

## A Relevant proofs

### A.1 Proof of Theorem 1

*Proof.* The proof is similar to that given by (Sutton et al. 2009) for TDC, but it is based on multi-time-scale stochastic approximation.

For the VMTDC algorithm, a new one-step linear TD solution is defined as:

$$0 = \mathbb{E}[(\phi - \gamma\phi' - \mathbb{E}[\phi - \gamma\phi'])\phi^\top] \mathbb{E}[\phi\phi^\top]^{-1} \mathbb{E}[(\delta - \mathbb{E}[\delta])\phi] = \mathbf{A}^\top \mathbf{C}^{-1}(-\mathbf{A}\boldsymbol{\theta} + \mathbf{b}).$$

The matrix  $\mathbf{A}^\top \mathbf{C}^{-1} \mathbf{A}$  is positive definite. Thus, the VMTD's solution is  $\boldsymbol{\theta}_{\text{VMTDC}} = \mathbf{A}^{-1} \mathbf{b}$ .

First, note that recursion (5) and (6) can be rewritten as, respectively,

$$\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \zeta_k \mathbf{x}(k),$$

$$\mathbf{u}_{k+1} \leftarrow \mathbf{u}_k + \beta_k \mathbf{y}(k),$$

where

$$\begin{aligned} \mathbf{x}(k) &= \frac{\alpha_k}{\zeta_k} [(\delta_k - \omega_k)\phi_k - \gamma\phi'_k(\phi_k^\top \mathbf{u}_k)], \\ \mathbf{y}(k) &= \frac{\zeta_k}{\beta_k} [\delta_k - \omega_k - \phi_k^\top \mathbf{u}_k]\phi_k. \end{aligned}$$

Recursion (5) can also be rewritten as

$$\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \beta_k z(k),$$

where

$$z(k) = \frac{\alpha_k}{\beta_k} [(\delta_k - \omega_k)\phi_k - \gamma\phi'_k(\phi_k^\top \mathbf{u}_k)],$$

Due to the settings of step-size schedule  $\alpha_k = o(\zeta_k)$ ,  $\zeta_k = o(\beta_k)$ ,  $\mathbf{x}(k) \rightarrow 0$ ,  $\mathbf{y}(k) \rightarrow 0$ ,  $z(k) \rightarrow 0$  almost surely as  $k \rightarrow \infty$ . That is that the increments in iteration (7) are uniformly larger than those in (6) and the increments in iteration (6) are uniformly larger than those in (5), thus (7) is the fastest recursion, (6) is the second fast recursion and (5) is the slower recursion. Along the fastest time scale, iterations of (5), (6) and (7) are associated to ODEs system as follows:

$$\dot{\boldsymbol{\theta}}(t) = 0, \tag{A-1}$$

$$\dot{\mathbf{u}}(t) = 0, \tag{A-2}$$

$$\dot{\omega}(t) = \mathbb{E}[\delta_t | \mathbf{u}(t), \boldsymbol{\theta}(t)] - \omega(t). \tag{A-3}$$

Based on the ODE (A-1) and (A-2), both  $\boldsymbol{\theta}(t) \equiv \boldsymbol{\theta}$  and  $\mathbf{u}(t) \equiv \mathbf{u}$  when viewed from the fastest timescale. By the Hirsch lemma (Hirsch 1989), it follows that  $\|\boldsymbol{\theta}_k - \boldsymbol{\theta}\| \rightarrow 0$  a.s. as  $k \rightarrow \infty$  for some  $\boldsymbol{\theta}$  that depends on the initial condition  $\boldsymbol{\theta}_0$  of recursion (5) and  $\|\mathbf{u}_k - \mathbf{u}\| \rightarrow 0$  a.s. as  $k \rightarrow \infty$  for some  $\mathbf{u}$  that depends on the initial condition  $\mathbf{u}_0$  of recursion (6). Thus, the ODE pair (A-1)-(refomegavmtdcFastest) can be written as

$$\dot{\omega}(t) = \mathbb{E}[\delta_t | \mathbf{u}, \boldsymbol{\theta}] - \omega(t). \tag{A-4}$$

Consider the function  $h(\omega) = \mathbb{E}[\delta | \boldsymbol{\theta}, \mathbf{u}] - \omega$ , i.e., the driving vector field of the ODE (A-4). It is easy to find that the function  $h$  is Lipschitz with coefficient  $-1$ . Let  $h_\infty(\cdot)$  be the function defined by  $h_\infty(\omega) = \lim_{r \rightarrow \infty} \frac{h(r\omega)}{r}$ . Then  $h_\infty(\omega) = -\omega$ , is well-defined. For (A-4),  $\omega^* = \mathbb{E}[\delta | \boldsymbol{\theta}, \mathbf{u}]$  is the unique globally asymptotically stable equilibrium. For the ODE

$$\dot{\omega}(t) = h_\infty(\omega(t)) = -\omega(t), \tag{A-5}$$

apply  $\vec{V}(\omega) = (-\omega)^\top (-\omega)/2$  as its associated strict Liapunov function. Then, the origin of (A-5) is a globally asymptotically stable equilibrium.

