



Aalto University
School of Electrical
Engineering

ELEC-E8125 Reinforcement Learning Model-based RL

Joni Pajarinen

24.10.2023

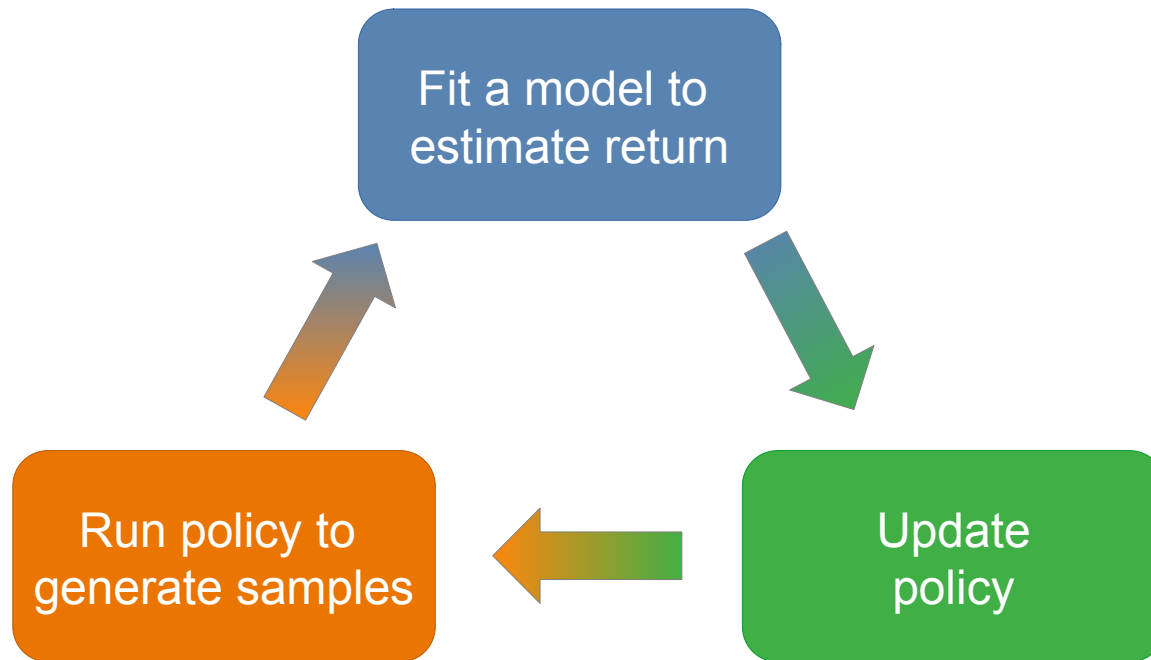
Learning goals

- Understand how optimal control relates to model-based reinforcement learning
- Understand basic sampling based policy optimization

