# Introduction to Reinforcement Learning
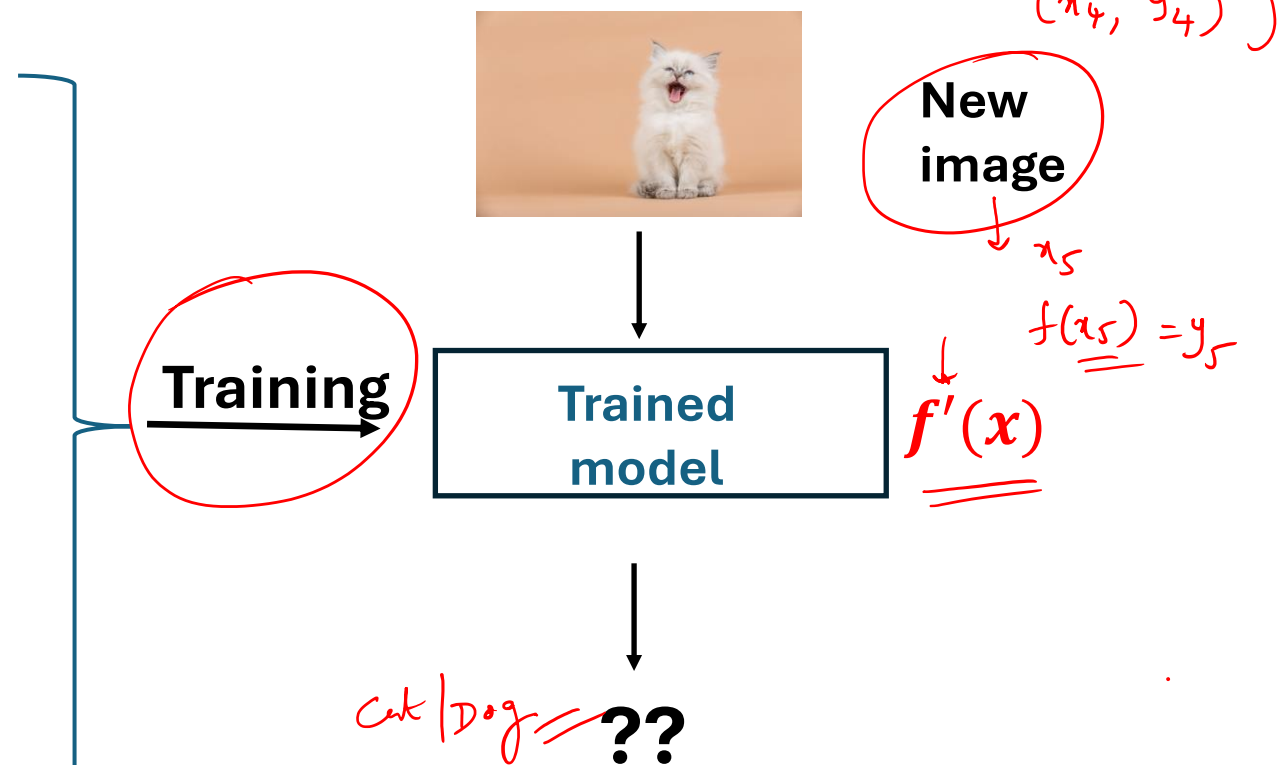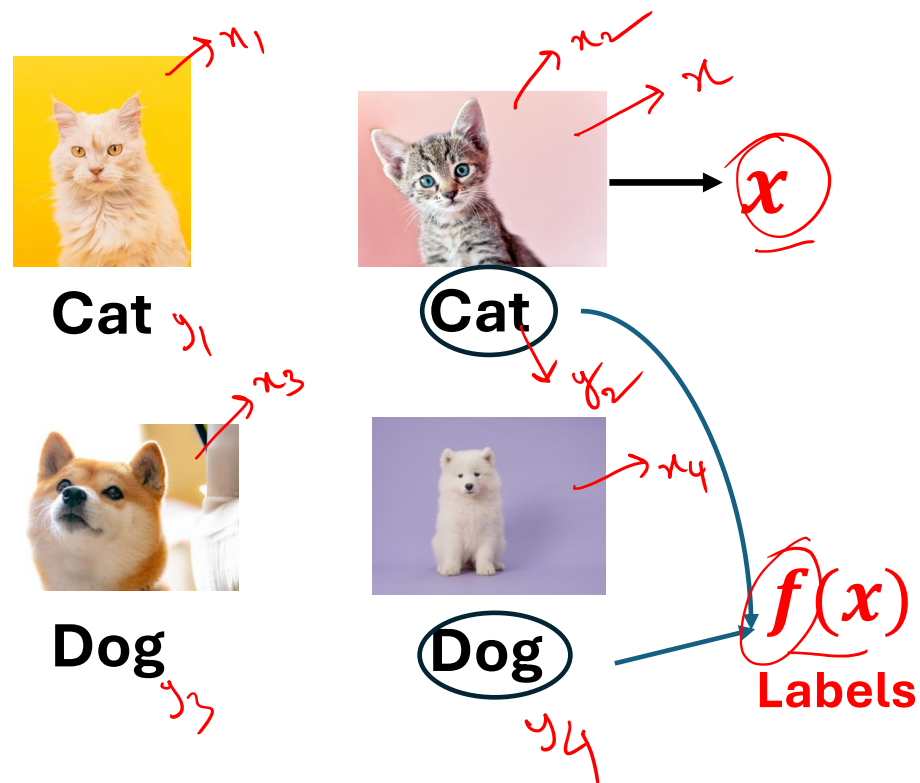
Subrahmanya Swamy Peruru