Consider now the recursion (7). Let  $M_{k+1} = (\delta_k - \omega_k) - \mathbb{E}[(\delta_k - \omega_k) | \mathcal{F}(k)]$ , where  $\mathcal{F}(k) = \sigma(\omega_l, \mathbf{u}_l, \boldsymbol{\theta}_l, l \leq k; \phi_s, \phi'_s, r_s, s < k)$ ,  $k \geq 1$  are the sigma fields generated by  $\omega_0, \mathbf{u}_0, \boldsymbol{\theta}_0, \omega_{l+1}, \mathbf{u}_{l+1}, \boldsymbol{\theta}_{l+1}, \phi_l, \phi'_l, 0 \leq l < k$ . It is easy to verify that  $M_{k+1}, k \geq 0$  are integrable random variables that satisfy  $\mathbb{E}[M_{k+1} | \mathcal{F}(k)] = 0, \forall k \geq 0$ . Because  $\phi_k, r_k$ , and  $\phi'_k$  have uniformly bounded second moments, it can be seen that for some constant  $c_1 > 0, \forall k \geq 0$ ,

$$\mathbb{E}[\|M_{k+1}\|^2 | \mathcal{F}(k)] \leq c_1(1 + \|\omega_k\|^2 + \|\mathbf{u}_k\|^2 + \|\boldsymbol{\theta}_k\|^2).$$

Now Assumptions (A1) and (A2) of (Borkar and Meyn 2000) are verified. Furthermore, Assumptions (TS) of (Borkar and Meyn 2000) is satisfied by our conditions on the step-size sequences  $\alpha_k, \zeta_k, \beta_k$ . Thus, by Theorem 2.2 of (Borkar and Meyn 2000) we obtain that  $\|\omega_k - \omega^*\| \rightarrow 0$  almost surely as  $k \rightarrow \infty$ .

Consider now the second time scale recursion (6). Based on the above analysis, (6) can be rewritten as

$$\dot{\boldsymbol{\theta}}(t) = 0, \tag{A-6}$$

$$\dot{\mathbf{u}}(t) = \mathbb{E}[(\delta_t - \mathbb{E}[\delta_t|\mathbf{u}(t), \boldsymbol{\theta}(t)])\boldsymbol{\phi}_t|\boldsymbol{\theta}(t)] - \mathbf{C}\mathbf{u}(t). \quad (\text{A-7})$$

The ODE (A-6) suggests that  $\boldsymbol{\theta}(t) \equiv \boldsymbol{\theta}$  (i.e., a time invariant parameter) when viewed from the second fast timescale. By the Hirsch lemma (Hirsch 1989), it follows that  $\|\boldsymbol{\theta}_k - \boldsymbol{\theta}\| \rightarrow 0$  a.s. as  $k \rightarrow \infty$  for some  $\boldsymbol{\theta}$  that depends on the initial condition  $\boldsymbol{\theta}_0$  of recursion (5).

Consider now the recursion (6). Let  $N_{k+1} = ((\delta_k - \mathbb{E}[\delta_k]) - \boldsymbol{\phi}_k \boldsymbol{\phi}_k^\top \mathbf{u}_k) - \mathbb{E}[(\delta_k - \mathbb{E}[\delta_k]) - \boldsymbol{\phi}_k \boldsymbol{\phi}_k^\top \mathbf{u}_k]|\mathcal{I}(k)]$ , where  $\mathcal{I}(k) = \sigma(\mathbf{u}_l, \boldsymbol{\theta}_l, l \leq k; \boldsymbol{\phi}_s, \boldsymbol{\phi}'_s, r_s, s < k)$ ,  $k \geq 1$  are the sigma fields generated by  $\mathbf{u}_0, \boldsymbol{\theta}_0, \mathbf{u}_{l+1}, \boldsymbol{\theta}_{l+1}, \boldsymbol{\phi}_l, \boldsymbol{\phi}'_l, 0 \leq l < k$ . It is easy to verify that  $N_{k+1}, k \geq 0$  are integrable random variables that satisfy  $\mathbb{E}[N_{k+1}|\mathcal{I}(k)] = 0, \forall k \geq 0$ . Because  $\boldsymbol{\phi}_k, r_k$ , and  $\boldsymbol{\phi}'_k$  have uniformly bounded second moments, it can be seen that for some constant  $c_2 > 0, \forall k \geq 0$ ,

$$\mathbb{E}[\|N_{k+1}\|^2|\mathcal{I}(k)] \leq c_2(1 + \|\mathbf{u}_k\|^2 + \|\boldsymbol{\theta}_k\|^2).$$

Because  $\boldsymbol{\theta}(t) \equiv \boldsymbol{\theta}$  from (A-6), the ODE pair (A-6)-(A-7) can be written as

$$\dot{\mathbf{u}}(t) = \mathbb{E}[(\delta_t - \mathbb{E}[\delta_t|\boldsymbol{\theta}])\boldsymbol{\phi}_t|\boldsymbol{\theta}] - \mathbf{C}\mathbf{u}(t). \quad (\text{A-8})$$

