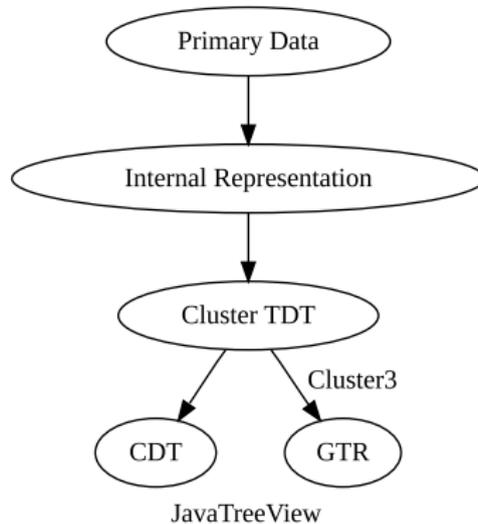


Practical Bioinformatics

Mark Voorhies

4/25/2011

Clustering protocol



How can we discover how we should format a file?

Whiteboard Image

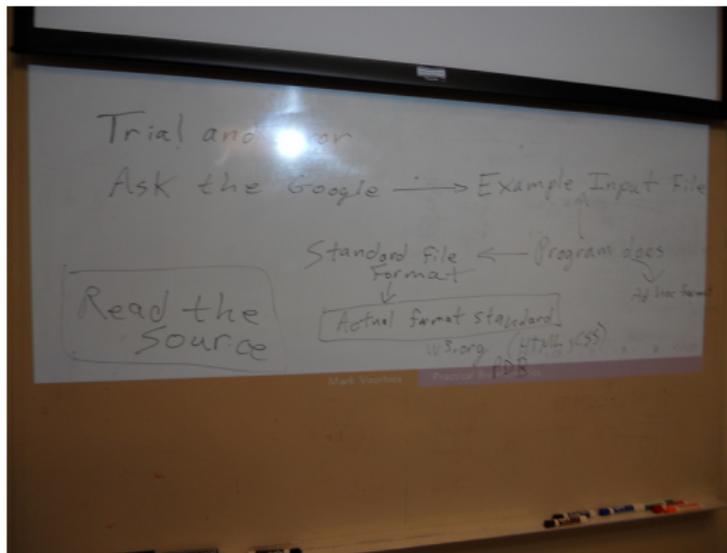
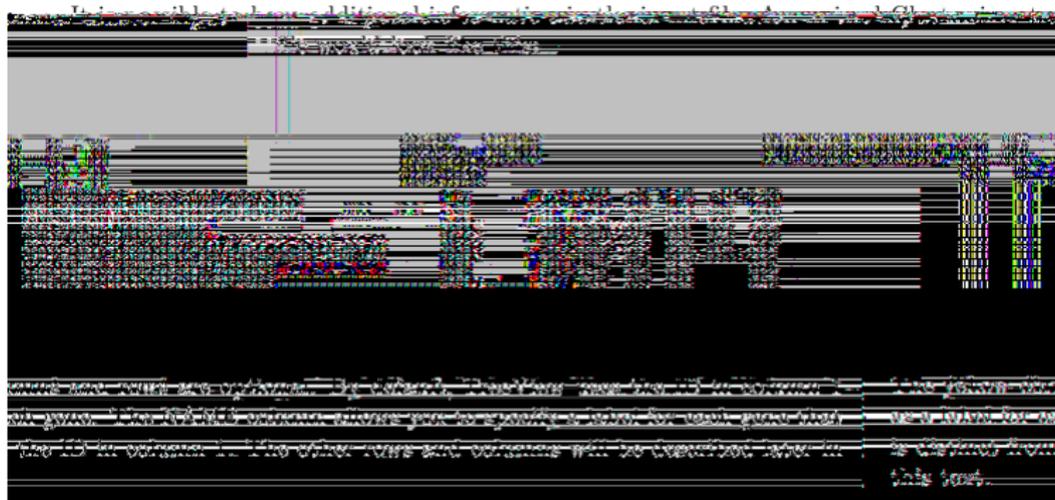


Figure 2.16. Screenshot of Generalized CDT file in Excel

	A	B	C	D	E	F	G	H	I	J
1	GID	YORF	NAME	PlateID	Spot	FGCOLOR	Sector	GWEIGHT	yB4n059 5 h	yB4n059 7 h
2	FGCOLOR					#000000			#999900	#999905
3	EWEIGHT					#000000		1	1	1
4	GENE5998X	YML051W	1709 GAL80	25570	7590	#000000	40	1	0.07	-0.27
5	GENE6516X	YNL064C	2407 YDJ1	25571	8170	#000000	43	1	0.94	0.6
6	GENE6229X	YBR035C	3774 PDX3	25558	7853	#000000	41	1	0.38	-0.54
7	GENE6073X	YCL145W	5500 TIR06	25564	7632	#000000	40	1	0.95	0



How does Cluster3 handle newlines?