# Paradigms of Machine Learning

- Supervised Learning ✓ ①

- Unsupervised Learning ②

- Reinforcement Learning ③

# Supervised Learning

**Labeled Training Data**

$x_1$

$x_2$

$x \rightarrow x$

**Cat**   $y_1$

**Cat**   $y_2$

$x_3$

$x_4$

**Dog**   $y_3$

**Dog**   $y_4$

$f(x)$

**Labels**

function fitting

$x, f(x)$

$(x, y)^{y}$

$(x_1, y_1)$
$(x_2, y_2)$
$(x_3, y_3)$
$(x_4, y_4)$

**Training**

**Trained model**

**New image**

$x_5$

$f(x_5) = y_5$

$f'(x)$

Cat / Dog **??**

# Unsupervised Learning

$(x_1, y_1)$
$\downarrow \quad \downarrow$
i/p label

**Unlabeled Data**

$\rightarrow x_1$

$\downarrow x_2$

$x_3$

$x_4$

$x_5$

$x$

**Identify patterns**

# Reinforcement Learning

**Feedback:**
**Score,**
**new display**

**State:**
**Display**

**Actions:**
**UP / Left /**
**Right / Down**

+100

+1
-5

**Learn by Trial and Error**

$S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, \dots S_T$

Terminal

Episode

Programs

**Next State, Reward**
$(S_{t+1}, R_{t+1})$

**Environment**

$\pi : S \rightarrow A$

$\pi(s)$

**Agent**

*Action* $(A_t)$

*State* $(S_t)$

1. Agent observes the state and takes action
2. Environment puts the agent in a new state &
3. Gives a reward based on the action taken

+1  -1
Win / lose

**GOAL:** **Learn policy to maximize**
**the cumulative reward**

$\sum_t R\_t$

$R_1 + R_2$

$S_0, A_0, R_1, S_1, A_1, R_2, S_2 \dots$

0

0

# Paradigms of Machine Learning

- ## Supervised Learning
  - Fitting a function for the given labeled data $(x, y)$
  - $y \approx f(x)$

- ## Unsupervised Learning
  - Identifying patterns in unlabeled data
  - E.g. Clustering

- ## Reinforcement Learning
  - Learning sequential tasks through trial and error
  - Feedback through reward/penalty

# RL Demonstrations

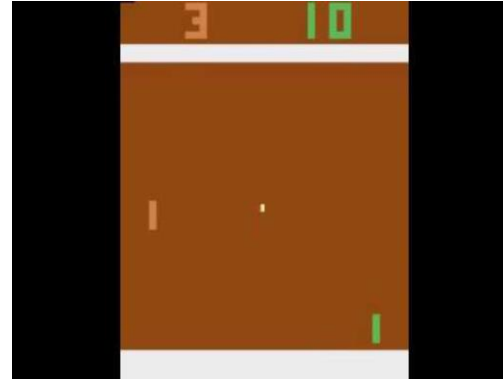DeepRL

MCTS

GO

→ Finance

→ Wireless Networks

→ Robotics

2016



"Autonomous Helicopter"

