

Data Analysis using the R Project for Statistical Computing

Daniela Ushizima

NERSC Analytics/Visualization and Math Groups
Lawrence Berkeley National Laboratory



Outline

I. R-programming

- Why to use R
- R in the scientific community
- Extensible
- Graphics
- Profiling

II. Exploratory data analysis

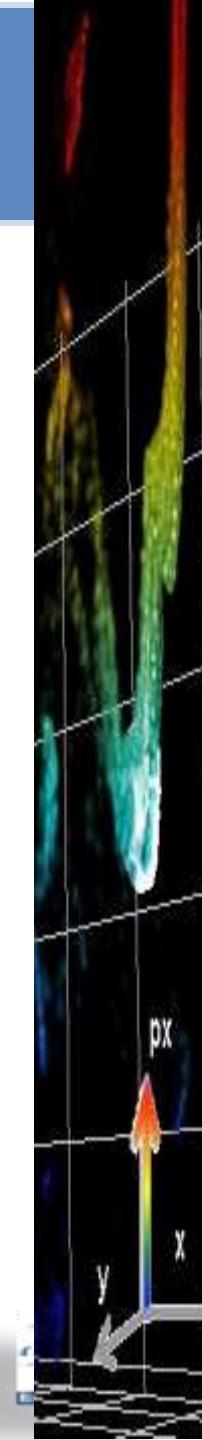
- Regression
- Clustering algorithms

III. Case study

- Accelerated laser-wakefield particles

IV. HPC

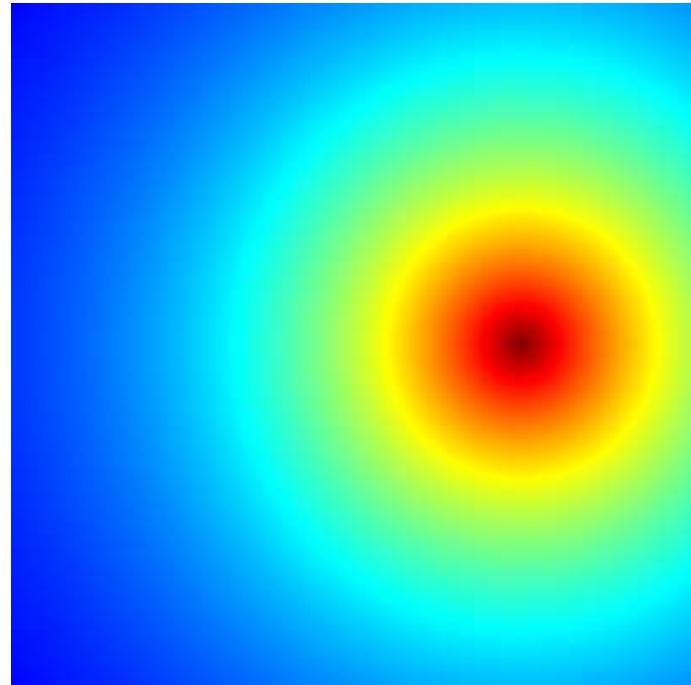
- State-of-the-art





Packages, data visualization and examples

R-PROGRAMMING



Search Technology

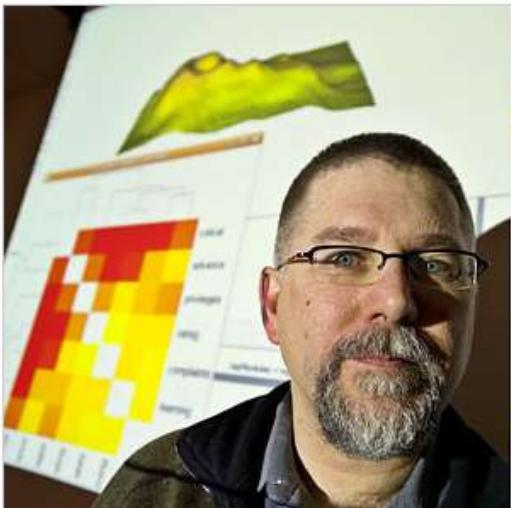
Inside Technology

Internet | Start-Ups | Business Computing

Go



Data Analysts Captivated by R's Power



Stuart Isett for The New York Times

R first appeared in 1996, when the statistics professors Robert Gentleman, left, and Ross Ihaka released the code as a free software package.

By ASHLEE VANCE

Published: January 6, 2009

To some people R is just the 18th letter of the alphabet. To others, it's the rating on racy movies, a measure of an attic's insulation or what pirates in movies say.

Related

[Bits: R You Ready for R?](#)[The R Project for Statistical Computing](#)

R is also the name of a popular programming language used by a growing number of data analysts inside corporations and academia. It is becoming their lingua franca partly

is a language and environment for statistical computing and graphics, a GNU project.

R provides a wide variety of statistical (linear and nonlinear modeling, classical statistical tests, time-series analysis, classification, clustering, ...) and graphical techniques, and is highly extensible.

Download:

<http://www.r-project.org>

Recommended tutorial:

http://cran.r-project.org/doc/contrib/Paradis-rdebuts_en.pdf

1.Why to use R?

- Open-source, multiplatform, extensible;
- Easy on users with familiarity with S/S+, Matlab, Python or IDL;
- Active and growing community:
 - Google, Pfizer, Merck, Bank of America, Boeing, the InterContinental Hotels Group and Shell.

2.R in the scientific community

- Google summer of code and projects using R-project to mine large datasets:

<http://www.r-project.org/SoC08/ideas.html>



- With Pfizer:
 - predict the safety of compounds, specifically carcinogenic side effects in potential drugs.
 - models eliminate the expensive and time-consuming process of studying a large number of potential compounds in the physical laboratory..."

<http://www.bio-medicine.org/medicine-news-1/Pfizer-Partners-with-REvolution-Computing-to-Improve-Medicine-Production-Pipeline-17917-2/>



2.1. You R with NERSC

- Get started with R:

> module load R

> R

>help()

>demo()

>help.start()

>source('your_function.R')

>library(*package_name*)



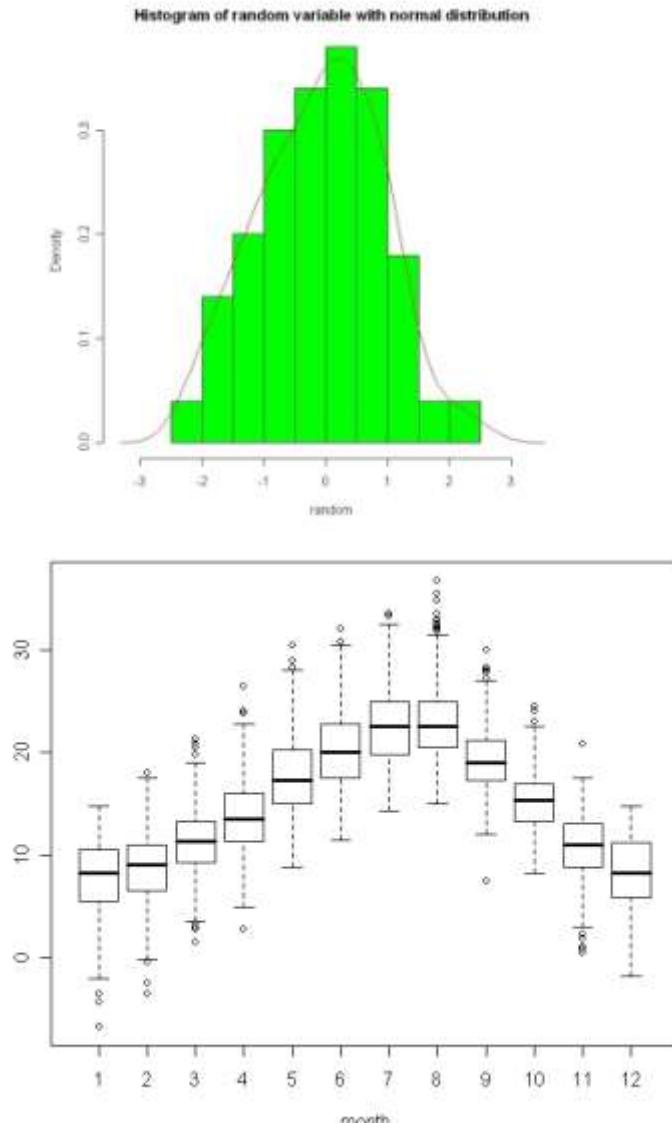
3. Extensible

- Add-on packages:
 - Data input/output: hdf5, Rnetcdf, DICOM, etc.
 - Graphics: trellis, gplot, RGL, fields, etc.
 - Multivariate analysis: MASS, mclust, ape, etc.
 - Other languages: Rcpp, Rpy, R.matlab, etc.



