# Pairwise Alignment

Mark Voorhies

3/27/2012

## Review: Tips and tricks

Making a file executable:

chmod "a+x" pydotter.py

Handling file/directory names with spaces:

**cd** My\ Directory\ with\ Spaces

or

**cd** "My Directory with Spaces"

## Review: Tips and tricks

Killing a process on OS X:

- Try ctrl-c
- If that doesn't work:

    - **ps** −awx | grep name_of_process
    - First column in ps output is PID (process ID)
    - **kill** PID
    - If that doesn't work:

        **kill** −KILL PID

- On Linux:

    **ps** −ealf | grep name_of_process

- FASTA files

  ```
  >Name Free-form annotation
  MGCLLIMKEGGPGRKHKLIVMLYLDENQ
  EHELPIMTRAPPEDINADNAMACHINEW
  NQEDLYMNILKHGPPGEDEDRKHEDEDG
  ```
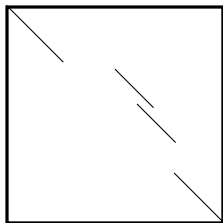
- FASTA files

  ```
  >Name Free-form annotation
  MGCLLIMKEGGPGRKHKLIVMLYLDENQ
  EHELPIMTRAPPEDINADNAMACHINEW
  NQEDLYMNILKHGPPGEDEDRKHEDEDG
  ```

- Dotplots: unbiased plot of all possible ungapped alignments of two sequences.

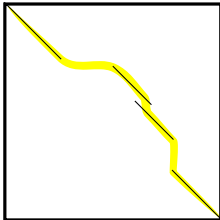How can we automate our dotplot protocol to find the "best" gapped alignment of our sequences?

How can we automate our dotplot protocol to find the "best" gapped alignment of our sequences?

What do we mean by best?

How can we automate our dotplot protocol to find the "best" gapped alignment of our sequences?

What do we mean by best?

- Residues with equivalent functional roles are paired

How can we automate our dotplot protocol to find the "best" gapped alignment of our sequences?

What do we mean by best?



- Residues with equivalent functional roles are paired
- Residues that derive from the same position in the common ancestor are paired (homology)

How can we automate our dotplot protocol to find the "best" gapped alignment of our sequences?

What do we mean by best?



- Residues with equivalent functional roles are paired
- Residues that derive from the same position in the common ancestor are paired (homology)
- The sequence alignment maximizes a similarity function

Frequency of residue $i$:

$$p_i$$

# Deriving scores from alignments

Frequency of residue $i$:

$$p_i$$

Frequency of residue $i$ aligned to residue $j$:

$$q_{ij}$$

# Deriving scores from alignments

Frequency of residue $i$:

$$p_i$$

Frequency of residue $i$ aligned to residue $j$:

$$q_{ij}$$

Expected frequency if $i$ and $j$ are independent:

$$p_i p_j$$

# Deriving scores from alignments

Frequency of residue $i$:

$$p_i$$

Frequency of residue $i$ aligned to residue $j$:

$$q_{ij}$$

Expected frequency if $i$ and $j$ are independent:

$$p_i p_j$$

Ratio of observed to expected frequency:

$$\frac{q_{ij}}{p_i p_j}$$

## Deriving scores from alignments

Frequency of residue $i$:

$$p_i$$

Frequency of residue $i$ aligned to residue $j$:

$$q_{ij}$$

Expected frequency if $i$ and $j$ are independent:

$$p_i p_j$$

Ratio of observed to expected frequency:

$$\frac{q_{ij}}{p_i p_j}$$

Log odds (LOD) score:

$$s(i,j) = \log \frac{q_{ij}}{p_i p_j}$$

- PAM1 matrix originally calculated from manual alignments of highly conserved sequences (myoglobin, cytochrome C, etc.)

- PAM1 matrix originally calculated from manual alignments of highly conserved sequences (myoglobin, cytochrome C, etc.)
- We can think of a PAM matrix as evolving a sequence by one unit of time.

## PAM (Dayhoff) and BLOSUM matrices

- PAM1 matrix originally calculated from manual alignments of highly conserved sequences (myoglobin, cytochrome C, etc.)
- We can think of a PAM matrix as evolving a sequence by one unit of time.
- If evolution is uniform over time, then PAM matrices for larger evolutionary steps can be generated by multiplying PAM1 by itself (so, higher numbered PAM matrices represent greater evolutionary distances).

# PAM (Dayhoff) and BLOSUM matrices

- PAM1 matrix originally calculated from manual alignments of highly conserved sequences (myoglobin, cytochrome C, etc.)
- We can think of a PAM matrix as evolving a sequence by one unit of time.
- If evolution is uniform over time, then PAM matrices for larger evolutionary steps can be generated by multiplying PAM1 by itself (so, higher numbered PAM matrices represent greater evolutionary distances).
- The BLOSUM matrices were determined from automatically generated ungapped alignments. Higher numbered BLOSUM matrices correspond to *smaller* evolutionary distances. BLOSUM62 is the default matrix for BLAST.

## Fun with logarithms

In log space, multiplication and division become addition and subtraction:

$$\begin{aligned} \log(xy) &= \log(x) + \log(y) \\ \log(x/y) &= \log(x) - \log(y) \end{aligned}$$

Therefore, exponentiation becomes multiplication:

$$\log(x^y) = y \log(x)$$

Also, we can change of the base of a logarithm like so:

$$\log_A(x) = \log(x)/\log(A)$$

# Scoring an alignment

Log odds (LOD) score:

$$s(i,j) = \log \frac{q_{ij}}{p_i p_j}$$

## Scoring an alignment

Log odds (LOD) score:

$$s(i,j) = \log \frac{q_{ij}}{p_i p_j}$$

Multiplying independent probabilities is equivalent to adding independent log probabilities.

