

Reinforcement Learning

Lecture 6: Deep Q-Networks (DQN)

Taehoon Kim

Sogang University MIMIC Lab
<https://mimic-lab.com>

Fall Semester 2025

Session Goals

By the end of this lecture, you will be able to:

- ➊ Explain why tabular Q-learning fails with high-dimensional state spaces
- ➋ Understand how DQN addresses divergence risks with experience replay and target networks
- ➌ Derive the DQN learning target and implement the Huber loss objective
- ➍ Implement a production-grade DQN agent in PyTorch with:
 - Replay buffer, target network, and ϵ -greedy exploration
 - Mixed precision (AMP) and `torch.compile()` acceleration
 - Checkpoint save/restore and TensorBoard logging
- ➎ Run ablations on hyperparameters and interpret results

Prerequisites: Q-learning (Lecture 5), PyTorch basics (Lecture 2)

Focus to kick things off

- Diagnose why tabular Q-learning breaks in large/continuous spaces
- Build up the DQN recipe: replay buffer, target network, loss
- Reason about stability tricks (Double DQN, Huber loss, AMP)

Flow

- 1 Revisiting tabular Q-learning pitfalls
- 2 Introducing neural function approximation
- 3 Deriving the DQN update and variants

Goal: Bridge tabular intuition to deep RL foundations before coding.

The Curse of Dimensionality

Tabular Q-Learning Limitations:

- State space explosion: CartPole with 10 bins/dimension $\rightarrow 10^4 = 10,000$ states
- Memory requirements: $O(|S| \times |A|)$
- No generalization between similar states
- Must visit every state-action pair

Real-world Examples:

- Atari games: 210×160 pixels \times 128 colors $\approx 10^{120,000}$ states
- Robotics: Continuous joint angles and velocities
- Autonomous driving: High-dimensional sensor inputs

Solution: Function Approximation with Neural Networks

Function Approximation in Q-Learning

From Table to Function:

Tabular Q-Learning:

- $Q : S \times A \rightarrow \mathbb{R}$
- Stored as table $Q[s, a]$
- Update: $Q[s, a] \leftarrow \text{target}$
- Exact values per state

Benefits of Neural Networks:

- Automatic feature extraction
- Generalization to similar states
- Handles high-dimensional inputs (images, etc.)
- Compact representation

Deep Q-Learning:

- $Q_\theta : S \times A \rightarrow \mathbb{R}$
- Neural network with parameters θ
- Update: $\theta \leftarrow \theta - \alpha \nabla_\theta L$
- Approximate values

Why Naive Neural Q-Learning Fails

Deadly Triad of Instability:

- 1 **Function Approximation:** Non-tabular representation
- 2 **Bootstrapping:** Using estimates to update estimates
- 3 **Off-policy Learning:** Learning from old experiences

Specific Problems:

- **Moving Targets:** $Q_{\theta}(s', a')$ changes as we update θ
- **Correlation:** Sequential samples are highly correlated
- **Feedback Loops:** Updates affect future targets
- **Overestimation:** Max operator causes positive bias

Result: Divergence, oscillation, or poor performance

DQN Innovation 1: Experience Replay

Breaking Correlation with Memory:

- Store transitions $(s, a, r, s', done)$ in replay buffer \mathcal{D}
- Sample random mini-batches for training
- Breaks temporal correlation
- Improves sample efficiency (reuse experiences)

Replay Buffer Implementation:

```
1 class ReplayBuffer:
2     def __init__(self, capacity):
3         self.buffer = deque(maxlen=capacity)
4
5     def push(self, state, action, reward, next_state, done):
6         self.buffer.append((state, action, reward, next_state, done))
7
8     def sample(self, batch_size):
9         return random.sample(self.buffer, batch_size)
```

Experience Replay: Why It Works

Without Replay:

- Samples: $s_t, s_{t+1}, s_{t+2}, \dots$
- High correlation
- Recent bias
- Catastrophic forgetting
- Unstable gradients

With Replay:

- Samples: $s_{17}, s_{203}, s_5, \dots$
- I.I.D.-like sampling
- Balanced experience
- Better coverage
- Stable gradients

Key Parameters:

- Buffer size: Typically 10^4 to 10^6 transitions
- Batch size: Usually 32-256
- Start learning after: 1000+ transitions (warmup)

DQN Innovation 2: Target Network

Stabilizing the Learning Target:

- Maintain two networks: Online Q_θ and Target Q_{θ^-}
- Use target network for computing TD targets
- Update target network periodically (every C steps)

The Key Insight:

$$\text{Without target network: } y = r + \gamma \max_{a'} Q_\theta(s', a') \quad (1)$$

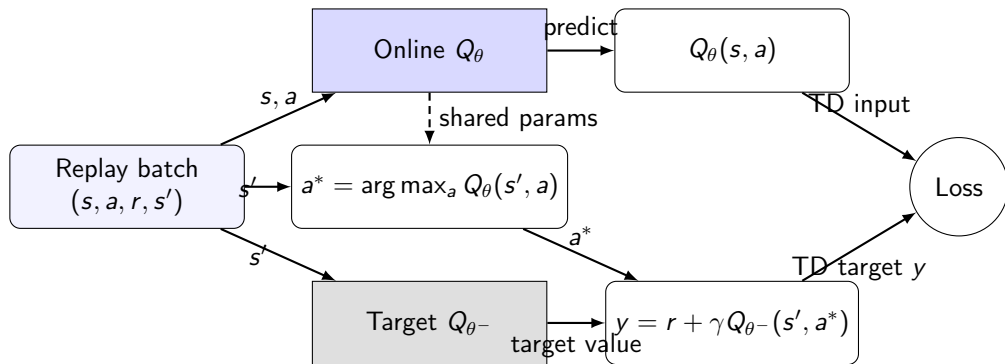
$$\text{With target network: } y = r + \gamma \max_{a'} Q_{\theta^-}(s', a') \quad (2)$$

Target θ^- remains fixed during updates, breaking harmful feedback loops!

Update Strategies:

- Hard update: $\theta^- \leftarrow \theta$ every C steps
- Soft update: $\theta^- \leftarrow \tau\theta + (1 - \tau)\theta^-$ each step

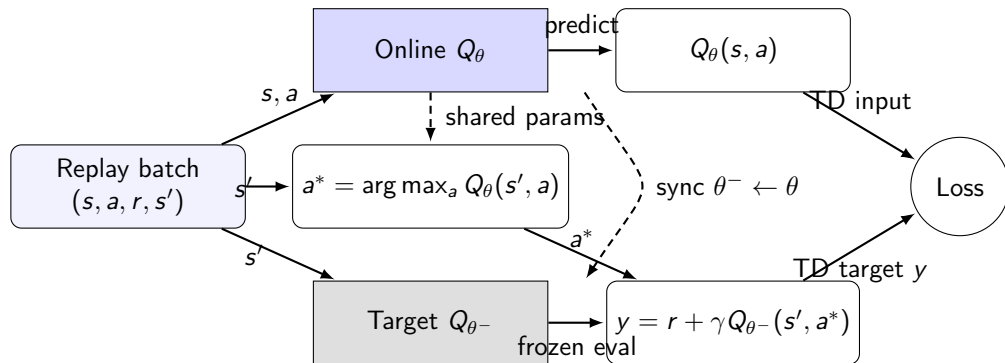
Target Network Flow Forward



Forward pass:

- 1 Sample a mini-batch from replay to evaluate $Q_\theta(s, a)$.
- 2 Use the online network to select a^* for each next state s' .
- 3 Combine a^* with the frozen target network to construct TD targets y .

Target Network Flow Updates



Updates and stability:

- 1 Backpropagate the loss through Q_θ only. Keep Q_{θ^-} frozen.
- 2 Periodically copy parameters $(\theta^- \leftarrow \theta)$ to refresh the target.

DQN Loss Function

The DQN Objective:

$$L(\theta) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[(y - Q_{\theta}(s, a))^2 \right] \quad (3)$$

where the target is:

$$y = r + \gamma(1 - \text{done}) \cdot \max_{a'} Q_{\theta-}(s', a') \quad (4)$$

Huber Loss (Smooth L1):

$$\mathcal{L}_{\delta}(x) = \begin{cases} \frac{1}{2}x^2 & \text{if } |x| \leq \delta \\ \delta(|x| - \frac{1}{2}\delta) & \text{otherwise} \end{cases} \quad (5)$$

