

## 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 if 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

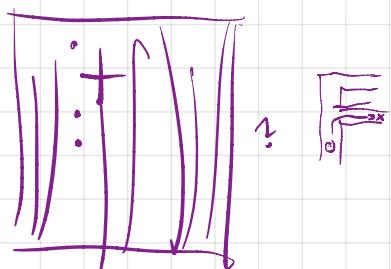
more precisely:

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?



→ self. (---)



Idea: Learn

Pseudo-algorithm:

if not in\_safety:  
    not\_in\_safety = within\_bomb\_expllosion\_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:  
    Idea: 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 crates,  
multi-bomb strategies?  
expected coin speed =  $\frac{\text{remaining\_reachablecoins} / \text{remaining\_crates}}{\text{total\_time}}$   
 $\times \text{coin\_density} \times \text{crates\_destroyed}$   
 $\rightarrow \text{time to reach good spot + bomb cooldown}$   
 $\rightarrow \text{expected\_time\_to\_reach\_coins}$

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  $\Rightarrow$  a matrix  $\in \mathbb{R}^{6 \times 6}$  new features?
- modes maybe individually trainable and possibly individually analysable, 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 approx. algorithms
- may prevent agent from employing simultaneous or hybrid strategies \*
- (quite like 2019 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 / incentives / inform agent when to prioritize one set of features over another?
- big Q matrix (maybe unnecessary)
  - seems not to increase likelihood of advanced strategies  
↳ why not?

- Learn when to switch between modes
    - + let the low level stuff be handled by algorithms
    - + focus thinking / learning on deciding between high level strategies
- ? How to do that?

Idea for info categories / naming convention

Idea: Provide medium level info such that  
agent knows just the things relevant to high level decisions.

low level: info directly from game-state (or almost directly)  
medium level: info from algo such as look-for-targets  
high level: info about strategies and plans

- e.g.
- where never to go : into walls, crates, bomb, agents
  - where to go for coins (direction?, proximity?, what about multiple ones?)
  - where to go to evade danger (direction?, threat level?)
  - where the next best place to place bomb and destroy crate (direction?, gain?)

How to decide between strategies?

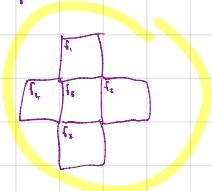
↳ maybe learn from long-term experience

$$k_c = s(4, 3) = 15$$

$$k_m = s(4, 5) \cdot s(1, 3) = 210 \text{ at most}$$

$$k_h$$

Ausprobieren:



- $\Rightarrow$
- |           |   |   |
|-----------|---|---|
| $f_{i-n}$ | $0: \text{blockiert}$ (nicht passabel oder tödlich gefährlich)<br>$1: \text{frei aber nicht besonders interessant}$<br>$2: \text{richtung für nächster coin}$<br>$3: \text{richtung für beste bomb-stelle (good-spot)}$<br>$4: \text{sofort bomb oder auch bomb}$ | } |
|           | $\rightarrow \text{nicht dort hin}$<br>$\} \text{Lerne entscheidung}$<br>$\} \text{zwischen diesen alternativen}$   |   |

$$s^+(5, 5)$$

( $i$ )

- $\Rightarrow$
- |  |   |
|--|---|
| $0: \text{auf bomben bew. ist man leben}$<br>$1: \text{uninteressant}$       | $\rightarrow \text{nicht wait}$   |
| $2: \text{auf coin (nicht good-spot)}$<br>$3: \text{auf good-spot (analog)}$ | $\rightarrow \text{wait}$ (was sehr selten, unwichtig. In ande policies fast gar nix gute Ergebnisse liefern.)<br>$\rightarrow \text{bomb}$ |
- $\Rightarrow$

Herausforderungen:

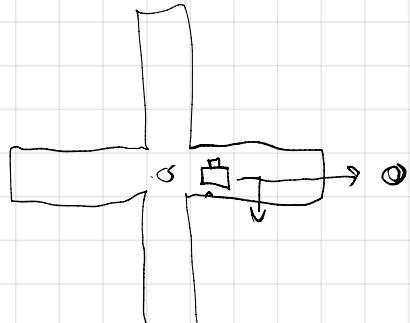
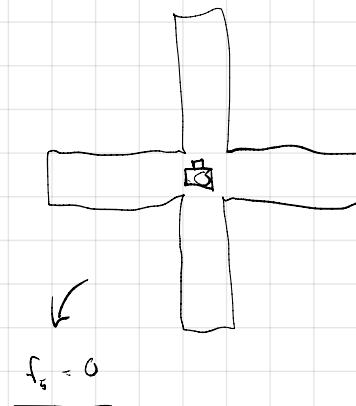
wie priorisieren zwischen nearest\_coin und good\_spot

