# Pong from Pixels

Deep RL Bootcamp
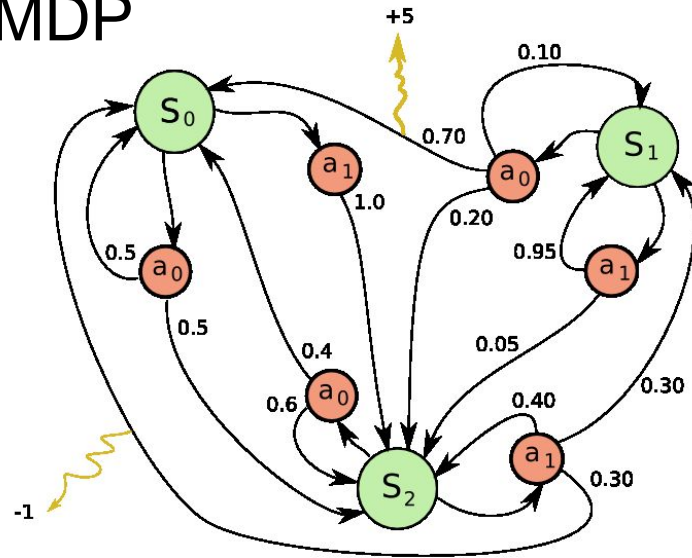
Andrej Karpathy, Aug 26, 2017
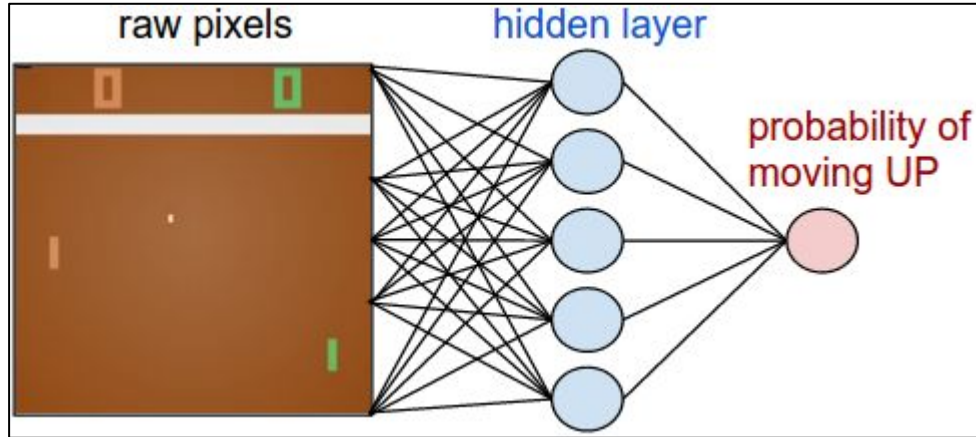
MDP

