

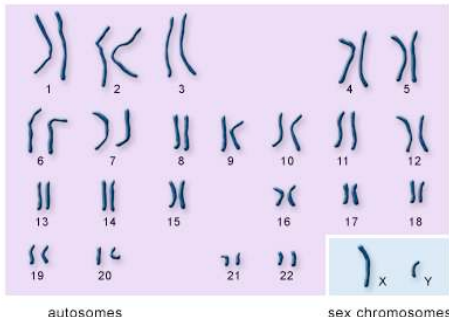
Inference for SNPchip Data in the Presence of Genotype and Copy Number Uncertainty

Ingo Ruczinski

Department of Biostatistics
Johns Hopkins Bloomberg School of Public Health

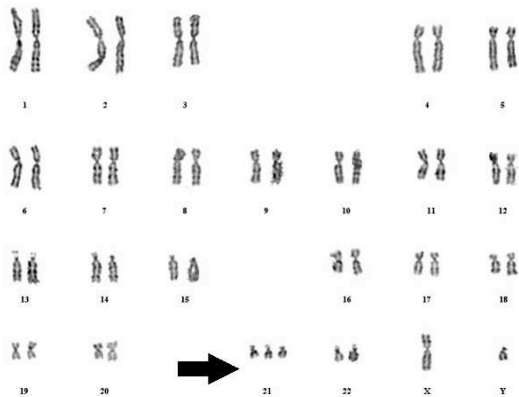
July 27, 2009

Karyotypes

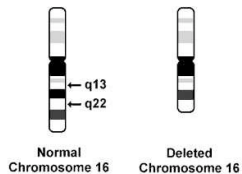
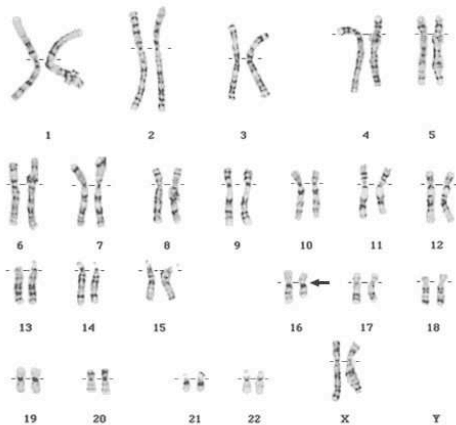


