

Contents

1 Chapter 3: SAC on Dense Reach -- Validating the Off-Policy Stack	2
1.1 What This Chapter Is Really About	2
1.1.1 Early Training: Exploration (30k steps, entropy = 0.47)	2
1.1.2 Final Training: Exploitation (1M steps, entropy = 0.0004)	2
1.2 Part 0: Setting the Stage	3
1.2.1 0.1 Why SAC After PPO?	3
1.2.2 0.2 The Diagnostic Mindset (Again)	3
1.3 Part 1: WHY -- Maximum Entropy Reinforcement Learning	3
1.3.1 1.1 The Standard RL Objective (Review)	3
1.3.2 1.2 Why Determinism Is Actually a Problem	3
1.3.3 1.3 The Maximum Entropy Objective	4
1.3.4 1.4 Automatic Temperature Tuning	4
1.3.5 1.5 Why This Matters for Robotics	4
1.4 Part 2: HOW -- The SAC Algorithm	4
1.4.1 2.1 The Components	4
1.4.2 2.2 The Training Loop	5
1.4.3 2.3 Key Differences from PPO	6
1.4.4 2.4 Key Hyperparameters	6
1.5 Part 2.5: BUILD IT -- From Equations to Code	6
1.5.1 2.5.1 The Replay Buffer	6
1.5.2 2.5.2 The Twin Q-Network	7
1.5.3 2.5.3 The Squashed Gaussian Policy	7
1.5.4 2.5.4 Twin Q-Network Loss (Critic Update)	7
1.5.5 2.5.5 Actor Loss with Entropy	8
1.5.6 2.5.6 Automatic Temperature Tuning	8
1.5.7 2.5.7 Wiring It Together: The SAC Update	9
1.5.8 2.5.8 Verify the Full Lab	9
1.5.9 2.5.9 Verify vs SB3 (Optional)	10
1.5.10 2.5.10 Exercises: Modify and Observe	10
1.5.11 2.5.11 Demo: SAC Solves Pendulum	11
1.6 Part 3: WHAT -- Running the Experiment (Run It)	12
1.6.1 3.1 The One-Command Version	12
1.6.2 3.2 What to Expect	12
1.6.3 3.3 What the Diagnostics Callback Logs	12
1.6.4 3.4 What to Watch For	13
1.6.5 3.5 Throughput Scaling	13
1.6.6 3.6 Actual Results	13
1.6.7 3.7 Understanding GPU Utilization	13
1.6.8 3.8 Generating Demo Videos	14
1.7 Part 4: Understanding What You Built	14
1.7.1 4.1 The Replay Buffer	14
1.7.2 4.2 The Min-Q Trick	15
1.7.3 4.3 The Squashed Gaussian	15
1.8 Part 5: Exercises	15
1.8.1 Exercise 3.1: Reproduce the Baseline	15
1.8.2 Exercise 3.2: Compare Learning Curves	15
1.8.3 Exercise 3.3: Throughput Analysis	15
1.8.4 Exercise 3.4: Q-Value Analysis (Written)	16
1.8.5 Exercise 3.5: Temperature Ablation	16
1.9 Part 6: Common Failures and Solutions	16
1.9.1 "Q-values explode (>1000)"	16
1.9.2 "Entropy coefficient goes to zero immediately"	16

1.9.3 "Success rate stalls below 50%"	16
1.9.4 "Training is much slower than PPO"	17
1.10 Part 7: Looking Ahead	17
1.11 References	17

1 Chapter 3: SAC on Dense Reach -- Validating the Off-Policy Stack

1.1 What This Chapter Is Really About

Chapter 2 validated our training pipeline with PPO, achieving 100% success on dense Reach and confirming that the infrastructure works. But PPO has a fundamental limitation: **it discards data after every update**.

The result: A SAC policy that learned to reach any point in 3D space, matching PPO's performance while building the off-policy machinery we need for HER.

1.1.1 Early Training: Exploration (30k steps, entropy = 0.47)

SAC early exploration

At 30k steps, the policy is still exploring. The high entropy coefficient (0.47) encourages random-ish actions. The gripper moves erratically, rarely reaching the target. Success rate: 0%.

1.1.2 Final Training: Exploitation (1M steps, entropy = 0.0004)

SAC trained policy

At 1M steps, the policy has converged. The entropy coefficient dropped to 0.0004 (nearly deterministic). The gripper moves directly to the target. Success rate: 100%.

This is the entropy bonus in action: Early training explores broadly (high entropy), late training exploits what was learned (low entropy). SAC auto-tunes this tradeoff.

PPO's data inefficiency is wasteful: each transition took real simulation time to collect, yet PPO uses it for a few gradient steps and then throws it away. For dense-reward tasks where signal is plentiful this inefficiency is tolerable, but for sparse-reward tasks (coming in Chapter 4) it becomes catastrophic.

Off-policy methods solve this by storing transitions in a **replay buffer** and reusing them across many updates, which is dramatically more sample-efficient -- though it introduces new failure modes: stale data, value overestimation, and training instability.

This chapter validates the off-policy machinery on dense Reach -- the same "easy" task where we know success is achievable. If SAC fails here, the bug is in our off-policy implementation, not the task difficulty.