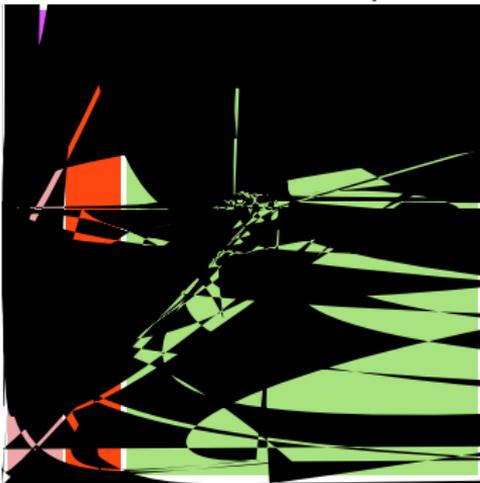
```
static char* GetLine(FILE* inputfile)
/* The function GetLine reads one line from the inputfile , and returns it as a
 * null-terminated string. If inputfile is at EOF, a null pointer is returned.
 * The calling routine should free the char* returned by GetLine.
 */
{ int c;
  int n = 0;
  int size = 1023;
  char* line = malloc((size+1)*sizeof(char));
  while ((c = getc(inputfile))!=EOF && c!='\r' && c!='\n')
  { if (n == size)
    { size *= 2;
      line = realloc(line ,(size+1)*sizeof(char));
    }
    line[n] = (char)c;
    n++;
  }
  if (c=='\r')
  { c = getc(inputfile);
    if (c!='\n' && c!=EOF) ungetc(c, inputfile);
  }
  if (n==0 && c==EOF)
  { free(line);
    return 0;
  }
  line[n] = '\0';
  line = realloc(line ,(n+1)*sizeof(char));
  return line;
}
```

How does Cluster3 handle delimiters?

```
static char* tokenize(char* s)
{
    char* p = s;
    while (1)
    {
        if (*p=='\0') return NULL;
        if (*p=='\t')
        {
            *p = '\0';
            return p+1;
        }
        p++;
    }
    /* Never get here */
    return NULL;
}
```

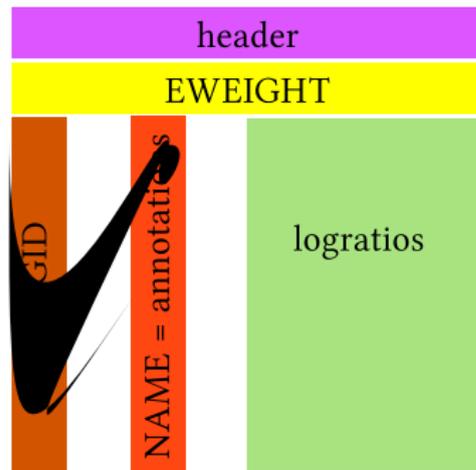
The CDT file format

Minimal CLUSTER input



- Tab delimited (`\t`)
- UNIX newlines (`\n`)
- Missing values \rightarrow empty cells

Cluster3 CDT output



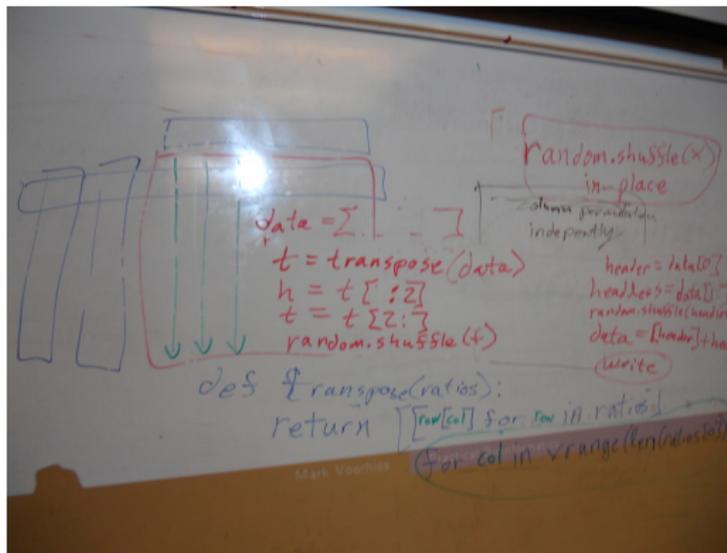
Be liberal in what you accept, strict in what you emit

```
def writerow(out, row):
    out.write(
        # Use tab delimiters
        "\t".join(
            # Normalize all whitespace (including \t, \r, and \n) and "
            # to "normal" spaces.
            [re.sub(r'[\s]', " ", str(i)) for i in row])
        # Use UNIX-style newlines
        +"\n")

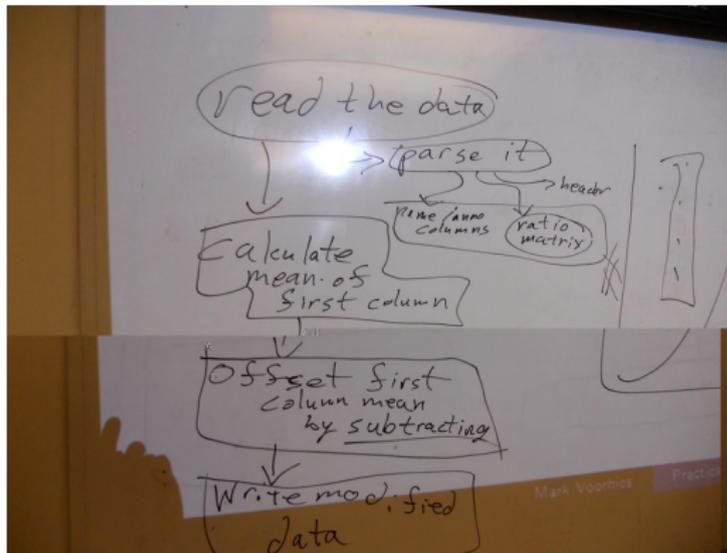
def writectd(filename, matrix,
             row_names, row_annotations, columns):
    out = open(filename, "w")
    writerow(out, ["UNIQID", "NAME"] + columns)
    for (uniqid, name, row) in zip(row_names, row_annotations, matrix):
        writerow(out, [uniqid, name] + row)
    out.close()
```

