

Deep Reinforcement Learning in R with rIR package¹

Xudong Sun

July 4, 2018

¹<https://github.com/smilesun/rIR>

Teach computer to play games with computer

rlR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rlR package

API

```
library(rlR)
env = makeGymEnv("Pong-v0")
env$overview()

##
## action cnt: 6
## state dim: 210, 160, 3
## discrete action
```

Teach computer to play games with computer

rlR

Xudong Sun

Reinforcement
Learning

Problem

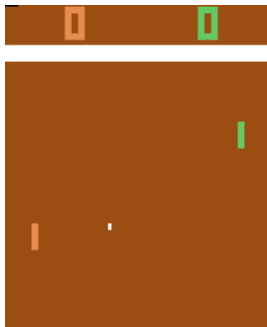
Concept

Theory

the rlR package

API

```
env$snapshot(steps = 25)
```



Available Environments in rIR

rIR

Xudong Sun

Reinforcement
Learning

Problem

Concept
Theory

the rIR package
API

```
library(rIR)
```

```
listGymEnvs()[1:30]
```

```
## [1] "CartPole-v0" "CartPole-v1"
## [3] "MountainCar-v0" "MountainCarContinuous-v0"
## [5] "Pendulum-v0" "Acrobot-v1"
## [7] "LunarLander-v2" "LunarLanderContinuous-v2"
## [9] "BipedalWalker-v2" "BipedalWalkerHardcore-v2"
## [11] "CarRacing-v0" "Blackjack-v0"
## [13] "KellyCoinflip-v0" "KellyCoinflipGeneralized-v0"
## [15] "FrozenLake-v0" "FrozenLake8x8-v0"
## [17] "CliffWalking-v0" "NChain-v0"
## [19] "Roulette-v0" "Taxi-v2"
## [21] "GuessingGame-v0" "HotterColder-v0"
## [23] "Reacher-v2" "Pusher-v2"
## [25] "Thrower-v2" "Striker-v2"
## [27] "InvertedPendulum-v2" "InvertedDoublePendulum-v2"
## [29] "HalfCheetah-v2" "Hopper-v2"
```

RL as a Functional Optimization process

rlR

Xudong Sun

Reinforcement
Learning

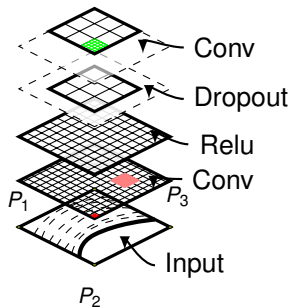
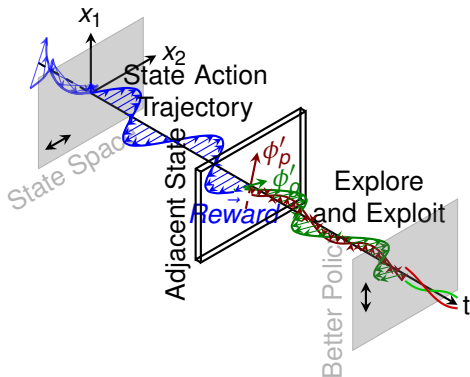
Problem

Concept

Theory

the rlR package

API



Classical Control Problem

rlR

Xudong Sun

Reinforcement
Learning

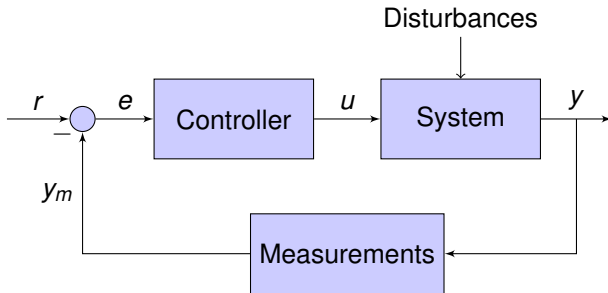
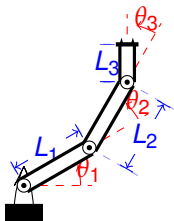
Problem

Concept

Theory

the rlR package

API



CliffWalker²

rIR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rIR package

API

##												
##	V	V	V	>	V	V	V	V	V	V	>	V
##	V	V	V	V	V	V	V	V	V	V	>	V
##	>	>	>	>	>	>	>	>	>	>	>	V
##	^	^	^	^	^	^	^	^	^	^	^	^
##												
##	>	>	>	>	>	>	>	>	>	>	>	V
##	>	>	>	>	>	>	>	>	>	>	>	V
##	>	>	>	>	>	>	>	>	>	>	>	V
##	^	^	^	^	^	^	^	^	^	^	^	^