## Scoring an alignment

Log odds (LOD) score:

$$s(i,j) = \log \frac{q_{ij}}{p_i p_j}$$

Multiplying independent probabilities is equivalent to adding independent log probabilities.

Therefore, for an ungapped alignment can be scored as:

$$S(x,y) = \log \prod_i^N \frac{q_{x_i y_i}}{p_{x_i} p_{y_i}} = \sum_i^N s(x_i, y_i)$$

## Scoring an alignment

Log odds (LOD) score:

$$s(i,j) = \log \frac{q_{ij}}{p_i p_j}$$

Multiplying independent probabilities is equivalent to adding independent log probabilities.

Therefore, for an ungapped alignment can be scored as:

$$S(x,y) = \log \prod_i^N \frac{q_{x_i y_i}}{p_{x_i} p_{y_i}} = \sum_i^N s(x_i, y_i)$$

What about gaps?

## Scoring an alignment

Log odds (LOD) score:

$$s(i,j) = \log \frac{q_{ij}}{p_i p_j}$$

Multiplying independent probabilities is equivalent to adding independent log probabilities.

Therefore, for an ungapped alignment can be scored as:

$$S(x,y) = \log \prod_i^N \frac{q_{x_i y_i}}{p_{x_i} p_{y_i}} = \sum_i^N s(x_i, y_i)$$

What about gaps?

- Probability of an insertion/deletion event (gap opening, $G$)
- Length distribution of insertions/deletions (gap extension, $E$)

## Scoring an alignment

Log odds (LOD) score:

$$s(i,j) = \log \frac{q_{ij}}{p_i p_j}$$

Multiplying independent probabilities is equivalent to adding independent log probabilities.

Therefore, for an ungapped alignment can be scored as:

$$S(x,y) = \log \prod_i^N \frac{q_{x_i y_i}}{p_{x_i} p_{y_i}} = \sum_i^N s(x_i, y_i)$$

What about gaps?

- Probability of an insertion/deletion event (gap opening, $G$)
- Length distribution of insertions/deletions (gap extension, $E$)

$$S_{gapped}(x,y) = S(x,y) + \sum_i^{gaps} (G + E * L_i)$$

## Scoring an alignment

Log odds (LOD) score:

$$s(i,j) = \log \frac{q_{ij}}{p_i p_j}$$

Multiplying independent probabilities is equivalent to adding independent log probabilities.

Therefore, for an ungapped alignment can be scored as:

$$S(x,y) = \log \prod_i^N \frac{q_{x_i y_i}}{p_{x_i} p_{y_i}} = \sum_i^N s(x_i, y_i)$$

What about gaps?

- Probability of an insertion/deletion event (gap opening, $G$)
- Length distribution of insertions/deletions (gap extension, $E$)

$$S_{gapped}(x,y) = S(x,y) + \sum_i^{gaps} (G + E * L_i)$$

We find an optimal alignment by finding $x$ and $y$ that maximize $S$.

# How many ways can we align two sequences?

# How many ways can we align two sequences?

Binomial formula:

$$\binom{k}{r} = \frac{k!}{(k-r)!r!}$$

$$\binom{2n}{n} = \frac{(2n)!}{n!n!}$$

Stirling's approximation:

$$x! \approx \sqrt{2\pi} \left( x^{x+\frac{1}{2}} \right) e^{-x}$$

$$\binom{2n}{n} \approx \frac{2^{2n}}{\sqrt{\pi n}}$$

$\frac{2^{2n}}{\sqrt{\pi n}}$ is too expensive.

$\frac{2^{2n}}{\sqrt{\pi n}}$ is too expensive.

$$S_{gapped}(x, y) = S(x, y) + \sum_{i}^{gaps} (G + E * L_i)$$

$\frac{2^{2n}}{\sqrt{\pi n}}$ is too expensive.

$$S_{gapped}(x,y) = S(x,y) + \sum_{i}^{gaps} (G + E * L_i)$$

The best alignment of any pair of subsequences is independent of the global alignment.

ATGC
ATGC

ATGC
ATGC

ATGCTTCG
ATGC...G

- DOTTER: $O(n^2)$
- Exhaustive search: $\frac{2^{2n}}{\sqrt{\pi n}}$
- Dynamic programming: $O(n^2)$ to $O(n^3)$

# Annotating features in JALVIEW

# Annotating features in JALVIEW

# Annotating features in JALVIEW

# Annotating features in JALVIEW

```
# List all differences between two text files
# (empty output for identity)
diff HvSs.gap0.0.both.aln HvSs.gap0.0.mult.aln
# Report only whether the files differ
# (empty output for identity)
diff -q HvSs.gap0.0.both.aln HvSs.gap0.0.mult.aln
```

(*NIX = *BSD, OS X, Solaris, Linux, Windows with Cygwin, ...)

- Use a text or sequence editor to create a spliced variant of HvPhytepsin that can be aligned to the full HsSaposinC sequence

- Find the GenBank entries for HvPhytepsin and SsPepsinogen (tip: use the identifiers from the FASTA files) and find the corresponding transcript sequences.

  - How easy is it to align the proteins vs. the transcripts?
  - Can you tell if you are getting equivalent results from the two alignments?

- Try repeating this exercise for a pair of sequences where genomic sequence is available; *e.g.*, *A. nidulans* VosA (ABQ18268.1), and *H. capsulatum* Ryp2 (ACB59236.1).