# Function Approximation
## Control using VFA

Marius Lindauer

**tnt**

Leibniz
Universität
Hannover

Automated
Machine Learning
Hannover

# Control using Value Function Approximation

- Use value function approximation to represent state-action values $\hat{Q}^{\pi}(s, a; \mathbf{w}) \approx Q^{\pi}$
- Interleave
  - Approximate policy evaluation using value function approximation
  - Perform $\epsilon$-greedy policy improvement
- Can be unstable. Generally involves intersection of the following:
  - Function approximation
  - Bootstrapping
  - Off-policy learning

# Action-Value Function Approximation with an Oracle

- $\hat{Q}^\pi(s, a; \mathbf{w}) \approx Q^\pi$
- Minimize the mean-squared error between the true action-value function $Q^\pi(s, a)$ and the approximate action-value function:

$$J(\mathbf{w}) = \mathbb{E}_\pi[(Q^\pi(s, a) - \hat{Q}^\pi(s, a; \mathbf{w}))^2]$$

- Use stochastic gradient descent to find a local minimum

$$
\begin{aligned}
-\frac{1}{2}\nabla_\mathbf{w} J(\mathbf{w}) &= \mathbb{E}\left[(Q^\pi(s, a) - \hat{Q}^\pi(s, a; \mathbf{w}))\nabla_\mathbf{w}\hat{Q}^\pi(s, a; \mathbf{w})\right] \\
\Delta\mathbf{w} &= -\frac{1}{2}\alpha\nabla_\mathbf{w} J(\mathbf{w})
\end{aligned}
$$

- Stochastic gradient descent (SGD) samples the gradient

# Linear State Action Value Function Approximation with an Oracle

- Use features to represent both the state and action

$$\mathbf{x}(s,a) = \begin{pmatrix} \mathbf{x}_1(s,a) \\ \mathbf{x}_2(s,a) \\ \dots \\ \mathbf{x}_n(s,a) \end{pmatrix}$$

- Represent state-action value function with a weighted linear combination of features

$$\hat{Q}(s,a;\mathbf{w}) = \mathbf{x}(s,a)^T \mathbf{w} = \sum_{j=1}^{n} x_j(s,a)w_j$$

- Stochastic gradient descent update

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \nabla_{\mathbf{w}} \mathbb{E}_\pi [(Q^\pi(s,a) - \hat{Q}^\pi(s,a;\mathbf{w}))^2]$$

# Incremental Model-Free Control Approaches

- Similar to policy evaluation, true state-action value function for a state is unknown and so substitute a target value
- In Monte Carlo methods, use a return $G_t$ as a substitute target

$$\Delta \mathbf{w} = \alpha(G_t - \hat{Q}(s_t, a_t; \mathbf{w}))\nabla_{\mathbf{w}}\hat{Q}(s_t, a_t; \mathbf{w})$$

- For SARSA instead use a TD target $r + \gamma\hat{Q}(s', a'; \mathbf{w})$ which leverages the current function approximations value

$$\Delta \mathbf{w} = \alpha(r + \gamma\hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w}))\nabla_{\mathbf{w}}\hat{Q}(s, a; \mathbf{w})$$

- For Q-learning instead use a TD target $r + \gamma\max_{a'}\hat{Q}(s', a'; \mathbf{w})$ which leverages the max of the current function approximations value

$$\Delta \mathbf{w} = \alpha(r + \gamma\max_{a'}\hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w}))\nabla_{\mathbf{w}}\hat{Q}(s, a; \mathbf{w})$$