4. Statistical analysis and graphs

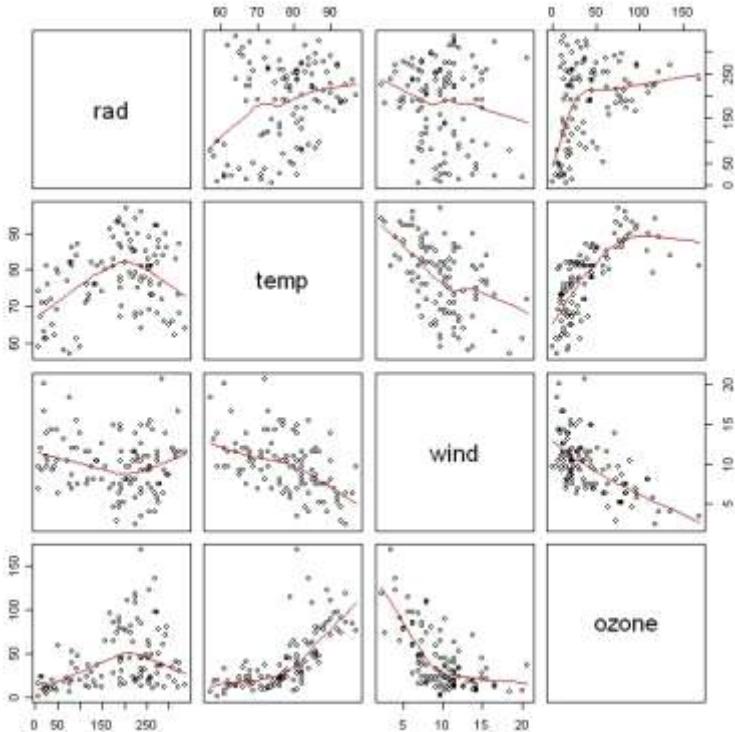
- Histogram
- Density
- Boxplot
- Multivariate plot
- Conditioning plot
- Contour plot



4.1. Multivariate plots

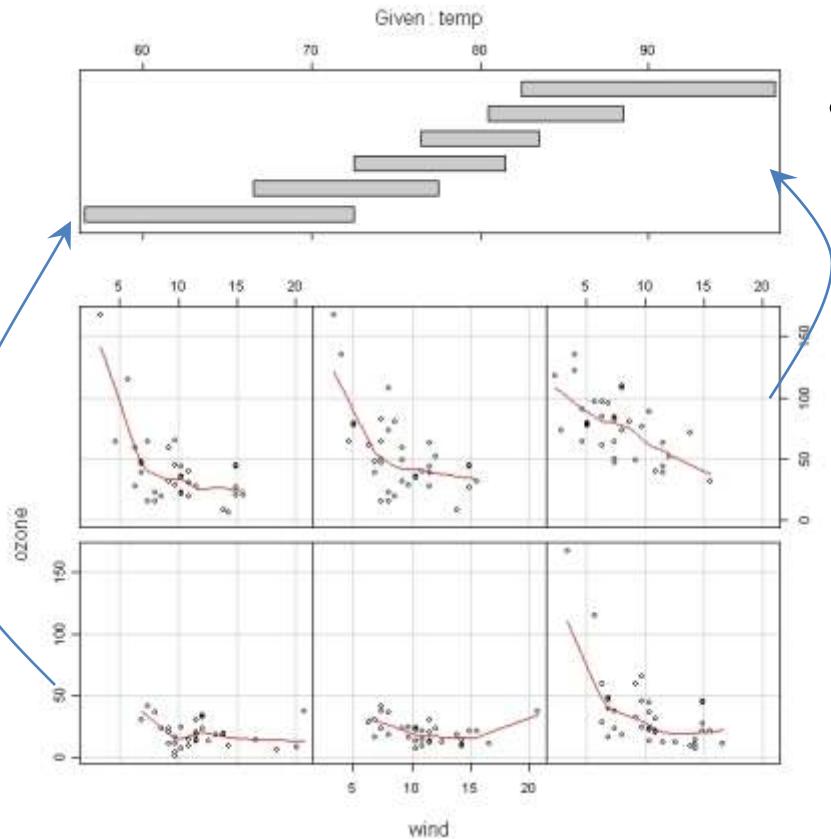
Ex: Explanatory variables: solar radiation, temperature, wind and the response variable ozone;

- use of pairs() with dataframes to check for dependencies between the variables.



```
> data=read.table('ozone.txt',  
+ header=T)  
> names(data)  
[1] "rad" "temp" "wind" "ozone"  
> pairs(data,panel.smooth)  
#panel.smooth = locally-weighted polynomial regression
```

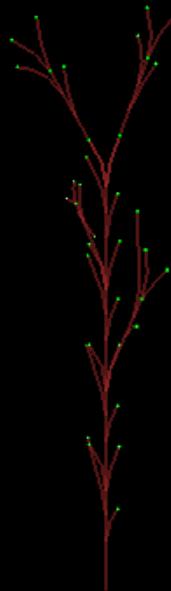
4.2. Conditional plots



- Check the relation of the two explanatory variables *wind*, *temp* and the response variable *ozone*:
 - Low temp: no influence of wind on levels of zone;
 - High temp: negative relationship between wind speed and ozone concentraton

```
>coplot(ozone~wind|temp,panel=panel.smooth)
```

4.3. Package RGL for 3D visualization



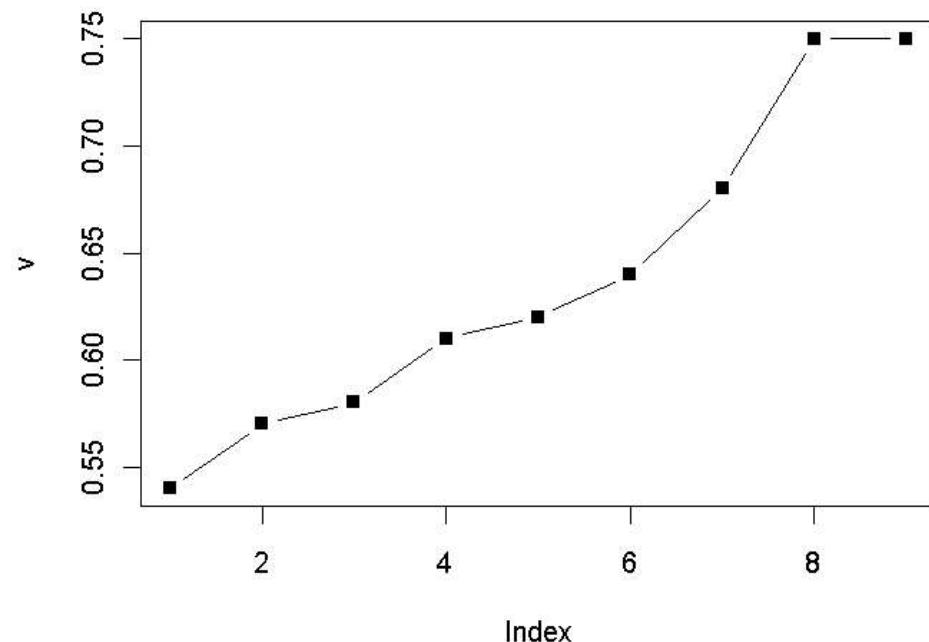
- OpenGL
- > `rgl.demo.lsystem()` - kernel density estimation

5. Profiling

Variable number of arguments

- Where does your program spend more time?

Matrix product computation time



Also try packages:
profr and *prof-tools*

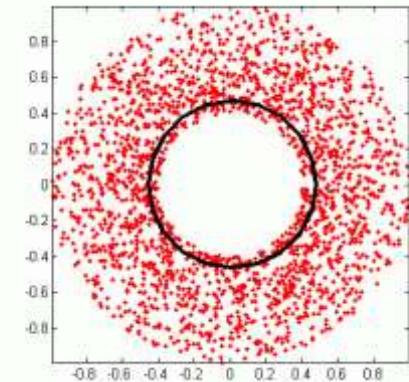
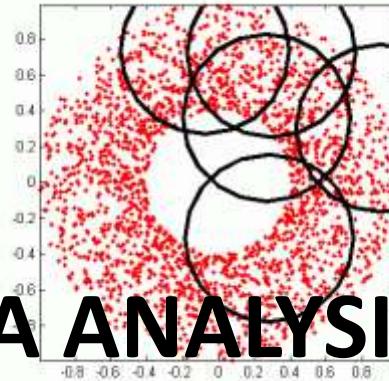
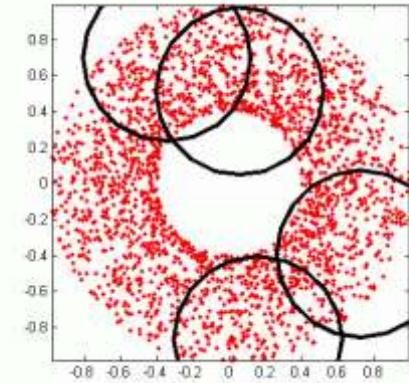
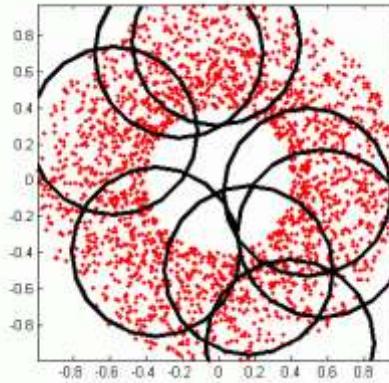
```
several.times <- function (n, f, ...) {
  for (i in 1:n) {
    f(...)
  }
}

matrix.multiplication <- function (s) {
  A <- matrix(1:(s*s), nr=s, nc=s)
  B <- matrix(1:(s*s), nr=s, nc=s)
  C <- A %*% B
}

v <- NULL
for (i in 2:10) {
  v <- append(
    v,
    system.time(
      several.times(
        10000,
      matrix.multiplication,
      i
    )
  )[1]
}
plot(v, type = 'b', pch = 15,
      main = "Matrix product computation time")
```

