

Deep Q-learning. \rightarrow perform action a .

Risk environment: first state

Q-learning or function approximation

what we do:

Deep Q-Learning (DQN) Review

- Basic Q-Learning:

```
s = env.reset()  $\rightarrow$  env.reset. Step 1  
while not done:  
    a = argmax(Q[s, :]) # can also use epsilon-greedy /  
    s', r, done = env.step(a)  $\begin{matrix} \text{action} \\ \text{giving open play} \end{matrix}$  across all actions  
    td_err = (r + gamma * max { Q[s', :] } - Q[s, a])^2  
    params = params - learning_rate * grad(td_err, params)  
    s = s'
```

TD-error: $\text{Jeni gradient} - \text{target} \mid \text{square target}$

true target \rightarrow estimate

diff of once

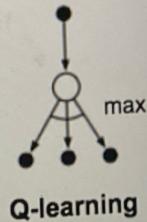
Details (semi gradient)

Estimated target

$$(r + \gamma \max \{ Q[s', :] \}) - Q[s, a]^2$$

Prediction

- It's not true gradient descent (we call it "semi-gradient")
- In supervised learning, we normally have $(\text{target} - \text{prediction})^2$
- In Q-Learning, the target depends on the prediction! *only estimate* *not true* *(target - pred)*²
- The target does not tell us the *true* Q value
- Even though Q shows up twice (once in target, once in prediction) we only differentiate wrt prediction
- We pretend the target is fixed
- Q-learning is off-policy
- The target is still $r + \gamma Q(s', a')$ where $a' = \arg\max\{Q(s', :)\}$
- Even if we don't take the action a' in the next step



- *for semi-gradient*
- *Target depends on the prediction*
- *Estimate based on own network* / model predictions are off-policy learning algorithm.
- off-policy learning algorithm
- diff. once wrt Q
- *Q-learning - off policy learning.*
- $$Q = r + \gamma \max \{ Q[s', a'] \}$$
 → *target*

$$a' = \arg\max \{ Q[s', :] \}$$
 same policy regardless of action I take.

→ fixed feature extract?
Q learning feature approximator
 N.N unstable → cannot work with fixed feature extraction with linear model on top.

Deep Q-learning tricks

N.N → deep learning
 $(s, a, r, s', \text{done})$ - tuple } experienced replay memory / buffer
 Input \leftrightarrow tuple form

$$\text{Target: } r + \gamma \times (\text{done}) \times \max(Q(s', :))$$

$$\text{Pred: } Q(s, a)$$

Why $(\text{done}) = 1$, then target = r ~~fit~~
 $(\text{done}) = 0$, then target = 0 $\sum \text{reward from}$
 terminal state $s' = 0$

$$Q(s, a) \rightarrow 0$$

Dr learning tricks
 replay buffer = 1 million | most recent sample
 perform sampling on correlated-
 $(s, a, r, s', \text{done})$ or the batch from the replay buffer
 cal (input (tuple)) → $\boxed{\text{grad}(s, a, r, s', \text{done})}$

experience replay buffer Strategy

Trick 1: Replay buffer | size 1 million
gradients $\{s, a, r, s', \text{done}\}$ but $\begin{cases} \text{sample from batch} \\ \text{sequential states episodes are correlated} \end{cases}$
most recent

Experience replay pseudocode

```
s = env.reset()  
while not done:  
    a = argmax(Q[s, :]) # can also use epsilon-greedy  
    s', r, done = env.step(a) optimal action (a)  
    replay_buffer.add((s, a, r, s', done)) action perform add replay buffer  
    inputs, targets = sample data from replay_buffer  
    params = params - lr * grad((targets - predictions)^2, params)  
    s = s'
```

Sampling Data: 1 million past samples to choose from

Trick 2: Use separate target network

$$Old: r + \gamma(1 - d) \max_{a'} \{Q(s', a')\}$$

$$New: r + \gamma(1 - d) \max_{a'} \{Q_{targ}(s', a')\}$$

Targt is not from target / semi-gradient).

$$Old: r + \gamma(1 - d) \max_{a'} \{Q(s', a')\}$$

$$New: r + \gamma(1 - d) \max_{a'} \{Q_{targ}(s', a')\}$$

copy

target steps

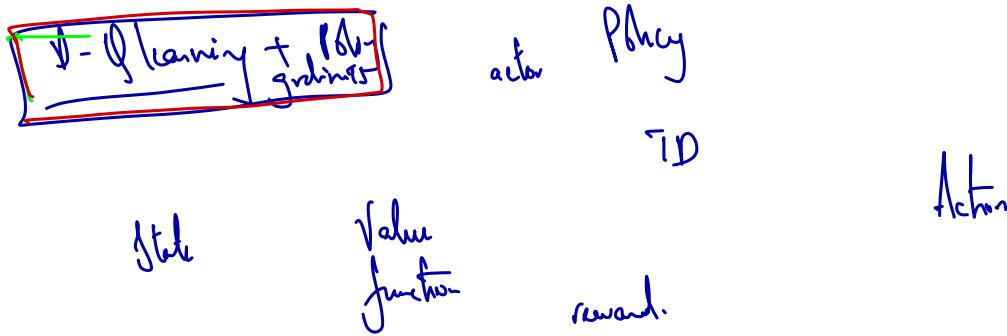
Pseudocode

```
s = env.reset()
step = 0
while not done:
    a = argmax(Q[s, :]) # can also use epsilon-greedy
    s', r, done = env.step(a)
    replay_buffer.add((s, a, r, s', done))
    inputs, targets = sample data from replay_buffer
    params = params - lr * grad((targets - predictions)^2, params)
    s = s'
    step++
    if step % 10000 == 0:
        copy params to target network
```

if steps are multiples of 10000, then copy main
parameters to target network.

