

VanillaICE: Hidden Markov Models for the Assessment of Chromosomal Alterations using High-throughput SNP Arrays

Robert Scharpf

May 12, 2008

1 Introduction

Chromosomal DNA is characterized by variation between individuals at the level of entire chromosomes (e.g. aneuploidy in which the chromosome copy number is altered), segmental changes (including insertions, deletions, inversions, and translocations), and changes to small genomic regions (including single nucleotide polymorphisms). A variety of alterations that occur in chromosomal DNA, many of which can be detected using high density single nucleotide polymorphism (SNP) microarrays, are linked to normal variation as well as disease and therefore of particular interest. These include changes in copy number (deletions and duplications) and genotype (e.g. the occurrence of regions of homozygosity). Hidden Markov models (HMM) are particularly useful for detecting such abnormalities, modeling the spatial dependence between neighboring SNPs. Here, we extend previous approaches that utilize HMM frameworks for inference in high throughput SNP arrays by integrating copy number, genotype calls, and the corresponding measures of uncertainty when available. Using simulated and real data, we demonstrate how confidence scores control smoothing in a probabilistic framework. The goal of this vignette is to provide a simple interface for fitting HMMs and plotting functions to help visualize the predicted states alongside the experimental data.

2 Simple Usage

2.1 Vanilla HMM

```
> library(VanillaICE)
```

Before fitting the HMM, the data must be organized into one of the following classes for high-throughput SNP data:

- `SnpCallSet` (genotype calls)
- `SnpCopyNumberSet` (copy number estimates)
- `oligoSnpSet` (genotype calls and copy number estimates)
- `RatioSnpSet` (genotype calls and ratios of A and B allele intensities)

When pre-processing Affymetrix SNP chips with the R package *oligo*, an object of one of the above classes is created and can be used directly with the HMMs described in this vignette. These classes can also be created from Illumina data as described in the *IlluminaHowTo* vignette. The simulated data provided with this package is an instance of class `oligoSnpSet`

```

> data(chromosome1)
> annotation(chromosome1)

[1] "pd.mapping50k.hind240,pd.mapping50k.xba240"

> chromosome1

oligoSnpSet (storageMode: lockedEnvironment)
assayData: 9165 features, 1 samples
  element names: calls, callsConfidence, cnConfidence, copyNumber
experimentData: use 'experimentData(object)'
Annotation: pd.mapping50k.hind240,pd.mapping50k.xba240
phenoData
An object of class "AnnotatedDataFrame"
  sampleNames: NA06993
  varLabels and varMetadata description:
    family: trio variable
    upd: uniparental isodisomy indicator
featureData
An object of class "AnnotatedDataFrame"
  rowNames: SNP_A-1677174, SNP_A-1718890, ..., SNP_A-1677548 (9165 total)
  varLabels and varMetadata description:
    dbsnp_rs_id: dbsnp_rs_id
    chromosome: chrom
    ...: ...
    enzyme: enzyme
    (8 total)
Annotation [1] "pd.mapping50k.hind240,pd.mapping50k.xba240"

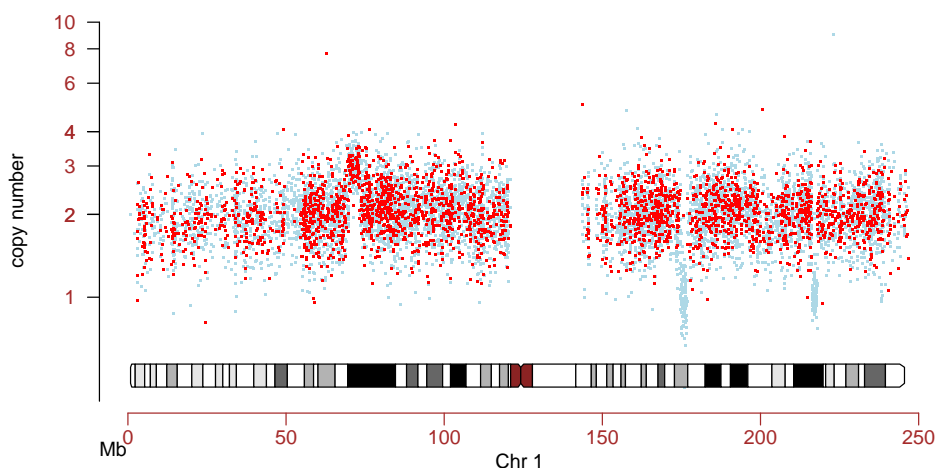
```

Before deciding whether to fit a HMM, plot the data:

```

> gp <- plotSnp(chromosome1)
> show(gp)

```



The HMM assumes that the copy number estimates, conditional on the hidden state, are approximately Gaussian. A log transformation will be performed automatically provided that the copy number estimates are all positive. See the documentation pages in the R package *VanillaICE* for more information about the `chromosome1` example dataset. Fitting a hidden Markov model requires the following components:

