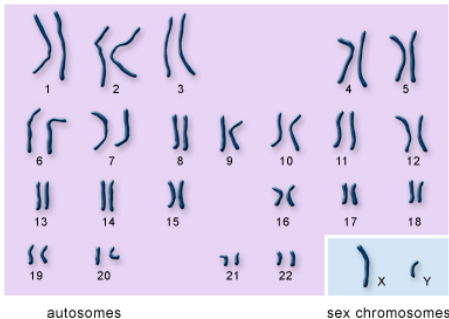# An Integrated Approach for the Assessment of Chromosomal Abnormalities

Ingo Ruczinski

Department of Biostatistics
Johns Hopkins Bloomberg School of Public Health
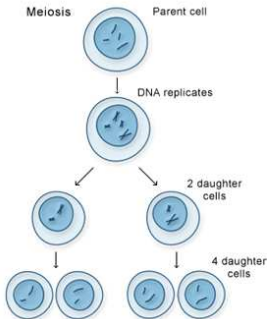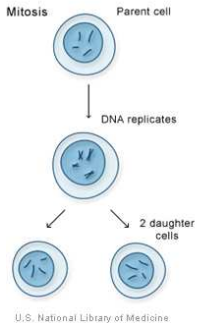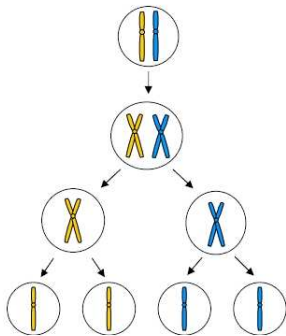
June 6, 2007

autosomes                    sex chromosomes

U.S. National Library of Medicine

U.S. National Library of Medicine

Missense mutation

Original DNA code for an amino acid sequence.

DNA → C A T C A T C A T C A T C A T C A T C A T
bases

His His His His His His His

Amino acid

Replacement of a single nucleotide.

C A T C A T C A T C C T C A T C A T C A T

His His His Pro His His His

Incorrect amino acid, which may produce a malfunctioning protein.

U.S. National Library of Medicine

# Cancer samples

# SNPchip S4 classes and methods

```
> merged <- subset(merged, samples = 1)
> summary(merged)
$NA06985
                     Chr 1 Chr 2 Chr 3 Chr 4 Chr 5 Chr 6 Chr 7 Chr 8 Chr 9
mean copy number      2.06  2.07  2.09  2.15  2.08  2.09  2.09  2.06  2.06
sd copy number        0.45  0.45  0.43  0.46  0.45  0.46  0.45  0.43  0.45
% heterozygous calls  0.26  0.26  0.27  0.27  0.27  0.29  0.28  0.28  0.27
% homozygous calls    0.73  0.72  0.72  0.73  0.72  0.69  0.71  0.71  0.72
% no calls            0.01  0.01  0.01  0.01  0.01  0.01  0.01  0.01  0.01
                     Chr 10 Chr 11 Chr 12 Chr 13 Chr 14 Chr 15 Chr 16 Chr 17
mean copy number       2.04   2.06   2.06   2.11   2.09   2.02   2.00   1.98
sd copy number         0.45   0.45   0.46   0.46   0.46   0.43   0.44   0.44
% heterozygous calls   0.30   0.28   0.25   0.27   0.26   0.25   0.26   0.29
% homozygous calls     0.69   0.70   0.74   0.72   0.73   0.74   0.72   0.70
% no calls             0.01   0.01   0.01   0.01   0.01   0.01   0.01   0.01
                     Chr 18 Chr 19 Chr 20 Chr 21 Chr 22 Chr X Total (autosomes)
mean copy number       2.09   1.98   1.98   2.10   1.92  2.13               2.07
sd copy number         0.46   0.44   0.43   0.47   0.44  0.44               0.45
% heterozygous calls   0.26   0.24   0.30   0.27   0.26  0.26               0.27
% homozygous calls     0.73   0.74   0.68   0.72   0.72  0.73               0.72
% no calls             0.01   0.02   0.01   0.02   0.02  0.01               0.01

>
```

**1** By SNP:

Estimate genotype and copy number for each SNP.

