

III Bayesian Learning

Bayesian learning Provides a probabilistic approach to inference derive it is based on the assumption that the quantities of interest are governed by Probability distributions.

It is important machine learning because it provides a quantitative approach to weighing the evidence supporting alternative hypothesis.

Features of Machine learning Methods.

- 1) Each observed training example can incrementally decrease or increase the estimated probability that a hypothesis is correct
- 2) Prior knowledge can be combined with the observed data to determine the final probability of a hypotheses
- 3) Bayesian methods can accommodate hypotheses that make Probabilistic Predictions.
- 4) New instances can be classified by combining the predictions of multiple hypotheses, weighted by a probability.

Bayes Theorem :-

In machine learning we are often interested in determining the best hypotheses from some space ' H ', given in the observed training data ' D '.

Notations :-

$P(h)$ → initial probability that hypothesis ' h ' holds
(Prior Probability of h)

$P(D)$ → The prior probability that training data ' D '

will be observed

$P(D|h)$ → Probability of observing data 'D' given some word in which hypothesis 'h' holds.

$P(x|y)$ → Probability of x given y

In machine learning we are interested in probability $P(h|D)$ that 'h' holds given observed training data 'D', is called "Posterior Probability" of 'h'

$$\therefore \text{Bayes Theorem: } P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

In many learning scenarios the learner interested in finding the most probable hypothesis. Any search maximally probable hypothesis is called a "Maximum a Posteriori (MAP) hypothesis".

$$h_{\text{MAP}} \equiv \underset{h \in H}{\operatorname{argmax}} P(h|D)$$

$$= \underset{h \in D}{\operatorname{argmax}} \frac{P(D|h)P(h)}{P(D)}$$

$$= \underset{h \in H}{\operatorname{argmax}} P(D|h)P(h)$$

In some cases we assume that every hypothesis in 'H' is equally probable a priori

$$(P(h_i) = P(h_j))$$

∴ Maximum Likelihood (ML) hypothesis, h_{ML}

$$h_{\text{ML}} \equiv \underset{h \in H}{\operatorname{argmax}} P(D|h)$$

Probability Formulas

Product Rule: Probability $P(A \cap B)$ of a conjunction of two events

$$A \in B \text{ is } P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

Sum Rule: Probability of a disjunction of two events $A \in B$ is

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Bayes Theorem: The posterior probability $P(h|D)$ of 'h given D'

$$\text{is } P(h|D) = \frac{P(D|h) \cdot P(h)}{\sum_i P(D|h_i) \cdot P(h_i)}$$

Theorem of Total Probability: If events A_1, A_2, \dots, A_n are mutually exclusive with $\sum_{i=1}^n P(A_i) = 1$, then

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i)$$

Bayes Theorem & Concept Learning

Bayes theorem provides a principle way to calculate for posterior probability of each hypothesis given the training data, we can use it as the basis for a learning algorithm that calculates the probability for each possible hypothesis, then outputs the most probable hypothesis.

i) Brute-Force MAP Learning Algorithm

step 1: For each hypothesis h in H , calculate the posterior probability

$$P(h|D) = \frac{P(D|h) \cdot P(h)}{\sum_i P(D|h_i) \cdot P(h_i)}$$

step 2: Output the hypothesis h_{MAP} with highest posterior probability

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(h|D)$$

- || | | |

In order to specify a learning problem for the brute-force MAP learning algo we must specify what values are to be used for $p(h)$ & for $P(D|h)$. It is reasonable to assign the same prior probability to every hypothesis 'h' in H when there is no prior knowledge, sum to one.

$$\therefore p(h) = \frac{1}{|H|} \text{ for all } h \in H$$

$$P(D|h) = \begin{cases} 1 & \text{if } d_i = h(x_i) \text{ for all } d_i \in D \\ 0 & \text{if } d_i \neq h(x_i) \end{cases}$$

In other words the probability of data D given hypothesis 'h' is 1 if D is consistent with 'h', & 0 otherwise

Case 1: h is inconsistent

$$P(h|D) = \frac{0 \cdot P(h)}{P(D)} = 0, \text{ if } h \text{ is inconsistent with } D$$

Case 2: h is consistent

$$P(h|D) = \frac{\frac{1}{|H|}}{\frac{P(D)}{|V_{SH,D}|}} = \frac{\frac{1}{|H|}}{\frac{|V_{SH,D}|}{|H|}} = \frac{1}{|V_{SH,D}|} \text{ if } h \text{ is consistent with } D.$$

$V_{SH,D}$ — subsets of hypothesis from H that are consistent with D .

