

Reinforcement Learning

Lecture 1: Course Overview and Environment Setup

Taehoon Kim

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

Fall Semester 2025

Learning Objectives

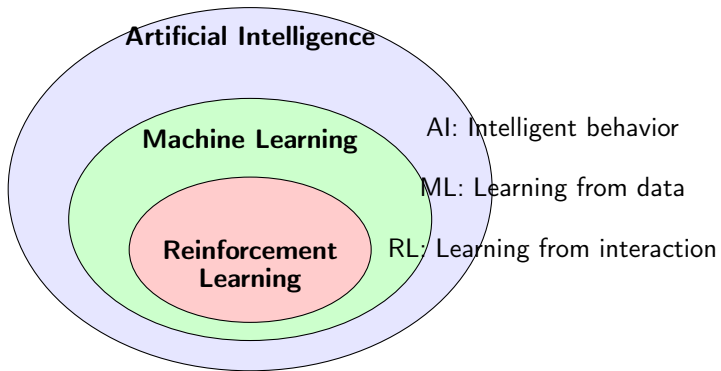
By the end of this lecture, you will:

- Understand the relationship among AI, ML, and RL
- Master the MDP formalism and core RL notation
- Set up a reproducible PyTorch 2.x environment
- Implement the standard code header for the course
- Complete 9 hands-on experiments
- Pass the integrated smoke test

Prerequisites

- Python programming experience
- Basic linear algebra and calculus
- Familiarity with neural networks (helpful)

The AI-ML-RL Hierarchy

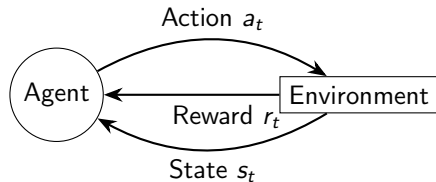


What is Reinforcement Learning?

Definition

RL learns optimal behavior through **trial and error** interaction with an environment

- Agent takes actions
- Environment provides rewards
- Goal: maximize cumulative reward
- No explicit supervision



RL vs Supervised Learning

Aspect	Supervised	Reinforcement
Feedback	Immediate labels	Delayed rewards
Data	i.i.d. samples	Sequential, correlated
Exploration	Not needed	Essential
Goal	Minimize error	Maximize return
Training	Offline, batch	Online, interactive

Key Insight

RL faces the **exploration-exploitation dilemma**: Should the agent try new actions (explore) or stick with known good actions (exploit)?

13-Week Course Structure

Foundations (Weeks 1-4)

- Week 1: Environment Setup
- Week 2: Deep Learning Essentials
- Week 3: RL Fundamentals
- Week 4: Mathematical Foundations

Value-Based (Weeks 5-7)

- Week 5: Q-Learning
- Week 6: Deep Q-Networks
- Week 7: DQN Project

Policy-Based (Weeks 8-10)

- Week 8: Policy Gradients
- Week 9: Actor-Critic Methods
- Week 10: PPO

Advanced (Weeks 11-13)

- Week 11: Current Trends
- Week 12: Project Development
- Week 13: Final Presentations

Required Tools and Resources

Software

- Python 3.10-3.12
- PyTorch 2.x
- Gymnasium
- TensorBoard
- Jupyter/Colab
- Git

Hardware

- CPU: Any modern processor
- GPU: Optional but recommended
- RAM: 8GB minimum
- Storage: 20GB free space

Cloud Alternative

Google Colab provides free GPU access - all experiments will run there

Understanding the MDP Framework

Topics

- Markov Decision Processes (MDPs)
- States, Actions, and Rewards
- Policies and Value Functions
- Bellman Equations
- Optimality Conditions

Markov Decision Process (MDP)

An MDP is defined as a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$:

- \mathcal{S} : State space
- \mathcal{A} : Action space
- $\mathcal{P}(s'|s, a)$: Transition probability
- $r(s, a)$: Reward function
- $\gamma \in [0, 1)$: Discount factor

Markov Property

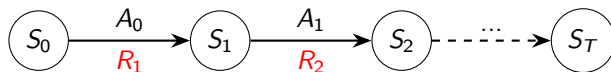
The future depends only on the current state, not the history:

$$P(S_{t+1}|S_t, A_t, S_{t-1}, A_{t-1}, \dots) = P(S_{t+1}|S_t, A_t)$$

Episode Structure

An episode is a sequence of interactions:

$$(S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T, S_T)$$



Terminal state S_T ends the episode

Return and Discounting

The **return** G_t is the cumulative discounted reward:

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Why Discounting?

