

## UP - UP, distance between points

### Tracks

1. Track 1: unmarked points
2. Track 2: unmarked points

### Questions

Where in the genome are the points in track 1 closer to/further apart from points in track 2 than expected by chance?

Comment:

- We assume points in track 2 as fixed and want to find out whether points in track 1 are closer to or further apart from the closest point in track 2 than expected. The test may indicate that the two tracks are independent. The test is not symmetric in the two tracks.
- Significance is determined by means of p-values. Small p-values identify regions where the points in track 1 are closer to or further apart from the closest point in track 2 than expected. P-values are computed as explained below, where the null hypothesis is explained in detail.
- The p-values may be found by simulation or by an approximate calculation. It is necessary to specify a distribution of the unmarked points in track 1.

### Bins

The genome (or the areas of the genome under study) are divided into small regions, called bins. The tests are performed in each bin.

### Hypothesis tested

For each bin  $i$  we have the four different null hypotheses corresponding to each of the four alternative preservation rules given below:

**H<sub>0</sub>:** *Assume points in track 1 are independent of points in track 2*

with the following alternative hypotheses:

**H<sub>1</sub>:** *Points in track 1 are closer to points in track 2 than expected or*

**H<sub>2</sub>:** *Points in track 1 are further apart from points in track 2 than expected.*

Let  $g(r)$  be the point in track 2 that is closest to the point  $r$  in track 1. Define the distance  $d(r)$  as the distance between the position of  $r$  and the position of  $g(r)$  (see Figure 1). Let  $r_1, \dots, r_n$  be the points in track 1 in bin  $i$ , and let  $\hat{\mu} = \frac{1}{n} \sum_{j=1}^n d(r_j)$  be the mean distance between points the tracks 1 and it's nearest point in track 2. In all tests the points in track 2 will be considered as fixed, the points in track 1 as random and  $\hat{\mu}$  will be used as test statistic.

The  $H_0$  hypothesis is rejected for each bin  $i$  if:  $\hat{\mu}_i > c_{\alpha,i}$  or  $\hat{\mu}_i < d_{\alpha,i}$  or  $c_{\alpha/2,i} < \hat{\mu}_i < d_{\alpha/2,i}$  corresponding to the average distance is significantly larger/smaller/different than expected. The critical values  $c_{\alpha,i}$  and  $d_{\alpha,i}$  are found by simulation and depend on the threshold  $\alpha$  and the bin  $i$ .

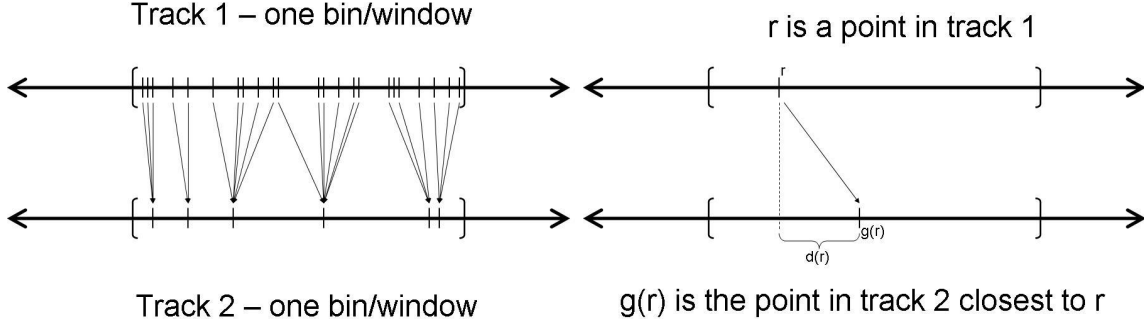


Figure 1: *Comparing positions in track 1 and 2.*

We may assume four different preservation rules for the distribution of points in track 1. These give different null distributions for  $\hat{\mu}$  and hence different test results. In all four cases we use Monte-Carlo simulation for obtaining samples of track 1 under the null hypothesis. For each sample of track 1, the corresponding  $\hat{\mu}$  is computed and the distribution of  $\hat{\mu}$  under the null hypothesis is obtained. How to sample the points of track 1 under each of the four different preservation rules is described below.

