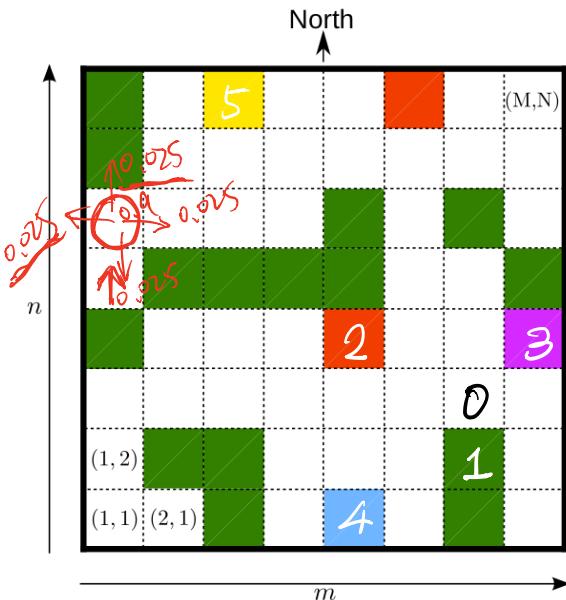


## ① Problem Setup

map = double  $m \times n$



mapSize = [ , ]

random word or not

plot : policy / cost = T/F

$\gamma$  AMMA R NC P-WIND

# FREE TREE SHOOTER PLCK\_UP

## DROP-OFF BASE

input = { ↑ ↓ → ← ⊗ }

NORTH SOUTH EAST WEST HOVER

1 2 3 4 5

stateSpace : Kx1 Matrix

## index

$\begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 \end{pmatrix}$  ! not continuous in  $N'$ , because no TREE

2 | 1 1 1

## Some important states :

$$3 \mid 120$$

① terminal state :  $(m_{drop} \ n_{drop} \ 1)$

$$\begin{array}{c|cc} 4 & & \\ & | & | \\ & 2 & \end{array}$$

② base state after crash : ( $m_b$   $n_b$  0)

1

⚠ index  $\leftrightarrow (m, n, p)$

✓

1

1

1

1

1 m n

R m n l

How to debug : switch xxx implemented from "false" to "true"

## ② Help functions

<1>  $bool = \text{invalid\_input}(m, n, u, \text{map})$

→ when position is  $(m, n)$ , input  $u$  is invalid ( $=1$ ) or valid ( $=0$ ) depends on boundary and trees

<2>  $idx = \text{index}(m, n, p, \text{stateSpace})$

→ find the index of the state  $(m, n, p)$  in stateSpace  
"for"

<3>  $pr = p\_not\_shot(m, n, \text{map})$

→ return the Pr. of not get shot by residents when position is  $(m, n)$

→ for every residents,  $pr = pr \times (1 - \text{getshot})$

<4>  $tpm = \text{tempM}(m, n, u, \text{map})$

→ if  $u$  is valid input for  $(m, n)$ , return  $m$  after  $u$ .  
else return NaN with error displayed

<5>  $tpn = \text{tempN}(m, n, u, \text{map})$

→ if  $u$  is valid input for  $(m, n)$ , return  $n$  after  $u$ .  
else return NaN with error displayed

### ③ Compute Transition Probabilities

$P(i, j, u) = \text{zeros}(k, k, 5)$ : the Pr from state  $i$  to state  $j$  with input  $u$ .

for  $i = 1 : K$

if  $i = \text{Terminal}$  then  $P(i, i, :) = 1$ ,  $P(i, \neq i, :) = 0$   
continue;

end

for  $u = 1 : 5$

if  $\text{invalid}(m, n, u, \text{map}) == 1$   $P(i, j, u) = \underline{\text{NaN}}$ .  
continue;