U.S. National Library of Medicine

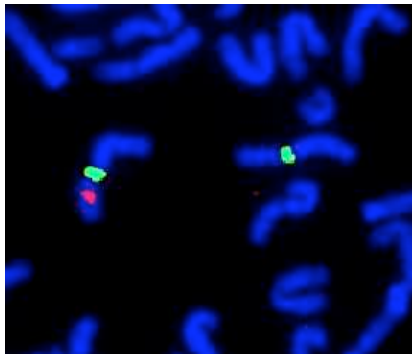
Trisomy



Karyotypes

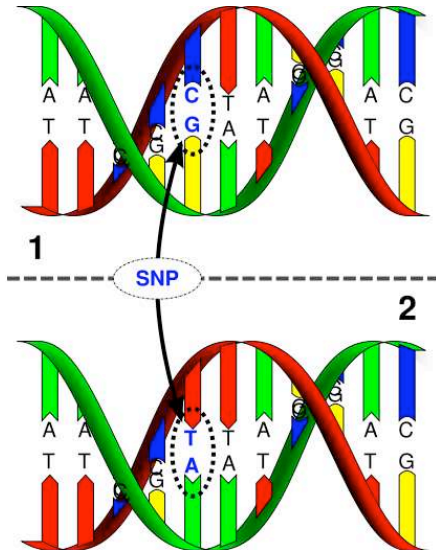


General Cytogenetics Information <http://members.aol.com/chrominfo/>

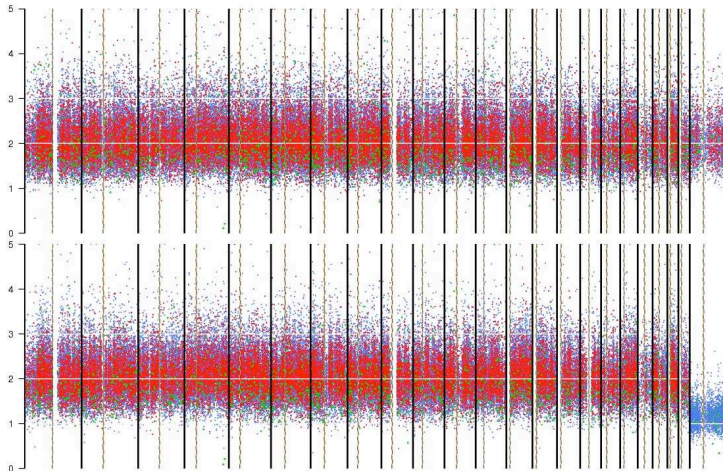


Courtesy of the Pevsner Laboratory

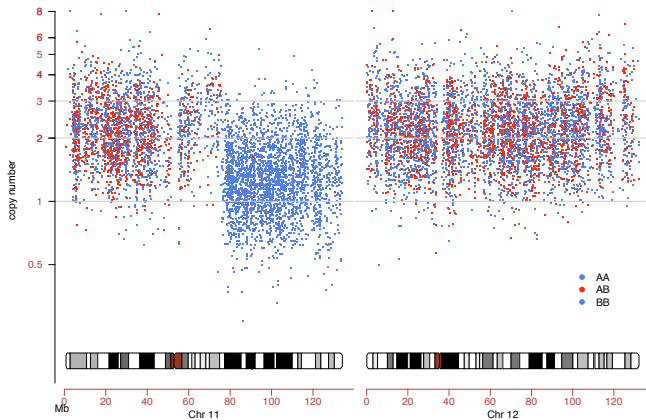
Single nucleotide polymorphisms



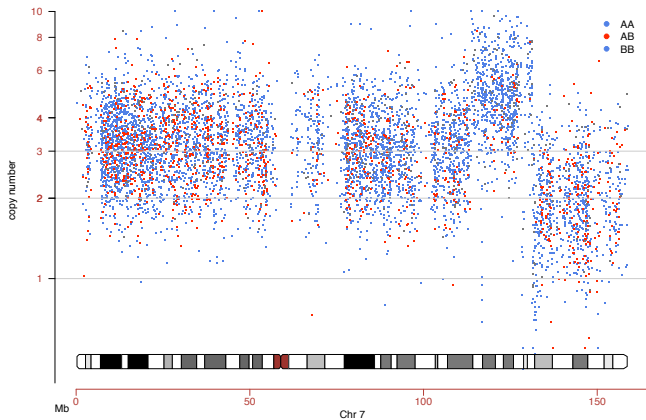
SNP chip data



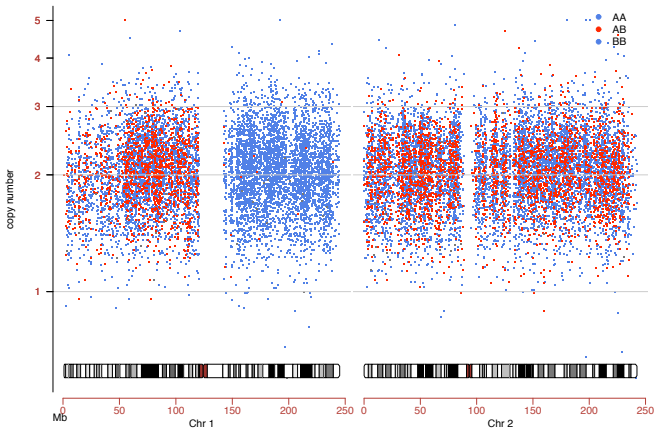
Deletion



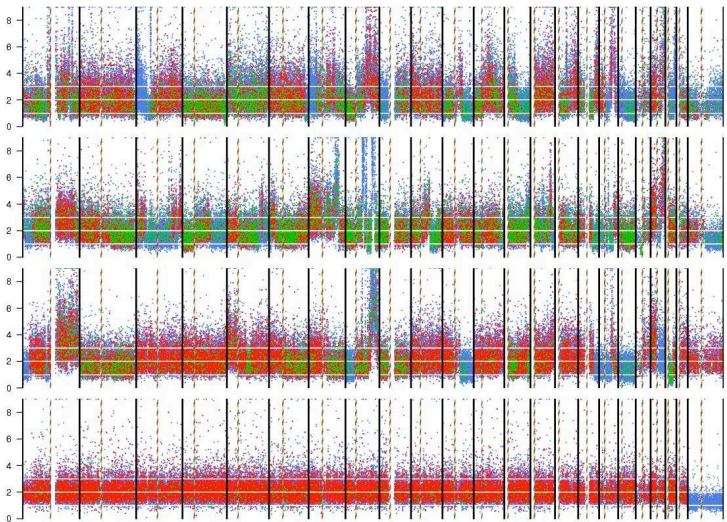
Amplification



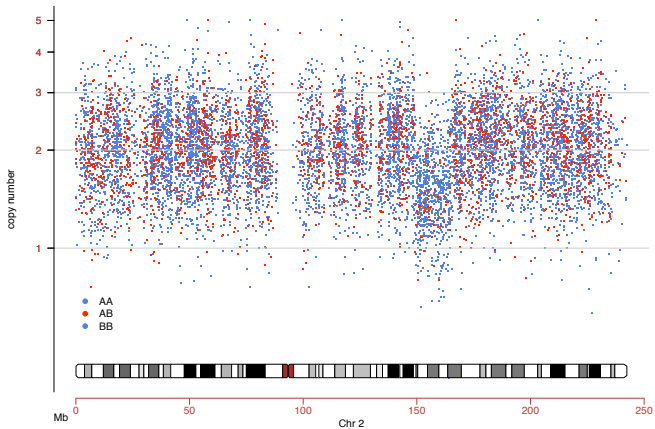
Uniparental Isodisomy



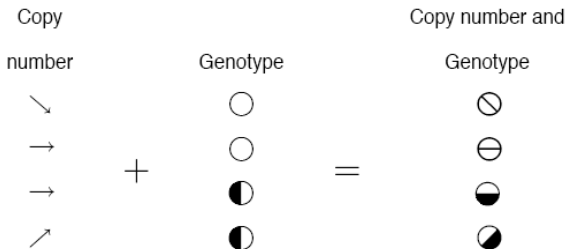
Cancer samples



Mosaicism

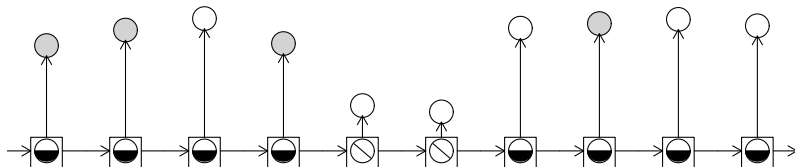


Hidden states



The structure of the data we observe

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



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)

QuantiSNP and *PennCNV* also model genotype and copy number jointly!

Colella et. al. (2007), *QuantiSNP: an objective Bayes hidden-Markov model...*, Nucleic Acids Res 35(6): 2013-25.

Wang et. al. (2008), *PennCNV: An integrated hidden-Markov model designed for...*, Genome Research 17: 1665-74.

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:

$$\begin{aligned}\tau_{\ominus|\ominus}(d) &= P\{\ominus_{i+1} | \ominus_i, d\} \\ &= P\{\ominus_{i+1}, I | \ominus_i, d\} + P\{\ominus_{i+1}, I^c | \ominus_i, d\} \\ &= P\{\ominus_{i+1} | I, \ominus_i, d\} \times P\{I | \ominus_i, d\} + \\ &\quad P\{\ominus_{i+1} | I^c, \ominus_i, d\} \times P\{I^c | \ominus_i, d\} \\ &= P\{\ominus\} \times \theta(d).\end{aligned}$$

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

For example:

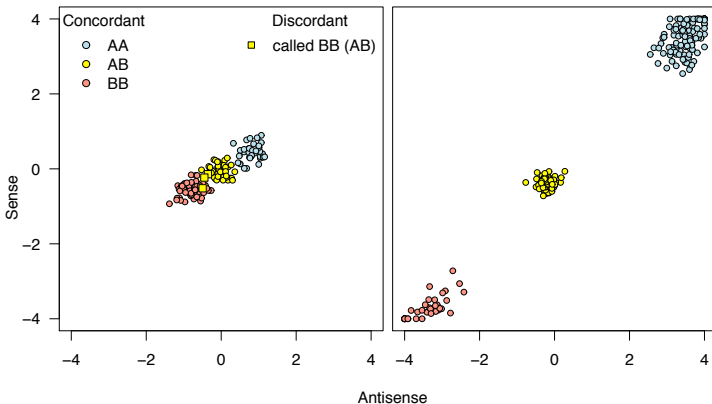
$$\begin{aligned}f(\widehat{\text{CN}}, \widehat{\text{GT}} | \circ) &= f(\widehat{\text{CN}} | \circ) \times f(\widehat{\text{GT}} | \circ) \\&= f\{\widehat{\text{CN}} | \searrow\} \times f\{\widehat{\text{GT}} | \circ\} \\&= \beta_{\searrow}\{\widehat{\text{CN}}\} \times \beta_{\circ}\{\widehat{\text{GT}}\}.\end{aligned}$$

Some geezer down the hallway...



More information

- The confidence in genotype calls can differ substantially between SNPs!



Integrating confidence estimates for genotype calls

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

We can estimate from Hapmap the following densities:

$$f\{S_{\widehat{HOM}} | \widehat{HOM}, \text{HOM}\}, f\{S_{\widehat{HOM}} | \widehat{HOM}, \text{HET}\}, f\{S_{\widehat{HET}} | \widehat{HET}, \text{HOM}\}, f\{S_{\widehat{HET}} | \widehat{HET}, \text{HET}\}.$$

→ Note:

$$f\{S_{\widehat{HOM}} | \widehat{HOM}, \circ\} \approx f\{S_{\widehat{HOM}} | \widehat{HOM}, \text{HOM}\}$$

$$f\{S_{\widehat{HET}} | \widehat{HET}, \circ\} \approx f\{S_{\widehat{HET}} | \widehat{HET}, \text{HOM}\}.$$

Recall that

$$\begin{aligned}f(\widehat{\text{CN}}, \widehat{\text{GT}} | \circ) &= f(\widehat{\text{CN}} | \circ) \times f(\widehat{\text{GT}} | \circ) \\&= f\{\widehat{\text{CN}} | \searrow\} \times f\{\widehat{\text{GT}} | \circ\} \\&= \beta_{\searrow}\{\widehat{\text{CN}}\} \times \beta_{\circ}\{\widehat{\text{GT}}\}.\end{aligned}$$