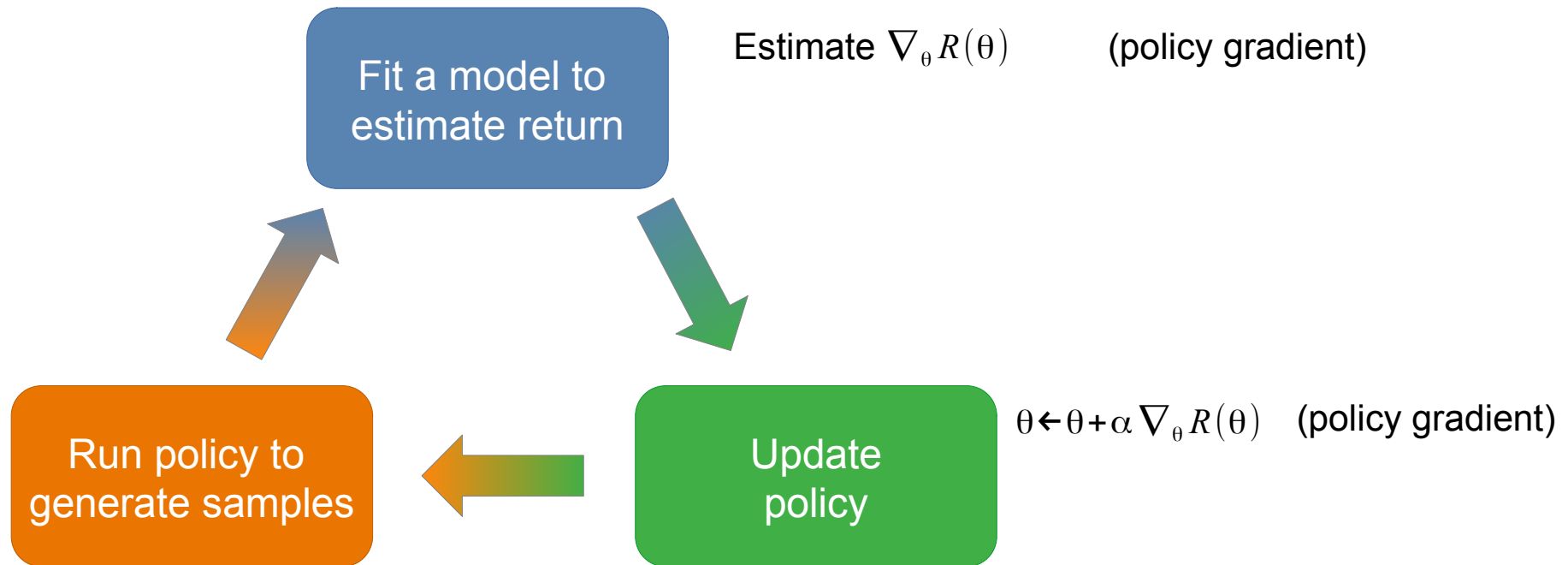
Motivation from two perspectives

- Reinforcement learning has limited sample efficiency
 - Locally optimal control can control complex systems
 - For example, whole body control of a humanoid robot
<https://www.youtube.com/watch?v=vl-8xgJ6ct0>
 - Caveat: optimal control requires knowing the system dynamics
- Learned policies are task, that is, reward-function-specific, learned knowledge cannot be reused