- Mathematical convenience (convergence)
- Uncertainty about the future
- Preference for immediate rewards

Example

If $\gamma = 0.9$ and rewards are $[1, 2, 3, \dots]$:

$$G_0 = 1 + 0.9 \cdot 2 + 0.81 \cdot 3 + \dots = 10$$

Policy

A **policy** π defines the agent's behavior:

$$\pi(a|s) = P(A_t = a | S_t = s)$$

Deterministic Policy

$$a = \pi(s)$$

One action per state

Stochastic Policy

$$\pi(a|s) \in [0, 1]$$

Probability distribution over actions

Goal

Find the optimal policy π^* that maximizes expected return:

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \right]$$

Value Functions

State Value Function

Expected return starting from state s following policy π :

$$v^{\pi}(s) = \mathbb{E}_{\pi}[G_t \mid S_t = s]$$

Action Value Function (Q-function)

Expected return starting from state s , taking action a , then following π :

$$q^{\pi}(s, a) = \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a]$$

Relationship:

The policy $\pi(a|s)$ is a probability distribution over actions given state s . Therefore, the state value is the expected action value under this distribution:

$$v^{\pi}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \cdot q^{\pi}(s, a)$$

Bellman Equations

Value functions satisfy recursive relationships:

Bellman Expectation Equation

$$v^{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} \mathcal{P}(s'|s, a) [r(s, a) + \gamma v^{\pi}(s')]$$

Bellman Optimality Equation

$$v^*(s) = \max_a \sum_{s'} \mathcal{P}(s'|s, a) [r(s, a) + \gamma v^*(s')]$$

These equations are the foundation for RL algorithms!

Bellman Equations: Detailed Explanation

Bellman Expectation Equation

- Under a given policy π , the value of a state s is
- the weighted sum of possible actions, where weights are given by $\pi(a|s)$.
- Each action leads probabilistically to a next state s' according to $\mathcal{P}(s'|s, a)$, yielding an immediate reward $r(s, a)$.
- Therefore, $v^\pi(s)$ is the expected immediate reward plus the discounted value of the next state.

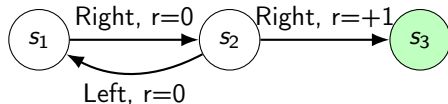
Bellman Optimality Equation

- For the optimal policy, the value of state s is
- the maximum expected return achievable over all possible actions.
- Instead of averaging with $\pi(a|s)$, we take \max_a .
- Solving this recursive relation gives the optimal value function $v^*(s)$ and the optimal policy.

Key idea: The Bellman equations express the value of a state as a recursive relationship between immediate reward and future value.

Bellman: Expectation vs Optimality (GridWorld)

Expectation (policy evaluation)



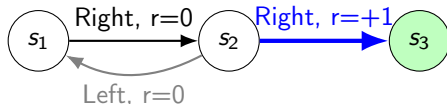
$$\pi(\cdot | s_2): \pi(\text{Right} | s_2) = 0.7, \pi(\text{Left} | s_2) = 0.3$$

$$\gamma = 0.9, \quad v^\pi(s_3) = 0 \text{ (terminal state)}$$

$$v^\pi(s_2) = 0.7 [1 + \gamma v^\pi(s_3)] + 0.3 [0 + \gamma v^\pi(s_1)]$$

$$v^\pi(s_1) = 0 + \gamma v^\pi(s_2)$$

Optimality (control)



$$\gamma = 0.9, \quad v^*(s_3) = 0$$

$$v^*(s_2) = \max\{ 1 + \gamma v^*(s_3), \gamma v^*(s_1) \}$$

$$v^*(s_1) = \gamma v^*(s_2)$$

Greedy policy at s_2 : choose Right

Expectation averages action values using $\pi(a|s)$, while Optimality takes the maximum over actions.

Optimal Policy and Value Functions

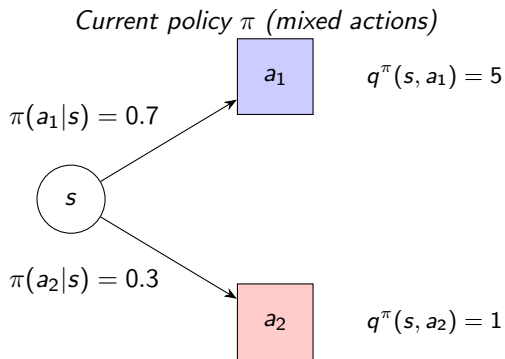
- Optimal value function: $v^*(s) = \max_{\pi} v^{\pi}(s)$
- Optimal Q-function: $q^*(s, a) = \max_{\pi} q^{\pi}(s, a)$
- Optimal policy: $\pi^*(a|s) = \arg \max_a q^*(s, a)$