Basics and beyond

EXPLORATORY DATA ANALYSIS

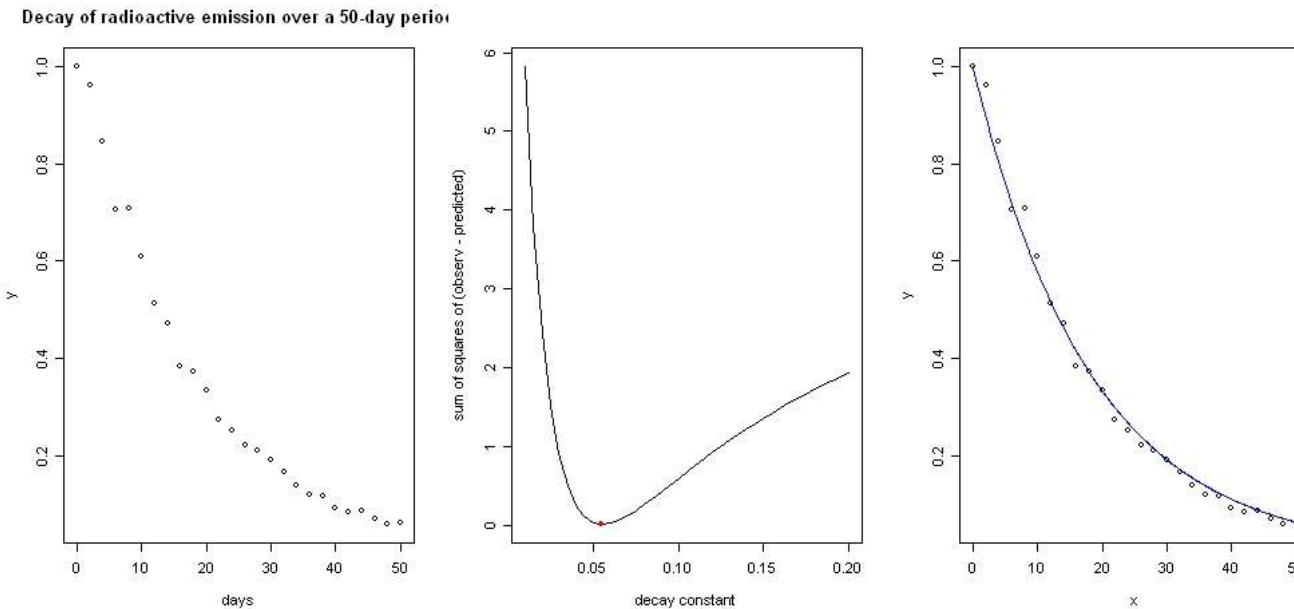


1. Statistical analysis

- Statistical modeling: check for variations in the response variable given explanatory variables;
 - Linear regression
- Multivariate statistics: look for structure in the data;
 - Clustering:
 - Hierarchical
 - Dendrograms
 - Partitioning
 - Kmeans (stats)
 - Mixture-models (mclust)

2. Linear regression

- Ex: Find the equation that best fit the data, given the decay of radioactive emission over a 50-day period



- Linear regression: variables expected to be linearly related;
- Maximum likelihood estimates of parameters = least squares;

2.1.Linear regression

```
data = read.table('sapdecay.txt',header=T)
attach(data)
# the log(y) gives a rough idea of the decay constant, a, for these data by linear regression of log(y) against x
mylm = lm(log(y)~x)
print(mylm$coefficients)
# sum of squares of the difference between the observed yv and predicted yp values of y, given a specific value of parameter a
sumsq <-function(a,xv=x,yv=y)
{
  yp = exp(-a*xv) #predicted model for y
  sum((yv-yp)^2)
}
a=seq(0.01,0.2,.005)
sq=sapply(a,sumsq)
decayK=a[min(sq)==sq] #this is the least-squares estimate for the decay constant
days=seq(0,50,0.1)

par(mfrow=c(1,3))
plot(x,y,main='Decay of radioactive emission over a 50-day period',xlab='days')
plot(a,sq,type='l',xlab='decay constant',ylab='sum of squares of (observ - predicted)')
matplot(decayK,min(sq),pch=19,col='red',add=T)
plot(x,y); lines(days,exp(-decayK*days),col='blue')
detach()
```

3. Cluster analysis

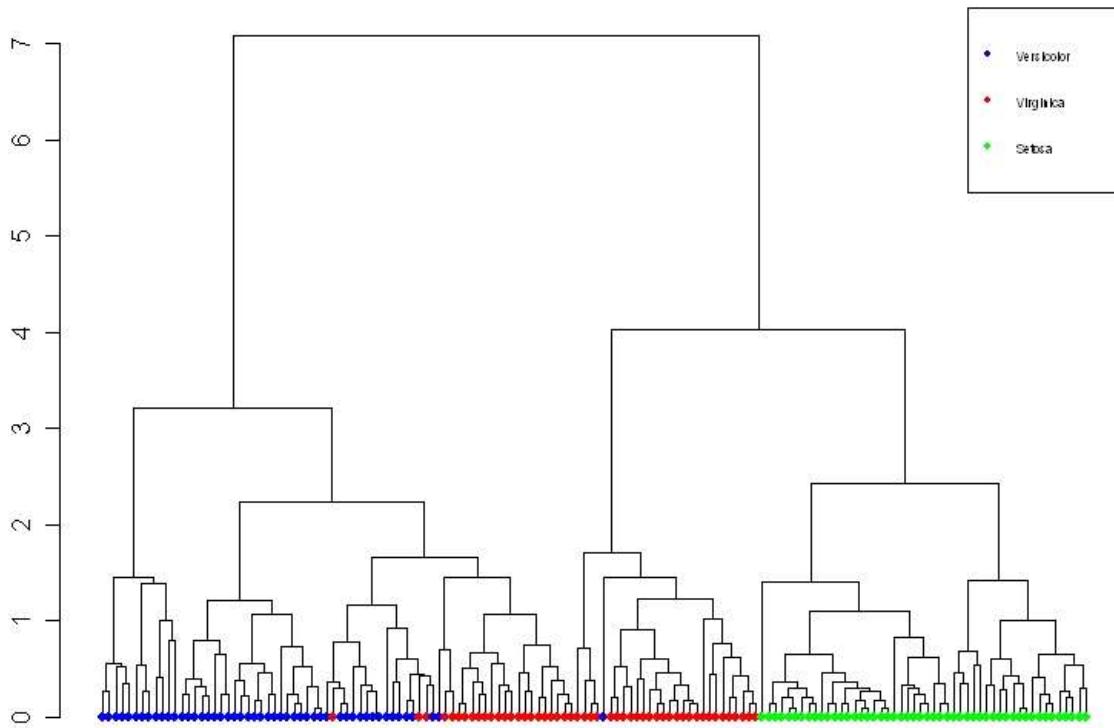
- Hierarchical
 - dendrogram(stats)
- Partitioning
 - kmeans (stats)
- Mixture-models:
 - Mclust (mclust)



Iris dataset: 150 samples of Iris flowers described in terms of its petal and sepal length and width

3.1.Hierarchical clustering

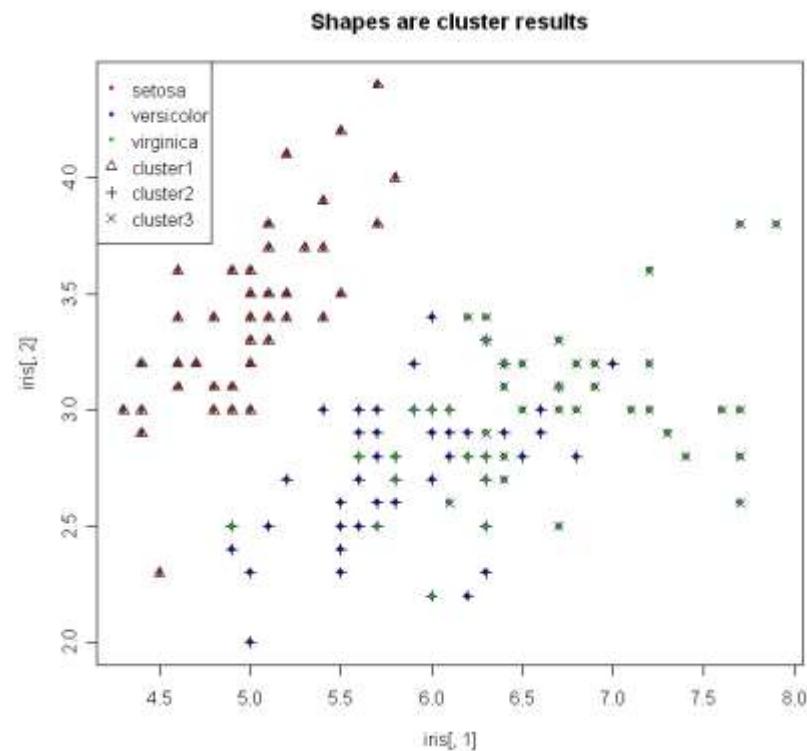
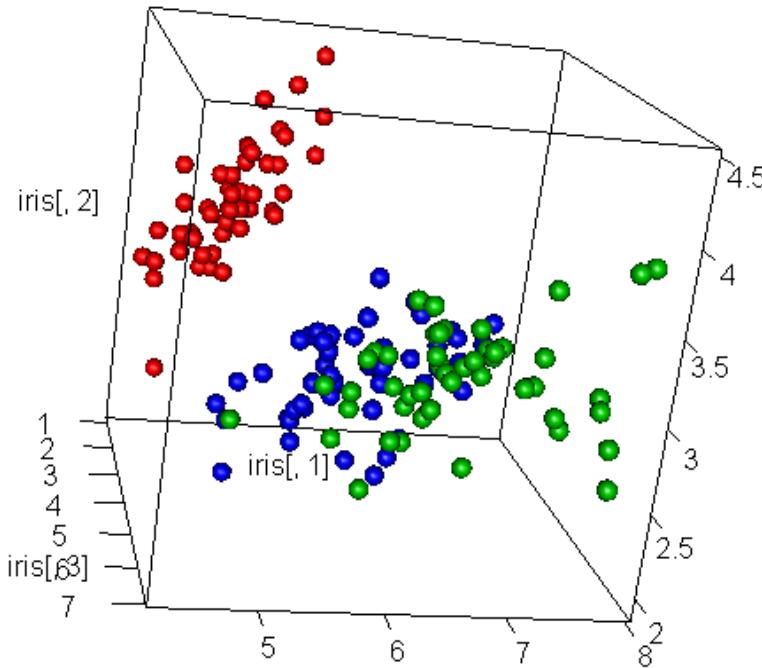
Iris database



- Analysis on a set of dissimilarities, combined to agglomeration methods for analyzing it:
- Dissimilarities: Euclidean, Manhattan, ...
- Methods:
 - ward, single, complete, average, mcquitty, median or centroid.



