Adaptive Bounding Subgradient Method for Convex Functions

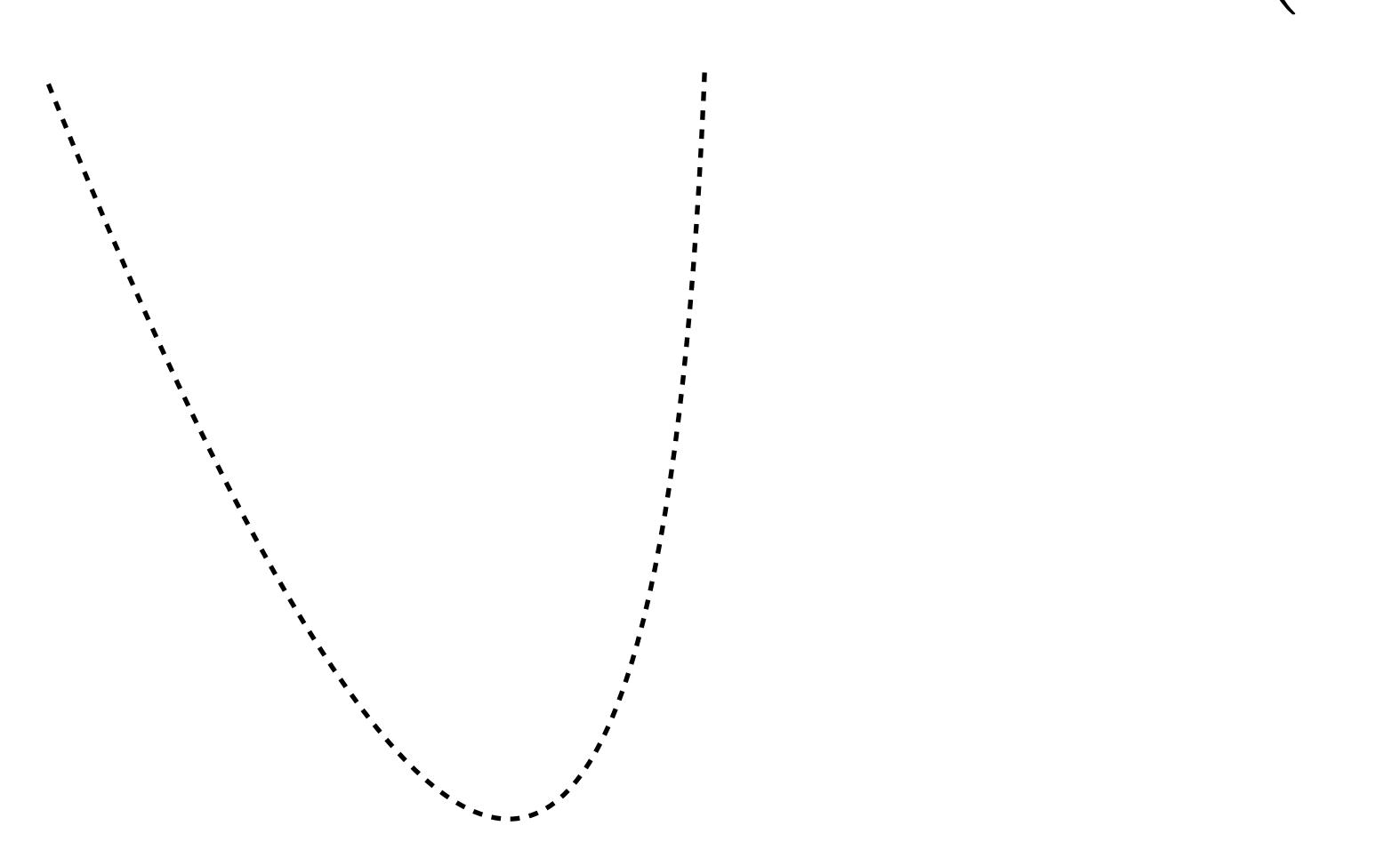
Nicholas Dwork, Charles Tsao, Cody Coleman, John M. Pauly Radiology, University of California in San Francisco Electrical Engineering, Stanford University



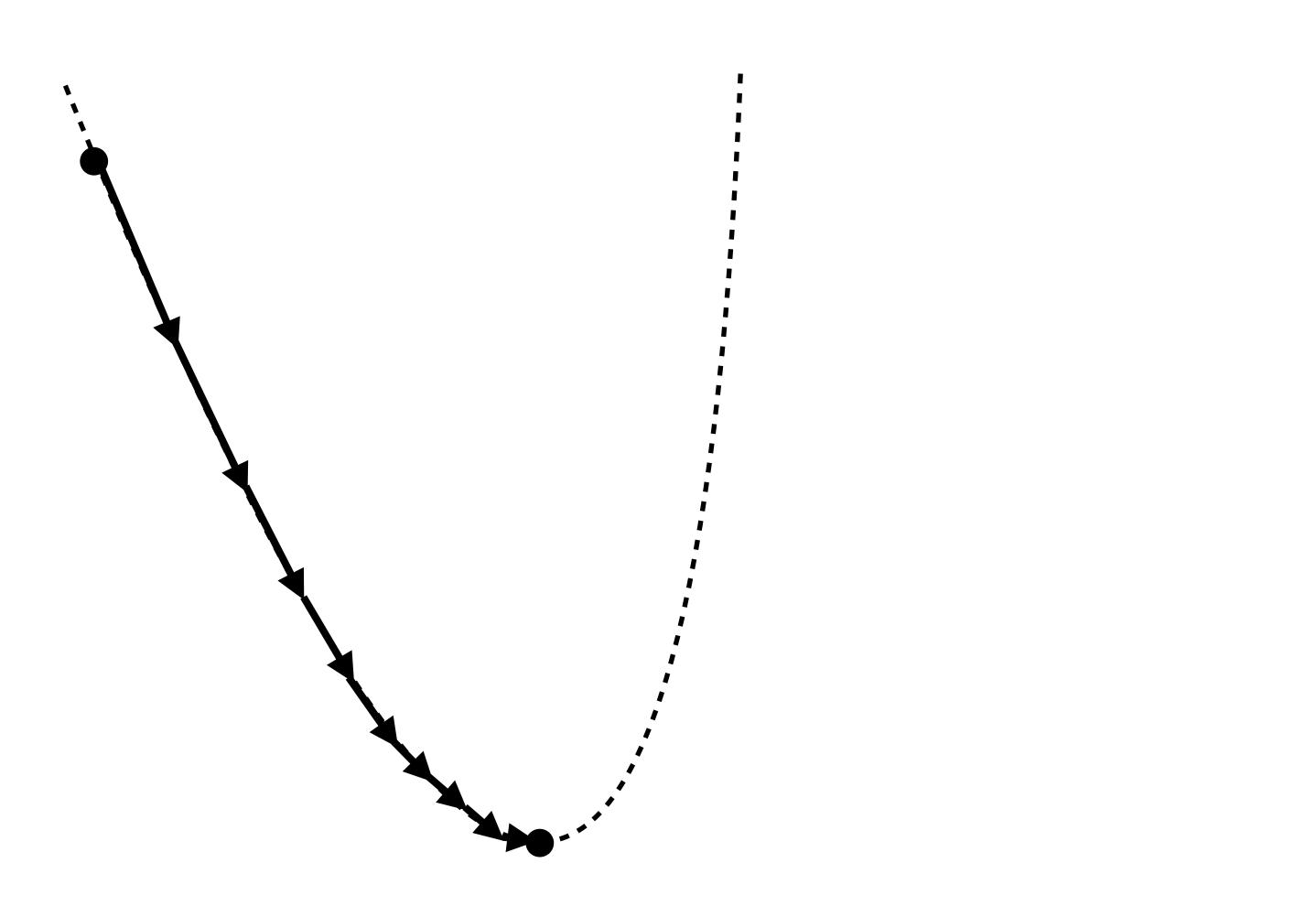


$$x^{(k+1)} = x^{(k)} - \alpha g(x^{(k)})$$

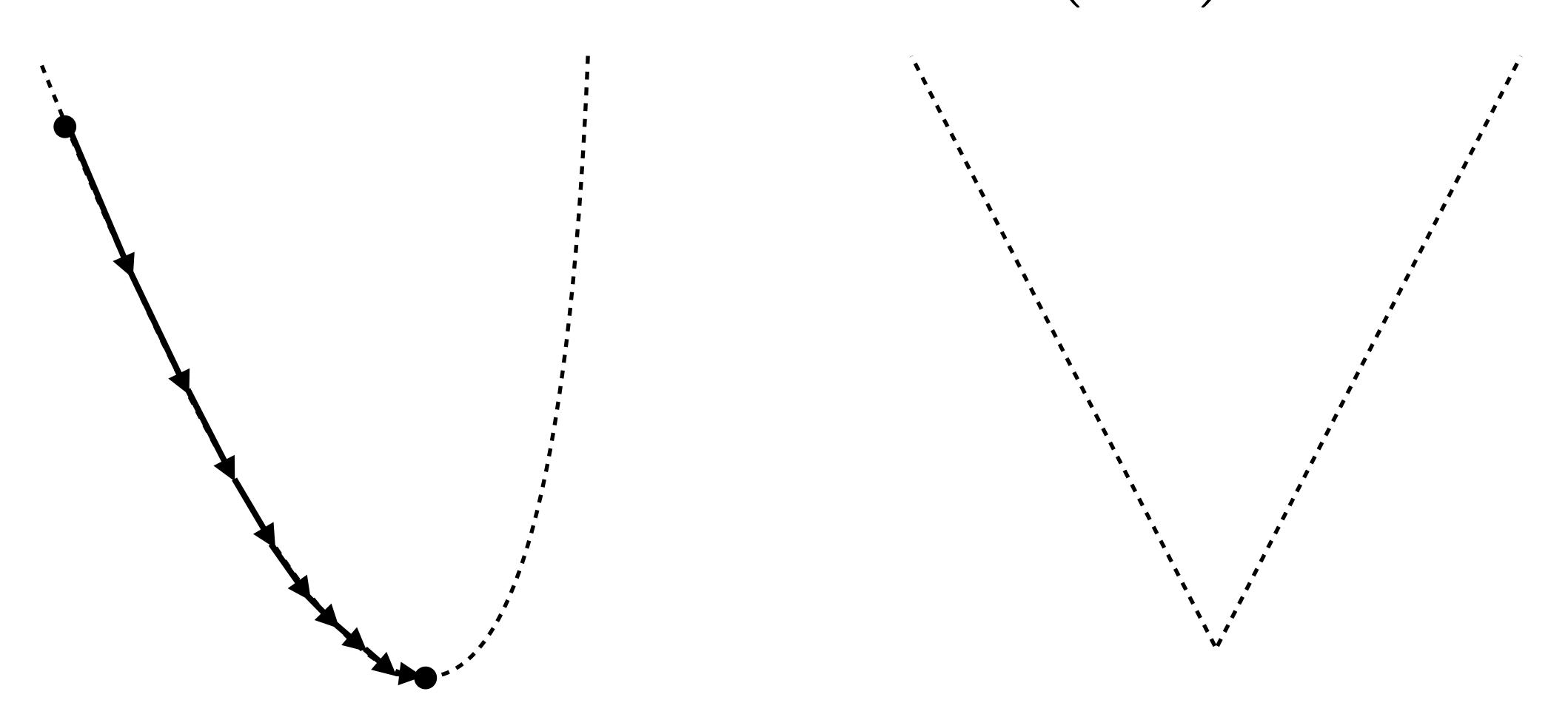
$$x^{(k+1)} = x^{(k)} - \alpha g(x^{(k)})$$



