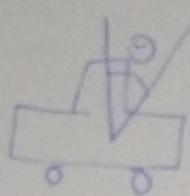


Lecture-18

- ④ Let's say you want to build a driving car.



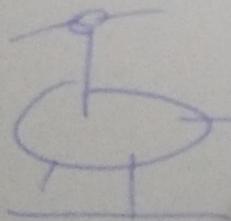
→ How do you model the state of the car:-

→ $x, y \rightarrow$ Position of the car
 $\rightarrow \theta \rightarrow$ Orientation of the car

\leftarrow
Six dimensions
state space

→ $\dot{x}, \dot{y}, \dot{\theta} \rightarrow$ Velocity in different direction.

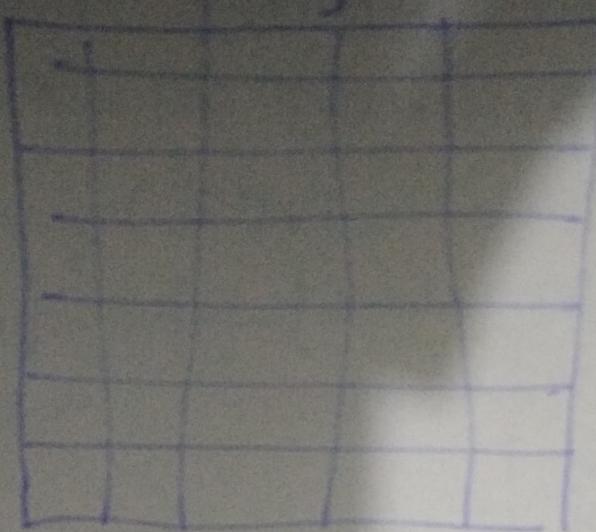
- ④ Let's say you want to build a helicopter.



→ How do you model the state of the helicopter?

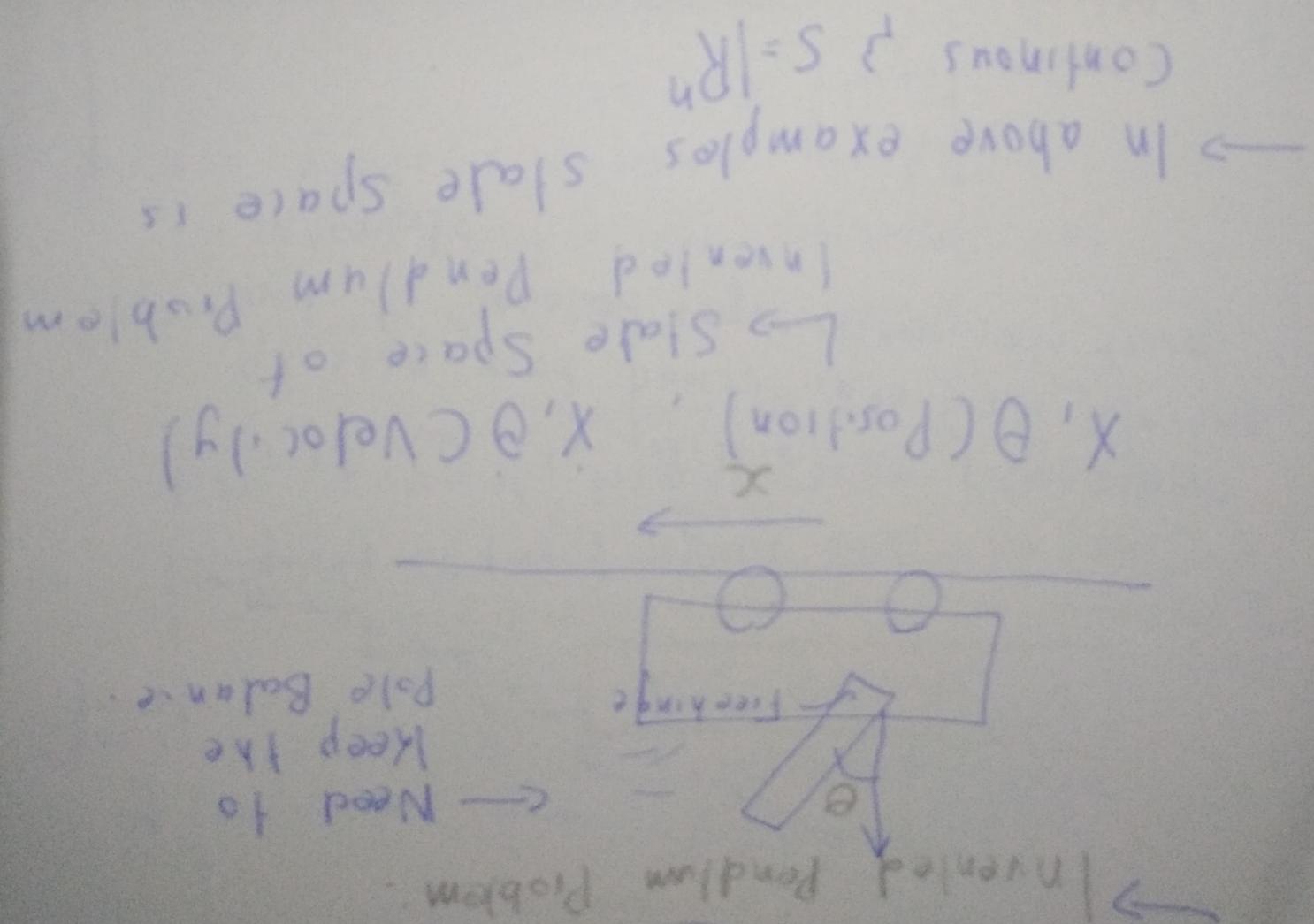
→ $x, y, z, \beta, \theta, \psi \leftarrow$ Position
 \downarrow Roll Pitch Yaw
 \leftarrow Position