(Version space of H wrt D)

$$\therefore P(h|D) = \begin{cases} \frac{1}{|V_{SH,D}|}, & \text{if } h \text{ is consistent with } D \\ 0, & \text{otherwise} \end{cases}$$

$|V_{SH,D}| \rightarrow$ no. of hypotheses from H consistent with D .

∴ under our choice for $p(h)$ and $P(D|h)$, every consistent hypothesis has Posterior Probability $\left[\frac{1}{|V_{SH,D}|} \right]$ and every inconsistent

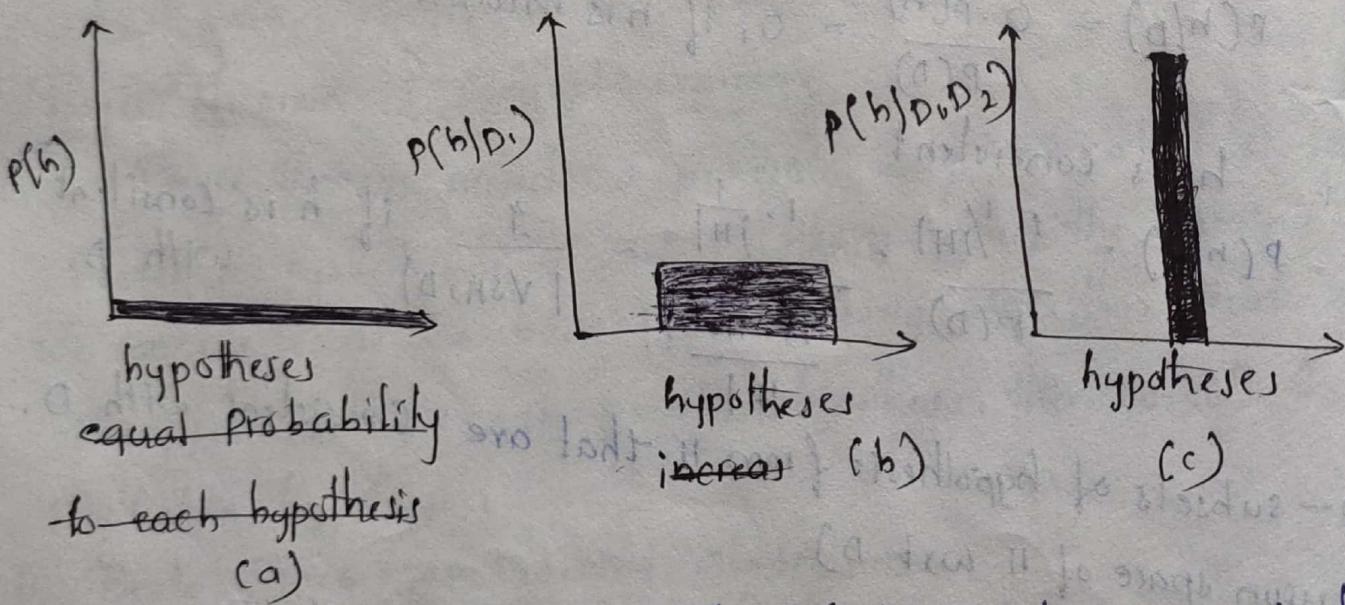
hypothesis has posterior probability $P(h|D)$
→ every consistent hypothesis is therefore a MAP hypothesis

(ii) MAP hypotheses & consistent learners :-

A learning algo is a "consistent learner" provided it outputs a hypothesis that contains zero errors over the training examples.

Every consistent learner upto a MAP hypothesis if we assume a uniform real probability distribution over H ($P(h_i) = P(h_j)$)

→ Concept learning algo find-S searches the hypothesis space H from specific to general hypothesis, outputting a maximally specific consistent hypothesis.



evolution of Posterior Probability $p(h|D)$ with increase in training data.

- a. Uniform Priors assign equal probability to each hypothesis.
- b. as training data increases first to D_1 .
- b. Then to $D_1 \cup D_2$
- c. The posterior Prob of inconsistent hypothesis becomes zero while posterior Prob increase for hypotheses remaining in the version space.

