

Safe Zones for Tracking the Median Function

Assaf Schuster, Technion

Daniel Keren, Haifa University

Tsachi Sharfman, Technion

Minos Garofalakis, Technical University of Crete

Vasilis Samoladas, Technical University of Crete

Goal

- Construct safe zones that ***contain*** the safe zones defined by the geometric method
 - Guarantees that resulting constraints remain effective for longer periods of time
 - Can potentially reduce communications by orders of magnitude

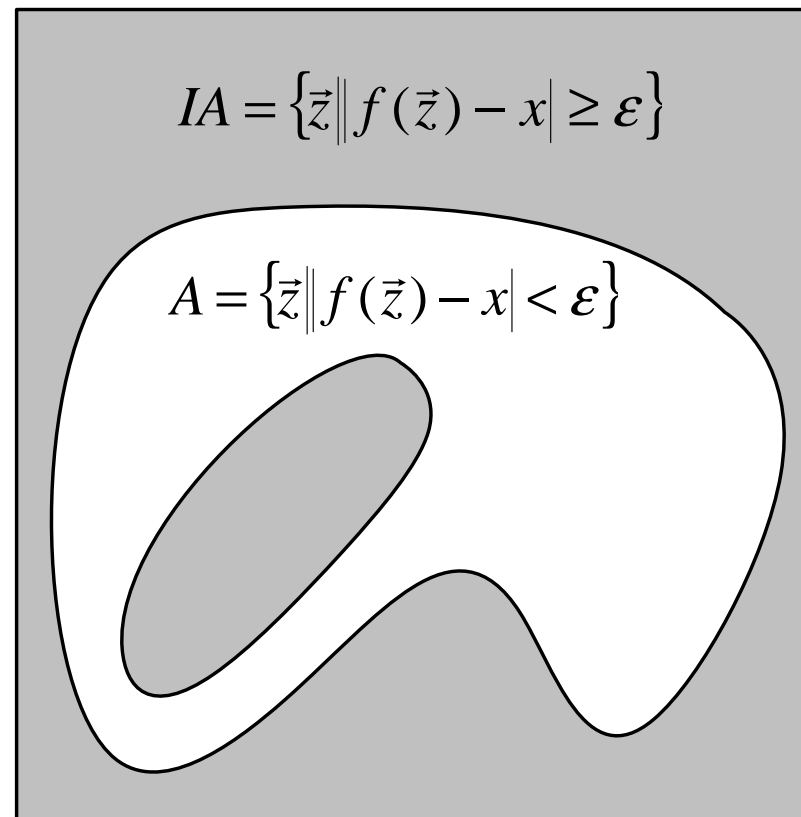
Tracking the Median Function

- We have n nodes, each holding a time varying d dimensional vector denoted by \mathbf{v}_i .
- We denote their average by \mathbf{v} .
- We are given the median function $med(\mathbf{v})$ and an approximation margin ε .
- At any time we would like to hold a number x , referred to as the approximation value, such that $|med(\mathbf{v}) - x| < \varepsilon$.

