

Pong from Pixels

Deep RL Bootcamp

Andrej Karpathy, Aug 26, 2017

1

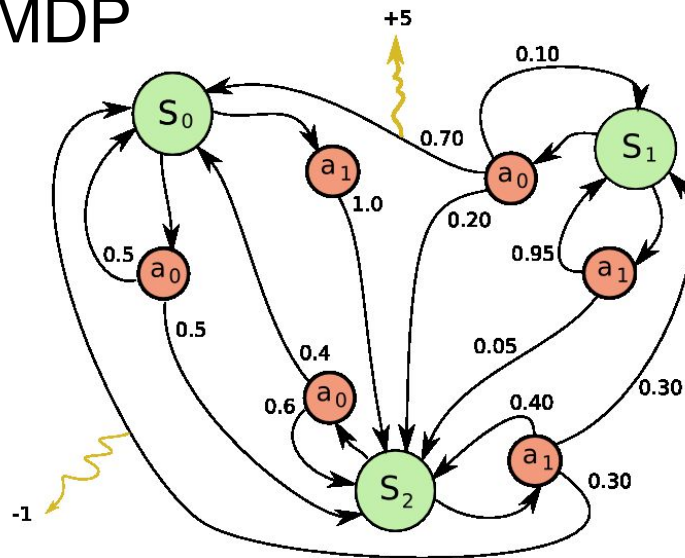
1

1

MDP:

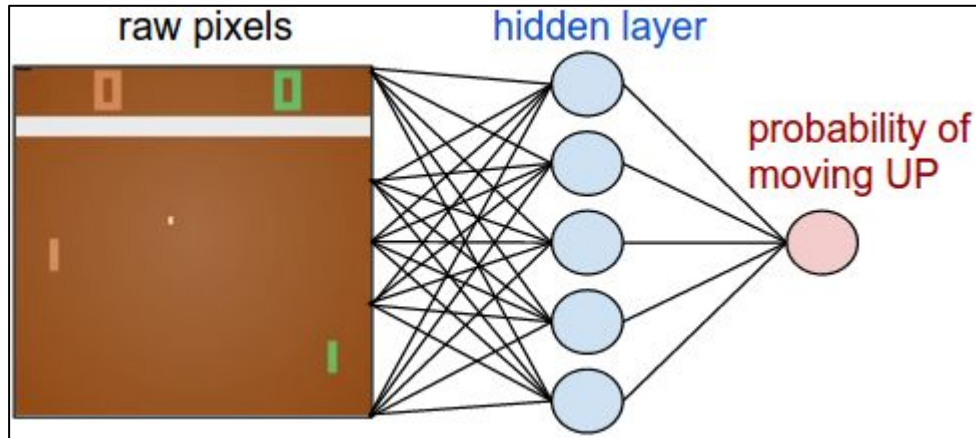
- Reward: +1 if we get the ball past the opponent, -1 if they get it past us
- State: 1) location of ball, 2) location of our paddle, 3) location of opponent's paddle
- Actions: 1) Move paddle up, 2) Move paddle down

MDP



Policy network

- single-layer MLP
- we predict the
probability of moving up



Policy network

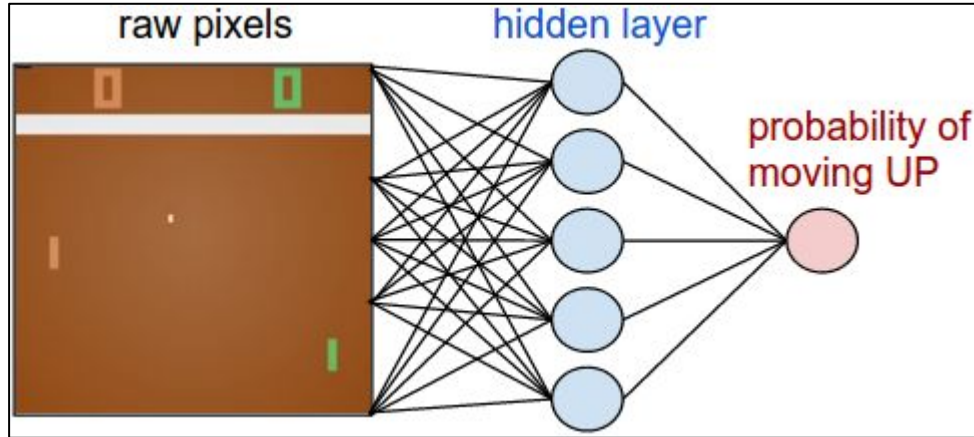
e.g.,

height width

[80 x 80]

array of

Grey-scale

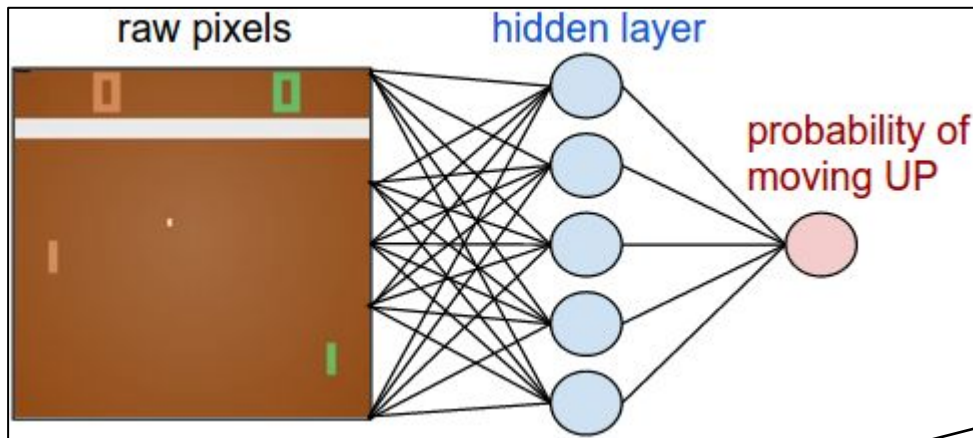


Policy network

height width

[80 x 80]
array

input layer is 80*80



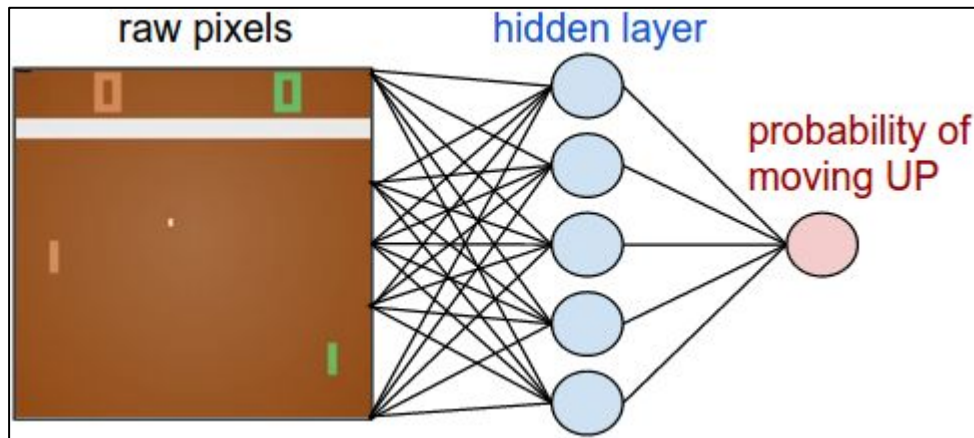
get the logit

```
h = np.dot(W1, x) # compute hidden layer neuron activations
h[h<0] = 0 # ReLU nonlinearity: threshold at zero
logp = np.dot(W2, h) # compute log probability of going up
p = 1.0 / (1.0 + np.exp(-logp)) # sigmoid function (gives probability of going up)
```

apply sigmoid to squash this between
0 and 1

Policy network

height width
[80 x 80]
array



biases

E.g. 200 nodes in the hidden network, so:

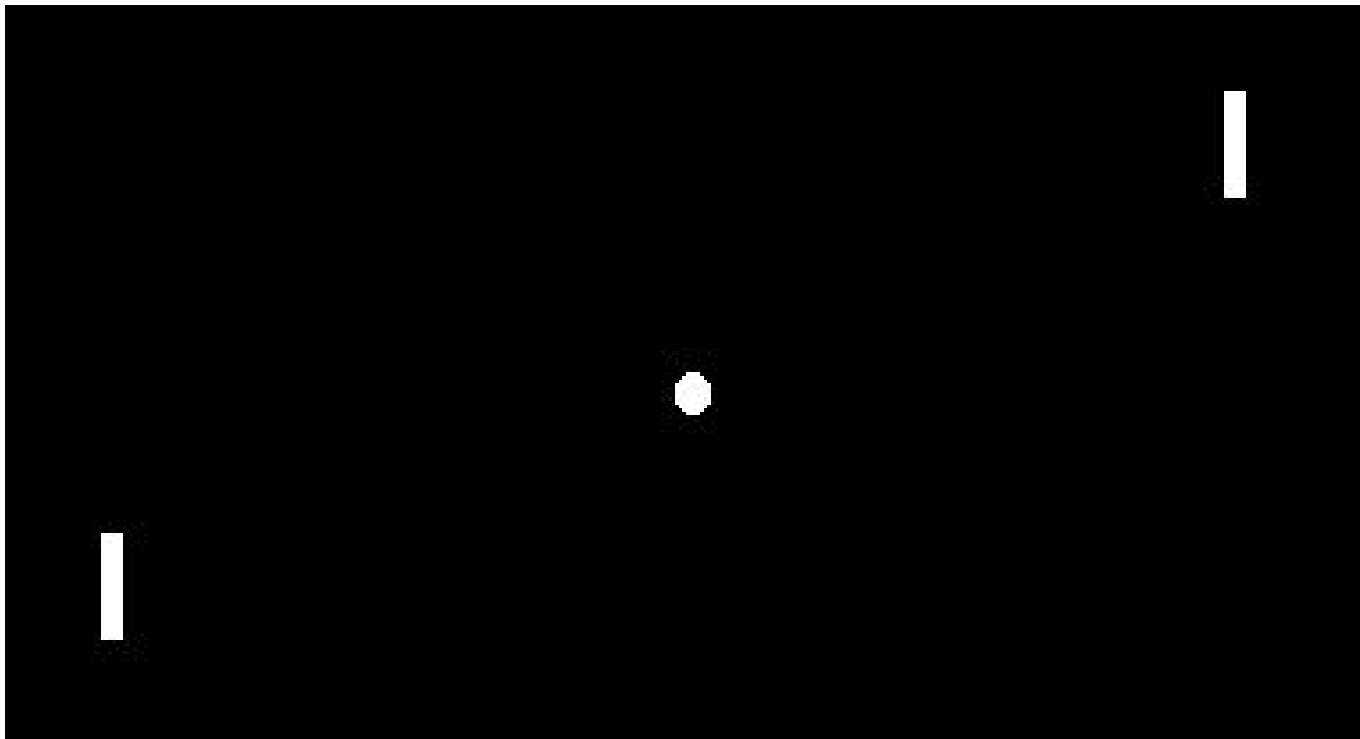
$$[(80*80)*200 + 200] + [200*1 + 1] = \sim 1.3\text{M parameters}$$

Layer 1

Layer 2

Note that we can't just take in a single frame because that would only be a static input and the game is dynamic.

Instead we could 1) concatenate/stack multiple frames in the input, or 2) take the difference between two frames as the input



Network does not see this. Network sees $80 \times 80 = 6,400$ numbers.
It gets a reward of +1 or -1, some of the time.
Q: How do we efficiently find a good setting of the 1.3M parameters?

Problem is easy if you want to be inefficient...

random search

1. **Repeat Forever:**
2. Sample 1.3M random numbers
3. Run the policy for a while
4. If the performance is best so far, save it
5. Return the best policy

Problem is easy if you want to be inefficient...

1. **Repe**
2. Sa
3. Ru
4. If th
5. Retur



Problem is easy if you want to be inefficient...

1. **Repeat**
2. **Save**
3. **Run**
4. **If th**
5. **Return**



Policy Gradients

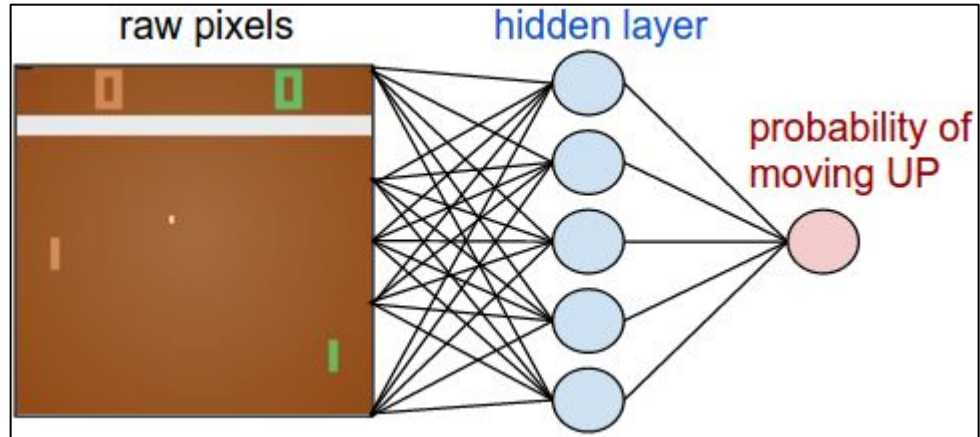


Suppose we had the training labels...
(we know what to do in any state)

(x1,UP)
(x2,DOWN)
(x3,UP)
...

Suppose we had the training labels...
(we know what to do in any state)

(x1,UP)
(x2,DOWN)
(x3,UP)
...



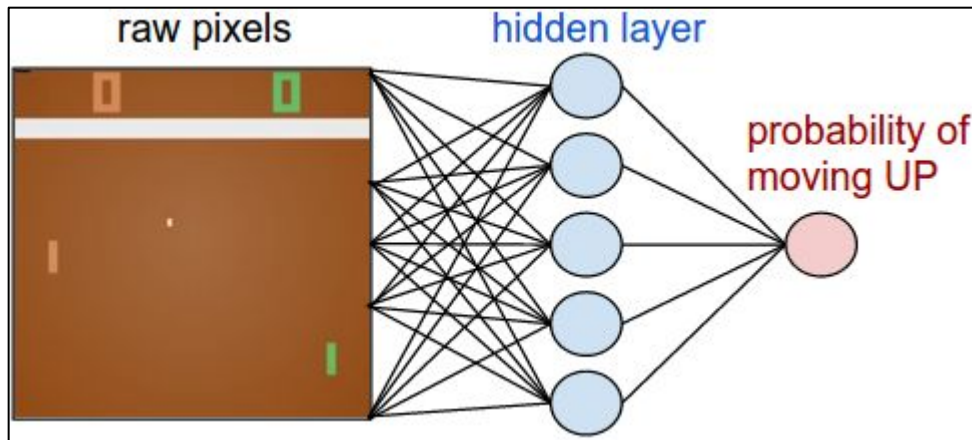
Suppose we had the training labels...
(we know what to do in any state)

(x1,UP)
(x2,DOWN)
(x3,UP)
...

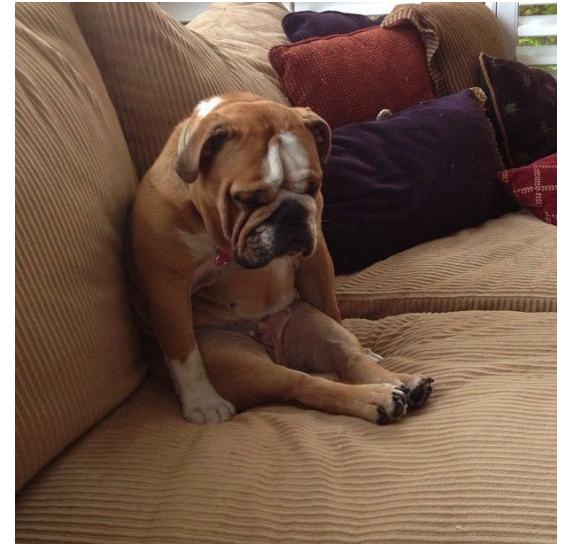
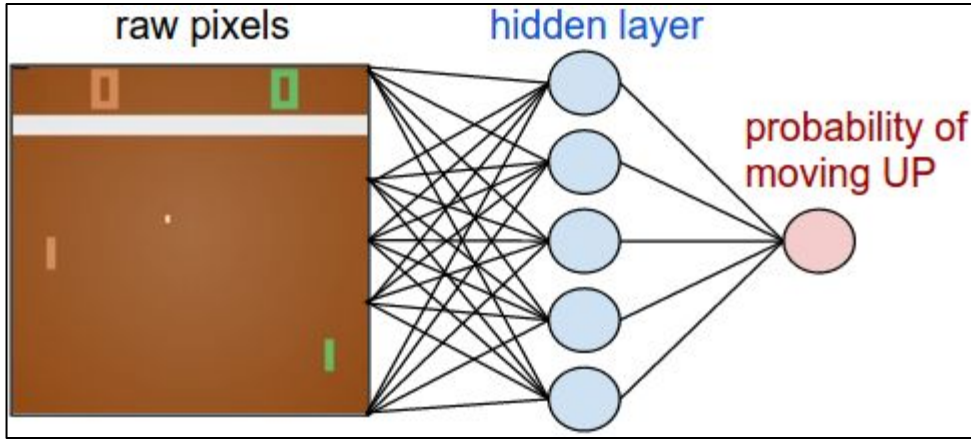
so if we were told what to do
at every state, we could just
train a NN and this becomes
a supervised learning
problem where we try to
maximize the log probability
of getting the answer correct

maximize:

$$\sum_i \log p(y_i | x_i)$$



Except, we don't have labels...