- the hidden states
- the emission probabilities
- transition probabilities

Here, we illustrate how one may specify the hidden states and calculate the emission and transition probabilities below. We first create an instance of the class `HmmOptions` that contains the hidden states and, if copy number estimates are available, the location parameters for the Gaussian distribution. Though we will model the logarithm of the copy number estimates as approximately normally distributed, it is more natural to specify the location parameters on the original scale (e.g., 1 copy is deletion (D), 2 copies is normal (N), etc.). Emission probabilities for the genotype calls are a homozygous genotype call (`probHomCall` argument). Again, the ordering must correspond to the order of the hidden states. computed from the supplied probabilities of Note that the ordering of the numeric vector for `copyNumber.location` must correspond to the ordering of the states. A number of validity checks are performed on this object and it is a good idea to check whether the instantiated object is valid.

```
> options <- new("HmmOptions", states = c("D", "N", "L",
+     "A"), snpset = chromosome1, copyNumber.location = c(1,
+     2, 2, 3), probHomCall = c(0.99, 0.7, 0.99, 0.7))
> validObject(options)
```

```
[1] TRUE
```

From an object of class `HmmOptions`, we can compute the emission probabilities. The emission probabilities are returned on the log scale. Conditional on the hidden state, we assume that the copy number and genotype are independent. Therefore, the emission probabilities for an HMM that models the copy number and genotypes jointly are computed by adding the emission probabilities for copy number and genotype.

```
> params <- new("HmmParameter", states = states(options),
+     initialStateProbability = 0.99)
> cn.emission <- copyNumber.emission(options)

[1] "Calculating emission probabilities on the log(copy number)"

> gt.emission <- calls.emission(options)
> emission(params) <- cn.emission + gt.emission

> genomicDistance(params) <- exp(-2 * physicalDistance(options)/(100 *
+     1e+06))
> transitionScale(params) <- scaleTransitionProbability(options)
> class(params)

[1] "HmmParameter"
attr(,"package")
[1] "VanillaICE"
```

We may then fit the HMM by

```

> fit <- hmm(options, params)

> fit

HmmPredict (storageMode: lockedEnvironment)
assayData: 9165 features, 1 samples
  element names: predictions
phenoData
  sampleNames: NA06993
  varLabels and varMetadata description:
    family: trio variable
    upd: uniparental isodisomy indicator
featureData
  featureNames: SNP_A-1677174, SNP_A-1718890, ..., SNP_A-1677548 (9165 total)
  fvarLabels and fvarMetadata description:
    dbsnp_rs_id: dbsnp_rs_id
    chromosome: chrom
    ...: ...
    arm: NA
    (9 total)
experimentData: use 'experimentData(object)'
Annotation: pd.mapping50k.hind240,pd.mapping50k.xba240
hidden states: D N L A
breakpoints:
'data.frame':      11 obs. of  9 variables:
 $ id   : chr  "NA06993" "NA06993" "NA06993" "NA06993" ...
 $ chr  : chr  "1" "1" "1" "1" ...
 $ state: chr  "N" "L" "N" "A" ...
 $ size : num  48.71  4.81 15.34  3.30 101.64 ...
 $ N    : num  997 102 902 202 3797 ...
 $ start: num  0.837 49.598 54.499 69.854 73.174 ...
 $ last : num  49.5 54.4 69.8 73.2 174.8 ...
 $ prev : num    NA 49.5 54.4 69.8 73.2 ...
 $ next : num  49.6 54.4 69.9 73.2 174.8 ...

> breakpoints(fit)

      id chr state      size      N      start      last
1.1 NA06993  1    N 48.708312  997   0.836727 49.54504
1.2 NA06993  1    L  4.811945  102  49.597810 54.40975
1.3 NA06993  1    N 15.339118  902  54.498950 69.83807
1.4 NA06993  1    A  3.299434  202  69.854466 73.15390
1.5 NA06993  1    N 101.640222 3797  73.174070 174.81429
1.6 NA06993  1    D  1.985684  100 174.815096 176.80078
1.7 NA06993  1    N 39.439073 1900 176.800844 216.23992
1.8 NA06993  1    D  1.586808   99 216.286002 217.87281
1.9 NA06993  1    N 20.410491  901 217.892177 238.30267
1.10 NA06993  1    D  0.097921   5 238.319943 238.41786
1.11 NA06993  1    N  8.431130  160 238.429864 246.86099
      prev      next
1.1      NA 49.59781
1.2 49.54504 54.40975

```

```

1.3  54.40975  69.85447
1.4  69.83807  73.17407
1.5  73.15390 174.81429
1.6 174.81429 176.80084
1.7 176.80078 216.28600
1.8 216.23992 217.89218
1.9 217.87281 238.31994
1.10 238.30267 238.41786
1.11 238.41786 246.86099

