

## VII.2. Efficient Implementations

### Sparse Matrix Representation

Recall that we used the matrix reduction of  $\delta$  to compute the persistence homology of the filtration  $K_0 \subseteq K_1 \subseteq \dots \subseteq K_n = k$ , where  $K_i = \{s_0, \dots, s_i\}$ .

The resulting reduced matrix  $R$  consists of the information of the persistence diagram. (See VII.1)

Since the columns of  $\delta$  consist of simplices of all dimension and the boundary of a simplex is a codimension-1 chain, the matrix  $\delta$  is sparse (by convention, a matrix is "sparse" if its number of nonzero entries is approximately the number of columns / rows).

In most cases,  $R$  is also sparse. There are more efficient ways to store sparse matrices.

- ① Dictionary of keys (DOK): Store a sparse matrix by collections of dictionaries that map (row, column)-triples to the values. In the case of 0-1 matrices, only (row, column)-triples. This method is convenient for construction but poor for iteration.
- ② List of list (LIL): Store one list per column/row, with each list consisting of tuples of (row/column index, value)-pairs. In the case of 0-1 matrices, only 0-1 matrices. Typically, the inner lists are sorted by increasing/decreasing row/column indices.

### Matrix Reduction using sparse matrix representation

The boundary matrix  $\Delta_{nn}$  is given.

Step 0: Store the boundary matrix by LIL, one list per column, with row indices sorted decreasingly (the most recently added simplex is at the top). Denote each list  $\tilde{\delta}[j]$ .  $\tilde{\delta}[j] = [\text{low}_j, \dots]$

Step  $j$ : Create an empty list  $\tilde{R}[j]$ . Put the  $j$ th column of the reduced matrix  $R$  on  $\tilde{R}[j]$ , where  $i = \text{low}_j$ . If  $\text{low}_j$  undefined, do nothing. We do this by the loop

$L = \tilde{R}[j]$

while  $L \neq \text{numh}$ . &&  $\tilde{R}[L(0)] \neq \text{numh}$ , do

$L = L + R[L(0)]$  By "+" we mean merge the two lists and discard duplicate terms.

endwhile

If  $L \neq \text{numh}$ , then set  $\tilde{R}[L(0)] = L$ . endif.

The condition " $L \neq \text{numh}$ " means  $L$  is not all 0 after previous reduction. Because of the last command, " $\tilde{R}[L(0)] \neq \text{numh}$ " means there is no previous reduced column  $L'$  with  $L'(0) = L(0)$ , i.e.  $\text{low}(L') = \text{low}(L)$ .

If  $L$  is empty after some addition, then the  $j$ -th iteration is end.

If the "if" is executed, then the "while" loop is end, i.e.,

$\tilde{R}[L(0)] = \text{numh}$  is available, we set its value to be  $L$ .

As a result, the  $i$ -th array  $\tilde{R}_{ij}$  of  $\tilde{R}$  stores the  $j$ -th column of the reduced matrix  $R$ , where  $i = \text{low}(j)$ . If such  $j$  do not exist,  $\tilde{R}_{ij} = \text{numh}$ .

What's good:

- ① There are fast algorithms for sparse matrices.

② At the  $j$ -th iteration, we need not to compare  $\text{low}(j')$  with  $\text{low}(j)$  for  $j' < j$ . It suffices to check the emptiness of  $R[L(0)]$ .

Time complexity: The time complexity of the  $j$ -th iteration is  $(p+1)(j-i)$ , where  $p$  is the dimension of  $\sigma_j$ . Indeed, in the while loop,  $L(0)$  is strictly decreasing ( $\text{low}(j)$  is strictly increasing), until it reach  $i = \text{low}(j)$ . Thus the largest number of addition is  $j-i$ . (we try to store  $L$  in  $\tilde{R}_{ik}$  for some  $i \leq k < j$  every time. If it fails, an addition occur.)