# Policy network
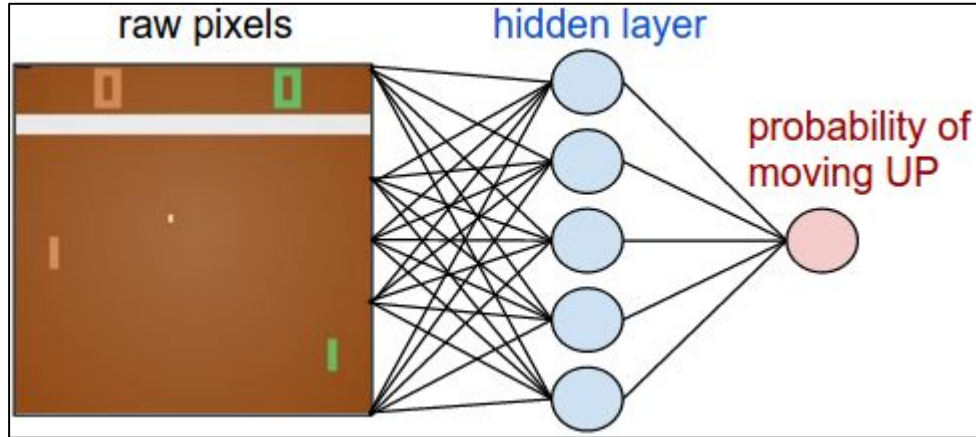
# Policy network

e.g.,

height  width

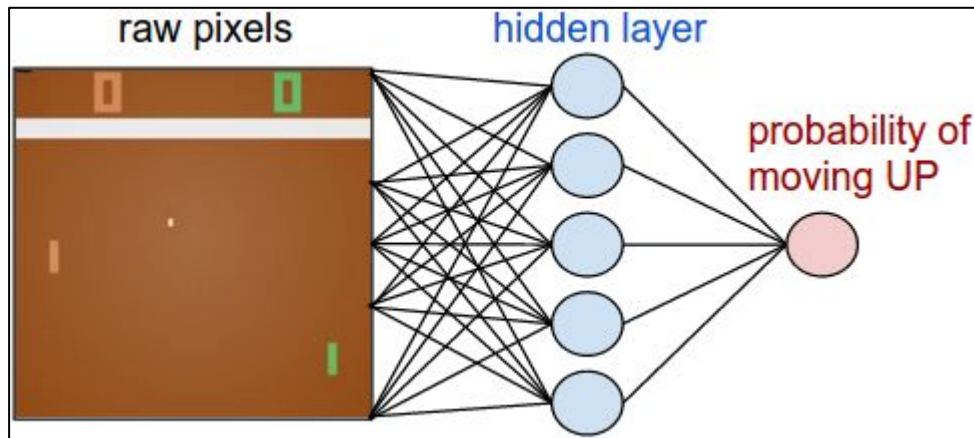[**80** x **80**]
array of

# Policy network

height width

[**80** x **80**]
array



```
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)
```

# Policy network



height  width

[**80** x **80**]
array

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

[(80*80)*200 + 200] + [200*1 + 1] = **~1.3M parameters**

Layer 1                    Layer 2

Network does not see this. Network sees 80*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...

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. **Repe**
2.    Sa
3.    Ru
4.    If t
5. Retur

# 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)

...

maximize:

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



raw pixels · hidden layer · probability of moving UP

Except, we don't have labels...



raw pixels    hidden layer    probability of moving UP

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...



Rollout the policy and collect an episode

# Collect many rollouts...

**4 rollouts:**

# 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: $\log p(y_i \mid x_i)$

Pretend every action we took here was the wrong label.

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

# 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:

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

# 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:

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

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

# 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:

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

+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.



Reward +1.0

Reward -1.0

# Discounting

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

Discounted rewards

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

0.21   0.24   0.27   -0.81   -0.9   -1   0   0



Reward +1.0

Reward -1.0

\gamma = 0.9

time

**discounted reward**

episode start

ball gets past our paddle

-1 reward

we screwed up
somewhere here

this was all hopeless and only
contributes noise to the gradient

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

```python
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()

  # 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

  # 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!

  # 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 (see http://cs231n.github.io/neural-networks-2/#los

  # 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)

  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

    # 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)

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

    # 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

    # 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 ' !!!!!!!!')
```

```python
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()
```

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    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/#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      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 ' !!!!!!!!')
```

```
# 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
```

```
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()
```

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 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       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 ' !!!!!!!!')
```

```python
# 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!
```

```python
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
```

```python
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/#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      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 ' !!!!!!!!')
```

```
# 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.