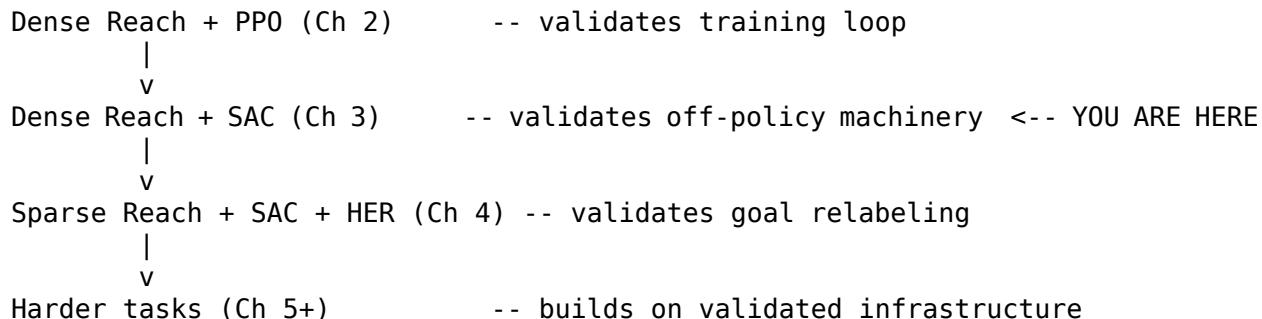
By the end, you will have:

1. A trained SAC policy matching or exceeding PPO's performance
 2. Understanding of maximum-entropy RL and why it matters
 3. Replay buffer diagnostics to detect common off-policy pathologies
 4. Throughput benchmarks for scaling to larger experiments
-

1.2 Part 0: Setting the Stage

1.2.1 0.1 Why SAC After PPO?

We are not choosing SAC because it is "better" -- we are building toward HER (Hindsight Experience Replay), which requires off-policy learning. The dependency chain makes this concrete:



Each step isolates one new component: if SAC fails on dense Reach, we know the bug is in SAC, not the environment; if HER fails on sparse Reach, we know the bug is in HER, not SAC.

1.2.2 0.2 The Diagnostic Mindset (Again)

Dense Reach with SAC should achieve similar performance to PPO, since both are solving the same task with the same amount of signal. If SAC performs much worse, the cause is likely value overestimation (Q-values growing unbounded), entropy collapse (exploration stopped too early), or replay buffer issues (stale data, wrong sampling). The diagnostics callback we use catches these early.

1.3 Part 1: WHY -- Maximum Entropy Reinforcement Learning

1.3.1 1.1 The Standard RL Objective (Review)

Recall from Chapter 2 that standard RL maximizes expected return:

$$J(\theta) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t R_t \right]$$

This finds a policy that gets high reward, but there is a subtle problem: **the optimal policy is deterministic**. Once you know the best action for each state, why ever do anything else?

1.3.2 1.2 Why Determinism Is Actually a Problem

In theory, a deterministic optimal policy is fine. In practice, it causes three related problems.

Problem 1: Exploration dies. A deterministic policy exploits what it knows -- if the current best action gets reward +10, the policy commits to it. But what if there is an action that would get +20, which the policy has never tried because it stopped exploring?

Problem 2: Brittleness. A policy that commits fully to one action per state is fragile, since small perturbations (observation noise, model mismatch) can push it into unfamiliar states where it has no idea what to do.

Problem 3: Training instability. When the policy is nearly deterministic, small value estimate changes cause large behavioral changes (the "winning" action flips), which amplifies noise in the training process.

1.3.3 1.3 The Maximum Entropy Objective

SAC modifies the objective to include an **entropy bonus**:

$$J(\theta) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t \left(R_t + \alpha \mathcal{H}(\pi(\cdot|s_t)) \right) \right]$$

where $\mathcal{H}(\pi(\cdot|s))$ is the **entropy** of the policy at state s (measuring how "spread out" the action distribution is) and $\alpha > 0$ is the **temperature** parameter controlling the exploration-exploitation tradeoff.

Definition (Entropy). For a continuous policy outputting a Gaussian distribution over actions, the entropy is:

$$\mathcal{H}(\pi(\cdot|s)) = \frac{1}{2} \ln((2\pi e)^d |\Sigma|)$$

where d is the action dimension and $|\Sigma|$ is the determinant of the covariance matrix. Higher entropy means the policy is more random; lower entropy means more deterministic.

The key insight: The optimal policy no longer commits fully to one action. Instead, it prefers actions proportionally to their Q-values:

$$\pi^*(a|s) \propto \exp(Q^*(s, a) / \alpha)$$

This is a **Boltzmann distribution** (or "softmax" over continuous actions). High-Q actions are more likely, but low-Q actions still have some probability. The temperature α controls the sharpness: as $\alpha \rightarrow 0$ the policy becomes deterministic (argmax), while as $\alpha \rightarrow \infty$ it approaches uniform random.

