

Trust Region Policy Optimization (TRPO)

Shusen Wang

<http://wangshusen.github.io/>

Optimization Basics

Gradient Ascent

Problem: Find $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$.

Gradient Ascent

Problem: Find $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$.

- **Gradient ascent** repeats:

1. At $\boldsymbol{\theta}_{\text{old}}$, compute gradient $\mathbf{g} = \frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{\text{old}}}$.
2. Gradient ascent: $\boldsymbol{\theta}_{\text{new}} \leftarrow \boldsymbol{\theta}_{\text{old}} + \alpha \cdot \mathbf{g}$.

Stochastic Gradient Ascent

Problem: Find $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$.

- Assume $J(\boldsymbol{\theta}) = \mathbb{E}_S[V(S; \boldsymbol{\theta})]$.

Stochastic Gradient Ascent

Problem: Find $\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} J(\boldsymbol{\theta})$.

- Assume $J(\boldsymbol{\theta}) = \mathbb{E}_S[V(S; \boldsymbol{\theta})]$.
- **Stochastic gradient ascent** repeats:
 1. $s \leftarrow$ random sampling.
 2. At $\boldsymbol{\theta}_{\text{old}}$, compute gradient $\mathbf{g} = \frac{\partial V(s; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_{\text{old}}}$.
 3. Gradient ascent: $\boldsymbol{\theta}_{\text{new}} \leftarrow \boldsymbol{\theta}_{\text{old}} + \alpha \cdot \mathbf{g}$.

Trust Region

Problem: Find $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$.

- Let $\mathcal{N}(\boldsymbol{\theta}_{\text{old}})$ be a neighborhood of $\boldsymbol{\theta}_{\text{old}}$

Trust Region

Problem: Find $\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} J(\boldsymbol{\theta})$.

- Let $\mathcal{N}(\boldsymbol{\theta}_{\text{old}})$ be a neighborhood of $\boldsymbol{\theta}_{\text{old}}$, e.g.,

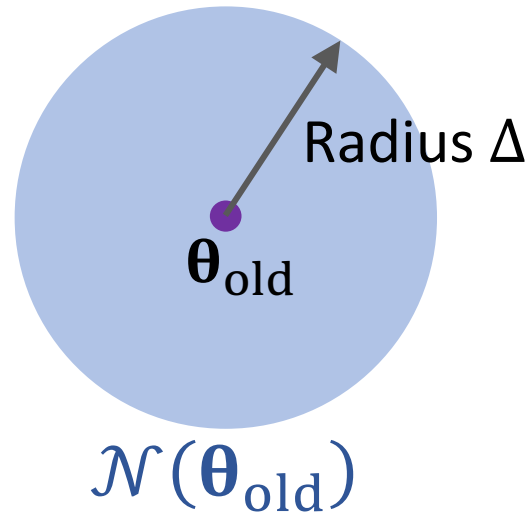
$$\mathcal{N}(\boldsymbol{\theta}_{\text{old}}) = \left\{ \boldsymbol{\theta} \mid \|\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{old}}\|_2 \leq \Delta \right\}.$$

Trust Region

Problem: Find $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$.

- Let $\mathcal{N}(\boldsymbol{\theta}_{\text{old}})$ be a neighborhood of $\boldsymbol{\theta}_{\text{old}}$, e.g.,

$$\mathcal{N}(\boldsymbol{\theta}_{\text{old}}) = \left\{ \boldsymbol{\theta} \mid \|\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{old}}\|_2 \leq \Delta \right\}.$$



Trust Region

Problem: Find $\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} J(\boldsymbol{\theta})$.

- Let $\mathcal{N}(\boldsymbol{\theta}_{\text{old}})$ be a neighborhood of $\boldsymbol{\theta}_{\text{old}}$, e.g.,

$$\mathcal{N}(\boldsymbol{\theta}_{\text{old}}) = \left\{ \boldsymbol{\theta} \mid \|\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{old}}\|_2 \leq \Delta \right\}.$$

- If we have a function, $L(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}})$, that well approximates $J(\boldsymbol{\theta})$ in $\mathcal{N}(\boldsymbol{\theta}_{\text{old}})$, then $\mathcal{N}(\boldsymbol{\theta}_{\text{old}})$ is called “trust region”.

Trust Region Algorithms

Problem: Find $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$.

Trust region algorithms

Trust Region Algorithms

Problem: Find $\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$.

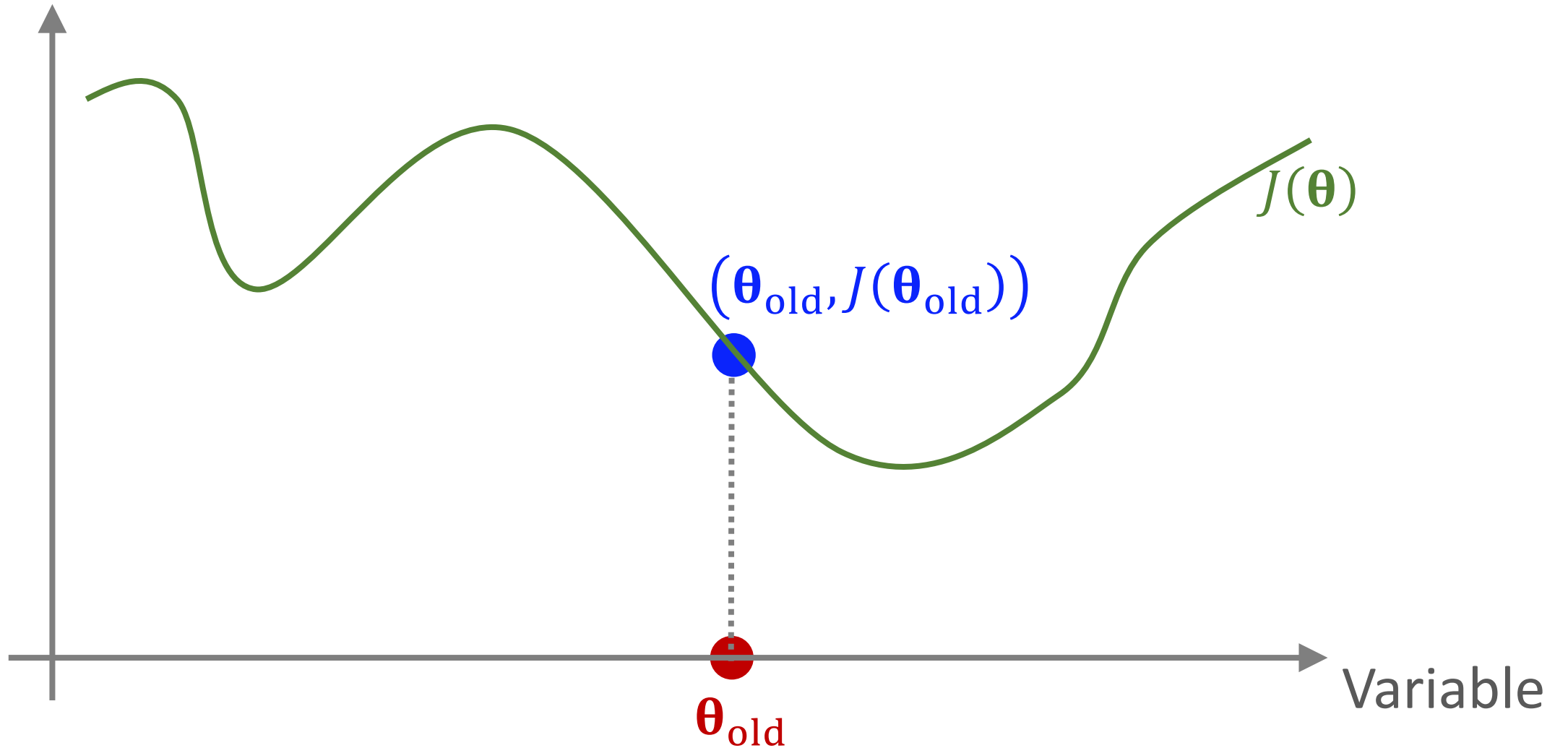
Trust region algorithms repeat:

- 1. Approximation:** Given $\boldsymbol{\theta}_{\text{old}}$, construct $L(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}})$, which is an approximation to $J(\boldsymbol{\theta})$ in $\mathcal{N}(\boldsymbol{\theta}_{\text{old}})$.
- 2. Maximization:** In the trust region, find $\boldsymbol{\theta}_{\text{new}}$ by:

$$\boldsymbol{\theta}_{\text{new}} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta} \in \mathcal{N}(\boldsymbol{\theta}_{\text{old}})} L(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}).$$

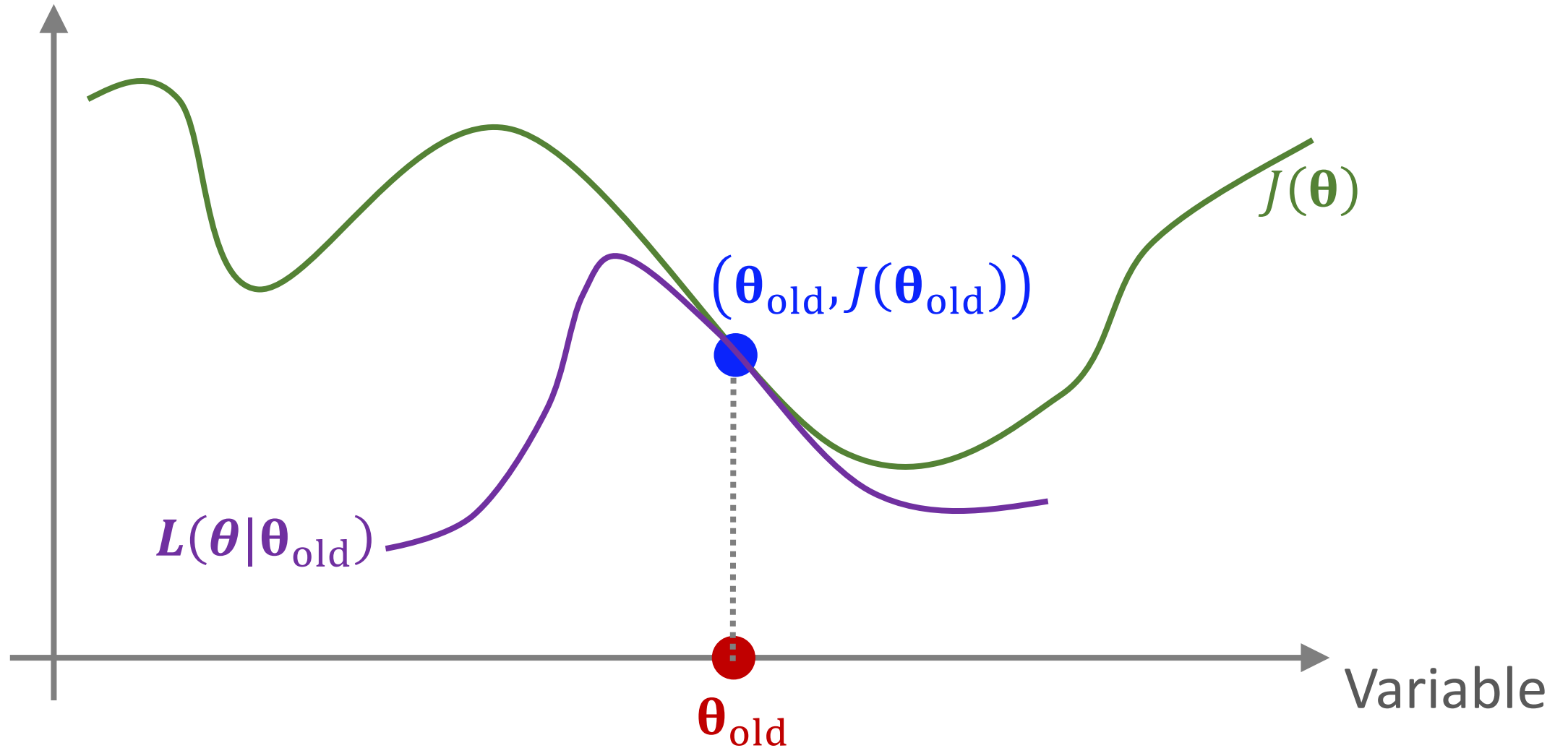
Trust Region Algorithms

Function Value



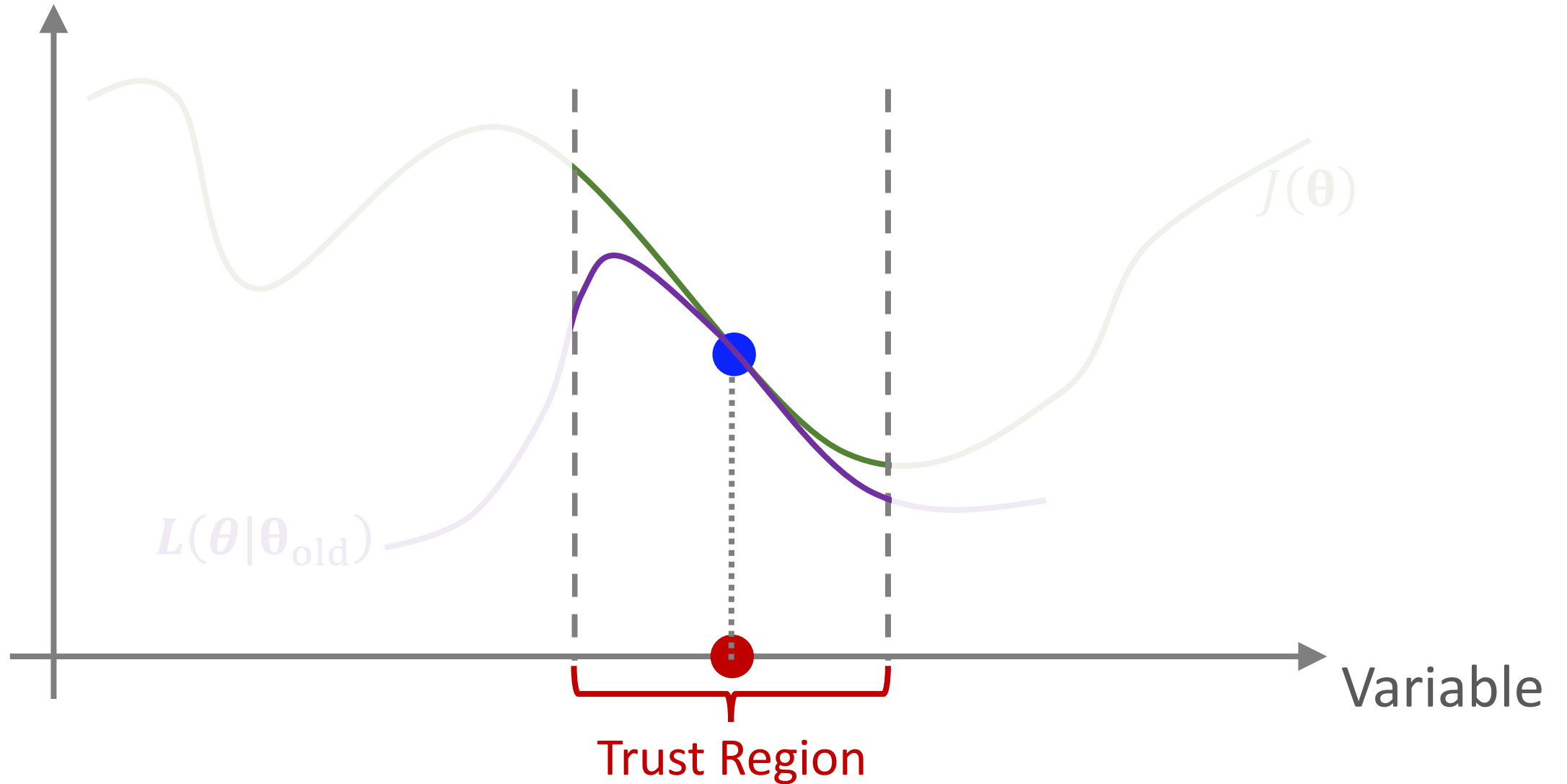
Trust Region Algorithms

Function Value



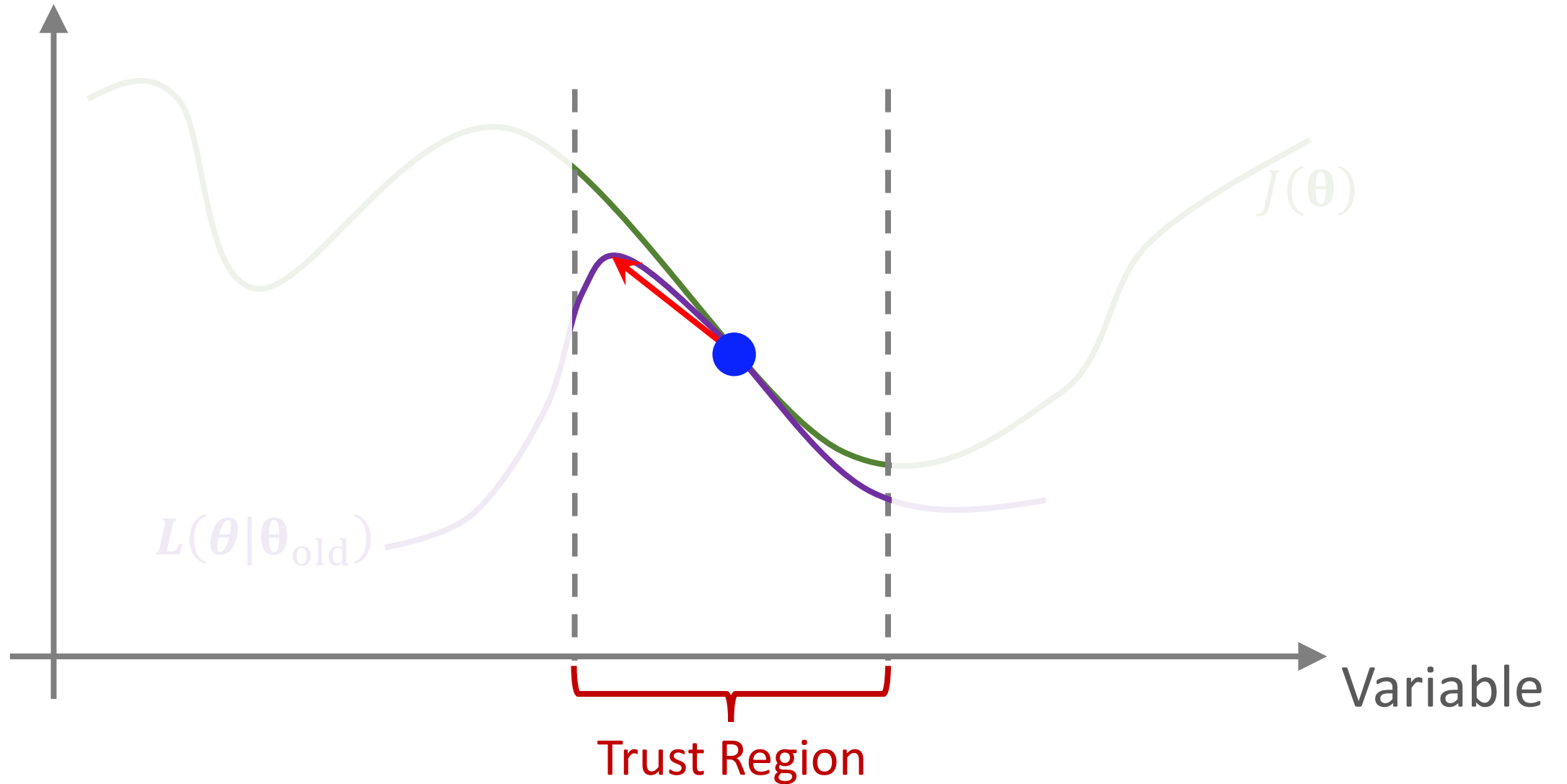
Trust Region Algorithms

Function Value



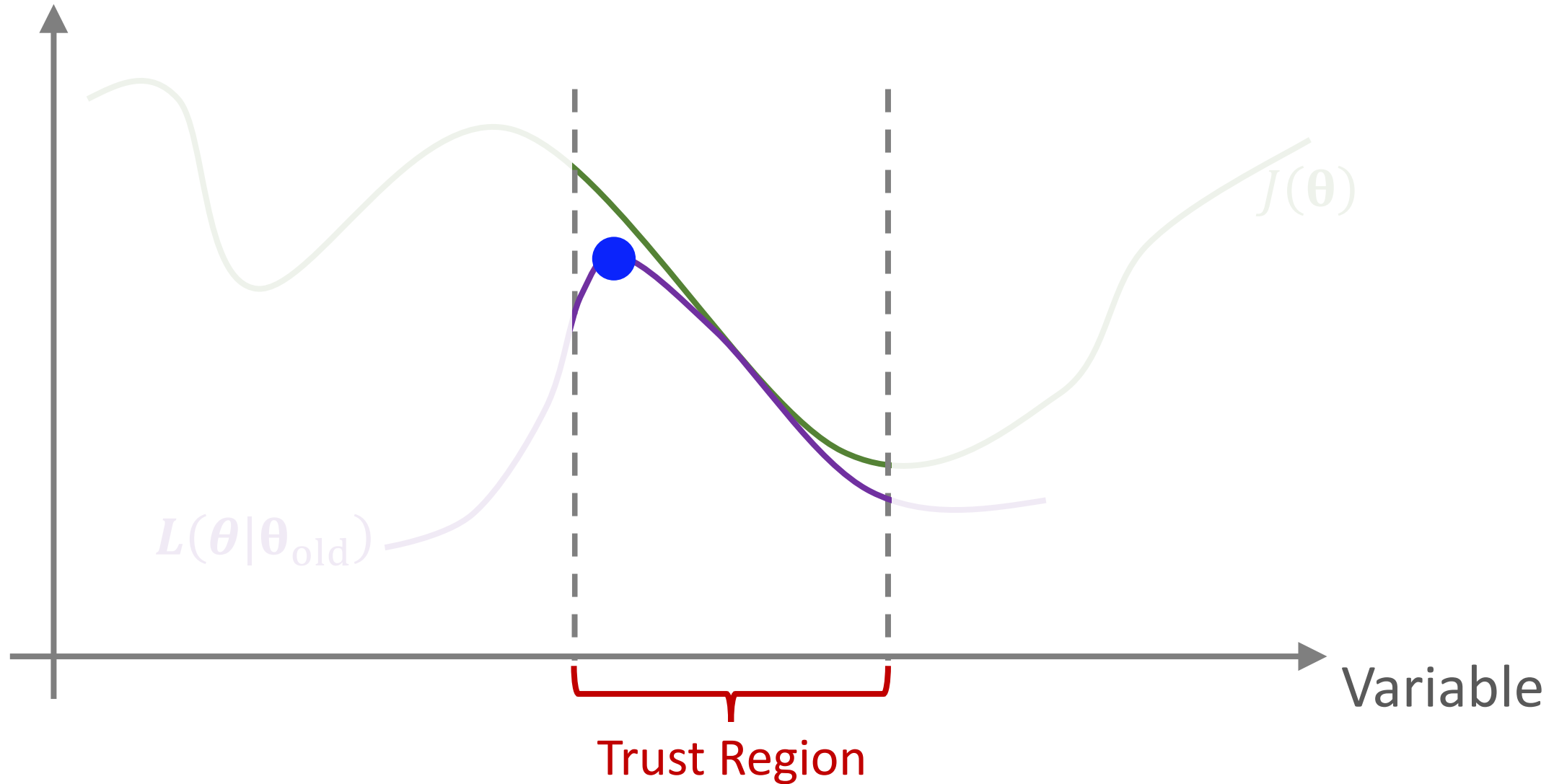
Trust Region Algorithms

Function Value



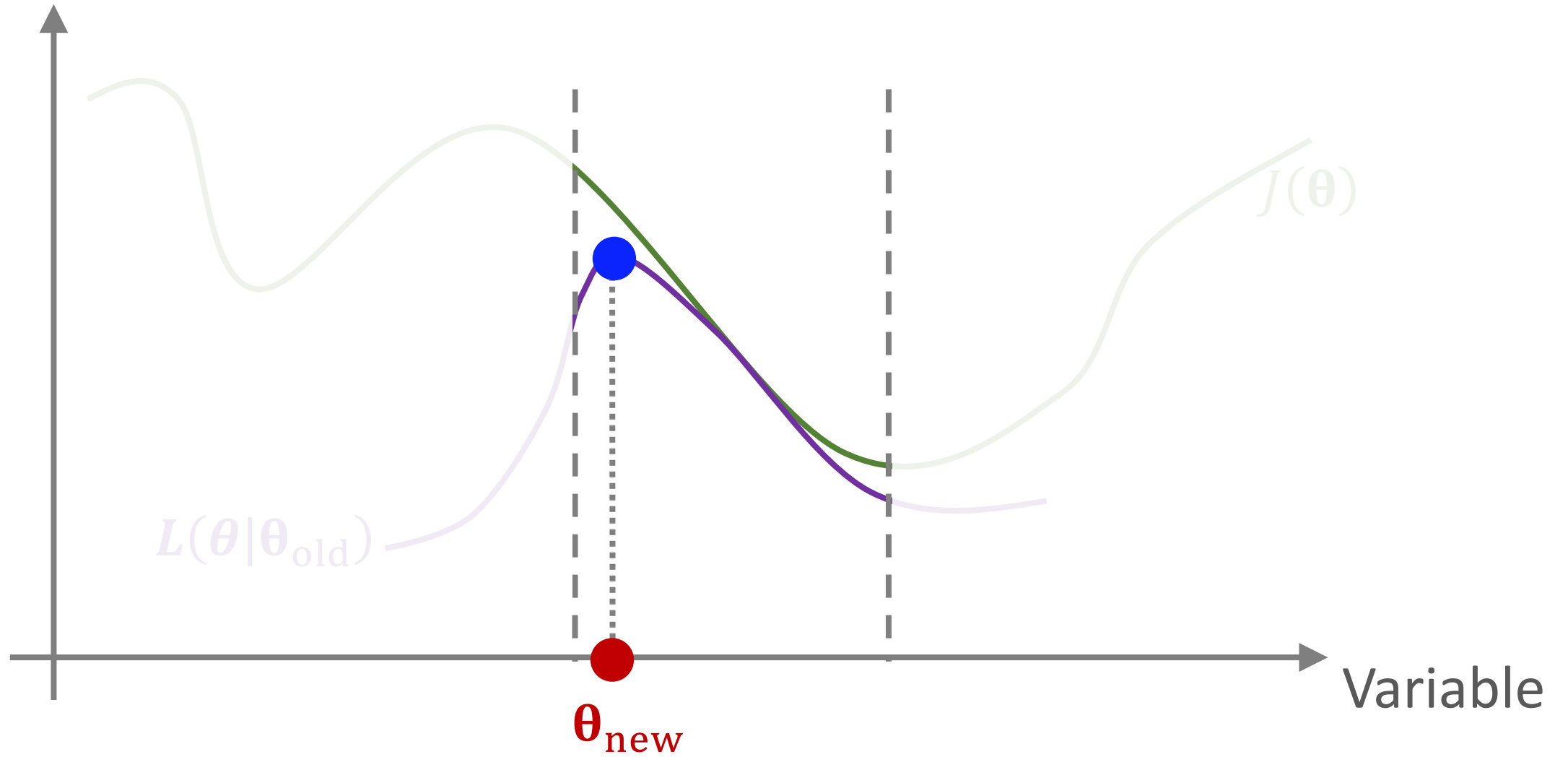
Trust Region Algorithms

Function Value



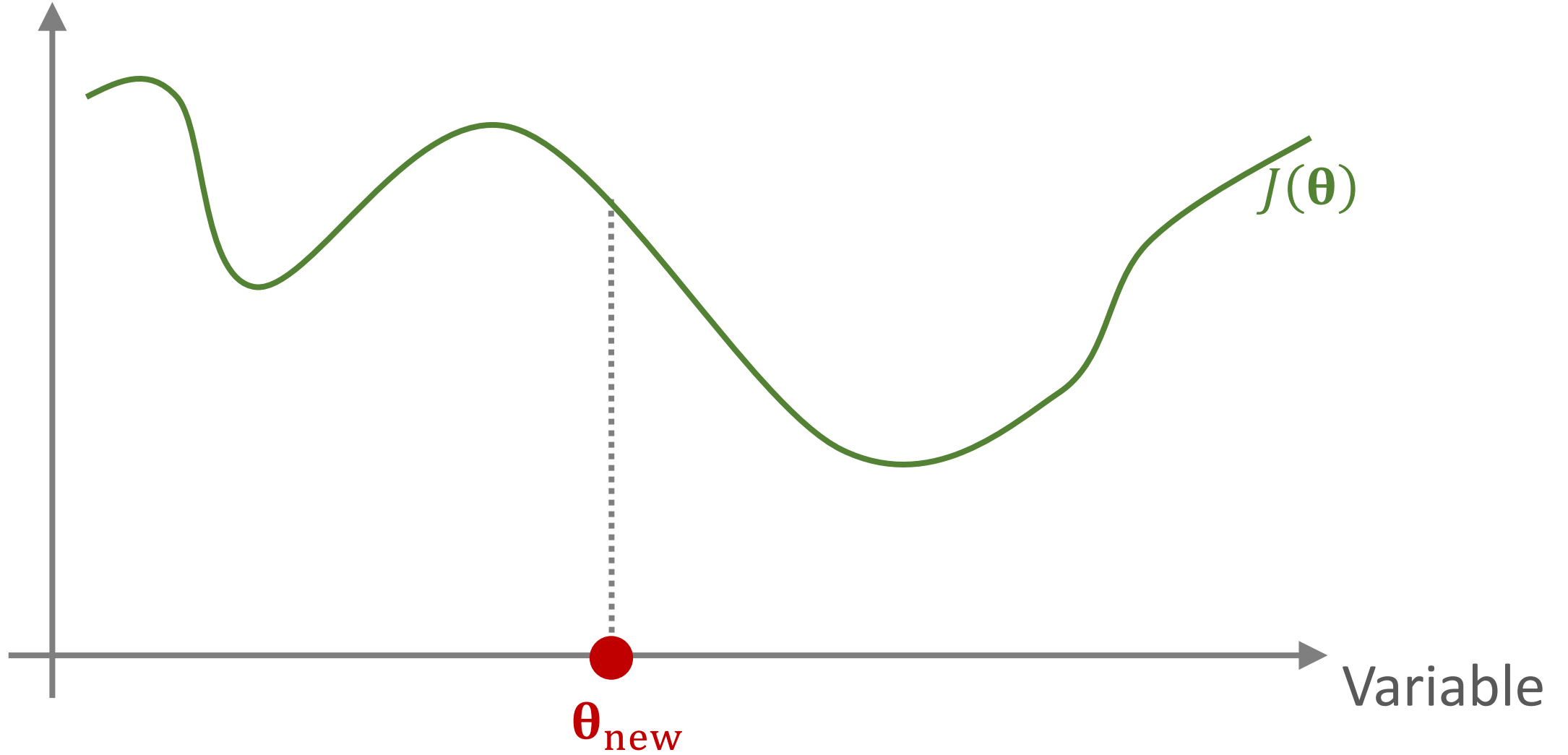
Trust Region Algorithms

Function Value