Strategies for Shuffling

Whiteboard Image

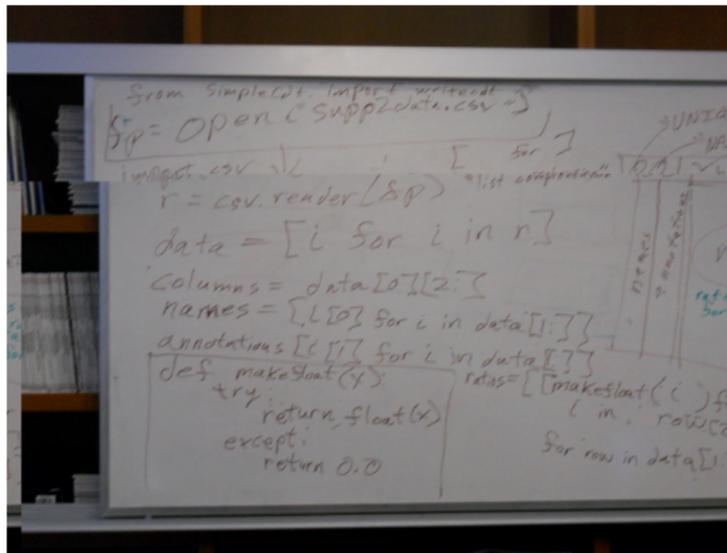


Whiteboard Image



Whiteboard Image

```
from SimpleCsv import writecsv
f = open('Supp2data.csv')
import csv
r = csv.reader(f)
data = [i for i in r]
columns = data[0][1:]
names = [L[0] for L in data[1:]]
annotations = [L[1] for L in data[1:]]
def makefloat(x):
    try:
        return float(x)
    except:
        return 0.0
ratas = [makefloat(i) for i in row[2:] for row in data[1:]]
```



The whiteboard contains handwritten Python code for reading a CSV file. The code defines a list 'data' from a CSV reader, extracts 'columns' and 'names' from the first row, and uses list comprehensions to extract 'annotations' and 'ratas'. A 'makefloat' function is defined to handle non-numeric values. To the right, a diagram shows a CSV file structure with columns labeled 'UNIA', 'NPA', and 'V'. The first row is 'UNIA NPA V' and subsequent rows contain numerical values. A 'ratas' list is shown as a list of floats derived from the 'V' column.

Whiteboard Image

import writecsv
supp2data.csv

header(['SP'])
for i in n]

