

Coordinated Sampling

Ok, so lets sample..

- Lets start with sampling 1 element..
- We get a stream of items want to stay with 1 uniformly random



Ok, so lets sample..

- Lets start with sampling 1 element..
- We get a stream of items want to stay with 1 uniformly random

1

Sample:

1

Ok, so lets sample..

- Start with 1 element..
- We get a stream of items want to stay with 1 uniformly random

1 2

Sample: 1

Replace with probability $1/2$

Ok, so lets sample..

- Start with 1 element..
- We get a stream of items want to stay with 1 uniformly random

1 2 3

Sample: 1

Replace with probability $1/3$

Ok, so lets sample..

- Start with 1 element..
- We get a stream of items want to stay with 1 uniformly random

1 2 3

Sample: 3

Replace with probability $1/3$

Reservoir sampling (Vitter 85)

- Start with 1 element..
- We get a stream of items want to stay with 1 uniformly random

1 2 3 4

Sample: 3

Replace with probability $1/4$

If an item appears more than once ?



Min hash



Choose at random a hash function h from some family H

Your sample is the item x with $\min h(x)$

Need h to be min-wise independent..

“Permanent random numbers”

We shall assume that $h(x) \in U[0,1]$

$h(x_i)$ are independent

Ok, so lets **sample k items**

- **k-mins:** Repeat the experiment k times with h_1, h_2, \dots, h_k
- **Bottom-k:** Take the k smallest

Ok, so lets **sample k items**

- **k-mins:** Repeat the experiment k times with h_1, h_2, \dots, h_k **Sampling with replacement**

- **Bottom-k:** Take the k smallest **Sampling without replacement**

Develop estimators

Develop estimators

- For the size **of the set = # of distinct items ?**
- For the size of selected subsets (say, the number of **green** items) ?

Estimators (# greens: $|G|$)

- Let n = #distinct items (assume it is known)
- Estimate the total # of green items to be $g \cdot (n/k)$

Adjusted weights (HT 52)

- Let p ($=k/n$) be the probability that item x is sampled
- Let the adjusted weight of x be

$$a(x) = \begin{cases} \frac{1}{p} = \frac{n}{k} & \text{if } x \text{ is sampled} \\ 0 & \text{otherwise} \end{cases} \quad E[a(x)] = 1$$

- Estimate $|G|$ by

$$g \frac{n}{k} = \sum_{x \in S \cap G} \frac{1}{p} = \sum_{x \in G} a(x)$$

Adjusted weights (HT 52)

- Let p ($=k/n$) be the probability that item x is sampled
- Let the adjusted weight of x be

$$a(x) = \begin{cases} \frac{1}{p} = \frac{n}{k} & \text{if } x \text{ is sampled} \\ 0 & \text{otherwise} \end{cases} \quad E[a(x)] = 1$$

- Unbiased:

$$E\left[\frac{n}{k}\right] = E\left[\sum_{x \in S \cap G} \frac{1}{p}\right] = E\left[\sum_{x \in G} a(x)\right] = \sum_{x \in G} E[a(x)] = |G|$$

Caveats

- We do not really know n = #distinct items
- So we do not know the adjusted weight (n/k)
- In fact, often we want to estimate n itself ?

Terminology

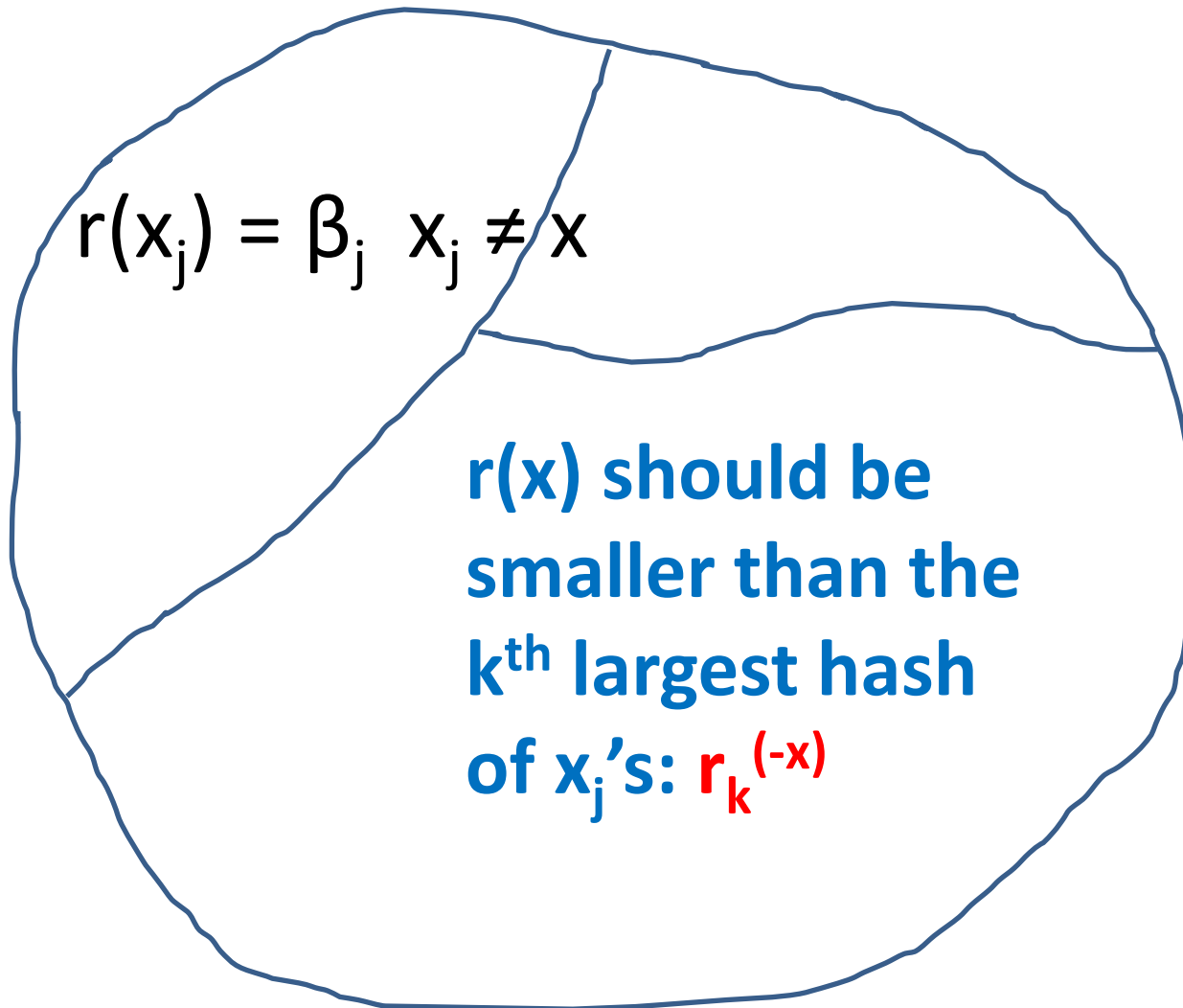
- When items are unweighted I refer to $h(x)=r(x)$ as the rank of x

Partition the sample space

Partition the sample space

- Recall: a “point” in our sample space is an assignment of ranks $\in U[0,1]$ to the elements
- Lets condition on the ranks assigned to all elements but **x**
- **What is the probability that we sample **x**?**

Partition the sample space



Conditioned adjusted weights

- Let p be the probability that item x is sampled conditioned on $r(x_j) = \beta_j$, $x_j \neq x$
- $\mathbf{p} = \mathbf{r}_k^{(-x)}$

$$a(x) = \begin{cases} \cancel{1/p} = \cancel{1/r_k^{(-x)}} & \text{if } x \text{ is sampled} \\ 0 & \text{otherwise} \end{cases} \quad E[a(x)] = 1$$

- Estimate $|G|$ by

$$g \frac{1}{r_k^{(-x)}} = \sum_{x \in S \cap G} \frac{1}{p} = \sum_{x \in G} a(x)$$

Partition the sample space

- How do we know $r_k^{(-x)}$?
- Lets keep with the sample the $k+1$ smallest hash value, r_{k+1}
- $r_k^{(-x)} = r_{k+1}$ if x is in the sample

Conditioned adjusted weights

- Let p be the probability that item x is sampled conditioned on $r(x_j) = \beta_j$, $x_j \neq x$

- $\mathbf{p} = \mathbf{r}_{k+1}$

$$a(x) = \begin{cases} \frac{1}{p} = \frac{1}{r_{k+1}} & \text{if } x \text{ is sampled} \\ 0 & \text{otherwise} \end{cases} \quad E[a(x)] = 1$$

- Estimate $|G|$ by

$$g \frac{1}{r_k} = \sum_{x \in S \cap G} \frac{1}{p} = \sum_{x \in G} a(x)$$

Coordinated sampling

Coordinated sampling

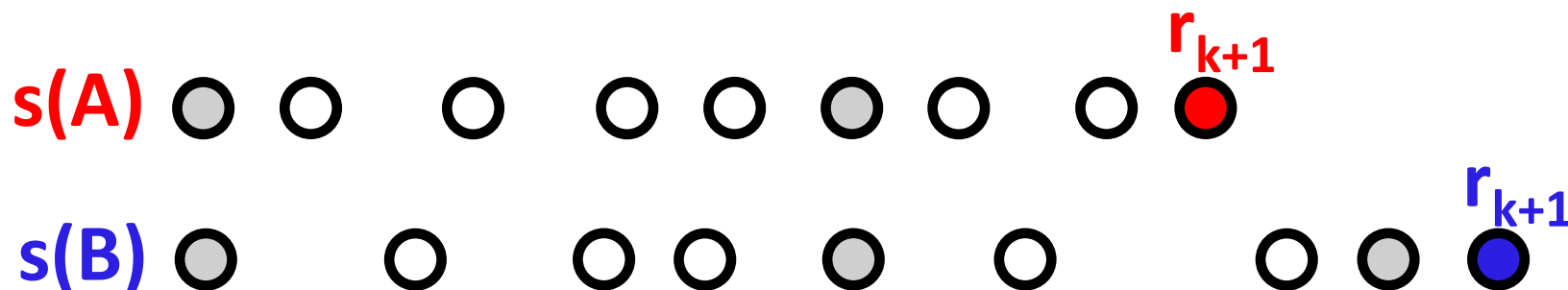
- Say we have two sets **A** and **B**
- We have a bottom-k sketches **s(A)** of **A** and **s(B)** of **B**, prepared with the same hash function !
- Want to estimate $|\mathbf{A} \cap \mathbf{B}|$?

Multiple set scenario

item	Router 1	Router 2
132.169.1.1	✓	
132.66.235.47		✓
157.166.238.17	✓	✓
128.112.132.86	✓	
170.149.172.130		
72.52.4.119	✓	✓
192.115.76.44		
128.139.199.7		✓
107.23.224.136	✓	✓

Coordinated sampling

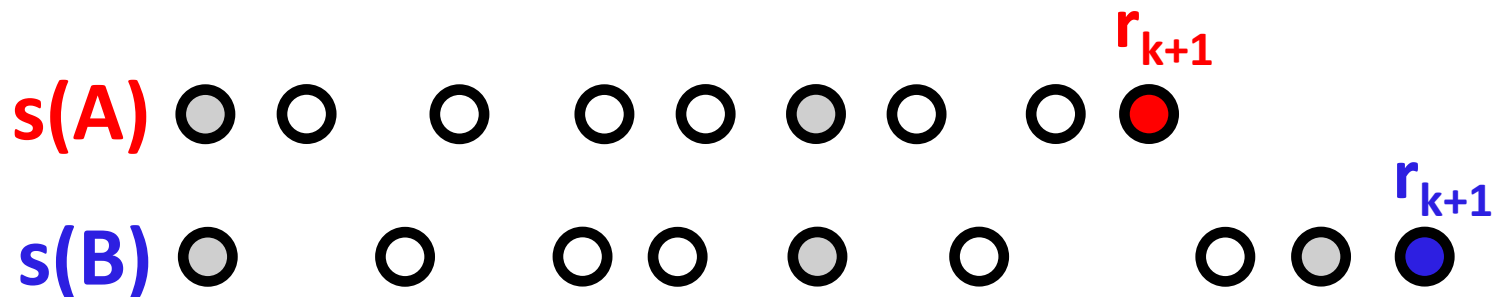
Give adjusted weights to the elements in $A \cap B$:



$$a(x) = \begin{cases} \frac{1}{p} & \text{if } x \in s(A) \cap s(B) \\ 0 & \text{otherwise} \end{cases}$$

p is the probability that $x \in s(A) \cap s(B)$

Estimating $A \cap B$ with coordination

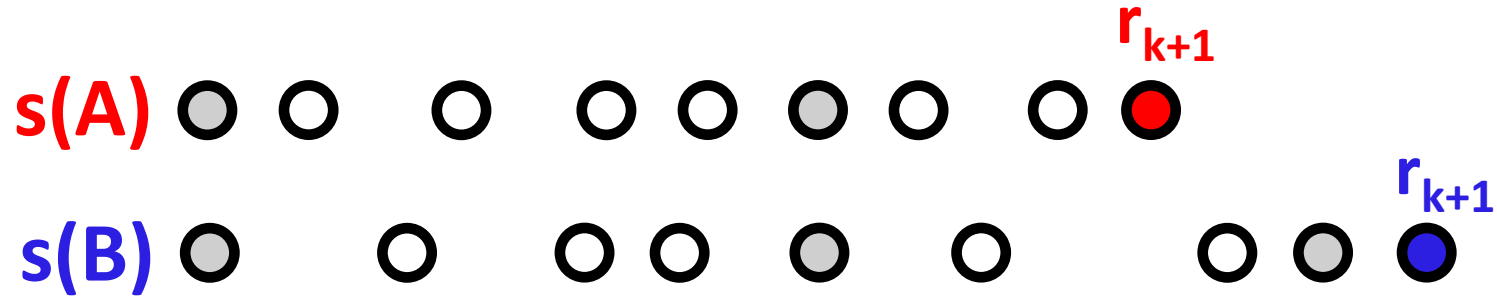


$$a(x) = \begin{cases} \frac{1}{p} & \text{if } x \in s(A) \cap s(B) \\ 0 & \text{otherwise} \end{cases}$$

p is the probability that $x \in s(A) \cap s(B)$
conditioned on all ranks of $y \neq x$

What is p ??

Estimating $A \cap B$ with coordination



$$a(x) = \begin{cases} \frac{1}{p} & \text{if } x \in s(A) \cap s(B) \\ 0 & \text{otherwise} \end{cases}$$

$$p = \min\{r_{k+1}, r_{k+1}\}$$

rather than $p = r_{k+1} \cdot r_{k+1}$ without coordination !

Homework

- Develop an estimate for $|\mathbf{B}-\mathbf{A}|$?

Weighted items

- Each item has a weight $w(x)$
- Want to estimate total weight of a selected subset: $w(G)$ =sum of the weights of the greens

Multiple set scenario

item	Router 1	Router 2
132.169.1.1	30	
132.66.235.47		40
157.166.238.17	100	100
128.112.132.86	1000	1000
170.149.172.130		
72.52.4.119	440	440
192.115.76.44		
128.139.199.7		515
107.23.224.136	w	w

Weighted items

- Assign adjusted weights:

$$a(x) = \begin{cases} \frac{w(x)}{p} = \frac{w(x)}{r_{k+1}} & \text{if } x \text{ is sampled} \\ 0 & \text{otherwise} \end{cases} \quad E[a(x)] = w(x)$$

- Estimate $w(G)$ =weight of greens by

$$\sum_{x \in S \cap G} \frac{w(x)}{r_{k+1}} = \sum_{x \in S \cap G} \frac{1}{p} = \sum_{x \in G} a(x)$$

Weighted items

- Assign adjusted weights:

$$a(x) = \begin{cases} \frac{w(x)}{p} = \frac{w(x)}{r_{k+1}} & \text{if } x \text{ is sampled} \\ 0 & \text{otherwise} \end{cases} \quad E[a(x)] = w(x)$$

- This is an unbiased estimator

$$E\left[\sum_{x \in S \cap G} \frac{w(x)}{r_{k+1}}\right] = E\left[\sum_{x \in S \cap G} \frac{1}{p}\right] = E\left[\sum_{x \in G} a(x)\right] = \sum_{x \in G} E[a(x)] = \sum_{x \in G} w(x)$$

HT Variance

$$\text{Var}[a(x)] = p \left(\frac{w(x)}{p} \right)^2 - w^2(x) = w^2(x) \left(\frac{1}{p} - 1 \right)$$

- Want p to be large
- In particular if $w(x)$ is large

Weighted items

➔ Draw the random rank of x proportionally to $w(x)$

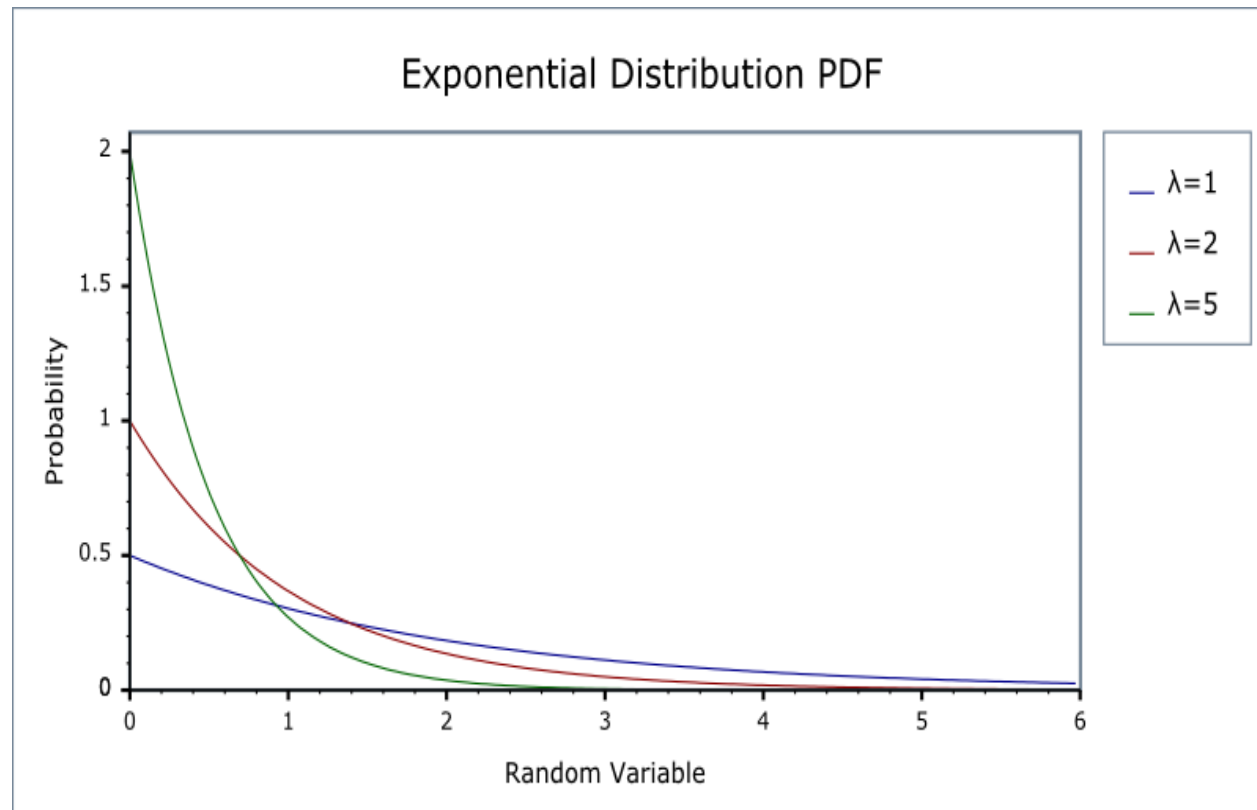
There are different ways to do it

Priority ranks

- Set $r(x) = h(x)/w(x)$, $h(x) \in U[0,1]$

Draw from: $\text{Exp}(w)$

- PDF: $w e^{-wx}, x \geq 0$
- CDF: $1 - e^{-wx}$
- $\mu = \sigma = \frac{1}{w}$



Exp(w) ranks

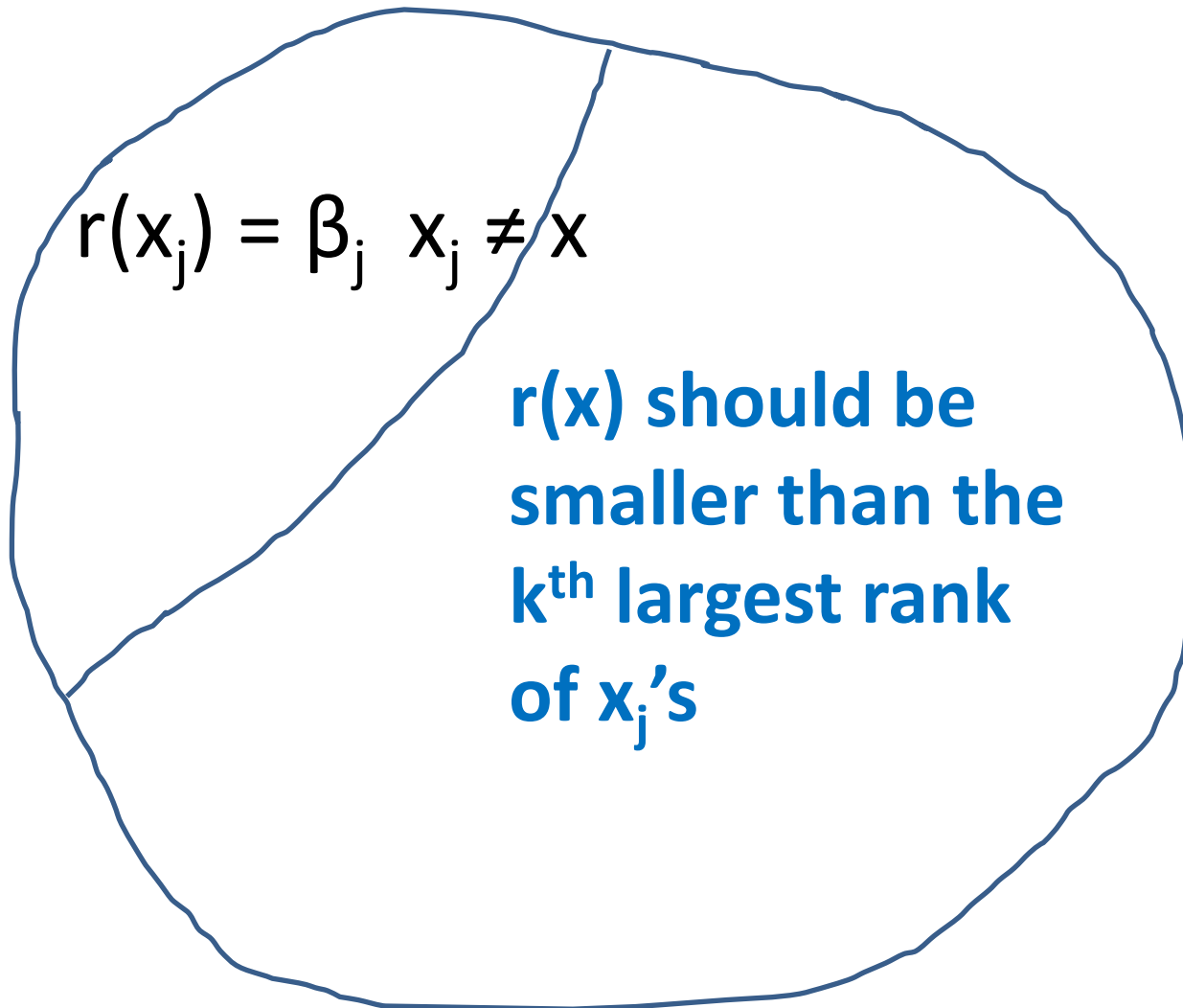
- Larger weight \rightarrow smaller rank
- Equivalent to weighted sampling with replacement
- Can draw using the following:

$$h(x) \in U[0,1], \quad r(x) = -\frac{\ln(1-h(x))}{w} \sim \text{Exp}(w)$$

Partition the sample space

- A “point” in our sample space is an assignment of ranks $\in U[0,1]/w(x)$ to each element x
- Lets condition on the ranks assigned to all elements but one, say x
- **What is the probability that we sample x ?**

