

Computational Intelligence Notes

Nature Inspired Computing \rightarrow inspired by nature

↳ to simulate / emulate nature

↳ bring properties to our systems such as:

↳ robustness, fault tolerance, self-repair.

Combinatorial Optimization \rightarrow finding optimal object from a set of finite objects \rightarrow Exhaustive search not tractable.

Evolutionary Computing (EC)

\rightarrow Evolution \rightarrow change in inherited traits of a population of organisms through successive generations.

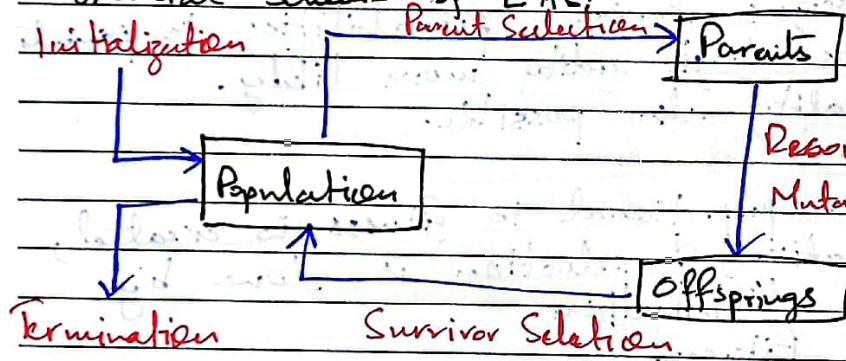
\rightarrow EC is the study of comp sys. inspired by natural evolution.

↳ an EA is a generic population-based meta-heuristic optimization algorithm.

↳ uses mechanisms such as:

reproduction \leftrightarrow Mutation \rightarrow Recombination
 \downarrow Selection

General Schema of EAs



Variation & Selection \rightarrow Reduces diversity

↳ some alleles survive

↳ Builds up genetic Diversity

↳ produces new alleles.

Exploitation vs. Exploration

↳ a good alg must find a good tradeoff b/w the two.

↳ Exploration investigates new & unknown areas in the search space

↳ Exploitation makes use of knowledge found at points previously visited to find better points.

Typical EA Cycle:

(2)

- 1 → Initialize random samples.
 - 2 → While no convergence (or terminating condition).
 - (a) → evaluate fitness
 - (b) → Perform crossover
 - (i) Select two parents
 - (ii) Produce offsprings
 - (c) perform mutation (Select one individual & mutate it)
 - (d) select the new generations
 - (e) evolve the next generation.

→ Crossover (merging two parents) is explorative since it makes a jump b/w two parent areas.

→ Mutation is exploitative since it creates small diversities thereby staying near the parent.

* Only crossover can combine info

*only mutation can introduce new info.

Selection Schemes

Parent Selection → Assigns variable probabilities to parents based on fitness → high quality more likely.
↳ bad quality also possible.

1. Fitness Proportional Scheme:

b A distribution proportional to fitness is created, by normalization; & selection is done by sampling the distribution.

b Selection & Fitness.

b) High selective pressure \rightarrow strong individuals dominate.

2. Rank Based Solution

- ↳ Rank populations acc. to fitness & base selection properties on rank: where worst has 1, best has μ .
- ↳ Selection is the first step.

b Selection index: of actual fitness

↳ best want dominate

↳ preserves constant selection pressure

3. Tournament Selection

- ↳ If population v. large, might be too time consuming for FP or P, or if universal fitness not def.
- ↳ In TS, pick n_{ts} members at random & then select the best of those.
- ↳ inherits adv. of rank.
- ↳ For crossover, TS done twice (one for each parent).
- ↳ If n_{ts} not too large, prevents the best from dominating → lower selective pressure.
- ↳ If n_{ts} too small, bad chances increase.
- ↳ Pressure \propto n_{ts}. (at n_{ts} = 1 \rightarrow random, n_{ts} = s \rightarrow best does)

4. Random Selection

- ↳ simplest sol. ($P(x) = 1/n$ for each indiv.)

Replacement:

Age-Based \rightarrow remove certain indivs that have not been replaced in a set time.

Fitness-Based \rightarrow On fitness.

↳ Replace Worst / Truncation.

↳ worst 1 members are replaced.