Eg: Consider a medical diagnosis pblm in which there are 2 alternative hypotheses:

- i. That the patient has a particular form of Cancer
- ii. That the patient does not

* Two Possible outcomes: \oplus Positive
 \ominus Negative

* Prior knowledge that over the entire population of people, Only 0.008 have this disease

* Correct positive result $\rightarrow 98\%$ } Test
Correct negative result $\rightarrow 97\%$.

In other case the test returns the opposite result

$$\therefore P(\text{cancer}) = 0.008 \quad P(\neg \text{cancer}) = 0.992$$

$$P(\oplus/\text{cancer}) = 0.98, \quad P(\ominus/\text{cancer}) = 0.02$$

$$P(\oplus/\neg \text{cancer}) = 0.03, \quad P(\ominus/\neg \text{cancer}) = 0.97$$

Suppose we have observe a new patient for whom the lab test return a positive result. Should we diagnose the patient as having cancer or not?

The maximum a posterior (MAP) hypothesis can be found using the following equation.

$$h_{\text{MAP}} = \underset{h \in H}{\operatorname{argmax}} \quad P(D|h)P(h)$$

$$P(\oplus/\text{cancer})P(\text{cancer}) = (0.98)(0.008) = 0.0078$$

$$P(\oplus/\neg \text{cancer})P(\neg \text{cancer}) = (0.03)(0.992) = 0.0298$$

$$\text{Thus } h_{\text{MAP}} = \neg \text{cancer}$$

Minimum Description length (MDL) Principle

The minimum Description length MDL principle is motivated by interpreting the definition of h-map.

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} P(D|h)P(h)$$

which can be equivalently expressed in terms of maximizing the \log_2 .

$$h_{MAP} = \underset{h \in H}{\operatorname{argmax}} \log_2 P(D|h) + \log_2 P(h)$$

or alternative, minimizing the negative of this quantity

$$h_{MAP} = \underset{h \in H}{\operatorname{argmin}} -\log_2 P(D|h) - \log_2 P(h)$$

It can be interpreted as a statement that short hypotheses are preferred

We can rewrite the above eqn as

$$h_{MAP} = \underset{h}{\operatorname{argmin}} L_C(h) + L_{C_0}(D|h)$$

C_H and C_0/h are optimal encoding for H and for D given h .

∴ Minimum Description length principle

$$h_{MDL} = \underset{h \in H}{\operatorname{argmin}} L_C(h) + L_{C_2}(D|h)$$

C & C_2 are optimal encodings of hypotheses

Bayes optimal classification

What is the most probable classification hypothesis given the training data?

What is the most probable classification of the new instance given the training data.

Eg: Posterior probabilities of 3 hypothesis h_1, h_2, h_3 are 0.4, 0.3, 0.3

Thus h_1 is the MAP hypothesis
Suppose a new instance x is encountered which is classified positive by h_1 , but negative by $h_2 \& h_3$.
Taking all hypotheses into account the probability that x is positive is 0.4, and the probability that x is negative is 0.6 ($0.3 + 0.3$)

The most probable classification (Gre) with 0.6
in this case is different from classification generated by the MAP hypothesis (0.4)

Bayes optimal classification is

$$\operatorname{argmax}_{v_j \in V} \sum_{h_i \in H} P(v_j | h_i) P(h_i | D)$$

$$P(h_1 | D) = 0.4, P(\ominus | h_1) = 0, P(\oplus | h_1) = 1$$

$$P(h_2 | D) = 0.3, P(\ominus | h_2) = 1, P(\oplus | h_2) = 0$$

$$P(h_3 | D) = 0.3, P(\ominus | h_3) = 1, P(\oplus | h_3) = 0$$

Therefore $\sum_{h_i \in H} P(\oplus | h_i) P(h_i | D) = 0.4$

$$\leq \sum_{h_i \in H} P(\ominus | h_i) P(h_i | D) = 0.6$$

and argmax

$$v_j \in \{\oplus, \ominus\} \leq P(v_j | h_i) P(h_i | D) = \ominus$$

Based Instance Bayes learning

Learning in instance bayes algo consists of simply storing the results of training presented data. When a new query instance is encountered a set of similar related instances is retrieved from memory & used to classify the new query instance.

k-Nearest Neighbour (kNN) Learning

This algo assumes all instances correspond to points in the n dimensional space \mathbb{R}^n .

The nearest neighbours of an instance are defined in terms of standard Euclidian distance.

Let an arbitrary instance x be described by the feature vector

$$\langle a_1(x), a_2(x), \dots, a_n(x) \rangle$$

where $a_r(x) \rightarrow$ value of r th attribute of instance x .