Theorem (Policy Improvement)

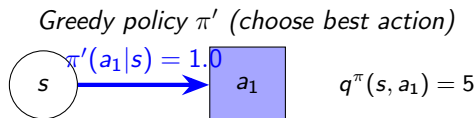
For any policy π , the greedy policy with respect to v^{π} is at least as good as π

This leads to **policy iteration** and **value iteration** algorithms

Policy Improvement Intuition



$$v^\pi(s) = 0.7 \times 5 + 0.3 \times 1 = 3.8$$



$$v^{\pi'}(s) = \max_a q^\pi(s, a) = 5$$

Greedy policy always achieves at least as high value as the original policy.

Types of RL Problems

Dimension	Types	Examples
State space	Discrete/Continuous	Grid/Robot control
Action space	Discrete/Continuous	Chess/Driving
Observation	Full/Partial	Go/Poker
Model	Model-based/free	Planning/Q-learning
Policy	On-policy/Off-policy	SARSA/Q-learning

This Course Focus

- Start with discrete spaces (tabular methods)
- Move to continuous (function approximation)
- Both model-free and model-based approaches

Mathematical Prerequisites

Linear Algebra

- Vector operations
- Matrix multiplication
- Eigenvalues (optional)

Calculus

- Derivatives
- Chain rule
- Gradients

Probability

- Expectations
- Conditional probability
- Distributions

Optimization

- Gradient descent
- Convexity (optional)
- Convergence

Building Your RL Development Environment

Components to Install

- Python environment (Anaconda/Miniconda)
- PyTorch 2.x with CUDA support
- Essential libraries
- Reproducibility tools
- Version control (Git)

Python Environment Setup

```
1 # Create conda environment
2 conda create -n rl2025 python=3.10
3 conda activate rl2025
4
5 # Install PyTorch (with CUDA 11.8)
6 conda install pytorch torchvision torchaudio \
7     pytorch-cuda=11.8 -c pytorch -c nvidia
8
9 # Install essential packages
10 pip install numpy matplotlib pandas tqdm
11 pip install tensorboard jupyterlab
12 # gymnasium will be installed in Lecture 3
```

Important

Python 3.10-3.12 required for compatibility

Device Detection Logic

```
1  # Proper device selection (CUDA > MPS > CPU)
2  device = torch.device(
3      'cuda' if torch.cuda.is_available()
4      else 'mps' if hasattr(torch.backends, 'mps')
5                      and torch.backends.mps.is_available()
6      else 'cpu'
7  )
8
9  print(f"Using device: {device}")
10
11 # Check CUDA details if available
12 if torch.cuda.is_available():
13     print(f"GPU: {torch.cuda.get_device_name(0)}")
14     print(f"CUDA: {torch.version.cuda}")
```

Ensuring Reproducibility

```
1 def setup_seed(seed=42):
2     random.seed(seed)
3     np.random.seed(seed)
4     torch.manual_seed(seed)
5     if torch.cuda.is_available():
6         torch.cuda.manual_seed_all(seed)
7
8     # Deterministic algorithms
9     torch.use_deterministic_algorithms(True)
10    torch.backends.cudnn.benchmark = False
11    torch.backends.cudnn.deterministic = True
12
13    # Always call at start of experiments
14    setup_seed(42)
```

Critical for reproducing results!

Recommended Project Structure

```
1  rl2025/  
2  |-- envs/           # Environment configs  
3  |   |-- environment.yml  
4  |   +-- requirements.txt  
5  |-- experiments/    # Experiment scripts  
6  |   +-- exp01_setup.py  
7  |-- runs/           # Logs and checkpoints  
8  |   |-- checkpoints/  
9  |   +-- tensorboard/  
10 |-- notebooks/      # Jupyter notebooks  
11 +-- README.md
```

Keep code, data, and results organized!

Automatic Mixed Precision (AMP)

What is AMP?