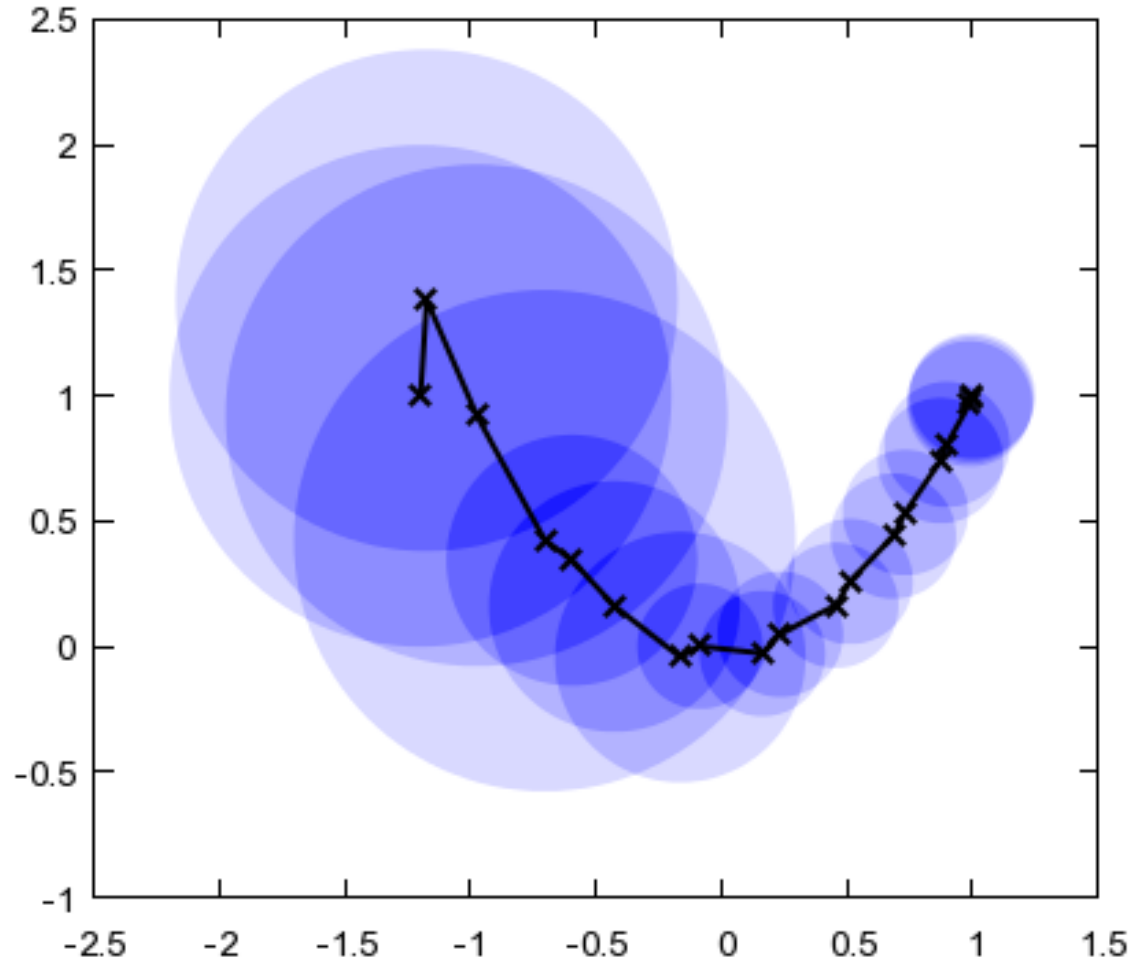


Trust Region Algorithms

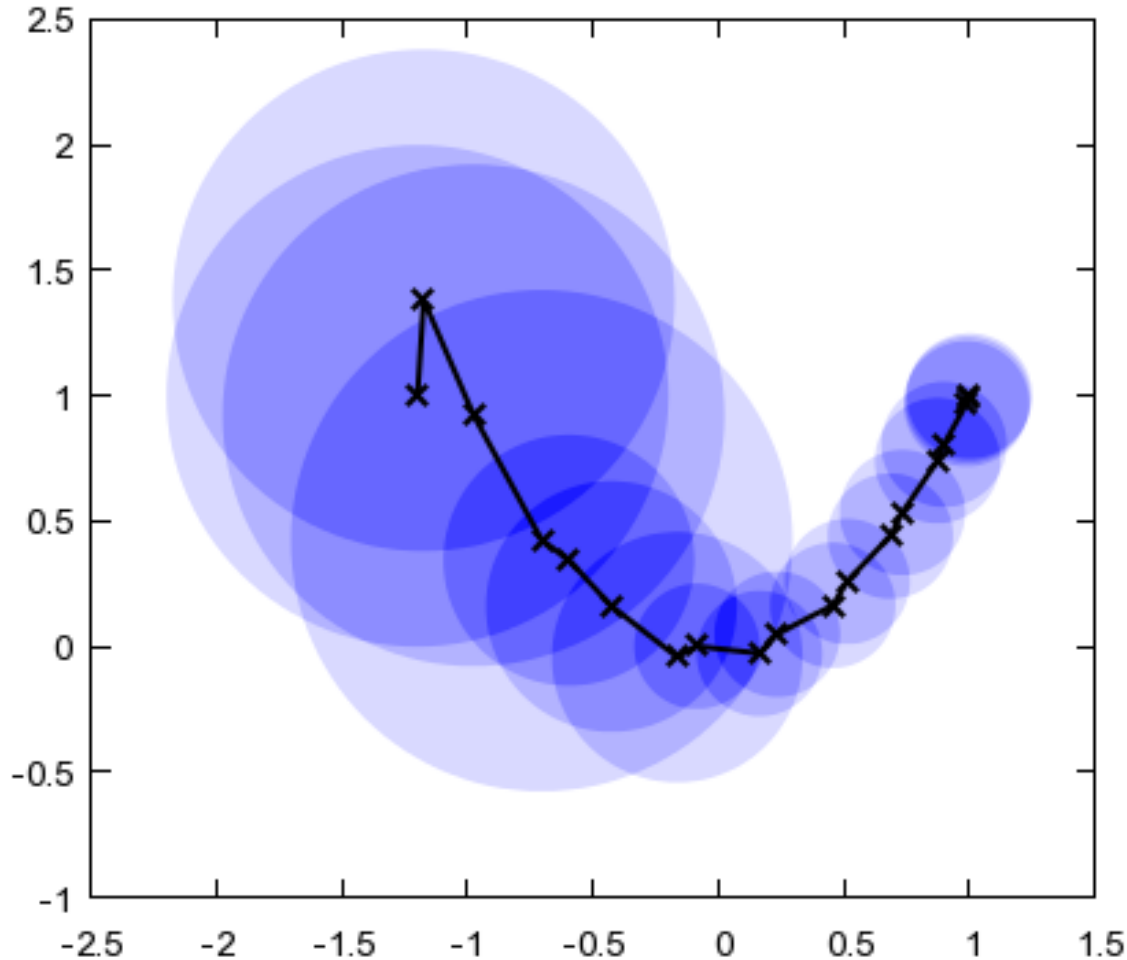
Function Value



Trust Region Algorithms



Trust Region Algorithms



Trust region algorithms repeat:

- 1. Approximation:** Given θ_{old} , construct L which approximates J in the neighborhood of θ_{old} .
- 2. Maximization:** In the trust region (i.e., the neighborhood of θ_{old}), find θ_{new} by maximizing L .

Policy-Based Reinforcement Learning

Policy-Based Reinforcement Learning

- Use policy network, $\pi(a|s; \theta)$, for controlling the agent.

