

Association of Genomic Features with Integration

Charles C. Berry

July 5, 2007

Contents

1	Introduction	2
2	Preference for Genes	4
2.1	refGene Genes	4
2.2	ensGenes	6
2.3	genScan Genes	8
2.4	oncogenes	10
3	CpG Island Neighborhoods	12
3.1	1 kilobase neighborhoods	12
3.2	5 kilobase neighborhoods	13
3.3	10 kilobase neighborhoods	14
3.4	25 kilobase neighborhoods	14
3.5	50 kilobase neighborhoods	15
4	Gene Density, Expression 'Density', and CpG Island Density	17
4.1	25 kilobase Window	17
4.2	50 kilobase Window	22
4.3	100 kilobase Window	27
4.4	250 kilobase Window	32
4.5	500 kilobase Window	37
4.6	1 megabase Window	42
4.7	2 megabase Window	47
4.8	4 megabase Window	52
4.9	8 megabase Window	57
4.10	16 megabase Window	62
4.11	32 megabase Window	67
5	Juxtaposition with Gene Start and End Positions	73
5.1	Refseq Annotations	73
5.2	genScan Annotations	77
6	Oncogenes	81

7 GC content	83
8 Cytobands	100

1 Introduction

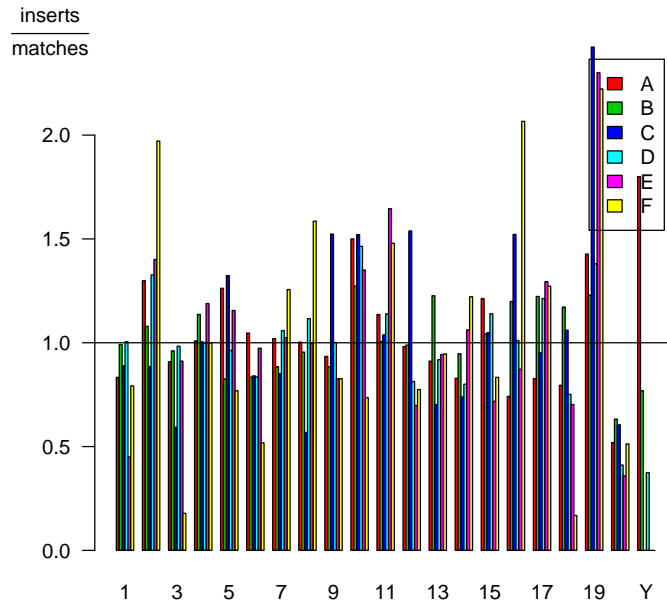
In this document, I examine the association of integration sites with various genomic features.

The data consist of both actual integration sites and sets of control sites, each set chosen to match the spacing (in bases) from the nearest restriction site (according to the direction in which the sequence was read) to an integration site. The numbers of insertion and matching sites for several data sets are shown below:

Origin.of.data.set	type	
	insertion	match
A	2441	7308
B	1554	15440
C	374	3740
D	4828	14475
E	478	4780
F	320	3200

The advantage of choosing 'control' sites that match the spacing from the nearest restriction site is that biases due to location and density of restriction sites are eliminated by applying the classical multinomial logit model (reviewed in [2]). This model allows regression procedures to be applied to the study of integration intensity as a function of genomic features. The `clogit` function of the R `survival` library) implements estimation and fitting for such models along with the usual likelihood ratio and Wald tests.

The distribution of relative frequency of insertions across the chromosomes is given in this barplot:

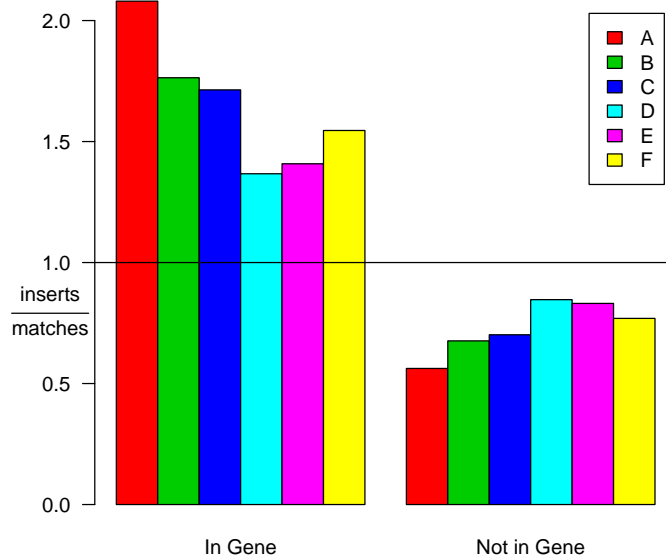


It seems evident that there are some chromosomes that are particularly favored for integration. This is reinforced by a test of statistical significance. The test performed used the likelihood ratio statistic for the multinomial logit model (reviewed in [2]) as implemented by the `clogit` function of the R `survival` library). The null hypothesis tested is that the ratio of true integration events to matched control sites is constant across all chromosomes. This test attains a p-value of $< 2.22e - 16$.

2 Preference for Genes

2.1 refGene Genes

Here we examine the preference that integration events have for genes. In the following plot we show the relative frequency of integrations in genes according to the 'refGene' annotation. The bars grouped over the label "In Gene" give the relative frequency of integration events (compared to control sites) between bases located within refGene gene annotations, while the label "Not in Gene" give the relative frequency of integration events (compared to control sites) between bases not located within refGene gene annotations.



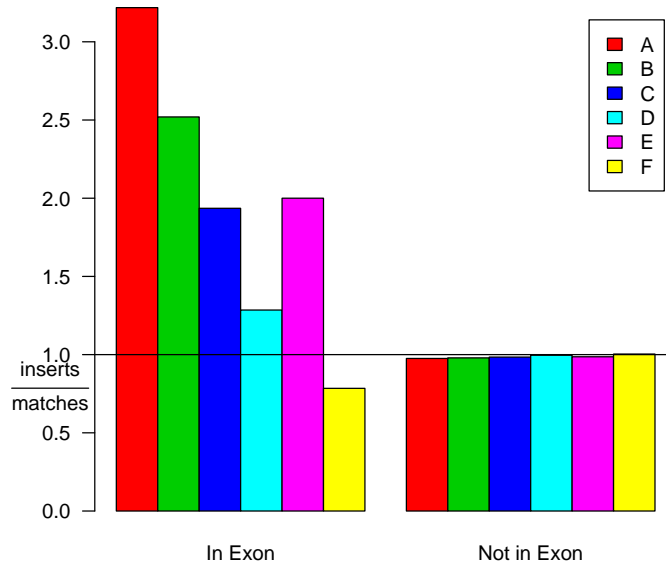
It seems evident that there is a strong tendency for insertions to occur in genes. A formal test of significance bears this out with a p-value of $< 2.22e-16$. Also, it appears that the tendency of different viruses to integrate into genes varies, and a test for this hypothesis attains $< 2.22e-16$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

	coef	se	z	p
A	1.320	0.0506	26.10	5.52e-150
B	0.967	0.0540	17.90	1.29e-71
C	0.882	0.1090	8.10	5.62e-16

D 0.475 0.0345 13.80 4.11e-43
 E 0.527 0.0983 5.37 8.06e-08
 F 0.700 0.1190 5.89 3.95e-09

As is evident, there are some differences in the coefficients. The largest coefficient is seen in the A data set, while the smallest is seen in the D data set.

In the following plot we show the relative frequency of insertions in exons according to the 'refGene' annotation. The bars grouped over the label "In Exon" give the relative frequency of integration events (compared to control sites) between bases located in exons according to the Acembly annotation, while the label "Not in Exon" give the relative frequency of integration events (compared to control sites) between bases not located in exons according to the Acembly gene annotation.



Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

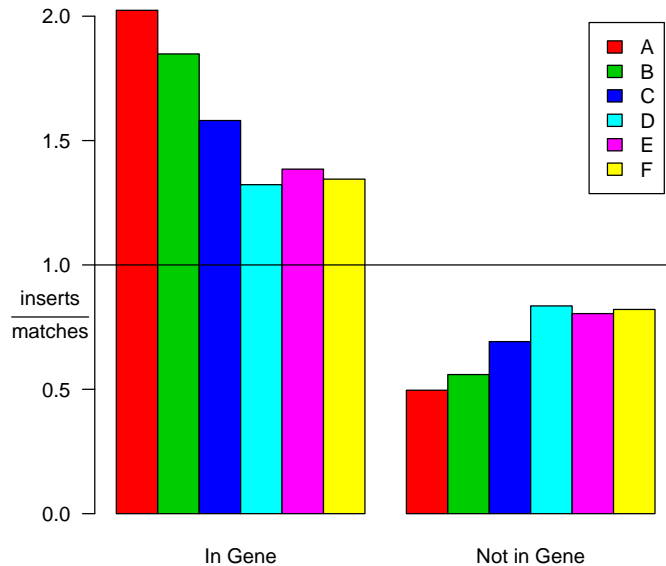
	coef	se	z	p
A	0.5660	0.163	3.470	0.000522
B	0.3640	0.159	2.290	0.021800
C	0.1790	0.326	0.549	0.583000
D	-0.0655	0.131	-0.499	0.618000
E	0.3850	0.315	1.220	0.222000

F -0.7310 0.529 -1.380 0.167000

The model on which these coefficients are based include terms for whether the site is in a gene or not. Thus, the effect shown as "In Exon" is net of that due to being in a gene. Note that in the barplot above the 'Not in Exon' bars include both the introns and intergenic regions, so the impression given by the table may differ from that for the barplot.

2.2 ensGenes

Here we examine the preference that insertions have for genes. In the following plot we show the relative frequency of insertions in genes according to the 'ensGene' annotation.



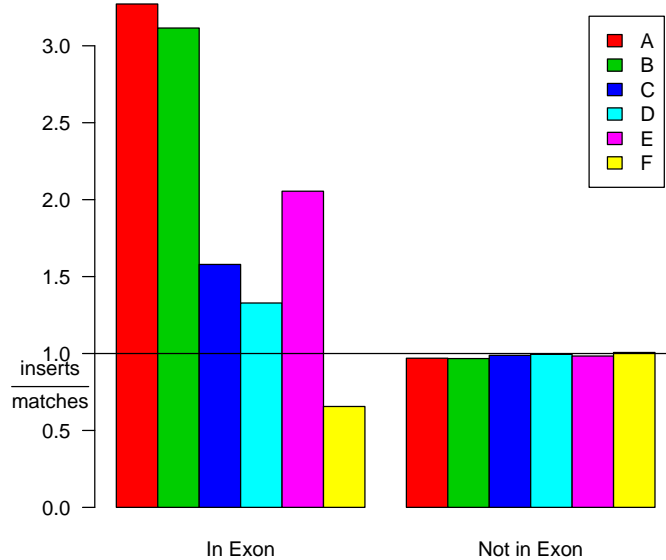
It seems evident that there is a strong tendency for insertions to occur in genes. A formal test of significance bears this out with a p-value of $< 2.22e-16$. Also, it appears that the tendency of different viruses to integrate into genes varies, and a test for this hypothesis attains $< 2.22e-16$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

	coef	se	z	p
A	1.420	0.0517	27.50	3.28e-166

B	1.210	0.0558	21.70	2.05e-104
C	0.827	0.1100	7.54	4.78e-14
D	0.456	0.0338	13.50	1.90e-41
E	0.545	0.0968	5.63	1.78e-08
F	0.492	0.1180	4.16	3.13e-05

As is evident, there are some differences in the coefficients. The largest coefficient is seen in the A data set, while the smallest is seen in the D data set.

In the following plot we show the relative frequency of insertions in exons according to the 'ensGene' annotation.



Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

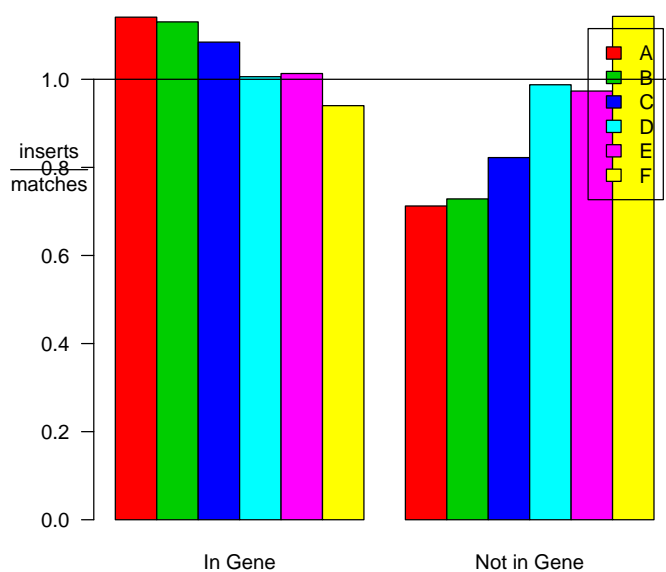
	coef	se	z	p
A	0.55400	0.147	3.7600	1.69e-04
B	0.55100	0.138	3.9900	6.57e-05
C	0.03140	0.320	0.0979	9.22e-01
D	0.00932	0.120	0.0780	9.38e-01
E	0.42900	0.294	1.4600	1.45e-01
F	-0.75800	0.525	-1.4400	1.49e-01

The model on which these coefficients are based include terms for whether

the site is in a gene or not. Thus, the effect shown as "In Exon" is net of that due to being in a gene.

2.3 genScan Genes

Here we examine the preference that insertions have for genes. In the following plot we show the relative frequency of insertions in genes according to the 'genScan' annotation.

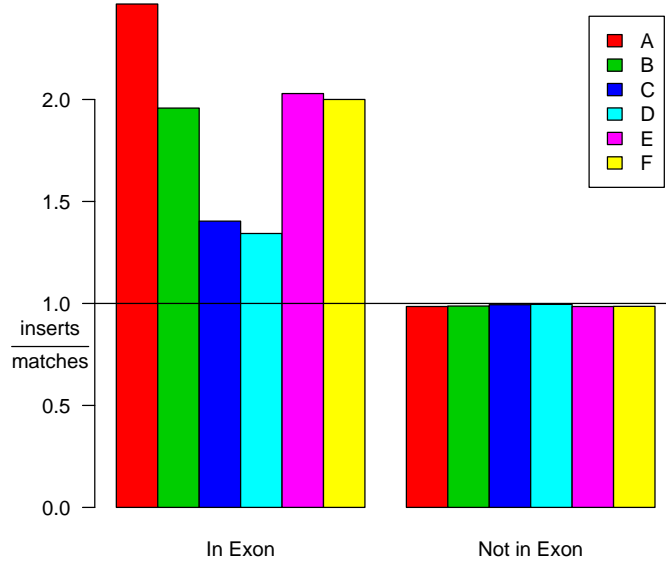


It seems evident that there is a strong tendency for insertions to occur in genes. A formal test of significance bears this out with a p-value of $1.3959e-15$. Also, it appears that the tendency of different viruses to integrate into genes varies, and a test for this hypothesis attains $4.9227e-16$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

	coef	se	z	p
A	0.4790	0.0542	8.840	$9.86e-19$
B	0.4460	0.0623	7.160	$8.32e-13$
C	0.2790	0.1230	2.270	$2.31e-02$
D	0.0202	0.0356	0.569	$5.70e-01$
E	0.0404	0.1030	0.392	$6.95e-01$
F	-0.1960	0.1250	-1.570	$1.16e-01$

As is evident, there are some differences in the coefficients. The largest coefficient is seen in the A data set, while the smallest is seen in the F data set.

In the following plot we show the relative frequency of insertions in exons according to the 'genScan' annotation.



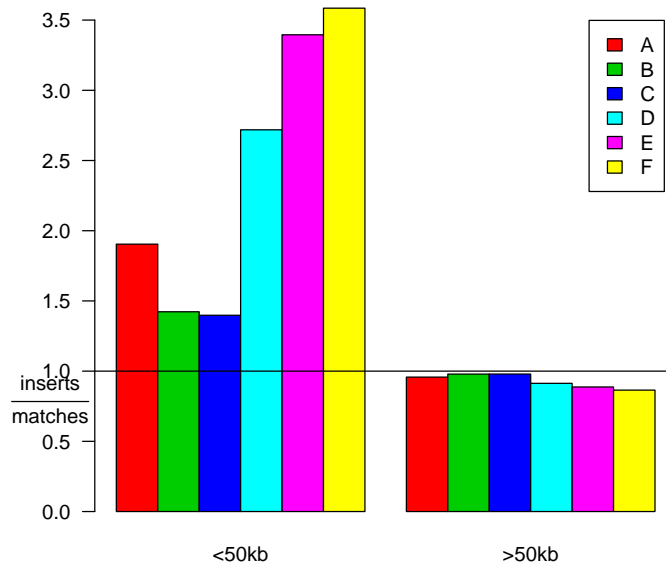
Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

	coef	se	z	p
A	0.828	0.179	4.630	3.68e-06
B	0.566	0.176	3.220	1.27e-03
C	0.263	0.381	0.692	4.89e-01
D	0.294	0.131	2.250	2.47e-02
E	0.733	0.303	2.420	1.56e-02
F	0.771	0.370	2.080	3.73e-02

The model on which these coefficients are based include terms for whether the site is in a gene or not. Thus, the effect shown as "In Exon" is net of that due to being in a gene.

2.4 oncogenes

Here we examine the preference that insertions have for oncogenes. In the following plot we show the relative frequency of insertions with 50kb of an oncogene 5' end.



A formal test of oncogenic insertion returns p-value of $< 2.22e - 16$. The tendency of different viruses to integrate near oncogenes may vary, and a test for this hypothesis attains $7.3792e - 14$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

	coef	se	z	p
A	-0.678	0.0917	-7.40	1.35e-13
B	-0.372	0.1070	-3.47	5.26e-04
C	-0.355	0.2170	-1.64	1.01e-01
D	-1.090	0.0584	-18.70	3.02e-78
E	-1.370	0.1480	-9.26	2.14e-20
F	-1.460	0.1710	-8.51	1.81e-17
A	NA	0.0000	NA	NA
B	NA	0.0000	NA	NA
C	NA	0.0000	NA	NA
D	NA	0.0000	NA	NA
E	NA	0.0000	NA	NA

F NA 0.0000 NA NA

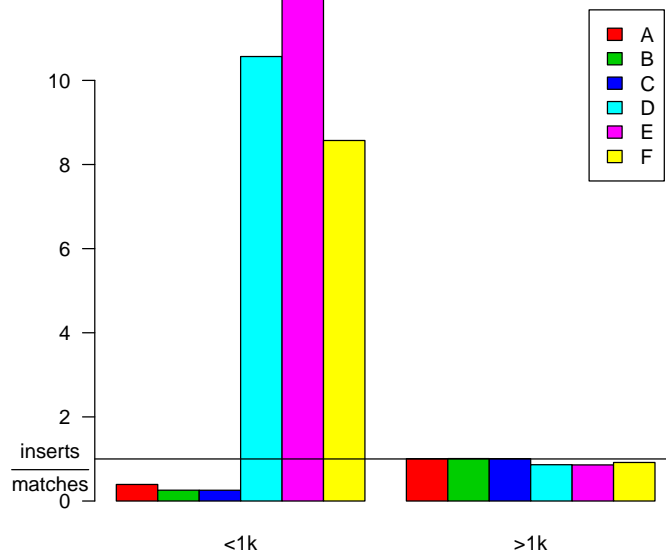
As is evident, there are some differences in the coefficients. The largest coefficient is seen in the C data set, while the smallest is seen in the F data set.

3 CpG Island Neighborhoods

Here we study the effect of being in the neighborhood of CpG Islands. Following Wu et al [3], who found that the neighborhoods within $\pm 1\text{kb}$ of CpG islands are enriched for MLV insertions, we study such neighborhoods.

3.1 1 kilobase neighborhoods

The following plot shows the effect of being in or within $\pm 1\text{kb}$ of a CpG island:



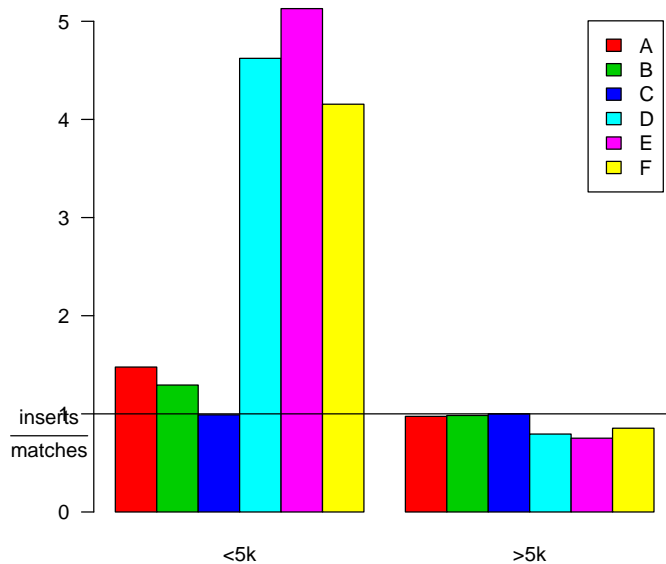
A formal test of significance comparing the difference attains a p-value of $< 2.22e - 16$. A test for differences between viruses attains $< 2.22e - 16$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

	coef	se	z	p
A	-0.947	0.3390	-2.80	5.19e-03
B	-1.360	0.4540	-3.00	2.71e-03
C	-1.360	1.0100	-1.35	1.78e-01
D	2.490	0.0844	29.50	1.09e-190
E	2.640	0.1870	14.10	1.96e-45
F	2.220	0.2570	8.63	6.35e-18

The largest coefficient is seen in the E data set, while the smallest is seen in the C data set.

