# Learning about (Deep) Reinforcement Learning

#### Scott Rome

romescott@gmail.com

## About Me



- 2 Cats
- Currently the Lead Data Scientist at an ad tech firm
  - Previously a Data Scientist in healthcare
- Ph.D. in Mathematics
- Run a ML blog: <u>http://srome.github.io</u>

# You May Also Know Me From...



About 1,370,000 results (0.46 seconds)

#### python BeautifulSoup parsing table - Stack Overflow

https://stackoverflow.com/questions/23377533/python-beautifulsoup-parsing-table Mar 21, 2016 - Here you go: data = [] table = soup.find('table', attrs={'class':'lineItemsTable'}) table\_body = table.find('tbody') rows = table\_body.find\_all('tr') for row in rows: cols = row.find\_all('td') cols = [ele.text.strip() for ele in cols] data.append([ele for ele in cols if ele]) # Get rid of empty values. This gives you: [ [u'1359711259', u'SRF', ...

beautifulsoup - Python Beautiful Soup parse a table with a	Oct 25, 2013
How to parse html table with python and beautifulsoup and write to	Mar 6, 2013
python - BeautifulSoup: Get the contents of a specific table	Nov 16, 2011
python - parsing table with BeautifulSoup and write in text file	Feb 8, 2010
More results from stackoverflow.com	





srome.github.io/Parsing-HTML-Tables-in-Python-with-BeautifulSoup-and-pandas/ -

May 30, 2016 - Something that seems daunting at first when switching from R to Python is replacing all the ready-made functions R has. For example, R has a nice CSV reader out of the box. Python users will eventually find pandas, but what about other R libraries like their HTML **Table** Reader from the xml package?

# Sadly, I'm A One Hit Wonder

What pages do your users visit?

Page	Pageviews	Page Value
/Parsing-HTML-TablesfulSoup-and-pandas/	1,462	\$0.00
/An-Annotated-Proofplementation-Notes/	92	\$0.00
/Train-A-Neural_Net-TJack-With-Q-Learning/	54	\$0.00
/A-Tour-Of-Gotchas-Wras-And-OpenAi-Gym/	51	\$0.00
/Build-Your-Own-Eventacktester-In-Python/	51	\$0.00
/Async-SGD-in-Pythonplementing-Hogwild!/	50	\$0.00
/	44	\$0.00
/Eigenvesting-I-LinearYour-Stock-Portfolio/	35	\$0.00
/archive/	30	\$0.00
/Visualizing-the-Learetwork-Geometrically/	25	\$0.00
Last 7 days 🔻	PAGES REPORT >	

**PSA: Volunteers Wanted!** 



# You should leave with...

- Some ideas about what (deep) reinforcement learning is
- An understanding of the basic training process for reinforcement learning
- Having seen more than enough code snippets from <a href="https://github.com/srome/ExPyDQN">https://github.com/srome/ExPyDQN</a>
- Not having heard all the annoying gotchas when implementing a reinforcement learning framework (see the post <u>"A Tour of Gotchas When</u> <u>Implementing Deep Q Networks with Keras and OpenAi Gym"</u> from my blog for that).

# **Reinforcement Learning: A Brief History**

- (1950s) Roots in Optimal Control / dynamic programming, developed by Bellman at the RAND Corporation
  - Called "dynamic programming" because the Secretary of Defense "actually had a pathological fear and hatred of the word, research"
- (1977) Formal study of Temporal Difference Learning began
- (1989) Q-Learning was published
- (1992) TD-Gammon developed
- (2013) Deep Q Learning published (Neural Networks used to play Atari)
- (2016) AlphaGo beats Lee Sidol
- (2017) AlphaGo Zero beats Alpha Go

Notice, the activity starting around the early 2010's... that is in part due to...



# **Goal of Reinforcement Learning**

# Training an agent (e.g., a model) to interact with an environment



Follow ~

Deep RL is popular because it's the only area in ML where it's socially acceptable to train on the test set.