$\rightarrow$  kann kann, aber mit gegebener info was good was besser ist (braucht aber langfristiges denken, sprich & hoch, bzw. n-step qmb)

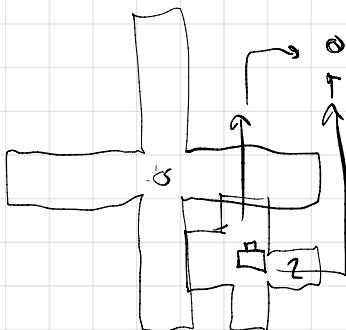
Benötigt folgende Algorithmen:

- $\rightarrow$  • direction\_nearest\_coins()
- $\rightarrow$  • direction\_good\_spot()
- $\rightarrow$  • direction\_lethal\_danger()

Outputs:  
(alle)  
[self-position,  
neighbor-positions]



$$f_{i-n} = 0$$



Basic feature architecture:

state-to-features()      Input: game-state      Output:  $[f_1, f_2, f_3, f_4, f_5]$

coins =  $[(x, y)] \leftarrow \text{find\_nearest\_coins}(\dots)$  ← would have to be tweaked

spots =  $[(x, y)] \leftarrow \text{find\_crate\_bombing\_spots}()$

danger =  $[(x, y)] \leftarrow \text{find\_danger}()$  ↗

[ directions(coins)  
[ directions(spots)  
[ directions(danger)  
] make into features

Algorithm design:      find-crate-bombing-spots()

Input: "game-state"

Output:  $[(x, y)]$  (multiple coordinates possible)

Aim: Solve at what spot a bomb placed results in the highest coin speed.

Task: calculate the best spot(s) to place the next bomb, such that  $\frac{\text{expected collected coins}}{\text{total time}}$  is maximised

Assumptions:

- Only value of bombing crates is coins

revealed maybe?  
After, and if, coins are revealed, the coin and crate algorithms can decide on further action.

Potential issues:

- What about the safety aspect of bombing crates?  
Isn't there risk of trapping oneself between walls/crates and bombs?
- May be blind to currently non-reachable coins.

Mathematically:

$$\text{return } \underset{\vec{s}}{\operatorname{argmax}} (V_c(\vec{s}, \vec{p}, F))$$

$$V_c(\vec{s}, \vec{p}, F) = \frac{c(\vec{s}, F)}{t_{\text{tot}}(\vec{s}, \vec{p}, F)}$$

$$c(\vec{s}, F) = k_d(\vec{s}, F) \cdot p_c(F)$$

$$p_c = c_h / k(F)$$

$$t_{\text{tot}}(\vec{s}, \vec{p}, F) = t_r(\vec{s}, \vec{p}, F) + t_b$$

$\vec{s} = (x, y)$  bomb placing spot

$\vec{p} = (x, y)$  own current position

$F = \text{game\_state}['field']$

$V_c(\vec{s})$  expected coin collection speed

$c(\vec{s})$  expected number of coins revealed by bomb  
good measure?

$t_{\text{tot}}$  expected time until next bomb can be placed

$k_d(\vec{s})$  crates destroyed by bomb at spot  $\vec{s}$

$p_c$  density of remaining hidden coins

$c_h$  number of remaining hidden coins

$k$  number of remaining crates

$t_r(\vec{s}, \vec{p})$  time to reach spot  $\vec{s}$  from  $\vec{p}$

$t_b$  time until next bomb can be placed

Ideas for implementation:

$t_r(\vec{s}, \vec{p}, F)$  could be a map  $\in \mathbb{N}^{B \times B}$  based on own position  $\vec{p}$   
calculated with a breadth-first search  
Non reachable values could be NaN or None

board length = 17 ?

$t_b = 7$  empirical value

$k(F)$  count all 1's in  $\text{game\_state}['field']$

$c_h$  set up some coin counter that tracks appearing and disappearing coins

$k_d(\vec{s})$  map bomb explosion radius over  $F$   
(maybe update only if  $F$  changes?)