end

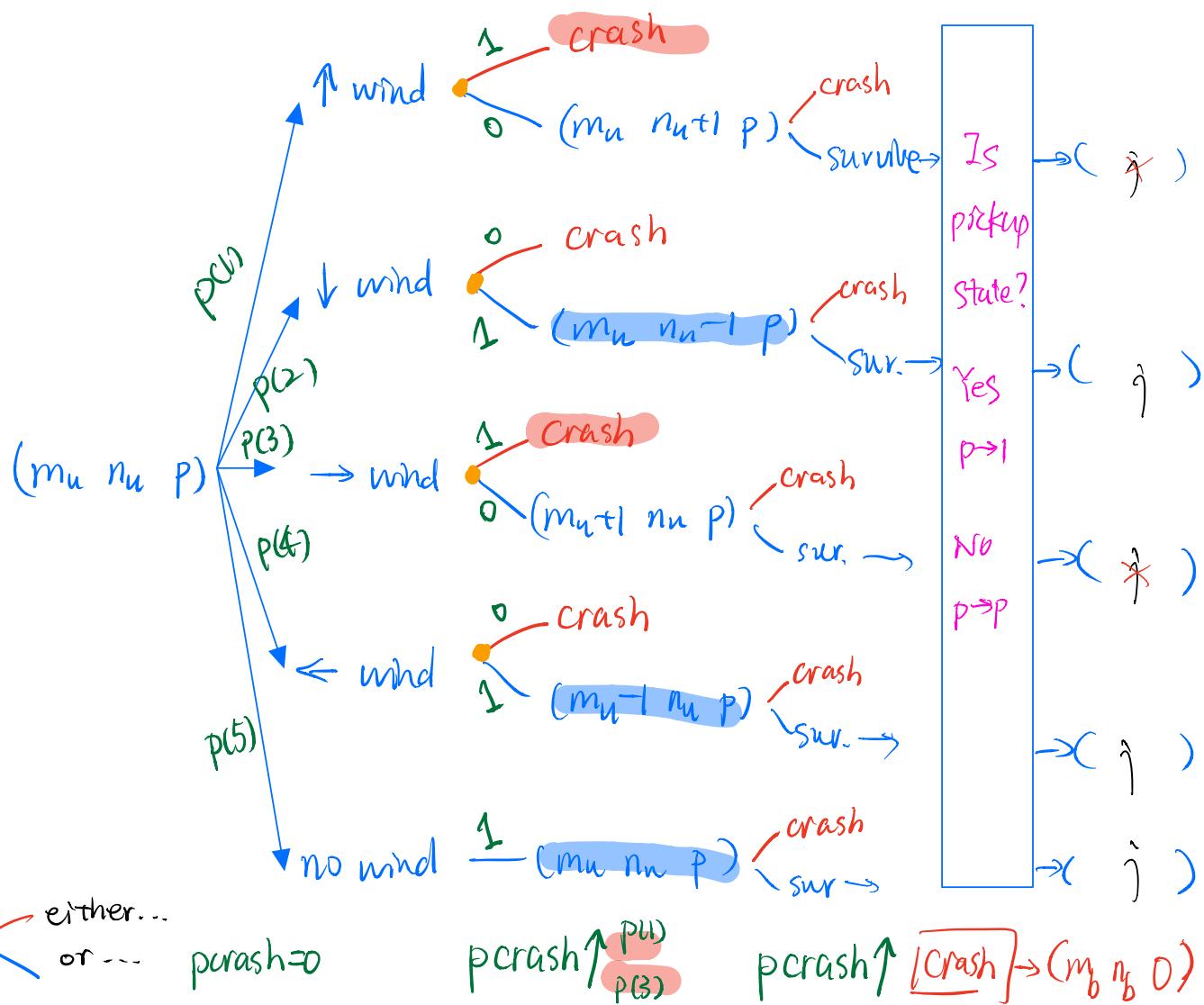
$$(m, n, p) \xrightarrow{u} (m_u, n_u, p)$$