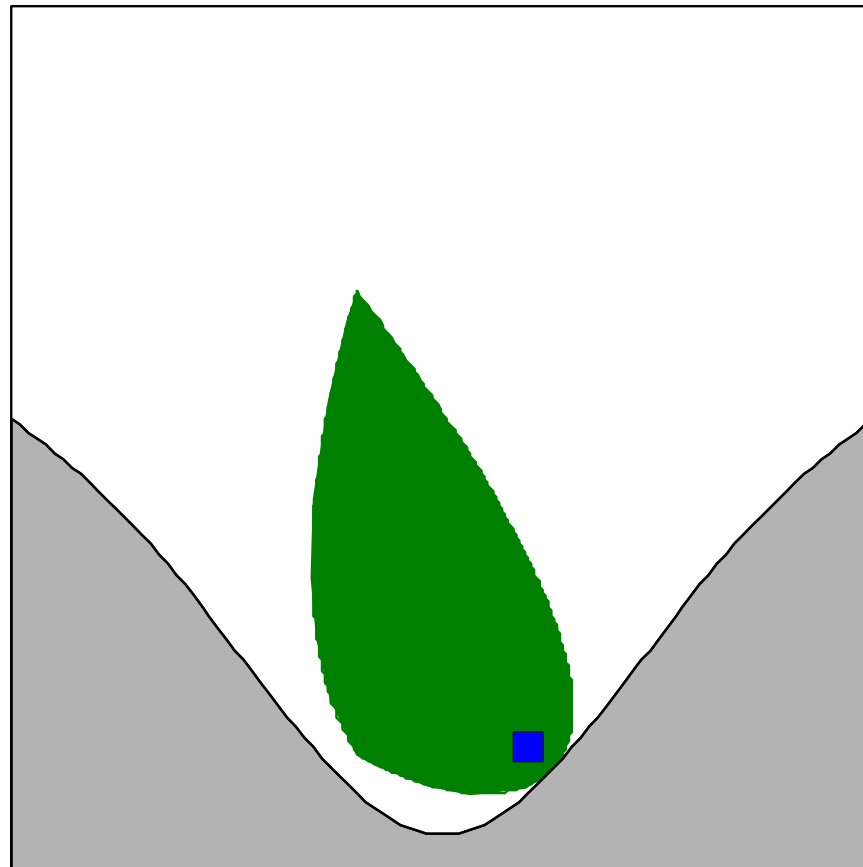
Admissible and Inadmissible Regions

- Given a function f , an estimate value x , and an approximation margin ε , we define the following:
 - Let the *admissible region*, denoted by $A(f,x,\varepsilon)$ be the set of vectors for which the value of the function is within the approximation range, i.e.
 $A(f,x,\varepsilon) = \{\mathbf{z} \mid |f(\mathbf{z}) - x| < \varepsilon\}$
 - Similarly, let the *inadmissible region*, denoted by $IA(f,x,\varepsilon)$ be the set of vectors for which the value of the function is not within the approximation range, i.e. $IA(f,x,\varepsilon) = \{\mathbf{z} \mid |f(\mathbf{z}) - x| \geq \varepsilon\}$

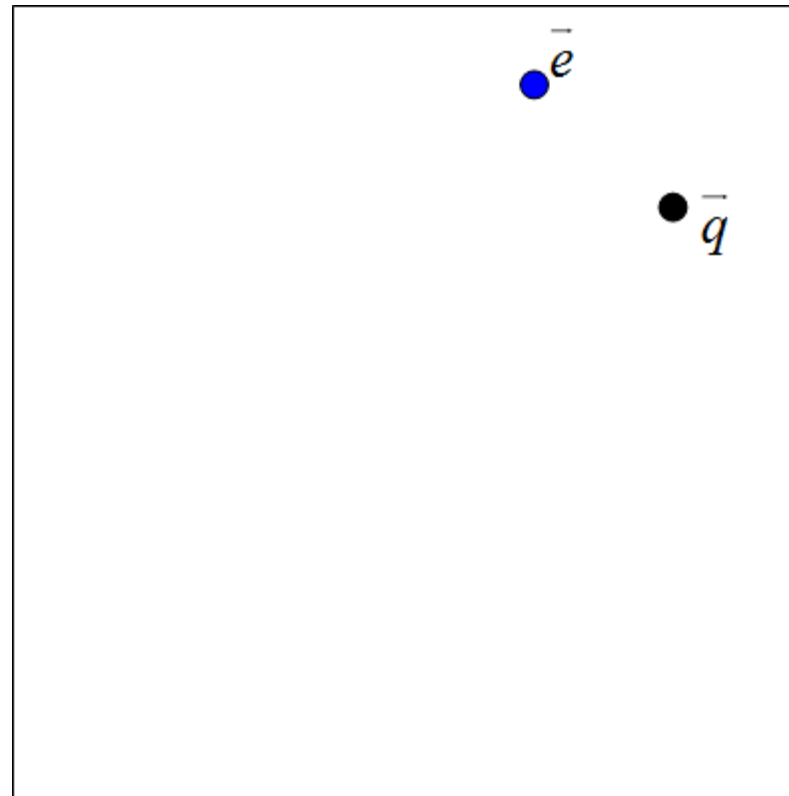
Admissible and Inadmissible Regions (Cont.)



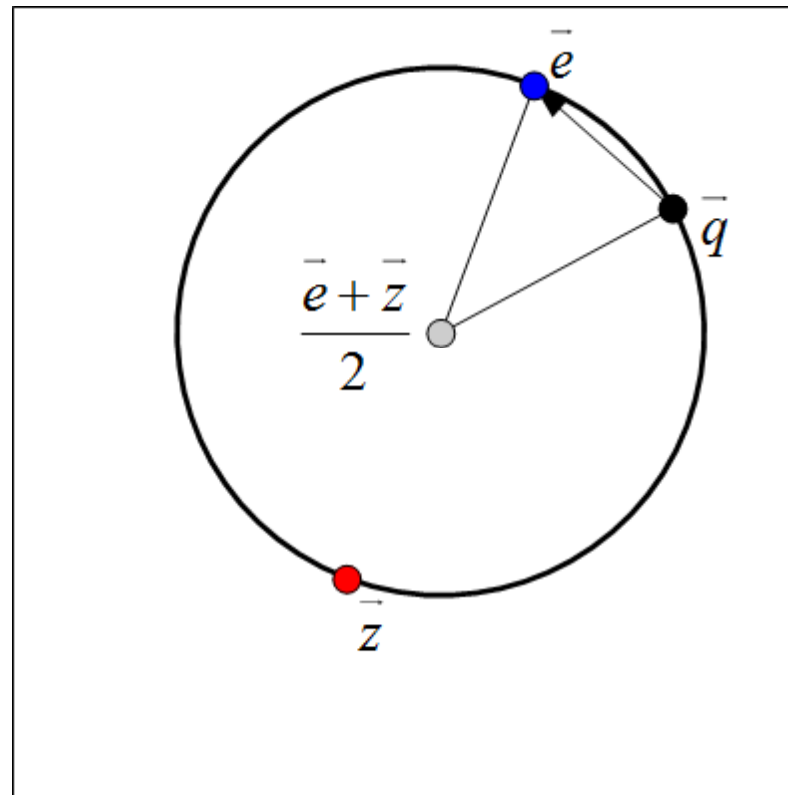
Safe Zone Induced by the Geometric Method