```

Summary statistics for the breakpoints (more useful when multiple chromosomes and samples are in the `fit` object)

```
> summary(fit)
```

```
$L
```

```

      Freq med(length) avg(length) sd(length) med(n.snp) avg(n.snp)
1      1    4811.945    4811.945         NA         102         102
      sd(n.snp)
1             NA

```

```
$A
```

```

      Freq med(length) avg(length) sd(length) med(n.snp) avg(n.snp)
1      1    3299.434    3299.434         NA         202         202
      sd(n.snp)
1             NA

```

```
$D
```

```

      Freq med(length) avg(length) sd(length) med(n.snp) avg(n.snp)
1      3    1586.808    1223.471    994.949         99         68
      sd(n.snp)
1      54.562

```

See [2] for a more complete description of the simulated dataset and the features detected by this HMM. We may plot the data along with the predictions as follows:

```
> gp <- plotSnp(snpset(options), fit)
```

```

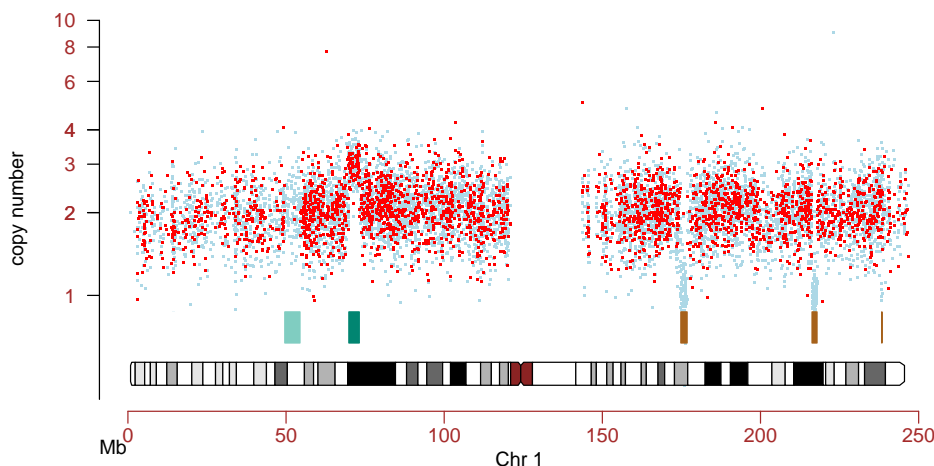
[1] "col.predict not specified in list of graphical parameters. Using the following colors:"
[1] "#A6611A" "white"   "#80CDC1" "#018571"

```

```

> gp$abline.v <- TRUE
> show(gp)

```



2.2 Integrating Confidence Estimates (ICE)

In this section, we illustrate how one may fit an HMM that incorporates confidence estimates of the SNP-level summaries for genotype calls and copy number. Confidence scores (inverse of standard errors) are available for this object (see Section 2.3 for how confidence scores were derived).

This information is incorporated into the HMM emission probabilities. Probably the easiest way to do that is recreate the options object, and then recalculate the emission probabilities.

```
> options <- new("HmmOptions", snpset = chromosome1, states = c("D",
+   "N", "L", "A"), copyNumber.location = c(1, 2, 2,
+   3), copyNumber.ICE = TRUE, probHomCall = c(0.99,
+   0.75, 0.99, 0.75))
> cn.emission <- copyNumber.emission(options)

[1] "Calculating emission probabilities on the log(copy number)"
[1] "Using 1/cnConfidence(object) as standard errors for the copy number"

> emission(params) <- cn.emission + gt.emission
> fit.ice <- hmm(options, params)

[1] "Transforming copy number to log2 scale."
[1] "Fitting HMM to sample 1"
```