Environment Agent Interaction through Policy

rlR

Xudong Sun

Reinforcement
Learning

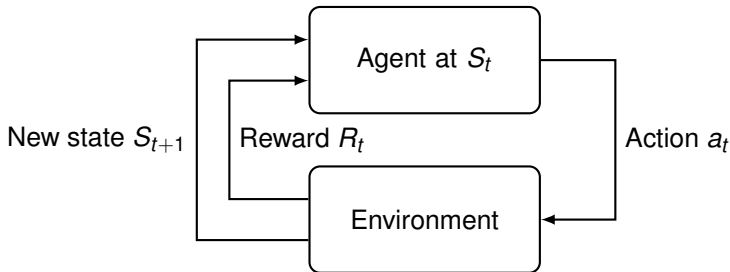
Problem

Concept

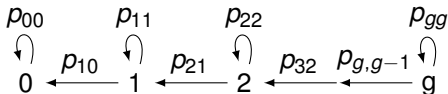
Theory

the rlR package

API



- Agent Environment Interaction
- State s , Action a , Transition(Environment MDP), Reward R
- Learning reaction Policy(e.g. Look Up Table) $\pi(a|s)$
- Returns: accumulated gain of reward $G_t = \sum_{i=0}^{\infty} \gamma^i R_{t+i}$



rIR sequence diagram

rIR

Xudong Sun

Reinforcement
Learning

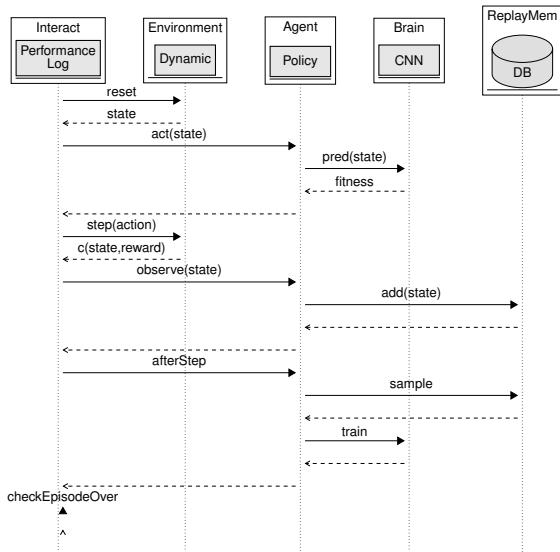
Problem

Concept

Theory

the rIR package

API



BackUp Diagram and Policy Environment Uncertainty

rlR

Xudong Sun

Reinforcement
Learning

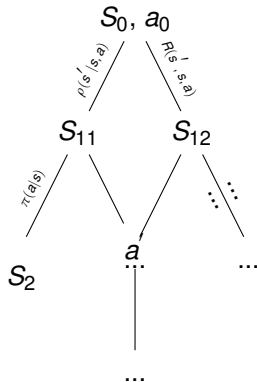
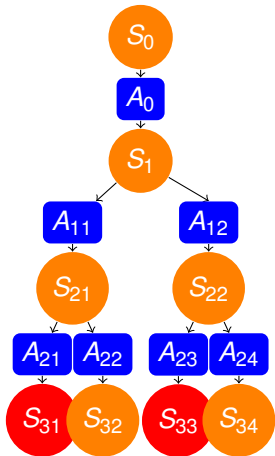
Problem

