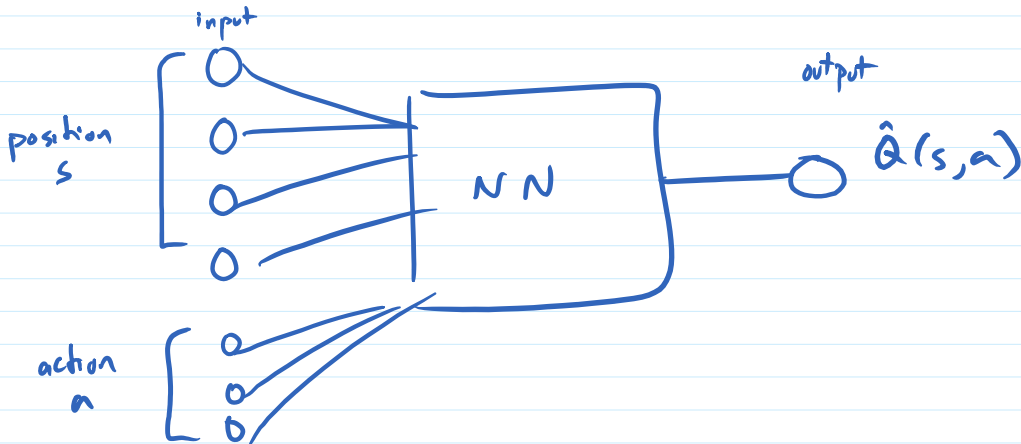initialize learning, target networks

for each iteration
$\rightarrow$ for each episode
for each event
add $(s, a, s', r)$ to replay database
sample replay database
train learning network towards $r + \gamma \cdot \max\limits_{a} \hat{Q}(s', a')$

*some smallish number*

*computed by target network*

if **enough time has passed**
copy learning network to target network

ANN for approximating $Q(s,a)$

input

position
s

N N

output

$\hat{Q}(s,a)$

action
a

$Q(\,(80, 4, 10, 2\cancel{4}), 1) = 0.53\ldots$    $\hat{Q}(\,(80, 4, 10, 74), 1) = 0.30$

observe    $r + \gamma V(s')$    expected $\hat{Q}(s,a)$
*output of neural network*

and record in
replay database
(experience replay)

$\boxed{\max\limits_{a} \hat{Q}(s', a)}$    *use as training input*

use 2 separate networks    1) to estimate future reward
"target" network

2) to train
learning or "prediction" network