data[0][2:]
[0] for i in data[1:]]
[1] for i in data[1:]]
t(x):
n_float(x)
0.0

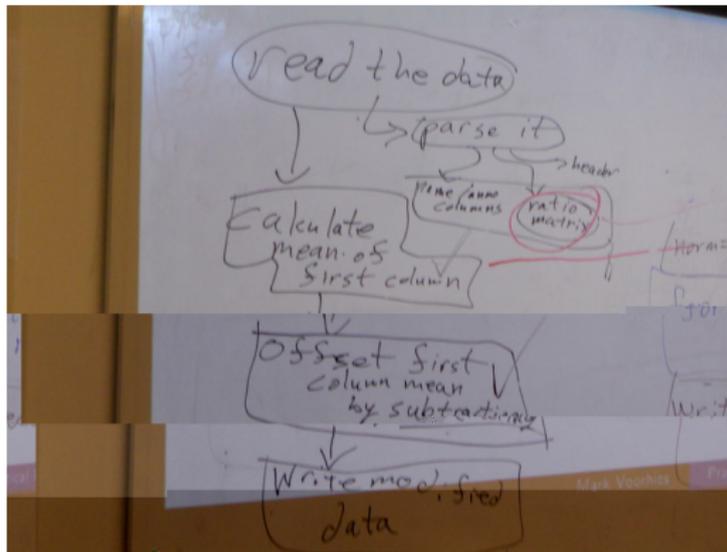
UNIQID
NAME
headers
ratios

ratios = []
for row in data[1:]:
ratios.append(
for i in row[2:]:
ratios.append(
makefloat(i))

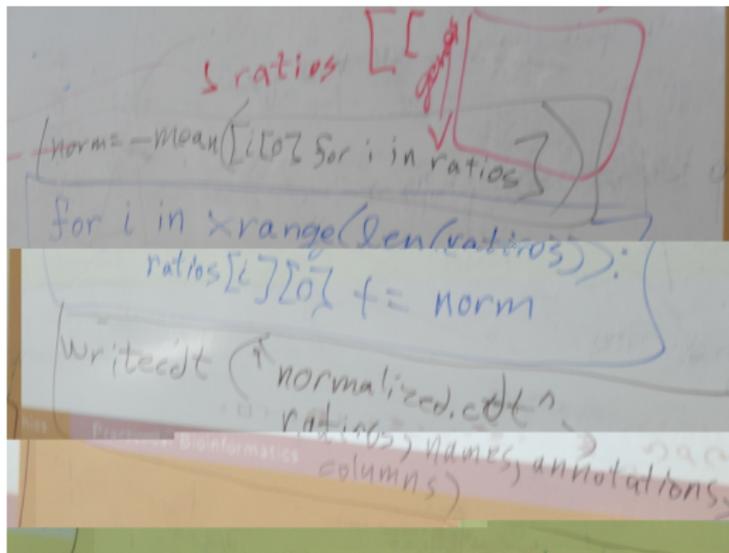
ratios = [makefloat(i) for i in row[2:]]
for row in data[1:]]

names
no. to ratios

Whiteboard Image



Whiteboard Image



Using the Cluster3 GUI

Gene Cluster 3.0

File Help

File loaded

Job name

Data set has Rows Columns

Filter Genes

- % Present >= 80
- SD (Gene Vector) 2,0
- At least 1 observations with abs(Val) >= 2,0
- MaxVal - MinVal >= 2,0