1.3.4 1.4 Automatic Temperature Tuning

Choosing α manually is tricky -- too low and you lose exploration, too high and you never exploit. SAC can learn α automatically by targeting a desired entropy level:

$$\alpha^* = \arg\min_{\alpha} \mathbb{E}_{a \sim \pi} \left[-\alpha \log \pi(a|s) - \alpha \bar{\mathcal{H}} \right]$$

where $\bar{\mathcal{H}}$ is the **target entropy** (typically set to $-\dim(\mathcal{A})$, the negative action dimension).

Translation: "Adjust α so that the policy's entropy stays near $\bar{\mathcal{H}}$."

This means SAC adapts its exploration automatically: early in training, when the policy is uncertain, entropy is high naturally, while late in training, when the policy is confident, α decreases to allow more exploitation.

1.3.5 1.5 Why This Matters for Robotics

In robotics, brittleness kills -- a policy that works perfectly in simulation but fails on real hardware is useless. Maximum entropy helps on several fronts: it encourages **diverse training data** (the policy explores many actions, seeing more of the state space), produces **robust behaviors** (the policy does not commit fully to any single action, making it more tolerant to noise), and leads to **smoother training** (the Boltzmann-style policy changes gradually as Q-values change rather than snapping between discrete action choices).

1.4 Part 2: HOW -- The SAC Algorithm

1.4.1 2.1 The Components

SAC maintains five networks (in Stable Baselines 3's implementation):

Network	Purpose	Updates
Actor π_θ	Maps states to action distributions	Policy gradient
Critic 1 Q_{ϕ_1}	Estimates Q-values	Bellman backup
Critic 2 Q_{ϕ_2}	Second Q estimate (reduces overestimation)	Bellman backup
Target Critic 1 $Q_{\bar{\phi}_1}$	Stable target for critic updates	Polyak averaging
Target Critic 2 $Q_{\bar{\phi}_2}$	Stable target for critic updates	Polyak averaging

Why two critics? Q-learning tends to overestimate values because we take a max over noisy estimates. Using two critics and taking the minimum reduces this bias. This is called **clipped double Q-learning**.

Why target networks? If we update the critic using its own predictions as targets, we get a moving target problem. Target networks are slow-moving copies that provide stable targets:

$$\bar{\phi} \leftarrow \tau \phi + (1 - \tau) \bar{\phi}$$

with $\tau = 0.005$ typically (slow updates).

1.4.2 2.2 The Training Loop

repeat:

1. Collect transitions using current policy, store in replay buffer
2. Sample minibatch from replay buffer
3. Update critics using Bellman backup (target uses min of two Q-targets)
4. Update actor to maximize Q + entropy
5. Update temperature alpha (if auto-tuning)
6. Soft-update target networks

Step 1: Collect Data

Unlike PPO, SAC collects data continuously. Each step, it samples an action from the current policy (with stochasticity for exploration), executes it in the environment, and stores the resulting transition $(s, a, r, s', done)$ in the replay buffer. The replay buffer is circular -- when full, new transitions overwrite the oldest.

Step 2: Sample Minibatch

Uniformly sample `batch_size` transitions from the buffer. This is where sample reuse happens -- the same transition might be sampled many times across training.

Step 3: Update Critics

For each critic, minimize the Bellman error:

$$L(\phi_i) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{B}} \left[\left(Q_{\phi_i}(s, a) - y \right)^2 \right]$$

where the target y is:

$$y = r + \gamma (1 - done) \min_{j=1,2} Q_{\bar{\phi}_j}(s', a') - \alpha \log \pi(a|s)$$

with $a' \sim \pi(\cdot|s')$.

Translation: "The Q-value of (s, a) should equal the reward plus the (discounted, entropy-adjusted) value of the next state."

Step 4: Update Actor

Maximize expected Q-value plus entropy:

$$L(\theta) = \mathbb{E}_{s \sim \mathcal{B}, a \sim \pi_\theta} \left[\alpha \log \pi_\theta(a|s) - Q_{\phi_1}(s, a) \right]$$

Translation: "Choose actions that have high Q-values while maintaining high entropy."

1.4.3 2.3 Key Differences from PPO

Aspect	PPO	SAC
Data reuse	None (on-policy)	Extensive (replay buffer)
Exploration	Via entropy bonus (optional)	Via entropy bonus (core)
Sample efficiency	Low	High
Stability	High (clipped updates)	Medium (moving targets)
Complexity	Lower	Higher (more networks)

1.4.4 2.4 Key Hyperparameters

Parameter	Default	What It Controls
buffer_size	1,000,000	Replay buffer capacity
batch_size	256	Minibatch size for updates
learning_starts	10,000	Steps before training begins
tau	0.005	Target network update rate
ent_coef	"auto"	Entropy temperature (auto-tuned)
learning_rate	3e-4	Gradient step size

1.5 Part 2.5: BUILD IT -- From Equations to Code

This section builds SAC piece by piece, verifying each component before moving to the next. We use pedagogical implementations from `scripts/labs/sac_from_scratch.py` -- these are for understanding, not production.

