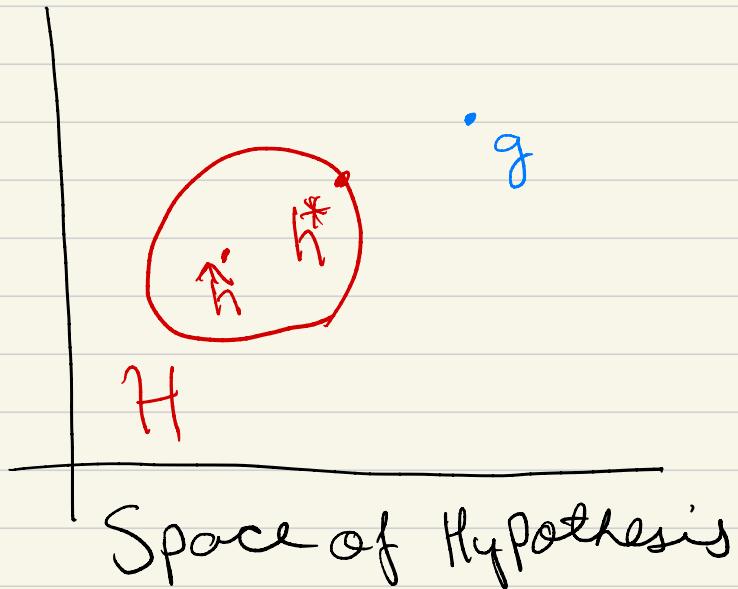


4.7 Approximation & Estimation Error



g : Best Possible Hypothesis

H : Set of all hypothesis of a particular type
Eg: Set of all SVMs

h^* : Best possible hypothesis of
that class: Eg: Most optimal SVM

\hat{h} : Learned from finite Data

$E(h)$ = Risk / Generalization Error

$$= E_{(x,y) \sim D} [1 \{ h(x) \neq y \}]$$

$E_s(h)$: Empirical Risk

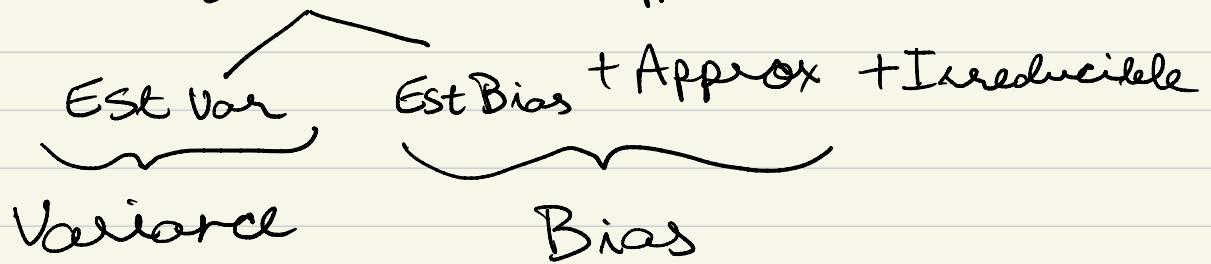
$$= \frac{1}{m} \sum_{i=1}^m 1 \{ h(x^{(i)}) \neq y^{(i)} \}$$

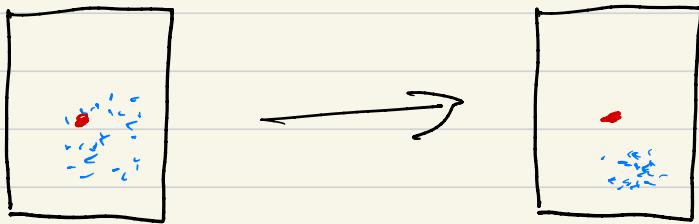
$E(g)$ = Bayes Error / Irreducible (unavoidable)

$E(h^*) - E(g)$ = Approximation Error (Class)

