# Homework 3

## CS 4364/5364
## Spring 2022

Due: 23 February 2022

1. **(25 points)** *Profile Alignment Problem* Given two sequence profiles $T$ and $S$, both of size $\sigma \times n$ (that is each represents a sequence of length $n$, but with probabilities of each character from the alphabet at each position), determine the optimal alignment (i.e. which columns of $S$ align with which columns of $T$) under the scoring scheme $\delta$.

   Your task: **modify** the Needlman-Wunch global alignment algorithm to consider these profiles rather than sequences. You can assume that the replacement costs are defined in a function $\delta(a,b) \to \mathbb{Z}, \forall a,b \in \Sigma \cup \{'-'\}$. Give the algorithm, an explanation of correctness, and analysis of it's running time.

   An example alignment is shown below over the alphabet $\Sigma = \{A, C, T, G\}$, as well as it's alignment score. Note that the score for a column is now no longer the value of $\delta$ for the two characters being aligned, but the weighted sum of these values.



Figure 1: Alignment of two profiles



Figure 2: Scoring scheme