1.5.1 2.5.1 The Replay Buffer

Off-policy learning (Section 2.2) requires storing transitions for reuse. Unlike PPO, which discards data after each update, SAC stores every transition in a circular buffer and samples from it repeatedly:

```
--8<-- "scripts/labs/sac_from_scratch.py:replay_buffer"
!!! lab "Checkpoint" Add transitions and verify sampling works:
```python
buf = ReplayBuffer(obs_dim=10, act_dim=4, capacity=1000)
for i in range(100):
 buf.add(np.random.randn(10), np.random.randn(4), -0.5, np.random.randn(10), False)

batch = buf.sample(32, torch.device("cpu"))
print(f"obs shape: {batch['obs'].shape}") # (32, 10)
print(f"actions shape: {batch['actions'].shape}") # (32, 4)
print(f"rewards shape: {batch['rewards'].shape}") # (32,)
print(f"buffer size: {buf.size}") # 100
```
```

1.5.2 2.5.2 The Twin Q-Network

SAC uses two Q-networks to reduce overestimation bias (Section 2.1). Each takes a state-action pair and outputs a scalar Q-value:

```
--8<-- "scripts/labs/sac_from_scratch.py:twin_q_network"
```

!!! lab "Checkpoint" Verify forward pass shapes and that initial Q-values are near zero (random initialization):

```
```python
q_net = TwinQNetwork(obs_dim=10, act_dim=4)
obs = torch.randn(32, 10); actions = torch.randn(32, 4)
q1, q2 = q_net(obs, actions)
print(f"Q1 shape: {q1.shape}") # (32,)
print(f"Q1 mean: {q1.mean():.4f}") # near 0 at init
print(f"Q2 mean: {q2.mean():.4f}") # near 0 at init
```
```

1.5.3 2.5.3 The Squashed Gaussian Policy

SAC's policy outputs a squashed Gaussian -- a normal distribution passed through tanh to bound actions to $[-1, 1]$. The log probability includes a Jacobian correction for the tanh transformation:

```
--8<-- "scripts/labs/sac_from_scratch.py:gaussian_policy"
```

!!! lab "Checkpoint" Verify actions are bounded and log probabilities are finite:

```
```python
policy = GaussianPolicy(obs_dim=10, act_dim=4)
obs = torch.randn(32, 10)
actions, log_probs = policy(obs)
print(f"Actions in [{actions.min():.2f}, {actions.max():.2f}]") # within [-1, 1]
print(f"Log probs finite: {torch.isfinite(log_probs).all()}") # True
print(f"Log probs mean: {log_probs.mean():.2f}") # negative (probabilities < 1)
```
```

1.5.4 2.5.4 Twin Q-Network Loss (Critic Update)

The critic update from Section 2.2 minimizes the Bellman error:

$$L(\phi_i) = \mathbb{E} \left[\left(Q_{\phi_i}(s, a) - y \right)^2 \right] \quad \text{where} \quad y = r + \gamma \left[\min_j Q_{\bar{\phi}_j}(s', a') - \alpha \log \pi(a'|s') \right]$$

In code:

```
--8<-- "scripts/labs/sac_from_scratch.py:twin_q_loss"
```

Key mapping:

| Math | Code | Meaning |
|---------------------------|---------------------------------|---|
| $Q_{\phi_i}(s, a)$ | q1, q2 | Twin Q-network outputs |
| $\min_j Q_{\bar{\phi}_j}$ | torch.min(target_q1, target_q2) | Pessimistic target (reduces overestimation) |

| Math | Code | Meaning |
|-------------------|-------------------------------------|-------------------------|
| $\alpha \log \pi$ | <code>alpha * next_log_probs</code> | Entropy bonus in target |

!!! lab "Checkpoint" The initial Q-loss magnitude tells you the critic is learning. With random networks, Q-values and targets are both near zero, so the loss should be small initially:

```
```python
After creating q_network, target_q_network, policy, and a random batch:
q_loss, info = compute_q_loss(q_network, target_q_network, policy, batch)
print(f"Q1 loss: {info['q1_loss']:.4f}") # typically < 1.0 at init
print(f"Q1 mean: {info['q1_mean']:.4f}") # near 0 at init
print(f"Target Q mean: {info['target_q_mean']:.4f}") # near 0 at init
```
```

1.5.5 2.5.5 Actor Loss with Entropy

The policy update from Section 2.2 maximizes Q-values while maintaining entropy:

$$L(\theta) = \mathbb{E} \left[\alpha \log \pi_{\theta}(a|s) - \min_i Q_{\phi_i}(s, a) \right]$$

In code:

```
--8<-- "scripts/labs/sac_from_scratch.py:actor_loss"
```

Key insight: The actor wants high Q-values (good actions) *and* high entropy (exploration). The α parameter balances these goals.