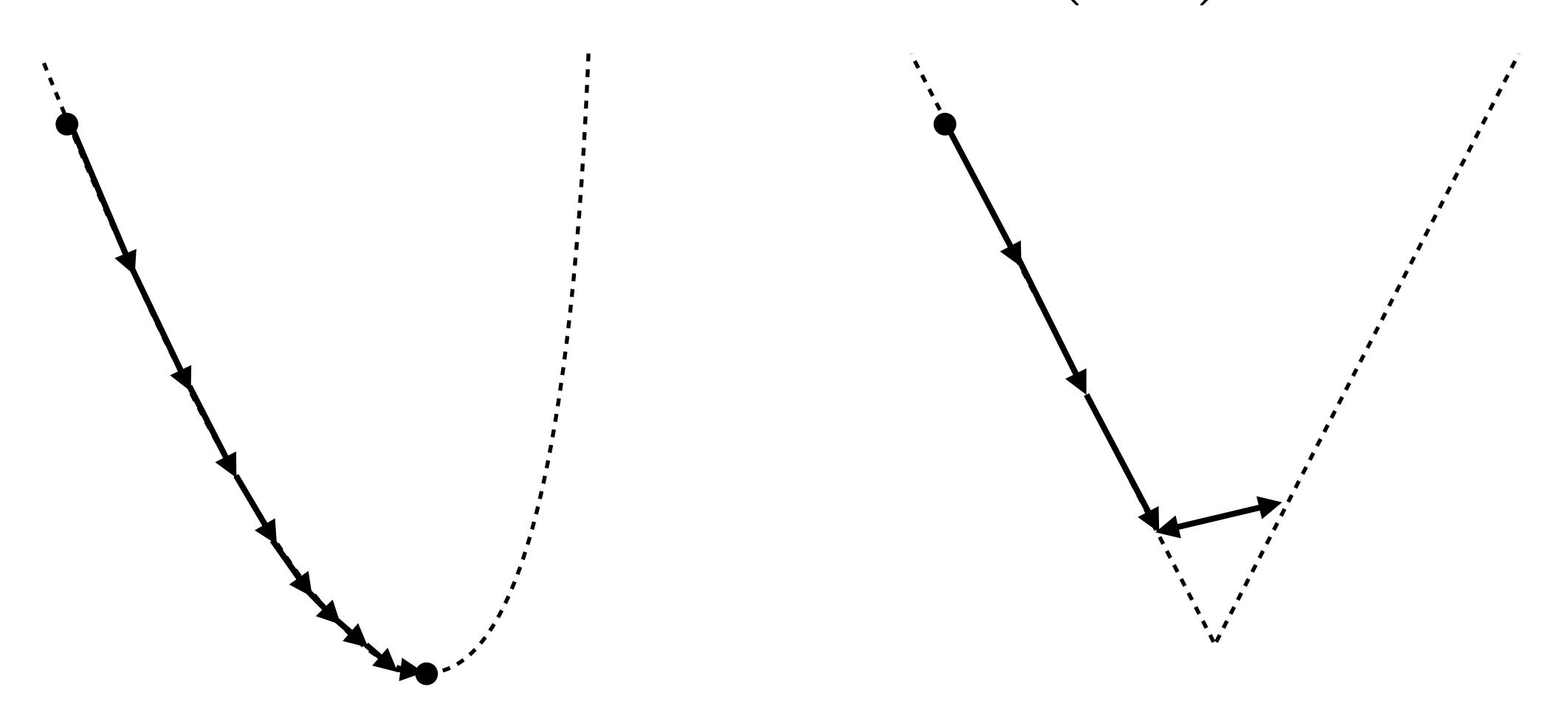
$$x^{(k+1)} = x^{(k)} - \alpha g(x^{(k)})$$



$$x^{(k+1)} = x^{(k)} - \alpha g(x^{(k)})$$

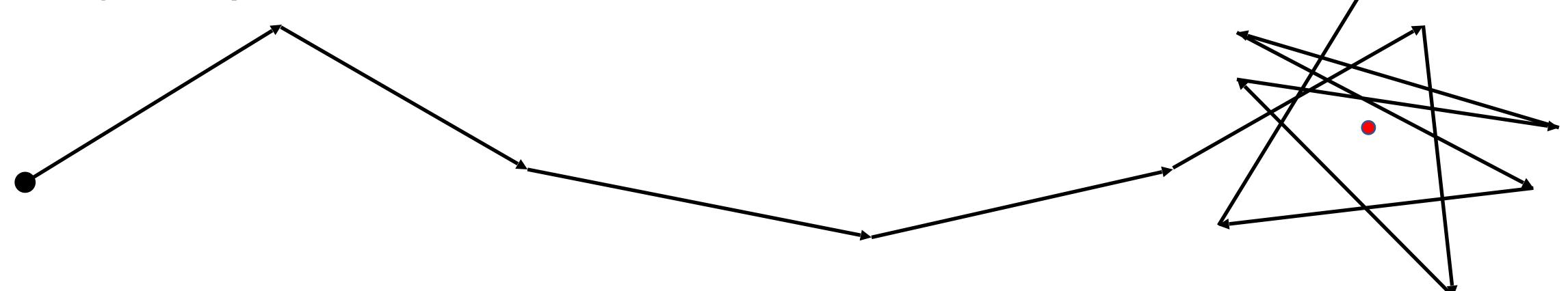


$$x^{(k+1)} = x^{(k)} - \alpha g(x^{(k)})$$

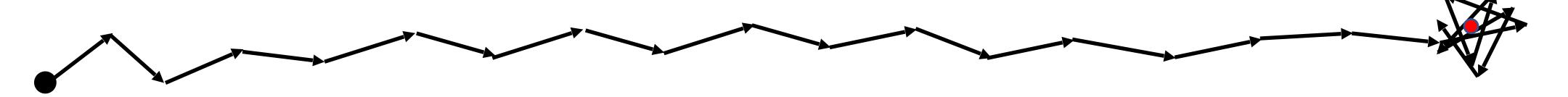


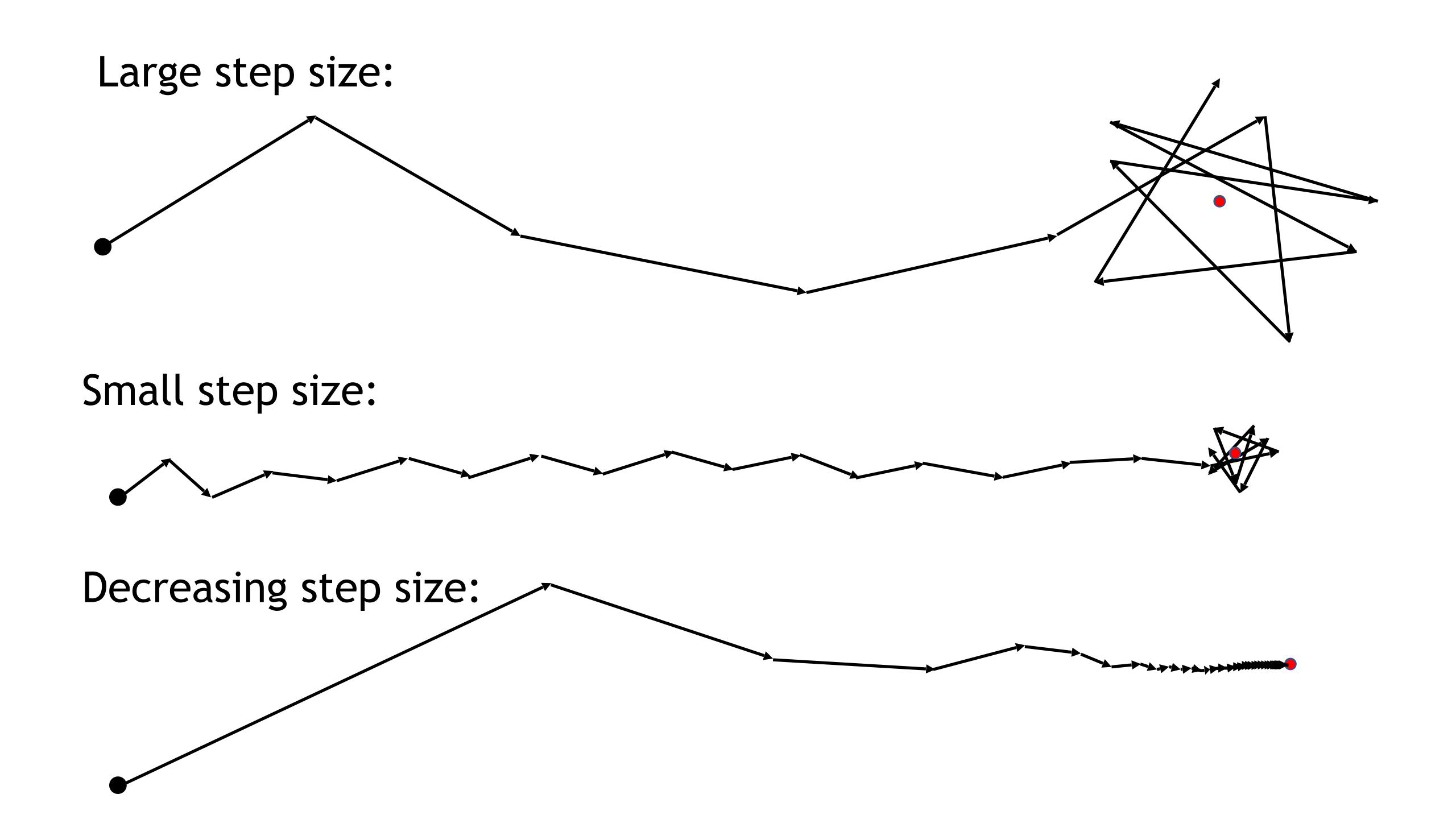
Large step size:

Large step size:



Small step size:





Proof:

Proof:
$$x^{(k+1)} = x^{(k)} - \alpha g_k$$

Proof:
$$x^{(k+1)} = x^{(k)} - \alpha g_k$$

$$\|x^{(k+1)} - x^*\|_2^2 = \|x^{(k)} - \alpha g_k - x^*\|_2^2$$

Proof:
$$x^{(k+1)} = x^{(k)} - \alpha g_k$$

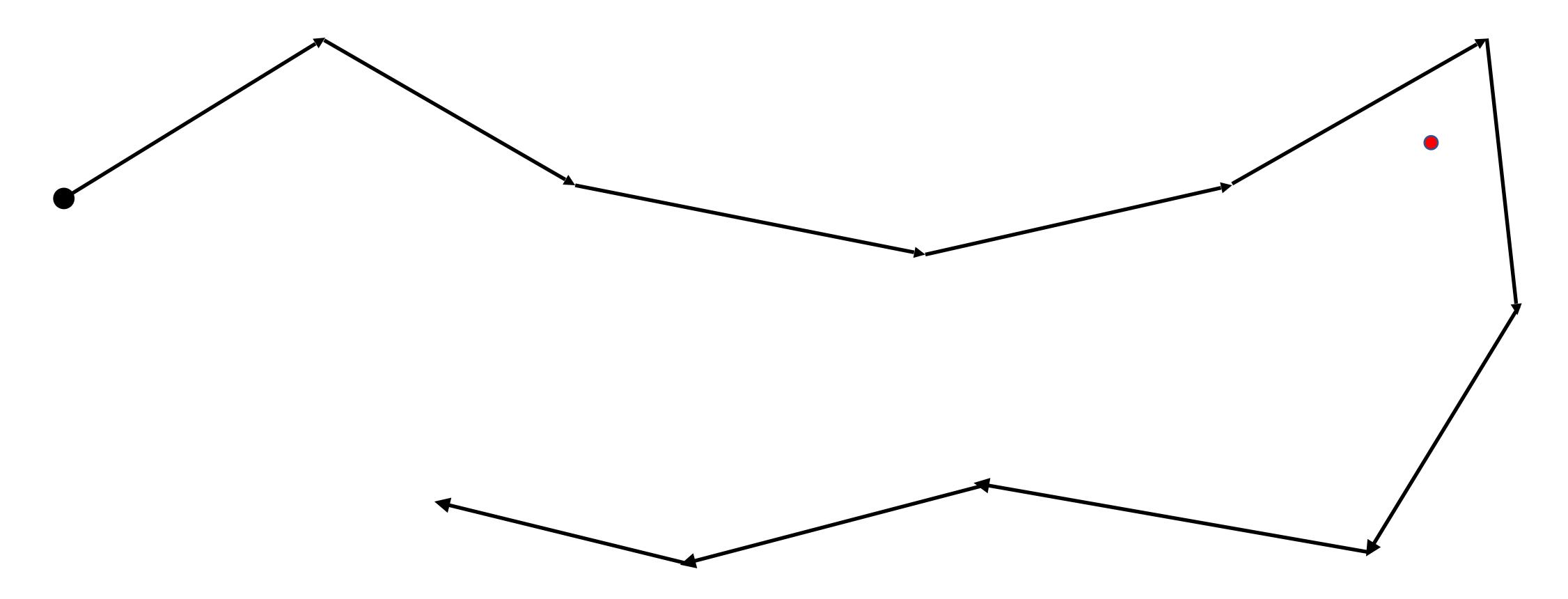
 $\|x^{(k+1)} - x^*\|_2^2 = \|x^{(k)} - \alpha g_k - x^*\|_2^2$
 $= \|x^{(k)} - x^*\|_2^2 - 2\alpha \langle g_k, x^{(k)} - x^* \rangle + \alpha^2 \|g_k\|_2^2$

$$\begin{array}{ll} \text{Proof:} & x^{(k+1)} = x^{(k)} - \alpha \, g_k \\ & \|x^{(k+1)} - x^\star\|_2^2 = \|x^{(k)} - \alpha g_k - x^\star\|_2^2 \\ & = \|x^{(k)} - x^\star\|_2^2 - 2\alpha \langle g_k, x^{(k)} - x^\star \rangle + \alpha^2 \|g_k\|_2^2 \end{array}$$
 By definition of g_k , $f(x^\star) \geq f\left(x^{(k)}\right) + \langle g_k, x^\star - x^{(k)} \rangle$

$$\begin{aligned} \text{Proof:} \quad & x^{(k+1)} = x^{(k)} - \alpha \, g_k \\ & \|x^{(k+1)} - x^\star\|_2^2 = \|x^{(k)} - \alpha g_k - x^\star\|_2^2 \\ & = \|x^{(k)} - x^\star\|_2^2 - 2\alpha \langle g_k, x^{(k)} - x^\star \rangle + \alpha^2 \|g_k\|_2^2 \\ & \text{By definition of } g_k, \ f\left(x^\star\right) \geq f\left(x^{(k)}\right) + \langle g_k, x^\star - x^{(k)}\rangle \\ & \|x^{(k+1)} - x^\star\|_2^2 \leq \|x^{(k)} - x^\star\|_2^2 - 2\alpha \left(f\left(x^{(k)}\right) - f(x^\star)\right) + \alpha^2 \|g_k\|_2^2 \end{aligned}$$