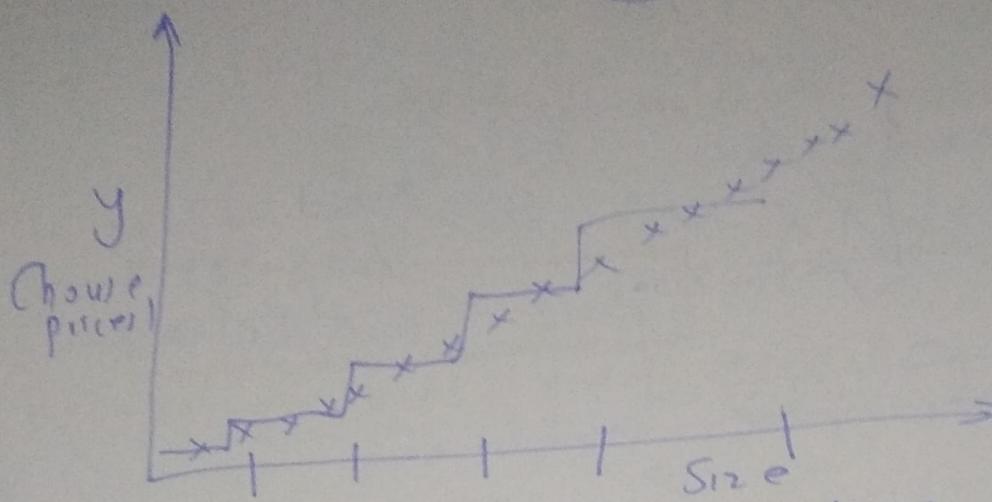
b pen direction
 slide square has
 common
 space direction
 Two stable
 dimensions