Then the distance b/w the two instances x_i and x_j is defined to be

$$d(x_i, x_j)$$

where,

$$d(x_i, x_j) = \sqrt{\sum_{r=1}^n (a_r(x_i) - a_r(x_j))^2}$$

Algorithm :-

Training algorithm

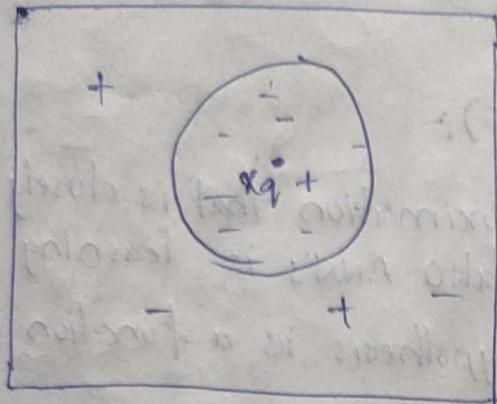
For each training example $\langle x_i, f(x_i) \rangle$ add the example to the list training examples

Classification Algorithm

- Given a query instance x_q to be classified,
 - let x_1, \dots, x_k denote the k instances from training examples that are nearest to x_q .
 - Return

$$f(x_q) \leftarrow \operatorname{argmax}_{v \in V} \sum_{i=1}^k \delta(v_i, f(x_i))$$

where $\delta(a, b) = 1$ if $a = b$ and where $\delta(a, b) = 0$, otherwise



A set of +ve and -ve training examples is shown on the graph along with new query instance x_q to be classified.

One-nearest neighbour algo classified x_q as +ve, whereas 5-nearest neighbour classifies x_q as -ve.

Locally Weighted Regression :-

The nearest neighbour approach is approximating the target function $f(x)$ at the single query point $x_q = x$.

whereas locally weighted regression constructs an explicit approximation to $f(x)$ over a local region surrounding x_q .

Locally weighted regression uses near by or distance weighted training examples to form this local approximation to the target function $f(x)$.

Let us consider the case of locally weighted regression in which the target func f is approximated near x_q using the linear function of the form $f(x) = w_0 + w_1 a_1(x) + \dots + w_n a_n(x)$ where, $a_i(x) \rightarrow$ The value of i th attribute of the instance x we derived methods to choose weights that minimize the squared error

Summed over, the set D of training examples.

$$E = \frac{1}{2} \sum_{x \in D} (f(x) - \hat{f}(x))^2$$

which let us to gradient descent training rule

$$\Delta w_j = \eta \sum_{x \in D} (f(x) - \hat{f}(x)) a_j(x)$$

* Radial Basis Functions (RBF)

One approach to function approximation that is closely related to distance weighted regression & also ANN's is learning with RBF's.
→ In this approach the learned hypothesis is a function of the form

$$\hat{f}(x) = w_0 + \sum_{u=1}^k w_u k_u(d(x_u, x)) \quad \text{①}$$

$x_u \rightarrow$ instance from X

Kernal function $k_u(d(x_u, x))$ is decreases as the distance $d(x_u, x)$ increases.

Even though $\hat{f}(x)$ is a global approximation to $f(x)$, the contribution from each of the $k_u(d(x_u, x))$ term is localized to a region nearby the point x_u .

It is common to choose each function $k_u(d(x_u, x))$ to be Gaussian function, centered at the point x_u with some

Variance, σ_u^2

$$k_u(d(x_u, x)) = e^{-\frac{1}{2\sigma_u^2} d^2(x_u, x)}$$

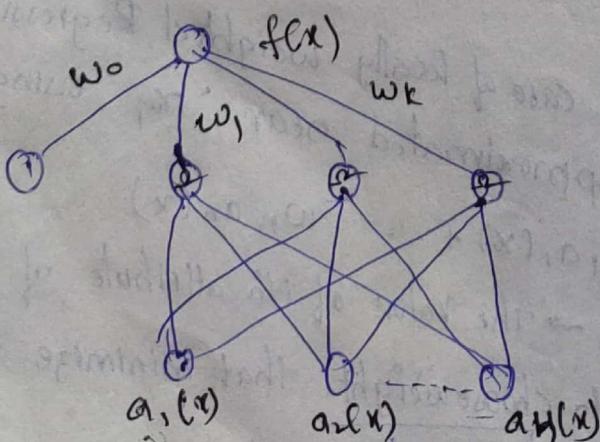


fig: A Radial Basis Functions (RBF) Network.