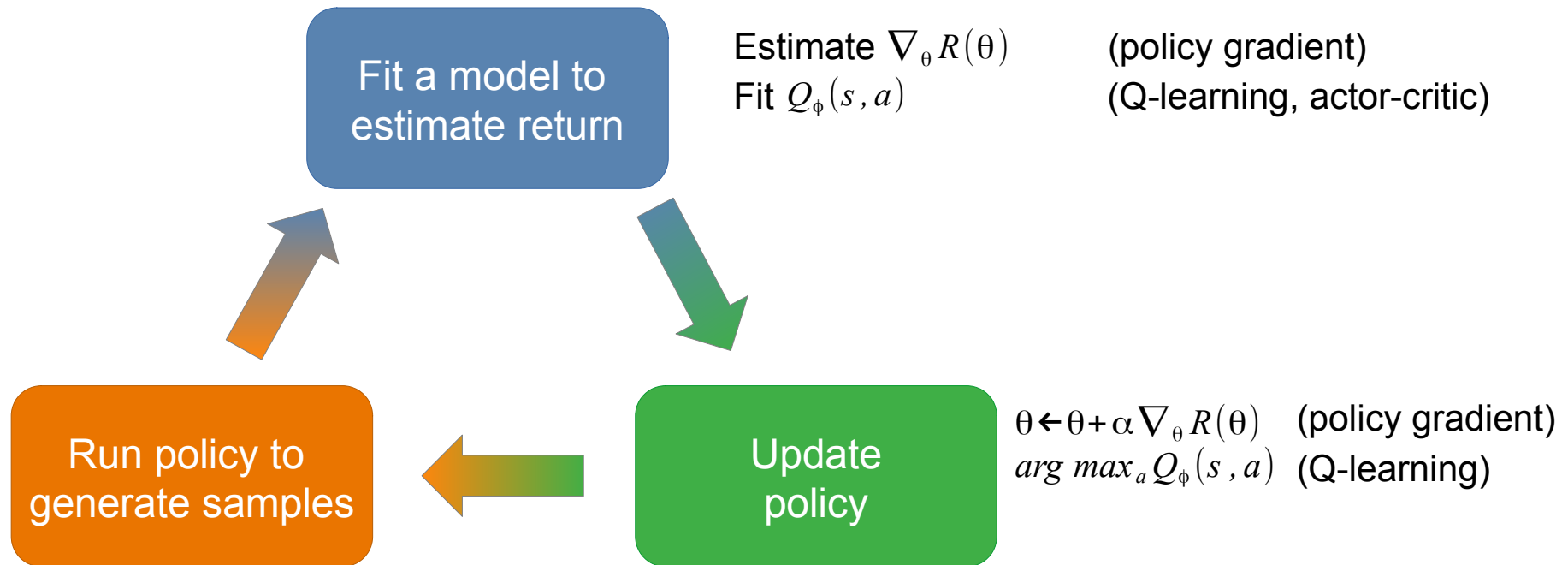
Anatomy of reinforcement learning



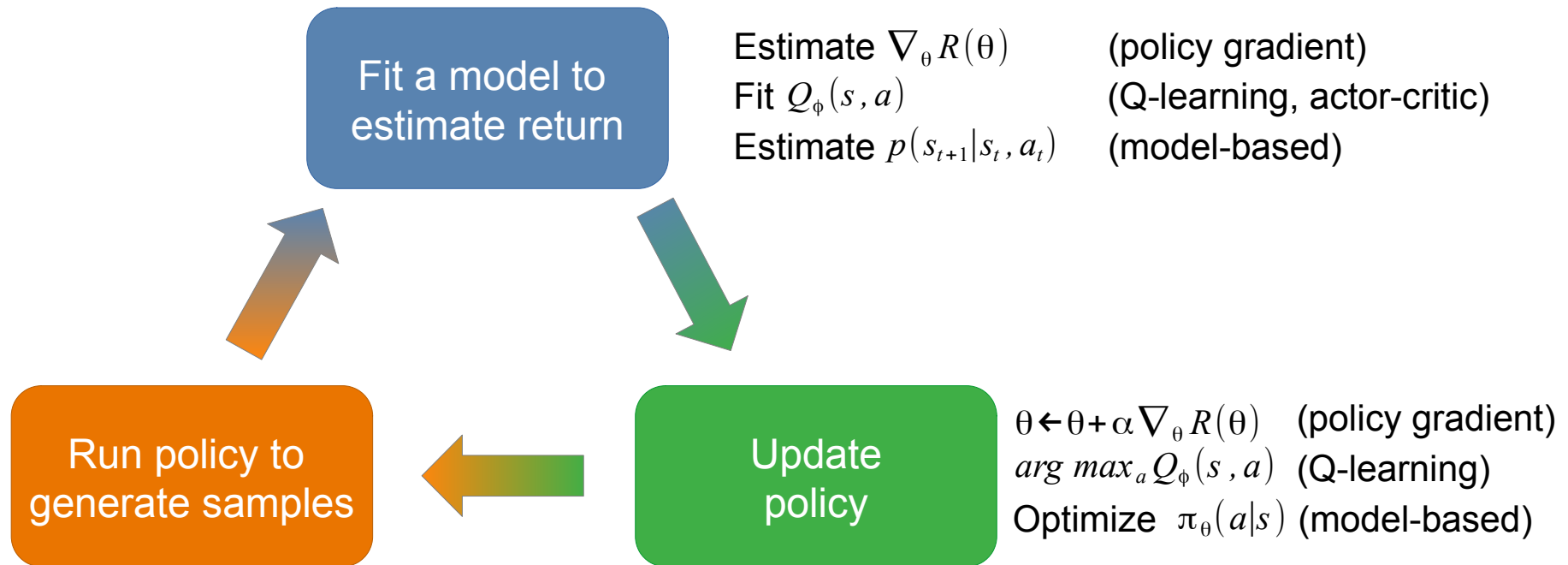
Anatomy of reinforcement learning: *Policy gradient*



