# Practical Bioinformatics

Mark Voorhies

4/6/2017

```python
# Use import the first time you load a module
#   (And keep using import until it loads
#     successfully)
import my_module

my_module.my_function(42)

# Once a module has been loaded, use reload to
#   force python to read your new code
from importlib import reload
reload(my_module)
```

## Pearson distances

Pearson similarity

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

# Pearson distances

Pearson similarity

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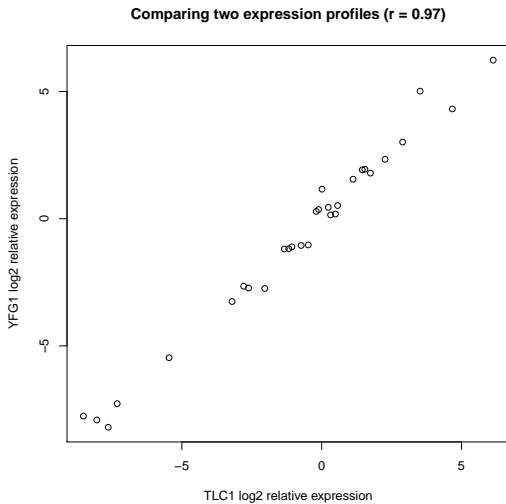
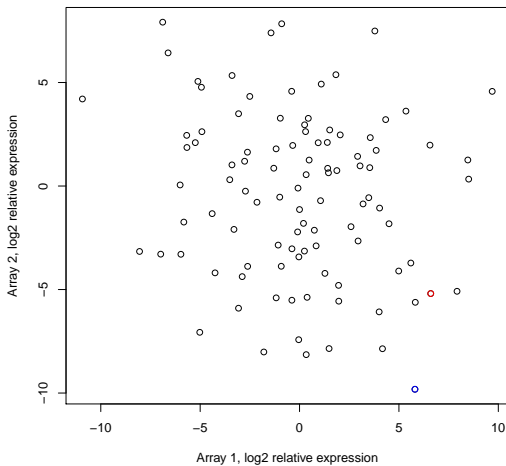Pearson distance

$$d(x, y) = 1 - s(x, y)$$

# Pearson distances

Pearson similarity

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

Pearson distance

$$d(x, y) = 1 - s(x, y)$$

Euclidean distance

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

Comparing two expression profiles (r = 0.97)

# Comparing all genes for two measurements

**Euclidean Distance**

Array 2, log2 relative expression

Array 1, log2 relative expression

**Uncentered Pearson**

$$\frac{(N-1)^2}{2}$$

It's hard work at times, but you have to be realistic. If you have a large database with many variables and your goal is to get a good understanding of the interrelationships, then, unless you get lucky, this complex structure is bound to require some hard work to understand.

Bill Cleveland and Rick Becker
http://stat.bell-labs.com/project/trellis/interview.html

# Using JavaTreeView

# Adjust pixel settings for global view

# Adjust pixel settings for global view

# Select annotation columns

# Select annotation columns

# Select URL for gene annotations

# Select URL for gene annotations

# Activate and detach annotation window

# Activate and detach annotation window

# Activate and detach annotation window

Write functions to reproduce the shuffling controls in figure 3 of the Eisen paper (removing correlations among genes and/or arrays).

Write functions to reproduce the shuffling controls in figure 3 of the Eisen paper (removing correlations among genes and/or arrays).

```python
def shuffleGenes(self, seed = None):
    """Shuffle expression matrix by row."""
    import random
    if(seed != None):
        random.seed(seed)
    indices = range(len(self.genes))
    random.shuffle(indices)
    genes = [self.geneName[i] for i in indices]
    self.geneName = genes
    annotations = [self.geneAnn[i] for i in indices]
    self.geneAnn = genes
    num = [self.num[i] for i in indices]
    self.num = num
```

Write functions to reproduce the shuffling controls in figure 3 of the Eisen paper (removing correlations among genes and/or arrays).

Write functions to reproduce the shuffling controls in figure 3 of the
Eisen paper (removing correlations among genes and/or arrays).

```python
def shuffleRows(self, seed = None):
    """Permute ratio values within rows."""
    import random
    if(seed != None):
        random.seed(seed)
    for i in self.num:
        random.shuffle(i)
```

# Clustering exercises – Negative controls

Write functions to reproduce the shuffling controls in figure 3 of the Eisen paper (removing correlations among genes and/or arrays).

```python
def shuffleRows(self, seed = None):
    """Permute ratio values within rows."""
    import random
    if(seed != None):
        random.seed(seed)
    for i in self.num:
        random.shuffle(i)


def shuffleCols(self, seed = None):
    """Permute ratio values within columns."""
    import random
    if(seed != None):
        random.seed(seed)
    # Transpose the expression matrix
    cols = []
    for col in xrange(len(self.num[0])):
        cols.append([row[col] for row in self.num])
    # Shuffle
    for i in cols:
        random.shuffle(i)
    # Transpose back to original orientation
    self.num = []
    for row in xrange(len(cols)):
        self.num.append([col[row] for col in row])
```

1. Explore different clustering methods and/or distance methods
2. Try additional shufflings of the data: how do they affect your ability to cluster the data? *C.f. figure 3 the Eisen paper*
   - Permute the columns
   - Independently permute the columns of each row