Should we go UP or DOWN?

Except, we don't have labels...

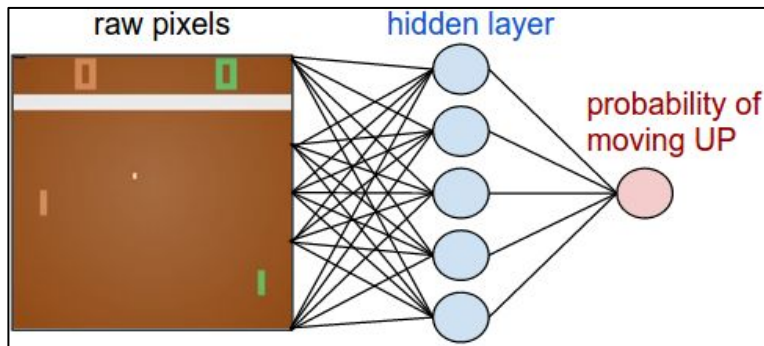


“Try a bunch of stuff and see what happens. Do more of the stuff that worked in the future.”

-RL

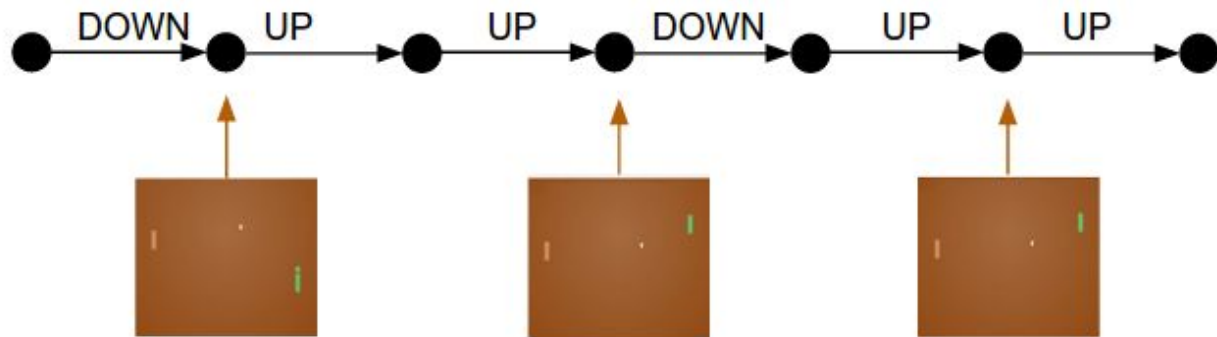
Let's just act according to our current policy...

A rollout is a single run/episode of the game with the current policy.



Start with a randomly initialized policy, then take actions for each state

So this is what actions we end up doing...At first this is pretty much just random actions



WIN

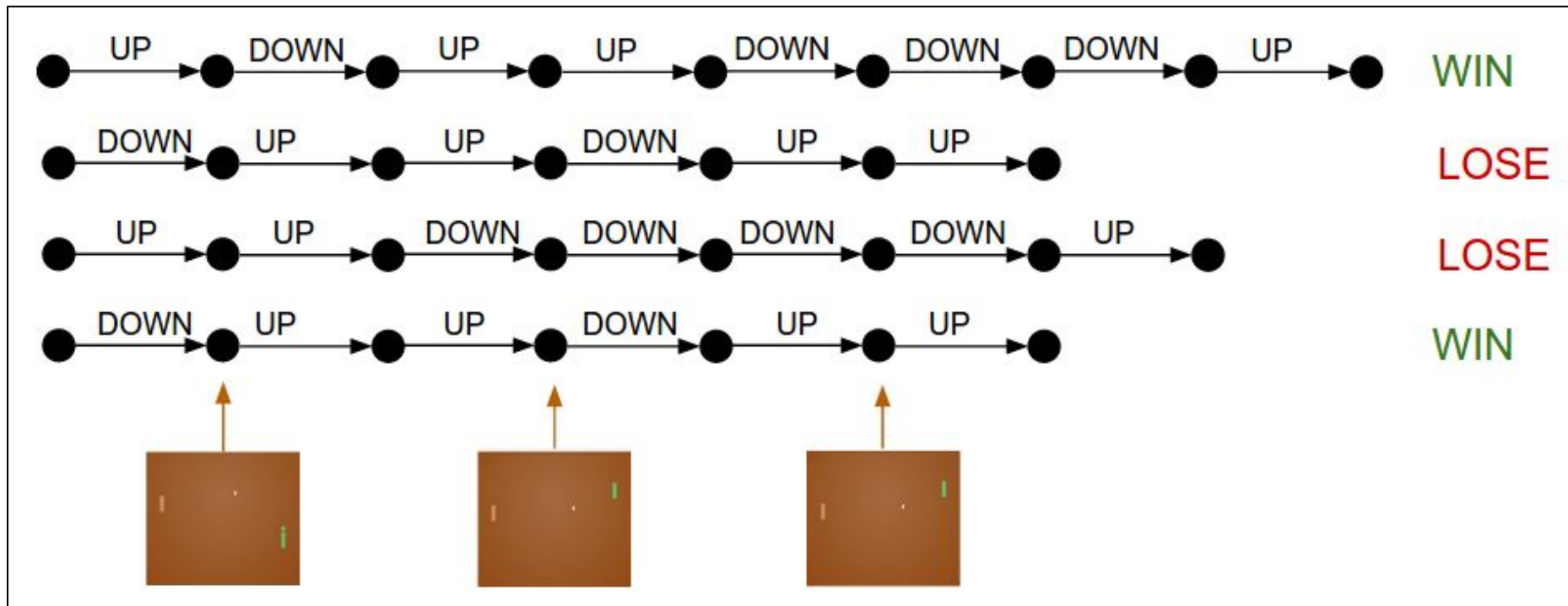
Rollout the policy and collect an episode

Collect many rollouts...

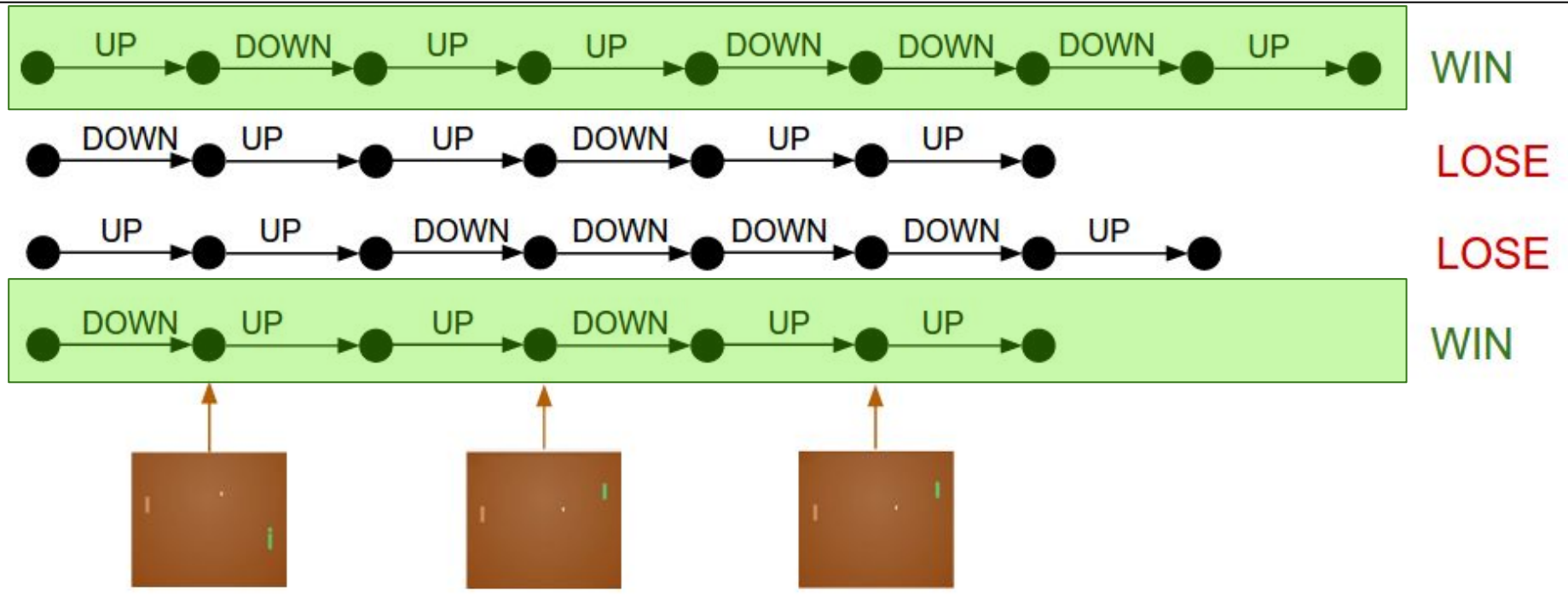
4 rollouts:

