# Controlling the false discovery rate at detection of biological aberrations in -omics data

Djork-Arné Clevert, Andreas Mayr, Andreas Mitterecker, Günter Klambauer, and Sepp Hochreiter

Institute of Bioinformatics, JKU Linz

July 15th, 2011

# Sources of -omics data

- Measurement techniques
  - Next generations sequencing
  - Mass Spec
  - Microarrays

- Application
  - mRNA and miRNA (gene expression profiling/transcriptomics)
  - Copy numbers (structural variant determination/genomics)
  - SNPs (genotyping/genomics)
  - Proteins and Metabolites (proteomics/metabolomics)

- Aberrations
  - Differentially expressed
  - Loss and gain of DNA segments, loss of heterozygosity
  - SNP frequencies
  - Different concentration

# Characteristics of -omics data

- High-dimensional
  - # Genes, #miRNAs
  - # Loci
  - # SNPs
  - # Proteins

- Noisy
  - Measurement noise
  - Cross-hybridization, GC-content bias

Many falsely discovered aberrations or false positives (high FDR)

## Problems caused by false discoveries

**BIOINF**

- FPs are not associated with an experimental condition
  while correction for multiple testing must take them into account
  - Decreases the study's discovery power
  - Decreases the significance of discoveries
- FPs misguide researchers

Demand for **methods with a low false discovery** rate at detection of biological aberrations in -omics data

# Latent variable models

- Decompose observation into noise and signal by a generative model
  - Remove noise
    $\Rightarrow$ aberration detection in noise-free data
  - Decrease dimensionality
    $\Rightarrow$ signal variance for filtering (Informative/ Non-Informative calls)
- Model across samples for each gene, locus, SNP, or protein
- Use latent variable to represent the gene expression level, DNA copy number, genotype, or protein concentration

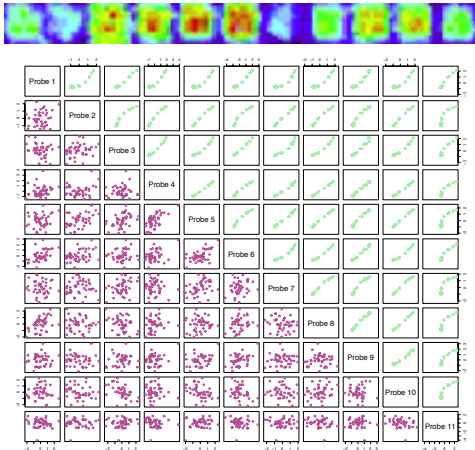# Our latent variable models

Next generations sequencing:

- cn.MOPS (Copy Number estimation by a Mixture Of PoissonS)
  - $\rightarrow$ DNA copy numbers

Microarray:

- FARMS (Factor Analysis for Robust Microarray Summarization)
  - $\rightarrow$ gene expression and miRNA
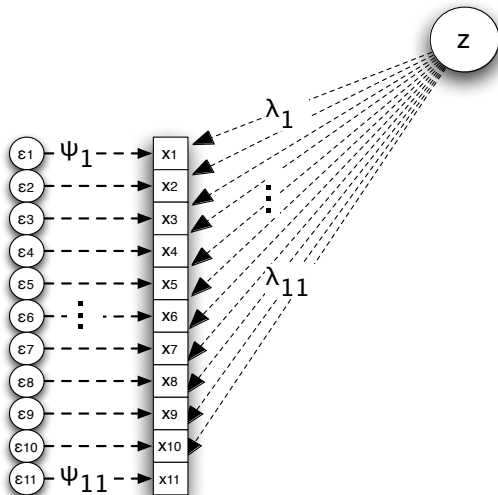- cn.FARMS ( Copy Number estimation by FARMS)
  - $\rightarrow$ DNA copy numbers

# FARMS: Facts and assumptions

- Gene measured by different probes
- Goal: summarize probe intensities to an expression value
- Noise-free probes are positively correlated
  - Variable probe qualities
  - High quality probes are linear dependent
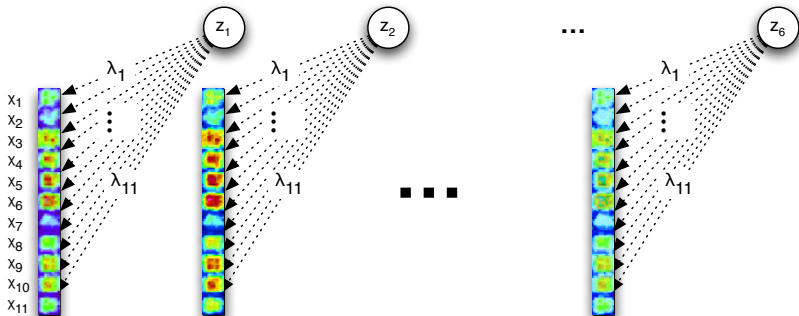- Replicate probe intensities are Gaussian distributed