$$\text{NaN} + \text{NaN} = \text{NaN}$$

$$\text{NaN} \times \text{NaN} = \text{NaN}$$

$$\text{number} + \text{NaN} = \text{NaN}$$

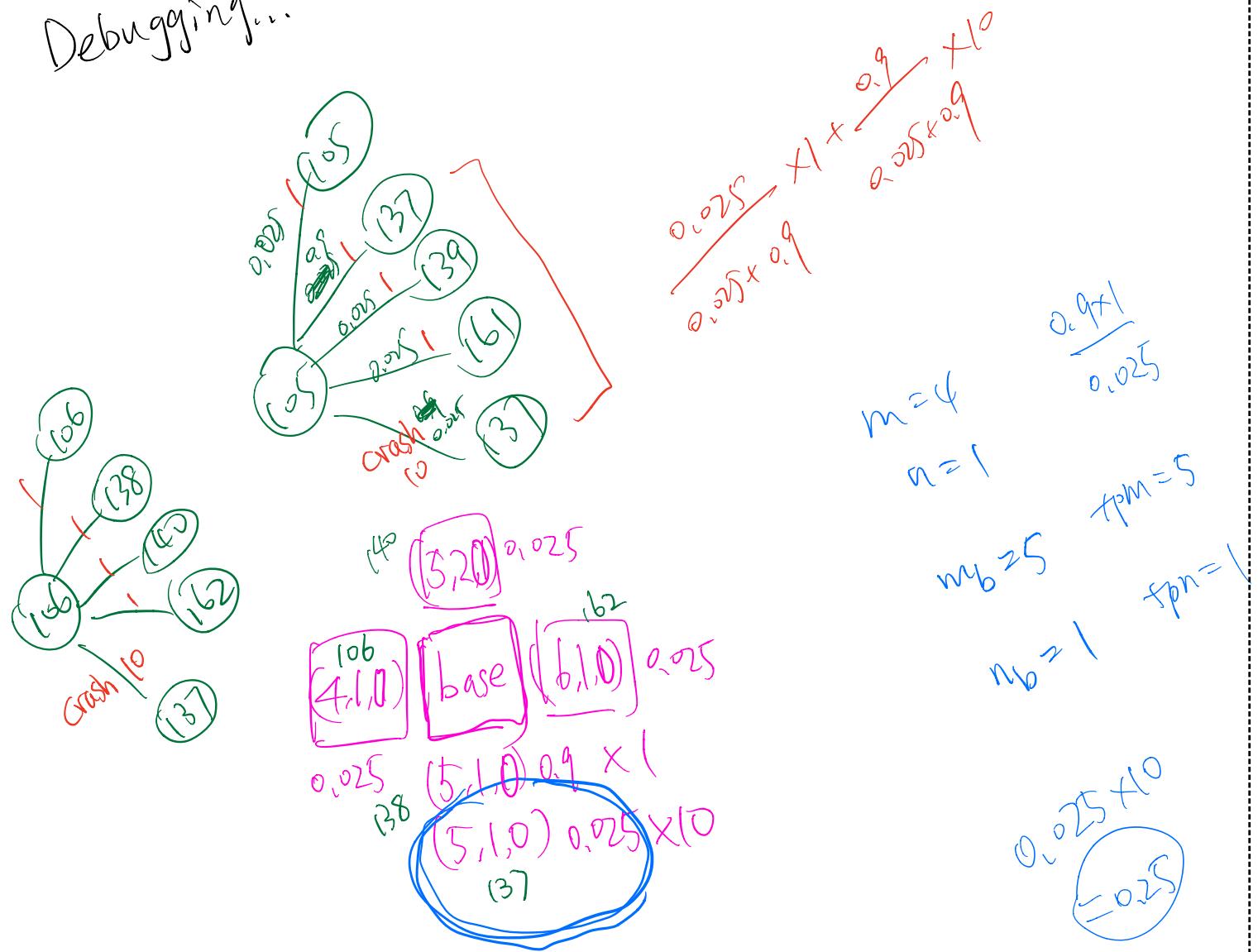
$$P(i, j, u) = \underline{\text{NaN}}$$



Check out:  $P(i, \cdot, u)$  : with stateSpace / Map

- ① if  $u$  is invalid ,  $P(i, \cdot, u) = \text{NaN}$
  - ② if no shooter ,  $P(i, \cdot, u) = \text{some 0.025}$   
Base 137 ↑
  - ③ if shooter ,  $P(i, \cdot, u) =$

# Debugging...



④ Compute Stage Costs.

$$\begin{aligned} \mathbb{E}(i, u) &= q(i, u) = \underset{(w|x=i, u=u)}{\mathbb{E}} [g(x_k, u_k, w)] \\ &= \sum_j P(i, j, u) \cdot \underline{\text{cost}(i, j, u)} \end{aligned}$$