So this becomes our dataset

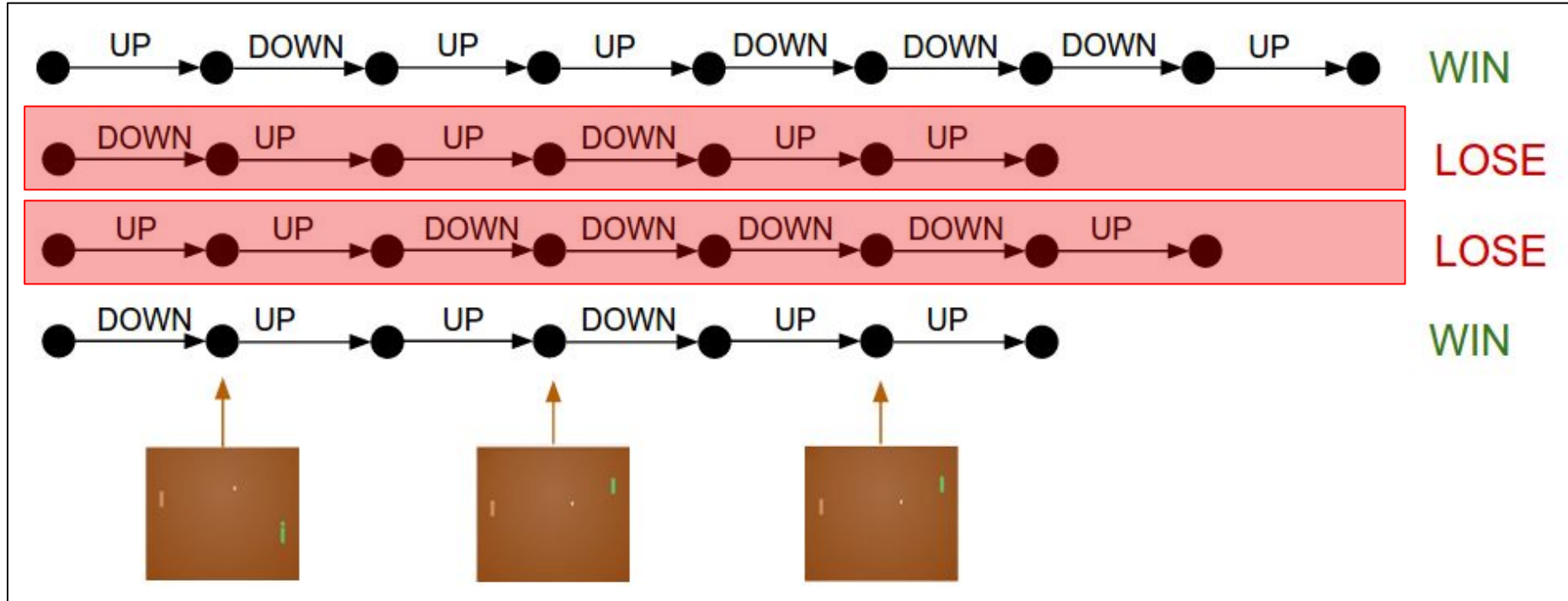
The labels are the actions that we happened to do



Not sure whatever we did here, but
apparently it was good.



Not sure whatever we did here, but it was bad.

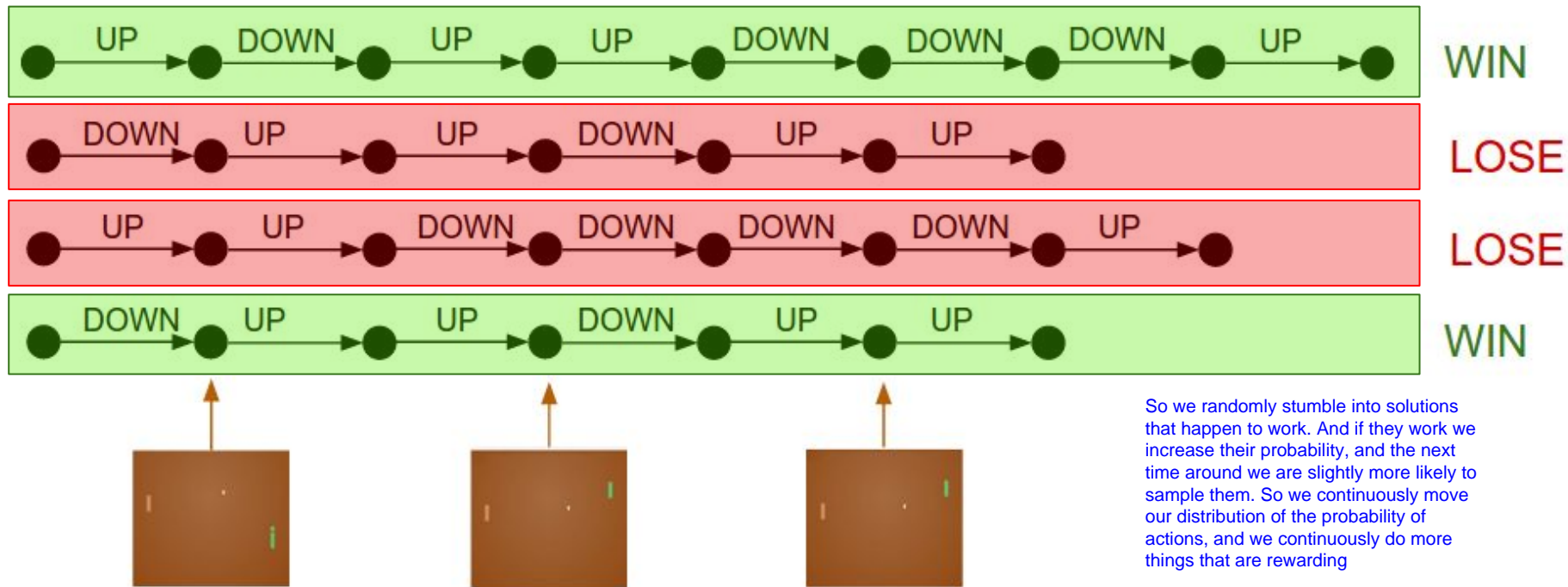


Pretend every action we took here was the correct label.

maximize: $\log p(y_i | x_i)$

Pretend every action we took here was the wrong label.

maximize: $(-1) * \log p(y_i | x_i)$



Supervised Learning

maximize:

$$\sum_i \log p(y_i | x_i)$$

For images x_i and their labels y_i .

Supervised Learning

maximize:

$$\sum_i \log p(y_i | x_i)$$

For images x_i and their labels y_i .

Reinforcement Learning

Supervised Learning

maximize:

$$\sum_i \log p(y_i | x_i)$$

For images x_i and their labels y_i .

Reinforcement Learning

1) we have no labels so we sample:

$$y_i \sim p(\cdot | x_i)$$

we rollout our policy and the actions we happen to take become our data

Supervised Learning

maximize:

$$\sum_i \log p(y_i | x_i)$$

For images x_i and their labels y_i .

Reinforcement Learning

1) we have no labels so we sample:

$$y_i \sim p(\cdot | x_i)$$

2) once we collect a batch of rollouts:

maximize:

$$\sum_i A_i * \log p(y_i | x_i)$$

we want this scalar to be high if we want to encourage that action in the future, and we want it to be small if we want to discourage that action in the future

Supervised Learning

maximize:

$$\sum_i \log p(y_i | x_i)$$

For images x_i and their labels y_i .

So at each state/frame,

- 1) Pass state into NN as input
- 2) NN produces probability of 1 (i.e. going up).
So for example, it could output 0.8
- 3) Pick a random number between 0-1. If it's > than 0.8, take that action (i.e. go up), and otherwise go down.
- 4) That action becomes our y_i

Once the game is over,

- 5) Find the advantage. If you won, $A=1$, if you lost $A=-1$.
- 6) Apply this action to all of the state/actions in that rollout/trajectory. So now we have (s, a, A) for each step in a rollout. Based on this advantage, we'll either increase or decrease the probability of doing that action
- 7) Now we want to alter our policy to maximize the sum of $A * \log p(y|x)$ for everything we did in our batch of rollouts.

Reinforcement Learning

1) we have no labels so we sample:

$$y_i \sim p(\cdot | x_i)$$

the labels are the actions we happened to take

given an action

2) once we collect a batch of rollouts:
maximize:

$$\sum_i A_i * \log p(y_i | x_i)$$

we want to increase the logprob of actions if it turned out good, and decrease it if it turned out bad

We call this the **advantage**, it's a number, like +1.0 or -1.0 based on how this action eventually turned out.

so A_i will be the same for all actions/data within the same rollout. So we equally blame all the actions.

