

Genetic Algo

Population based Algo.
Extremely slow convergence

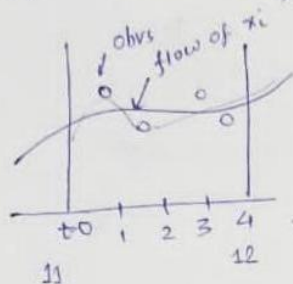
* Derivative-free methods

(a) 4D - Data Assimilation problem:

$$\Psi(x_0) = \|x_0 - x_0^b\|_{B^{-1}}^2 + \sum_{i=1}^n \|H_i(x_i) - y_i\|_{R^{-1}}^2 \equiv \begin{pmatrix} \text{initial value} \\ \text{PDE} \end{pmatrix}$$

obj fn at $x_i = M_{i,i-1}(x_{i-1})$ - forward model
 x_0 = initial value

y_i = observation



can read
obs at all
pts might
be possible

M = governing eqn. (phys sys)
(phys space)
 H = mapping from obs. loc to
entire domain. (obs space)

x_0 = decision variable (to be determined)
 \Rightarrow Reverse engg problem to determine x_0 & match
it with obs value. based on final values

$$\Psi(x_0) = \|x_0 - x_0^b\|_{B^{-1}}^2 + \sum_{i=1}^n \|H_i(x_i) - y_i\|_{R^{-1}}^2 + \sum_{i=1}^n \lambda_i (x_i - M_{i,i-1}(x_{i-1}))$$

Applying SGD, we need $\nabla \Psi(x_0)$

$$\nabla \Psi(x_0) = B^{-1}x_0 + \left[\frac{\partial H_i(x_i)}{\partial x_0} \right] x_{i-1} \dots \frac{\partial x_i}{\partial x_0}$$

to be done for $i=1$ to n

This to be done for each iter.

Hence computational cost is \uparrow in cal. the derivatives
Even storing the derivatives is expensive as
it will take a lot of storage (2000 x 2000 grid pts)

(b) Derivative can't be cal.
Ex in LASSO regression

$$L + \sum_{i=1}^n \lambda_i |O_i| \rightarrow \text{can't be differentiated}$$

(c) No. of sep hyperparameter tuning.
can't be learnt from data, done by grid search
method. No fn relates the hyperparameter to
loss fn. ~~No~~ No direct reln

Derivative free Methods \rightarrow local search - Gen. pattern search

\rightarrow Global opti - Evolutionary Algo - Genetic Algo.

local search
4 all algo are learnt tell now

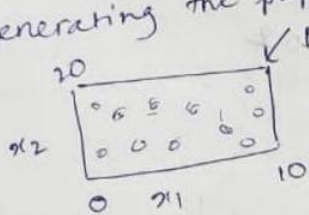
global optima due to stochasticity involved.

Evol. Algo

1. Generating the population
 2. Recombination
 3. New population for next iteration
- } Depends on method

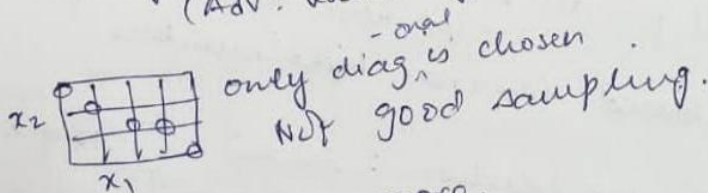
Genetic Algo.

- (1) Generating the population (candidate soln).



Domain
randomly generate samples from domain
samples = $4 \times \# \text{variables (gen)}$

Any random dist method (Normal / Uniform)
to generate samples.
They should span the entire set. well.
(Adv: Latin Hypercube sampling).



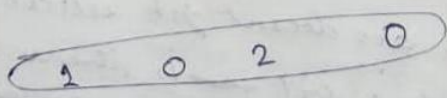
only diag. is chosen
Not good sampling.