Now consider the function  $h(\mathbf{u}) = \mathbb{E}[\delta_t - \mathbb{E}[\delta_t|\boldsymbol{\theta}]]\boldsymbol{\phi}|\boldsymbol{\theta}] - \mathbf{C}\mathbf{u}$ , i.e., the driving vector field of the ODE (A-8). For (A-8),  $\mathbf{u}^* = \mathbf{C}^{-1}\mathbb{E}[(\delta - \mathbb{E}[\delta|\boldsymbol{\theta}])\boldsymbol{\phi}|\boldsymbol{\theta}]$  is the unique globally asymptotically stable equilibrium. Let  $h_\infty(\mathbf{u}) = -\mathbf{C}\mathbf{u}$ . For the ODE

$$\dot{\mathbf{u}}(t) = h_\infty(\mathbf{u}(t)) = -\mathbf{C}\mathbf{u}(t), \quad (\text{A-9})$$

the origin of (A-9) is a globally asymptotically stable equilibrium because  $\mathbf{C}$  is a positive definite matrix (because it is non-negative definite and nonsingular). Now Assumptions (A1) and (A2) of (Borkar and Meyn 2000) are verified. Furthermore, Assumptions (TS) of (Borkar and Meyn 2000) is satisfied by our conditions on the step-size sequences  $\alpha_k, \zeta_k, \beta_k$ . Thus, by Theorem 2.2 of (Borkar and Meyn 2000) we obtain that  $\|\mathbf{u}_k - \mathbf{u}^*\| \rightarrow 0$  almost surely as  $k \rightarrow \infty$ .

Consider now the slower timescale recursion (5). In the light of the above, (5) can be rewritten as

$$\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \alpha_k(\delta_k - \mathbb{E}[\delta_k|\boldsymbol{\theta}_k])\boldsymbol{\phi}_k - \alpha_k\gamma\boldsymbol{\phi}'_k(\boldsymbol{\phi}_k^\top \mathbf{C}^{-1}\mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\boldsymbol{\theta}_k])\boldsymbol{\phi}|\boldsymbol{\theta}_k]). \quad (\text{A-10})$$

Let  $\mathcal{G}(k) = \sigma(\boldsymbol{\theta}_l, l \leq k; \boldsymbol{\phi}_s, \boldsymbol{\phi}'_s, r_s, s < k)$ ,  $k \geq 1$  be the sigma fields generated by  $\boldsymbol{\theta}_0, \boldsymbol{\theta}_{l+1}, \boldsymbol{\phi}_l, \boldsymbol{\phi}'_l, 0 \leq l < k$ . Let

$$\begin{aligned} Z_{k+1} &= ((\delta_k - \mathbb{E}[\delta_k|\boldsymbol{\theta}_k])\boldsymbol{\phi}_k - \gamma\boldsymbol{\phi}'_k\boldsymbol{\phi}_k^\top \mathbf{C}^{-1}\mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\boldsymbol{\theta}_k])\boldsymbol{\phi}|\boldsymbol{\theta}_k]) \\ &\quad - \mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\boldsymbol{\theta}_k])\boldsymbol{\phi}_k - \gamma\boldsymbol{\phi}'_k\boldsymbol{\phi}_k^\top \mathbf{C}^{-1}\mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\boldsymbol{\theta}_k])\boldsymbol{\phi}|\boldsymbol{\theta}_k])|\mathcal{G}(k)] \\ &= ((\delta_k - \mathbb{E}[\delta_k|\boldsymbol{\theta}_k])\boldsymbol{\phi}_k - \gamma\boldsymbol{\phi}'_k\boldsymbol{\phi}_k^\top \mathbf{C}^{-1}\mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\boldsymbol{\theta}_k])\boldsymbol{\phi}|\boldsymbol{\theta}_k]) \\ &\quad - \mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\boldsymbol{\theta}_k])\boldsymbol{\phi}_k|\boldsymbol{\theta}_k] - \gamma\mathbb{E}[\boldsymbol{\phi}'_k\boldsymbol{\phi}_k^\top]\mathbf{C}^{-1}\mathbb{E}[(\delta_k - \mathbb{E}[\delta_k|\boldsymbol{\theta}_k])\boldsymbol{\phi}_k|\boldsymbol{\theta}_k]). \end{aligned}$$

It is easy to see that  $Z_k, k \geq 0$  are integrable random variables and  $\mathbb{E}[Z_{k+1}|\mathcal{G}(k)] = 0, \forall k \geq 0$ . Further,

$$\mathbb{E}[\|Z_{k+1}\|^2|\mathcal{G}(k)] \leq c_3(1 + \|\boldsymbol{\theta}_k\|^2), k \geq 0$$

for some constant  $c_3 \geq 0$ , again because  $\boldsymbol{\phi}_k, r_k$ , and  $\boldsymbol{\phi}'_k$  have uniformly bounded second moments, it can be seen that for some constant.