Supervised Learning

maximize:

$$\sum_i \log p(y_i | x_i)$$

For images x_i and their labels y_i .

So the only difference between RL and supervised learning is that we have the Advantage factor, which might be negative. In supervised learning, the "advantage" is always 1.

the advantage is a scalar that tells you whether or not you want to encourage or discourage the action to happen to take

Reinforcement Learning

1) we have no labels so we sample:

$$y_i \sim p(\cdot | x_i)$$

2) once we collect a batch of rollouts:
maximize:

$$\sum_i A_i * \log p(y_i | x_i)$$

+ve advantage will make that action more likely in the future, for that state.

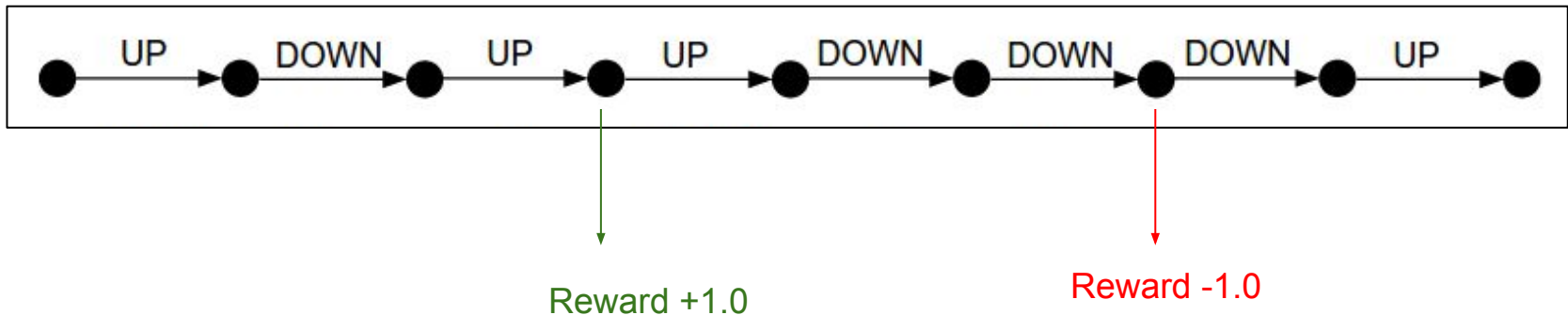
-ve advantage will make that action less likely in the future, for that state.



Discounting

Blame each action assuming that its effects have exponentially decaying impact into the future.

It's kind of odd to assign equal blame to every single action in the rollout, because some of them could have actually been good, while others could have been really bad. So in practice, we discount to modulate the blame. We blame each action in an exponentially decreasing action. The further away the action was from the reward, the less it's blamed



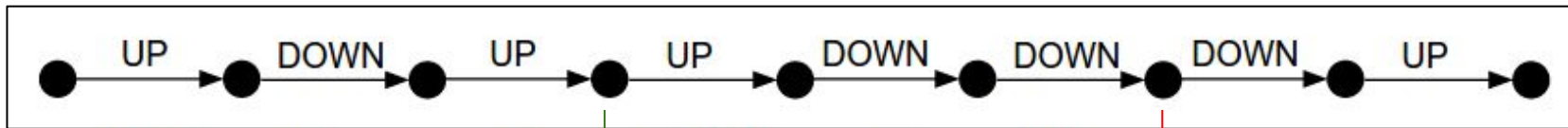
Discounting

Blame each action assuming that its effects have exponentially decaying impact into the future.

Discounted rewards

$$\sum_i A_i * \log p(y_i | x_i)$$

0.21 0.24 0.27 -0.81 -0.9 -1 0 0



Reward +1.0

Reward -1.0

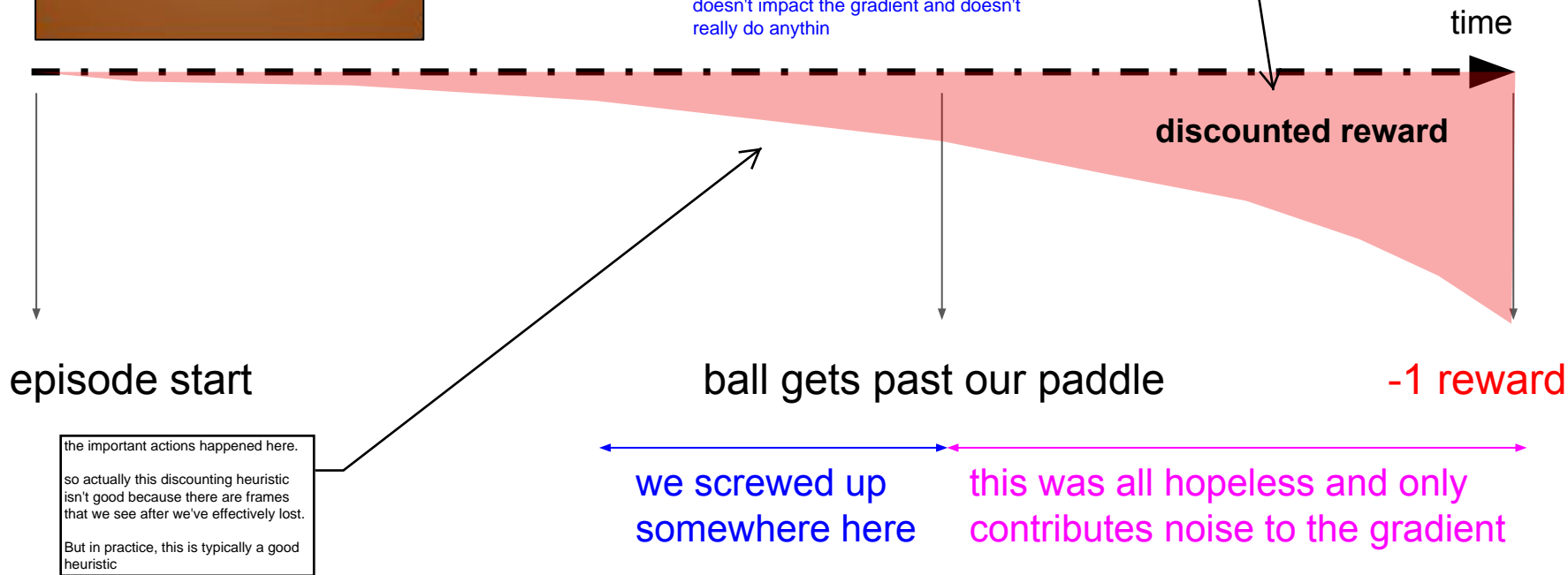
$\gamma = 0.9$



if you were to use a Value Function, and you have a Value Function estimator, it could learn that when the ball is past you, it knows you're about to lose. So it learns that the expected return is -1. So the baseline would become -1. So in this case the Advantage is return minus baseline. So then the empirical return minus the baseline is 0, so it doesn't contribute the gradient.

This is would be good because the actions after the ball gets past you doesn't impact the gradient and doesn't really do anything

there's several frames where the ball is past us, but there's nothing we can do



**YEAH, IF YOU COULD JUST SHOW ME
SOME CODE**



THAT'D BE GREAT

memegenerator.net

130 line gist, numpy as the only dependency.

<https://gist.github.com/karpathy/a4166c7fe253700972fcbc77e4ea32c5>

```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprobs, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprobs else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - aprobs) # grad that encourages the action that was taken to be taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     ddrs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99         eps = np.vstack(xs)
100         eph = np.vstack(hs)
101         epdlogp = np.vstack(dlogps)
102         epr = np.vstack(drs)
103         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105         # compute the discounted reward backwards through time
106         discounted_epr = discount_rewards(epr)
107         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108         discounted_epr -= np.mean(discounted_epr)
109         discounted_epr /= np.std(discounted_epr)
110
111         epdlogp *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
112         grad = policy_backward(eph, epdlogp)
113         for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
114
115         # perform rmsprop parameter update every batch_size episodes
116         if episode_number % batch_size == 0:
117             for k,v in model.iteritems():
118                 g = grad_buffer[k] # gradient
119                 rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
120                 model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k]) + 1e-5)
121                 grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
122
123         # boring book-keeping
124         running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
125         print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
126         if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
127         reward_sum = 0
128         observation = env.reset() # reset env
129         prev_x = None
130
131     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
132         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

```

env = gym.make("Pong-v0")
observation = env.reset()

prev_x = None # used in computing the difference frame
xs,hs,dlogps,drs = [],[],[],[]

running_reward = None
reward_sum = 0
episode_number = 0
while True:
    if render: env.render()

```

training loop
- we iterate through the learning
algorithm

Nothing too scary over here.

We use OpenAI Gym.
And start the main training loop.