!!! lab "Checkpoint" The actor loss should be positive at initialization (minimizing negative Q + alpha * log_pi):

```
```python
actor_loss, info = compute_actor_loss(policy, q_network, obs, alpha=0.2)
print(f"Actor loss: {info['actor_loss']:.4f}") # positive (minimizing)
print(f"Entropy: {info['entropy']:.4f}") # positive (H = -E[log pi])
```
```

1.5.6 2.5.6 Automatic Temperature Tuning

Instead of manually choosing the entropy coefficient α , SAC can learn it by targeting a desired entropy level:

$$L(\alpha) = \mathbb{E} \left[-\alpha (\log \pi(a|s) + \bar{\mathcal{H}}) \right]$$

where $\bar{\mathcal{H}} = -\dim(\mathcal{A})$ is the target entropy (Section 1.4).

```
--8<-- "scripts/labs/sac_from_scratch.py:temperature_loss"
```

!!! lab "Checkpoint" Alpha starts at 1.0 (since log_alpha is initialized to 0) and adjusts based on the gap between current entropy and target:

```
```python
log_alpha = torch.tensor(0.0, requires_grad=True)
print(f"Initial alpha: {log_alpha.exp().item():.3f}") # 1.000
```



```
After a few updates, alpha adjusts toward the target entropy
```
```

1.5.7 2.5.7 Wiring It Together: The SAC Update

The individual components above are combined into a single update step. SAC updates three components in sequence, then soft-updates the target networks:

1. Q-networks: minimize Bellman error
2. Policy: maximize Q-value + entropy
3. Temperature: maintain target entropy
4. Target networks: Polyak averaging $\bar{\phi} \leftarrow \tau\phi + (1 - \tau)\bar{\phi}$

```
--8<-- "scripts/labs/sac_from_scratch.py:sac_update"
```

!!! lab "Checkpoint" After 20 updates on random data, alpha should have changed and Q-values should have shifted:

```
```python
initial_alpha = log_alpha.exp().item()
for _ in range(20):
 info = sac_update(policy, q_network, target_q_network, log_alpha,
 policy_optimizer, q_optimizer, alpha_optimizer,
 batch, target_entropy)
print(f"Alpha: {initial_alpha:.4f} -> {info['alpha']:.4f}") # has changed
print(f"Q1 loss: {info['q1_loss']:.4f}") # finite
print(f"Actor loss: {info['actor_loss']:.4f}") # finite
```
```

1.5.8 2.5.8 Verify the Full Lab

Run the from-scratch implementation's sanity checks -- this exercises all the components above end-to-end:

```
bash docker/dev.sh python scripts/labs/sac_from_scratch.py --verify
```

Expected output:

```
=====
SAC From Scratch -- Verification
=====
Verifying Q-network...
  Q1 mean: X.XXXX                # near 0 at init
  Q2 mean: X.XXXX                # near 0 at init
[PASS] Q-network OK

Verifying policy...
  Actions in [-X.XX, X.XX]       # within [-1, 1]
  Log prob mean: -X.XXXX         # negative (probabilities < 1)
[PASS] Policy OK

Verifying SAC update...
  Q1 loss: X.XXXX                # finite
  Actor loss: X.XXXX             # finite
  Alpha: 1.0000 -> X.XXXX        # has changed
[PASS] SAC update OK
```

```
=====
[ALL PASS] SAC implementation verified
=====
```

This lab is **not** how we train policies -- that's what SB3 is for. The lab shows *what* SB3 is doing internally, with every tensor operation visible.

1.5.9 2.5.9 Verify vs SB3 (Optional)

SB3 is the scaling engine we use in the Run It track. This optional check confirms that our squashed Gaussian log-probability matches SB3's distribution implementation:

```
bash docker/dev.sh python scripts/labs/sac_from_scratch.py --compare-sb3
```

Expected output:

```
=====
SAC From Scratch -- SB3 Comparison
=====
```

```
Max abs log_prob diff: 2.055e-02
Tolerance (atol):      5.0e-02
```

```
[PASS] Our squashed Gaussian log_prob matches SB3
```

The ~ 0.02 nat difference comes from SB3's epsilon safety term: SB3 computes $\log(1 - a^2 + 1e-6)$ while our numerically stable formula computes the exact $\log(1 - \tanh^2(u))$. For non-saturated actions the formulas agree; the epsilon only matters when $|\tanh(u)| \rightarrow 1$.

1.5.10 2.5.10 Exercises: Modify and Observe

Exercise 2.5.1: Twin Q-Network Ablation

In `compute_q_loss()`, the target uses `torch.min(target_q1, target_q2)`. Try using just one Q-network:

```
# Change: target_q = torch.min(target_q1, target_q2)
# To:     target_q = target_q1 # Use only Q1
```

Question: What happens to `q1_mean` over training? Does it grow unbounded? This demonstrates why twin Q-networks matter.

Exercise 2.5.2: Entropy Coefficient Effect

In `verify_sac_update()`, observe how `alpha` changes:

```
bash docker/dev.sh python scripts/labs/sac_from_scratch.py --verify
```

Question: Alpha starts at 1.0 and trends toward ~ 0.99 . What does this tell you about the target entropy? Try changing `target_entropy = -float(act_dim)` to a smaller value and observe.

Exercise 2.5.3: Fixed vs Auto Alpha

Modify `sac_update()` to use a fixed alpha instead of auto-tuning:

```
# Replace: alpha = log_alpha.exp().item()
# With:    alpha = 0.2 # Fixed temperature
```

Question: With fixed `alpha=0.2`, how does training differ? When might you prefer fixed over auto-tuned temperature?

