

Presentations and posters

Tools for Reproducible Research

Karl Broman

Biostatistics & Medical Informatics, UW–Madison

kbroman.org

github.com/kbroman

@kbroman

Course web: kbroman.org/Tools4RR

Powerpoint/Keynote

- + Standard
- + Easy to share slides
- + WYSIWYG (mostly)
- + Fancy animations
- Font problems
- Lots of copy-paste
- Hard to get equations
- Not reproducible

On the Complexity of SNP Block Partitioning Under the Perfect Phylogeny Model

Jens Gramm¹ Tzvika Hartman² Till Nierhoff³
Roded Sharan⁴ **Till Tantau⁵**

¹Universität Tübingen, Germany

²Bar-Ilan University, Ramat-Gan, Israel

³International Computer Science Institute, Berkeley, USA

⁴Tel-Aviv University, Israel

⁵Universität zu Lübeck, Germany

Workshop on Algorithms in Bioinformatics, 2006

Get rid of the junk

```
\usetheme{default}
```

```
\beamertemplatenavigationsymbolsempy
```

Change colors

```
\definecolor{foreground}{RGB}{255,255,255}
\definecolor{background}{RGB}{24,24,24}
\definecolor{title}{RGB}{107,174,214}
\definecolor{subtitle}{RGB}{102,255,204}
\definecolor{hilit}{RGB}{102,255,204}
\definecolor{lolit}{RGB}{155,155,155}

\setbeamercolor{titlelike}{fg=title}
\setbeamercolor{subtitle}{fg=subtitle}
\setbeamercolor{institute}{fg=lolit}
\setbeamercolor{normal text}{fg=foreground,bg=background}
\setbeamercolor{item}{fg=foreground} % color of bullets
\setbeamercolor{subitem}{fg=lolit}
\setbeamercolor{itemize/enumerate subbody}{fg=lolit}
\setbeamerfont{itemize subitem}{\textendash}
\setbeamerfont{itemize/enumerate subbody}{size=\footnotesize}
\setbeamerfont{itemize/enumerate subitem}{size=\footnotesize}

\newcommand{\hilit}{\color{hilit}}
\newcommand{\lolit}{\color{lolit}}
```

Also, slide numbers and fonts

```
% slide number
\setbeamertemplate{footline}{%
  \raisebox{5pt}{\makebox[\paperwidth]{\hfill\makebox[20pt]{\lolit
    \scriptsize\insertframenumbers}}}\hspace*{5pt}}

% font
\usepackage{fontspec}
% http://www.gust.org.pl/projects/e-foundry/tex-gyre/
% ... heros/qhv2.004otf.zip
\setsansfont
  [ ExternalLocation = ../fonts/ ,
    UprightFont = *-regular ,
    BoldFont = *-bold ,
    ItalicFont = *-italic ,
    BoldItalicFont = *-bolditalic ]{texgyreheros}
% Palatino for notes
\setbeamerfont{note page}{family*=pplx,size=\footnotesize}
```

Title slide

```
\title{Put title here}
\subtitle{And maybe a subtitle}
\author{Author name}
\institute{Biostatistics \& Medical Informatics,
  UW{\textendash}Madison}
\date{\tt \scriptsize biostat.wisc.edu/{\textasciitilde}kbroman}

\begin{document}

{
\setbeamertemplate{footline}{} % no slide number here
\frame{
  \titlepage

\note{
  Summary of the talk, as a note.
}
} }
```

Typical slide

```
\begin{frame}{Title of slide}

\bbi
  \item Bullet 1
  \item Bullet 2
  \item Bullet 3
\ei

\note{
  Put a note here
}
\end{frame}
```


Typical slide

```
\begin{frame}{Title of slide}

\vspace{24pt} \begin{itemize} \itemsep8pt
  \item Bullet 1
  \item Bullet 2
  \item Bullet 3
\end{itemize}

\note{
  Put a note here
}
\end{frame}
```

Slide with a figure

```
\begin{frame}{Title of slide}

\figh{Figs/a_figure.png}{0.75}

\note{
  Put a note here
}
\end{frame}
```

Slide with a figure

```
\begin{frame}{Title of slide}

\centerline{\includegraphics[height=0.75\textheight]{%
    Figs/a_figure.png}}

\note{
  Put a note here
}
\end{frame}
```