Higher mRNA concentration → larger intensities

# FARMS: The idea

# FARMS: The data

## FARMS: The model

BIOINF

$$x = \boldsymbol{\lambda} \, z + \boldsymbol{\epsilon}$$

- $x, \boldsymbol{\lambda} \in \mathbb{R}^n$ and $z \sim \mathcal{N}(0, 1)$, $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \boldsymbol{\Psi})$.
- Probe intensities: $\{x\} = \{x_1, \ldots, x_N\}$ (log-transformed, standardized)
- Hidden factor $z$ represents the gene expression level
- $\boldsymbol{\epsilon}$ accounts for the independent noise in the probes intensities
- $\boldsymbol{\epsilon}$ and $z$ are independent

Model selection: maximize the likelihood $\left(x \sim \mathcal{N}\left(\boldsymbol{\lambda}\boldsymbol{\lambda}^T + \boldsymbol{\Psi}\right)\right)$ with respect to $\boldsymbol{\Psi}$ and $\boldsymbol{\lambda}$ by an EM algorithm

# FARMS: Bayes framework

**BIOINF**

## Posterior

$$p(\boldsymbol{\lambda}, \boldsymbol{\Psi} \,|\, \{x\}) \;\propto\; p(\{x\} \,|\, \boldsymbol{\lambda}, \boldsymbol{\Psi}) \; p(\boldsymbol{\lambda})$$

### Prior knowledge

- Positive $\boldsymbol{\lambda}$ ensure positive probe correlation

- Most genes show no or small signal (large signals are of interest in a study)

### Rectified Gaussian



$\lambda_j = \max\{y_j, 0\}$ with
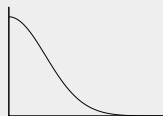$y_j \sim \mathcal{N}(\mu_\lambda, \sigma_\lambda)$

# FARMS: Bayes framework

BIOINF

## Posterior

$$p(\boldsymbol{\lambda}, \boldsymbol{\Psi} \,|\, \{x\}) \;\propto\; p(\{x\} \,|\, \boldsymbol{\lambda}, \boldsymbol{\Psi}) \; p(\boldsymbol{\lambda})$$

## Prior knowledge

- Positive $\boldsymbol{\lambda}$ ensure positive probe correlation
- Most genes show no or small signal (large signals are of interest in a study)

## Rectified Gaussian

$\lambda_j = \max\{y_j, 0\}$ with
$y_j \sim \mathcal{N}(\mu_\lambda, \sigma_\lambda)$

# FARMS: Bayes framework

## Posterior

$$p(\boldsymbol{\lambda}, \boldsymbol{\Psi} \,|\, \{x\}) \;\propto\; p(\{x\} \,|\, \boldsymbol{\lambda}, \boldsymbol{\Psi}) \; p(\boldsymbol{\lambda})$$

## Prior knowledge

- Positive $\boldsymbol{\lambda}$ ensure positive probe correlation
- Most genes show no or small signal (large signals are of interest in a study)

## Rectified Gaussian



$\lambda_j = \max\{y_j, 0\}$ with
$y_j \sim \mathcal{N}(\mu_\lambda, \sigma_\lambda)$

# FARMS: EM updates

**BIOINF**

## E-step:

$$E_{z_i|x_i}(z_i) = \mu_{z_i|x_i} \quad \text{and} \quad E_{z_i|x_i}(z_i^2) = \mu_{z_i|x_i}^2 + \sigma_{z_i|x_i}^2$$
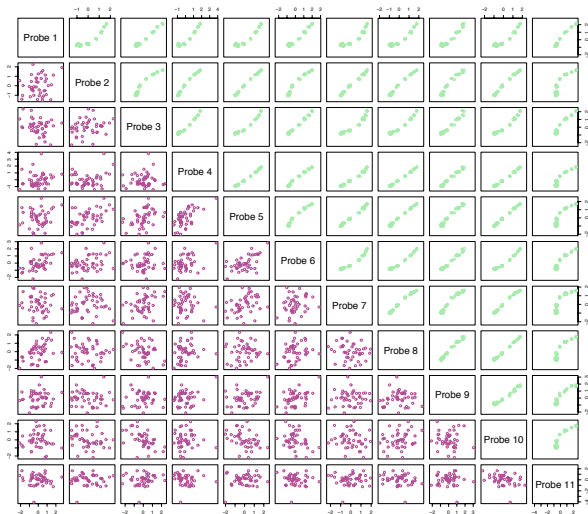
## M-step:

