# Graphical Data and Data Graphics

Paul Murrell
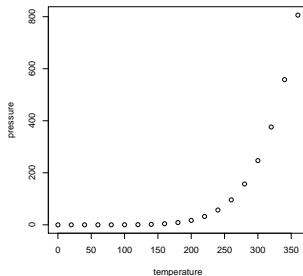
The University of Auckland

July 12 2007

# Graphical Statistics

```
> pressure
   temperature pressure
1            0   0.0002
2           20   0.0012
3           40   0.0060
4           60   0.0300
5           80   0.0900
6          100   0.2700
7          120   0.7500
8          140   1.8500
...
```
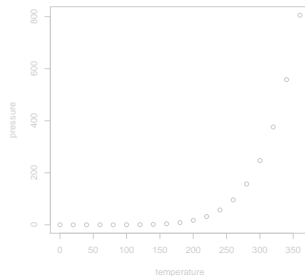
$\longrightarrow$

# Statistical Graphics

```
> pressure
    temperature  pressure
1            0    0.0002
2           20    0.0012
3           40    0.0060
4           60    0.0300
5           80    0.0900
6          100    0.2700
7          120    0.7500
8          140    1.8500
...
```
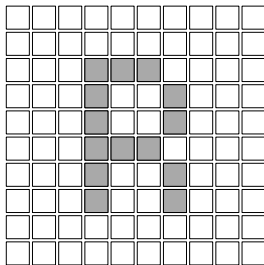
$\longrightarrow$

# Graphical Data and Data Graphics

- Graphical Statistics: *data* → *plot*
- Statistical Graphics: *data* → *plot*

- Graphical Data: *plot* → *data*
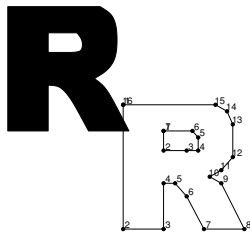- Data Graphics: *plot* → *data*

# Graphical Formats
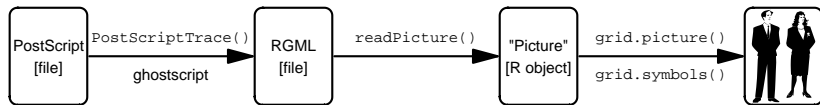
Raster



**pixmap** package
**EBimage** package

Vector



**grImport** package

# The grImport Package

# The PostScript Bezier Tiger

```
%!PS-Adobe-2.0 EPSF-1.2
%%Creator: Adobe Illustrator(TM)
%%For: OpenWindows Version 2
%%Title: tiger.eps
...
.8 setgray
clippath fill
-110 -300 translate
1.1 dup scale

0 g
0 G
0 i
0 J
0 j
0.172 w
10 M
[]0 d
0 0 0 0 k
...
```

# Converting the Tiger to Data

```
PostScriptTrace("tiger.ps")

tiger <-
  readPicture("tiger.ps.xml")
```

# Using the Tiger in a Plot

`grid.picture(tiger)`

# A Chess Board

```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE svg PUBLIC "-//W3C//DTD SVG"
"http://www.w3.org/TR/2001/REC-SVG...">
<!-- Created with Sodipodi -->
<svg version="1.0">
...
  <g
     style="font-size:12;"
     id="g874">
    <path
       d="M 0 437 L 437 0 "
       style="fill:none;fill-opacity:1"
       id="path616" />
...
```

```
# Convert SVG to PostScript
# using InkScape

PostScriptTrace("chess.ps")

chess <-
  readPicture("chess.ps.xml")
```

# The Paths in the Chess Board

```
picturePaths(chess[125:136])
```

## A Chess Piece as a Plotting Symbols

The number of moves required to complete chess games for
different opening gambits. From the career of Louis Charles Mahe
De La Bourdonnais (circa 1830).

```
grid.symbols(
  chess[205:206],
  x=games$num.moves,
  y=1:ngames,
  "native",
  size=unit(0.5, "cm"))
```

# Statistical Data Graphics

- Graphical Statistics: *data* → *plot*
- Statistical Graphics: *data* → *plot*

- Graphical Data: *plot* → *data*
- Data Graphics: *plot* → *data*

- Statistical Data Graphics: *data* → *plot* → *data*

# A Published Plot

```
# Extract just page 2
# and convert to PostScript

PostScriptTrace("Fig1.ps")

Fig1 <-
  readPicture("Fig1.ps.xml")

grid.picture(Fig1)
```

picturePaths(Fig1)

