*"And always, he fought the temptation to choose a clear, safe course, warning "That path leads ever down into stagnation.""*

Frank Herbert, *Dune*

# CMPUT 365
# Introduction to
# Sequential-Decision Making

# Plan

- Motivation

- *Non-comprehensive* overview of Intro to Sequential-Decision Making in Coursera (Bandits, Chapter 2 of the textbook)

3

# Reminder

You **should be enr** for CMPUT 365.

I **cannot** use marks

You **need** to **check** ou are submitting quizzes and assign

The deadlines in the oursera.

If you have any que

cmput365@ualbe



Marlos C. Machado

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

# Let's play a game!



Marlos C. Machado

# Bandits

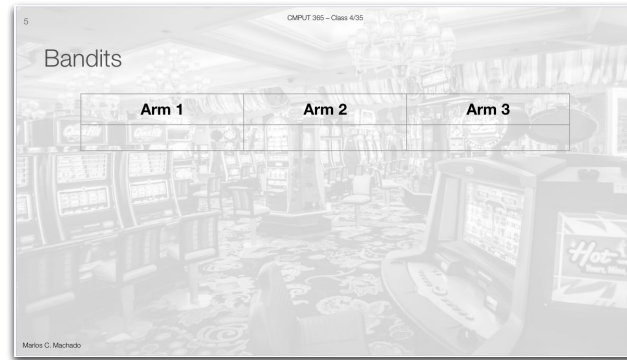| Arm 1 | Arm 2 | Arm 3 |
|---|---|---|
| 9, 7, 11, 12, 7, 6 | 6, 4, 5 | 8, 9, 9, 10 |

# Reinforcement learning (RL)

- RL is about learning from *evaluative* feedback (an evaluation of the taken actions) rather than *instructive* feedback (being given the correct actions).
  - Exploration is essential in reinforcement learning.

- It is not necessarily about online learning, as said in the videos, but more generally about sequential decision-making.

- Reinforcement learning potentially allows for continual learning but in practice, quite often we deploy our systems.

Marlos C. Machado

# Why study bandits?

- Bandits are the simplest possible reinforcement learning problem.
  - Actions have no delayed consequences.

- Bandits are deployed in so many places! [Source: Csaba's slides]
  - Recommender systems (Microsoft paper):
    - News,
    - Videos,
    - …
  - Targeted COVID-19 border testing (Deployed in Greece, paper).
  - Adapting audits (Being deployed at IRS in the USA, paper).
  - Customer support bots (Microsoft paper).
  - … and more.