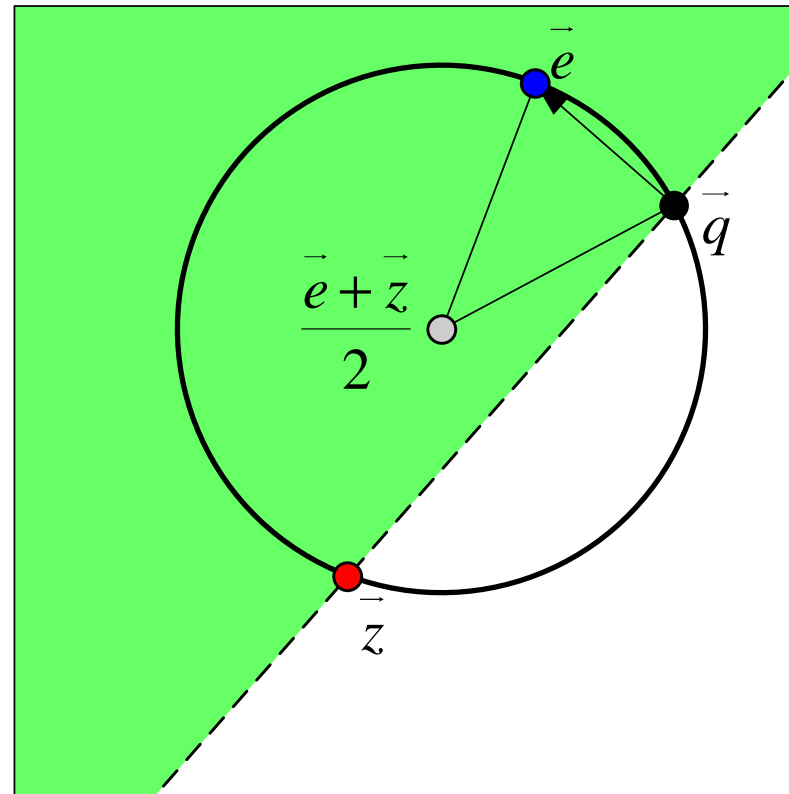
Safe Zone Induced by a Singular Vector



Safe Zone Induced by a Singular Vector



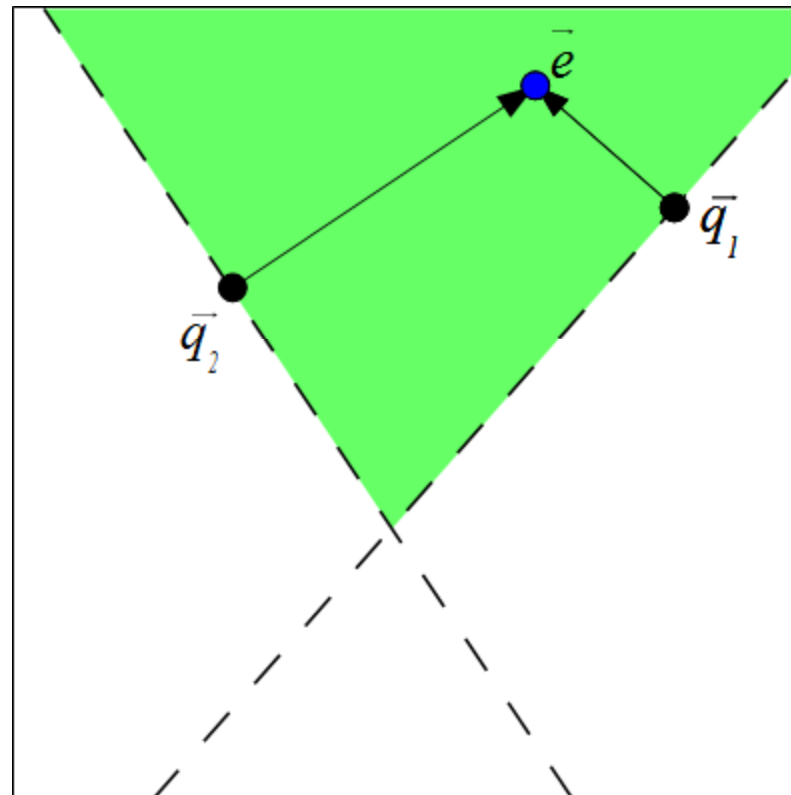
Safe Zone Induced by a Singular Vector



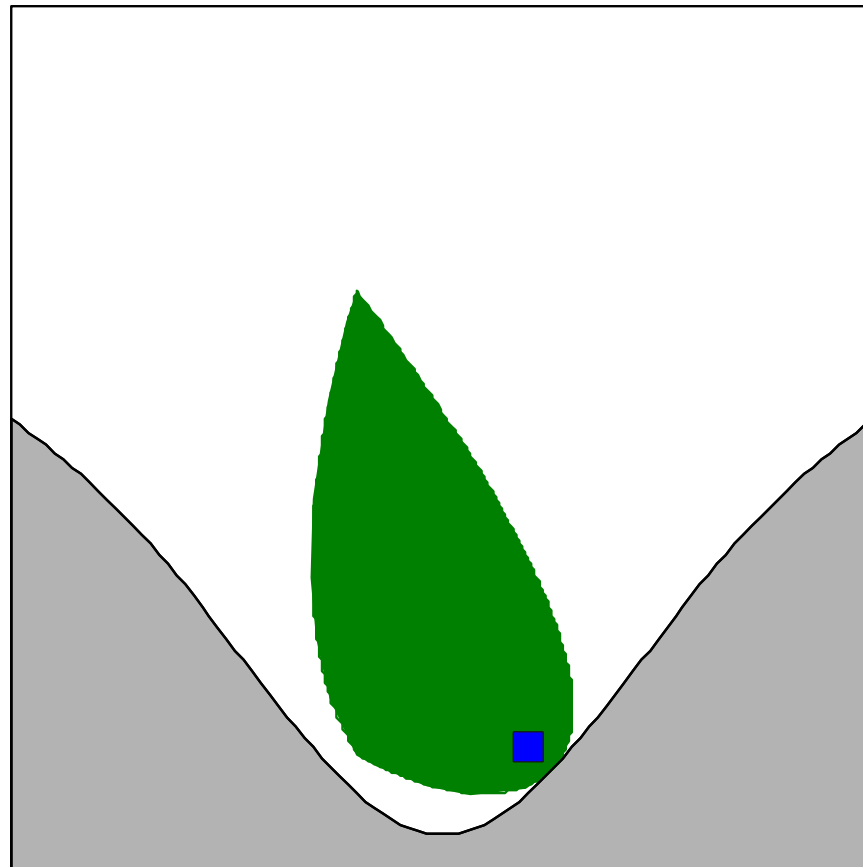
An Induced Half-Space

- Given an inadmissible vector \mathbf{q} and an estimate vector \mathbf{e} , let $H(\mathbf{e}, \mathbf{q})$ be the half-space defined by the hyper-plane passing through \mathbf{q} , orthogonal to $\mathbf{e} - \mathbf{q}$, and containing \mathbf{e} . We refer to $H(\mathbf{e}, \mathbf{q})$ as the half-space induced by \mathbf{q} .

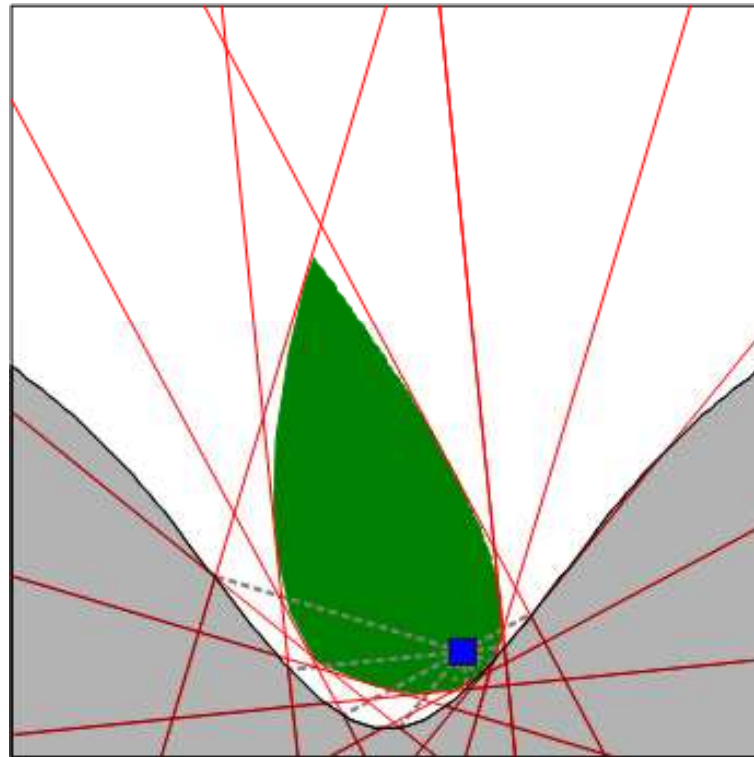
Safe Zone Induced by Two Singular Vectors