3.2 5 kilobase neighborhoods

The following plot shows the effect of being in or within ± 5 kb of a CpG island:



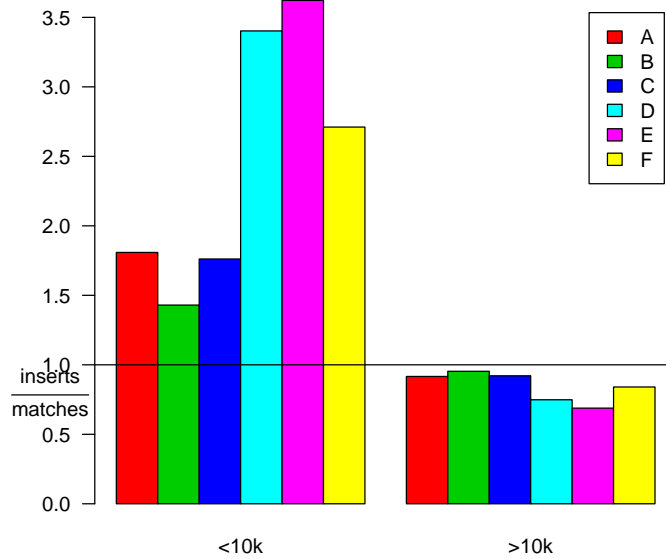
A formal test of significance comparing the difference attains a p-value of $< 2.22e - 16$. A test for differences between viruses attains $< 2.22e - 16$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

	coef	se	z	p
A	0.4120	0.0929	4.4400	9.11e-06
B	0.2820	0.1070	2.6400	8.23e-03
C	-0.0111	0.2470	-0.0448	9.64e-01
D	1.7800	0.0523	34.0000	8.54e-254
E	1.9400	0.1220	15.9000	3.34e-57
F	1.5700	0.1680	9.3500	8.55e-21

The largest coefficient is seen in the E data set, while the smallest is seen in the C data set.

3.3 10 kilobase neighborhoods

The following plot shows the effect of being in or within ± 10 kb of a CpG island:



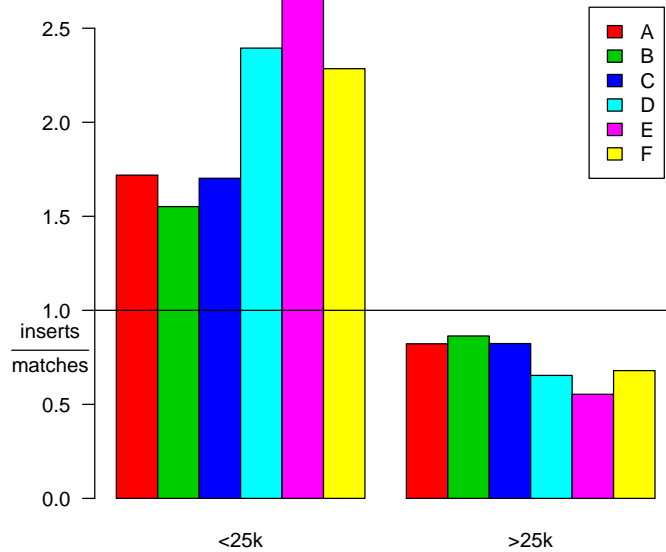
A formal test of significance comparing the difference attains a p-value of $< 2.22e - 16$. A test for differences between viruses attains $< 2.22e - 16$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

	coef	se	z	p
A	0.681	0.0676	10.10	7.28e-24
B	0.402	0.0782	5.14	2.73e-07
C	0.644	0.1490	4.31	1.61e-05
D	1.530	0.0439	34.80	8.27e-266
E	1.670	0.1060	15.70	2.04e-55
F	1.170	0.1480	7.94	2.07e-15

The largest coefficient is seen in the E data set, while the smallest is seen in the B data set.

3.4 25 kilobase neighborhoods

The following plot shows the effect of being in or within ± 25 kb of a CpG island:



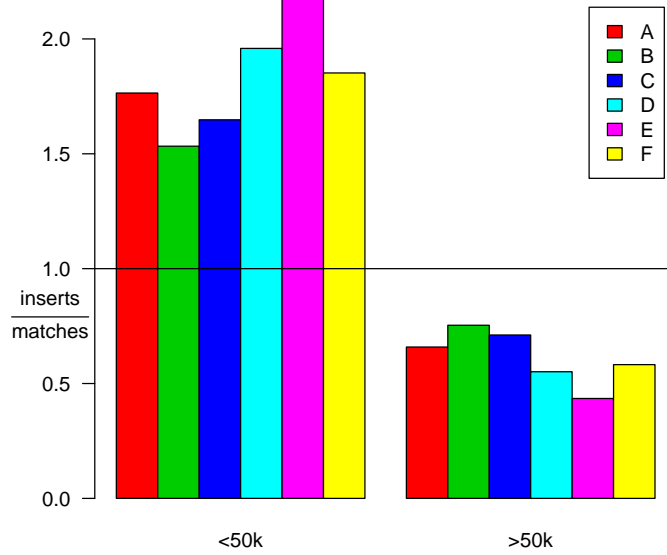
A formal test of significance comparing the difference attains a p-value of $< 2.22e - 16$. A test for differences between viruses attains $< 2.22e - 16$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

	coef	se	z	p
A	0.734	0.0521	14.10	4.33e-45
B	0.587	0.0586	10.00	1.32e-23
C	0.730	0.1170	6.25	4.12e-10
D	1.310	0.0370	35.30	1.04e-272
E	1.580	0.1000	15.80	3.48e-56
F	1.230	0.1220	10.10	8.67e-24

The largest coefficient is seen in the E data set, while the smallest is seen in the B data set.

3.5 50 kilobase neighborhoods

The following plot shows the effect of being in or within ± 50 kb of a CpG island:



A formal test of significance comparing the difference attains a p-value of $< 2.22e - 16$. A test for differences between viruses attains $< 2.22e - 16$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites along with their standard errors, z statistics, and p-values for each data set:

	coef	se	z	p
A	0.985	0.0488	20.20	8.86e-91
B	0.716	0.0538	13.30	1.81e-40
C	0.856	0.1110	7.74	9.83e-15
D	1.280	0.0359	35.50	1.74e-276
E	1.620	0.1060	15.30	5.57e-53
F	1.170	0.1210	9.60	8.08e-22

The largest coefficient is seen in the E data set, while the smallest is seen in the B data set.

4 Gene Density, Expression 'Density', and CpG Island Density

In this section the association with gene density is examined. For expression analysis, the 'genes' that are counted are the genes represented on the microarray. In addition, we the number of such genes expressed at various levels. The levels are

low.ex Count genes whose expression is in the upper half and divide by number of bases

med.ex Count genes whose expression is in the upper $1/8^{th}$ and divide by number of bases

high.ex Count genes whose expression is in the upper $1/16^{th}$ and divide by number of bases

The bolded terms are used as abbreviations in what follows. The abbreviation **dens** is used to indicate gene density as number of genes per base.

4.1 25 kilobase Window

In the barplot that follows we examine the association of insertion sites with gene density in a 25 kilobase window surrounding each locus. More such plots will follow and the method of their construction is always to try to divide the data according to the deciles of density. However, it often happens that there is a very skewed distribution of density and even the 90^{th} percentile is zero. In that case, the barplots simply show the sites for which the density is zero and those for which it is non-zero. If there are fewer than ten groups of bars, the groupings contain ten percent of the sites each except for the leftmost grouping which will contain all of the remaining sites.

Also note that the title of the plot contains clues as to its content; the prefix indicates the type of variable studied while the suffix indicates the window width in the number of bases. The p-value given is the result of fitting a cubic polynomial to the gene density values.

The following expression data and probe set were used for this report:

```
[1] "MEF-GSE3400-MGU74Av2"
```

```
[1] "MG_U74"
```

Density data too sparse for barplot

	coef	se	z	p
A	1.010	0.0694	14.50	7.21e-48
B	0.736	0.0764	9.63	5.77e-22
C	0.521	0.1660	3.14	1.71e-03

D 0.943 0.0483 19.50 1.03e-84
E 1.270 0.1250 10.10 5.32e-24
F 0.803 0.1680 4.79 1.70e-06

Here are the results for expression density. First, we count just genes that are in the upper half.

Density data too sparse for barplot

	coef	se	z	p
A	1.220	0.0856	14.20	1.02e-45
B	0.879	0.0934	9.41	4.79e-21
C	0.730	0.1980	3.68	2.31e-04
D	1.130	0.0593	19.00	2.02e-80
E	1.200	0.1530	7.84	4.36e-15
F	0.706	0.2170	3.25	1.15e-03

Now we count genes in the upper $1/8^{th}$:

Density data too sparse for barplot

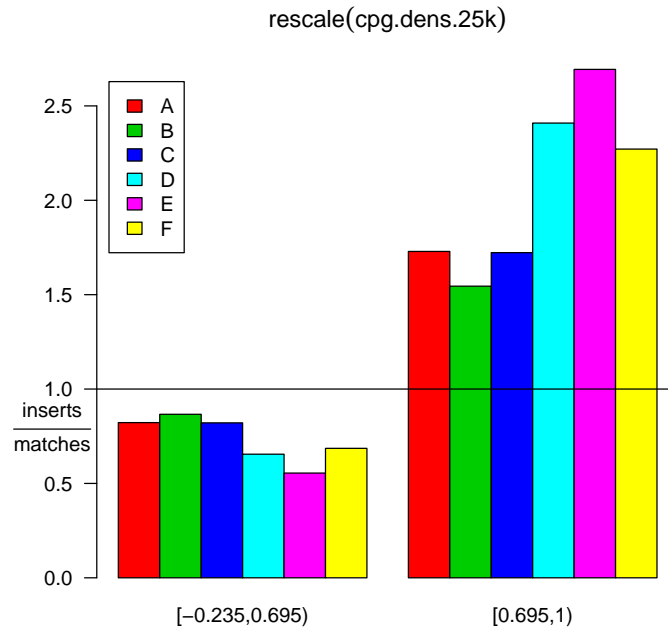
	coef	se	z	p
A	1.270	0.1120	11.40	7.35e-30
B	0.699	0.1290	5.41	6.33e-08
C	0.797	0.2530	3.15	1.65e-03
D	1.350	0.0768	17.60	1.80e-69
E	1.190	0.1980	6.01	1.87e-09
F	0.554	0.3000	1.85	6.42e-02

And here we count genes in the upper $1/16^{th}$:

Density data too sparse for barplot

	coef	se	z	p
A	1.120	0.149	7.52	5.61e-14
B	0.629	0.176	3.57	3.56e-04
C	1.180	0.303	3.88	1.04e-04
D	1.530	0.102	14.90	2.37e-50
E	1.430	0.248	5.77	7.80e-09
F	0.713	0.395	1.81	7.07e-02

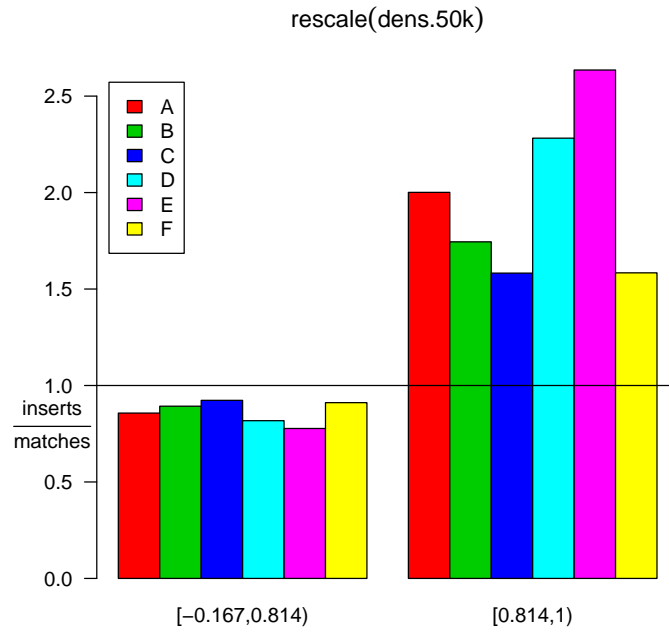
Here the effect of density of CpG islands is studied:



	coef	se	z	p
A	0.740	0.0523	14.10	2.05e-45
B	0.577	0.0588	9.81	1.05e-22
C	0.744	0.1170	6.37	1.88e-10
D	1.310	0.0371	35.30	1.89e-273
E	1.580	0.0999	15.80	1.72e-56
F	1.210	0.1220	9.93	3.21e-23

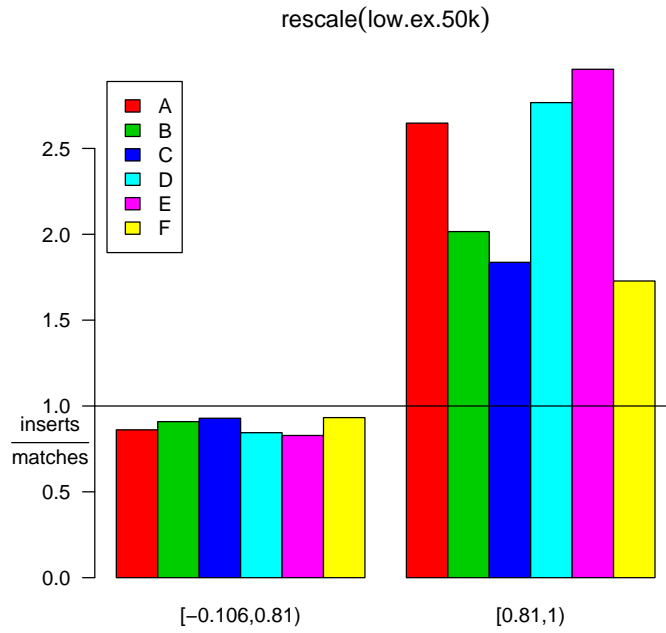
4.2 50 kilobase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 50 kilobase window surrounding each locus. First, we count just the number of genes represented on the chip.



	coef	se	z	p
A	0.924	0.0573	16.10	1.68e-58
B	0.700	0.0628	11.10	8.55e-29
C	0.567	0.1340	4.24	2.27e-05
D	1.040	0.0400	26.00	1.52e-148
E	1.190	0.1060	11.20	4.63e-29
F	0.638	0.1400	4.57	4.97e-06

Here are the results for expression density. First, we count just genes that are in the upper half.



	coef	se	z	p
A	1.130	0.0674	16.80	2.55e-63
B	0.786	0.0733	10.70	8.74e-27
C	0.704	0.1540	4.56	5.00e-06
D	1.170	0.0462	25.30	8.00e-141
E	1.290	0.1190	10.90	1.29e-27
F	0.612	0.1690	3.62	2.93e-04

Now we count genes in the upper $1/8^{th}$:

Density data too sparse for barplot

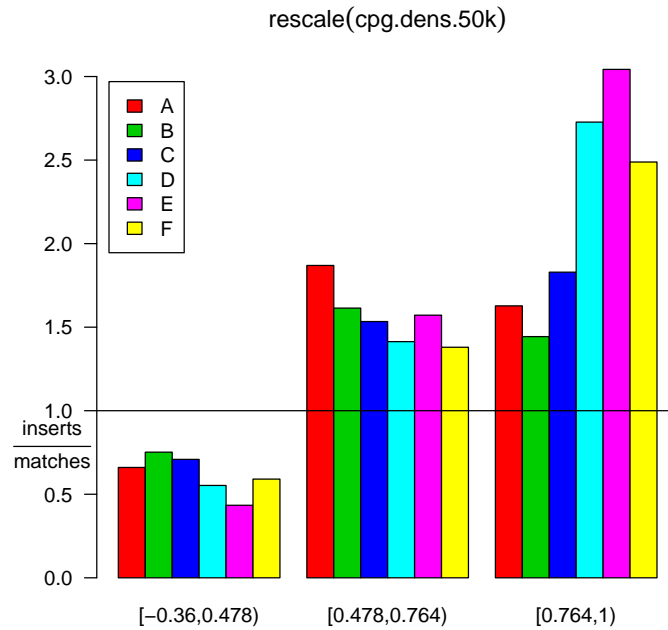
	coef	se	z	p
A	1.130	0.0841	13.50	3.02e-41
B	0.604	0.0986	6.13	8.83e-10
C	0.649	0.1940	3.34	8.39e-04
D	1.270	0.0574	22.10	2.02e-108
E	1.300	0.1440	9.05	1.39e-19
F	0.772	0.2030	3.81	1.40e-04

And here we count genes in the upper $1/16^{th}$:

Density data too sparse for barplot

	coef	se	z	p
A	1.030	0.1120	9.18	4.29e-20
B	0.531	0.1340	3.95	7.72e-05
C	0.944	0.2400	3.94	8.28e-05
D	1.400	0.0748	18.70	3.05e-78
E	1.340	0.1800	7.44	1.01e-13
F	1.110	0.2440	4.56	5.11e-06

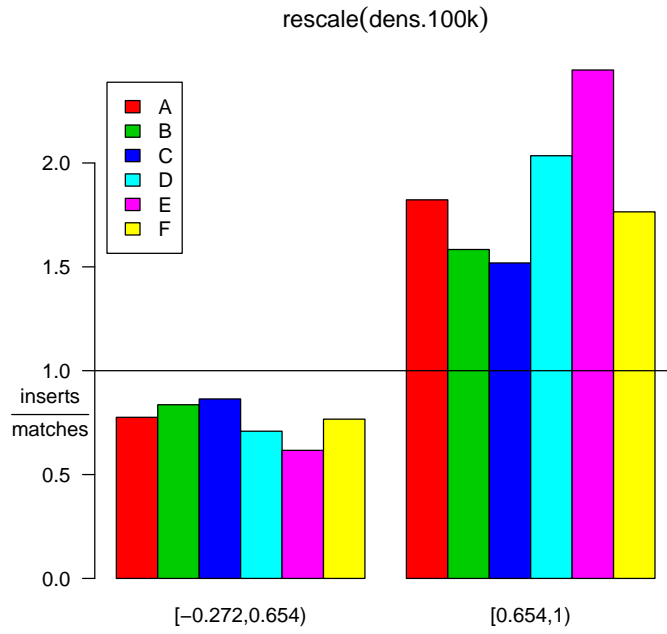
Here the effect of density of CpG islands is studied:



	coef	se	z	p
A	0.981	0.0487	20.10	4.00e-90
B	0.716	0.0538	13.30	2.10e-40
C	0.866	0.1110	7.83	4.72e-15
D	1.270	0.0359	35.40	3.45e-275
E	1.620	0.1060	15.40	2.59e-53
F	1.140	0.1210	9.44	3.67e-21

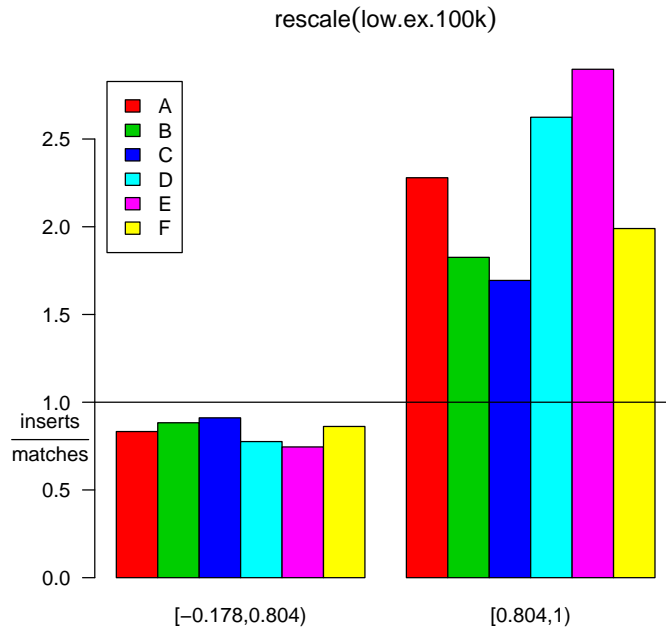
4.3 100 kilobase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 100 kilobase window surrounding each locus. First, we count just the number of genes represented on the chip.