3.2.K-means



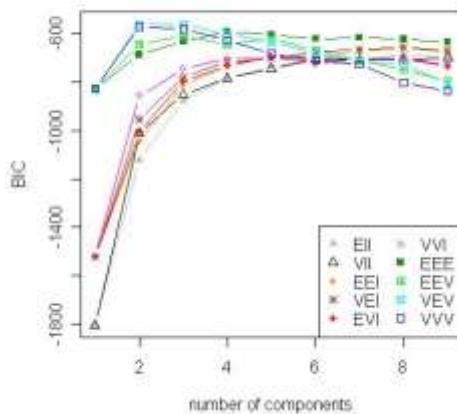
- Split n observations into k clusters;
 - each observation belongs to the cluster with the nearest mean.

setosa versicolor virginica

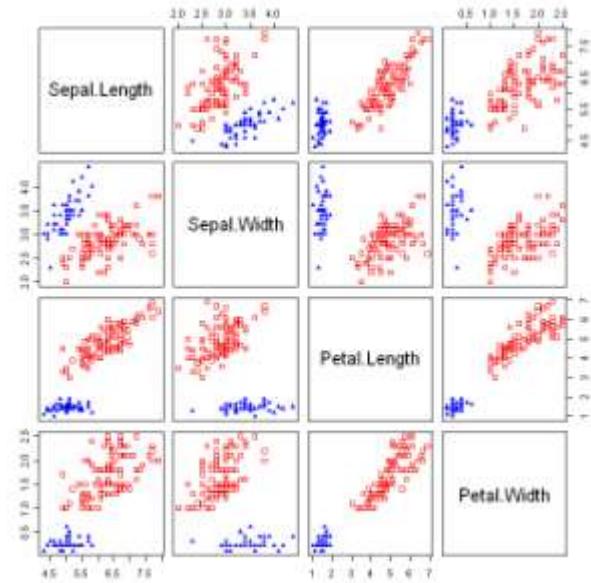
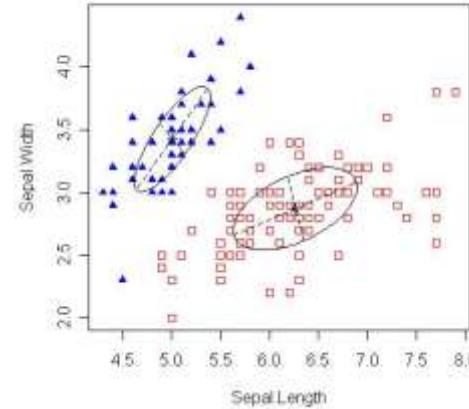
1	0	48	14
2	0	2	36
3	50	0	0



3.3. Model-based clustering

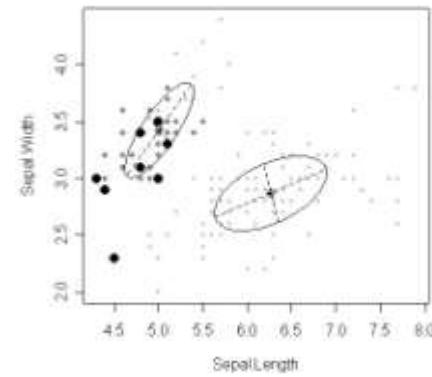


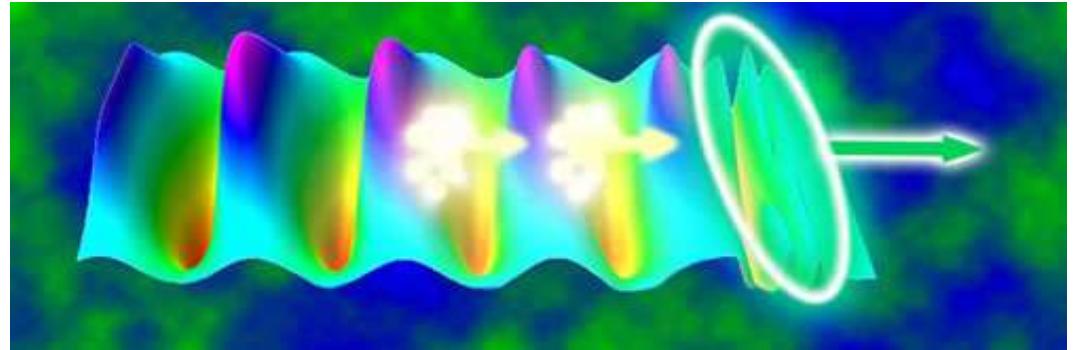
1.2 Coordinate Projection showing Classification



- Mixture Models
 - Each cluster is mathematically represented by a **parametric distribution**;
 - Set of k distributions is called a *mixture*, and the overall model is a *finite mixture model*;
 - Each probability distribution gives the probability of an instance being in a given cluster.

1.2 Coordinate Projection showing Uncertainty





<http://www.lbl.gov/publicinfo/newscenter/features/2008/apr/af-bella.html>

Accelerated laser-wakefield particles

Case study

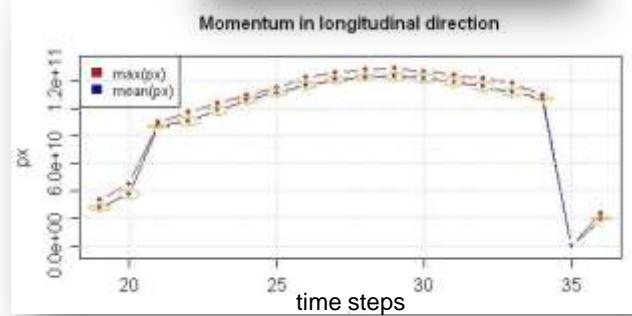
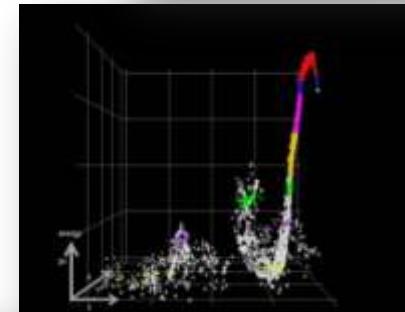
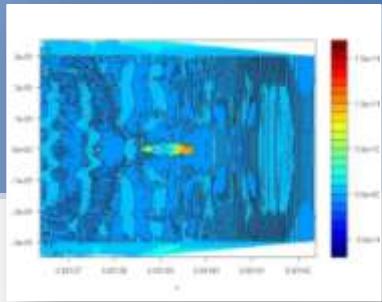
Knowledge discovery in LWFA science via machine learning

- C. Geddes (LBNL) in LOASIS program headed by W. Leemans and SciDAC COMPASS project.

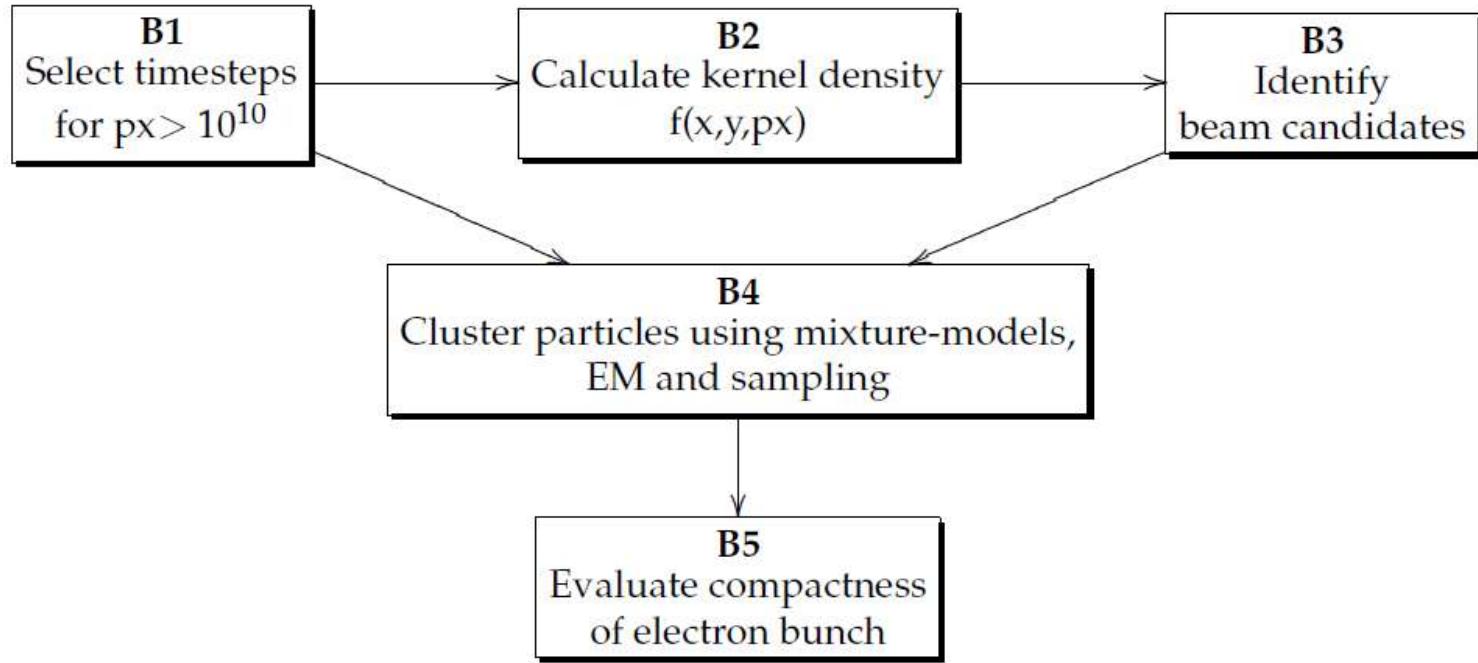
- Highlights:
 - Described compact electron clouds using minimum enclosing ellipsoids;
 - Developed algorithms to adapt mixture model clustering to large datasets;