```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprobs, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprobs else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - aprobs) # grad that encourages the action that was taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99         eps = np.vstack(xs)
100         eph = np.vstack(hs)
101         epdlogps = np.vstack(dlogps)
102         epr = np.vstack(drs)
103         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105         # compute the discounted reward backwards through time
106         discounted_epr = discount_rewards(epr)
107         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108         discounted_epr -= np.mean(discounted_epr)
109         discounted_epr /= np.std(discounted_epr)
110
111         epdlogps *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
112         grad = policy_backward(eph, epdlogps)
113         for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
114
115         # perform rmsprop parameter update every batch_size episodes
116         if episode_number % batch_size == 0:
117             for k,v in model.iteritems():
118                 g = grad_buffer[k] # gradient
119                 rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
120                 model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k]) + 1e-5)
121                 grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
122
123         # boring book-keeping
124         running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
125         print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
126         if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
127         reward_sum = 0
128         observation = env.reset() # reset env
129         prev_x = None
130
131     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
132         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

```

# preprocess the observation, set input to network to be difference image
cur_x = prepro(observation)
x = cur_x - prev_x if prev_x is not None else np.zeros(D)
prev_x = cur_x

```

To get temporal difference, we subtract the previous preprocessed image from the current reprocessed image. We then pass this difference to the NN so we have some inherent motion in the network. If we instead passed the raw frame, then we wouldn't have any temporal information in the NN

```

def prepro(I):
    """ prepro 210x160x3 uint8 frame into 6400 (80x80) 1D float vector """
    I = I[35:195] # crop
    I = I[:,::2,::2,0] # downsample by factor of 2
    I[I == 144] = 0 # erase background (background type 1)
    I[I == 109] = 0 # erase background (background type 2)
    I[I != 0] = 1 # everything else (paddles, ball) just set to 1
    return I.astype(np.float).ravel()

```

Preprocessing:

- Take a frame from gym (i.e. observation)
- Crop out unimportant stuff
- Downsample by only taking every even pixel
- Make pixels binary
- Then return its float

Get the current image and preprocess it.

```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprob, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprob else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - aprob) # grad that encourages the action that was taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98     # stack together all inputs, hidden states, action gradients, and rewards for this episode
99     eps = np.vstack(xs)
100     eph = np.vstack(hs)
101     epdlogp = np.vstack(dlogps)
102     epr = np.vstack(drs)
103     xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105     # compute the discounted reward backwards through time
106     discounted_epr = discount_rewards(epr)
107     # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108     discounted_epr -= np.mean(discounted_epr)
109     discounted_epr /= np.std(discounted_epr)
110
111     epdlogp *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
112     grad = policy_backward(eph, epdlogp)
113     for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
114
115     # perform rmsprop parameter update every batch_size episodes
116     if episode_number % batch_size == 0:
117         for k,v in model.iteritems():
118             g = grad_buffer[k] # gradient
119             rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
120             model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k] + 1e-5))
121             grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
122
123     # boring book-keeping
124     running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
125     print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
126     if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
127     reward_sum = 0
128     observation = env.reset() # reset env
129     prev_x = None
130
131     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
132         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

```

# forward the policy network and sample an action from the returned probability
aprob, h = policy_forward(x)
action = 2 if np.random.uniform() < aprob else 3 # roll the dice!

```

```
def policy_forward(x):
```

```
    h = np.dot(model['W1'], x)
```

```
    h[h<0] = 0 # ReLU nonlinearity
```

```
    logp = np.dot(model['W2'], h)
```

```
    p = sigmoid(logp)
```

```
    return p, h # return probability of taking action 2, and hidden state
```

We pass in the difference image to the NN
We return the probability of going up, and the hidden state so we can cache this and do back prop later. This is because we aren't using any NN library.

in the gym environment action 2 = up
and action 3 = down



```
def sigmoid(x):
```

```
    return 1.0 / (1.0 + np.exp(-x)) # sigmoid "squashing" function to interval [0,1]
```

```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprob, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprob else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - aprob) # grad that encourages the action that was taken to be taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99         eps = np.vstack(xs)
100         eph = np.vstack(hs)
101         epdlogp = np.vstack(dlogps)
102         epr = np.vstack(drs)
103         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105         # compute the discounted reward backwards through time
106         discounted_epr = discount_rewards(epr)
107         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108         discounted_epr -= np.mean(discounted_epr)
109         discounted_epr /= np.std(discounted_epr)
110
111         epdlogp *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
112         grad = policy_backward(eph, epdlogp)
113         for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
114
115         # perform rmsprop parameter update every batch_size episodes
116         if episode_number % batch_size == 0:
117             for k,v in model.iteritems():
118                 g = grad_buffer[k] # gradient
119                 rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
120                 model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k]) + 1e-5)
121                 grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
122
123         # boring book-keeping
124         running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
125         print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
126         if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
127         reward_sum = 0
128         observation = env.reset() # reset env
129         prev_x = None
130
131     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
132         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

```

# record various intermediates (needed later for backprop)
xs.append(x) # observation
hs.append(h) # hidden state
y = 1 if action == 2 else 0 # a "fake label"
dlogps.append(y - aprob) # grad that encourages the action that was taken to be taken

```

Bookkeeping so that we can do backpropagation later. If you were to use PyTorch or something, this would not be needed.


```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprobs, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprobs else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - aprobs) # grad that encourages the action that was taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     ddrs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99         eps = np.vstack(xs)
100         eph = np.vstack(hs)
101         epdlogps = np.vstack(dlogps)
102         epr = np.vstack(drs)
103         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105         # compute the discounted reward backwards through time
106         discounted_epr = discount_rewards(epr)
107         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108         discounted_epr -= np.mean(discounted_epr)
109         discounted_epr /= np.std(discounted_epr)
110
111         epdlogps *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
112         grad = policy_backward(eph, epdlogps)
113         for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
114
115         # perform rmsprop parameter update every batch_size episodes
116         if episode_number % batch_size == 0:
117             for k,v in model.iteritems():
118                 g = grad_buffer[k] # gradient
119                 rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
120                 model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k]) + 1e-5)
121                 grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
122
123         # boring book-keeping
124         running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
125         print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
126         if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
127         reward_sum = 0
128         observation = env.reset() # reset env
129         prev_x = None
130
131     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
132         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

```

# record various intermediates (needed later for backprop)
xs.append(x) # observation
hs.append(h) # hidden state
y = 1 if action == 2 else 0 # a "fake label"
dlogps.append(y - aprobs) # grad that encourages the action that was taken to be taken

```

A small piece of backprop:

Derivative of the [log probability of the taken action given this image] with respect to the [output of the network (before sigmoid)]

recall: loss:

$$\sum_i A_i * \log p(y_i | x_i)$$

$$s = W_2 f(W_1 x)$$

$$p = 1 / (1 + e^{-s})$$

$$y \sim p$$

```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprobs, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprobs else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - aprobs) # grad that encourages the action that was taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99         eps = np.vstack(xs)
100         eph = np.vstack(hs)
101         epdlogps = np.vstack(dlogps)
102         epr = np.vstack(drs)
103         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105         # compute the discounted reward backwards through time
106         discounted_epr = discount_rewards(epr)
107         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108         discounted_epr -= np.mean(discounted_epr)
109         discounted_epr /= np.std(discounted_epr)
110
111         epdlogps *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
112         grad = policy_backward(eph, epdlogps)
113         for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
114
115         # perform rmsprop parameter update every batch_size episodes
116         if episode_number % batch_size == 0:
117             for k,v in model.iteritems():
118                 g = grad_buffer[k] # gradient
119                 rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
120                 model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k]) + 1e-5)
121                 grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
122
123         # boring book-keeping
124         running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
125         print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
126         if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
127         reward_sum = 0
128         observation = env.reset() # reset env
129         prev_x = None
130
131     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
132         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

```

# record various intermediates (needed later for backprop)
xs.append(x) # observation
hs.append(h) # hidden state
y = 1 if action == 2 else 0 # a "fake label"
dlogps.append(y - aprobs) # grad that encourages the action that was taken to be taken

```

A small piece of backprop:

Derivative of the [log probability of the taken action given this image] with respect to the [output of the network (before sigmoid)]

recall: loss:

$$\sum_i A_i * \log p(y_i | x_i)$$

output of network
before the sigmoid

$$s = W_2 f(W_1 x)$$

$$p = 1 / (1 + e^{-s})$$

$$y \sim p$$

KL divergence!