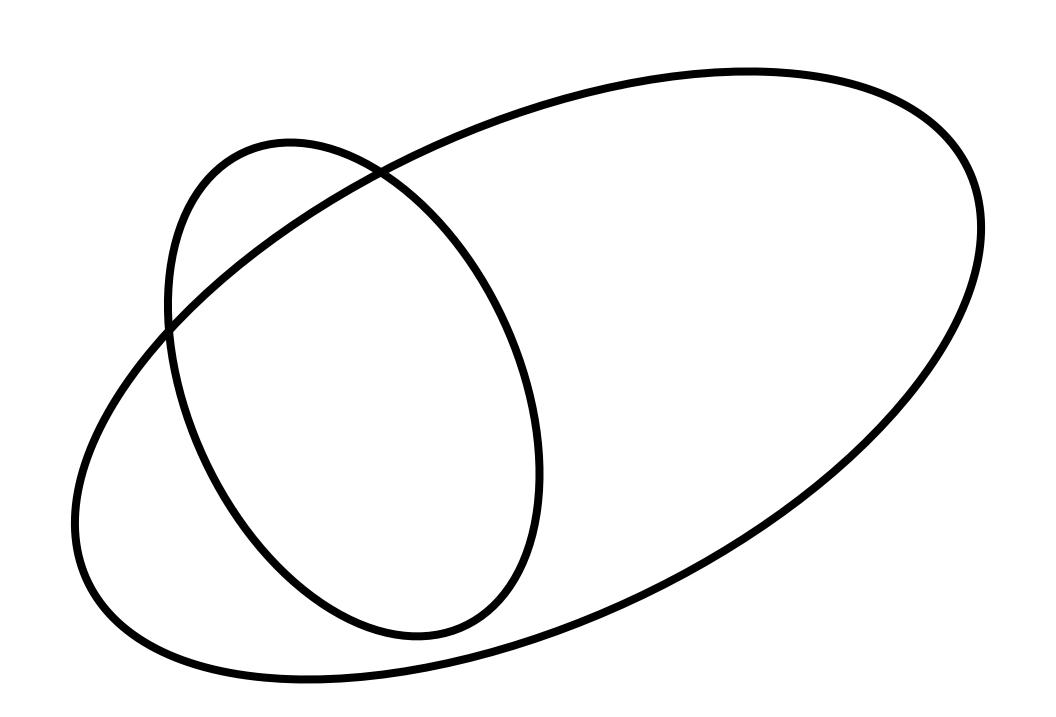
$$\begin{aligned} \text{Proof:} \quad & x^{(k+1)} = x^{(k)} - \alpha \, g_k \\ & \| x^{(k+1)} - x^\star \|_2^2 = \| x^{(k)} - \alpha g_k - x^\star \|_2^2 \\ & = \| x^{(k)} - x^\star \|_2^2 - 2\alpha \langle g_k, x^{(k)} - x^\star \rangle + \alpha^2 \| g_k \|_2^2 \\ & \text{By definition of } g_k, \ f\left(x^\star\right) \geq f\left(x^{(k)}\right) + \langle g_k, x^\star - x^{(k)}\rangle \\ & \| x^{(k+1)} - x^\star \|_2^2 \leq \| x^{(k)} - x^\star \|_2^2 - 2\alpha \left(f\left(x^{(k)}\right) - f(x^\star)\right) + \alpha^2 \| g_k \|_2^2 \\ & \therefore \ \| x^{(k+1)} - x^\star \|_2^2 < \| x^{(k)} - x^\star \|_2^2 \ \text{ when } \ \alpha < 2 \left(f\left(x^{(k)}\right) - f^\star\right) / \| g_k \|_2^2 \end{aligned}$$

Does the trajectory stay close, or can it move far away?

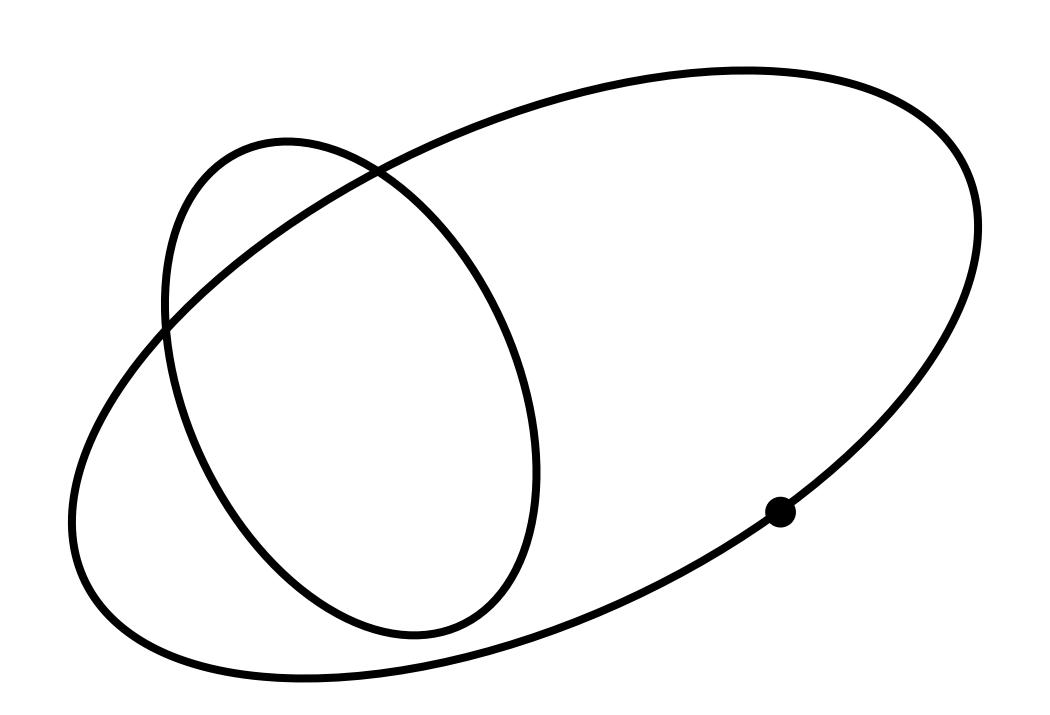


$$H(X,Y) = \max \left(\sup_{x \in X^*} \inf_{y \in S_k} ||x - y||_2, \sup_{y \in S_k} \inf_{x \in X^*} ||x - y||_2 \right)$$

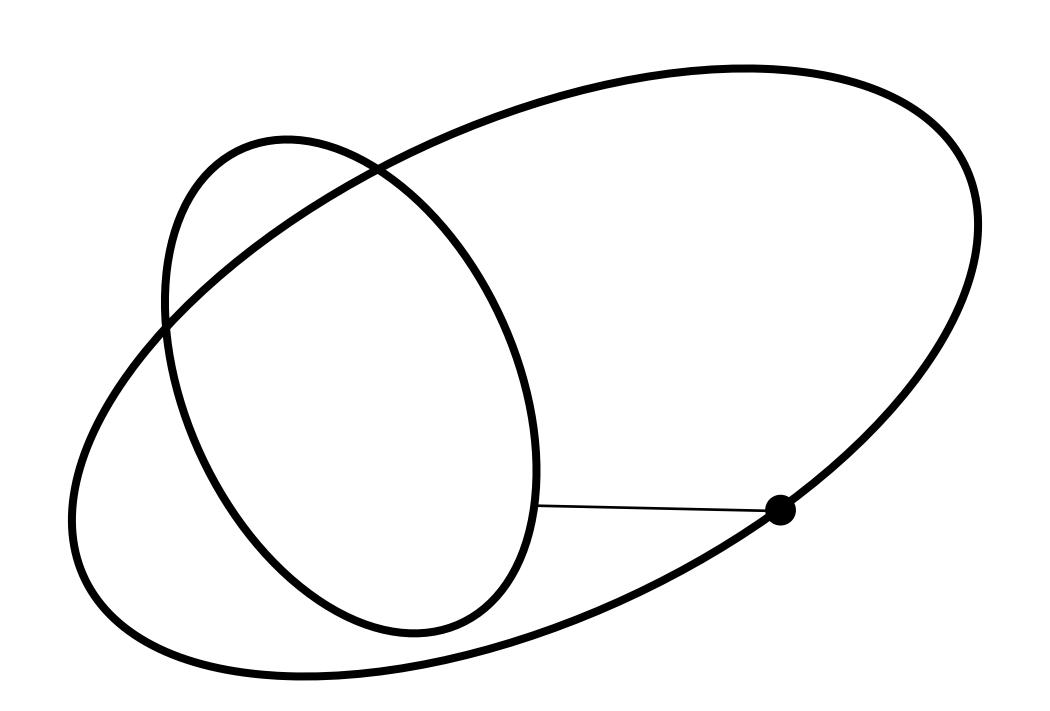
$$H(X,Y) = \max \left(\sup_{x \in X^*} \inf_{y \in S_k} ||x - y||_2, \sup_{y \in S_k} \inf_{x \in X^*} ||x - y||_2 \right)$$



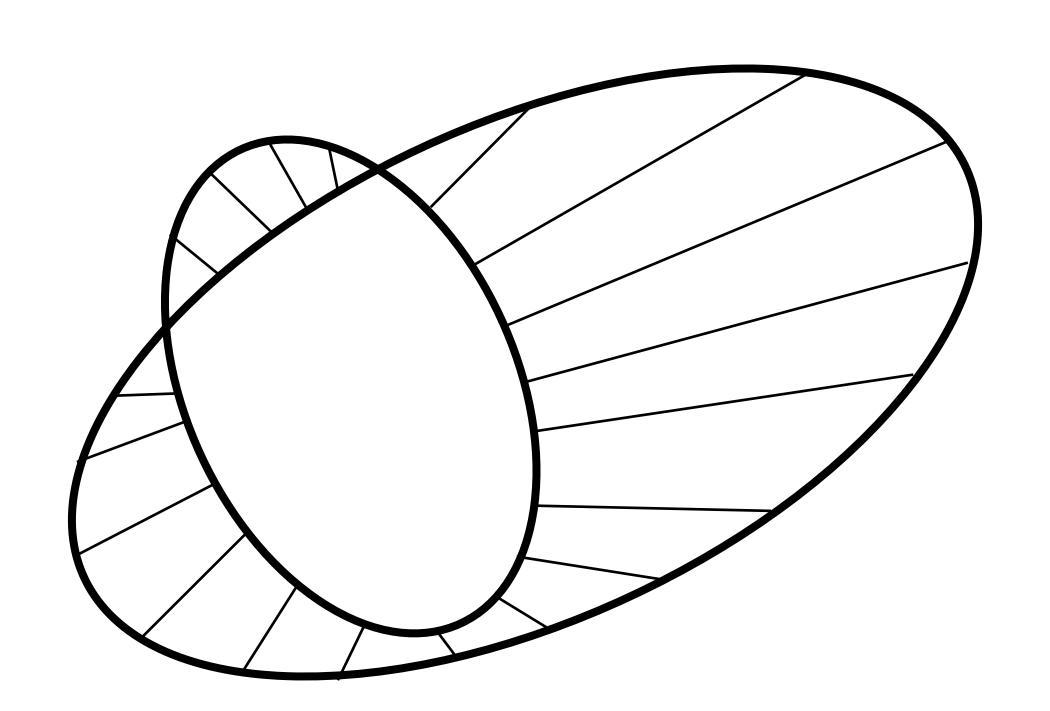
$$H(X,Y) = \max \left(\sup_{x \in X^*} \inf_{y \in S_k} ||x - y||_2, \sup_{y \in S_k} \inf_{x \in X^*} ||x - y||_2 \right)$$



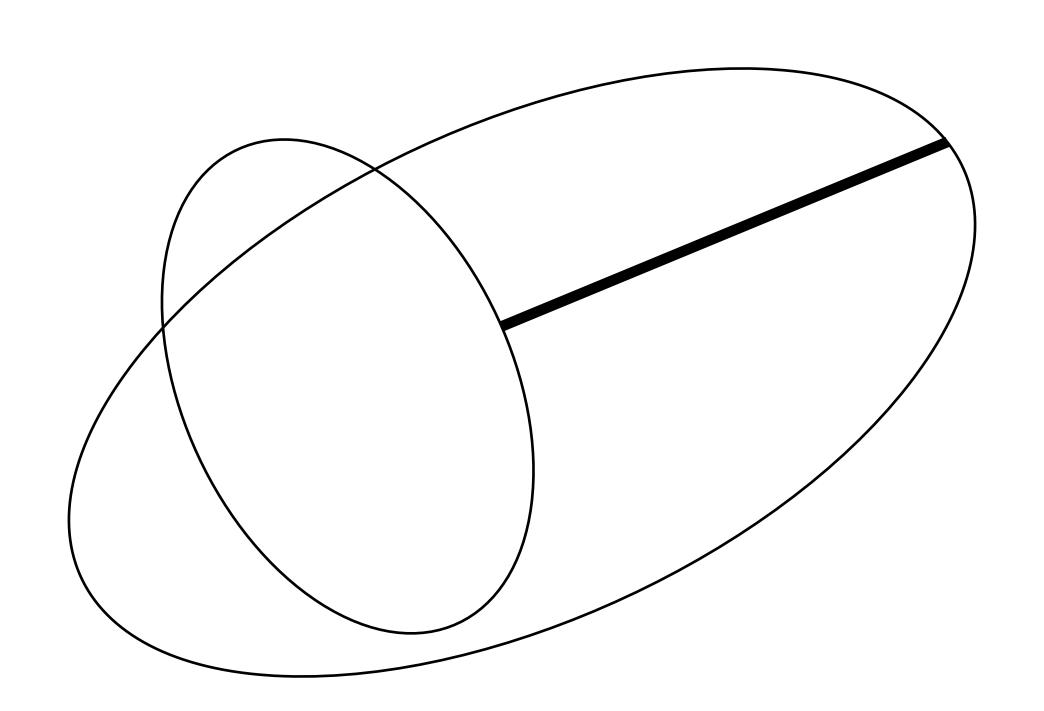
$$H(X,Y) = \max \left(\sup_{x \in X^*} \inf_{y \in S_k} ||x - y||_2, \sup_{y \in S_k} \inf_{x \in X^*} ||x - y||_2 \right)$$



$$H(X,Y) = \max \left(\sup_{x \in X^*} \inf_{y \in S_k} \|x - y\|_2, \sup_{y \in S_k} \inf_{x \in X^*} \|x - y\|_2 \right)$$



$$H(X,Y) = \max \left(\sup_{x \in X^*} \inf_{y \in S_k} ||x - y||_2, \sup_{y \in S_k} \inf_{x \in X^*} ||x - y||_2 \right)$$



Sublevel set: $S_k = \{x \in \mathbb{R}^N \mid f(x) \leq f(x^{(k)})\}$

Sublevel set:
$$S_k = \{x \in \mathbb{R}^N \mid f(x) \leq f(x^{(k)})\}$$

$$H_k = H\left(S_k, X^{\star}\right)$$

Sublevel set: $S_k = \{x \in \mathbb{R}^N \mid f(x) \leq f(x^{(k)})\}$

$$H_k = H\left(S_k, X^{\star}\right)$$

Theorem: the trajectory always stays within a bounded distance of \mathcal{H}_k .

