A detailed digital illustration of a desert canyon. The scene is characterized by layered, reddish-brown rock formations and deep, shadowed crevices. Sunlight filters through the upper levels, creating a hazy, golden atmosphere. In the foreground, a lone figure in a dark, hooded robe stands on a rocky ledge, looking out over the vast, layered landscape. The overall mood is one of solitude and grandeur.

“One learns from books and example only that certain things can be done. Actual learning requires that you do those things.”

Frank Herbert, *Children of Dune*

# CMPUT 365

## Introduction to RL

# Coursera Reminder

You **should be enrolled in the private session** we created in Coursera for CMPUT 365.

I **cannot** use marks from the public repository for your course marks.

You **need to check, every time**, if you are in the private session and if you are submitting quizzes and assignments to the private section.

The deadlines in the public session **do not align** with the deadlines in Coursera.

If you have any questions or concerns, **talk with the TAs** or email us `cmput365@ualberta.ca`.

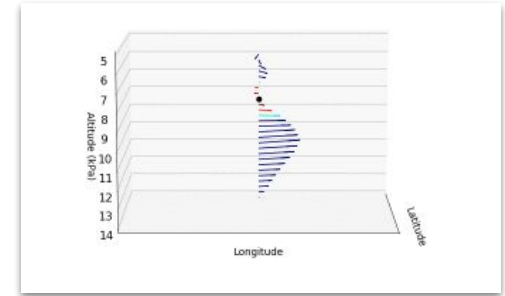
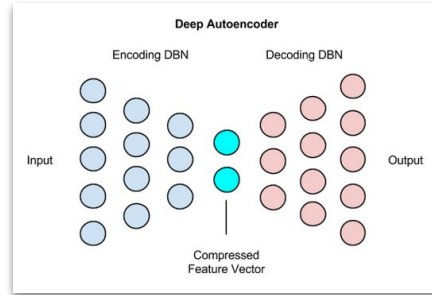
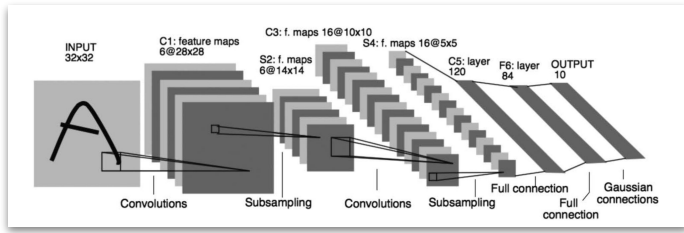
# Reminders and Notes

- The practice quiz for Control with FA is due today (only 4 to go!).
- The programming assignment for Control with FA is due Wednesday.
- Rich Sutton will give a guest lecture Dec 9th, Monday. Spread the word.
- A note on the final exam:
  - The required reading from the syllabus does not mean that's what will be covered in the final exam. There are some mismatches. Anything we discussed in class is fair game, including Maximization Bias and Double Learning (Section 6.7), and Nonlinear Function Approximation: Artificial Neural Networks (Section 9.7).
- The Student Perspectives of Teaching (SPOT) Survey is now available.

**Please, interrupt me at any time!**



# Last Class: Neural Networks and Feature/Input Construction



<https://wiki.pathmind.com/deep-autoencoder>

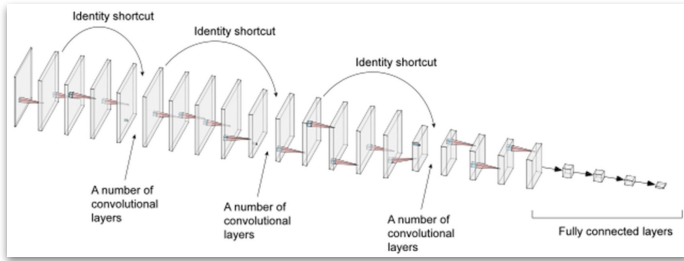
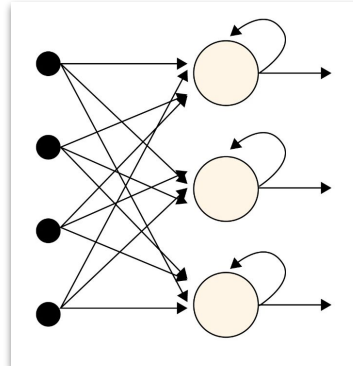


Image by Yu, Miao, and Wang (2022)



<https://www.scaler.com/topics/deep-learning/rnn/>

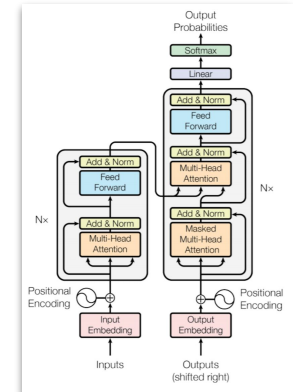


Image by Vaswani et al. (2017)



# Episodic Semi-gradient Control

- We need to approximate the action-value function now,  $\hat{q} \approx q_\pi$ , that is represented as a parameterized function form with weight vector  $\mathbf{w}$ .
- Before (until last class):  $S_t \mapsto U_t$ .  
Now:  $S_t, A_t \mapsto U_t$ .