↳ Rapid improvements in the mean

↳ can lead to premature convergences.

↳ Elitism \rightarrow ensuring that the best survive.

↳ best copied without mutation.

↳ less diversity

↳ ~~consistently~~ guarantees convergence (risk of local minima)

↳ Hall of Fame \rightarrow best individual of each generation is selected to be inserted into H.o.F.

↳ parent pool for crossover.

Variants of EA

EC vs Classical Optimization [CO]

- * CO uses deterministic rules to move from one point to another. EC uses prob.
- * COs more successful in linear, quadratic, strongly convex, unimodal & other specialized problems.
- ECs more efficient in discontinuous & multimodal pr_{gs}.
- * EC starts from diverse set of points,
- CO starts from one point.
- * CO uses derivative info, ECs use fitness value.

Biology

Swarm Intelligence

Ant Colony Optimization (ACO).

Swarm Behaviours in Animals:

- food foraging
- flocking / schooling
- clustering
- Natural Construction

Key Principles - social insects work without supervision

- individuals largely homogeneous
- act asynchronously in parallel
- self-organization
- flexibility → colony can adapt to changing environment
- Robustness → group can perform even if an individual fails.
- Communication → stigmergy.

Stigmergy: * Communication via signs or ~~dead~~ cues placed by one entity which affects the behaviour of other entities who encounter them

* Stigmergy is the indirect communication that seems to underpin cooperation among social insects.

↳ form of self-organization → produces complex, seemingly intelligent structures without planning or control b/w agents.

↳ supports efficient collaboration

→ ants lay down pheromones on their way back.

When they find food.

Termites example → deposit mudballs, positive feedback

① Ant Colony Optimization (ACO)

is model food foraging behaviour of ants

② Particle Swarm Optimization (PSO)

↳ flocking & fish schooling

with the (fish) scales are best for protection not catching

most birds will respond to threat to their nest

① Ant Colony Optimization (ACO):

- ↳ While initially random, once a food source is located, activity patterns become more organized.
- ↳ This behaviour emerges from individuals following 2 simple rules:
 - lay pheromones
 - follow the trails of others.

- * Working of ACO: Two ants leave at the same time, taking separate paths to the food source, marking their trail with pheromone.
- * The ant taking the shorter path returns first, thus this trail is worked by twice as much pheromones.
- * Since shorter path has more concentrations of pheromones, next ant is attracted to that.
- * As more ants take that route, they further lay their pheromones, further amplifying their attractiveness.
- * As time passes, pheromones also evaporate, so longer paths vanish.
- (i) Also prevents convergence to a sub optimal solution.
- * Without evaporation, exploration is restricted.

* ACO Algorithm:

* Let $B_i(t)$ be number of ants in town i at time t ,

* Let $m = \sum B_i(t)$ the total number of ants.

* Each ant is a simple agent with the following characteristics:

↳ chooses the town to go to with a probability that is a function of the town distance & of the amount of trail present on the connecting edge.

↳ To form legal tours, transitions to already visited towns are disallowed until a tour is completed.

↳ When a trail is completed, a "trail" is laid on each edge (i, j) visited.

* Let $\tau_{ij}(t)$ be intensity of trail on edge (i, j) at time t .

* Each ant at time t chooses the next town $(t+1)$.

- * In n iterations, each ant completes a tour.
- * Trail intensity updated:

$$\tau_{ij}(t+n) = \rho \cdot \tau_{ij}(t) + \Delta \tau_{ij}$$

$\therefore \rho$ is a coefficient such that $(1-\rho)$ is the evaporation of trail between time t & $t+n$.

$$\therefore \Delta \tau_{ij} = \sum_{k=1}^m \Delta \tau_{ij}^k$$

$$\Delta \tau_{ij}^k = \begin{cases} \alpha/l_k & \text{if } k\text{-th ant uses edge } (i,j) \\ 0 & \text{in its tour.} \\ 0 & \text{otherwise.} \end{cases}$$

α is constant, l_k is tour length of k -th ant.

- * Tabu List Data Structure associated with each ant.
- ↳ saves already visited tours upto time ' t '

- * ACO has an advantage over evolutionary algorithms when the graph may change dynamically.

↳ Can adapt in real time

↳ Interest in network routing & urban transportation systems.