Figures with KnitR

```
<<knitr_options, echo=FALSE>>=
opts_chunk$set(echo=FALSE, fig.height=7, fig.width=10)
change_colors <-
function(bg=rgb(24,24,24, maxColorValue=255), fg="white")
  par(bg=bg, fg=fg, col=fg, col.axis=fg, col.lab=fg,
      col.main=fg, col.sub=fg)
@

<<pdf_figure>>=
change_colors()
par(las=1)
n <- 100
x <- rnorm(n)
y <- 2*x + rnorm(n)
plot(x, y, pch=16, col="slateblue")
@
```

Figures with KnitR

```
% << >>= all on one line!  
<<png_figure, dev="png", fig.align="center",  
  dev.args=list(pointsize=30),  
  fig.height=15, fig.width=15, out.height="0.75\\textheight",  
  out.width="0.75\\textheight">>=  
change_colors(bg=rgb(32,32,32,maxColorValue=255))  
par(las=1)  
n <- 251  
x <- y <- seq(-pi, pi, len=n)  
z <- matrix(ncol=n, nrow=n)  
for(i in seq(along=x))  
  for(j in seq(along=y))  
    z[i,j] <- sin(x[i]) + cos(y[j])  
image(x,y,z)  
@
```

Slides with notes

```
\documentclass[12pt,t]{beamer}  
\setbeameroption{hide notes}  
\setbeamertemplate{note page}[plain]
```

```
\documentclass[12pt,t,handout]{beamer}  
\setbeameroption{show notes}  
\setbeamertemplate{note page}[plain]  
\def\notescolors{1}
```

```
\ifx\notescolors\undefined % slides  
  \definecolor{foreground}{RGB}{255,255,255}  
  \definecolor{background}{RGB}{24,24,24}  
\else % notes  
  \definecolor{background}{RGB}{255,255,255}  
  \definecolor{foreground}{RGB}{24,24,24}  
\fi
```

Simple animations

```
\begin{frame}{Bullets entering one at a time}

\bbi
\item Bullet 1
\onslide<2->{\item Bullet 2}
\onslide<3->{\item Bullet 3}
\onslide<4->{\item Bullet 4}
\ei

\note{
  Do this sparingly.
}
\end{frame}
```

Simple animations

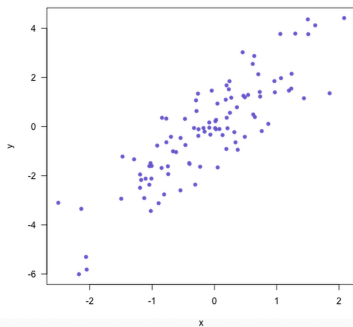
```
\begin{frame}{Bullets entering one at a time}

\bbi
\item {\lollit \only<1>{\color{foreground}} Bullet 1}
\item {\lollit \only<2>{\color{foreground}} Bullet 2}
\item {\lollit \only<3>{\color{foreground}} Bullet 3}
\item {\lollit \only<4>{\color{foreground}} Bullet 4}
\ei

\note{
  Do this sparingly.
}
\end{frame}
```


Slidify and R Markdown

A figure



Slidify and R Markdown

```
## Slide title

- Bullet 1
- Bullet 2
- Bullet 3
- Bullet 4

---

## A figure

```{r a_figure, echo=FALSE, fig.align="center"}
par(las=1)
n <- 100
x <- rnorm(n)
y <- 2*x + rnorm(n)
plot(x, y, pch=16, col="slateblue")
```
```

Using slidify

```
library(devtools)
install_github("slidify", "ramnathv")
install_github("slidifyLibraries", "ramnathv")

library(slidify)
setwd("~/Docs/Talks/")
author("slidify_example")

# edit ~/Docs/Talks/slidify_example/index.Rmd

slidify("index.Rmd")
browseURL("index.html")
```

YAML header

```
---
title       : Slidify example
subtitle    : Tools for reproducible research
author      : Karl Broman
job         : Biostatistics & Medical Informatics, UW-Madison
framework   : io2012           # {io2012, html5slides, shower, ...}
highlighter  : highlight.js    # {highlight.js, prettify, highlight}
hitheme     : tomorrow         #
widgets     : [mathjax]        # {mathjax, quiz, bootstrap}
mode        : standalone       # {selfcontained, standalone, draft}
---
```

