



Aalto University
School of Electrical
Engineering

ELEC-E8125 Reinforcement Learning

Model-based reinforcement learning

Ville Kyrki

12.11.2019

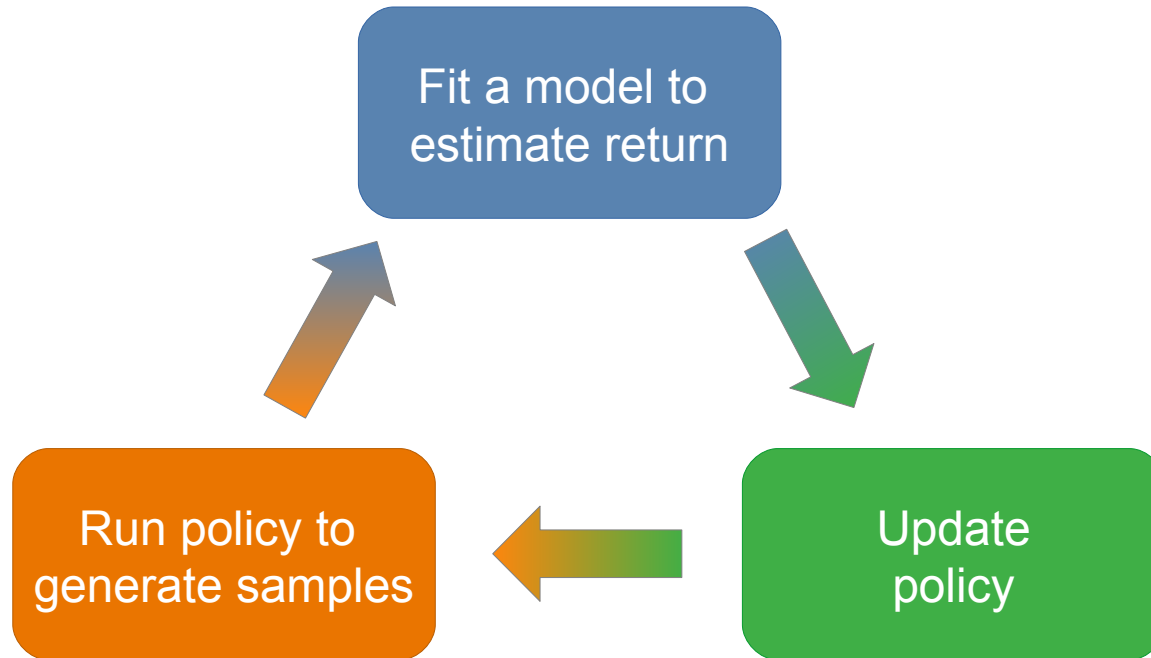
Learning goals

- Understand basic approaches of model-based reinforcement learning.

Motivation from two perspectives

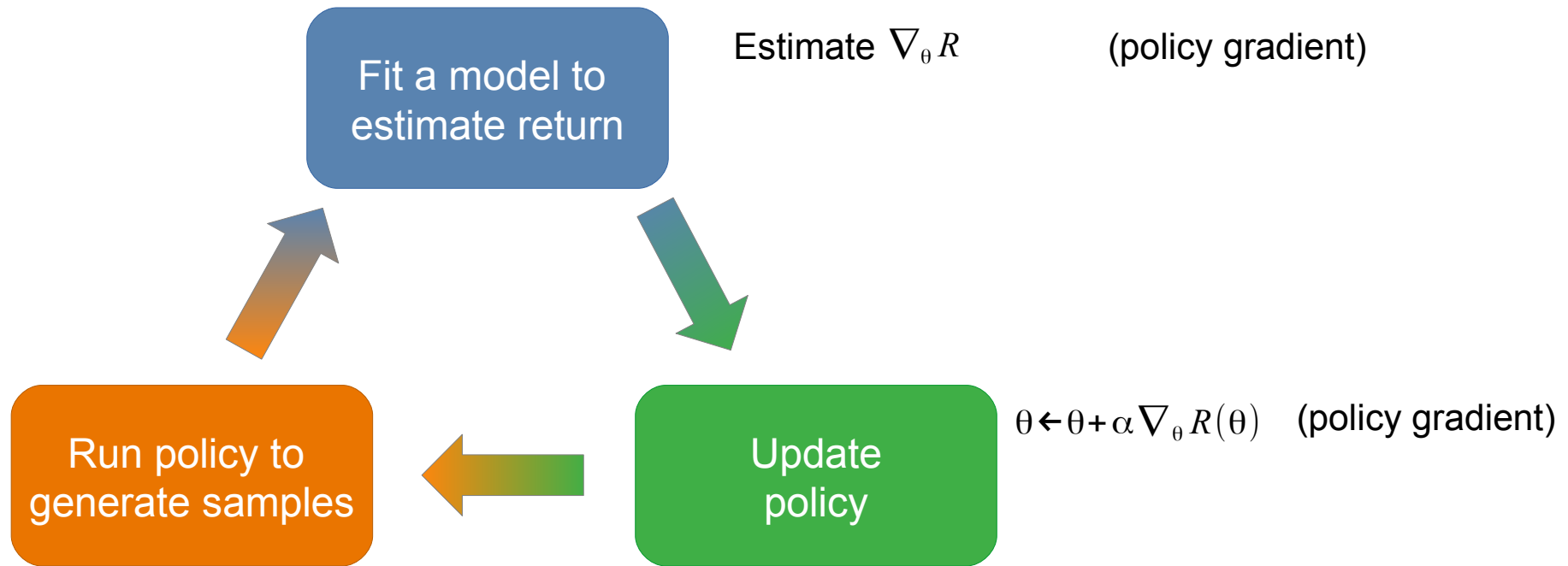
- Reinforcement learning has limited sample efficiency.
 - Locally optimal control has recently shown progress for control of complex systems.
 - For example, whole body control of a humanoid robot
<https://www.youtube.com/watch?v=vl-8xgJ6ct0>
 - But optimal control requires knowing the system dynamics.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