ACO has been applied to various Combinatorial Optimization Problems (COPs)

- Assignment Problem
- Closure Problem
- Constraint Satisfaction Problem
- Cutting Stock Problem
- Integer Programming
- Knapsack Problem
- Minimum Spanning Tree
- Nurse Scheduling Problem
- Vehicle Routing Problem

Particle Systems & Particle Swarm Optimisation (PSO)

- PSO → inspired by bird flocking & fish schooling
- population based stochastic optimization technique
- no evolution operators
- initially a random solution, searches for optima by updating generations.

* 1 Particle → current position
↳ velocity

Velocity Update: → current velocity }
 { personal best [PB]
 { global best [GB] } influence velocity update

- * At each step, particles are displaced from their current position by a vector gradient
- Magnitude & direction of their velocity at each step is influenced by their velocity in the previous iteration, simulated momentum, location of the particle relative to the PB & GB.
- * Particle ~~update~~ updates its velocity & position.

$$v[t] = \varphi v[t] + c_1 * rand() * (pbest[t] - present[t]) \\ + c_2 * rand() * (gbest[t] - present[t])$$

$$\text{present} = \text{present}[t] + v[t]$$

c_1 & c_2 are learning factors (usually 2)

φ is inertia weight.

* PSO: → Swarm, a set of particles (S)

→ Particle, a potential solution

• position: $x_i \in (x_{i,1}, x_{i,2}, \dots, x_{i,n}) \in \mathbb{R}^n$

• Velocity: $v_i \in (v_{i,1}, v_{i,2}, \dots, v_{i,n}) \in \mathbb{R}^n$

→ each particle maintains its PB

→ Swarm maintains its GB

PSO Algorithm: 1 → Randomly initialize the swarm to form the solution space.

2 → evaluate the fitness of each particle

3 → Update individual (p_{best}) & global best.

4 → update velocity & position of each particle.

5 → Repeat from step 2 until termination condition.

* Initially p_{best} is the starting location & g_{best} is the best out of all particles (best p_{best} for all particles).

* In each generation, a particle is accelerated towards its previous best position & towards a global best position, forcing particles to continually search in the most-promising regions.

Date: _____

Artificial Neural Networks (ANNs)

- Neural Networks modelled after our understanding of how the brain works.
- Because of massive number of neurons & interconnections, how the brain works remains a mystery.

Neural Network → Made up of highly connected network of individual computing elements (mimicking neurons)

→ Collectively used to solve interesting & difficult problems.

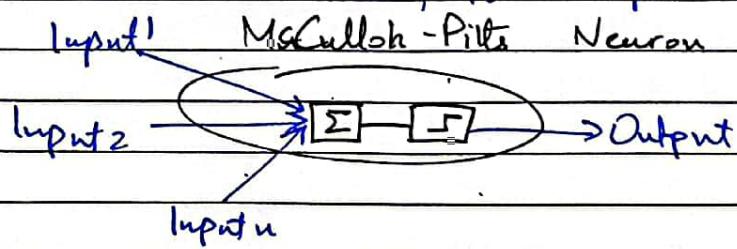
Neurons: responsive cells in the nervous system that process & transmit information by chemical signals within the neuron.

- number of different neurons exist.
- respond to stimulus & communicate the presence of that stimuli to the central nervous system which processes that info & sends responses to the other parts of the body

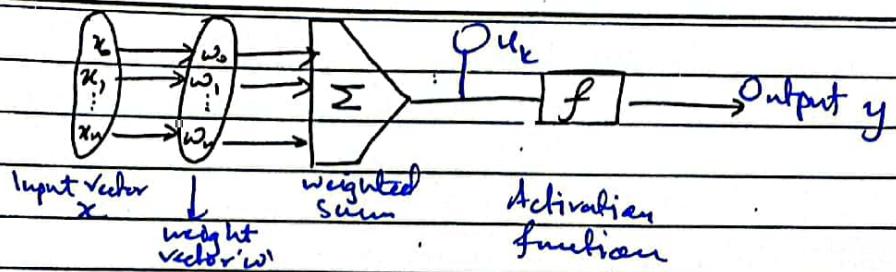
Artificial Neural Network:

- highly connected network of individual computing elements (mimicking neurons) that collectively can be used to solve interesting & difficult problems.
- An "artificial neuron" is an information processing unit that is fundamental to the operation of an ANN.
- Artificial Neuron → Set of input links

↳ to Output
 ↳ An adder for \sum input signals weighted by respective synaptic strengths.
 ↳ Activation function for limiting the amplitude of the output of a neuron.