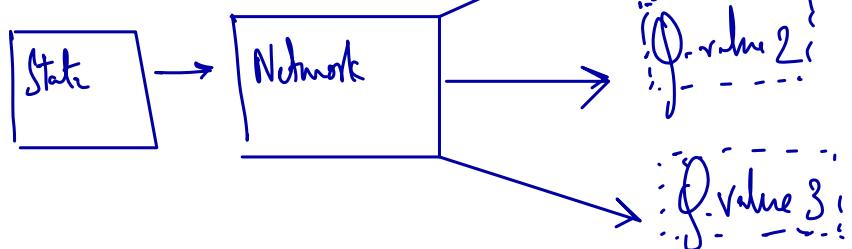
Replay buffer
 Target network } $\xrightarrow{\text{Holds}} \text{Stability}$ → } allow to use C.N.N with
 Q-learning

D.D.P.G : Deep Deterministic Policy Gradient



Cannot handle continuous
 action spaces

Limitations of Q & QN
 $Q(s, a)$



continuous scale
continuous vector

Binning:

$$a_0 = \{1, 0\}$$

$$a_1 = \{0, 1\}$$

$$a_L = \{1, 2\}$$

0.9
1.1

Similar but
Binning

Cost of dimensionality

- $\{-1, 0, +1\}^7 \rightarrow 7$ degrees of freedom.

- not expressive enough

$$|A| = 3^7 = 2187$$

} neural network output

Policy as Probability

$$\pi(a|s) : \text{softmax } (W^T s) \rightarrow \text{sample from distibn}$$

policy parameters : W (shape \mathcal{W})

flat actions, rewards Theoretically nice rewards nice.

— Discrete
— Continuous

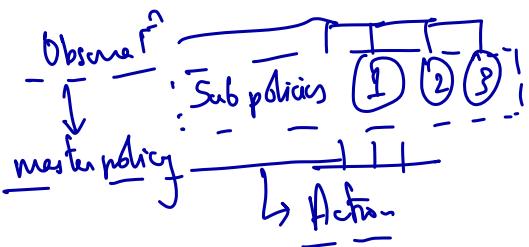
$0, 1, \dots, k-1$ $s \mid s \text{ states} : 0 \dots s-1 $ $0 = \text{pos}$ $1 = \text{act}$ $2 = \text{value}$	<u>targets</u> continuous \rightarrow tensors with >1 dimension
---	--

Current state — formal state
Current action $\pi(a|s) : \text{softmax } (W^T s)$

maximizing rewards in future

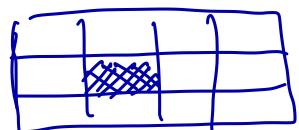
dictionary mapping

flat (s) $\xrightarrow{\text{values}}$ Action (a).



Action ($\uparrow \downarrow \leftarrow \rightarrow$)

policy $\{(0,0) :$



Indexing array
faster than indexing
dictionary.

Policy as Prob.

rand. no. $\rightarrow (\epsilon) < 0.1$
action from A

Els.

$a = \text{fixed policy}(s)$

value =

greedy Epsilon algo:

convert categorical to

action space

$\pi(a|s)$ probability
which action to take

$$\pi(a|s) = \text{softmax}(w^T s)$$

continuous state space: s^1 is a

D. dimensional vector

w

$(w^T s)$

gives rise to softmax function.

for each action space. $\pi(a|s)$

from env. $\text{argmax}(\pi(a|s))$

Merton predictor

The quick brown fox jumps over the lazy dog

Only predict what happens
before.

Markov Decision Process (MDP)

$$p(s_{t+1}, r_{t+1} | s_t, a_t)$$

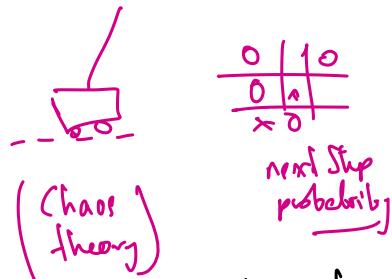
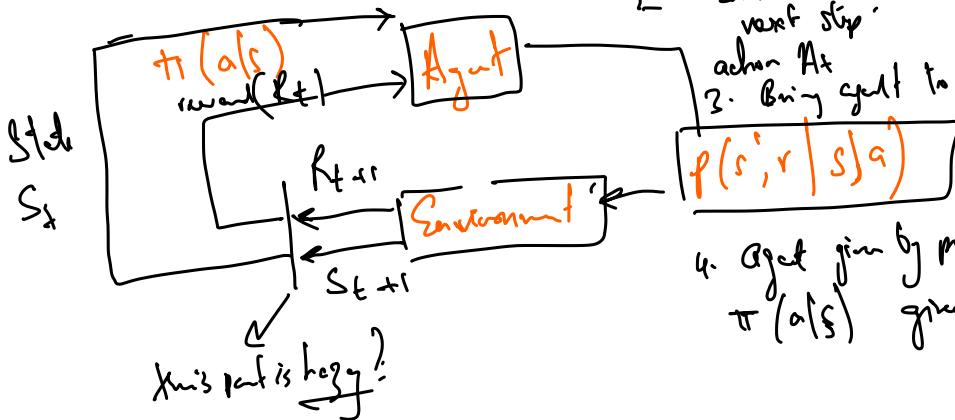
probability of getting to state $s_{t+1}(t+1)$ & get reward (r) at $t+1$ when current state is s_t at t & taking action a_t

$$\underbrace{p(s', r | s, a)}_{\text{Deterministic: } R(s, a, s')} \quad \begin{matrix} \text{state transition probability.} \\ \text{or} \\ \underbrace{R(s), R(s') \text{ reward}} \end{matrix}$$

$$\boxed{p(s' | s, a)}$$

Useful of soft transition probability:

- Built on Q-learning
- Law of physics: Deterministic? (No)
- Inverse law of physics: inaccu-
- Tautology: Environment dynamics



1. next state takes from environment updated from last step.
2. action a_t
3. Bring a_t to next step & next reward.
4. Agent given by prob. $\pi(a|s)$ given s .

- Environment state transition probability.
- well defined problem. → find solution.

maximizing reward.

Agent Goal: $\sum (\text{future rewards})$ agent: wants to max

Dilemma:

- next step
- last part rewards
- plan future max.
- actions to maximize the reward.
- Immediate gratification } long-term planning.
- miss T: V show

$$\sum_{t=1}^T (\text{future rewards}) = \text{Return}$$

$$G(t) = \sum_{z=1}^{\infty} R(t+z) = R(t+1) + R(t+2) + R(t+3) \dots R(T)$$

sum upto true terminal state.

Infinite Horizon MDP:

- Discounted for episodic | infinite
- ↳ Main theor.
- ↳ focus about rewards now | then future.

$$G(t) = \sum_{z=1}^{T-t} \gamma^{z-1} R(t+z) = R(t+1) + \gamma R(t+2) + \gamma^2 R(t+3) \dots + \gamma^{T-t-1} R(T).$$

$$G(t) = \underbrace{r(t+1)}_{\text{return at } (t+1)} + \gamma \underbrace{G(t+1)}_{\text{1. reward}} \quad \text{recursively,}$$

↓ ↓

1. reward 2. return

Value functions & Bellman curve-

Expressing problem as que?

Expected value = Mean = (μ)
Probability $\text{Prob.} \begin{pmatrix} 1 \\ 0 \end{pmatrix}$

$$(\mu = 2.6)$$

Expected value = 0.5
value $\left(\text{avg. of the values you can expect} \right)$

Discrete : $E(x) = \sum_{k=-\infty}^{+\infty} p(k)k$

Continuous : $E(x) = \int_{-\infty}^{\infty} x p(x) dx$

Expected value is weighted sum of all the possible values of rand. variable where weights = prob of values.

Expected value } Biased coin: $P_-(0.6)$ heads.
 $E(x) = 0.6 \times 1 + 0.4 \times 0 = \boxed{0.6}$ Expected value.

Environment dynamics
 \hookrightarrow probabilistic

playing
 ←
 maximise expected sum of future rewards

Final sum of future rewards

S

: Value function

$$V(s) = E(G_t | S_t = s)$$

Since

$$G_t(t) = R(t+1) + \gamma G(t+1)$$

repeating

(condition): we want state s at time t

$$V'(s) = E(R_{t+1} + \gamma V(s')) | S_t = s$$

Problem to solve?

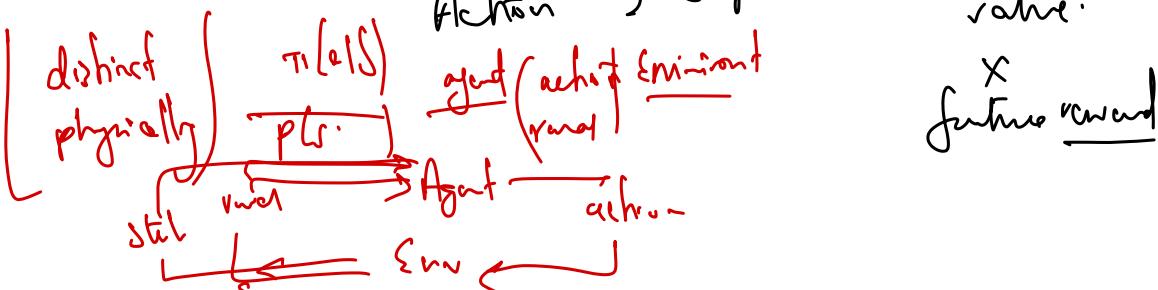
Consume building up problem.

$$V(s) = \sum_a \sum_{s'} \sum_r \pi(a|s) p(s', r | s, a) [r + \gamma V(s')]$$

Sum over all
possible actions (a),
possible next states (s'),
possible reward (r).
↓
probability
of performing
action a ,
given we
are in state s

probability of
landing in state
(s') & getting
reward r

Action \rightarrow Step \rightarrow reward \rightarrow reward
value.



problem & soln

multiple policies \rightarrow Good
 Bad.
 value function

$V_{\pi}(s) = \text{value function given policy} = \boxed{\text{prediction problem}}$

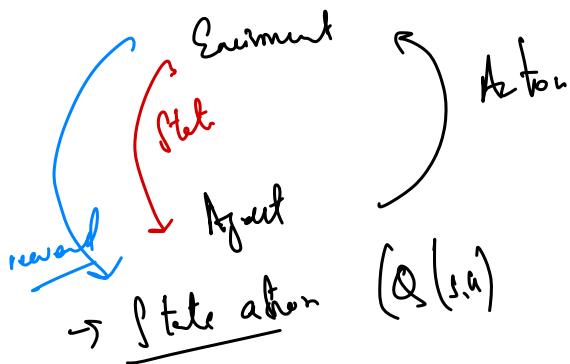
$\pi(a|s) \rightarrow \text{known (policy). function}$
 $p(s', r | s, a)$ fixed transition prob | Gridworld. environment.

learn meaning: Reinforcement learning problem.

