March 23rd (T) yy, 7L.

1. Demo $\begin{cases} \underline{a} \text{ Photo/Image} \\ \underline{b} \text{ Video} \end{cases}$
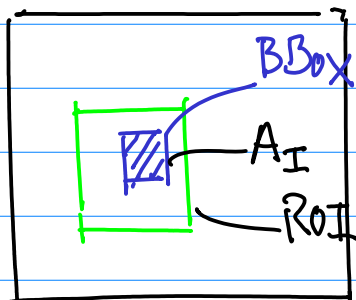
~ roi.py ⟍ main
~ Box.py ⟋ Program
Bounding

Testing Program : rectangle_intersection.py.

To do : Send readme.txt,

Threshold T = 0.25



Area (Intersection)
Blue

$A_I$ V.S. T

$A_I \geqslant T \rightarrow$ Alarm
$A_I < T$ No Action

Define 2 Thresholds. $\begin{cases} \text{Upper Threshold} \\ \quad T_u \\ \text{Lower Threshold} \\ \quad T_L \end{cases}$

$\begin{cases} A_I \\ T_u \simeq 50\% \\ \quad A_I \text{ Interested} \\ T_L \simeq 10\% \\ \text{Discarded} \\ \quad A_I \text{ is not to be Considered} \end{cases}$

Discard
$A_I$ if it is Bigger than $T_u$

DRL-SAC March 24 (Wed)

Soft Actor Critic Algorithm & ATPS.
UC Berkeley & Google      Pp. 4
Standard RL Objective to Maximize
Reward

$$\text{Max} \left\{ \sum_{t} E_{(s_t, a_t) \sim \rho_\pi} \left[ r(\vec{s_t}, \vec{a_t}) \right] \right\} \quad \cdots (1)$$

To Learn Policy

$$\underline{\pi(\vec{a_t} | \vec{s_t})} \quad \cdots (2)$$

to maximize Eqn (1).

Augment Eqn (1) with
Entropy to maximize
entropy at each visited
state.

So

$$\alpha \mathcal{H}\left( \pi(a_t | s_t) \right) \quad \cdots (3)$$

Weight ⟍ Entropy  Independent
            Function  Variable π
("Temperature" parameter  Policy
determines the relative
importance of Entropy term)

Hence Reward in Eqn(1) becomes

$$\max \left\{ \sum_t E_{(s_t, a_t) \sim \rho_\pi} \left[ r\left(\vec{s_t}, \vec{a_t}\right) \right] + \alpha \mathcal{H} \left( \pi\left(a_t | s_t\right) \right) \right\}$$

Soft Policy Iteration

Q-value (Reward)

$$Q : S \times A \to R \qquad \cdots (1)$$

Discount factor

$$T^\pi Q\left(\vec{s_t}, \vec{a_t}\right) \overset{\triangle}{=} r\left(\vec{s_t}, \vec{a_t}\right) + \gamma E_{\vec{s}_{x+1} \sim p} \left[ V\left(\vec{s}_{t+1}\right) \right] \qquad \cdots (1*)$$

Bellman
Operator
for Each Action
leads to Reward

Reward Function

Q-Value function Computed
when $\pi\left(\vec{s}_{t+1} | \vec{a}_{t+1}\right)$ is
applied, and reward is
generated

where

$$V\left(\vec{s}_{t+1}\right) = E_{\vec{a_t} \sim \pi} \left[ Q\left(\vec{s_t}, \vec{a_t}\right) - \alpha \log \pi\left(\vec{a_t} | \vec{s_t}\right) \right] \qquad \cdots (2)$$

Policy introduced Q-Value
(Reward)



Note: Based on (2)

$\pi(\cdot) \quad \alpha \log \pi\left(\vec{a_t} | \vec{s_t}\right)$, So $\pi\left(\vec{a_t} | s_t\right)$ is
Reward Based Function

Lemma 1. Bellman Operator T   "0" initial time t=0.