Date: _____



- * The ' n ' dimensional input vector ' x ' is mapped onto variable ' y ' by means of the scalar product & a nonlinear function mapping.

Multi-Layer Feed Forward Networks

- one of most imp & popular classes of ANNs
- referred to as multilayer perceptrons.

Multi Layer Perceptron (MLP)

- ① Network contains one or more layer of hidden neurons not part of input. These hidden nodes enable the network to learn complex & nonlinear tasks by extracting progressively more meaningful features from input patterns.
- ② Model of each neuron usually includes a ~~at~~ non-linear activation function, sigmoid or hyperbolic.
- ③ High degree of connectivity from one layer to the next.

How a multi-layer neural network works?

- Inputs correspond to attributes for each training type
- Inputs fed simultaneously to the input layer
- Weighted & fed simultaneously to (a) hidden layer(s).
- Weighted outputs are inputs to the output layer which emits the networks prediction.
- Non linear Regression is performed

Specification of an ANN:

- Number of input attributes found within the individual instances determines the number of input layer nodes.
- The user specifies the number of hidden layers as well as the number of nodes within a specified hidden layer.

* Input Format: → Numeric $\in [0, 1]$

- Appropriate way to numerically represent categorical data
 - Attribute Color: {R, G, B, Y}

→ Conversion method for numerical data falling outside the [0, 1] range $\in (100, 200, 420)$.

* Output Format: [0, 1].

* Activation Functions:

⇒ Sigmoid Function

(1) → Must output $\in [0, 1]$

↳ Criterion 1

(2) → Value output should be close to 1 when function is sufficiently excited.

Sigmoid meets the above 2 criteria!

$$f(x) = \frac{1}{1 + e^{-x}} = \text{Sigmoid}(x)$$

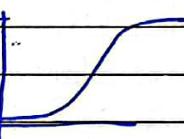
$$\tanh\left(\frac{x}{2}\right) = \frac{1 - e^{-x}}{1 + e^{-x}}$$

learning a Neural Network

* Learning is accomplished by modifying network connection weights while a set of input instances is repeatedly passed through the network

* Once trained, an unknown instance is passed through the model/network & classified according to the values at the output layer.

Working of an ANN



↓
Assign initial random weights

Select random test from tot set

↓
Propagate inputs forward

↓
Back Propag. Error Backwards

↓
Reached pre defined MSE

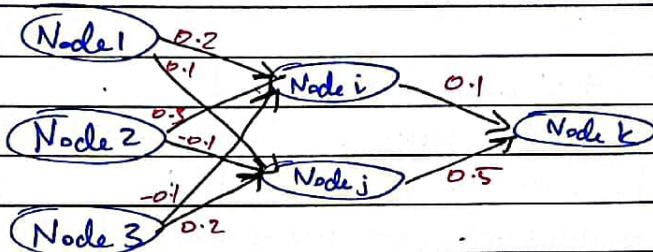
↓
Done

Feed-Forwarding Data:

* Input = {1.0, 0.4, 0.7}

* Input to node $i = 0.2(1.0) + 0.3(0.4) - 0.1(0.7)$

* Now apply the sigmoid function: $= 0.562$

Gradient Descent * Alg that minimizes function

* Starts with an initial set of ^{param} values & iteratively moves towards a set of parameter values that minimize the function.

* Achieved using Calculus \rightarrow taking steps in the negative direction of the function gradient.

Gradient: \rightarrow grad of each weight gives an indication on how to modify the weight to get the expected output / reduce error.

\rightarrow Each weight has gradient \rightarrow slope of the error function

$\hookrightarrow 0 \Rightarrow$ weight not contributing to error

\hookrightarrow -ve \Rightarrow weight should be increased for lower error

\hookrightarrow +ve \Rightarrow weight should be decreased for lower error.

