

Dynamic Programming in Haskell

Christian Neukirchen*

March 2006

The purpose of this Literate Haskell program is to implement a function that does global sequence alignment using Needleman/Wunsch techniques.¹

The algorithm is based on two steps: first, filling a matrix with the maximal alignment scores for each element and then tracing a path connecting the top-left and the bottom-right cell. Note that the matrix is $O(n \cdot m)$ memory-wise and therefore pretty inefficient, you don't want to use this on bigger sequences.

In Haskell, a good way to implement Dynamic Programming like this is an array that will memoize a lazy stream of scores per cell. This allows $O(1)$ -lookup of formerly calculated values without losing referential transparency and (to an extent) lazy evaluation.

```
import Array
```

align is the function that wraps all the functions below and calls them in correct order. It takes two strings and returns them aligned:

```
align      :: String → String → [String]
align da db = format $ reverse $ traceback lena lenb
where
  lena = length da
  lenb = length db
```

The algorithm is easier to express when the sequences to align are one-indexed, since the borders of the matrix are used as special values. An easy way to achieve this is prepending a space:

```
a = '␣' : da
b = '␣' : db
```

*The author can be reached at <http://chneukirchen.org>.

¹More about these techniques, graphics helpful for understanding, and a codeless step-by-step explanation can be found at <http://www.sbc.su.se/~pjk/molbioinfo2001/dynprog/dynamic.html>.

memscore is the array that contains the actual matrix. It is filled using a lazy stream of scores for each element.

```
memscore = listArray ((0,0), (lena, lenb))
             [score x y | x ← [0..lena], y ← [0..lenb]]
```

The scoring function looks very confusing since Haskell's array access operator is not very elegant. I'll introduce an infix operator $i @@ j$ that corresponds to $M_{i,j}$:

```
infix 5 @@
(@@) i j = memscore ! (i,j)
```

The score $M_{i,j}$ of each element is determined in below code as follows, the borders of the matrix with $i = 0$ and $j = 0$ are initialized to zero. (More complex scoring algorithms could be added easily.)

$$M_{i,j} = \text{maximum of } \begin{cases} M_{i-1,j-1} + S_{i,j} \\ M_{i,j-1} + w \\ M_{i-1,j} + w \end{cases}$$

The gap penalty w is zero here for reasons of simplicity.

```
score 0 _ = 0
score _ 0 = 0
score x y = maximum [(x - 1 @@ y - 1) + difference x y,
                    x - 1 @@ y,
                    x      @@ y - 1]
```

$S_{i,j}$ is a mismatch penalty defined here like this:

$$S_{i,j} = \begin{cases} 0 & \text{if the symbols at position } i \text{ and position } j \text{ match} \\ 1 & \text{otherwise} \end{cases}$$

```
where difference x y | a !! x ≡ b !! y = 1
                    | otherwise      = 0
```

traceback now finds the path connecting both corners of the matrix and collects the appropriate symbols (or spaces for gaps).

```
traceback :: Int → Int → [(Char, Char)]
traceback 0 0 = []
traceback x y | x      ≡ 0      = ('␣' , b !! y) : traceback 0      (y - 1)
              | y      ≡ 0      = (a !! x, '␣' ) : traceback (x - 1) 0
              | x @@ y ≡ x @@ y - 1 = ('␣' , b !! y) : traceback x      (y - 1)
              | x @@ y ≡ x - 1 @@ y = (a !! x, '␣' ) : traceback (x - 1) y
              | otherwise          = (a !! x, b !! y) : traceback (x - 1) (y - 1)
```

The resulting list of tuples like [('a', 'd'), ('b', 'e'), ('c', 'f')] gets converted by *format* into ["abc", "def"].

```
format l = [map fst l, map snd l]
```

Finally, a small main program to test the algorithm:

```
dna1 = "GAATTCAGTTA"  
dna2 = "GGATCGA"  
main = mapM_ putStrLn $ align dna1 dna2
```

Expected output:

```
G_AATTCAGTTA  
GGA_T_C_G_A
```

As you can see, corresponding symbols are aligned, with appropriate gaps in between. Implementing more complex rules for alignment is left as an exercise for the reader.

A run where bigger gaps are needed:

```
dna1 = "ATGGCTTCTACC"  
dna2 = "TATCAAAAGCCG"  
  
_ATGGCTTCTA_____CC_  
TAT__C_____AAAAGCCG
```