Anatomy of reinforcement learning: *Value-function based*

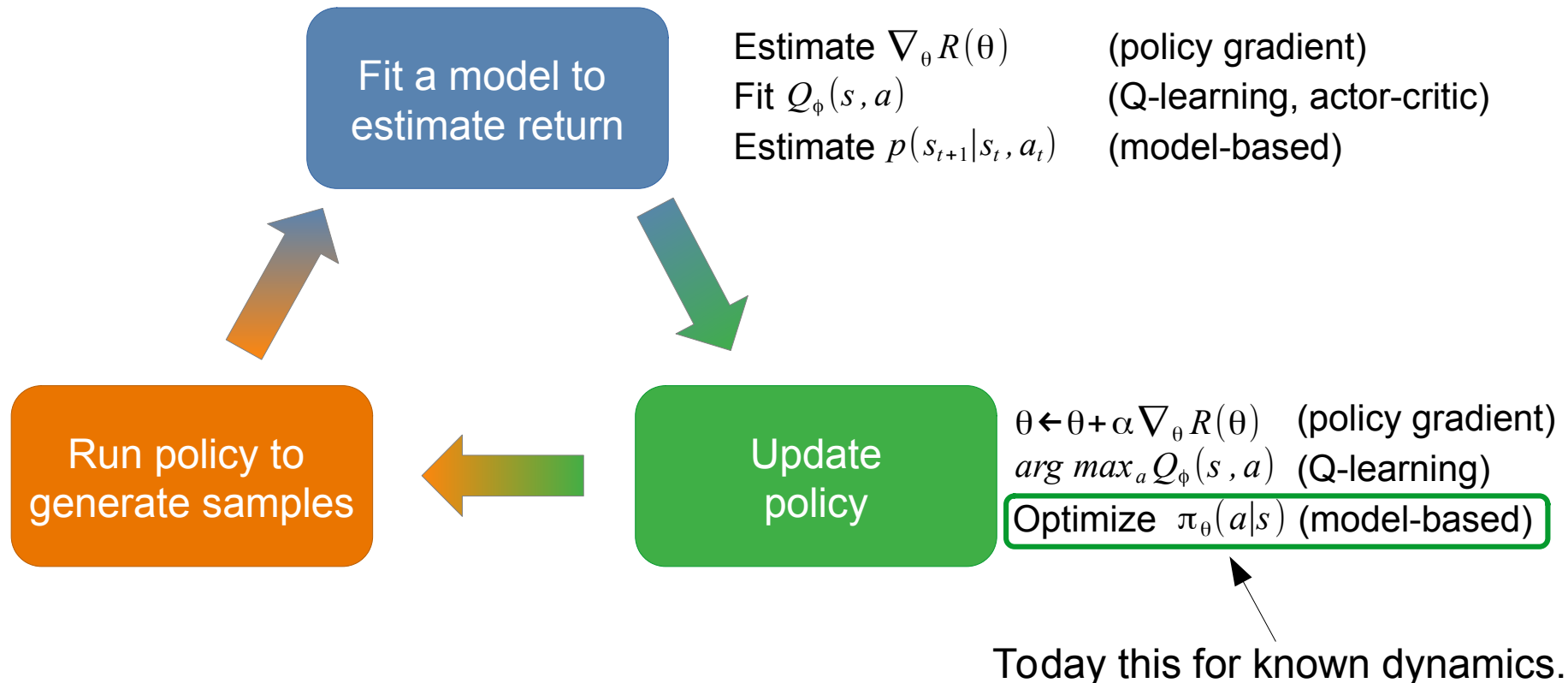


Anatomy of reinforcement learning: *Model-based*



Anatomy of reinforcement learning

Model-based



Solving optimal control problems

Optimal control
optimization objective

$$\min \sum_t c(s_t, a_t)$$

↑
cost
function

Reinforcement learning
optimization objective

$$\max \sum_t r(s_t, a_t)$$

↑
reward
function

$$c(s_t, a_t) = -r(s_t, a_t)$$

Solving (deterministic, finite-horizon) optimal control problems

$$\min_{a_1, \dots, a_T} \sum_t c(s_t, a_t) \quad s.t. \quad s_{t+1} = f(s_t, a_t)$$

cost function system dynamics

Can also be written as:

$$\min_{a_1, \dots, a_T} c(s_1, a_1) + c(f(s_1, a_1), a_2) + \dots + c(f(f(\dots)), a_T)$$

Shooting vs collocation

Shooting methods: Optimize actions

$$\min_{a_1, \dots, a_T} c(s_1, a_1) + c(f(s_1, a_1), a_2) + \dots + c(f(f(\dots)), a_T)$$

Collocation methods: Optimize actions and states
(constrained optimization)

$$\min_{a_1, \dots, a_T, s_1, \dots, s_T} \sum_t c(s_t, a_t) \quad s.t. \quad s_{t+1} = f(s_t, a_t)$$

LQR (linear-quadratic regulator)

Problem definition (finite horizon)

$$\min_{a_1, \dots, a_T} c(s_1, a_1) + c(f(s_1, a_1), a_2) + \dots + c(f(f(\dots)), a_T)$$

$$f(s_t, a_t) = \begin{pmatrix} A_t & B_t \end{pmatrix} \begin{pmatrix} s_t \\ a_t \end{pmatrix} + f_t = F_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + f_t$$

$$c_t(s_t, a_t) = \frac{1}{2} \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T C_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T c_t$$

↑
Note: costs for different time steps may vary.
For example, different costs for final time step.