Assume:

$$||g||_2 < G \quad \forall \quad g \in \partial f$$

$$\alpha < \left(f(x^{(k)}) - f^* \right) / G^2$$

Assume:

$$||g||_2 < G \quad \forall \quad g \in \partial f$$

$$\alpha < \left(f(x^{(k)}) - f^* \right) / G^2$$

Theorem:

$$\operatorname{dist}\left(x^{(m)}, X^{\star}\right) < H_k + 2\left(f\left(x^{(k)}\right) - f^{\star}\right)/G$$

Assume: $||g||_2 < G \quad \forall \quad g \in \partial f$

$$\alpha < \left(f(x^{(k)}) - f^* \right) / G^2$$

Theorem: dist $\left(x^{(m)}, X^{\star}\right) < H_k + 2\left(f\left(x^{(k)}\right) - f^{\star}\right)/G$

Proof (by induction):

Assume: $||g||_2 < G \quad \forall \quad g \in \partial f$

$$\alpha < \left(f(x^{(k)}) - f^* \right) / G^2$$

Theorem: dist $\left(x^{(m)}, X^{\star}\right) < H_k + 2\left(f\left(x^{(k)}\right) - f^{\star}\right)/G$

Proof (by induction):

Base step: since $m=k,\ x^{(k)}\in S_k$. Therefore the theorem is true by the previous theorem.

Case 1: $x^{(k)} \in S_k$

Case 1:
$$x^{(k)} \in S_k$$
 Since $\alpha < \left(f(x^{(k)}) - f^\star\right)/G^2$,
$$\|x^{(m+1)} - x^{(m)}\|_2 = \alpha \|g^{(m)}\|_2$$

Case 1:
$$x^{(k)} \in S_k$$
 Since $\alpha < \left(f(x^{(k)}) - f^\star\right)/G^2$,
$$\|x^{(m+1)} - x^{(m)}\|_2 = \alpha \|g^{(m)}\|_2 \le aG$$

Case 1:
$$x^{(k)} \in S_k$$
 Since $\alpha < \left(f(x^{(k)}) - f^\star\right)/G^2$,
$$\|x^{(m+1)} - x^{(m)}\|_2 = \alpha \|g^{(m)}\|_2 \le aG < 2\left(f(x^{(k)}) - f^\star\right)/G$$

Case 1:
$$x^{(k)} \in S_k$$

Since
$$\alpha < \left(f(x^{(k)}) - f^* \right) / G^2$$
,

$$||x^{(m+1)} - x^{(m)}||_2 = \alpha ||g^{(m)}||_2 \le aG < 2\left(f(x^{(k)}) - f^*\right)/G$$

By the triangle inequality,

$$\operatorname{dist}\left(x^{(m+1)}, X^{\star}\right) \leq \operatorname{dist}\left(x^{(m+1)}, x^{(m)}\right) + \operatorname{dist}\left(x^{(m)}, X^{\star}\right)$$

Case 1:
$$x^{(k)} \in S_k$$

Since
$$\alpha < \left(f(x^{(k)}) - f^* \right) / G^2$$
,

$$||x^{(m+1)} - x^{(m)}||_2 = \alpha ||g^{(m)}||_2 \le aG < 2\left(f(x^{(k)}) - f^*\right)/G$$

By the triangle inequality,

$$\operatorname{dist}\left(x^{(m+1)}, X^{\star}\right) \leq \operatorname{dist}\left(x^{(m+1)}, x^{(m)}\right) + \operatorname{dist}\left(x^{(m)}, X^{\star}\right)$$
$$< 2\left(f(x^{(k)}) - f^{\star}\right) / G + H_{k}$$

Case 1:
$$x^{(k)} \in S_k$$

Since
$$\alpha < \left(f(x^{(k)}) - f^* \right) / G^2$$
,

$$||x^{(m+1)} - x^{(m)}||_2 = \alpha ||g^{(m)}||_2 \le aG < 2\left(f(x^{(k)}) - f^*\right)/G$$

By the triangle inequality,

$$\operatorname{dist}\left(x^{(m+1)}, X^{\star}\right) \leq \operatorname{dist}\left(x^{(m+1)}, x^{(m)}\right) + \operatorname{dist}\left(x^{(m)}, X^{\star}\right)$$
$$< 2\left(f(x^{(k)}) - f^{\star}\right) / G + H_{k}$$

Case 2: $x^{(k)} \notin S_k$.

Case 1:
$$x^{(k)} \in S_k$$
 Since $\alpha < \left(f(x^{(k)}) - f^\star \right) / G^2$,

$$||x^{(m+1)} - x^{(m)}||_2 = \alpha ||g^{(m)}||_2 \le aG < 2\left(f(x^{(k)}) - f^*\right)/G$$

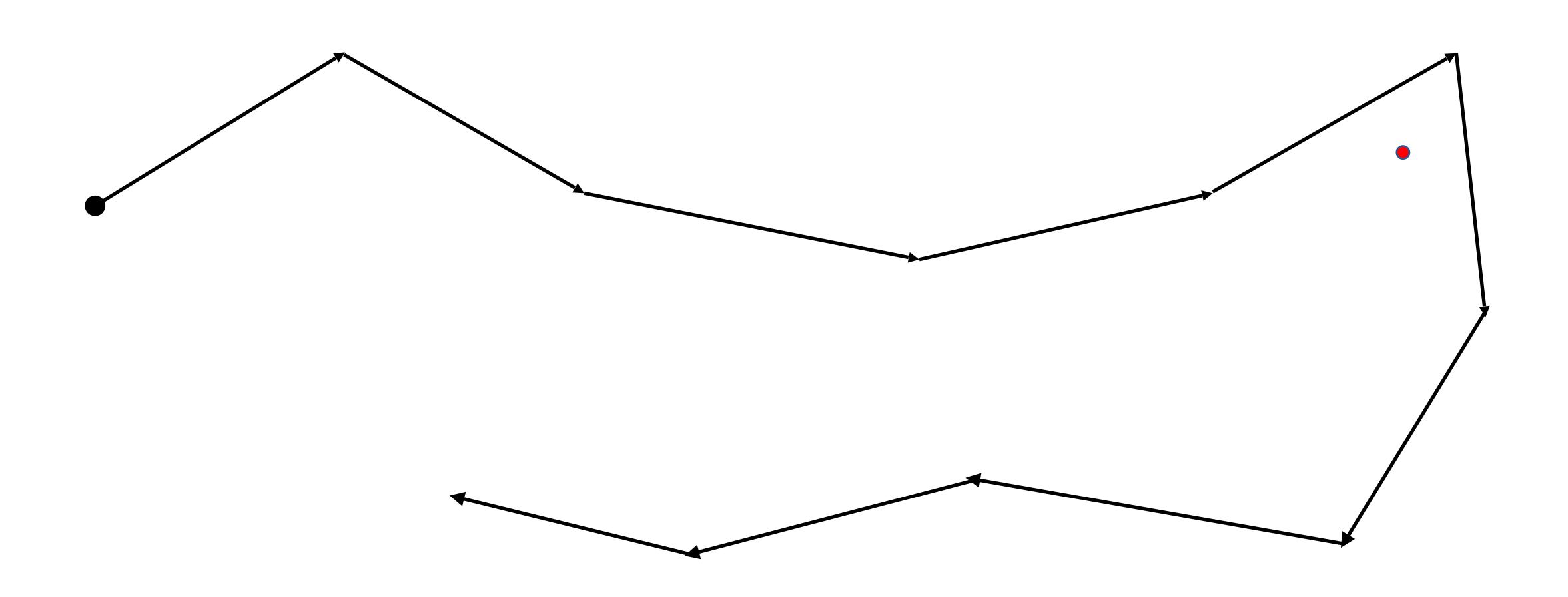
By the triangle inequality,

$$\operatorname{dist}\left(x^{(m+1)}, X^{\star}\right) \leq \operatorname{dist}\left(x^{(m+1)}, x^{(m)}\right) + \operatorname{dist}\left(x^{(m)}, X^{\star}\right)$$
$$< 2\left(f(x^{(k)}) - f^{\star}\right) / G + H_{k}$$

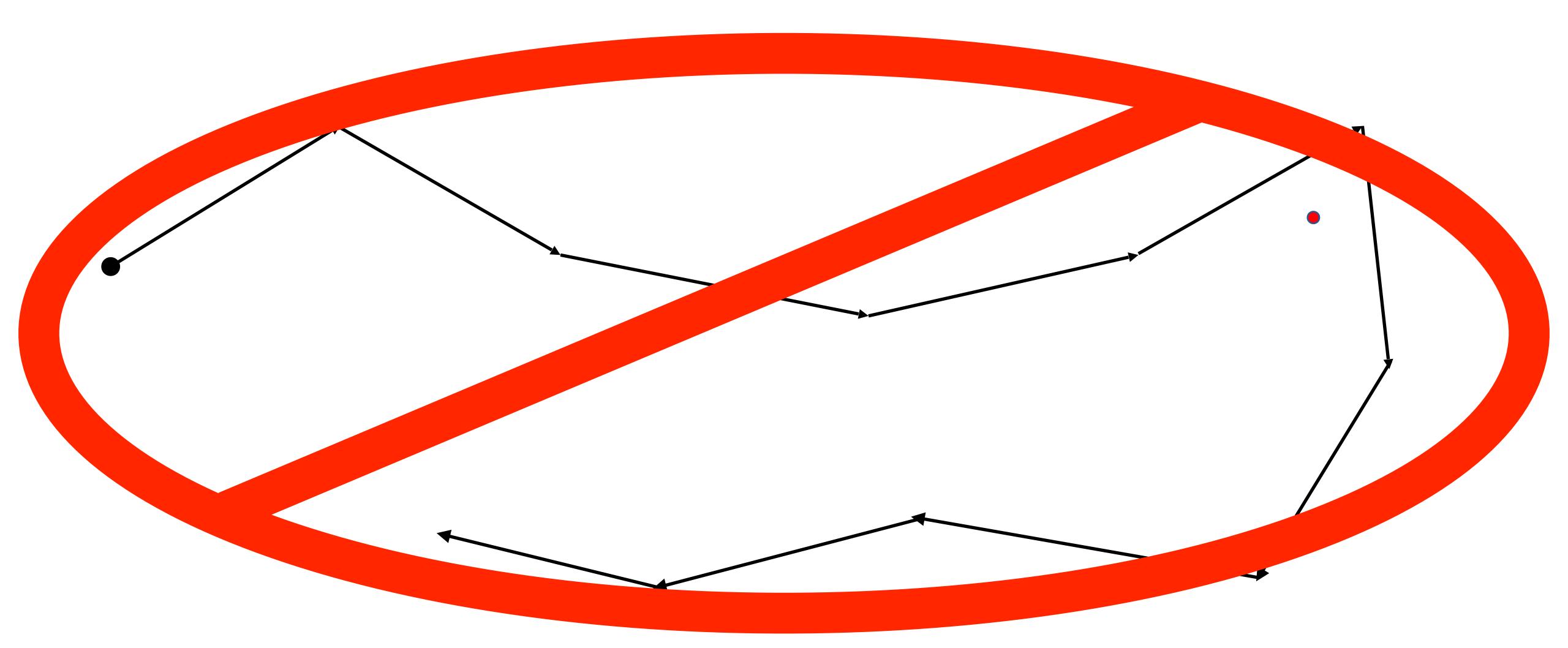
Case 2: $x^{(k)} \notin S_k$.

Then we move closer to the optimal set (as before).

