

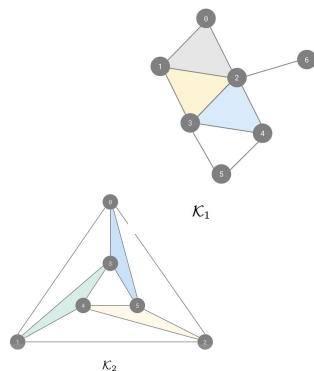
# Approximating Simplex Frequency Distribution for Simplicial Complexes

Hamid Beigy<sup>1</sup>, Mohammad Mahini<sup>2</sup>, Salman Qadami<sup>3</sup>, and Morteza Saghafian<sup>4</sup>

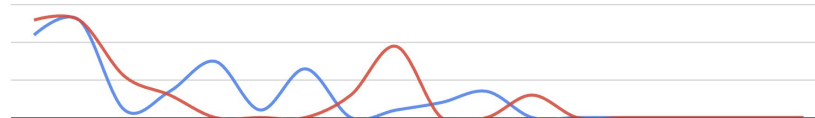
- 1 Sharif University of Technology  
beigy@sharif.edu
- 2 Sharif University of Technology  
m\_mahini@ce.sharif.edu
- 3 Amirkabir University of Technology  
salmanqadami@gmail.com
- 4 Institute of Science and Technology Austria  
morteza.saghafian@ist.ac.at



# Approximating Simplex Frequency Distribution for Simplicial Complexes



Simplex Type	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
$ \mathcal{S}_{\mathcal{K}_1} $	10	12	1	3	7	1	6	0	1	2	3	0	0	0	0	0	0	0
$ \phi_{\mathcal{K}_1} $	0.22	0.26	0.02	0.07	0.15	0.02	0.13	0	0.02	0.04	0.07	0	0	0	0	0	0	0
$ \mathcal{S}_{\mathcal{K}_2} $	12	12	5	3	0	0	0	3	9	0	0	3	0	0	0	0	0	0
$ \phi_{\mathcal{K}_2} $	0.26	0.26	0.11	0.06	0	0	0	0.06	0.19	0	0	0.06	0	0	0	0	0	0

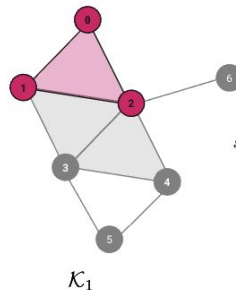


Creating a vector from a Simplicial Complex (SC)  
 Based on its **Small Building Blocks (Simplexes)**

That can be helpful for ML applications such as classification.

# Simplet

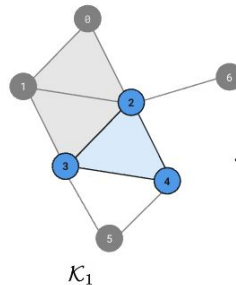
- Simplets are **small induced connected sub-complexes** of a massive complex that appear at any frequency.
- Every simplet can be identified by its vertices.
- *Simplet Types* are isomorphic classes of simplets.
- We denoted  $\mathcal{S}_{\mathcal{K}}(i)$  as a set of all simplets of type  $i$  in  $\mathcal{K}$ .



$$s_1 = \{(0), (1), (2),$$

(vertices)  
(0, 1), (0, 2), (1, 2),  
(edges)  
(0, 1, 2)\}

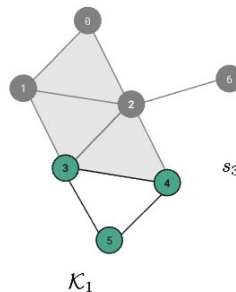
(triangles)



$$s_2 = \{(2), (3), (4),$$

(vertices)  
(2, 3), (2, 4), (3, 4),  
(edges)  
(2, 3, 4)\}

(triangles)



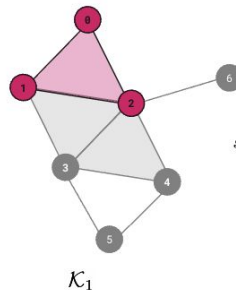
$$s_3 = \{(3), (4), (5),$$

(vertices)  
(3, 4), (3, 5), (4, 5)\}

(edges)

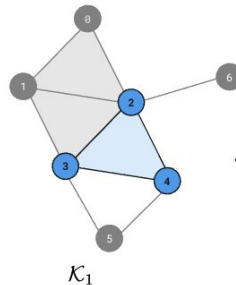
# Simplet

- Simplets are **small induced connected sub-complexes** of a massive complex that appear at any frequency.
- Every simplet can be identified by its **vertices**.
- *Simplet Types* are isomorphic classes of simplets.
- We denoted  $\mathcal{S}_{\mathcal{K}}(i)$  as a set of all simplets of type  $i$  in  $\mathcal{K}$ .