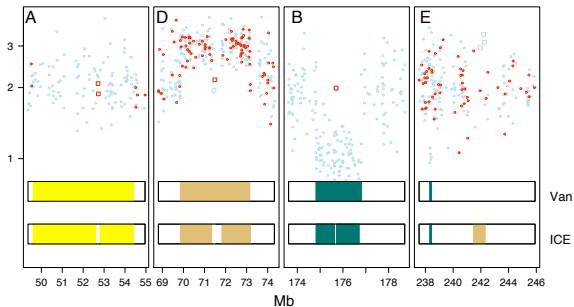
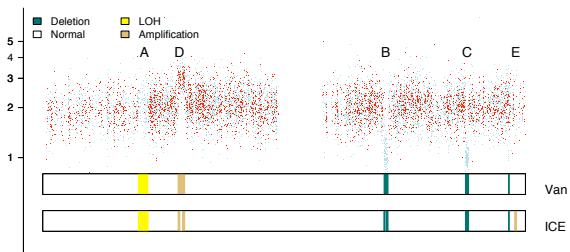
If the state for a particular SNP is *Loss*, we have

$$\beta_{\circ}\{\widehat{\text{GT}}, \mathbf{s}_{\widehat{\text{GT}}}\} = f\{\widehat{\text{GT}} | \circ\} \times f\{\mathbf{s}_{\widehat{\text{GT}}} | \widehat{\text{GT}}, \circ\}.$$

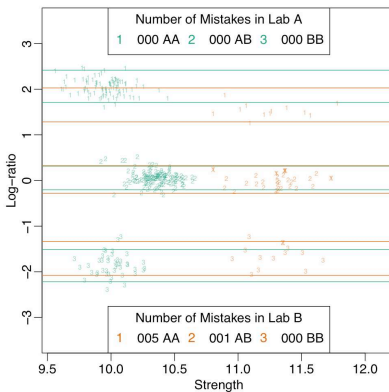
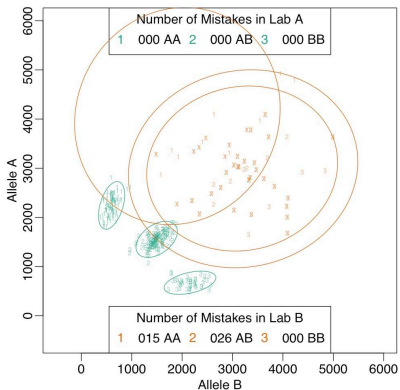
For retention, the true genotype can be HET or HOM:

$$\begin{aligned} & \beta_{\bullet} \{ \widehat{GT}, S_{\widehat{GT}} \} \\ = & f \{ \widehat{GT} | \bullet \} f \{ S_{\widehat{GT}} | \widehat{GT}, \bullet \} \\ = & f \{ \widehat{GT} | \bullet \} (f \{ S_{\widehat{GT}, \text{HOM}} | \widehat{GT}, \bullet \} + f \{ S_{\widehat{GT}, \text{HET}} | \widehat{GT}, \bullet \}) \\ = & f \{ \widehat{GT} | \bullet \} (f \{ S_{\widehat{GT}} | \text{HOM}, \widehat{GT}, \bullet \} f \{ \text{HOM} | \widehat{GT}, \bullet \} + f \{ S_{\widehat{GT}} | \text{HET}, \widehat{GT}, \bullet \} f \{ \text{HET} | \widehat{GT}, \bullet \}) \\ = & f \{ \widehat{GT} | \bullet \} (f \{ S_{\widehat{GT}} | \text{HOM}, \widehat{GT} \} f \{ \text{HOM} | \widehat{GT}, \bullet \} + f \{ S_{\widehat{GT}} | \text{HET}, \widehat{GT} \} f \{ \text{HET} | \widehat{GT}, \bullet \}) \end{aligned}$$

Vanilla and ICE HMMs for genotype and copy number

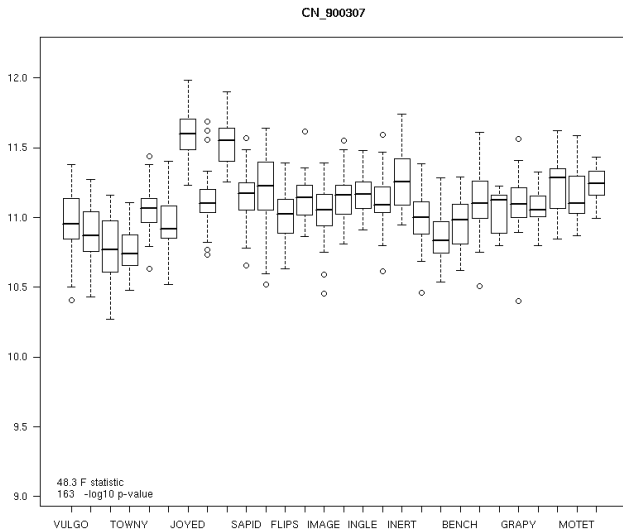


Genotypes and copy numbers



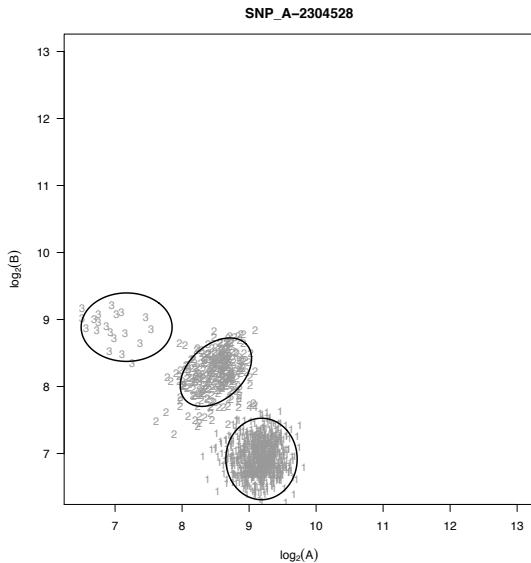
From Benilton Carvalho and Rafa Irizarry

