# In-Context RL

## 1 Data Generation

### 1.1 5-Armed Bandit Problem

State space $S = \{0, 1, 2, 3, 4\}$; action space $A = \{0, 1, 2, 3, 4\}$

```
Initialize empty pretraining dataset B
for i in [N] do
    p₁ ~ Dirchlet distribution(𝟙)
    p₂ ~ point-mass distribution
    ω ~ Unif(0.1[10])
    action distribution p = (1 − ω)p₁ + ωp₂
    action means μ ~ Unif[0,1]⁵

    for h in [H = 500] do
        action aₕ ~ p (OHE)
        reward rₕ ~ N(μₐ, σ²) where σ = 0.3
        append (aₕ, rₕ) to goal g
        append aₕ, rₕ to B
```

## 2 Value Function Approximation

```
for i in [#iterations] do
    sample offline data {sᵗⁱ, aᵗⁱ, sᵗ₊₁ⁱ, gᵗⁱ}ⁱ₌₁ᴺ ~ B, {s₀ⁱ}ⁱ₌₁ᴹ ~ μ₀
    estimate the reward function R̂ for each task in the current offline
data
    Value objective: Lᵥ(θ) = (1−γ)/M Σᵢ₌₁ᴹ [Vθ(s₀ⁱ; g₀ⁱ)] + 1/N Σᵢ₌₁ᴺ [f⋆(Rᵗⁱ + γV(sᵗ₊₁ⁱ; gᵗⁱ) − V(sᵗⁱ; gᵗⁱ))]
    update Vθ:   Vθ ← Vθ − αᵥ∇Lᵥ(θ)
```

For bandit data, condition $V$ on $a$ and $s$ instead, and there is no $\gamma V$ term

## 3 DT Training

```
for i in [#iterations] do
    sample offline data {sᵗⁱ, aᵗⁱ, sᵗ₊₁ⁱ, gᵗⁱ}ⁱ₌₁ᴺ ~ B
```

estimate the reward function $\hat{R}$ for each task in the current offline
data

Policy objective: $L_\pi(\phi) = \sum_{i=1}^{N} \left[ \left( f'_\star \left( R_t^i + \gamma V_\theta(s_{t+1}^i; g_t^i) - V_\theta(s_t^i; g_t^i) \right) \right) \log \pi(a \mid s, g) \right]$

Update $\pi_\phi$: $\pi_\phi \leftarrow \pi_\phi - \alpha_\pi$

For bandit data, condition $V$ on $a$ and $s$ instead, and there is no $\gamma V$ term

# 4 Test

## 4.1 Offline Test

```
# Bandit version
subopt = []
for i in [500] do
```
  sample dataset $D$ with number of $i$ data $\sim \mathcal{B}_{\text{test}}$

  $s = s_0$

  $a^* = \arg\max_a \mu$

  $\hat{a} = \arg\max_{a \in \mathcal{A}} \pi_\phi(\cdot | s, D)$

  suboptimality = $\mu_{a^*} - \mu_{\hat{a}}$

  append suboptimality to subopt

## 4.2 Online Test

```
# Bandit version
suboptimality = 0
subopt = []
Initialize D = {}
for ep in [max_eps=500] do
```
  sample dataset $D \sim \mathcal{B}_{\text{test}}$

  $s = s_0 \sim \text{Unif}[0, 1]$

  $\hat{a} = \pi_\phi(\cdot | s, D)$

  suboptimality += $\mu_{a^*} - \mu_{\hat{a}}$

  append suboptimality to subopt

  add $(a, r)$ to D