Note: We will follow notation that clumps together state and action, opposite to traditional control literature, because most recent RL papers use that. We also include the bias term from the beginning.

$$C_t = \begin{pmatrix} C_{s_t, s_t} & C_{s_t, a_t} \\ C_{a_t, s_t} & C_{a_t, a_t} \end{pmatrix}$$

$$c_t = \begin{pmatrix} c_{s_t} \\ c_{a_t} \end{pmatrix}$$

Example system: 1-D particle motion

$$f(s_t, a_t) = F_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + f_t$$

$$c_t(s_t, a_t) = \frac{1}{2} \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T C_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T c_t$$

LQR partial derivation, final step

$$\min_{a_1, \dots, a_T} c(s_1, a_1) + c(f(s_1, a_1), a_2) + \dots + \underbrace{c(f(f(\dots)), a_T)}_{\text{Only cost depending on } a_T}$$

$$f(s_t, a_t) = F_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + f_t$$

$$c_t(s_t, a_t) = \frac{1}{2} \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T C_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T c_t$$

Action-value function:

$$\left. \begin{aligned} Q(s_T, a_T) &= \text{const} + \frac{1}{2} \begin{pmatrix} s_T \\ a_T \end{pmatrix}^T C_T \begin{pmatrix} s_T \\ a_T \end{pmatrix} + \begin{pmatrix} s_T \\ a_T \end{pmatrix}^T c_T \\ \nabla_{a_t} Q(s_T, a_T) &= C_{a_T, s_T} s_T + C_{a_T, a_T} a_t + c_{a_t} = 0 \\ a_T &= -C_{a_T, a_T}^{-1} (C_{a_t, s_t} s_t + c_{a_t}) \end{aligned} \right\}$$



$$\begin{aligned} a_T &= K_T s_T + k_T \\ K_T &= -C_{a_T, a_T}^{-1} C_{a_t, s_t} \\ k_T &= -C_{a_T, a_T}^{-1} c_{a_t} \end{aligned}$$

$$C_t = \begin{pmatrix} C_{s_t, s_t} & C_{s_t, a_t} \\ C_{a_t, s_t} & C_{a_t, a_t} \end{pmatrix}$$

$$c_t = \begin{pmatrix} c_{s_t} \\ c_{a_t} \end{pmatrix}$$

LQR partial derivation, final step

$$\min_{a_1, \dots, a_T} c(s_1, a_1) + c(f(s_1, a_1), a_2) + \dots + c(f(f(\dots)), a_T)$$

$$a_T = K_T s_T + k_T \quad K_T = -C_{a_T, a_T}^{-1} C_{a_i, s_i} \quad k_T = -C_{a_T, a_T}^{-1} c_{a_i}$$

State-value function (by substitution):

$$V(s_T) = \text{const} + \frac{1}{2} \begin{pmatrix} s_T \\ K_T s_T + k_T \end{pmatrix}^T C_T \begin{pmatrix} s_T \\ K_T s_T + k_T \end{pmatrix} + \begin{pmatrix} s_T \\ K_T s_T + k_T \end{pmatrix}^T c_T$$

State value function is quadratic in s_T !

$$V(s_T) = \text{const} + \frac{1}{2} s_T^T V_T s_T + s_T^T v_T$$

$$V(s_T) = \text{const} + \frac{1}{2} s_T^T V_T s_T + s_T^T v_T$$

LQR partial derivation, other steps

$$\begin{aligned}
 Q(s_t, a_t) &= \text{const} + \overbrace{\frac{1}{2} \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T C_t \begin{pmatrix} s_t \\ a_t \end{pmatrix}}^{\text{quadratic}} + \overbrace{\begin{pmatrix} s_t \\ a_t \end{pmatrix}^T c_t + V(f(s_t, a_t))}^{\text{quadratic}} \\
 &= \text{const} + \frac{1}{2} \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T Q_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T q_t
 \end{aligned}$$

$$Q_t = C_t + F_t^T V_{t+1} F_t$$

$$q_t = c_t + F_t^T V_{t+1} f_t + F_t^T v_{t+1}$$

Note: We skip here the derivation of V_t, v_t

$$V(s_T) = \text{const} + \frac{1}{2} s_T^T V_T s_T + s_T^T v_T$$

LQR partial derivation, other steps

$$\begin{aligned}
 Q(s_t, a_t) &= \text{const} + \overbrace{\frac{1}{2} \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T C_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T c_t}^{\text{quadratic}} + \overbrace{V(f(s_t, a_t))}^{\text{quadratic}} \\
 &= \text{const} + \frac{1}{2} \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T Q_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T q_t
 \end{aligned}$$