Plate effects

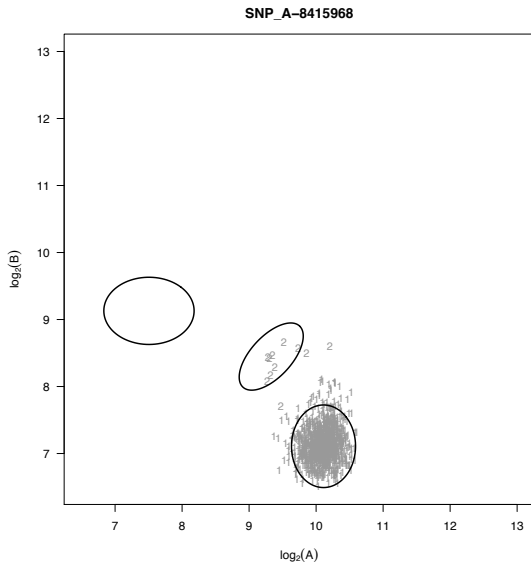


Bipolar GWAS (EA controls) from dbGap

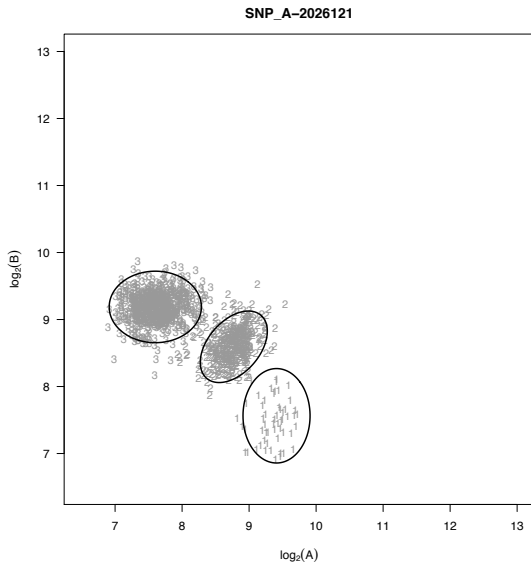
A versus B plots



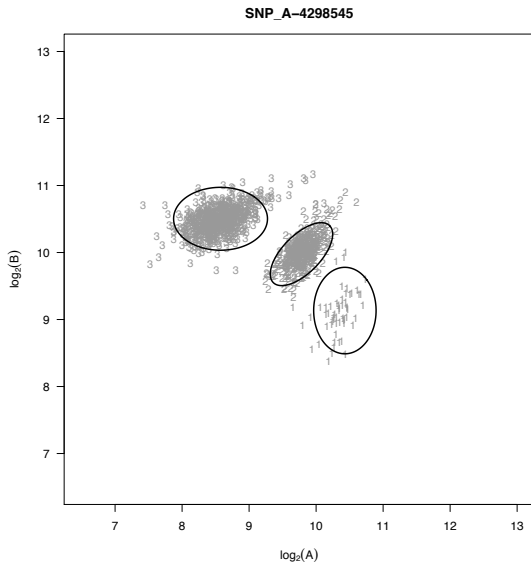
A versus B plots



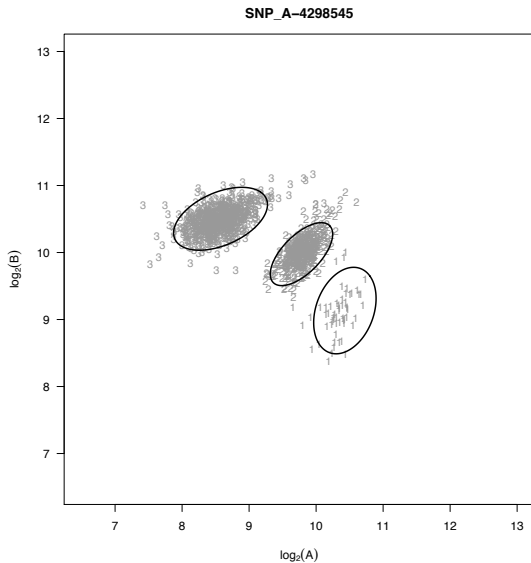
A versus B plots



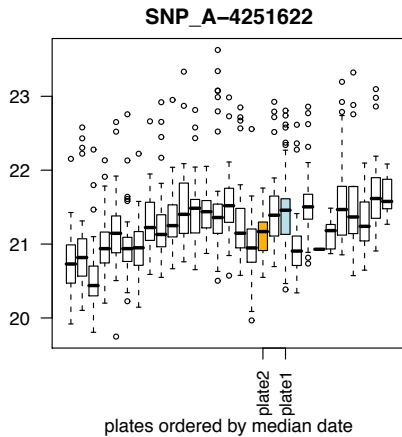
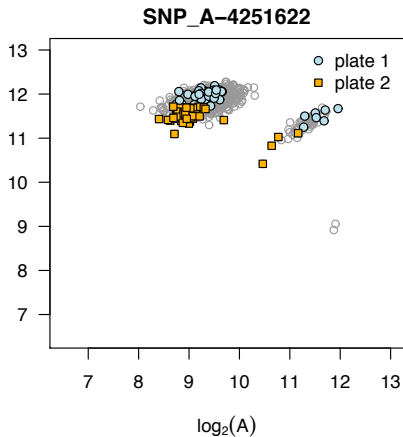
A versus B plots