Benefits: Less sensitive to outliers, prevents gradient explosion

The Complete DQN Algorithm

- 1: Initialize replay buffer \mathcal{D} with capacity N
- 2: Initialize Q-network Q_θ with random weights
- 3: Initialize target network Q_{θ^-} with $\theta^- = \theta$
- 4: **for** episode = 1 to M **do**
- 5: Initialize state s_1
- 6: **for** $t = 1$ to T **do**
- 7: Select action $a_t = \begin{cases} \text{random action} & \text{with probability } \epsilon \\ \arg \max_a Q_\theta(s_t, a) & \text{otherwise} \end{cases}$
- 8: Execute a_t , observe r_t, s_{t+1}
- 9: Store (s_t, a_t, r_t, s_{t+1}) in \mathcal{D}
- 10: Sample mini-batch from \mathcal{D}
- 11: Compute targets $y_i = r_i + \gamma \max_{a'} Q_{\theta^-}(s'_i, a')$
- 12: Update θ by minimizing $L = \frac{1}{|B|} \sum_i (y_i - Q_\theta(s_i, a_i))^2$
- 13: Every C steps: $\theta^- \leftarrow \theta$
- 14: **end for**
- 15: **end for**

Double DQN: Reducing Overestimation

The Overestimation Problem:

- Max operator in Q-learning causes positive bias
- $\mathbb{E}[\max_a Q(s, a)] \geq \max_a \mathbb{E}[Q(s, a)]$ (Jensen's inequality)
- Errors accumulate through bootstrapping

Double DQN Solution:

- Decouple action selection from evaluation
- Use online network to select actions
- Use target network to evaluate them

$$\text{DQN: } y = r + \gamma \max_{a'} Q_{\theta^-}(s', a') \quad (6)$$

$$\text{Double DQN: } y = r + \gamma Q_{\theta^-}(s', \arg \max_{a'} Q_{\theta}(s', a')) \quad (7)$$

Exploration Strategy: ϵ -Greedy

Balancing Exploration and Exploitation:

```
1 def select_action(state, epsilon):  
2     if random.random() < epsilon:  
3         return env.action_space.sample() # Explore  
4     else:  
5         q_values = q_network(state)  
6         return q_values.argmax() # Exploit
```

Epsilon Scheduling:

- Linear decay: $\epsilon_t = \epsilon_{start} - t \cdot \frac{\epsilon_{start} - \epsilon_{end}}{T}$
- Exponential decay: $\epsilon_t = \epsilon_{end} + (\epsilon_{start} - \epsilon_{end}) \cdot e^{-t/\tau}$
- Step decay: Reduce by factor at milestones

Typical values: $\epsilon_{start} = 1.0$, $\epsilon_{end} = 0.01$

DQN Hyperparameters

Parameter	Typical Value	Description
Learning rate	10^{-4} to 10^{-3}	Gradient descent step size
Batch size	32-256	Mini-batch size
Buffer size	10^4 to 10^6	Replay buffer capacity
Target update	1000-10000 steps	Hard update frequency
γ	0.99	Discount factor
ϵ_{start}	1.0	Initial exploration
ϵ_{end}	0.01-0.1	Final exploration
Hidden layers	[128, 128]	Network architecture
Optimizer	Adam	Gradient optimizer
Loss	Huber	Loss function

Common Pitfalls and Solutions

1 Insufficient Warmup

- Problem: Learning from small buffer
- Solution: Start after 1000+ transitions

2 Exploration Collapse

- Problem: ϵ decays too quickly
- Solution: Longer decay schedule

3 Gradient Explosion

- Problem: Unstable training
- Solution: Gradient clipping, Huber loss

4 Reward Scale

- Problem: Rewards too large/small
- Solution: Reward clipping or normalization

5 Target Update Frequency

- Problem: Too frequent or too rare
- Solution: Tune based on environment

DQN Variants and Extensions

Major DQN Improvements:

- **Double DQN (2015)**: Reduce overestimation bias
- **Prioritized Replay (2015)**: Sample important transitions more
- **Dueling DQN (2016)**: Separate value and advantage streams
- **Rainbow (2017)**: Combine all improvements
- **C51 (2017)**: Distributional Q-learning
- **Noisy Networks (2017)**: Parameter noise for exploration

Rainbow Components:

- 1 Double Q-learning
- 2 Prioritized replay
- 3 Dueling networks
- 4 Multi-step learning
- 5 Distributional RL
- 6 Noisy networks

When to Use DQN

Good fit for DQN:

- Discrete action spaces
- High-dimensional observations (images)
- Need sample efficiency
- Off-policy learning beneficial
- Deterministic environments

Consider alternatives when:

- Continuous actions → DDPG, TD3, SAC
- On-policy preferred → PPO, A2C
- Simple state space → Tabular Q-learning
- Safety critical → Conservative algorithms

Success stories: Atari games, resource allocation, trading

Key Equations for DQN

Bellman Optimality with Function Approximation:

$$Q^*(s, a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right] \quad (8)$$

Semi-gradient Update:

$$\theta_{t+1} = \theta_t + \alpha [y_t - Q_\theta(s_t, a_t)] \nabla_\theta Q_\theta(s_t, a_t) \quad (9)$$

Convergence Conditions (Tabular):

- All state-action pairs visited infinitely
- Learning rate satisfies Robbins-Monro conditions
- $\sum_t \alpha_t = \infty, \sum_t \alpha_t^2 < \infty$

Note: Neural Q-learning convergence not guaranteed!

Efficient Replay Buffer

```
1 class ReplayBuffer:
2     def __init__(self, capacity, obs_dim):
3         # Pre-allocate arrays for efficiency
4         self.observations = np.zeros((capacity, obs_dim),
5                                     dtype=np.float32)
6         self.actions = np.zeros(capacity, dtype=np.int64)
7         self.rewards = np.zeros(capacity, dtype=np.float32)
8         self.next_observations = np.zeros((capacity, obs_dim),
9                                          dtype=np.float32)
10        self.dones = np.zeros(capacity, dtype=np.float32)
11
12        self.position = 0
13        self.size = 0
14        self.capacity = capacity
15
16    def push(self, obs, action, reward, next_obs, done):
17        idx = self.position
18        self.observations[idx] = obs
19        self.actions[idx] = action
20        # ... store other components
21        self.position = (self.position + 1) % self.capacity
22        self.size = min(self.size + 1, self.capacity)
```

Q-Network Implementation

```
1 class QNetwork(nn.Module):
2     def __init__(self, obs_dim, n_actions,
3                 hidden_sizes=(128, 128)):
4         super().__init__()
5
6         layers = []
7         input_size = obs_dim
8
9         # Build hidden layers
10        for hidden_size in hidden_sizes:
11            layers.append(nn.Linear(input_size, hidden_size))
12            layers.append(nn.ReLU())
13            input_size = hidden_size
14
15        # Output layer (Q-values for each action)
16        layers.append(nn.Linear(input_size, n_actions))
17
18        self.network = nn.Sequential(*layers)
19
20    def forward(self, x):
21        return self.network(x) # Shape: [batch_size, n_actions]
```

Target Network Mechanism

```
1 def hard_update(target_net, online_net):
2     """Copy weights from online to target network"""
3     target_net.load_state_dict(online_net.state_dict())
4
5 def soft_update(target_net, online_net, tau=0.005):
6     """Polyak averaging update"""
7     with torch.no_grad():
8         for target_param, param in zip(
9             target_net.parameters(),
10            online_net.parameters()
11        ):
12            target_param.data.copy_(
13                tau * param.data + (1 - tau) * target_param.data
14            )
15
16 # Usage in training loop
17 if step % target_update_freq == 0:
18     hard_update(target_network, q_network)
19 # OR for soft updates:
20 soft_update(target_network, q_network, tau=0.005)
```

DQN Loss Computation

```
1 def compute_dqn_loss(batch, q_network, target_network,
2                       gamma=0.99):
3     obs, actions, rewards, next_obs, dones = batch
4
5     # Current Q-values for taken actions
6     q_values = q_network(obs) # [B, n_actions]
7     q_values = q_values.gather(1, actions.unsqueeze(1))
8     q_values = q_values.squeeze() # [B]
9
10    # Compute targets with target network
11    with torch.no_grad():
12        next_q_values = target_network(next_obs) # [B, n_actions]
13        next_q_max = next_q_values.max(1)[0] # [B]
14
15    # TD targets
16    targets = rewards + gamma * (1 - dones) * next_q_max
17
18    # Huber loss
19    loss = F.huber_loss(q_values, targets)
20    return loss
```