Concept

Theory

the rlR package

API



Long Term Consideration

rlR

Xudong Sun

Reinforcement
Learning

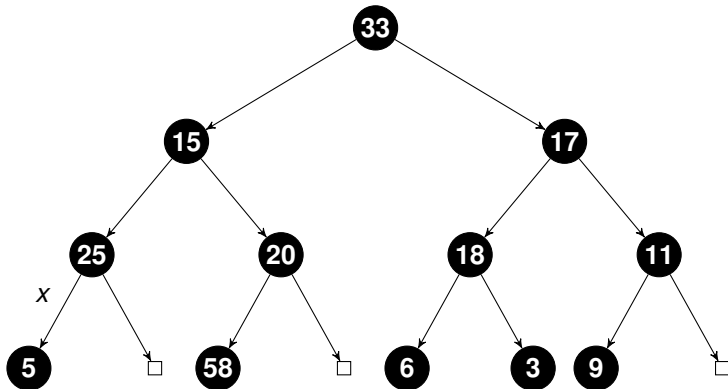
Problem

Concept

Theory

the rlR package

API



State Value Function and Bellman Equation

rlR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rlR package

API

$$\begin{aligned} V^\pi(s) &= E_\pi\left[\sum_{i=0}^{\infty} \gamma^i R_{t+i} \mid S_t = s\right] \\ &= E_\pi\left[R_t + \sum_{i=1}^{\infty} \gamma^i R_{t+i} \mid S_t = s\right] \\ &= E_\pi\left[R_t + \gamma \sum_{i=1}^{\infty} \gamma^{(i-1)} R_{t+i} \mid S_t = s\right] \\ &= E_\pi\left[R_t + \gamma \sum_{i'=0}^{\infty} \gamma^{i'} R_{t+1+i'} \mid S_t = s\right] \\ &= E_\pi\left[R_t + \gamma V^\pi(S_{t+1}) \mid S_t = s\right] \end{aligned}$$

State (Action) Value Function and Bellman Equation

rlR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rlR package

API

- $Q^\pi(s, a) = E_{\pi, \varepsilon}[R_t + \sum_{i=1}^{\infty} \gamma^i R_{t+i} | S_t = s, A_t = a]$
- $V(s_t) = \sum_a \pi(a|s_t) Q(s_t, a)$ since
 $\sum_a \pi(a|s_t) V(s_{t+1}) = V(s_{t+1})$
- $\pi^* = \operatorname{argmax}_{\pi} V^\pi(s) = \operatorname{argmax}_{\pi} E_{\pi}[R_t + \gamma V^\pi(S_{t+1}) | S_t = s], \forall s \in S$ (Otherwise replace to the better action at step t)
- $V^{\pi^*}(s) = E_{\pi^*}[R_t + \gamma V^{\pi^*}(s_{t+1}) | S_t = s]$ (optimal act each step)
- $Q^{\pi^*}(s, a) = E_{\pi^*, \varepsilon}[R_t + \sum_{i=1}^{\infty} \gamma^i R_{t+i} | S_t = s, A_t = a]$
- $V^{\pi^*}(s) = \max_a Q^{\pi^*}(s, a)$
- $Q^{\pi^*}(s, a) = E_{\varepsilon, \pi^*}[R_t] + \gamma \max_a \{Q^{\pi^*}(s_{t+1}, a)\}$

GPI

rlR

Xudong Sun

Reinforcement
Learning

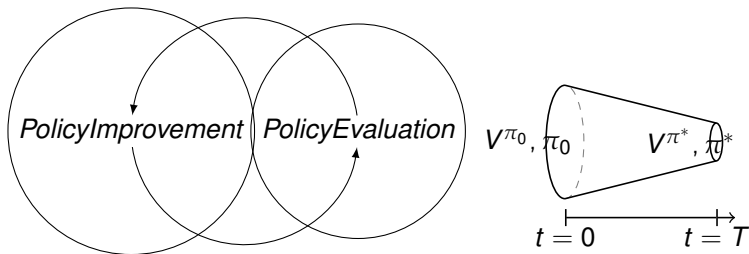
Problem

Concept

Theory

the rlR package

API



Dynamic Programming and Monte Carlo

rlR

Xudong Sun

Reinforcement
Learning

Problem

Concept

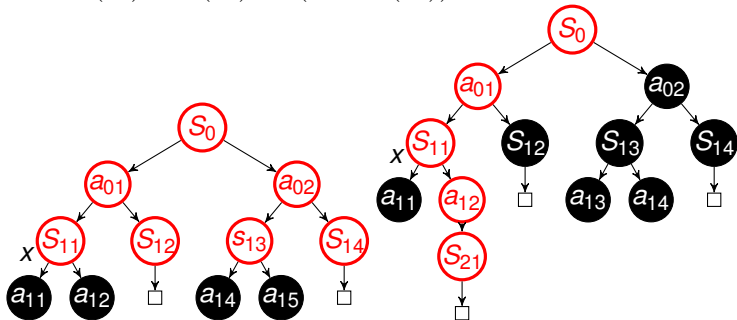
Theory

the rlR package

API

$$DP : V(S_t) = E_{\pi}[R_t + \alpha V(S_{t+1})]$$

$$MC : V(S_t) = V(S_t) + \alpha(R_t - V(S_t))$$



TD(λ) algorithm and (stochastic or deterministic) Policy Gradient

rlR

Xudong Sun

Reinforcement

Learning

Problem

Concept

Theory

the rlR package

API

$$Q^{\pi^*}(s, a) = E_{\varepsilon, \pi^*}[R_t] + \gamma \max_a \{Q^{\pi^*}(s_{t+1}, a)\}$$