1. Prediction problem.
 given π , find $\boxed{V(s)}$
2. Control problem.
 find π^* which gives $\max \boxed{V(s)}$

Greedy distribution

$$\Rightarrow a^* = \underset{a}{\operatorname{argmax}} Q^*(s, a)$$

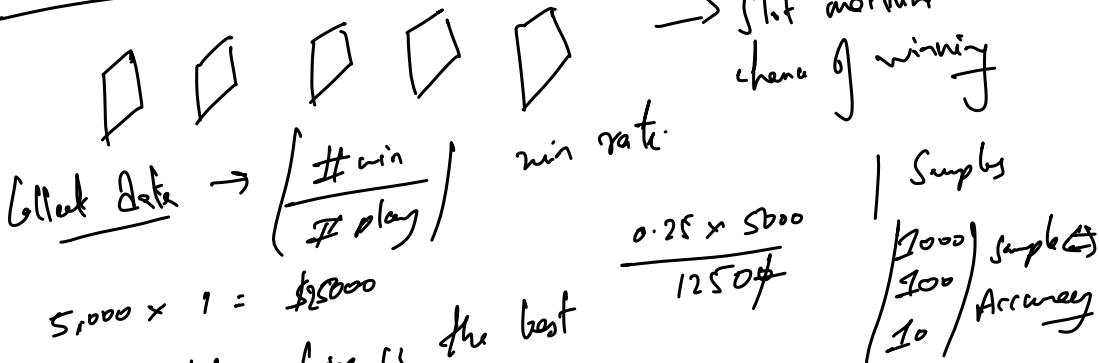


3 Action

$$\begin{cases} Q(s, a_1) \rightarrow \begin{cases} Q(s, a_1) \text{ initialized to } 0 \\ Q(s, a_1) = 0 \end{cases} \\ Q(s, a_2) \rightarrow \begin{cases} Q(s, a_2) \text{ initialized to } 0 \\ Q(s, a_2) = 0 \end{cases} \\ Q(s, a_3) \rightarrow \begin{cases} Q(s, a_3) \text{ initialized to } 1 \\ Q(s, a_3) = 1 \end{cases} \end{cases}$$

Explore - Exploit Dilemma

initialized in bad way



$$\frac{0.25 \times 5000}{12500}$$

Samples	5000
Accuracy	100%
Accuracy	10%

Epsilon Greedy problem

random no. $\alpha \varepsilon \leq 1$

$\sum = \text{random value}$
 Q -table.

1) $a = \underset{a}{\operatorname{argmax}}$ Standard policy

2) choose random no. $\leq \varepsilon$ if
if not met,
then random action

else:
 $a = \underset{a}{\operatorname{argmax}} (Q(s, a))$

Greedy
Epsilon

Q-learning algorithm

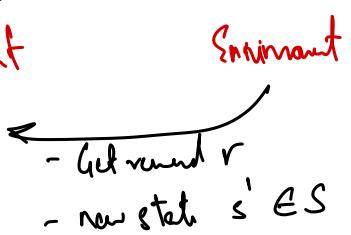
MDP - reward agent, action, environment
 - State $s \in S$
 - take action $a \in A$.

Prediction:

Control:

$$\text{Imitation: } (\sum \text{returns}) \text{ not known before episodes / rewards.}$$

$$(\sum \text{returns}) = \text{Reward.}$$



- long time for episodes.

- MC sub optimal learning

Temporal difference

$$\Rightarrow G(t) = R(t+1) + \gamma G(t+1)$$

return \rightarrow reward

learning much slower

diminishing returns

expected value problem.

Temporal diff $\xrightarrow{\text{approx}}$ M. Carlo method

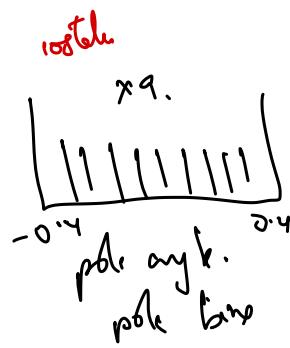
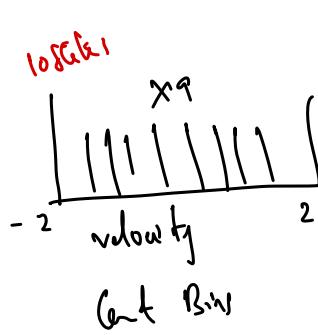
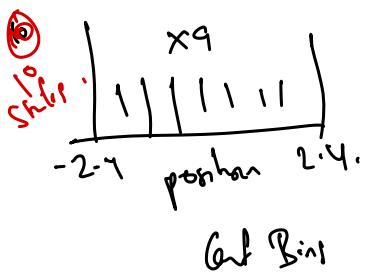
$$Q(s,a) = Q(s,a) + \alpha(g - Q(s,a))$$

$$J = (g - V(s))^2$$

free height &
sample (g)

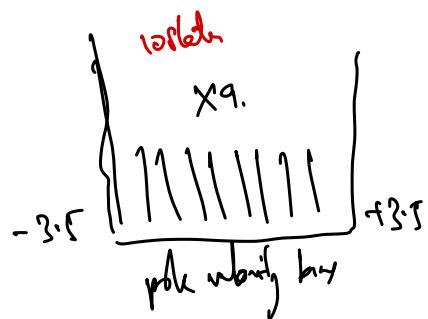
$V(c)$
predicted value for state 's'

Models from Binning.



