

---

## Q-Learning with Stocks: Definitions for our seminar’s example

---

### *Immediate view and choices*

---

Agent	The trader who takes actions.
Long vs. Flat	Business jargon for buying (long) or not trading (flat) an asset. However, we use the term flat here to indicate that we are not invested.
State	All the information the agent <i>sees right now</i> when deciding (e.g. last 10 returns and whether he/she was already long or flat). The state should be “enough” to predict what happens next (Markov property).
Action	The choice available in that state – here: <i>stay/turn flat</i> (0) or <i>stay/turn long</i> (1).
Policy	The agent’s rule (possibly random) that maps each state to an action, written $\pi(a   s)$ .

### *World the agent interacts with*

---

Environment	The outside world producing the next state and reward after an action. In this lesson: the SPY price stream plus our trading rules.
Episode	One complete run until the environment signals <i>done</i> (here: when we reach the last price bar).
Reward	The immediate payoff after an action (daily return minus transaction cost).
Markov Decision Process (MDP)	The formal “state, action, reward, next state, done” framework underlying the whole problem.

### *Value estimation*

---

Q-function (action-value)	$Q^\pi(s, a)$ = expected <i>total future</i> reward if you start in state $s$ , take action $a$ once, then follow policy $\pi$ . It is <i>computed with</i> the Bellman equation.
Bellman equation	The recursion tying a Q-value to its own next-step reward: $Q(s, a) = \mathbb{E}[r + \gamma \max_{a'} Q(s', a')]$ . Foundation of Q-learning.
Discount factor $\gamma$	Number in $(0, 1)$ telling how strongly we value future rewards compared with immediate ones.

### *Neural networks used*

---

Online (Q-)network	The neural net we keep updating so its outputs approximate true Q-values; it chooses actions during training.
Target (frozen) network	A periodically copied snapshot of the online network. Its stable Q-values are used as learning targets until the next copy.

### *Experience storage and training data*

---

Replay buffer	Rolling memory of past transitions $(s, a, r, s', \text{done})$ so batches are less correlated.
Mini-batch	A random subset of transitions drawn from the replay buffer for one gradient update.
Warm-up memory	Minimum buffer size required before training starts, to ensure batches are diverse.

### *Exploration mechanisms*

---

$\varepsilon$ -greedy	With probability $\varepsilon$ pick a random action (explore); with $1 - \varepsilon$ pick the best estimated one (exploit). $\varepsilon$ decays over time.
Exploration vs. exploitation	The tension between “try something new to learn” and “use what I already know to earn reward now”. $\varepsilon$ -greedy is one simple strategy to balance the two.

---

Table 1: Glossary of key reinforcement-learning terms for the trading Q-learning example.