- (2) Evaluate fitness score.

$$\min f(x) = f(x_1, \dots, x_4) \quad (\text{unconstrained})$$

suppose we have 8 samples

Fitness Score (FS)	x_1	x_2	x_3	x_4	fit
1	1	3	2	5	
2	2	1	9	10	
3					
4					
5					
6					
7					
8					



candidate soln. } Chromosome
Suppose $f(x) = x_1^2 + x_2^2 + x_3^2 + x_4^2$
Each sample gives a soln
soln = 0, 0, 0, 0

Fitness score may be any other fn of the fit.
lowest score is the soln. (3) selection

~~(2) selection~~

Encoding in (2):

$$\min f(x_1, x_2)$$

(a) bit encoding — good for small int, higher int
higher numbers of bits

candidate som: 1 3 5 2
↓ ↓ ↓ ↓
0001 0011 0101 0010

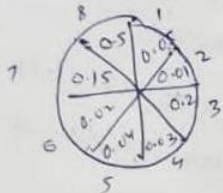
random cs
cs:

1011	1000	0000	1110
------	------	------	------

we have fitness score for each cs.

Roulette wheel selection /
Tournament

Roulette wheel: Fitness score in each sectors



Arrange →

0.01
0.02
0.03
0.04
0.05
0.15
0.5

say Random no. = 0.6.

Choose the cs where the cumsum is crossed the Random no. 0.8

Mate in pool

all parents are kept here

pop size is fixed.

This is biased towards higher numbers.

Tournament selection: for min for max
40-50% of the pop. 70-80%
hyperparameter

Select randomly & make smallest mating pool.

↓ do this 8 times

Hence largest one doesn't get selected

Repetition in mating pool is there
If we do 70-80%, then lowest value will be repeated many no. of times
it will be there in most of the 70-80% selections

(4) Crossover: Generate offspring from mating pool.
Choose any two ^{diff} parents from mating pool.
↓
constrained. to generate new popⁿ.

* 1 point / 2 point crossover / N point / uniform crossover

4 point !

1 point
 a 1011 10000
 c 0111 1001
 0000 1110 → 0.2 randomly → say 1
 1100 1111 → 0.2 randomly → say 1
 0.3 fs → FS may increase be careful

α 1011 1000
 α 0111 1001

$\times 0.3$ FS \rightarrow FS may increase
 be careful

offs: 1011, 1001, 1001, 11 $\rightarrow 0.3$ fs \rightarrow FS may be careful
half of parents may join length stochastic involved

Any half of population
2 point: Choose 2 random length
N " " "

stochasticity
involved
↓
goes to
global opt.

Uniform Crossover:

Generate random mask

101	0011	1001	1000
1000	0100	1111	

$\therefore \begin{cases} 1 \rightarrow \text{from 1st parent} \\ 0 \rightarrow \text{2nd "} \end{cases}$

→ 1011 1000 0100 1111

000 0100 1111
Helps in generating offspring w/ Φ FS
 \hookleftarrow parents FS

How to know if children doing better than

(5) parents ?
Do mutilation → avoids trapping in local minima.

Randomly change children string

offs: \rightarrow 1011 1000 1000 1110

Mark: 1000 0000 0001 0100
0001 1000 1001 1010
(8) & after

FS for offs after crossover & after mutation

we have 24 CS \rightarrow $\left. \begin{array}{l} 8 \text{ Curr pop} + \\ 8 \text{ crossover} + \\ 8 \text{ Mut} \end{array} \right\} 24 \text{ CS}$

8 Mutⁿ.
choose lowest 8 FS probⁿ $\leq 20\%$
as New popⁿ (6) Replacement: New
genⁿ for next itⁿ

with prob of 1:
($\frac{3}{16}$)
= 20%

of entire
length is
changed

usually
"mutation
prob" $\leq 20\%$.

Since all the f_1 eval are independent tasks,
hence high no. of f_1 eval is not expensive
as it can be parallelized.

Repeat 3-6 until converged.

Convergence criteria:-

- No. of itr
 - No. of itr w/o improvement
 - ~~Tolerance~~ Tolerance on obj fn (if we know the min f_1 value).
- ↓
or tolerance of obj fn from
1 step to another.

Population should be diversified at each itr.
If New Genⁿ has FS b/w 0.01 and 0.02, it is not
diversified and no improvement of soln may occur.
So we may choose percent of popⁿ from top, middle &
bottom FS

no guarantee that crossover will produce offspring which
performs better than both parents.

so we have a strategy: Elitist strategy

Elite Ratio $\sim 5\%$: preserve the best candidate from
current popⁿ & rest from crossover +
mutation
(for new Genⁿ)

Crossover Probability: How often a chromosome is being changed.