$$\text{if } y = 1, L = \log p, dL/ds = 1 - p$$

$$\text{if } y = 0, L = \log(1 - p), dL/ds = -p$$

More compact:

$$L = y \log(p) + (1 - y) \log(1 - p)$$

$$dL/ds = y - p$$


```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprobs, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprobs else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - aprobs) # grad that encourages the action that was taken to be taken (see http://cs231n.github.io/neural-networks-2/#los
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99         eps = np.vstack(xs)
100         eph = np.vstack(hs)
101         epdlogps = np.vstack(dlogps)
102         epr = np.vstack(drs)
103         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105         # compute the discounted reward backwards through time
106         discounted_epr = discount_rewards(epr)
107         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108         discounted_epr -= np.mean(discounted_epr)
109         discounted_epr /= np.std(discounted_epr)
110
111         epdlogps *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
112         grad = policy_backward(eph, epdlogps)
113         for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
114
115         # perform rmsprop parameter update every batch_size episodes
116         if episode_number % batch_size == 0:
117             for k,v in model.iteritems():
118                 g = grad_buffer[k] # gradient
119                 rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
120                 model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k]) + 1e-5)
121                 grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
122
123         # boring book-keeping
124         running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
125         print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
126         if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
127         reward_sum = 0
128         observation = env.reset() # reset env
129         prev_x = None
130
131     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
132         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

```

# step the environment and get new measurements
observation, reward, done, info = env.step(action)
reward_sum += reward

drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)

```

Step the environment

(execute the action, get new state and record the reward)

```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprobs, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprobs else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - aprobs) # grad that encourages the action that was taken to be taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99         epx = np.vstack(xs)
100         eph = np.vstack(hs)
101         epdlogp = np.vstack(dlogps)
102         epr = np.vstack(drs)
103         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105         # compute the discounted reward backwards through time
106         discounted_epr = discount_rewards(epr)
107         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108         discounted_epr -= np.mean(discounted_epr)
109         discounted_epr /= np.std(discounted_epr)
110
111         epdlogp *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
112         grad = policy_backward(eph, epdlogp)
113         for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
114
115         # perform rmsprop parameter update every batch_size episodes
116         if episode_number % batch_size == 0:
117             for k,v in model.iteritems():
118                 g = grad_buffer[k] # gradient
119                 rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
120                 model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k]) + 1e-5)
121                 grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
122
123         # boring book-keeping
124         running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
125         print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
126         if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
127         reward_sum = 0
128         observation = env.reset() # reset env
129         prev_x = None
130
131     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
132         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

```
if done: # an episode finished
```

```
episode_number += 1
```

```
# stack together all inputs, hidden states, action gradients, and rewards for this episode
```

```
epx = np.vstack(xs)
```

```
eph = np.vstack(hs)
```

```
epdlogp = np.vstack(dlogps)
```

```
epr = np.vstack(drs)
```

```
xs,hs,dlogps,drs = [],[],[],[] # reset array memory
```

Once a rollout is done,
Concatenate together all images, hidden
states, etc. that were seen in this batch.

Again, if using PyTorch, no need to do this.

```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprobs, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprobs else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(- aprobs) and grad that encourages the action that was taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99         eps = np.vstack(xs)
100         eph = np.vstack(hs)
101         epdlogps = np.vstack(dlogps)
102         epr = np.vstack(drs)
103         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105         # compute the discounted reward backwards through time
106         discounted_epr = discount_rewards(epr)
107         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108         discounted_epr -= np.mean(discounted_epr)
109         discounted_epr /= np.std(discounted_epr)
110
111         # boring book-keeping
112         running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
113         print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
114         if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
115         reward_sum = 0
116         observation = env.reset() # reset env
117         prev_x = None
118
119     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
120         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

compute the discounted reward backwards through time

discounted_epr = discount_rewards(epr)

standardize the rewards to be unit normal (helps control the gradient estimator variance)

discounted_epr -= np.mean(discounted_epr)

discounted_epr /= np.std(discounted_epr)

epr = episode rewards

We want to compute an Advantage. We want to the discounted reward to be the Advantage.

def discount_rewards(r):

""" take 1D float array of rewards and compute discounted reward """

discounted_r = np.zeros_like(r)

running_add = 0

for t in reversed(xrange(0, r.size)):

if r[t] != 0: running_add = 0 # reset the sum, since this was a discountable action. And we also

running_add = running_add * gamma + r[t]

discounted_r[t] = running_add

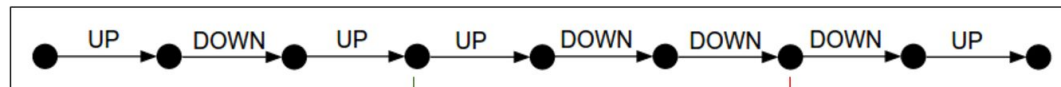
return discounted_r

But to reduce the variance in the Advantage we will get it into zero-mean unit variance. So we are going to subtract the mean and divide by the standard deviation. This will center and scale our rewards. This will allow us to encourage half the actions and discourage half the actions. And we also make sure that the Advantage is on the same scale across rollouts/episodes. So we do this over a batch of rollouts

Discounted rewards

$$\sum_i A_i * \log p(y_i | x_i)$$

0.21 0.24 0.27 -0.81 -0.9 -1 0 0



Reward +1.0

Reward -1.0

```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprobs, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprobs else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - aprobs) # grad that encourages the action that was taken to be taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99         eps = np.vstack(xs)
100         eph = np.vstack(hs)
101         epdlogps = np.vstack(dlogps)
102         epr = np.vstack(drs)
103         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105         # compute the discounted reward backwards through time
106         discounted_epr = discount_rewards(epr)
107         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108         discounted_epr -= np.mean(discounted_epr)
109         discounted_epr /= np.std(discounted_epr)
110
111         epdlogps *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
112         grad = policy_backward(eph, epdlogps)
113         for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
114
115         # perform rmsprop parameter update every batch_size episodes
116         if episode_number % batch_size == 0:
117             for k,v in model.iteritems():
118                 g = grad_buffer[k] # gradient
119                 rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
120                 model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k]) + 1e-5)
121                 grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
122
123         # boring book-keeping
124         running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
125         print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
126         if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
127         reward_sum = 0
128         observation = env.reset() # reset env
129         prev_x = None
130
131     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
132         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

here we get the weight matrix for W1 and W2. We can then update the network. By doing this the network will be more likely to take actions that increase reward, and less likely to take actions that decrease reward

```

epdlogps *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
grad = policy_backward(eph, epdlogps)
for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch

```

$$\sum_i A_i * \log p(y_i | x_i)$$

Advantage modulation

```

def policy_backward(eph, epdlogps):
    """ backward pass. (eph is array of intermediate hidden states) """
    dw2 = np.dot(eph.T, epdlogps).ravel()
    dh = np.outer(epdlogps, model['W2'])
    dh[eph <= 0] = 0 # backpro prelu
    dw1 = np.dot(dh.T, epx)
    return {'W1':dw1, 'W2':dw2}

```

backprop!!!!!!1

```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     apro, h = policy_forward(x)
81     action = 2 if np.random.uniform() < apro else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - apro) # grad that encourages the action that was taken to be taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99
100         eps = np.vstack(xs)
101         eph = np.vstack(hs)
102         epdlogp = np.vstack(dlogps)
103         epr = np.vstack(drs)
104         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
105
106         # compute the discounted reward backwards through time
107         discounted_epr = discount_rewards(epr)
108         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
109         discounted_epr -= np.mean(discounted_epr)
110         discounted_epr /= np.std(discounted_epr)
111
112         epdlogp *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
113         grad = policy_backward(eph, epdlogp)
114         for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
115
116     # perform rmsprop parameter update every batch_size episodes
117     if episode_number % batch_size == 0:
118         for k,v in model.iteritems():
119             g = grad_buffer[k] # gradient
120             rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
121             model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k]) + 1e-5)
122             grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
123
124     # boring book-keeping
125     running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
126     print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
127     if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
128     reward_sum = 0
129     observation = env.reset() # reset env
130     prev_x = None
131
132     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
133         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' ' if reward == -1 else '!!!!!!!')

```

```

# perform rmsprop parameter update every batch_size episodes
if episode_number % batch_size == 0:
    for k,v in model.iteritems():
        g = grad_buffer[k] # gradient
        rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
        model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k]) + 1e-5)
        grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer

```

Use RMSProp for the parameter update.

