

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

Robert Scharpf

April 22, 2010

Abstract

This package provides an implementation of a hidden Markov Model for high throughput SNP arrays. Users of this package should already have available locus-level estimates of copy number. Copy number estimates can be relative or absolute.

1 Overview

This vignette requires that you have

- an absolute estimate of the *total* copy number organized such that rows correspond to loci and columns correspond to samples
and / or
- a matrix of genotype calls (1=AA, 2 = AB, 3= BB): rows correspond to loci and columns correspond to samples

Additional options that can improve the HMM predictions include

- a CRLMM confidence score of the genotype call
- standard errors of the copy number estimates

Other HMM implementations are available for the joint analysis of copy number and genotype, including QuantiSNP [1] and PennCNV [5].

Data considerations. The HMM implemented in this package is most relevant for heritable diseases for which integer copy numbers are expected. For somatic cell diseases such as cancer, we suggest circular binary segmentation, as implemented in the R package DNACopy [3].

Citing this software. Robert B Scharpf, Giovanni Parmigiani, Jonathan Pevsner, and Ingo Ruczinski. Hidden Markov models for the assessment of chromosomal alterations using high-throughput SNP arrays. *Annals of Applied Statistics*, 2(2):687–713, 2008.

2 Organizing the locus-level data

This package includes simulated genotype and copy number data for approximately 9165 SNPs on chromosome 1 and 100 SNPs on chromosome 2.

```
> library(VanillaICE)
```

```

Loading package ff 2.1-2
-getOption("fftempdir")==="/tmp/RtmpeE5qrS"
-getOption("ffextension")==".ff"
-getOption("ffdrop")==TRUE
-getOption("fffinonexit")==TRUE
-getOption("ffpagesize")==65536
-getOption("ffcaching")== "mmnoflush" -- consider "ffeachflush" if your system stalls on large writes
-getOption("ffbatchbytes")==16777216 -- consider a different value for tuning your system

> data(locusLevelData)

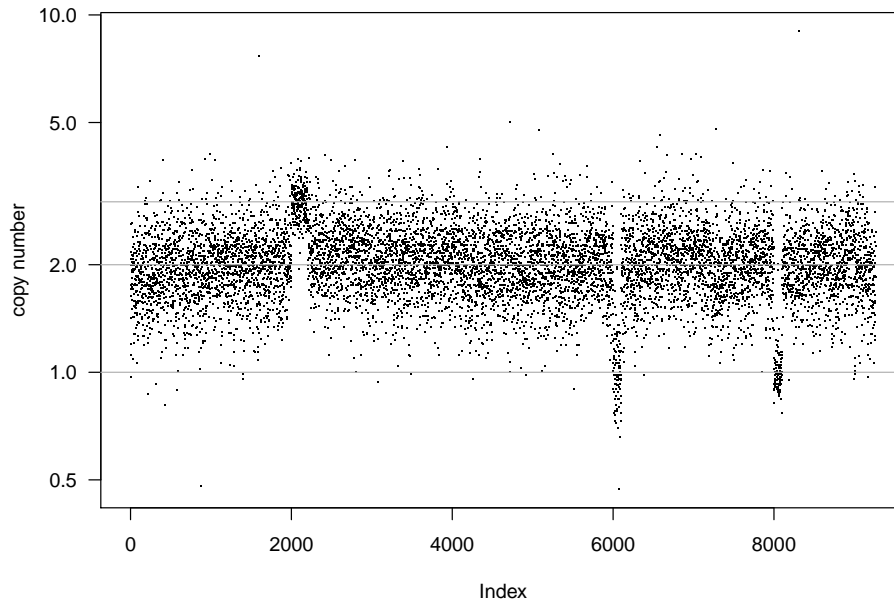
```