`grid.picture(Fig1[4:48])`

```
> barePlot <- Fig1[seq(4, 38, 2)]
```

```
> grid.picture(barePlot)
```

```
> slotNames(barePlot)

[1] "paths"    "summary"

> barePlot@summary

An object of class "PictureSummary"
Slot "numPaths":
[1] 18

Slot "xscale":
[1] 2563 5046

Slot "yscale":
[1] 6108 7371
```

```
> class(barePlot@paths)

[1] "list"

> barePlot@paths[[1]]

An object of class "PictureFill"
Slot "x":
move line line line line
2563 5046 5046 2563 2563

Slot "y":
move line line line line
6109 6109 7371 7371 6109

Slot "rgb":
[1] "#E6E6E6"

Slot "lwd":
[1] 1.33
```

```
> scaledMax <- function(x, summary) {
      (max(x@y) - summary@yscale[1]) /
      diff(range(summary@yscale))
  }

> barProportions <- sapply(barePlot@paths[-1],
                           scaledMax,
                           barePlot@summary)

> barProportions * 45

 [1] 26.8 28.8 29.1 29.6 30.5 31.9 32.3 34.3 34.6 35.1 35.1
[12] 35.4 35.5 35.9 36.2 36.4 39.2
```

picturePaths(Fig1)

```
> grid.picture(Fig1[39:41])
```

```
> errorBars <- explodePaths(Fig1[39:41])
> grid.picture(errorBars)
```

```
> picturePaths(errorBars)
```

```r
> topBars <- errorBars[seq(3, 35, 2)]
> bottomBars <- errorBars[seq(37, 69, 2)]
> scaledMin <- function(x, summary) {
      (min(x@y) - summary@yscale[1]) /
       diff(range(summary@yscale))
  }
> barMaxProp <- sapply(topBars@paths,
                       scaledMax,
                       barePlot@summary)
> barMinProp <- sapply(bottomBars@paths,
                       scaledMin,
                       barePlot@summary)
```
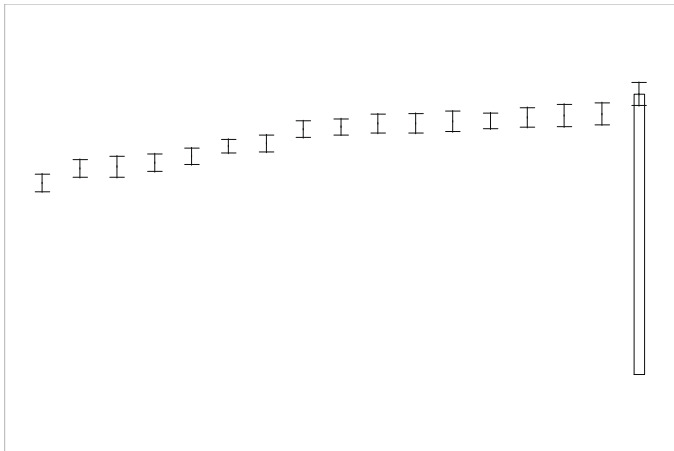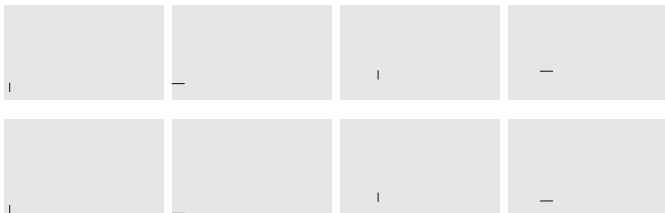
```
> barMaxProp * 45

 [1] 28.0 30.0 30.5 30.8 31.6 32.8 33.4 35.4 35.7 36.3 36.4
[12] 36.8 36.5 37.2 37.7 37.9 40.8

> barMinProp * 45

 [1] 25.5 27.5 27.5 28.4 29.3 30.9 31.1 33.1 33.4 33.7 33.7
[12] 33.9 34.3 34.5 34.6 34.8 37.6
```
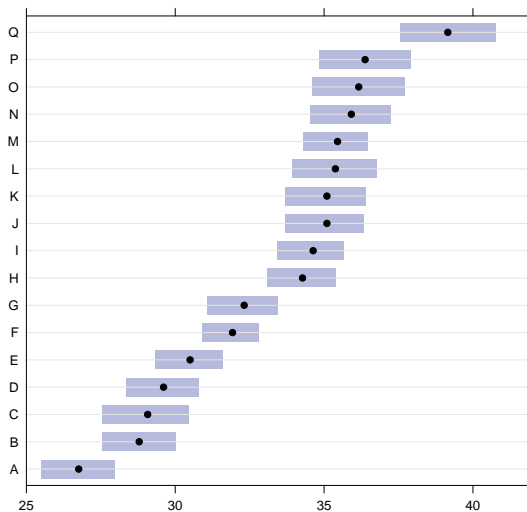
# Graphical Data Graphical Statistics

- Graphical Statistics: *data* → *plot*
- Statistical Graphics: *data* → *plot*

- Graphical Data: *plot* → *data*
- Data Graphics: *plot* → *data*

- Statistical Data Graphics: *data* → *plot* → *data*

- Graphical Data Graphical Statistics:
  *data* → *plot* → *data* → *plot*

dotplot(LETTERS[1:17] ~ barProportions*45)

# Acknowledgements