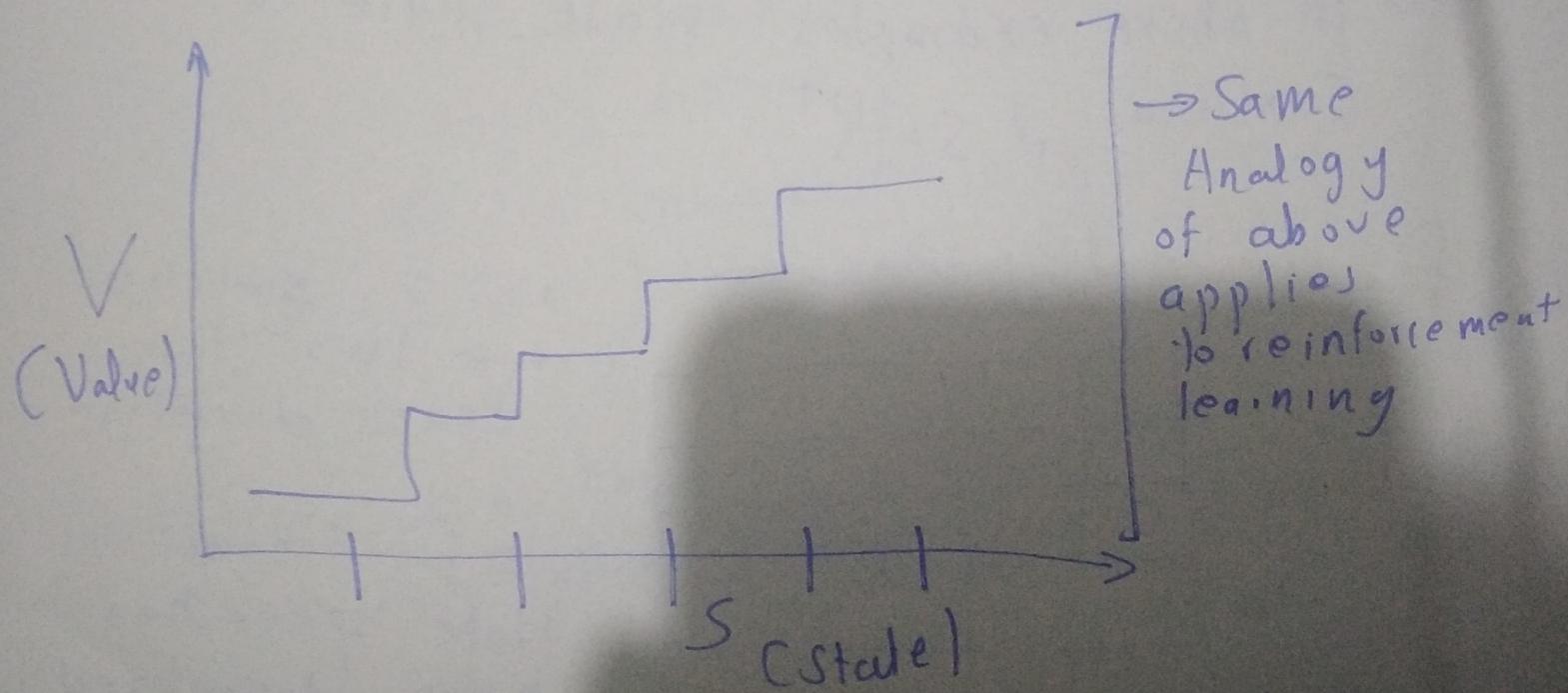
5²

④ Discrete solution.





↳ If we to discretize regression of size and house into five interval and fit separate regression into five intervals.



↳ State space is discretized to map the value function.

↳ Function is not smooth and not a good representation.

- Curse of dimensionality:

④ $S = \mathbb{R}^n$ and discretize each dimension into k values, get k^n discrete states.

→ Example of curse of dimensionality:-

④ One of area where people have applied reinforcement learning is factory optimization:

④ So we have a factory with 100 machines in a factory.

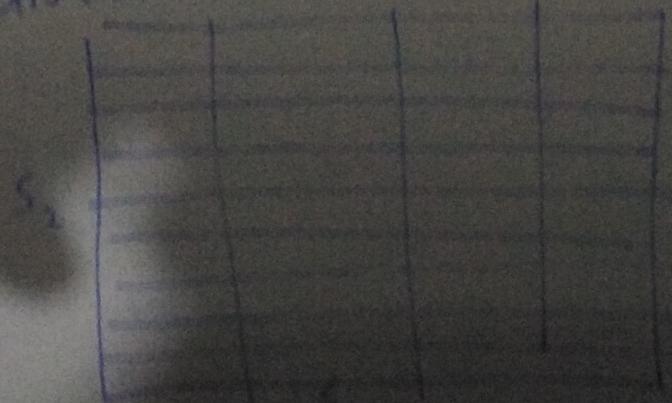
④ Every machine is doing something slightly different

④ 100 machines in k states

④ Then total number of states is k^{100}

④ Discretization however is effective in small state space (\mathbb{R}^3 state space, 3 dimension state space)

④ As you reach four to six dimension choose your discretization more carefully.



→ Small interval for because we believe the state is more important than time

- Seven to eight dimensions are not suitable for discretization
- Approximate V^* directly without resorting to discretization
 - In linear regression you say:

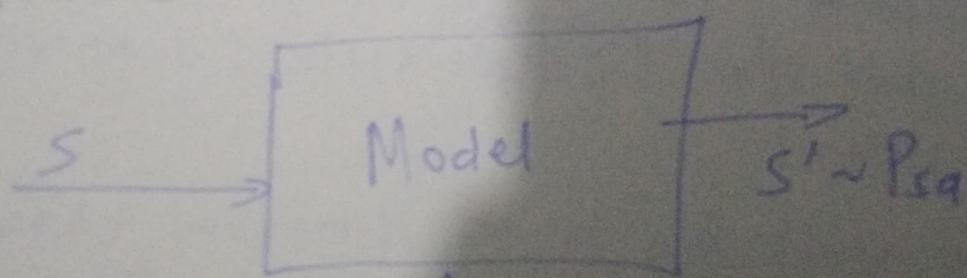
$$X \approx \Theta^T y \quad y = \Theta^T X \text{ or } \Theta^T \phi(x)$$

$\phi(x)$ features of x . $\phi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_{12} \\ x_1 x_2 \\ \vdots \end{bmatrix}$

$$\rightarrow V(s) \approx \Theta^T \phi(s) \quad (\text{fitted value})$$

↳ where $\phi(s)$ is features of state

- Model(Simulation) of MDP.