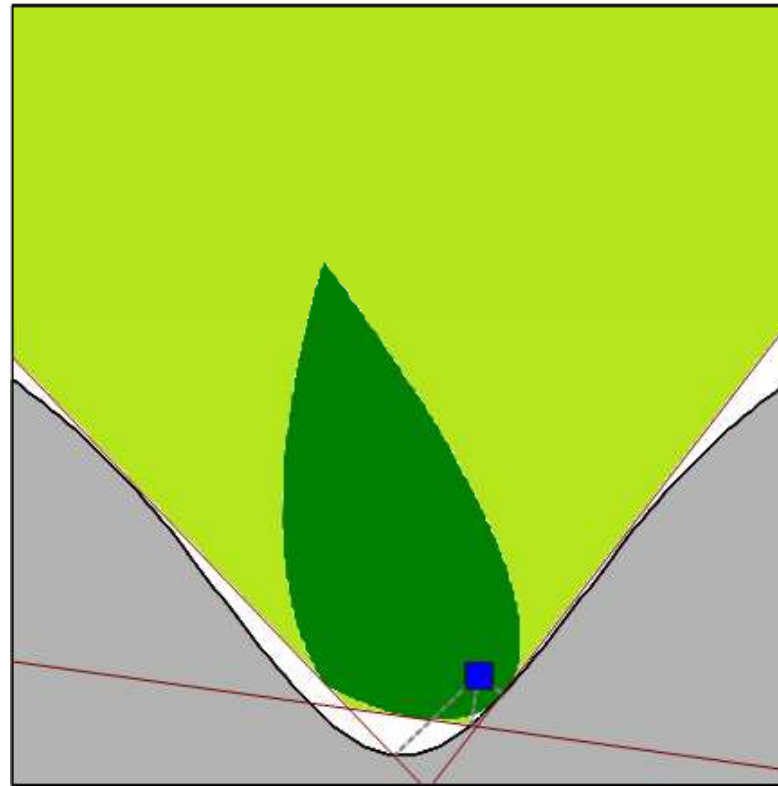
Safe Zone Induced by the Geometric Method



Safe Zone as an Intersection of Half-Spaces



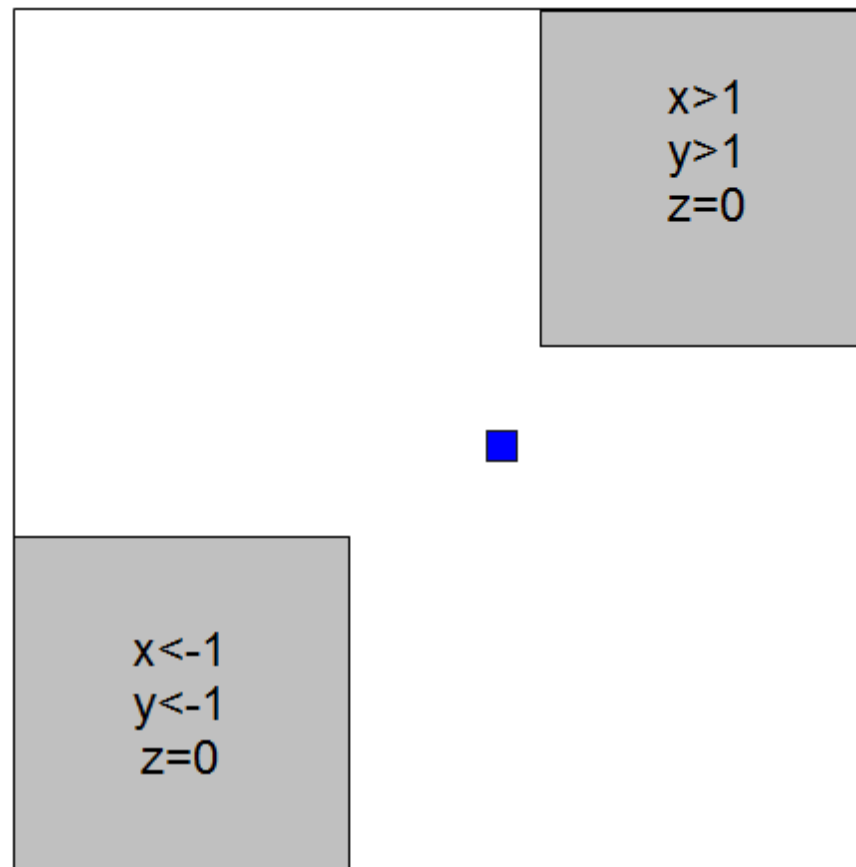
Safe Zone as an Intersection of Half-Spaces