We may also incorporate the confidence scores for the genotype calls by specifying `calls.ICE=TRUE`. This feature is only available for the Affymetrix 100k and 500k platforms. The slot `probHomCall` stores user-specified probabilities of $P(\text{call is AA or BB} \mid \text{state is LOH})$ and $P(\text{call is AA or BB} \mid \text{state is normal})$. These probabilities must be specified in this order. The emission probabilities for the genotype calls will only be calculated for the states LOH (LOH is defined as a stretch of homozygous genotype calls longer than what one would expect by chance) and Normal refers to typical ratios of heterozygous to homozygous genotype calls. The slot `term5` contains user-specified probabilities for the $P(\text{true genotype is HET} \mid \text{genotype call is AB, hidden state is Normal})$ and $P(\text{true genotype is HET} \mid \text{genotype call is AA or BB, hidden state is Normal})$, respectively. Default values are provided when not specified, as the following example illustrates.

```

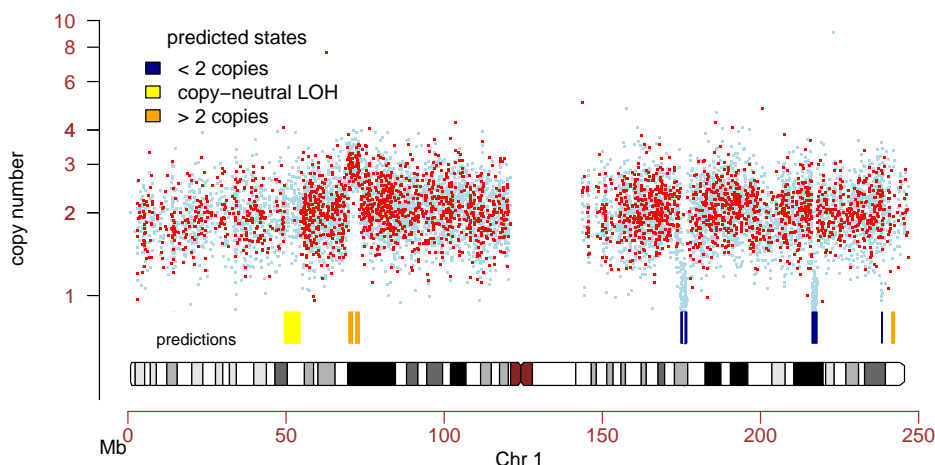
> options <- new("HmmOptions", snpset = chromosome1, states = c("D",
+   "N", "L", "A"), copyNumber.location = c(1, 2, 2,
+   3), copyNumber.ICE = TRUE, calls.ICE = TRUE, probHomCall = c(0.99,
+   0.75))
> params <- new("HmmParameter", states = states(options),
+   initialStateProbability = 0.99)
> cn.emission <- copyNumber.emission(options)
> genomicDistance(params) <- exp(-2 * physicalDistance(options)/(100 *
+   1e+06))
> transitionScale(params) <- scaleTransitionProbability(options)
> gt.emit <- calls.emission(options)
> gt.emission <- array(NA, dim(cn.emission))
> gt.emission[, , 1:2] <- gt.emit
> gt.emission[, , 3:4] <- gt.emit
> emission(params) <- cn.emission + gt.emission
> fit.ice <- hmm(options, params)

> gp <- plotSnp(snpset(options), fit.ice)

[1] "col.predict not specified in list of graphical parameters. Using the following colors:"
[1] "#A6611A" "white"    "#80CDC1" "#018571"

> gp$abline.v <- TRUE
> gp$col.predict <- c("darkblue", "white", "yellow", "orange")
> show(gp)
> legend(-0.05, 10, fill = gp$col.predict[c(1, 3, 4)],
+   legend = c("< 2 copies", "copy-neutral LOH", "> 2 copies"),
+   bty = "n", title = "predicted states")
> legend(0, 0.8, legend = "predictions", bty = "n", cex = 0.8,
+   adj = 0)

```



Note that the ICE HMM correctly identifies the simulated normal segments in features B and C (the normal segments were simulated to have high confidence scores). Additionally, the ICE HMM detects the micro-amplification in region E (also simulated to have high confidence scores).

2.3 Confidence scores

Confidence scores for genotype calls We suggest using the CRLMM algorithm [1] for genotype calls. CRLMM (in the R package *oligo*) provides confidence scores ($S_{\widehat{GT}}$) of the genotype estimates (\widehat{GT}). From 269 HapMap samples assayed on the Affymetrix 50k Xba and Hind chips, we have a gold standard of the true genotype defined by the consensus of the HapMap centers. We use kernel based density estimates to obtain

$$f\{S_{\widehat{HOM}} | \widehat{HOM}, HOM\}, f\{S_{\widehat{HOM}} | \widehat{HOM}, HET\}, f\{S_{\widehat{HET}} | \widehat{HET}, HOM\}, \text{ and } f\{S_{\widehat{HET}} | \widehat{HET}, HET\} \quad (1)$$

separately for the Xba and Hind 50k chips. The first term in (1), for example, denotes the density of the scores when the genotype is correctly called homozygous (\widehat{HOM}) and the true genotype is homozygous (HOM). See [2] for a more complete description of the methods. The data needed to estimate these densities is stored in the experiment data package *callsConfidence*. *callsConfidence* is available from the author's website.

Confidence scores for copy number estimates To illustrate how standard errors of the copy number estimate could be integrated in the HMM, the R object `chromosome1` contains standard errors simulated from a shifted Gamma: $\text{Gamma}(1, 2) + 0.3$, where 1 is the shape parameter and 2 is the rate parameter. To ascertain the effect of qualitatively high confidence scores on the ICE HMM, we scaled a robust estimate of the copy number standard deviation by $\frac{1}{2}$. Similarly, to simulate less precise \widehat{CN} we scaled ϵ by 2. For more detailed information about how the data in the `chromosome1` was generated, see the documentation for this object in the R package *VanillaICE*.