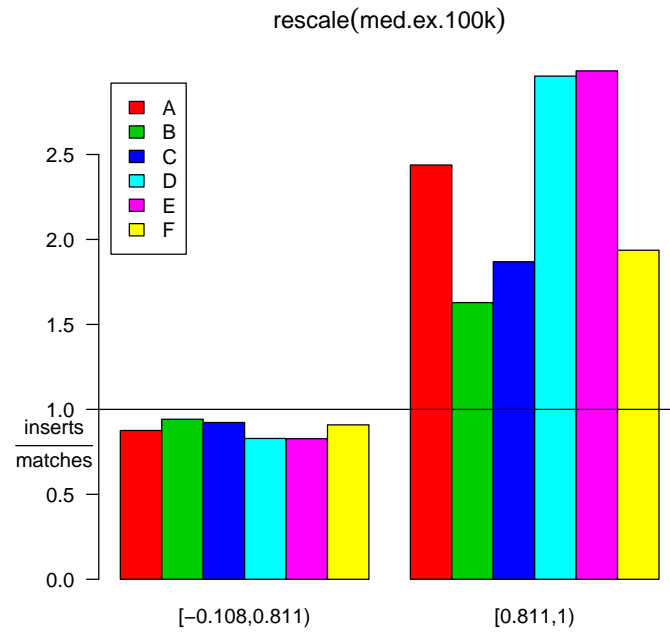
	coef	se	z	p
A	0.928	0.0504	18.40	1.15e-75
B	0.647	0.0559	11.60	5.71e-31
C	0.593	0.1160	5.13	2.91e-07
D	1.070	0.0358	29.90	7.97e-196
E	1.370	0.0989	13.90	6.66e-44
F	1.020	0.1200	8.48	2.25e-17

Here are the results for expression density. First, we count just genes that are in the upper half.



	coef	se	z	p
A	1.080	0.0561	19.30	1.29e-82
B	0.771	0.0615	12.50	3.99e-36
C	0.711	0.1270	5.60	2.14e-08
D	1.210	0.0389	31.20	1.39e-213
E	1.320	0.1040	12.70	4.99e-37
F	1.060	0.1300	8.18	2.88e-16

Now we count genes in the upper $1/8^{th}$:



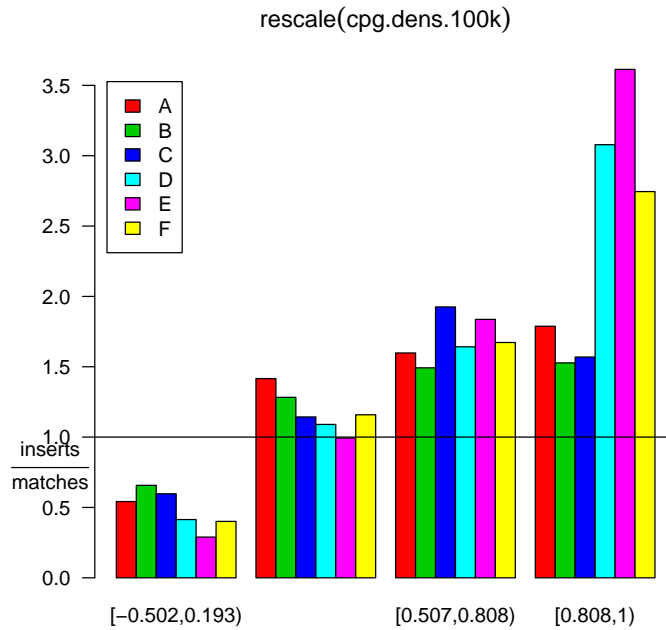
	coef	se	z	p
A	1.040	0.0662	15.80	5.24e-56
B	0.534	0.0782	6.82	8.85e-12
C	0.764	0.1510	5.07	4.08e-07
D	1.250	0.0455	27.50	4.47e-167
E	1.260	0.1200	10.50	1.20e-25
F	0.803	0.1590	5.04	4.63e-07

And here we count genes in the upper $1/16^{th}$:

Density data too sparse for barplot

	coef	se	z	p
A	0.922	0.0874	10.50	5.29e-26
B	0.471	0.1040	4.53	5.85e-06
C	0.787	0.1940	4.06	4.84e-05
D	1.300	0.0569	22.90	2.71e-116
E	1.160	0.1460	7.94	2.03e-15
F	1.030	0.1960	5.25	1.50e-07

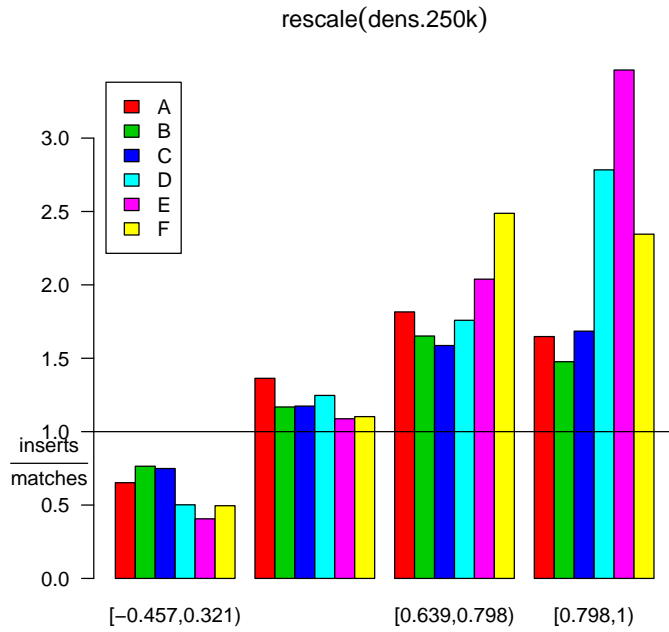
Here the effect of density of CpG islands is studied:



	coef	se	z	p
A	0.779	0.0491	15.90	1.24e-56
B	0.611	0.0549	11.10	9.38e-29
C	0.888	0.1110	7.97	1.61e-15
D	1.290	0.0358	36.00	3.27e-283
E	1.640	0.1020	16.10	2.70e-58
F	1.240	0.1220	10.20	1.40e-24

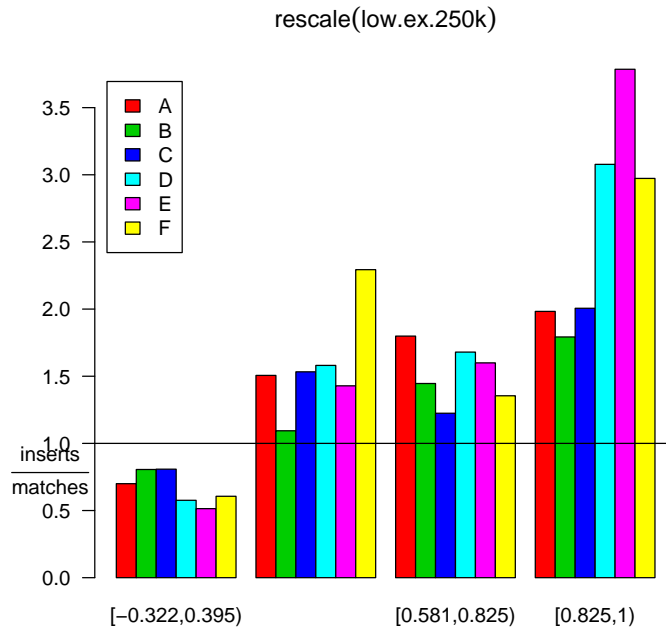
4.4 250 kilobase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 250 kilobase window surrounding each locus. First, we count just the number of genes represented on the chip.



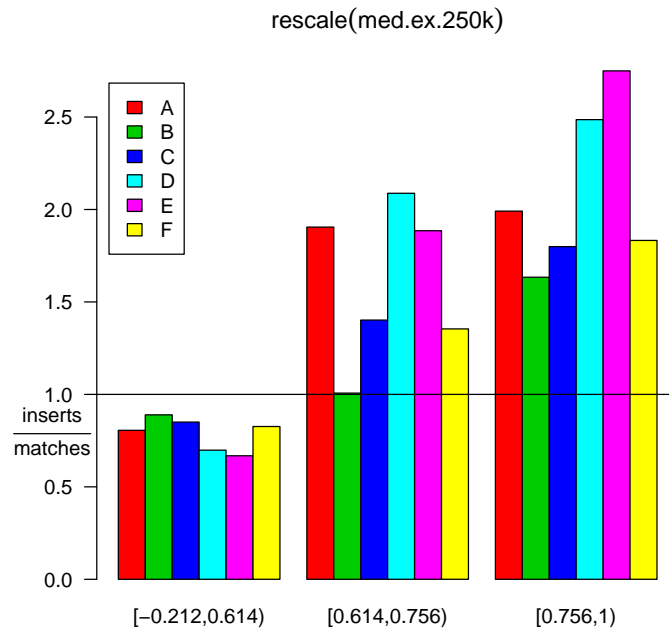
	coef	se	z	p
A	0.865	0.0490	17.70	8.22e-70
B	0.586	0.0538	10.90	1.12e-27
C	0.619	0.1090	5.65	1.60e-08
D	1.280	0.0371	34.40	1.31e-259
E	1.550	0.1130	13.70	6.61e-43
F	1.260	0.1310	9.61	7.13e-22

Here are the results for expression density. First, we count just genes that are in the upper half.



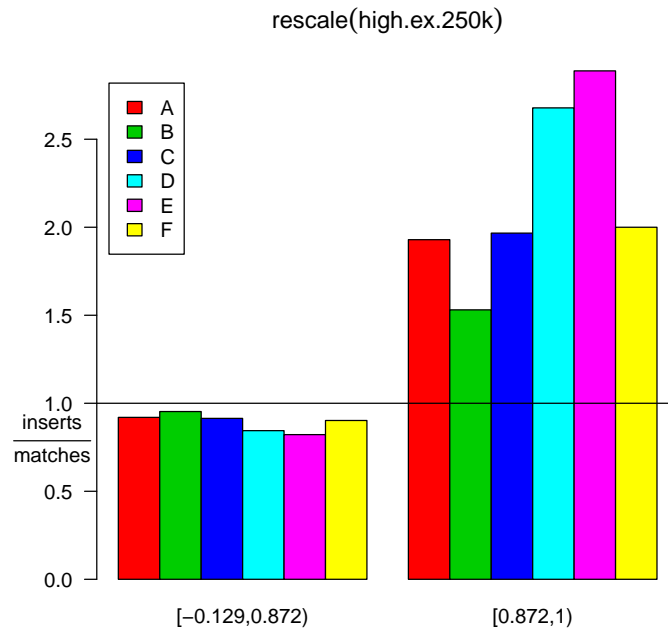
	coef	se	z	p
A	0.978	0.0494	19.80	2.87e-87
B	0.638	0.0546	11.70	1.90e-31
C	0.637	0.1120	5.69	1.24e-08
D	1.300	0.0358	36.40	1.32e-289
E	1.490	0.1010	14.80	1.60e-49
F	1.170	0.1200	9.68	3.50e-22

Now we count genes in the upper $1/8^{th}$:



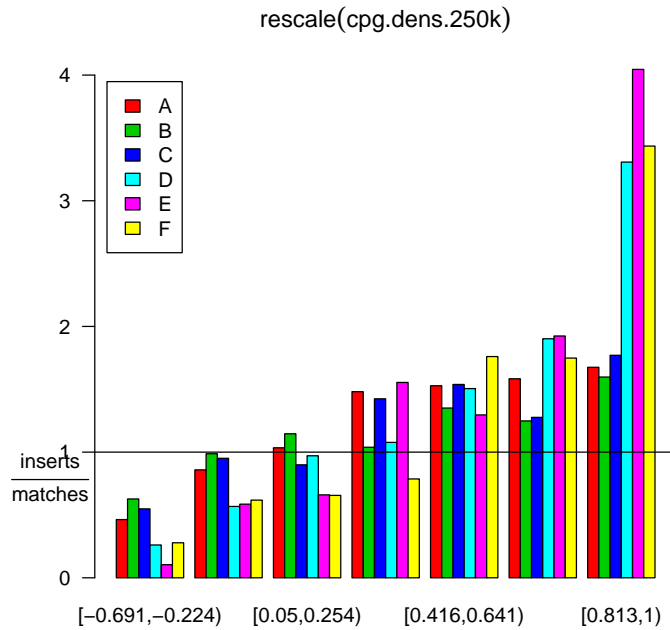
	coef	se	z	p
A	0.907	0.0536	16.90	2.82e-64
B	0.541	0.0611	8.85	8.85e-19
C	0.756	0.1200	6.30	3.06e-10
D	1.250	0.0371	33.70	1.76e-248
E	1.350	0.1010	13.40	9.80e-41
F	0.724	0.1270	5.70	1.19e-08

And here we count genes in the upper 1/16th:



	coef	se	z	p
A	0.817	0.0649	12.60	2.41e-36
B	0.425	0.0756	5.62	1.92e-08
C	0.768	0.1400	5.50	3.80e-08
D	1.220	0.0427	28.50	3.61e-178
E	1.230	0.1130	10.90	1.04e-27
F	0.857	0.1420	6.02	1.70e-09

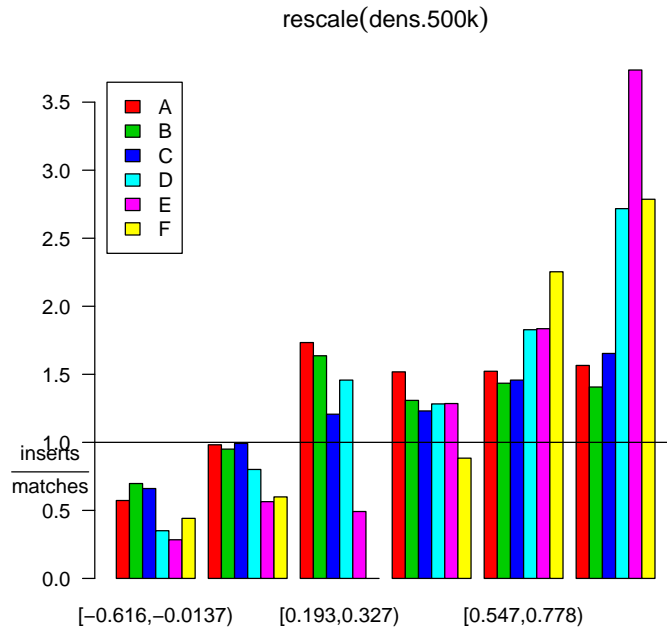
Here the effect of density of CpG islands is studied:



	coef	se	z	p
A	0.840	0.0481	17.50	2.06e-68
B	0.472	0.0536	8.80	1.31e-18
C	0.739	0.1090	6.75	1.44e-11
D	1.380	0.0372	37.10	5.83e-302
E	1.840	0.1170	15.80	3.23e-56
F	1.460	0.1310	11.20	5.41e-29

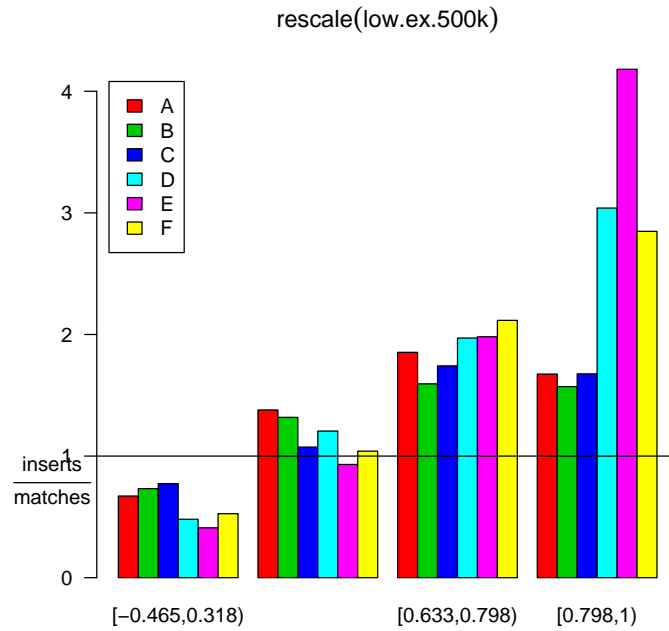
4.5 500 kilobase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 500 kilobase window surrounding each locus. First, we count just the number of genes represented on the chip.



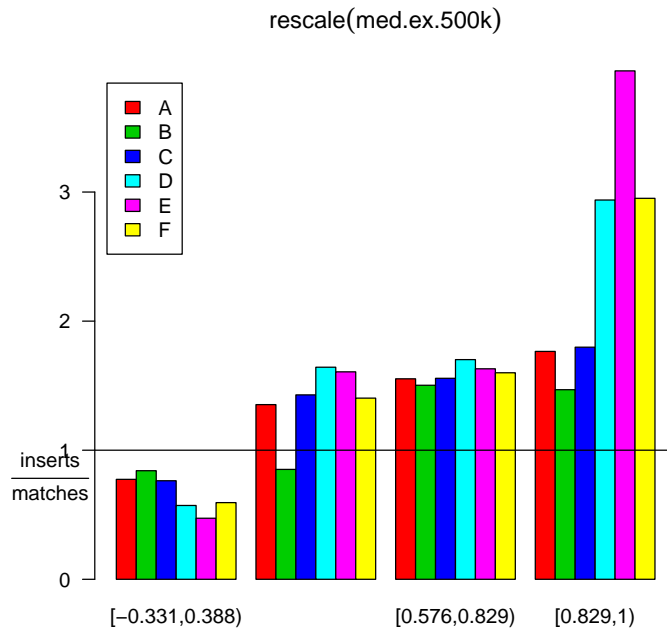
	coef	se	z	p
A	0.801	0.0484	16.60	1.24e-61
B	0.572	0.0535	10.70	1.11e-26
C	0.682	0.1100	6.21	5.38e-10
D	1.350	0.0367	36.70	5.83e-295
E	1.730	0.1120	15.40	1.83e-53
F	1.300	0.1270	10.30	6.87e-25

Here are the results for expression density. First, we count just genes that are in the upper half.