⊕ { 1000 : Crossover prob = 50% as 2 chromosomes are changed
2101 : " " = 25% " 1 " " "
= 2100

we define this probability for mating, $P_c = 0.8 = 80\%$ (say)

⊗ $\begin{matrix} \text{dom} \\ 1001 & 1100 & 1010 & 1111 \\ & & & 1001 \end{matrix}$: $\begin{matrix} 0.3 & 0.55 \\ 0 & 0 \end{matrix}$

⊗ 0110 0000 1101

Per Parent: (taking each as dominant parent once)
For each bit, generate a random no. \leftarrow take from dom parent
If random no. $> P_c$: don't change the bit } So we get 2 diff children from 2 diff parent
If random no. $< P_c$: don't take bit from other parent (non dominant)

We had fixed in this strategy.

But we see in case of ① 1000 \rightarrow Less change / P_c
① 2101 \rightarrow More change / P_c

So we change P_c adaptively.

say we put $P_{c1} = 0.9$, $P_{c2} = 0.1$, $P_{c1} + P_{c2} \neq 1$.

FS for ① : f_1
FS " ② : f_2
 $f_{avg} = \text{Avg FS of Curr. Pop}$
 $f_{min} = \min \{ \text{FS of curr pop} \}$

$P_c = \begin{cases} \frac{(P_{c1} - P_{c2})(f_{min} - f_{avg})}{f_{min} - f_{avg}}, & f < f_{avg} \\ P_{c1} = \text{highest prob.}, & \text{otherwise} \end{cases} \Rightarrow \text{any } f_n \text{ can be defined.}$
 \downarrow
then go for prev. stat with this

P_{c2} = How much we need to change the candidate soln

$P_c = 1 \rightarrow$ No change, $P_c = 0 \rightarrow$ Full change.

If the 2 solns from crossover are close to the parents then the f_n is not changing much in each it? & it gets stuck in the local minima. So we go for mutation (as less as possible mut? prob). This will also create diversified pop?

In Derivative based method, converge to the soln from initial pt.
" Genetic Algo " , we jump discretely anywhere.
Convergence is hypothesized \rightarrow "micro Schema Theorem, simulated annealing, macro mutation hypothesis"

Rate of convergence : empirically

\rightarrow based on population being generated.

\rightarrow should span entire search space : diversified.

⊗ If only 1 child is desired, take only 1 dominant parent (one w/ less FS)

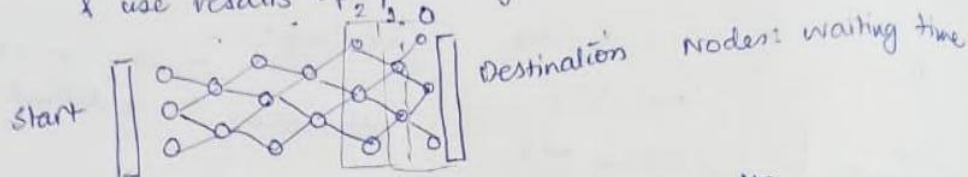
If 2 child are desired, take both as dominant by turn.

Dynamic Programming

Road Network:

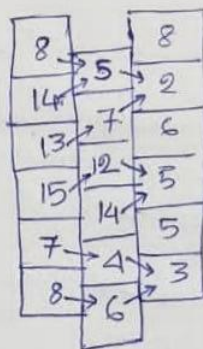
Brute force method — infeasible for large no. of layers & nodes.

So recursive way is used to solve layer wise & use results of prev. layers / subproblems.



	3	2	1	0 ← stage No.
5		3		8
		9	3	2
		6	5	6
		6	7	5
2		3	9	5
1		3	1	3
0		2	3	
↑				
state No.				

Stage 0:



Stage 0: All states have 1 soln.

Stage 1: solve for min^m taking lower sum for stage (1+0)

$$v_n(s_n) = \min \{ t_n(s_n) + v_{n-1}(s_{n-1}) \}$$

(cost to go)

n = Stage No.

s_n = State No. of stage n

$$s_{n-1} = \begin{cases} s_n + 1, & \text{choose up if } n \text{ is even} \\ s_n - 1, & \text{" down if } n \text{ is odd} \\ s_n, & \text{o.w.} \end{cases}$$

State transition

t_n = 1 step time duration at stage n } for each state.

