

Solving Bomberman

WITH Q-TABLE AND NEURAL NETWORKS

BEHROOZ MONTAZERAN

JANNIS HEISING

TOMÁŠ SLÁMA

20. 12. 2023

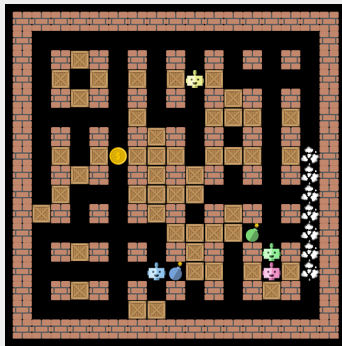


Introduction

PROBLEM STATEMENT

PROBLEM STATEMENT

- **6 actions** – 4 movement directions + wait + place bomb
- destructible crates, indestructible walls
- **+1** for coin, **+5** for killing opponent



METHODS USED

METHODS USED

Q-learning – optimizing action-value Q-function:

- **Q-table:** store each input/output in a table
- **Neural Network:** learn function via a neural network

METHODS USED

Q-learning – optimizing action-value Q-function:

- **Q-table:** store each input/output in a table
- **Neural Network:** learn function via a neural network

TD update – Use best possible predicted value in next step:

$$Y_t \leftarrow R_t + \gamma \cdot \max_a Q(s_{t+1}, a),$$
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \mu \cdot (Y_t - Q(s_t, a_t)).$$

METHODS USED

Q-learning – optimizing action-value Q-function:

- **Q-table:** store each input/output in a table
- **Neural Network:** learn function via a neural network

TD update – Use best possible predicted value in next step:

$$Y_t \leftarrow R_t + \gamma \cdot \max_a Q(s_{t+1}, a),$$
$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \mu \cdot (Y_t - Q(s_t, a_t)).$$

Stability can be improved with policy/target model:

- **target:** trained on the results of actions
- **policy:** updated slowly / replaced periodically

Evaluation with `elo.py`:

- comparing relative strengths of agents using the elo system
- agents start with 1 000 elo, updated after win/loss/draw

Evaluation with `elo.py`:

- comparing relative strengths of agents using the elo system
- agents start with 1 000 elo, updated after win/loss/draw

Training with `train.py`:

- subcommands for training traits (coin picking, agent hunting)
- periodically calculates strength with `elo.py`
- useful training flags: `--infinite <n>`, `--continue`

TRAINING COMMAND EXAMPLE

TRAINING COMMAND EXAMPLE

```
TASKS = {  
    ...  
    "complete": [  
        # first learn to pick up coins  
        (["--scenario", "coin-heaven", "--n-rounds", "100"], False),  
        # then learn to hunt/evade the rule-based agent  
        (["rule_based_agent", "--scenario", "empty", "--n-rounds", "1000"], False),  
        # then learn to work with crates  
        (["--scenario", "classic", "--n-rounds", "1000"], False),  
        # then learn to work with crates + hard agent  
        (["rule_based_agent", "--scenario", "classic", "--n-rounds", "200"], True),  
        # then learn to work with crates + really hard agent  
        (["binary_agent_v3", "--scenario", "classic", "--n-rounds", "200"], True),  
        # finally play against itself (None gets substituted)  
        ([None, "--scenario", "classic", "--n-rounds", "200"], True),  
    ],  
    ...  
}
```

Neural Networks

Network structure: fully connected layers with ReLU

- uses TD for Q-value updating

Network structure: fully connected layers with ReLU

- uses TD for Q-value updating

Main idea – carefully crafted binary state vector:

- 1... 5: direction to the nearest **coin**
- 6...10: direction to the nearest **crate**
- 11...15: direction to where a bomb will **kill opponent**
- 16...20: direction to **safety** (if in danger of dying)
- 21: **can place a bomb** and there is a way to escape it

Network structure: fully connected layers with ReLU

- uses TD for Q-value updating

Main idea – carefully crafted binary state vector:

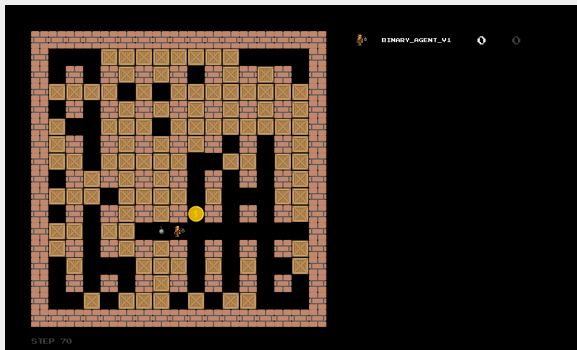
- 1... 5: direction to the nearest **coin**
- 6...10: direction to the nearest **crate**
- 11...15: direction to where a bomb will **kill opponent**
- 16...20: direction to **safety** (if in danger of dying)
- 21: **can place a bomb** and there is a way to escape it