Policy-Based Reinforcement Learning

- Use policy network, $\pi(a|s; \theta)$, for controlling the agent.
- State-value function:

$$\begin{aligned} V_{\pi}(s) &= \mathbb{E}_{A \sim \pi}[Q_{\pi}(s, A)] \\ &= \sum_a \pi(a|s; \theta) \cdot Q_{\pi}(s, a). \end{aligned}$$

Policy-Based Reinforcement Learning

- Use policy network, $\pi(a|s; \theta)$, for controlling the agent.
- State-value function:

$$\begin{aligned} V_{\pi}(s) &= \mathbb{E}_{A \sim \pi}[Q_{\pi}(s, A)] \\ &= \sum_a \pi(a|s; \theta) \cdot Q_{\pi}(s, a). \end{aligned}$$

- Objective function:

$$J(\theta) = \mathbb{E}_S[V_{\pi}(S)].$$

Deriving Objective Function

- State-value function:

$$V_{\pi}(s) = \sum_a \pi(a \mid s; \theta) \cdot Q_{\pi}(s, a)$$

Deriving Objective Function

- State-value function:

$$\begin{aligned} V_{\pi}(s) &= \sum_a \pi(a \mid s; \theta) \cdot Q_{\pi}(s, a) \\ &= \sum_a \pi(a \mid s; \theta_{\text{old}}) \cdot \frac{\pi(a \mid s; \theta)}{\pi(a \mid s; \theta_{\text{old}})} \cdot Q_{\pi}(s, a) \end{aligned}$$