- Science Impact:
 - Automated detection and analysis of compact electron clouds;
 - Derived dispersion features of electron clouds;
 - Extensible algorithms to other science problems;

- Collaborators:
 - Tech-X
 - Math Group, LBNL
 - UC Davis, University of Kaiserlautern

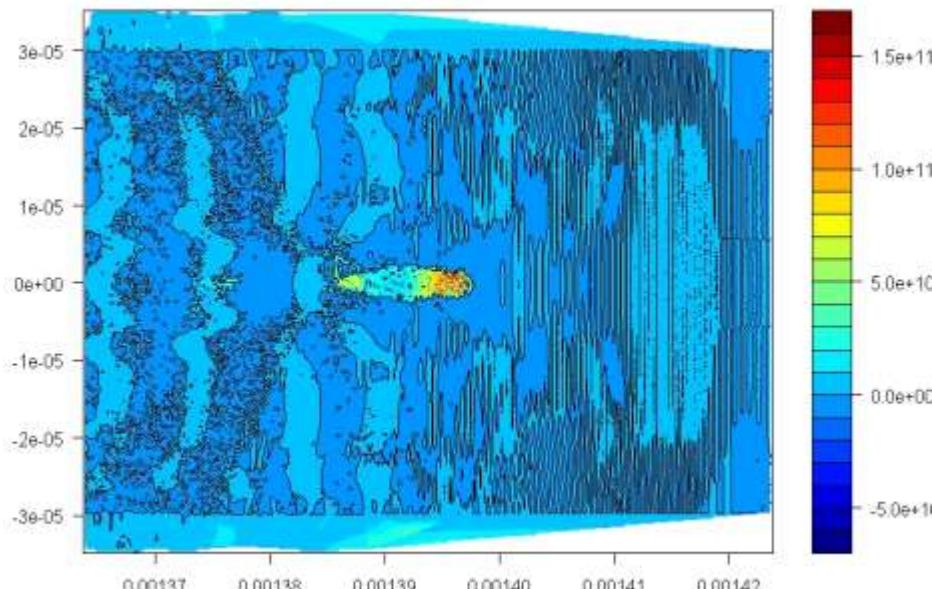


Framework



- Goal: automate the analysis of electron bunches by detecting compact groups of particles, with similar momentum and spatio-temporal coherence.

B1. Select relevant particles

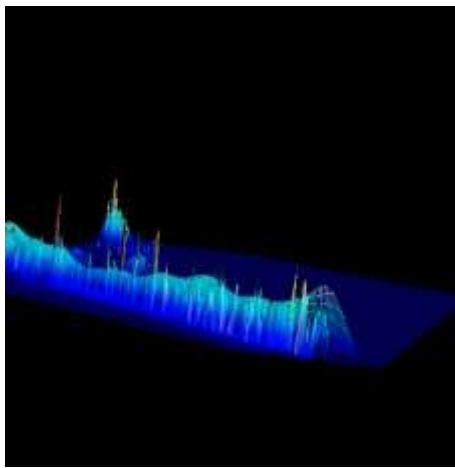


Representation of particle momentum in one time step: spline interpolation onto a grid for visualization of irregularly spaced input data.

- Beams of interest are characterized by high density of high-energy particles:
1. Elimination of low energy particles ($px < 1e10$)
 - Wake oscillation: $px \leq 1e9$
 - Excludes particles of the background
 2. Calculation of the simulation average number of particles (μ_s);
 3. Elimination of timesteps with number of particles inferior to μ_s ;

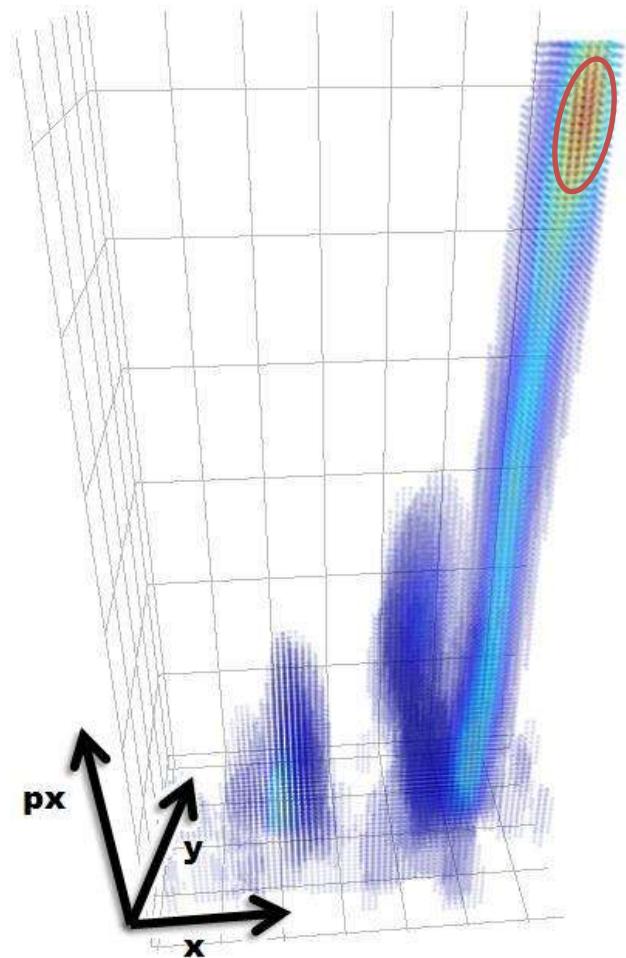
Packages:
akima, hdf5, fields

B2. Kernel-based estimation

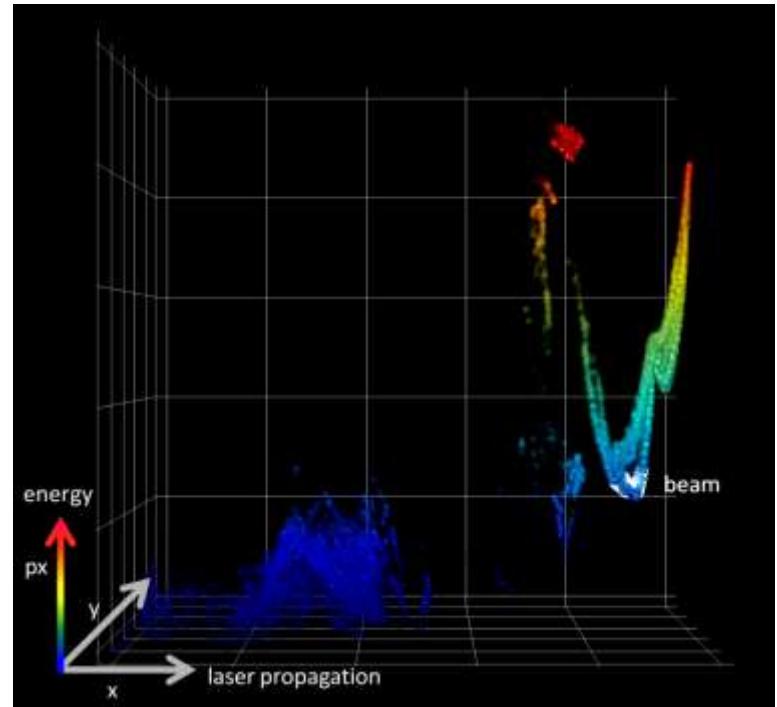
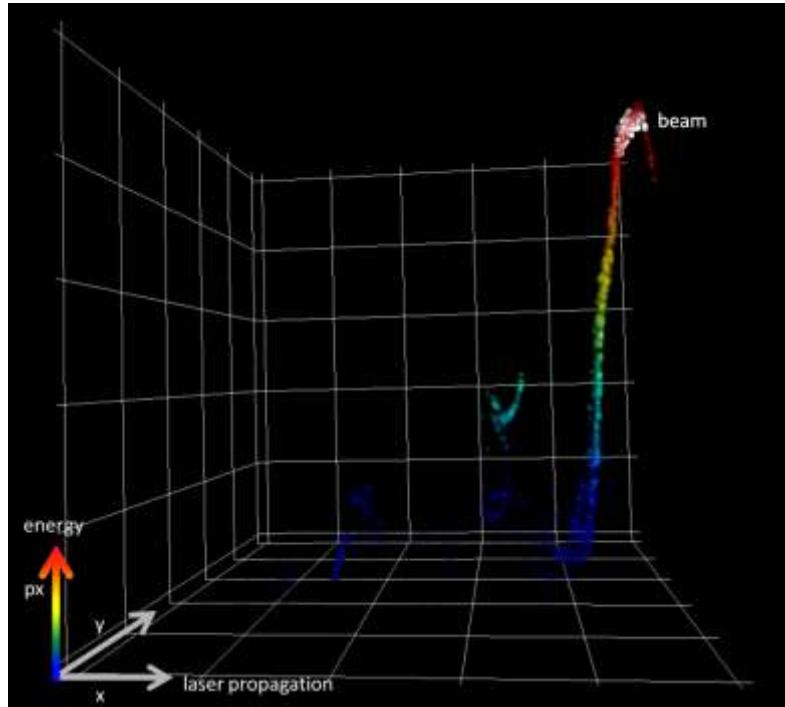


- Kernel density estimators are less sensitive to the placement of the bin edges;
- Goal: retrieve a dense group of particles with similar spatial and momentum characteristics:
 - argmax $f(x,y,px)$,
 - Neighborhood: 2 μm