Double DQN: Reducing Overestimation

```
1 def compute_double_dqn_loss(batch, q_network, target_network,
2                             gamma=0.99):
3     obs, actions, rewards, next_obs, dones = batch
4
5     # Current Q-values
6     q_values = q_network(obs).gather(1, actions.unsqueeze(1))
7     q_values = q_values.squeeze()
8
9     with torch.no_grad():
10         # Action selection with online network
11         next_q_online = q_network(next_obs)
12         next_actions = next_q_online.argmax(1, keepdim=True)
13
14         # Action evaluation with target network
15         next_q_target = target_network(next_obs)
16         next_q_values = next_q_target.gather(1, next_actions)
17         next_q_values = next_q_values.squeeze()
18
19         targets = rewards + gamma * (1 - dones) * next_q_values
20
21     return F.huber_loss(q_values, targets)
```

Main Training Loop

```
1 def train_dqn(env, n_episodes=1000):
2     q_network = QNetwork(obs_dim, n_actions).to(device)
3     target_network = QNetwork(obs_dim, n_actions).to(device)
4     target_network.load_state_dict(q_network.state_dict())
5
6     optimizer = optim.Adam(q_network.parameters(), lr=1e-3)
7     buffer = ReplayBuffer(capacity=10000, obs_dim=obs_dim)
8
9     for episode in range(n_episodes):
10         obs, _ = env.reset()
11         episode_reward = 0
12
13         while not done:
14             # Select action (epsilon-greedy)
15             action = select_action(obs, epsilon)
16             next_obs, reward, done, _, _ = env.step(action)
17
18             # Store and learn
19             buffer.push(obs, action, reward, next_obs, done)
20             if len(buffer) >= batch_size:
21                 batch = buffer.sample(batch_size)
22                 loss = compute_dqn_loss(batch, q_network,
23                                         target_network)
24
25                 # ... optimize
```

TensorBoard Integration

```
1 from torch.utils.tensorboard import SummaryWriter
2
3 writer = SummaryWriter('runs/dqn_experiment')
4
5 # Log scalars
6 writer.add_scalar('loss/td_error', loss.item(), step)
7 writer.add_scalar('metrics/episode_reward', reward, episode)
8 writer.add_scalar('metrics/epsilon', epsilon, step)
9
10 # Log hyperparameters and metrics
11 writer.add_hparams(
12     {'lr': 1e-3, 'batch_size': 32, 'gamma': 0.99},
13     {'final_reward': best_reward}
14 )
15
16 # Visualize in terminal:
17 # tensorboard --logdir runs/
```

Key Metrics to Track: Loss, episode rewards, Q-values, epsilon, learning rate

Saving and Loading Models

```
1 def save_checkpoint(agent, filepath):
2     checkpoint = {
3         'q_network': agent.q_network.state_dict(),
4         'target_network': agent.target_network.state_dict(),
5         'optimizer': agent.optimizer.state_dict(),
6         'episode': agent.episode,
7         'epsilon': agent.epsilon,
8         'replay_buffer': agent.buffer, # Optional
9         'config': agent.config
10    }
11    torch.save(checkpoint, filepath)
12
13 def load_checkpoint(filepath):
14     checkpoint = torch.load(filepath, map_location=device)
15
16     agent.q_network.load_state_dict(checkpoint['q_network'])
17     agent.target_network.load_state_dict(
18         checkpoint['target_network'])
19     agent.optimizer.load_state_dict(checkpoint['optimizer'])
20     # ... restore other components
21
22     return checkpoint['episode']
```

Proper Evaluation

```
1 def evaluate(agent, env, n_episodes=10):
2     """Evaluate agent without exploration"""
3     eval_rewards = []
4
5     for episode in range(n_episodes):
6         obs, _ = env.reset(seed=seed + episode)
7         episode_reward = 0
8         done = False
9
10        while not done:
11            # Greedy action selection (no exploration)
12            with torch.no_grad():
13                q_values = agent.q_network(
14                    torch.FloatTensor(obs).unsqueeze(0))
15                action = q_values.argmax().item()
16
17            obs, reward, terminated, truncated, _ = env.step(action)
18            done = terminated or truncated
19            episode_reward += reward
20
21        eval_rewards.append(episode_reward)
22
23    return np.mean(eval_rewards), np.std(eval_rewards)
```

Hyperparameter Tuning Strategy

Grid Search Example:

```
1 param_grid = {  
2     'lr': [1e-3, 3e-4, 1e-4],  
3     'batch_size': [32, 64, 128],  
4     'target_update': [100, 500, 1000],  
5     'buffer_size': [10000, 50000]  
6 }
```

