

# TD-CD-SAC Online Training Algorithm

## Algorithm Implementation

February 5, 2026

## 1 Notation

- $Q_{\theta_j}$ :  $j$ -th critic network with parameters  $\theta_j$
- $\pi_\phi$ : Actor network with parameters  $\phi$
- $\bar{\theta}_j$ : Target network parameters for critic  $j$
- $\alpha$ : Temperature parameter for entropy regularization
- $\mathcal{D}$ : Replay buffer
- $c_t$ : Constraint violation at time  $t$
- $\delta_t$ : Soft termination probability at time  $t$
- $\gamma_t$ : Discount factor at time  $t$
- $c_{\max}^{seen}$ : Maximum constraint violation seen so far
- $c_{\max}^{ema}$ : Exponential moving average of maximum constraint violation

## 2 Algorithm Description

This algorithm implements TD-CD-SAC (Time-Dependent Constraint-Directed Soft Actor-Critic), an online reinforcement learning method that handles constraints through time-dependent discount factors. The key innovation is the use of soft termination probabilities  $\delta_t$  that modify the effective discount factor  $\gamma_t$  based on constraint violations.

The algorithm operates in two main phases:

1. **Data Collection:** Actions are sampled from the policy and executed in the environment to collect transitions.
2. **Network Updates:** The actor and critic networks are updated using the collected data, with modified target values that account for constraint violations.

The constraint handling mechanism works by:

- Detecting constraint violations through environment information
- Computing soft termination probabilities based on violation severity
- Modifying discount factors to penalize constraint-violating trajectories
- Updating the policy to avoid constraint violations while maximizing rewards

---

**Algorithm 1** TD-CD-SAC Online Training (cd\_sac\_ball)

---

**Require:** Target  $g$ , constrained environment  $\mathcal{E}$  with bounds  $(v_{\max}, a_{\max})$   
**Require:** Discount factor  $\gamma \in (0, 1)$ , soft-update parameter  $\tau$ , batch size  $B$   
**Require:** TD-CD hyperparameters  $p^{\max} \in [0, 1]$ , EMA factor  $\tau_c \in [0, 1]$   
**Require:** Switch `use_amount`  $\in \{0, 1\}$  (binary/continuous constraint signal)

- 1: Initialize SAC networks  $(Q_{\theta_1}, Q_{\theta_2}, \pi_\phi)$  and targets  $\bar{\theta}_1 \leftarrow \theta_1$ ,  $\bar{\theta}_2 \leftarrow \theta_2$
- 2: Initialize replay buffer  $\mathcal{D} \leftarrow \emptyset$
- 3: Initialize TD-CD stats:  $c_{\max}^{\text{seen}} \leftarrow 0$ ,  $c_{\max}^{\text{ema}} \leftarrow 1$
- 4: Sample initial state  $s_0 \sim \mathcal{E}.\text{reset}()$
- 5: **for**  $t = 0, 1, 2, \dots, T - 1$  **do**
- 6:     Select action  $u_t \sim \begin{cases} \mathcal{U}([-1, 1]^m) & t < t_{\text{start}} \\ \pi_\phi(\cdot | s_t) & \text{otherwise} \end{cases}$
- 7:     Execute in environment:  $(s_{t+1}, r_t, d_t, \text{info}_t) \leftarrow \mathcal{E}.\text{step}(u_t)$
- 8:         ▷ Time-limit truncation is treated as non-terminal
- 9:      $d_t^{\text{buf}} \leftarrow d_t$ ; **if**  $\text{info}_t.\text{time\_limit} = 1$  **then**  $d_t^{\text{buf}} \leftarrow 0$
- 10:    **if**  $d_t^{\text{buf}} = 1$  **then**
- 11:       $\gamma_t \leftarrow 0$
- 12:    **else**
- 13:         ▷ TD-CD Equation (7): build soft termination  $\delta_t$
- 14:         **if** `use_amount`=1 **then**
- 15:            $c_t \leftarrow \text{info}_t.\text{vel\_violation\_amount}$  ▷  $c_t \geq 0$
- 16:            $c_{\max}^{\text{seen}} \leftarrow \max(c_{\max}^{\text{seen}}, |c_t|)$
- 17:            $\delta_t \leftarrow p^{\max} \cdot \text{clip}\left(\frac{|c_t|}{c_{\max}^{\text{ema}}}, 0, 1\right)$
- 18:         **else**
- 19:            $c_t \leftarrow \mathbb{I}[\text{info}_t.\text{constraint\_violation} = 1]$
- 20:            $\delta_t \leftarrow p^{\max} \cdot c_t$
- 21:         **end if**
- 22:         ▷ TD-CD Equation (9): per-step discount
- 23:          $\gamma_t \leftarrow \gamma(1 - \delta_t)$
- 24:     **end if**
- 25:     Store transition:  $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, u_t, r_t, s_{t+1}, d_t^{\text{buf}}, \gamma_t)\}$
- 26:      $s_t \leftarrow s_{t+1}$
- 27:     **if**  $t \geq t_{\text{update}}$  and  $|\mathcal{D}| \geq B$  **then**
- 28:         **for**  $k = 1$  to  $K$  **do** ▷  $K$  updates per environment step
- 29:             Sample minibatch  $\{(s_i, u_i, r_i, s'_i, d_i, \gamma_i)\}_{i=1}^B \sim \mathcal{D}$
- 30:             Sample next actions  $a'_i \sim \pi_\phi(\cdot | s'_i)$  and log-probs  $\log \pi_\phi(a'_i | s'_i)$
- 31:             Compute soft value target:
- 32:              $V(s'_i) = \min_{j \in \{1, 2\}} Q_{\bar{\theta}_j}(s'_i, a'_i) - \alpha \log \pi_\phi(a'_i | s'_i)$  ▷ Key change vs vanilla SAC: use per-transition  $\gamma_i$
- 33:             Critic target:  $y_i \leftarrow r_i + \gamma_i V(s'_i)$
- 34:             Update critics by MSE:  $\theta_j \leftarrow \arg \min_{\theta_j} \frac{1}{B} \sum_i (Q_{\theta_j}(s_i, u_i) - y_i)^2$
- 35:             **for**  $j = 1, 2$  **do** Optional: Update temperature  $\alpha$  by entropy tuning
- 36:                 Soft-update targets:  $\bar{\theta}_j \leftarrow \tau \theta_j + (1 - \tau) \bar{\theta}_j$  for  $j = 1, 2$
- 37:             **end for**
- 38:         **end if**
- 39:         **if**  $t > 0$  and  $t \bmod N_{\text{eval}} = 0$  **then** ▷ TD-CD Equation (8): EMA update once per window
- 40:             Update actor (standard SAC):
- 41:              $\phi \leftarrow \arg \min_{\phi} \frac{1}{B} \sum_i (\alpha \log \pi_\phi(a_i | s_i) - \min_j Q_{\theta_j}(s_i, a_i))$
- 42:              $c_{\max}^{\text{ema}} \leftarrow \tau_c c_{\max}^{\text{ema}} + (1 - \tau_c) \max(c_{\max}^{\text{seen}}, \epsilon)$
- 43:              $c_{\max}^{\text{seen}} \leftarrow 0$
- 44:             (Optional) evaluate policy and save checkpoint
- 45:         **end if**