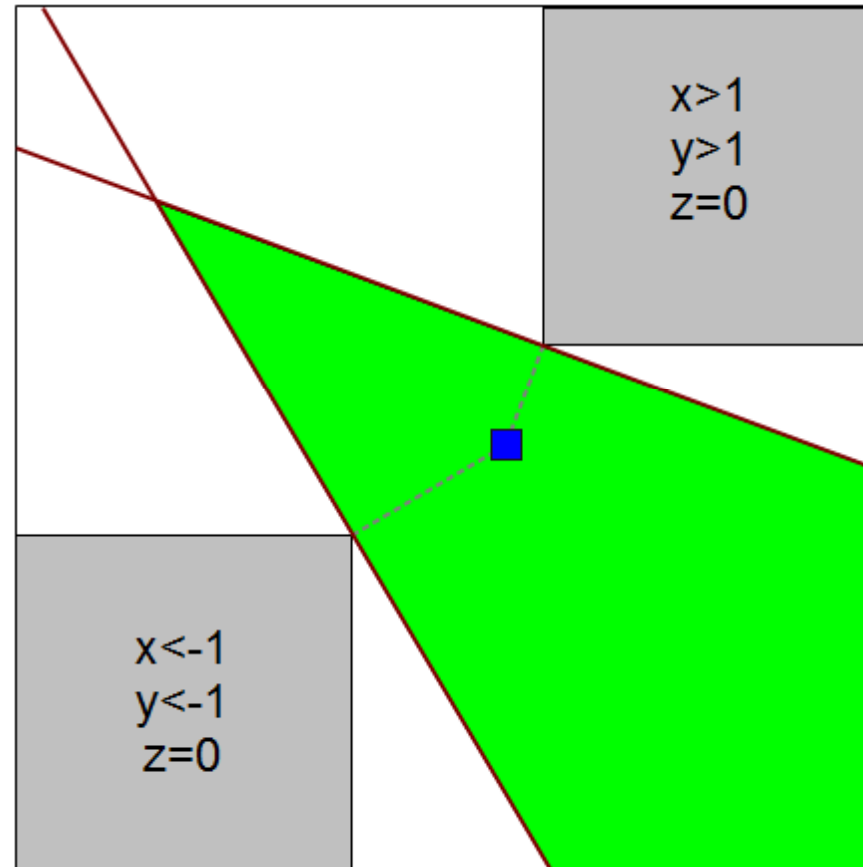
Support Vectors

- Given an inadmissible region and an estimate vector, we would like to select a small set of vectors from the set such that the intersection of the half spaces they induce is contained in the admissible region.
- We refer to these vectors as support vectors

The Median Function



The Median Function

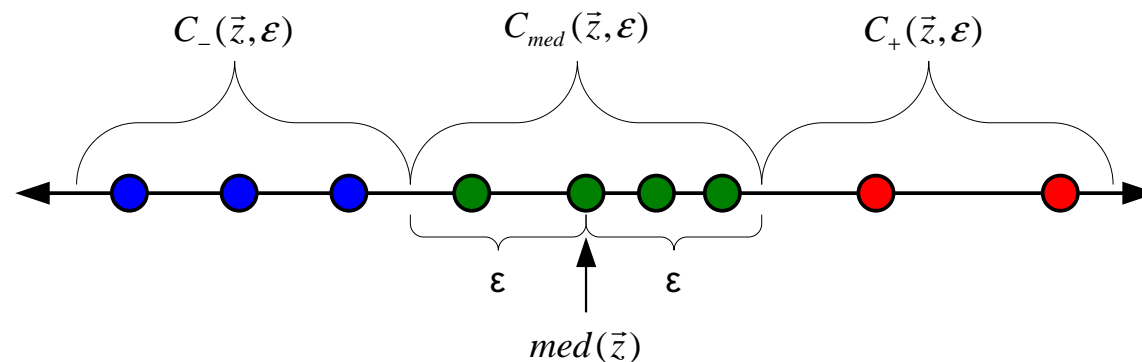


Notations

- Let I denote a subset of the natural numbers in the range $\{1, \dots, d\}$.
- Let r be a real value, and let \mathbf{z} be a d -dimensional vector. Let $\mathbf{z}^{(I,r)}$ denote the vector \mathbf{z} with the value of the components whose index is in I set to r .
- Let $\binom{I}{m}$ be the set of all m sized subsets of I .