Given initial Condition $\vec{Q^0} : S \times A \to R$ with $|A| < \infty$

Define Bellman Operator T as

"All" total Number of finite
Actions

$$Q^{K+1} = T^\pi Q^K \qquad \cdots (3)$$

Time Index

$$Q^{k+1} = T^{\pi} Q^{k}$$

Q-Value
(Reward) Function  Bellman
                    Operator
                    under Policy $\pi$

For k=0, we have

$$Q^1 = T^{\pi} Q^0$$

$$K=1, \quad Q^2 = T^{\pi} Q^1 = T^{\pi}(T^{\pi} Q^0)$$
$$= (T^{\pi})^2 Q^0 \qquad \ldots(4)$$

$\ldots$

$$K=i \quad Q^{i+1} = (T^{\pi})^{i+1} Q^0 \qquad \ldots(5)$$

$\begin{cases} \text{Update Q-Value (Reward)} \\ \text{Update Policy } \pi \end{cases}$

$$\pi_{New} = \underset{\pi' \in \Pi}{\text{argmin}} \; D_{KL}\left(\pi'(\cdot|\vec{S_t}) \,\Big\|\, \frac{e^{\frac{1}{2} Q^{\pi_{old}}(\vec{S_t}|\cdot)}}{Z^{\pi_{old}}(\vec{S_t})}\right) \quad \ldots(6)$$

Note: $1^o$ $\pi \in \Pi$

Policy from (belongs to)
all collections of the
Policies     Collection of All
             Polices

Note: "Argmin" find minimum
from a Collection of elements,
functions, etc.

Example $\text{Argmin}\{1, -1.2, -110, 212\}$
$= -110.$

$2^o$ $\pi(\vec{a_t}|\vec{S_x})$ Policy

$\pi(\cdot|\vec{S_t})$ All Polices
        for each/every $a_x$
        from $\vec{a_t}$.

$\pi'(\cdot|\vec{S_t})$ one selected from
        the All polices (
all action, e.g. $\cdot$ "for $\vec{a_t}$)

$2^o$ $D_{KL}$ KL: Kull Back—Leibler
                Divergence

$3^o$ $\pi'$ One $\pi$ realization
from the collection of All
Polices, A particular Policy

Generalize Eqn(b), we have

$$\pi_{New} = f\left(\pi_{old} \text{ from one of} \mid \text{Conditions}\right)$$
the all possible
Policies
$$\cdots (6b)$$

$f$ to be minimization problem, so

$$\pi_{New} = \text{argmin}\, \partial_{KL}\left(\pi_{old} \text{ from One of} \mid \text{Condition}\right)$$
the all polices

$$\cdots (6c) \qquad \text{Probability or likelihood}$$

Note: Parameterized Policy as Gaussian Distribution

$$\frac{e^{\frac{1}{\alpha} Q^{\pi_{old}}\left(\vec{S_t}, \cdot\right)}}{Z^{\pi_{old}}\left(\vec{S_t}\right)} \qquad \cdots (7)$$

Note: In NN Softmax Activation Function

We map Neuron Output in $(-\infty, +\infty)$ to Probabilistic distribution $[0,1]$ e.g,

$$f : Z \in (-\infty, +\infty) \rightarrow f(z) \in [0,1]$$

where

$$f(z) = \frac{e^{z_i}}{\sum_{j=1}^{N} e^{z_j}} \qquad \cdots (8)$$

Lemma 2.

$$Q^{\pi^*}\left(\vec{S_t}, \vec{a_t}\right) \geq Q^{\pi}\left(\vec{S_t}, \vec{a_t}\right)$$

where $\pi^*$ is from Eqn(b)

Policy $\pi \in \Pi$, under Eqn(b) converges to

$$\pi^*$$

Note: Eqn(b) needs Better Explanation, as why it is formulated as

$D_{KL}$ minimization Problem.

Now, introduce parameterized Q function, and Policy $\pi$.