Feature transform

- transform one obs \rightarrow at one time.
- not ideal, max probability & isn't unk addressed.



Model:

$$\begin{aligned} \text{env.} &= \text{con.} \\ \text{state} &= \text{10 bins} \\ \text{action} &= 2 \\ Q &= \text{min} \end{aligned}$$

10 bins for state variables (x_9)
num states, \times num-actions

$$\begin{aligned} (\text{actions} = 2) &\quad \text{lift (light)} \\ \text{Build state} &= \left[\begin{array}{cccc} \text{get pos.}, \text{cont vel}, \text{pole angle}, \text{pole vel} \\ 1 \quad 2 \quad 3 \quad 4 \end{array} \right] \end{aligned}$$

4 states
actions
space

```

50 class Model:
51     def __init__(self, env, feature_transformer):
52         self.env = env
53         self.feature_transformer = feature_transformer
54
55         num_states = 10*env.observation_space.shape[0]
56         num_actions = env.action_space.n
57         self.Q = np.random.uniform(low=-1, high=1, size=(num_states, num_actions))
58
59     def predict(self, s):
60         x = self.feature_transformer.transform(s)
61         return self.Q[x]
62
63     def update(self, s, a, G):
64         x = self.feature_transformer.transform(s)
65         self.Q[x, a] += 10e-3*(G - self.Q[x, a])
66
67     def sample_action(self, s, eps):
68         if np.random.random() < eps:
69             return self.env.action_space.sample()
70         else:
71             p = self.predict(s)
72             return np.argmax(p)
73

```

Next we have a sample action function.

Chassis → Chassis

Stra → Stra

Inly → Inblect

(best answer)

$G = R + \gamma \arg\max_{(a)} Q(s, a)$

Target return
Q-learning goal

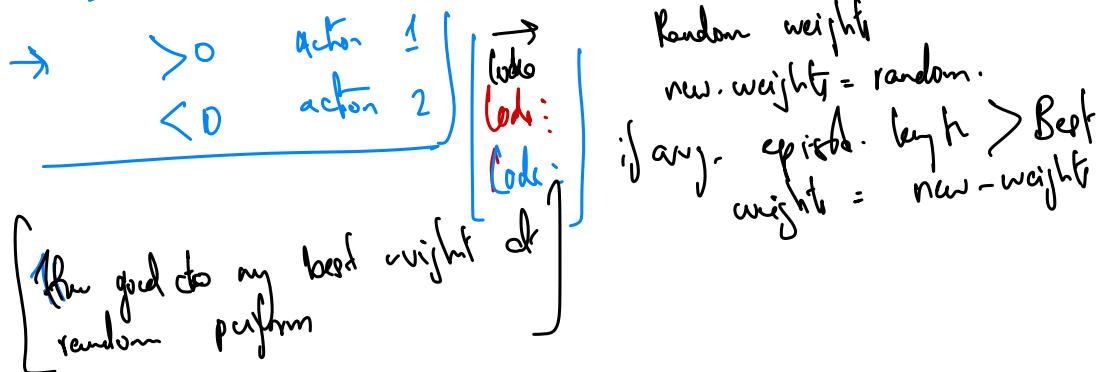
```

def play_one(model, eps, gamma):
    observation = env.reset() - initial State
    done = False
    totalreward = 0
    iters = 0
    while not done and iters < 10000:
        action = model.sample_action(observation, eps) - sample action = greedy option
        prev_observation = observation - initial state
        observation, reward, done, info = env.step(action)

        totalreward += reward
        if done and iters < 199: - (if falls down then
            reward = -300 - reward is -300)
            # update the model
            G = reward + gamma * np.max(model.predict(observation)) - G = reward + gamma * np.max(model.predict(observation))
            model.update(prev_observation, action, G) - model.update(prev_observation, action, G)
        iters += 1
    return totalreward