Calculated by **searching state space** (other players stand still)

- ⇒ finds path even in ongoing explosions
- ⇒ shortest on-board path \neq shortest path stepwise

Q-VECTOR EXAMPLE (1)

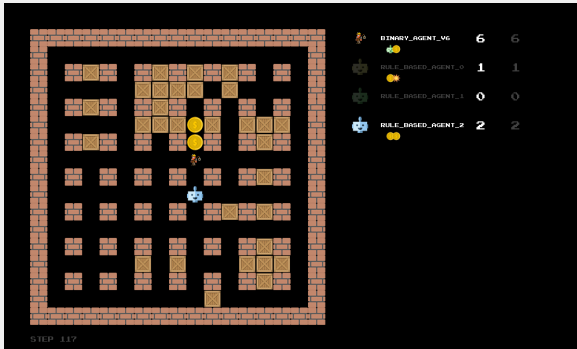
Q-VECTOR EXAMPLE (1)



| | | | | |
|---|---|---|---|---|
| 0 | 1 | 0 | 0 | 0 |
| ↓ | | | | |
| Nearest coin to the right | | | | |
| 0 | 1 | 0 | 0 | 0 |
| ↓ | | | | |
| Nearest crate to the right (by state space search) | | | | |
| 0 | 0 | 0 | 0 | 0 |
| ↓ | | | | |
| No players reachable | | | | |
| 0 | 1 | 0 | 0 | 0 |
| ↓ | | | | |
| Safety to the right | | | | |
| 0 | | | | |
| ↓ | | | | |
| Can't place a bomb | | | | |

Q-VECTOR EXAMPLE (2)

Q-VECTOR EXAMPLE (2)



1 0 0 0 0 | 1 0 0 0 0 | 0 0 0 0 1 | 0 0 0 0 0 | 1

↓ ↓ ↓ ↓ ↓

Nearest coin up Nearest crate up Place bomb to endanger player Safe Can place a bomb

TRAINING PARAMETERS

TRAINING PARAMETERS

Hyperparameters

```
BATCH_SIZE = 256
MEMORY_SIZE = 1000

# chances of picking random action
# happens per-round, not per-epoch!
EPS_START = 0.99
EPS_END = 0.05
EPS_DECAY = 10

# (soft) update rate of target network
TAU = 1e-3

# LR of the optimizer
LR = 1e-4

# discount for future states
GAMMA = 0.99

OPTIMIZER = optim.Adam
LAYER_SIZES = [1024, 1024]
```

Event Rewards

```
# hunt coins
MOVED_TOWARD_COIN: 50
DID_NOT_MOVE_TOWARD_COIN: -100

# hunt people
MOVED_TOWARD_PLAYER: 25
DID_NOT_MOVE_TOWARD_PLAYER: -10

# blow up crates
MOVED_TOWARD_CRATE: 1

# basic stuff
e.INVALID_ACTION: -100
DID_NOT_MOVE_TOWARD_SAFETY: -500

# be active!
USELESS_WAIT: -100

# meaningful bombs
PLACED_USEFUL_BOMB: 60
PLACED_SUPER_USEFUL_BOMB: 150
DID_NOT_PLACE_USEFUL_BOMB: -500
```

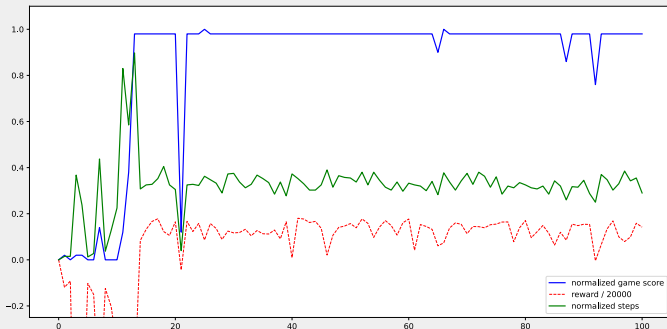

TRAINING LOOP

TRAINING LOOP

1. `train.py` with the "complete" task and `--infinite` 2
2. pick the best-performing model based on calculated elo
3. decrease `EPS_START`, `EPS_END`, `TAU` and `LR` hyperparameters
4. go to step 1 (with `--continue` to not overwrite the model)

