, DNA polymorphism analyses by the coalescent and other methods. *Bioinformatics*, 19:2496–2497, 2003.
- [4] F. Tajima. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics*, 123:585–595, 1989.