Exercise 2.5.4: Soft Target Updates

The tau parameter controls how quickly target networks track the main networks. In `sac_update()`:

Try: tau = 0.001 (slower), tau = 0.01 (faster), tau = 1.0 (hard update)

Question: What happens with tau=1.0 (hard update every step)? Why does slow tracking help stability?

1.5.11 2.5.11 Demo: SAC Solves Pendulum

The from-scratch implementation can actually solve a continuous control task. Run:

Quick demo (~30 seconds): shows learning but doesn't solve
`bash docker/dev.sh python scripts/labs/sac_from_scratch.py --demo`

Full demo (~8 minutes on CPU): actually solves Pendulum
`bash docker/dev.sh python scripts/labs/sac_from_scratch.py --demo --steps 50000`

Actual results (50k steps on CPU):

| Step | Avg Return | Alpha | What's Happening |
|-------|-------------|-------|-------------------------------------|
| 5000 | -843 | 0.415 | Learning started, still random-ish |
| 10000 | -190 | 0.197 | Solved! (threshold: -200) |
| 25000 | -141 | 0.048 | Refining swing-up behavior |
| 50000 | -134 | 0.017 | Stable, nearly deterministic policy |

Key observations: Automatic temperature tuning works as expected -- alpha drops from 0.415 to 0.017, meaning the policy transitions from exploratory (high entropy) to exploitative (low entropy) without manual intervention. The task is solved by step 10k (crossing the -200 threshold), after which additional training refines the policy further. Most importantly, this is the same algorithm SB3 uses; the from-scratch implementation matches SB3's SAC, with the difference being pedagogical clarity rather than algorithmic substance.

Exercise 2.5.5: Record a GIF of the Trained Policy

After training, you can record a GIF showing the learned behavior:

`bash docker/dev.sh python scripts/labs/sac_from_scratch.py --demo --steps 50000 --record`

This trains SAC on Pendulum for 50k steps and saves a GIF to `videos/sac_pendulum_demo.gif`.

Note: If running outside of `dev.sh` (e.g., in scripts or CI), use the full Docker command:

```
docker run --rm \
  -e MUJOCO_GL=egl \
  -e PYOPENGL_PLATFORM=egl \
  -e PYTHONUNBUFFERED=1 \
  -e HOME=/tmp \
  -e XDG_CACHE_HOME=/tmp/.cache \
  -e TORCH_HOME=/tmp/.cache/torch \
  -e TORCHINDUCTOR_CACHE_DIR=/tmp/.cache/torch_inductor \
  -e MPLCONFIGDIR=/tmp/.cache/matplotlib \
  -e USER=user \
  -e LOGNAME=user \
  -v "$PWD:/workspace" \
  -w /workspace \
  --gpus all \
```

```
--ipc=host \
robotics-rl:latest \
bash -c 'source .venv/bin/activate && python scripts/labs/sac_from_scratch.py --demo --s
```

Trained policy swinging up and balancing:

SAC Pendulum Demo

1.6 Part 3: WHAT -- Running the Experiment (Run It)

1.6.1 3.1 The One-Command Version

```
bash docker/dev.sh python scripts/ch03_sac_dense_reach.py all --seed 0
```

This runs training (with diagnostics), evaluation, and comparison to PPO.

For faster iteration (~2 minutes):

```
bash docker/dev.sh python scripts/ch03_sac_dense_reach.py train --total-steps 100000
```

1.6.2 3.2 What to Expect

Training progress milestones for SAC (note: `learning_starts=10000` means no training for first 10k steps):

| Timesteps | Success Rate | What's Happening |
|-----------|--------------|--|
| 0-10k | 0% | Collecting random data (no training yet) |
| 10k-50k | 0-5% | Training begins, Q-values learning |
| 50k-150k | 20-60% | Rapid improvement, entropy decreasing |
| 150k-400k | 80-99% | Policy converging |
| 400k-1M | 100% | Fine-tuning, entropy near zero |

Our test run achieved **100% success rate** after roughly 300k steps, with an **18.6mm average goal distance** (well within the 50mm threshold) and **~594 steps/second** throughput on NVIDIA GB10. The entropy coefficient dropped from 0.47 to 0.0004 over the course of training, meaning the policy became nearly deterministic as it converged.

1.6.3 3.3 What the Diagnostics Callback Logs

Our custom SACDiagnosticsCallback logs to TensorBoard:

| Metric | What It Means | Healthy Range |
|------------------------------|--------------------------|-------------------------------|
| replay/q1_mean, q2_mean | Average Q-values | Should stabilize, not explode |
| replay/q_min_mean | Min of two Q-values | Tracks actual value estimates |
| replay/ent_coef | Temperature α | Should decrease over training |
| replay/reward_mean | Average reward in buffer | Task-dependent |
| replay/goal_distance_mean | How far from goals | Should decrease |
| replay/goal_within_threshold | Fraction within success | Should increase |

1.6.4 3.4 What to Watch For

Q-Value Behavior: Healthy Q-values stabilize in a reasonable range (we observed $q_{min_mean} \sim -1.5$), while unhealthy Q-values grow unbounded (>100 and continuing to increase). Our run showed Q-values starting positive (~ 18 at 30k steps during random exploration) then settling to ~ -1.5 as the policy learned -- this is healthy, indicating the critic learned accurate value estimates.