At each addition, the lists added contain only nonzero terms at faces of simplices with index in  $\llbracket i, j \rrbracket$ . There are at most  $(j-i)(p+1)$  such simplices. Therefore, the time complexity is  $(j-i)^2(p+1)$ . In practice it is usually much faster than this.

Thm (Persistence equivalence theorem) Let  $\phi_i : U_i \rightarrow V_i$  be isomorphism of vector spaces. If the diagram

$$\begin{array}{ccccccc} V_0 & \rightarrow & V_1 & \rightarrow & \cdots & \rightarrow & V_{n-1} \rightarrow V_n \\ \uparrow & \uparrow & & & & & \uparrow \\ U_0 & \rightarrow & U_1 & \rightarrow & \cdots & \rightarrow & U_{n-1} \rightarrow U_n \end{array}$$

commute, then

the persistence diagram of the  $\{V_i\}$ 's are the same as the  $\{U_i\}$ 's.

Pf: Omitted.

Rmk: In practice,  $V_i = H_p(k_i)$  for some filtration  $k_i$ .  $\phi_i$  are isomorphisms commuting with  $f_{i,j}^{i,j}: H_p(k_i) \rightarrow H_p(k_j)$ . A sequence of such isomorphisms gives the equivalence of "persistence vector space", that is, a seq of vector space  $\{V_i\}$  with linear maps  $f^{i,j}: V_i \rightarrow V_j$  s.t. the composition law holds.

## Surfaces

Let  $K$  be a simplicial complex that triangulates a 2-manifold.

Recall that we can define a piecewise linear function  $f: |K| \rightarrow \mathbb{R}$  by specifying the values on the vertices  $f(v_1) < f(v_2) < \dots < f(v_n)$  and extend to the whole  $|K|$  using barycentric coordinates.  $f(x) = \sum_i b_i(x) f(v_i)$ .

Then we obtain the lower star filtration of  $f$ .  $K_0 \subseteq K_1 \subseteq \dots \subseteq K_n = K$ ,

where  $K_i = \bigcup_{\substack{\text{f}(v_j) \leq i \\ \text{f}(v_k) > i}} \overline{s(v_i)}$ . If we define  $g: K \rightarrow \mathbb{R}$  by  $g(v) = \max_{x \in v} f(x)$ , then

for  $f(v_i) \leq a < f(v_{i+1})$ , we have  $K_i = g^{-1}(-\infty, a]$ . Let  $K_a := g^{-1}(-\infty, a]$ .

As discussed in VI.3.  $|K|$  is a deformation retract to  $f^{-1}(-\infty, a]$ ;

Also observe that  $|K_a| \subseteq |K|_a := f^{-1}(-\infty, a)$  for all  $a \in \mathbb{R}$ .

Hence the map  $H_p(K_a) \xrightarrow{\sim} H_p(|K|_a)$  induced by inclusion is an isomorphism.

The commuting diagram  $H_p(|K|_a) \rightarrow H_p(|K|_b)$ ,  $a \leq b$  together with

$$\begin{array}{ccc} H_p(K_a) & \longrightarrow & H_p(K_b) \\ \uparrow & & \uparrow \\ H_p(K_a) & \longrightarrow & H_p(K_b) \end{array}$$

the persistence equivalence theorem implies that calculating the persistence homology of the lower star filtration of  $K$  is equivalent to that of  $f$ .

## Zeroth Diagram

There is a much more efficient method to computing the 0-th diagram. Given a 1-simplex  $\sigma$ . Suppose  $\partial\sigma = u + v$  and  $u < v$ . If  $u$  and  $v$  are already in the same connected component, then the addition of  $\sigma$  does not affect the zeroth diagram. Otherwise, it kills a 0-cycle that is

Created by the younger one of the oldest simplex of the two connected components.