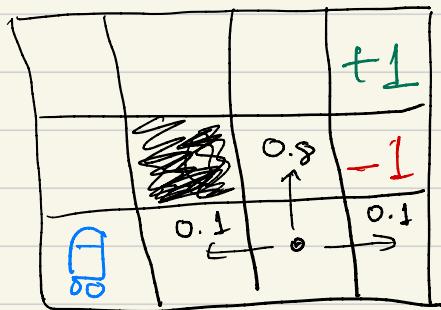
$E(h^*) - E(h)$ = Estimation Error (Data)

Total Error = Estimation + Approximation + Irreducible





— — —
RL Intro



Total Payoff :

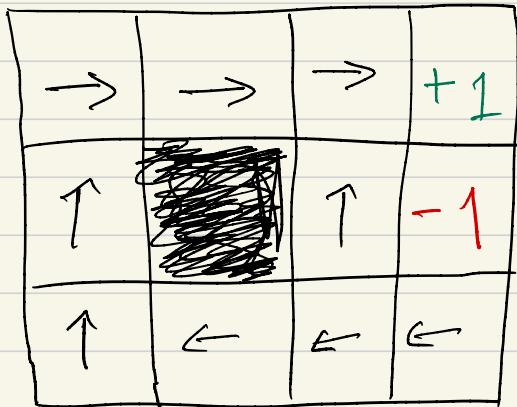
$$R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2) + \dots$$

$$\gamma < 1 \text{ (usually 0.99)}$$

Goal of RL is to maximise

$$E[R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2) + \dots]$$

Optimal Policy



Define: V^π, V^*, π^*

$$V^\pi(s) = E[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots | \pi, s_0=s]$$

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{\pi(s)}(s') V^\pi(s')$$

Immediate Reward
 ↑
 Computes Expectation
 ↑
 Transition Probability
 ↑
 Recursive Value Function
 ↓

$$V^*(s) = \max V^\pi(s)$$

Expected future reward

$$V^*(s) = R(s) + \max_a \underbrace{\gamma \sum_{s'} P_{sa} V^*(s')}_{\substack{\uparrow \\ \text{if we take action } a}}$$

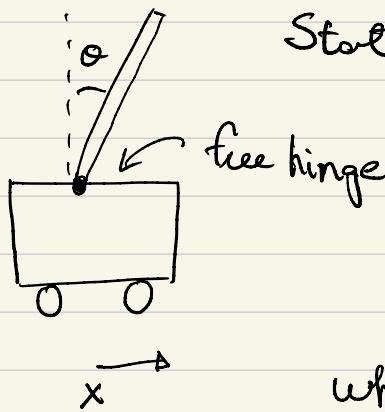
$$\pi^*(s) = \operatorname{argmax}_{a \in A} \sum_{s'} P_{sa} V^*(s')$$

if we take action a

What if I don't know P_{sa} ?

$$P_{sa}(S') = \frac{\# \text{ times I took action } a \text{ in } S \text{ and got } S'}{\# \text{ times I took action } a \text{ in } S}$$

(or $\frac{1}{|S|}$ if above is 0)



$$\text{State Space} = [x, \theta, \dot{x}, \dot{\theta}]$$

$s \in \mathbb{R}^N$
where N is dimensionality
of state space

In Linear Regression, we try to
approximate

$$y \approx \theta^T \phi(x)$$

↑ ↑
Parameters (Model) State
Features

$$V^*(s) \approx \theta^T \phi(s)$$

$$S_0 \xrightarrow{a_0^{(1)}} S_1 \xrightarrow{a_1^{(1)}} S_2 \xrightarrow{a_2^{(1)}} S_3 \xrightarrow{a_3^{(1)}} \dots \xrightarrow{a_T^{(1)}} S_T$$

$$S_0 \xrightarrow{a_0^{(2)}} S_1 \xrightarrow{a_1^{(2)}} S_2 \xrightarrow{a_2^{(2)}} S_3 \xrightarrow{a_3^{(2)}} \dots \xrightarrow{a_T^{(2)}} S_T$$