$$\lambda_j^{\text{Gauss}} = \left( \frac{1}{N} \sum_{i=1}^{N} x_{ij} \, E_{z_i|x_i}(z_i) + \frac{1}{N} \frac{\mu_\lambda \, \Psi_{jj}^{\text{old}}}{\sigma_\lambda^2} \right) \left( \frac{1}{N} \sum_{i=1}^{N} E_{z_i|x_i}(z_i^2) + \frac{1}{N} \frac{\Psi_{jj}^{\text{old}}}{\sigma_\lambda^2} \right)^{-1}$$

$$\lambda_j^{\text{new}} = \begin{cases} \lambda_j^{\text{Gauss}} & \text{for} \quad \lambda_j^{\text{Gauss}} > 0 \\ 0 & \text{for} \quad \lambda_j^{\text{Gauss}} \le 0 \end{cases},$$

$$\Psi_{jj}^{\text{new}} = \left[ \text{diagvect} \left( \frac{1}{N} \sum_{i=1}^{N} x_i x_i^T \right) \right]_j - \lambda_j^{\text{new}} \left[ \frac{1}{N} \sum_{i=1}^{N} E_{z_i|x_i}(z_i) \, x_i \right]_j +$$

$$\frac{1}{N} \frac{\Psi_{jj}^{\text{old}}}{\sigma_\lambda^2} \lambda_j^{\text{new}} (\mu_\lambda - \lambda_j^{\text{new}})$$

## FARMS: EM updates

**BIOINF**

### E-step:

$$E_{z_i|x_i}(z_i) = \mu_{z_i|x_i} \quad \text{and} \quad E_{z_i|x_i}(z_i^2) = \mu_{z_i|x_i}^2 + \sigma_{z_i|x_i}^2$$

### M-step:

$$\lambda_j^{\text{Gauss}} = \left(\frac{1}{N}\sum_{i=1}^{N} x_{ij}\, E_{z_i|x_i}(z_i) + \frac{1}{N}\frac{\mu_\lambda\,\Psi_{jj}^{\text{old}}}{\sigma_\lambda^2}\right)\left(\frac{1}{N}\sum_{i=1}^{N} E_{z_i|x_i}(z_i^2) + \frac{1}{N}\frac{\Psi_{jj}^{\text{old}}}{\sigma_\lambda^2}\right)^{-1}$$

$$\lambda_j^{\text{new}} = \begin{cases} \lambda_j^{\text{Gauss}} & \text{for} \quad \lambda_j^{\text{Gauss}} > 0 \\ 0 & \text{for} \quad \lambda_j^{\text{Gauss}} \leq 0 \end{cases},$$

$$\Psi_{jj}^{\text{new}} = \left[\text{diagvect}\left(\frac{1}{N}\sum_{i=1}^{N} x_i x_i^T\right)\right]_j - \lambda_j^{\text{new}}\left[\frac{1}{N}\sum_{i=1}^{N} E_{z_i|x_i}(z_i)\, x_i\right]_j +$$

$$\frac{1}{N}\frac{\Psi_{jj}^{\text{old}}}{\sigma_\lambda^2}\lambda_j^{\text{new}}(\mu_\lambda - \lambda_j^{\text{new}})$$

# FARMS: Filtering by signal variance

# FARMS: $z$-posterior

### Variance of $z \mid x$

Model

$$x = \boldsymbol{\lambda}\, z + \boldsymbol{\epsilon}$$

and Gaussian $z$-prior $\mathcal{N}(0,1)$ results in the $z$-posterior $p(z \mid x)$:

$$z \mid x \sim \mathcal{N}\left(\mu_{z|x}, \ \sigma_{z|x}^2\right)$$

$$\mu_{z|x} = (x)^T\, \boldsymbol{\Psi}^{-1}\boldsymbol{\lambda}\, \left(1 + \boldsymbol{\lambda}^T\boldsymbol{\Psi}^{-1}\boldsymbol{\lambda}\right)^{-1}$$

$$\sigma_{z|x}^2 = \left(1 + \boldsymbol{\lambda}^T\boldsymbol{\Psi}^{-1}\boldsymbol{\lambda}\right)^{-1}$$

## FARMS: The I/NI call

The variance of $z$ is decomposed into a signal and a noise part:

$$1 = \mathrm{var}(z) = \frac{1}{N} \sum_{i=1}^{N} \mathrm{E}_{z_i|x_i} \left( z_i^2 \right) = \frac{1}{N} \sum_{i=1}^{N} \left( \mu_{z_i|x_i}^2 + \sigma_{z_i|x_i}^2 \right)$$

$$\frac{1}{N} \sum_{i=1}^{N} \sigma_{z_i|x_i}^2 = 1 - \frac{1}{N} \sum_{i=1}^{N} \mu_{z_i|x_i}^2$$

$$\sigma_{z|x}^2 = 1 - \frac{1}{N} \sum_{i=1}^{N} \mu_{z_i|x_i}^2 = \left( 1 + \boldsymbol{\lambda}^T \boldsymbol{\Psi}^{-1} \boldsymbol{\lambda} \right)^{-1}$$

$\sigma_{z|x}^2$ is called the "Informative/NonInformative (I/NI) call" and is one minus the signal variance. We see that large $\boldsymbol{\lambda}$ (going with low noise $\boldsymbol{\Psi}$) leads to low variance of $z \,|\, x$ which means a precise conditional $z$.