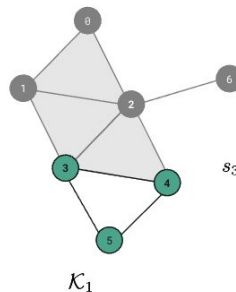
$$s_1 = \{(0), (1), (2), \\ (0, 1), (0, 2), (1, 2), \\ (0, 1, 2)\}$$

(vertices)  
(edges)  
(triangles)



$$s_2 = \{(2), (3), (4), \\ (2, 3), (2, 4), (3, 4), \\ (2, 3, 4)\}$$

(vertices)  
(edges)  
(triangles)

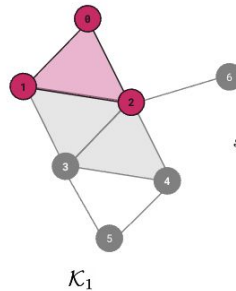


$$s_3 = \{(3), (4), (5), \\ (3, 4), (3, 5), (4, 5)\}$$

(vertices)  
(edges)

# Simplet

- Simplets are **small induced connected sub-complexes** of a massive complex that appear at any frequency.
- Every simplet can be identified by its vertices.
- **Simplet Types** are isomorphic classes of simplets.
- We denoted  $\mathcal{S}_{\mathcal{K}}(i)$  as a set of all simplets of type  $i$  in  $\mathcal{K}$ .



$$s_1 = \{(0), (1), (2),$$

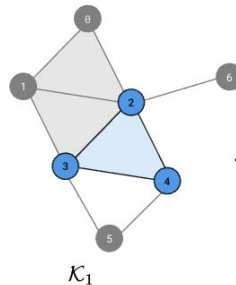
(vertices)

$$(0, 1), (0, 2), (1, 2),$$

(edges)

$$(0, 1, 2)\}$$

(triangles)



$$s_2 = \{(2), (3), (4),$$

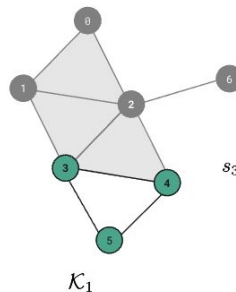
(vertices)

$$(2, 3), (2, 4), (3, 4),$$

(edges)

$$(2, 3, 4)\}$$

(triangles)



$$s_3 = \{(3), (4), (5),$$

(vertices)

$$(3, 4), (3, 5), (4, 5)\}$$

(edges)

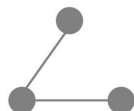
# Simplet Types

- Simplet types with two to four vertices

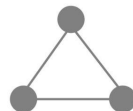
- Example:



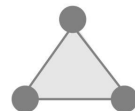
1



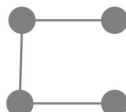
2



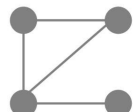
3



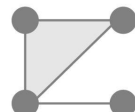
4



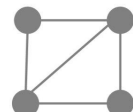
5



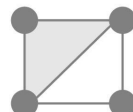
6



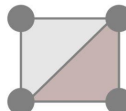
7



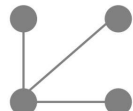
8



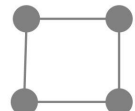
9



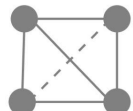
10



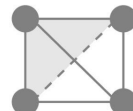
11



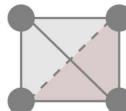
12



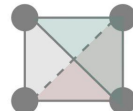
13



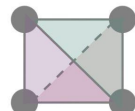
14  
1 triangle



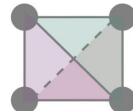
15  
2 triangles



16  
3 triangles



17  
4 triangles

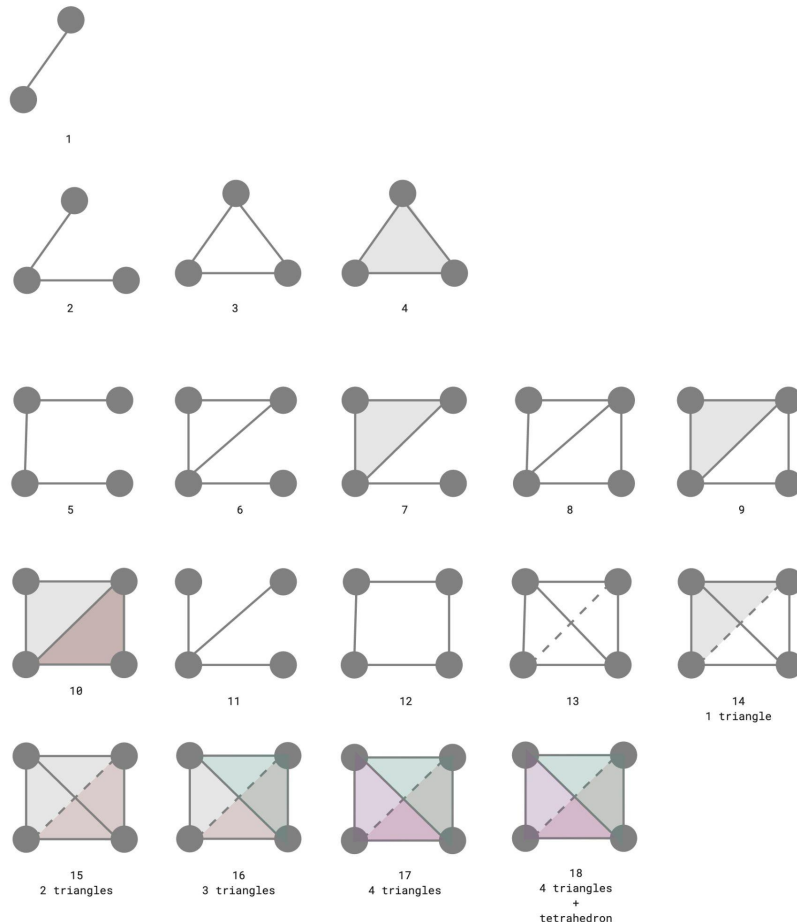
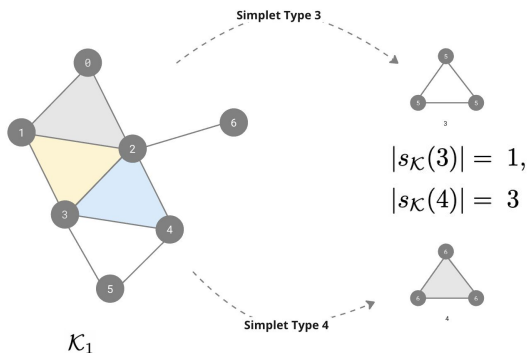


18  
4 triangles  
+  
tetrahedron

# Simplet Types

- Simplet types with two to four vertices

- Example:



# Simplet Frequency Distribution (SFD) Vector

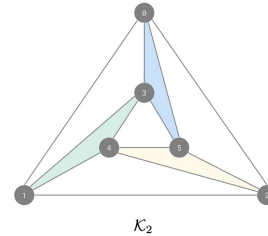
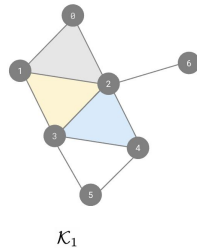
- The relative frequencies of various simplet types in  $\mathcal{K}$
- The frequency denoted by  $\phi_{\mathcal{K}}(i)$  is obtained by dividing  $|\mathcal{S}_{\mathcal{K}}(i)|$  by  $\sum_{j=1}^{N_m} |\mathcal{S}_{\mathcal{K}}(j)|$
- The vector  $(\phi_{\mathcal{K}}(1), \dots, \phi_{\mathcal{K}}(N_m))$  is called the SFD vector of the  $\mathcal{K}$
  
- Example:



# Simplet Frequency Distribution (SFD) Vector

- The relative frequencies of various simplet types in  $\mathcal{K}$
- The frequency denoted by  $\phi_{\mathcal{K}}(i)$  is obtained by dividing  $|\mathcal{S}_{\mathcal{K}}(i)|$  by  $\sum_{j=1}^{N_m} |\mathcal{S}_{\mathcal{K}}(j)|$
- The vector  $(\phi_{\mathcal{K}}(1), \dots, \phi_{\mathcal{K}}(N_m))$  is called the SFD vector of the  $\mathcal{K}$

