

Topics in machine learning

Siddharth Bhat

Monsoon 2019

Contents

1	Policy iteration
---	------------------

5

Chapter 1

Policy iteration

$$\pi_{k+1}(s) = \arg \max_{a \in A(s)} r(s, a) + \gamma \sum_s P(s'|s, a) v_{\pi_k}(s')$$

Theorem 1 *The policy iteration algorithm generates a sequence of policies with non-decreasing state values. That is, $V^{\pi_{k+1}} \geq V^{\pi_k}$, $V^\pi \in \mathbb{R}^n$, is the vector of state values for state π*

Proof 1 F^{π_k} is the bellman expectation operator (?)

Since V^{π_k} is a fixed point of F^{π_k} ,

$$V^{\pi_k} = F^{\pi_k}(V^{\pi_k}) \leq F(V^{\pi_k}) \quad (\text{upper bounded by max value})$$

$$F(V^{\pi_k}) = F^{\pi_{k+1}}(V^{\pi_k}) \quad (\text{By defn of policy improvement step})$$

$$V^{\pi_k} \leq F^{\pi_{k+1}}(V^{\pi_k}) \quad (\text{eqn 1})$$

$$F^{\pi_{k+1}}(V^{\pi_k}) \leq (F^{\pi_{k+1}})^2(V^{\pi_k}) \quad (\text{Monotonicity of } F^{\pi_{k+1}})$$

$$\forall t \geq 1, F^{\pi_{k+1}}(V^{\pi_k}) \leq (F^{\pi_{k+1}})^t(V^{\pi_k}) \quad (\text{Monotonicity of } F^{\pi_{k+1}})$$

$$F^{\pi_{k+1}}(V^{\pi_k}) \leq (F^{\pi_{k+1}})^t(V^{\pi_k}) \leq V^{\pi_{k+1}} \quad (\text{Contraction mapping, } V^{\pi_{k+1}} \text{ is fixed point})$$