

Evaluating a Policy

Policy : function π : observations \rightarrow action
 position s set of actions $A(s)$

$v(\pi)$ = expected reward following policy π
 can compute exactly for finite game (slow)

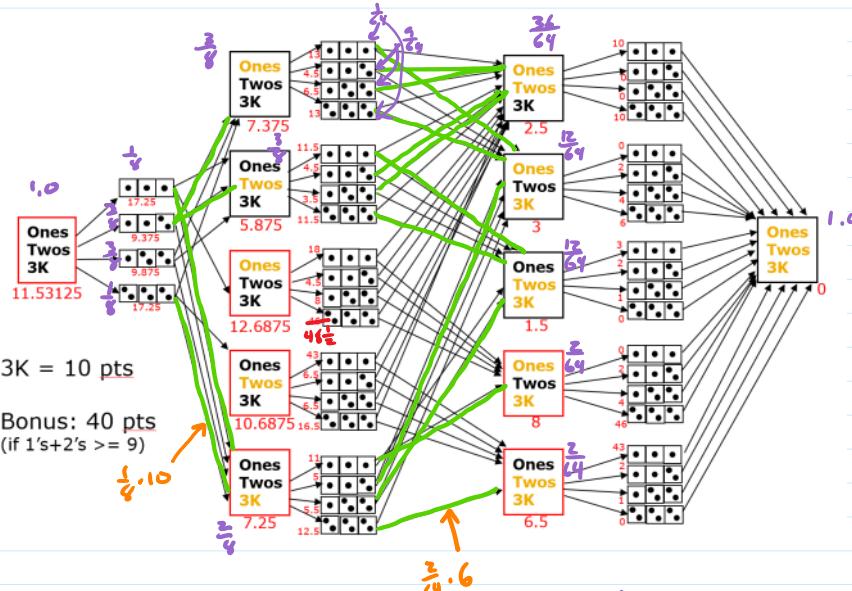
$P_{\pi}(s) = \text{prob reach position } s \text{ when following policy } \pi$

$$= \begin{cases} 1 & \text{if } s = \overbrace{s_0}^{\text{initial position}} \\ \sum_{g \text{ leads to } s} P_{\pi}(g) \cdot P_{\pi(g)}(g \rightarrow s) & \text{otherwise} \end{cases}$$

prob π
 leads to g prob action selected
 at g results in s

(special case of
Markov decision process)
- finite, no discount

$$v(\pi) = \sum_{\pi(g)=s} P_\pi(g) \cdot \underbrace{r_{\pi(g)}(g \rightarrow s)}_{\text{reward earned by chosen action}}$$



find them showing most
keep them

{ ones, twos, ..., sixes, straight, yahtzee, full house, chance }

$$\text{Yah+tree : } A(S) = \begin{cases} \text{up to } 32 \text{ contours} \\ \text{of size to keep} \end{cases} \quad \text{for scroll positions} \\ \left\{ c \mid \text{category } c \text{ is open} \right\} \quad \text{if } s \text{ is end of tree}$$

if results == 2
 keep 5's, 6's
else
 keep 4's, 5's, 6's

else
keep 4's, 5's, 6's

normal computations for computing $\pi_{\text{opt}}(s)$: (for finite process)

$V_a(s) = \text{value of action } a \text{ at position } s$

$$= \sum_{P_a(s \rightarrow s') \neq 0} P_a(s \rightarrow s') \cdot \underbrace{(R_a(s \rightarrow s') + V_{\pi_{\text{opt}}}(s'))}_{\substack{\text{reward for} \\ \text{action}}} \underbrace{V_{\pi_{\text{opt}}}(s')}_{\substack{\text{value at} \\ \text{resulting position}}}$$

$$\pi_{\text{opt}}(s) = \underset{a}{\operatorname{argmax}} V_a(s)$$

works even with a restricted set of actions
to determine opt policy wrt those actions

$$V_{\pi_{\text{opt}}}(s) = V_{\pi_{\text{opt}}}(s)$$

Classifiers

Classifier : function $\text{positions} \rightarrow \text{actions}$
 Classifier : function $\text{attributes} \rightarrow \text{class}$

Iris flower data set

MNIST database

<http://yann.lecun.com/exdb/mnist/>

Learning : supervised — examples available

reinforcement — reward observable

unsupervised
(clustering)

Methods : k-nearest neighbors treat example inputs as pts in d-dim space,
 classify new input according to class of majority
 of the k nearest neighbors

decision trees

nested ifs

if sepal-len > 0.2
 if petal-wid < 0.4
 versicolor virginica
 setosa

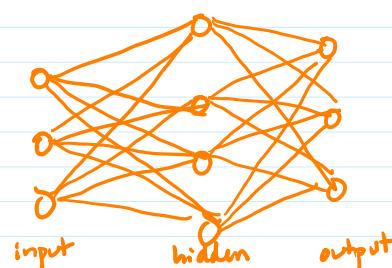
perceptron



$$\text{output} = f \left(\sum_{k=1}^K i_k \cdot w_k \right)$$

$$\text{ex: threshold } f(y) = \begin{cases} 1 & \text{if } x > 5 \\ 0 & \text{otherwise} \end{cases}$$

multi-layer perceptron
(a type of artificial neural network)



ReLU



Deep Q network learning to play Pong