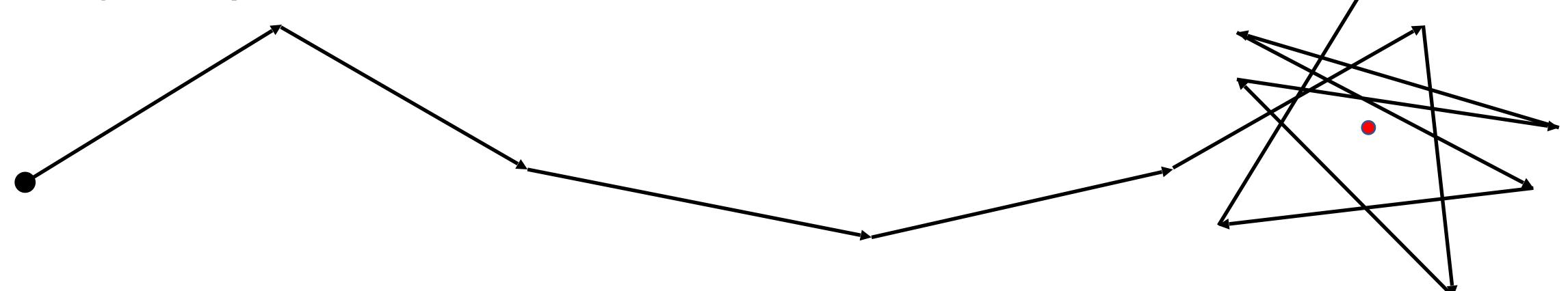
Theorem: the trajectory always stays within a bounded distance of \mathcal{H}_k .



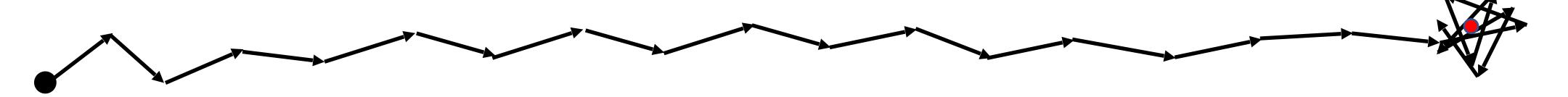
Theorem: the trajectory always stays within a bounded distance of \mathcal{H}_k .

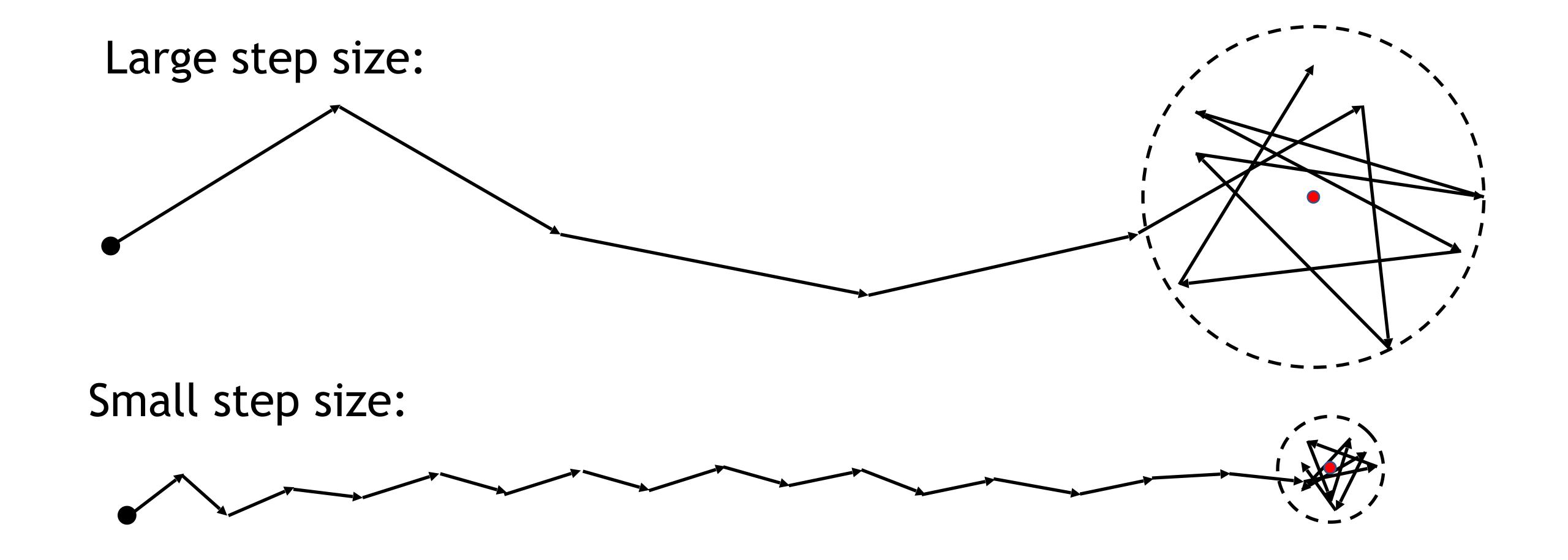


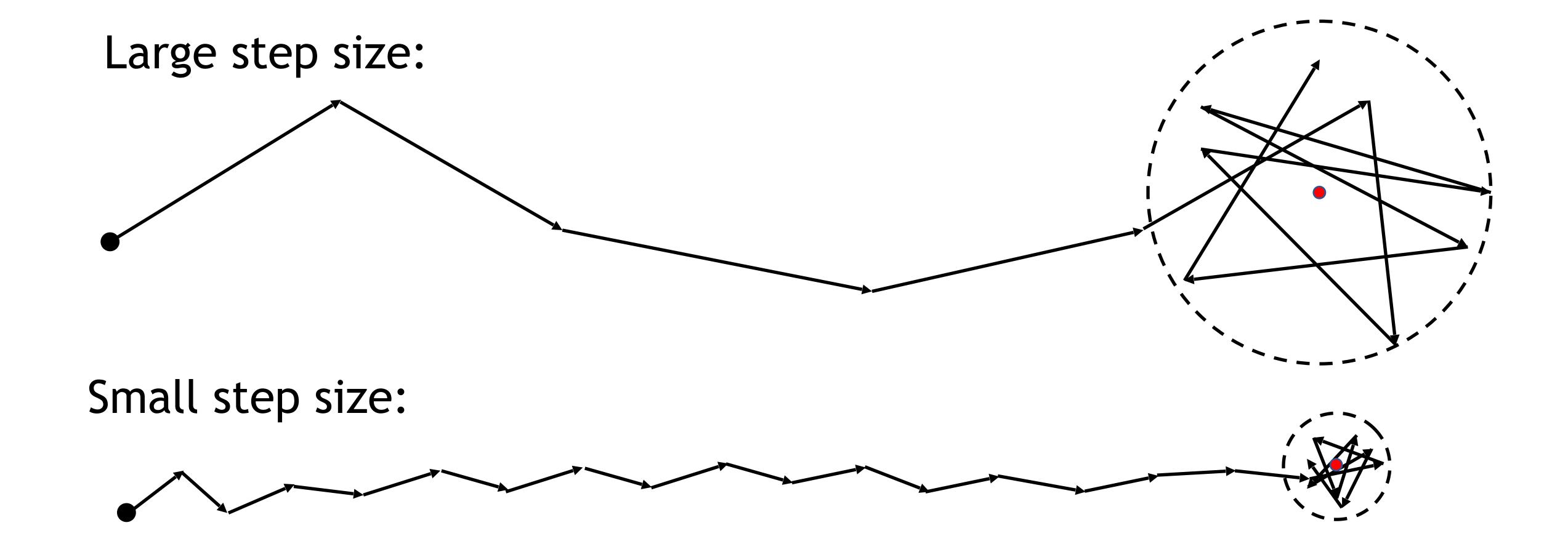
Large step size:



Small step size:

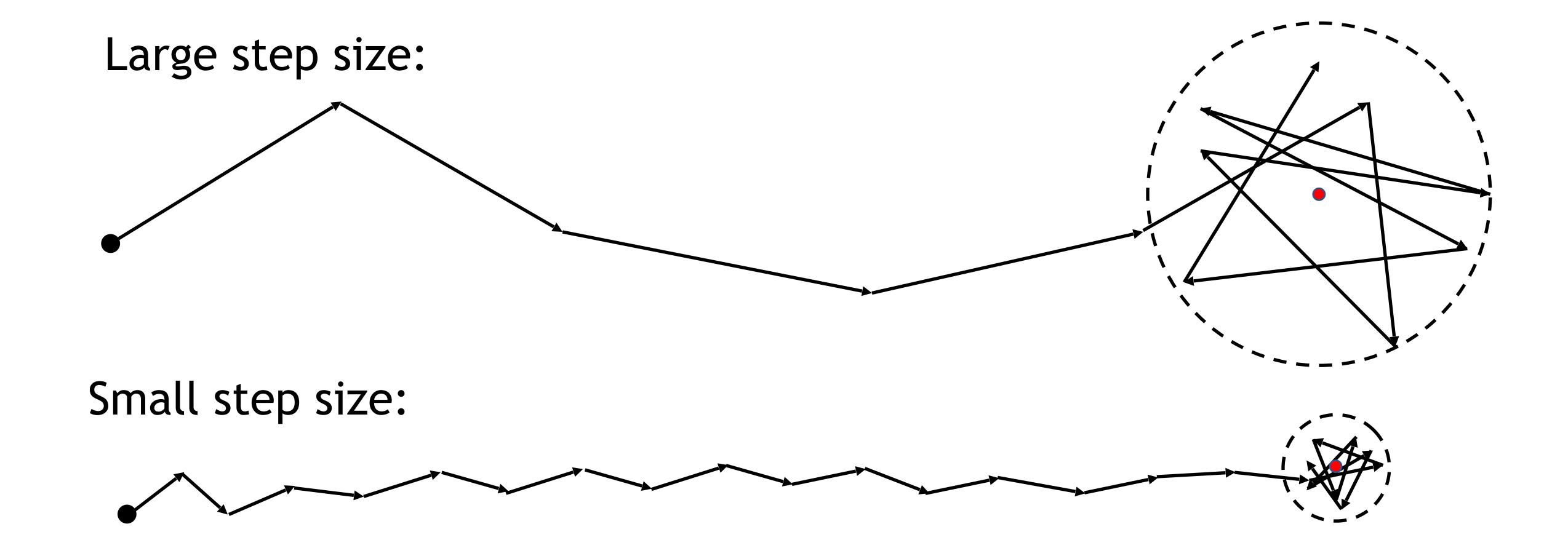




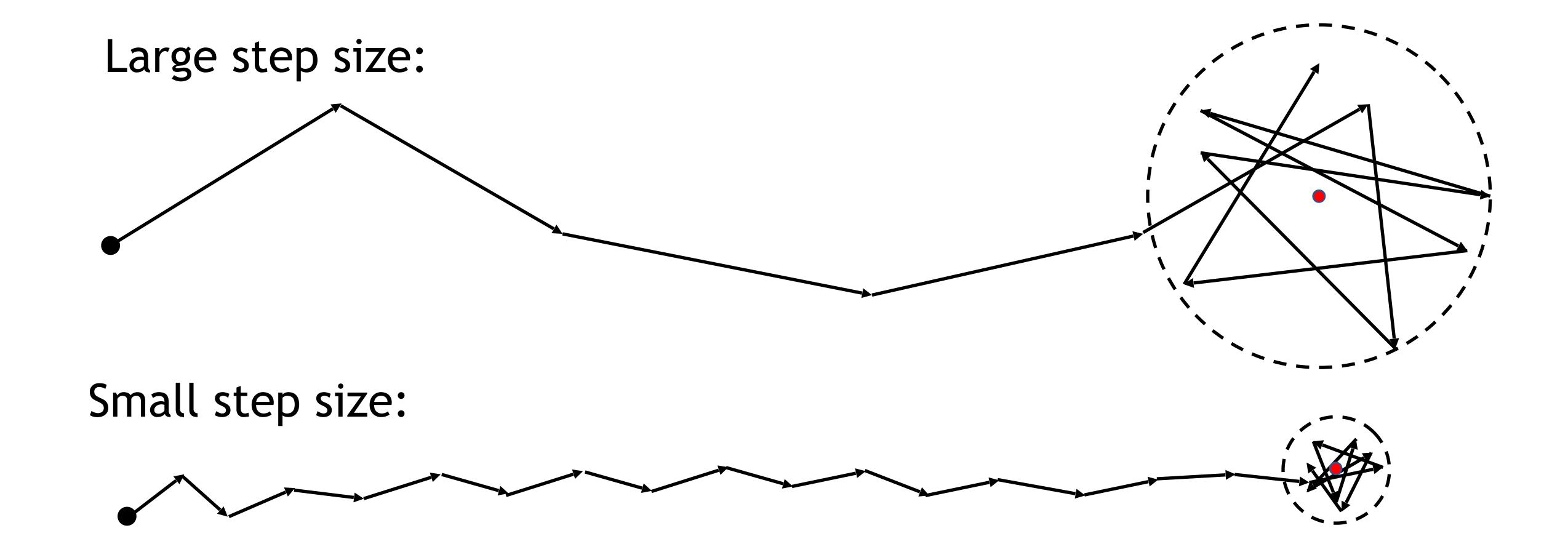


Two possibilities

- 1) We're converging to a solution
- 2) We're bouncing around the solution



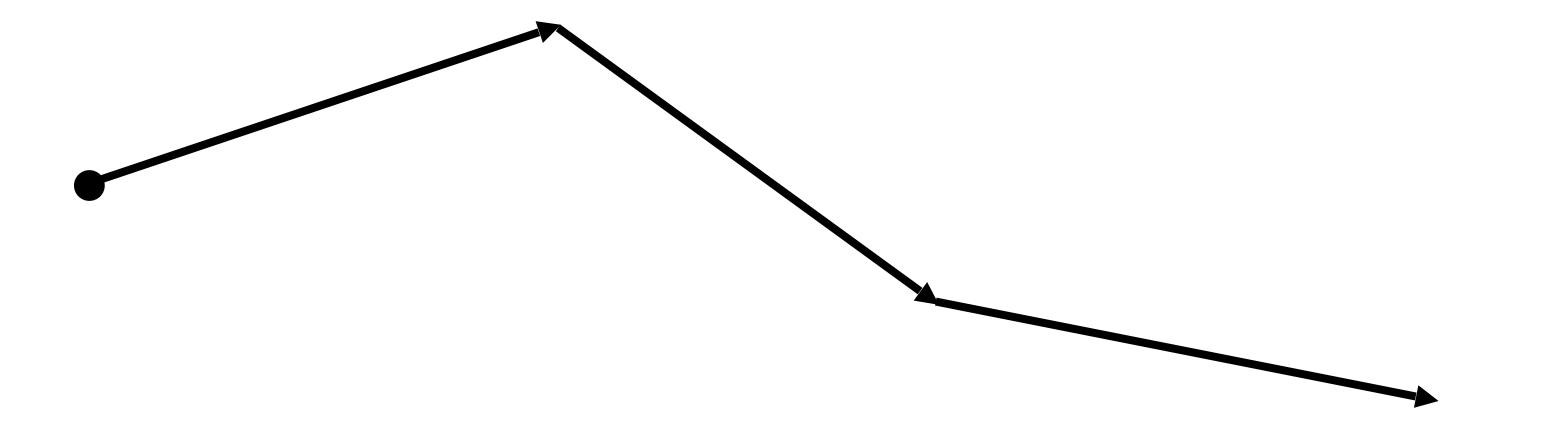
Can we detect when the trajectory is bouncing around the optimal set?