Partition the sample space



Weighted items

- Adjusted weights for weighted sampling

$$a(x) = \begin{cases} \frac{w(x)}{p(h(x) < r_{k+1} w(x))} & \text{if } x \text{ is sampled} \\ 0 & \text{otherwise} \end{cases}$$

- Estimate $w(G)$ =weight of greens by

$$\sum_{x \in S \cap G} \frac{w(x)}{p(h(x) < r_{k+1} w(x))} = \sum_{x \in S \cap G} \frac{1}{p} = \sum_{x \in G} a(x)$$

Multiple weights

Multiple set scenario

item	Router 1	Router 2	Router 3
132.169.1.1	30		200
132.66.235.47		40	
157.166.238.17	100	200	90
128.112.132.86	1000	500	
170.149.172.13	102		9999
72.52.4.119		440	450
192.115.76.44			330
128.139.199.7		515	111
107.23.224.136	w_1	w_2	w_3

Weighted sampling and Coordination

- We want the samples to be **weighted**
- We also want **coordination**, that is if we sampled x at router 1 then x is more likely to be sampled also at router 2
- How could we achieve this ?

Coordinated priority ranks

$$h(x) \sim U[0, 1]$$

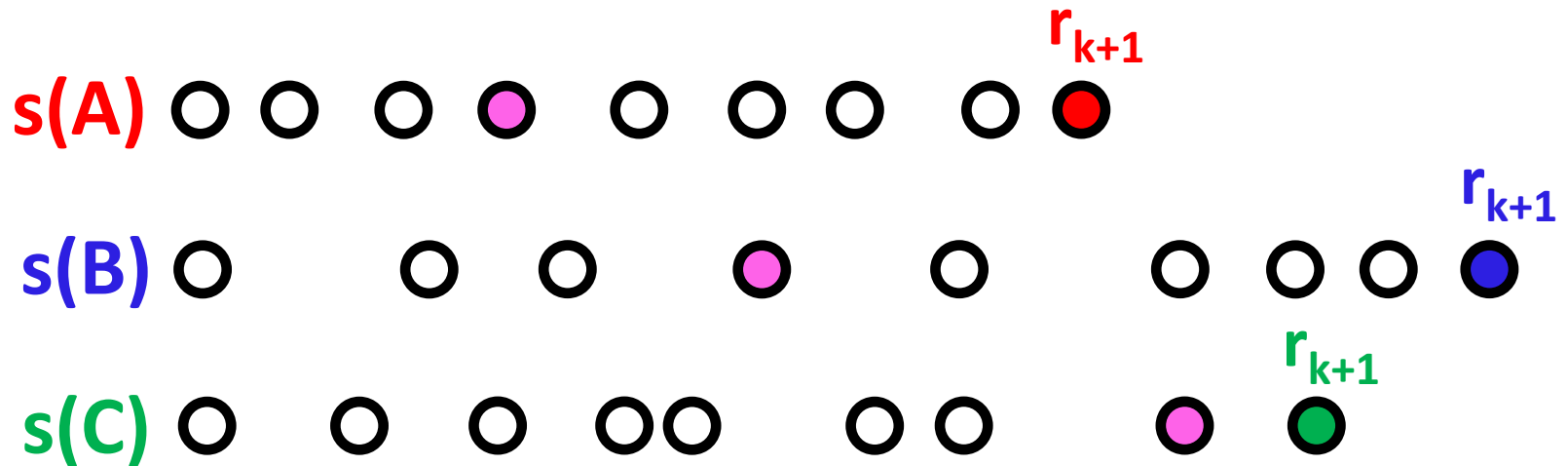
$$r_1 = h(x)/w_1$$

$$r_2 = h(x)/w_2$$

$$r_3 = h(x)/w_3$$

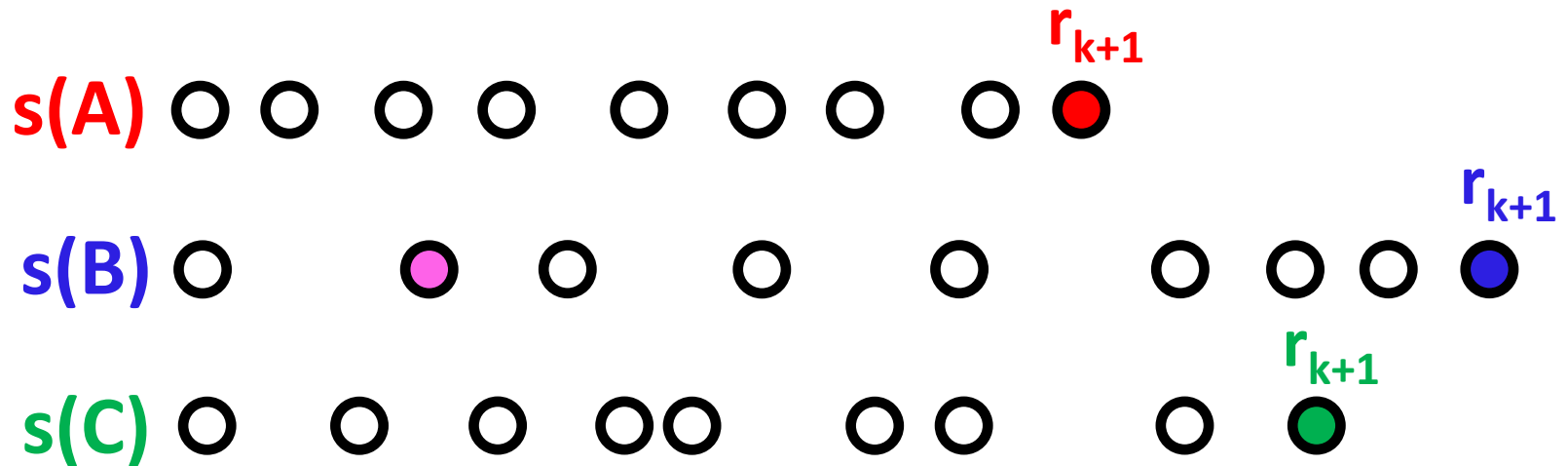
Coordinated weighted samples

- Estimate $\sum_x \max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x}), \mathbf{w}_3(\mathbf{x})\}$?



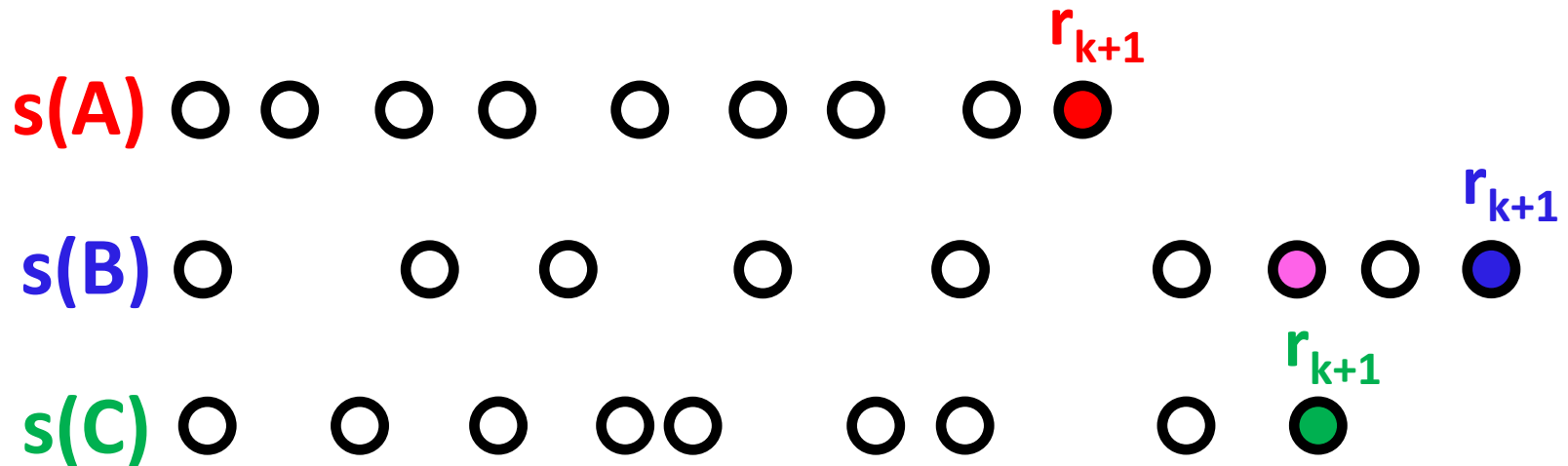
Coordinated weighted samples

- Estimate $\sum_x \max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x}), \mathbf{w}_3(\mathbf{x})\}$?



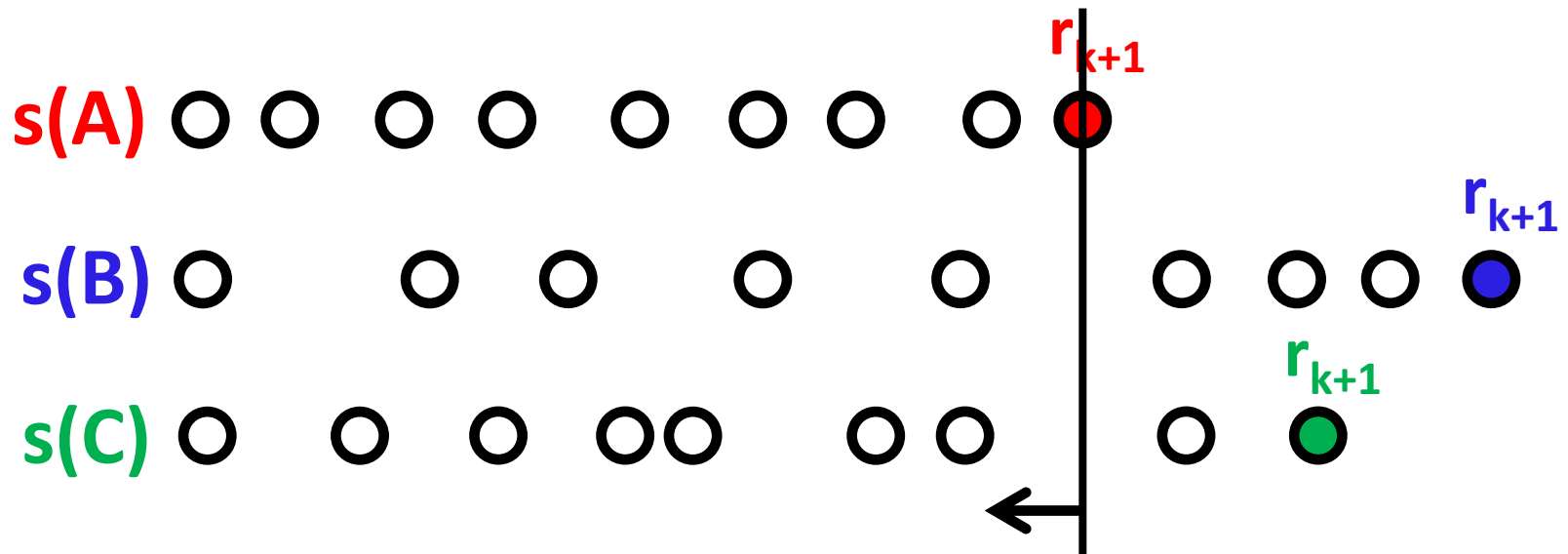
Coordinated weighted samples

- Estimate $\sum_x \max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x}), \mathbf{w}_3(\mathbf{x})\}$?



Coordinated weighted samples

- Estimate $\sum_x \max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x}), \mathbf{w}_3(\mathbf{x})\}$?



- We know $\max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x}), \mathbf{w}_3(\mathbf{x})\}$ for all x such that $\exists i \ r_i(x) < \mathbf{r}_{k+1}$

Coordinated sampling

$$a(x) = \begin{cases} \frac{\max\{w_i(x)\}}{p} & \min\{\mathbf{r}_1(x), \mathbf{r}_2(x), \mathbf{r}_3(x)\} < \min\{\mathbf{r}_{k+1}, \mathbf{r}_{k+1}, \mathbf{r}_{k+1}\} \\ 0 & \text{otherwise} \end{cases}$$

$$p = P(\min\{\mathbf{r}_1(x), \mathbf{r}_2(x), \mathbf{r}_3(x)\} < \min\{\mathbf{r}_{k+1}, \mathbf{r}_{k+1}, \mathbf{r}_{k+1}\})$$

$$p = P(h(x) < \min\{\mathbf{r}_{k+1}, \mathbf{r}_{k+1}, \mathbf{r}_{k+1}\} \max\{\mathbf{w}_1(x), \mathbf{w}_2(x), \mathbf{w}_3(x)\})$$

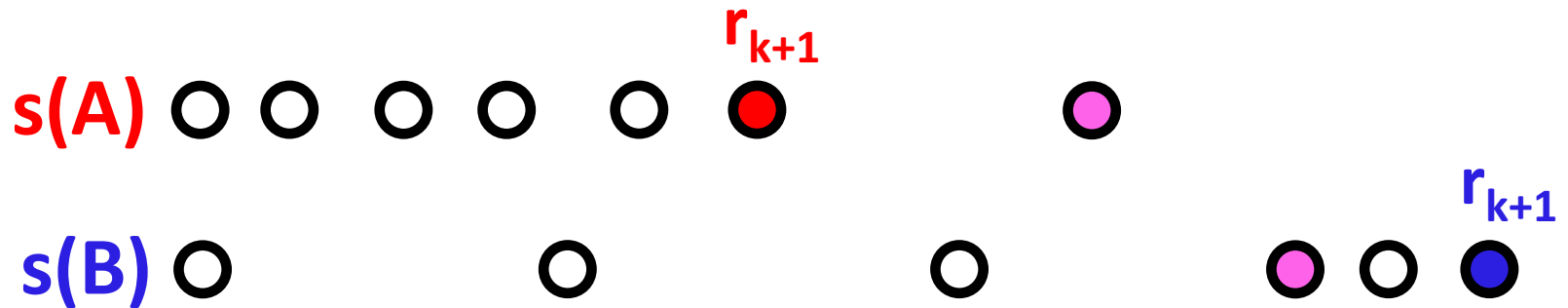
Homework

- Develop an estimator for $\sum_x \text{median}\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x}), \mathbf{w}_3(\mathbf{x})\} ?$

Using partial information

Partial information

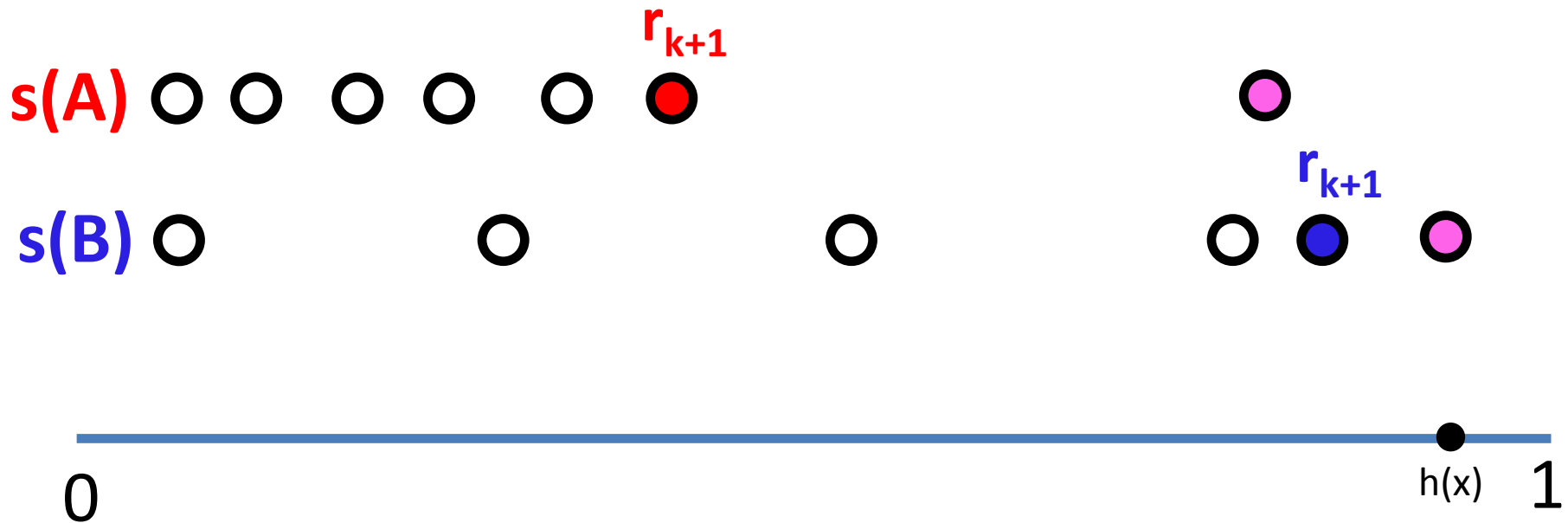
- Estimate $\sum_x \max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x})\}$?



- We may not know $\max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x})\}$ but we do know $\mathbf{w}_2(\mathbf{x})$ which is a lower bound

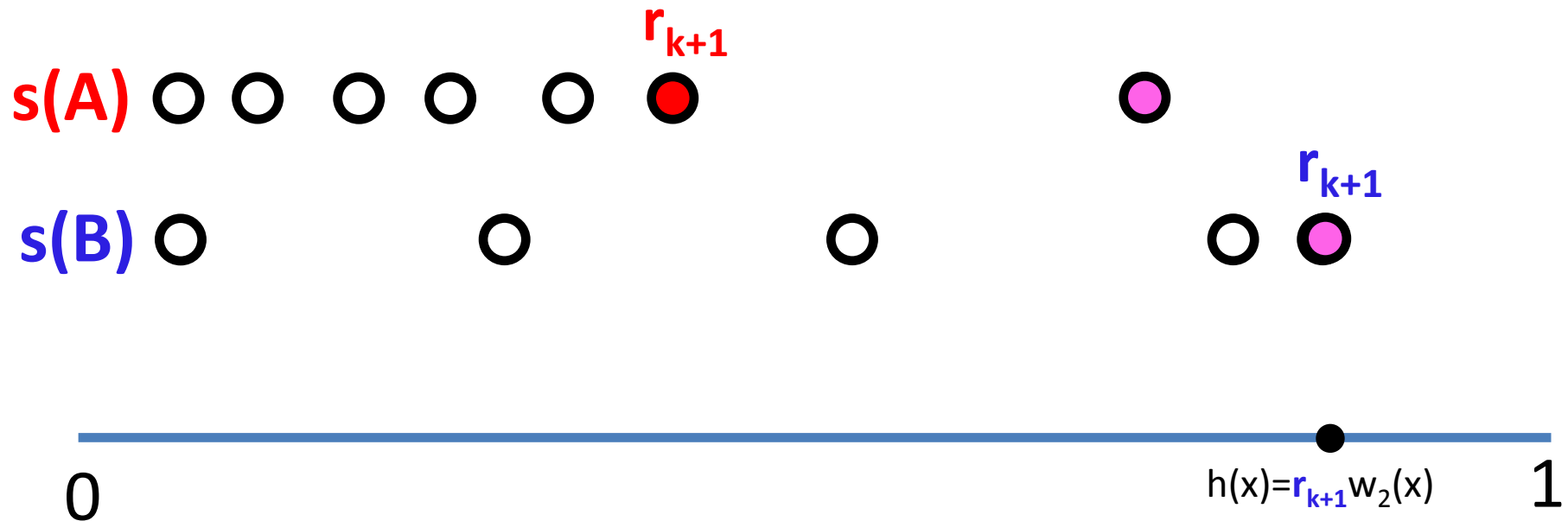
Partial information

- Estimate $\sum_x \max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x})\}$?