# FARMS: Independent I/NI calls filtering

## Independent filtering increases detection power for high-throughput experiments

**Richard Bourgon[a], Robert Gentleman[b], and Wolfgang Huber[c,1]**

[a]European Bioinformatics Institute, Cambridge CB10 1SD, United Kingdom; [b]Genentech, Inc., 1 DNA Way, South San Francisco, CA 94080-4990; and [c]European Molecular Biology Laboratory, 69117 Heidelberg, Germany

- For permutation invariant test statistics and for the *t*-test statistic $T$ (only for Gaussian $z$-prior), the I/NI call filter applied to null hypotheses is independent of the statistic

- This guarantees type I error rate control if first filtering by I/NI calls, then using these statistics, and finally applying correction for multiple testing.

- http://www.bioinf.jku.at/software/cnfarms/proof_ini.pdf

# FARMS: I/NI calls distribution



### Bimodal distribution

- Enforced by the parameter prior
- Modes clearly separated (insensitive for filtering threshold)
- Works for unbalanced data (few samples contain a signal) in contrast to variance filtering (Bourgon et al. (2010))
- Works for few genes with a signal

# A pipeline for gene expression analysis

BIOINF



**Figure:** Probe-level modeling is a mandatory step

# Receiver Operator Characteristics (ROC)

Affycomp II / GoldenSpike Benchmark (AUC - area under the curve):

|  | Intensity | FARMS | RMA | GCRMA | MAS 5.0 | MBEI |
|---|---|---|---|---|---|---|
| HGU133 | Low | **0.94** | 0.51 | 0.62 | 0.07 | 0.21 |
|  | Med | **0.99** | 0.91 | 0.94 | 0.00 | 0.43 |
|  | High | **1.00** | 0.64 | 0.59 | 0.00 | 0.16 |
|  | Mean | **0.95** | 0.60 | 0.69 | 0.05 | 0.26 |
| HGU95 | Low | **0.91** | 0.57 | 0.45 | 0.09 | - |
|  | Med | **1.00** | 0.91 | 0.91 | 0.00 | - |
|  | High | **0.98** | 0.96 | 0.92 | 0.00 | - |
|  | Mean | **0.93** | 0.65 | 0.57 | 0.06 | - |
| GoldenSpike |  | **0.85** | 0.76 | 0.78 | 0.28 | 0.39 |

Computational costs for processing 60 arrays

|  | FARMS | RMA | MAS 5.0 | MBEI |
|---|---|---|---|---|
| Computational time [s] | **92** | 384 | 851 | 591 |

# Results I/NI call

BIOINF

- Leads on average to 84 ($\pm$1.5)% exclusion rate
  - Applied on 30 real life studies
  - A/P calls excluded only 33 ($\pm$1)%
- Validation was carried out on spiked-in data:

Exclusion rate on spiked-in data sets:

|  | INFORMATIVE | NON-INFORMATIVE | EXCLUSION RATE | DETECTED SPIKED-INS | DETECTED PSEUDO SPIKED-INS |
|---|---|---|---|---|---|
| HGU133A | 81 | 22219 | 99.63% | 42/42 | 28/28* |
| HGU95_V2 | 56 | 12570 | 99.56% | 14/14 | 5/5** |
| HU. GENE 1.0 ST | 40 | 19,753 | 99.80% | 15/15*** | - |

*McGee et al. 2006; **Wolfinger and Chu 2002; Cope et al. 2004; ***long spiked-in fragments

# I/NI call vs. A/P call



**Figure:** Variance and mean of genes selected by A/P calls and I/NI calls.

# A pipeline for copy number analysis



**Figure:** Copy number analysis for (Affymetrix) DNA genotyping arrays as a three-step pipeline: (1) Normalization, (2) Modeling, and (3) Segmentation.