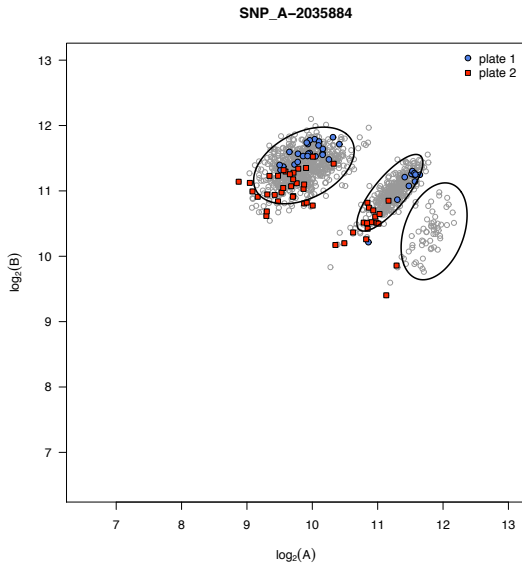
A versus B plots



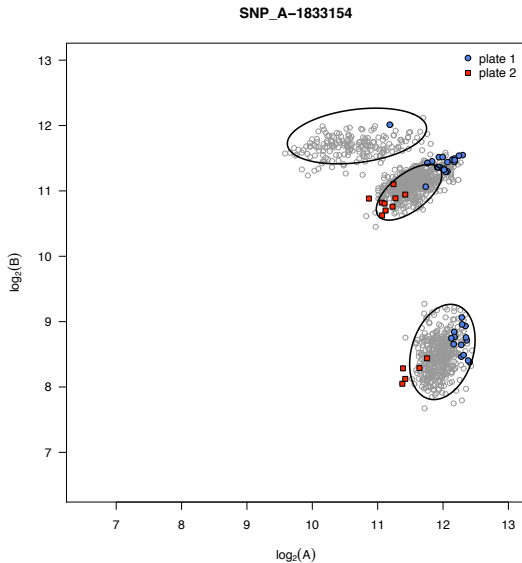
A versus B plots



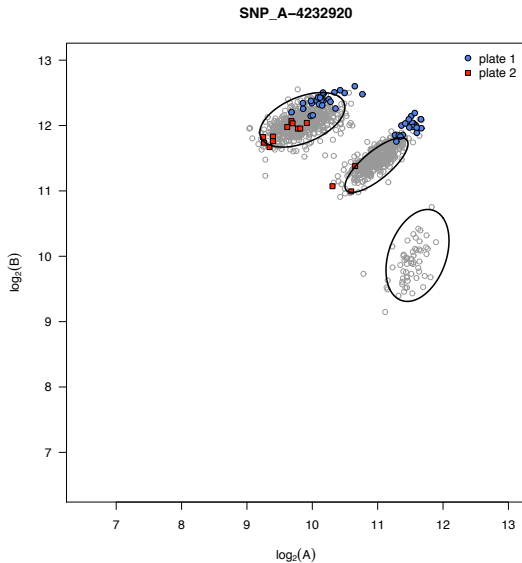
A versus B plots



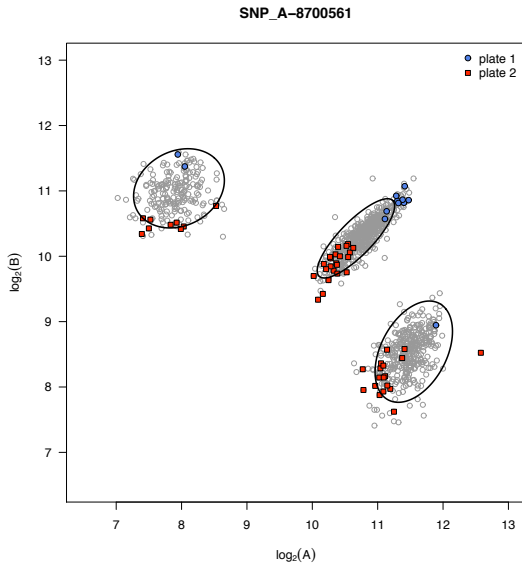
A versus B plots



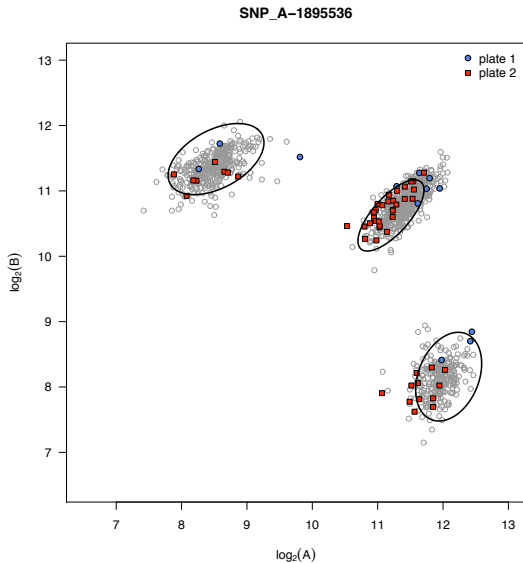
A versus B plots



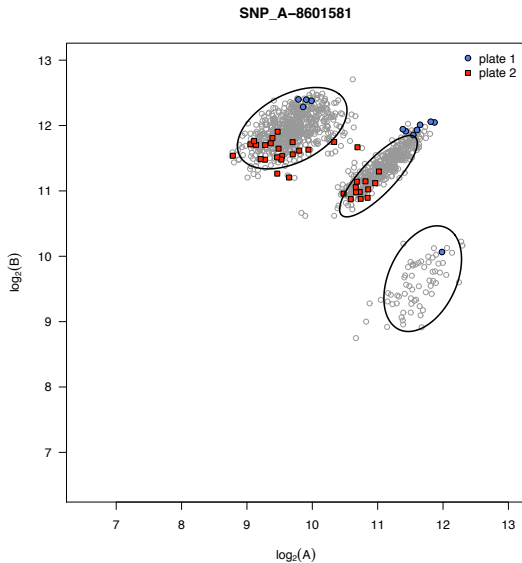
A versus B plots



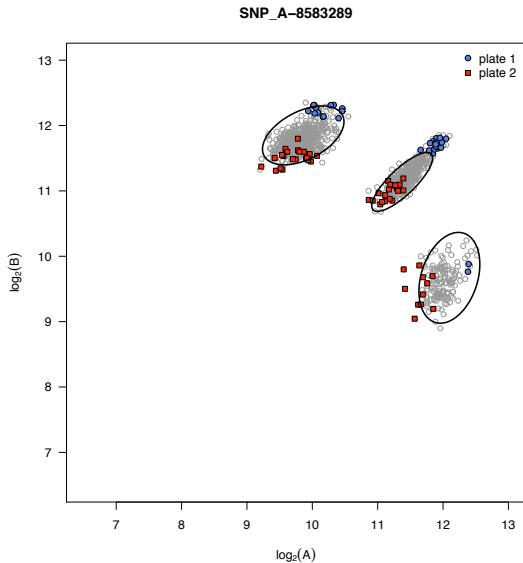
A versus B plots



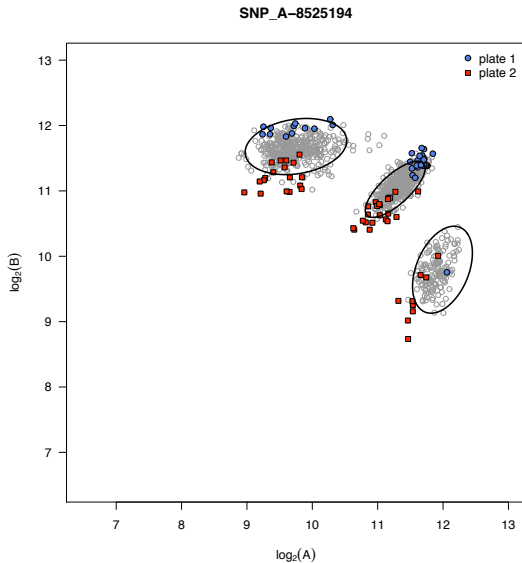
A versus B plots



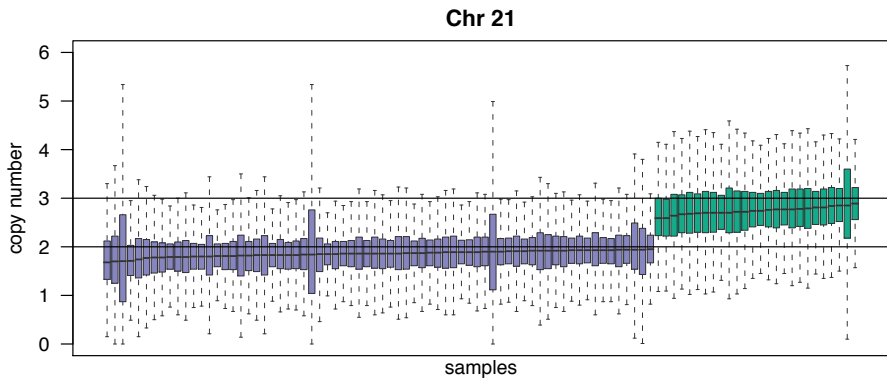
A versus B plots



A versus B plots