3 The HmmParameter class

An instance of the class is created by the method `new`:

```
> new("HmmParameter")
```

The object `params` contains all of the parameters needed for fitting the HMM, including an estimate of the genomic distance between SNPs (used for calculating SNP-specific transition probabilities), emission probabilities (slot: `emission`), and initial state probabilities.

Emission probabilities. The emission probabilities are stored as an array in the `params` object. The emission probability array has dimension $R \times C \times S$, where S is the number of hidden states, R is the number of rows (SNPs), and C is the number of samples. One may use `[` to subset object of class `HmmParameter`.

```
> params[5, 1, ]
```

```
Formal class 'HmmParameter' [package "VanillaICE"] with 5 slots
..@ states                : chr [1:4] "D" "N" "L" "A"
..@ initialStateProbability: num [1:4] 0.00333 0.99000 0.00333 0.00333
..@ emission              : num [1, 1, 1:4] -1.20 -1.54 -1.19 -1.86
.. ..- attr(*, "dimnames")=List of 3
.. .. ..$ : chr "SNP_A-1662392"
.. .. ..$ : chr "NA06993"
.. .. ..$ : chr [1:4] "D" "N" "L" "A"
..@ genomicDistance       : num [1:3] 1 1 1
```



```

..@ transitionScale      : num [1:4, 1:4] 1 1 0.75 0.75 1.5 1 1.5 1.5 0.75 1 ...
.. ..- attr(*, "dimnames")=List of 2
.. .. ..$ : chr [1:4] "D" "N" "L" "A"
.. .. ..$ : chr [1:4] "D" "N" "L" "A"

```

Transition probabilities. The probability of remaining in the same state, $P(S_t = S_{t+1})$ (the diagonal of the transition probability matrix) is a function of the distance (d) between SNPs: $e^{-2d(100*1e6)}$. This value is stored in the slot **tau** of the **params** object. The probability of leaving a state is ϵ , where $\epsilon = 1 - P(S_t = S_{t+1})$. The ϵ is split among $S - 1$ states. By default, the probability of transitioning from an altered state back to the normal state is twice as likely as the probability of transitioning between two altered states. The weights for ϵ are provided in the **tau.scale** matrix in the R object **params**

```
> transitionScale(params)
```

```

      D    N    L    A
D 1.00 1.5 0.75 0.75
N 1.00 1.0 1.00 1.00
L 0.75 1.5 1.00 0.75
A 0.75 1.5 0.75 1.00

```

and can be adjusted by the **SCALE** argument. For illustration, one could make the probability of transitioning from an altered state to a normal state 10 times as likely as the probability of transitioning between two altered states by the following command:

```

> transitionScale(params) <- scaleTransitionProbability(params,
+   SCALE = 10)
> transitionScale(params)

```

```

      D    N    L    A
D 1.00 2.5 0.25 0.25
N 1.00 1.0 1.00 1.00
L 0.25 2.5 1.00 0.25
A 0.25 2.5 0.25 1.00

```

4 The HmmOptions Class

To be completed ...

5 The HmmPredict Class

The output from the HMM is an instance of the **HmmPredict** class and contains the predicted states as well as the breakpoints for the different states.

```
> fit
```

```

HmmPredict (storageMode: lockedEnvironment)
assayData: 9165 features, 1 samples
  element names: predictions
phenoData
  sampleNames: NA06993
  varLabels and varMetadata description:

```

```

    family: trio variable
    upd: uniparental isodisomy indicator
featureData
  featureNames: SNP_A-1677174, SNP_A-1718890, ..., SNP_A-1677548 (9165 total)
  fvarLabels and fvarMetadata description:
    dbsnp_rs_id: dbsnp_rs_id
    chromosome: chrom
    ...: ...
    arm: NA
    (9 total)
experimentData: use 'experimentData(object)'
Annotation: pd.mapping50k.hind240,pd.mapping50k.xba240
hidden states: D N L A
breakpoints:
'data.frame':      11 obs. of  9 variables:
 $ id   : chr  "NA06993" "NA06993" "NA06993" "NA06993" ...
 $ chr  : chr  "1" "1" "1" "1" ...
 $ state: chr  "N" "L" "N" "A" ...
 $ size : num  48.71  4.81 15.34  3.30 101.64 ...
 $ N    : num  997 102 902 202 3797 ...
 $ start: num  0.837 49.598 54.499 69.854 73.174 ...
 $ last : num  49.5 54.4 69.8 73.2 174.8 ...
 $ prev : num    NA 49.5 54.4 69.8 73.2 ...
 $ next : num  49.6 54.4 69.9 73.2 174.8 ...