12:27 PM - 28 Oct 2017

# Training Algorithm at a High Level





In []: env # emulator environment
state # initial state from environment
agent # the thing we want to train
while game\_still\_going:
 action = agent.get\_action(state) # get action from the agent
 next\_state, reward, game\_still\_going = env.step(action) # pass action to emulat
or, get next state
 agent.update(state, next\_state, action, reward) # update the agent

# Differences vs. "Typical" Machine Learning Problems

- Training data is usually generated during training
- The agent (model) is not independent from the data it is trained on:
  - The agent affects the generation of new training data.
  - The target variable at each update step (as we will see) depends on a version of the agent.
- It's possible to miss part of your possible training set due to data collection decisions

In [ ]: while game\_still\_going:

action = agent.get\_action(state)
next\_state, reward, game\_still\_going = env.step(action)
agent.update(state, next state, action, reward)

# **Big Questions**



- What libraries do we use?
- What (who?) is the agent?
- How do we select  $a_t$ ?
- What is *s*<sub>*t*</sub>?
- What is our target  $y_t$ ?
  - What is the minibatch?

# **Python Libraries**



\* Disclaimer: Not a graphic artist.

#### The Agent

From (3),





symbolizes sliding of each filter across input image) and two fully connected layers with a single output for each valid action. Each hidden layer is followed by a rectifier nonlinearity (that is, max(0,x)).

#### How Do We Select $a_t$ : Exploration / Exploitation

The "Multi-Armed Bandit" Problem



"Originally considered by Allied scientists in World War II, it proved so intractable that, according to Peter Whittle, the problem was proposed to be dropped over Germany so that German scientists could also waste their time on it." - Wikipedia

# Approximate Solution to MAB: $\epsilon$ - Greedy Strategy

- For each state, select a random action (level) with probability  $\epsilon$ .
- Otherwise, allow the agent to choose the action.
- $\epsilon$  annealing is used to allow the agent to select more actions over time.

```
In [ ]: def get_action(self, image=None, epsilon_override=None):
    epsilon = self.get_epsilon(epsilon_override) # Annealing code is hidden in here
    if np.random.uniform(0, 1) < 1 - epsilon:
        rewards = self._get_q_values(image=image)
        action = self._actions[np.argmax(rewards)]
    else:
        action = np.random.choice(self._actions)
    return action</pre>
```

# What is $s_t$ ?

The emulator returns an image after each action. To form  $s_t$ , we must

- Drop image to grayscale
- Frame skip
  - Every *n*-th frame is considered for the definition of the current state
  - n is sometimes referred to as  $\phi$  length
- Frames are stacked as an input
  - For a (84,84) grayscale image and frame skip of 4: an instance state has the dimension (4,84,84,1)
- Consecutive max
  - Take pixel-wise max of (*n* − 1)-th and *n*-th image from the emulator as the *n*-th frame.

From (5),



# Frame Skip / Consecutive Max Code during Emulator Step

```
In [6]: | def step(self, action):
             """ This relies on the fact that the underlying environment creates new images
         for each frame.
                 By default, opengym uses atari py's getScreenRGB2 which creates a new array
        for each frame."""
             total_reward = 0
             obs = None
             for k in range(self.frame skip):
                 last obs = obs
                 # Make sure you are using a "NoFrameskip-v4" ROM
                 obs, reward, is terminal, info = self. env.step(action)
                 total reward += reward
                 if is terminal:
                     # End episode if is terminal
                     if k == 0 and last obs is None:
                         last obs = obs
                     break
             if self.consecutive max and self.frame skip > 1:
                 obs = np.maximum(last obs, obs)
             # Store observation, greyscale and cropping are applied in here
             self.store observation(obs)
             return self. get current state(), total reward, is terminal
```

#### As Promised in the Abstract...



# What is $y_t$ ?

Let  $r' \in \mathbb{R}$  be a reward from performing action a at state s. The action-value function Q provides a mapping from  $(s, a) \rightarrow r'$ .

- Using Q, you can define a policy for the agent.
- Greedy Policy:
  - At each time step, the agent selects the action associated with the highest reward.

#### **More Definitions**

Assume  $r \in \mathbb{R}$  is the resulting reward from the emulator for performing *a* at state *s*, and let  $\gamma \in \mathbb{R}$  be a discount.

We define the optimal policy as the actions chosen via  $Q^*$ , which satisfies the *Bellman Equation*:

$$Q^*(s,a) = \mathbb{E}_{s'\sim\mathcal{E}}\left[r + \gamma \max_{a'} Q^*(s',a') \middle| s,a\right],$$

i.e., the reward plus the discounted future rewards for following the policy.