Packages:
misc3d, rgl, fields

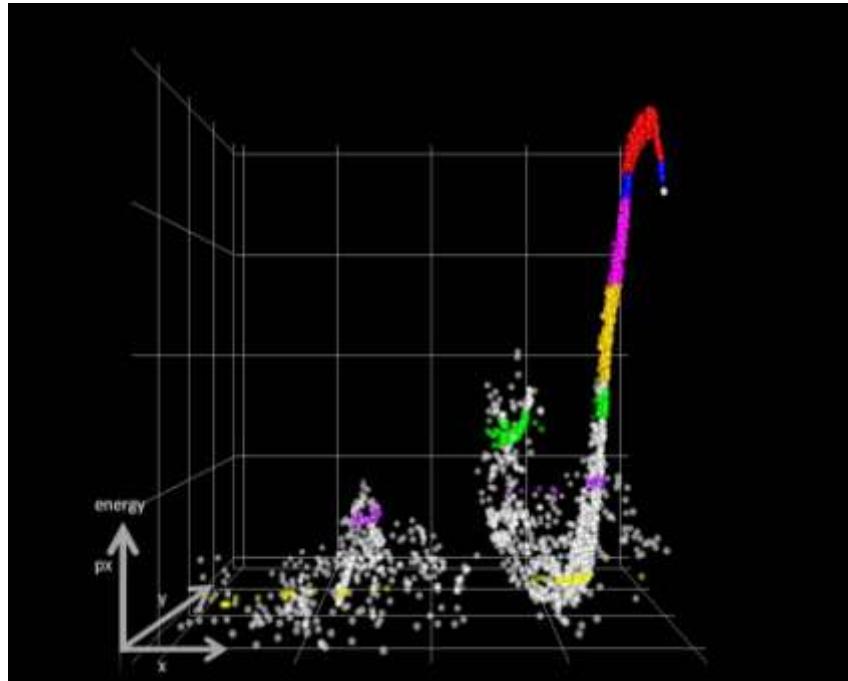


B3. Identify beam candidates



- Detection of compact groups of particles independent of being a maximum in one of the variables;

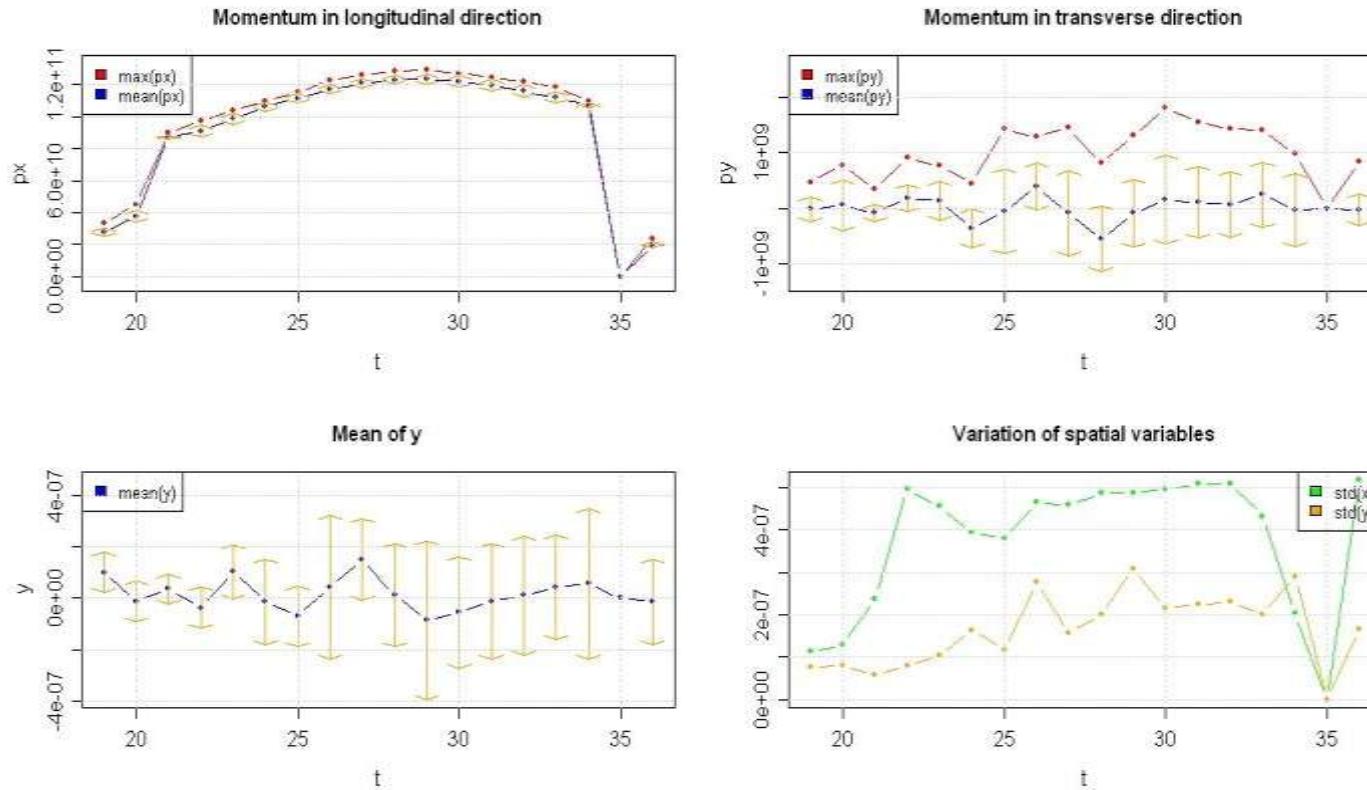
B4. Cluster using mixture models



- Model and number of clusters can be selected at run time (*mclust*);
- Partition of multidimensional space;
- Assume that the functional form of the underlying probability density follows a mixture of normal distributions;

Packages:
mclust, rgl

B5. Evaluation of compactness



- Bunches of interest move at speed $\approx c$, hence are nearly stationary in the moving simulation window;
- Moving averages smoothes out short-term fluctuations and highlights longer-term trends.

Packages, challenges and new businesses

High performance computing



1. Improve performance/reusability

- Good coding: avoid loops, vectorization;
- Extend R using C++ compiled code:
 - packages: Rcpp, inline
- Reuse your Python codes:
 - Package: Rpython
- Parallelism:
 - Explicit: packages Rmpi, Rpvm, nws
 - Implicit: packages pnmath, pnmath0 for multithreaded math functions
- Use out-of-memory processing with
 - packages bigmemory and ff

2. What is going on HPC in R?

- Parallelism:
 - Multicore: multicore, pnmath, ...
 - Computer cluster: snow, Rmpi, rpvm, ...
 - Grid computing: GRIDR, ...
- GPU:
 - gputools: parallel algorithms using CUDA + CULA
- Extremely large data:
 - ff: memory mapped pages of binary flat files.

3. Nothing is perfect...



- Limits on individual objects: on all versions of R, the maximum number of elements of a vector is $2^{31} - 1$;
- R will take all the RAM it can get (Linux only);
- More information, type:

```
>help('Memory-limits')
```

```
>gc() #garbage collector
```

```
>object.size(your_obj) #size of your object
```



Take home

- **Everything is an object.** This means that your variables are objects, but so are output from analyses. Everything that can possibly be an object by some stretch of the imagination... is an object.
- **R works in columns, not rows.** R thinks of variables first, and when you line them up as columns, then you have your dataset. Even though it seems fine in theory (we analyze variables, not rows), it becomes annoying when you have to jump through hoops to pull out specific rows of data with all variables.
- **R likes lists.** If you aren't sure how to give data to an R function, assume it will be something like this: `c("item 1", "item 2")` meaning "concatenate into a list the 2 objects named Item 1, Item 2". Also, "list" is different to R from "vector" and "matrix" and "dataframe" etc.
- **Its open source.** It won't work the way you want. It has far too many commands instead of an optimized core set. It has multiple ways to do things, none of them really complete. People on the mailing lists revel in their power over complexity, lack of patience, and complete inability to forgive a novice. We just have to get used to it, grit our teeth, and help them become better people.

References

- Michael J. Crawley. *Statistics: An Introduction using R*. Wiley, 2005. ISBN 0-470-02297-3.
 - data: <http://www.bio.ic.ac.uk/research/mjcraw/therbook/>
- Robert H. Shumway and David S. Stoffer. *Time Series Analysis and Its Applications With R Examples*. Springer, New York, 2006. ISBN 978-0-387-29317-2
- Basics
 - <http://cran.r-project.org/doc/contrib/Short-refcard.pdf>
 - <http://cran.r-project.org/doc/contrib/refcard.pdf>
 - http://cran.r-project.org/doc/contrib/Paradis-rdebut_en.pdf
 - [http://www.manning.com/kabacoff/Kabacoff MEAPCH1.pdf](http://www.manning.com/kabacoff/Kabacoff_MEAPCH1.pdf)
- Intermediate
 - <http://math.acadiau.ca/ACMMaC/Rmpi/basics.html>
 - User-lists

} Cheat sheets

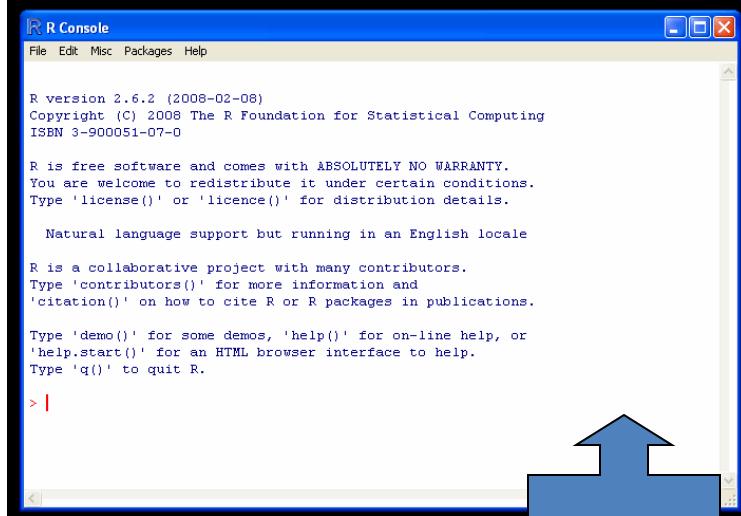
Basic but fundamental

EXTRA SLIDES

1. Install R

1. Download R from: <http://www.r-project.org>
2. Install the binary
3. Start R
4. Print one of the “cheat sheets”
5. Check the packages you may need
6. Customize by typing the cmd in your R session:
`install.packages('<name_pkg>')`