# Benchmark data sets



- 30 male and 30 female CEU founders
  - SNP 6.0 and 250K NSP Arrays
  - Classification task: distinguish males from females by their copy number on the X chromosome

- Evaluation on:
  - Single-locus / multi-loci classification (window mode)
  - Multi-loci summarization with
    - cn.FARMS
    - Median locus for dChip and CRMA_v2

# ROC-Curve (250K arrays)

**single-locus**



**multi-loci, 3 markers**



## TPR / FPR

True positive rate (TPR) = TP/(TP+FN)

False positive rate (FPR) = FP/(FP+TN)

# ROC-Curve (250K arrays)

## single-locus



## multi-loci, 3 markers



## TPR / FPR

True positive rate (TPR) = TP/(TP+FN)

False positive rate (FPR) = FP/(FP+TN)

# ROC-Curve (SNP 6.0 arrays)

**single-locus**



**multi-loci, 3 markers**



## TPR / FPR

True positive rate (TPR) $= TP/(TP+FN)$

False positive rate (FPR) $= FP/(FP+TN)$

# ROC-Curve (SNP 6.0 arrays)

**single-locus**



**multi-loci, 3 markers**



**TPR / FPR**

True positive rate (TPR) = TP/(TP+FN)

False positive rate (FPR) = FP/(FP+TN)

# Results cn.FARMS

BIOINF

| Loci | Criteria | Affymetrix Mapping250K_NSP | | | Affymetrix SNP 6.0 | | |
|---|---|---|---|---|---|---|---|
| | | **cn.FARMS** | CRMA_v2 | dChip | **cn.FARMS** | CRMA_v2 | dChip |
| 1 | AUC | **0.9852** | 0.9820 | 0.9819 | **0.9838** | 0.9807 | 0.9721 |
| | FP | **8472** | 9106 | 9018 | 56145 | 68593 | 77438 |
| | P-VALUE | – | 1.8e-65 | 3.1e-26 | – | 1e-1160 | 1e-6049 |
| 2 | AUC | **0.9983** | 0.9974 | 0.9969 | **0.9983** | 0.9963 | 0.9894 |
| | FP | **1375** | 1449 | 1611 | 9777 | 11705 | 18039 |
| | P-VALUE | – | 2.7e-4 | 2.5e-12 | – | 1e-317 | 1e-3713 |
| 3 | AUC | **0.9998** | 0.9995 | 0.9992 | **0.9998** | 0.9990 | 0.9953 |
| | FP | **240** | 366 | 440 | 1573 | 3462 | 6625 |
| | P-VALUE | – | 2.6e-38 | 7.2e-58 | – | 1e-896 | 1e-3455 |

**Table:** AUC values at the sex classification task for 59 HapMap CEU founders based on the X chromosome copy numbers:

# CNV detection benchmark

- "The International HapMap Project" phase 2 data set with Affymetrix SNP 6.0 arrays
  - Goal is to identify true rare CNV regions with a low FDR
  - "True CNV regions" are those regions which were detected and verified by different bio-technologies
    - NimbleGen tiling arrays, Agilent CGH arrays, Illumina Infinium genotyping (Human660W)
  - 2,515 true CNV regions as reference
- CNV calling criteria:
  - I/NI call for cn.FARMS
  - Variance of the raw copy numbers on the samples for dChip and CRMA_v2

# CNV detection plot



**Figure:** CNV calling plots across chromosome 4 for 3-loci regions (each point in the plot summarizes 3 loci).

# CNV detection on HapMap (multi-loci 3)

**BIOINF**

## Chromosome 8



## Whole genome



## Precision / Recall

Recall = TP/(TP+FN)

Precision = TP/(TP+FP) = 1 - FDR

# CNV detection on HapMap (multi-loci 3)

**BIOINF**

## Chromosome 8



## Whole genome



## Precision / Recall

Recall = TP/(TP+FN)

Precision = TP/(TP+FP) = 1 - FDR

# CNV detection on HapMap (multi-loci 5)

**BIOINF**

## Chromosome 8



## Whole genome



## Precision / Recall

Recall = TP/(TP+FN)

Precision = TP/(TP+FP) = 1 - FDR

# CNV detection on HapMap (multi-loci 5)

**BIOINF**

## Chromosome 8



## Whole genome



## Precision / Recall

Recall = TP/(TP+FN)
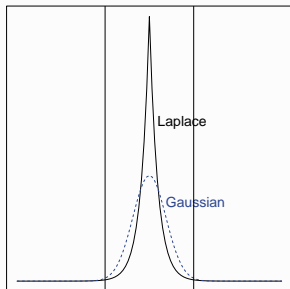
Precision = TP/(TP+FP) = 1 - FDR

# Interim results

- cn.FARMS outperforms aroma.affymetrix, *dChip*, *CNAG* and *CNAT* in terms of sensitivity and specificity
  - Shows good signal detection while being robust against measurement noise
- I/NI call correctly prioritizes CNV regions of interest
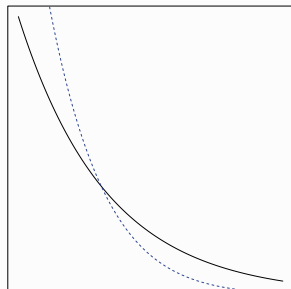  - Reduces the FDR at CNV detection