Partial information

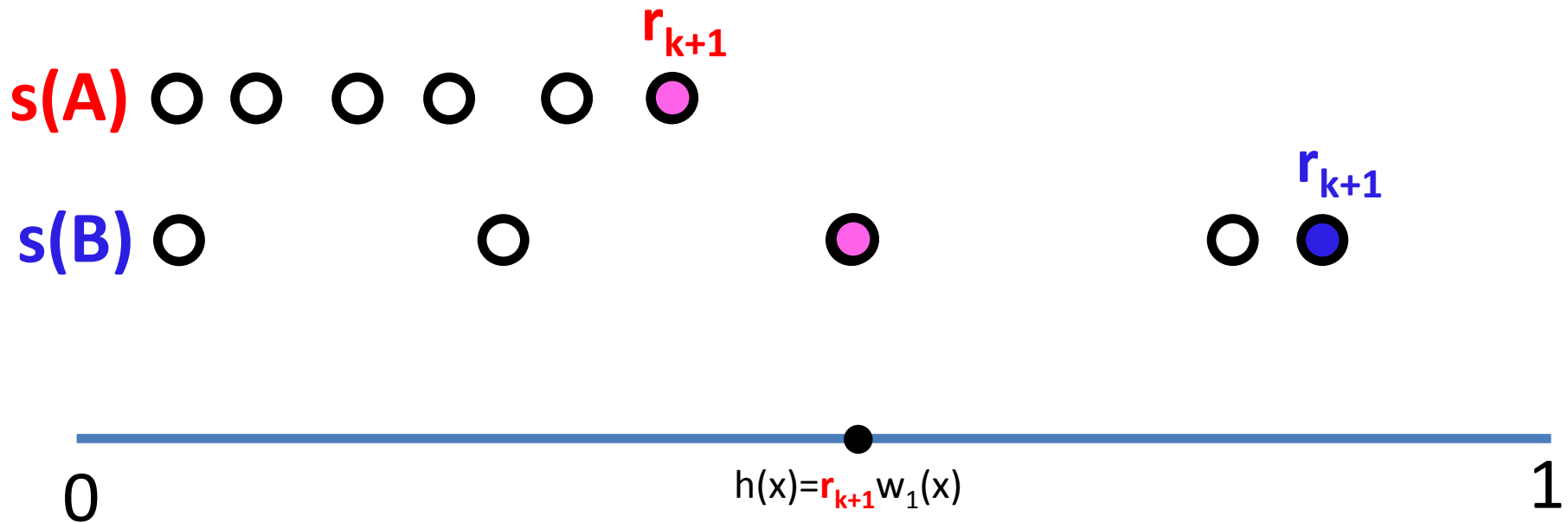
- Estimate $\sum_x \max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x})\}$?



Estimate: $w_2(x) / (r_{k+1} w_2(x)) = 1 / r_{k+1}$

Partial information

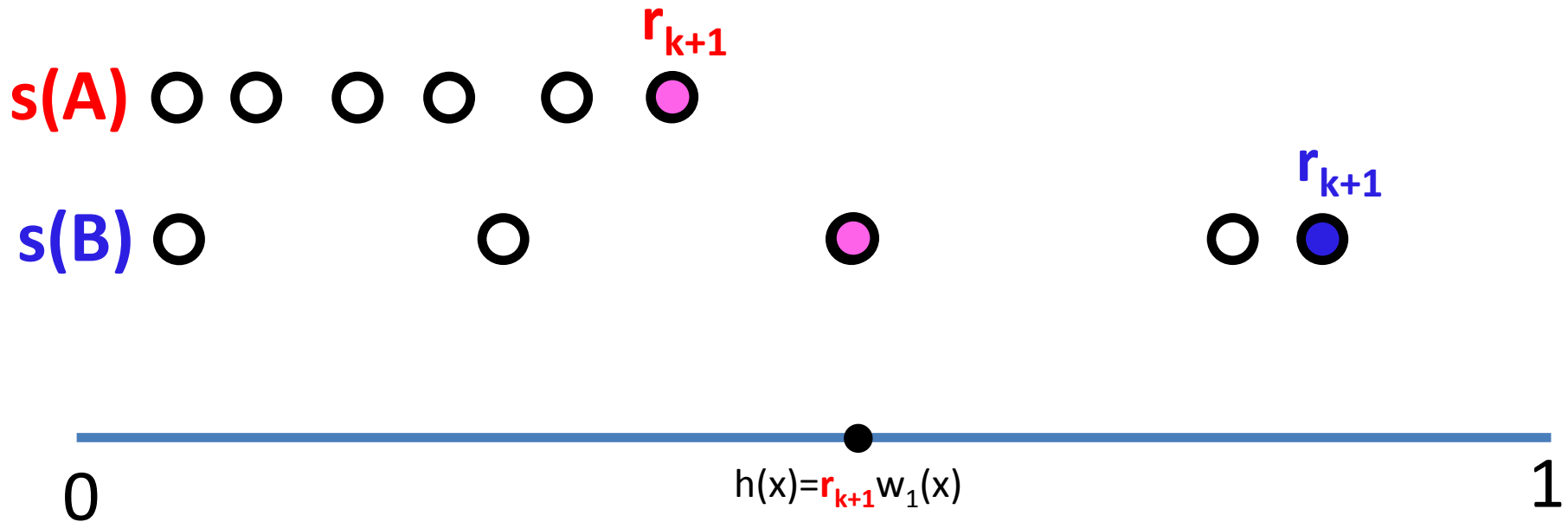
- Estimate $\sum_x \max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x})\}$?



Estimate: ??

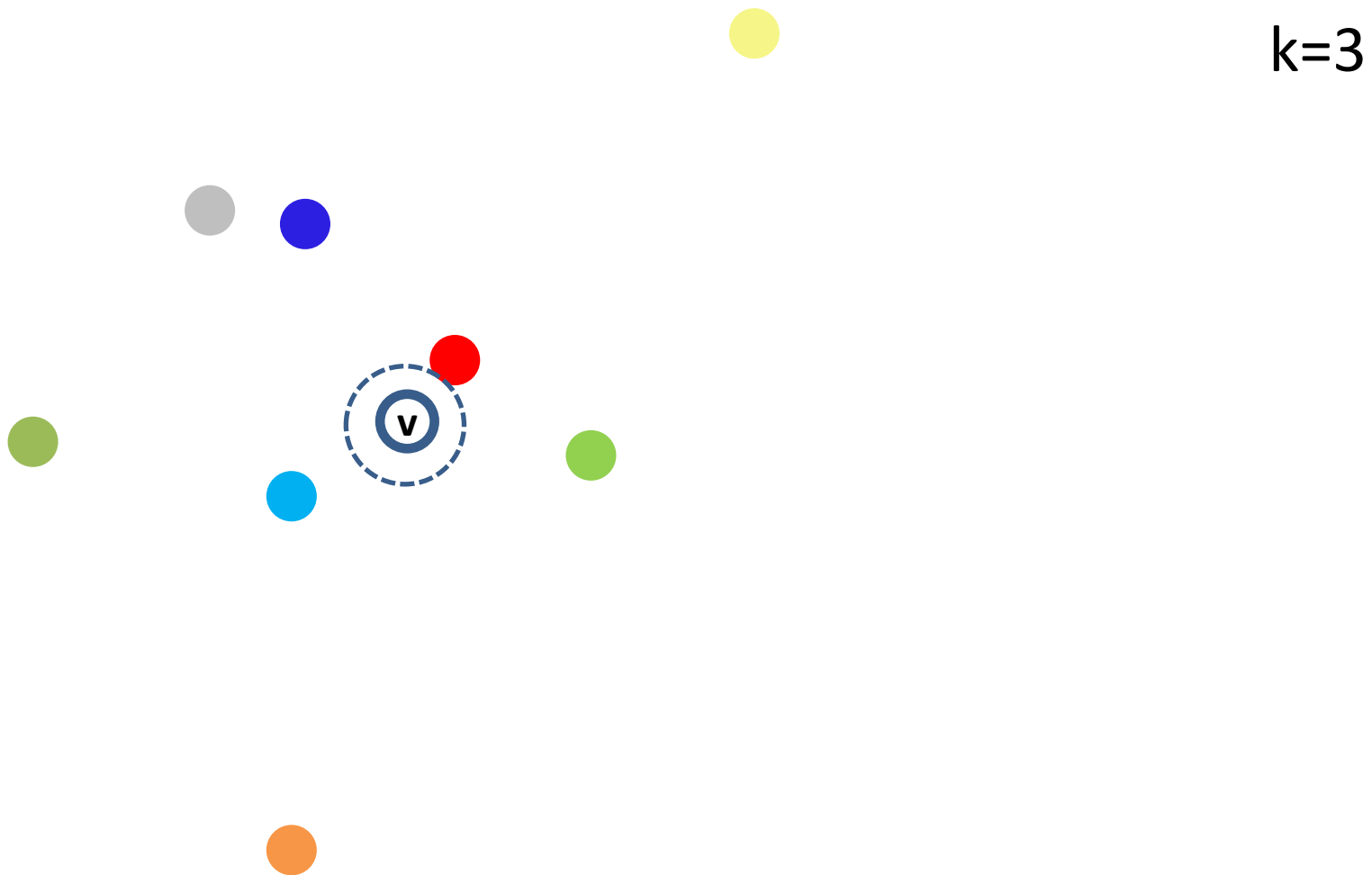
Partial information

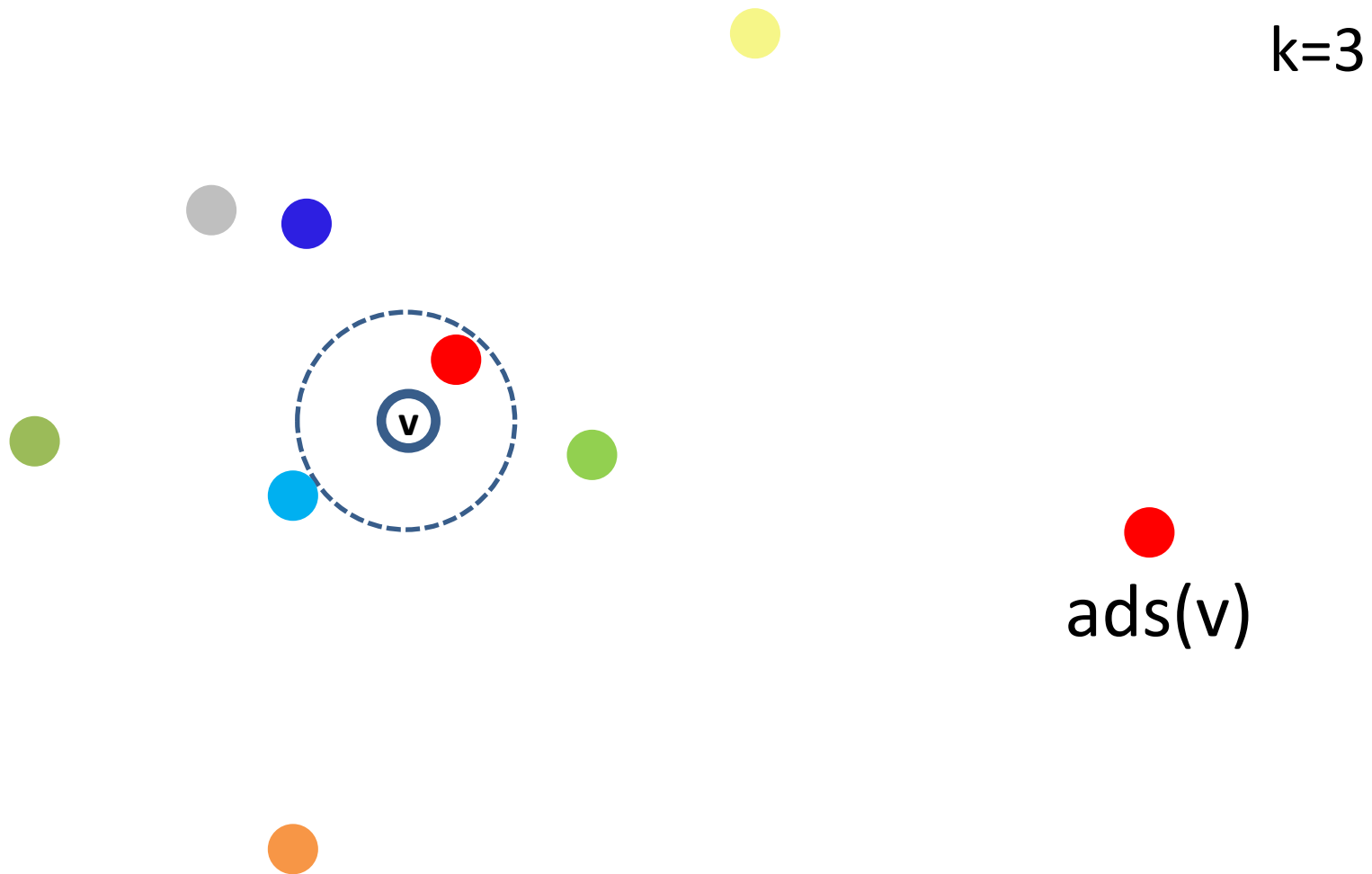
- Estimate $\sum_x \max\{\mathbf{w}_1(\mathbf{x}), \mathbf{w}_2(\mathbf{x})\}$?

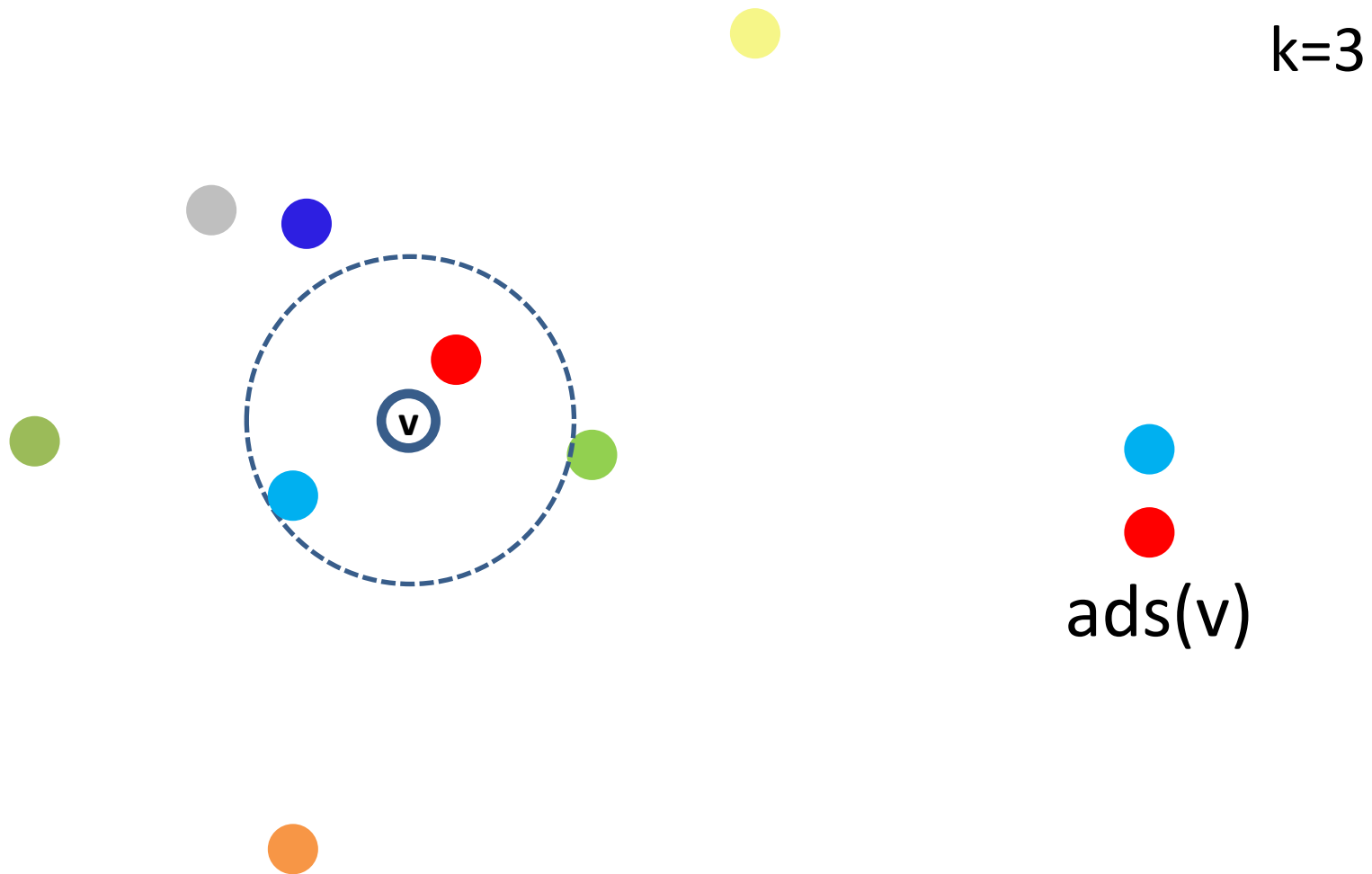


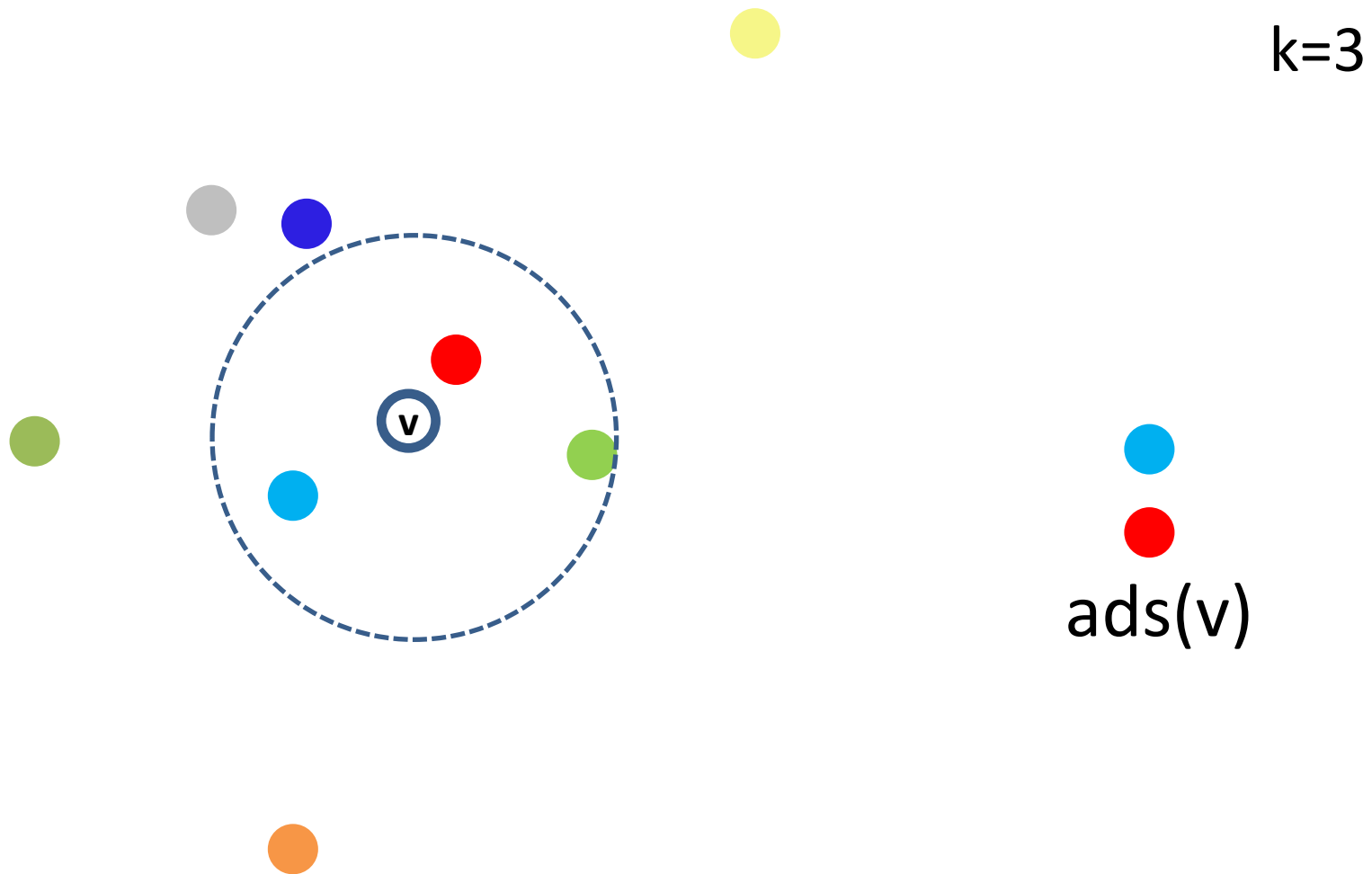
$$(1/\mathbf{r}_{k+1})(\mathbf{r}_{k+1} w_2(x) - \mathbf{r}_{k+1} w_1(x)) + \mathbf{E} \times \mathbf{r}_{k+1} w_1(x) = w_1(x)$$