Consider now the following ODE associated with (5):

$$\dot{\boldsymbol{\theta}}(t) = (\mathbf{I} - \mathbb{E}[\gamma\boldsymbol{\phi}'\boldsymbol{\phi}^\top]\mathbf{C}^{-1})\mathbb{E}[(\delta - \mathbb{E}[\delta|\boldsymbol{\theta}(t)])\boldsymbol{\phi}|\boldsymbol{\theta}(t)]. \quad (\text{A-11})$$

Let

$$\begin{aligned} \vec{h}(\boldsymbol{\theta}(t)) &= (\mathbf{I} - \mathbb{E}[\gamma\boldsymbol{\phi}'\boldsymbol{\phi}^\top]\mathbf{C}^{-1})\mathbb{E}[(\delta - \mathbb{E}[\delta|\boldsymbol{\theta}(t)])\boldsymbol{\phi}|\boldsymbol{\theta}(t)] \\ &= (\mathbf{C} - \mathbb{E}[\gamma\boldsymbol{\phi}'\boldsymbol{\phi}^\top])\mathbf{C}^{-1}\mathbb{E}[(\delta - \mathbb{E}[\delta|\boldsymbol{\theta}(t)])\boldsymbol{\phi}|\boldsymbol{\theta}(t)] \\ &= (\mathbb{E}[\boldsymbol{\phi}\boldsymbol{\phi}^\top] - \mathbb{E}[\gamma\boldsymbol{\phi}'\boldsymbol{\phi}^\top])\mathbf{C}^{-1}\mathbb{E}[(\delta - \mathbb{E}[\delta|\boldsymbol{\theta}(t)])\boldsymbol{\phi}|\boldsymbol{\theta}(t)] \\ &= \mathbf{A}^\top \mathbf{C}^{-1}(-\mathbf{A}\boldsymbol{\theta}(t) + \mathbf{b}), \end{aligned}$$

because  $\mathbb{E}[(\delta - \mathbb{E}[\delta|\boldsymbol{\theta}(t)])\boldsymbol{\phi}|\boldsymbol{\theta}(t)] = -\mathbf{A}\boldsymbol{\theta}(t) + \mathbf{b}$ , where  $\mathbf{A} = \text{Cov}(\boldsymbol{\phi}, \boldsymbol{\phi} - \gamma\boldsymbol{\phi}')$ ,  $\mathbf{b} = \text{Cov}(r, \boldsymbol{\phi})$ , and  $\mathbf{C} = \mathbb{E}[\boldsymbol{\phi}\boldsymbol{\phi}^\top]$

Therefore,  $\boldsymbol{\theta}^* = \mathbf{A}^{-1}\mathbf{b}$  can be seen to be the unique globally asymptotically stable equilibrium for ODE (A-11). Let  $\vec{h}_\infty(\boldsymbol{\theta}) = \lim_{r \rightarrow \infty} \frac{\vec{h}(r\boldsymbol{\theta})}{r}$ . Then  $\vec{h}_\infty(\boldsymbol{\theta}) = -\mathbf{A}^\top \mathbf{C}^{-1}\mathbf{A}\boldsymbol{\theta}$  is well-defined. Consider now the ODE

$$\dot{\boldsymbol{\theta}}(t) = -\mathbf{A}^\top \mathbf{C}^{-1}\mathbf{A}\boldsymbol{\theta}(t). \quad (\text{A-12})$$

Because  $\mathbf{C}^{-1}$  is positive definite and  $\mathbf{A}$  has full rank (as it is nonsingular by assumption), the matrix  $\mathbf{A}^\top \mathbf{C}^{-1}\mathbf{A}$  is also positive definite. The ODE (A-12) has the origin as its unique globally asymptotically stable equilibrium. Thus, the assumption (A1) and (A2) are verified.

The proof is given above. In the fastest time scale, the parameter  $w$  converges to  $\mathbb{E}[\delta|\mathbf{u}_k, \boldsymbol{\theta}_k]$ . In the second fast time scale, the parameter  $u$  converges to  $\mathbf{C}^{-1}\mathbb{E}[(\delta - \mathbb{E}[\delta|\boldsymbol{\theta}_k])\boldsymbol{\phi}|\boldsymbol{\theta}_k]$ . In the slower time scale, the parameter  $\boldsymbol{\theta}$  converges to  $\mathbf{A}^{-1}\mathbf{b}$ .  $\square$

## A.2 Proof of Theorem 2

*Proof.* The proof of VMETD's convergence is also based on Borkar's Theorem for general stochastic approximation recursions with two time scales (Borkar 1997).

The VMTD's solution is  $\theta_{\text{VMETD}} = \mathbf{A}_{\text{VMETD}}^{-1} \mathbf{b}_{\text{VMETD}}$ .

First, note that recursion (19) can be rewritten as

$$\theta_{k+1} \leftarrow \theta_k + \beta_k \xi(k),$$

where

$$\xi(k) = \frac{\alpha_k}{\beta_k} (F_k \rho_k \delta_k - \omega_{k+1}) \phi_k$$

Due to the settings of step-size schedule  $\alpha_k = o(\beta_k)$ ,  $\xi(k) \rightarrow 0$  almost surely as  $k \rightarrow \infty$ . That is the increments in iteration (13) are uniformly larger than those in (12), thus (13) is the faster recursion. Along the faster time scale, iterations of (12) and (13) are associated to ODEs system as follows:

$$\dot{\theta}(t) = 0, \quad (\text{A-13})$$

$$\dot{\omega}(t) = \mathbb{E}_\mu[F_t \rho_t \delta_t | \theta(t)] - \omega(t). \quad (\text{A-14})$$

Based on the ODE (A-13),  $\theta(t) \equiv \theta$  when viewed from the faster timescale. By the Hirsch lemma (Hirsch 1989), it follows that  $\|\theta_k - \theta\| \rightarrow 0$  a.s. as  $k \rightarrow \infty$  for some  $\theta$  that depends on the initial condition  $\theta_0$  of recursion (12). Thus, the ODE pair (A-13)-(A-14) can be written as

$$\dot{\omega}(t) = \mathbb{E}_\mu[F_t \rho_t \delta_t | \theta] - \omega(t). \quad (\text{A-15})$$