# Episodic Semi-gradient Control

- We need to approximate the action-value function now,  $\hat{q} \approx q_\pi$ , that is represented as a parameterized function form with weight vector  $\mathbf{w}$ .
- Before (until last class):  $S_t \mapsto U_t$ .  
Now:  $S_t, A_t \mapsto U_t$ .

- Action-value prediction:

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha \left[ U_t - \hat{q}(S_t, A_t, \mathbf{w}_t) \right] \nabla \hat{q}(S_t, A_t, \mathbf{w}_t)$$



# Episodic Semi-gradient Control

- We need to approximate the action-value function now,  $\hat{q} \approx q_\pi$ , that is represented as a parameterized function form with weight vector  $\mathbf{w}$ .
- Before (until last class):  $S_t \mapsto U_t$ .  
Now:  $S_t, A_t \mapsto U_t$ .

- Action-value prediction:

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha \left[ U_t - \hat{q}(S_t, A_t, \mathbf{w}_t) \right] \nabla \hat{q}(S_t, A_t, \mathbf{w}_t)$$

- Episodic semi-gradient one-step *Sarsa*:

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha \left[ R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}_t) - \hat{q}(S_t, A_t, \mathbf{w}_t) \right] \nabla \hat{q}(S_t, A_t, \mathbf{w}_t)$$

# Episodic Semi-gradient Sarsa

## Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step size  $\alpha > 0$ , small  $\varepsilon > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Loop for each episode:

$S, A \leftarrow$  initial state and action of episode (e.g.,  $\varepsilon$ -greedy)

Loop for each step of episode:

Take action  $A$ , observe  $R, S'$

If  $S'$  is terminal:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

Go to next episode

Choose  $A'$  as a function of  $\hat{q}(S', \cdot, \mathbf{w})$  (e.g.,  $\varepsilon$ -greedy)

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$$

$S \leftarrow S'$

$A \leftarrow A'$



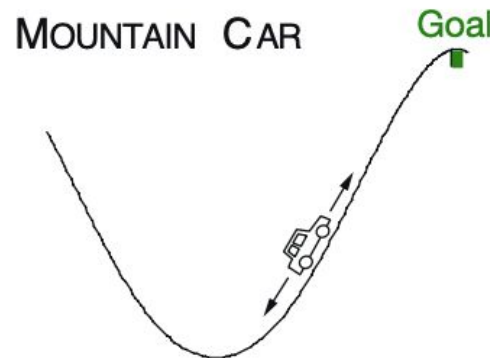
# Example: Mountain Car Task

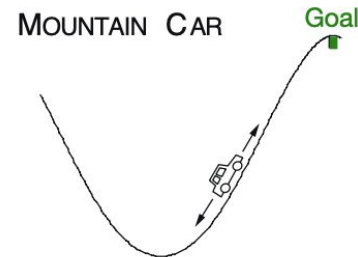
- Observations:  $(x, \dot{x})$
- Actions:
  - Full throttle forward: +1
  - Full throttle reverse: -1
  - Zero throttle: 0
- Rewards: -1 at every time step, until end of episode.
- Dynamics:

$$x_{t+1} \doteq \text{bound}[x_t + \dot{x}_{t+1}]$$

$$\dot{x}_{t+1} \doteq \text{bound}[\dot{x}_t + 0.001A_t - 0.0025 \cos(3x_t)]$$

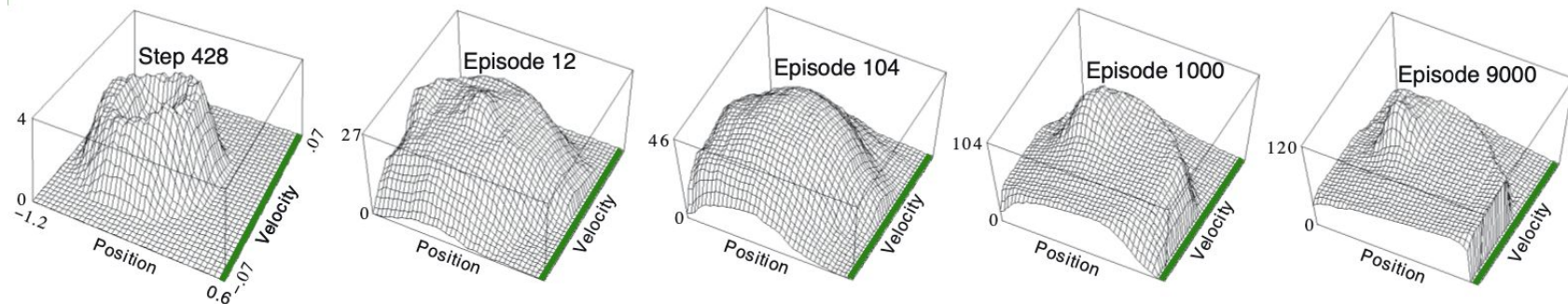
$$-1.2 \leq x_{t+1} \leq 0.5 \text{ and } -0.07 \leq \dot{x}_{t+1} \leq 0.07$$