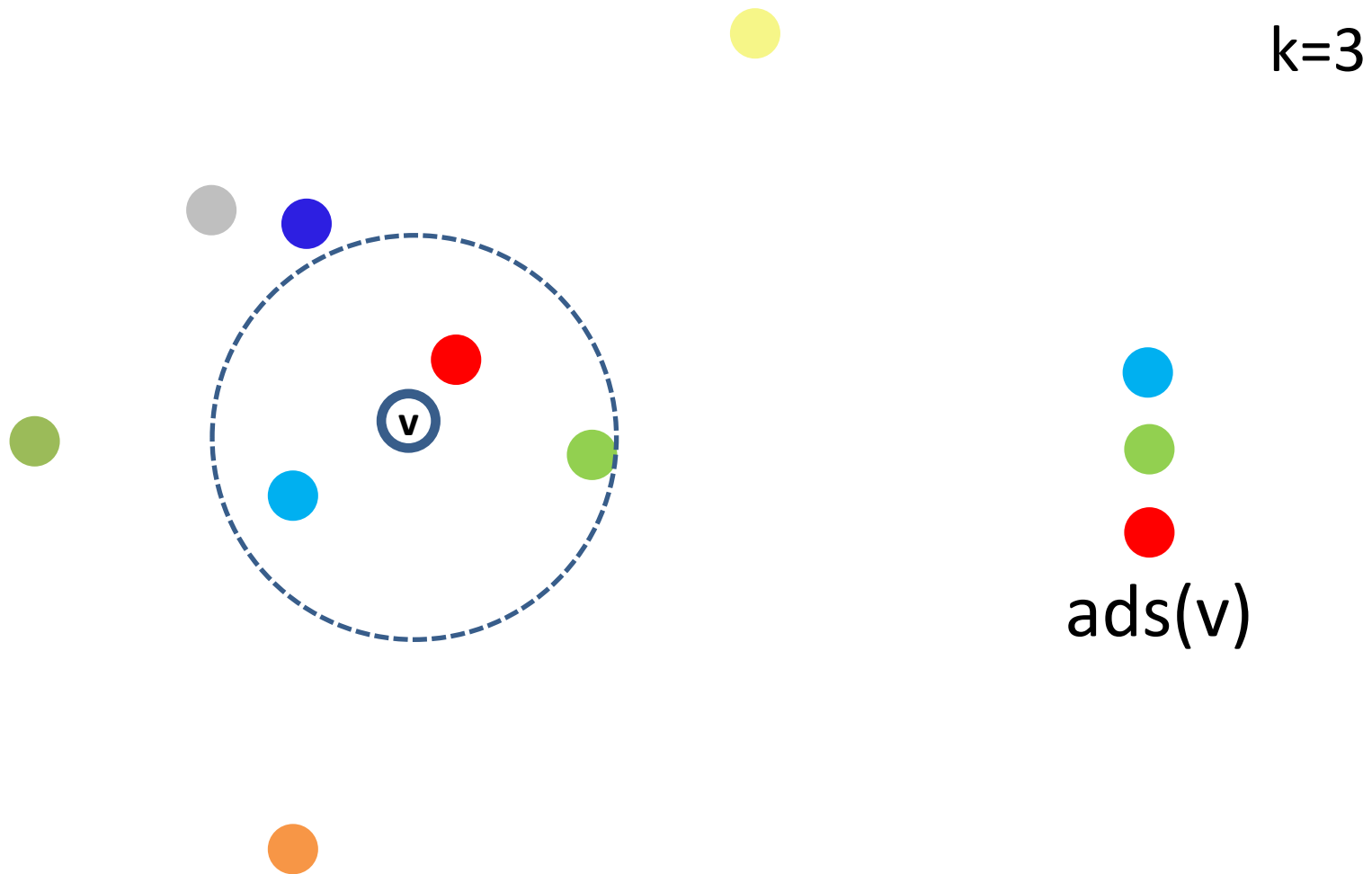
All Distances (coordinated) Sketches (ADS)

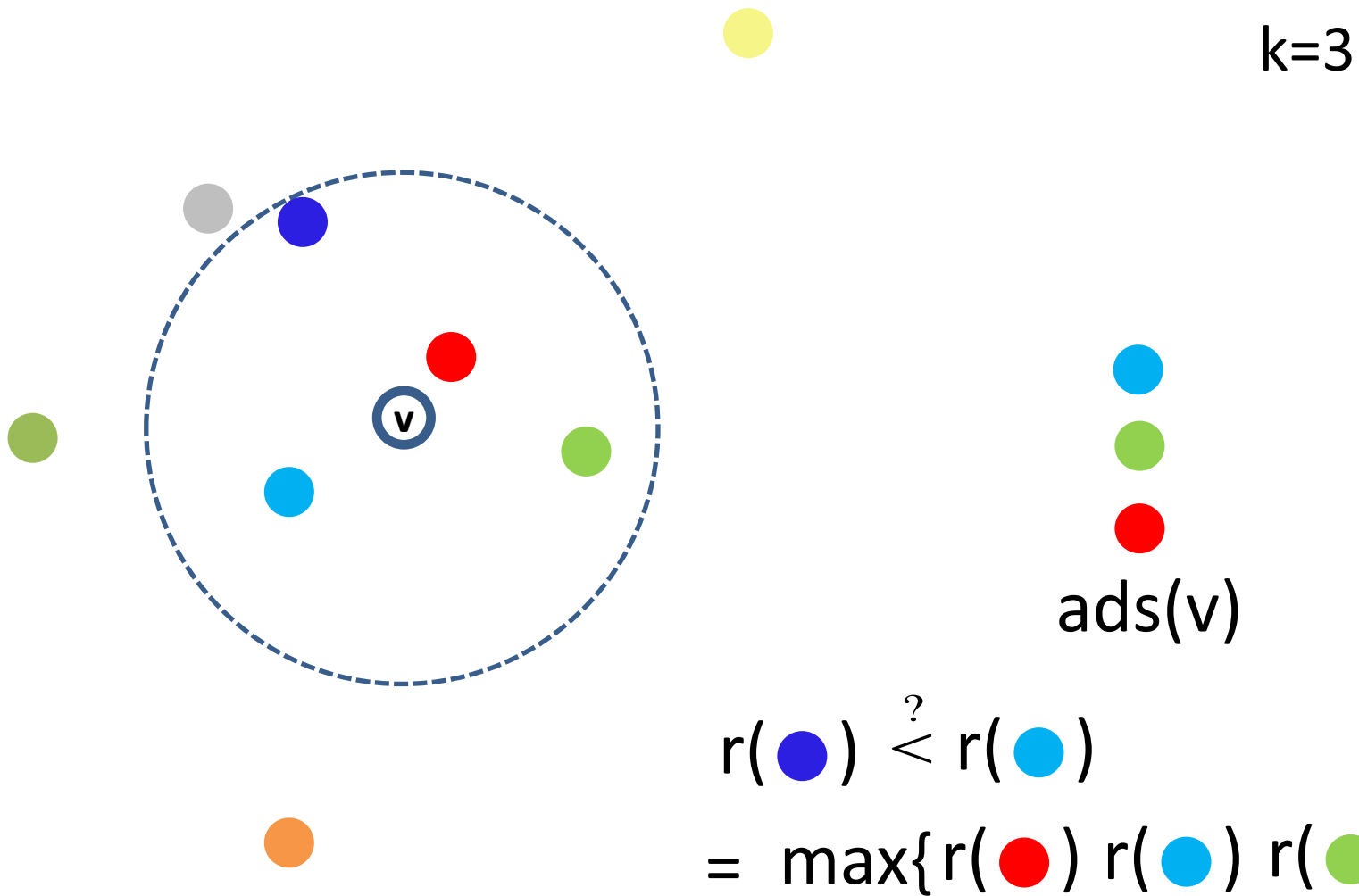


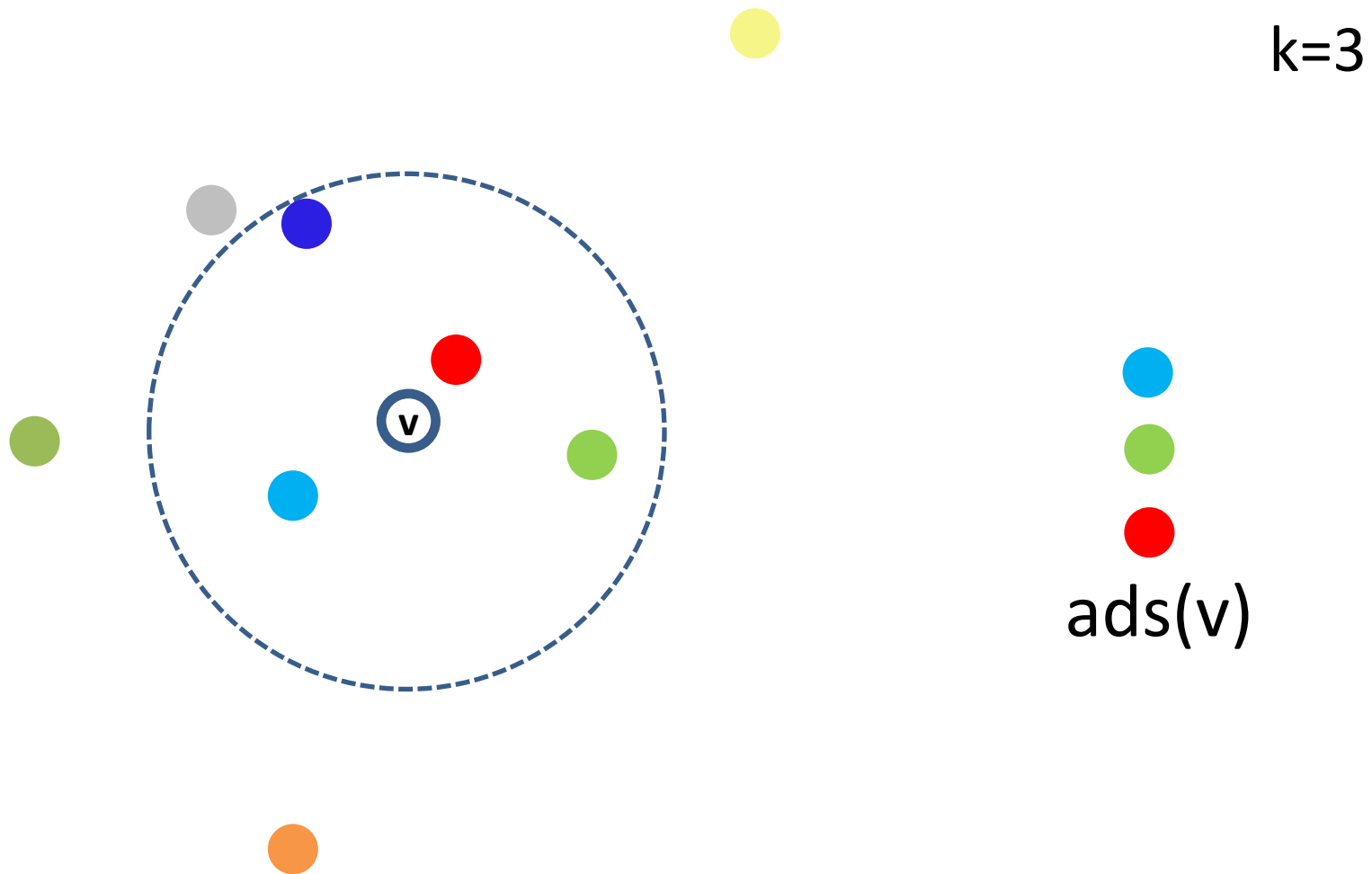


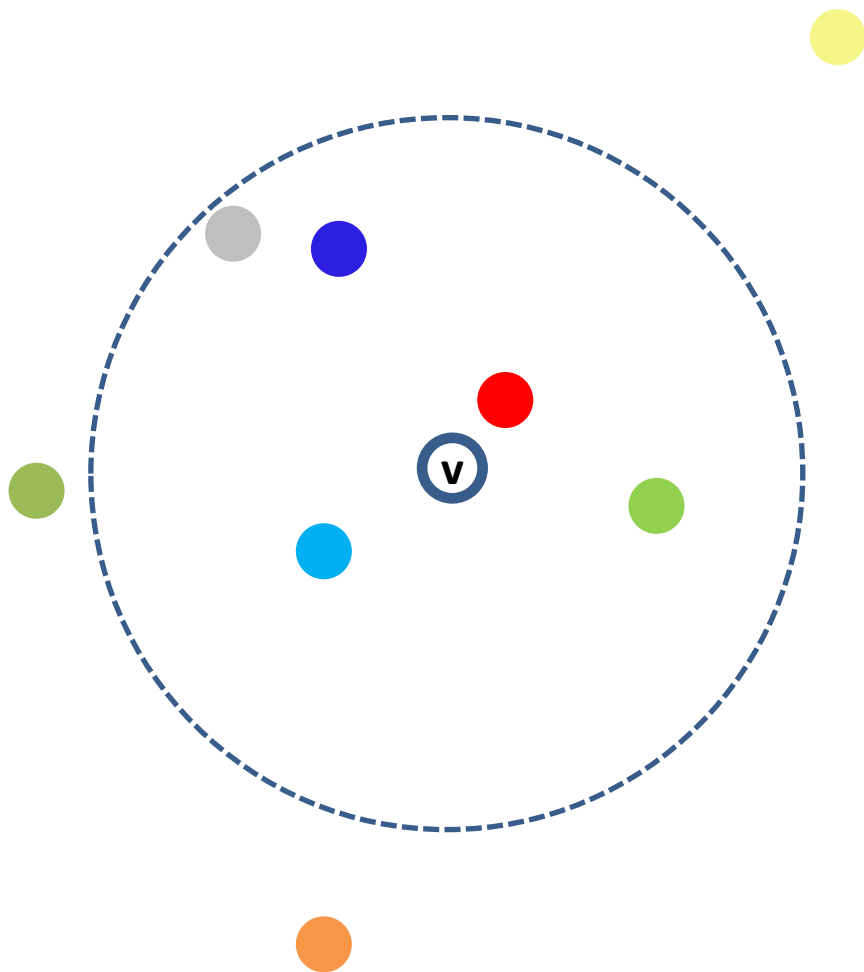








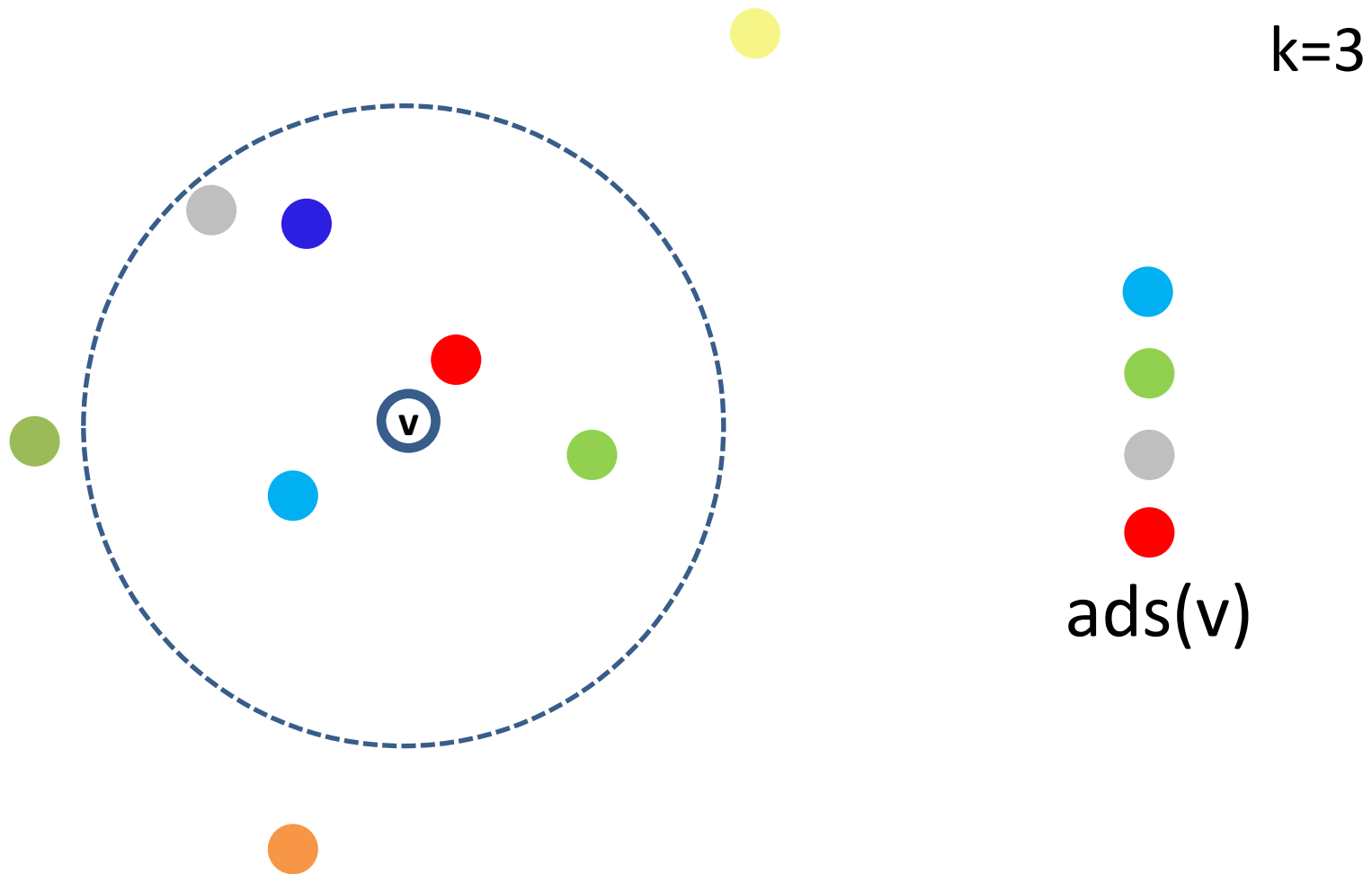


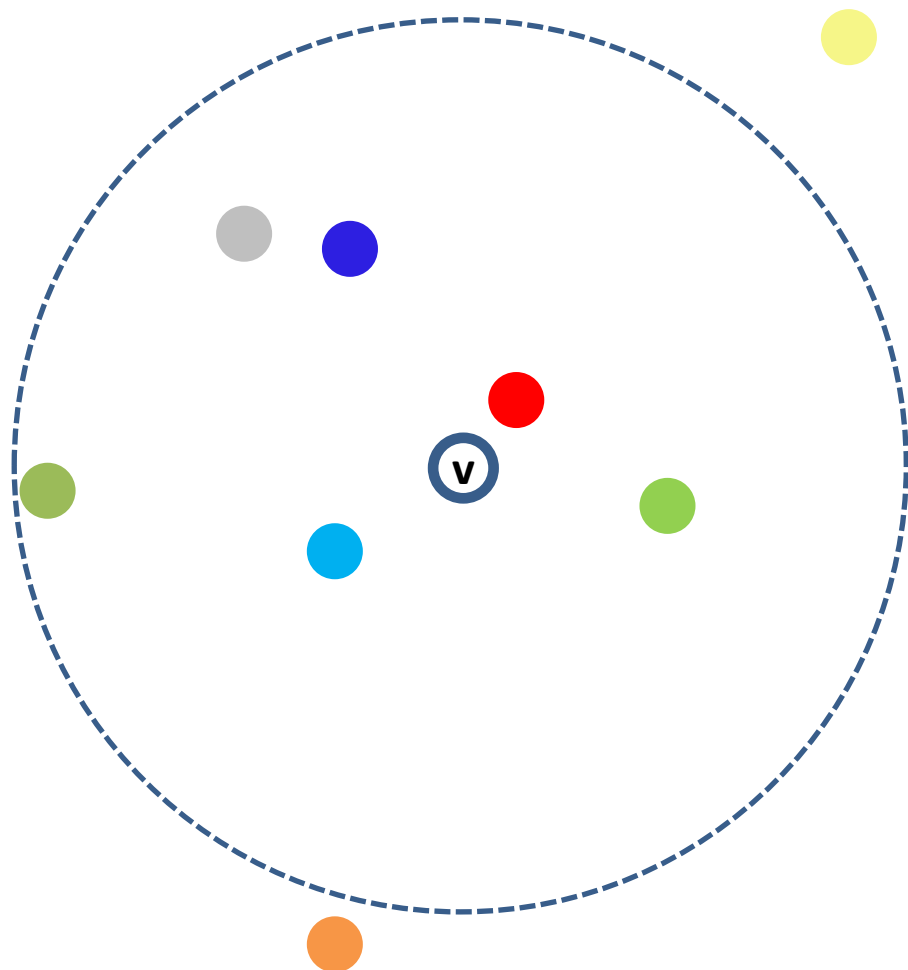


$k=3$

ads(v)