Training with mixed float16/float32 precision for:

- 2-3x speedup on modern GPUs
- 50% memory reduction
- Maintained accuracy

Benefits

- Larger batch sizes
- Faster training
- More complex models

Requirements

- CUDA-capable GPU
- PyTorch 1.6+
- Volta architecture or newer

AMP Implementation

```
1 from torch.cuda.amp import autocast, GradScaler
2
3 # Initialize scaler
4 scaler = GradScaler()
5
6 # Training step with AMP
7 optimizer.zero_grad()
8
9 # Forward pass with autocast
10 with autocast():
11     output = model(input) # FP16 computation
12     loss = criterion(output, target)
13
14 # Backward pass with scaling
15 scaler.scale(loss).backward()
16 scaler.step(optimizer)
17 scaler.update()
```

PyTorch 2.x Compilation

```
1  # Compile model for optimized execution
2  model = torch.compile(model, mode='default')
3
4  # Different compilation modes:
5  # 'default': Balanced optimization
6  # 'reduce-overhead': Minimize kernel launches
7  # 'max-autotune': Maximum performance
8
9  # Fallback for older PyTorch
10 def compile_if_available(module):
11     if hasattr(torch, 'compile'):
12         return torch.compile(module)
13     return module
```

Up to 2x speedup with torch.compile!

TensorBoard Integration

```
1 from torch.utils.tensorboard import SummaryWriter
2
3 # Initialize writer
4 writer = SummaryWriter('runs/experiment_1')
5
6 # Log scalars
7 writer.add_scalar('loss/train', loss, step)
8
9 # Log histograms
10 writer.add_histogram('weights', model.fc.weight, step)
11
12 # Log model graph
13 writer.add_graph(model, sample_input)
14
15 # Close when done
16 writer.close()
```

View with: `tensorboard -logdir runs`

Checkpoint Management

```
1 # Save checkpoint
2 checkpoint = {
3     'epoch': epoch,
4     'model': model.state_dict(),
5     'optimizer': optimizer.state_dict(),
6     'loss': loss,
7     'rng_states': {
8         'torch': torch.get_rng_state(),
9         'cuda': torch.cuda.get_rng_state_all()
10    }
11 }
12 torch.save(checkpoint, 'checkpoint.pt')
13
14 # Load checkpoint
15 checkpoint = torch.load('checkpoint.pt')
16 model.load_state_dict(checkpoint['model'])
17 optimizer.load_state_dict(checkpoint['optimizer'])
```

Version Control for RL

```
1 # Initialize repository
2 git init
3 git config user.name "Your Name"
4 git config user.email "email@example.com"
5
6 # Create .gitignore
7 echo "runs/" >> .gitignore
8 echo "__pycache__/" >> .gitignore
9 echo "*.pt" >> .gitignore
10
11 # Track experiment
12 git add experiment.py
13 git commit -m "Experiment: DQN baseline
14     Config: lr=0.001, batch=32
15     Result: 195.3 avg reward"
```

Google Colab Setup

```
1 # Colab bootstrap cell
2 import sys
3 IN_COLAB = 'google.colab' in sys.modules
4
5 if IN_COLAB:
6     # Install packages
7     !pip install -q torch tensorboard
8
9 # Mount Google Drive
10 if IN_COLAB:
11     from google.colab import drive
12     drive.mount('/content/drive')
13
14 # Check GPU
15 !nvidia-smi # Should show Tesla T4 or better
```

Free GPU access for experiments!

Unified Starting Point for All Experiments

Components

- Reproducibility (seeds)
- Device management
- AMP support
- Logging utilities
- Checkpoint handling
- Common RL functions

All experiments will import from this header!

Testing Your Setup

```
1  # Run integrated test
2  python exp09_integrated_test.py
3
4  # Expected output:
5  # =====
6  # Test 1: Environment Setup      [PASS]
7  # Test 2: Reproducibility       [PASS]
8  # Test 3: Model Training        [PASS]
9  # Test 4: DQN Components        [PASS]
10 # Test 5: Checkpointing         [PASS]
11 # Test 6: Logging               [PASS]
12 # =====
13 # All tests passed!
```

9 Progressive Experiments

- 1 Environment verification (exp01)
- 2 PyTorch basics (exp02)
- 3 Reproducibility (exp03)
- 4 AMP benchmarks (exp04)
- 5 Standard header (exp05)
- 6 Logging setup (exp06)
- 7 Checkpointing (exp07)
- 8 Git integration (exp08)
- 9 Integration test (exp09)

Experiment 1: Environment Verification

Goal: Verify Python and package installations

Tasks:

- Check Python version (3.10-3.12)
- Verify PyTorch installation
- List installed packages
- Create environment files
- Save system information

Run: `python exp01_setup.py`

Exp1: Key Code

```
1 def check_python_version():
2     version = sys.version_info
3     if not (3, 10) <= (version.major, version.minor) <= (3, 12):
4         print("Warning: Python 3.10-3.12 recommended")
5         return False
6     return True
7
8 def check_package_installations():
9     required = ['torch', 'numpy', 'matplotlib']
10    for package in required:
11        try:
12            __import__(package)
13            print(f"[OK] {package}")
14        except ImportError:
15            print(f"[MISSING] {package}")
```