(The copy number estimates in the locusLevelData object were multiplied by 100 and saved as an integer.) Verify that it is reasonable to assume integer copy number for the HMM by plotting the locus-level estimates as a function of the physical position.

```

> par(las = 1)
> plot(locusLevelData[["copynumber"]][, 1]/100, pch = ".",
+       ylab = "copy number", log = "y")
> abline(h = 1:3, col = "grey70")

```



Next, create an object of class oligoSnpSet from the simulated data:

```

> oligoSet <- new("oligoSnpSet", copyNumber = locusLevelData[["copynumber"]]/100,
+               call = locusLevelData[["genotypes"]], callProbability = locusLevelData[["crlmmConfidence"]],
+               annotation = locusLevelData[["platform"]])

```

If confidence scores or inverse standard errors for the copy number estimates are available, these should be supplied to the `cnConfidence` slot in the `assayData`. For illustration, in the following code chunk we transform the copy number estimates to the log scale and calculate a robust estimate of the standard deviation. If uncertainty estimates are not available for copy number, the HMM will calculate the median absolute deviation (MAD). See the the function `robustSds`.

```
> sds <- robustSds(log2(locusLevelData[["copynumber"]]/100))
```

The inverse of the `sds` object can be assigned to the `cnConfidence` slot.

3 Fitting the HMM

3.1 Vanilla HMM

When jointly modeling the copy number and genotype data, we assume that the genotype estimates and copy number estimates are independent conditional on the underlying hidden state. The emission probabilities for the genotypes are then calculated using either (i) assumptions of the probability of observing a homozygous genotype call given the underlying state. Note that the SNPs should be ordered by chromosome and physical position.

```
> copyNumber(oligoSet) <- log2(copyNumber(oligoSet))
> oligoSet <- oligoSet[order(chromosome(oligoSet), position(oligoSet)),
+   ]
> hmmOpts <- hmm.setup(oligoSet, copynumberStates = log2(c(1,
+   2, 2, 3)), states = c("hem-del", "ROH", "normal",
+   "amp"), normalIndex = 3, log.initialP = rep(log(1/4),
+   4), prGenotypeHomozygous = c(0.99, 0.99, 0.7, 0.7))
```

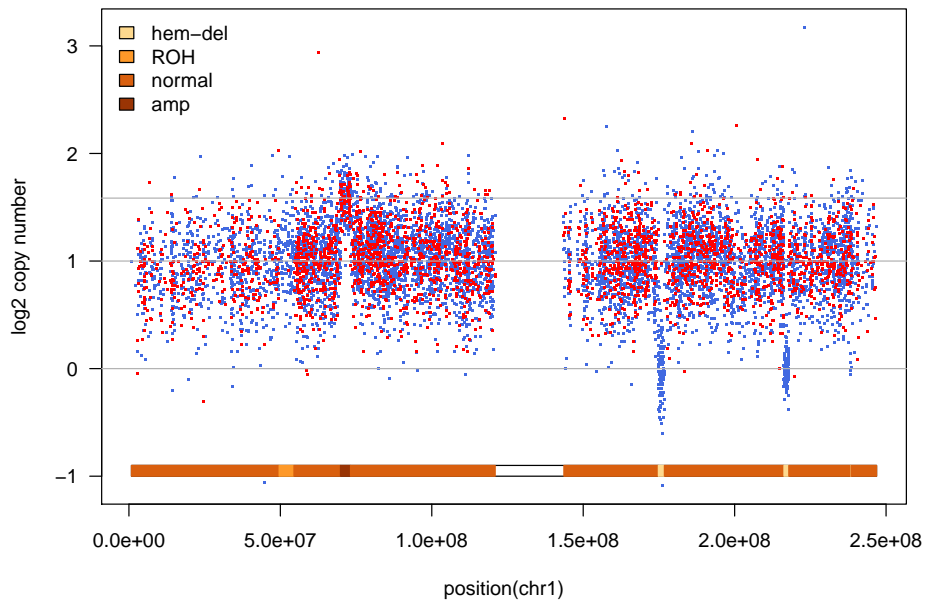
The log-scale emission probabilities:

```
> dim(hmmOpts$log.emission)
```

```
[1] 9265    1    4
```