$$r(\text{grey dot}) \stackrel{?}{<} r(\text{blue dot})$$



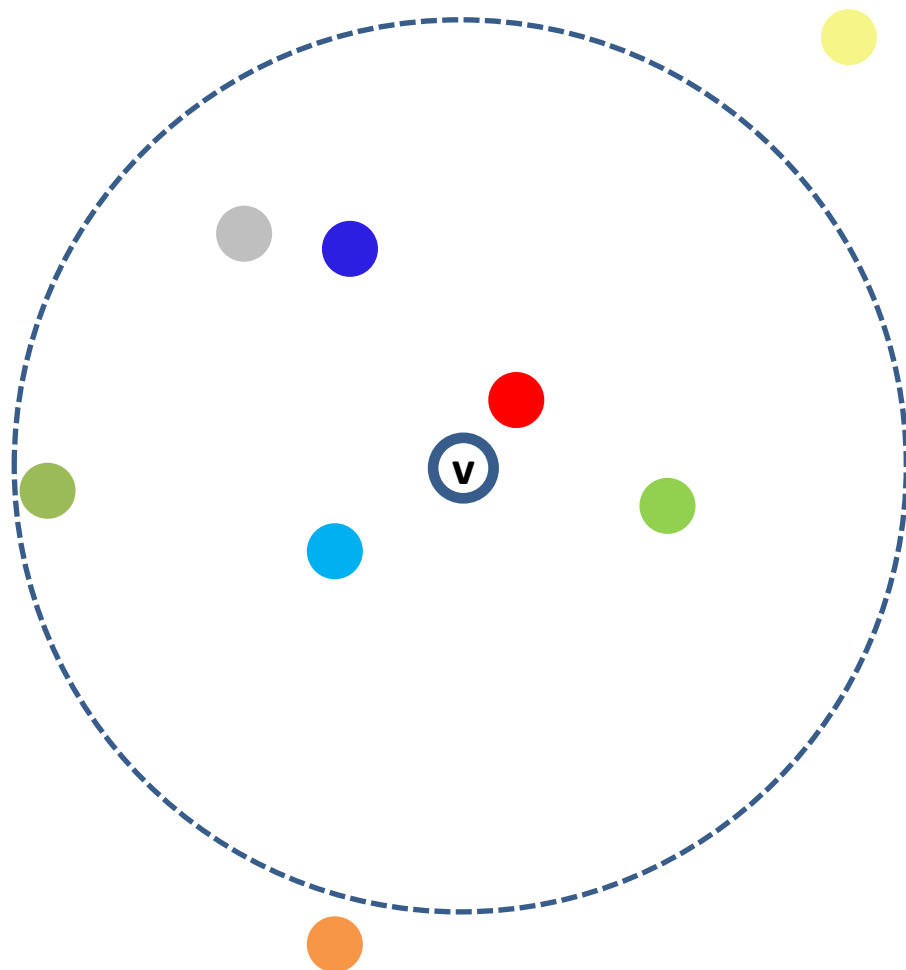


$k=3$

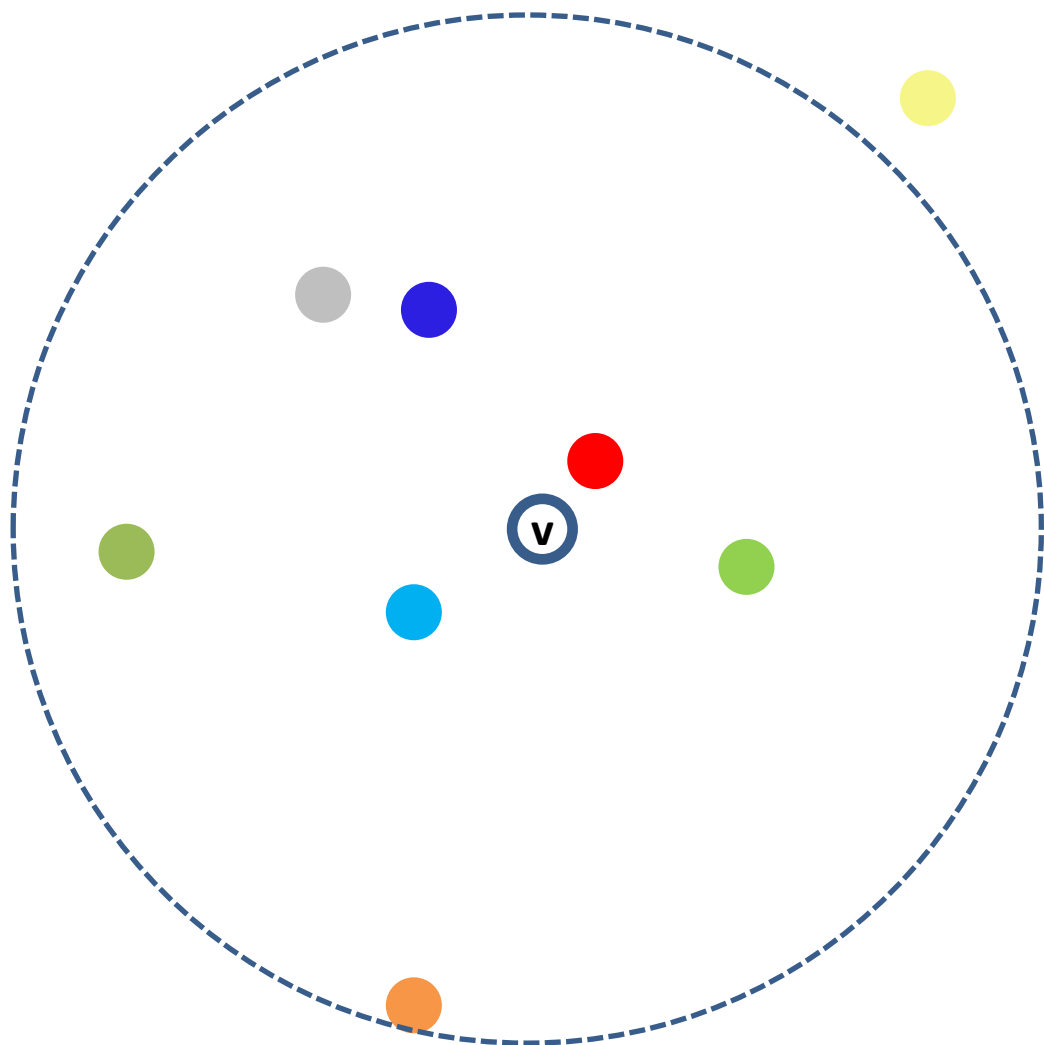
ads(v)

$r(\bullet) \stackrel{?}{<} r(\bullet)$

$k=3$



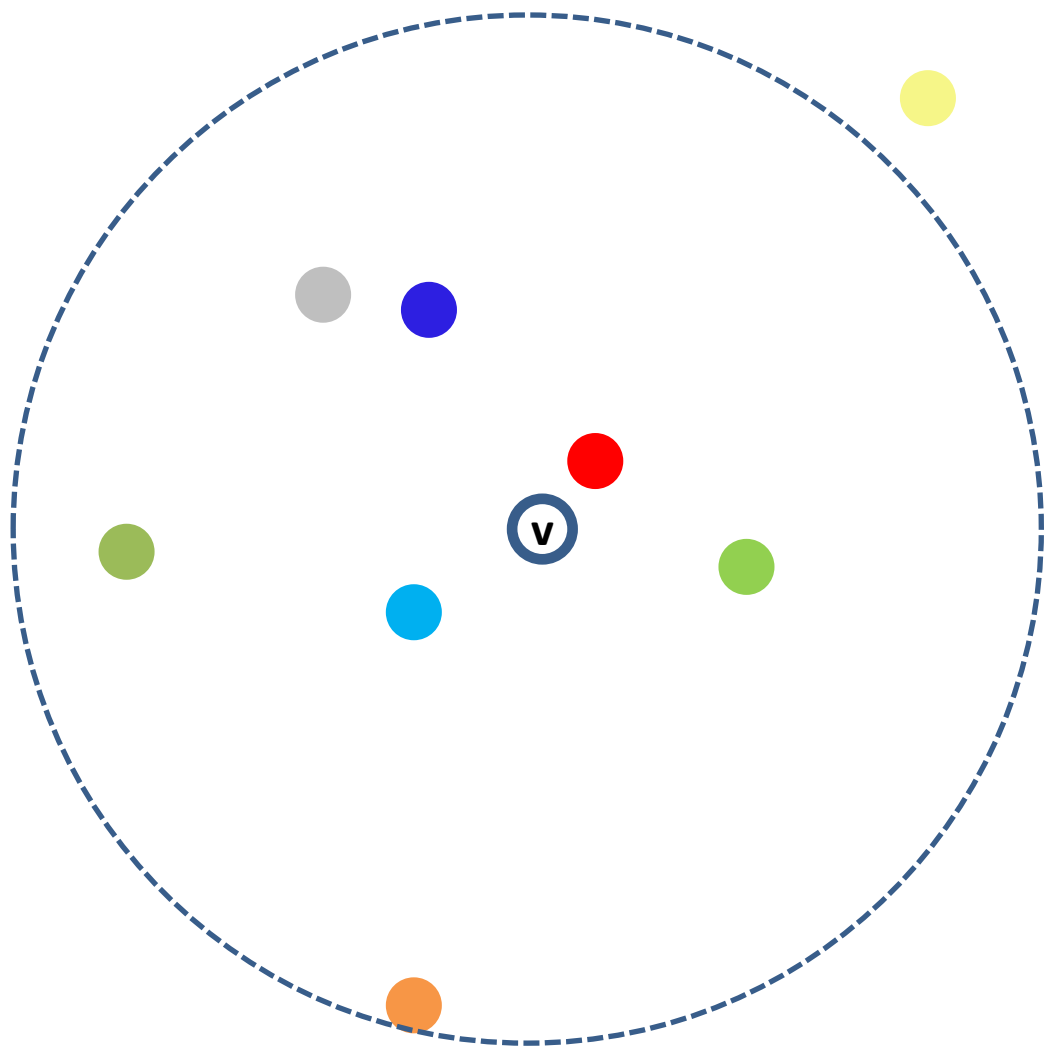
ads(v)