$$S_0 \xrightarrow{a_0^{(3)}} S_1 \xrightarrow{a_1^{(3)}} S_2 \xrightarrow{a_2^{(3)}} S_3 \xrightarrow{a_3^{(3)}} \dots \xrightarrow{a_T^{(3)}} S_T$$

$$\vdots \quad \vdots \quad \vdots$$
$$S_0 \xrightarrow{a_0^{(P)}} S_1 \xrightarrow{a_1^{(P)}} S_2 \xrightarrow{a_2^{(P)}} S_3 \xrightarrow{a_3^{(P)}} \dots \xrightarrow{a_T^{(P)}} S_T$$

S_{t+1} as function of S_t , a_t
(S')

Eg : Linear Regression version

$$S_{t+1} = A S_t + B a_t$$
$$\underbrace{A}_{\mathbb{R}^{N \times N}} \quad \underbrace{B}_{\mathbb{R}^{N \times K}}$$

(tall skinny matrix)

Objective :

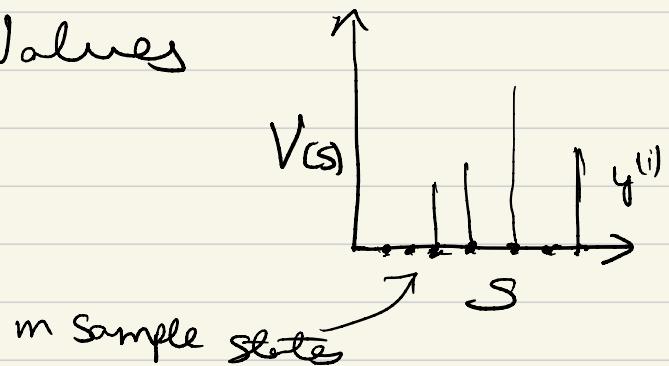
$$\min_{A, B} \sum_{i=0}^P \sum_{t=0}^T \| S_{t+1}^{(i)} - (A S_t^{(i)} + B a_t^{(i)}) \|$$

$$S_{t+1} = AS_t + BS_t + \epsilon_t \quad (\text{Stochastic Model})$$

↑
noise

Learn a mapping from $S \rightarrow V(S)$

States $\xrightarrow{\text{to}}$ Values



1. Randomly sample n states $s^{(1)}, s^{(2)}, \dots, s^{(n)} \in S$.

2. Initialize $\theta := 0$.

3. Repeat {

For $i = 1, \dots, n$ {

For each action $a \in A$ {

Sample $s'_1, \dots, s'_k \sim P_{s^{(i)}a}$ (using a model of the MDP).

Set $q(a) = \frac{1}{k} \sum_{j=1}^k R(s^{(i)}) + \gamma \underbrace{V(s'_j)}_{\text{old estimate of } V = \theta^T \phi(s)}$

// Hence, $q(a)$ is an estimate of $R(s^{(i)}) + \gamma \mathbb{E}_{s' \sim P_{s^{(i)}a}} [V(s')]$.

}

Set $y^{(i)} = \max_a q(a)$.

// Hence, $y^{(i)}$ is an estimate of $R(s^{(i)}) + \gamma \max_a \mathbb{E}_{s' \sim P_{s^{(i)}a}} [V(s')]$.

}

// In the original value iteration algorithm (over discrete states)

// we updated the value function according to $V(s^{(i)}) := y^{(i)}$.

// In this algorithm, we want $V(s^{(i)}) \approx y^{(i)}$, which we'll achieve

// using supervised learning (linear regression).