Important Relationships:

- Larger buffer \rightarrow more stable but slower
- Frequent target updates \rightarrow less stable
- Larger batch \rightarrow more stable gradients
- Higher $\gamma \rightarrow$ longer-term planning

Start with:

- Published hyperparameters for similar tasks
- Conservative values (small LR, large buffer)
- Tune one parameter at a time

Debugging Checklist

Common Issues and Solutions:

1 Q-values exploding

- Check reward scale
- Verify target network updates
- Use gradient clipping

2 No learning progress

- Increase exploration (ϵ)
- Check replay buffer filling
- Verify loss computation

3 Unstable training

- Reduce learning rate
- Increase batch size
- Use Huber loss instead of MSE

4 Poor final performance

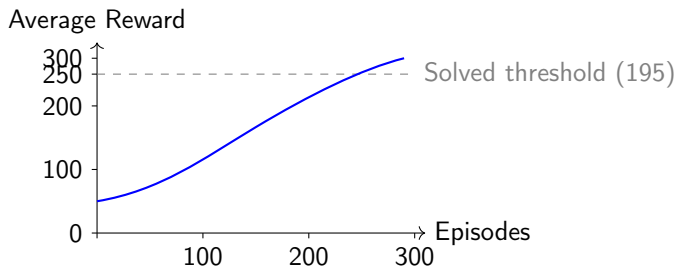
- Tune exploration schedule
- Increase network capacity
- Try Double DQN

Expected Results: CartPole-v1

Training Progression:

- Episodes 1-50: Random behavior (reward ≈ 20)
- Episodes 50-100: Learning begins (reward ≈ 50 -100)
- Episodes 100-150: Rapid improvement (reward ≈ 150 -200)
- Episodes 150+: Solved (reward = 500, maximum)

Typical Training Curve:



Component Impact Analysis:

Configuration	Episodes to Solve	Final Reward
Full DQN	150	500
No target network	Diverges	N/A
No replay buffer	300+	200
Small buffer (1000)	250	450
No exploration decay	400+	300
MSE loss (not Huber)	200	480
Double DQN	130	500

Key Insights:

- Target network is essential for stability
- Replay buffer significantly improves efficiency
- Proper exploration schedule crucial
- Double DQN provides modest improvement

Real-world DQN Applications

Game Playing:

- Atari 2600 games (original DQN paper)
- StarCraft II micro-management
- Poker and card games

Robotics:

- Robotic grasping
- Navigation in discrete spaces
- Task scheduling

Resource Management:

- Data center cooling (Google)
- Network routing
- Inventory management

Finance:

- Portfolio optimization
- Trading strategies
- Risk management

Performance Optimization Tips

Memory Efficiency:

- Use numpy arrays in replay buffer
- Store observations as uint8 if possible
- Implement circular buffer
- Clear gradients with `set_to_none=True`

Computation Speed:

- Batch environment steps (vectorized envs)
- Use `torch.compile()` for JIT optimization
- Enable AMP on compatible GPUs
- Profile with `torch.profiler`

Training Stability:

- Normalize observations
- Clip rewards to `[-1, 1]` range
- Use learning rate scheduling
- Monitor gradient norms

Key Mathematical Results

TD Error:

$$\delta_t = r_t + \gamma \max_{a'} Q_{\theta^-}(s_{t+1}, a') - Q_{\theta}(s_t, a_t) \quad (10)$$

Gradient Update:

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} L(\theta_t) \quad (11)$$

Loss Gradient:

$$\nabla_{\theta} L = -\mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} [\delta_t \nabla_{\theta} Q_{\theta}(s, a)] \quad (12)$$

Optimal Q-function satisfies:

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) \quad \forall s, a \quad (13)$$

Summary

Key Theoretical Contributions:

- ➊ **Problem:** Q-learning fails with function approximation
- ➋ **Solution 1:** Experience replay breaks correlation
- ➌ **Solution 2:** Target network stabilizes learning
- ➍ **Enhancement:** Double DQN reduces overestimation

Mathematical Foundation:

- Semi-gradient TD learning
- Bellman optimality with approximation
- Convergence not guaranteed but works in practice

Critical Insights:

- Deadly triad: FA + bootstrapping + off-policy
- Replay enables i.i.d. sampling assumption
- Target network breaks feedback loops