# Optimizing a Finite-state controller

Goal: Optimize a policy $\Pi_\theta$ to maximize average rewards

POLICY/STRATEGY : $\Pi_\theta$ (parametrized by $\theta$)
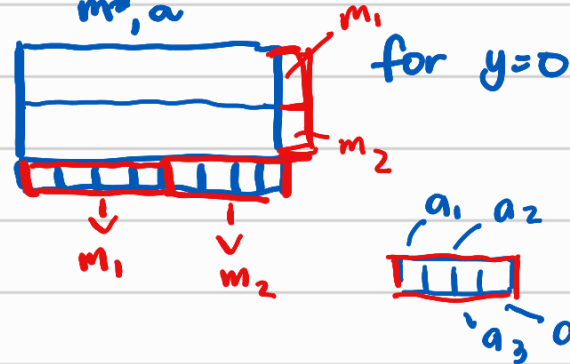
$\Pi_\theta = \text{Prob}(a, m^* | m, y)$     in code: $Pi(O, M, M^*A)$

i.e. $pi[0, :, :]$

$\boxed{M=2}$   $\xi$    for $y=0$

markov chain illustration



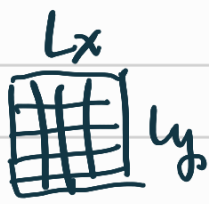$\rightarrow y=1$

$\dashrightarrow y=0$

(memory state transitions only)

*To couple w/ state actions, above the arrow, add the dominant action:

e.g. $\text{Prob}(a=\uparrow, m^* = m_1 | m = m_1, y = 0) = 0.8$



ENVIRONMENT SET-UP



size: $L_y \times L_x$

target: $(L_{x_0}, L_{y_0})$

information / signal :    in code: $PObs(O, L_y \times L_x)$

$\text{Prob}(y | s) \, z_\theta$

compute from particle counts from part 1

Options: (1) time-averaged

(2) time-averaged, symmetrical

(3) snapshot at certain time : $c/c_{max}$

THRESHOLD!!

$P(y=1|s) = 1$

if $\bar{c}(s) > th$

Some detection model we could use :

(1) conical distribution ( w/ tanh functions)

(2) From paper of infotaxis

$$C(S|S_0) = \frac{Cmax}{\|S-S_0\|_2 + 0.01} e^{-\frac{\|S-S_0\|_2}{\lambda} - \frac{\gamma V D}{2}}$$

$$\text{where} \quad \lambda = \sqrt{\frac{D\tau}{1 + \frac{v^2\tau}{4D}}}$$

■ Rewards function
in our set-up : reward = 0 if target found
reward = ↓ if target not
− (1·γ) found

Average reward given state⁻ $S, m$ :

$$r(s,m) = -(1-\gamma) \sum_{a,s',y} p(s'|s, a) f(y|s) \pi_\theta(a^*, m'|m, y) \, \mathbb{1}(x \neq x_s)$$

$\eta(s,m) \Rightarrow$ occupancy of state - memory pair

$$\eta = (1-\gamma T)^{-1} \rho \qquad T : S \times M \to S \times M$$

$$T(s', m' | s, m)$$

$V(s,m) \Rightarrow$ value of the policy

$$V^T = r^T (1- \gamma T)^{-1}$$

expected return : $\quad G = r^T \eta = V^T \rho \quad \bigg|$ for $\gamma \Rightarrow 1$

$G \approx -(1-\gamma) *$
av. time

# Optimization: NPG (Natural Policy Gradient)

pseudocode:

init $\theta$ [ 0, M, M*A]
pi = softmax ($\theta$)          $\rightarrow$ Boltzman param
value_prev = 0
value_new = 1

until  value_new - value_prev > tol :
    calculate $\eta$, Q
    grad = f $\times$ $\eta$ $\times$ Q
    $\theta$ = $\theta$ + (learning_rate) grad