Set $\theta := \arg \min_{\theta} \frac{1}{2} \sum_{i=1}^n (\theta^T \phi(s^{(i)}) - y^{(i)})^2$

}

↑
update V by learning new mapping $s \rightarrow V(s)$

(Can also be a neural network or any supervised learning algo)

$$V_t^*(s) = \max_a R(s, a) + \sum_{s'} P_{sa}(s') V_{t+1}^*(s')$$

$$V_T^*(s) = \max_a R(s, a)$$

LOR:

Applies to an MDP: $(S, A, \{P_{sa}\}, T, R)$

$$S_{t+1} = A s_t + B a_t + w_t$$

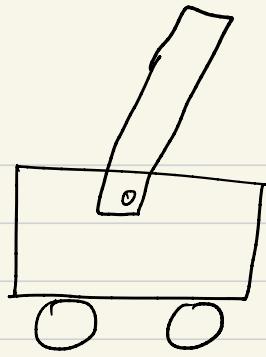
$\begin{matrix} R^{n \times n} & R^{n \times d} \\ \hookdownarrow & \hookdownarrow \\ S_t: \mathbb{R}^n & A: \mathbb{R}^d \end{matrix}$
 Matrices
 Linear System

$$R(s, a) = -(s^T U s + a^T V a)$$

Reward where $U \in \mathbb{R}^{n \times n}$, $V \in \mathbb{R}^{d \times d}$

$U, V \geq 0$ (positive semi-definite)

Quadratic cost function



$$\alpha \approx \pi_\theta(s) = \frac{1}{1 + e^{-\theta^T s}}$$

Let us assume that we train a logistic

$$\pi_\theta(s, "R") = \frac{1}{1 + e^{-\theta^T s}}$$

$$\pi_\theta(s, "L") = 1 - \frac{1}{1 + e^{\theta^T s}}$$

$$s = \begin{pmatrix} 1 \\ x \\ \dot{x} \\ \phi \\ \dot{\phi} \end{pmatrix}$$

Suppose Policy $\theta = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

$$\pi_\theta(s, "R") = \frac{1}{1 + e^{-\phi}}$$

Intuitively, policy aggressively pushes the cart to the right as the pole angle ϕ increases

Goal : Find θ so that when we execute $\Pi_\theta(s, a)$ we maximize

$$\max_{\theta} E[R(s_0, a_0) + \dots + R(s_T, a_T) | \Pi_\theta]$$

Observation $y = \begin{pmatrix} x \\ \phi \end{pmatrix} + \text{noise}$

Only a partial state is measured.

Cannot use $V^*(s)$ because s is not fully observable.

Policy search attempts to learn $\Pi^*(y)$

Assume $T=1$

$$E[R(s_0, a_0) + R(s_1, a_1) | \Pi_\theta]$$

$$= \sum_{s_0, a_0, s_1, a_1} P(s_0) \Pi_\theta(s_0, a_0) P_{s_0, a_0}(s_1) \Pi_\theta(s_1, a_1) \cdot [\text{payoff}]$$

Reinforce Algorithm

Loop : {

Run your MDP (using current policy)

Sample $s_0, a_0, s_1, a_1, \dots, s_T, a_T$

Compute Payoff = $R(s_0) + R(s_1) + \dots + R(s_T)$

Update

$$\theta := \theta + \alpha \left[\frac{\nabla_{\theta} \pi_{\theta}(s_0, a_0)}{\pi_{\theta}(s_0, a_0)} + \frac{\nabla_{\theta} \pi_{\theta}(s_1, a_1)}{\pi_{\theta}(s_1, a_1)} + \dots \right] \cdot \text{Payoff}$$

$$\pi(s,a) = \frac{e^{\theta^T s}[a]}{\sum_{a'} e^{\theta^T s}[a']}$$

indexing start action
with a certain index

$$\frac{e^{\theta^T s}}{\sum_{a'} e^{\theta^T s}}$$

} (num_states, num_actions)

$$\nabla_\theta \ln (\pi_\theta(s_t, a_t))$$

$$\nabla_\theta \ln \left(\frac{e^{\theta^T s_t}[a_t]}{\sum_{a'} e^{\theta^T s_t}[a'_t]} \right)$$