TRAINING LOOP

1. `train.py` with the "complete" task and `--infinite`
2. pick the best-performing model based on calculated elo
3. decrease `EPS_START`, `EPS_END`, `TAU` and `LR` hyperparameters
4. go to step 1 (with `--continue` to not overwrite the model)

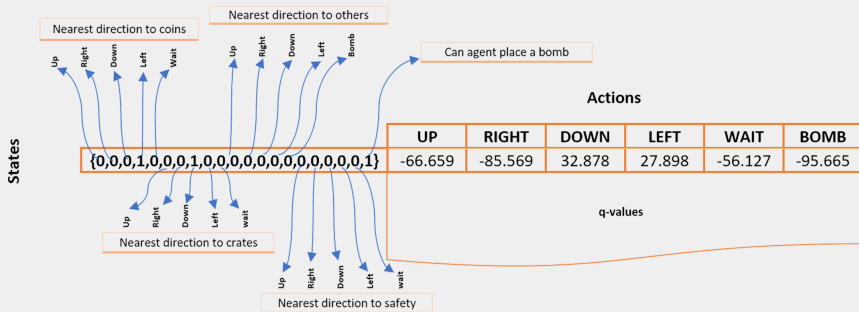


Q-table

STORING Q-TABLE

STORING Q-TABLE

- Theoretically $207 \times 2^{21} \times 6 = 2.60466278e9$ possible states
- 207 free tiles for the agent to go
- 2^{21} possible actions-state pairs
- 6 actions



STORING Q-TABLE

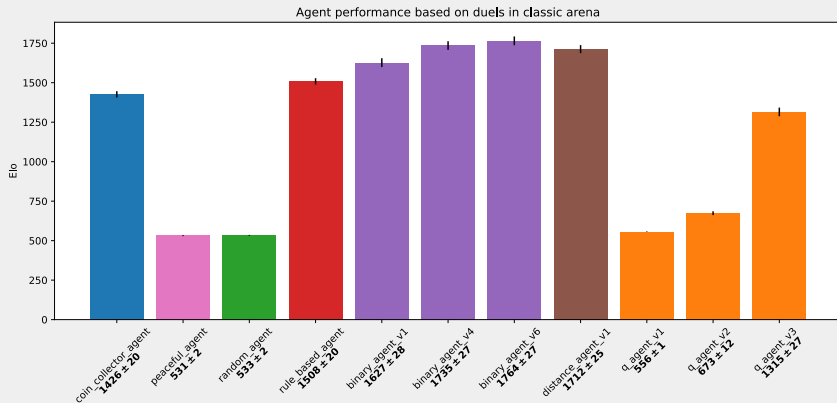
In reality, agent encounters much fewer states

| Scenario | Opponent(s) | Number of Rounds | Maximum new encountered states | Encountered states | Run time | Total run time |
|---------------------|--|------------------|--------------------------------|--------------------|-----------|----------------|
| Coin-heaven | - | 5,000 | ~100 | 100 | ~00:03:40 | 00:03:40 |
| Sparse-crate | - | 30,000 | ~560 | 660 | ~00:20:30 | 00:24:10 |
| Sparse-crate | Peaceful-agent | 40,000 | ~2290 | 2950 | ~00:47:35 | 01:11:45 |
| classic | Peaceful-agent | 40,000 | 0 | 2950 | ~00:27:28 | 01:39:13 |
| classic | Peaceful-agent & rule-based-agent | 40,000 | ~2530 | 5480 | ~01:15:17 | 02:54:30 |
| classic | rule-based-agent & rule-based-agent & rule-based-agent | 100,000 | ~5315 | 10795 | ~05:44:10 | 08:38:40 |
| Doubled | Above in order | 255,000 | ~16 | 10811 | ~04:20:40 | 12:59:20 |

Conclusion

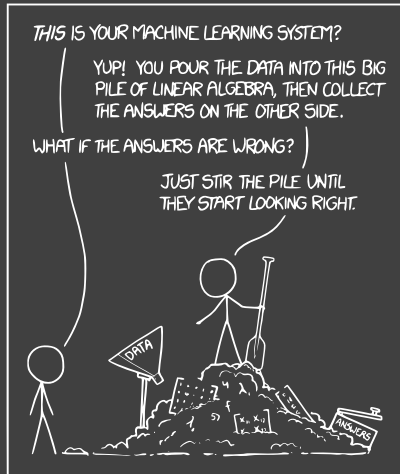
RESULTS

RESULTS



THANKS FOR YOUR ATTENTION!

SPECIAL THANKS TO ULLRICH KÖTHE AND TUTORS



REFERENCES