$Q_{\theta}$, $\pi_{\phi}$, $\theta$ (Theta) $\phi$ (phi) are

Theta  phi

$\phi$: Policy Parameters

$$J_Q(\theta) = E_{(\vec{s}_t, \vec{a}_t) \sim \theta} \left[ \frac{1}{2} \left( Q_\theta(\vec{s}_t, \vec{a}_t) - \left( r(\vec{s}_t, \vec{a}_t) + \gamma E_{s_{t+1} \sim P} [V_\theta(\vec{s}_{t+1})] \right) \right)^2 \right]$$

Objective Function $\qquad \dots (9)$

$\qquad\qquad\qquad$ Eqn (1*) Bellman Iteration

$Q(\theta)$ Q-value (Reward) Function with Parameter $\theta$ (Theta)

$$\nabla J_Q(\theta) = E_{(\vec{s}_t, \vec{a}_t) \sim \theta} \left\{ \left[ Q_\theta(\vec{s}_t, \vec{a}_t) - \left( r(\vec{s}_t, \vec{a}_t) + \gamma E_{s_{t+1} \sim P}[V_\theta(\vec{s}_{t+1})] \right) \right] \cdot \nabla Q_\theta(\vec{s}_t, \vec{a}_t) \right\} \qquad \dots (10)$$

$$J_\pi(\phi) = E_{\vec{s}_t \sim D} \left[ E_{\vec{a} \sim \pi_\phi} \left[ \alpha \log(\pi_\phi(\vec{a}_t | \vec{s}_t)) - Q_\theta(\vec{s}_t, \vec{a}_t)) \right] \right] \qquad \dots (11)$$

$$J_\pi(\phi) = E_{\vec{s}_t \sim D} \left[ E_{\vec{a}_t \sim \pi_\phi} [\cdot] \right]$$

$\qquad$ for all states $\quad$ For all actions $\vec{a}_t$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\pi(\vec{a}_t, \vec{s}_t)$

Note:

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ Eqn(1), PP. 7

$\frac{d}{dx} \ln(x) = \frac{1}{x}, \quad \frac{d}{dx} \log_a(x) = \frac{1}{x} \frac{1}{\ln(a)} \qquad \dots (12)$ $\quad$ 2. id.

$\vec{a}_t = f_\phi(\epsilon_\pi; \vec{s}_t) \qquad \dots (13)$ Parameterized Policy?

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ Action

Action Vector

multi-dimensional $\qquad\quad J_\pi(\phi) = E_{\vec{s}_t \sim D, \epsilon_t \sim N} \left[ \alpha \log \pi_\phi(f_\phi(\epsilon_t; \vec{s}_t) | \vec{s}_t) - \right.$

$\vec{a}_t = (a_{1_t}, a_{2_t}, \dots, a_{N_t})$ Parameter $\qquad\qquad\qquad\qquad \left. Q_\theta(\vec{s}_t, f_\phi(\epsilon_t; \vec{s}_t)) \right]$

$\qquad\qquad\qquad\qquad$ for $\pi$ $\qquad\qquad \dots (14)$

$\pi_\phi \left( f_\phi(\epsilon_t | \vec{s}_t) | \vec{s}_t \right)$ is from

$\quad \pi_\phi \left( \vec{a}_{\phi t} | \vec{s}_t \right)$, or is from

$\quad \pi_\phi (\vec{a}_t | \vec{s}_t)$ is from Eqn(2) PP.2.

$\hat{\nabla}_\phi J_\pi(\phi) = \nabla_\phi \alpha \log \left( \pi_\phi (\vec{a}_t | \vec{s}_t) \right) + \nabla_{\vec{a}_t} \alpha \log \left( \pi_\phi (\vec{a}_t | \vec{s}_t) \right) -$

$\qquad\qquad\qquad \nabla_{\vec{a}_t} Q(\vec{s}_t, \vec{a}_t)) \nabla_\phi f_\phi (\epsilon_t ; \vec{s}_t)$

$\qquad\qquad\qquad\qquad\qquad\qquad \cdots (15)$

Note: Jacobian

Suppose $\vec{f} : \mathbb{R}^n \to \mathbb{R}^m$

$\quad \vec{x} \in \mathbb{R}^n, \vec{f}(\vec{x}) \in \mathbb{R}^m \quad \cdots (16a)$

$\nabla f_1(\vec{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} \\ \vdots \\ \frac{\partial f_1}{\partial x_n} \end{bmatrix}$