Consider the function  $h(\omega) = \mathbb{E}_\mu[F \rho \delta | \theta] - \omega$ , i.e., the driving vector field of the ODE (A-15). It is easy to find that the function  $h$  is Lipschitz with coefficient  $-1$ . Let  $h_\infty(\cdot)$  be the function defined by  $h_\infty(\omega) = \lim_{x \rightarrow \infty} \frac{h(x\omega)}{x}$ . Then  $h_\infty(\omega) = -\omega$ , is well-defined. For (A-15),  $\omega^* = \mathbb{E}_\mu[F \rho \delta | \theta]$  is the unique globally asymptotically stable equilibrium. For the ODE

$$\dot{\omega}(t) = h_\infty(\omega(t)) = -\omega(t), \quad (\text{A-16})$$

apply  $\bar{V}(\omega) = (-\omega)^\top (-\omega)/2$  as its associated strict Liapunov function. Then, the origin of (A-16) is a globally asymptotically stable equilibrium.

Consider now the recursion (13). Let  $M_{k+1} = (F_k \rho_k \delta_k - \omega_k) - \mathbb{E}_\mu[(F_k \rho_k \delta_k - \omega_k) | \mathcal{F}(k)]$ , where  $\mathcal{F}(k) = \sigma(\omega_l, \theta_l, l \leq k; \phi_s, \phi'_s, r_s, s < k)$ ,  $k \geq 1$  are the sigma fields generated by  $\omega_0, \theta_0, \omega_{l+1}, \theta_{l+1}, \phi_l, \phi'_l$ ,  $0 \leq l < k$ . It is easy to verify that  $M_{k+1}, k \geq 0$  are integrable random variables that satisfy  $\mathbb{E}[M_{k+1} | \mathcal{F}(k)] = 0$ ,  $\forall k \geq 0$ . Because  $\phi_k, r_k$ , and  $\phi'_k$  have uniformly bounded second moments, it can be seen that for some constant  $c_1 > 0$ ,  $\forall k \geq 0$ ,

$$\mathbb{E}[\|M_{k+1}\|^2 | \mathcal{F}(k)] \leq c_1(1 + \|\omega_k\|^2 + \|\theta_k\|^2).$$

Now Assumptions (A1) and (A2) of (Borkar and Meyn 2000) are verified. Furthermore, Assumptions (TS) of (Borkar and Meyn 2000) is satisfied by our conditions on the step-size sequences  $\alpha_k, \beta_k$ . Thus, by Theorem 2.2 of (Borkar and Meyn 2000) we obtain that  $\|\omega_k - \omega^*\| \rightarrow 0$  almost surely as  $k \rightarrow \infty$ .

Consider now the slower time scale recursion (12). Based on the above analysis, (12) can be rewritten as

$$\begin{aligned} \theta_{k+1} &\leftarrow \theta_k + \alpha_k (F_k \rho_k \delta_k - \omega_k) \phi_k - \alpha_k \omega_{k+1} \phi_k \\ &= \theta_k + \alpha_k (F_k \rho_k \delta_k - \mathbb{E}_\mu[F_k \rho_k \delta_k | \theta_k]) \phi_k \\ &= \theta_k + \alpha_k F_k \rho_k (R_{k+1} + \gamma \theta_k^\top \phi_{k+1} - \theta_k^\top \phi_k) \phi_k - \alpha_k \mathbb{E}_\mu[F_k \rho_k \delta_k] \phi_k \\ &= \theta_k + \alpha_k \underbrace{\{(F_k \rho_k R_{k+1} - \mathbb{E}_\mu[F_k \rho_k R_{k+1}]) \phi_k\}}_{\mathbf{b}_{\text{VMETD},k}} - \underbrace{\{(F_k \rho_k \phi_k (\phi_k - \gamma \phi_{k+1})^\top - \phi_k \mathbb{E}_\mu[F_k \rho_k (\phi_k - \gamma \phi_{k+1})]^\top) \theta_k\}}_{\mathbf{A}_{\text{VMETD},k}} \end{aligned}$$

Let  $\mathcal{G}(k) = \sigma(\theta_l, l \leq k; \phi_s, \phi'_s, r_s, s < k)$ ,  $k \geq 1$  be the sigma fields generated by  $\theta_0, \theta_{l+1}, \phi_l, \phi'_l$ ,  $0 \leq l < k$ . Let  $Z_{k+1} = Y_k - \mathbb{E}[Y_k | \mathcal{G}(k)]$ , where

$$Y_k = (F_k \rho_k \delta_k - \mathbb{E}_\mu[F_k \rho_k \delta_k | \theta_k]) \phi_k.$$

Consequently,

$$\begin{aligned} \mathbb{E}_\mu[Y_k | \mathcal{G}(k)] &= \mathbb{E}_\mu[(F_k \rho_k \delta_k - \mathbb{E}_\mu[F_k \rho_k \delta_k | \theta_k]) \phi_k | \mathcal{G}(k)] \\ &= \mathbb{E}_\mu[F_k \rho_k \delta_k \phi_k | \theta_k] - \mathbb{E}_\mu[\mathbb{E}_\mu[F_k \rho_k \delta_k | \theta_k] \phi_k] \\ &= \mathbb{E}_\mu[F_k \rho_k \delta_k \phi_k | \theta_k] - \mathbb{E}_\mu[F_k \rho_k \delta_k | \theta_k] \mathbb{E}_\mu[\phi_k] \\ &= \text{Cov}(F_k \rho_k \delta_k | \theta_k, \phi_k), \end{aligned}$$

where  $\text{Cov}(\cdot, \cdot)$  is a covariance operator.