$$\delta = Q^w(s, a) - ([R_t] + \gamma \max_a \{Q^w(s_{t+1}, a)\})$$

$$\nabla_{\theta} v_{\pi}(s) =$$

$$\sum_a [\nabla_{\theta} \pi_{\theta}(a|s) q_{\pi}(s, a) + \pi_{\theta}(a|s) [\sum_{s'} p(s'|s, a) + (\gamma \nabla_{\theta} v_{\pi(\theta)}(s'))]]$$

Implemented Algorithms in rIR

rIR

Xudong Sun

Reinforcement
Learning

Problem
Concept
Theory

the rIR package

API

```
env = makeGymEnv("CartPole-v0")
rIR::listAvailAgent(env)

## $AgentDQN
## [1] "Deep Q learning"
##
## $AgentFDQN
## [1] "Frozen Target Deep Q Learning"
##
## $AgentDDQN
## [1] "Double Deep QLearning"
##
## $AgentPG
## [1] "Policy Gradient Monte Carlo"
##
## $AgentPGBaseline
## [1] "Policy Gradient with Baseline"
##
## $AgentActorCritic
## [1] "Actor Critic Method"
```

Experiment reproducibility in rIR

rIR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rIR package

API

```
rIR::showDefaultConf()
```

```
##
## render FALSE
## log FALSE
## console FALSE
## agent.gamma 0.99
## agent.flag.reset.net TRUE
## agent.lr.decay 0.999000499833375
## agent.lr 0.001
## agent.store.model FALSE
## agent.clip.td FALSE
## policy.maxEpsilon 0.01
## policy.minEpsilon 0.01
## policy.decay 1
## policy.decay.type geometric1
## policy.aneal.steps 1e+06
## policy.softmax.magnify 1
```

rIR on Inverted Pendulum

rIR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rIR package

API

```
env = makeGymEnv("CartPole-v0")  
env$snapshot()
```



rlR on Inverted Pendulum

rlR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rlR package

API

```
library(rlR)
env = makeGymEnv("CartPole-v0")
conf = getDefaultConf("AgentDQN")
agent = makeAgent("AgentDQN", env, conf)
perf = agent$learn(200)
perf$plot()
```

Deep Q Learning on CartPole-v0

rIR

Xudong Sun

Reinforcement
Learning

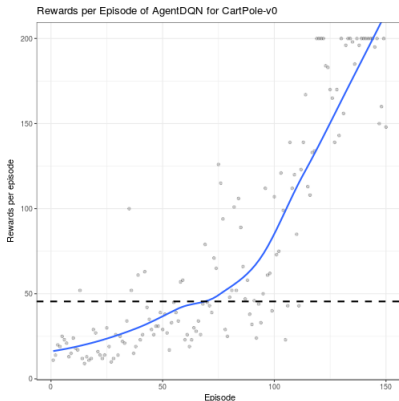
Problem

Concept

Theory

the rIR package

API



Custom Neural Network

rlR

Xudong Sun

Reinforcement

Learning

Problem

Concept

Theory

the rlR package

API

```
env = makeGymEnv("MountainCar-v0",
  act_cheat = c(0, 2))
conf = getDefaultConf("AgentDQN")
conf$set(console = TRUE, render = TRUE,
  policy.maxEpsilon = 0.15, policy.minEpsilon = 0,
  policy.decay = 1.0 / 1.01, replay.batchsize = 10,
  replay.epochs = 4, agent.lr.decay = 1,
  agent.gamma = 0.95)
agent = makeAgent("AgentDQN", env, conf)
env$overview()

##
## action cnt: 2
## state dim: 2
## discrete action
```