**2** Within a sample:

Borrow strength between SNPs to infer regions of LOH and copy number changes.

**3** Between samples:

Comparison between normal and disease populations to find chromosomal alterations associated with disease.

$\longrightarrow$ The confidence in genotype calls can differ substantially between SNPs!

$\longrightarrow$ At each SNP, we observe a noisy measure of the true copy number and genotype (and possibly also measures of confidence in those estimates).

Novel (and we believe, important) HMM features:

1. Model the observation sequence of genotype calls and copy number jointly (Vanilla)

2. Integrate confidence estimates of the genotype calls and copy number estimates (ICE)

# The Vanilla HMM components

- Observations $\widehat{CN}$ and $\widehat{GT}$

- Hidden states

- Initial state probability distribution

- Transition probabilities

- Emission probabilities

## Transition probabilities

Following suggestions in the literature, we model the transition probabilities as a function of the distance $d$ between SNPs.

Specifically, let $\theta(d) \equiv 1 - e^{-2d}$ denote the probability that SNP $i$ is not informative ($I^c$) for SNP at $i + 1$.

For example:

$$
\begin{aligned}
\tau_{\ominus \mid \oslash}(d) &= P\left\{ \ominus_{i+1} \mid \oslash_i, d \right\} \\[1mm]
&= P\left\{ \ominus_{i+1}, I \mid \oslash_i, d \right\} + P\left\{ \ominus_{i+1}, I^c \mid \oslash_i, d \right\} \\[1mm]
&= P\left\{ \ominus_{i+1} \mid I, \oslash_i, d \right\} \times P\left\{ I \mid \oslash_i, d \right\} + \\
&\quad P\left\{ \ominus_{i+1} \mid I^c, \oslash_i, d \right\} \times P\left\{ I^c \mid \oslash_i, d \right\} \\[1mm]
&= P\left\{ \ominus \right\} \times \theta(d).
\end{aligned}
$$

We assume conditional independence between copy number estimates and the genotype calls.

For example:

$$
\begin{aligned}
f(\widehat{\mathrm{CN}}, \widehat{\mathrm{GT}} | \oslash) &= f(\widehat{\mathrm{CN}} | \oslash) \times f(\widehat{\mathrm{GT}} | \oslash) \\
&= f\left\{ \widehat{\mathrm{CN}} \mid \searrow \right\} \times f\left\{ \widehat{\mathrm{GT}} \mid \bigcirc \right\} \\
&= \beta_{\searrow}\left\{ \widehat{\mathrm{CN}} \right\} \times \beta_{\bigcirc}\left\{ \widehat{\mathrm{GT}} \right\}.
\end{aligned}
$$

Let $S_{\widehat{GT}}$ be the confidence score for the genotype estimate.

We can estimate from Hapmap the following densities:

$$f\left\{ S_{\widehat{HOM}} \mid \widehat{HOM}, HOM \right\}, f\left\{ S_{\widehat{HOM}} \mid \widehat{HOM}, HET \right\}, f\left\{ S_{\widehat{HET}} \mid \widehat{HET}, HOM \right\}, f\left\{ S_{\widehat{HET}} \mid \widehat{HET}, HET \right\}.$$

$\longrightarrow$ Note:

$$f\left\{ S_{\widehat{HOM}} \mid \widehat{HOM}, \circ \right\} \approx f\left\{ S_{\widehat{HOM}} \mid \widehat{HOM}, HOM \right\}$$

$$f\left\{ S_{\widehat{HET}} \mid \widehat{HET}, \circ \right\} \approx f\left\{ S_{\widehat{HET}} \mid \widehat{HET}, HOM \right\}.$$

Recall that

$$
\begin{aligned}
f(\widehat{CN}, \widehat{GT} | \oslash) &= f(\widehat{CN} | \oslash) \times f(\widehat{GT} | \oslash) \\
&= f\left\{ \widehat{CN} \mid \searrow \right\} \times f\left\{ \widehat{GT} \mid \bigcirc \right\} \\
&= \beta_{\searrow}\left\{ \widehat{CN} \right\} \times \beta_{\bigcirc}\left\{ \widehat{GT} \right\}.
\end{aligned}
$$