Experiment 2: PyTorch Basics

Goal: Master PyTorch fundamentals and device management

Tasks:

- Device detection (CUDA > MPS > CPU)
- Tensor operations
- Automatic differentiation
- Performance benchmarking

Key Learning: Proper device selection is critical

Exp2: Device Selection

```
1 def get_device():
2     if torch.cuda.is_available():
3         device = torch.device('cuda')
4         print(f"Using CUDA: {torch.cuda.get_device_name(0)}")
5     elif hasattr(torch.backends, 'mps') and \
6         torch.backends.mps.is_available():
7         device = torch.device('mps')
8         print("Using MPS (Apple Silicon)")
9     else:
10        device = torch.device('cpu')
11        print("Using CPU")
12    return device
```

Experiment 3: Reproducibility

Goal: Ensure experiments are reproducible

Tasks:

- Set seeds for all RNGs
- Test reproducibility
- Handle DataLoader workers
- Save RNG states

Critical: Same seed → Same results

Exp3: Complete Seeding

```
1 def setup_seed(seed=42, deterministic=True):
2     # Python RNG
3     random.seed(seed)
4     # NumPy RNG
5     np.random.seed(seed)
6     # PyTorch RNG
7     torch.manual_seed(seed)
8     # CUDA RNG
9     if torch.cuda.is_available():
10         torch.cuda.manual_seed_all(seed)
11     # Deterministic mode
12     if deterministic:
13         torch.use_deterministic_algorithms(True)
```

Experiment 4: AMP and Compilation

Goal: Benchmark performance optimizations

Configurations tested:

- Baseline (FP32, no compile)
- AMP only (FP16/BF16)
- Compile only (torch.compile)
- AMP + Compile

Expected speedup: 2-4x on GPU

Exp4: Benchmark Results

Configuration	Time (ms/step)	Speedup
Baseline (FP32)	100	1.0x
AMP only	60	1.7x
Compile only	55	1.8x
AMP + Compile	35	2.9x

Combining AMP with compilation gives best performance!

Experiment 5: Standard Code Header

Goal: Implement reusable components

Components:

- Seeding functions
- Device management
- AMP context manager
- DQN training step
- Policy evaluation
- Model compilation

This becomes your toolkit for the course!

Exp5: DQN Training Step

```
1 def dqn_td_step(q_net, target_q_net, batch,
2                 gamma=0.99, optimizer=None):
3     states, actions, rewards, next_states, dones = batch
4
5     # Current Q-values
6     q_values = q_net(states).gather(1, actions.unsqueeze(1))
7
8     # Target Q-values
9     with torch.no_grad():
10         next_q = target_q_net(next_states).max(1)[0]
11         targets = rewards + gamma * (1 - dones) * next_q
12
13     loss = F.smooth_l1_loss(q_values.squeeze(), targets)
14
15     if optimizer:
16         optimizer.zero_grad()
17         loss.backward()
18         optimizer.step()
19
20     return loss.item()
```

Experiment 6: Logging and TensorBoard

Goal: Set up experiment tracking

Features:

- Automatic system info logging
- Scalar and histogram tracking
- Hyperparameter logging
- Model graph visualization

View results: `tensorboard -logdir runs`

Experiment 7: Checkpointing

Goal: Save and restore training state

What to save:

- Model weights
- Optimizer state
- Learning rate scheduler
- Training step/epoch
- RNG states
- Loss history

Enable training continuation after interruption!

Experiment 8: Git Integration

Goal: Version control for experiments

Best practices:

- Commit before experiments
- Include config hash in commits
- Track results with Git LFS
- Use meaningful commit messages
- Tag successful experiments

Experiment 9: Integration Test

Goal: Validate complete setup

Tests performed:

- 1 Environment check
- 2 Reproducibility verification
- 3 Model training
- 4 DQN components
- 5 Checkpoint save/load
- 6 Logging functionality

Must pass all tests before proceeding!

Key Takeaways

- 1 **RL is different:** Sequential decisions, delayed rewards
- 2 **MDP framework:** Foundation for all RL algorithms
- 3 **Reproducibility matters:** Always set seeds
- 4 **Device awareness:** CUDA > MPS > CPU
- 5 **Use optimizations:** AMP and compilation
- 6 **Track everything:** Logs, checkpoints, versions