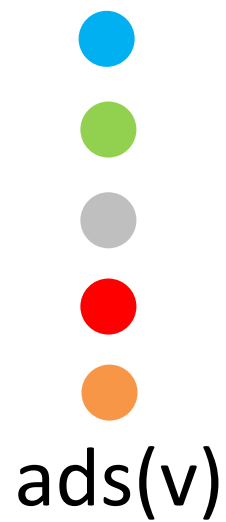
$k=3$

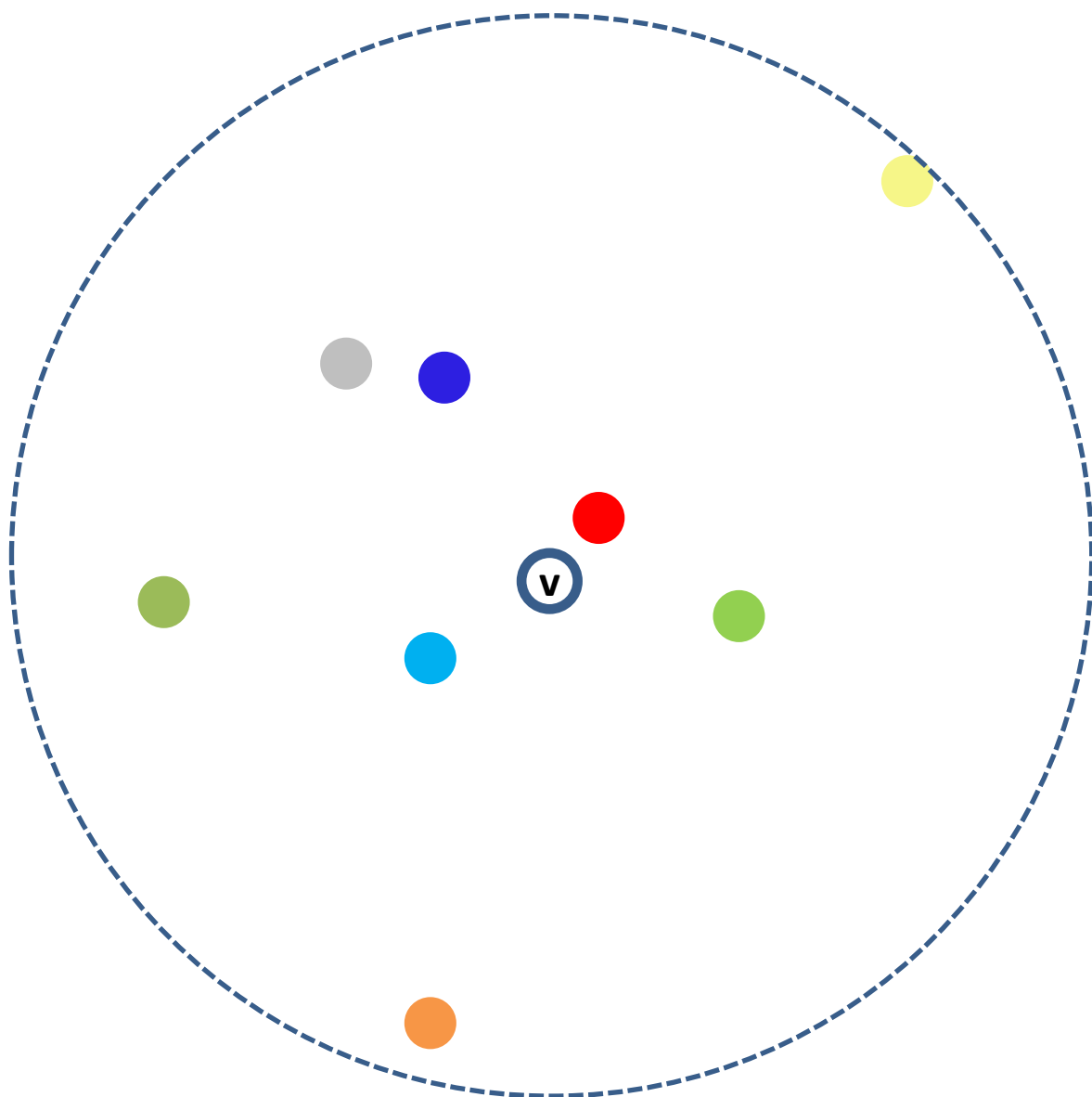
$\text{ads}(v)$

$r(\text{orange dot}) <? r(\text{green dot})$



$k=3$



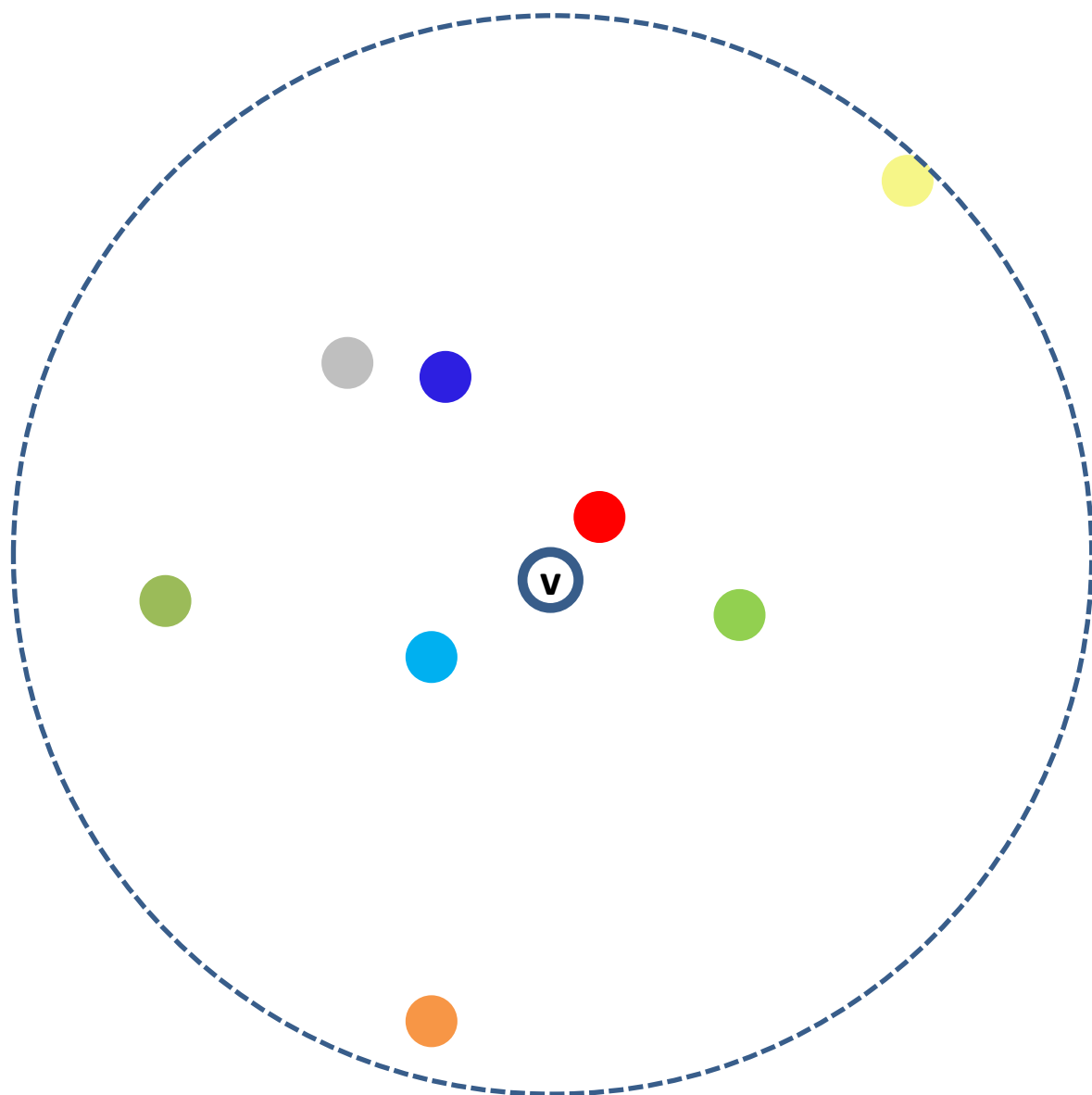


$k=3$



$\text{ads}(v)$

$$r(\text{yellow}) \stackrel{?}{<} r(\text{grey})$$

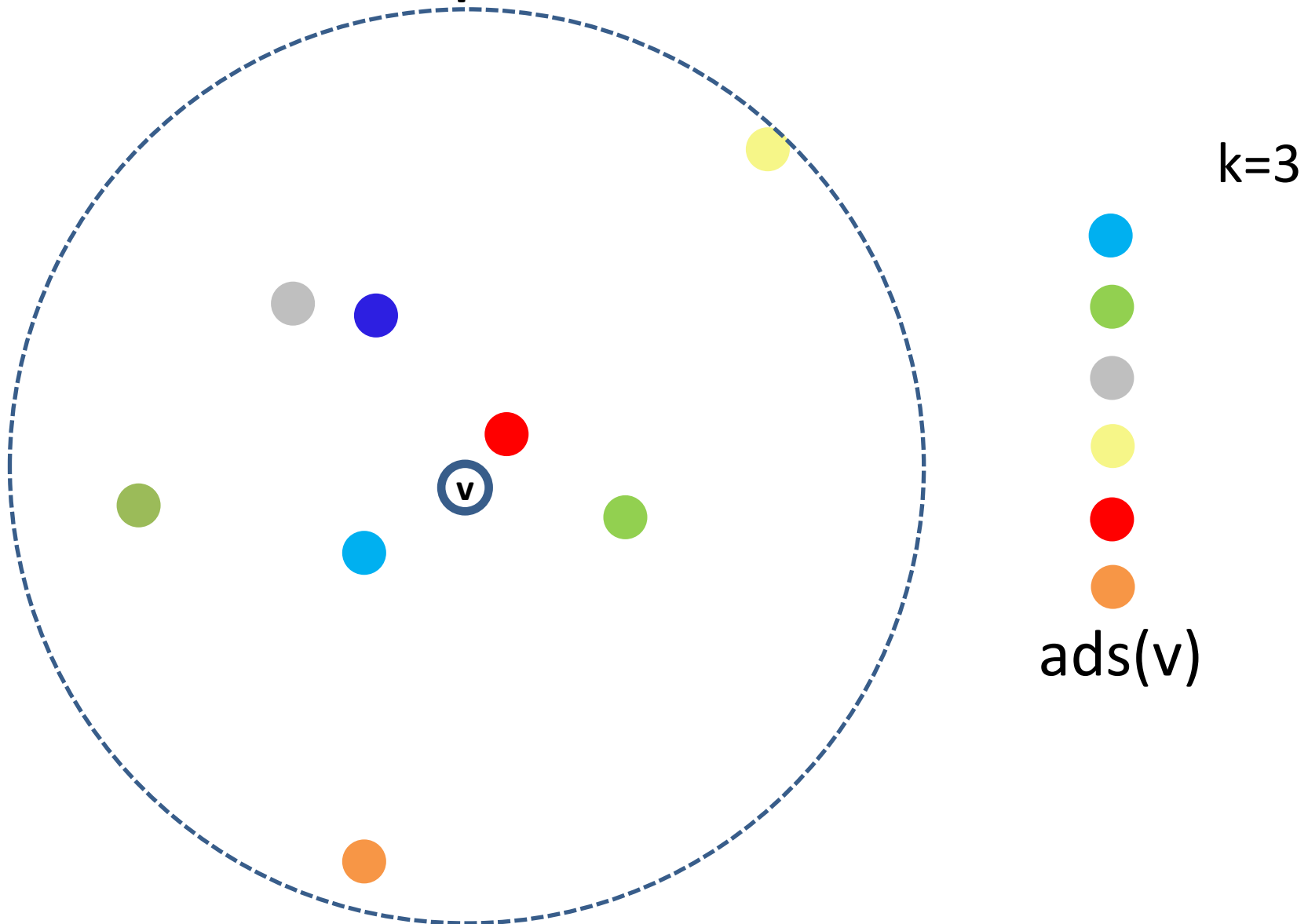


$k=3$



$\text{ads}(v)$

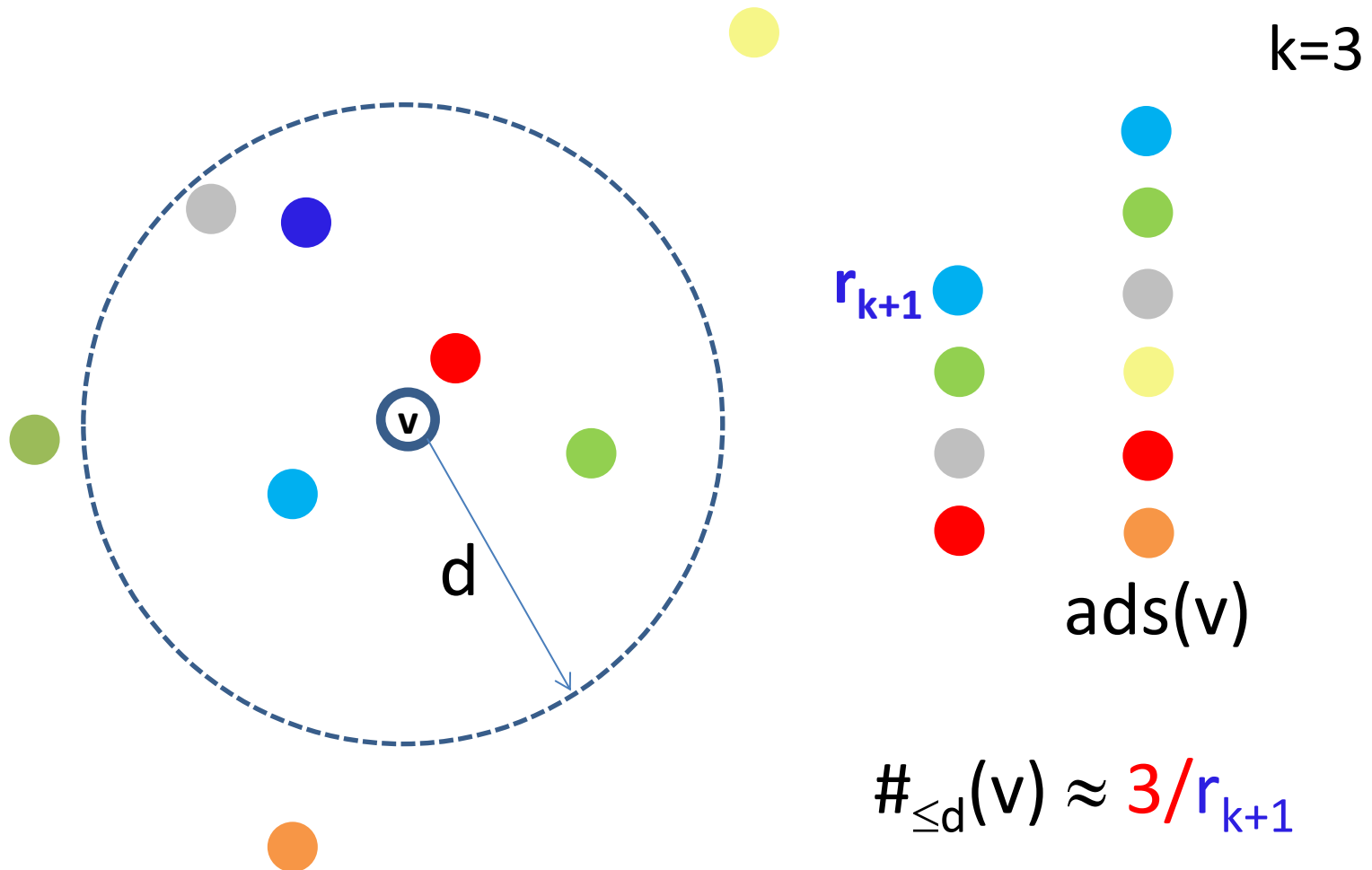
Properties of the ADS



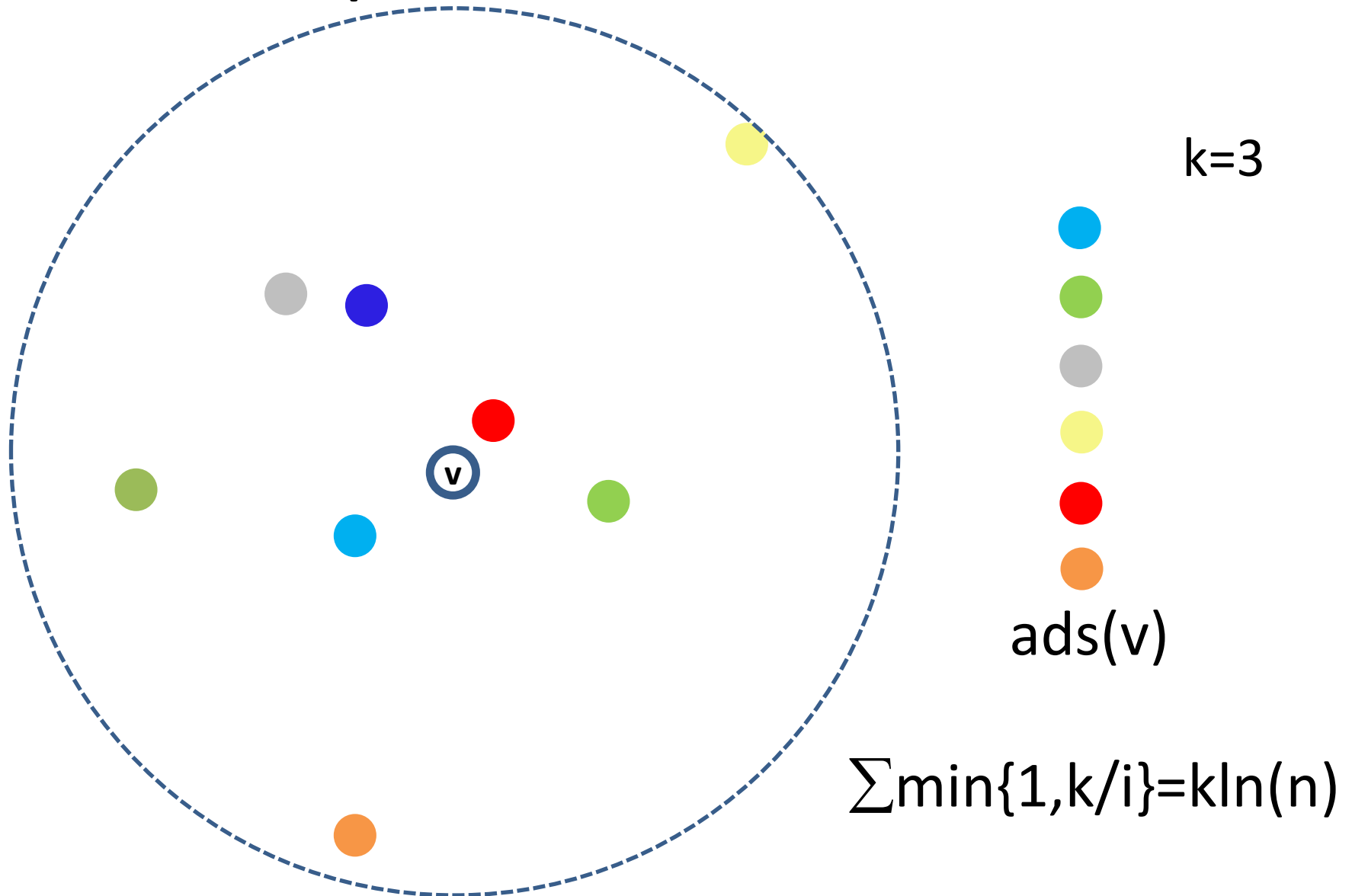
Properties of the ADS

- Includes a bottom-k sketch for every ball around v
 - ➔ Estimate how many vertices are there in every ball around v

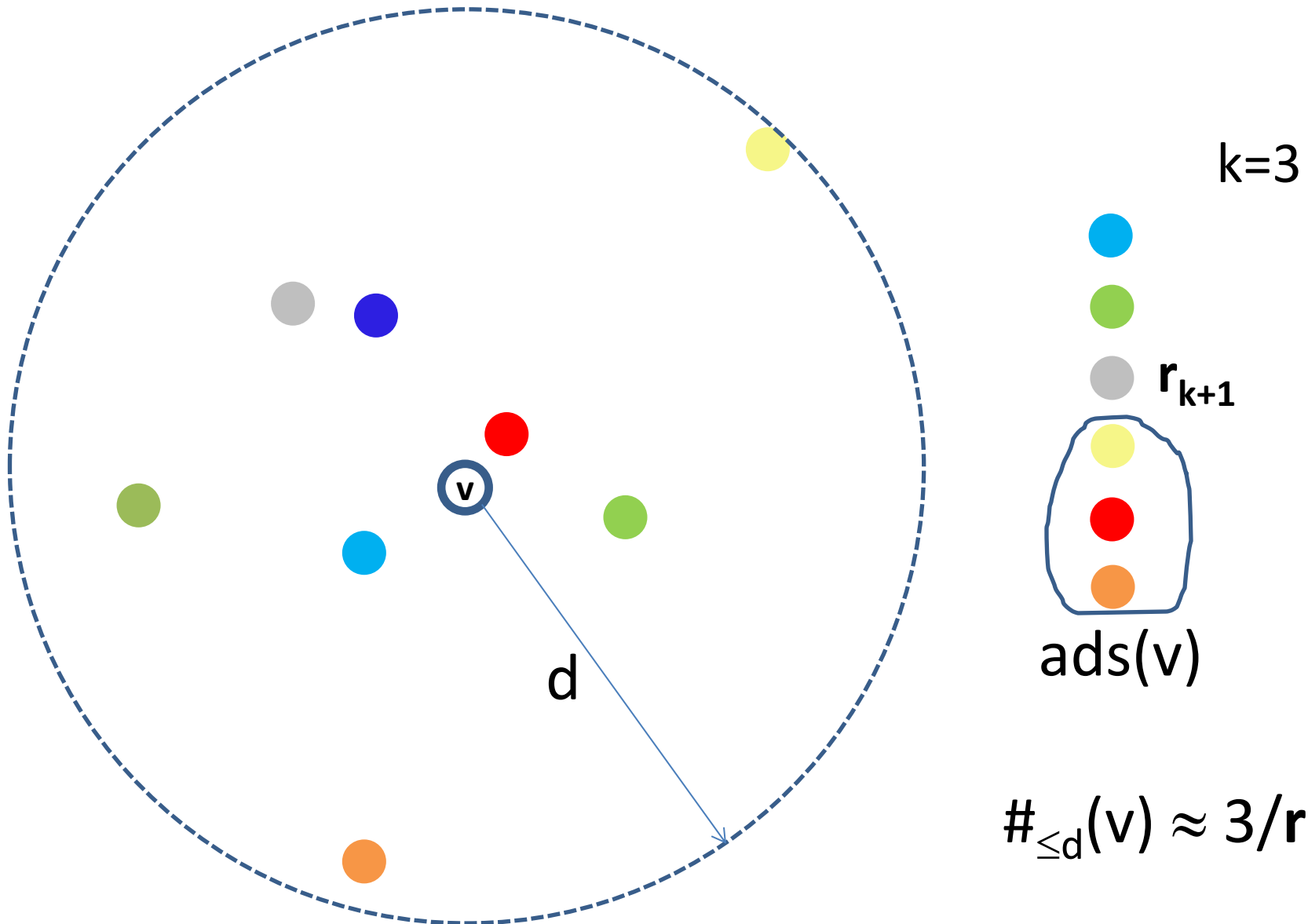
Estimating using the ADS



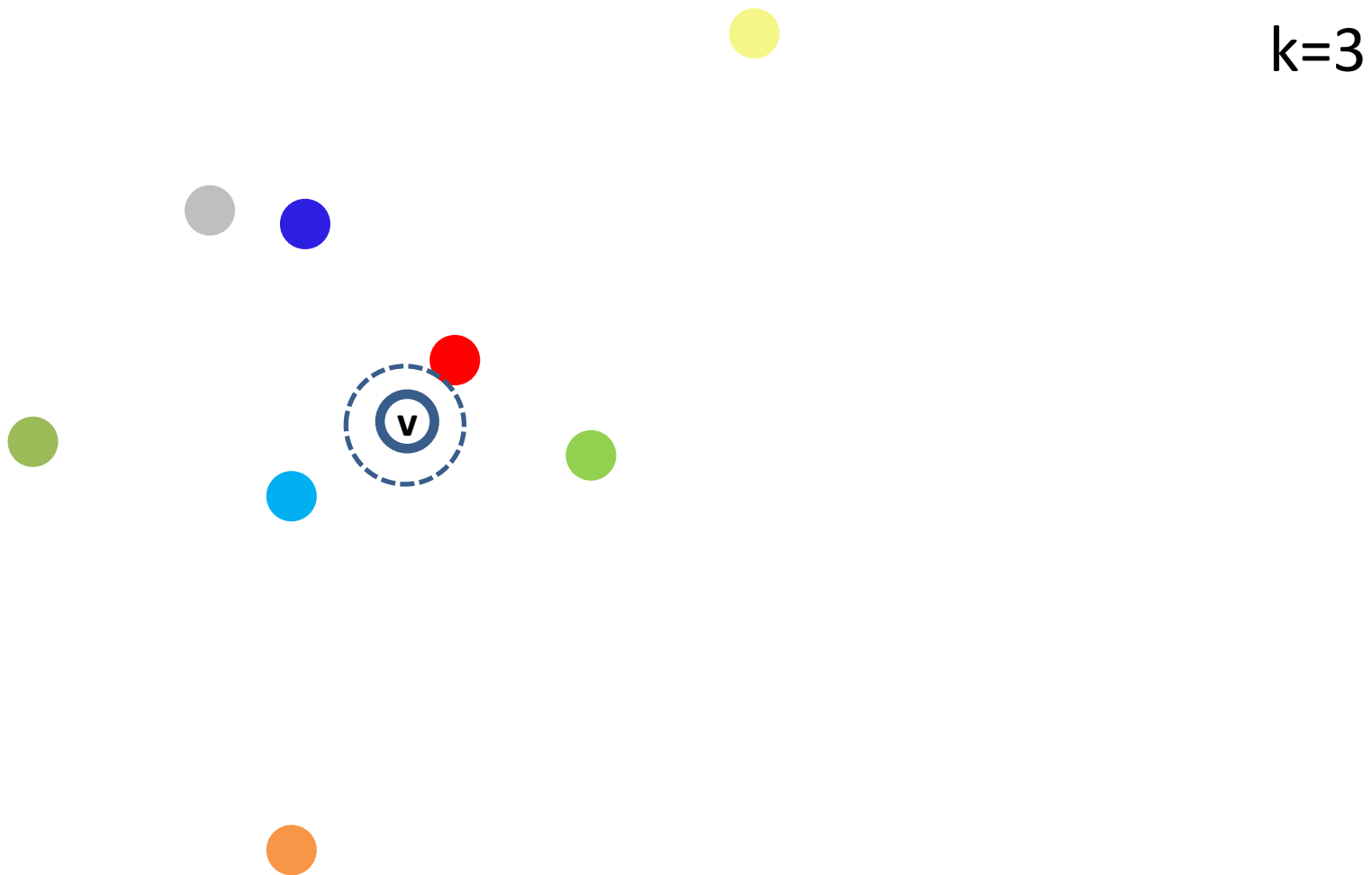
Expected size of the ADS



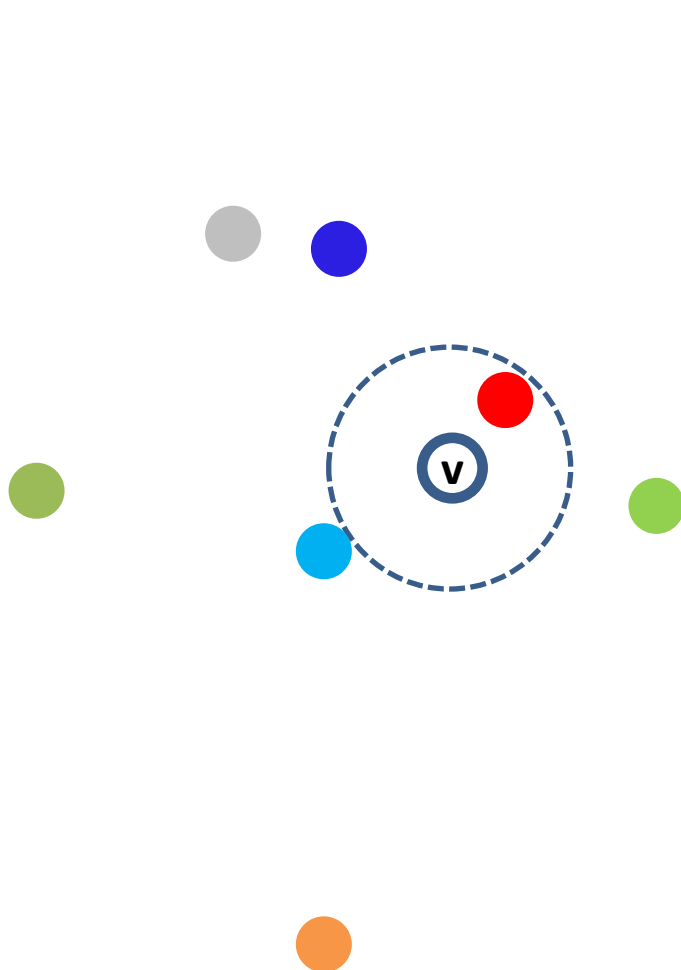
Do we use the most out of the ADS ?