Entropy Coefficient: A healthy entropy coefficient starts high and gradually decreases (we observed $0.47 \rightarrow 0.0004$), whereas an unhealthy one drops to near-zero within the first 10k steps (exploration collapsed). Our run showed gradual decrease over 1M steps, indicating the auto-tuning worked correctly. The final value of 0.0004 means the policy became nearly deterministic -- appropriate for a solved task.

Goal Distance (in replay buffer): Note that this metric shows the *average* distance across the entire buffer, not current policy performance. Our run showed $\sim 0.20m$ mean distance (a historical average including early random data). Current policy performance is better measured by rollout/success_rate and the evaluation metrics.

1.6.5 3.5 Throughput Scaling

Off-policy methods can benefit from parallel environments, but the relationship is complex (more envs = more data, but also more stale):

```
bash docker/dev.sh python scripts/ch03_sac_dense_reach.py throughput --n-envs-list 1,2,4,8
```

This measures steps/second for different `n_envs` values. Expect near-linear scaling up to some point, followed by diminishing returns as CPU becomes the bottleneck, with the sweet spot typically around 4-16 envs for SAC.

1.6.6 3.6 Actual Results

On FetchReachDense-v4 with 1M steps (NVIDIA GB10):

| Metric | PPO (Ch 2) | SAC (Actual) |
|-------------------|-----------------|----------------|
| Success Rate | 100% | 100% |
| Mean Return | -0.40 | -1.06 |
| Final Distance | 4.6mm | 18.6mm |
| Action Smoothness | 1.40 | 1.68 |
| Training Time | ~ 6 min | ~ 28 min |
| Throughput | ~ 1300 fps | ~ 594 fps |

Analysis: Both algorithms achieve 100% success, which validates the off-policy stack. SAC has a higher final distance (18.6mm vs 4.6mm) but still well within the 50mm success threshold, and it is slower due to more network updates (actor + 2 critics + 2 targets) per step. SAC's slightly less smooth actions reflect the entropy bonus encouraging exploration even late in training.

1.6.7 3.7 Understanding GPU Utilization

You may notice low GPU utilization ($\sim 5-10\%$) during training:

```
$ nvidia-smi dmon -s u
# gpu      sm      mem
   0        7      0
```

This is expected for RL training. The bottleneck is CPU, not GPU:

| Component | Runs On | Time Fraction |
|---------------------------|---------|---------------|
| MuJoCo physics simulation | CPU | ~60-70% |
| Environment step/reset | CPU | ~10-15% |
| Neural network forward | GPU | ~5-10% |
| Neural network backward | GPU | ~10-15% |
| Replay buffer operations | CPU | ~5% |

With small batch sizes (256) and simple MLPs, GPU operations complete in microseconds, so the GPU idles while waiting for CPU-bound simulation. This is why throughput (~600 fps) is limited by CPU, not GPU.

To increase GPU utilization, one could use larger batch sizes (though with diminishing returns for small networks), run multiple parallel training workers (which adds complexity), or switch to vision-based policies with CNNs (Week 8+). For our curriculum, ~600 fps is sufficient -- there is no reason to optimize prematurely.

1.6.8 3.8 Generating Demo Videos

To visualize your trained policy, use the video generation script:

```
# Generate videos from SAC checkpoint
bash docker/dev.sh python scripts/generate_demo_videos.py \
  --ckpt checkpoints/sac_FetchReachDense-v4_seed0.zip \
  --env FetchReachDense-v4 \
  --n-episodes 5 \
  --out videos/sac_reach_demo \
  --grid \
  --gif
```

This creates three output files: `videos/sac_reach_demo.mp4` (a compilation of all episodes), `videos/sac_reach_demo_grid.mp4` (a 2x2 grid showing 4 episodes simultaneously), and `videos/sac_reach_demo_grid.gif` (a GIF version for documentation or social media).

Video annotations: Each frame shows the real-time distance to target in centimeters, a green "TARGET REACHED!" indicator when within the 5cm threshold, and a legend (red dot for target position, green dot for gripper position).

The script automatically loads the trained model (works with PPO, SAC, or TD3), enlarges the target sphere for visibility, adds text overlays with goal distance, and creates both MP4 and GIF formats.

1.7 Part 4: Understanding What You Built

1.7.1 4.1 The Replay Buffer

The replay buffer is a circular array storing transitions:

```
[t_0, t_1, t_2, ..., t_{n-1}, t_n, t_{n+1}, ...]
  ^-- oldest                               ^-- newest
```

When full, new transitions overwrite the oldest, which means recent experience is always available while very old experience is eventually forgotten -- in effect, buffer size controls the "memory horizon."

For 1M buffer and 1M training steps, every transition is seen roughly once on average. Increase buffer size for longer memory; decrease for faster adaptation.

