Sequence Alignment - Dynamic Programming Method

Antonio Osamu Katagiri Tanaka (https://www.katagiri-mx.com/about-me) - MNT16 - A01212611

```
# Clear all objects (from the workspace)
rm(list = ls())
# Increase the limit for max.print in R
options(max.print=10000)
```

```
f(\mathrm{seqX},\mathrm{seqY}) = egin{cases} & \mathrm{score\_match} & \mathrm{if}\ \mathrm{seqX} = \mathrm{seqY} \ & \mathrm{score\_mismatch} & \mathrm{if}\ \mathrm{seqX} 
eq \mathrm{Y} \end{cases}
```

```
f <- function(seqX, seqY, score_match, score_mismatch)
{
   if (seqX == seqY)
   {
      score <- score_match
   }
   else
   {
      score <- score_mismatch
   }
   return(score)
}</pre>
```

```
f_0 = s_{i-1,j-1} + f(\text{seqX}, \text{seqY})
```

```
# Score function
f0 <- function(seqX, seqY, score_match, score_mismatch, score_global_prev)
{
    score_global_curr = score_global_prev + f(seqX, seqY, score_match, score_mismatch)
    return(score_global_curr)
}</pre>
```

$$f_{1_x} = \max_{x \geq 1} \left\{ s_{i-1,j} - w_x
ight\} \, \, \& \, f_{1_y} = \max_{y \geq 1} \left\{ s_{i,j-1} - w_y
ight\}$$

```
# Score function if seqX (or seqY) is aligned
f1 <- function(score_global_prev, weight)
{
   score_global_curr <- score_global_prev + weight
   return(score_global_curr)
}</pre>
```

$$f_{score} = \max egin{cases} & f_0 \ & f_{1_x} \ & f_{1_{-}} \end{cases}$$

```
f_score <- function(seqX, seqY, score_match, score_mismatch, score_global_prev, score_global_prevX, score_g
lobal_prevY, weight)
{
    score0 = f0(seqX, seqY, score_match, score_mismatch, score_global_prev)
    score1 = f1(score_global_prevX, weight)
    score2 = f1(score_global_prevY, weight)
    score = max(c(score0, score1, score2))
    # TODO: but what if score0=score1, score0=score2, score1=score2, or score0=score1=score2?
    return(score)
}</pre>
```

Let's write a function to print the alignment

```
f_alignment <- function(seqx, seqy, score, score_match, score_mismatch, weights)</pre>
  ## Traceback
  i <- length(seqx)</pre>
  j <- length(seqy)</pre>
  ax <- character()</pre>
  ay <- character()</pre>
  # Add a 'dash' if a gap is required
  while (i>1 && j>1){
    sc <- score[i-1,j-1]</pre>
    sc <- f0(seqx[i], seqy[j], score_match, score_mismatch, sc)</pre>
    if (sc == score[i,j]) # best score is in the upper diagonal
      ax <- c(seqx[i], ax)</pre>
      ay <- c(seqy[j], ay)</pre>
      i <- i-1
      j <- j-1
    else if (f1(score[i,j-1], weights) == score[i,j]) # best score is up
      ax <- c("-",
      ay <- c(seqy[j], ay)</pre>
      j <- j-1
    else #if ((score[i-1,j] + weights) == score[i,j]) # best score is to the Left
      ax <- c(seqx[i], ax)</pre>
      ay <- c("-",
                        ay)
      i <- i-1
    # TODO: consider more paths if the best score appears in more than one 'cell'
  }
  # Add a 'bar' if both are equal
  seq.x <- c('',unlist(strsplit(paste(ax, collapse=''), '')))</pre>
  seq.y <- c('',unlist(strsplit(paste(ay, collapse=''), '')))</pre>
  bar = character()
  for (i in 2:length(seq.x))
    if (seq.x[i]==seq.y[i])
      bar <- c(bar, '|')
    else
      bar <- c(bar, ' ')</pre>
  }
  cat(paste(ax, collapse=''), "\n")
  cat(paste(bar, collapse=''), "\n")
  cat(paste(ay, collapse=''), "\n")
  return (c(ax, ay))
}
```