- Example:



Simplet Type	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
$ \mathcal{S}_{\mathcal{K}_1} $	10	12	1	3	7	1	6	0	1	2	3	0	0	0	0	0	0	0
$ \phi_{\mathcal{K}_1} $	0.22	0.26	0.02	0.07	0.15	0.02	0.13	0	0.02	0.04	0.07	0	0	0	0	0	0	0
$ \mathcal{S}_{\mathcal{K}_2} $	12	12	5	3	0	0	0	3	9	0	0	3	0	0	0	0	0	0
$ \phi_{\mathcal{K}_2} $	0.26	0.26	0.11	0.06	0	0	0	0.06	0.19	0	0	0.06	0	0	0	0	0	0

# Calculating The SFD Vector for Large SCs

- Simple Simplet counting algorithm is in  $\Theta(n^k)$

Our approach:

- Instead of calculating the exact counts we use an approximation on Simplet frequencies.
- Our algorithm is sublinear in the size of  $\mathcal{K}$

# Calculating The SFD Vector for Large SCs

- Simple Simplet counting algorithm is in  $\Theta(n^k)$

Our approach:

- Instead of calculating the exact counts we use an approximation on Simplet frequencies.
- Our algorithm is sublinear for large and sparse SCs in the size of  $\mathcal{K}$

# Approximating The SFD Vector with Sampling

1. **[Number of Samples]** With a set of  $\frac{c}{\epsilon^2}(1 + \ln \frac{1}{\delta})$  simplexes sampled uniformly from SC  $\mathcal{K}$ , we can have an  $(\epsilon, \delta)$ -approximation on the SFD vector of  $\mathcal{K}$ .
2. **[Sampling Algorithm]** We propose a uniform sampling algorithm for simplexes in a connected simplicial complex that find a sample with complexity in  $O(\log(n) \cdot \Delta \cdot \text{diam}(\mathcal{K})^2)$ .

The time complexity of  $(\epsilon, \delta)$ -approximation of SFD vector of  $\mathcal{K}$  is

$$O\left(\frac{1}{\epsilon^2} \cdot (1 + \ln \frac{1}{\delta}) \cdot \log(n) \cdot \Delta \cdot \text{diam}(\mathcal{K})^2\right)$$

# Approximating The SFD Vector with Sampling

1. Number of Samples With a set of  $\frac{c}{\epsilon^2}(1 + \ln \frac{1}{\delta})$  simplets sampled uniformly from SC  $\mathcal{K}$ , we can have an  $(\epsilon, \delta)$ -approximation on the SFD vector of  $\mathcal{K}$ .
2. Sampling Algorithm We propose a uniform sampling algorithm for simplets in a connected simplicial complex that find a sample with complexity in  $O(\log(n) \cdot \Delta \cdot \text{diam}(\mathcal{K})^2)$ .

The time complexity of  $(\epsilon, \delta)$ -approximation of SFD vector of  $\mathcal{K}$  is

$$O\left(\frac{1}{\epsilon^2} \cdot \left(1 + \ln \frac{1}{\delta}\right) \cdot \log(n) \cdot \Delta \cdot \text{diam}(\mathcal{K})^2\right)$$

# Number of Samples

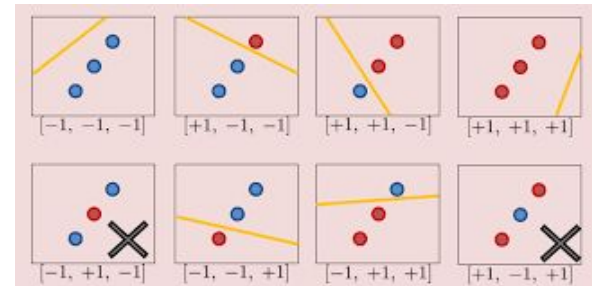
[Number of Samples] With a set of  $\frac{c}{\epsilon^2} (1 + \ln \frac{1}{\delta})$  simplets sampled uniformly from SC  $\mathcal{K}$ , we can have an  $(\epsilon, \delta)$ -approximation on the SFD vector of  $\mathcal{K}$ .

$$\frac{c}{\epsilon^2} \left( 1 + \ln \frac{1}{\delta} \right)$$

# Number of Samples

[Number of Samples] With a set of  $\frac{c}{\epsilon^2}(1 + \ln \frac{1}{\delta})$  simplets sampled uniformly from SC  $\mathcal{K}$ , we can have an  $(\epsilon, \delta)$ -approximation on the SFD vector of  $\mathcal{K}$ .