Example: $f(x_1, x_2) = f(\vec{x}), \vec{x} \in \mathbb{R}^2$

Then $\qquad\qquad\qquad \vec{x} = (x_1, x_2)$

$\quad \vec{f}(\vec{x}) = \left( f_1(\vec{x}), f_2(\vec{x}) \right), \qquad \cdots(16b)$

$\qquad\qquad\qquad\qquad\qquad\qquad \cdots (16e)$

$\quad \vec{f}(\vec{x}) \in \mathbb{R}^2. \qquad \cdots(16c)$

$\nabla^T f_1(\vec{x})$

$= \left( \frac{\partial f_1}{\partial x_1} \cdots \frac{\partial f_1}{\partial x_n} \right)$

$\qquad\qquad\qquad \cdots(16f)$

Jacobian: Matrix of All its 1st order

$\qquad$ Partial Derivatives.

Hence, Jacobian

$\vec{J} \triangleq \left[ \frac{\partial \vec{f}}{\partial x_1}, \cdots, \frac{\partial \vec{f}}{\partial x_n} \right] = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$

$\vec{J} = \begin{bmatrix} \nabla^T f_1 \\ \vdots \\ \nabla^T f_m \end{bmatrix} \quad \cdots(16g)$

gradient of $f_1(\vec{x})$

$\qquad\qquad\qquad \cdots(16d)$

Note: $\nabla f_1(\vec{x})$, if

$\qquad f(\vec{x})$ changed to $f(\vec{x}; \phi)$

$\nabla f_1 =$

$\nabla_\phi f_1(\vec{x}; \phi)$, so we are

taking Partial derivative wrt $\phi$.

$\phi = (\phi_1, \phi_2, \cdots, \phi_n)$, $\vec{\phi}$ is written

Simplified as $\phi$.

Automating Entropy Adjustment

$$\max_{\pi_{0:T}} E_{\rho_\pi} \left[ \sum_{t=0}^{T} r(\vec{s}_t, \vec{a}_t) \right] \quad \text{s.t.} \quad E_{(\vec{s}_t, \vec{a}_t) \sim \rho_\pi} \left[ -\log(\pi_t(\vec{a}_t | \vec{s}_t)) \right] \geq \mathcal{H}, \forall x$$

$$\cdots (17) \qquad \cdots (18)$$

$\underset{=}{a}$ $r(\vec{s}_t, \vec{a}_t)$ : reward function

$\underset{=}{b}$ $\sum_{t=0}^{T} r(\vec{s}_t, \vec{a}_t)$ : Summation of

reward function from $t=0$
to $t=T$ ;

$\underset{=}{c}$ $\max \left[ \sum_{t=0}^{T} r(\vec{s}_t, \vec{a}_t) \right]$

Maximize the reward function
from $t \in [0, T]$ period.

$\underset{=}{d}$ $\max E \left[ \sum_{t=0}^{T} r(\vec{s}_t, \vec{a}_t) \right]$

maximize expected reward

function from $[0, T]$

Now Add Notation to
detail it up

$$\max E \left[ \sum_{t=0}^{T} r(\vec{s}_t, \vec{a}_t) \right]$$

$\pi_{0:T}$     $\rho_\pi$ under

Policy $\pi$ from   Policy, Policy
time 0 to T.    $\pi$

So, we have Eqn (17).

$-\log(\pi_t(\vec{a}_t | \vec{s}_t)) = \log \dfrac{1}{\pi_t(\vec{a}_t | \vec{s}_t)}$

Dynamic Programming, Solving
for the policy backward through time

Rewrite Objective Function, Eqn (17)

$$\max_{\pi_0}\left\{E[r(\vec{s}_0,\vec{a}_0)] + \max_{\pi_1}\left\{E[\cdots] + \max_{\pi_T}E(r(\vec{s}_t,\vec{a}_t))\right\}\right\} \cdots \text{(18)}$$

$$\max_{\pi_0}\left\{E[r(\vec{s}_0,\vec{a}_0)] + \max_{\pi_1}\left\{E[r(\vec{s}_1,\vec{a}_1)] + \max_{\pi_2}\left\{E[r(\vec{s}_2,\vec{a}_2)]\right\}\right\}\right\} \quad \text{for } t=2$$

$$\max_{\pi_2}\left\{E[r(\vec{s}_2,\vec{a}_2)] + \max_{\pi_3}\left\{E[r(\vec{s}_3,\vec{a}_3)]\right\}\right\} \quad \text{for } t=3$$

$$\max_{\pi_3}\left\{E[r(\vec{s}_3,\vec{a}_3)] + \max\left\{E[r(\vec{s}_4,\vec{a}_4)]\right\}\right\} \quad \text{for } t=4$$

$$\ddots$$

Note :