Change the title slide colors

```
<style>
.title-slide {
  background-color: #EEE;
}

.title-slide hgroup > h1,
.title-slide hgroup > h2 {
  color: #005;
}
</style>
```

Beamer-based posters

Identifying and correcting sample mix-ups in eQTL data

Karl W Broman¹, Mark P Keller², Ameer Teo Boman³, Daniele M Greenawald⁴,

Christina Kendziora⁵, Eric E Schadt⁶, Štursnik Sen⁷, Brian S Yandell^{8,9}, and Alan D Attie⁶

¹Department of Medical Informatics, ²Biostatistics, ³Statistics, ⁴Heriotskupa, UW Madison, ⁵Merck & Co., Inc., ⁶Pacific Biosciences, ⁷UC-San Francisco

Abstract

In a recent interview with more than 100 authors and genome-wide gene expression data on six tissues, we identified a high proportion of sample mix-ups in the genotype data, on the order of 17%.

Local eQTL genetic test following genome expression with increased false effect may be used to have a classifier for predicting an individual's eQTL genotype from its gene expression levels. By considering multiple eQTLs and their related transcripts, we identified numerous individuals whose predicted eQTL genotypes based on their expression data did not match their observed genotypes, and then used an to identify other individuals whose genotypes did match the predicted eQTL genotypes.

The concordance of predictions across six tissues indicated that the problem was due to mix-ups in the genotypes. Consideration of the joint probabilities of the samples indicated a number of mix-up by one and six-by-two cases, likely the result of pipetting errors.

Such sample mix-ups are a problem to any genetic study. As we show, eQTL data often can be identified, and even correct, such problems.

Data

- ~500 Mb x 100M reference data, all 46/46
- Genotype at 267 SNPs (Adyze/10k chip)
- Gene expression in six tissues (Affymetrix arrays)
- Before, genome-wide results. Significant gene-tissue hits (see)
- Numerous clinical phenotypes in a local single tissue and across tissues

Initial observation: Sex swaps



We should have:
Y: females, B: Y or B: Y
X: males, heterozygous in B

But 30 male had X chromosome genotype that combined with sex.

Which are correct: genotypes or sexes?

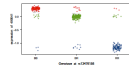
We could look for a transcripting, Xist whose expression level is diagnostic for sex.

Even better, we can look at transcripts with strong local eQTL, for which genotype is strongly associated with expression levels.

Transcripts with strong local eQTL can distinguish the genotypes. By considering multiple such transcripts across the genome, we can learn a DNA language.

- eQTL - quantitative trait locus: a genomic region that influences a quantitative trait
- eQTL - expression QTL, a QTL, that influences the level of expression of a gene

A diagnostic transcript



Colors indicate the inferred eQTL genotype according to a nearest neighbor classifier, with gray points not called.

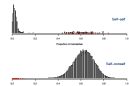
The method

Identify expression traits with strong local eQTL (that is, for which genotype of the transcript's genomic position is strongly associated with its expression level).

For each trait, create a classifier for predicting eQTL genotype from expression phenotype.

In each pair of sites, calculate the proportion of mismatches between the observed eQTL genotypes of one tissue and the inferred eQTL genotypes of the other.

Proportions of mismatches in eQTL genotypes



Decisions



There were ~100 tissue with genotypes and ~500 with expression data.

For each tissue, we did the proportion of mismatches between its observed genotype data and the genotypes inferred from the corresponding gene expression data, against the estimated false proportion of mismatches, comparing that observed genotype data to each set of inferred genotypes.

Inferred genotype mix-ups

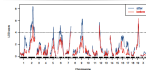


The gene expression data from the multiple tissues were concordant and both called that the problems were in the genotype data.

(This did, however, identify and correct a small number of simple mix-ups within each of the six sets of gene expression data. This was done by crossing single pairs of tissues and measuring the correlation in a tissue's gene expression across tissues.)

The bulk of the problems concerned apparent genotyping errors in the genotyping plates: a series of all-by-one and all-by-two errors covering half of each of two plates. (This did not happen in MassMut).