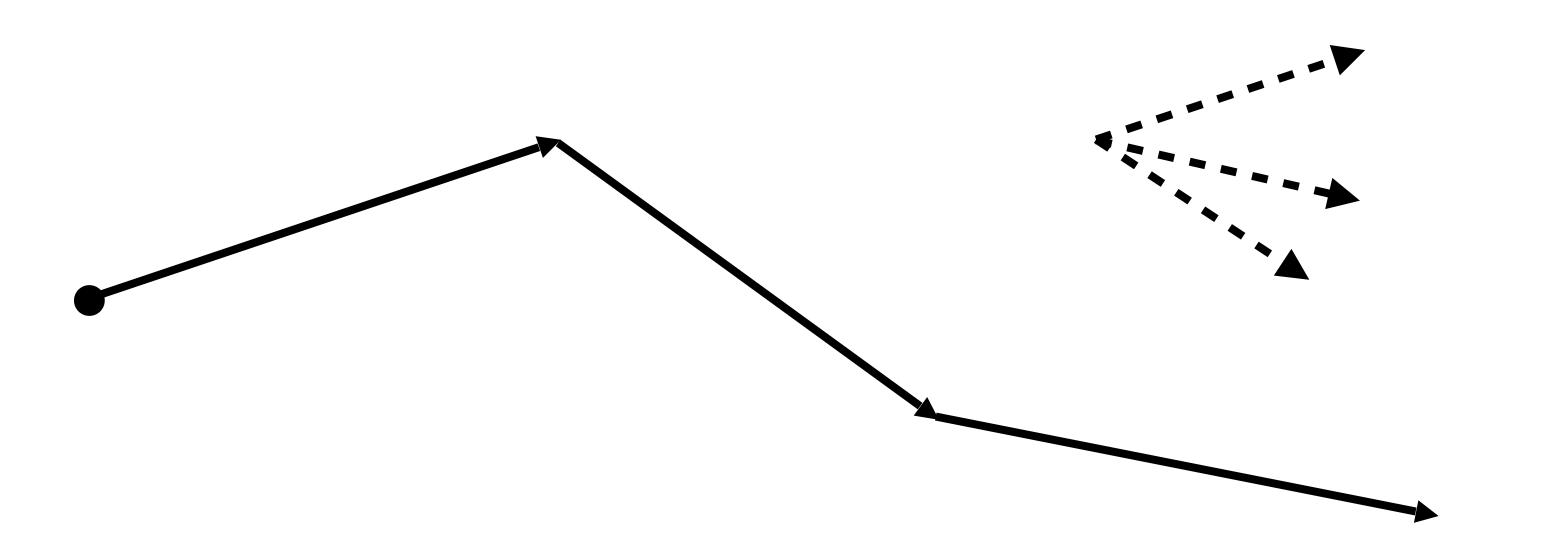
Can we detect when the trajectory is bouncing around the optimal set?

If so, can we make use of this knowledge?

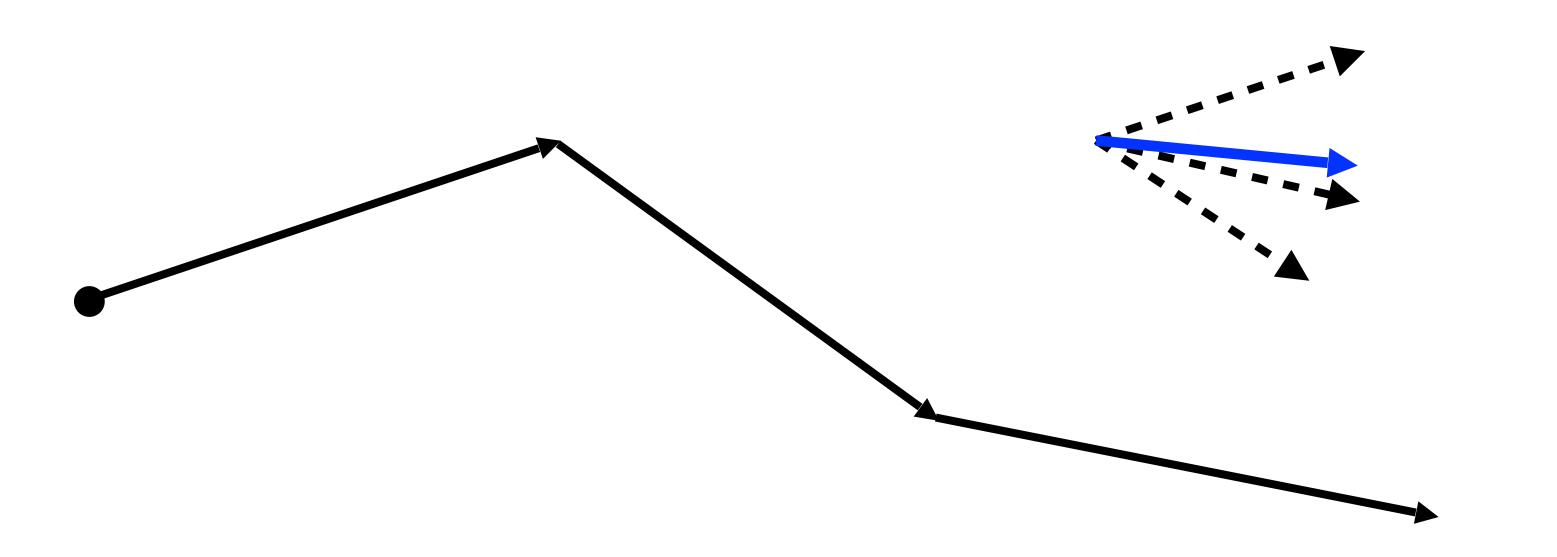
$$u_k = \frac{-g_k}{\|g_k\|_2} \qquad m = \left\|\frac{\sum u_k}{n}\right\|_2$$



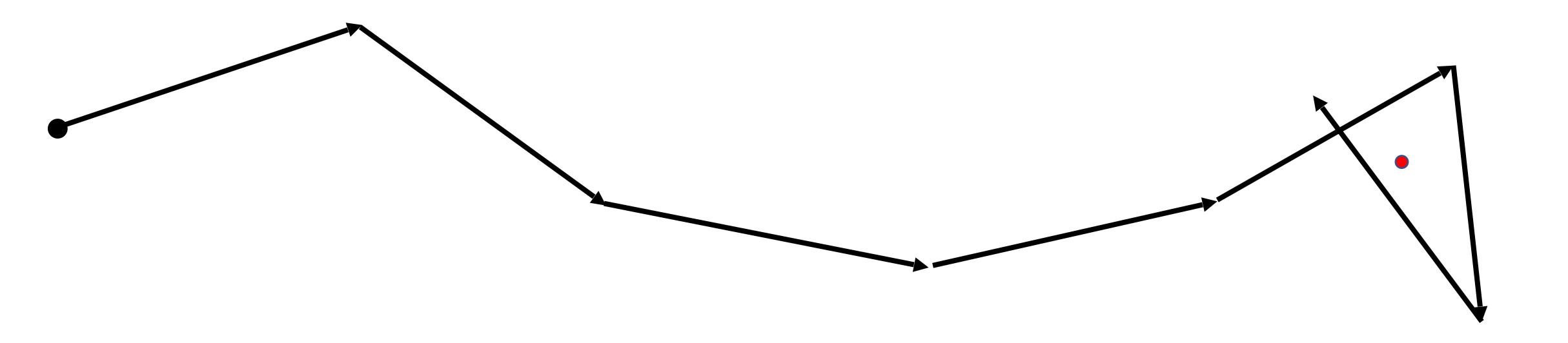
$$u_k = \frac{-g_k}{\|g_k\|_2} \qquad m = \left\|\frac{\sum u_k}{n}\right\|_2$$



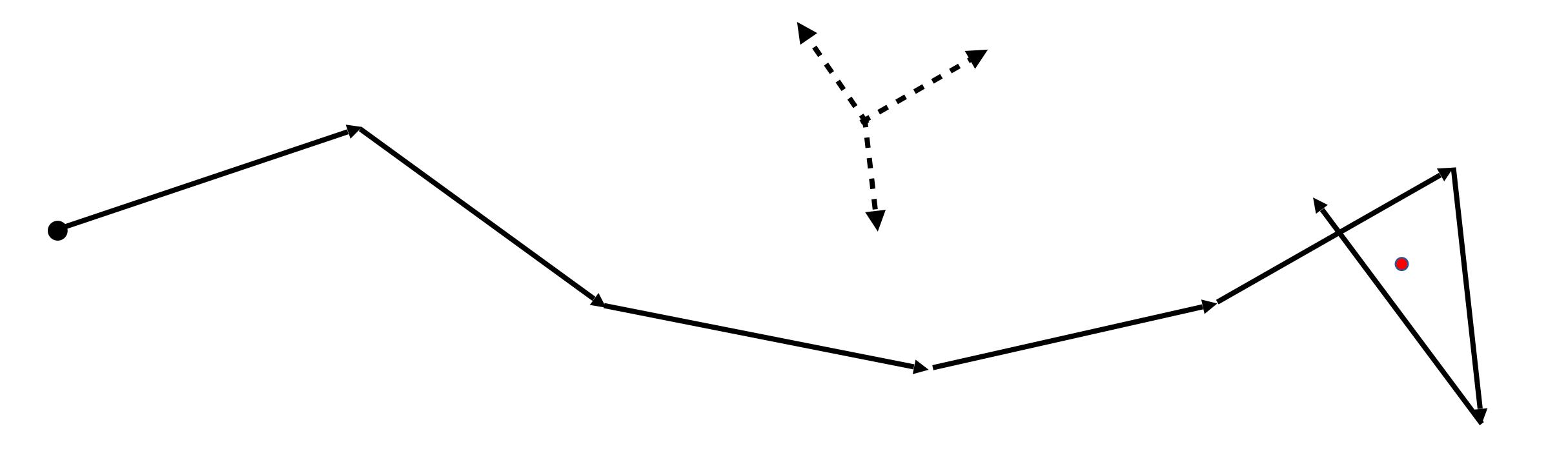
$$u_k = \frac{-g_k}{\|g_k\|_2} \qquad m = \left\|\frac{\sum u_k}{n}\right\|_2$$



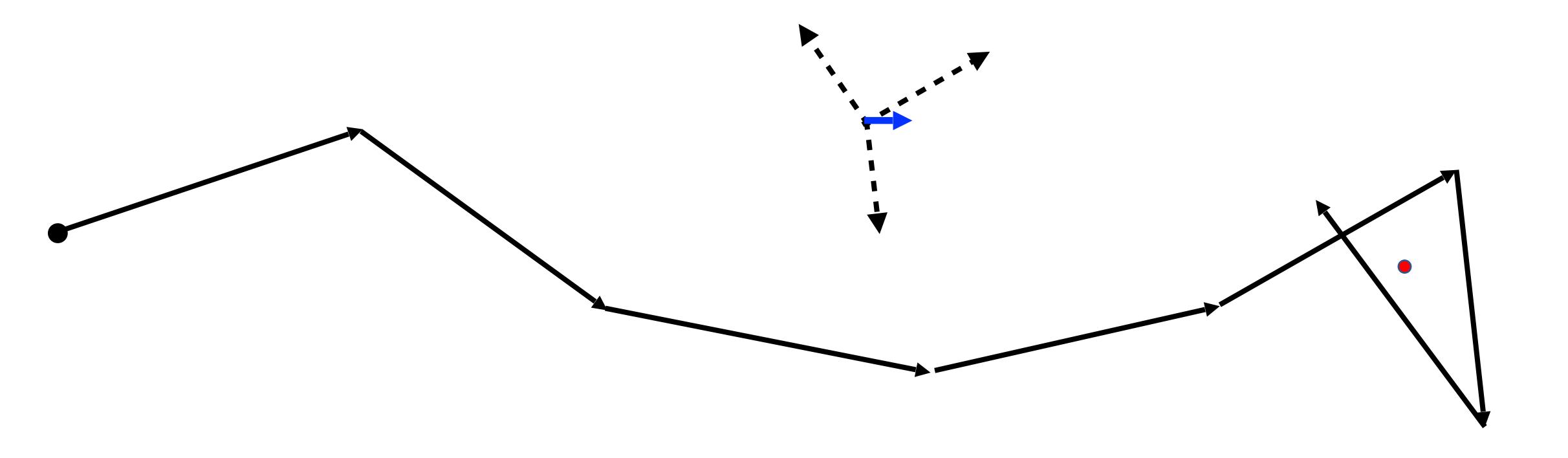
$$u_k = \frac{-g_k}{\|g_k\|_2} \qquad m = \left\|\frac{\sum u_k}{n}\right\|_2$$



$$u_k = \frac{-g_k}{\|g_k\|_2} \qquad m = \left\|\frac{\sum u_k}{n}\right\|_2$$



$$u_k = \frac{-g_k}{\|g_k\|_2} \qquad m = \left\|\frac{\sum u_k}{n}\right\|_2$$



We know that if the step size is small enough, then the next trajectory point is closer to the solution than the current point.