1° Stochastic Gradient

$$\hat{\nabla}_\theta J_Q(\theta) \qquad \text{from Eqn (6)} \quad \text{Ref. SAC.}$$

2° Policy Gradient

$$\hat{\nabla}_\phi J_\pi(\phi) \qquad \text{from Eqn (10), Ref}$$

Ref. Fujimoto, 2018

21-Dimensional Humanoid

3° Single Q function → 2 Soft Q-function
Speed up training.

SAC ref.

Algorithm 1. (PP.8 )

Input: $\theta_1, \theta_2, \phi$     $\theta_{1,2}$ for Reward

$\bar{\theta}_1 \leftarrow \theta_1, \bar{\theta}_2 \leftarrow \theta_2$    $\phi$ for Policy $\pi$

$\mathcal{D} \leftarrow 0$       Replay Pool

for each iteration do

    for environment Step do

      $\vec{a}_t \sim \pi_\phi(\vec{a}_t | \vec{s}_t)$    Sample action from Policy

      $\vec{s}_{t+1} \sim \mathcal{P}(\vec{s}_{t+1} | \vec{s}_t, \vec{a}_t)$   Sample Transition from Enviro.

      $\mathcal{D} \leftarrow \mathcal{D} \cup \left\{ (\vec{s}_t, \vec{a}_t, r(\vec{s}_t, \vec{a}_t), \vec{s}_{t+1}) \right\}$   Store transition in Replay pool

    end for

    for each gradient Step do

      $\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i)$ for $i \in \{1,2\}$   update Q-func Parameter

      $\phi \leftarrow \phi - \lambda_\pi \hat{\nabla}_\phi J_\pi(\phi)$    update Policy weights

      $\alpha \leftarrow \alpha - \lambda \hat{\nabla}_\alpha J(\alpha)$    update temperature

      $\bar{\theta}_i \leftarrow \tau \theta_i + (1-\tau) \bar{\theta}_i$ for $i \in \{1,2\}$   Update network Weights

    end for

end for

Output : $\theta_1, \theta_2, \phi$

Note : 4 Dual Objectives  $\longrightarrow$ Dual Gradient Descent
                                   Alternating Between
                                   Lagrangian.

$$\alpha_t^* = \underset{\alpha_t}{\text{argmin}} \, E_{\vec{a}_t \sim \pi_t^*}[\cdot] \quad \cdots (17)$$
                                          Ref.

5. "Off-Policy" Data ?

April 2nd, 2021

Reward Functions Based On Difference

For Autonomous Driving

{ Micro-Scale: $\overrightarrow{P_p}$ to $\overrightarrow{P_q}$ within One Image $(X_w - Y_w - Z_w)$
  I(x,y)
medium Scale: Within One Room.
Global Scale: One Environment Multiple Rooms. }
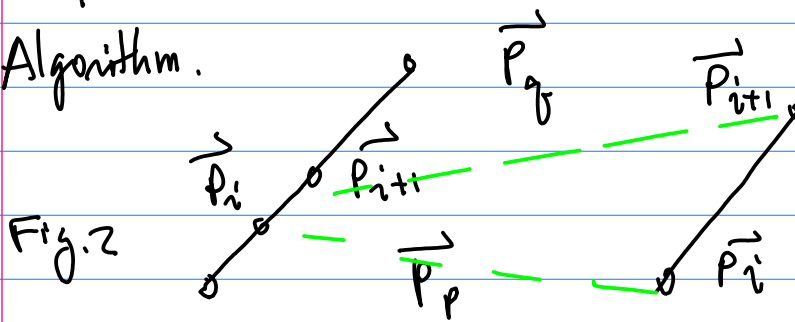
Micro-Scale