Custom Neural Network

rIR

Xudong Sun

Reinforcement
Learning

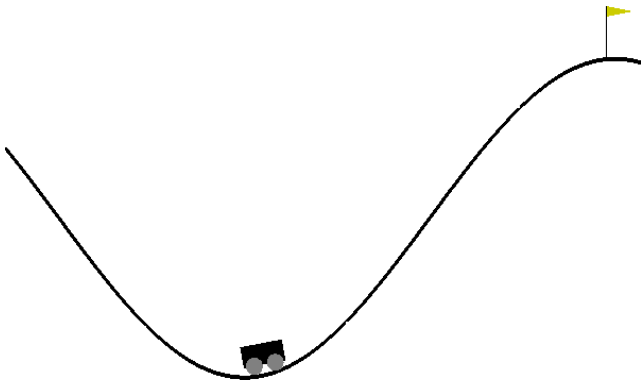
Problem

Concept

Theory

the rIR package

API



Custom Neural Network

rlR

Xudong Sun

Reinforcement
Learning

Problem
Concept
Theory

the rlR package
API

```
mfun = function(state_dim, act_cnt) {  
  requireNamespace("keras")  
  model = keras::keras_model_sequential()  
  model %>%  
    layer_dense(units = 10, activation = "relu"  
      input_shape = c(state_dim)) %>%  
    layer_dropout(rate = 0.25) %>%  
    layer_dense(units = act_cnt,  
      activation = "linear")  
  model$compile(loss = "mse",  
    optimizer = optimizer_rmsprop(lr = 0.001))  
  model  
}  
agent$customizeBrain(value_fun = mfun)  
agent$learn(500L)
```


Custom Neural Network

rlR

Xudong Sun

Reinforcement
Learning

Problem
Concept
Theory

the rlR package
API

```
conf = getDefaultConf("AgentFDQN")
conf$set(replay.batchsize = 32, replay.freq = 4L,
  console = TRUE,
  agent.lr.decay = 1, agent.lr = 0.00025,
  replay.memname = "UniformStack", render = FALSE,
  policy.decay = exp(-2.2 / 1e6),
  policy.minEpsilon = 0.1,
  agent.start.learn = 5e4, replay.mem.size = 1e6,
  log = FALSE,
  agent.update.target.freq = 10000L, agent.clip.td
  policy.decay.type = "linear")
env = makeGymEnv("KungFuMaster-v0", observ_stack_le
agent = makeAgent("AgentFDQN", env, conf)
```

Custom Neural Network

rIR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rIR package

API



Custom Neural Network

rlR

Xudong Sun

Reinforcement
Learning

Problem
Concept
Theory

the rlR package
API

```
pong_fun = function (state_dim, act_cnt) {  
  model <- keras_model_sequential()  
  model%>%  
    layer_conv_2d(filter = 32, kernel_size = c(8,8),  
      padding = "same", input_shape = state_dim) %>%  
    layer_conv_2d(filter = 64, kernel_size = c(4,4),  
    layer_activation("relu") %>%  
    layer_conv_2d(filter = 64, kernel_size = c(3,3),  
      strides = c(1,1), padding = "same") %>%  
    layer_activation("relu") %>%  
    layer_flatten() %>% layer_dense(512) %>%  
    layer_activation("relu") %>% layer_dense(act_cnt) %>%  
    layer_activation("linear")  
  opt = optimizer_rmsprop(lr = 0.00025)  
  model %>% compile(loss = "mse", optimizer = opt,  
  return(model)  
}
```

Custom Neural Network

rIR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rIR package

API

```
agent$customizeBrain(value_fun = pong_fun)  
agent$learn(5000)
```

Custom Neural Network

rlR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rlR package

API

```
rlR::listAvailConf()
```

```
##      [1] "render"                                "log"
##      [3] "console"                              "agent.gamma"
##      [5] "agent.flag.reset.net"                 "agent.lr.decay"
##      [7] "agent.lr"                             "agent.store.mod"
##      [9] "agent.update.target.freq"             "agent.start.lear"
##     [11] "agent.clip.td"                        "policy.maxEpsilon"
##     [13] "policy.minEpsilon"                   "policy.decay"
##     [15] "policy.decay.type"                   "policy.aneal.st"
##     [17] "policy.softmax.magnify"               "replay.batchsiz"
##     [19] "replay.memname"                      "replay.mem.size"
##     [21] "replay.epochs"                       "replay.freq"
```

Thanks for your attention!

rlR

Xudong Sun

Reinforcement
Learning

Problem

Concept

Theory

the rlR package

API

- URL: <https://github.com/smilesun/rlR>
- BugReports: <https://github.com/smilesun/rlR/issues>
- `devtools::install_github("smilesun/rlR", dependencies=TRUE)`
- Install, try it out, and have fun!