Deriving Objective Function

- State-value function:

$$\begin{aligned} V_{\pi}(s) &= \sum_a \pi(a \mid s; \theta) \cdot Q_{\pi}(s, a) \\ &= \sum_a \pi(a \mid s; \theta_{\text{old}}) \cdot \frac{\pi(a \mid s; \theta)}{\pi(a \mid s; \theta_{\text{old}})} \cdot Q_{\pi}(s, a) \\ &= \mathbb{E}_{A \sim \pi(\cdot \mid s; \theta_{\text{old}})} \left[\frac{\pi(A \mid s; \theta)}{\pi(A \mid s; \theta_{\text{old}})} \cdot Q_{\pi}(s, A) \right]. \end{aligned}$$

Deriving Objective Function

- Objective function:

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathcal{S}}[V_{\pi}(\mathcal{S})]$$

Deriving Objective Function

- Objective function:

$$\begin{aligned} J(\boldsymbol{\theta}) &= \mathbb{E}_{\mathcal{S}}[V_{\pi}(\mathcal{S})] \\ &= \mathbb{E}_{\mathcal{S}} \left[\mathbb{E}_{\mathcal{A}} \left[\frac{\pi(\mathcal{A} \mid \mathcal{S}; \boldsymbol{\theta})}{\pi(\mathcal{A} \mid \mathcal{S}; \boldsymbol{\theta}_{\text{old}})} \cdot Q_{\pi}(\mathcal{S}, \mathcal{A}) \right] \right]. \end{aligned}$$

Trust Region Policy Optimization (TRPO)

Reference:

- Schulman, Levine, Abbeel, Jordan, & Moritz. [Trust region policy optimization](#). In *ICML*, 2015.

Why TRPO?

- More robust than policy gradient algorithms.
- More sample efficient than policy gradient algorithms.

Overview of TRPO

Repeat the two steps:

1. **Approximation:** Given θ_{old} , we construct $L(\theta | \theta_{\text{old}})$ that approximates $J(\theta)$ in the neighborhood of θ_{old} .
2. **Maximization:** In the trust region, $\mathcal{N}(\theta_{\text{old}})$, find θ_{new} by:

$$\theta_{\text{new}} \leftarrow \underset{\theta \in \mathcal{N}(\theta_{\text{old}})}{\operatorname{argmax}} L(\theta | \theta_{\text{old}}).$$

Step 1: Approximation


Approximate $J(\boldsymbol{\theta})$ in the neighborhood of $\boldsymbol{\theta}_{\text{old}}$.

Step 1: Approximation

Objective function: $J(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{S}, \boldsymbol{A}} \left[\frac{\pi(\boldsymbol{A} | \boldsymbol{S}; \boldsymbol{\theta})}{\pi(\boldsymbol{A} | \boldsymbol{S}; \boldsymbol{\theta}_{\text{old}})} \cdot Q_{\pi}(\boldsymbol{S}, \boldsymbol{A}) \right].$

Step 1: Approximation

Objective function: $J(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{s}, \boldsymbol{a}} \left[\frac{\pi(\boldsymbol{a} | \boldsymbol{s}; \boldsymbol{\theta})}{\pi(\boldsymbol{a} | \boldsymbol{s}; \boldsymbol{\theta}_{\text{old}})} \cdot Q_{\pi}(\boldsymbol{s}, \boldsymbol{a}) \right].$



- \boldsymbol{s} is sampled from the **state transition** of the environment.
- \boldsymbol{a} is sampled from the **policy** $\pi(\boldsymbol{a} | \boldsymbol{s}; \boldsymbol{\theta}_{\text{old}})$.

Step 1: Approximation

Objective function: $J(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{S}, \boldsymbol{A}} \left[\frac{\pi(\boldsymbol{A} \mid \boldsymbol{S}; \boldsymbol{\theta})}{\pi(\boldsymbol{A} \mid \boldsymbol{S}; \boldsymbol{\theta}_{\text{old}})} \cdot Q_{\pi}(\boldsymbol{S}, \boldsymbol{A}) \right].$

- Controlled by the policy, $\pi(\boldsymbol{A} \mid \boldsymbol{S}; \boldsymbol{\theta}_{\text{old}})$, the agent collects a trajectory:

$$\boldsymbol{s}_1, \boldsymbol{a}_1, \boldsymbol{r}_1, \boldsymbol{s}_2, \boldsymbol{a}_2, \boldsymbol{r}_2, \dots, \boldsymbol{s}_n, \boldsymbol{a}_n, \boldsymbol{r}_n.$$

Step 1: Approximation

$$\text{Objective function: } J(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{s}, \boldsymbol{A}} \left[\frac{\pi(\boldsymbol{A} | \boldsymbol{s}; \boldsymbol{\theta})}{\pi(\boldsymbol{A} | \boldsymbol{s}; \boldsymbol{\theta}_{\text{old}})} \cdot Q_{\pi}(\boldsymbol{s}, \boldsymbol{A}) \right].$$

- Controlled by the policy, $\pi(\boldsymbol{A} | \boldsymbol{s}; \boldsymbol{\theta}_{\text{old}})$, the agent collects a trajectory:

$$\boldsymbol{s}_1, \boldsymbol{a}_1, \boldsymbol{r}_1, \boldsymbol{s}_2, \boldsymbol{a}_2, \boldsymbol{r}_2, \dots, \boldsymbol{s}_n, \boldsymbol{a}_n, \boldsymbol{r}_n.$$

- Monte Carlo approximation to the expectation:

$$L(\boldsymbol{\theta} | \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(\boldsymbol{a}_i | \boldsymbol{s}_i; \boldsymbol{\theta})}{\pi(\boldsymbol{a}_i | \boldsymbol{s}_i; \boldsymbol{\theta}_{\text{old}})} \cdot Q_{\pi}(\boldsymbol{s}_i, \boldsymbol{a}_i).$$

Step 1: Approximation

$$\text{Approximate } J(\boldsymbol{\theta}) \text{ by } L(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_i \mid s_i; \boldsymbol{\theta})}{\pi(a_i \mid s_i; \boldsymbol{\theta}_{\text{old}})} \cdot Q_{\pi}(s_i, a_i).$$

Step 1: Approximation

$$\text{Approximate } J(\boldsymbol{\theta}) \text{ by } L(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_i \mid s_i; \boldsymbol{\theta})}{\pi(a_i \mid s_i; \boldsymbol{\theta}_{\text{old}})} \cdot Q_{\pi}(s_i, a_i).$$

- Observed rewards in an episode: $r_1, r_2, r_3, \dots, r_n$.
- Discounted return:

$$u_i = r_i + \gamma \cdot r_{i+1} + \gamma^2 \cdot r_{i+2} + \dots + \gamma^{n-i} \cdot r_n.$$

Step 1: Approximation

$$\text{Approximate } J(\boldsymbol{\theta}) \text{ by } L(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_i \mid s_i; \boldsymbol{\theta})}{\pi(a_i \mid s_i; \boldsymbol{\theta}_{\text{old}})} \cdot Q_{\pi}(s_i, a_i).$$

- Observed rewards in an episode: $r_1, r_2, r_3, \dots, r_n$.