Anatomy of reinforcement learning



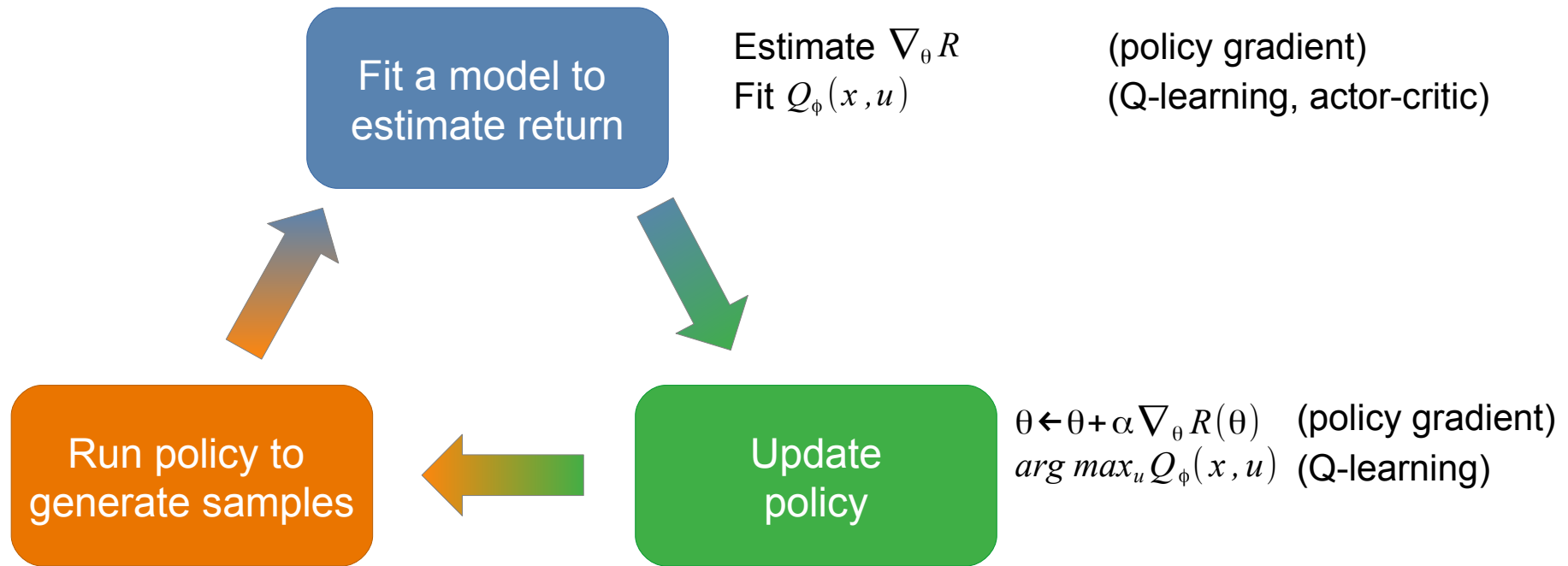
Anatomy of reinforcement learning

Policy gradient



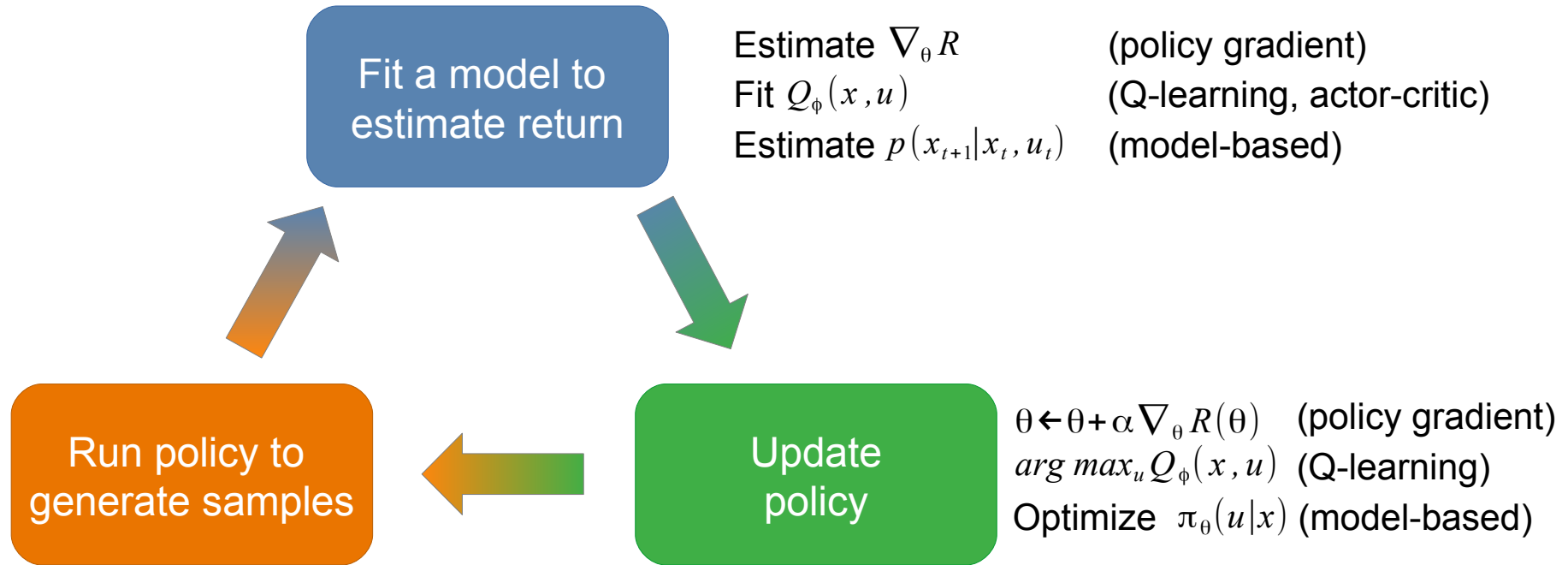
Anatomy of reinforcement learning

Value-function based



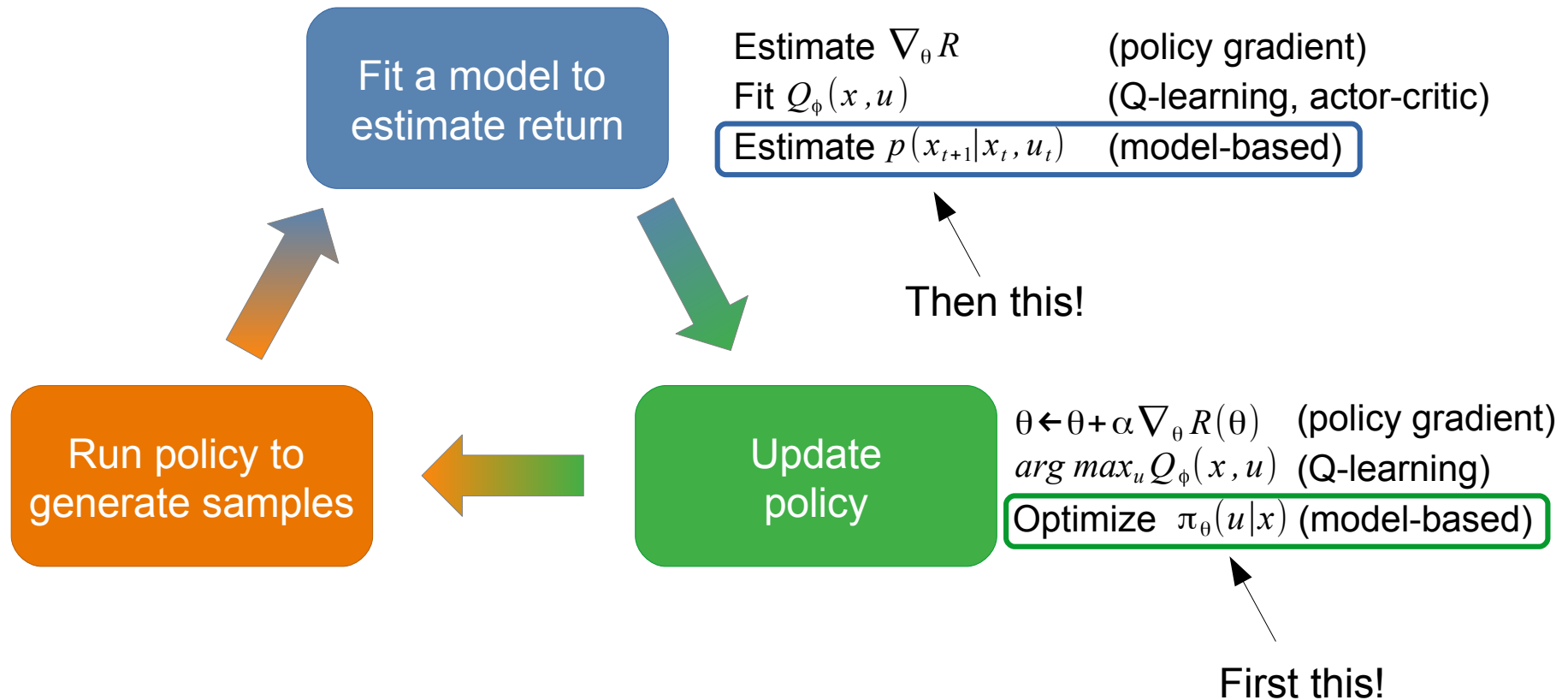
Anatomy of reinforcement learning

Model-based



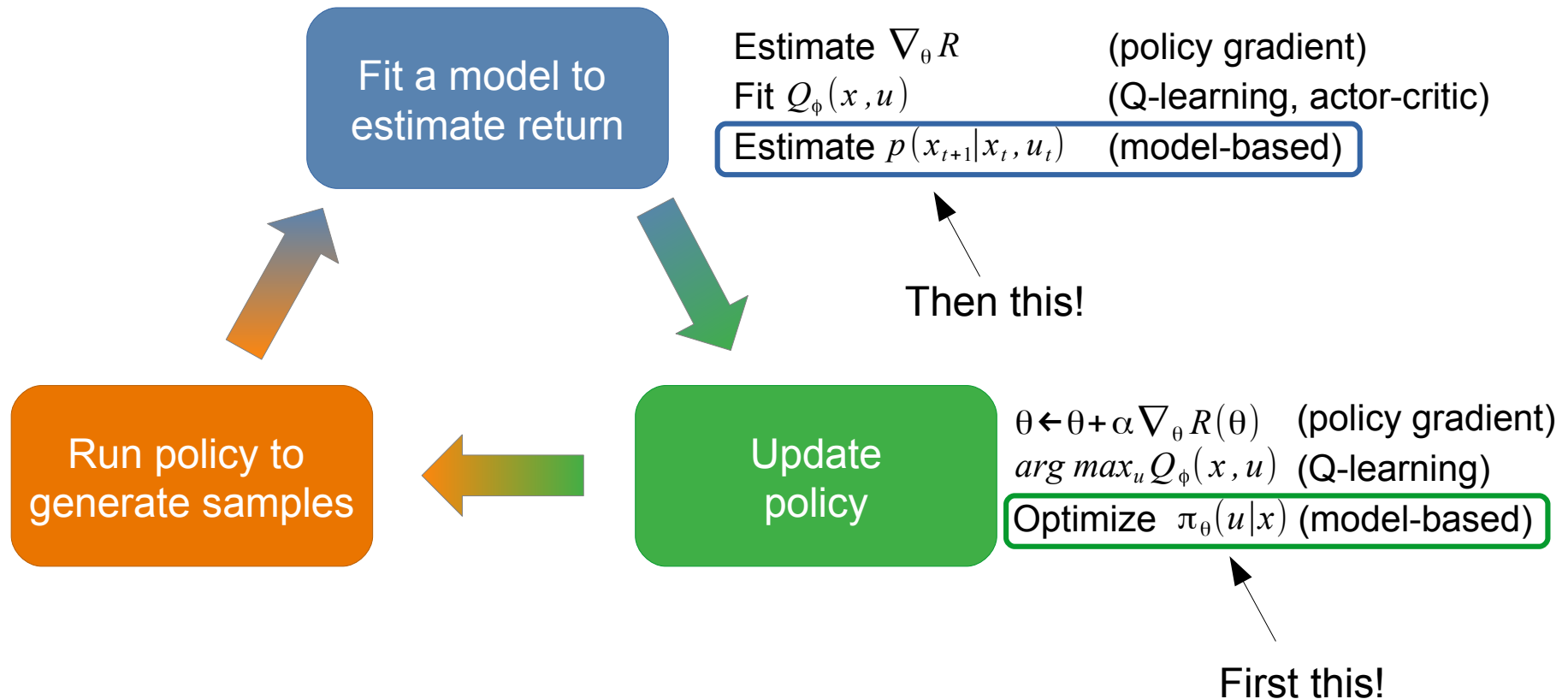
Anatomy of reinforcement learning

Model-based



Anatomy of reinforcement learning

Model-based



Solving (deterministic, finite-horizon) optimal control problems

$$\min_{u_1, \dots, u_T} \sum_t c(\mathbf{x}_t, \mathbf{u}_t) \quad s.t. \quad \mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t)$$

$$\min_{u_1, \dots, u_T} c(\mathbf{x}_1, \mathbf{u}_1) + c(f(\mathbf{x}_1, \mathbf{u}_1), \mathbf{u}_2) + \dots + c(f(f(\dots)), \mathbf{u}_T)$$

