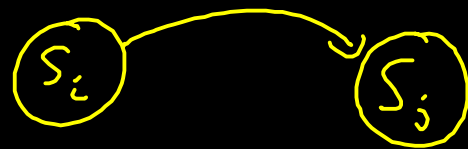


Theorem

Assuming a constant size of the alphabet above greedy algorithm given input $S = \{s_1, \dots, s_n\}$ takes running time in $O(N \cdot (n + \log(N)))$, $N = \sum_{1 \leq i \leq n} |s_i|$.

Proof:

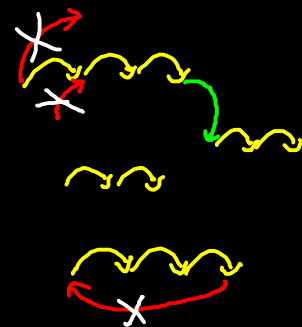
Under all remaining edges find one with maximal weight which "can be used":



$\langle s_i, s_j \rangle$



a) Row of start vertex
or column of goal vertex
labeled (constant time)



b) FIND on incident vertices gives
same partition. (constant time)

a) v b) \rightarrow edge not used.

At most $2n-1$ edges cannot be used (per iteration)

\implies all iterations together in time $O(n^2) = O(n \cdot N)$

If we "insert" the edge $s_i \rightarrow s_j$ we do the following:

- label row i and column j of adjacency matrix

→ excluding $2n-2$ edges (constant time)

• UNION (S_1, S_2)

→ scan shortest list element by element and update $R[v]$; merge the two lists (time $O(n)$)

⇒ iteration over all edges selected $O(n^2) = O(n \cdot N)$.

□

Performance guarantees?

Example: Let $S = \{c(ab)^m, (ba)^m, (ab)^m c\}$, $m \in \mathbb{N}$, the input for our greedy algorithm. The pair with maximal overlap is given by $(c(ab)^m, (ab)^m c)$ so the algorithm does merge $\langle (c(ab)^m, (ab)^m c) \rangle = c(ab)^m c$.

The next iteration then processes set $\{(ba)^m, c(ab)^m c\}$.

As both possible pairings do not yield an overlap concatenating is our only choice, resulting in the superstring $(ba)^m c(ab)^m c$ or $c(ab)^m c(ba)^m$ both of length $4m + 2$.

It would have been better to first place the string $(ba)^m$ between the words $c(ab)^m$ and $(ab)^m c$ leading to the optimal superstring $ca(ba)^m bc$ of length $2m + 4$.

So the approximation rate of the greedy method for this example is $\lim_{m \rightarrow \infty} \frac{4m+2}{2m+4} = 2$.

We have hence shown that for the greedy algorithm there can be no performance guarantee **better** than 2. It is widely believed that this is the actual performance guarantee but none has yet been able to prove this.

The following is proven:

Theorem

The greedy algorithm to compute a superstring is a 4-approximation algorithm for SCSP.

□

As mentioned above an optimal solution for SCSP is also an optimal solution for MCCSP. However we do not yet know any performance guarantee for the greedy algorithm wrt. maximizing compression.

Theorem

The greedy algorithm to compute a superstring is a 3-approximation algorithm for MCCSP.

Proof:

$W_0 =$ optimal superstring for input $\Sigma = \{S_1, S_2, \dots, S_n\}$

We order the S_i according to their first (left-most) appearance within W_0

↪ $(S_{i_1}, S_{i_2}, \dots, S_{i_m})$, i.e.

$W_0 = \langle \dots \langle \langle S_{i_1}, S_{i_2} \rangle, S_{i_3} \rangle \dots \rangle, S_{i_m} \rangle$



$$\text{comp}(W_0) = \sum_{m \in \text{merge}} |O_v(m)| = \sum_{1 \leq k < l} O_v(S_{i_k}, S_{i_{k+1}})$$

Greedy algorithm generates an ordering

$(S_{j_1}, S_{j_2}, \dots, S_{j_n})$ analogously.

We show: merge $\langle S_{j_k}, S_{j_{k+m}} \rangle$ can make at most three merges of optimal solution impossible.