1. Within One Image $I(x,y)$



$\overrightarrow{P_q}(X_q, Y_q)$

$\overrightarrow{P_p}(X_p, Y_p)$  Fig.1

2. the Path from $\overrightarrow{P_p}$ to $\overrightarrow{P_q}$ is accomplished Step by Step Starting from $\overrightarrow{P_p}$, then Update each Step $\overrightarrow{P_i}$ to $\overrightarrow{P_{i+1}}$, ..., $i = 1, 2, ..., N$

$\overrightarrow{P_1} = \overrightarrow{P_p}$ (Starting Point)

$\overrightarrow{P_N} = \overrightarrow{P_q}$ (ending Point)

3. $\overrightarrow{P_i}$ for $i = 1, 2, ..., N$ is

Computed Based On DDA

Algorithm.



Fig.2

a $\overrightarrow{P_i}$ to $\overrightarrow{P_{i+1}}$ is one pixel equivalent distance in the world Coordinate $(X_w - Y_w - Z_w)$

b Actual AGV Travels during that Step i to i+1, e.g. time $t_i$ to $t_{i+1}$, is $\overrightarrow{P_{ACT}}$, or $\overrightarrow{P_a}(X_a, Y_a)$

c Case I    Actual Displacement within $\overrightarrow{P_{i+1}}$



Fig 3

Case II    Actual Displacement Outside $\overrightarrow{P_{i+1}}$



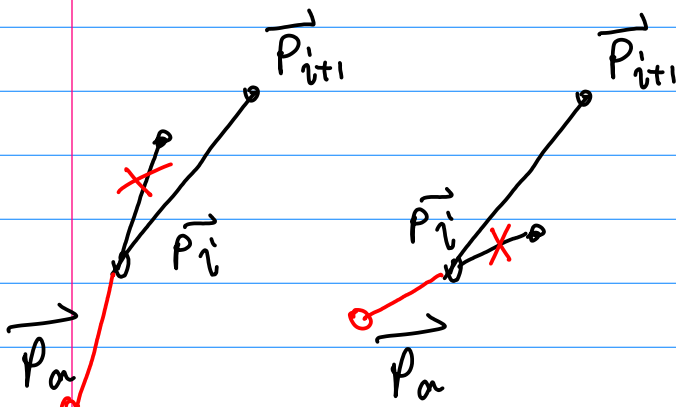Fig 4.

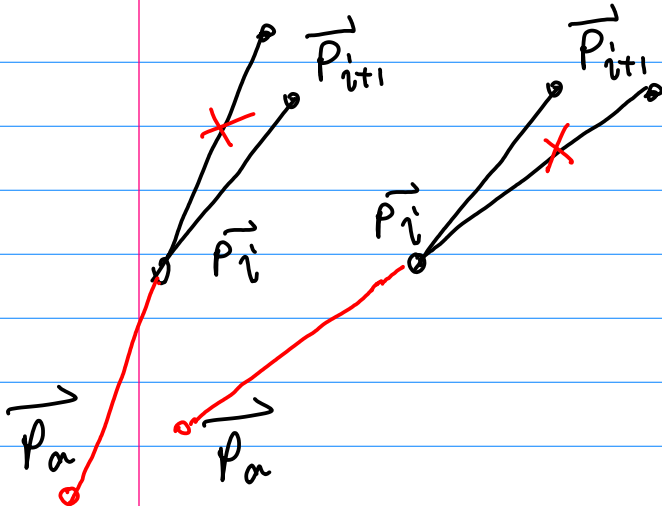Case III Opposite Directions as in Case I & II.