V_n = Σ time till stage n

Decision Making : minimise Σ cost OR
maximise Σ rewards at each stage.
and minimize the risk / uncertainty at each stage.

Markov Prop: Future depends on present, not past.

MDP: < Markov Decision Process >

$\langle S, A, P, R \rangle$ S = Finite set of states

state trans prob.

(π, u, f, g, J)

A = " " " actions

$$p^a(s, s') = P[S_{t+1} = s' | S_t = s, A_t = a]$$

$S = \{1, 2, \dots, 9\}$

R = Reward fn

$$= R^a(s) = E[R_{t+1} | S_t = s, A_t = a]$$

(scalar)

$A = \{u, l, d, r\}$

$$R(s, a, s') = E[R(s, a, s')]$$

Terminal state

Deterministic

Stochastic

$$P(2|5, u) = 1$$

$$P(2|5, u) = 0.5$$

$$P(1|5, u) = 0$$

$$P(3|5, u) = 0.5$$

supposing wind blows towards rt w.p 0.5

1	2	3
4	5	6
7	8	9

9, 5, 4A

$$R(s, a) = \left. \begin{aligned} R(2, l) &= -1 \\ R(2, d) &= -1 \\ R(2, u) &= -2 \end{aligned} \right\}$$

G_t (Return) = Sum of future rewards

$$G_t = R_{t+1} + R_{t+2} + R_{t+3} + \dots$$

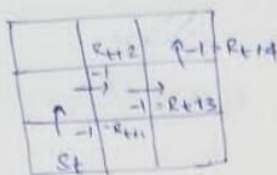
$$= R_{t+1} + G_{t+1}$$

Policy, π : $\pi(a|s) = P[A_t = a | S_t = s]$
 $\pi(s) = a$

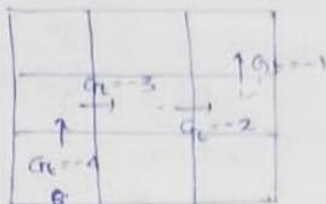
T = Terminal state

fn of state you're in, prob dist of actions at each state
 $s_t \rightarrow (A) \rightarrow \text{deterministic}$
(stochastic) ~~example~~

In Policy grid given, all actions are given
: Deterministic policy



Imm rewards



Returns

Value functions :

State-value fn, $v_{\pi}(s) = E[G_t | S_t = s]$

Action " fn, $q_{\pi}(s, a) = E[G_t | S_t = s, A_t = a]$

-4	-3	0
-3	-2	-1
-4	-3	-2

v_{π}

~~200501~~

$$q_{\pi}(G_t | 7, r) = -4$$

$$q_{\pi}(G_t | 7, a) = -4$$

followed
rt & then
follows
policy

followed

$$v_{\pi}(s) = \max_a q_{\pi}(s, a)$$

For deterministic,

$v_{\pi}(s) = G_t$ at
each grid

$$q_{\pi}(7, l) = -6 \quad \text{left wall : (stay & get -2)}$$

$$q_{\pi}(7, r) = -4$$

Bellman Expectation Eqn

In deterministic env,

$$v_{\pi}(s) = R^a(s) + v_{\pi}(s')$$

~~200501~~

$$v_{\pi}(4) = (-1) + (-4) = -5$$

$$v_{\pi}(1) = (-1) + v_{\pi}(4)$$

$$v_{\pi}(2) = (-1) + v_{\pi}(5)$$

$$v_{\pi}(4) = -1 + v_{\pi}(5)$$

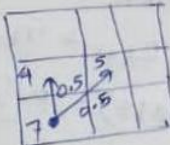
} 9 eqn, 9 variables

In stoch env,

$$v_{\pi}(s) = R^a(s) + \sum_{s' \in S} p(s, s') v_{\pi}(s')$$

$$v_{\pi}(7) = (-1) + [v_{\pi}(4) \times 0.5 + v_{\pi}(5) \times 0.5]$$

9 eqn, 9 var



$$V_{\pi}^{k+1}(s) = R^{\pi}(s) + \sum_{s' \in S} P^{\pi}(s, s') V_{\pi}^k(s') \quad (*)$$

start with $V_{\pi}^{(0)}(s')$ as say $\{0, 0, \dots, 0\}$
 → Solve w/ iterative scheme

→ Policy Evaluation.