```
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()

  # 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

  # 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!

  # 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 (see http://cs231n.github.io/neural-networks-2/#los

  # 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)

  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

    # 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)

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

    # 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

    # 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 ' !!!!!!!!')
```

```
# 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
```

# 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    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/#los

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 ' !!!!!!!!')
```

```
# 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
```

# 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:

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

$$s = W_2 f(W_1 x)$$
$$p = 1/(1 + e^{-s})$$
$$y \sim p$$

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

More compact:

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

```python
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()

  # 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

  # 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!

  # 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 (see http://cs231n.github.io/neural-networks-2/#los

  # 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)

  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

    # 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)

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

    # 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

    # 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 ' !!!!!!!!')
```

```python
# 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)

```python
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()

  # 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

  # 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!

  # 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 (see http://cs231n.github.io/neural-networks-2/#los

  # 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)

  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

    # 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)

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

    # 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

    # 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 ' !!!!!!!!')
```

```python
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    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/#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      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 ' !!!!!!!!')
```

```
# 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)
```
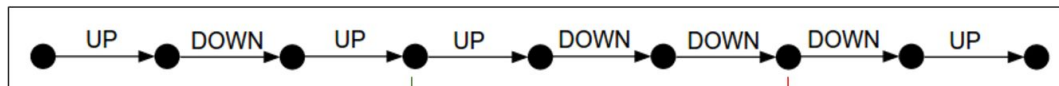
```
def discount_rewards(r):
  """ take 1D float array of rewards and compute discounted reward """
  discounted_r = np.zeros_like(r)
  running_add = 0
  for t in reversed(xrange(0, r.size)):
    if r[t] != 0: running_add = 0 # reset the sum, since this was a game boundary (pong specific!)
    running_add = running_add * gamma + r[t]
    discounted_r[t] = running_add
  return discounted_r
```



Discounted rewards

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

| 0.21 | 0.24 | 0.27 | -0.81 | -0.9 | -1 | 0 | 0 |

UP → DOWN → UP → UP → DOWN → DOWN → DOWN → UP

Reward +1.0

Reward -1.0

```python
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()

  # 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

  # 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!

  # 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 (see http://cs231n.github.io/neural-networks-2/#los

  # 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)

  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

    # 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)

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

    # 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

    # 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 ' !!!!!!!!')
```

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

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

Advantage modulation

```python
def policy_backward(eph, epdlogp):
  """ backward pass. (eph is array of intermediate hidden states) """
  dW2 = np.dot(eph.T, epdlogp).ravel()
  dh = np.outer(epdlogp, model['W2'])
  dh[eph <= 0] = 0 # backpro prelu
  dW1 = np.dot(dh.T, epx)
  return {'W1':dW1, 'W2':dW2}
```

backprop!!!!!!1

```python
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()