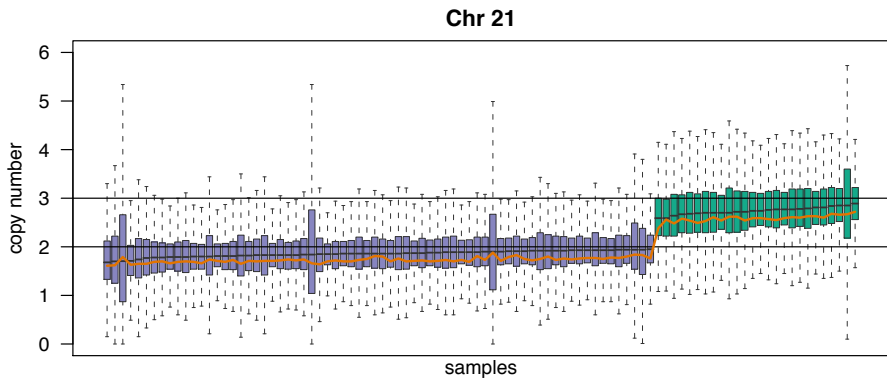


Trisomy 21



Samples from Aravinda Chakravarti and Betty Doan

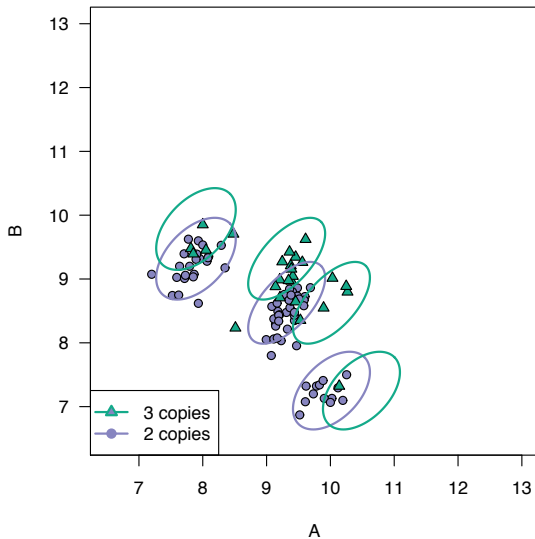
Trisomy 21



Samples from Aravinda Chakravarti and Betty Doan

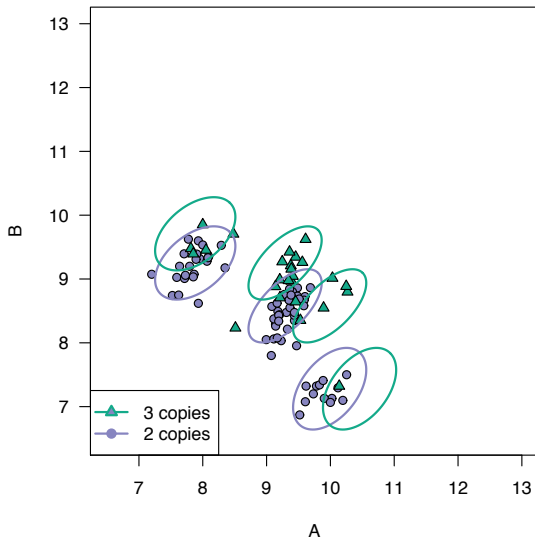
A versus B plots

SNP_A-8348190

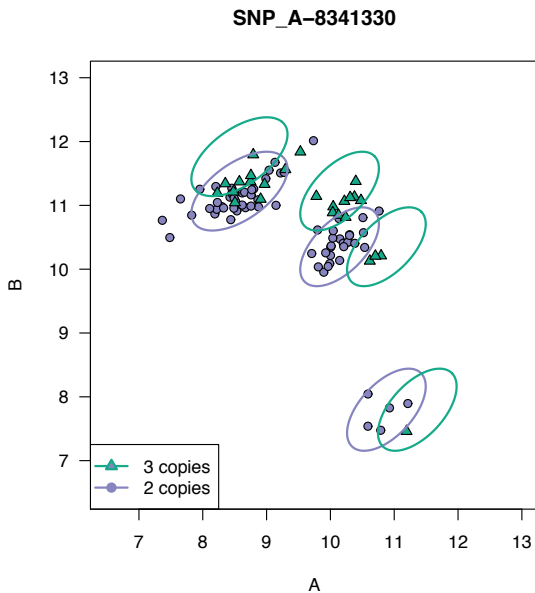


A versus B plots

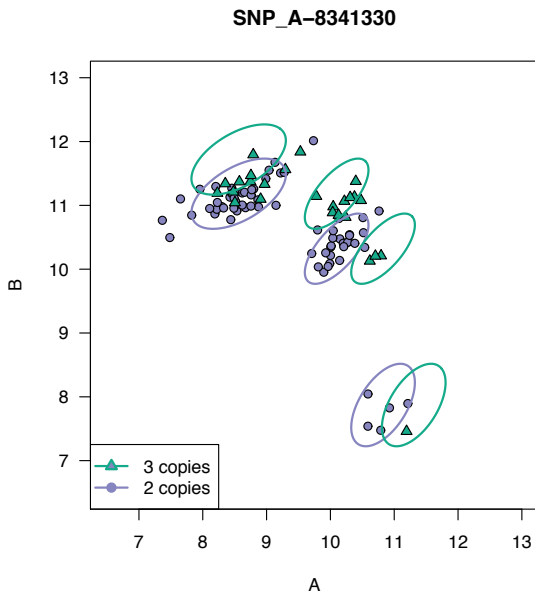
SNP_A-8348190



A versus B plots

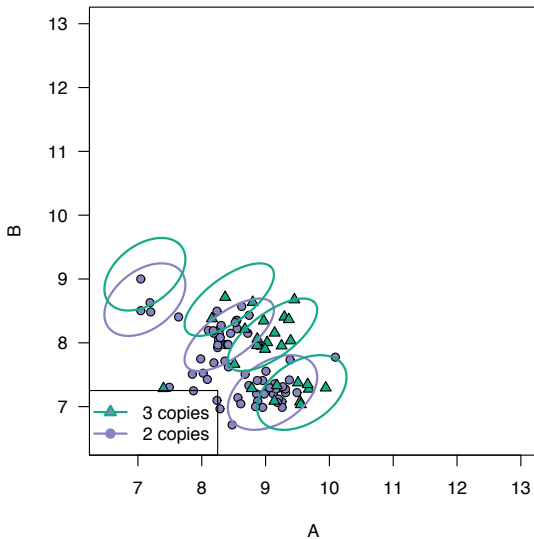


A versus B plots



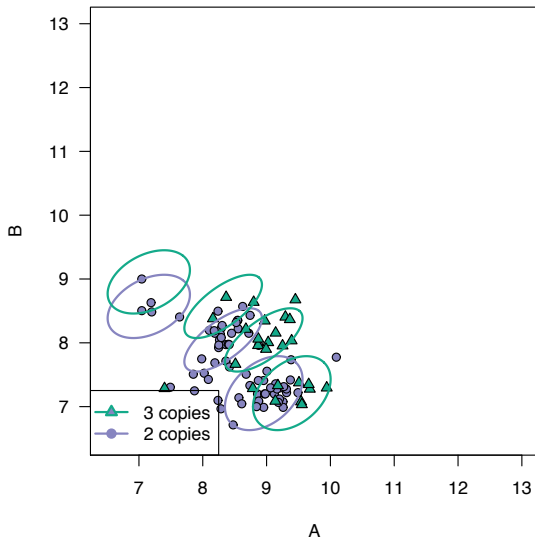
A versus B plots

SNP_A-8339372



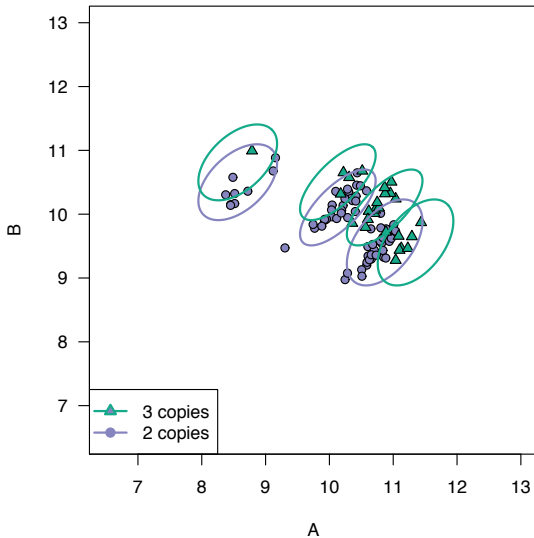