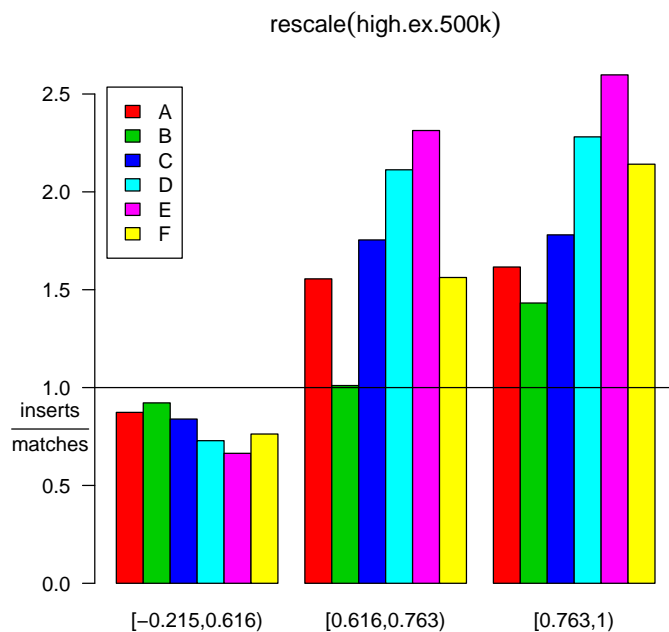
	coef	se	z	p
A	0.904	0.0489	18.50	3.23e-76
B	0.661	0.0542	12.20	3.56e-34
C	0.736	0.1110	6.65	2.90e-11
D	1.450	0.0387	37.30	3.74e-305
E	1.720	0.1190	14.40	4.68e-47
F	1.280	0.1310	9.76	1.61e-22

Now we count genes in the upper $1/8^{th}$:



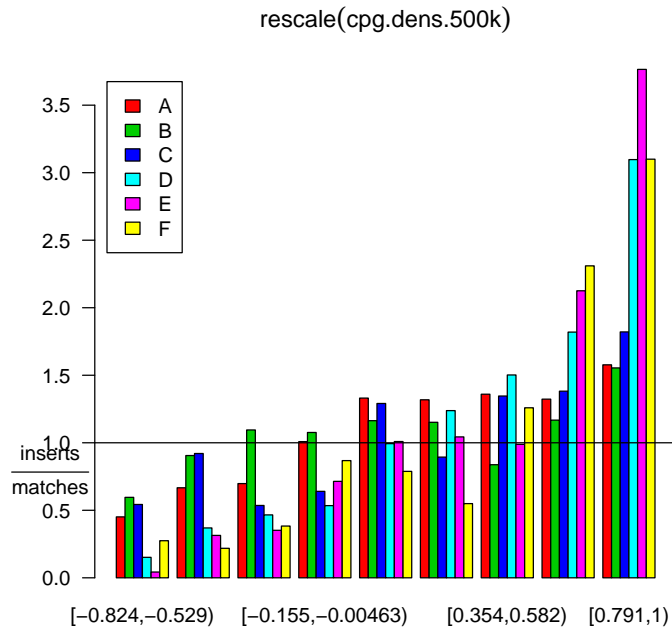
	coef	se	z	p
A	0.738	0.0487	15.20	5.43e-52
B	0.515	0.0549	9.38	6.34e-21
C	0.742	0.1100	6.73	1.65e-11
D	1.300	0.0357	36.30	1.13e-288
E	1.610	0.1030	15.60	1.13e-54
F	1.210	0.1210	10.00	1.46e-23

And here we count genes in the upper $1/16^{th}$:



	coef	se	z	p
A	0.623	0.0542	11.50	1.30e-30
B	0.407	0.0622	6.55	5.61e-11
C	0.824	0.1180	6.98	2.89e-12
D	1.150	0.0371	31.00	2.66e-211
E	1.310	0.0990	13.20	6.16e-40
F	1.010	0.1230	8.25	1.56e-16

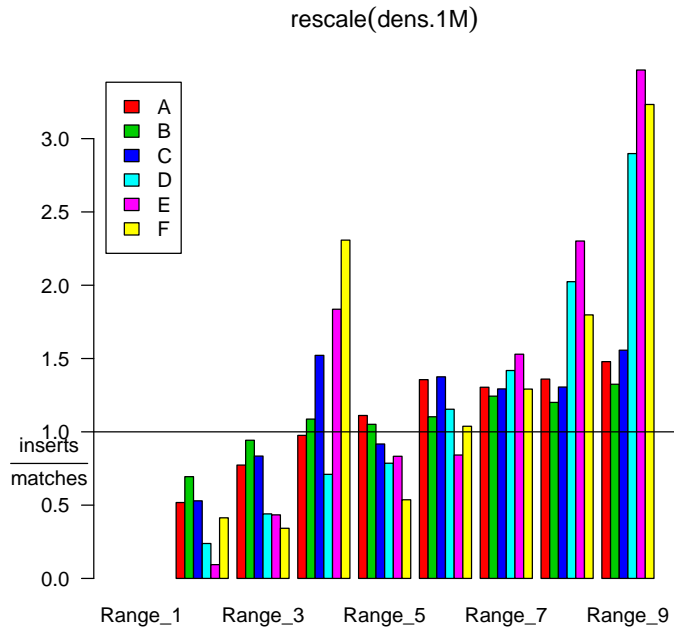
Here the effect of density of CpG islands is studied:



	coef	se	z	p
A	0.706	0.0479	14.70	3.49e-49
B	0.287	0.0535	5.36	8.17e-08
C	0.736	0.1100	6.67	2.57e-11
D	1.520	0.0393	38.80	0.00e+00
E	1.570	0.1160	13.50	1.04e-41
F	1.370	0.1350	10.10	3.87e-24

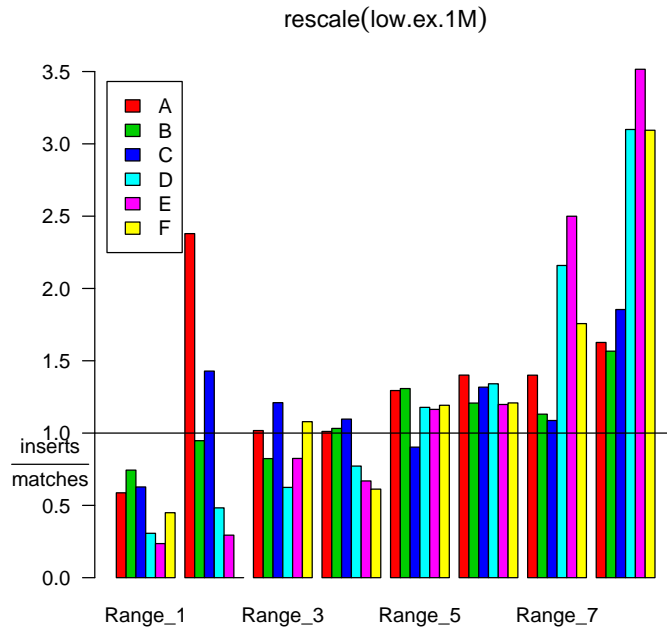
4.6 1 megabase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 1 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.



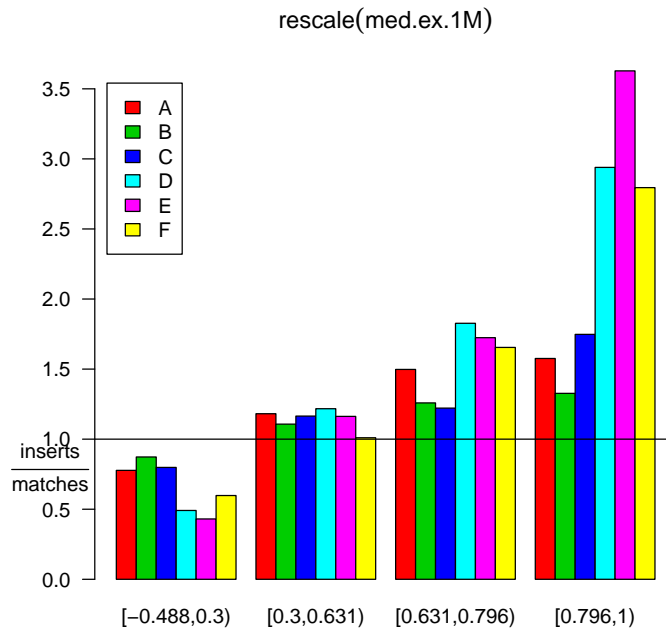
	coef	se	z	p
A	0.671	0.0481	13.90	3.86e-44
B	0.342	0.0535	6.39	1.70e-10
C	0.626	0.1110	5.64	1.66e-08
D	1.450	0.0396	36.60	3.76e-293
E	1.550	0.1170	13.20	4.81e-40
F	1.260	0.1360	9.26	2.07e-20

Here are the results for expression density. First, we count just genes that are in the upper half.



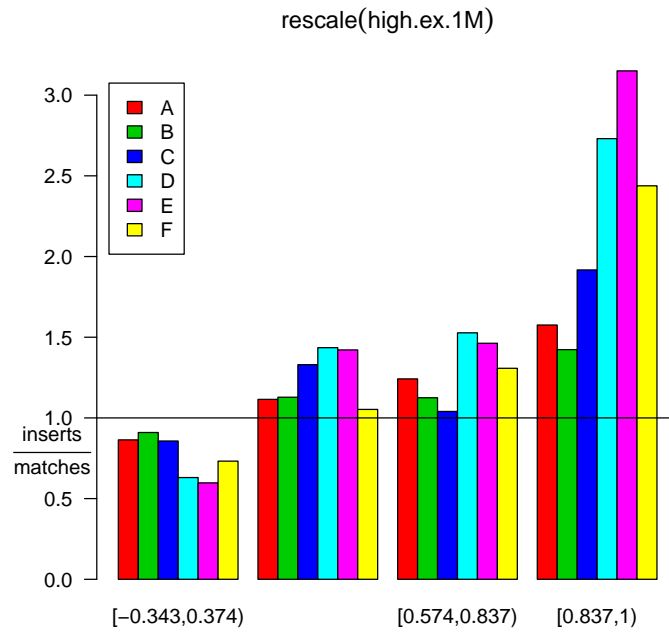
	coef	se	z	p
A	0.667	0.0476	14.00	1.61e-44
B	0.434	0.0535	8.10	5.37e-16
C	0.507	0.1090	4.66	3.09e-06
D	1.410	0.0374	37.70	0.00e+00
E	1.580	0.1100	14.30	3.67e-46
F	1.190	0.1270	9.37	7.20e-21

Now we count genes in the upper $1/8^{th}$:



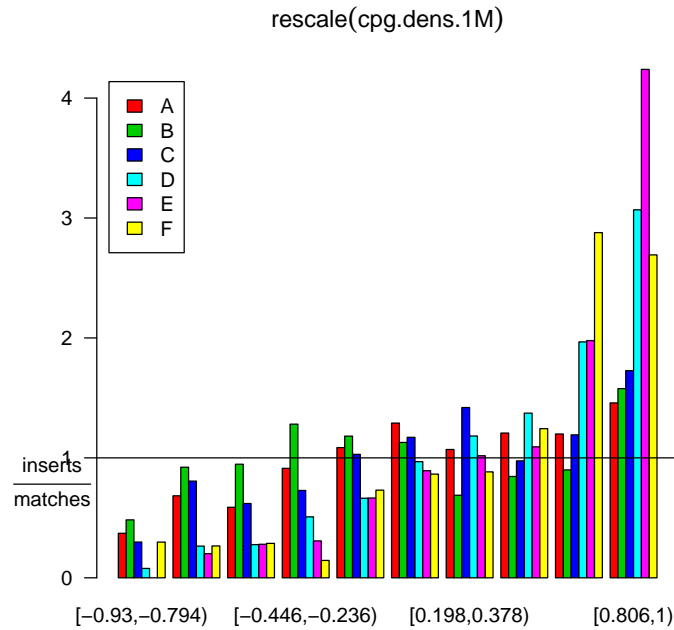
	coef	se	z	p
A	0.592	0.0476	12.40	1.52e-35
B	0.308	0.0535	5.75	8.70e-09
C	0.533	0.1100	4.83	1.33e-06
D	1.350	0.0386	35.00	4.97e-269
E	1.600	0.1190	13.50	1.91e-41
F	1.140	0.1310	8.67	4.49e-18

And here we count genes in the upper $1/16^{th}$:



	coef	se	z	p
A	0.454	0.0486	9.34	9.85e-21
B	0.296	0.0553	5.35	8.64e-08
C	0.427	0.1120	3.82	1.33e-04
D	1.120	0.0353	31.80	1.35e-221
E	1.150	0.0977	11.80	5.18e-32
F	0.802	0.1180	6.79	1.10e-11

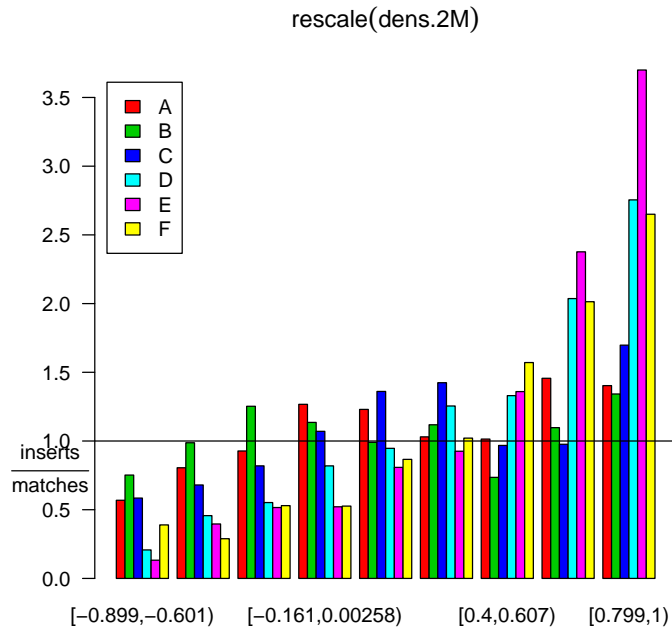
Here the effect of density of CpG islands is studied:



	coef	se	z	p
A	0.4490	0.0473	9.49	2.42e-21
B	-0.0171	0.0534	-0.32	7.49e-01
C	0.5260	0.1100	4.79	1.64e-06
D	1.4300	0.0395	36.20	2.96e-287
E	1.7400	0.1260	13.80	3.00e-43
F	1.4200	0.1420	9.97	1.99e-23

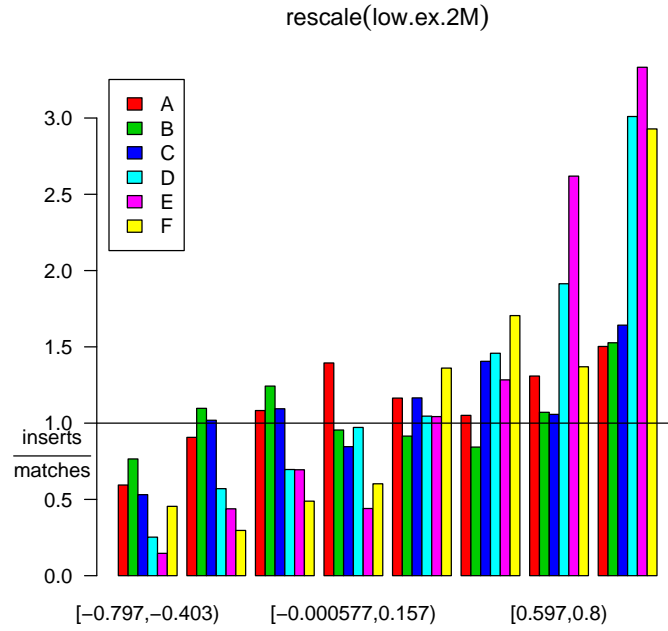
4.7 2 megabase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 2 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.



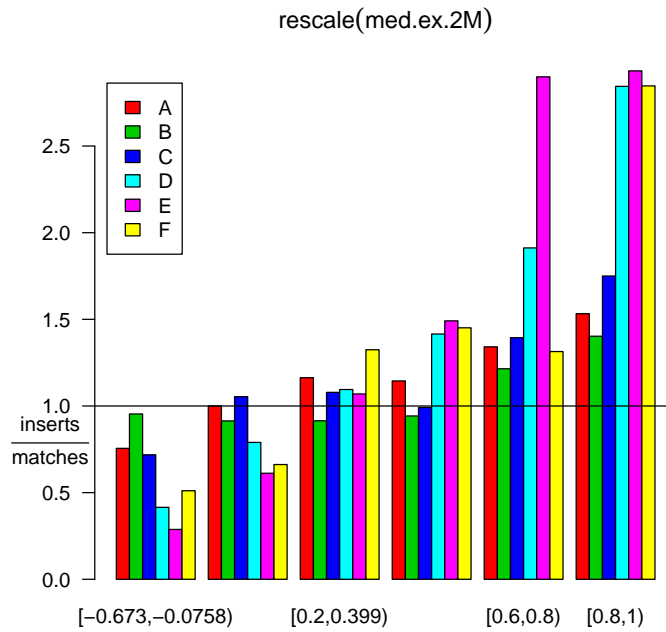
	coef	se	z	p
A	0.341	0.0468	7.28	3.45e-13
B	0.114	0.0535	2.13	3.33e-02
C	0.522	0.1090	4.78	1.80e-06
D	1.260	0.0375	33.70	7.30e-249
E	1.610	0.1180	13.60	2.14e-42
F	1.350	0.1370	9.79	1.26e-22

Here are the results for expression density. First, we count just genes that are in the upper half.



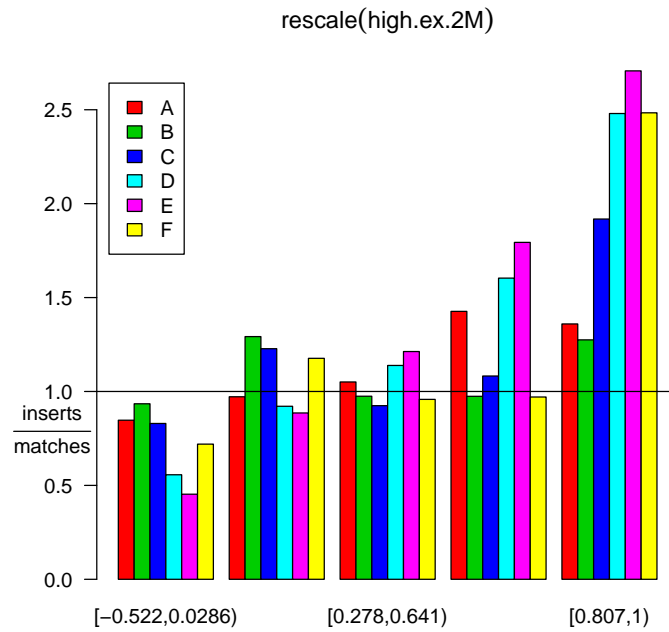
	coef	se	z	p
A	0.473	0.0474	9.97	2.03e-23
B	0.083	0.0534	1.56	1.20e-01
C	0.454	0.1090	4.15	3.31e-05
D	1.310	0.0385	34.10	3.66e-255
E	1.620	0.1220	13.30	3.33e-40
F	1.410	0.1430	9.88	5.26e-23

Now we count genes in the upper $1/8^{th}$:



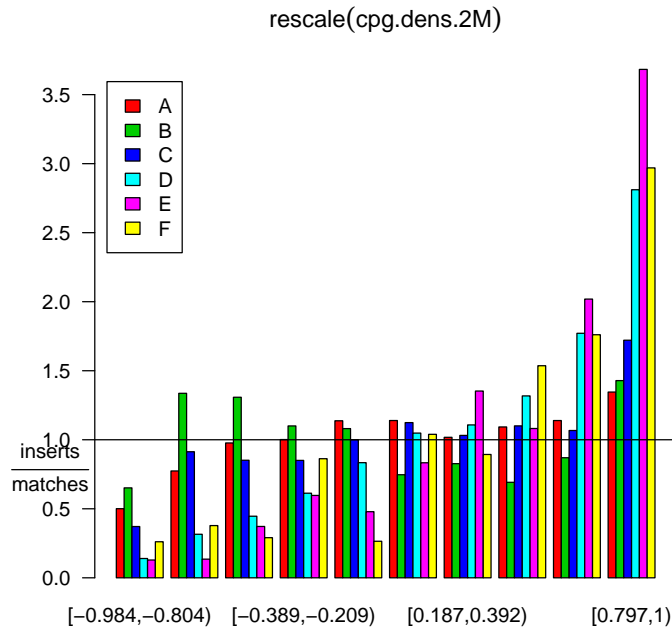
	coef	se	z	p
A	0.466	0.0471	9.89	4.71e-23
B	0.115	0.0535	2.15	3.19e-02
C	0.482	0.1090	4.43	9.45e-06
D	1.210	0.0366	33.20	4.30e-241
E	1.620	0.1170	13.90	1.12e-43
F	1.060	0.1270	8.36	6.43e-17

And here we count genes in the upper $1/16^{th}$:



	coef	se	z	p
A	0.3270	0.0470	6.96	3.52e-12
B	0.0838	0.0535	1.57	1.17e-01
C	0.2940	0.1090	2.71	6.79e-03
D	0.9780	0.0356	27.40	1.41e-165
E	1.2300	0.1070	11.50	1.16e-30
F	0.4990	0.1190	4.17	3.00e-05

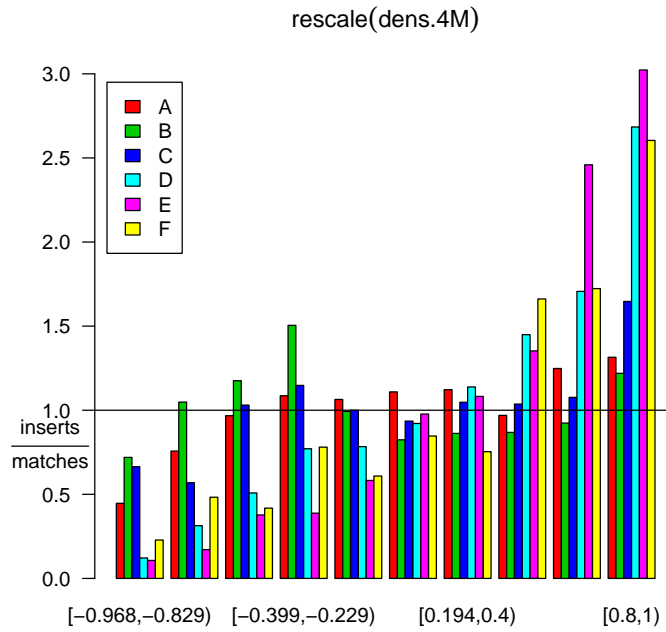
Here the effect of density of CpG islands is studied:



	coef	se	z	p
A	0.250	0.0469	5.33	9.63e-08
B	-0.146	0.0539	-2.70	6.84e-03
C	0.415	0.1090	3.80	1.46e-04
D	1.140	0.0370	30.90	1.71e-209
E	1.630	0.1210	13.50	1.80e-41
F	1.380	0.1390	9.95	2.48e-23

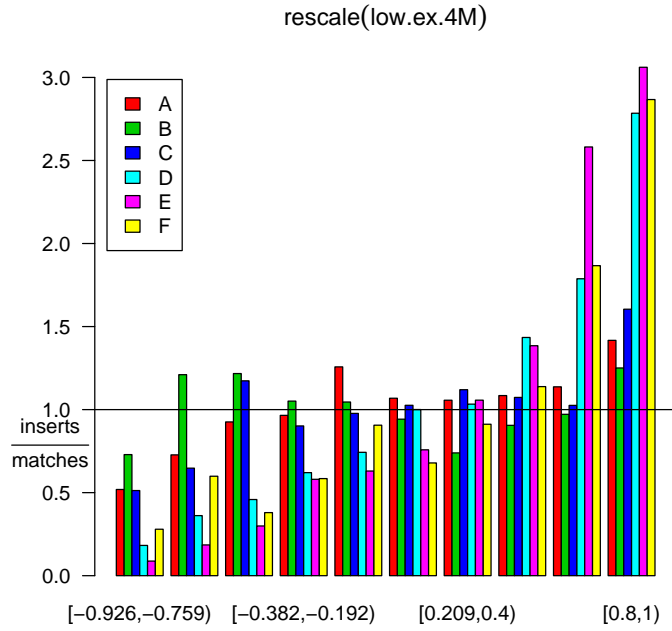
4.8 4 megabase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 4 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.