1) Given any initial conditions  $V_0(1), \dots, V_0(n)$ , the sequence  $V_l(i)$  generated by the iteration

$$V_{l+1}(i) = \min_{u \in \mathcal{U}(i)} \left( q(i, u) + \sum_{j=1}^n P_{ij}(u) V_l(j) \right), \quad \forall i \in \mathcal{S}^+ \quad (4.7)$$

where  $\mathcal{S}^+ := \mathcal{S} \setminus \{0\}$  and

$$q(i, u) := \underset{(w|x=i, u=u)}{\mathbb{E}} [g(x, u, w)],$$

converges to the optimal cost  $J^*(i)$  for all  $i \in \mathcal{S}^+$ ;

$\text{cost}(i, j, u, \text{stateSpace}, \text{map})$ :  
 average cost when state  $i \xrightarrow{u} j$  with input  $u$ .

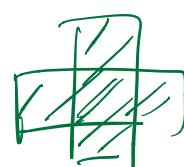
⚠ crash  $\Rightarrow$  "j=basestate" but "j=base state"  $\not\Rightarrow$  crash  
 so apply average cost

Solution:

$$\text{cost}(i, \text{base}, u) = N_c \times (1 - p(\text{not crash} | i, \text{base}, u)) + \\ 1 \times p(\text{not crash} | i, \text{base}, u)$$

$$\textcircled{2} \quad p(\text{not crash} | i \xrightarrow{u} \text{base}) = \frac{p(\text{not crash} \& i \xrightarrow{u} \text{base})}{p(i \xrightarrow{u} \text{base})}$$

$$= \begin{cases} p(\text{good wind}) \times p(\text{not get shot}) / p(i, \text{base}, u) & \text{if "next to"} \\ 0 & \text{if not "next to"} \end{cases}$$



next to base.

## ⑤ Value Iteration

- 1) Given any initial conditions  $V_0(1), \dots, V_0(n)$ , the sequence  $V_l(i)$  generated by the iteration

$$V_{l+1}(i) = \min_{u \in \mathcal{U}(i)} \left( q(i, u) + \sum_{j=1}^n P_{ij}(u) V_l(j) \right), \quad \forall i \in \mathcal{S}^+ \quad (4.7)$$

where  $\mathcal{S}^+ := \mathcal{S} \setminus \{0\}$  and

$$q(i, u) := \underset{(w|x=i, u=u)}{\mathbb{E}} [g(x, u, w)],$$

converges to the optimal cost  $J^*(i)$  for all  $i \in \mathcal{S}^+$ ;

- 2) The optimal costs satisfy the Bellman Equation (BE):

$$J^*(i) = \min_{u \in \mathcal{U}(i)} \left( q(i, u) + \sum_{j=1}^n P_{ij}(u) J^*(j) \right) \quad \forall i \in \mathcal{S}^+;$$

- 3) The solution to the BE is unique;

- 4) The minimizing  $u$  for each  $i \in \mathcal{S}^+$  of the BE gives an optimal policy, which is proper.

stateSpace :  $K \times 1$   
index

1	1	1	0
2	1	1	1
3	1	2	0
4	1	2	1
⋮	⋮	⋮	⋮

for iterator = 1 : K

if iter = T, continue

$v(\text{iter}) = \min \{ \dots \}$

M = [ Inf ]

if not invalid u

Question ?

How to speed up VI?

terminal state {0}

$i \in \mathcal{S}^+ = \mathcal{S} \setminus \{0\}$ .

m	n	0
m	n	1

# ⑥ Policy Iteration

## Policy Iteration

**Initialization:** Initialize with a proper policy  $\mu^0 \in \Pi$ .

**Stage 1 (Policy Evaluation):** Given a policy  $\mu^h$ , solve for the corresponding cost  $J_{\mu^h}$  by solving the linear system of equations

$$J_{\mu^h}(i) = q(i, \mu^h(i)) + \sum_{j=1}^n P_{ij}(\mu^h(i)) J_{\mu^h}(j), \quad \forall i \in \mathcal{S}^+.$$

(5.3)

$J = Q + P \cdot J$

**Stage 2 (Policy Improvement):** Obtain a new stationary policy  $\mu^{h+1}$  as follows

$$\mu^{h+1}(i) = \arg \min_{u \in \mathcal{U}(i)} \left( q(i, u) + \sum_{j=1}^n P_{ij}(u) J_{\mu^h}(j) \right), \quad \forall i \in \mathcal{S}^+.$$

$q(i, u) + \sum_{j=1}^n P_{ij}(u) J_{\mu^h}(j)$

(5.4)

Iterate between Stage 1 and 2 until  $J_{\mu^{h+1}}(i) = J_{\mu^h}(i)$  for all  $i \in \mathcal{S}^+$ .