If the state for a particular SNP is *Loss*, we have

$$
\beta_{\bigcirc}\left\{ \widehat{GT}, S_{\widehat{GT}} \right\} = f\left\{ \widehat{GT} \mid \bigcirc \right\} \times f\left\{ S_{\widehat{GT}} \mid \widehat{GT}, \bigcirc \right\}.
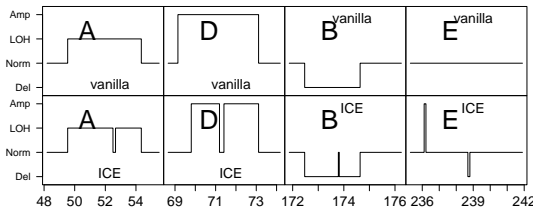$$

For retention, the true genotype can be HET or HOM:

$$\beta_{\bullet} \left\{ \widehat{GT}, s_{\widehat{GT}} \right\}$$

$$= f\left\{ \widehat{GT} \mid \bullet \right\} f\left\{ s_{\widehat{GT}} \mid \widehat{GT}, \bullet \right\}$$

$$= f\left\{ \widehat{GT} \mid \bullet \right\} \left( f\left\{ s_{\widehat{GT}}, HOM \mid \widehat{GT}, \bullet \right\} + f\left\{ s_{\widehat{GT}}, HET \mid \widehat{GT}, \bullet \right\} \right)$$

$$= f\left\{ \widehat{GT} \mid \bullet \right\} \left( f\left\{ s_{\widehat{GT}} \mid HOM, \widehat{GT}, \bullet \right\} f\left\{ HOM \mid \widehat{GT}, \bullet \right\} + f\left\{ s_{\widehat{GT}} \mid HET, \widehat{GT}, \bullet \right\} f\left\{ HET \mid \widehat{GT}, \bullet \right\} \right)$$

$$= f\left\{ \widehat{GT} \mid \bullet \right\} \left( f\left\{ s_{\widehat{GT}} \mid HOM, \widehat{GT} \right\} f\left\{ HOM \mid \widehat{GT}, \bullet \right\} + f\left\{ s_{\widehat{GT}} \mid HET, \widehat{GT} \right\} f\left\{ HET \mid \widehat{GT}, \bullet \right\} \right)$$

Bioconductor package: ICE

| Father | Mother | Child | iUPI-P (1) | hUPI-P (1) | BPI (2) | hUPI-M (3) | iUPI-M (3) | MI-S (4) | MI-D (5) |
|--------|--------|-------|------------|------------|---------|------------|------------|----------|----------|
| AA | AA | AA | X | X | X | X | X | | |
| | | AB | | | | | | X | |
| | | BB | | | | | | | X |
| | AB | AA | X | X | X | | X | | |
| | | AB | | | X | X | | | |
| | | BB | | | | | X | | |
| | BB | AA | X | X | | | | | |
| | | AB | | | X | | | | |
| | | BB | | | | X | X | | |
| AB | AA | AA | X | | X | X | X | | |
| | | AB | | X | X | | | | |
| | | BB | X | | | | | | |
| | AB | AA | X | | X | | X | | |
| | | AB | | X | X | X | | | |
| | | BB | X | | X | | X | | |
| | BB | AA | X | | | | | | |
| | | AB | | X | X | | | | |
| | | BB | | | X | X | X | | |
| BB | AA | AA | | | | X | X | | |
| | | AB | | | X | | | | |
| | | BB | X | X | | | | | |
| | AB | AA | | | | | X | | |
| | | AB | | | X | X | | | |
| | | BB | X | X | X | | X | | |
| | BB | AA | | | | | | | X |
| | | AB | | | | | | X | |
| | | BB | X | X | X | X | X | | |

$\longrightarrow$ Given the number of SNPs for a particular event region type ($M$), among all uploaded informative autosomal SNPs ($X$) of that trio, what is the probability to have the observed number ($N$) or more informative non-BPI SNPs of that type clustered together solely by chance?

In other words, if the $M$ SNPs of a particular event region type were randomly dispersed among the $X$ informative autosomal SNPs, how probable is it to observe an event of the same or larger magnitude (defined by the number of consecutive SNPs of that event region)?