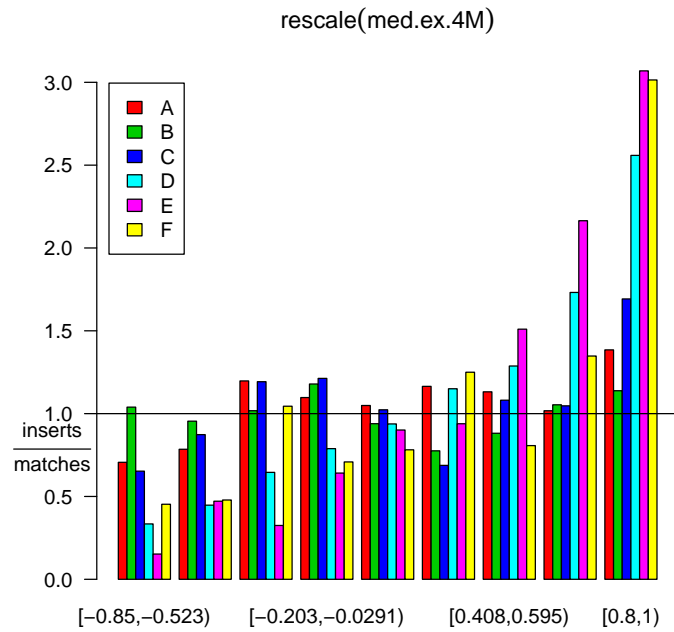
	coef	se	z	p
A	0.280	0.0473	5.93	3.01e-09
B	-0.095	0.0536	-1.77	7.63e-02
C	0.236	0.1090	2.18	2.96e-02
D	1.130	0.0373	30.20	8.06e-201
E	1.560	0.1200	13.00	8.45e-39
F	1.200	0.1360	8.85	8.69e-19

Here are the results for expression density. First, we count just genes that are in the upper half.



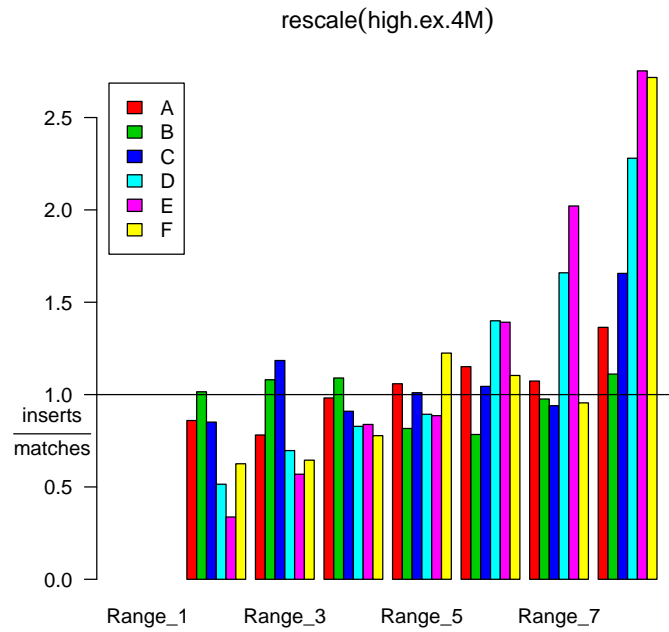
	coef	se	z	p
A	0.274	0.0471	5.83	5.67e-09
B	-0.102	0.0538	-1.90	5.77e-02
C	0.290	0.1080	2.67	7.52e-03
D	1.140	0.0370	30.80	1.36e-208
E	1.550	0.1180	13.10	3.23e-39
F	0.985	0.1280	7.67	1.71e-14

Now we count genes in the upper $1/8^{th}$:



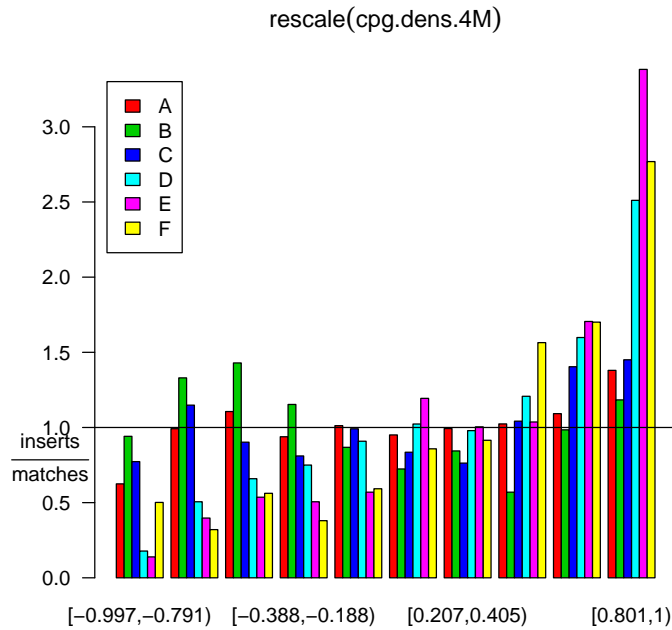
	coef	se	z	p
A	0.276	0.0471	5.87	4.35e-09
B	-0.108	0.0537	-2.02	4.38e-02
C	0.235	0.1080	2.17	3.03e-02
D	1.070	0.0369	29.00	8.82e-185
E	1.510	0.1170	12.90	4.93e-38
F	0.945	0.1270	7.43	1.11e-13

And here we count genes in the upper $1/16^{th}$:



	coef	se	z	p
A	0.2680	0.0471	5.68	1.33e-08
B	-0.0868	0.0537	-1.62	1.06e-01
C	0.2330	0.1080	2.15	3.14e-02
D	0.8910	0.0357	25.00	1.75e-137
E	1.1900	0.1090	11.00	6.55e-28
F	0.6880	0.1220	5.63	1.77e-08

Here the effect of density of CpG islands is studied:

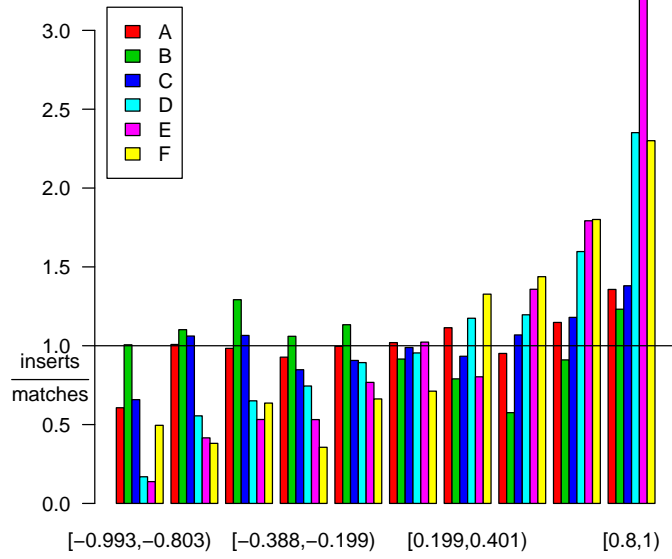


	coef	se	z	p
A	0.140	0.0470	2.98	2.93e-03
B	-0.289	0.0545	-5.30	1.15e-07
C	0.186	0.1080	1.72	8.51e-02
D	0.863	0.0355	24.30	8.83e-131
E	1.280	0.1110	11.50	1.56e-30
F	1.160	0.1330	8.73	2.48e-18

4.9 8 megabase Window

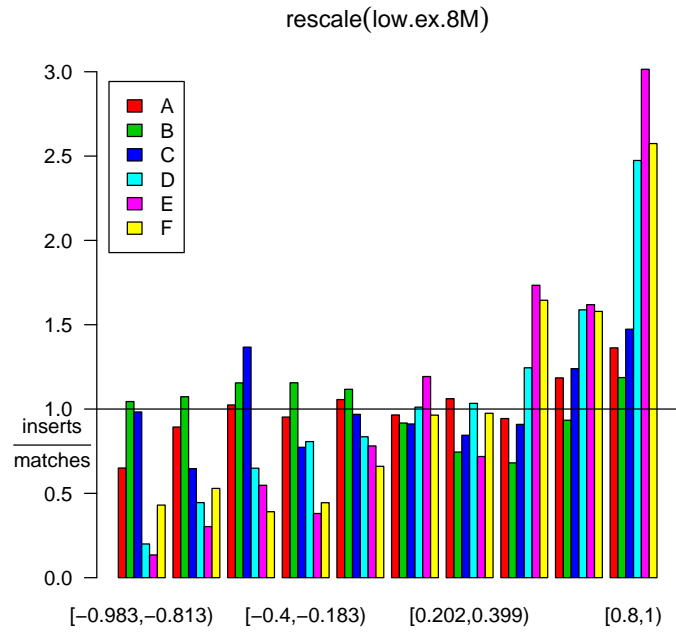
In the barplot that follows we examine the association of insertion sites with expression density in a 8 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.

rescale(dens.8M)



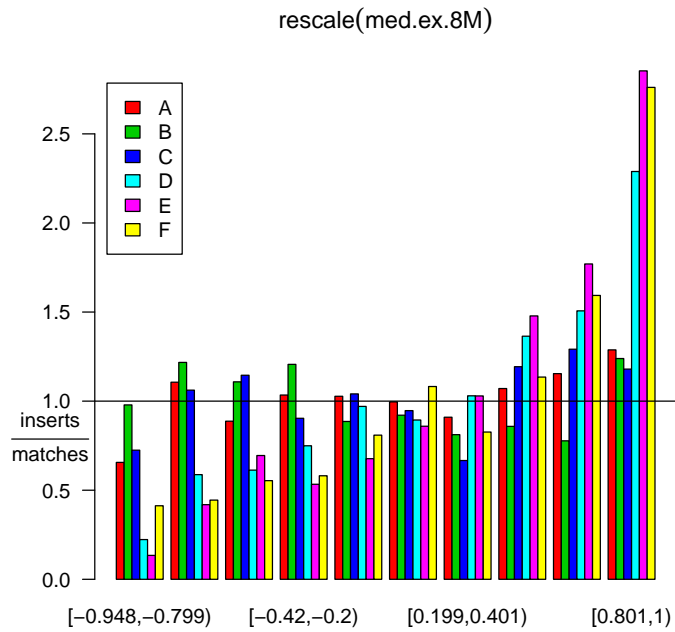
	coef	se	z	p
A	0.206	0.0469	4.39	1.14e-05
B	-0.243	0.0541	-4.50	6.78e-06
C	0.201	0.1090	1.85	6.45e-02
D	0.865	0.0357	24.20	1.04e-129
E	1.160	0.1090	10.60	2.80e-26
F	1.110	0.1330	8.35	6.91e-17

Here are the results for expression density. First, we count just genes that are in the upper half.



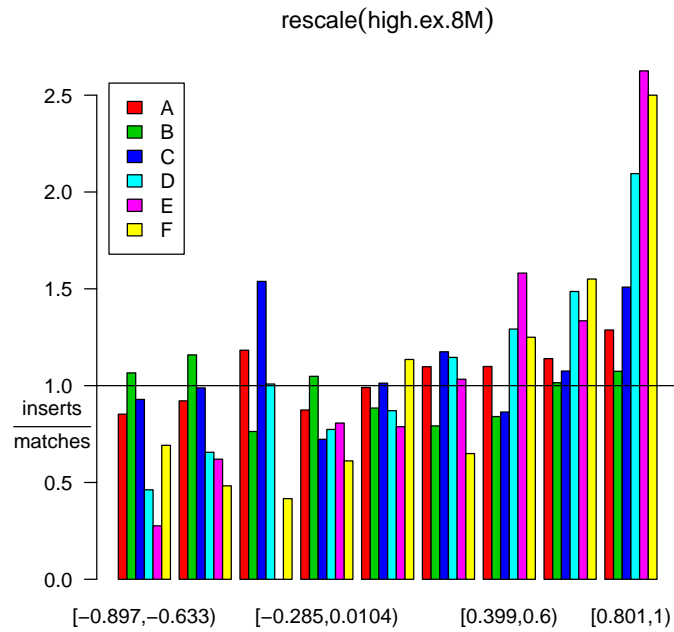
	coef	se	z	p
A	0.165	0.0470	3.51	4.52e-04
B	-0.221	0.0541	-4.07	4.64e-05
C	0.126	0.1090	1.16	2.47e-01
D	0.883	0.0356	24.80	1.47e-135
E	1.240	0.1110	11.20	4.31e-29
F	1.150	0.1330	8.67	4.46e-18

Now we count genes in the upper $1/8^{th}$:



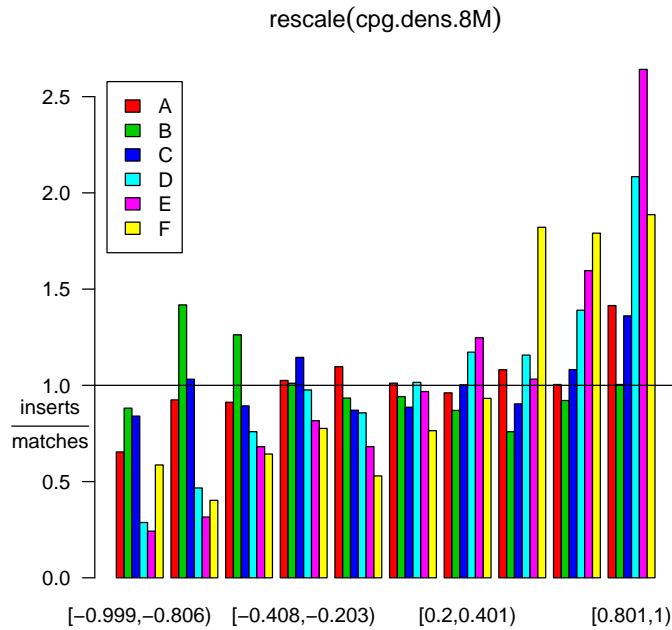
	coef	se	z	p
A	0.1420	0.0469	3.020	2.53e-03
B	-0.1660	0.0538	-3.080	2.05e-03
C	0.0811	0.1080	0.749	4.54e-01
D	0.8000	0.0354	22.600	7.90e-113
E	1.1500	0.1090	10.600	4.04e-26
F	0.9450	0.1290	7.320	2.53e-13

And here we count genes in the upper 1/16th:



	coef	se	z	p
A	0.221	0.0471	4.69	2.79e-06
B	-0.136	0.0542	-2.51	1.22e-02
C	0.244	0.1080	2.25	2.46e-02
D	0.742	0.0348	21.30	9.04e-101
E	0.930	0.1020	9.15	5.63e-20
F	0.677	0.1210	5.60	2.09e-08

Here the effect of density of CpG islands is studied:

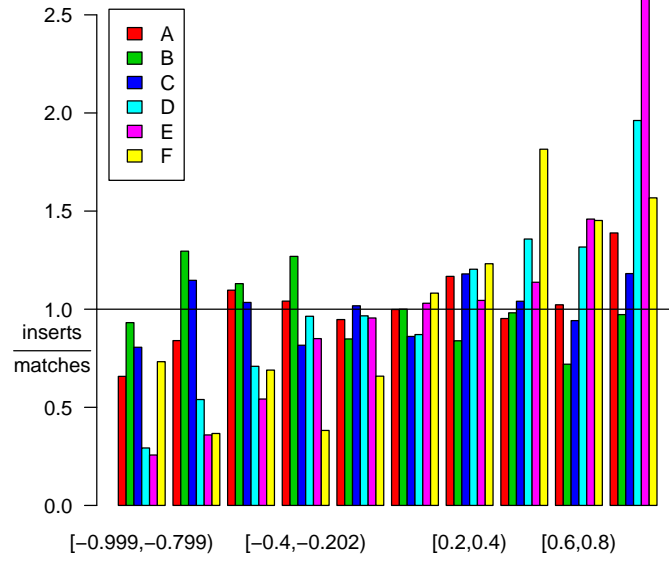


	coef	se	z	p
A	0.162	0.0470	3.440	5.83e-04
B	-0.223	0.0540	-4.120	3.71e-05
C	0.106	0.1080	0.984	3.25e-01
D	0.701	0.0346	20.200	5.20e-91
E	0.940	0.1040	9.020	1.94e-19
F	0.899	0.1270	7.090	1.39e-12

4.10 16 megabase Window

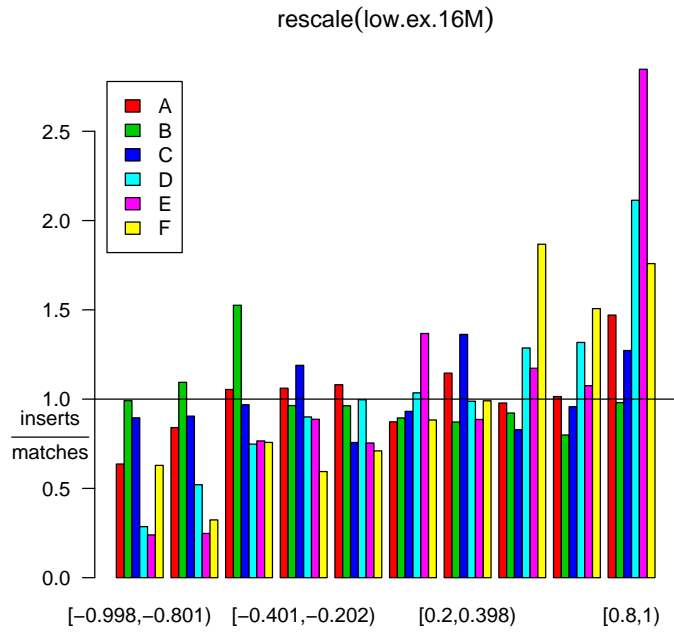
In the barplot that follows we examine the association of insertion sites with expression density in a 16 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.

rescale(dens.16M)



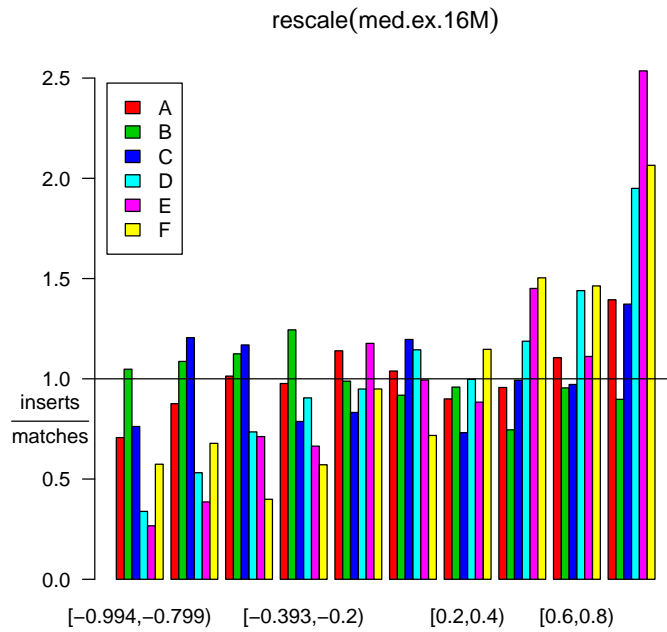
	coef	se	z	p
A	0.1880	0.0469	4.020	5.92e-05
B	-0.1940	0.0539	-3.590	3.27e-04
C	0.0832	0.1080	0.769	4.42e-01
D	0.6560	0.0345	19.000	2.21e-80
E	0.8950	0.1040	8.600	8.01e-18
F	0.9170	0.1290	7.100	1.28e-12

Here are the results for expression density. First, we count just genes that are in the upper half.