The viterbi algorithm is used to obtain the most likely sequence of hidden states given the observed data. For efficiency, we return an object of class `RangedData` with genomic coordinates of the normal and altered regions. We also return the log-likelihood ratio (LLR) of the predicted sequence in an interval versus the null of normal copy number. For intervals with typical copy number (2) and percent heterozygosity (the 3rd state in the above codechunk), the LLR is zero.

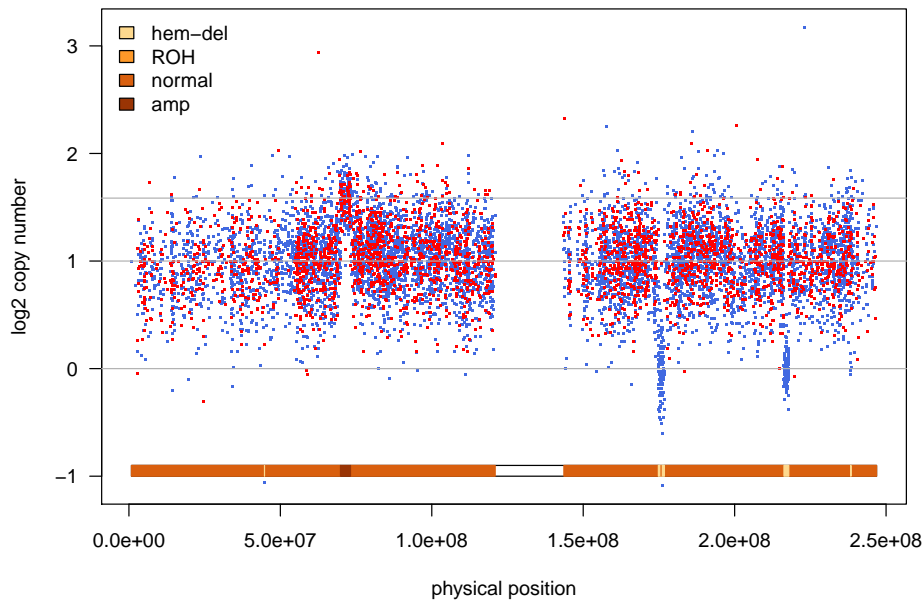
```
> fit.van <- hmm(oligoSet, hmmOpts)
```



3.2 ICE HMM

To compute emission probabilities that incorporate the *crmm* genotype confidence scores, (i) set `ICE` to `TRUE` in the `hmm.setup` function and (ii) indicate which of the states are expected to be largely homozygous (`rohStates`).

```
> hmmOpts <- hmm.setup(oligoSet, ICE = TRUE, copynumberStates = log2(c(1,
+ 2, 2, 3)), states = c("hem-del", "ROH", "normal",
+ "amp"), normalIndex = 3, log.initialP = rep(log(1/4),
+ 4), rohStates = c(TRUE, TRUE, FALSE, FALSE))
> fit.ice <- hmm(oligoSet, hmmOpts)
```



3.3 Other options

Copy number. A HMM for copy number only (e.g., if genotypes are ignored or are unavailable) can be fit as follows.

```
> cnSet <- new("CopyNumberSet", copyNumber = log2(locusLevelData[["copynumber"]]/100),
+   annotation = locusLevelData[["platform"]])
> cnSet <- cnSet[order(chromosome(cnSet), position(cnSet)),
+   ]
> hmmOpts <- hmm.setup(cnSet, copynumberStates = log2(c(0,
+   1, 2, 3)), states = c("hom-del", "hem-del", "normal",
+   "amp"), normalIndex = 3, log.initialP = rep(log(1/4),
+   4))
> fit.cn <- hmm(cnSet, hmmOpts)
```