```

The breakpoints are provided as a `data.frame`:

```

> breaks <- breakpoints(fit)
> breaks <- breaks[breaks[, "state"] != "N", ]

```

One may order the altered states from biggest to smallest for each chromosome as follows:

```

> breaks[order(breaks[, "chr"], breaks[, "size"], decreasing = TRUE),
+ ]

```

| | id | chr | state | size | N | start | last | prev |
|------|-----------|-----|-------|----------|-----|-----------|-----------|-----------|
| 1.2 | NA06993 | 1 | L | 4.811945 | 102 | 49.59781 | 54.40975 | 49.54504 |
| 1.4 | NA06993 | 1 | A | 3.299434 | 202 | 69.85447 | 73.15390 | 69.83807 |
| 1.6 | NA06993 | 1 | D | 1.985684 | 100 | 174.81510 | 176.80078 | 174.81429 |
| 1.8 | NA06993 | 1 | D | 1.586808 | 99 | 216.28600 | 217.87281 | 216.23992 |
| 1.10 | NA06993 | 1 | D | 0.097921 | 5 | 238.31994 | 238.41786 | 238.30267 |
| | next | | | | | | | |
| 1.2 | 54.40975 | | | | | | | |
| 1.4 | 73.17407 | | | | | | | |
| 1.6 | 176.80084 | | | | | | | |
| 1.8 | 217.89218 | | | | | | | |
| 1.10 | 238.41786 | | | | | | | |

The `summary` method returns a list where each element in the list provides statistics for an altered states. For each chromosome, the mean, median, and standard deviation of the size of the features and number of SNPs involved are reported. For instance, in this HMM the altered states were loss of heterozygosity (L), amplification of copy number (A), and deletion of copy number (D). Because we only have one sample and one chromosome in the `fit` example, the `summary` method is not that useful:

```
> summary(fit)

$L
  Freq med(length) avg(length) sd(length) med(n.snp) avg(n.snp)
1    1    4811.945    4811.945         NA         102         102
  sd(n.snp)
1          NA

$A
  Freq med(length) avg(length) sd(length) med(n.snp) avg(n.snp)
1    1    3299.434    3299.434         NA         202         202
  sd(n.snp)
1          NA

$D
  Freq med(length) avg(length) sd(length) med(n.snp) avg(n.snp)
1    3    1586.808    1223.471    994.949         99         68
  sd(n.snp)
1    54.562
```

6 HMMs for different classes of data

6.1 Copy number

The method `hmm` has a different set of underlying hidden states depending on whether copy number estimates, genotype calls, or both are available. When only copy number estimates are available, the hidden states (for autosomes) are hemizygous or homozygous deletion (one or fewer copies), normal (two copies), and amplification (three or more copies). The corresponding data class is `SnpcopyNumberSet`. To illustrate, we convert the `chromosome1` example to an object of this class and fit the HMM.

```
> chr1.cn <- as(chromosome1, "SnpcopyNumberSet")
> options <- new("HmmOptions", snpset = chr1.cn, states = c("D",
+   "N", "A"), copyNumber.location = 1:3)
> params.cn <- new("HmmParameter", states = c("D", "N",
+   "A"))
> emission(params.cn) <- copyNumber.emission(options)

[1] "Calculating emission probabilities on the log(copy number)"

> transitionScale(params.cn) <- scaleTransitionProbability(options)
> genomicDistance(params.cn) <- exp(-2 * physicalDistance(options)/(100 *
+   1e+06))
> fit.cn <- hmm(options, params.cn)

[1] "Transforming copy number to log2 scale."
[1] "Fitting HMM to sample 1"

> breakpoints(fit.cn)

      id chr state      size    N      start      last
1.1 NA06993   1    N 43.865725  877   0.836727 44.70245
1.2 NA06993   1    D  0.040126   2 44.722116 44.76224
```

| | | | | | | | |
|------|---------|---|---|------------|------|------------|-----------|
| 1.3 | NA06993 | 1 | N | 25.058717 | 1122 | 44.779351 | 69.83807 |
| 1.4 | NA06993 | 1 | A | 3.299434 | 202 | 69.854466 | 73.15390 |
| 1.5 | NA06993 | 1 | N | 101.640222 | 3797 | 73.174070 | 174.81429 |
| 1.6 | NA06993 | 1 | D | 1.985748 | 101 | 174.815096 | 176.80084 |
| 1.7 | NA06993 | 1 | N | 39.313361 | 1899 | 176.926556 | 216.23992 |
| 1.8 | NA06993 | 1 | D | 1.586808 | 99 | 216.286002 | 217.87281 |
| 1.9 | NA06993 | 1 | N | 20.410491 | 901 | 217.892177 | 238.30267 |
| 1.10 | NA06993 | 1 | D | 0.097921 | 5 | 238.319943 | 238.41786 |
| 1.11 | NA06993 | 1 | N | 8.431130 | 160 | 238.429864 | 246.86099 |

| | prev | next |
|------|-----------|-----------|
| 1.1 | NA | 44.70245 |
| 1.2 | 44.70245 | 44.76224 |
| 1.3 | 44.76224 | 69.85447 |
| 1.4 | 69.83807 | 73.17407 |
| 1.5 | 73.15390 | 174.81429 |
| 1.6 | 174.81429 | 176.80084 |
| 1.7 | 176.80084 | 216.28600 |
| 1.8 | 216.23992 | 217.89218 |
| 1.9 | 217.87281 | 238.31994 |
| 1.10 | 238.30267 | 238.41786 |
| 1.11 | 238.41786 | 246.86099 |

```
> graph.par <- plotSnp(snpset(options), fit.cn)
```