```

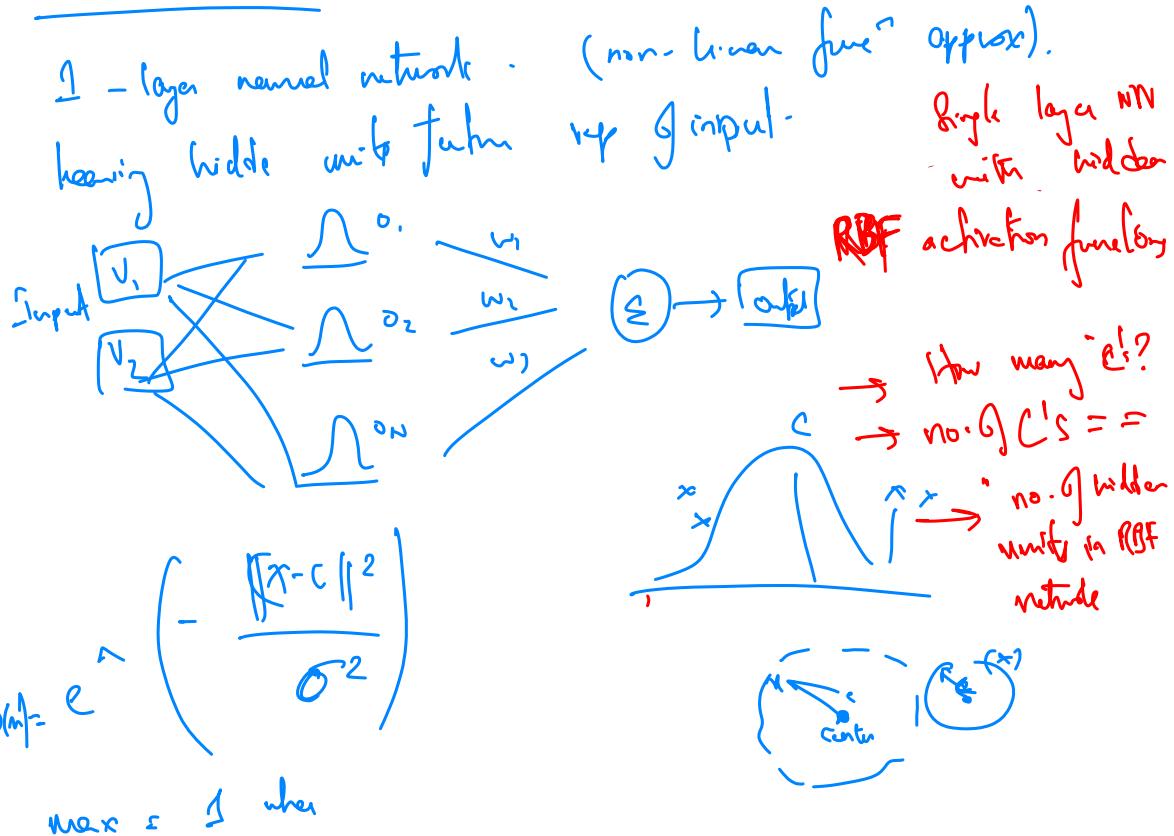
Random Search in parameter space for linear model.



```
random_search.py *
1 from __future__ import print_function, division
2 from builtins import range
3 # Note: you may need to update your version of future
4 # sudo pip install -U future
5
6 import gym
7 import numpy as np
8 import matplotlib.pyplot as plt
9
10
11 def get_action(s, w):
12     return 1 if s.dot(w) > 0 else 0
13
14
15 def play_one_episode(env, params):
16     observation = env.reset()
17     done = False
18     t = 0
19
20     while not done and t < 10000:
21         # env.render()
22         t += 1
23         action = get_action(observation, params)
24         observation, reward, done, info = env.step(action)
25         if done:
```

```
graph TD
    N1[If dot product > 0  
then val = 1  
else val = 0] --- W1[while not done and t < 10000]
    W1 --- N2[if done:]
```

R-B-F networks - Radial Basis func?



max = 1 when

no. of exemplars = no. of hidden units in Network.
 hyperparameters to be tuned.

sklearn Sci-kit learn. usage:

from sklearn.kernel_approx import RBFSampler.

sampler = RBFSampler()

sampler.fit(raw_data)

features = sampler.transform(raw_data).

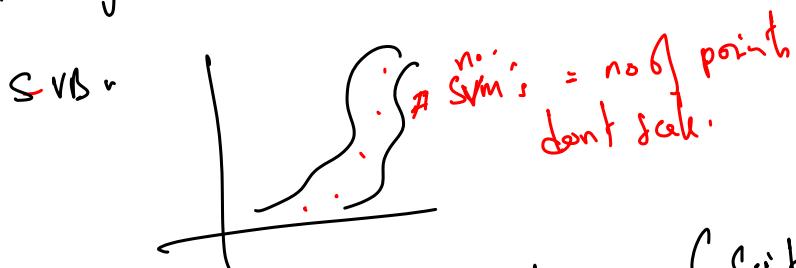
from sklearn.neighbors import KNeighborsClassifier.

SGD Regressor:

partial fit (X, y) -> one step of grad. desc.

predict (x)

partial fit \rightarrow dummy values / initializ.



exemplary := hyperparameters {Scikit learn implementation}

- * full partial fit SGD regressor target = 0
- * optimistic initial values method

instead of linear model.

use $x \xrightarrow{\text{transform}(s)}$

for every model, $Q(s)$ represents different action-

3 actions $Q(s) \vdash \text{Left}$

$Q(s) \vdash \text{Right}$

$Q(s) \vdash \text{Nothing}$

Policy gradient
Optimal policy π^*

Modeling
 \rightarrow (final π^{**})

1. find value given current policy
2. improve policy given the current value (gradient).

\rightarrow policy iteration

parameterized policy : score for action $a(j)$
 $\text{score}_j = f(a_j, s, \theta) + \phi(s)^T \theta_j$ if linear

softmax : $\pi(a_j | s) = \frac{\exp(\text{score}_j)}{\sum_j \exp(\text{score}_j)}$

initial state (s_0)
maximize the performance over the state s_0 = $V_{\pi}(s_0)$ = $\eta(\theta_p)$

parameterize (1) Policy (2) Value function

1. Policy : $\pi(a | s, \theta_p)$ = $f(s; \theta_p)$

policy response

action given state s | give policy θ_p

↓
performance
etc.

fn

2. Value function = $V_{\pi}(s) = f(s; \theta_v)$

value approximation

(θ_v) not included in the future optimal approximation?

