

Building a game player  $\rightarrow$  Evaluation function

$$f(b) = \sum_i w_i \Phi_i(b)$$

↑                          domain knowledge  
learned  
by playing games

Set of actions: for each  $a$  there is a reward  $R(a)$

$R(a)$  might be random (in general, drawn from a distribution)

Ad placement:  $\Phi_i \rightarrow$  Ads/content in ad

Find ad  $a$  s.t. we obtain max reward  $R(a)$

$a_i \rightarrow$  try it  $n_i$  times:  $R_1 \dots R_{n_i}$

Find  $\vec{w}_i$  (weight or parameter vector)  $\rightarrow$  array

s.t.  $\underbrace{\vec{w}_i \cdot \vec{\Phi}_i}_{\sum_j w_{ij} \Phi_{ij}(a)}$   $\rightarrow$  correct-ish estimate of  $R(a_i)$

$$\sum_j w_{ij} \Phi_{ij}(a)$$

Minimize the MSE between our estimates and what we get from the environment

[Idea 1]

num n

(237)

Error function:  
 $\min_{\mathbf{w}} \sum_{j=1}^m (h(\mathbf{x}) - R_j)^2$

all data    prediction on jth try

E.g.  $h = \sum w_i \phi_i$  in the previous case.

$$h(\mathbf{x}) \quad \min_{\mathbf{w}} \sum_{j=1}^m (h(\mathbf{x}, j) - R_j)^2$$

take a derivative w.r.t  $w$   
and set them to 0 to find  $w$

In linear case: we can solve this in closed form (nice!)

Look for action  $a_i$  that maximizes reward (according to our estimates)

$$\arg \max h(w_i),$$

$\hookrightarrow$  best action

Exploration - Exploitation dilemma:

- try out something new? or
- use the best guess you have so far?

(137-3)  
A1 → payoff of 100 | Big difference → A1  
A2 → payoff of 10

Regret: how much loss do I get from not doing the optimal thing

A1 → payoff of 100 →? ] A2 not much

A2 → payoff of 99 ] worse

A1 → payoff of 100 ] Both tried 10<sup>6</sup> times

A2 → payoff of 99 ] lots of data → more "greedy"

Randomized strategies for picking actions

→ Best actions according to current picked more

→  $a_i$  and  $a_j$  have close reward est  $\rightarrow p(a_i) \approx p(a_j)$

→ Action that is best according to our est should see lots prob  $\rightarrow 1$  as data  $\rightarrow \infty$

Boltzmann distribution

$R^*$ -best est. so far  $(R^* - R_i)/\tau_{mi}$

$p(a_i)$  prop

to

↑

"temperature"

Bandit problem → really used for ad placement

(37-)

Delayed rewards: reward for some action is seen much later!

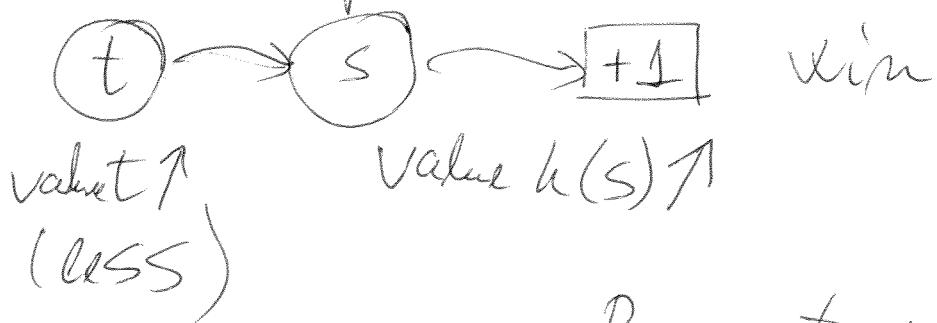
Suppose we are in some situation  $s$  (eg board position)

$s \xrightarrow[\text{action}]{a} s'$   
new situation (random)

$h(s) + r + h(s') \rightarrow$  using our val fn

intermediate rewards (eg. capture a piece)

Train  $h$  to be similar to  $r + h(s')$   
↳ adjust its parameters



Propagate rewards  
back to previous  
situation  
adjust  $h$

Similar idea to how classical conditioning works.

## Example applications

- ↳ Game playing
  - ↳ Automatically flying helicopters
- state/situation: set of pos/velocities +  
sensors

action: how much torque on rotor?

reward:

+100	for reaching destination
-1000	for crashing
0	otherwise

=> Train! → simulations

Choosing actions → max preferred (innate  
bias of alg)

Robots → same

→ Stock market!