

to the repeated application of the chain-rule, the computation of the gradient itself is often mathematically more involved than a sampling-based estimate.

Let us consider an example where the immediate reward r only depends on the state (generalizations to control-dependent rewards are straightforward) and the system dynamics are deterministic, such that $\mathbf{x}_{t+1} = f(\mathbf{x}_t, \mathbf{u}_t) = f(\mathbf{x}_t, \pi_{\boldsymbol{\theta}}(\mathbf{x}_t, \boldsymbol{\theta}))$, where f is a (nonlinear) transition function, $\pi_{\boldsymbol{\theta}}$ is the (deterministic) policy, and $\boldsymbol{\theta}$ are the policy parameters. The gradient of the long-term reward $J_{\boldsymbol{\theta}} = \sum_t \gamma^t r(\mathbf{x}_t)$ with respect to the policy parameters is obtained by applying the chain-rule repeatedly:

$$\frac{dJ_{\boldsymbol{\theta}}}{d\boldsymbol{\theta}} = \sum_t \gamma^t \frac{dr(\mathbf{x}_t)}{d\boldsymbol{\theta}} = \sum_t \gamma^t \frac{\partial r(\mathbf{x}_t)}{\partial \mathbf{x}_t} \frac{d\mathbf{x}_t}{d\boldsymbol{\theta}} \quad (3.20)$$

$$= \sum_t \gamma^t \frac{\partial r(\mathbf{x}_t)}{\partial \mathbf{x}_t} \left(\frac{\partial \mathbf{x}_t}{\partial \mathbf{x}_{t-1}} \frac{d\mathbf{x}_{t-1}}{d\boldsymbol{\theta}} + \frac{\partial \mathbf{x}_t}{\partial \mathbf{u}_{t-1}} \frac{d\mathbf{u}_{t-1}}{d\boldsymbol{\theta}} \right). \quad (3.21)$$

From these equations we observe that the total derivative $d\mathbf{x}_t/d\boldsymbol{\theta}$ depends on the total derivative $d\mathbf{x}_{t-1}/d\boldsymbol{\theta}$ at the previous time step. Therefore, the derivative $dJ_{\boldsymbol{\theta}}/d\boldsymbol{\theta}$ can be computed iteratively.

Extension to Probabilistic Models and Stochastic MDPs. For the extension to derivatives in stochastic MDPs and/or probabilistic models, we have to make a few adaptations to the gradients in Equation (3.20)–(3.21): When the state \mathbf{x}_t is represented by a probability distribution $p(\mathbf{x}_t)$, we have to compute the *expected* reward $\mathbb{E}[r(\mathbf{x}_t)] = \int r(\mathbf{x}_t)p(\mathbf{x}_t) d\mathbf{x}_t$. Moreover, we need to compute the derivatives with respect to the parameters of the state distribution, assuming that $p(\mathbf{x}_t)$ has a parametric representation.

For example, if $p(\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_t^x, \boldsymbol{\Sigma}_t^x)$, we compute the derivatives of $\mathbb{E}[r(\mathbf{x}_t)]$ with respect to the mean $\boldsymbol{\mu}_t^x$ and covariance $\boldsymbol{\Sigma}_t^x$ of the state distribution and continue applying the chain-rule similarly to Equation (3.20)–(3.21): With the definition $\mathcal{E}_t := \mathbb{E}_{\mathbf{x}_t}[r(\mathbf{x}_t)]$, we obtain the