RMSProp

Update rule:

$$R_t = \gamma R_{t-1} + (1 - \gamma) \nabla L_t(W_{t-1})^2$$

$$W_t = W_{t-1} - \alpha \frac{\nabla L_t(W_{t-1})}{\sqrt{R_t}}$$

Similar to AdaGrad but with an exponential moving average controlled by $\gamma \in [0, 1]$ (smaller $\gamma \implies$ more emphasis on recent gradients).


```

64 env = gym.make("Pong-v0")
65 observation = env.reset()
66 prev_x = None # used in computing the difference frame
67 xs,hs,dlogps,drs = [],[],[],[]
68 running_reward = None
69 reward_sum = 0
70 episode_number = 0
71 while True:
72     if render: env.render()
73
74     # preprocess the observation, set input to network to be difference image
75     cur_x = prepro(observation)
76     x = cur_x - prev_x if prev_x is not None else np.zeros(D)
77     prev_x = cur_x
78
79     # forward the policy network and sample an action from the returned probability
80     aprobs, h = policy_forward(x)
81     action = 2 if np.random.uniform() < aprobs else 3 # roll the dice!
82
83     # record various intermediates (needed later for backprop)
84     xs.append(x) # observation
85     hs.append(h) # hidden state
86     y = 1 if action == 2 else 0 # a "fake label"
87     dlogps.append(y - aprobs) # grad that encourages the action that was taken (see http://cs231n.github.io/neural-networks-2/#loss)
88
89     # step the environment and get new measurements
90     observation, reward, done, info = env.step(action)
91     reward_sum += reward
92
93     drs.append(reward) # record reward (has to be done after we call step() to get reward for previous action)
94
95     if done: # an episode finished
96         episode_number += 1
97
98         # stack together all inputs, hidden states, action gradients, and rewards for this episode
99         eps = np.vstack(xs)
100         eph = np.vstack(hs)
101         epdlogps = np.vstack(dlogps)
102         epr = np.vstack(drs)
103         xs,hs,dlogps,drs = [],[],[],[] # reset array memory
104
105         # compute the discounted reward backwards through time
106         discounted_epr = discount_rewards(epr)
107         # standardize the rewards to be unit normal (helps control the gradient estimator variance)
108         discounted_epr -= np.mean(discounted_epr)
109         discounted_epr /= np.std(discounted_epr)
110
111         epdlogps *= discounted_epr # modulate the gradient with advantage (PG magic happens right here.)
112         grad = policy_backward(eph, epdlogps)
113         for k in model: grad_buffer[k] += grad[k] # accumulate grad over batch
114
115         # perform rmsprop parameter update every batch_size episodes
116         if episode_number % batch_size == 0:
117             for k,v in model.iteritems():
118                 g = grad_buffer[k] # gradient
119                 rmsprop_cache[k] = decay_rate * rmsprop_cache[k] + (1 - decay_rate) * g**2
120                 model[k] += learning_rate * g / (np.sqrt(rmsprop_cache[k] + 1e-5))
121                 grad_buffer[k] = np.zeros_like(v) # reset batch gradient buffer
122
123         # boring book-keeping
124         running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
125         print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
126         if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
127         reward_sum = 0
128         observation = env.reset() # reset env
129         prev_x = None
130
131     if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
132         print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' if reward == -1 else ' !!!!!!!!')

```

```
# boring book-keeping
```

```

running_reward = reward_sum if running_reward is None else running_reward * 0.99 + reward_sum * 0.01
print 'resetting env. episode reward total was %f. running mean: %f' % (reward_sum, running_reward)
if episode_number % 100 == 0: pickle.dump(model, open('save.p', 'wb'))
reward_sum = 0
observation = env.reset() # reset env
prev_x = None

```

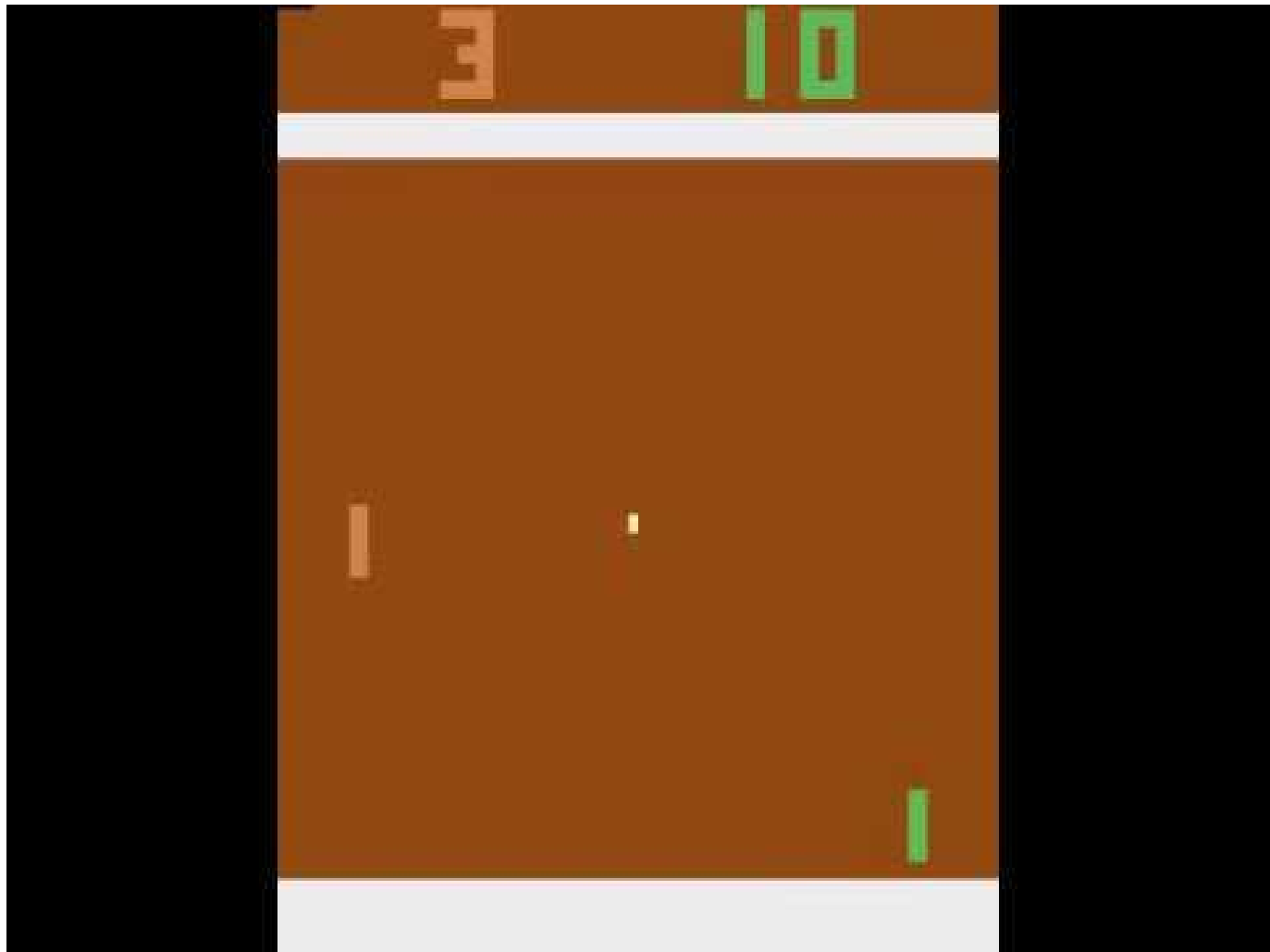
```
if reward != 0: # Pong has either +1 or -1 reward exactly when game ends.
```

```
print ('ep %d: game finished, reward: %f' % (episode_number, reward)) + (' if reward == -1 else ' !!!!!!!!')
```

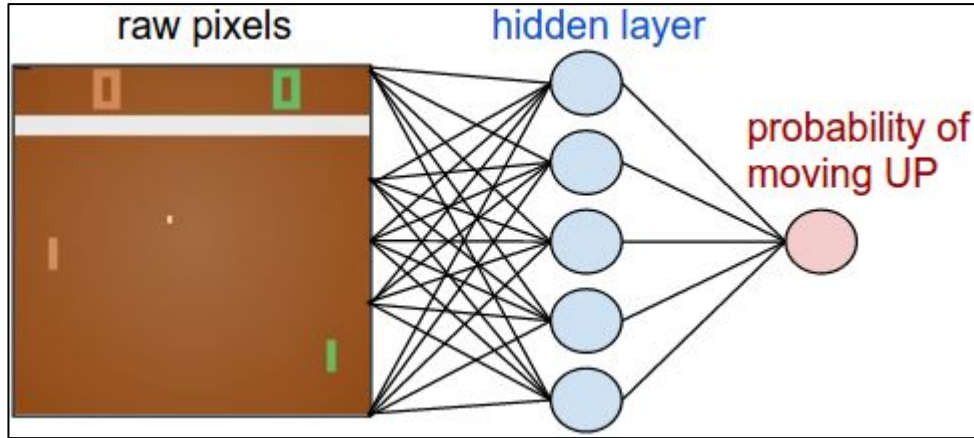
prints etc

In summary

1. Initialize a policy network at random
2. **Repeat Forever:**
3. Collect a bunch of rollouts with the policy
4. Increase the probability of actions that worked well
5. ???
6. Profit.



Thank you! Questions?



$$\sum_i A_i * \log p(y_i | x_i)$$