A versus B plots

SNP_A-8339372

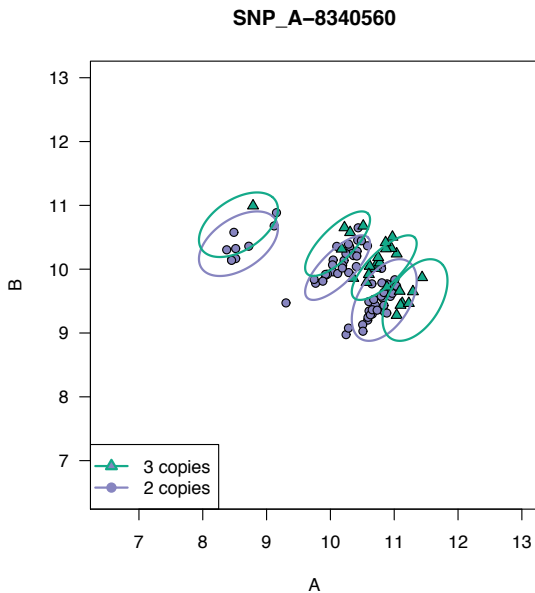


A versus B plots

SNP_A-8340560

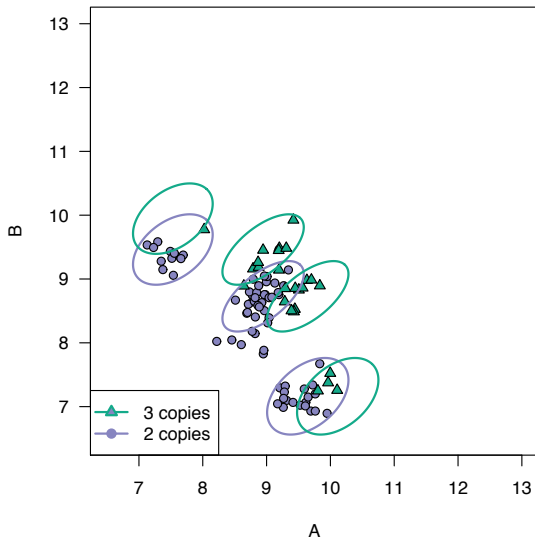


A versus B plots



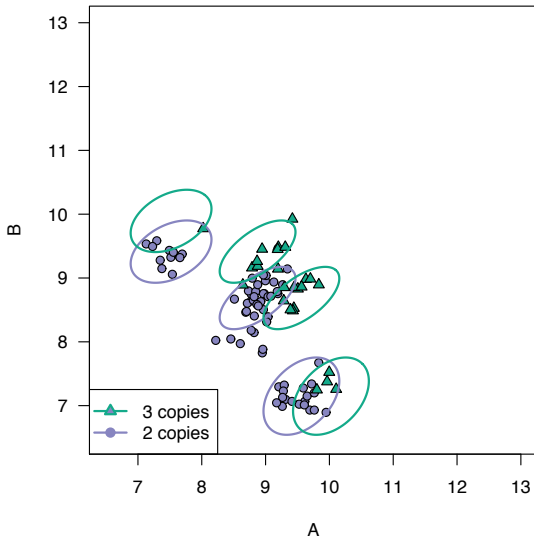
A versus B plots

SNP_A-1969323

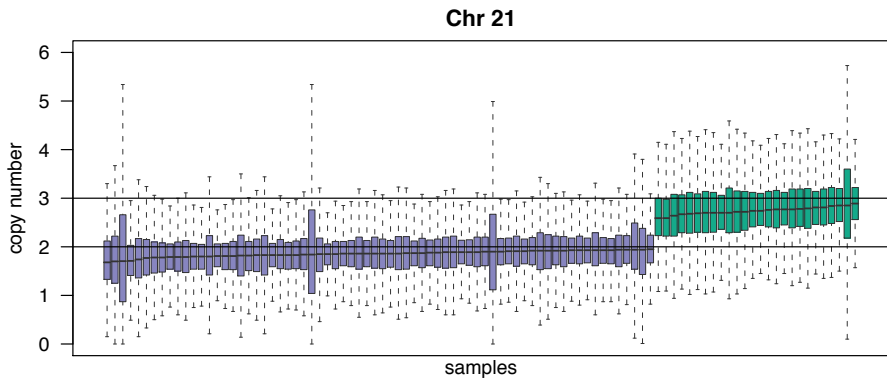


A versus B plots

SNP_A-1969323

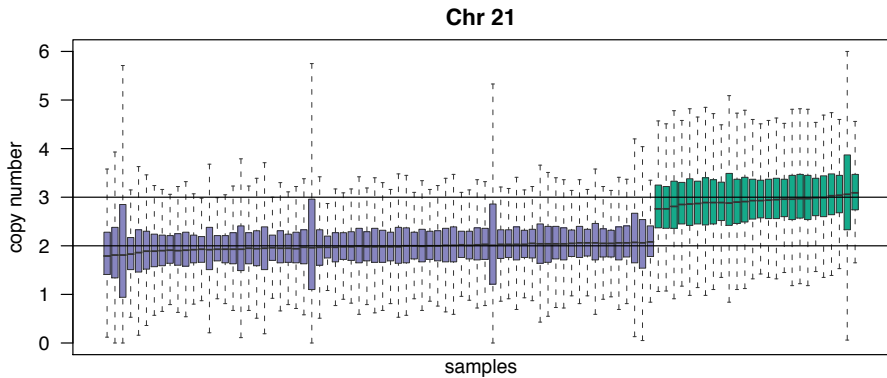


Trisomy 21



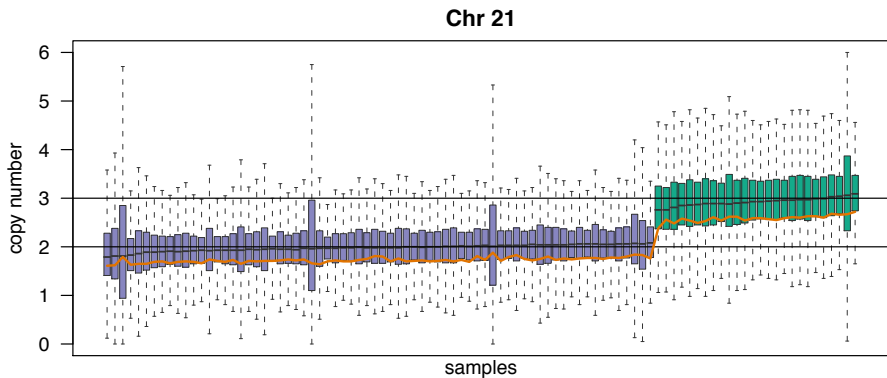
Samples from Aravinda Chakravarti and Betty Doan

Trisomy 21



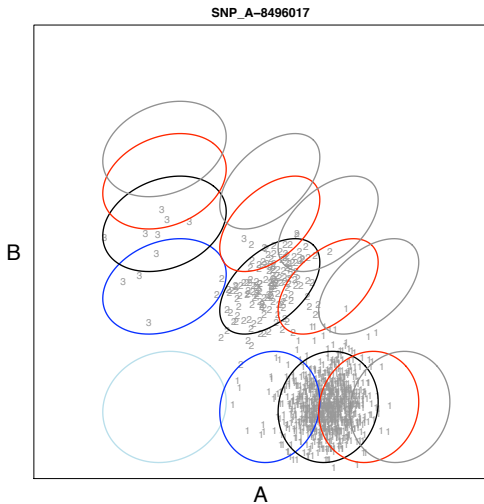
Samples from Aravinda Chakravarti and Betty Doan