$$Q_t = C_t + F_t^T V_{t+1} F_t$$

$$q_t = c_t + F_t^T V_{t+1} f_t + F_t^T v_{t+1}$$

$$\nabla_{a_t} Q(s_t, a_t) = Q_{a_t, s_t} s_t + Q_{a_t, a_t} a_t + q_t^T = 0$$

$$a_t = K_t s_t + k_t \quad K_t = -Q_{a_t, a_t}^{-1} Q_{a_t, s_t} \quad k_t = -Q_{a_t, a_t}^{-1} q_{a_t}$$

LQR algorithm

Backward recursion:

For $t = T$ down to 1

$$Q_t = C_t + F_t^T V_{t+1} F_t$$

$$q_t = c_t + F_t^T V_{t+1} f_t + F_t^T v_{t+1}$$

$$K_t = -Q_{a_t, a_t}^{-1} Q_{a_t, s_t}$$

$$k_t = -Q_{a_t, a_t}^{-1} q_{a_t}$$

$$V_t = Q_{s_t, s_t} + Q_{s_t, a_t} K_t + K_t^T Q_{a_t, s_t} + K_t^T Q_{a_t, a_t} K_t$$

$$v_t = q_{s_t} + Q_{s_t, a_t} k_t + K_t^T q_{a_t} + K_t^T Q_{a_t, a_t} k_t$$



First: compute the gains.

Forward recursion:

For $t = 1$ to T

$$a_t = K_t s_t + k_t$$

$$s_{t+1} = f(s_t, a_t)$$



Then: apply the law to compute controls.

System uncertainty / stochastic dynamics

Gaussian noise

$$f(s_t, a_t) = F_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + f_t + w_t \quad w_t \sim N(\mathbf{0}, \Sigma_t)$$

$$p(s_{t+1} | s_t, a_t) \sim N \left(F_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + f_t, \Sigma_t \right)$$

- A linear system with Gaussian noise can be controlled optimally using *separation principle*:
 - Use optimal observer (Kalman filter) to observe state
 - Control system using LQR with mean predicted state
- No change in algorithm!

Non-linear systems - Iterative LQR

- Approximate a non-linear system as a linear-quadratic

$$f(s_t, a_t) = F_t \begin{pmatrix} s_t \\ a_t \end{pmatrix}$$

$$c_t(s_t, a_t) = \frac{1}{2} \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T C_t \begin{pmatrix} s_t \\ a_t \end{pmatrix} + \begin{pmatrix} s_t \\ a_t \end{pmatrix}^T c_t$$

$$f(s_t, a_t) \approx f(\hat{s}_t, \hat{a}_t) + \nabla_{s_t, a_t} f(\hat{s}_t, \hat{a}_t) \begin{pmatrix} s_t - \hat{s}_t \\ a_t - \hat{a}_t \end{pmatrix}$$

$$c_t(s_t, a_t) \approx c(\hat{s}_t, \hat{a}_t) + \frac{1}{2} \begin{pmatrix} s_t - \hat{s}_t \\ a_t - \hat{a}_t \end{pmatrix}^T \nabla_{s_t, a_t}^2 c(\hat{s}_t, \hat{a}_t) \begin{pmatrix} s_t - \hat{s}_t \\ a_t - \hat{a}_t \end{pmatrix} + \nabla_{s_t, a_t} c(\hat{s}_t, \hat{a}_t) \begin{pmatrix} s_t - \hat{s}_t \\ a_t - \hat{a}_t \end{pmatrix}$$

Non-linear systems - Iterative LQR

$$f(s_t, a_t) \approx f(\hat{s}_t, \hat{a}_t) + \nabla_{s_t, a_t} f(\hat{s}_t, \hat{a}_t) \begin{pmatrix} s_t - \hat{s}_t \\ a_t - \hat{a}_t \end{pmatrix}$$

$$c_t(s_t, a_t) = c(\hat{s}_t, \hat{a}_t) + \frac{1}{2} \begin{pmatrix} s_t - \hat{s}_t \\ a_t - \hat{a}_t \end{pmatrix}^T \nabla_{s_t, a_t}^2 c(\hat{s}_t, \hat{a}_t) \begin{pmatrix} s_t - \hat{s}_t \\ a_t - \hat{a}_t \end{pmatrix} + \nabla_{s_t, a_t} c(\hat{s}_t, \hat{a}_t) \begin{pmatrix} s_t - \hat{s}_t \\ a_t - \hat{a}_t \end{pmatrix}$$

$$\bar{f}(\delta s_t, \delta a_t) = F_t \begin{pmatrix} \delta s_t \\ \delta a_t \end{pmatrix}$$