→ There always exists a soln which is unique given some mild conditions.

Finding the Optimal Policy:

Control Problem

Given an MDP $\langle S, A, P, R \rangle$, And opt. pol. π^*

$$\pi^* \text{ s.t., } V_{\pi^*}(s) \geq V_{\pi}(s) \quad \forall s \in S \quad \forall \pi \text{ where } \pi \neq \pi^*$$

In given policies, $V_{\pi_2}(s) \geq V_{\pi_1}(s) \quad \forall s$

$\pi_1 \rightarrow$ Sub optimal

$\pi_2, \pi_3 \rightarrow$ optimal

* Multiple optimal policies (maybe),

But optimal fn value is same for all

Opti Problem : — {best / optimal policy}
 $\max_{\pi} V_{\pi}(s) / \left. \begin{aligned} V_{\pi^*}(s) &= \max_{\pi} V_{\pi}(s) \\ V^*(s) &= \max_{\pi} V_{\pi}(s) \end{aligned} \right\}$

rewards \uparrow
 if cost, then min.

Bellman Optimality Eqn.

$$V^*(s) = \max_x E[R(s,a) + V^*(s')]$$

$$V_n(s_n) = \min \{ t_n(s_n) + V_{n-1}(s_{n-1}) \} \quad \left\{ \begin{array}{l} \text{converted} \\ \text{to update} \\ \text{rule} \end{array} \right.$$

$$(*) V_{\pi}^{k+1}(s) \leftarrow \max_{a \in A} \left\{ R^a(s) + \sum_{s' \in S} P^{\pi}(s, s') V_{\pi}^k(s') \right\}$$

Value Iteration Algo, $k \rightarrow$ stages
 Spl case of
 GPI in
 RL

Initialize

0	0	0
0	0	0
0	0	0
0	0	0

-1	-1	-1
-1	-1	-1
-1	-1	-1
-1	-1	-1

V_1

$$V^*(s) = \max \{ -1+0, -1+0, \dots \}$$

$$= -1$$

$R(a,s) = -1, V(s') = 0$
 state I land up in

-1	-2	-2
-1	-2	-2
-1	-2	-2
-1	-2	-2

-1	-2	-2
-1	-2	-2
-1	-2	-2
-1	-2	-2

In each stage, one diag entries get converged.

Based on domain knowledge, we can prune out some nodes which don't get updated further.

DP methods are model based;

Model of world has to be known a priori

$P^a(s, s')$ has to be known a priori

This is offline planning.

What if model is not known?

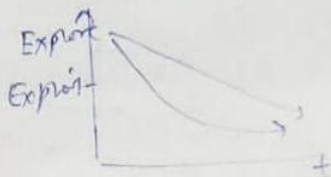
→ Reinforcement Learning: learning to decide through interactions w/ world.

Learning from experience

Regret: not knowing a priori & hence missing out on exploiting

Tradeoff: Explore & Exploit

initially do this lot (so that we get a good model of the world)



W/o Exploration, we may learn the suboptimal policy.

Note " , more better / optimal solution.

Prediction & Control Problem: (Model is not known)

Policy Evaluation
Find value for given policy

Find optimal policy
Gen policy itr.
sample mean

Prediction: $q_{\pi}(s, a) = E[G_t | s_t = s, A_t = a]$

$$V_{\pi}^{(k+1)}(s) = R^a(s) + \sum P^a(s, s') V_{\pi}^{(k)}(s')$$

not known

Find: Action value for by Estimating Expectation

↓	↓	×
→	→	↑
↑	↑	↑

$$G^{(1)} = -4$$

$$G^{(2)} = -4$$

$$G^{(3)} = 4$$

$$\begin{aligned}
 \bar{G}(\text{start}, \uparrow) &= (-4) + (-4) + (-4) \\
 &= \frac{-12}{3} \\
 &= -4
 \end{aligned}$$

Take sample mean of return

thru episodes.