We know that if the step size is small enough, then the next trajectory point is closer to the solution than the current point.

We know that the trajectory will bounce around the solution if the step size is too large.

We know that if the step size is small enough, then the next trajectory point is closer to the solution than the current point.

We know that the trajectory will bounce around the solution if the step size is too large.

We know that we can detect when the trajectory is bouncing around the solution.

Adaptive Bounding Method

Adaptive Bounding Method

If the bouncing metric is below a threshold

Reduce the step size: $\alpha_k = r \alpha_{k-1}$ 0 < r < 1

Adaptive Bounding Method

If the bouncing metric is below a threshold

Reduce the step size: $\alpha_k = r \alpha_{k-1}$ 0 < r < 1

Perform a subgradient descent update

$$x^{(k+1)} = x^{(k)} - \alpha_k g(x^{(k)})$$

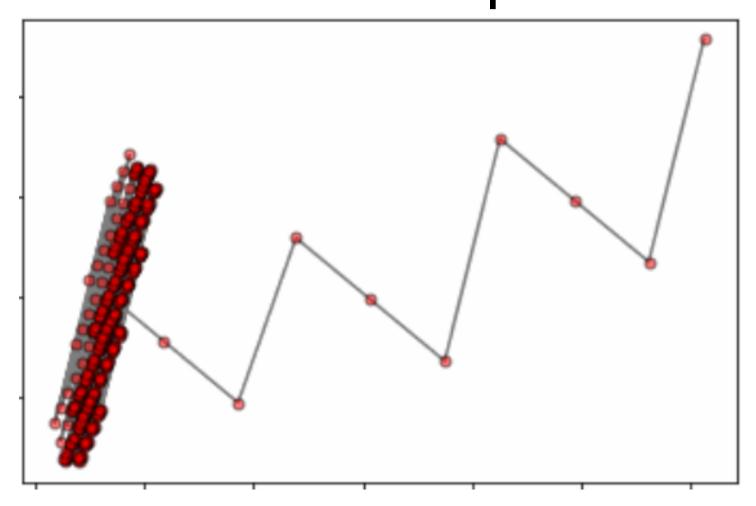
$$f(x) = \max(A_i x + b_i) \qquad A \in \mathbb{R}^{500 \times 2} \quad b \in \mathbb{R}^{500}$$

$$f(x) = \max(A_i x + b_i)$$

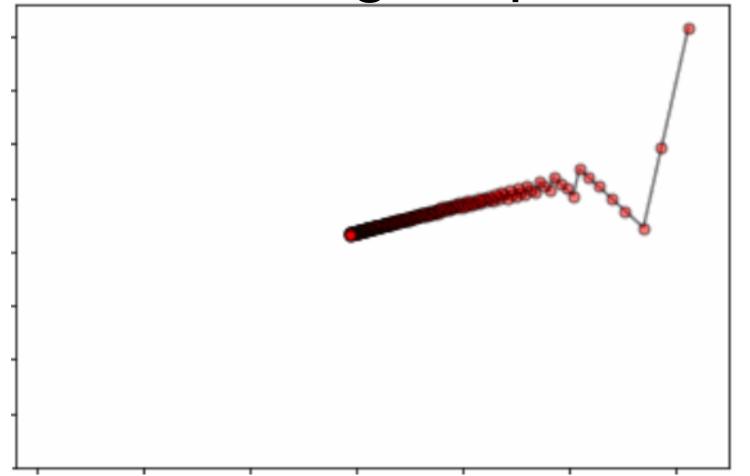
$$A \in \mathbb{R}^{500 \times 2}$$

$$b \in \mathbb{R}^{500}$$

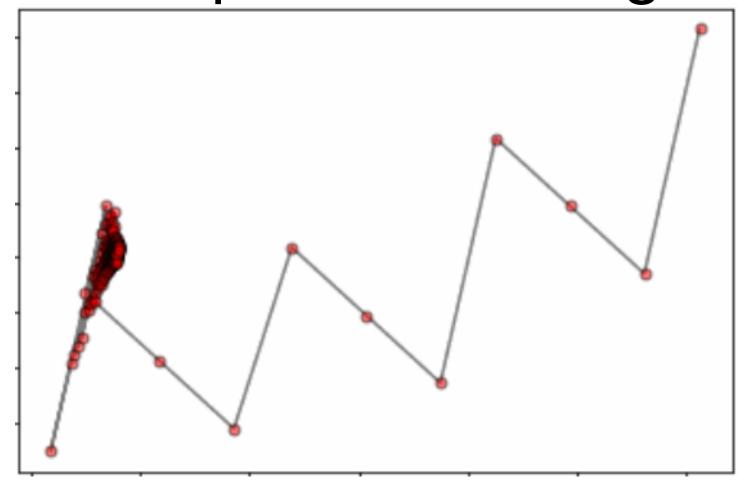
Constant Step Size



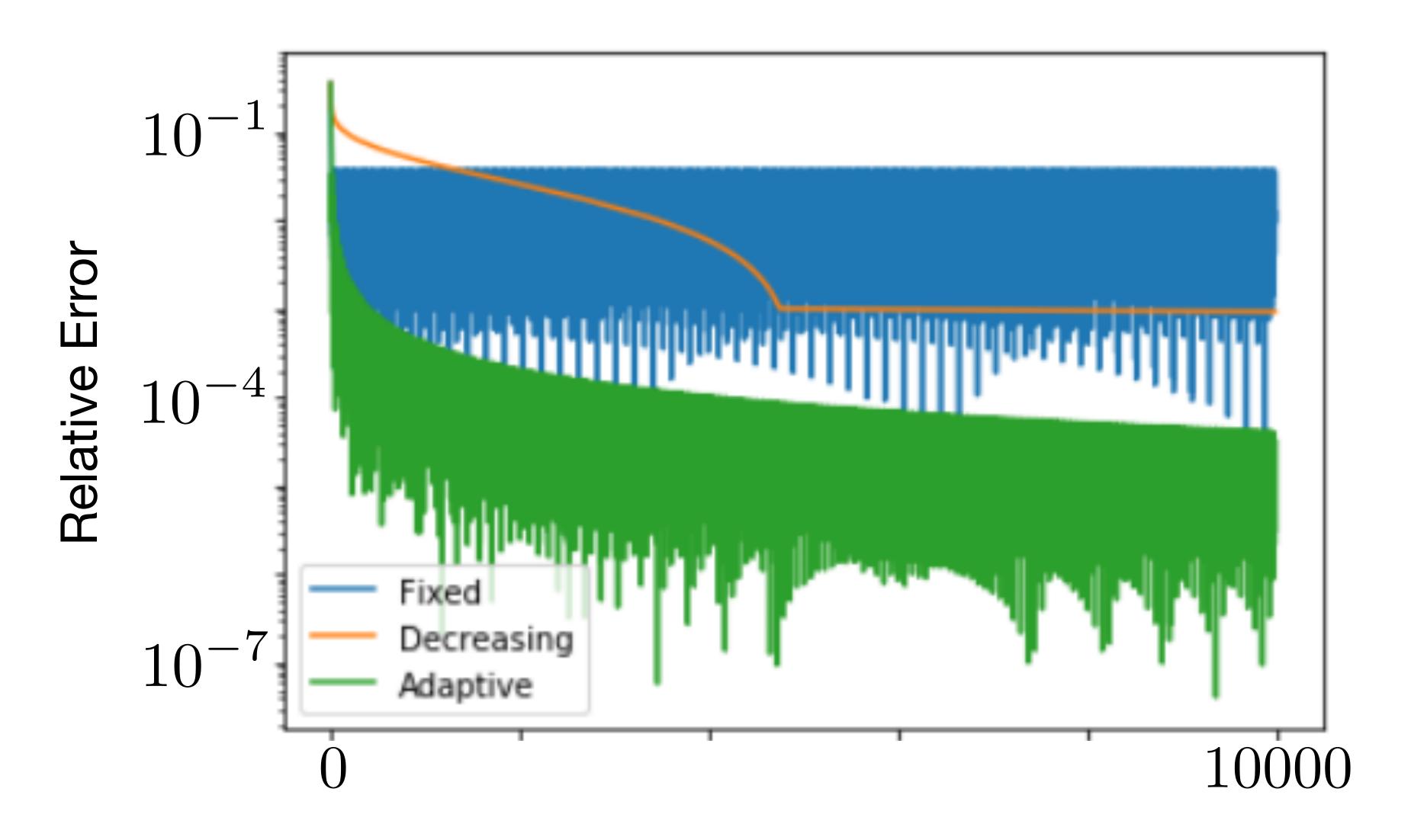
Decreasing Step Size



Adaptive Bounding



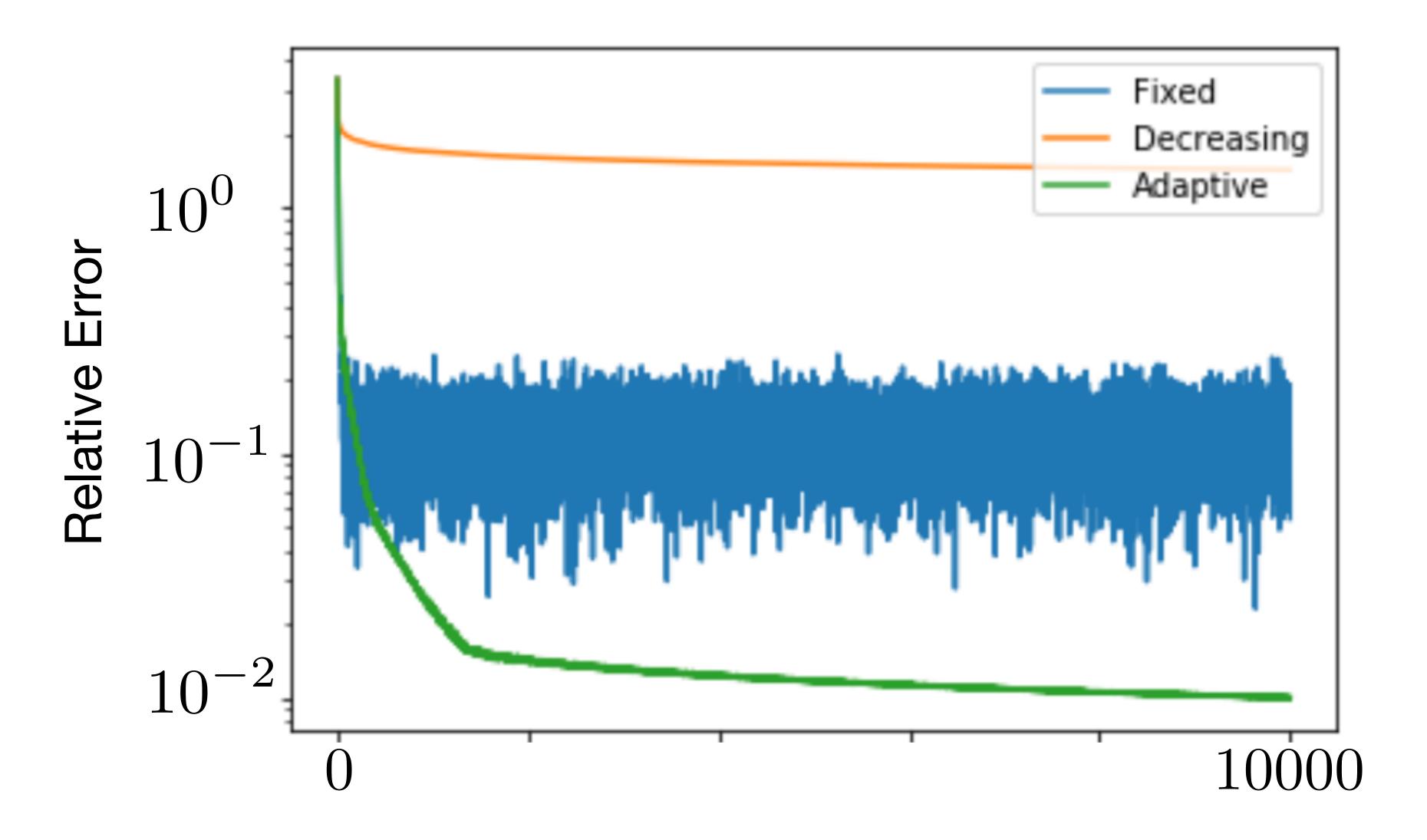
$$f(x) = \max(A_i x + b_i) \qquad A \in \mathbb{R}^{500 \times 2} \quad b \in \mathbb{R}^{500}$$



$$f(x) = \max(A_i x + b_i) \qquad A \in \mathbb{R}^{200 \times 10} \quad b \in \mathbb{R}^{200}$$