Explanation of Back Propagation Algorithm

$$w_{1i} = 0.2, w_{1j} = 0.1, w_{2i} = 0.3, w_{2j} = -0.1, w_{3i} = -0.1, w_{3j} = 0.2$$

$$w_{ik} = 0.1, w_{jk} = 0.5, T = 0.65$$

$$* \text{Input} = \{1.0, 0.4, 0.7\}$$

$$* \text{Input to } N_i = 0.2(1.0) + 0.3(0.4) - 0.1(0.7) = 0.25$$

$$* \text{Apply Sigmoid: } f(0.25) = \underline{\underline{0.562}}$$

$$* \text{Input to } N_j = 1.0(0.1) + 0.4(-0.1) + 0.7(0.2) = 0.2$$

$$f(0.2) = \underline{\underline{0.549}}$$

$$* \text{Input to } N_k = 0.582$$

$$* \text{Error}(k) = (T - O_k) O_k (1-O_k) = (0.65 - 0.582)(0.582)(1-0.582) = \underline{\underline{0.418}}$$

Date: _____

* Error(i) = Error(k) $w_{ik} O_i (1-O_i)$

$$= 0.418 (0.1)(0.562)(1-0.562) = 0.01$$

* Error(j) = Error(k) $w_{jk} O_j (1-O_j)$

$$= 0.418 (0.5)(0.549)(1-0.549) = 0.0517$$

* Next step is to update the weights ~~based on associated~~ with individual node connectives.

* Adjustments made using delta rule.

↳ Minimize \sum square errors [error = diff(actual, actual output)]

$\Rightarrow w_{ik} = w_{ik}(\text{current}) + \Delta w_{ik}$

* $\Delta w_{ik} = r \times \text{Error}(k) \times O_i$

↳ r is learning rate or cl

* Compute Δw_{ik} , Δw_{j1} , Δw_{j2} , Δw_{j3}

$$\begin{aligned} * w_{ik} &= \overset{w_{ik}(\text{current})}{0.5} + (\overset{\text{error}}{0.8} \times \overset{O_i}{0.549} \times \overset{\Delta w_{ik}}{0.0517}) \\ &= 0.683 \end{aligned}$$

*

Algorithm:

1 * Initialize the network

i - create topology by choosing the number of nodes for the input, hidden, & output layers.

ii - initialize weights for all node connectives ~~randomly~~ to arbitrary values $[-1.0, 1.0]$

iii - choose $r \in [0, 1]$

iv - choose terminating condition.

2 * For all training instances:

i - feed training ~~data~~ instance through network

ii - determine the output error

iii - update network weights.

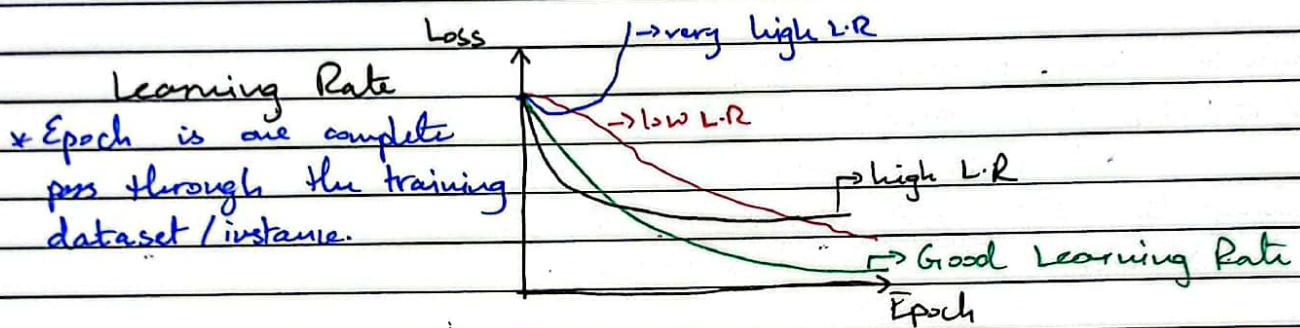
3 * If terminating condition not met, repeat step 2.

4 * Test accuracy of the network on a test data set

If accuracy is less than optimal, change one or more parameters & start over.