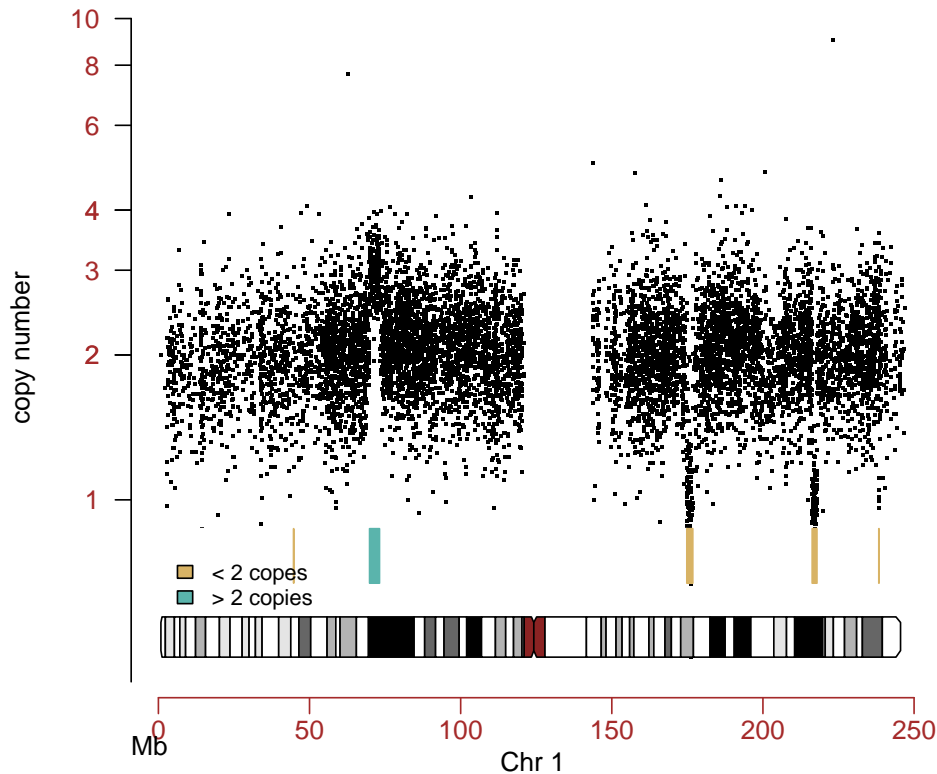
```
[1] "col.predict not specified in list of graphical parameters. Using the following colors:"
```

```
[1] "#D8B365" "white"    "#5AB4AC"
```

```
> graph.par$abline.v <- FALSE
```

```
> print(graph.par)
```

```
> legend(0, 0.8, fill = graph.par$col.predict[c(1, 3)],
+       legend = c("< 2 copes", "> 2 copies"), bty = "n",
+       cex = 0.8)
```



