April 13 (Tue)
Reward Function Example: Ref.
  a github 105d ; b ML-Robot-Arm

Example: Reward Function, pp.1,
    Section.

$$S \times A = (S_1 a_1, S_1 a_2, \cdots,$$
$$S_n a_n)$$
$$S = \{S_1, S_2\}$$
$$A = \{a_1, a_1\}$$
$$S \times A = \{S_1, S_2\} \times \{a_1, a_2\}$$
$$= \{S_1 a_1, S_1 a_2, S_2 a_1, S_2 a_2\}$$

Reward $r : S \times A \rightarrow \underline{\underline{R}}$

$$R = \{r_1, r_2, r_3, r_4\}$$

whole collection
of all reward

$S_1 a_1 \quad S_1 a_2 \quad \cdots \quad S_2 a_2 \cdots (1)$

Denote a reward at time $t$ as

$$R(s_t, a_t) = R\left(\text{Hit the Ground}, a_t\right) = -1 \quad \cdots (2)$$

$$R(s_t, a_t) = R\left(\text{Reachs the target}, a_t\right) = 1 \quad \cdots (3)$$

$\overrightarrow{P(x, y, z)}$ in $X_w - Y_w - Z_w$ World Coordinate:



$\overrightarrow{P_{tgt}}, \overrightarrow{P_{END}}$

End Effector

Move $\overrightarrow{P_{END}}$ to $\overrightarrow{P_{tgt}}$ By $\underline{N}$ steps

$$(S_t, a_t) \text{ for } t = 1, 2, \ldots, N \quad \cdots (3b)$$

if $\overrightarrow{P_{END}} = \overrightarrow{P_{tgt}}$ then End the Episod.

$$\overrightarrow{P_{END}} \simeq \overrightarrow{P_{tgt}}, \quad \overrightarrow{P_{END}} - \overrightarrow{P_{tgt}} \simeq 0 \quad \cdots (4)$$

$$\| \overrightarrow{P_{END}} - \overrightarrow{P_{tgt}} \| \ll \epsilon \quad \cdots (5)$$

$$\overrightarrow{P_{END}}(x, y, z) = \overrightarrow{P_{END}}(x(t), y(t), z(t))$$

for $t = 1$,

$$\overrightarrow{P_{END}}(x(1), y(1), z(1))$$

for $t = 2$

$$\overrightarrow{P_{END}}(x(2), y(2), z(2))$$

for $t = i$

$$\overrightarrow{P_{END}}(x(i), y(i), z(i))$$

$$\vdots$$

$$\overrightarrow{P_{END}}(x(N), y(N), z(N))$$

for $t-1$ to $t$, the distance (movement of $\overrightarrow{P_{END}}$)

$$\overrightarrow{P_{END}}(x(t), y(t), z(t)) - \overrightarrow{P_{END}}(x(t-1), y(t-1), z(t-1)) \quad \cdots (6)$$

Suppose

$$\overrightarrow{P_{END}}(x(0), y(0), z(0)) = (1, 1, 2)$$

$$\overrightarrow{P_{END}}(x(1), y(1), z(1)) = (0.5, -1, 2)$$

Find the distance

$$\overrightarrow{P_{END}}(x(1), y(1), z(1)) - \overrightarrow{P_{END}}(x(0), y(0), z(0)) =$$

$$(1 - 0.5, -1 - 1, 2 - 2) = (0.5, -2, 0)$$

$$\| \overrightarrow{P_{END}}(x(1), y(1), z(1)) - \overrightarrow{P_{END}}(x(0), y(0), z(0)) \|_2$$

$$= \sqrt{(x(1) - x(0))^2 + (y(1) - y(0))^2 + (z(1) - z(0))^2}$$

$$= \sqrt{0.5^2 + (-2)^2 + 0^2} = \sqrt{4.25} \quad \text{Distance to the target}$$

$$d = \| \overrightarrow{P_{END}}(t-1) - \overrightarrow{P_{tgt}} \|_2, @ t-1$$

$$\downarrow d(t-1)$$

$$d(t) = \| \overrightarrow{P_{END}}(t) - \overrightarrow{P_{tgt}} \|_2, @ t$$

$$R(S_t, a_t) = \frac{d(\overrightarrow{P_{END}}(t-1) - \overrightarrow{P_{tgt}}) - d(\overrightarrow{P_{END}}(t) - \overrightarrow{P_{tgt}})}{d(\overrightarrow{P_{END}}(0) - \overrightarrow{P_{tgt}})} \quad \cdots (9*)$$

Difference of the distance

Assume $d(\overrightarrow{P_{END}}(t-1) - \overrightarrow{P_{tgt}}) = \sqrt{6.25}$    Positive Reward.

$$d(\overrightarrow{P_{END}}(t) - \overrightarrow{P_{tgt}}) = \sqrt{4.25}$$

$$\widetilde{R}(S_t, a_t) = \sqrt{6.25} - \sqrt{4.25}$$

$\downarrow$ Normalization

$$R(S_t, a_t) = \widetilde{R} \Big/ d(\overrightarrow{P_{END}^{(0)}} - \overrightarrow{P_{tgt}})$$

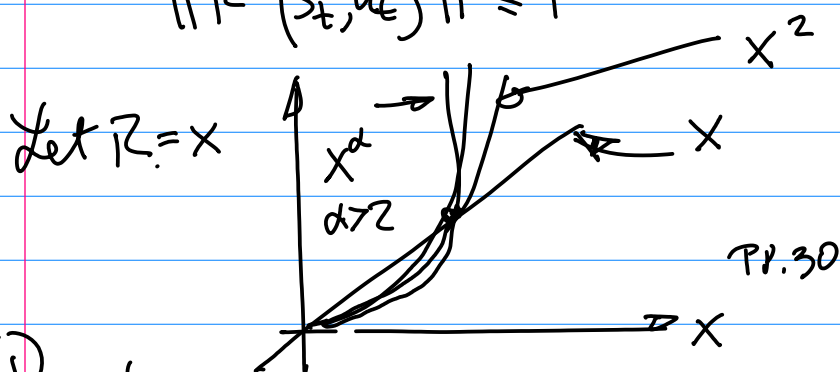For Negative Reward from Eqn $(9*)$ holds good as well, Since Eqn $(9*)$ will give negative value.

$$R(s_t, a_t) = \frac{d(\vec{P}_{END}(t-1) - \vec{P}_{tgt}) - d(\vec{P}_{END}(t) - \vec{P}_{tgt})}{d(\vec{P}_{END}(0) - \vec{P}_{tgt})}$$

$$-1 \leq R(s_t, a_t) \leq 1$$

$$-1 \leq R^2(s_t, a_t) \leq 1 \Rightarrow 0 \leq R^2(s_t, a_t) \leq 1$$

$$\| R^\alpha(s_t, a_t) \| \leq 1$$

Let $R = x$



$x^\alpha$
$\alpha > 2$

$x^2$

$x$

PP.30

① Code/Source Walk-Through ② Dynamic Programming

Non-Linear Accelerated Reward Function
with $\alpha > 2$

Question: How About $e^{\pm R}$ or $\log_a R$



$\log_a x \quad a > 1$