- We use VC dimension to prove this bound.
- For a domain  $D$  and collection  $\mathcal{R}$  of subsets of  $D$ , the VC dimension represents the maximum size of a set  $X \subseteq D$  that can be shattered by  $\mathcal{R}$  which means  $\{r \cap X | \forall r \in R\} = 2^{|X|}$ .



# Number of Samples

[Number of Samples] With a set of  $\frac{c}{\epsilon^2}(1 + \ln \frac{1}{\delta})$  simplets sampled uniformly from SC  $\mathcal{K}$ , we can have an  $(\epsilon, \delta)$ -approximation on the SFD vector of  $\mathcal{K}$ .

- **VC Dimension of Simplets** Let  $\mathcal{R} = \{\mathcal{S}_i \mid 1 \leq i \leq N_m\}$  be a family of all simplet sets where  $N_m$  is the number of simplet types with at most  $m$  vertices, and  $D$  is all simplets of SC  $\mathcal{K}$ , Then we have  $VC(D, \mathcal{R}) = 1$ .

Proof. Let  $\{s_1, s_2\} \subseteq D$

- If  $s_1$  and  $s_2$  are belong to the same simplet type,  $\{s_1\}$  can't be shattered.
- Otherwise,  $\{s_1, s_2\}$  can't be shattered.



# Number of Samples

$$\frac{c}{\epsilon^2} \left(1 + \ln \frac{1}{\delta}\right)$$

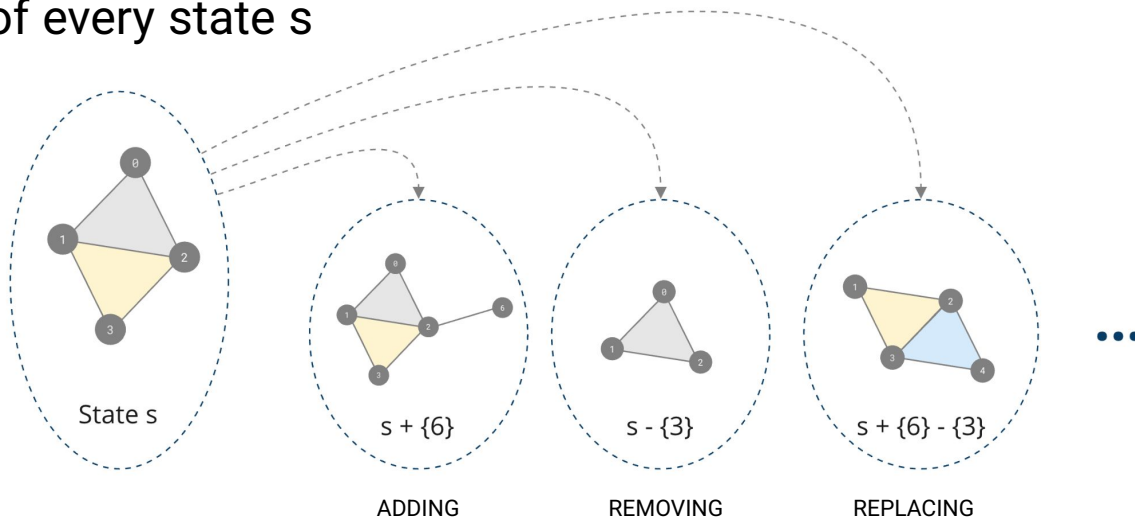
[Number of Samples] With a set of  $\frac{c}{\epsilon^2} \left(1 + \ln \frac{1}{\delta}\right)$  simplets sampled uniformly from SC  $\mathcal{K}$ , we can have an  $(\epsilon, \delta)$ -approximation on the SFD vector of  $\mathcal{K}$ .

- **VC Dimension of Simplets** Let  $\mathcal{R} = \{\mathcal{S}_i \mid 1 \leq i \leq N_m\}$  be a family of all simplet sets where  $N_m$  is the number of simplet types with at most  $m$  vertices, and  $D$  is all simplets of SC  $\mathcal{K}$ , Then we have  $VC(D, \mathcal{R}) = 1$ .