**Pong game**

**AlphaGo by DeepMind**

3    10

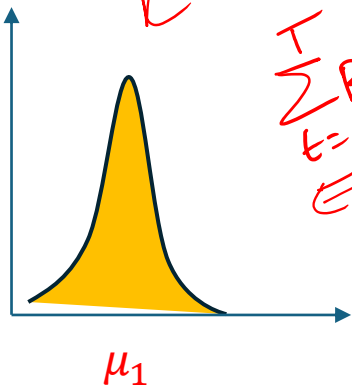"Multi-arm Bandits"

# One State RL: Multi-arm Bandits

- Simplified version of RL problem: "Multi-arm Bandit" problem.
  - Only one state
  - Multiple actions (a.k.a. arms)
    - $\mathcal{A}$ — Action set

$a_1$

$a_2$

$a_3$

- A reward distribution corresponding to each arm
  - $\mathcal{R}_a$ — Reward distrution for action a
  - $\mu_a = \mathbb{E}[\mathcal{R}_a]$ — Expected reward for action $a$

- Applications: Recommendation systems, Ad placement, ...
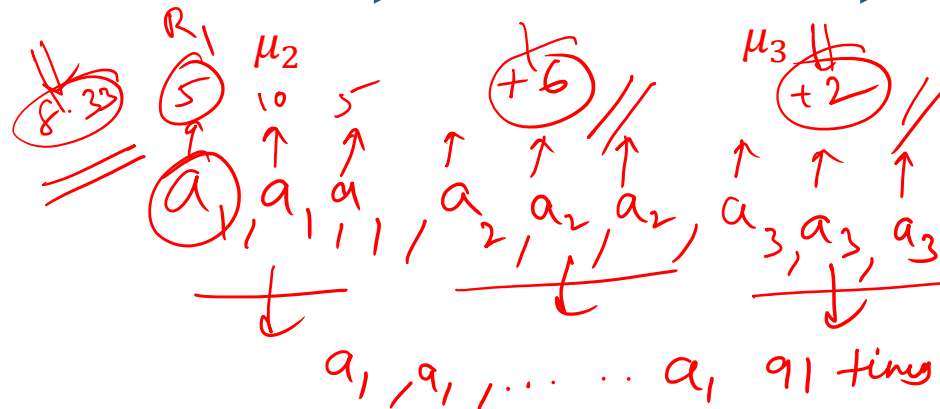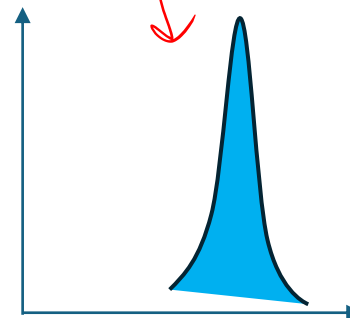
"Cognitive Radios"

# Multi-arm Bandits



Reward Distributions

$\sum_{t=1}^{T} R_t$

$\mu_1$  $\mu_2$  $\mu_3$

**Problem:**
- Reward distributions are **unknown**
- Given **T chances** to pull the arms
- Which arms should be pulled to **maximize the total reward** in those T rounds

Exploration
Vs
Exploitation dilemma

# ETC (Explore-Then-Commit)

*handwritten annotations (red):*
$\mu(a)$
$\mu(b)$ $\mu(c)$  — b)

$+5$ $+8$

$E[x] = \int_x x \, P_X(x) \, dx = \mu(a)$ — unknown

$\mu^* = \max\{\mu(a), \mu(b), \mu(c)\}$ — b

---

1. **Explore:** Play each arm $N$ times    $T - KN$

2. Compute the sample average rewards $\bar{\mu}(a) = \frac{1}{N}\sum_{t=1}^{KN} R_t \, 1\{a_t = a\}$   for each arm $a \in \mathcal{A}$   94

3. **Commit:** Play the arm with the highest sample average for the remaining $T - KN$ rounds

---

$\mu^*$ - Optimal arm's expected reward    $R_t$ - Sample reward obtained in round $t$

$a_t$ - Arm played in round $t$    $T$   - Total number of rounds

$K$   - Number of arms

*handwritten:* $\mu^* + \mu^* + \mu^* + \cdots \mu^*$   $E\left[R_1 + R_2 + \cdots R_T\right]$