How to solve these?

Shooting vs collocation

Shooting methods: Optimize actions

$$\min_{u_1, \dots, u_T} c(\mathbf{x}_1, \mathbf{u}_1) + c(f(\mathbf{x}_1, \mathbf{u}_1), \mathbf{u}_2) + \dots + c(f(f(\dots)), \mathbf{u}_T)$$

Collocation methods: Optimize actions and states
(constrained optimization)

$$\min_{u_1, \dots, u_T} \sum_t c(\mathbf{x}_t, \mathbf{u}_t) \quad s.t. \quad \mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t)$$

LQR (linear-quadratic regulator)

Problem definition (finite horizon)

$$\min_{u_1, \dots, u_T} c(\mathbf{x}_1, \mathbf{u}_1) + c(f(\mathbf{x}_1, \mathbf{u}_1), \mathbf{u}_2) + \dots + c(f(f(\dots)), \mathbf{u}_T)$$

$$f(\mathbf{x}_t, \mathbf{u}_t) = \begin{pmatrix} A_t & 0 \\ 0 & B_t \end{pmatrix} \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix} + \mathbf{f}_t = F_t \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix} + \mathbf{f}_t$$

$$c_t(\mathbf{x}_t, \mathbf{u}_t) = \frac{1}{2} \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{C}_t \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix} + \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{c}_t$$

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.

$$\mathbf{C}_t = \begin{pmatrix} \mathbf{C}_{x_t, x_t} & \mathbf{C}_{x_t, u_t} \\ \mathbf{C}_{u_t, x_t} & \mathbf{C}_{u_t, u_t} \end{pmatrix} \quad \mathbf{c}_t = \begin{pmatrix} \mathbf{c}_{x_t} \\ \mathbf{c}_{u_t} \end{pmatrix}$$