Improved results



LOD curves for loci, indicating the evidence for QTL before and after correcting the sample mix-ups. The corrected data give stronger evidence and more QTL.

Summary

- Sample mix-ups happen
- With eQTL data, we can both identify and correct mix-ups
- The general idea here has wide application for high throughput data
- R package: <https://github.com/kwbroman/samplemixup>
- Key article: <https://doi.org/10.1093/bioinformatics/btt101>, 2010

Contact



Karl Broman
karlb@pacbio.com or kb@u.wisc.edu
<http://www.kbroman.com>

This work was supported by grants from National Institutes of Health and Pacific Biosciences.

Beamer-based posters

Data visualizations should be more interactive

Karl W Broman

Biostatistics & Medical Informatics, University of Wisconsin-Madison

Introduction

- High-dimensional data can be bewildering.
 - With 2000 gene expression arrays, you'd think you'd make a lot of graphs, but we need to make 9 graphs. We can't look at 2000 histograms, so why look at any?
 - Interactive graphics provide a solution to this problem.
 - For context to the visualization:
 - This visualization is also more important than formal statistics.
 - All graphs could be improved with some interactivity.
- Download: bit.ly/9waz2114

Opportunities

- Exploration
 - Being persistent
 - Identifying outliers
 - The fancy plot is 100x faster than
 - Using histograms
- Reports for collaboration
 - Using histograms
 - Allow deeper exploration of the results
 - Can you do sample questions?
- Big Data
 - Don't just do an summary statistics
 - Graphs compressed information, but with access to the details
 - Consistent data display
 - More exploration, more connections
- Teaching
 - Can't always look at all plots with it
 - Animated distributions of key concepts
 - Demonstrates data exploration
 - Enable users to explore on their own

Barriers

- We never learned how
- It's a hassle
- No consistent platforms
- Journal articles are static and obscure methods
- Most statisticians are still creating terrible static plots (even some interactive tables)

But... many exciting new tools

- HTML5 + Stable vector graphics (SVG)
- Incredible power of modern web browsers
- JavaScript based web tools
- Resizable fonts

DS

- JavaScript library for manipulating HTML and SVG elements
- Connects data to elements
- Low level, but flexible