2. Getting started



R version 2.6.2 (2008-02-08)
 Copyright (C) 2008 The R Foundation for Statistical Computing
 ISBN 3-900051-07-0

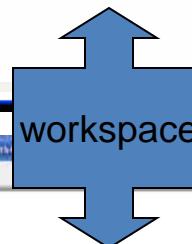
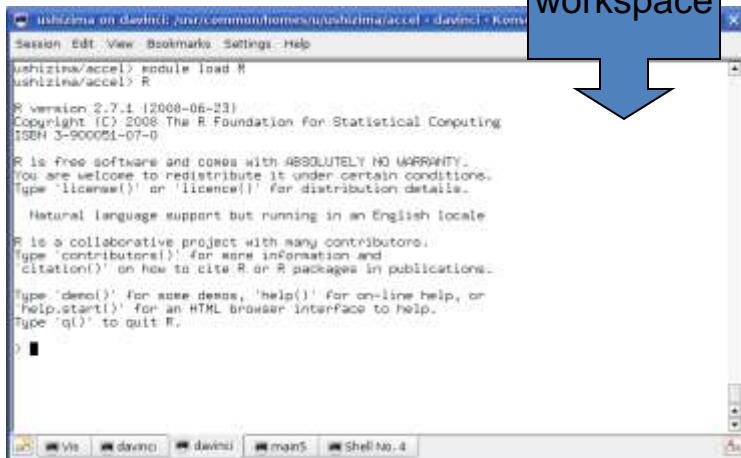
R is free software and comes with ABSOLUTELY NO WARRANTY.
 You are welcome to redistribute it under certain conditions.
 Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
 Type 'contributors()' for more information and
 'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
 'help.start()' for an HTML browser interface to help.
 Type 'q()' to quit R.

```
> |
```

ushizima@davinci:~\$ R
 Session Edit View Bookmarks Settings Help
 ushizima@accel1: module load R
 ushizima@accel1> R

R version 2.7.1 (2008-06-23)
 Copyright (C) 2008 The R Foundation for Statistical Computing
 ISBN 3-900051-07-0

R is free software and comes with ABSOLUTELY NO WARRANTY.
 You are welcome to redistribute it under certain conditions.
 Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
 Type 'contributors()' for more information and
 'citation()' on how to cite R or R packages in publications.

Type 'demo()', 'help()', 'for on-line help, or
 'help.start()' for an HTML browser interface to help.
 Type 'q()' to quit R.

```
> |
```

- 1) your question can be a valid package name or valid command:

```
> help(graphics) or ?plot
```
- 2) this will search anything that contain your query string:

```
> help.search('fourier')
```
- 3) which package contains the cmd?

```
> find("plot")
```
- 4) get working directory:

```
> getwd()
```
- 5) set working directory:

```
> setwd()
```
- 6) variables in your R-session:

```
> ls()
```
- 7) remove your variable:

```
> rm(mytrash_var)
```
- 8) List the objects which contain 'n'

```
> ls(pat='n')
```
- 9) Source a function:

```
> source('myfunction.R')
```
- 10) Load a library

```
> library(fields)
```

3. Simple syntax

- Basic types: numeric, character, complex or logical

```
> v1=c(7,33,1,7) #this is a vector  
> v2=1:4      #this is also a vector  
> v3=array(1,c(4,4,3)) #create a multidimensional array  
> i=complex(real=1,imag=3) #this is a complex number
```

- Functions:

```
> n=11; print(n); sqrt(n);  
> ifelse(n>11,n+1,n%%2)  
[1] 1
```

- Operators: + * / - ^ < <= > >= == !=

%/%, ^, %% , sqrt(): integer division, power, modulo, square root

```
>A%*%B #matrix multiplication
```

- Packages

```
> install.packages()  
> library(stats)
```

4.Type of objects representing data

Object	Modes	Allow mode heterogeneity ?
vector	numeric, character, complex or logical,	No
factor	numeric or character	No
array	numeric, character, complex or logical	No
matrix	numeric, character, complex or logical	No
data.frame	numeric, character, complex or logical	YES
ts	numeric, character, complex or logical	No
list	numeric, character, complex or logical, function, expression,	YES

1) Test type of object/mode: `is.type()`

2) Coerce: `as.type()`

Ex:

```
x=c(8,3,6,3)
is.character(x)
m=as.character(x)
```

R is a collaborative project with many contributors.
 Type 'contributors()' for more information and
 'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
 'help.start()' for an HTML browser interface to help.

Type 'q()' to quit R.

```
> v=array(1:10)
> v
[1] 1 2 3 4 5 6 7 8 9 10
> m=matrix(v,2,5)
> m
[,1] [,2] [,3] [,4] [,5]
[1,] 1 3 5 7 9
[2,] 2 4 6 8 10
> summary(m)
      V1        V2        V3        V4        V5 
Min. :1.00  Min. :3.00  Min. :5.00  Min. :7.00  Min. : 9.00 
1st Qu.:1.25 1st Qu.:3.25 1st Qu.:5.25 1st Qu.:7.25 1st Qu.: 9.25 
Median :1.50 Median :3.50 Median :5.50 Median :7.50 Median : 9.50 
Mean   :1.50 Mean   :3.50 Mean   :5.50 Mean   :7.50 Mean   : 9.50 
3rd Qu.:1.75 3rd Qu.:3.75 3rd Qu.:5.75 3rd Qu.:7.75 3rd Qu.: 9.75 
Max.   :2.00 Max.   :4.00 Max.   :6.00 Max.   :8.00 Max.   :10.00
> dim(m)
[1] 2 5
> mydata.frame=data.frame(x=v,label=letters[1:10])
> mydata.frame
  x label
1 1 a
2 2 b
3 3 c
4 4 d
5 5 e
6 6 f
7 7 g
8 8 h
9 9 i
10 10 j
> attributes(mydata.frame)
$names
[1] "x"      "label"

$row.names
[1] 1 2 3 4 5 6 7 8 9 10

$class
[1] "data.frame"
> |
```

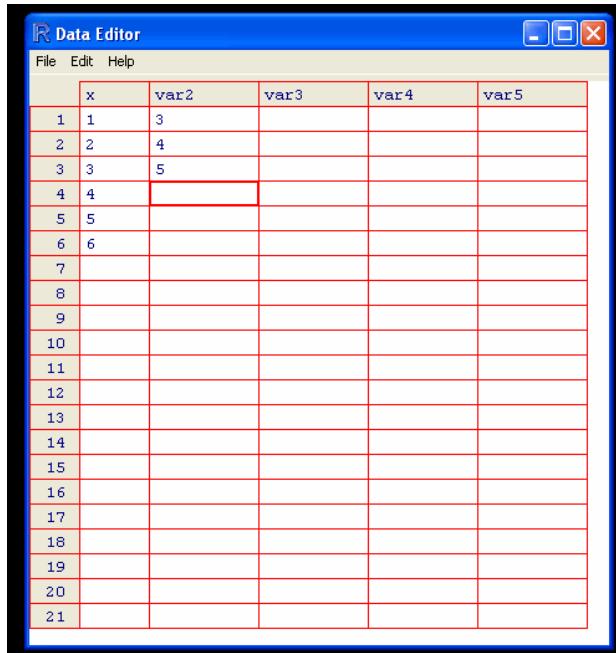
4.1.Creating objects

- Arrays;
- Matrices;
- Data frame: set of vectors of the same length;
- Factor: ‘category’, ‘enumerated type’

> **summary(.)**

> **attributes(.)**

4.2.Data input/output



- Graphical spreadsheet-like editor:

```
>data.edit(x) #open editor
```

```
>x=c(5,7,2,33,9,14)
```

```
>x=scan()
```

```
>data=read.table("data.txt",header=T)
```

- Ex output:

```
>write.table(d,"new_file.txt")
```

5. Functions

```
>myfun=function(x=1,y){  
+ z=x+y  
+ z}  
> myfun(2,3)  
[1] 5
```

- Several mathematical, statistical and graphical functions;
- The arguments can be: “data”, formulae, expressions, . . .
- Functions always need to be written with parentheses in order to be executed, even if there are no parameters;
 - Type the function without parentheses: R will display the content of the function.

5.1. Built-in functions

1. Basic functions

- `sin()`, `cos()`, `exp()`, `log()`, ...

2. Distributions

- `rnorm()`, `beta()`, `gamma()`, `binom()`, `cauchy()`, `mvrnorm()`, ...

3. Matrix algebra

- `sum()`, `diag()`, `var()`, `det()`, `ginv()`, `eigen()`, ...