- Discounted return:

$$u_i = r_i + \gamma \cdot r_{i+1} + \gamma^2 \cdot r_{i+2} + \dots + \gamma^{n-i} \cdot r_n.$$

- Monte Carlo approximation: $Q_{\pi}(s_i, a_i) \approx u_i$.

Step 1: Approximation

Approximate $J(\boldsymbol{\theta})$ by $L(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_i \mid s_i; \boldsymbol{\theta})}{\pi(a_i \mid s_i; \boldsymbol{\theta}_{\text{old}})} \cdot Q_{\pi}(s_i, a_i).$

$\approx u_i$



Approximate $J(\boldsymbol{\theta})$ by $\tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_i \mid s_i; \boldsymbol{\theta})}{\pi(a_i \mid s_i; \boldsymbol{\theta}_{\text{old}})} \cdot u_i.$

Step 1: Approximation

$$\text{Approximate } J(\boldsymbol{\theta}) \text{ by } \tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(\textcolor{red}{a}_i \mid \textcolor{green}{s}_i; \boldsymbol{\theta})}{\pi(\textcolor{red}{a}_i \mid \textcolor{green}{s}_i; \boldsymbol{\theta}_{\text{old}})} \cdot \textcolor{blue}{u}_i.$$

Step 2: Maximization

$$\text{Approximate } J(\boldsymbol{\theta}) \text{ by } \tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(\textcolor{red}{a}_i \mid \textcolor{green}{s}_i; \boldsymbol{\theta})}{\pi(\textcolor{red}{a}_i \mid \textcolor{green}{s}_i; \boldsymbol{\theta}_{\text{old}})} \cdot \textcolor{blue}{u}_i.$$

- In the trust region, $\mathcal{N}(\boldsymbol{\theta}_{\text{old}})$, find $\boldsymbol{\theta}_{\text{new}}$ by:

$$\boldsymbol{\theta}_{\text{new}} \leftarrow \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}); \quad \text{s.t. } \boldsymbol{\theta} \in \mathcal{N}(\boldsymbol{\theta}_{\text{old}}).$$

Step 2: Maximization

$$\text{Approximate } J(\boldsymbol{\theta}) \text{ by } \tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_i \mid s_i; \boldsymbol{\theta})}{\pi(a_i \mid s_i; \boldsymbol{\theta}_{\text{old}})} \cdot u_i.$$

- In the trust region, $\mathcal{N}(\boldsymbol{\theta}_{\text{old}})$, find $\boldsymbol{\theta}_{\text{new}}$ by:

$$\boldsymbol{\theta}_{\text{new}} \leftarrow \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}); \quad \text{s.t. } \boldsymbol{\theta} \in \mathcal{N}(\boldsymbol{\theta}_{\text{old}}).$$

- **Option 1:** $||\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{old}}|| < \Delta.$

Step 2: Maximization

$$\text{Approximate } J(\boldsymbol{\theta}) \text{ by } \tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_i \mid s_i; \boldsymbol{\theta})}{\pi(a_i \mid s_i; \boldsymbol{\theta}_{\text{old}})} \cdot u_i.$$

- In the trust region, $\mathcal{N}(\boldsymbol{\theta}_{\text{old}})$, find $\boldsymbol{\theta}_{\text{new}}$ by:

$$\boldsymbol{\theta}_{\text{new}} \leftarrow \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}); \quad \text{s.t. } \boldsymbol{\theta} \in \mathcal{N}(\boldsymbol{\theta}_{\text{old}}).$$

- **Option 1:** $\|\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{old}}\| < \Delta.$
- **Option 2:** $\frac{1}{n} \sum_{i=1}^n \text{KL}[\pi(\cdot \mid s_i; \boldsymbol{\theta}_{\text{old}}) \parallel \pi(\cdot \mid s_i; \boldsymbol{\theta})] < \Delta.$

TRPO: Summary

1. Controlled by the policy, $\pi(\cdot \mid \mathbf{s}; \boldsymbol{\theta}_{\text{old}})$, the agent plays a game to the end and observes a trajectory:

$$s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n.$$

2. For $i = 1, \dots, n$, compute discounted returns: $u_i = \sum_{k=i}^n \gamma^{k-i} \cdot r_k$.

TRPO: Summary

1. Controlled by the policy, $\pi(\cdot \mid \mathbf{s}; \boldsymbol{\theta}_{\text{old}})$, the agent plays a game to the end and observes a trajectory:

$$s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_n, a_n, r_n.$$

2. For $i = 1, \dots, n$, compute discounted returns: $u_i = \sum_{k=i}^n \gamma^{k-i} \cdot r_k$.

3. **Approximation:** $\tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_i \mid s_i; \boldsymbol{\theta})}{\pi(a_i \mid s_i; \boldsymbol{\theta}_{\text{old}})} \cdot u_i$.

TRPO: Summary

1. Controlled by the policy, $\pi(\cdot \mid \mathbf{s}; \boldsymbol{\theta}_{\text{old}})$, the agent plays a game to the end and observes a trajectory:

$$\mathbf{s}_1, \mathbf{a}_1, \mathbf{r}_1, \mathbf{s}_2, \mathbf{a}_2, \mathbf{r}_2, \dots, \mathbf{s}_n, \mathbf{a}_n, \mathbf{r}_n.$$

2. For $i = 1, \dots, n$, compute discounted returns: $u_i = \sum_{k=i}^n \gamma^{k-i} \cdot \mathbf{r}_k$.

3. **Approximation:** $\tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}) = \frac{1}{n} \sum_{i=1}^n \frac{\pi(\mathbf{a}_i \mid \mathbf{s}_i; \boldsymbol{\theta})}{\pi(\mathbf{a}_i \mid \mathbf{s}_i; \boldsymbol{\theta}_{\text{old}})} \cdot u_i$.

4. **Maximization:**

$$\boldsymbol{\theta}_{\text{new}} \leftarrow \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \tilde{L}(\boldsymbol{\theta} \mid \boldsymbol{\theta}_{\text{old}}); \quad \text{s. t.} \quad \|\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{old}}\| < \Delta.$$

Policy Gradient versus TRPO

- They are both policy-based reinforcement learning; they have the same objective function:

$$J(\boldsymbol{\theta}) = \mathbb{E}_{\mathcal{S}}[V_{\pi}(\mathcal{S})].$$

- Policy gradient algorithms maximize $J(\boldsymbol{\theta})$ by stochastic gradient ascent.
- TRPO maximizes $J(\boldsymbol{\theta})$ by trust-region algorithm.

Thank you!

<http://wangshusen.github.io/>