LQR partial derivation, final step

$$\min_{u_1, \dots, u_T} c(x_1, u_1) + c(f(x_1, u_1), u_2) + \dots + c(\underbrace{f(f(\dots))}_{\text{Only cost depending on } u_T}, u_T)$$

$$f(x_t, u_t) = F_t \begin{pmatrix} x_t \\ u_t \end{pmatrix} + f_t$$

Only cost depending on u_T

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

Action-value function:

$$Q(x_T, u_T) = \text{const} + \frac{1}{2} \begin{pmatrix} x_T \\ u_T \end{pmatrix}^T C_T \begin{pmatrix} x_T \\ u_T \end{pmatrix} + \begin{pmatrix} x_T \\ u_T \end{pmatrix}^T c_T$$

$$\nabla_{u_t} Q(x_T, u_T) = C_{u_t, x_t} x_T + C_{u_t, u_t} u_t + c_{u_t} = 0$$

$$u_T = -C_{u_T, u_T}^{-1} (C_{u_t, x_t} x_t + c_{u_t})$$



$$\begin{aligned} u_T &= K_T x_T + k_T \\ K_T &= -C_{u_T, u_T}^{-1} C_{u_t, x_t} \\ k_T &= -C_{u_T, u_T}^{-1} c_{u_t} \end{aligned}$$

$$C_t = \begin{pmatrix} C_{x_t, x_t} & C_{x_t, u_t} \\ C_{u_t, x_t} & C_{u_t, u_t} \end{pmatrix} \quad c_t = \begin{pmatrix} c_{x_t} \\ c_{u_t} \end{pmatrix}$$

LQR partial derivation, final step

$$\min_{u_1, \dots, u_T} c(\mathbf{x}_1, \mathbf{u}_1) + c(f(\mathbf{x}_1, \mathbf{u}_1), \mathbf{u}_2) + \dots + c(f(f(\dots)), \mathbf{u}_T) \dot{}$$

$$\mathbf{u}_T = \mathbf{K}_T \mathbf{x}_T + \mathbf{k}_T \quad \mathbf{K}_T = -\mathbf{C}_{u_T, u_T}^{-1} \mathbf{C}_{u_T, x_T} \quad \mathbf{k}_T = -\mathbf{C}_{u_T, u_T}^{-1} \mathbf{c}_{u_T}$$

State-value function (by substitution):

$$V(\mathbf{x}_T) = \text{const} + \frac{1}{2} \begin{pmatrix} \mathbf{x}_T \\ \mathbf{K}_T \mathbf{x}_T + \mathbf{k}_T \end{pmatrix}^T \mathbf{C}_T \begin{pmatrix} \mathbf{x}_T \\ \mathbf{K}_T \mathbf{x}_T + \mathbf{k}_T \end{pmatrix} + \begin{pmatrix} \mathbf{x}_T \\ \mathbf{K}_T \mathbf{x}_T + \mathbf{k}_T \end{pmatrix}^T \mathbf{c}_T$$

State value function is quadratic in \mathbf{x}_T !

$$V(\mathbf{x}_T) = \text{const} + \frac{1}{2} \mathbf{x}_T^T \mathbf{V}_T \mathbf{x}_T + \mathbf{x}_T^T \mathbf{v}_T$$

$$V(\mathbf{x}_T) = \text{const} + \frac{1}{2} \mathbf{x}_T^T V_T \mathbf{x}_T + \mathbf{x}_T^T \mathbf{v}_T$$

LQR partial derivation, other steps

$$\begin{aligned}
 Q(\mathbf{x}_t, \mathbf{u}_t) &= \text{const} + \overbrace{\frac{1}{2} \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{C}_t \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}}^{\text{quadratic}} + \overbrace{\begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{c}_t + V(f(\mathbf{x}_t, \mathbf{u}_t))}_{\text{quadratic}} \\
 &= \text{const} + \frac{1}{2} \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{Q}_t \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix} + \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{q}_t
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{Q}_t &= \mathbf{C}_t + \mathbf{F}_t^T V_{t+1} \mathbf{F}_t \\
 \mathbf{q}_t &= \mathbf{c}_t + \mathbf{F}_t^T V_{t+1} \mathbf{f}_t + \mathbf{F}_t^T \mathbf{v}_{t+1}
 \end{aligned}$$

Note: We skip here the derivation of V_t, \mathbf{v}_t

$$V(\mathbf{x}_T) = \text{const} + \frac{1}{2} \mathbf{x}_T^T V_T \mathbf{x}_T + \mathbf{x}_T^T \mathbf{v}_T$$

LQR partial derivation, other steps

$$\begin{aligned}
 Q(\mathbf{x}_t, \mathbf{u}_t) &= \text{const} + \overbrace{\frac{1}{2} \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{C}_t \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}}^{\text{quadratic}} + \overbrace{\begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{c}_t + V(f(\mathbf{x}_t, \mathbf{u}_t))}_{\text{quadratic}} \\
 &= \text{const} + \frac{1}{2} \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{Q}_t \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix} + \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{q}_t
 \end{aligned}$$

$$\mathbf{Q}_t = \mathbf{C}_t + \mathbf{F}_t^T V_{t+1} \mathbf{F}_t$$

$$\mathbf{q}_t = \mathbf{c}_t + \mathbf{F}_t^T V_{t+1} \mathbf{f}_t + \mathbf{F}_t^T \mathbf{v}_{t+1}$$

$$\nabla_{\mathbf{u}_t} Q(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{Q}_{\mathbf{u}_t, \mathbf{x}_t} \mathbf{x}_t + \mathbf{Q}_{\mathbf{u}_t, \mathbf{u}_t} \mathbf{u}_t + \mathbf{q}_t^T = 0$$

$$\mathbf{u}_t = \mathbf{K}_t \mathbf{x}_t + \mathbf{k}_t \quad \mathbf{K}_t = -\mathbf{Q}_{\mathbf{u}_t, \mathbf{u}_t}^{-1} \mathbf{Q}_{\mathbf{u}_t, \mathbf{x}_t} \quad \mathbf{k}_t = -\mathbf{Q}_{\mathbf{u}_t, \mathbf{u}_t}^{-1} \mathbf{q}_{\mathbf{u}_t}$$

LQR algorithm

Backward recursion:

For $t = T$ downto 1

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

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

$$K_t = -Q_{u_t, u_t}^{-1} Q_{u_t, x_t}$$

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

$$V_t = Q_{x_t, x_t} + Q_{x_t, u_t} K_t + K_t^T Q_{u_t, x_t} + K_t^T Q_{u_t, u_t} K_t$$

$$v_t = q_{x_t} + Q_{x_t, u_t} k_t + K_t^T q_{u_t} + K_t^T Q_{u_t, u_t} k_t$$



First: compute the gains.