4. Calculus

- Ex: `D()`, `integrate()`

5. Differential equation

- `Rk4()` #library(odesolve)

6. Manipulation of objects – step1

- Querying data:

```
> f=rep(2:4,each=3) #repeats each element of the 1st parameter 3X  
> which(f==3) #indexes of where f==3 holds  
[1] 4 5 6
```

- Related commands:

```
> seq(), unique(), sort(), rank(), order(), rev()
```

- NaN is not NA:

```
> 0/0  
[1] NaN  
  
> is.nan(0/0) #this is not a number  
[1] TRUE  
  
>names=c('mary','john',NA) # use of not available
```

different



6.1. Manipulation of objects – step2

- Faster operations: `apply()`, `lapply()`, `sapply()`, `tapply()`
 - `apply` = for applying functions to the rows or columns of matrices or dataframes

```
>apply(M,2,max) #max of col
```

– `lapply` = for lists

```
>lapply(list(x=1:10,y=1:30), max)
```

– `sapply` = for vectors

```
>sapply(m=sapply(rnorm(2000),(function(x){x^2}))
```

– `tapply` = for tables

```
> mylist <- list(c(1, 2, 2,1), c("A", "A", "B","C"))
```

```
>tapply(1:length(mylist), mylist)
```



6.2. Create your own function

- Use the command line or an editor to create a function:
- Editor:

```
>fact <- function (x){ ifelse (x>1,x*fat(x-1),1)}
```

- Save in a file name fact.R

```
>source('fact.R')
```

```
>fact(3)
```

- You can also save the history:

```
>savehistory('facthistory')
```

```
>loadhistory('facthistory')
```

```
>history(5) #see last 5 commands
```



- *Tip: save filename = function name*

7. Graphics

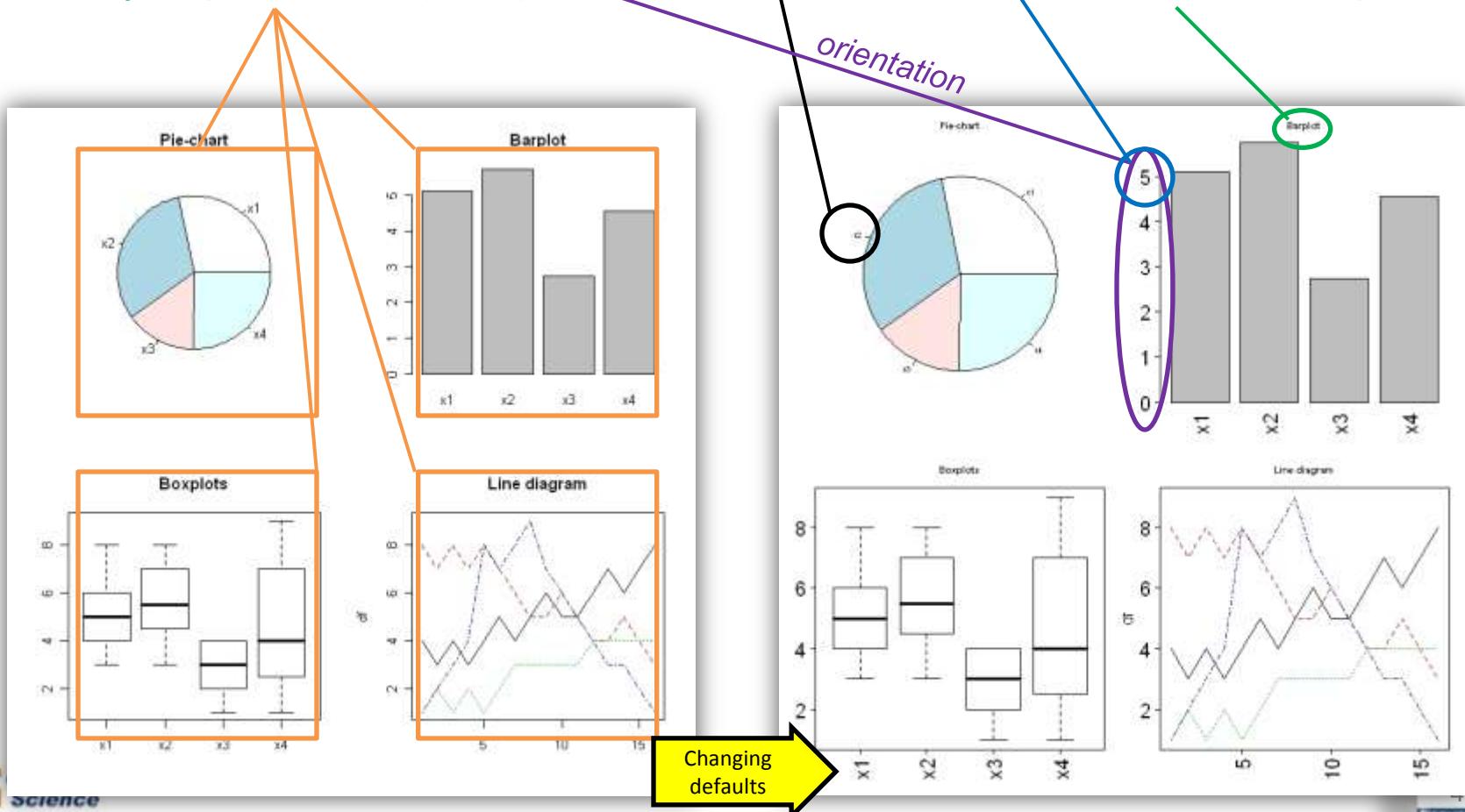
- Get a sense of what R can do: `>demo(graphics)`
- The graphical windows are called X11 under Unix/Linux and windows under Windows
- Other graphical devices: pdf, ps, jpg, png

```
>x11()  
  
>windows()  
  
>png()  
  
>pdf()
```

7.1. Plot structure

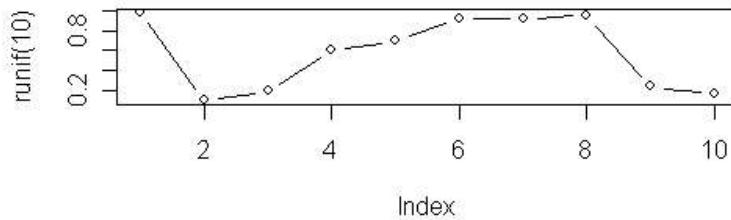
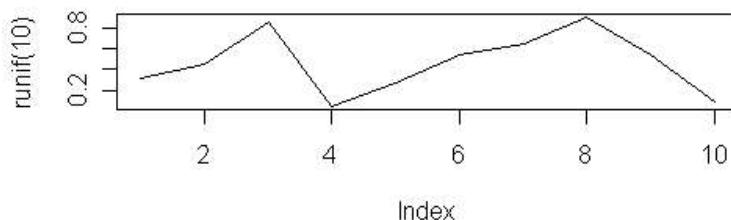
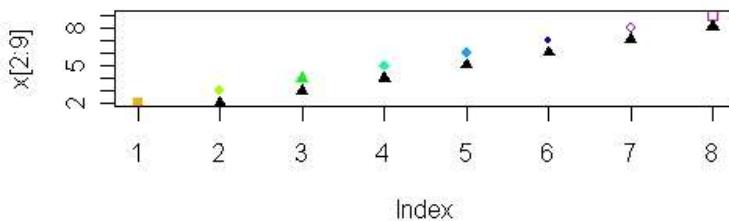
- Graph parameter function:

```
> par(mfrow=c(2, 2), las=2, cex=.5, cex.axis=2.5, cex.lab=2)
```



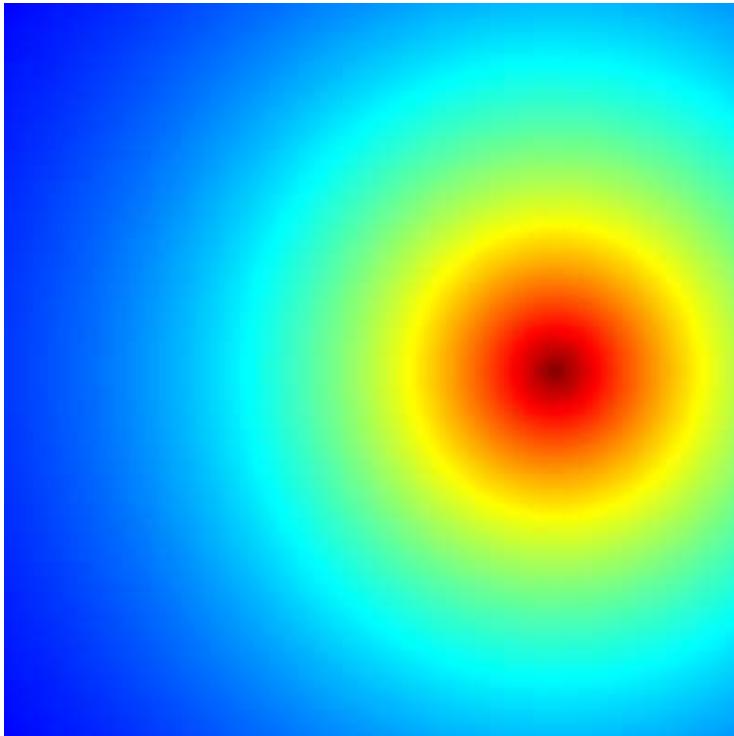
7.2.Example of data plots

Using layout cmd



```
> x=1:10
> layout(matrix(1:3, 3, 1))
> par(cex=1)
> plot(x[2:9],col=rainbow(9-2+2)[2:9], pch=15:(15+9-2+1)) #add to the plot
> matplot(x,pch=17,add=T);
> title('Using layout cmd')
> plot(runif(10),type='l')
> plot(runif(10),type='b')
```

7.3.Creating a gif animation



```
library(fields) # for tim.colors
library(caTools) # for write.gif
m <- 400 # grid size
C <- complex(real=rep(seq(-1.8,0.6, length.out=m),
  each=m), imag=rep(seq(-1.2,1.2, length.out=m),
  m))
C <- matrix(C, m, m)
Z <- 0
X <- array(0, c(m, m, 20))
for (k in 1:20)
{
  Z <- Z^2+C
  X[,,k] <- exp(-abs(Z))
}
col <- tim.colors(256)
col[1] <- "transparent"
write.gif(X, "rplot-mandelbrot.gif", col=col, delay=100)
image(X[,,k], col=col) # show final image in R
```

8.Data input/output – special formats

- Availability of several libraries. Ex:
 - Rnetcdf: netcdf functions and access to udunits calendar conversions;
 - DICOM: provides functions to import and manipulate medical imaging data via the Digital Imaging and Communications in Medicine (DICOM) Standard.
 - hdf5: interface to NCSA HDF5 library; read and write (only) the entire data