θ_p = parameterize with 'p'

Policy gradient theorem:

$$\nabla \pi(\theta_P) = E \left[\sum_a Q_\pi(s, a) \nabla_{\theta_P} \pi(a | s, \theta_P) \right]$$

? Expected value.

$$\nabla \eta(\theta_P) = E \left[\sum_a Q_\pi(s, a) \nabla_{\theta_P} \pi(a | s, \theta_P) \right]$$

$$\begin{aligned} \left(\text{gradient of performance} \right) &= \left(\text{Sum of rewards} \right) \\ &\quad \left(\text{of action over states} \right) \end{aligned}$$

↓

$$\left(\text{gradient of policy parameters w.r.t p} \right)$$

$$\nabla \eta(\theta_P) = E \left[Q_\pi(s, a) \nabla_{\theta_P} \pi(a | s, \theta_P) \frac{1}{\pi(a | s, \theta_P)} \right]$$

$$\nabla \log f(x) = \frac{\nabla f(x)}{f(x)}$$

$$\nabla \eta(\theta_P) = E \left[Q_\pi(s, a) \nabla_{\theta_P} \log \pi(a | s, \theta_P) \right]$$

$$\nabla \eta(\theta_P) = E \left[G \nabla_{\theta_P} \log \pi(a | s, \theta_P) \right]$$

$$\begin{aligned} \left(\text{return we get from playing an episode} \right) &\rightarrow \text{parametrized pol.} \\ &\quad \text{pol.} \end{aligned}$$

- Steps
1. Play episode |
calc. the return.
 2. perform grad-
-ascent
- Much different than
gradient descent.
- Since goal is to max.
not minimize reward.

Batch gradient descent
By the time batch is over, have the rewards.

$$E(X) = \frac{1}{N} \sum_{n=1}^N x_n \rightarrow \begin{cases} \text{Expected value: mean } \bar{x} \\ \text{sample} \end{cases}$$

$$\nabla \eta(\theta_p) \approx \frac{1}{T} \sum_{t=1}^T G_t \nabla_{\theta_p} \log \pi(a_t | s_t, \theta_p)$$

$$E(X) \approx \frac{1}{N} \sum_{n=1}^N X_n$$

$$\nabla \eta(\theta_p) \approx \frac{1}{T} \sum_{t=1}^T G_t \nabla_{\theta_p} \log \pi(a_t | s_t, \theta_p)$$

$$\frac{1}{T} \sum_{t=1}^T G_t \nabla_{\theta_p} \log \pi(a_t | s_t, \theta_p) = \frac{1}{T} \sum_{t=1}^T \nabla_{\theta_p} G_t \log \pi(a_t | s_t, \theta_p)$$

constant

$$\frac{1}{T} \sum_{t=1}^T \nabla_{\theta_p} G_t \log \pi(a_t | s_t, \theta_p) = \frac{1}{T} \nabla_{\theta_p} \sum_{t=1}^T G_t \log \pi(a_t | s_t, \theta_p)$$

\rightarrow derivative of sum = sum of derivatives

minimize : $-\sum_{t=1}^T G_t \log \pi(a_t | s_t, \theta_p)$
of an episode because this involves the sum

minimize is opp of maximize
 $T = \text{length of episode}$
 $t = \text{time step of episode}$

Stochastic Grad. method.

$$\theta_{P,t+1} = \theta_{P,t} + \alpha G_t \frac{\nabla \pi(a_t | s_t)}{\pi(a_t | s_t)} \quad (= \theta_{P,t} + \alpha G_t \nabla \log \pi(a_t | s_t))$$

$$\Theta(a_t | s_t)$$

minimizing
reward
 \propto value of G_t

$\nabla \pi$ grad $\nabla(\pi)$

Direction of max reward.

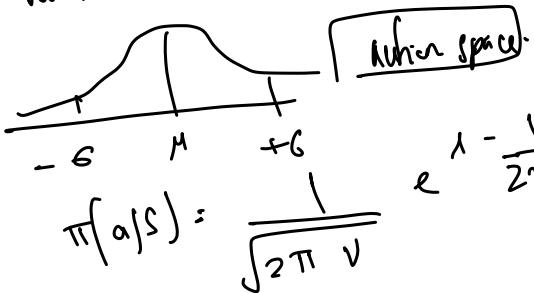
$(\pi) \rightarrow$: prob of choosing action (a_t)
grad $(\pi) \rightarrow$ direction of gradient increase
in (π)

$$\theta_{V,t+1} = \theta_{V,t} + \alpha(G_t - V_t) \nabla V_t$$

faster learning rate.
Baseline of $V(s)$

$V(s_t)$

Policy gradient method:
 Continuous action spaces.
 What distribution appropriate?
 $\pi(a|s)$ vs parameters (μ, σ) ($wx + b$)

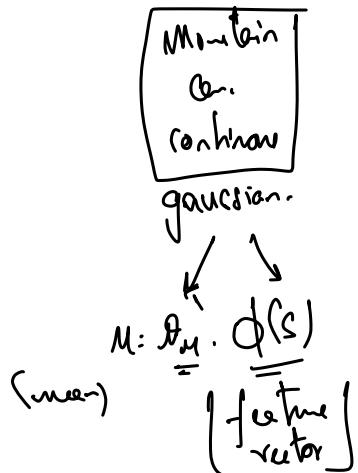


$$N(a; \mu, \sigma^2)$$



variance: $\exp(\Theta_v \cdot \phi(s))$
 [feature vector]

variance: $\text{softmax}(\Theta^\top \cdot \phi(s))$

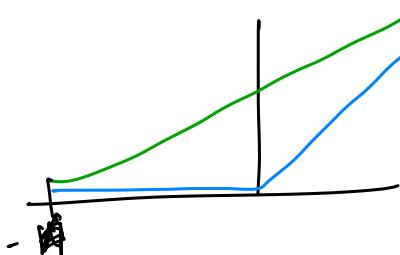


$$= (+ve)$$

make sure variance is always positive.