**Preservation rule 1: Preserve the number of points in the bin in track 1** Assume track 1 has  $n$  points. The locations of the  $n$  points are drawn independently and uniformly in the bin.

**Preservation rule 2: Preserve the number of points and also the interpoint distances in the bin in track 1** The points in track 1 are sampled by permuting the interpoint distances of the original track 1.

**Preservation rule 3: Preserve the distribution of the interpoint distances in the bin in track 1** The leftmost point might be drawn by drawing a distance  $d$  from the distribution  $D$  of the interpoint distances, and then draw the distance from the bin start to the first point from the uniform distribution  $U[0,d]$ . The next points in track 1 are sampled one by one from left to right by drawing the interpoint distances from the distribution  $D$ . We stop drawing new points when the next point would have been placed outside the bin.

If a control track is available the four sampling procedure above might be extended as indicated in the note "Sampling MC-locations from the candidate track".

### Approximation under preservation rule 1

For preservation rule 1 we may, alternatively, use an approximation for the null distribution of  $\hat{\mu}$  as described below.

**Assume that the number of points in the bin in track 1 is preserved.** Let  $D_1, \dots, D_n$  be independently, identically distributed random variables for the distances of the points in track 1,  $d(r_1), \dots, d(r_n)$ .

The locations for the  $n$  points in track 1 are independent. Let  $f$  be the prior on possible locations for one point. In the special case that  $f$  is uniform, we observe that the distribution of each of  $D_1, \dots, D_n$  is a mixture of non-overlapping uniform distributions i.e. the

distribution  $f_D$  is a piecewise constant distribution (Figure 2):

$$f_D(d) = \sum_{i=1}^m c_i \cdot U[b_{i-1}, b_i],$$

where  $m$ ,  $a_i$  and  $b_i$ ,  $i = 1, \dots, m$ , are as indicated in Figure 2,  $b_0 = 0$  and  $c_i$  is the fraction of the bin that is covered by  $a_i$  segments.  $b_1$  is the shortest half-distance,  $b_2$  the next shortest etc. and  $a_i = b_i - b_{i-1}$ .

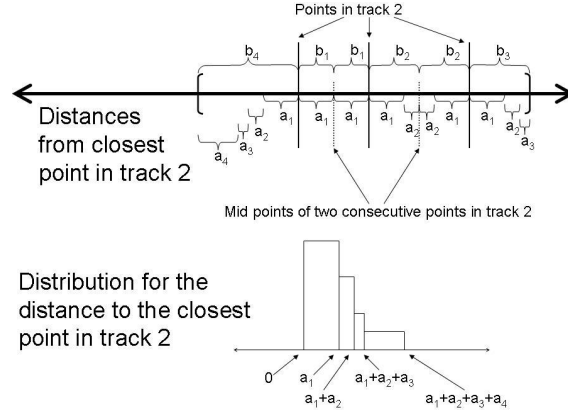


Figure 2: *Distribution of the distances of the points in track 1 for a bin with three points in track 2 assuming a uniform prior.*

When the prior on possible locations is not a uniform, the distribution  $f_D$  for the distance to the closest point in track 2 is obtained as follows: Let  $m$  be the largest possible distance and define a function  $g : \{0, \dots, m\} \rightarrow R$  by  $g(i) = \sum_{\text{location } l \text{ with } d(l)=i} f(l)$ . Then  $f_D(i) = \frac{g(i)}{\sum_{j=0}^m g(j)}$ . Also in this case  $f_D$  is a piecewise constant distribution, but each interval with constant values is very short. To obtain longer intervals we might approximate  $f_D(i)$  with another piecewise constant distribution, f.ex. by repeatedly merging some neighbour intervals with quite similar values into a new interval with constant value equal to the mean of the original values.

The null distribution for  $\mu$ , the mean distance for the points in track 1, may be approximated by a piecewise constant density. The density may be found by adding one and one of the terms in the sum (this should be done in such a way that as few additions as possible are performed). When the number of points in track 1 is large, the null distribution for  $\mu$  might be approximated by a normal distribution. Small p-values are obtained when  $\hat{\mu}$  computed from the data occurs in the left/right/left or right tail of the null-distribution, corresponding to the tests of whether the average distance is significantly larger/smaller/different than expected.