Definitions

- Given a d -dimensional vector \mathbf{z} , an estimate value x and an approximation margin ε , we group the vector's components as follows:
 - $C_+(\mathbf{z}, \varepsilon)$: A set of the indices of the components of \mathbf{z} whose value is greater than or equal to $median(\mathbf{z}) + \varepsilon$.
 - $C_{med}(\mathbf{z}, \varepsilon)$: A set of the indices of the components of \mathbf{z} whose value is in the range $median(\mathbf{z}) \pm \varepsilon$.
 - $C_-(\mathbf{z}, \varepsilon)$: A set of the indices of the components of \mathbf{z} whose value is smaller than or equal to $median(\mathbf{z}) - \varepsilon$.



Support Vectors

- The safe zone consists of the half-spaces defined by the estimate vector and a set of vectors referred to as the support vectors.
- The set of support vectors is the union of the following two sets:

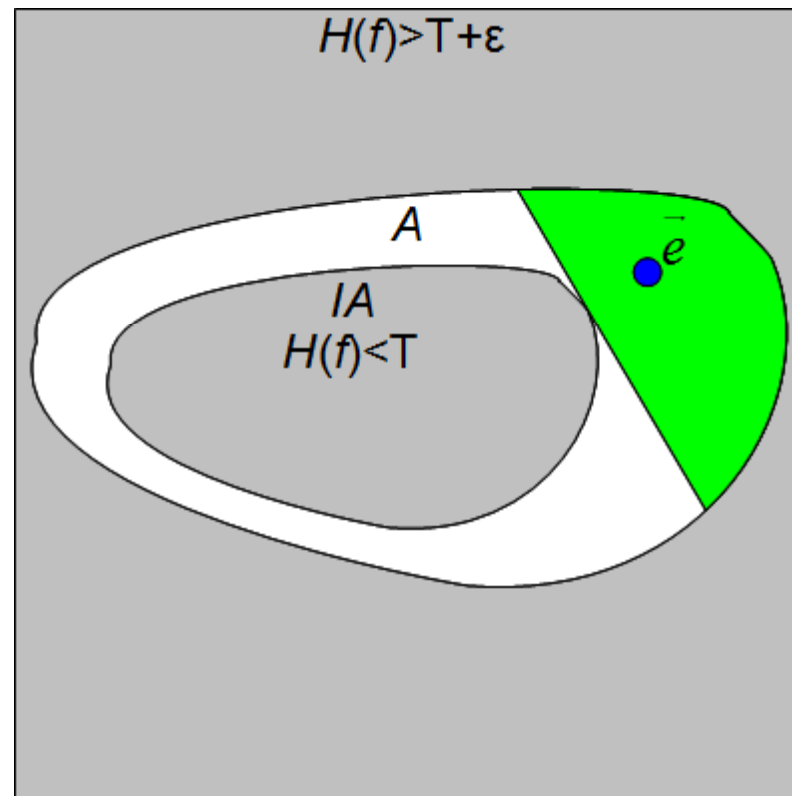
$$Support_+(\vec{e}, \varepsilon) = \left\{ \vec{e}^{(I, med(\vec{e}) - \varepsilon)} \mid I \in \left[\begin{array}{l} C_{med}(\vec{e}, \varepsilon) \cup C_+(\vec{e}, \varepsilon) \\ \frac{d+1}{2} - \|C_-(\vec{e}, \varepsilon)\| \end{array} \right] \right\}$$

$$Support_-(\vec{e}, \varepsilon) = \left\{ \vec{e}^{(I, med(\vec{e}) + \varepsilon)} \mid I \in \left[\begin{array}{l} C_{med}(\vec{e}, \varepsilon) \cup C_-(\vec{e}, \varepsilon) \\ \frac{d+1}{2} - \|C_+(\vec{e}, \varepsilon)\| \end{array} \right] \right\}$$

Monitoring Entropy

- We are given a distributed set of time-varying frequency vectors \mathbf{f}_i , let their \mathbf{f} denote their sum.
- We would like to monitor the entropy of the sum, $H(\mathbf{f})$, with an approximation error of ε .
- It is easy to show that the inadmissible region IA defined by $H(\mathbf{f}) > T$ is convex.
- Consequently, the admissible region A defined by $H(\mathbf{f}) < T + \varepsilon$ is convex.
- Given an estimate vector \mathbf{e} (that belongs to $A \setminus IA$) let \mathbf{e}^* be the vector in IA that is closest to \mathbf{e} .
- We can construct a safe zone by taking the intersection of $H(\mathbf{e}, \mathbf{e}^*)$ and A .