Forward recursion:

For $t = 1$ to T

$$u_t = K_t x_t + k_t$$

$$x_{t+1} = f(x_t, u_t)$$



Then: apply the law to compute controls.

System uncertainty / stochastic dynamics

Gaussian noise

$$f(\mathbf{x}_t, \mathbf{u}_t) = \mathbf{F}_t \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix} + \mathbf{f}_t + \mathbf{w}_t \quad \mathbf{w}_t \sim N(\mathbf{0}, \Sigma_t)$$

$$p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) \sim N \left(\mathbf{F}_t \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix} + \mathbf{f}_t, \Sigma_t \right)$$

- 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

$$f(\mathbf{x}_t, \mathbf{u}_t) = F_t \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}$$

$$c_t(\mathbf{x}_t, \mathbf{u}_t) = \frac{1}{2} \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{C}_t \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix} + \begin{pmatrix} \mathbf{x}_t \\ \mathbf{u}_t \end{pmatrix}^T \mathbf{c}_t$$

$$f(\mathbf{x}_t, \mathbf{u}_t) \approx f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) + \nabla_{x_t, u_t} f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{pmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{pmatrix}$$

$$c_t(\mathbf{x}_t, \mathbf{u}_t) = c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) + \frac{1}{2} \begin{pmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{pmatrix}^T \nabla_{x_t, u_t}^2 c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{pmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{pmatrix} + \nabla_{x_t, u_t} c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{pmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{pmatrix}$$

Non-linear systems - Iterative LQR

$$f(\mathbf{x}_t, \mathbf{u}_t) \approx f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) + \nabla_{\mathbf{x}_t, \mathbf{u}_t} f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{pmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{pmatrix}$$

$$c_t(\mathbf{x}_t, \mathbf{u}_t) = c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) + \frac{1}{2} \begin{pmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{pmatrix}^T \nabla_{\mathbf{x}_t, \mathbf{u}_t}^2 c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{pmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{pmatrix} + \nabla_{\mathbf{x}_t, \mathbf{u}_t} c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t) \begin{pmatrix} \mathbf{x}_t - \hat{\mathbf{x}}_t \\ \mathbf{u}_t - \hat{\mathbf{u}}_t \end{pmatrix}$$

$$\bar{f}(\delta \mathbf{x}_t, \delta \mathbf{u}_t) = \mathbf{F}_t \begin{pmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{pmatrix}$$

$$\nabla_{\mathbf{x}_t, \mathbf{u}_t} f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$$

$$\bar{c}_t(\delta \mathbf{x}_t, \delta \mathbf{u}_t) = \frac{1}{2} \begin{pmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{pmatrix}^T \mathbf{C}_t \begin{pmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{pmatrix} + \begin{pmatrix} \delta \mathbf{x}_t \\ \delta \mathbf{u}_t \end{pmatrix}^T \mathbf{c}_t$$

$$\nabla_{\mathbf{x}_t, \mathbf{u}_t}^2 c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$$

$$\nabla_{\mathbf{x}_t, \mathbf{u}_t} c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$$

Iterative LQR (iLQR) – Algorithm outline

Repeat

$$\mathbf{F}_t = \nabla_{x_t, u_t} f(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$$

$$\mathbf{C}_t = \nabla_{x_t, u_t}^2 c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$$

$$\mathbf{c}_t = \nabla_{x_t, u_t} c(\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t)$$