$\mu_0 = 5 \times \text{ones}(k, 1)$  (Assumption:  $\mu$  is valid)

stateSpace : $K \times 1$	
index	
1	1 1 0
2	1 1 1
3	1 2 0
4	1 2 1
⋮	⋮
T	⋮
m	n 0
K	m n 1

$\cup \mathcal{S}^t = \mathcal{S} / \{0\}$

$$J_{\mu^h} = \begin{bmatrix} J_{\mu^h}(1) \\ J_{\mu^h}(2) \\ \vdots \\ J_{\mu^h}(T) = 0 \\ \vdots \\ J_{\mu^h}(K) \end{bmatrix} \quad Q_{\mu^h} = \begin{bmatrix} q(1, \mu^h(1)) \\ q(2, \mu^h(2)) \\ \vdots \\ q(T, \mu^h(T)) = 0 \\ \vdots \\ q(K, \mu^h(K)) \end{bmatrix}$$

$$P = \begin{bmatrix} & & & \\ & & & \\ & & & \\ \cdots & P(i, j, \mu^h(i)) & \cdots & \\ & & & \end{bmatrix}$$

terminal state  $\{0\}$

$i \in \mathcal{S}^t = \mathcal{S} / \{0\}$ .

## ⑦ Linear Programming.

**Theorem 5.2.** The solution to the optimization problem

$$\begin{aligned} & \underset{V}{\text{maximize}} \quad \sum_{i \in \mathcal{S}^+} V(i) \\ & \text{subject to} \quad V(i) \leq \left( q(i, u) + \sum_{j=1}^n P_{ij}(u) V(j) \right), \quad \forall u \in \mathcal{U}(i), \quad \forall i \in \mathcal{S}^+. \end{aligned} \tag{5.12}$$

also solves the Bellman Equation to yield the optimal cost  $J^*$  for the SSP problem.

maximize  $V_1 + V_2 + \dots + V_{T-1} + V_{T+1} + \dots + V_K$

$\Leftrightarrow$  minimize  $-V_1 - V_2 - \dots - V_{T-1} - V_{T+1} - \dots - V_K$

$\Leftrightarrow$  minimize  $f^T V$ ,  $f = \begin{bmatrix} -1 \\ -1 \\ \vdots \\ -1 \\ 0 \\ \vdots \\ -1 \end{bmatrix}$   $V = \begin{bmatrix} V_1 \\ V_2 \\ \vdots \\ V_T \\ V_{T+1} \\ \vdots \\ V_K \end{bmatrix}$

for  $i = 1 : K$

if  $i = T$  continue.

for  $u = 1 : 5$

if  $\text{isnan}(G(i, u))$  continue.

$$\rightarrow V(i) \leq q(i, u) + P(i, 1, u)V(1) + P(i, 2, u)V(2) + \dots + P(i, i, u)V(i) + \dots + P(i, T, u)V(T) + \dots + P(i, K, u)V(K)$$

$$[-P(i, 1, u) \quad -P(i, 2, u) \quad \dots \quad -P(i, i, u) \quad \dots \quad -P(i, K, u)] \begin{bmatrix} V(1) \\ V(2) \\ \vdots \\ V(i) \\ V(T) \\ \vdots \\ V(K) \end{bmatrix} \leq q(i, u)$$

optimal policy?  $\mu^*(\cdot)$

2) The optimal costs satisfy the Bellman Equation (BE):

$$J^*(i) = \min_{u \in \mathcal{U}(i)} \left( q(i, u) + \sum_{j=1}^n P_{ij}(u) J^*(j) \right) \quad \forall i \in \mathcal{S}^+;$$

$$\begin{bmatrix} 2624 \times 1 \\ 4576 \times 1 \end{bmatrix}$$