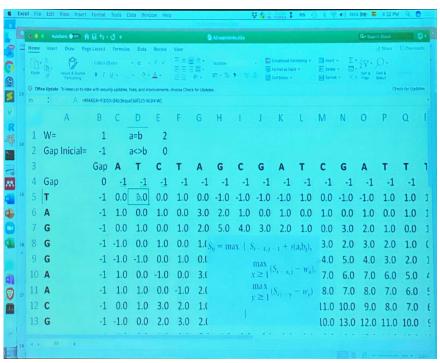
Let's write a function *f_global_score* to calculate the global score using the DPM algorithm

```
f_global_score <- function(Y, X, score_match, score_mismatch, weights, printTable, printAlignment)</pre>
  # Convert the strings into sequences (arrays) to access each letter in a loop
  # Also, contatenate an empty character at the begining of each sequence
  seq.x <- c('',unlist(strsplit(X, '')))</pre>
  seq.y <- c('',unlist(strsplit(Y, '')))</pre>
  # Initialize the DPM matrix
  tbl_score <- matrix(NA, length(seq.x), length(seq.y))</pre>
  # Add wy weights
  tbl_score[,1] <- sapply(1:length(seq.x)-1, function(x) weights)</pre>
  # Add wx weights
  tbl_score[1,] <- sapply(1:length(seq.y)-1, function(x) weights)</pre>
  # Make the top left corner eq to zero
  tbl_score[1,1] <- 0
  # Iterate tbl_score and do the DPM algorithm
  # NOTE: start looping from the 3rd element of the sequence as 1st is the "letter" and 2nd is the weight
  for (i in 2:length(seq.x))
    for (j in 2:length(seq.y))
      tbl_score[i,j] = f_score(
        seq.x[i],
                             # seqX
        seq.y[j],
                             # seqY
        score_match,
                             # score_match
        score mismatch,
                            # score mismatch
        tbl_score[i-1, j-1], # score_global_prev
        tbl_score[i-1,j], # score_global_prevX
        tbl_score[i,j-1],
                            # score_global_prevY
        weights)
                             # weight
    }
  }
  # Save table
  dimnames(tbl_score) = list(
      c('Gap', seq.x[-1]), # row names
      c('Gap', seq.y[-1])) # column names
  # Save global score
  globalScore = tbl_score[nrow(tbl_score),ncol(tbl_score)]
  # Print tbl_score and Global Score
  if (printTable)
    print(tbl_score)
    cat("Global Score =", globalScore, '\n')
  if (printAlignment)
    f_alignment(seq.x, seq.y, tbl_score, score_match, score_mismatch, weights)
  return(globalScore)
}
```

Let's test the f_global_score function with the class example

```
##
                 AGCGAT
                              C
                                G
                                   Α
                                     Т
     Gap
        ATC
              Т
                                        Т
## Gap
      0 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1
## T
                 0 -1 -1 -1 -1
                            1
## A
        1 0 1 0 3 2 1 0 1 0
                               1
      -1 0 1 0 1 2 5 4 3 2 1 0 3
## G
      -1 -1 0 1 0 1 4 5 6 5 4 3
## G
     -1 -1 -1 0 1 0 3 4 7 6 5 4 5 4 3
## G
## A
     -1 1 0 -1 0 3 2 3 6 9 8 7 6 7 6 5
      ## A
      -1 0 1 3 2 1 2 5 4 7 8 11 10 9 8 7
## C
## G
     -1 -1 0 2 3 2 3 4 7 6 7 10 13 12 11 10
## Global Score = 10
## TAGGGAACG---
## ||| || ||
## TAGCGATCGATT
```

[1] 10



Class Example

numbers match, so let's try a longer sequence ...

```
cat("Global Score =", sc_best, '\n')
```

```
## Global Score = 92
```

Class Example

Not equal, but close to the NCBI's BLAST

Generate score values with shuffled sequences

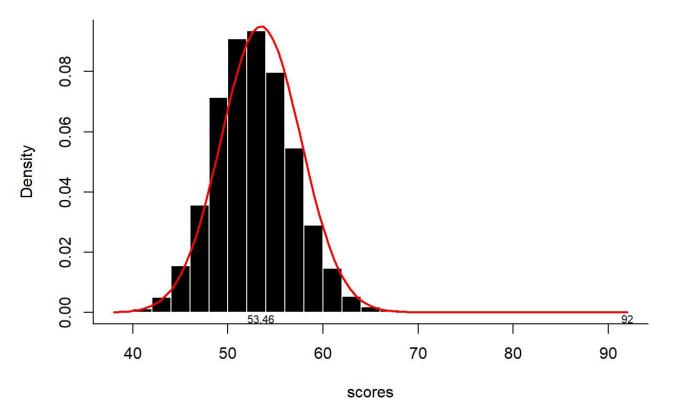
```
printAlignment <- FALSE

shuffled_scores = c(sc_best)
for (i in 0:10000)
{
    X_shuff = stringi::stri_rand_shuffle(X)
    Y_shuff = stringi::stri_rand_shuffle(Y)
    shuffled_scores = c(
        shuffled_scores,
        f_global_score(X_shuff, Y_shuff, score_match, score_mismatch, weights, printTable, printAlignment))
}</pre>
```

Analyze the shuffled score values

```
hist(shuffled_scores,
prob=TRUE, col="black", border="white", xlab="scores", breaks=25)
box(bty="1")
# Draw density function (assuming normal dist)
score_mean = mean(shuffled_scores)
score_sd = sd(shuffled_scores)
curve(dnorm(x,mean=score_mean,sd=score_sd), add=TRUE, col="red", lwd=2)
text(score_mean, -0.002, round(score_mean,2), cex = 0.7)
text(sc_best, -0.002, round(sc_best,2), cex = 0.7)
```

Histogram of shuffled_scores



Calculating a Single p Value From a Normal Distribution

```
t.test(shuffled_scores)
```

```
##
## One Sample t-test
##
## data: shuffled_scores
## t = 1274.2, df = 10001, p-value < 2.2e-16
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## 53.37996 53.54446
## sample estimates:
## mean of x
## 53.46221</pre>
```