# Rare CNV events

## Sparse data

CNV data is sparse with an kurtosis larger than 30 $\rightarrow$ change the model assumption to a Laplacian distributed hidden variable $z$.
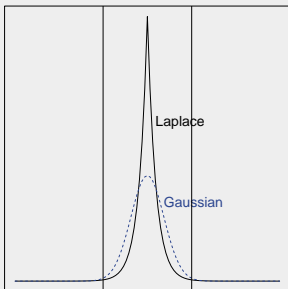


Gauss vs. Laplace



Close up Gauss vs. Laplace
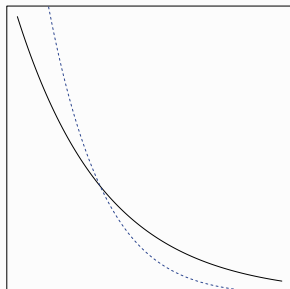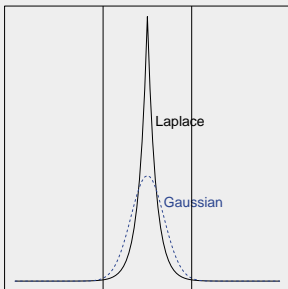
# Rare CNV events

## Sparse data

CNV data is sparse with an kurtosis larger than 30 $\rightarrow$ change the model assumption to a Laplacian distributed hidden variable $z$.

### Gauss vs. Laplace



### Close up Gauss vs. Laplace
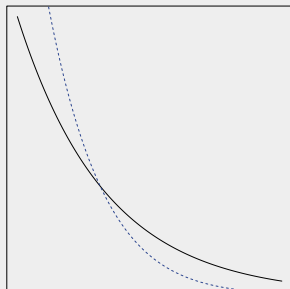
# Rare CNV events

**BIOINF**

### Sparse data

CNV data is sparse with an kurtosis larger than 30 $\rightarrow$ change the model assumption to a Laplacian distributed hidden variable $z$.

### Gauss vs. Laplace



Laplace

Gaussian

### Close up Gauss vs. Laplace

# Laplacian FARMS

## Data likelihood

$$p(\{x\} \mid \boldsymbol{\lambda}, \boldsymbol{\Psi}) = \int p(\{x\} \mid z, \boldsymbol{\lambda}, \boldsymbol{\Psi}) \; p(z) \; dz$$

## Problem

- The **likelihood is analytically intractable** for the non-Gaussian prior

## Solution

- Variational EM approach
- Based on a local Gaussian approximation to the mode

# Laplacian FARMS

## Data likelihood

$$p(\{x\} \mid \boldsymbol{\lambda}, \boldsymbol{\Psi}) = \int p(\{x\} \mid z, \boldsymbol{\lambda}, \boldsymbol{\Psi}) \ p(z) \ dz$$

## Problem

- The **likelihood is analytically intractable** for the non-Gaussian prior

## Solution

- Variational EM approach
- Based on a local Gaussian approximation to the mode

# Laplacian FARMS

**BIOINF**

## Data likelihood

$$p\left(\{x\}\,|\,\boldsymbol{\lambda},\boldsymbol{\Psi}\right) \,=\, \int p\left(\{x\}\,|\,z,\boldsymbol{\lambda},\boldsymbol{\Psi}\right)\,p\left(z\right)\,dz$$

## Problem

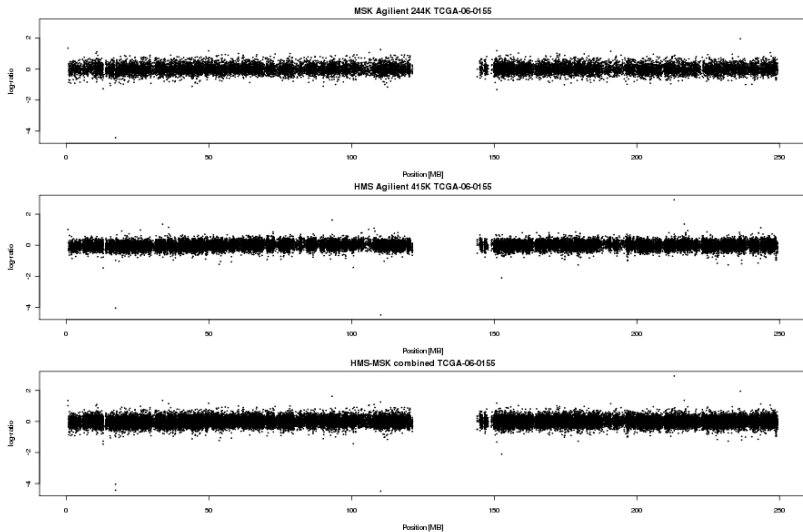- The **likelihood is analytically intractable** for the non-Gaussian prior