Fig 5a



Fig 5b

Actual Displacement Direction
Vector

$$\vec{d_a} \triangleq \vec{P_a} - \vec{P_i} = (x_a - x_i, y_a - y_i)$$
$$\cdots (3)$$

5. Define Deviation Index Based on
   $\begin{cases} \underline{a} & \text{Distance Difference} \\ \underline{b} & \text{Directional Difference.} \end{cases}$

Note:

$$\vec{d_d} \cdot \vec{d_a} = \| \vec{d_d} \| \| \vec{d_a} \| \cos \alpha \qquad \cdots (4)$$

Normalize it to make it an index

$$\frac{\vec{d_d} \cdot \vec{d_a}}{\| \vec{d_d} \| \| \vec{d_a} \|} = \cos \alpha \qquad \cdots (4-1)$$

4. Path Deviation $\vec{P_a}$ Calculation
   By Dot Product

$$\left( \vec{P_a} - \vec{P_i} \right) \cdot \left( \vec{P_{i+1}} - \vec{P_i} \right) \quad \cdots (1)$$

where
Desired Direction Vector

$$\vec{d_d} \triangleq \vec{P_{i+1}} - \vec{P_i} = (x_{i+1} - x_i, y_{i+1} - y_i)$$
$$\cdots (2)$$

Property 1:          Displacement
if $\vec{d_a} = \vec{d_d}$  and distance

$$\| \vec{P_a} - \vec{P_i} \| = \| \vec{P_{i+1}} - \vec{P_i} \|$$

then, the AGV is on the right
path direction and has the
exact desired displacement
distance to make $\vec{P_a} = \vec{P_{i+1}}$

Define

Definition 1. Direction Index $I_{dir}$, a Scalar is defined as the angle $\alpha$ Between $(\overrightarrow{P_a} - \overrightarrow{P_i})$ and $(\overrightarrow{P_{i+1}} - \overrightarrow{P_i})$;

e.g.

$$I_{dir} = \cos \alpha = \frac{\overrightarrow{d_a} \cdot \overrightarrow{d_d}}{\|\overrightarrow{d_a}\| \|\overrightarrow{d_d}\|} \quad \cdots (5)$$

Define Distance index $I_{dis}$, as Scalar as

$$I_{dis} = \begin{cases} \dfrac{\|\overrightarrow{P_a} - \overrightarrow{P_i}\|}{\|\overrightarrow{P_{i+1}} - \overrightarrow{P_i}\|}, & \cdots (6a) \\[4mm] \quad \text{if } \|\overrightarrow{P_a} - \overrightarrow{P_i}\| \le \|\overrightarrow{P_{i+1}} - \overrightarrow{P_i}\| \\[4mm] \dfrac{\|\overrightarrow{P_a} - \overrightarrow{P_i}\| - \|\overrightarrow{P_{i+1}} - \overrightarrow{P_i}\|}{\|\overrightarrow{P_{i+1}} - \overrightarrow{P_i}\|} \end{cases}$$

$$\text{if } \|\overrightarrow{P_a} - \overrightarrow{P_i}\| > \|\overrightarrow{P_{i+1}} - \overrightarrow{P_i}\| \quad \cdots (6b)$$

Note $I_{dis} \in [0,1]$; and $I_{dir} \in [0,1]$
$$(\cos \alpha)$$

Definition 2. Define Combo Index $I_\Sigma$ for AGV as

$$I_\Sigma = \beta I_{dir} + (1-\beta) I_{dis} \quad \cdots (7)$$

$$I_\Sigma = \beta I_{dir} + (1-\beta) I_{dis}$$
where $\beta \in [0,1]$
for Equal weights on direction and distance we have $\beta = \frac{1}{2}$. So
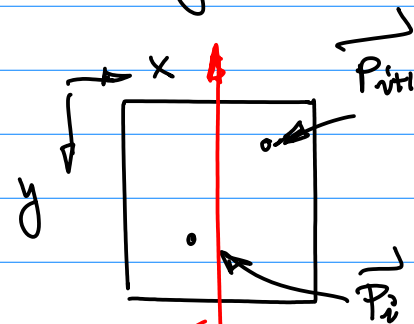
$$I_\Sigma = \frac{1}{2}\left(I_{dir} + I_{dis}\right) \quad \cdots (7a)$$

PART II. Consider Actuation.

Note 1. How to find $\overrightarrow{d_a}$?
magnetometer find $\overrightarrow{d_a}$

Note 2. How to find $\overrightarrow{d_d}$

See Eqn (2), Computer Vision
Technique plus $X_w - Y_w - Z_w$ mapping.



