

Thoughts on how to improve the Q-Learning algorithm:

Maybe also test models on a crate-free, coin-poor scenario (Task 1.5?)

Maybe test with reenabling bomb and wait to evaluate whether Q-learning can learn those two as well?

Things to maybe try out

- MC value estimation

$$Y_{\pi^*} = \sum_{t'=t+1}^{\infty} \gamma^{t'-t-1} R_{\text{max}}^{\pi^*} \tau_{t'}$$

- "offline" update:
 - either update π every K games, while accumulating a K -game TS
 - or update Q every game with a TS going back K games

- look into other learning approaches for Q-learning regression

- linear value approx

-

Feature Design Task 2:

Goal: Learn to collect all coins by blowing up crates safely

- maximise coin-collecting speed
- minimise deadly accidents

Approach: Develop features that

- make it as easy as possible for the agent to learn good behaviour.
- make it able to also handle complex/edge cases well.

What is required for good behaviour?

- If a bomb is placed, always move to safety in time
- Idea: place bombs where the expected number of destroyed crates is highest
or such that the crate destruction speed is maximised.
- Idea: Have a high priority on collecting accessible coins

more precisely:

is what is in proximity of agent more away to safety zone

dangerous zones
safe
How to show info about to agent
in time and space?



Ray tasks:

1. collect coins
2. when no coins reachable, place bomb where many crates are destroyed.
3. when bomb placed, move to safety (and stay away from explosion while it lasts)

Idea: Learn

Provide algorithm

if not in_safety:
not_in_safety = within bomb explosion radius
move to safety

else if non-reachable_coins:
non-reachable_coins = no visible coins or coins out of reach
(path blocked by crates)
place bomb at good-spot

else:
good_spot = Det. algorithmically
collect coins as in Task 1

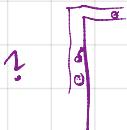
at what spot a bomb placed results in the highest coin collecting speed

[factor in: time to get to good_spot, time to place bomb, run to safety, time to return to new paths
similarity in routes,

multi-bomb strategies?

remaining_minedcoins / remaining_crates

expected_coins = total_time_to_reach * coin_density * crates_destroyed



Advantages:

- feature architecture with "recommended directions" can be used again, with additional "feature mode" and action "bomb" and "wait" enabled → do these actions require 3 nodes, 6 actions \Rightarrow matrix $\in \mathbb{R}^{6 \times 6}$ new features?
- modes maybe individually trainable and possibly individually analysing, debug and incentivise
- provided learning of task 1 has become good and "recommender algorithms" are good, this may actually be very easy to learn.

Disadvantages

- strategy heavily dependent on quality of algorithms
- may prevent agent from employing simultaneous or hybrid strategies *
- (quite like 2013 winner approach)
- learning not really necessary here

*possible advanced strategies

- Ignore coins that can be more easily reached by opponents
- Place one bomb. Then move in a direction to safety where you (after the first one is gone) can place another bomb. Move to safety exploring the hole of the first bomb while the second one explodes.
- When going after an accessible coin, place bombs along the way to find nominees.
- When a coin is inaccessible behind a few crates, place a bomb specifically to get to that coin.
- Further strategies likely exist.
- Hybrid strategies when opponents are present may become even more important.

Other ideas:

- Have the features of multiple modes separate from each other
 - + potentially allows for learned prioritizing between modes
- how to make/induce/inform agent when to prioritize one set of features over another?
 - big Q matrix (maybe unnecessary)
 - seems not to increase likelihood of advanced strategies
- Learn when to switch between modes
 - + let the low level stuff be handled by algorithms
 - + focus 'thinking' / learning on deciding between high level strategies