Regions of homozygosity. A HMM for genotype-only data can be used to find long stretches of homozygosity. Note that hemizygous deletions are also identified as 'ROH' when copy number is ignored (as the biallelic genotype call in a hemizygous deletions tends to be all homozygous calls).

```
> snpSet <- new("SnpSet", call = locusLevelData[["genotypes"]],
+   callProbability = locusLevelData[["crlmmConfidence"]],
+   annotation = locusLevelData[["platform"]])
> featureData(snpSet) <- addFeatureAnnotation(snpSet)
> fvarLabels(snpSet)

[1] "chromosome" "position"   "isSnp"

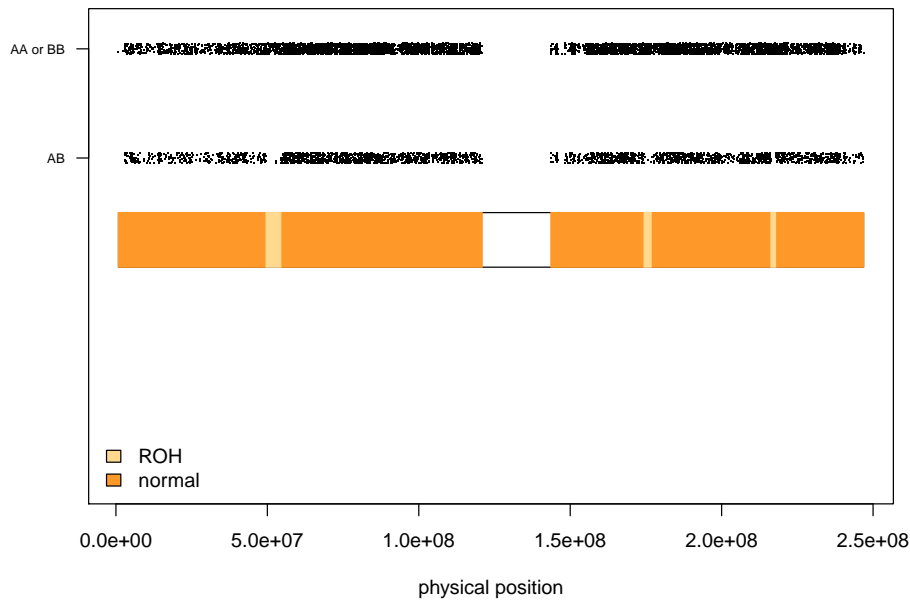
> snpSet <- snpSet[order(chromosome(snpSet), position(snpSet)),
+   ]
```

```

> hmmOpts <- hmm.setup(snpSet, states = c("ROH", "normal"),
+   normalIndex = 2, log.initialP = rep(log(1/2), 2),
+   prGenotypeHomozygous = c(0.99, 0.7), TAUP = 5e+07)
> fit.gt <- hmm(snpSet, hmmOpts)

```

A suggested visualization:



4 Quality control

4.1 Outliers

Copy number outliers can cause the HMM to become too jumpy. One approach to reduce the influence of outliers is to use some *light*-smoothing prior to fitting the HMM, as suggested in the R package `DNAcopy`. For instance, one could identify outliers by some criteria and then average the outliers using the estimates from neighboring probes. Here, we use the defaults in `smooth.CNA`.

```

> if (require("DNAcopy")) {
+   copyNumber(cnSet)[50] <- 10
+   copyNumber(cnSet)[45:55]
+   cnaObj <- CNA(genomdat = copyNumber(cnSet), chrom = chromosome(cnSet),
+     maploc = position(cnSet), data.type = "logratio",
+     sampleid = sampleNames(cnSet))
+   smoothed.cnaObj <- smooth.CNA(cnaObj)
+   copyNumber(cnSet) <- matrix(smoothed.cnaObj[[3]],
+     nrow(cnSet), 1)
+   copyNumber(cnSet)[50]
+ }

```

```
[1] 1.700717
```

One could also increase the value of TAUP in the viterbi algorithm to encourage a fit with fewer jumps. Note that with improved estimates of copy number uncertainty, many of these *post-hoc* approaches for addressing outliers would be less critical.