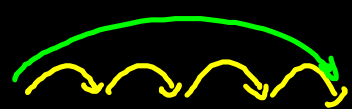
Case 1: $\langle S_{j_k}, S_{j_{k+m}} \rangle = \langle S_{i_k}, S_{i_{k+m}} \rangle, k \in [1:n-1]$

\leadsto both merges identical, i.e. no merge of optimal solution becomes impossible

Case 2: $\langle S_{j_k}, S_{j_{k+m}} \rangle = \langle S_{i_k}, S_{i_{k+m}} \rangle, m > 1$

\leadsto at most 2 merges become impossible, namely

$\langle S_{i_k}, S_{i_{k+m}} \rangle$ & $\langle S_{i_{k+m-1}}, S_{i_{k+m}} \rangle$



Case 3: $\langle S_{j_k}, S_{j_{k+m}} \rangle = \langle S_{i_{k+m}}, S_{i_k} \rangle, m \geq 1$

\leadsto at most 3 merges become impossible

$\langle S_{i_{k+m}}, S_{i_{k+m+1}} \rangle, \langle S_{i_{k-1}}, S_{i_k} \rangle$ and one

of $\langle S_{i_{k-1}}, S_{i_{k+1}} \rangle, \dots, \langle S_{i_{k+m-1}}, S_{i_{k+m}} \rangle$

Since the use of all of them would give rise to a cycle which is broken by not using one.

Notation: $M_x \hat{=}$ set of merges used by algorithm $x \in \{0, 1\}$.

$V(m) \subseteq M_0 \hat{=}$ set of merges being made impossible by merge $m = \langle S_{i_k}, S_{i_{k+1}} \rangle \in M_1$

$\hookrightarrow (\forall m \in M_1) (|V(m)| \leq 3)$.

For $ov(m), m \in M_1$, the length of the overlap we will have with merge m we have

$$(\forall m \in M_1): (ov(m) \geq \frac{1}{3} \left(\sum_{m' \in V(m)} ov(m') \right)) \quad (*)$$

Reason: greedy chooses among the remaining (possible) merges the one of maximal overlap

\hookrightarrow a merge made impossible by m never has overlap $> ov(m)$.

Furthermore:

$$(\forall m' \in M_0): (m' \in M_3) \vee \left(\overbrace{(\exists m \in M_3 \setminus M_0)}^{\checkmark}: \right. \\ \left. (m' \in V(m)) \right) \quad (**)$$

Reason: For $m' = \langle S, \bar{S} \rangle \notin M_3$ we must have merges with M_3 which merge S resp. \bar{S} with other strings. This however makes m' impossible

Consider

$$\sum_{m \in M_3} ov(m) = \sum_{m \in M_3 \cap M_0} ov(m) + \sum_{m \in M_3 \setminus M_0} ov(m)$$

We have

$$\sum_{m \in M_3 \cap M_0} ov(m) \geq \frac{1}{3} \sum_{m \in M_3 \cap M_0} ov(m) \quad \text{and}$$

$$\sum_{\substack{m \in M_3 \setminus M_0 \\ \subseteq M_3}} ov(m) \stackrel{(*)}{\geq} \sum_{m \in M_3 \setminus M_0} \frac{1}{3} \cdot \sum_{m' \in V(m)} ov(m')$$

$$\begin{aligned}
 & \overset{(\ast\ast)}{\geq} \frac{1}{3} \cdot \sum_{m' \in M_b \setminus M_g} \text{ov}(m') \\
 & \sim \exists m \in M_g \setminus M_0: m' \in V(m) \\
 \Rightarrow & \frac{\text{comp}(w_0)}{\text{comp}(w_g)} = \frac{\sum_{m' \in M_0} \text{ov}(m')}{\sum_{m \in M_g} \text{ov}(m)} \leq \frac{\sum_{m' \in M_0} \text{ov}(m')}{\frac{1}{3} \cdot \sum_{m' \in M_0} \text{ov}(m')} \\
 & = 3 \quad \square
 \end{aligned}$$

Using more elaborate reasoning one can show:

Theorem

The greedy algorithm to compute a superstring is a 2-approximation algorithm for MCCSP. □

Alternative Procedure

Idea: Do not use a single cycle to cover all vertices but a set of cycles

→ cycle-cover for which each vertex must be part of exactly one cycle.

Here: Cycle-cover of minimal cost (according to distance graph).

Remark: A minimal cycle-cover can be computed in time $O(n^3)$ (whereas the restriction to edges to be covered by at least one cycle leads to an \mathcal{NP} -complete problem).