# $y_t$ defined

We define our neural network with weights  $\theta$  as  $Q(s, a, \theta)$ . For each example  $(s_t, s_{t+1}, a_t, r_t)$ , we select weights  $\theta_t$  and define the training target to be:

$$y_t(\theta_t) := r_t + \gamma \max_{a'} Q(s_{t+1}, a', \theta_t)$$

- Notice, our target depends on  $\theta_t$ -- it's not fixed like traditional machine learning problems!
- Some out-of-scope math proves that this  $y_t$  allows us to learn  $Q^*$ 
  - (Technically, we don't optimize a single loss function, but a sequence of loss functions)

# Wait? $y_t$ depends on $\theta$ ? What $\theta$ ?

- New targets are generated every minibatch using a recent, fixed set of weights  $\hat{\theta}_t$ .
- $\hat{\theta}_t$  is typically updated with recent weights every ~10000 steps (for example).
  - i.e.  $\hat{\theta}_t = \hat{\theta}_{t_i}$  for some fixed older set of weights.
- This technique is referred to as the fixed target network.

# Other notes on $y_t$

- Notationally, it is easier to define the neural network as  $Q(s, a, \theta)$ .
- Practically, it is defined as  $Q(s, \theta)$  and  $y_t \in \mathbb{R}^{|A|}$ , i.e.  $y_t$  has a reward entry for every action.

# Calculating the $\max_{a'} Q(s_{t+1}, a', \theta_i)$

```
In [4]:
        def build training_variables(self, future_states, actions, is_terminal):
            future rewards = self. get model(fixed=True).predict(future states) # Get the f
        ixed target network!
            max rewards = np.max(future rewards, axis=1).reshape((len(future rewards), 1))
            # Code for an optimization around training
            # which updates reward for ONLY the completed action
            # training mask masks rewards for other actions so
            # we only update weights according to the action we chose
            training mask = np.zeros(future rewards.shape)
            action indexes = np.array([[np.where(self. actions == action)[0][0]] for action
        in actions])
            is terminal = is terminal * 1. # convert to float
            for index0, index1 in zip(range(self. minibatch), action indexes):
                training mask[index0, index1] = 1.
            return max rewards, training mask, is terminal
```

# What is the minibatch?

- Time correlation of examples causes divergence during training
- The solution is "Replay Memory"
  - Store *m* examples in memory *D* and sample minibatches from *D* for training!

```
In [ ]: def sample(self):
    # Code to sample from the replay memory
    ind = np.random.choice(self._size, size=self._mini_batch_size)
    # Avoiding a copy action as much as possible
    self._mini_batch_state[:] = self._memory_state[ind,:,:,:]
    self._mini_batch_future_state[:] = self._memory_future_state[ind,:,:,:]
    rewards = self._rewards[ind]
    is_terminal = self._is_terminal[ind]
    actions = self._actions[ind]
    return self._mini_batch_state, self._mini_batch_future_state, actions, rewards,
    is_terminal
```

# **Final Training Algorithm**

```
In []:
        def run episode(self, training=True):
             state = self.reset environment()
             is terminal = False
             total reward = 0
             steps = 0
            while not is terminal:
                 if not training:
                     self.render() # Speeds up training to not display it
                 action = self.get action(state, training)
                 next state, reward, is terminal = self.step(action)
                 if training:
                     clipped reward = np.clip(reward ,-self.reward clip ,self.reward clip)
                     self.replay memory.add example(state=state,
                                                    future state=next state,
                                                    action=action,
                                                    reward=clipped reward,
                                                    is terminal=is terminal)
                     self.training step()
                 state = next state
                 total reward += reward
                 steps += 1
             return total reward, steps # for logging
```

# **Training Results**



# Thank You

- Slides will be available on <a href="http://srome.github.io">http://srome.github.io</a>
- The (functional) framework code is available at <a href="https://github.com/srome/ExPyDQN">https://github.com/srome/ExPyDQN</a>.
- Check out TechGirlz!

# References

- 1. Playing Atari with Deep Reinforcement Learning
- 2. Wikipedia: Reinforcement Learning
- 3. <u>Human-level control through deep reinforcement learning</u>
- 4. Wikipedia: Banach fixed-point theorem
- 5. Frame Skipping and Preprocessing for Deep Q Networks on Atari 2600
- 6. Wikipedia: Multi-Armed Bandit
- 7. Google Image Search
- 8. Scholarpedia: Reinforcement Learning
- 9. Richard Bellman on the Birth of Dynamic Programming