Training / testing of ANN:

- * During training phase, training instances repeatedly passed through network while weights are modified.
- * Connection weights changed to minimize training set error rate.
- * Training repeat till terminating condition met.
- * Terminate \rightarrow convergence of the network to a minimum value, could be
 - \rightarrow time criterion
 - \rightarrow any
 - \rightarrow Max iterations.



- * A bias value allows you to shift the activation function to the left or right \rightarrow may be critical for successful learning.

Weakness of ANNs:

- * lack the ability to explain their behaviour
- * not guaranteed to converge to an optimal solution
- * NNs can be over trained on training data but work poorly on test data.

Self Organizing Maps (SOM)

- * A SOM falls under the domain of unsupervised learning in Neural Networks.
- * Grid of neurons - each denotes one cluster learned during training.

SOM vs ANN.

- * SOM is not a series of layers, but a 2D grid of neurons.
- * Don't error correct rather learn by competitive learning.
- * Deal with unsupervised machine learning problems.
- * Competitive Learning for SOMs refer to the fact that when an input is "presented", only one of the neurons will be activated. In a way, the neurons on the grid "compete" for each input.

Clustering in SOM.

- * Each neuron has a location.
- * Neurons close to each other represent a cluster \rightarrow similar properties.
- * Each neuron has a weightage vector
 - \hookrightarrow equal to the centroid of its particular cluster.

SOM Setup:

- * ~~Input~~ N-Dimensional Vector $x : x \in (x_1, x_2, x_3, \dots, x_n)$
- * N-dimensional weight vector $w : w \in (w_1, w_2, w_3, \dots, w_n)$
- * No lateral connections b/w neurons

Training SOM. \rightarrow Several Steps over many Iterations.

1. Weights randomly initialized for each node.
2. Random vector chosen from training data \rightarrow given to lattice.
3. Every node examined \rightarrow calculate node's weights most similar to input. Winning node commonly known as Best Matching Unit (BMU) ^{vector}.
4. Radius of the neighborhood of BMU ~~fixed~~ calculated. Nodes within this radius deemed inside BMU's neighborhood.
5. Each neighboring node's (from 4.) weights are adjusted to make them like the input vector. Closer to BMU, more the weight gets altered.
6. Repeat from (2.) for N-Iterations.

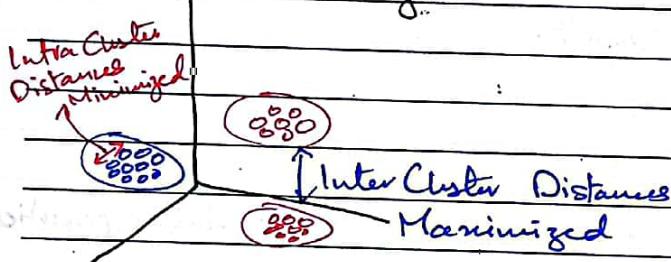
Neighborhood Function

- * Each Iteration after BMV, we find nodes within the BMV's neighborhood.
- ↳ All such nodes have their weights altered.
- * Kohonen Learning Algorithm → area of the neighborhood shrinks over time by reducing the radius.

$$\therefore w(t+1) = w(t) + \omega(t) L(t) [v(t) - w(t)]$$

- * L is the Learning Rate } Time Variant
- * ω → Theta is the neighborhood function. } Functions.

* Cluster Analysis:



- * A Self-Organizing Map (SOM) is a grid of neurons that adapt to the topological shape of the dataset.
 - ↳ Allows us to visualize large datasets & identify potential clusters.
- * SOM learns the shape of the dataset by repeatedly moving its neurons closer to the data points
 - ↳ Different groups of neurons thus reflect underlying clusters

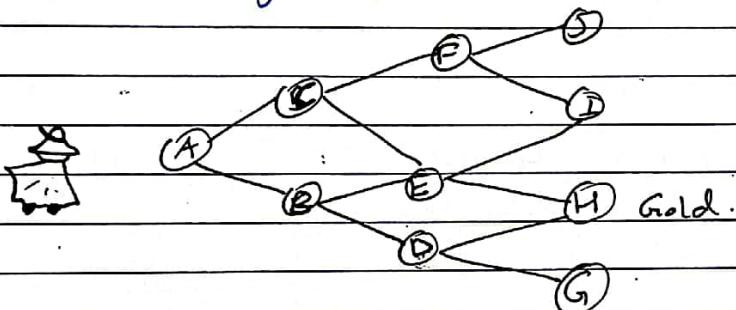
Date: _____

Reinforcement Learning (RL)