Vehicle Forward Looking Direction LSMSEN $\triangleq$ Physical World Caliberation in $X_w - Y_w - Z_w$

LSMSEN Sensor Input

$I(x,y)$
$m \times N$

i°. Angle $\theta = 132.615°$ for Example
(Theta), Vehicle forward direction angl

ii°. Red Line on $I(x,y)$ : Vehicle
forward Looking Direction

iii°. $\vec{P_i}$ Current time i Position

For $m \times N$ Image, $\vec{P_i}(x,y) = I_m\left(\frac{M}{2}, N-1\right)$

$\qquad \cdots (8)$

$\quad$ M : $\quad$ Col.
$\quad$ N : $\quad$ Row

iv°. mapping from $I(x,y)$ to $X_w - Y_w - Z_w$

$f : (x,y) \in \mathbb{I}^2 \longrightarrow (x_w, y_w, z_w) \in \mathbb{R}^3$
$\qquad \cdots (9)$

Definition 4. Define AGV forward
Looking Direction By LSMSEN Sensor
Input. e.g., Equal to the Angle
from LSMSEN.

Definition 5. Define AGV $\vec{P_i}$ Position
as the $(x, y) = \left(\frac{M}{2}, N-1\right)$, for
$m \times N$ $I(x,y)$.

Property 2. A feature Point $\vec{P_{i+1}}$
on Image $I(x,y)$'s location
$(x_{i+1}, y_{i+1})$ is mapped to the

$X_w - Y_w - Z_w$ By the mapping
function $f$ as follows

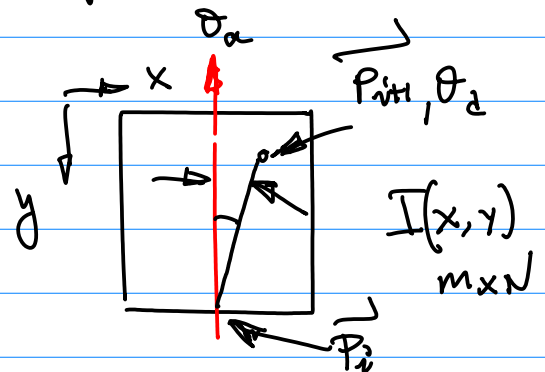$f : (x,y) \in \mathbb{I}^2 \longrightarrow (x_w, y_w, z_w) \in \mathbb{R}^3$

such that

$(x_{i+1}, y_{i+1}) \longrightarrow (x_{iw}, y_{iw}, z_i)\Big|_{\substack{w \\ z_i = C}}$

e.g. for Each $(x_{i+1}, y_{i+1})$
we have

$(x_{i+1}, y_{i+1}, C)_w$

5. Driving objectives:
To align the Red line with
$\vec{d_{dir}} = (\vec{P_{i+1}} - \vec{P_i})$

How much to Drive?

$\underline{\alpha}$ Angle $\theta_d - \theta_a = \Delta \theta$ to Be zero'd



$I(x,y)$
$m \times N$

distance to $\vec{P_{i+1}}$ to be zero'd.

| | map to $X_w$-$Y_w$-$Z_w$ | Result |
|---|---|---|
| $\theta_a$ | LSMSEN | |
| $\theta_d$ | $f(\cdot)$ | $X_{i+1,w}$ / $Y_{i+1,w}$ |

from $X_{i,w}, Y_{i,w}$ to $X_{i+1,w}, Y_{i+1,w}$

Left wheel $\begin{cases} \text{Displacement, dis} \\ \text{Angle} \end{cases}$

Right wheel $\begin{cases} \text{Displacement, dis} \\ \text{Angle} \end{cases}$

6 DoF Robot Arm Reward function.

1. 3 Types Movement
$\begin{cases} \text{Linear } \vec{P_i} \text{ to } \vec{P_{i+1}} \\ \text{Angular (Joint Angle)} \\ \text{Arc, } \vec{P_i} \text{ to } \vec{P_{i+1}} \text{ via Control pt.} \\ \vec{P_c} \end{cases}$