Trisomy 21

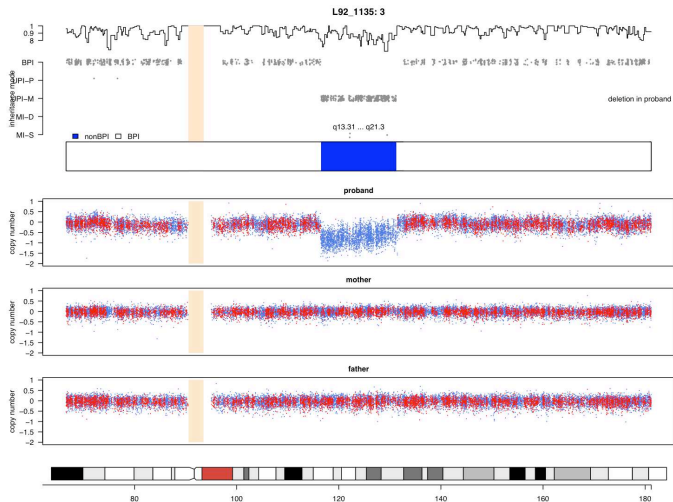


Samples from Aravinda Chakravarti and Betty Doan

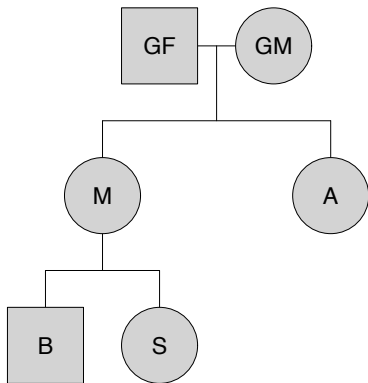
Prediction regions for copy number



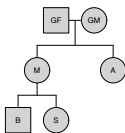
De novo deletion



Homozygous and hemizygous deletions

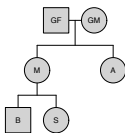


Homozygous and hemizygous deletions



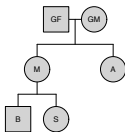
Grandfather		Grandmother		Mother		Aunt		Brother		Sister	
AB	0.01	BB	0.05	AB	-0.11	AB	-0.32	AB	0.06	BB	-0.02
AA	0.27	NC	-5.52	AA	-0.48	AA	-0.45	AB	0.12	BB	-0.42
AB	0.15	NC	-5.04	AA	-0.20	AA	-0.24	AB	-0.09	BB	-0.49
AB	-0.03	NC	-4.59	AA	-0.40	AA	-0.24	AA	0.30	AA	-0.72
BB	0.20	NC	-2.46	NC	-0.38	BB	-0.28	BB	0.22	BB	-0.45
AB	0.03	NC	-6.14	BB	-0.28	BB	-0.42	AB	0.09	AA	-0.70
AB	-0.05	NC	-5.02	BB	-0.17	BB	-0.34	AB	-0.22	AA	-1.06
BB	0.01	NC	-4.04	BB	0.04	BB	-0.68	BB	0.14	NC	-0.98
AB	0.17	NC	-4.06	AA	-0.27	AA	-0.33	AA	-0.03	AA	-0.76
AB	0.01	NC	-4.70	AA	-0.67	AA	-0.52	AA	0.16	AA	-0.80
AB	-0.10	NC	-4.42	BB	-0.25	BB	-0.62	AB	0.13	AA	-0.58
AB	0.01	NC	-8.29	BB	-0.17	BB	-0.15	AB	-0.15	AA	-0.29
BB	0.16	NC	-5.73	BB	-0.64	BB	-0.46	BB	0.10	BB	-0.52
AB	0.06	NC	-7.48	AA	-0.23	AA	-0.33	AB	0.07	BB	-0.47
AA	0.17	NC	-3.70	AA	-0.50	AA	-0.52	AB	-0.06	BB	-0.48
BB	0.02	NC	-5.00	BB	-0.34	BB	-0.45	AB	0.13	AA	-0.55
AA	0.21	NC	-6.10	AA	-0.43	AA	-0.40	AB	0.20	BB	-0.40
BB	0.05	BB	0.11	BB	0.15	BB	0.29	AB	0.13	AB	-0.01

Homozygous and hemizygous deletions



Grandfather		Grandmother		Mother		Aunt		Brother		Sister	
AB	0.01	BB	0.05	AB	-0.11	AB	-0.32	AB	0.06	BB	-0.02
AA	0.27	NC	-5.52	AA	-0.48	AA	-0.45	AB	0.12	BB	-0.42
AB	0.15	NC	-5.04	AA	-0.20	AA	-0.24	AB	-0.09	BB	-0.49
AB	-0.03	NC	-4.59	AA	-0.40	AA	-0.24	AA	0.30	AA	-0.72
BB	0.20	NC	-2.46	NC	-0.38	BB	-0.28	BB	0.22	BB	-0.45
AB	0.03	NC	-6.14	BB	-0.28	BB	-0.42	AB	0.09	AA	-0.70
AB	-0.05	NC	-5.02	BB	-0.17	BB	-0.34	AB	-0.22	AA	-1.06
BB	0.01	NC	-4.04	BB	0.04	BB	-0.68	BB	0.14	NC	-0.98
AB	0.17	NC	-4.06	AA	-0.27	AA	-0.33	AA	-0.03	AA	-0.76
AB	0.01	NC	-4.70	AA	-0.67	AA	-0.52	AA	0.16	AA	-0.80
AB	-0.10	NC	-4.42	BB	-0.25	BB	-0.62	AB	0.13	AA	-0.58
AB	0.01	NC	-8.29	BB	-0.17	BB	-0.15	AB	-0.15	AA	-0.29
BB	0.16	NC	-5.73	BB	-0.64	BB	-0.46	BB	0.10	BB	-0.52
AB	0.06	NC	-7.48	AA	-0.23	AA	-0.33	AB	0.07	BB	-0.47
AA	0.17	NC	-3.70	AA	-0.50	AA	-0.52	AB	-0.06	BB	-0.48
BB	0.02	NC	-5.00	BB	-0.34	BB	-0.45	AB	0.13	AA	-0.55
AA	0.21	NC	-6.10	AA	-0.43	AA	-0.40	AB	0.20	BB	-0.40
BB	0.05	BB	0.11	BB	0.15	BB	0.29	AB	0.13	AB	-0.01

Homozygous and hemizygous deletions



Grandfather		Grandmother		Mother		Aunt		Brother		Sister	
AB	0.01	BB	0.05	AB	-0.11	AB	-0.32	AB	0.06	BB	-0.02
AA	0.27	NC	-5.52	AA	-0.48	AA	-0.45	AB	0.12	BB	-0.42
AB	0.15	NC	-5.04	AA	-0.20	AA	-0.24	AB	-0.09	BB	-0.49
AB	-0.03	NC	-4.59	AA	-0.40	AA	-0.24	AA	0.30	AA	-0.72
BB	0.20	NC	-2.46	NC	-0.38	BB	-0.28	BB	0.22	BB	-0.45
AB	0.03	NC	-6.14	BB	-0.28	BB	-0.42	AB	0.09	AA	-0.70
AB	-0.05	NC	-5.02	BB	-0.17	BB	-0.34	AB	-0.22	AA	-1.06
BB	0.01	NC	-4.04	BB	0.04	BB	-0.68	BB	0.14	NC	-0.98
AB	0.17	NC	-4.06	AA	-0.27	AA	-0.33	AA	-0.03	AA	-0.76
AB	0.01	NC	-4.70	AA	-0.67	AA	-0.52	AA	0.16	AA	-0.80
AB	-0.10	NC	-4.42	BB	-0.25	BB	-0.62	AB	0.13	AA	-0.58
AB	0.01	NC	-8.29	BB	-0.17	BB	-0.15	AB	-0.15	AA	-0.29
BB	0.16	NC	-5.73	BB	-0.64	BB	-0.46	BB	0.10	BB	-0.52
AB	0.06	NC	-7.48	AA	-0.23	AA	-0.33	AB	0.07	BB	-0.47
AA	0.17	NC	-3.70	AA	-0.50	AA	-0.52	AB	-0.06	BB	-0.48
BB	0.02	NC	-5.00	BB	-0.34	BB	-0.45	AB	0.13	AA	-0.55
AA	0.21	NC	-6.10	AA	-0.43	AA	-0.40	AB	0.20	BB	-0.40
BB	0.05	BB	0.11	BB	0.15	BB	0.29	AB	0.13	AB	-0.01



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, Pevsner J, Ruczinski I (2008).
Hidden Markov models for the assessment of chromosomal alterations using high-throughput SNP arrays.
Annals of Applied Statistics, 2(2): 687-713.



Scharpf RB, Ruczinski I, Carvalho B, Doan B, Chakravarti A, Irizarry RA (2009).
A multilevel model to assess batch effects in copy number estimation using SNP arrays.
JHU Department of Biostatistics Working Papers, Working Paper 192.

→ Hopefully soon: more on plate effects and the trio HMM method.

<http://biostat.jhsph.edu/~iruczins/>