  # 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

  # 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!

  # 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 (see http://cs231n.github.io/neural-networks-2/#los

  # 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)

  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

    # 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)

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

    # 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

    # 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 ' !!!!!!!!')
```

```python
# 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).

```python
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()

  # 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

  # 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!

  # 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 (see http://cs231n.github.io/neural-networks-2/#los

  # 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)

  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

    # 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)

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

    # 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

    # 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 ' !!!!!!!!')
```
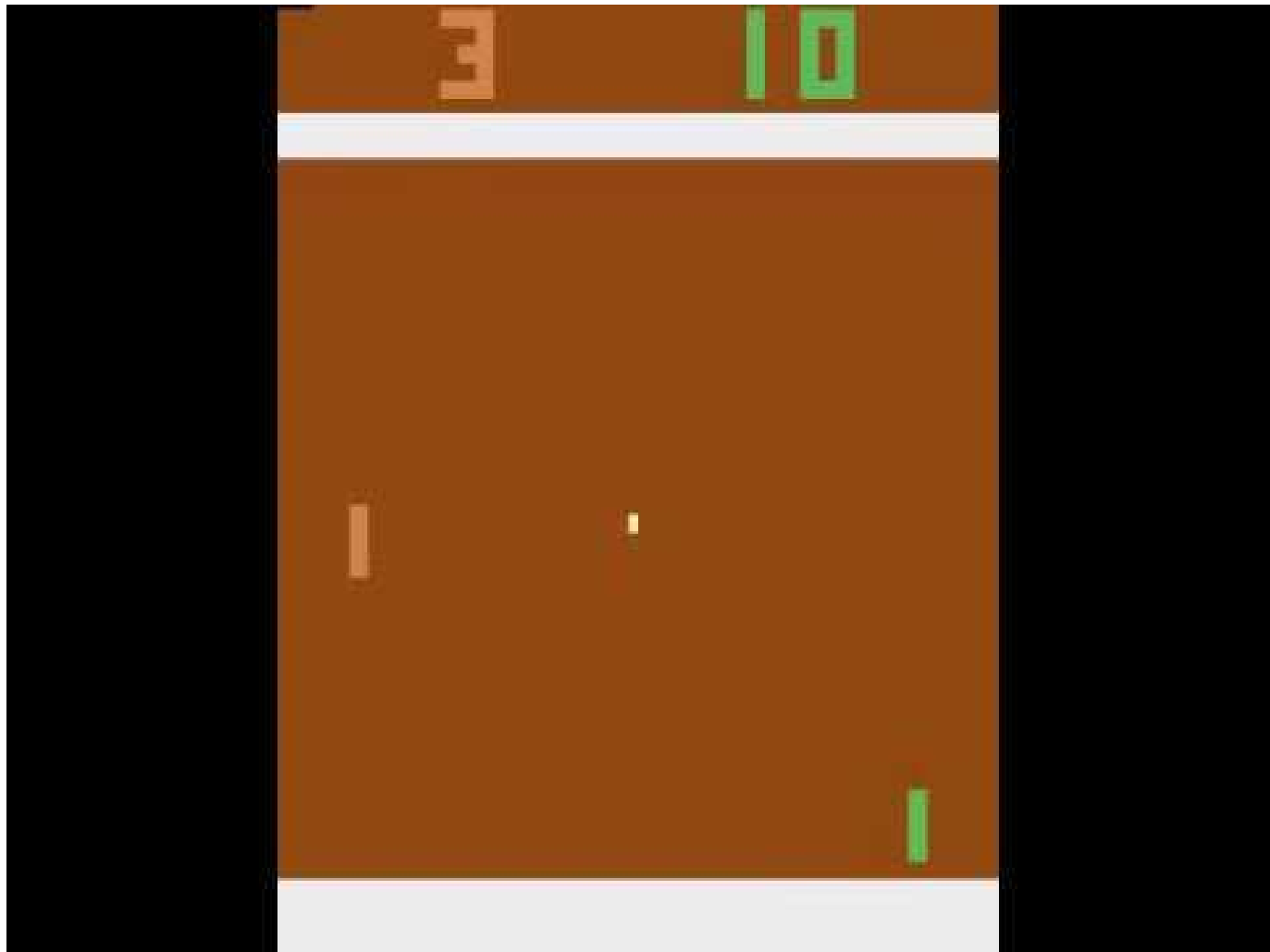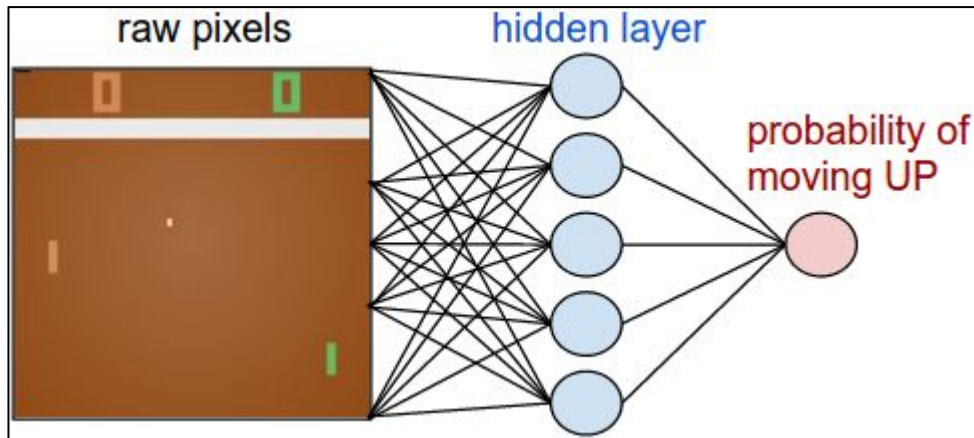
prints etc

# 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?



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