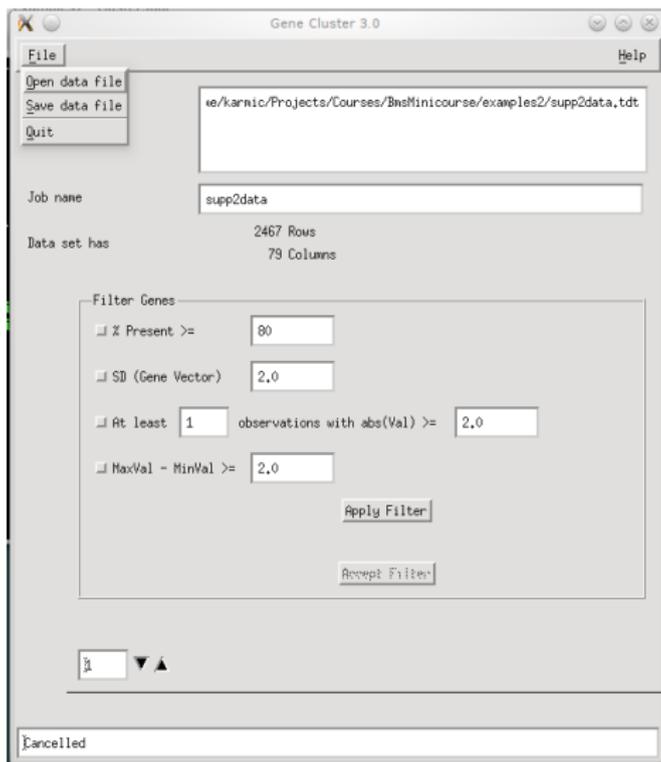
Apply Filter

Accept Filter

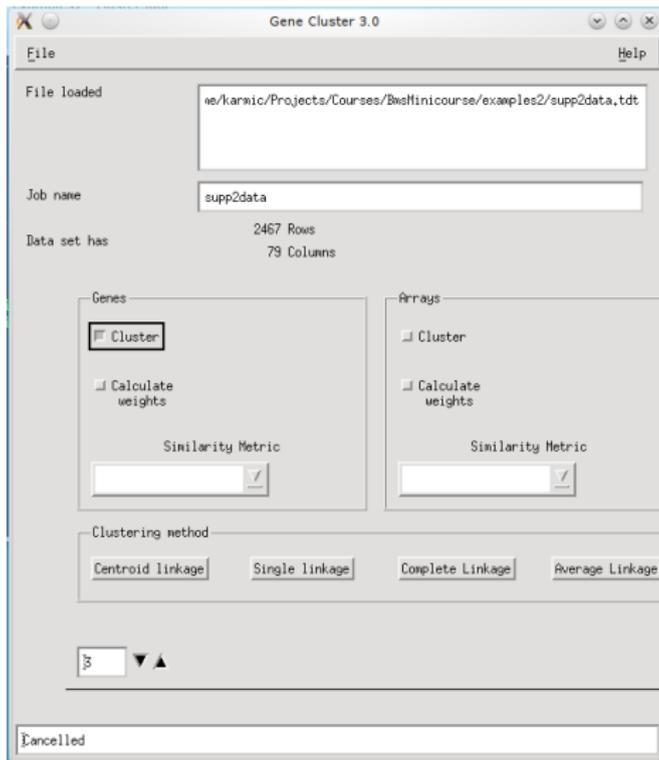
1

1

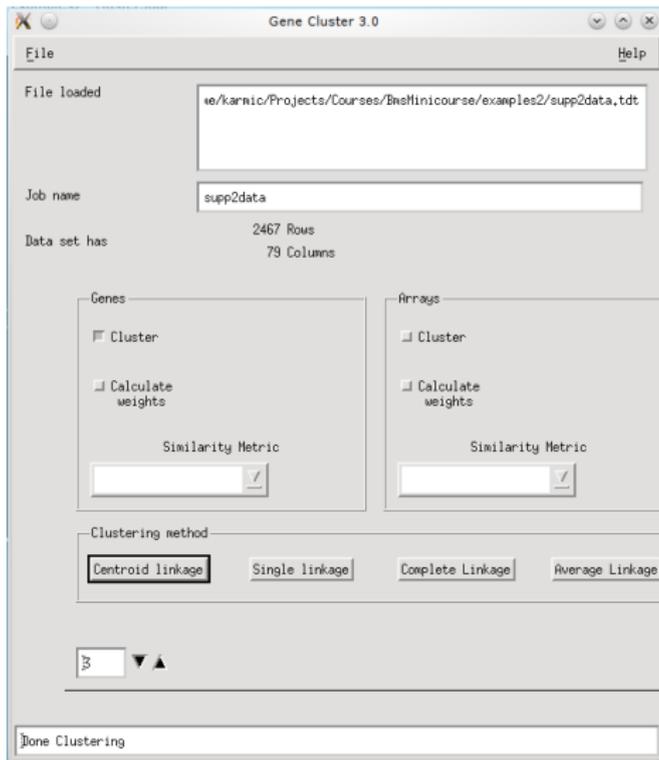
Load your data



Choose distance function



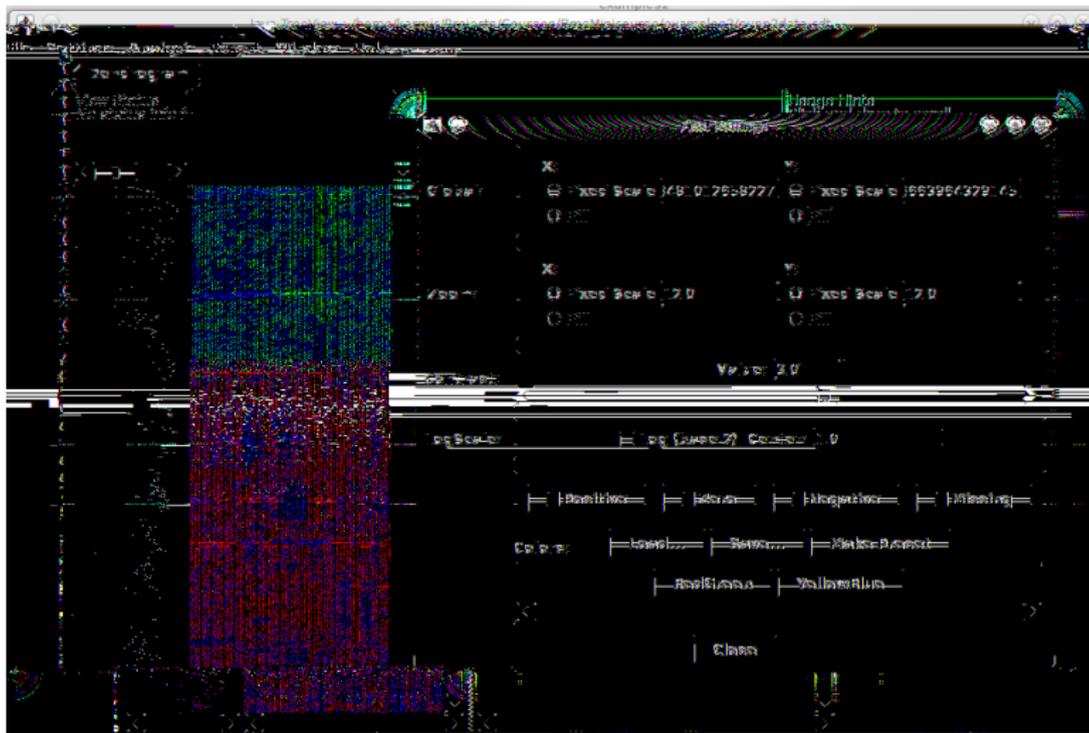
Choose linking method



Using JavaTreeView



Adjust pixel settings for global view



Select URL for gene annotations

The screenshot shows the Java TreeView application window titled "Java TreeView : /home/karmic/Projects/Courses/BmsMinicourse/examples2/supp2data.cdt". The menu bar includes "File", "Settings", "Analysis", "Export", "Window", and "Help". The "Settings" menu is open, showing options: "Annotations...", "Font Settings...", "Url Settings...", "Pixel Settings...", and "Presets". The "Presets" submenu is also open, displaying "Gene Url Presets..." (with a "Ctrl-P" shortcut) and "Array Url Presets...". A "Usage Hints" box in the top right corner of the application area contains the text "Click and drag to scroll". The main window area is filled with a complex, multi-colored tree view of data, likely representing a gene annotation hierarchy.

Select URL for gene annotations

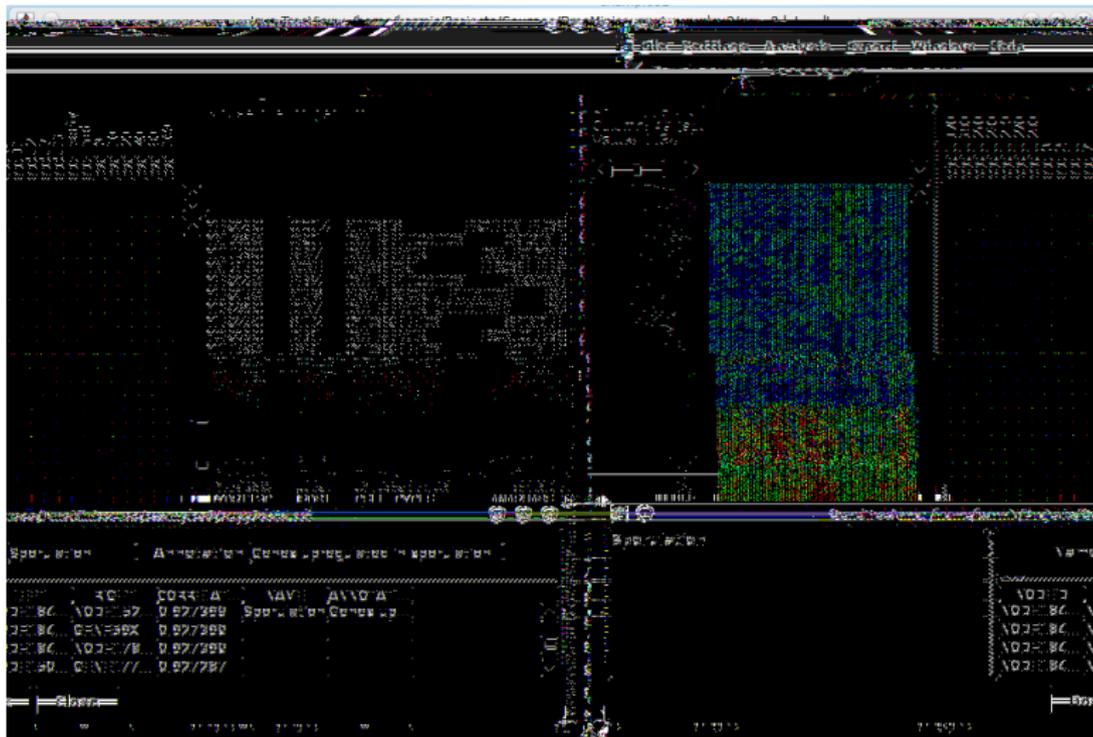
The screenshot shows the Java TreeView application interface. At the top, the title bar reads "Java TreeView - /home/karmic/Projects/Courses/BmsMinicourse/examples2/supp2data.cdt". The menu bar includes "File", "Settings", "Analysis", "Export", "Window", and "Help".