- Let $p$ be the probability of a positive outcome in a Bernoulli trial (e.g. a 1 for 0/1 outcomes).

- Assume that we have $Z$ trials.

- Let $LS$ be the length of the largest block of consecutive ones in those $Z$ trials.

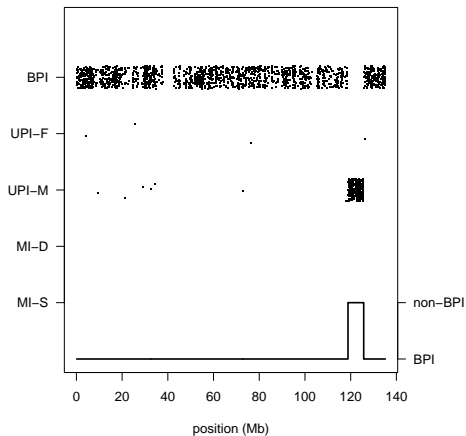- Let $P_Z(LS < N)$ be the probability that $LS$ is smaller than some number $N$.

- If $Z < N$, then $P_Z(LS < N) = 1$.

- If $Z = N$, then $P_Z(LS < N) = 1 - p^N$.

- If $Z > N$, then

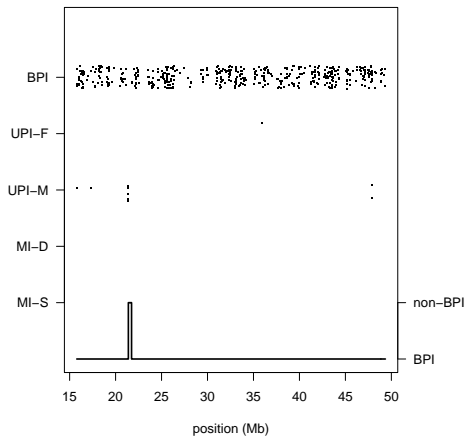$$P_Z(LS < N) = \sum_{k=0}^{N-1} p^k (1-p) P_{Z-k-1}(LS < N)$$

- Tabulate the probabilities $P_1(LS < N), \ldots, P_N(LS < N)$.

- Calculate $P_{N+1}(LS < N), P_{N+2}(LS < N), \ldots, P_X(LS < N)$, iteratively.

- Using $p = M/X$, the probability of having $N$ or more informative SNPs of a particular non-BPI type clustered together solely by chance is equal to $1 - P_X(LS < N)$.

# Acknowledgments

- Rob Scharpf

- Giovanni Parmigiani

- Rafael Irizarry & Gang

  Benilton Carvalho, Wenyi Wang

- Jonathan Pevsner & Lab

  Nate Miller, Eli Roberson, Jason Ting

📄 Carvalho B, Bengtsson H, Speed TP, Irizarry RA (2007)
Exploration, normalization, and genotype calls of high-density oligonucleotide
SNP array data. *Biostatistics*, 8(2):485-99.

📄 Scharpf RB, Ting JC, Pevsner J, Ruczinski I (2007).
SNPchip: R classes and methods for SNP array data. *Bioinformatics*, 23(5):
627-8.

📄 Scharpf RB, Parmigiani G, Ruczinski I (2007).
A hidden markov model for joint estimation of genotype and copy number in
high-throughput SNP chips. *JHU Biostatistics Working papers*, #136.

📄 Ting JC, Ye Y, Thomas GH, Ruczinski I, Pevsner J (2006).
Analysis and visualization of chromosomal abnormalities in SNP data with
SNPscan. *BMC Bioinformatics*, 7(1):25.

📄 Ting JC, Roberson ED, Miller N, et al, Ruczinski I, Thomas GH, Pevsner J (2007).
Visualization of uniparental inheritance, Mendelian inconsistencies, deletions and
parent of origin effects in single nucleotide polymorphism trio data with SNPtrio.
*Human Mutation*, (in press).

📄 Wang W, Caravalho B, Miller N, Pevsner J, Chakravarti A, Irizarry RA (2006)
Estimating genome-wide copy number using allele specific mixture models. *JHU
Biostatistics Working papers*, #122.