Other options

- Tabletop (1) and Identity D
- ggplot (2nd, 3rd) and various (4th, 5th, 6th, 7th, 8th, 9th, 10th, 11th, 12th, 13th, 14th, 15th, 16th, 17th, 18th, 19th, 20th, 21st, 22nd, 23rd, 24th, 25th, 26th, 27th, 28th, 29th, 30th, 31st, 32nd, 33rd, 34th, 35th, 36th, 37th, 38th, 39th, 40th, 41st, 42nd, 43rd, 44th, 45th, 46th, 47th, 48th, 49th, 50th, 51st, 52nd, 53rd, 54th, 55th, 56th, 57th, 58th, 59th, 60th, 61st, 62nd, 63rd, 64th, 65th, 66th, 67th, 68th, 69th, 70th, 71st, 72nd, 73rd, 74th, 75th, 76th, 77th, 78th, 79th, 80th, 81st, 82nd, 83rd, 84th, 85th, 86th, 87th, 88th, 89th, 90th, 91st, 92nd, 93rd, 94th, 95th, 96th, 97th, 98th, 99th, 100th)
- Node.js (34th, 35th, 36th, 37th, 38th, 39th, 40th, 41st, 42nd, 43rd, 44th, 45th, 46th, 47th, 48th, 49th, 50th, 51st, 52nd, 53rd, 54th, 55th, 56th, 57th, 58th, 59th, 60th, 61st, 62nd, 63rd, 64th, 65th, 66th, 67th, 68th, 69th, 70th, 71st, 72nd, 73rd, 74th, 75th, 76th, 77th, 78th, 79th, 80th, 81st, 82nd, 83rd, 84th, 85th, 86th, 87th, 88th, 89th, 90th, 91st, 92nd, 93rd, 94th, 95th, 96th, 97th, 98th, 99th, 100th)
- Anytime.js (41st, 42nd, 43rd, 44th, 45th, 46th, 47th, 48th, 49th, 50th, 51st, 52nd, 53rd, 54th, 55th, 56th, 57th, 58th, 59th, 60th, 61st, 62nd, 63rd, 64th, 65th, 66th, 67th, 68th, 69th, 70th, 71st, 72nd, 73rd, 74th, 75th, 76th, 77th, 78th, 79th, 80th, 81st, 82nd, 83rd, 84th, 85th, 86th, 87th, 88th, 89th, 90th, 91st, 92nd, 93rd, 94th, 95th, 96th, 97th, 98th, 99th, 100th)
- googleVis (46th, 47th, 48th, 49th, 50th, 51st, 52nd, 53rd, 54th, 55th, 56th, 57th, 58th, 59th, 60th, 61st, 62nd, 63rd, 64th, 65th, 66th, 67th, 68th, 69th, 70th, 71st, 72nd, 73rd, 74th, 75th, 76th, 77th, 78th, 79th, 80th, 81st, 82nd, 83rd, 84th, 85th, 86th, 87th, 88th, 89th, 90th, 91st, 92nd, 93rd, 94th, 95th, 96th, 97th, 98th, 99th, 100th)
- Skippy (48th, 49th, 50th, 51st, 52nd, 53rd, 54th, 55th, 56th, 57th, 58th, 59th, 60th, 61st, 62nd, 63rd, 64th, 65th, 66th, 67th, 68th, 69th, 70th, 71st, 72nd, 73rd, 74th, 75th, 76th, 77th, 78th, 79th, 80th, 81st, 82nd, 83rd, 84th, 85th, 86th, 87th, 88th, 89th, 90th, 91st, 92nd, 93rd, 94th, 95th, 96th, 97th, 98th, 99th, 100th)
- ggplot (49th, 50th, 51st, 52nd, 53rd, 54th, 55th, 56th, 57th, 58th, 59th, 60th, 61st, 62nd, 63rd, 64th, 65th, 66th, 67th, 68th, 69th, 70th, 71st, 72nd, 73rd, 74th, 75th, 76th, 77th, 78th, 79th, 80th, 81st, 82nd, 83rd, 84th, 85th, 86th, 87th, 88th, 89th, 90th, 91st, 92nd, 93rd, 94th, 95th, 96th, 97th, 98th, 99th, 100th)
- Richark (50th, 51st, 52nd, 53rd, 54th, 55th, 56th, 57th, 58th, 59th, 60th, 61st, 62nd, 63rd, 64th, 65th, 66th, 67th, 68th, 69th, 70th, 71st, 72nd, 73rd, 74th, 75th, 76th, 77th, 78th, 79th, 80th, 81st, 82nd, 83rd, 84th, 85th, 86th, 87th, 88th, 89th, 90th, 91st, 92nd, 93rd, 94th, 95th, 96th, 97th, 98th, 99th, 100th)

simple ↔ flexible

Choose one, I choose flexible.

Summary

- For high-dimensional data, good visualizations are critical
- Interactive graphics, responsive, beautiful
- Reduce exploration
- Use good information tools
- Enable connections with access to the details
- Visualizations must be tailored to the data and questions
- LOD on the top level, but it
 - is easily flexible (like a zoom graph)
 - to the bottom of exploration
 - Can provide other layers of resolution
- It's a glorious package under development (github.com, kbroman.com)

Acknowledgments

Example 1
Alisa Attar, Mark Keller, Alison De Broomer, Christina Kondratic, Brian Hankler, Eric Schmitt, Department of Biostatistics, Biostatistics & Medical Informatics, and Statistics, UW-Madison, Madison, WI

Example 2
Candace Moore, Roger Spalding, Logan Johnson, & Deep Kaula, Neuroscience Department of Biology, Statistics, and Computer Science, UW-Madison

Contact

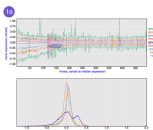
Karl Broman
kbroman@facstaff.wisc.edu
@kbroman
www.kbroman.com
www.kbroman.com/biostat
github.com/kbroman

The work was supported by grant GM087606

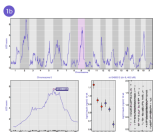
Example 1: Expression genetics

- Mouse inbreeds, 36 × 100K
- ~200 genes
- Genotypes of 200 SNPs

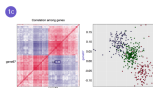
- Gene expression microarrays in six tissues
- Numerous clinical phenotypes



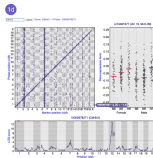
These are data from ~500 gene expression microarrays. The top panel is the 200 loci plot. Here are data on the 1, 10, and 99 percentiles for each of ~500 distributions. The distributions are sorted by their standard.