Monitoring Entropy (Cont')



Summary and Future Work

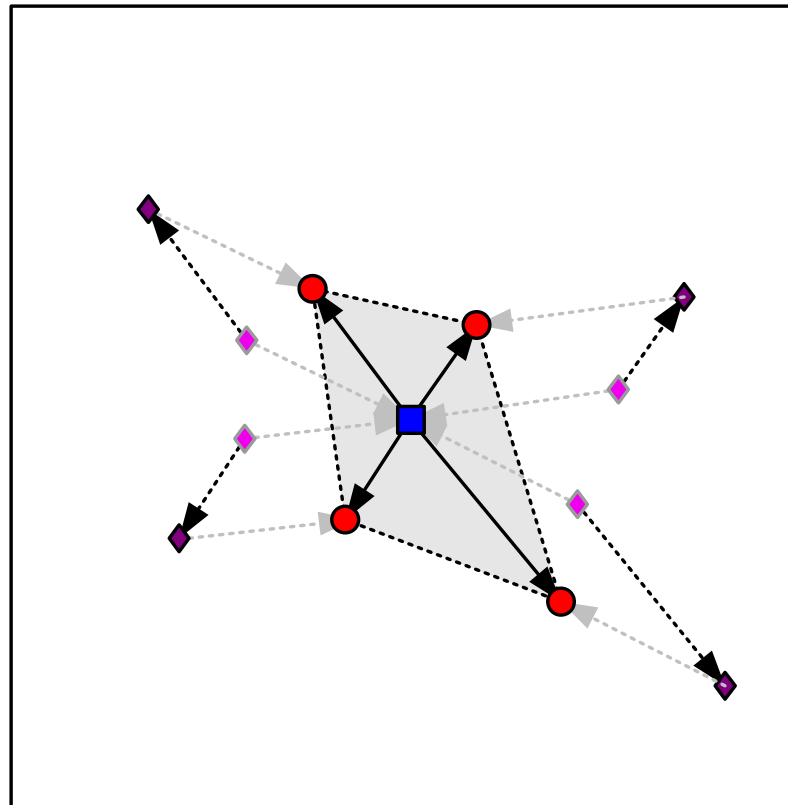
- Extend technique to general AMS sketches
- Extend technique to generic families of functions
- Refine and extend the notion of optimality

Backup Slides

Algorithmic Framework

- From time to time, as specified by the protocol, we will perform a synchronization procedure, which consists of:
 - Collecting local vectors from nodes. Value sent by the i^{th} vector is denoted by \mathbf{v}'_i .
 - Determining their average, denoted by \mathbf{e} .
 - Sending the average vector back to the nodes.
- After synchronization:
 - We set $x = f(\mathbf{e})$
 - At each node determine a vector $\delta_i = \mathbf{e} - \mathbf{v}'_i$. This value remains fixed until the next synchronization event.
 - Each node holds an additional vector \mathbf{u}_i referred to as the drift vector, $\mathbf{u}_i = \mathbf{v}_i + \delta_i$
- At any time, the average of the drift vectors is equal to the average of data vectors. After synchronization drift vectors are equal to the estimate vector.

Algorithmic Framework (cont.)



AMS Sketches

- Given two (very) large d vectors, \mathbf{v}_a and \mathbf{v}_b , their inner product $\langle \mathbf{v}_a, \mathbf{v}_b \rangle$ can be approximated as follows:
 - Given (ϵ, δ) , let $n = O(1/\epsilon^2)$ and $m = O(\ln(1/\delta))$
 - Each vector \mathbf{v} is represented by a $m \times n$ matrix $\mathbf{S}(\mathbf{v})$ referred as the vector's sketch:
 - We use $m \times n$ four-wise independent d dimensional random vectors \mathbf{r}_{ij} where each component receives either 1 or -1 with equal probability
 - Given a vector \mathbf{v} let $\mathbf{S}(\mathbf{v})_{ij} = \langle \mathbf{v}, \mathbf{r}_{ij} \rangle$
 - Let $\mathbf{S}(\mathbf{v}_a) \bullet \mathbf{S}(\mathbf{v}_b)$ denote the component-wise product of $\mathbf{S}(\mathbf{v}_a)$ and $\mathbf{S}(\mathbf{v}_b)$
 - The inner product of \mathbf{v}_a and \mathbf{v}_b is approximated by the median of the average of the rows of $\mathbf{S}(\mathbf{v}_a) \bullet \mathbf{S}(\mathbf{v}_b)$
- The AMS estimate has a error of $\epsilon \|\mathbf{v}_a\| \cdot \|\mathbf{v}_b\|$ with probability of are least $1 - \delta$