Run LQR backward pass with $\delta \mathbf{x}_t, \delta \mathbf{u}_t$

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

Update $\hat{\mathbf{x}}_t, \hat{\mathbf{u}}_t$ to results of forward pass

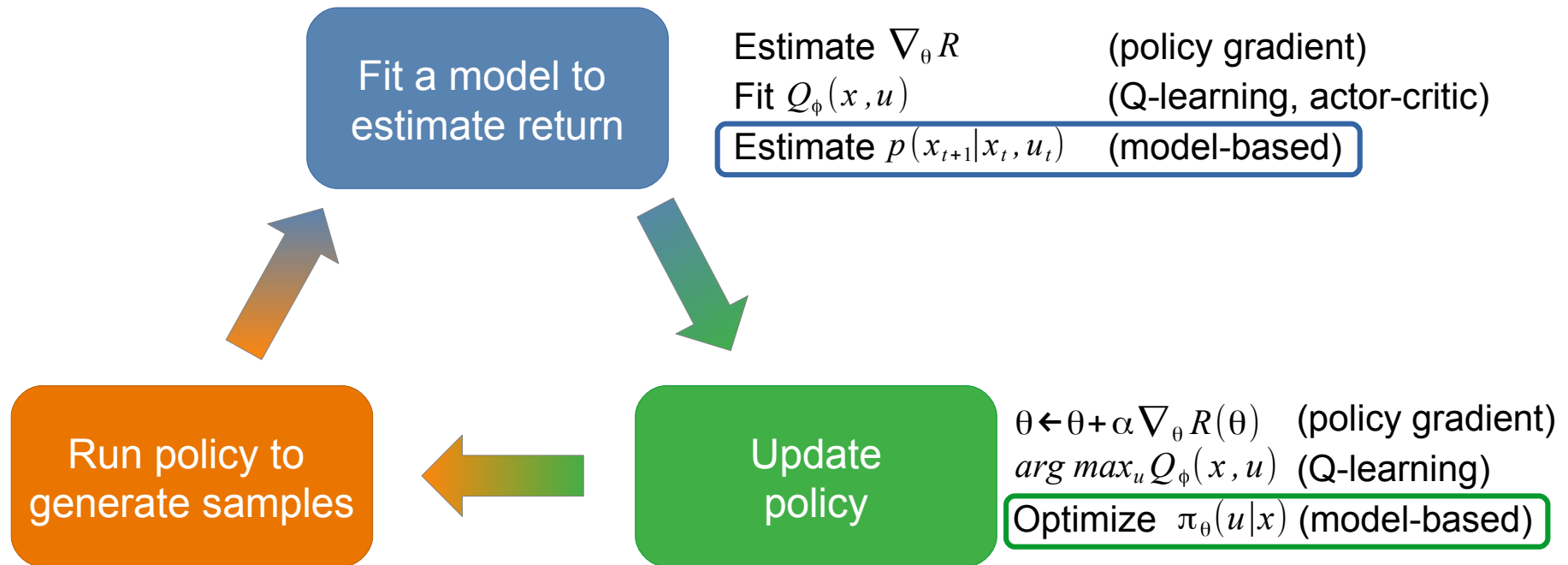
until convergence

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.

Anatomy of reinforcement learning

Model-based



Basic iterative model-based RL

Input: base policy π_0

Run base policy to collect data $D \leftarrow \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}$

Repeat

Fit dynamics model $f(\mathbf{x}, \mathbf{u})$ to minimize $\sum_i \|f(\mathbf{x}_i, \mathbf{u}_i) - \mathbf{x}'_i\|^2$

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

Execute first planned action, observe resulting state \mathbf{x}'

Update dataset $D \leftarrow D \cup \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')\}$

Basic iterative model-based RL

Input: base policy π_0

Run base policy to collect data $D \leftarrow \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}$

Repeat

Fit dynamics model $f(\mathbf{x}, \mathbf{u})$ to minimize $\sum_i \|f(\mathbf{x}_i, \mathbf{u}_i) - \mathbf{x}'_i\|^2$

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

Execute first planned action, observe resulting state \mathbf{x}'

Update dataset $D \leftarrow D \cup \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')\}$

Notes:

- Sample efficient.
- Computationally expensive for two reasons.
 - Dynamics fitting costly: model may be fitted only periodically (every n steps).
 - Planning costly for long horizons.
- Robust to moderate model errors.
- Choice of regression model is an important design parameter.

Basic iterative model-based RL

Input: base policy π_0

Run base policy to collect data $D \leftarrow \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}$

Repeat

Fit dynamics model $f(\mathbf{x}, \mathbf{u})$ to minimize $\sum_i \|f(\mathbf{x}_i, \mathbf{u}_i) - \mathbf{x}'_i\|^2$

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

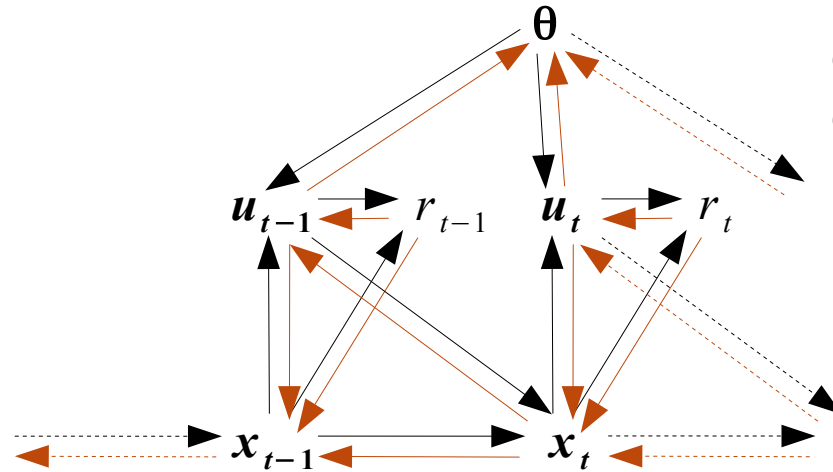
Execute first planned action, observe resulting state \mathbf{x}'

Update dataset $D \leftarrow D \cup \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')\}$

Notes:

- Sample efficient.
- Computationally expensive for two reasons.
 - Dynamics fitting costly: model may be fitted only periodically (every n steps).
 - **Planning costly for long horizons.**
- Robust to moderate model errors.
- Choice of regression model is an important design parameter.