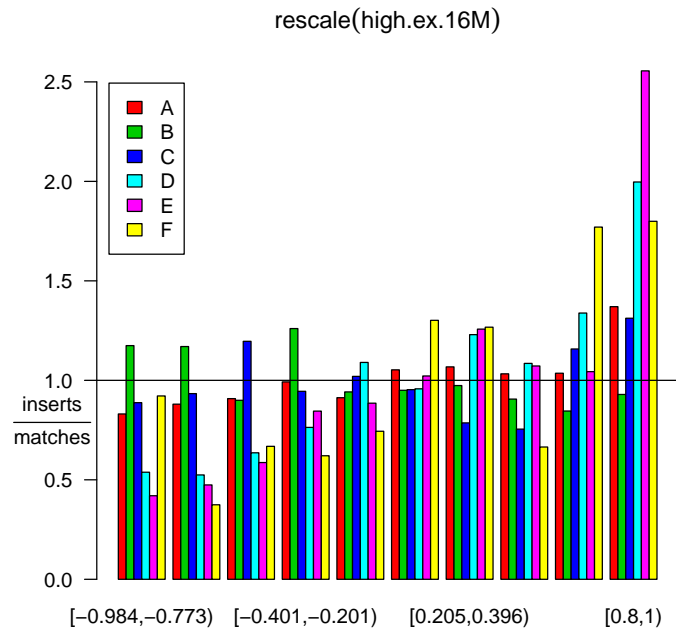
	coef	se	z	p
A	0.166	0.0470	3.54	4.07e-04
B	-0.215	0.0539	-3.99	6.51e-05
C	0.116	0.1080	1.07	2.83e-01
D	0.677	0.0347	19.50	7.07e-85
E	0.900	0.1040	8.68	4.05e-18
F	0.840	0.1270	6.63	3.41e-11

Now we count genes in the upper $1/8^{th}$:



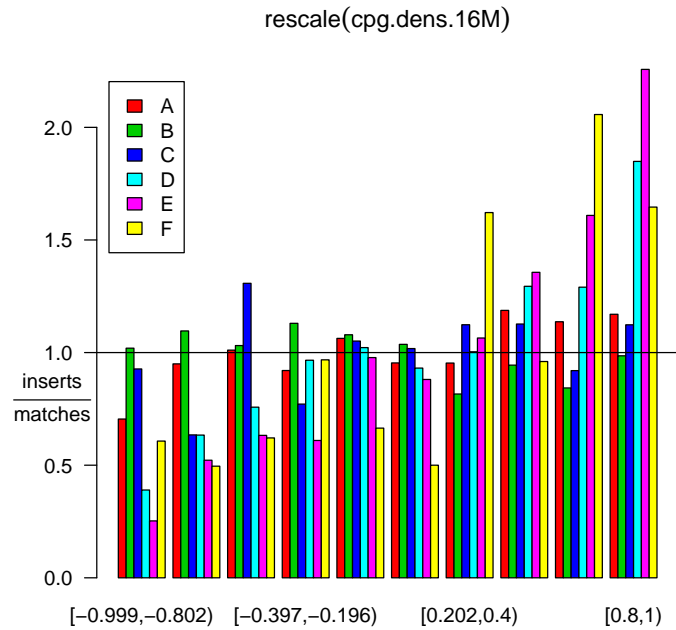
	coef	se	z	p
A	0.1320	0.0470	2.810	5.01e-03
B	-0.2050	0.0539	-3.800	1.42e-04
C	0.0777	0.1080	0.719	4.72e-01
D	0.6600	0.0347	19.000	8.39e-81
E	0.7730	0.1010	7.610	2.65e-14
F	0.7120	0.1240	5.730	9.96e-09

And here we count genes in the upper 1/16th:



	coef	se	z	p
A	0.2010	0.0469	4.280	1.85e-05
B	-0.1540	0.0538	-2.860	4.24e-03
C	-0.0222	0.1080	-0.206	8.37e-01
D	0.6140	0.0346	17.700	2.36e-70
E	0.7520	0.1020	7.390	1.42e-13
F	0.6820	0.1230	5.520	3.30e-08

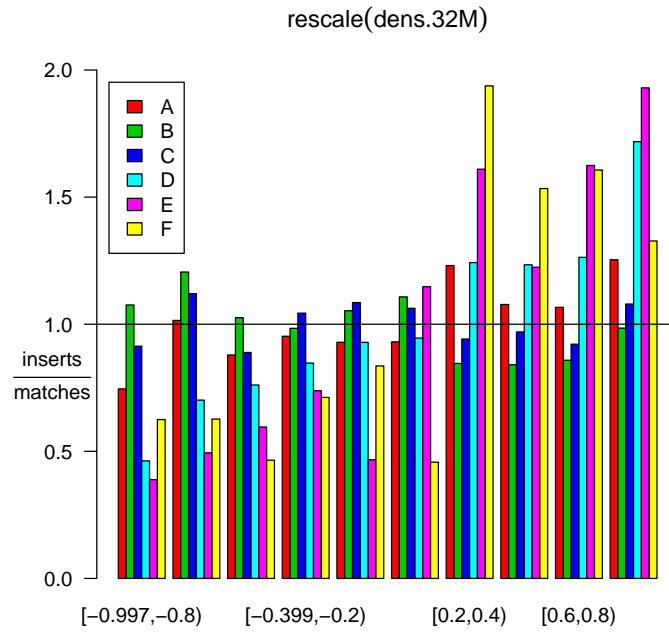
Here the effect of density of CpG islands is studied:



	coef	se	z	p
A	0.154	0.0469	3.28	1.05e-03
B	-0.132	0.0537	-2.45	1.42e-02
C	0.135	0.1080	1.25	2.12e-01
D	0.518	0.0340	15.20	2.52e-52
E	0.880	0.1040	8.48	2.30e-17
F	0.697	0.1230	5.65	1.63e-08

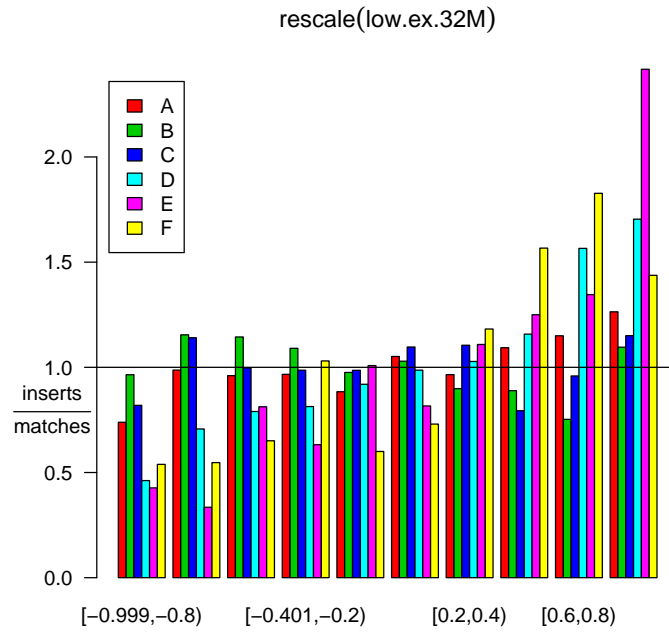
4.11 32 megabase Window

In the barplot that follows we examine the association of insertion sites with expression density in a 32 megabase window surrounding each locus. First, we count just the number of genes represented on the chip.



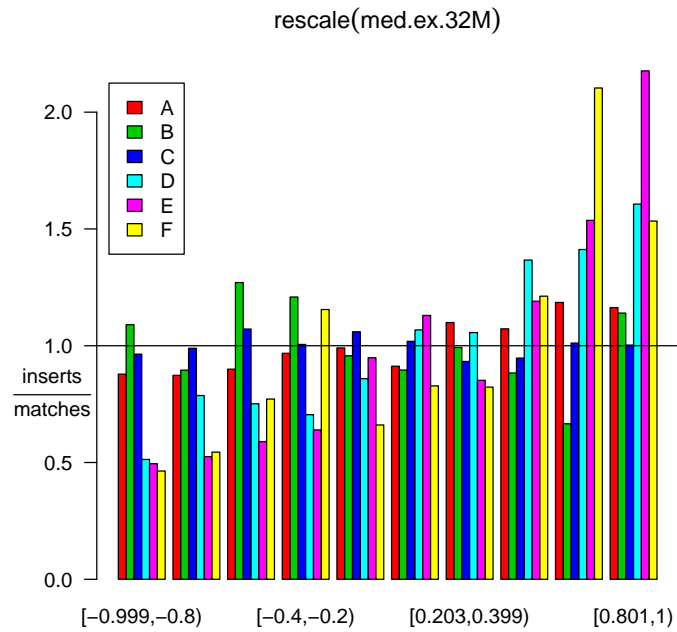
	coef	se	z	p
A	0.20400	0.0467	4.3800	1.21e-05
B	-0.13800	0.0537	-2.5800	1.00e-02
C	-0.00746	0.1080	-0.0689	9.45e-01
D	0.54700	0.0341	16.1000	5.69e-58
E	1.03000	0.1070	9.6900	3.25e-22
F	0.73200	0.1240	5.9100	3.37e-09

Here are the results for expression density. First, we count just genes that are in the upper half.



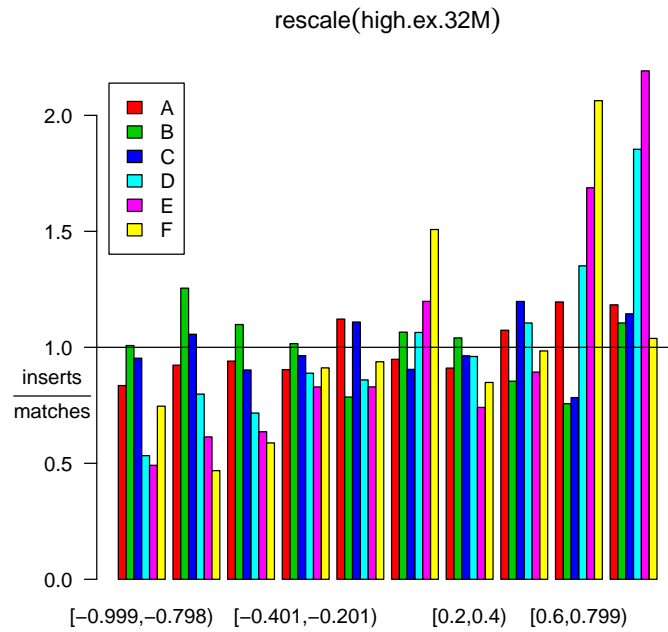
	coef	se	z	p
A	0.1980	0.0470	4.210	2.52e-05
B	-0.1330	0.0537	-2.470	1.34e-02
C	0.0363	0.1080	0.335	7.38e-01
D	0.5530	0.0342	16.200	6.40e-59
E	0.7600	0.1010	7.510	6.14e-14
F	0.6840	0.1220	5.590	2.31e-08

Now we count genes in the upper $1/8^{th}$:



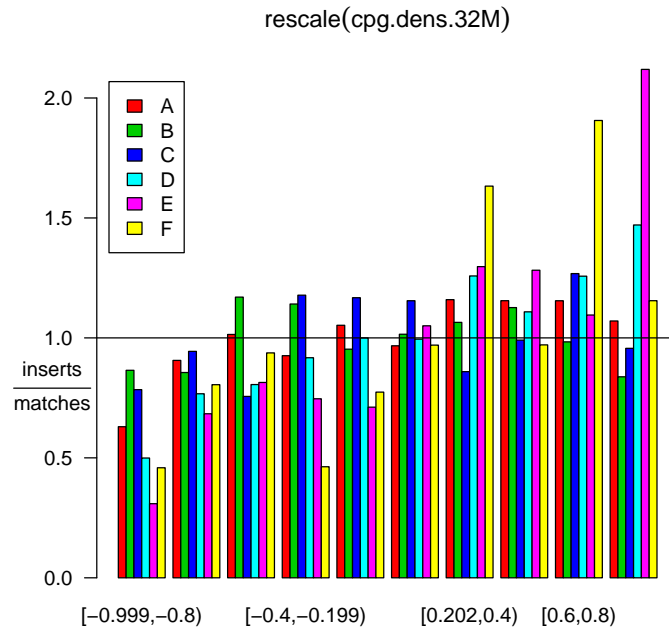
	coef	se	z	p
A	0.1620	0.0468	3.470	5.27e-04
B	-0.1690	0.0537	-3.140	1.66e-03
C	-0.0353	0.1090	-0.325	7.45e-01
D	0.5850	0.0342	17.100	2.11e-65
E	0.7610	0.1020	7.490	6.90e-14
F	0.5840	0.1220	4.800	1.62e-06

And here we count genes in the upper 1/16th:



	coef	se	z	p
A	0.1110	0.0468	2.370	1.77e-02
B	-0.0698	0.0537	-1.300	1.94e-01
C	-0.0139	0.1080	-0.128	8.98e-01
D	0.5030	0.0341	14.800	2.44e-49
E	0.6660	0.0999	6.670	2.58e-11
F	0.5520	0.1210	4.580	4.74e-06

Here the effect of density of CpG islands is studied:

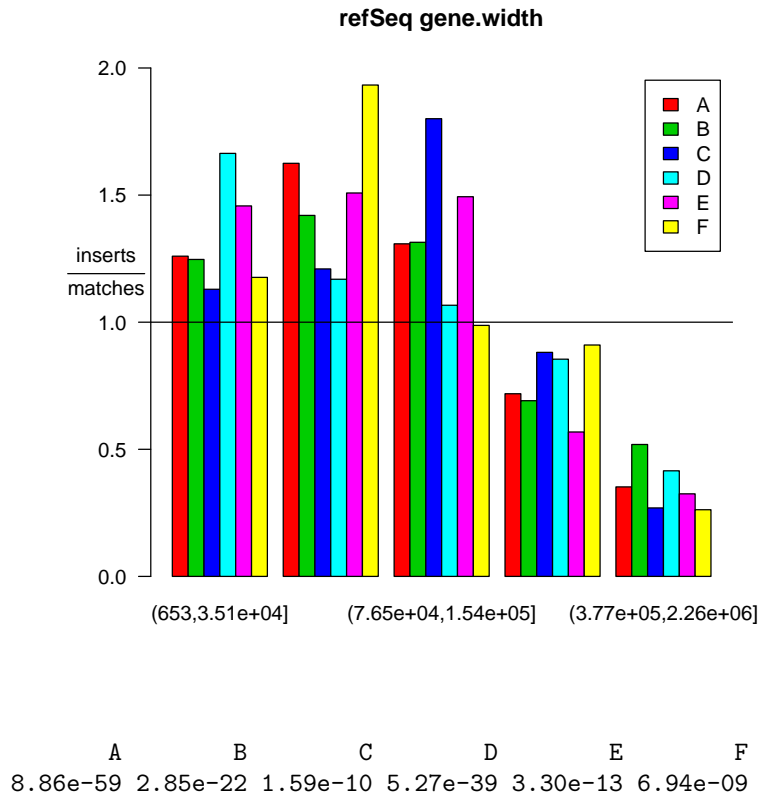


	coef	se	z	p
A	0.2030	0.0473	4.290	1.77e-05
B	0.0119	0.0534	0.223	8.23e-01
C	0.0748	0.1090	0.686	4.93e-01
D	0.4210	0.0338	12.500	1.40e-35
E	0.7510	0.1020	7.380	1.56e-13
F	0.6510	0.1220	5.350	8.79e-08

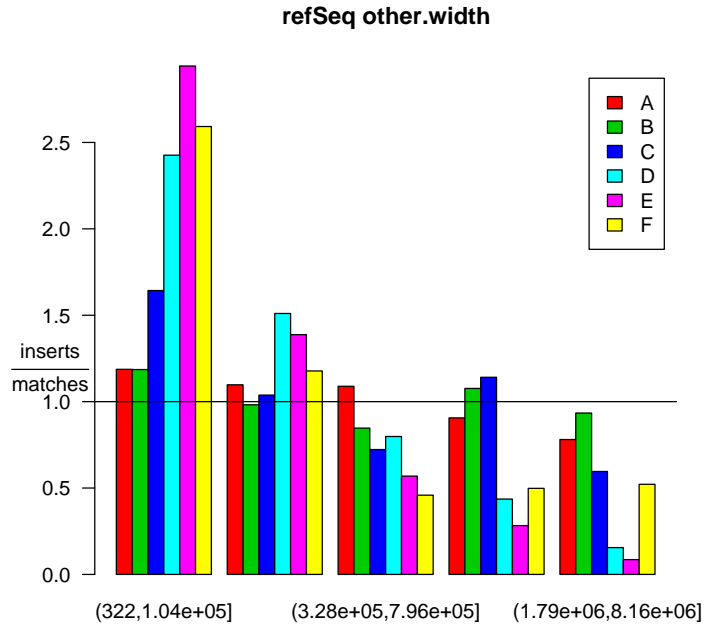
5 Juxtaposition with Gene Start and End Positions

5.1 Refseq Annotations

In this section we study the effect of juxtaposition in terms of gene start and end positions. The first barplot shows the effect of gene width for those insertions that are located within an RefSeq gene. The table following the barplot shows the p-values for a test of the hypothesis that the proportions in each of the categories that define the bars are equal in the insertions and their matches. This p-value is obtained from the $5 \times 2 \times k$ table of counts defined by gene width category, insertion/match status, and stratum (consisting of an insertion and its matched sites) using a likelihood ratio test for the hypothesis of no association between gene width category and insertion/match status. The test used compared the log-linear model [1] with all two-way configurations to that with no gene width category and insertion/match status configuration.

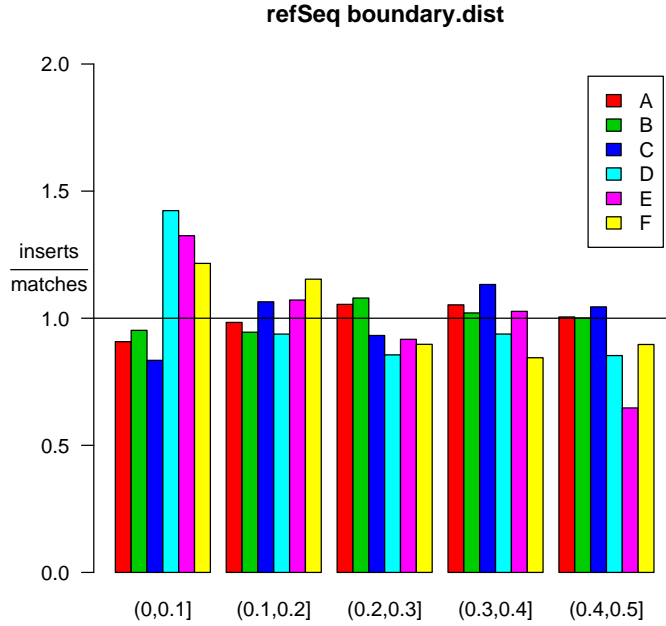


The next plot uses the width of a non-gene region for insertions that fall into such regions.



A	B	C	D	E	F
2.61e-05	4.01e-02	1.21e-03	0.00e+00	1.76e-54	2.40e-19

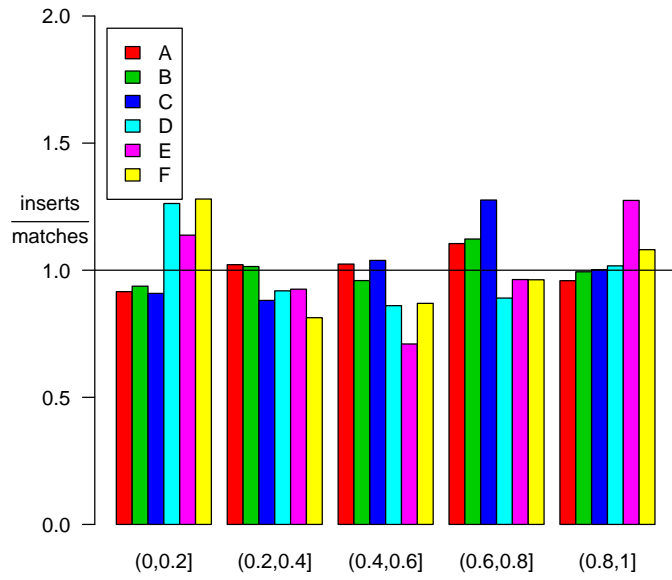
The next plot studies the distance to the nearest boundary between a gene and a non-gene region. The distance is expressed as a fraction of the length of the region. Thus, '0.25' refers to one quarter of the distance from the site to nearest boundary divided by the total width of the region.