Minimize the maximum of a set of affine functions (larger problem)

$$f(x) = \max(A_i x + b_i) \qquad A \in \mathbb{R}^{200 \times 10} \quad b \in \mathbb{R}^{200}$$

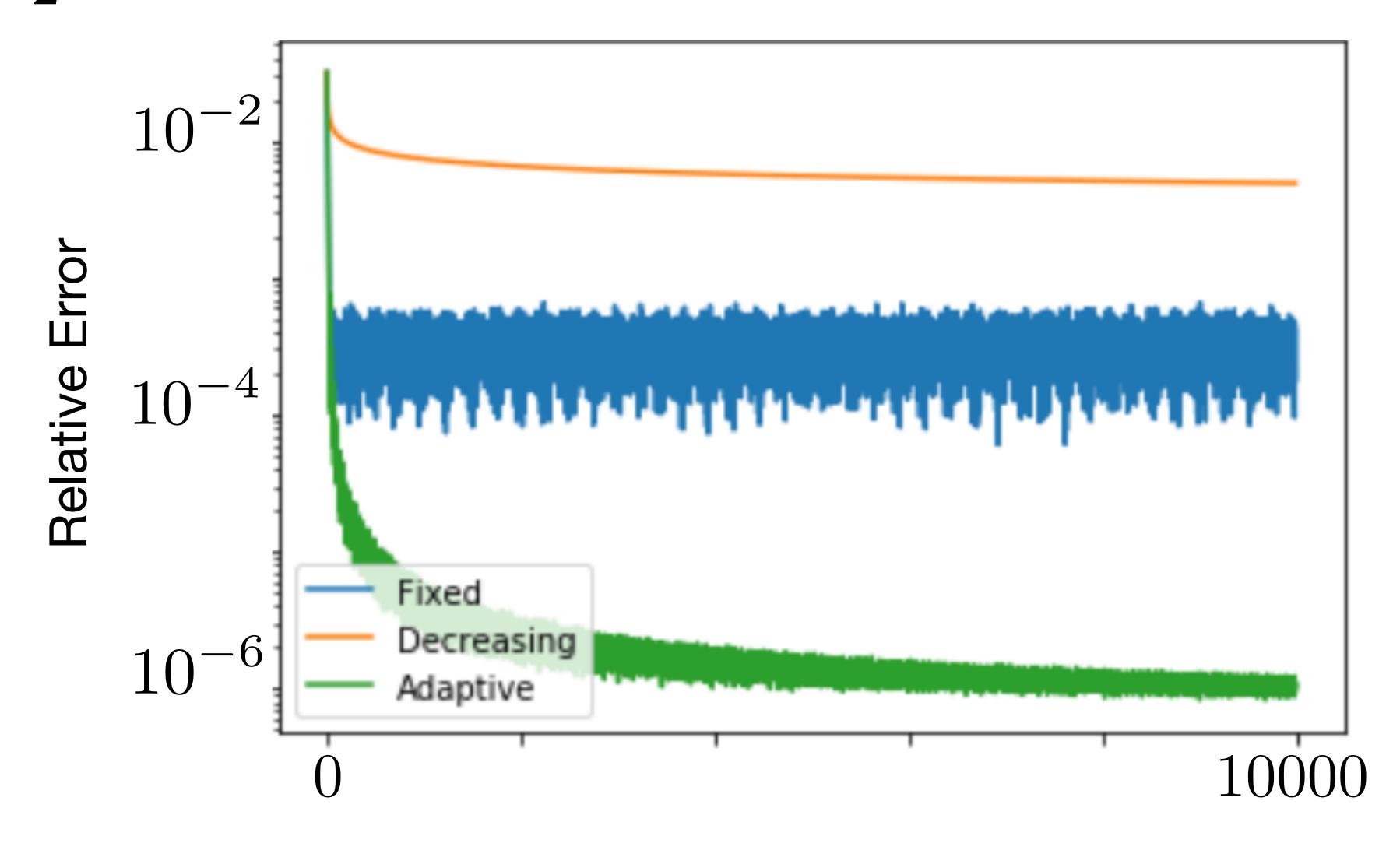


Lasso Problem

$$f(x) = \frac{1}{2} ||Ax - b||_2^2 + \lambda ||x||_1 \qquad A \in \mathbb{R}^{200 \times 20} \quad b \in \mathbb{R}^{200}$$

Lasso Problem

$$f(x) = \frac{1}{2} ||Ax - b||_2^2 + \lambda ||x||_1 \qquad A \in \mathbb{R}^{200 \times 20} \quad b \in \mathbb{R}^{200}$$

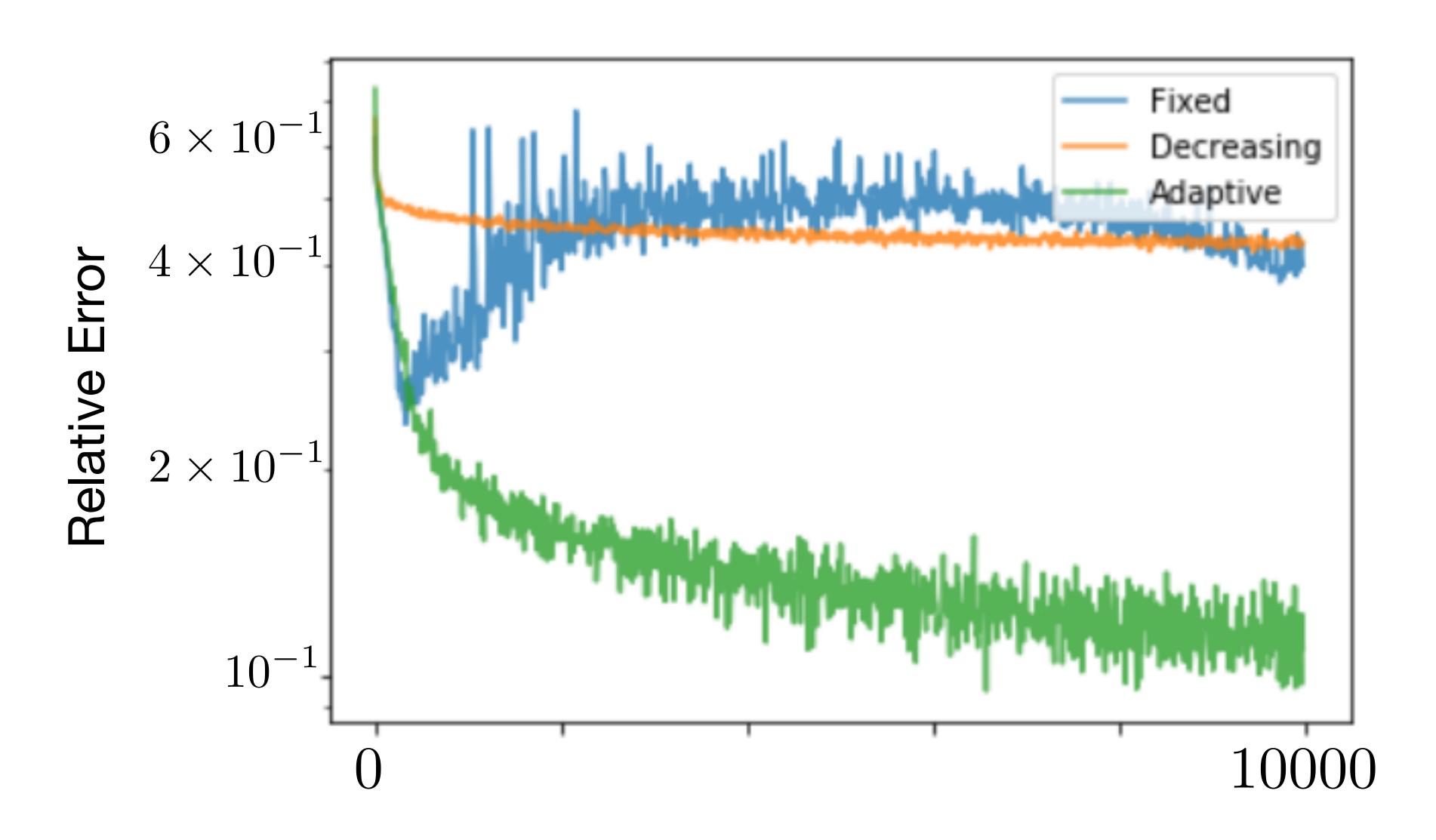


Neural Network

RESNET on CIFAR-10 dataset with L1 norm regularization

Neural Network

RESNET on CIFAR-10 dataset with L1 norm regularization



Conclusion

We have developed an apaptive method for determining the step size with subgradient descent.

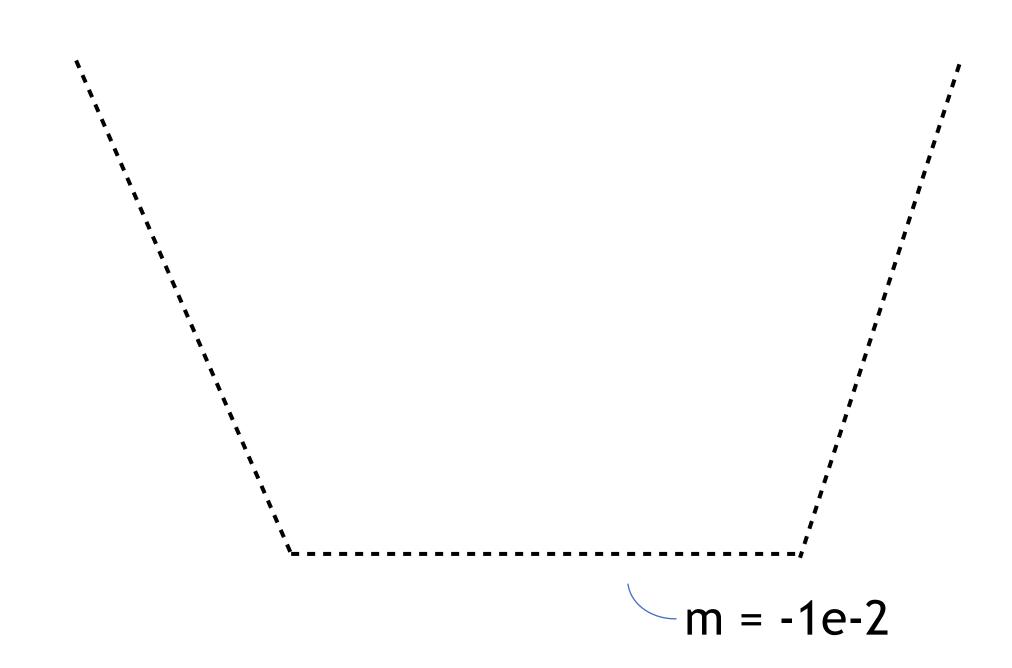
This increases the rate of convergence of the subgradient optimization algorithm.

Thank You

Backup

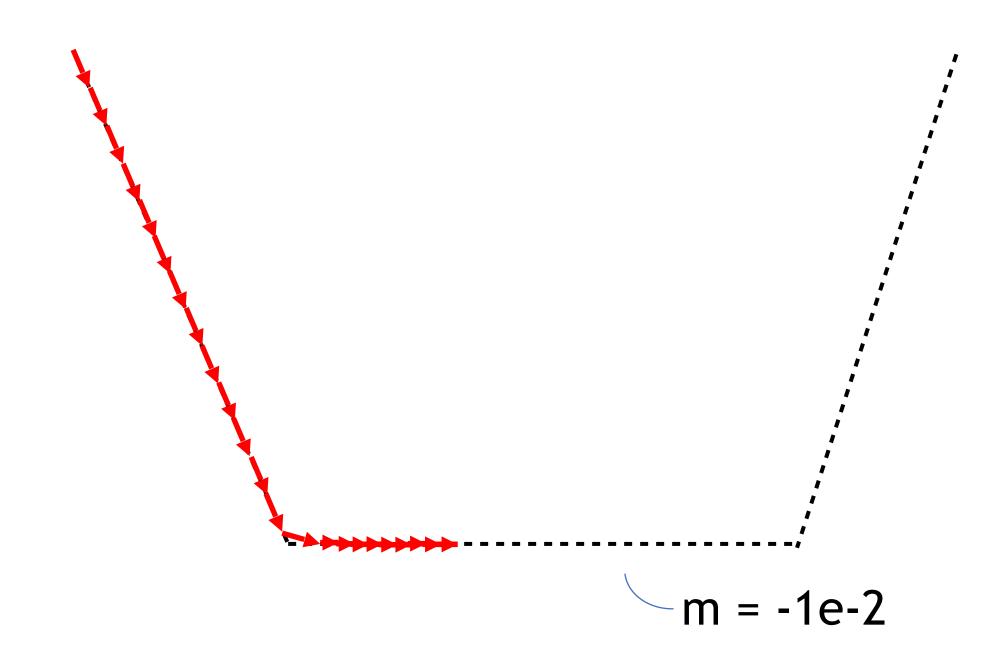
Existing Adaptive Method

Decrease the step size when the rate of change of the objective value becomes low



Existing Adaptive Method

Decrease the step size when the rate of change of the objective value becomes low



Existing Adaptive Method

Decrease the step size when the rate of change of the objective value becomes low

