

# Class Tutorial 11

---

## 1. Projection Operator

Consider the projection operator  $\Pi$  that projects a vector  $V \in \mathbb{R}^n$  to a linear subspace  $S$  that is the span of  $k$  features  $\phi_1(s), \dots, \phi_k(s)$  w.r.t. the  $d$ -weighted Euclidean norm:

$$\|X\|_d^2 = \sum_{j=1}^n d(j) X(j)^2$$

- Write down  $\Pi$  explicitly.
- Show that a projection is **non-expansive**:

$$\| \Pi V - \Pi \tilde{V} \|_d \leq \| V - \tilde{V} \|_d$$

### Solution

a. Let  $\Phi \in \mathbb{R}^{n \times k}$  denote a matrix with  $\phi(s)$  in its rows. Note that every vector in  $S$  can be written as  $\Phi w$  for some  $w$ .

For  $V \in \mathbb{R}^n$ , by definition, we have

$$\begin{aligned} \Pi V &= \Phi w^* \\ w^* &= \arg \min_{w \in \mathbb{R}^k} \|\Phi w - V\|_d^2 = \arg \min_{w \in \mathbb{R}^k} (\Phi w - V)^\top D (\Phi w - V) \end{aligned}$$

where  $D = \text{diag}(d)$ .

Taking a gradient and setting to zero gives:

$$\begin{aligned} \frac{\partial}{\partial w} (\Phi w - V)^\top D (\Phi w - V) &= 2\Phi^\top D \Phi w - 2\Phi^\top D V \\ w^* &= (\Phi^\top D \Phi)^{-1} \Phi^\top D V \end{aligned}$$

and

$$\Pi = \Phi (\Phi^\top D \Phi)^{-1} \Phi^\top D.$$

b. The error  $V - \Pi V$  is orthogonal to  $\Pi V$  (property of a projection to linear subspaces, see Random Signals course), therefore we have the Pythagorean Theorem:

$$\|X\|_d^2 = \|\Pi X\|_d^2 + \|X - \Pi X\|_d^2$$

Now, using the linearity of  $\Pi$  :

$$\| \Pi V - \Pi \tilde{V} \|_d = \| \Pi(V - \tilde{V}) \|_d = \| (I - \Pi)(V - \tilde{V}) \|_d = \| V - \tilde{V} \|_d$$

## 2. Projected Bellman Operator

Consider an MDP and some policy  $\pi$  and let denote  $P$  the state transition probability

matrix under  $\pi$ . Assume that  $P$  has a unique **stationary distribution**  $d$ , i.e.,  $\sum_{i=1}^n d_i P_{ij} = d_j$ .

Also, let  $T^\pi(v) \doteq r + \gamma P v$  denote the Bellman operator for policy  $\pi$ .

a. Show that  $\| P z \|_d^2 \leq \| z \|_d^2$ .

b. Show that  $\Pi T^\pi$  is a contraction in the  $\| \cdot \|_d$  norm.

## Solution

a.

$$\begin{aligned} \| P z \|_d^2 &= \sum_{i=1}^n d_i \left( \sum_{j=1}^n P_{ij} z_j \right)^2 \\ &\stackrel{\text{Jensen}}{\leq} \sum_{i=1}^n d_i \sum_{j=1}^n P_{ij} z_j^2 \\ &= \sum_{j=1}^n z_j^2 \sum_{i=1}^n d_i P_{ij} \\ &= \sum_{j=1}^n d_j z_j^2 \\ &\leq \| z \|_d^2 \end{aligned}$$

b. We have that  $T^\pi$  is a contraction in the  $\| \cdot \|_d$  norm:

$$\| T^\pi V - T^\pi \tilde{V} \|_d = \| \gamma P (V - \tilde{V}) \|_d \leq \gamma \| V - \tilde{V} \|_d$$

Since  $\Pi$  is non-expansive in the  $\| \cdot \|_d$  norm,  $\Pi T^\pi$  is a contraction.