A	B	C	D	E	F
1.06e-01	4.75e-01	3.59e-01	2.48e-41	8.21e-05	1.17e-01

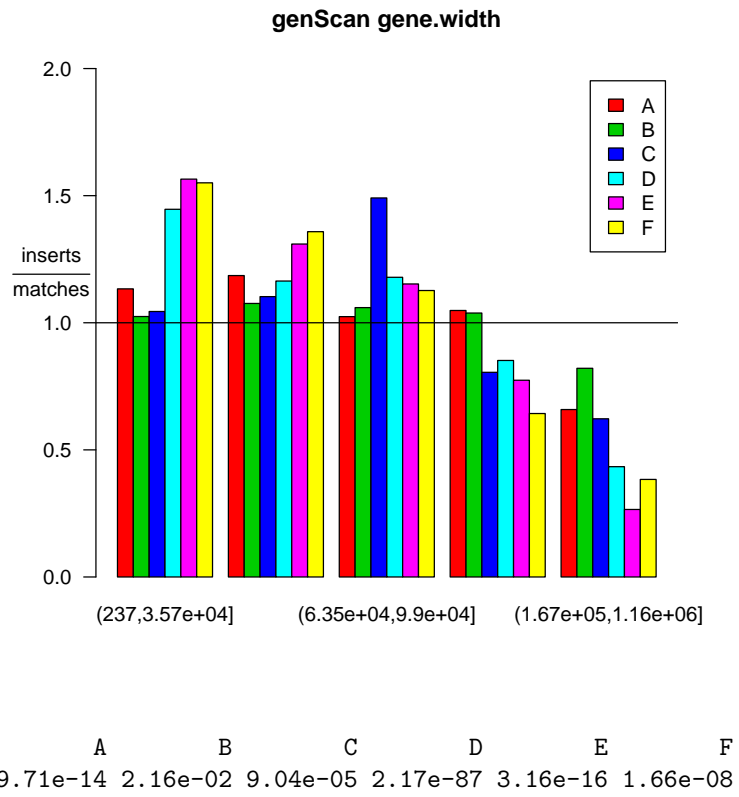
This plot studies the effect of nearness to the beginning of a transcript. For sites in genes, it is the distance to the start of the gene divided by the width of the gene. For other sites it is the distance from the site to the nearer gene if that gene boundary is also a transcription starting point. Locations near '0' are relatively near the beginning of transcription, while those near '1' are near the termination of the transcript.

refSeq start.dist

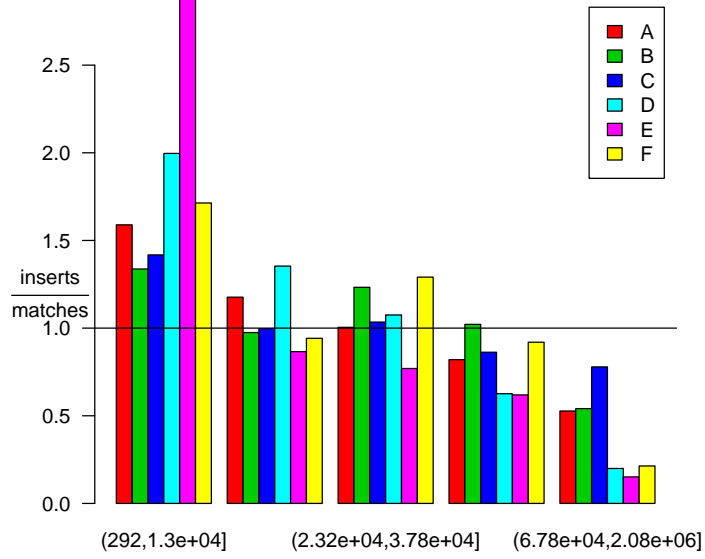


A	B	C	D	E	F
8.41e-02	3.03e-01	1.78e-01	2.33e-19	5.11e-03	1.02e-01

5.2 genScan Annotations

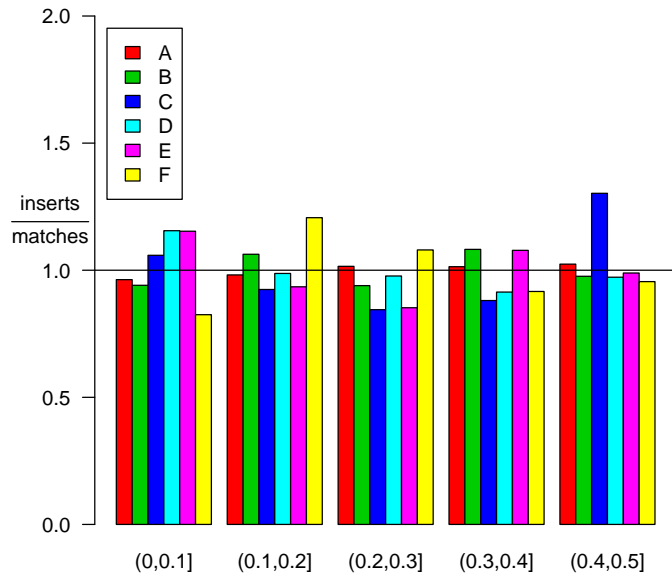


genScan other.width



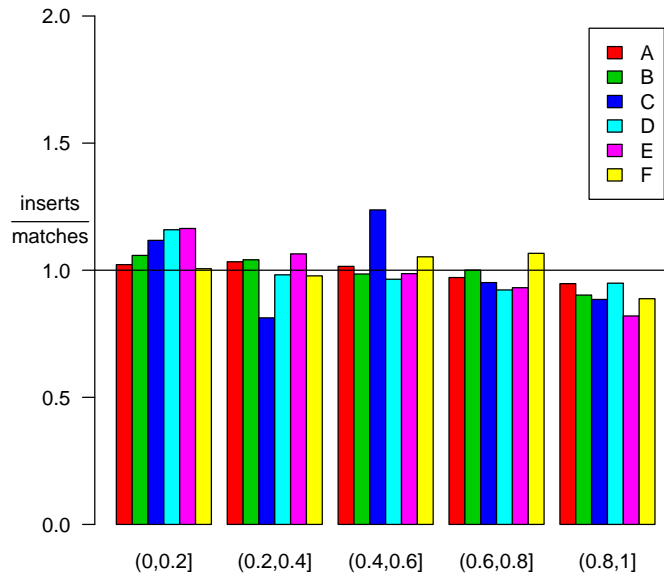
A 1.17e-05 B 7.71e-07 C 6.96e-03 D 1.02e-95 E 1.94e-22 F 1.62e-04

genScan boundary.dist



A 8.63e-01 B 2.27e-01 C 4.28e-02 D 4.64e-06 E 2.54e-01 F 2.34e-01

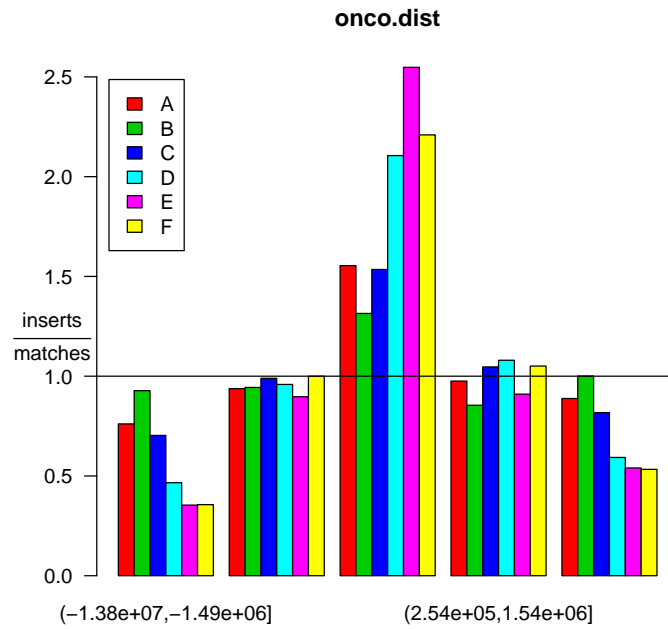
genScan start.dist



A B C D E F
7.09e-01 4.56e-01 1.01e-01 1.16e-06 1.82e-01 9.33e-01

6 Oncogenes

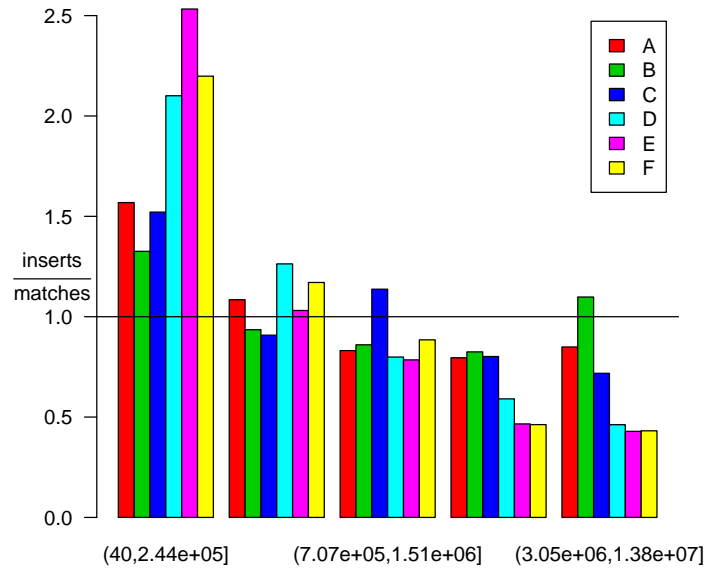
This plot studies the effect of nearness to the 5' end of an oncogene transcript. Positive values represent distances in which the integration site is upstream of the nearest oncogene 5' end, negative downstream.



A	B	C	D	E	F
$5.75e-31$	$8.45e-07$	$2.87e-05$	$7.42e-290$	$1.69e-47$	$3.66e-25$

Here is the same plot using absolute distance

abs abs.onco.dist

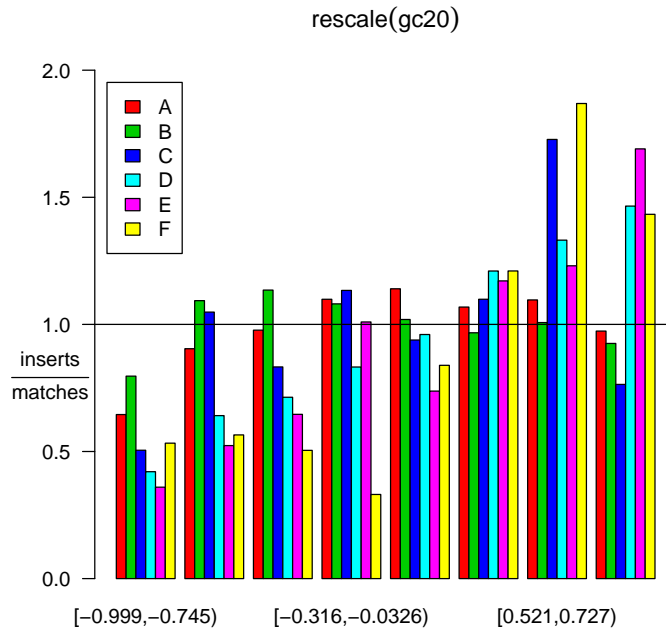


A 6.80e-35 B 1.36e-09 C 2.11e-05 D 7.74e-309 E 7.18e-47 F 5.79e-25

7 GC content

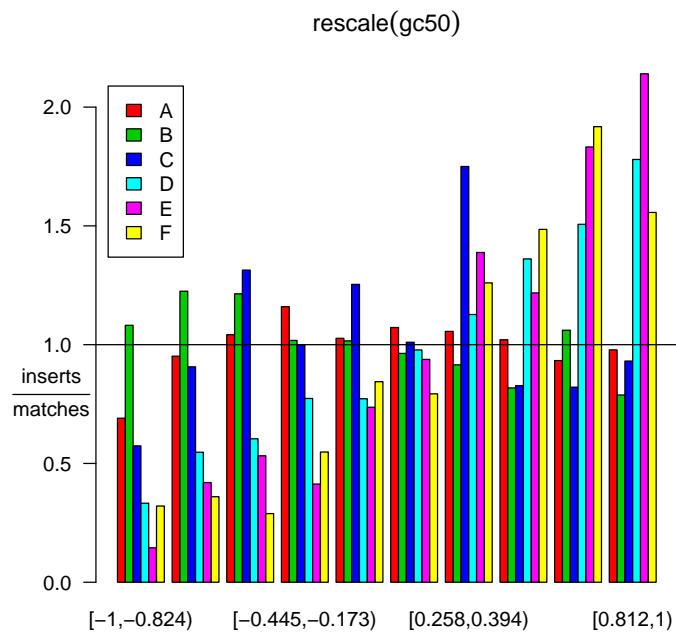
Here we study the effect of GC content on insertion. The GC content is taken from the Mouse Genome Draft at GoldenPath from the table

Following the plot is a table of fitted coefficients based on splitting the GC percent data at the median.

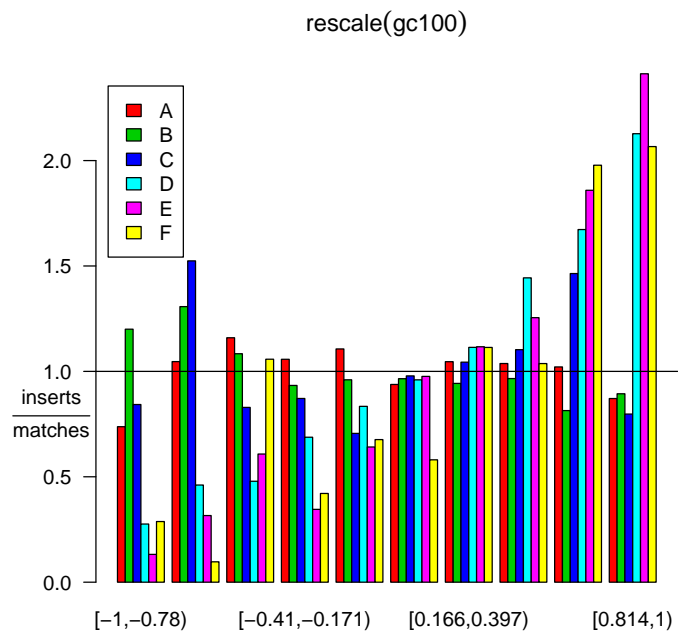


	coef	se	z	p
A	0.0651	0.0479	1.36	1.75e-01
B	-0.0656	0.0544	-1.21	2.28e-01
C	0.2280	0.1100	2.07	3.83e-02
D	0.5950	0.0340	17.50	1.45e-68
E	0.7080	0.0992	7.14	9.28e-13
F	0.9690	0.1280	7.59	3.17e-14

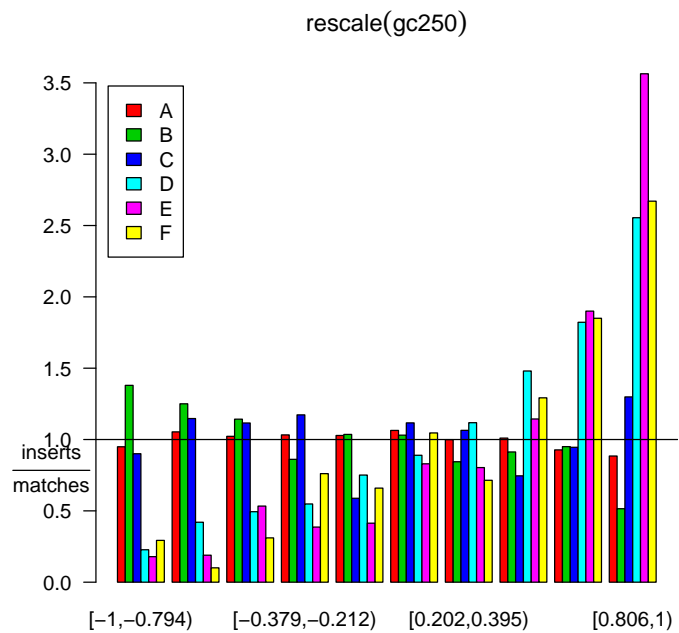
	coef	se	z	p
A	0.00906	0.0476	0.19	8.49e-01
B	-0.20600	0.0544	-3.79	1.49e-04
C	0.12900	0.1090	1.18	2.40e-01
D	0.76700	0.0349	22.00	8.50e-107
E	1.01000	0.1050	9.57	1.10e-21
F	0.95800	0.1310	7.30	2.86e-13



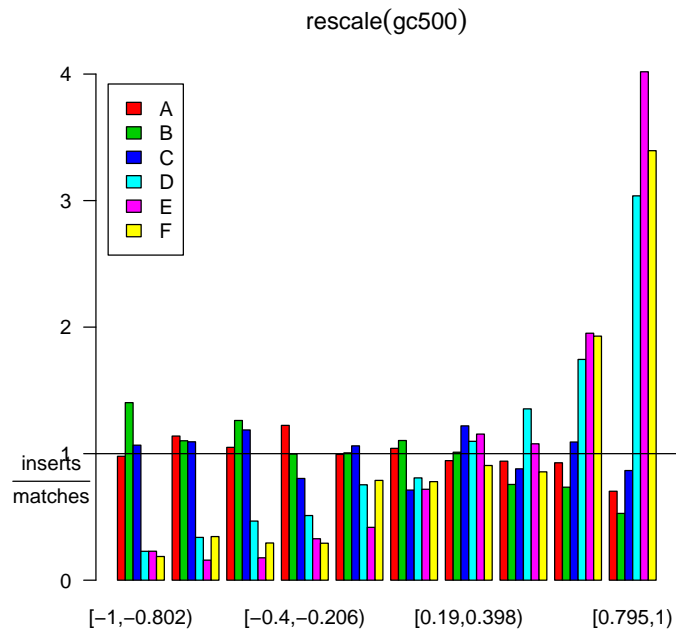
	coef	se	z	p
A	-0.0172	0.0479	-0.359	7.20e-01
B	-0.1800	0.0546	-3.310	9.45e-04
C	0.1610	0.1110	1.450	1.47e-01
D	1.0100	0.0365	27.600	7.04e-168
E	1.2400	0.1120	11.100	1.58e-28
F	1.0100	0.1330	7.610	2.82e-14



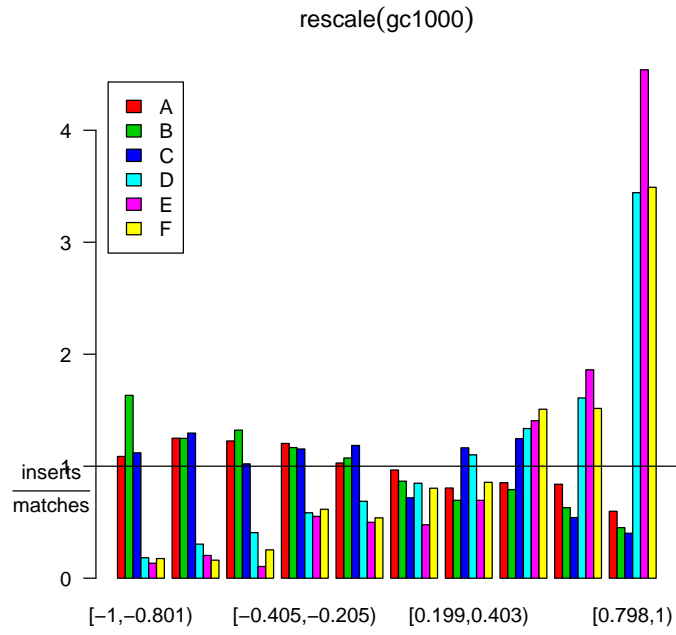
	coef	se	z	p
A	-0.0322	0.0472	-0.683	4.94e-01
B	-0.2670	0.0548	-4.880	1.06e-06
C	0.0142	0.1090	0.129	8.97e-01
D	1.1400	0.0371	30.700	4.87e-207
E	1.5200	0.1200	12.600	2.08e-36
F	1.3100	0.1420	9.230	2.82e-20



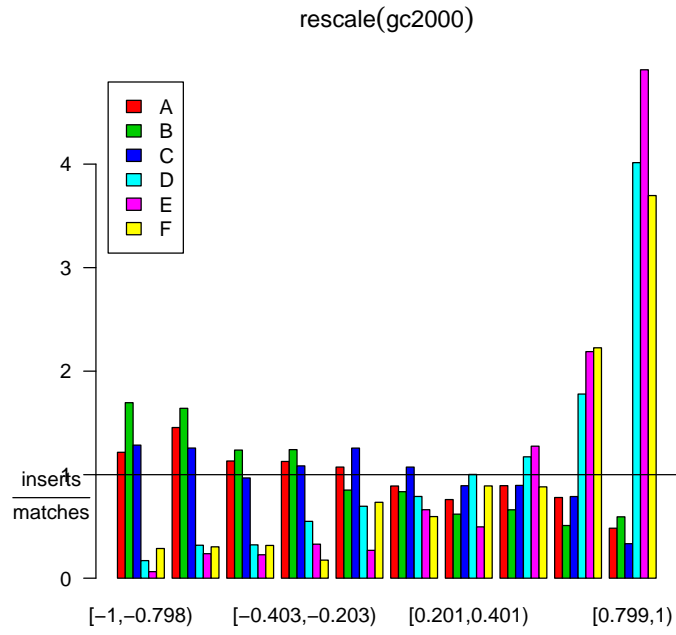
	coef	se	z	p
A	-0.1630	0.0475	-3.430	6.03e-04
B	-0.3380	0.0548	-6.170	6.78e-10
C	-0.0876	0.1100	-0.797	4.26e-01
D	1.2500	0.0384	32.600	4.34e-233
E	1.8400	0.1350	13.600	2.57e-42
F	1.4200	0.1460	9.690	3.41e-22