Better

grad-descent Θ : unbanded.
 Softmax.



Policy gradient

$$\theta := \theta + \alpha \sum (a_t - v(s_t)) \nabla \log \pi(a_t | s_t)$$

$$\pi(a | s) = \frac{1}{\sqrt{2\pi\nu}} e^{-\frac{1}{2\nu}(a-\mu)^2} = N(a; \mu, \nu)$$

$$\mu(s; \theta_\mu) = \theta_\mu^T \varphi(s) \rightarrow \text{parameterize } \mu$$

$$v(s; \theta_\nu) = \exp(\theta_\nu^T \varphi(s)) \rightarrow \text{the mean & variance}$$

$$\text{Alternative: } v(s; \theta_\nu) = \text{softplus}(\theta_\nu^T \varphi(s))$$

$$\text{softplus}(a) = \ln(1 + \exp(a))$$

$$\theta = \theta + \alpha \sum_{t=1}^T (G_t - V(s_t)) \nabla \log \pi(a_t | s_t)$$

$$\theta_{t+1} = \theta_t + \alpha(r_{t+1} + \gamma V(s_{t+1}) - V(s_t)) \nabla \log \pi(a_t | s_t)$$

] how to update
 → define model (1)
 → update model (2)

Q-updates:

→ update - \rightarrow degred, RMS prop.
 Any output distrib., parameterize \rightarrow (policy gradient)
 function approximat \rightarrow approx. / actor to π at step

Q-learning with func.
 → approximate
 → handle infinite spaces &
 action spaces

Policy gradient: Hill climbing



reward.

(-x action¹)

penalized for taking
illegal actions

HILL CLIMBING

→ Special type of random search.

2 parameters



initial params
(randomly)

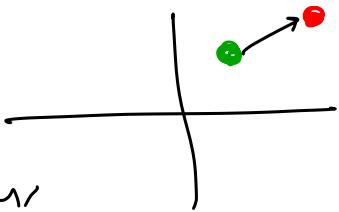
→ more episodes :

variance low

→ fewer episodes :

variance high

(now)



N.N.

Linear model, \rightarrow (RBF's)

parameters

Deep learning feelings

→ $\theta = \theta + \alpha(r + \gamma \max_a Q(s', a) - Q(s, a)) \nabla_{\theta} Q(s, a)$ → not converging to good sol?

→ automatic differentiation

Deep cf. networks

funcⁿ approx for $Q(s, a)$.

instead of $f(s, a)$, use $f(s)$.

→ State into input feature, each representing value for node.

Experience Replay

$(s_t, a_t) \leftarrow \text{Tuples}$

(s_t, a_t, r_t, s')

$(s_t, a_t, r_t, \text{reward}, \text{next step})$

→ Create env.

→ Main network

→ Target network

→ loop for predetermined # steps

$\begin{pmatrix} x \\ g \end{pmatrix}$ } shared context?

fully connected layer. hidden layer 'M'

Koef.

Actor critic method

making it parameterized.

Value - Value.

Actor \rightarrow Policy

~~Actor~~

$$G \sim V(s)$$

Return

Reward

- 2 neural networks to parameterize policy and value

$$\pi(a | s, \theta_p) = \text{NeuralNet}(input : s, weights : \theta_p) \rightarrow \text{parameterizing policy}$$

$$V(s, \theta_v) = \text{NeuralNet}(input : s, weights : \theta_v) \rightarrow \text{value.}$$

- For policy loss, work "backwards" from policy gradient, for value loss, use squared error

$$L_p = -(G - V(s)) \log \pi(a | s, \theta_p)$$

$$L_v = (G - V(s, \theta_v))^2$$

\rightarrow loss of policy

\rightarrow loss of value.

Pseudocode:

while not done :

$$a = \text{pi.sample}(s)$$

$$s', r, \text{done} = \text{env.step}(a)$$

$$G = r + \gamma * V(s')$$

$$L_p = -(G - V(s)) * \log(\text{pi}(s))$$

$$L_v = (G - V(s))^2$$

$$\theta_p: \theta_p - lr \times \frac{\partial L_p}{\partial \theta_p}$$

$$\theta_v: \theta_v - lr \times \frac{\partial L_v}{\partial \theta_v}$$

Entropy: (High loss)

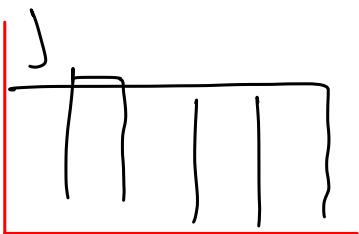
$$H$$

$$= - \sum_{k=1}^K \pi_k \log \pi_k$$

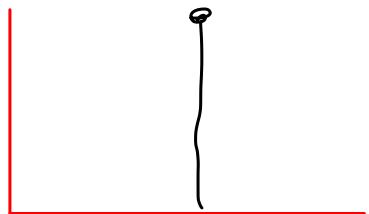
π_k → probability of event
 \propto log probability of event

new loss:

$$L'_p = L_p + C H$$



max. Entropy
= uniform
all same prob.



min. Entropy
= single event/
max prob.