## 3. TD(0)

Consider the following problem: given a sequence of states  $s_1, \dots, s_n$  and rewards  $r_1, \dots, r_n$  generated by an MDP with policy  $\pi$ , our goal is to estimate the value function  $V^\pi(s)$ . We approximate  $V^\pi(s)$  using linear function approximation, i.e.,

$$\tilde{V}^\pi(s) = \phi(s)^\top w$$

where  $\phi(s)$  denotes the features of state  $s$  and  $w$  are the approximation weights. Thus, our goal is to find  $w$ .

The **TD(0) algorithm** computes  $w$  iteratively as follows:

$$\begin{aligned} w_{t+1} &= w_t + \alpha_t \delta_t \phi(s_t) \\ \delta_t &\triangleq r(s_t) + \gamma \phi(s_{t+1})^\top w_t - \phi(s_t)^\top w_t \end{aligned} \quad (1)$$

We will now prove that (a simplified version of) TD(0) converges. Our approach is to represent TD(0) as a stochastic approximation of the form

$$w_{t+1} = w_t + \alpha_t (h(w_t) + m_t)$$

where  $h$  is a deterministic function and  $m_t$  are independent with  $\mathbb{E}[m_t | w_t] = 0$ .

Assume that the MDP under policy  $\pi$  has a unique **stationary distribution**  $d$ .

a. Let  $P_n(s) \triangleq \Pr(s_n = s)$  and  $P_n(s, s') \triangleq \Pr(s_n = s, s_{n+1} = s')$ . Calculate  $P_\infty(s) \triangleq \lim_{n \rightarrow \infty} P_n(s)$  and  $P_\infty(s, s') \triangleq \lim_{n \rightarrow \infty} P_n(s, s')$ .

We now make the following simplifying assumption: in the TD(0) algorithm (1) each pair  $s_t, s_{t+1}$  is drawn i.i.d. from  $P_\infty(s, s')$ .

b. Calculate  $\mathbb{E}[\delta_t \phi(s_t) | w_t]$ , and show that it may be written as a linear function  $A w_t + b$ .

c. Write the TD(0) algorithm as a stochastic approximation. What is the resulting ODE?

d. Show that the equilibrium point of the ODE corresponds to the solution of the projected fixed point equation  $V = \Pi T^\pi V$ . Conclude that the solution is unique.

e. Show that  $z^\top \gamma D P z \leq z^\top D z \quad \forall z \neq 0$ . Use the result to show that  $D(\gamma P - I)$  is negative definite in the sense that  $z^\top D(\gamma P - I) z < 0 \quad \forall z \neq 0$ . Show that if  $\Phi$  is full rank then  $A$  is also negative definite.

It can be shown that the negative definiteness of  $A$  means that the eigenvalues of  $A$  have a negative real part.

f. Using the stochastic approximation theorem learned in class, state a convergence theorem for TD(0).

g. Why did we need to assume that each pair  $s_t, s_{t+1}$  is drawn i.i.d. from  $P_\infty(s, s')$ ?

## Solution

a. We have  $\lim_{n \rightarrow \infty} P_n(s) = d(s)$  and  $\lim_{n \rightarrow \infty} P_n(s, s') = d(s) P_{s, s'}$ .

b. We have

$$\begin{aligned}\mathbb{E}[\delta_t \phi(s_t) | w_t] &= \mathbb{E}\left[\left(r(s_t) + \gamma \phi(s_{t+1})^\top w_t - \phi(s_t)^\top w_t\right) \phi(s_t) | w_t\right] \\ &= \mathbb{E}[r(s_t) \phi(s_t) | w_t] + \gamma \mathbb{E}[\phi(s_t) \phi(s_{t+1})^\top w_t | w_t] - \mathbb{E}[\phi(s_t) \phi(s_t)^\top w_t | w_t] \\ &= \mathbb{E}[r(s_t) \phi(s_t)] + \gamma w_t \mathbb{E}[\phi(s_t) \phi(s_{t+1})^\top] - w_t \mathbb{E}[\phi(s_t) \phi(s_t)^\top]\end{aligned}$$