1.7.2 4.2 The Min-Q Trick

When computing the target for critic updates, SAC uses:

$$\min_{j=1,2} Q_{\bar{\phi}_j}(s', a')$$

This seems pessimistic -- why take the minimum? Because Q-learning systematically overestimates:

$$\mathbb{E}[\max(Q_1, Q_2)] \geq \max(\mathbb{E}[Q_1], \mathbb{E}[Q_2])$$

Taking the min counteracts this bias, leading to more stable training.

1.7.3 4.3 The Squashed Gaussian

SAC uses a "squashed Gaussian" for the policy: first sample from a Gaussian ($z \sim \mathcal{N}(\mu, \sigma^2)$), then squash through \tanh ($a = \tanh(z)$). This ensures actions are bounded (important for robotics with joint limits), and the log probability calculation accounts for the squashing via a Jacobian correction.

1.8 Part 5: Exercises

1.8.1 Exercise 3.1: Reproduce the Baseline

Run SAC on FetchReachDense-v4 and verify >95% success rate:

```
bash docker/dev.sh python scripts/ch03_sac_dense_reach.py all --seed 0
```

1.8.2 Exercise 3.2: Compare Learning Curves

Open TensorBoard and compare PPO vs SAC:

```
bash docker/dev.sh tensorboard --logdir runs --bind_all
```

Questions to answer:

- Which reaches high success faster?
- How do the value loss curves differ?
- What happens to SAC's entropy coefficient over training?

1.8.3 Exercise 3.3: Throughput Analysis

Run the throughput experiment:

```
bash docker/dev.sh python scripts/ch03_sac_dense_reach.py throughput
```

Plot steps/sec vs n_{envs} . Where does scaling stop being linear?

1.8.4 Exercise 3.4: Q-Value Analysis (Written)

Watch the Q-value metrics in TensorBoard. Answer:

1. Do Q1 and Q2 diverge significantly, or stay close?
2. What is the typical range of Q-values for this task?
3. How would you detect Q-value overestimation if it occurred?

1.8.5 Exercise 3.5: Temperature Ablation

Train with fixed entropy coefficients instead of auto-tuning:

```
# In train.py, modify ent_coef to a fixed value:  
# ent_coef=0.1 (high exploration)  
# ent_coef=0.01 (medium)  
# ent_coef=0.001 (low exploration)
```

How does fixed vs auto-tuned entropy affect:

- Final performance?
 - Training stability?
 - Time to convergence?
-

1.9 Part 6: Common Failures and Solutions

1.9.1 "Q-values explode (>1000)"

Symptom: replay/q_min_mean keeps growing without bound.

Cause: An overestimation feedback loop -- Q-targets use overestimated Q-values, which train the critic to overestimate further.

Solutions:

1. Verify rewards are bounded (FetchReachDense: [-1, 0])
2. Check discount factor (should be <1, typically 0.99)
3. Reduce learning rate
4. Verify target networks are updating (tau should be small)

1.9.2 "Entropy coefficient goes to zero immediately"

Symptom: replay/ent_coef drops to <0.01 within first 10k steps.

Cause: Target entropy too low, or policy collapsed to near-deterministic.

Solutions:

1. Increase target entropy (default is -dim(action))
2. Check action space bounds (if very wide, entropy scale changes)
3. Use fixed entropy coefficient initially to debug

1.9.3 "Success rate stalls below 50%"

Symptom: Learning plateaus well below PPO baseline.

Cause: Often insufficient exploration or replay buffer issues.

Solutions:

1. Check entropy coefficient isn't too low
2. Verify `learning_starts` isn't too high (not using early data)
3. Check batch size (too small = high variance, too large = slow updates)

1.9.4 "Training is much slower than PPO"

Expected: SAC is slower per-step (more networks to update).

Unexpected slowness causes:

1. Buffer operations (should be $O(1)$ for sampling)
 2. GPU memory (buffer can be large)
 3. Too many gradient steps per env step
-

1.10 Part 7: Looking Ahead

With SAC validated on dense Reach, we have confirmed that the replay buffer correctly stores and samples transitions, that dual critics work without divergence, that entropy tuning behaves as expected, and that off-policy learning achieves good performance on this task.

Chapter 4 introduces HER (Hindsight Experience Replay) for sparse rewards. The key insight is that a failed trajectory to goal g is a successful trajectory to whatever goal we actually reached -- HER relabels transitions to manufacture success signal from failure. This requires off-policy learning (SAC or TD3), since on-policy methods like PPO cannot use HER because they cannot reuse relabeled data. Our validated SAC stack is the foundation.

1.11 References

1. Haarnoja, T., Zhou, A., Abbeel, P., & Levine, S. (2018). Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. ICML.
2. Haarnoja, T., Zhou, A., Hartikainen, K., Tucker, G., Ha, S., Tan, J., Kumar, V., Zhu, H., Gupta, A., Abbeel, P., & Levine, S. (2018). Soft Actor-Critic Algorithms and Applications. arXiv:1812.05905.
3. Fujimoto, S., van Hoof, H., & Meger, D. (2018). Addressing Function Approximation Error in Actor-Critic Methods. ICML. (TD3 paper, introduces clipped double-Q)
4. Stable Baselines3 SAC Documentation: <https://stable-baselines3.readthedocs.io/en/master/modules/sac.html>