Thus all we need to know is whether  $u$  and  $v$  belong to different connected components, and if they do, the oldest simplices of each component. This is exactly the Kruskal's algorithm for disjoint set systems we introduced in chapter I.1. A simplicial complex with only 0,1 simplices is just a simple graph!

Recall in section I.1. when  $u,v$  are not in the same connected component, we point the smaller component in the larger one. However in this case, we point the younger component to the older one because the younger component is the one being killed.

This algorithm has time complexity  $O(n^{\alpha})$  after doing parallelization, where  $\alpha$  is the inverse of the Ackermann function, which is a constant in practice.

### First diagram

If  $K$  is a triangulation of a 2-manifold (without boundary),  $f: K \rightarrow \mathbb{R}$  be an increasing function (i.e.  $0 \leq i \Rightarrow f(0) \leq f(i)$ ), we consider  $-f: K \rightarrow \mathbb{R}$ .

For  $a \in \mathbb{R}$ ,  $(-f)^{-1}(-\infty, -a] = K - f^{-1}(-\infty, a] = (f^{-1}(-\infty, a))^c$ . Let  $\tilde{K}_i = (-f)^{-1}(-\infty, -a_i)$ .

We obtain  $\tilde{K}_0 = K \supseteq \tilde{K}_1 \supseteq \dots \supseteq \tilde{K}_n = \emptyset$ . When adding a simplex  $\sigma$  to  $K_i$ ,

$K_{i+1} = K_i \cup \{\sigma\}$ .  $\sigma$  is deleted from  $\tilde{K}_i$  and  $\tilde{K}_i = \tilde{K}_{i+1} \cup \{\sigma\}$ .

Each  $\tilde{K}_i$  is not a simplicial complex, but we can still consider the (singular) homology of its underlying space, that is  $|K| - |K_i| \subseteq \mathbb{R}^n$ .

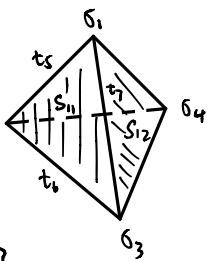
E.g.  $K = \{\sigma_1, \dots, \sigma_4, t_5 = (\sigma_1 \sigma_2), t_6 = (\sigma_2 \sigma_3), t_7 = (\sigma_1 \sigma_3), \dots, t_{10}, S_{11} = (\sigma_1 \sigma_2 \sigma_3), \dots, S_{14}\}$ .

a triangulation of  $S^2$ . The point  $(7, 11)$  is on the 1-st persistence diagram of the filtration  $K_0 \subseteq \dots \subseteq K_n$ . I.e. the cycle  $t_5 + t_6 + t_7$  is

created at  $K_7$  and dies at  $K_{11}$ . Now consider  $\tilde{K}_{11} = \{S_{11}, S_{12}, S_{13}, S_{14}\}$ .

having 4 disjoint 2-simplices, which we view as connected components.

$S_{11}$  is the one appear with the addition of  $S_{11}$ .  $S_{11}$  stays isolated until



the addition of  $t_7$  in  $\tilde{K}_7 = \{t_7, \dots, t_{10}, s_{11}, \dots, s_{14}\}$ , which connect  $S_{11}$  to  $S_{12}$ . Thus a zero cycle in  $\tilde{K}_n \subseteq \tilde{K}_{n-1} \subseteq \dots \subseteq \tilde{K}_0$  is created by  $s_{11}$  and killed by  $t_7$ . In other words,  $(-1, -7) \in Dgm_1(f)$  ( $\tilde{K}_i = (-f)^i(-\infty, -a_i)$ )

In general,  $(a, b) \in Dgm_1(f) \Leftrightarrow (-b, -a) \in Dgm_0(-f)$ . The rigorous construction and proof will be treated next chapter.

Therefore, the computation of  $Dgm_1(f)$  is reduced to that of  $Dgm_0(-f)$ , which can be done by Knuscal's algorithm whose time complexity is  $O(n) \approx \text{const.}$