- RL agent learns by interacting with the environment & observing the results of these interactions.
- Reward & Loss system
- Mimics how humans & animals learn

Treasure Journey Problem

- * traveler looking for gold
- * white & black stones at each vertex ~~at~~ that has a signpost.
- * white stone go up^{"u"}, black stone go down^{"d"}



- * After concluding a journey, travel back to A,
 - ↳ Put a stone with an additional stone if gold found
↳ Reward
 - ↳ Remove stone if gold not found → loss / Penalty

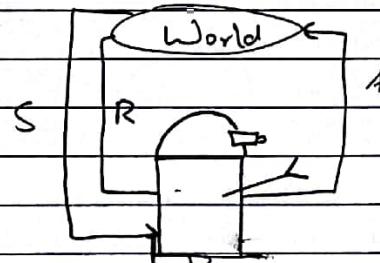
An RL Agent:

- observes an input state
- Action determined by a decision making function
- action performed
- agent receives a reward or penalty from environment
- information about reward is recorded
- * By observing the resulting reward, the policy for best action for a state can be fine-tuned.
- * Eventually the optimal decision policy will be generated & our agent can perform perfectly in that environment.

- ✓ RL is distinguished from other computational approaches
 - ↳ learning by the individual from interactions with the environment. → does not rely on exemplary supervision or complete models.

Basic RL model:

- 1 → Observe state, s_t
- 2 → Decide on action, a_t
- 3 → Perform action
- 4 → Observe new state, s_{t+1}
- 5 → Observe reward, r_{t+1}
- 6 → Learn from experience
- 7 → Repeat.



Markov Decision Process (MDP) 4-Tuple $(S, A, P(\cdot, \cdot), R(\cdot, \cdot))$

- * S is a finite set of states
- * A is a finite set of actions
- * $P_a(s, s')$
- * $R(s')$

Elements of RL → a policy

- optionally, a model of the environment
- a reward function
- a value function

Policy: * Mapping from perceived states of the environment to the actions to be taken when in those states.

Model = * Mimics the behaviour of the environment

- * might predict the resultant next state & next reward
- * RL models the environment in the form of MDPs

Reward Function: * defines goal in RL problem

- * maps each state to a single number
- * ↳ indicates intrinsic desirability of the state.

Value Functions:

- * value of a state is the amount of reward an agent can expect to accumulate over the future.
- * value function specifies what is good in the long-run.

Rewards vs Values: * without rewards, no values exists

- * purpose of estimating values is to get more rewards
- * we are mostly concerned with value when making & evaluating decisions.

- * \hookrightarrow we seek actions that bring about states of highest value, not highest reward

\hookrightarrow highest values yield highest rewards over the long run.

- * Efficiently calculating/estimating values is the most imp. component of almost all RL algs.

Cumulative Future Reward

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

$$G_t = \sum_{k=0}^{\infty} R_{t+k+1}$$

* G_t is total future rewards from time step t .

Time value of reward. \rightarrow Can't add rewards in reality.

\rightarrow rewards that come sooner are more probable to happen, since they are more predictable than the long term.

\hookrightarrow Discounting: * we define a discount rate 'gamma' $\in [0, 1]$.

* Larger gamma implies smaller discount.

\hookrightarrow learning agent cares more about the reward in long term

* smaller gamma means a bigger discount

\hookrightarrow learning agent care more about short-term reward.

Discounted Cumulative Reward

$$R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots$$

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad \text{where } \gamma \in [0, 1)$$

Date: _____

Value Iteration:Path: $s_0 \rightarrow s_1 \rightarrow s_2 \rightarrow s_3 \rightarrow \dots$

$$\text{Value}(s_0) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \gamma^3 R(s_3) + \dots$$

$$\text{Value}(s_0) \approx R(s_0) + \gamma V(s_1)$$

Since state transitions are stochastic, we have:

~~$$\text{Value}(s_0) \approx R(s_0) + \gamma E[V(s_1)]$$~~

* Value is the prediction of a series of rewards & is calculated as the expected sum of discounted rewards

* The agent will evaluate all the states in the state space, & the value function tells the agent how good each action is in a particular state.