A gene score for genetic loci (linked quantitative trait loci, QTL) influencing muscle level. The LOD score is a highly thresholded ratio measuring the strength of association between genotype and phenotype. Click on a chromosome on the top and a detailed view of the LOD score for that chromosome is shown on the bottom left. In the lower-left panel, hover over markers to see names; click to view an effect plot and phenotypic-genotype plot to the right.



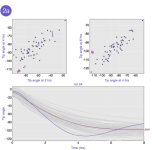
Association in gene expression among 200 genes that are influenced by a common genetic locus (QTL). The left panel is a heat map of the correlation matrix, with blue = -1 and red = +1. Hover over points in the correlation matrix on the left to see the values; click to see the corresponding heatmap on the right. Points in the heatmap are ordered by genotype of the underlying QTL.



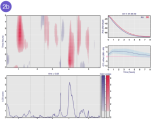
An investigation of genetic loci (QTL) influencing gene expression. In the top-left panel, the x-axis corresponds to marker location and the y-axis corresponds to the position of probes on a gene expression microarray. Each plotted point is an individual's eQTL. Hover over a point to see the probe ID and LOD score (measuring the strength of association), also highlighted as any other eQTL for that probe. Click on the point to see the LOD curve below. Hover over markers in the LOD curve plot to view marker names; click on a marker to see the phenotype-genotype plot to the right.

Example 2: Gravitropism

- Response to gravity in Arabidopsis seedlings
- Genetic variation of gravity and values over 4 hrs
- Measure the angle of the root by every 2 min



Average plot angle over time for 162 Arabidopsis lines. Hover over points in the top graphs or curves in the bottom, point to highlight the corresponding lines in the other panels.



The top-left panel is a heat map of a measure of association (LOD score) between genotype at a fixed position and the phenotype at a fixed time. Red lines indicate that 99.99% lines have larger phenotypes. When you hover over a point in the top-left plot, the LOD curve for the corresponding time is shown below, and the phenotype average and estimated genetic effect (area) time is shown to the right.

Beamer-based posters

```
\documentclass[final,plain]{beamer}
\usepackage[size=custom,width=152.4,height=91.44,scale=1.2]{%
  beamerposter}

\newlength{\sepwid}
\newlength{\onecolwid}
\newlength{\halfcolwid}
\newlength{\twocolwid}
\newlength{\threecolwid}

\setlength{\sepwid}{0.0192\paperwidth}
\setlength{\onecolwid}{0.176\paperwidth}
\setlength{\halfcolwid}{0.0784\paperwidth}
\setlength{\twocolwid}{0.3712\paperwidth}
\setlength{\threecolwid}{0.5664\paperwidth}
\setlength{\topmargin}{-0.5in}
\usetheme{confposter}
```


Basic code for a poster

```
\title{Data visualizations should be more interactive}
\author{Karl W Broman}
\institute{University of Wisconsin--Madison}

\begin{frame}[t]
\begin{columns}[t]
  \begin{column}{\sepwid}\end{column} % empty spacer column
  \begin{column}{\onecolwid}
    \begin{exampleblock}{\Large Introduction}{
      \begin{itemize} \itemsep18pt
        \item Bullet 1
        \item Bullet 2
      \end{itemize}
    }
  \colonevsep % between blocks
  \begin{block}{Barriers}{
  }
  \end{column}
\end{columns}
\end{frame}
```

Between-block spacing

```
\newcommand{\colonevsep}{\vspace{16mm}}  
\newcommand{\coltwovsep}{\vspace{35.5mm}}  
\newcommand{\colthreevsep}{\vspace{14mm}}  
\newcommand{\colfourvsep}{\vspace{16mm}}  
\newcommand{\colfivevsep}{\vspace{23mm}}
```

Summary

- ▶ Use LaTeX/Beamer or Slidify to create reproducible slides.
- ▶ Use LaTeX/Beamer to create reproducible posters.
- ▶ Include KnitR code chunks to create figures directly.
- ▶ Or keep the code for figures separate.