# “Solution”: Mountain Car Task

- Feature representation:
  - Grid-tilings with 8 tilings and asymmetrical offsets.
  - $\hat{q}(s, a, \mathbf{w}) \doteq \mathbf{w}^\top \mathbf{x}(s, a) = \sum_{i=1}^d w_i \cdot x_i(s, a)$
- Sarsa
  - Weights initialized at zero. Effectively optimistic initialization.





# Avg. Reward: A Problem Setting for Continuing Tasks

- Continuing problems without discounting.
  - The agent cares about all rewards equally.

# Avg. Reward: A Problem Setting for Continuing Tasks

- Continuing problems without discounting.
  - The agent cares about all rewards equally.
- Quality of a policy is defined by the average rate of reward,  $r(\pi)$ :

$$\begin{aligned}
 r(\pi) &\doteq \lim_{h \rightarrow \infty} \frac{1}{h} \sum_{t=1}^h \mathbb{E}[R_t \mid S_0, A_{0:t-1} \sim \pi] \\
 &= \lim_{t \rightarrow \infty} \mathbb{E}[R_t \mid S_0, A_{0:t-1} \sim \pi], \\
 &= \sum_s \mu_\pi(s) \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) r
 \end{aligned}$$

**If the MDP is *ergodic*:** the starting state and any early decision made by the agent can have only a temporary effect; in the long run the expectation of being in a state depends only on the policy and the MDP transition probabilities.



# Avg. Reward: A Problem Setting for Continuing Tasks

- (Differential) Return:

$$G_t \doteq R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + R_{t+3} - r(\pi) + \dots$$

# Avg. Reward: A Problem Setting for Continuing Tasks

- (Differential) Return:

$$G_t \doteq R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + R_{t+3} - r(\pi) + \dots$$

- Differential value functions:

$$v_\pi(s) = \sum_a \pi(a|s) \sum_{r,s'} p(s', r | s, a) \left[ r - r(\pi) + v_\pi(s') \right],$$

$$q_\pi(s, a) = \sum_{r,s'} p(s', r | s, a) \left[ r - r(\pi) + \sum_{a'} \pi(a'|s') q_\pi(s', a') \right],$$

$$v_*(s) = \max_a \sum_{r,s'} p(s', r | s, a) \left[ r - \max_\pi r(\pi) + v_*(s') \right], \text{ and}$$

$$q_*(s, a) = \sum_{r,s'} p(s', r | s, a) \left[ r - \max_\pi r(\pi) + \max_{a'} q_*(s', a') \right]$$

# Avg. Reward: A Problem Setting for Continuing Tasks

- Differential value functions:

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{r,s'} p(s', r | s, a) \left[ r - r(\pi) + v_{\pi}(s') \right],$$

$$q_{\pi}(s, a) = \sum_{r,s'} p(s', r | s, a) \left[ r - r(\pi) + \sum_{a'} \pi(a'|s') q_{\pi}(s', a') \right],$$

$$v_*(s) = \max_a \sum_{r,s'} p(s', r | s, a) \left[ r - \max_{\pi} r(\pi) + v_*(s') \right], \text{ and}$$

$$q_*(s, a) = \sum_{r,s'} p(s', r | s, a) \left[ r - \max_{\pi} r(\pi) + \max_{a'} q_*(s', a') \right]$$

- Differential TD error:

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \hat{v}(S_{t+1}, \mathbf{w}_t) - \hat{v}(S_t, \mathbf{w}_t),$$

$$\delta_t \doteq R_{t+1} - \bar{R}_t + \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}_t) - \hat{q}(S_t, A_t, \mathbf{w}_t).$$



# Differential semi-gradient Sarsa

## Differential semi-gradient Sarsa for estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step sizes  $\alpha, \beta > 0$ , small  $\varepsilon > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Initialize average reward estimate  $\bar{R} \in \mathbb{R}$  arbitrarily (e.g.,  $\bar{R} = 0$ )

Initialize state  $S$ , and action  $A$

Loop for each step:

    Take action  $A$ , observe  $R, S'$

    Choose  $A'$  as a function of  $\hat{q}(S', \cdot, \mathbf{w})$  (e.g.,  $\varepsilon$ -greedy)

$\delta \leftarrow R - \bar{R} + \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})$

$\bar{R} \leftarrow \bar{R} + \beta \delta$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta \nabla \hat{q}(S, A, \mathbf{w})$

$S \leftarrow S'$

$A \leftarrow A'$