↳ Agent chooses the best one.

Bellman Equation:

Expected Reward of executing action 'a' in 's'.

$$V(s) = \max_{a \in A(s)} \sum_{s' \in S} P_{ss'} [r(s, a, s') + \gamma V(s')]$$

Best Action from s For every state Immediate reward Discount factor Value of s'

Value Iteration:Initialize V arbitrarily, e.g., $V(s) = 0 \quad \forall s \in S^+$

Repeat:

$$\Delta \leftarrow 0$$

For each $s \in S$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V(s')]$$

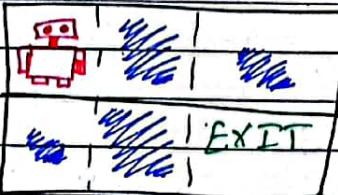
$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \delta$ (a small positive number)Output a deterministic Policy π^* s.t.

$$\pi^*(s) = \arg \max_a \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V(s')]$$

Date: _____

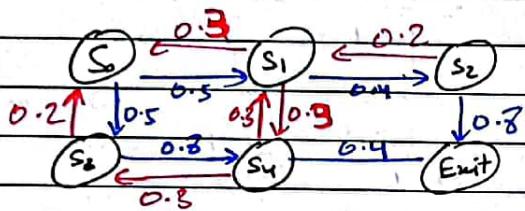
Baby Robot

* Small Puddles \rightarrow skid = 20%

↳ not skid = 80%

* Big Puddles \rightarrow skid = 60%

↳ not skid = 40%



Types of Learning

- * **Episodic** \Rightarrow rewards collected at the end of the episode, & then calculating maximum expected future reward.
- * **Temporal Difference learning** \Rightarrow Estimate the rewards at each step.

Q-Learning

- does not need environment → can be used online
- estimating values of state-action pairs
- Can perform adaptively in a world without understanding
 - ↳ sort out good actions from bad ones.
- $Q(s, a)$ value defined by the expected discount sum of future payoffs obtained by taking action 'a' from state 's' & following the current optimal policy.
 - ↳ Once values learned, optimal action from any state is one with the highest value.

Initialize $Q(s, a)$ arbitrarily

Repeat (for each episode)

 Initialize s

 Repeat (for each step of the episode)

 choose a from s using policy derived from Q

 Take ' a ', observe r, s'

$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_a Q(s', a) - Q(s, a)]$

$s \leftarrow s'$;

 until s is terminal.

Choosing an Action:

- function that assigns a probability of being chosen for each action in a given state.
- function should tend to choose actions with higher Q-values. should sometimes select lower Q-value actions
- prob of selecting highest-Q-value action should increase over time.

Action Selection Policy:

1. E-greedy:
 - $P(1 - e) \rightarrow$ choose max-value action
 - $P(e) \rightarrow$ randomly choose an action from all actions
- Big e at the start of Q-function
 - ↳ reduce e progressively as agent becomes better.

Exploration

e

Exploration

Exploitation

- know nothing abt env, more exploration.

2. Boltzmann Distribution.

→ ϵ can also take worst actions with equal chance even when it is unfavorable.

→ Vary action probs as a graded function

$$P(a|s) = \frac{e^{Q(s,a)/k}}{\sum_i e^{Q(s,a_i)/k}}$$

→ $k \rightarrow$ temperature → controls the prob of selecting non-optimal actions. $k \uparrow \rightarrow$ all actions selected uniformly
 $k \downarrow \& \lim k \rightarrow 0$, best always chosen.

→ Begin with large k & decrease it over time.

Calculating Error: \uparrow learning rate \uparrow Discount Factor

$$\text{New } Q(s,a) = Q(s,a) + \alpha \underbrace{[P(s,a) + \gamma \max\{Q'(s',a')\} - Q(s,a)]}_{\substack{\text{Reward} \\ \text{current val.} \\ \text{for } a' \\ \text{at } s'}} \quad \begin{array}{l} \text{Max expected} \\ \text{future reward} \\ \text{given new } s' \& \\ \text{all possible actions} \\ \text{at that new state.} \end{array}$$

$$\text{Loss} = (r + \gamma \max Q(s,a') - Q(s,a))^2$$