## Solution

- Variational EM approach
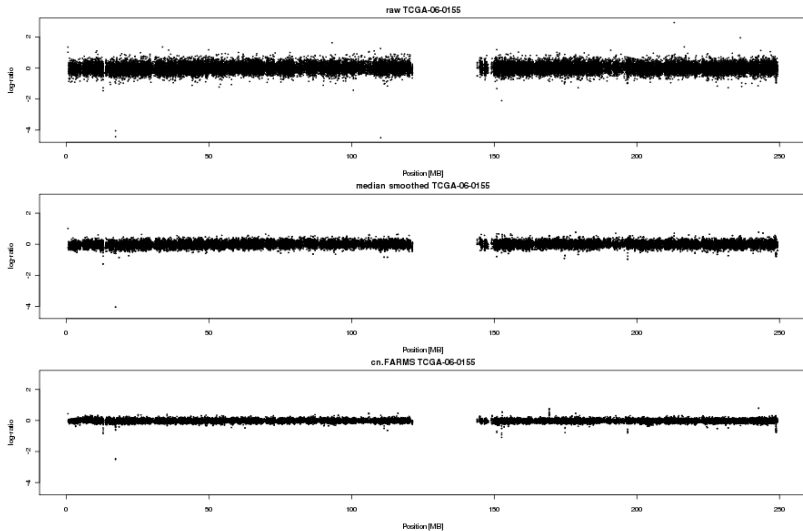- Based on a local Gaussian approximation to the mode

# CAMDA copy number data sets

- Glioblastoma multiforme data sets
  - 167 Agilent 415K CGH arrays from Harvard
  - 262 Agilent 244A CGH arrays from Harvard
  - 461 Agilent 244A CGH arrays from MSKCC
  - 533 Affymetrix SNP 6.0 arrays from Broad
  - 432 Illumina HumanHap 550 from Stanford
- CN data for SNP 6.0 and HumanHap 550 were not available
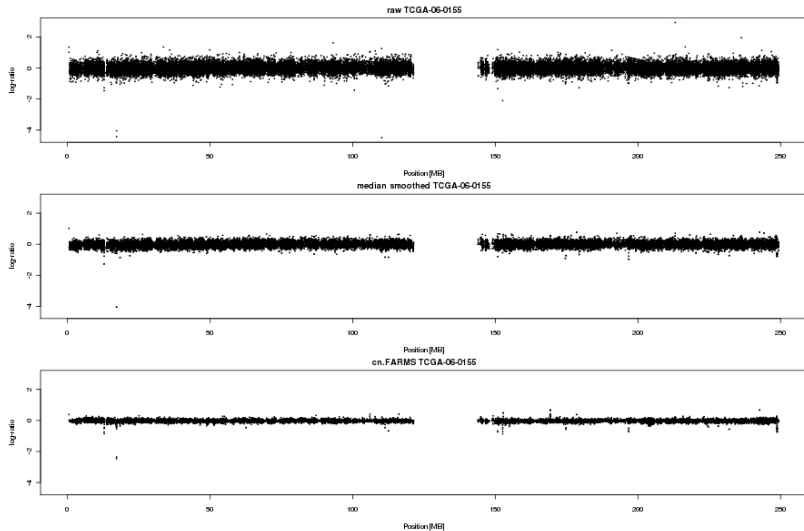- 167 matched arrays HMS 415K and MSKCC 244A remain
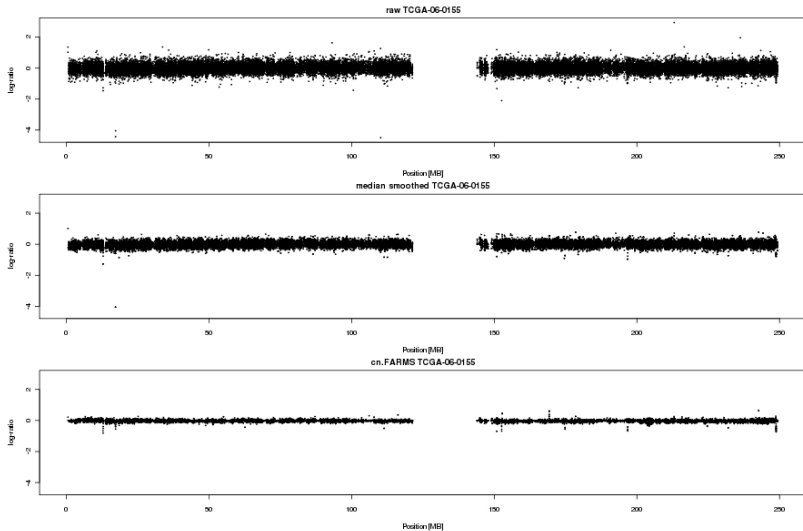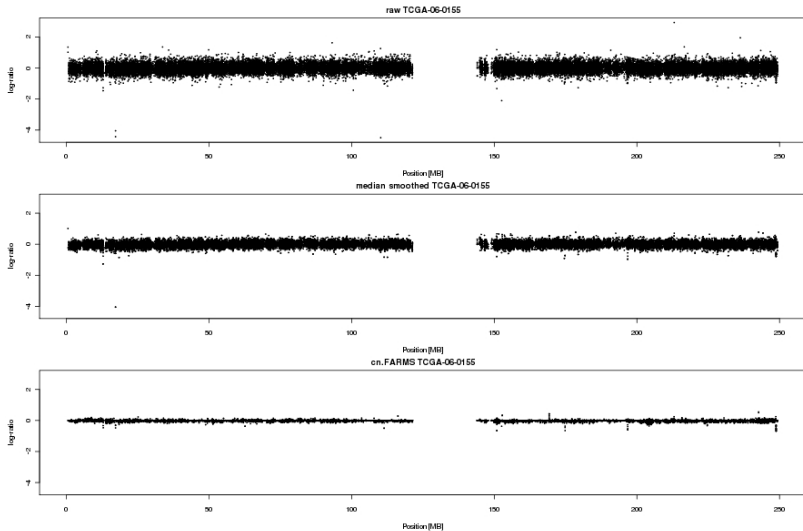
# Merged raw data (Chromosome 1)