Thus,

$$Z_{k+1} = (F_k \rho_k \delta_k - \mathbb{E}[\delta_k | \theta_k]) \phi_k - \text{Cov}(F_k \rho_k \delta_k | \theta_k, \phi_k).$$

It is easy to verify that  $Z_{k+1}, k \geq 0$  are integrable random variables that satisfy  $\mathbb{E}[Z_{k+1}|\mathcal{G}(k)] = 0, \forall k \geq 0$ . Also, because  $\phi_k, r_k$ , and  $\phi'_k$  have uniformly bounded second moments, it can be seen that for some constant  $c_2 > 0, \forall k \geq 0$ ,

$$\mathbb{E}[||Z_{k+1}||^2|\mathcal{G}(k)] \leq c_2(1 + ||\theta_k||^2).$$

Consider now the following ODE associated with (12):

$$\dot{\theta}(t) = -\mathbf{A}_{\text{VMETD}}\theta(t) + \mathbf{b}_{\text{VMETD}}. \quad (\text{A-17})$$

$$\begin{aligned} \mathbf{A}_{\text{VMETD}} &= \lim_{k \rightarrow \infty} \mathbb{E}[\mathbf{A}_{\text{VMETD},k}] \\ &= \lim_{k \rightarrow \infty} \mathbb{E}_\mu[F_k \rho_k \phi_k (\phi_k - \gamma \phi_{k+1})^\top] - \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k] \mathbb{E}_\mu[F_k \rho_k (\phi_k - \gamma \phi_{k+1})]^\top \\ &= \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k F_k \rho_k (\phi_k - \gamma \phi_{k+1})^\top] - \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k] \mathbb{E}_\mu[F_k \rho_k (\phi_k - \gamma \phi_{k+1})]^\top \\ &= \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k F_k \rho_k (\phi_k - \gamma \phi_{k+1})^\top] - \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k] \lim_{k \rightarrow \infty} \mathbb{E}_\mu[F_k \rho_k (\phi_k - \gamma \phi_{k+1})]^\top \\ &= \sum_s f(s) \phi(s) (\phi(s) - \gamma \sum_{s'} [\mathbf{P}_\pi]_{ss'} \phi(s'))^\top - \sum_s d_\mu(s) \phi(s) * \sum_s f(s) (\phi(s) - \gamma \sum_{s'} [\mathbf{P}_\pi]_{ss'} \phi(s'))^\top \quad (\text{A-18}) \\ &= \Phi^\top \mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) \Phi - \Phi^\top \mathbf{d}_\mu \mathbf{f}^\top (\mathbf{I} - \gamma \mathbf{P}_\pi) \Phi \\ &= \Phi^\top (\mathbf{F} - \mathbf{d}_\mu \mathbf{f}^\top) (\mathbf{I} - \gamma \mathbf{P}_\pi) \Phi \\ &= \Phi^\top (\mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{f}^\top (\mathbf{I} - \gamma \mathbf{P}_\pi)) \Phi \\ &= \Phi^\top (\mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{d}_\mu^\top) \Phi \end{aligned}$$

$$\begin{aligned} \mathbf{b}_{\text{VMETD}} &= \lim_{k \rightarrow \infty} \mathbb{E}[\mathbf{b}_{\text{VMETD},k}] \\ &= \lim_{k \rightarrow \infty} \mathbb{E}_\mu[F_k \rho_k R_{k+1} \phi_k] - \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k] \mathbb{E}_\mu[F_k \rho_k R_{k+1}] \\ &= \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k F_k \rho_k R_{k+1}] - \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k] \mathbb{E}_\mu[F_k \rho_k R_{k+1}] \\ &= \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k F_k \rho_k R_{k+1}] - \lim_{k \rightarrow \infty} \mathbb{E}_\mu[\phi_k] \lim_{k \rightarrow \infty} \mathbb{E}_\mu[F_k \rho_k R_{k+1}] \quad (\text{A-19}) \\ &= \sum_s f(s) \phi(s) r_\pi - \sum_s d_\mu(s) \phi(s) * \sum_s f(s) r_\pi \\ &= \Phi^\top (\mathbf{F} - \mathbf{d}_\mu \mathbf{f}^\top) \mathbf{r}_\pi \end{aligned}$$

Let  $\vec{h}(\theta(t))$  be the driving vector field of the ODE (A-17).

$$\vec{h}(\theta(t)) = -\mathbf{A}_{\text{VMETD}}\theta(t) + \mathbf{b}_{\text{VMETD}}.$$