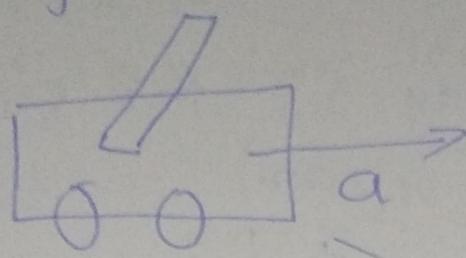


→ State is a q row value vector
 → For now let's assume action space is discrete

→ Usually state space is greater than action space

* How to get a model?

① - Physics simulation



$$\dot{s} = \begin{pmatrix} x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{pmatrix}$$

$$s = \begin{pmatrix} x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{pmatrix}$$

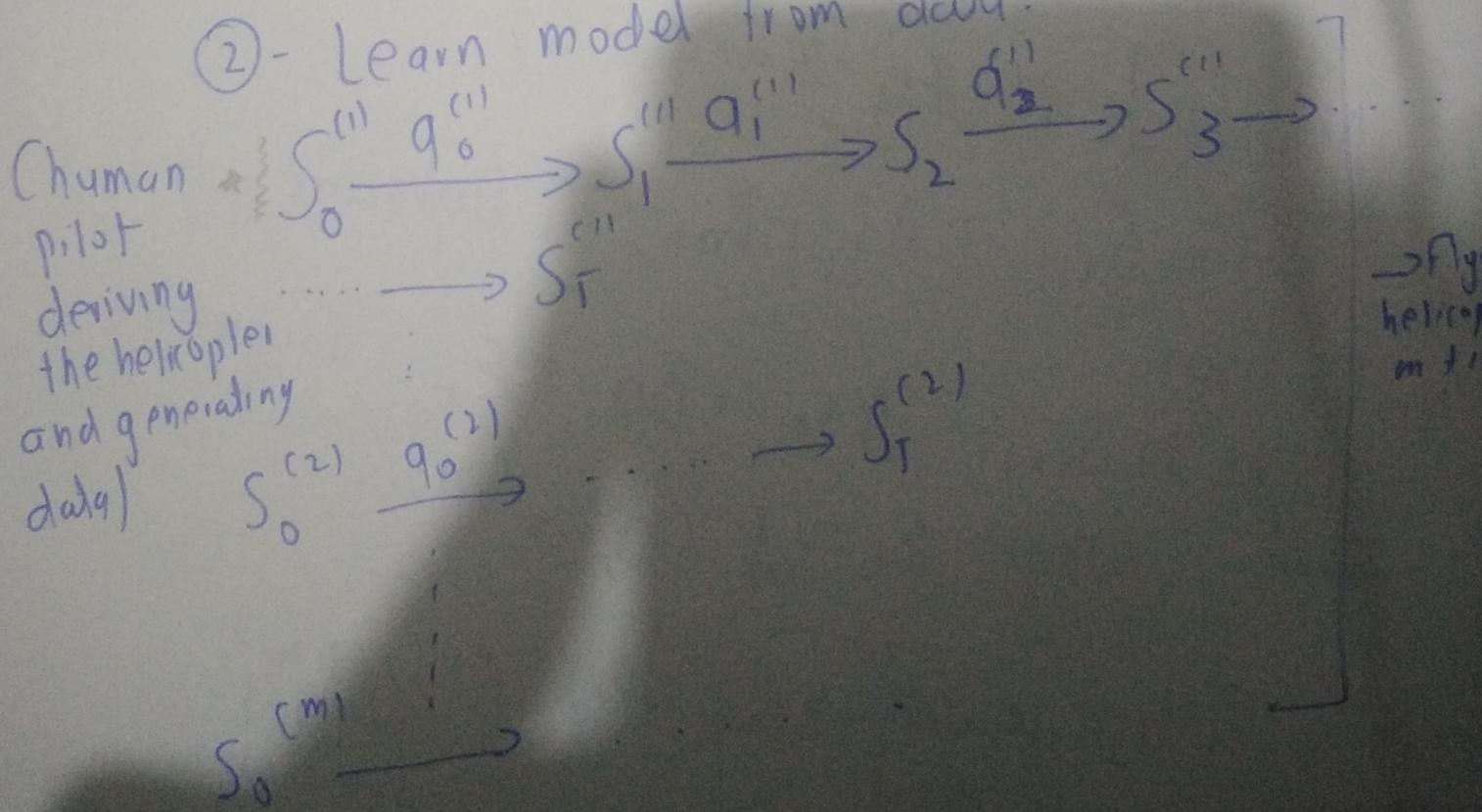
→ Conventional Physics law.