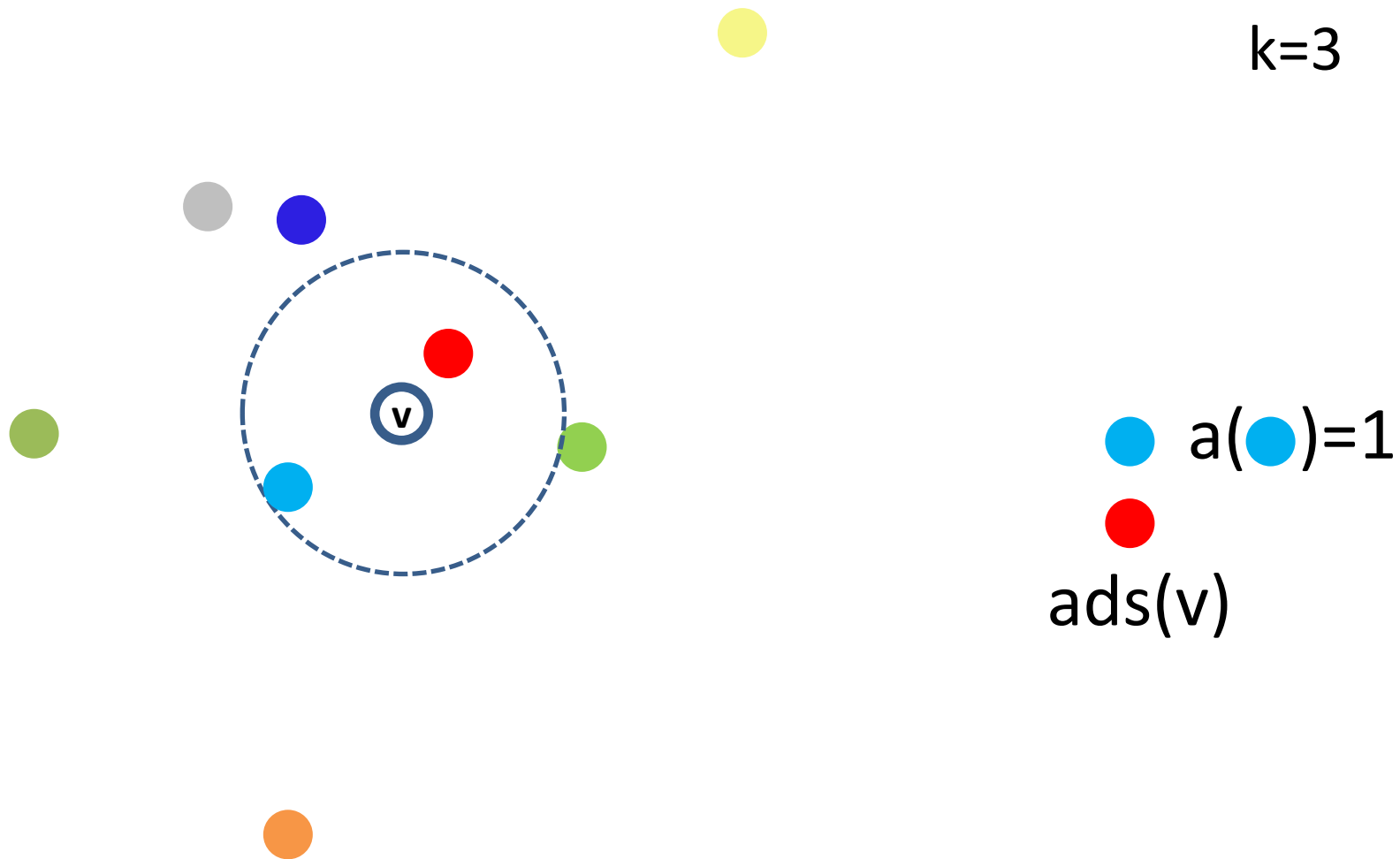
Give adjusted weights to all elements

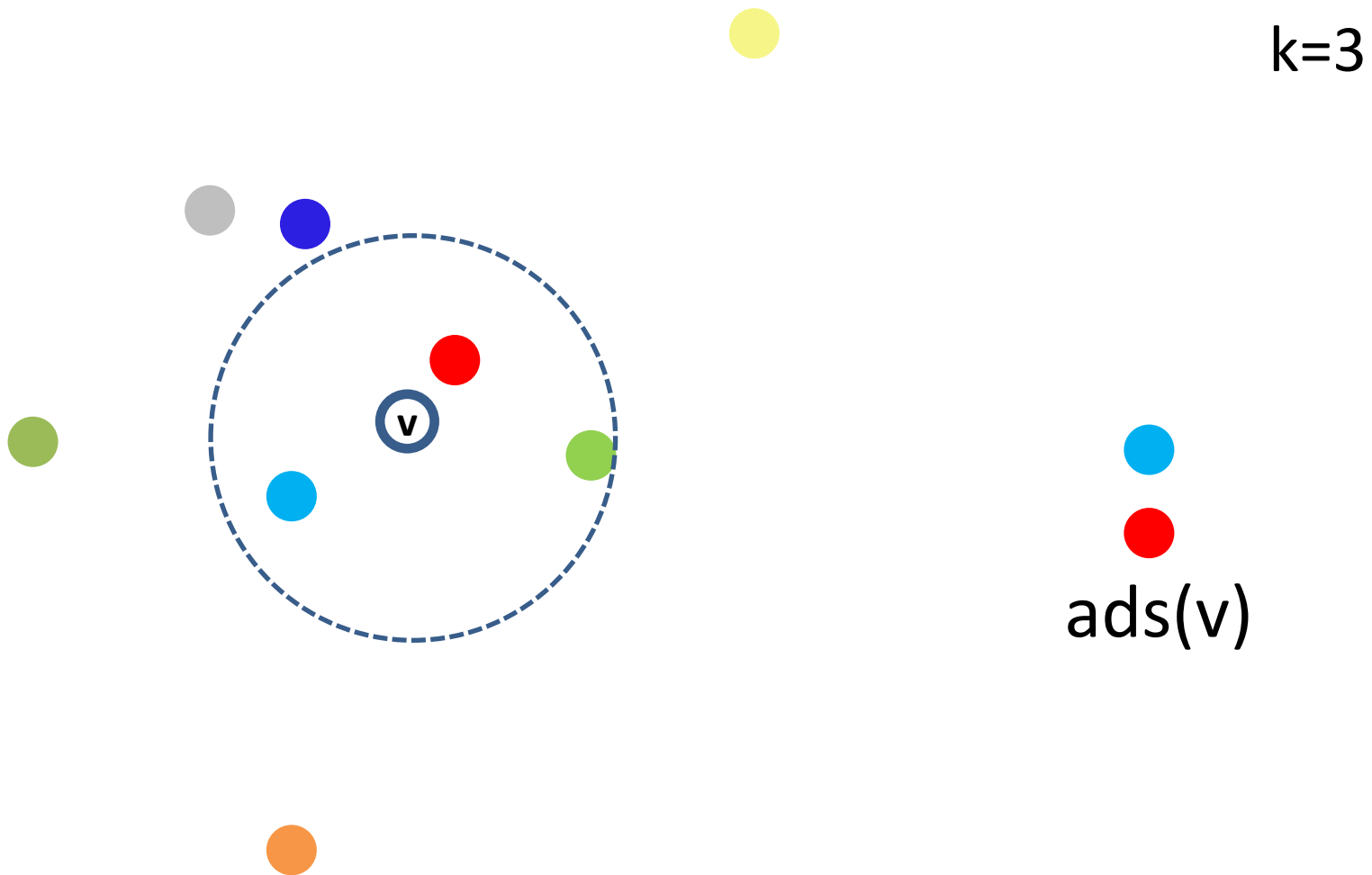


$k=3$

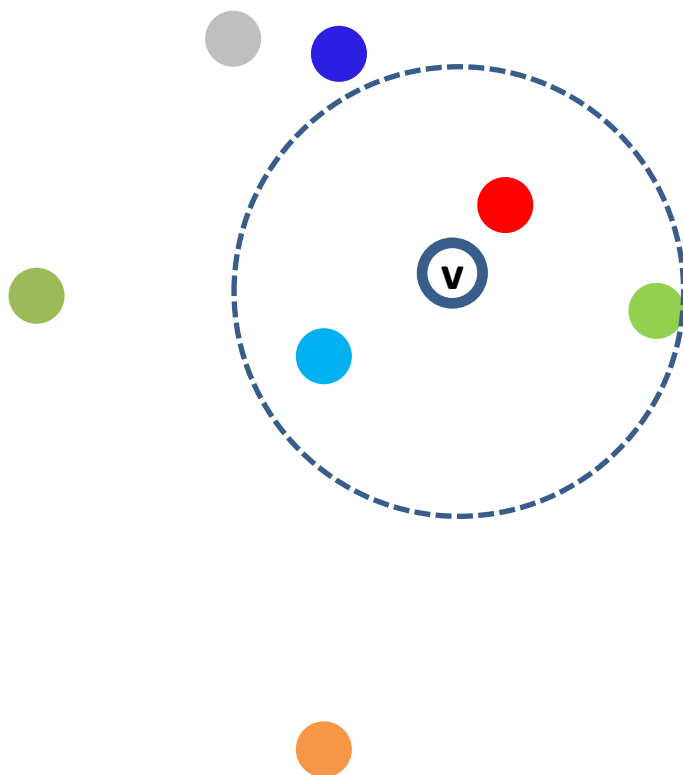


$a(\bullet)=1$
 $\text{ads}(v)$



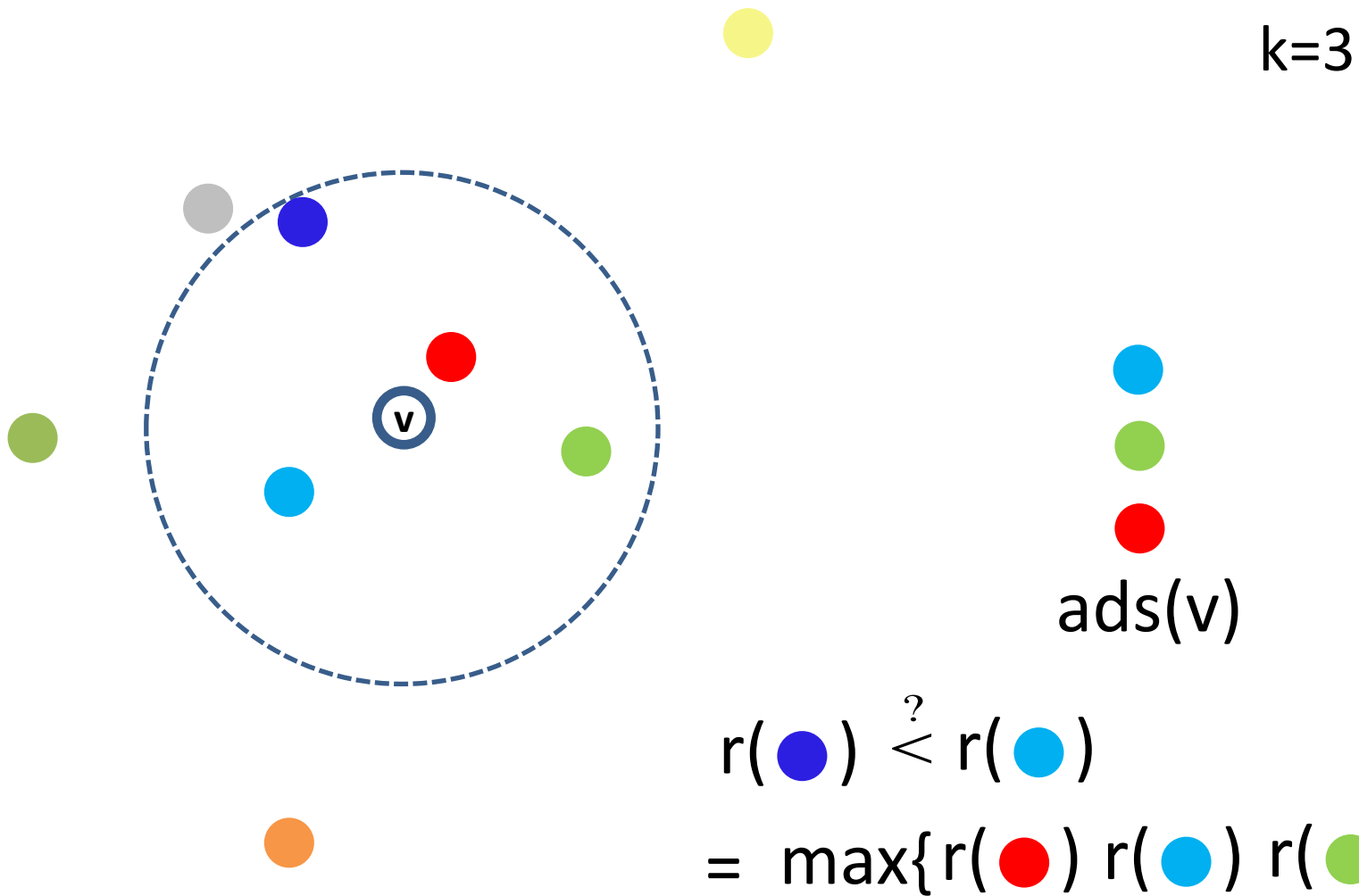


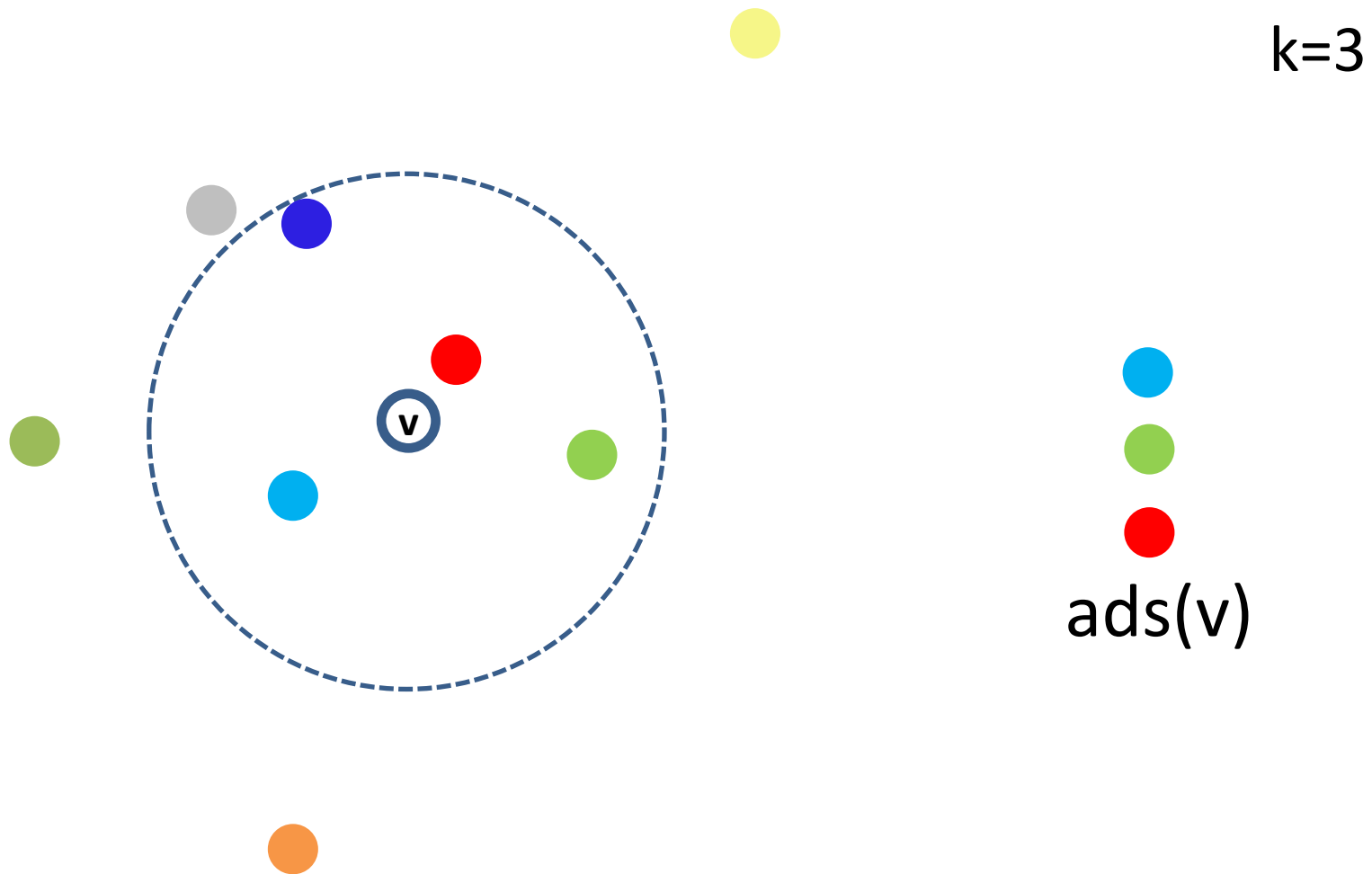
$k=3$

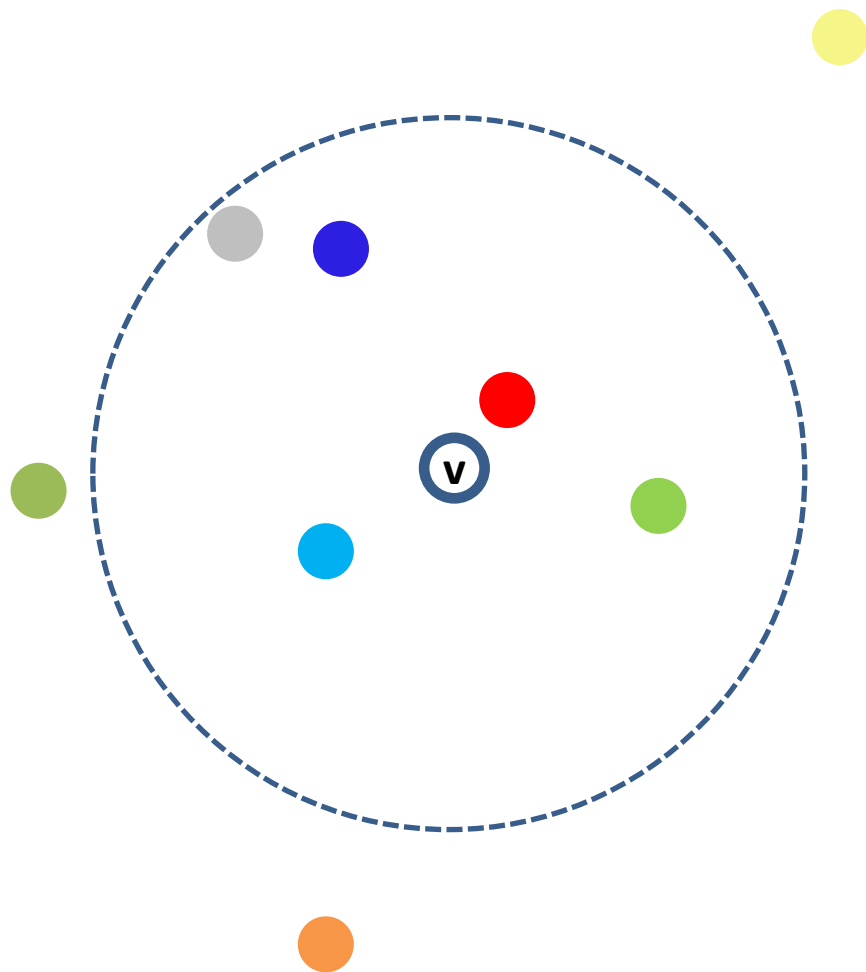


$a(\text{green node})=1$

$\text{ads}(v)$



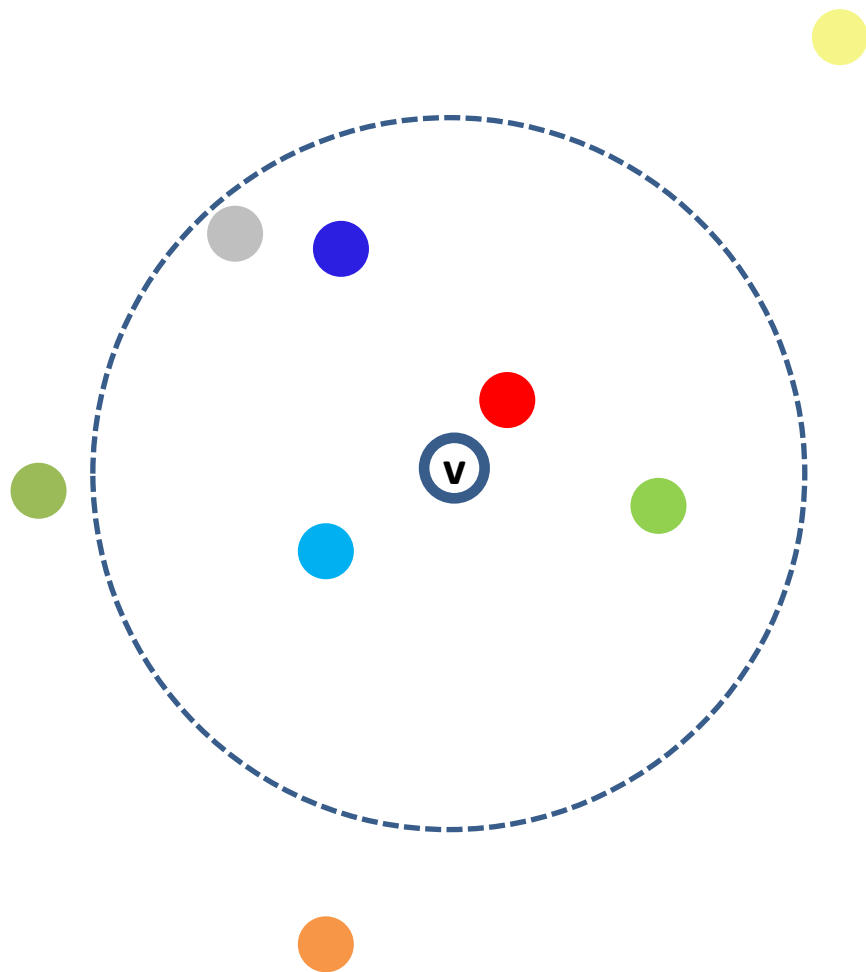




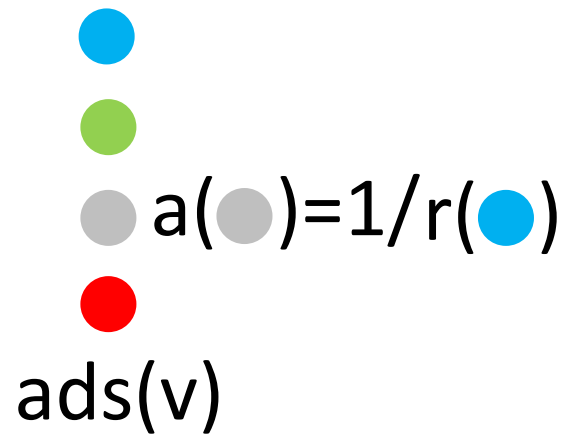
$k=3$

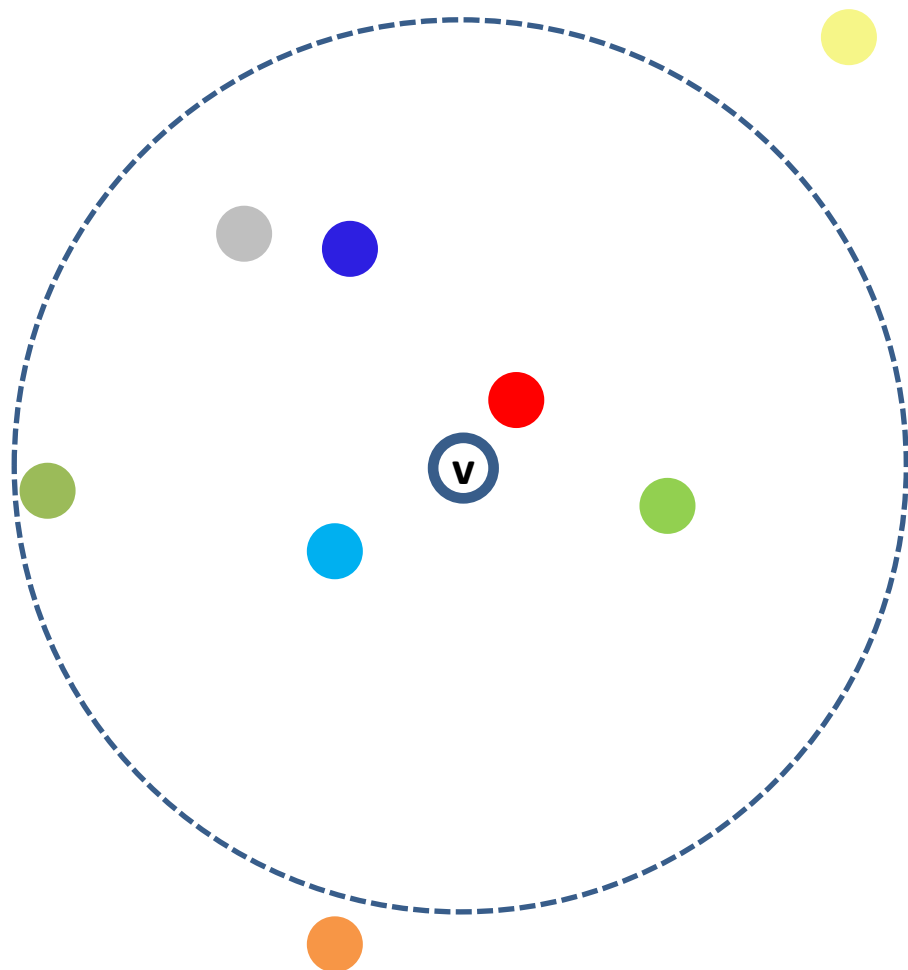
ads(v)