4.2 Batch effects

See the technical report available from the author's website [4].

5 Alternatives

See the `crlmmDownstream` vignette located in `inst/test` for an alternative approach that computes emission probabilities directly on the bivariate normal $\log(A)$ versus $\log(B)$ space [2]. A technical report describing the estimation procedure for copy number in `crlmm` is available from the author's website [4].

6 Session Information

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

- R version 2.11.0 (2010-04-22), x86_64-unknown-linux-gnu
- Locale: LC_CTYPE=en_US, LC_NUMERIC=C, LC_TIME=en_US, LC_COLLATE=en_US, LC_MONETARY=C, LC_MESSAGES=en_US, LC_PAPER=en_US, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=en_US, LC_IDENTIFICATION=C
- Base packages: base, datasets, graphics, grDevices, methods, stats, tools, utils
- Other packages: Biobase 2.8.0, crlmm 1.6.0, DBI 0.2-5, DNACopy 1.22.0, IRanges 1.6.0, oligo 1.12.0, oligoClasses 1.10.0, pd.mapping50k.hind240 1.0.0, pd.mapping50k.xba240 1.0.0, RColorBrewer 1.0-2, RSQLite 0.8-4, VanillaICE 1.10.0
- Loaded via a namespace (and not attached): affxparser 1.20.0, affyio 1.16.0, annotate 1.26.0, AnnotationDbi 1.10.0, Biostrings 2.16.0, bit 1.1-4, ellipse 0.3-5, ff 2.1-2, genefilter 1.30.0, mvtnorm 0.9-9, preprocessCore 1.10.0, SNPchip 1.12.0, splines 2.11.0, survival 2.35-8, xtable 1.5-6

References

- [1] Stefano Colella, Christopher Yau, Jennifer M Taylor, Ghazala Mirza, Helen Butler, Penny Clouston, Anne S Bassett, Anneke Seller, Christopher C Holmes, and Jiannis Ragoussis. QuantiSNP: an Objective Bayes Hidden-Markov Model to detect and accurately map copy number variation using SNP genotyping data. *Nucleic Acids Res*, 35(6):2013–2025, 2007.
- [2] Joshua M Korn, Finny G Kuruvilla, Steven A McCarroll, Alec Wysoker, James Nemes, Simon Cawley, Earl Hubbell, Jim Veitch, Patrick J Collins, Katayoon Darvishi, Charles Lee, Marcia M Nizzari, Stacey B Gabriel, Shaun Purcell, Mark J Daly, and David Altshuler. Integrated genotype calling and association analysis of snps, common copy number polymorphisms and rare cnvs. *Nat Genet*, 40(10):1253–1260, Oct 2008.
- [3] Adam B Olshen, E. S. Venkatraman, Robert Lucito, and Michael Wigler. Circular binary segmentation for the analysis of array-based DNA copy number data. *Biostatistics*, 5(4):557–72, Oct 2004.

- [4] Robert B Scharpf, Ingo Ruczinski, Benilton Carvalho, Betty Doan, Aravinda Chakravarti, and Rafael Irizarry. A multilevel model to address batch effects in copy number estimation using snp arrays. Submitted Revised Manuscript, Jan 2010.
- [5] Kai Wang, Mingyao Li, Dexter Hadley, Rui Liu, Joseph Glessner, Struan F A Grant, Hakon Hakonarson, and Maja Bucan. Penncnv: an integrated hidden markov model designed for high-resolution copy number variation detection in whole-genome snp genotyping data. *Genome Res*, 17(11):1665–1674, Nov 2007.