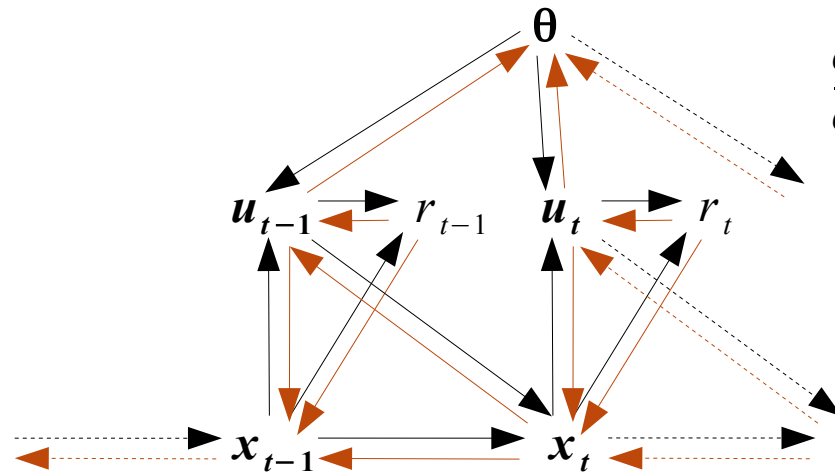
Combining parametric policy with learned dynamics by backpropagation



$$\frac{\partial r_t}{\partial \theta} = \frac{\partial r_t}{\partial u_t} \left(\frac{\partial u_t}{\partial \theta} + \frac{\partial u_t}{\partial x_t} \right) + \frac{\partial r_t}{\partial x_t} \left(\frac{\partial x_t}{\partial x_{t-1}} + \frac{\partial x_t}{\partial u_{t-1}} \right)$$

$$= \dots$$

Combining parametric policy with learned dynamics by backpropagation



$$\frac{\partial r_t}{\partial \theta} = \frac{\partial r_t}{\partial u_t} \left(\frac{\partial u_t}{\partial \theta} + \frac{\partial u_t}{\partial x_t} \right) + \frac{\partial r_t}{\partial x_t} \left(\frac{\partial x_t}{\partial x_{t-1}} + \frac{\partial x_t}{\partial u_{t-1}} \right)$$

$$= \dots$$

Run base policy to collect data $D \leftarrow \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')\}_i$

Repeat

Fit dynamics model $f_\phi(\mathbf{x}, \mathbf{u})$ to minimize $\sum_i \|f_\phi(\mathbf{x}_i, \mathbf{u}_i) - \mathbf{x}'_i\|^2$

Calculate policy gradient by backpropagating through dynamics

Execute updated policy (1 or more steps), collect data

Update dataset $D \leftarrow D \cup \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')\}$

Basic iterative model-based RL

Input: base policy π_0

Run base policy to collect data $D \leftarrow \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}$

Repeat

Fit dynamics model $f(\mathbf{x}, \mathbf{u})$ to minimize $\sum_i \|f(\mathbf{x}_i, \mathbf{u}_i) - \mathbf{x}'_i\|^2$

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

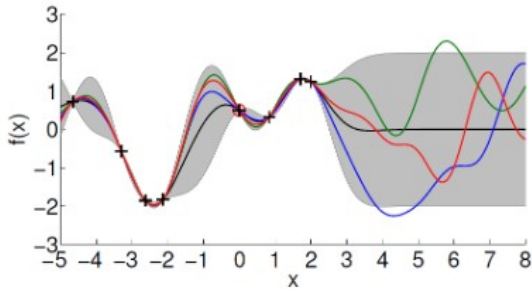
Execute first planned action, observe resulting state \mathbf{x}'

Update dataset $D \leftarrow D \cup \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')\}$

Notes:

- Sample efficient.
- Computationally expensive for two reasons.
 - Dynamics fitting costly: model may be fitted only periodically (every n steps).
 - Planning costly for long horizons.
- Robust to moderate model errors.
- **Choice of regression model is an important design parameter.**

Which model to use for dynamics?



Gaussian process (GP)

- Data-efficient
- Slow with big datasets
- May be too smooth for non-smooth dynamics

Neural networks

- Expressive
- Unpredictable with sparse data (overfit)

Linear models

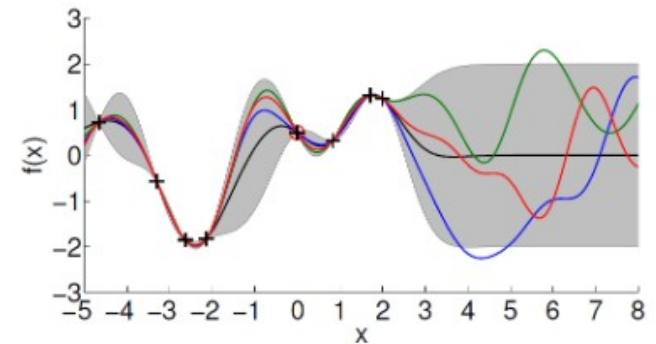
- May be used locally
- Do not overfit

Domain specific parametric models (e.g. physics parameters) can also be used.
→ Traditional control engineering approach of model identification + control.

Example

PILCO (Deisenroth&Rasmussen, 2011)

- Dynamics learning: Use Gaussian process models to include model uncertainty.
- Simulation: Simulate trajectory with learned model, including uncertainty.
- Policy: Radial basis function.
- Policy update: Calculate analytically policy gradient using learned dynamics and optimize with quasi-Newton optimizer (BFGS).



What about unknown rewards?

- Idea: Learn also regression function for rewards.
- BlackDROPS (2017) uses a Gaussian process to model reward function as well as dynamics.
- Uses CMA-ES (gradient free optimizer) for planning.

Summary

- Model-based RL requires typically less data than value-based or policy gradient approaches.
- Learned dynamics can be transferred across tasks.
- Potentially suboptimal: models do not optimize for task performance.
- Sometimes models are harder to learn than policy.
- Often require explicit choices (e.g. time horizon).

Next: Partial observability and POMDPs

- What changes if we cannot observe state directly?
- Reading: Tony Cassandra's on-line tutorial (see MyCourses for details)