$$r(\text{grey}) \stackrel{?}{<} r(\text{cyan})$$



$k=3$



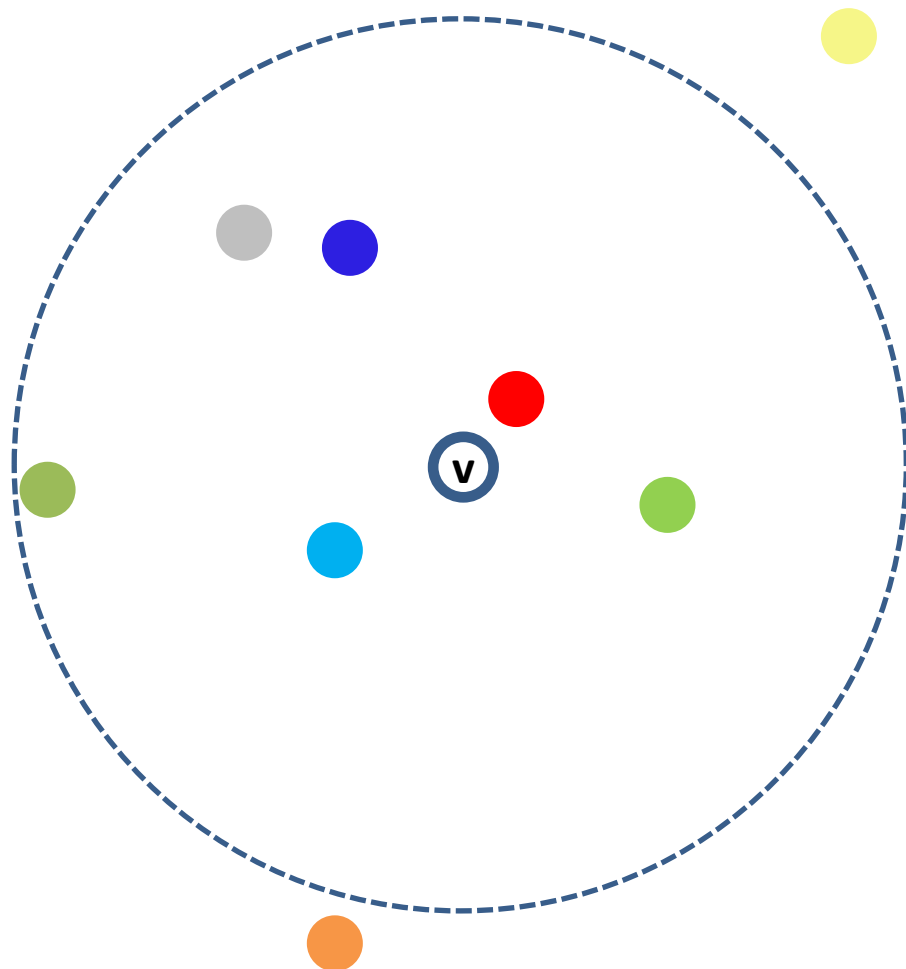


$k=3$

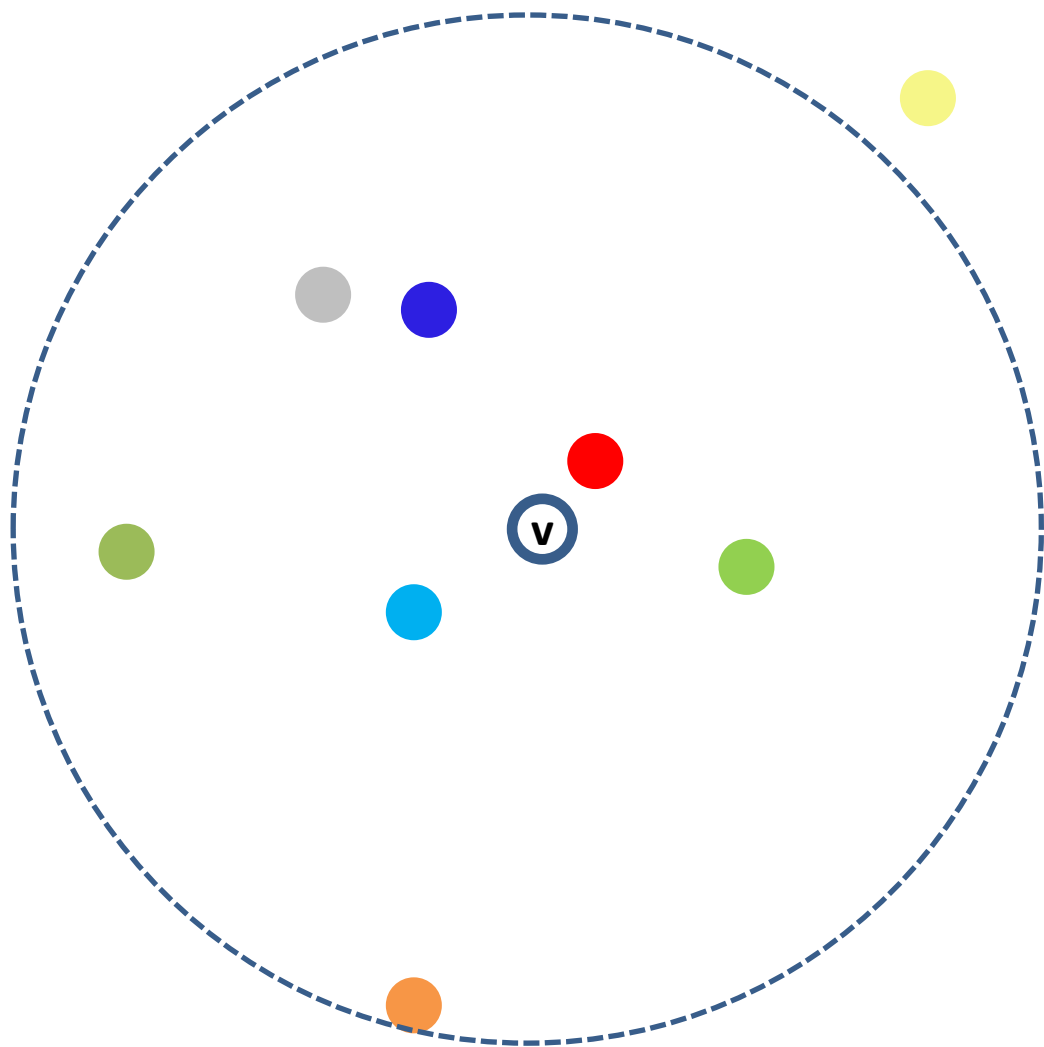
$\text{ads}(v)$

$r(\text{green point}) \stackrel{?}{<} r(\text{green point})$

$k=3$



ads(v)

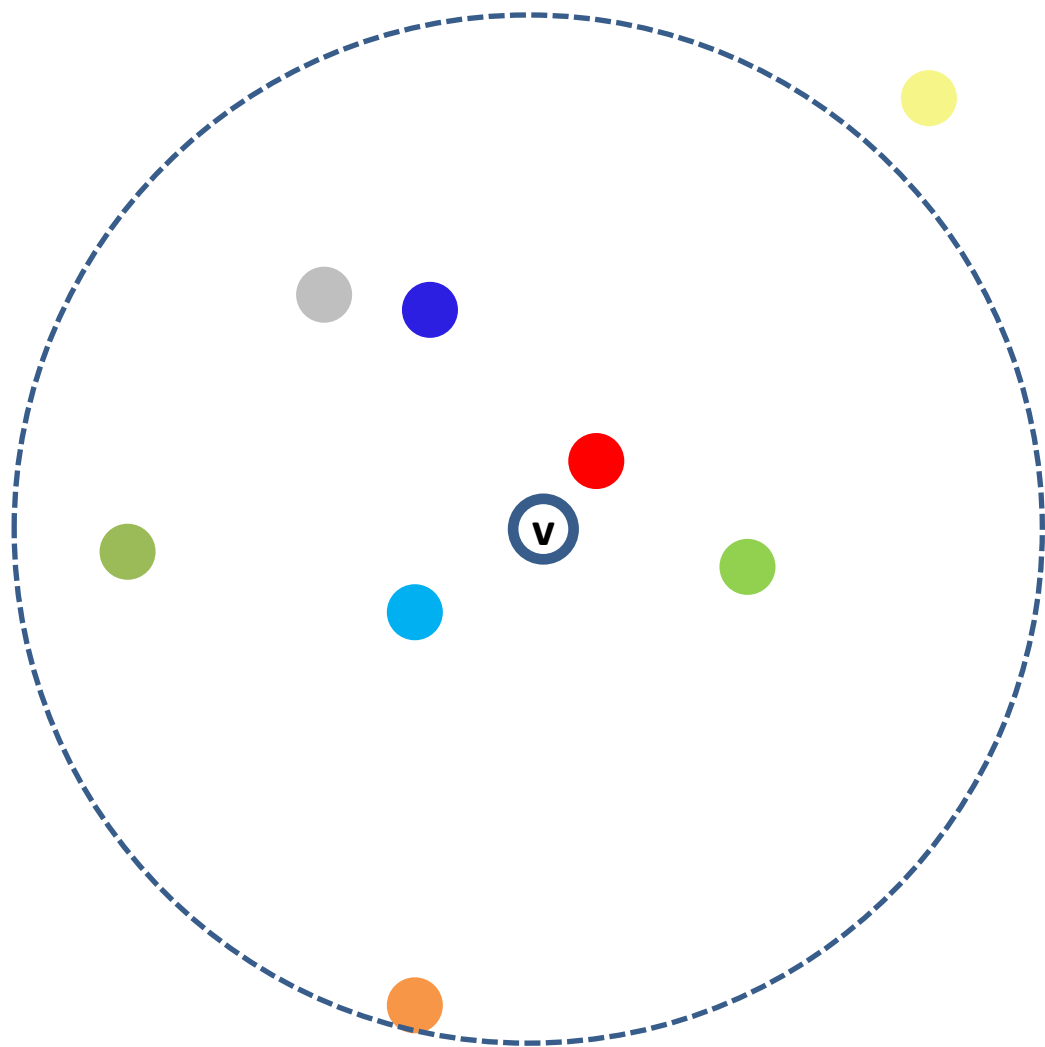


$k=3$



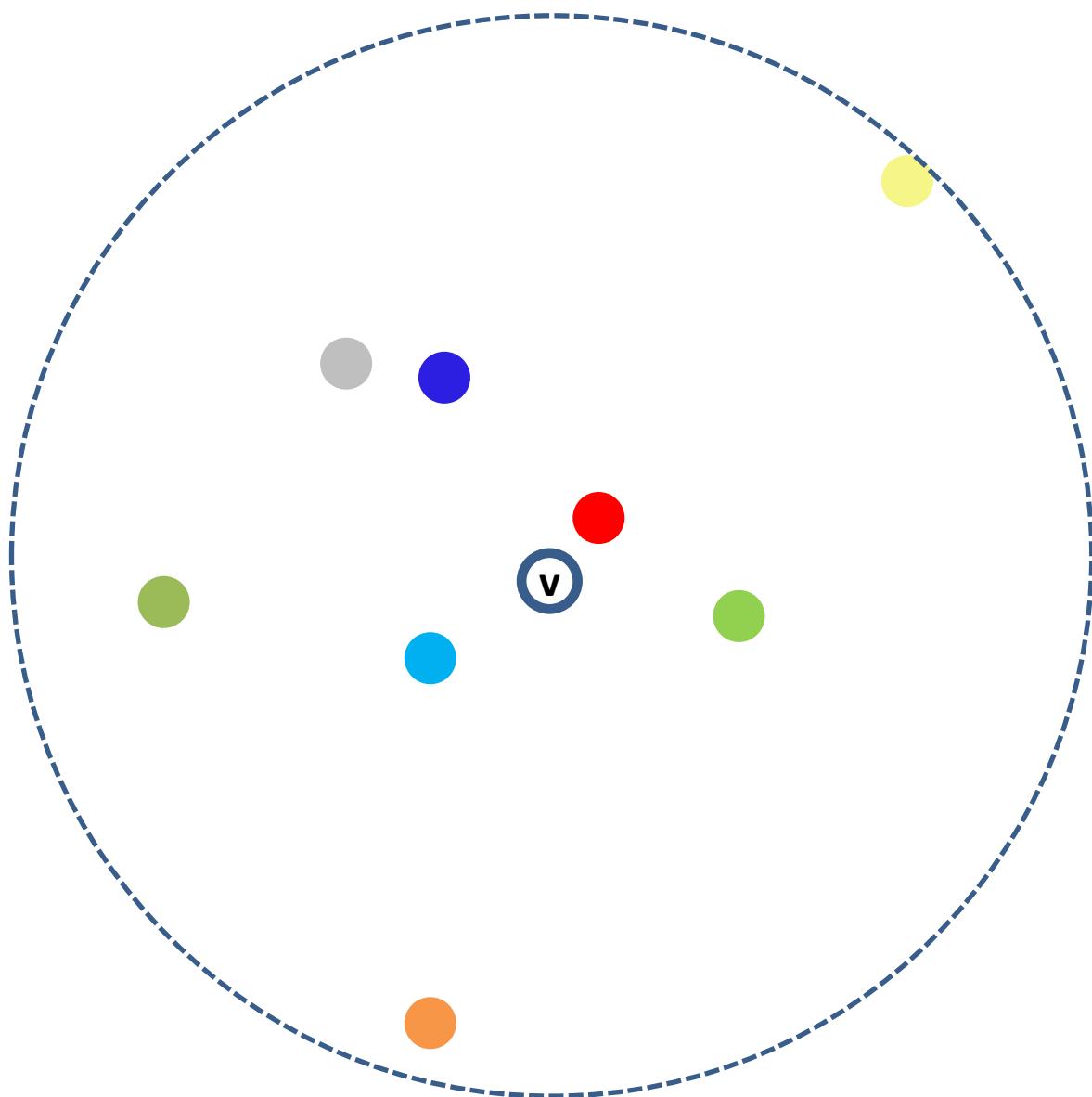
$\text{ads}(v)$

$$r(\text{orange circle}) \stackrel{?}{<} r(\text{green circle})$$



$k=3$

●
●
●
●
● $a(\text{orange}) = 1/r(\text{green})$
 $\text{ads}(v)$

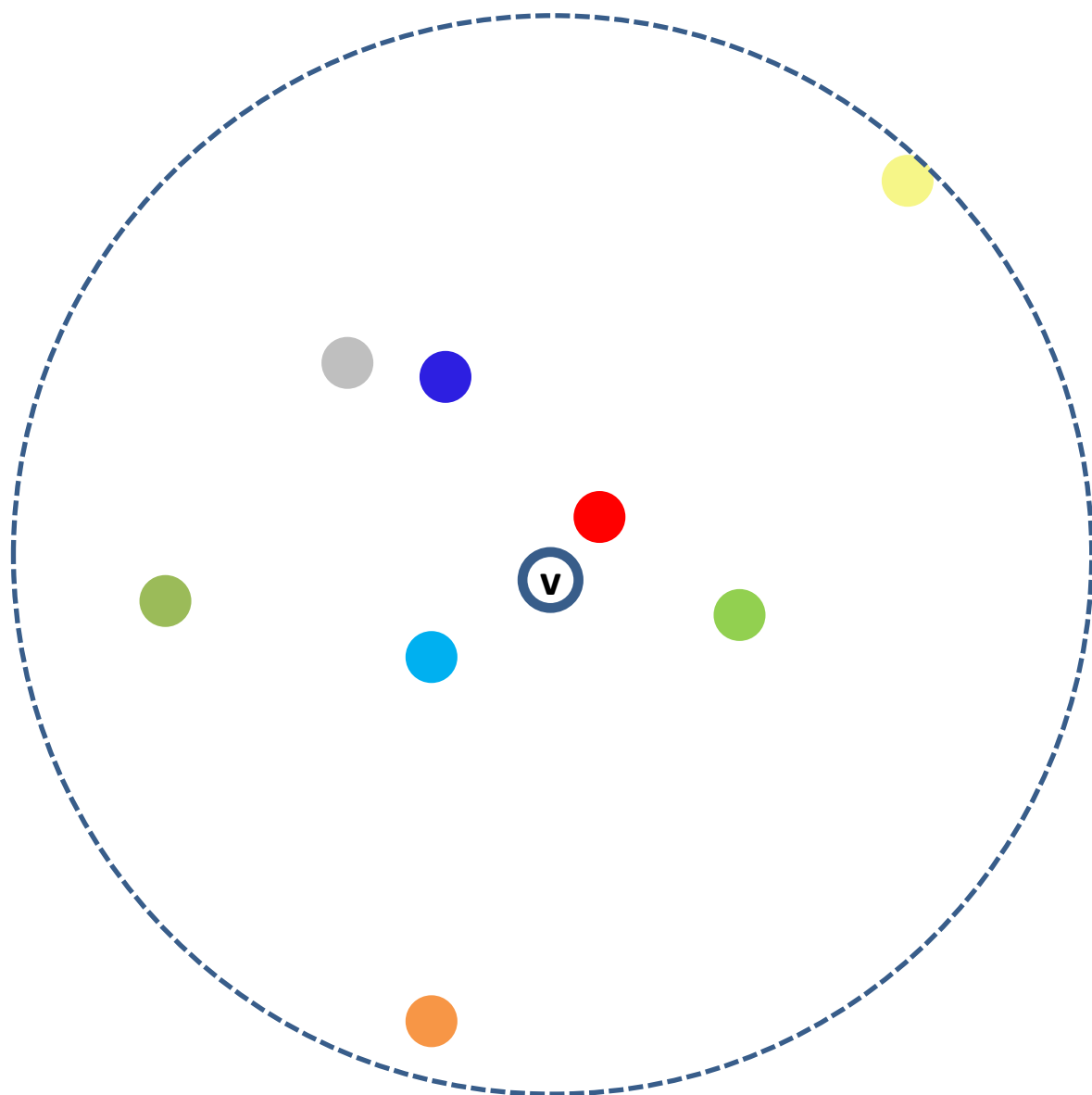


$k=3$



$\text{ads}(v)$

$$r(\text{yellow}) \stackrel{?}{<} r(\text{grey})$$



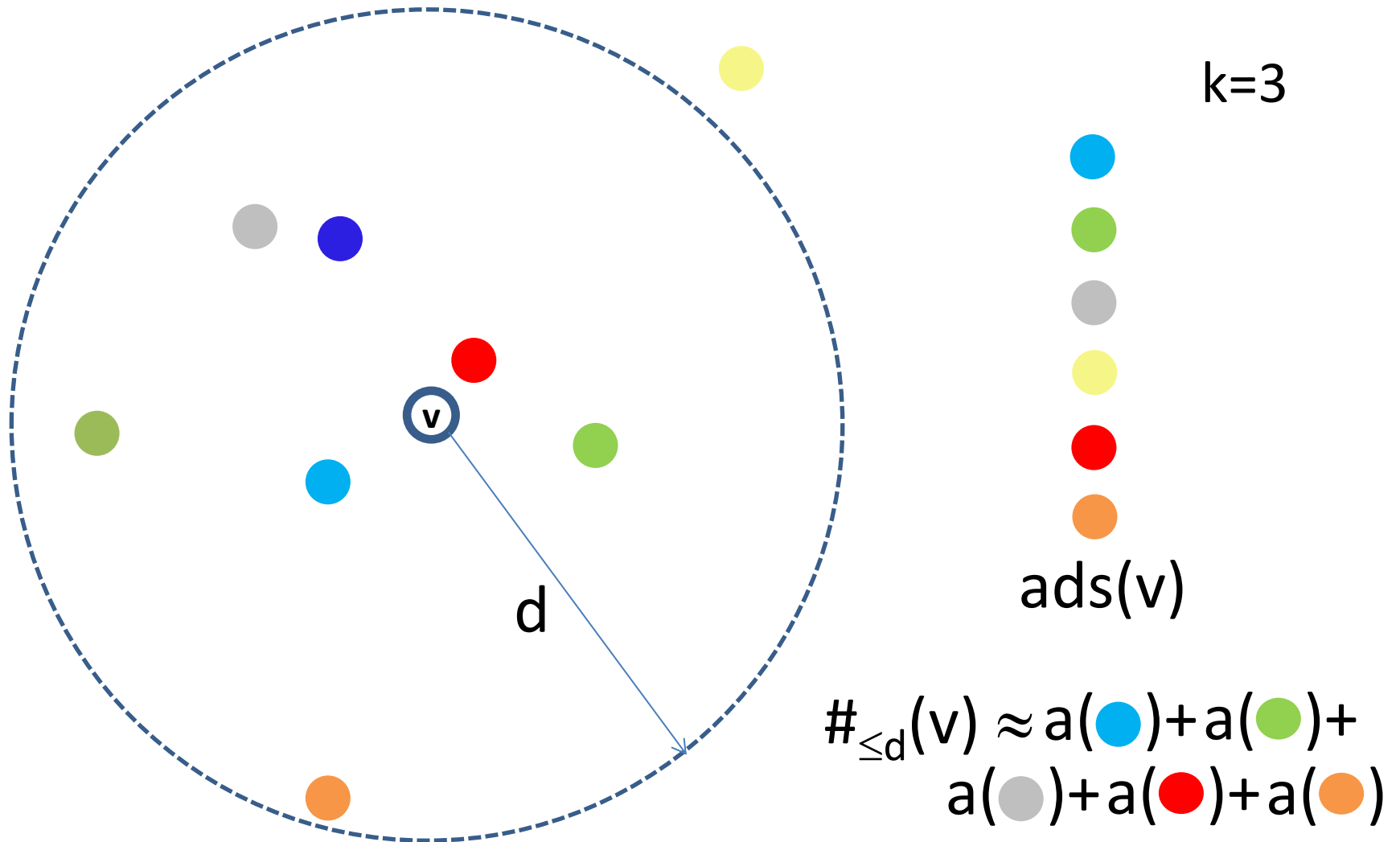
$k=3$



$$a(\text{yellow node}) = 1/r(\text{grey node})$$

$\text{ads}(v)$

Estimating



More material

- Edith Cohen: **All-distances sketches, revisited: HIP estimators for massive graphs analysis.** [PODS 2014](#)
- **Estimation for Monotone Sampling: Competitiveness and Customization.** [PODC 2014](#)
- Edith Cohen: **Distance Queries from Sampled Data: Accurate and Efficient.** [KDD 2014](#)
- Edith Cohen, Haim Kaplan: **What You Can Do with Coordinated Samples.** [APPROX-RANDOM 2013](#)
- Edith Cohen, Haim Kaplan: **Leveraging discarded samples for tighter estimation of multiple-set aggregates.** [SIGMETRICS/Performance 2009](#)
- **And earlier refs that you can find in the above**