The main window is divided into several sections:

- Dendrogram:** A tree structure on the left with a vertical scrollbar. The nodes are labeled "alpha" followed by numbers from 0 to 150.
- Heatmap:** A central area displaying a color-coded heatmap with red and green vertical bands.
- Usage Hints:** A text box on the right stating "Click to select node - use arrow keys to navigate tree".
- Presets:** A section below the heatmap with a "Presets" label and a list of gene names: YAL062W, GDH3, GLUTAMATE BIOSYNTHESIS, NAD1, YOR375C, GDH1, GLUTAMATE BIOSYNTHESIS, and GLUT.
- Gene List:** A table below the presets with columns: "Enabled", "Header", "Name", "Template", and "Default?". The "Template" column contains a URL: <http://genome-www1.stanford.edu/cgi-bin/CGI/locus.pl?locus=LEAD5>.

At the bottom of the window, there is a detailed genomic track visualization with various colored bars and annotations.

Activate and detach annotation window



Calculating the Distance Matrix

Pearson similarity

$$s(x, y) = \frac{\sum_i^N (x_i - x_{offset})(y_i - y_{offset})}{\sqrt{\sum_i^N (x_i - x_{offset})^2} \sqrt{\sum_i^N (y_i - y_{offset})^2}} \quad (1)$$

Pearson similarity

$$s(x, y) = \frac{\sum_i^N (x_i - x_{offset})(y_i - y_{offset})}{\sqrt{\sum_i^N (x_i - x_{offset})^2} \sqrt{\sum_i^N (y_i - y_{offset})^2}} \quad (1)$$

Pearson distance

$$d_{uncentered}(x, y) = 1 - s(x, y) \quad (2)$$

Pearson similarity

$$s(x, y) = \frac{\sum_i^N (x_i - x_{\text{offset}})(y_i - y_{\text{offset}})}{\sqrt{\sum_i^N (x_i - x_{\text{offset}})^2} \sqrt{\sum_i^N (y_i - y_{\text{offset}})^2}} \quad (1)$$