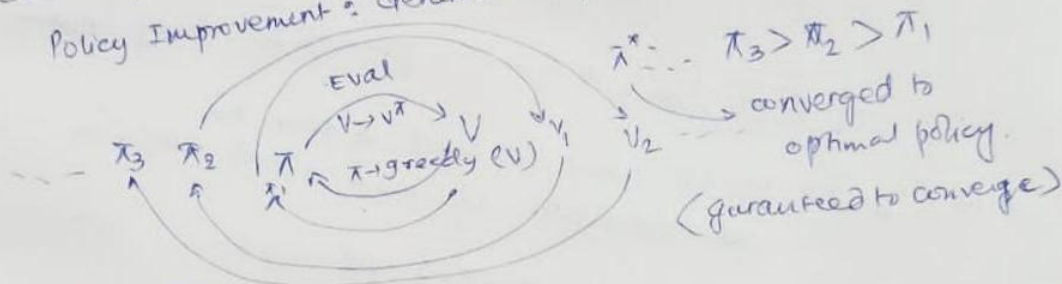
We know value fn of all states we visited

i.e. of $\{2, 3, 4, 5, 6, 7, 8\}$

Control: Gen Pol. Iter (start w/ random policy)

Calc: $v_{\pi}(s)$ for π (GPI)

Policy Improvement: Generate improved π' from $v_{\pi}(s)$



Greedy Pol. Imp: $\pi'(s) = \arg \max_a q_{\pi}(s, a) \forall s$

↓	↓	×
→	→	↑
↑	↑	↑

π_1

-4	-3	×
-3	-2	-1
-4	-3	-2

$V\pi_1$

-4	→	×
-3	-2	-1
-4	-3	-2

$V\pi_2$

↓	→	×
→	→	↑
↑	↑	↑

π_2

$$\begin{aligned}
 q_{\pi_1}(2, \downarrow) &= -3 \\
 q_{\pi_1}(2, \rightarrow) &= -1 \\
 q_{\pi_1}(2, \leftarrow) &= -5
 \end{aligned}$$

Epsilon Greedy Pol. Imp: $\epsilon \in (0, 1)$

$$\pi'(s) = \arg \max_a q_{\pi}(s, a) \text{ w.p. } 1-\epsilon$$

w.p. ϵ any other actions

$$\text{say } \pi'(a|a) = \begin{cases} \downarrow, \leftarrow, \rightarrow, \uparrow & \text{w.p. } \epsilon/4 \end{cases}$$

Thousands of interaction w/ π_1 to get the expectation of all actions to get $V\pi_1$

Theoretically $\rightarrow \infty$ times

Practically \rightarrow large no. of times

Monte Carlo Control

$$M_n = \frac{x_1 + \dots + x_n}{n} \Rightarrow M_n = M_{n-1} + \frac{1}{n} (x_n - M_{n-1})$$

$$Q_n = \frac{G_1 + \dots + G_n}{n} \Rightarrow Q_n = Q_{n-1} + \frac{1}{n} [G_n - Q_{n-1}]$$

Action value
 Q_n

return

$$Q_\pi(s, a)$$

for particular (s, a)

$$Q_\pi = E[G_t | S_t, A_t]$$

Policy Eval.

$\alpha_t = \frac{1}{n}$: for same state, we are updating the Q_n

$\epsilon \leftarrow 1/K, \pi \leftarrow \epsilon - \text{greedy}(Q)$ Policy Improve

$N(S_t, A_t)$ = NO. of times we have done the interaction w/ same state & action

$\epsilon \rightarrow$ more explore & start off w/ ϵ

$\epsilon \rightarrow \frac{1}{K}$ Exploration \uparrow $K \uparrow$ $\frac{1}{K} \downarrow$ $\epsilon \downarrow$

has more knowledge we already have about the world, so we chance of explore \downarrow as $K \uparrow$

$$\epsilon \propto (s, a)$$

\rightarrow Model free RL

* Model based RL: build the model within & then plan.

upper time limit set only T_i 's are shown, no V_i 's shown

$Q \rightarrow DNN$

Monte Carlo \rightarrow learning \rightarrow DQN Deep Q learning apply NN

moment of the bar is the policy

Two min paper.