**[Lemma]** For a domain  $D$  and collection  $\mathcal{R}$  of subsets of  $D$ , with  $VC(D, \mathcal{R}) \leq d$  and using  $\frac{c}{\epsilon^2} \left(d + \ln \frac{1}{\delta}\right)$  uniform samples, we can have an  $(\epsilon, \delta)$ -approximation on distribution of all subsets in  $\mathcal{R}$ .

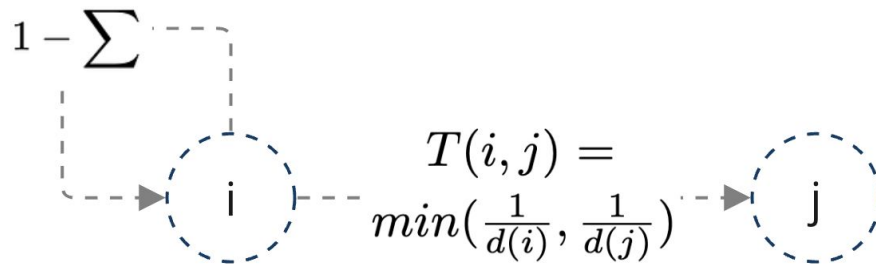
# Sampling Algorithm

- A Monte-Carlo Markov-Chain algorithm.
- Random walk on a directed graph  $\mathcal{P}_{\mathcal{K}}^m$  whose vertex set (states) is a set of all simplexes in complex  $\mathcal{K}$
- Out-neighbors of every state  $s$



# Sampling Algorithm

- $T(i, j)$  is Transition probability matrix  $T$  on  $\mathcal{P}_K^m$



- The random walk is
  - Irreducible
  - Aperiodic
  - Converges to the uniform stationary distribution

# Sampling Algorithm (Time Complexity)

- The mixing time of the markov chain on  $\mathcal{P}_{\mathcal{K}}^m$  is in

$$O(\log(n) \cdot \Delta \cdot \text{diam}(\mathcal{K})^2)$$

Proof.

- $\Delta(\mathcal{P}_{\mathcal{K}}^m) \in O(m^2 \cdot \Delta)$
- $\text{diam}(\mathcal{P}_{\mathcal{K}}^m) = \text{diam}(\mathcal{K}) + m$

( $m$  is the maximum number of simpleton vertices; and is constant in size of SC)

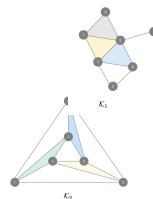
[Lemma] The mixing time  $t_{mix}^G$  of a random walk on graph  $G$  with  $n$  vertices is in

$$O(\log(n) \cdot \Delta(G) \cdot \text{diam}(G)^2)$$

# Approximating Simplex Frequency Distribution

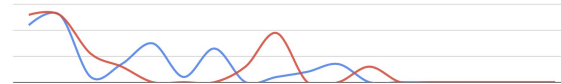
for  
Simplicial Complexes

Creating a vector from a Simplicial Complex (SC)  
using **local structures (Simplexes)**



## Conclusion

Simplex Type	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
$ SC_1 $	10	12	1	3	7	1	6	0	1	2	3	0	0	0	0	0	0	0
$ \phi_{SC_1} $	0.22	0.26	0.02	0.07	0.15	0.02	0.13	0	0.02	0.04	0.07	0	0	0	0	0	0	0
$ SC_2 $	12	12	5	3	0	0	0	3	9	0	0	3	0	0	0	0	0	0
$ \phi_{SC_2} $	0.26	0.26	0.11	0.06	0	0	0	0.06	0.19	0	0	0.06	0	0	0	0	0	0



$(\epsilon, \delta)$  - approximation  
algorithm



Upper-bound on  
the # of samples

$$\frac{c}{\epsilon^2} \left(1 + \ln \frac{1}{\delta}\right)$$



Uniform simplex  
sampling using random  
walk

$$O(\log(n) \cdot \Delta \cdot \text{diam}(\mathcal{K})^2)$$

**Sub-linear for Real World SCs**

# Future Directions

1. [The SFD vector for Simplexes]  $S_{\mathcal{K}}(s, i)$

$$SFD(s) = (\phi_{\mathcal{K}}(s, 1), \dots, \phi_{\mathcal{K}}(s, N_m))$$

2. [Centrality measure for simplexes]

$$\sum w_i \times \phi_{\mathcal{K}}(s, i)$$

3. [Alpha Complexes] - Filtering
4. [Simplicial Complex Similarity Metric]
5. [Classification Applications]



**Thanks for your attention**