$$s' = s + (\Delta t) \cdot \dot{s}$$

$\Delta t = 0.1 \text{ seconds}$

→ Simulation equations through physics

② - Learn model from data:-

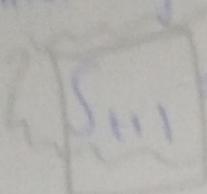


- Apply supervised learning to estimate

$\rightarrow S_{t+1}$ as function S_t, a_t

↳ Predicting S_{t+1} from S_t, a_t

Example regression version:



$$S_{t+1} = AS_t + Ba_t$$

A linear

Not a terrible model

You can easily make it nonlinear.

$$\min_{A, B} \sum_{i=1}^m ((S_{t+1}^{(i)}) - (AS_t^{(i)} + Ba_t^{(i)}))^2$$

$$\min_{A, B} \sum_{i=1}^m \sum_{j=0}^n \|S_{t+1}^{(i)} - (AS_t^{(i)} + Ba_t^{(i)})\|^2$$

Model B and RL

Model (Deterministic)

a) $S_{t+1} = AS_t + Ba_t$

(Deterministic)

or b) $S_{t+1} = AS_t + Ba_t + \epsilon_t$ (Stochastic)

You can choose either one of them.

where $\epsilon_t \sim N(0, \sigma^2)$

$S_t \rightarrow \boxed{\text{Model}} \rightarrow S_{t+1} \sim p_{SA}$ (Stochastic Model Simulation)

This diagram is introducing probability

Once you have build a model,
how do you approximate the
value of function:

- Fitted Value Iteration -

→ Choose features $\phi(s)$ of state

s

$$V(s) = \theta^T(s)$$

$$\phi(s) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ x_4 \\ \vdots \end{bmatrix} \rightarrow \text{Inverted Pendulum}$$

→ Designing features that conveys how well the robot is doing (Very important task!)

→ Feature Algo

→ fitted Value iteration Main Algorithm:

→ Previously (For discrete states)

① Value iteration:

$$\textcircled{1} \quad V(s) = R(s) + \gamma \max_a \sum_{s'} P_{sa}(s') V(s')$$

$$= R(s) + \gamma \max_a E_{s' \sim P_{sa}} [V(s')].$$

Going to
learn mapping
from
 $s \rightarrow V(s)$

② Let's generalize to fitted
value iteration:

③ Sample $\{s^{(1)}, s^{(2)}, \dots, s^{(m)}\}$ from S randomly

Initialize $\Theta := 0$

Algorithm
continues
on next page.

• Repeat {
 For $i = 1, \dots, m$ {
 For each action $a \in A$ {
 Sample $s'_1, s'_2, \dots, s'_K \sim P_{s^{(i)}|a}$ Using model
 Set $q(a) = \frac{1}{K} \sum_{j=1}^K [R(s'_j) + \gamma V(s'_j)]$

$$\text{Set } q(a) = \frac{1}{K} \sum_{j=1}^K [R(s'_j) + \gamma V(s'_j)].$$

$$R(s) + \sum_{s' \sim P_{sa}} [V(s')]$$

$$\text{Set } v^{(n)} = \max_a q(a)$$

→ Original Value Iteration :-

$$\{V(s)\} = y^{(1)}$$

→ Not part of the algorithm : Explanation

→ Fitted Value iteration:

Want $V(s^{(i)}) \approx y^{(i)}$

1

$$\mathbf{O}^T \phi(s^{(i)})$$

~~Post~~
~~post~~
of the
algorithm
→ ~~Def~~
concept

~~one explanation~~

→ Final Step : of Algorithm

π^* because

⑥ For fitted Value Iteration:

(211)

→ Then how do you choose $\pi^*(s)$ given a state:-

Say model is:-

$$S_{t+1} = f(S_t, a_t) + \epsilon_t$$

$$\text{e.g. } (AS_t + BS_t) + \epsilon_t$$

→ For your deployment (runtime,

④ Set $\epsilon_t = 0$ and $k=1$

→ When in state s ,

pick action

$$\pi^*(s) = \arg \max_q V(f(s, q))$$

✓ Simulator without noise.

$s' \sim P_{sa}$ but with deterministic model