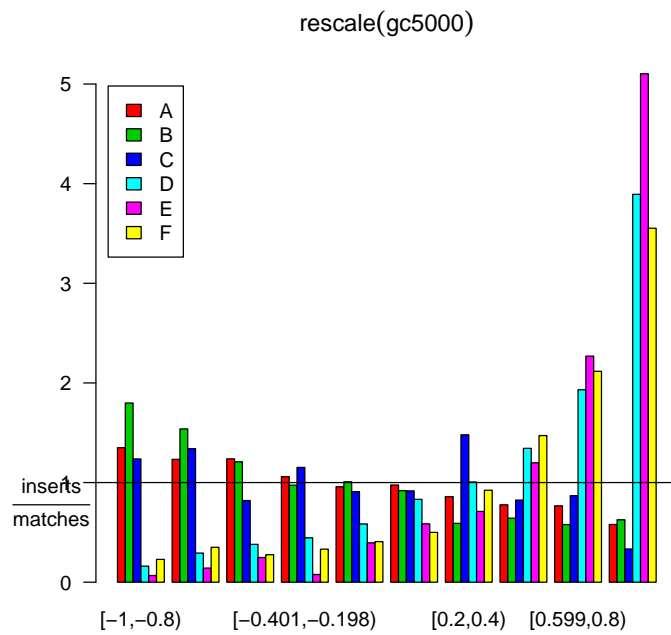
	coef	se	z	p
A	-0.325	0.0478	-6.80	1.07e-11
B	-0.618	0.0568	-10.90	1.24e-27
C	-0.311	0.1120	-2.78	5.47e-03
D	1.310	0.0390	33.60	6.32e-247
E	1.760	0.1320	13.40	7.47e-41
F	1.540	0.1520	10.20	2.91e-24



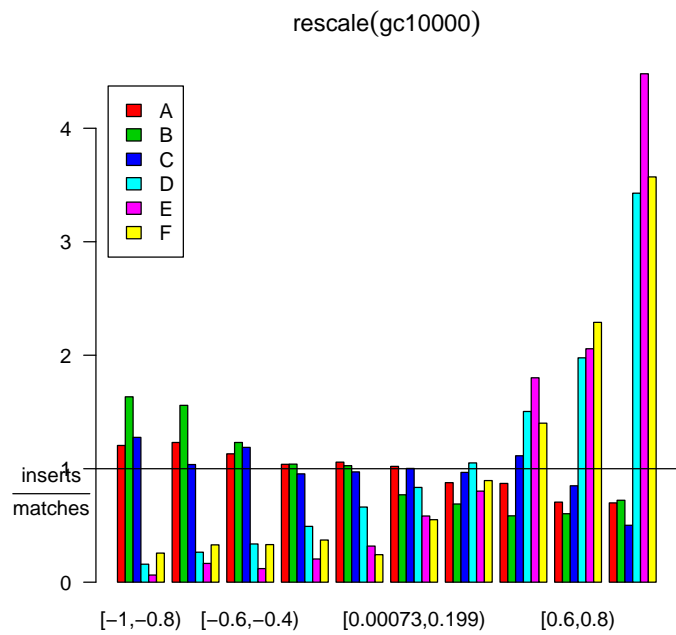
	coef	se	z	p
A	-0.435	0.0484	-9.00	2.33e-19
B	-0.715	0.0576	-12.40	1.93e-35
C	-0.376	0.1140	-3.31	9.17e-04
D	1.400	0.0399	35.00	3.78e-268
E	2.110	0.1480	14.20	7.36e-46
F	1.520	0.1510	10.10	7.05e-24



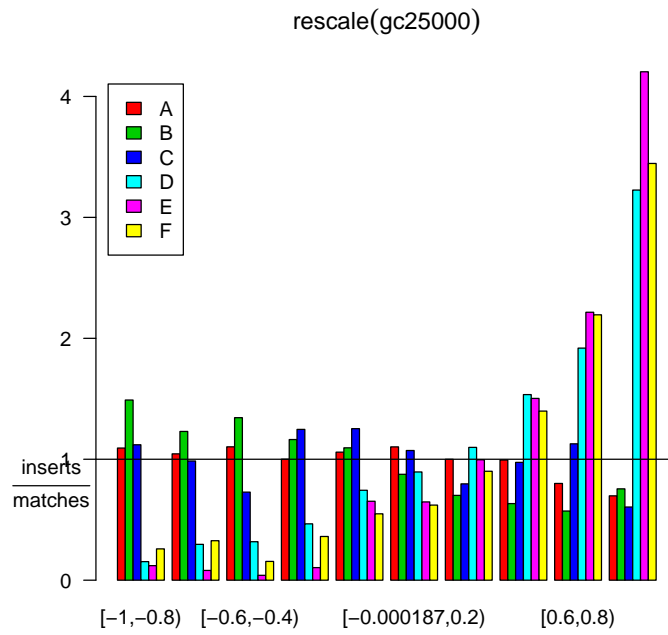
	coef	se	z	p
A	-0.366	0.0480	-7.63	2.33e-14
B	-0.654	0.0572	-11.40	2.81e-30
C	-0.191	0.1110	-1.72	8.61e-02
D	1.510	0.0411	36.80	7.34e-297
E	2.310	0.1600	14.50	1.91e-47
F	1.670	0.1580	10.60	4.65e-26



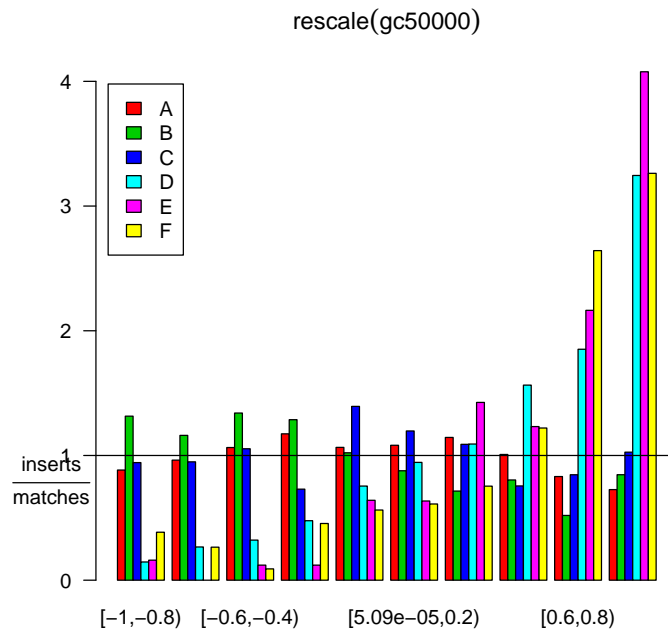
	coef	se	z	p
A	-0.290	0.0476	-6.10	1.09e-09
B	-0.655	0.0570	-11.50	1.58e-30
C	-0.191	0.1110	-1.72	8.49e-02
D	1.470	0.0407	36.20	2.18e-287
E	2.400	0.1630	14.70	9.99e-49
F	1.730	0.1600	10.80	4.97e-27



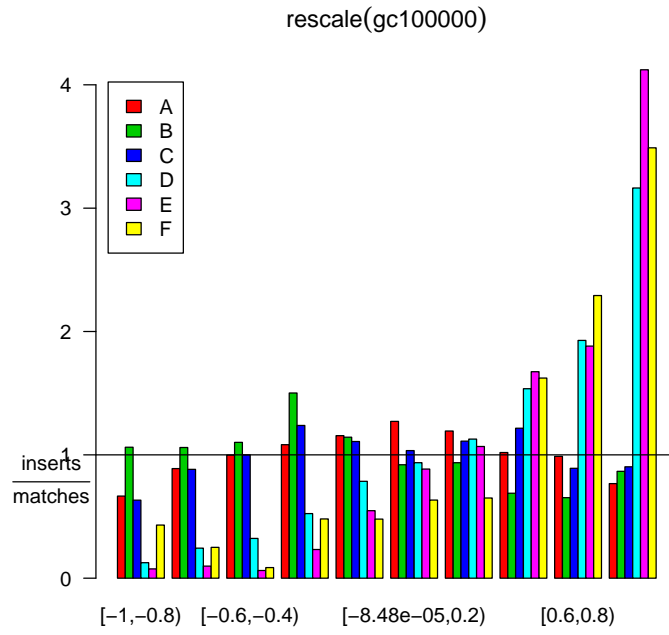
	coef	se	z	p
A	-0.131	0.0473	-2.78	5.49e-03
B	-0.576	0.0562	-10.20	1.21e-24
C	-0.142	0.1100	-1.30	1.95e-01
D	1.450	0.0403	35.90	7.04e-283
E	2.260	0.1550	14.60	5.71e-48
F	1.630	0.1540	10.60	2.17e-26



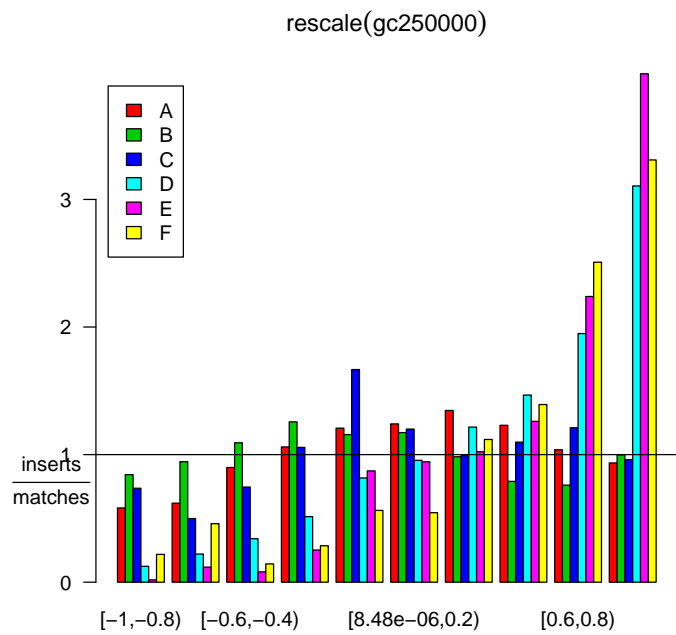
	coef	se	z	p
A	-0.0581	0.0471	-1.230	2.17e-01
B	-0.4880	0.0556	-8.770	1.71e-18
C	-0.0293	0.1090	-0.268	7.89e-01
D	1.4600	0.0405	36.100	1.53e-285
E	2.2500	0.1560	14.500	1.66e-47
F	1.5800	0.1510	10.500	1.36e-25



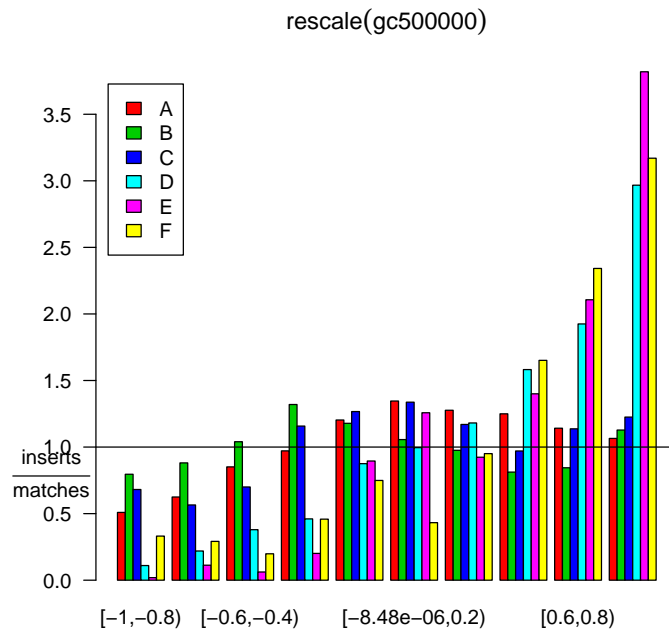
	coef	se	z	p
A	0.1030	0.0469	2.190	2.89e-02
B	-0.3550	0.0547	-6.490	8.68e-11
C	0.0602	0.1090	0.554	5.80e-01
D	1.4600	0.0405	36.000	2.90e-284
E	2.2500	0.1540	14.600	2.87e-48
F	1.6000	0.1510	10.700	1.61e-26



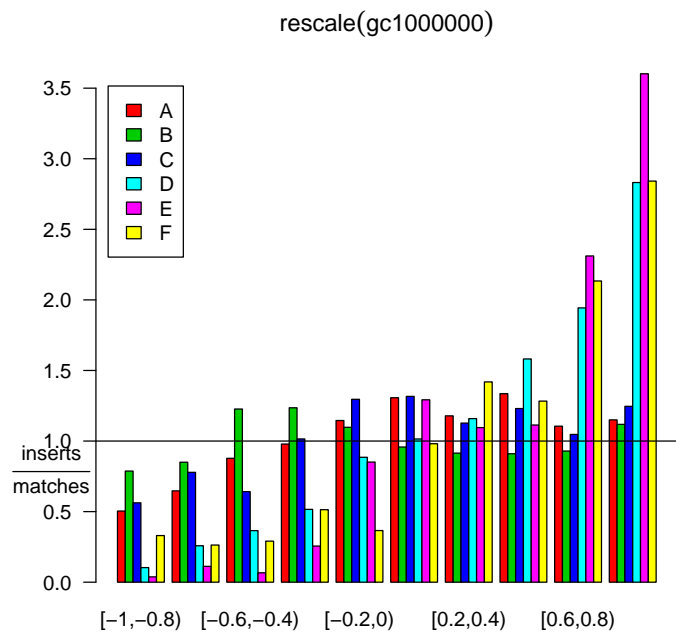
	coef	se	z	p
A	0.299	0.0467	6.41	1.45e-10
B	-0.108	0.0537	-2.02	4.38e-02
C	0.172	0.1080	1.59	1.12e-01
D	1.450	0.0402	36.00	1.55e-284
E	1.960	0.1380	14.20	6.94e-46
F	1.660	0.1510	10.90	6.67e-28



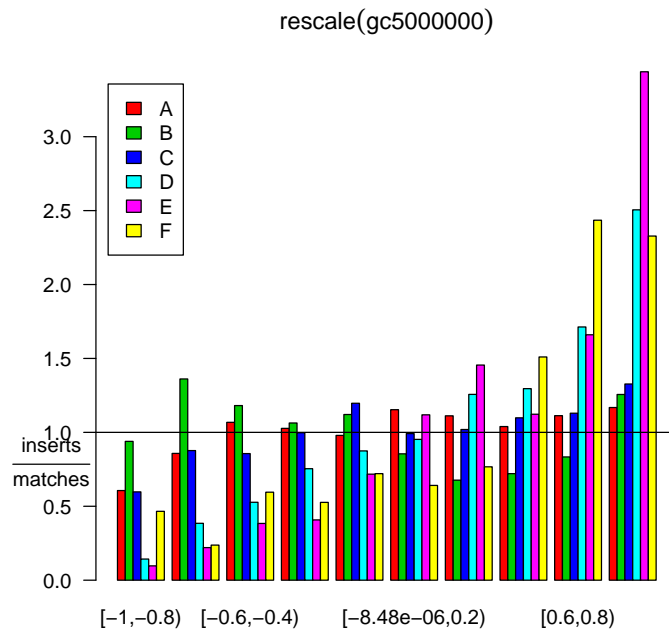
	coef	se	z	p
A	0.3980	0.0469	8.49	2.09e-17
B	-0.0657	0.0536	-1.22	2.21e-01
C	0.3030	0.1090	2.79	5.31e-03
D	1.4500	0.0403	36.00	1.32e-283
E	2.0200	0.1410	14.40	8.39e-47
F	1.4300	0.1420	10.10	7.54e-24



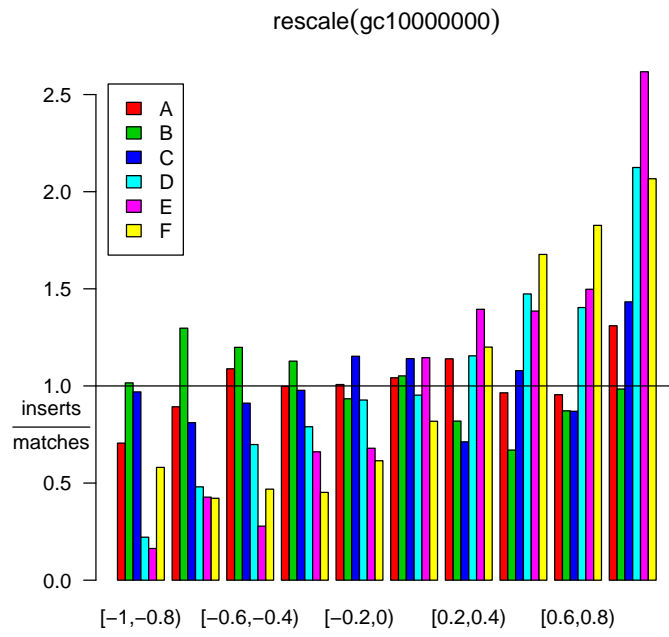
	coef	se	z	p
A	0.3940	0.0470	8.39	4.90e-17
B	-0.0633	0.0536	-1.18	2.37e-01
C	0.3490	0.1090	3.20	1.37e-03
D	1.4000	0.0397	35.20	1.07e-270
E	1.9600	0.1380	14.20	8.78e-46
F	1.5900	0.1490	10.70	1.42e-26



	coef	se	z	p
A	0.220	0.0469	4.69	2.74e-06
B	-0.272	0.0543	-5.01	5.55e-07
C	0.204	0.1080	1.88	6.02e-02
D	1.040	0.0366	28.50	7.66e-179
E	1.540	0.1210	12.80	2.54e-37
F	1.110	0.1330	8.33	8.42e-17



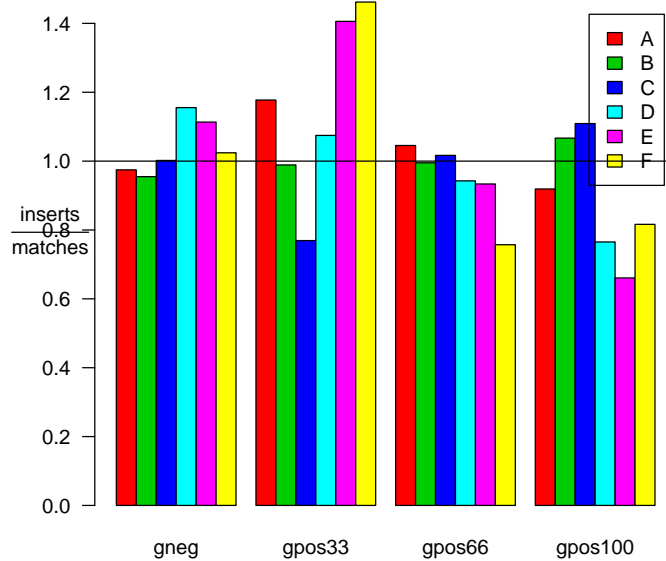
	coef	se	z	p
A	0.1410	0.0467	3.030	2.49e-03
B	-0.2330	0.0540	-4.300	1.68e-05
C	0.0812	0.1080	0.749	4.54e-01
D	0.8170	0.0353	23.200	1.03e-118
E	1.3000	0.1140	11.300	7.62e-30
F	1.0900	0.1330	8.220	2.08e-16



8 Cytobands

Here we study the association of cytoband with insertion intensity. The data are obtained from

<http://genome.ucsc.edu/goldenPath/hg17/database/cytoBand.txt.gz>.



A formal test of significance attains a p-value of $< 2.22e - 16$. Here is the table of coefficients of the log ratio of intensities for true insertion sites versus control insertion sites (comparing each category of Giemsa staining to 'gneg') along with their standard errors, z statistics, and p-values:

	coef	se	z	p
cyto.typegpos100	-0.2080	0.0276	-7.52	5.43e-14
cyto.typegpos33	0.0352	0.0319	1.10	2.69e-01
cyto.typegpos66	-0.0905	0.0372	-2.43	1.50e-02

References

- [1] Yvonne M.M. Bishop, Stephen E. Fienberg, and Paul W. Holland. *Discrete multivariate analyses: Theory and practice* (MIT Press, 1975).
- [2] P. McCullagh and John A. Nelder. *Generalized linear models*. (Chapman & Hall ltd, 1999).

- [3] Xiaolin Wu, Yuan Li, Bruce Crise, Shawn M. Burgess “Transcription Start Regions in the Human Genome Are Favored Targets for MLV Integration,” *Science*, **300**(5626), (June 2003): 1749-1751.