**Performance (**ETC Vs Best possible reward**) :** $T\mu^* - \sum_{t=1}^{T} \mathbb{E}[R_t]$

**How much to Explore?** $N \approx \left(\frac{T}{K}\right)^{\frac{2}{3}}$
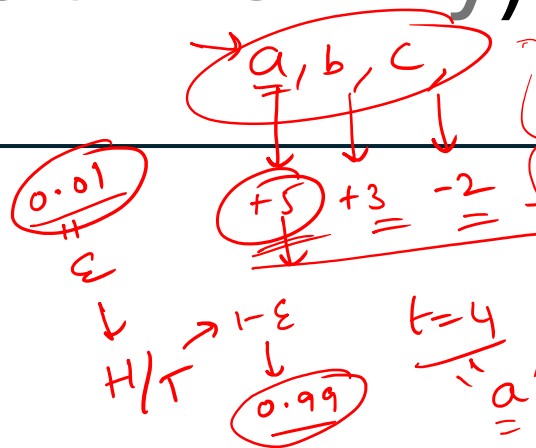
**$\epsilon$-Greedy** (Explore uniformly)

UCB

$a, b, c$   ETC

$4^{th}$ Exp NK   F-NK Exploit

"a" is best ← → NK

arm based on ← 

$a, b, c$

1. Play each arm once

0.01 #
$\epsilon$
↓
H/T → $1-\epsilon$
↓
0.99

$+5$  $+3$  $-2$   current estimates
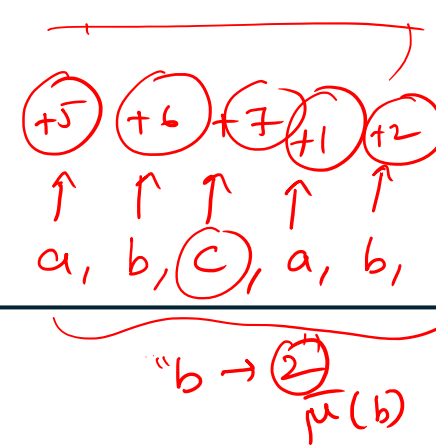
$t=4$
↓
"a"

2. In each round $t$:
   - Toss a coin with bias $\epsilon$.

   - If it lands in head: Explore - Play any arm randomly

   - Else: Exploit - Play the arm with the highest sample average so far

**What $\epsilon$ to choose?** $\epsilon \approx \left(\dfrac{K}{T}\right)^{\frac{1}{3}}$

# UCB (Upper Confidence Bound)

**Optimism under UnCertainty**

1. Play each arm once in the first $K$ rounds

2. For $t > K$:
   - Play the arm with the highest $UCB_t(a) = \overline{\mu_{t-1}}(a) + \sqrt{\dfrac{2 \log T}{n_{t-1}(a)}}$

   *UCB score* ↓   **Exploit**   **Explore**

   - Based on the observed sample reward $R_t$, update $n_t(a_t)$ and $\overline{\mu}_t(a_t)$
     - $n_t(a_t) = n_{t-1}(a_t) + 1$
     - $\overline{\mu}_t(a_t) = \dfrac{1}{n_t(a_t)}[(n_t(a_t) - 1)\,\overline{\mu_{t-1}}(a_t) + R_t]$

   $a_t = a$

**Exploit:** High sample reward arms are favoured
**Explore:** Least played arms are favoured

*Handwritten annotations:*

+5  +6  +7 +1 +2   $t = 7$

↑   ↑   ↑   ↑   ↑

a, b, ©, a, b,

$\mu_6(a) = \dfrac{5+1}{2} = 3$

$\mu_6(b) = 4$

"b → ②   $\mu(b)$   $\mu_6(c) = 7$

$n_{t-1}(a) = 2$
$n_{t-1}(b) = 2$
$n_{t-1}(c) = 1$

→ ETC
→ ε-greedy
→ ‖UCB‖

$\mu^* T - E[\sum_{t=1}^{T} R_t]$   3

Photos provided by Pexels

# Contextual Bandits – Multiple states

- ## News article Recommendation systems



- Articles – arms

- Like / Dislike – Reward

- User – State

Different users have different preferences to articles