6.2 Genotype calls

When only genotype calls are available, the hidden states are loss and retention (ret) of heterozygosity. We define *loss* to be a sequence of homozygous SNPs longer than what we would expect to observe by chance. Note that many long stretches of homozygosity may occur as a result of a population sharing a common underlying haplotype structure; loss predictions from an HMM fit to an individual do not necessarily reflect the 'loss' of an allele in that individual. For illustration, we convert the `chromosome1` example to an object of class `HmmSnpCallSet` and refit the HMM.

```
> chr1.calls <- as(chromosome1, "SnpCallSet")
> options <- new("HmmOptions", snpset = chr1.calls, states = c("L",
+   "N"), probHomCall = c(0.99, 0.7))
> params.calls <- new("HmmParameter", states = states(options))
> transitionScale(params.calls) <- scaleTransitionProbability(options)
> genomicDistance(params.calls) <- exp(-2 * physicalDistance(options)/(100 *
+   1e+06))
> emission(params.calls) <- calls.emission(options)
> fit.calls <- hmm(options, params.calls)
```

```
[1] "Fitting HMM to sample 1"
```

```

> breakpoints(fit.calls)

      id chr state      size      N      start      last      prev
1.1 NA06993 1      N 48.708312 997 0.836727 49.54504      NA
1.2 NA06993 1      L 4.811945 102 49.597810 54.40975 49.54504
1.3 NA06993 1      N 119.810058 4890 54.498950 174.30901 54.40975
1.4 NA06993 1      L 2.285950 109 174.418117 176.70407 174.30901
1.5 NA06993 1      N 39.439518 1902 176.800399 216.23992 176.70407
1.6 NA06993 1      L 1.485826 97 216.286002 217.77183 216.23992
1.7 NA06993 1      N 28.988981 1068 217.872013 246.86099 217.77183

      next
1.1 49.59781
1.2 54.40975
1.3 174.41812
1.4 176.80040
1.5 216.28600
1.6 217.87201
1.7 246.86099

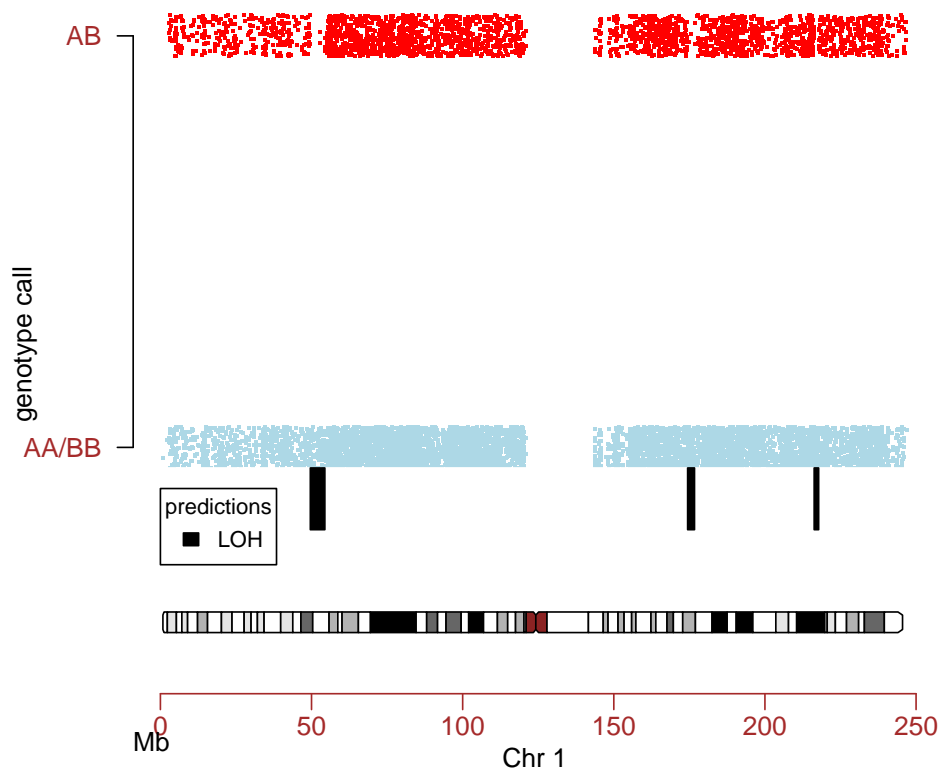
> gp <- plotSnp(snpset(options), fit.calls)

[1] "col.predict not specified in list of graphical parameters. Using the following colors:"
[1] "black" "white"

> gp$col.predict <- c("black", "white")
> gp$ylim <- c(-0.5, 1)
> gp$add.centromere <- FALSE
> gp$abline.v <- TRUE
> gp$cytoband.ycoords <- c(-0.45, -0.4)
> gp$hmm.ycoords <- c(-0.2, -0.05)

> show(gp)
> legend(0, -0.1, legend = "LOH", fill = "black", title = "predictions",
+       bty = "o", cex = 0.8)

```



6.3 Genotype calls and copy number

Section 2 illustrates how one may fit the HMM to objects of class `oligoSnpSet`.

More documentation about the classes can be found in the documentation for the R package *VanillaICE*.

7 Session Information

The version number of R and packages loaded for generating the vignette were:

- R version 2.7.0 (2008-04-22), i386-pc-mingw32
- Locale: LC_COLLATE=English_United States.1252;LC_CTYPE=English_United States.1252;LC_MONETARY=English_United States.1252
- Base packages: base, datasets, graphics, grDevices, methods, stats, tools, utils
- Other packages: Biobase 2.0.1, oligoClasses 1.2.0, RColorBrewer 1.0-2, SNPchip 1.4.0, VanillaICE 1.2.0

References

- [1] Benilton Carvalho, Henrik Bengtsson, Terence P Speed, and Rafael A Irizarry. Exploration, normalization, and genotype calls of high-density oligonucleotide SNP array data. *Biostatistics*, 8(2):485–499, Apr 2007.
- [2] Robert B Scharpf, Giovanni Parmigiani, Jonathan Pevsner, and Ingo Ruczinski. A hidden Markov model for joint estimation of genotype and copy number in high-throughput SNP chips. Technical Report Working Paper 136, Johns Hopkins University, February 2007.