# Why study bandits?



$$q^\star(a) \doteq \mathbb{E}[R_t \mid A_t = a]$$

$$A_t \doteq \text{argmax}_a\, Q_t(a)$$

We don't really know q*, so we use an estimate of it, $Q_t$

To exploit or to not exploit?
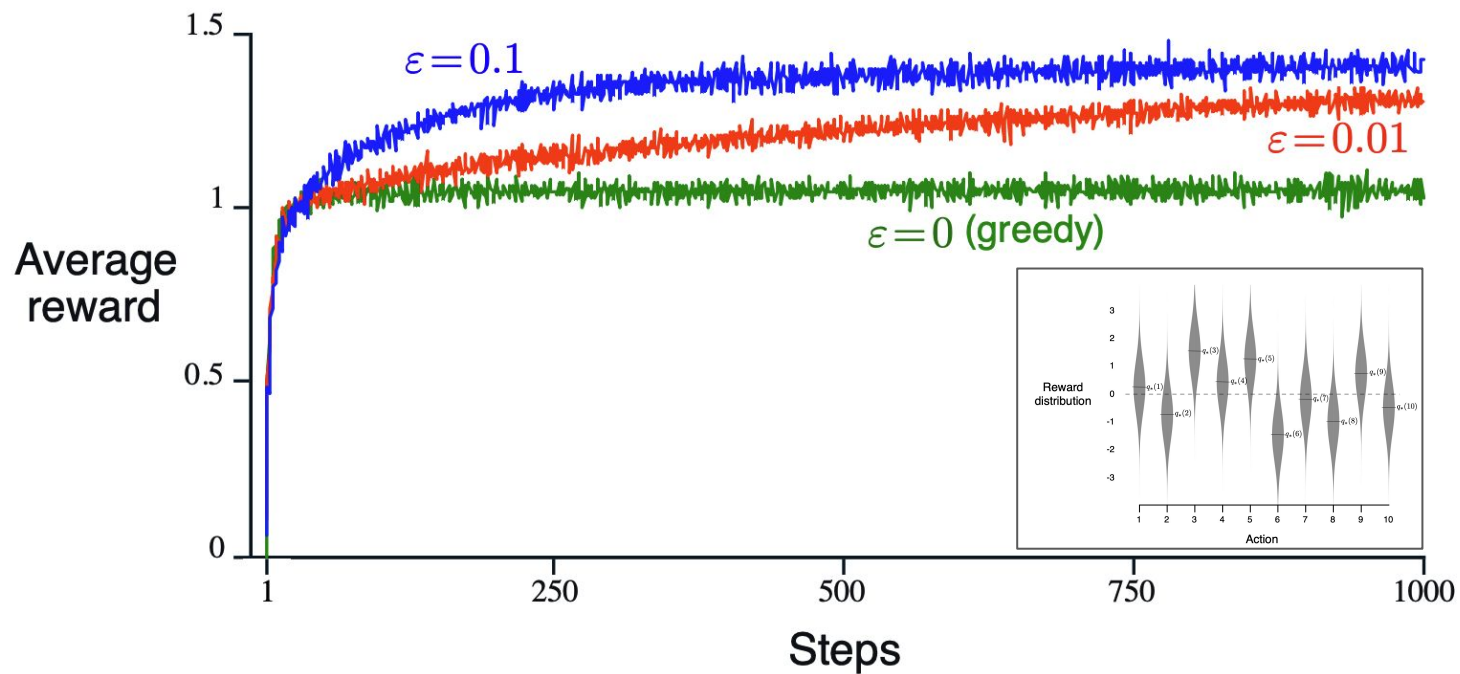
Greedy action

Marlos C. Machado

# Exploration

- Exploration is the opposite of exploitation.

- It is a whole, very active area of research, despite the textbook not focusing on it.

- How can we explore?
  - Randomly ($\in$ -greedy)
  - Optimism in the face of uncertainty
  - Uncertainty
  - Novelty / Boredom / Surprise
  - Temporally-extended exploration
  - …



To exploit or to not exploit?

# Exploration matters

# Incremental updates to estimate q*

$$Q_{n+1} \;=\; \frac{1}{n} \sum_{i=1}^{n} R_i$$

# Incremental updates to estimate q∗

$$
\begin{aligned}
Q_{n+1} \;&=\; \frac{1}{n} \sum_{i=1}^{n} R_i \\
&=\; \frac{1}{n} \left( R_n + \sum_{i=1}^{n-1} R_i \right) \\
&=\; \frac{1}{n} \left( R_n + (n-1)\frac{1}{n-1} \sum_{i=1}^{n-1} R_i \right) \\
&=\; \frac{1}{n} \left( R_n + (n-1)Q_n \right) \\
&=\; \frac{1}{n} \left( R_n + nQ_n - Q_n \right) \\
&=\; Q_n + \frac{1}{n} \left[ R_n - Q_n \right]
\end{aligned}
$$

# Next class

**Reminder: Practice Quiz for Coursera's Fundamentals of RL: Sequential decision-making is due today at midnight**.
**Programming Assignment for Coursera's Fundamentals of RL: Sequential decision-making is due on Wednesday.**

- What **I** plan to do: Wrap up Fundamentals of RL: An introduction to sequential decision-making (Bandits)
  - Go over some of your questions from Slack and eClass.
  - Time permitting, we'll work on some exercises in the classroom.