An  $\Phi^\top \mathbf{X} \Phi$  matrix of this form will be positive definite whenever the matrix  $\mathbf{X}$  is positive definite. Any matrix  $\mathbf{X}$  is positive definite if and only if the symmetric matrix  $\mathbf{S} = \mathbf{X} + \mathbf{X}^\top$  is positive definite. Any symmetric real matrix  $\mathbf{S}$  is positive definite if the absolute values of its diagonal entries are greater than the sum of the absolute values of the corresponding off-diagonal entries (Sutton, Mahmood, and White 2016).

$$\begin{aligned} (\mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{d}_\mu^\top) \mathbf{1} &= \mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) \mathbf{1} - \mathbf{d}_\mu \mathbf{d}_\mu^\top \mathbf{1} \\ &= \mathbf{F}(\mathbf{1} - \gamma \mathbf{P}_\pi \mathbf{1}) - \mathbf{d}_\mu \mathbf{d}_\mu^\top \mathbf{1} \\ &= (1 - \gamma) \mathbf{F} \mathbf{1} - \mathbf{d}_\mu \mathbf{d}_\mu^\top \mathbf{1} \\ &= (1 - \gamma) \mathbf{f} - \mathbf{d}_\mu \mathbf{d}_\mu^\top \mathbf{1} \\ &= (1 - \gamma) \mathbf{f} - \mathbf{d}_\mu \\ &= (1 - \gamma) (\mathbf{I} - \gamma \mathbf{P}_\pi^\top)^{-1} \mathbf{d}_\mu - \mathbf{d}_\mu \quad (\text{A-20}) \\ &= (1 - \gamma) [(\mathbf{I} - \gamma \mathbf{P}_\pi^\top)^{-1} - \mathbf{I}] \mathbf{d}_\mu \\ &= (1 - \gamma) \left[ \sum_{t=0}^{\infty} (\gamma \mathbf{P}_\pi^\top)^t - \mathbf{I} \right] \mathbf{d}_\mu \\ &= (1 - \gamma) \left[ \sum_{t=1}^{\infty} (\gamma \mathbf{P}_\pi^\top)^t \right] \mathbf{d}_\mu > 0 \end{aligned}$$

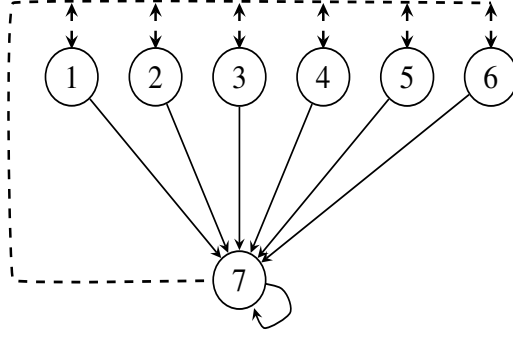


Figure 1: 7-state version of Baird's off-policy counterexample.

$$\begin{aligned}
\mathbf{1}^\top (\mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{d}_\mu^\top) &= \mathbf{1}^\top \mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{1}^\top \mathbf{d}_\mu \mathbf{d}_\mu^\top \\
&= \mathbf{d}_\mu^\top - \mathbf{1}^\top \mathbf{d}_\mu \mathbf{d}_\mu^\top \\
&= \mathbf{d}_\mu^\top - \mathbf{d}_\mu^\top \\
&= 0
\end{aligned} \tag{A-21}$$

(A-20) and (A-21) show that the matrix  $\mathbf{F}(\mathbf{I} - \gamma \mathbf{P}_\pi) - \mathbf{d}_\mu \mathbf{d}_\mu^\top$  of diagonal entries are positive and its off-diagonal entries are negative. So its each row sum plus the corresponding column sum is positive. So  $\mathbf{A}_{\text{VMETD}}$  is positive definite.

Therefore,  $\theta^* = \mathbf{A}_{\text{VMETD}}^{-1} \mathbf{b}_{\text{VMETD}}$  can be seen to be the unique globally asymptotically stable equilibrium for ODE (A-17). Let  $\vec{h}_\infty(\theta) = \lim_{r \rightarrow \infty} \frac{\vec{h}(r\theta)}{r}$ . Then  $\vec{h}_\infty(\theta) = -\mathbf{A}_{\text{VMETD}}\theta$  is well-defined. Consider now the ODE

$$\dot{\theta}(t) = -\mathbf{A}_{\text{VMETD}}\theta(t). \tag{A-22}$$

The ODE (A-22) has the origin as its unique globally asymptotically stable equilibrium. Thus, the assumption (A1) and (A2) are verified.  $\square$

## B Experimental details

The 2-state counterexample and the 7-state counterexample are well-known off-policy experimental environments. The 2-state counterexample is relatively simple, so next, I'll provide a detailed description of the 7-state counterexample environment.

**Baird's off-policy counterexample:** This task is well known as a counterexample, in which TD diverges (Baird et al. 1995; Sutton et al. 2009). As shown in Figure 1, reward for each transition is zero. Thus the true values are zeros for all states and for any given policy. The behaviour policy chooses actions represented by solid lines with a probability of  $\frac{1}{7}$  and actions represented by dotted lines with a probability of  $\frac{6}{7}$ . The target policy is expected to choose the solid line with more probability than  $\frac{1}{7}$ , and it chooses the solid line with probability of 1 in this paper. The discount factor  $\gamma = 0.99$ , and the feature matrix is defined in Appendix B (Baird et al. 1995; Sutton et al. 2009; Maei 2011).

