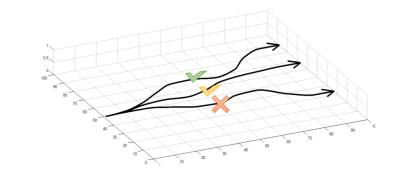
## Evaluating the policy gradient

recall: 
$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right] \approx \frac{1}{N} \sum_{i} \sum_{t} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$



$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[ \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) \right) \left( \sum_{t=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right) \right]$$

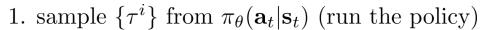
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

generate samples (i.e. run the policy)

## fit a model to estimate return





2. 
$$\nabla_{\theta} J(\theta) \approx \sum_{i} \left( \sum_{t} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t}^{i} | \mathbf{s}_{t}^{i}) \right) \left( \sum_{t} r(\mathbf{s}_{t}^{i}, \mathbf{a}_{t}^{i}) \right)$$

3. 
$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$



improve the policy