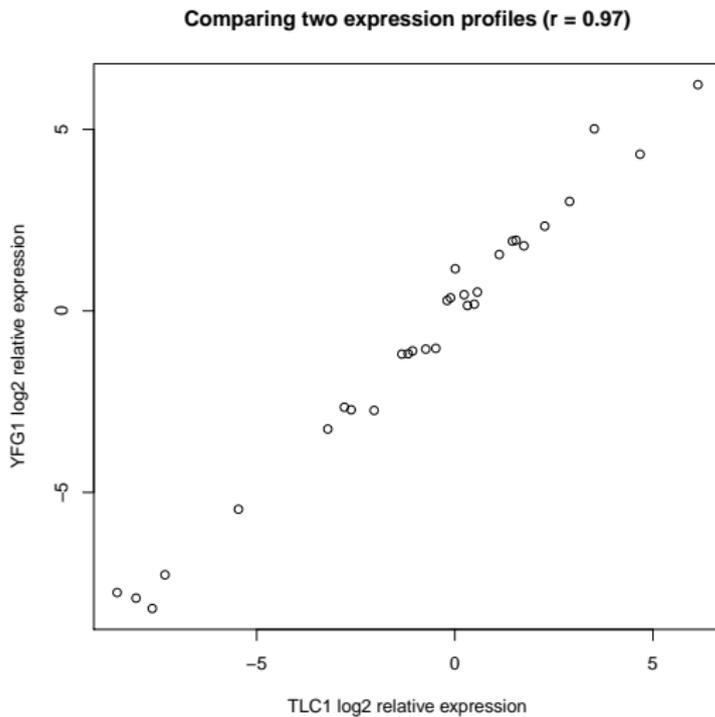
Pearson distance

$$d_{\text{uncentered}}(x, y) = 1 - s(x, y) \quad (2)$$

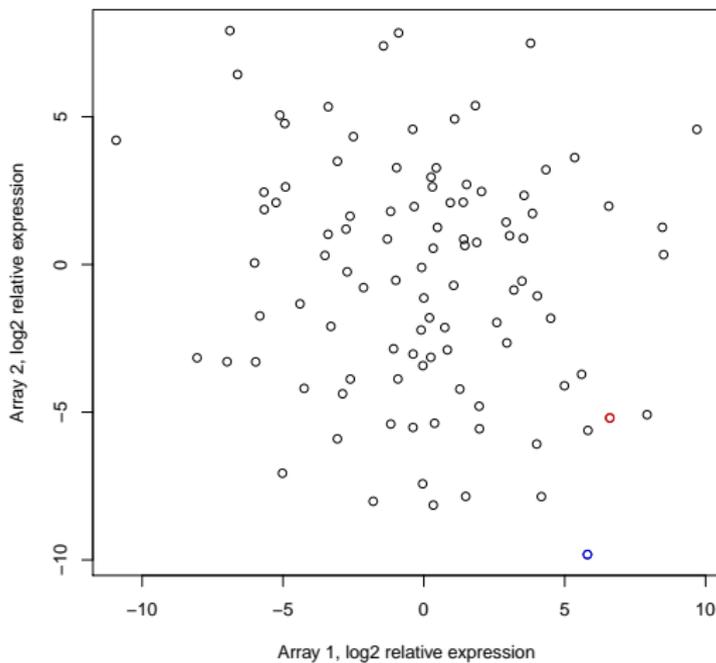
Euclidean distance

$$\frac{\sum_i^N (x_i - y_i)^2}{N} \quad (3)$$

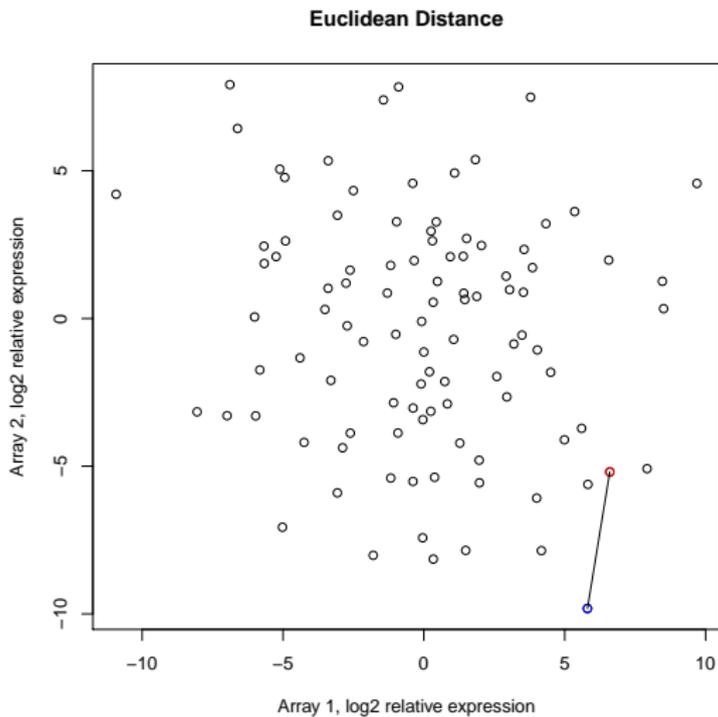
Comparing all measurements for two genes



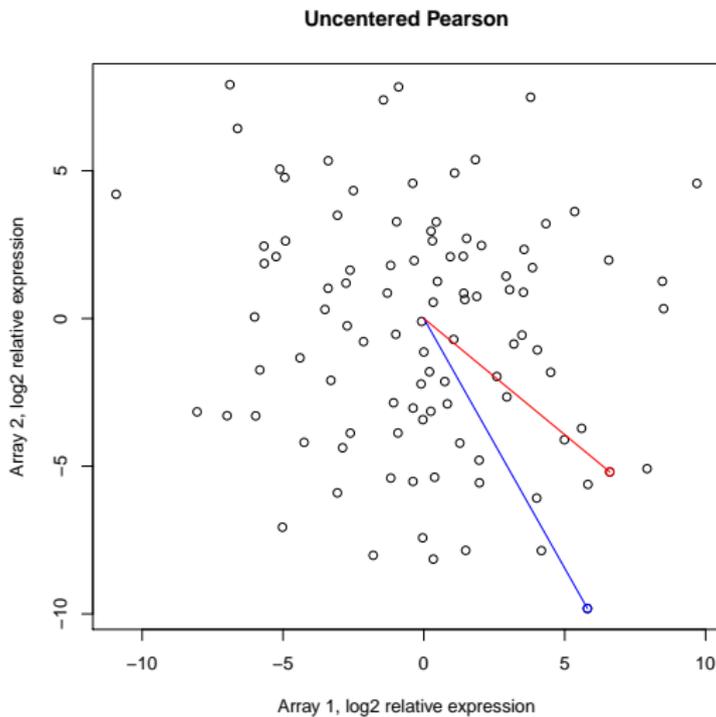
Comparing all genes for two measurements



Comparing all genes for two measurements



Comparing all genes for two measurements



Special cases for Pearson in Cluster3

```
int flag = 0;
/* flag will remain zero if no nonzero combinations of mask1 and mask2 are
 * found.
 */
int i;
for (i = 0; i < n; i++)
{ if (mask1[index1][i] && mask2[index2][i])
  { double term1 = data1[index1][i];
    double term2 = data2[index2][i];
    double w = weight[i];
    result += w*term1*term2;
    denom1 += w*term1*term1;
    denom2 += w*term2*term2;
    flag = 1;
  }
}
if (!flag) return 0.;
if (denom1==0.) return 1.;
if (denom2==0.) return 1.;
result = result / sqrt(denom1*denom2);
result = 1. - result;
return result;
```