\nwarrow
 $\nabla_{s_t, a_t} f(\hat{s}_t, \hat{a}_t)$

$$\bar{c}_t(\delta s_t, \delta a_t) = \frac{1}{2} \begin{pmatrix} \delta s_t \\ \delta a_t \end{pmatrix}^T C_t \begin{pmatrix} \delta s_t \\ \delta a_t \end{pmatrix} + \begin{pmatrix} \delta s_t \\ \delta a_t \end{pmatrix}^T c_t$$

\nwarrow
 $\nabla_{s_t, a_t}^2 c(\hat{s}_t, \hat{a}_t)$

\nwarrow
 $\nabla_{s_t, a_t} c(\hat{s}_t, \hat{a}_t)$

Iterative LQR (iLQR) – Algorithm outline

Repeat

$$\mathbf{F}_t = \nabla_{s_t, a_t} f(\hat{s}_t, \hat{a}_t)$$

$$\mathbf{C}_t = \nabla_{s_t, a_t}^2 c(\hat{s}_t, \hat{a}_t)$$

$$\mathbf{c}_t = \nabla_{s_t, a_t} c(\hat{s}_t, \hat{a}_t)$$

Run LQR backward pass with $\delta s_t, \delta a_t$

Run LQR forward pass with real dynamics and $\mathbf{a}_t = \mathbf{K}_t \delta s_t + \mathbf{k}_t + \hat{a}_t$

Update \hat{s}_t, \hat{a}_t to results of forward pass

until convergence

Practical considerations:

- Usually receding horizon is used: At every time-step, state is observed, iLQR is applied, and (only) first action is executed.
- On first iteration, gradients can be evaluated at starting point.

Planning by sampling – Shooting methods

Shooting methods: optimize actions

$$V(s_0) = \min_{a_0, \dots, a_{T-1}} c(s_0, a_0) + \dots + c(f(f(\dots)), a_{T-1})$$

$$V(s_0) = \max_{a_0, \dots, a_{T-1}} R(s_0, a_0) + \dots + R(f(f(\dots)), a_{T-1})$$

How to solve? Random shooting:

- Simulate multiple trajectories using random policy (remember Monte Carlo policy evaluation from lecture 3?)

$$Q_\pi(s_0, a_0) \approx \frac{1}{N} \sum_{i=0}^{N-1} \sum_{t=0}^{T-1} \gamma^t r_t$$

- Execute action with lowest cost / highest return
- Repeat

The cross-entropy method (CEM) general background

- Can be used to estimate $P_{\pi}(V(s) \geq d)$
- Monte Carlo estimation does not work well when $P_{\pi}(V(s) \geq d)$ is tiny
- CEM provides efficient estimation based on importance sampling (details in [De Boer 2005])
- The approach can be used also in optimization:
 - Select π to yield high probability for $P_{\pi}(V(s) \geq d)$
 - Increase d to reach higher $V(s)$ values
 - Repeat

CEM for optimization

- Goal: maximize $Q(a)$
- Choose sampling distribution. We choose a Gaussian
 $\pi_{\theta}(a) = \mathcal{N}(a | \mu, \sigma^2)$
- While not converged:
 - Sample N samples a^i from current sampling distribution $\pi_{\theta}(a)$
 - Evaluate objective function $Q(a^i)$ at each a^i
 - Select M ($M < N$) samples a^i with the highest $Q(a^i)$
 - Fit parameters $\theta = (\mu, \sigma^2)$ to the selected samples
 - Repeat

