# Practical Bioinformatics

Mark Voorhies

5/28/2015

# Needleman-Wunsch with g=0

|   | A | G | C | G | G | T | A |
|---|---|---|---|---|---|---|---|
|   |   |   |   |   |   |   |   |
| G |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |
| G |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |
| G |   |   |   |   |   |   |   |
| G |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |

|   | A | G | C | G | G | T | A |
|---|---|---|---|---|---|---|---|
| | 0 | | | | | | |
| G | | | | | | | |
| A | | | | | | | |
| G | | | | | | | |
| C | | | | | | | |
| G | | | | | | | |
| G | | | | | | | |
| A | | | | | | | |

```
def nw_fill(seq1, seq2, s, e):
    # Initialize dp matrix
    #    first dimension = seq1 positions
    #    second dimension = seq2 positions
    #    m[i][j] = best score for subalignment
    #              of seq1[:i], seq2[:j]

    m = [[0]]
    # Initialize pointer matrix, a two dimensional
    #    matrix of lists of (row, column) pointers
    p = [[None]]
```

|   | A | G | C | G | G | T | A |
|---|---|---|---|---|---|---|---|
| | 0 | -1 | -2 | -3 | -4 | -5 | -6 | -7 |
| G | | | | | | | | |
| A | | | | | | | | |
| G | | | | | | | | |
| C | | | | | | | | |
| G | | | | | | | | |
| G | | | | | | | | |
| A | | | | | | | | |

```
# Fill first row as leading gaps
for j in range(len(seq2)):
    m[-1].append(m[0][j]+e)
    p[-1].append([(0,j)])
```

|   | A | G | C | G | G | T | A |
|---|---|---|---|---|---|---|---|
|   | 0 | -1 | -2 | -3 | -4 | -5 | -6 | -7 |
| G | -1 | -1 | 0 | -1 | -2 | -3 | -4 | -5 |
| A | -2 | 0 | -1 | -1 | -2 | -3 | -4 | -3 |
| G | -3 | -1 | 1 | 0 | 0 | -1 | -2 | -3 |
| C | -4 | -2 | 0 | 2 | 1 | 0 | -1 | -2 |
| G | -5 | -3 | -1 | 1 | 3 | 2 | 1 | 0 |
| G | -6 | -4 | -2 | 0 | 2 | 4 | 3 | 2 |
| A | -7 | -5 | -3 | -1 | 1 | 3 | 3 | 4 |

```
for i in range(len(seq1)):
    # First column is leading gaps
    m.append([m[i][0]+e])
    p.append([[(i,0)]])
    for j in range(len(seq2)):
        # Score for aligning seq1[i] with seq2[j]
        match = m[i][j]+s[seq1[i]][seq2[j]]
        # Score for aligning seq1[i] with a gap
        hgap = m[i+1][j]+e
        # Score for aligning seq2[i] with a gap
        vgap = m[i][j+1]+e

        best = max(match, vgap, hgap)
        m[-1].append(best)
        p[-1].append([])

        if(match == best):
            p[-1][-1].append((i,j))
        if(hgap == best):
            p[-1][-1].append((i+1,j))
        if(vgap == best):
            p[-1][-1].append((i,j+1))
```