The feature matrix of 7-state version of Baird's off-policy counterexample is defined as follow:

$$\Phi_{\text{Counter}} = \begin{bmatrix} 1 & 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 2 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 2 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 2 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 2 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 2 & 0 \\ 2 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

2-state version of Baird's off-policy counterexample: All learning rates follow linear learning rate decay. For TDC algorithm,  $\frac{\alpha_k}{\zeta_k} = 5$  and  $\alpha_0 = 0.1$ . For ETD algorithm,  $\alpha_0 = 0.1$ . For VMTDC algorithm,  $\frac{\alpha_k}{\zeta_k} = 5$ ,  $\frac{\alpha_k}{\omega_k} = 4$ , and  $\alpha_0 = 0.1$ . For ETD algorithm,  $\frac{\alpha_k}{\omega_k} = 4$  and  $\alpha_0 = 0.1$ .

7-state version of Baird's off-policy counterexample: All learning rates follow linear learning rate decay. For TDC algorithm,  $\frac{\alpha_k}{\zeta_k} = 3$  and  $\alpha_0 = 0.1$ . For ETD algorithm,  $\alpha_0 = 0.1$ . For VMTDC algorithm,  $\frac{\alpha_k}{\zeta_k} = 3$ ,  $\frac{\alpha_k}{\omega_k} = 4$ , and  $\alpha_0 = 0.1$ . For ETD algorithm,  $\frac{\alpha_k}{\omega_k} = 4$  and  $\alpha_0 = 0.1$ .

For all policy evaluation experiments, each experiment is independently run 100 times.

For the four control experiments: The learning rates for each algorithm in all experiments are shown in Table 1. For all control experiments, each experiment is independently run 50 times.

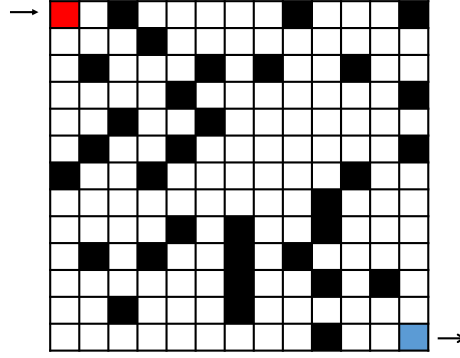


Figure 2: Maze.

**Maze:** The learning agent should find a shortest path from the upper left corner to the lower right corner. In each state, there are four alternative actions: *up*, *down*, *left*, and *right*, which takes the agent deterministically to the corresponding neighbour state, except when a movement is blocked by an obstacle or the edge of the maze. Rewards are  $-1$  in all transitions until the agent reaches the goal state. The discount factor  $\gamma = 0.99$ , and states  $s$  are represented by tabular features. The maximum number of moves in the game is set to 1000.

**The other three control environments:** Cliff Walking, Mountain Car, and Acrobot are selected from the gym official website and correspond to the following versions: “CliffWalking-v0”, “MountainCar-v0” and “Acrobot-v1”. For specific details, please refer to the gym official website. The maximum number of steps for the Mountain Car environment is set to 1000, while the default settings are used for the other two environments. In Mountain car and Acrobot, features are generated by tile coding.

Table 1: Learning rates ( $lr$ ) of four control experiments.

algorithms( $lr$ ) \ envs	Maze	Cliff walking	Mountain Car	Acrobot
GQ( $\alpha, \zeta$ )	0.1, 0.003	0.1, 0.004	0.1, 0.01	0.1, 0.01
EQ( $\alpha$ )	0.006	0.005	0.001	0.0005
VMGQ( $\alpha, \zeta, \beta$ )	0.1, 0.001, 0.001	0.1, 0.005, 1e-4	0.1, 5e-4, 1e-4	0.1, 5e-4, 1e-4
VMEQ( $\alpha, \beta$ )	0.001, 0.0005	0.005, 0.0001	0.001, 0.0001	0.0005, 0.0001

## References

- Baird, L.; et al. 1995. Residual algorithms: Reinforcement learning with function approximation. In *Proc. 12th Int. Conf. Mach. Learn.*, 30–37.
- Borkar, V. S. 1997. Stochastic approximation with two time scales. *Syst. & Control Letters*, 29(5): 291–294.
- Borkar, V. S.; and Meyn, S. P. 2000. The ODE method for convergence of stochastic approximation and reinforcement learning. *SIAM J. Control Optim.*, 38(2): 447–469.
- Hirsch, M. W. 1989. Convergent activation dynamics in continuous time networks. *Neural Netw.*, 2(5): 331–349.
- Maei, H. R. 2011. *Gradient temporal-difference learning algorithms*. Ph.D. thesis, University of Alberta.
- Sutton, R.; Maei, H.; Precup, D.; Bhatnagar, S.; Silver, D.; Szepesvári, C.; and Wiewiora, E. 2009. Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *Proc. 26th Int. Conf. Mach. Learn.*, 993–1000.
- Sutton, R. S.; Mahmood, A. R.; and White, M. 2016. An emphatic approach to the problem of off-policy temporal-difference learning. *The Journal of Machine Learning Research*, 17(1): 2603–2631.