Observe that

$$\mathbb{E}[r(s_t) \phi(s_t)] = \sum_s d(s) r(s) \phi(s) = \Phi^\top D r$$

where  $\Phi$  has  $\phi(s)$  in its rows, and  $D = \text{diag}(d)$ . Similarly, we have

$$\mathbb{E}[\phi(s_t) \phi(s_t)^\top] = \sum_s d(s) \phi(s) \phi(s)^\top = \Phi^\top D \Phi$$

$$\mathbb{E}[\phi(s_t) \phi(s_{t+1})^\top] = \sum_{s, s'} d(s) P_{s, s'} \phi(s) \phi(s')^\top = \Phi^\top D P \Phi$$

Thus, we may write

$$\mathbb{E}[\delta_t \phi(s_t) | w_t] = A w_t + b \text{ where } A = \Phi^\top D (\gamma P - I) \Phi \text{ and } b = \Phi^\top D r.$$

c. We write

$$w_{t+1} = w_t + \alpha_t \left( \mathbb{E}[\delta_t \phi(s_t) | w_t] + \delta_t \phi(s_t) - \mathbb{E}[\delta_t \phi(s_t) | w_t] \right),$$

Thus

$$h(w_t) = A w_t + b$$

And

$$m_t = \delta_t \phi(s_t) - \mathbb{E}[\delta_t \phi(s_t) | w_t].$$

The ODE is therefore

$$\dot{w} = A w + b.$$

d. Let us write  $V = \Pi T^\pi V$  explicitly. We first have  $T^\pi V = r + \gamma P V$ . We now write the equation explicitly as

$$\Phi w = \Phi \left( \Phi^\top D \Phi \right)^{-1} \Phi^\top D (r + \gamma P \Phi w) \quad (2)$$

The equilibrium point of the ODE satisfies

$$\Phi^\top D(\gamma P - I)\Phi w - \Phi^\top Dr = 0 \quad (3)$$

Multiplying Eq. (3) by  $\Phi(\Phi^\top D\Phi)^{-1}$  gives Eq. (2). Since  $\Pi T^\pi$  is a contraction the solution is unique.

e. Recall that  $\|Pz\|_d \leq \|z\|_d$ . We have

$$\begin{aligned} z^\top \gamma DPz &= \gamma z^\top D^{1/2} D^{1/2} Pz \\ &\stackrel{c.s.}{\leq} \gamma \|D^{1/2} z\| \|D^{1/2} Pz\| \\ &= \gamma \|z\|_d \|Pz\|_d \\ &\leq \gamma \|z\|_d \|z\|_d = \gamma z^\top Dz \end{aligned}$$

Thus

$$\begin{aligned} z^\top D(\gamma P - I)z &\leq \gamma z^\top Dz - z^\top Dz \\ &= (\gamma - 1) z^\top Dz < 0. \end{aligned} \quad (4)$$

Finally, assume by negation that  $z^\top \Phi^\top D(\gamma P - I)\Phi z \geq 0$  for some  $z \neq 0$ . If  $\Phi$  is full-rank then  $\Phi z \neq 0$ , which contradicts (4).

f. Let the step-sizes satisfy G1. It is easy to show that the noise satisfies N1. Also, clearly  $h$  is Lipschitz, and its unique equilibrium point is stable. The boundedness of  $w_t$  is not trivial, and may be shown using more advanced techniques (beyond the course material). Alternatively, one may use a projection approach to keep the iterates in some bounded set.

Assuming  $w_t$  is bounded w.p. 1, Theorem 1 of Lecture 10 guarantees that  $w_t$  converges to the unique fixed point of  $\Pi T^\pi$ .

g. In the regular TD(0) algorithm, consecutive states are dependent. Actually, the stochastic approximation theorem may be extended to such cases as well, but it needs to be shown that the dependence is short-ranged, i.e. – the Markov chain has a fast mixing property. This is beyond the scope of this course.