|   | A | G | C | G | G | T | A |
|---|---|---|---|---|---|---|---|
|   | 0 | -1 | -2 | -3 | -4 | -5 | -6 | -7 |
| G | -1 | -1 | 0 | -1 | -2 | -3 | -4 | -5 |
| A | -2 | 0 | -1 | -1 | -2 | -3 | -4 | -3 |
| G | -3 | -1 | 1 | 0 | 0 | -1 | -2 | -3 |
| C | -4 | -2 | 0 | 2 | 1 | 0 | -1 | -2 |
| G | -5 | -3 | -1 | 1 | 3 | 2 | 1 | 0 |
| G | -6 | -4 | -2 | 0 | 2 | 4 | 3 | 2 |
| A | -7 | -5 | -3 | -1 | 1 | 3 | 3 | 4 |

```python
# Start at bottom right corner
curpos = (len(seq1),len(seq2))
aligned1 = ""
aligned2 = ""

exitFlag = False
for i in range(len(seq1)+len(seq2)):
    plist = p[curpos[0]][curpos[1]]
    if(plist is None):
        exitFlag = True
        break
    nextpos = plist[0]
    # Check for vgap
    if(nextpos[0] == curpos[0]):
        aligned1 = "-"+aligned1
    else:
        aligned1 = seq1[nextpos[0]]+aligned1
    # Check for hgap
    if(nextpos[1] == curpos[1]):
        aligned2 = "-"+aligned2
    else:
        aligned2 = seq2[nextpos[1]]+aligned2

    curpos = nextpos

if(exitFlag == False):
    print "WARNING: Unexpected exit from traceback"
```

- For tools:
    - Read the manual
    - Read the paper
- Good general references:
    - The O'Reilly BLAST book
    - Durbin, Eddy, Krogh, and Mitcheson (HMMs)
    - Numerical Recipes
    - Branden & Tooze

Observation

Model

Hypothesis

Experiment

- Export, audit, edit, and import *independent* of a given program.
- Standard file formats for portability.
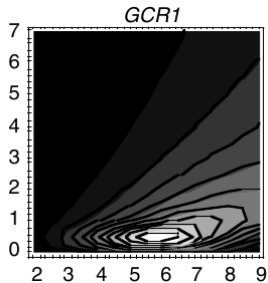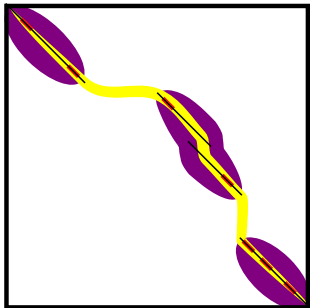- Don't be afraid to look inside and hack on *your* data files.

Iteration, clustering, dynamic programming, ...

CLUSTALW timings on Intel Core2 T7300@2.00GHz, 32bit

$y = (1.8e-9)x^\wedge(2.08)$

BLAST, HMMer3, MrBayes, ...

# Evolution is a rich source of information



- Infer homology from sequence similarity
- More sequences provide more information

# HMMs capture position and gap information



(Image from Sean Eddy, PLoS Comp. Biol. 4:e1000069)

- Make sure you know what you are measuring
- Nucleic acid sequences are especially easy to address
- Many phenotypes can be analyzed by common numerical methods

# Tools should support aggregation and annotation

# Groovy Packages

- numpy
- matplotlib
- scipy
- networkx
- h5py
- rpy
- pandas
- Pycluster (and BioPython)
- MySQLdb
- scikit-image
- scikit-learn

- Follow computational methods as they evolve (Web of Science, PubMed RSS, arXiv, Haldane's Sieve...)

- Follow computational methods as they evolve (Web of Science, PubMed RSS, arXiv, Haldane's Sieve...)
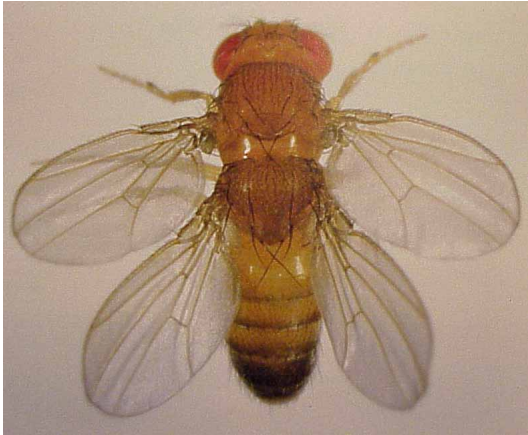- As a reviewer, insist on availability of source code

## Science is a Conversation

- Follow computational methods as they evolve (Web of Science, PubMed RSS, arXiv, Haldane's Sieve...)
- As a reviewer, insist on availability of source code
- Draw on your classmates' expertise

# We understand systems by breaking them