# Prior weight 0.5

# Prior weight 1.5

**BIOINF**

# Prior weight 2.0

# Prior weight 2.5

# Conclusion

BIOINF

- Latent variable models decompose observation into noise and signal
- Remove noise so that aberration detection take place in noise-free data
- Reduce dimensionality by filtering for signal variance

# Acknowledgments

**Johannes Kepler University**
Andreas Mayr
Andreas Mitterecker
Günter Klambauer
Martin Heusel
Ulrich Bodenhofer
Sepp Hochreiter

**Johnson & Johnson R&D**
Marianne Tuefferd
An De Bondt
Willem Talloen
Hinrich Göhlmann
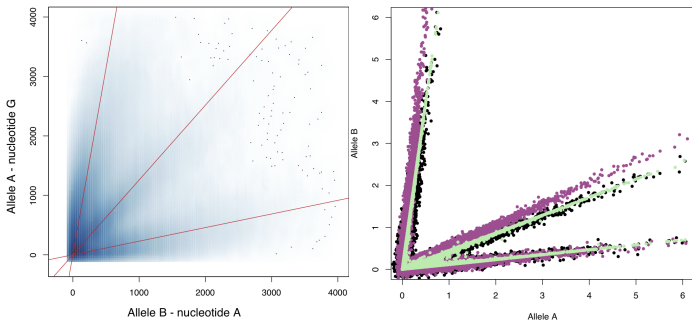
# Further information

- Clevert DA, Mitterecker A, Mayr A, Klambauer G, Tuefferd M, De Bondt A, Talloen W, Göhlmann H, and Hochreiter S: *cn.FARMS: a latent variable model to detect copy number variations in microarray data with a low false discovery rate* **Nucl. Acids Res.** (2011) 39: e79

- Talloen W, Hochreiter S, Bijnens L, Kasim A, Shkedy Z, Amaratunga D, and Göhlmann H: *Filtering data from high-throughput experiments based on measurement reliability* **PNAS** (2010) 107 (46) E173-E174

- Hochreiter S, Clevert DA, and Obermayer K: *A new summarization method for Affymetrix probe level data*. **Bioinformatics** (2006), 22: 943-949.

- Talloen W, Clevert DA, Hochreiter S, Amaratunga D, Bijnens L, Kass S and Göhlmann H: *I/NI-calls for the exclusion of non-informative genes: a highly effective filtering tool for microarray data*. **Bioinformatics** (2007) 23(21): 2897-2902

## Open source software



- FARMS, I/NI call and cn.FARMS are publicly available as Bioconductor R packages
- Software homepages:
  - http://www.bioinf.jku.at/software/farms/farms.html
  - http://www.bioinf.jku.at/software/cnfarms/cnfarms.html

# Sparse overcomplete representation



A sparse overcomplete representation of two-dimensional data $x_s \in \mathbb{R}^2$ can be modeled as: $x_s = \boldsymbol{\lambda}_s \, z_s + \boldsymbol{\epsilon}_s$ where $z_s \in \mathbb{R}^3$, $\boldsymbol{\lambda}_s \in \mathbb{R}^{2 \times 3}$. Sparseness is enforced by assuming a Laplacian prior for $z_s$:

$$p(z_s) = (2)^{-\frac{3}{2}} \prod_{l=1}^{3} \exp\left( \sqrt{2} \, |z_{sl}| \right)$$