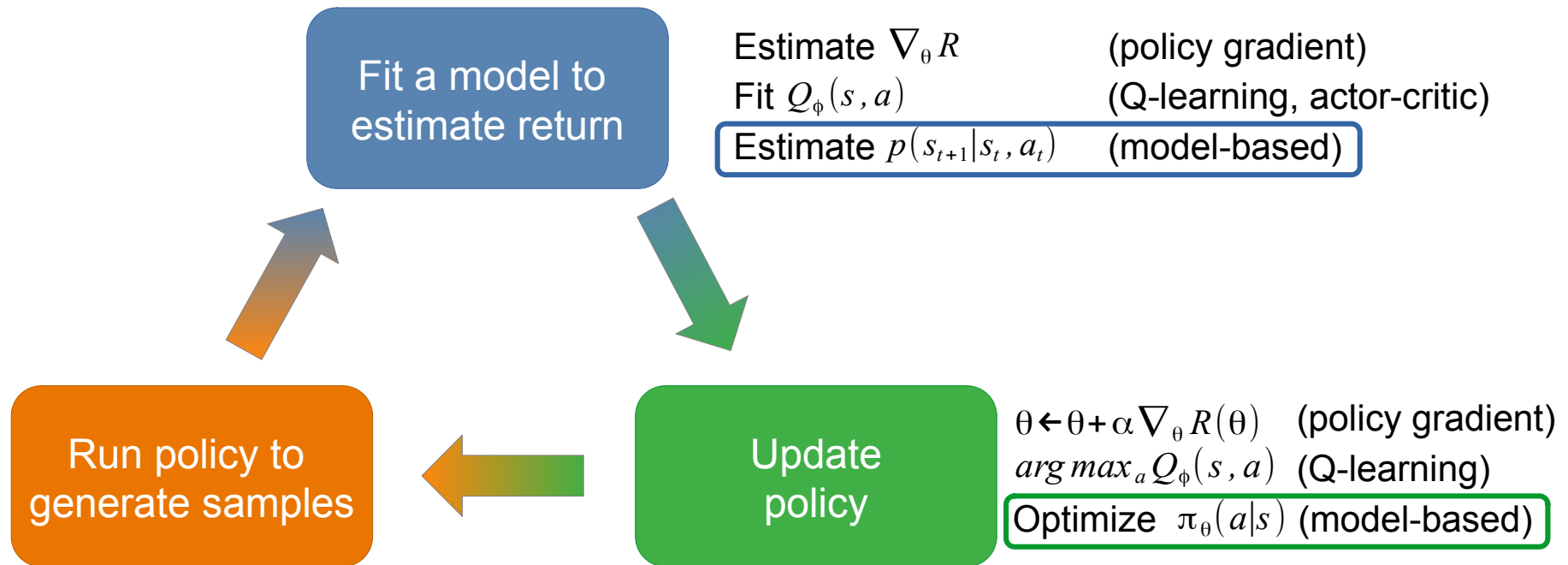
CEM in model-based RL

$$V(s_t^i) = \sum_{k=0}^H \gamma^k r_{t+k}^i$$

- Goal: maximize $V(s_0)$
- Sampling distribution at each time step t:
 $\pi_{\theta(t)}(a_t) = \mathcal{N}(a_t | \mu(t), \sigma^2(t))$
- While not converged:
 - Perform Monte Carlo evaluation (Lecture 3) over N trajectories using sampling distribution and dynamics model
 → We get for trajectory i at time step t sample s_t^i with value $V(s_t^i)$
 - For each time step t fit parameters $\theta(t) = (\mu(t), \sigma^2(t))$ of the sampling distribution to M ($M < N$) samples with the highest $V(s_t^i)$
 - Repeat

Anatomy of reinforcement learning

Model-based



Teaser: Basic iterative model-based RL

Input: base policy π_0

Run base policy to collect data $D \leftarrow \{(s, a, s')_i\}$

Repeat

Fit dynamics model $f(s, a)$ to minimize $\sum_i \|f(s_i, a_i) - s'_i\|^2$

Use model to plan (e.g. iLQR, CEM) actions

Execute first planned action, observe resulting state s'

Update dataset $D \leftarrow D \cup \{(s, a, s')\}$

Viewpoint: Use learned model as “simulator” that allows exploring various options to choose one that is (locally) optimal.

Summary

- Optimal control for linear systems with quadratic costs can be determined with LQR
- Locally optimal control for nonlinear systems can be performed using linearization of dynamics in iterative LQR
- CEM allows for sample based planning with arbitrary costs/reward and dynamics
- Model-based reinforcement learning aims especially to increase data efficiency

Next: Model-based RL – again – but with learned models

- What kind of dynamics model to use?
- How to optimize a general policy function using a dynamics model?
- Reading: Sutton & Barto, ch. 8-8.2. No quiz for next week's lecture.