Each hidden node produces an activation determined by a Gaussian func centered at some instance x_{ui} . Therefore its activation will be close to zero unless the input x is near to x_{ui} .

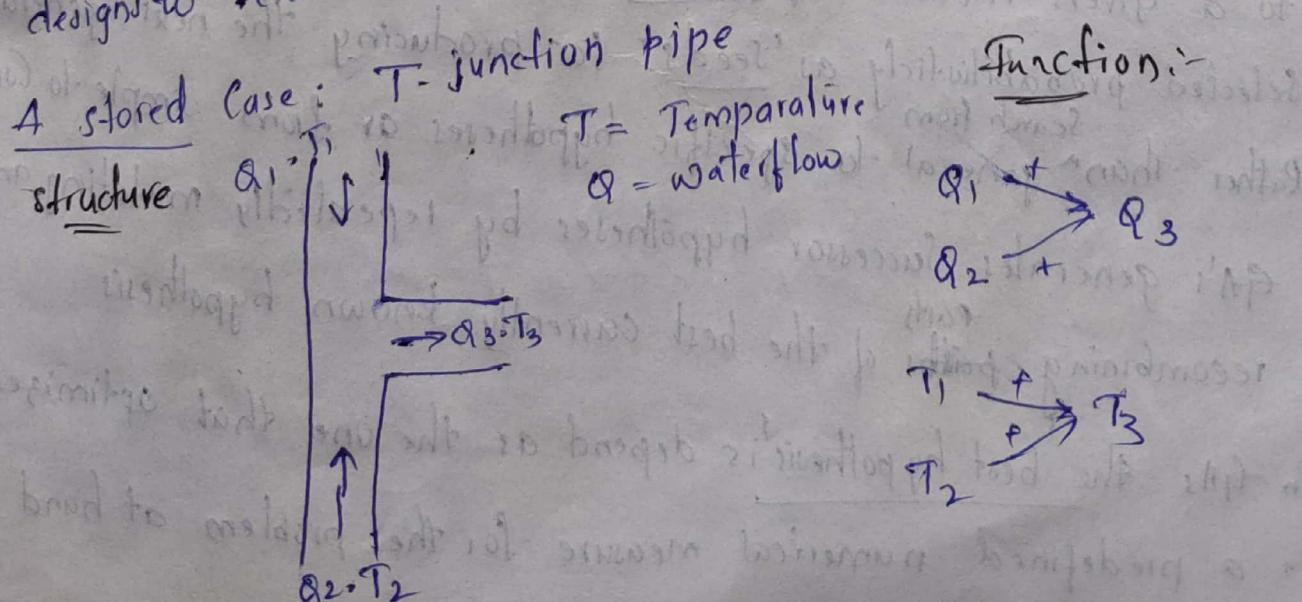
→ The o/p unit produces a linear combination of the hidden unit activations.

Case Based Reasoning (CBR):

In CBR instances are typically represented using rich symbolic descriptions & methods used to retrieve similar instances or correspondingly more elaborate.

→ CBR has been applied to problems such as conceptual design of mechanical devices based on a stored library of previous designs, reasoning about new legal cases based on previous rulings & solving planning and scheduling problems by reusing and combining portions of previous solutions to similar problems.

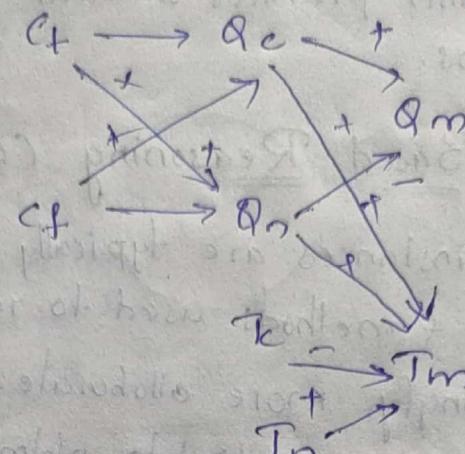
e.g:- CADET System
 The CADET system employs a case based reasoning to assist in the conceptual design of simple mechanical devices such as water faucets. It uses a library containing approximate previous designs & design fragments to suggest conceptual designs to meet the specifications of new design problems.



A Boolean Specification: Water-Paucet

Structure :-

Function:



* Genetic Algorithms:

Genetic algorithms provide an approach to learning that is based loosely on simulated evaluation.

The search for an appropriate hypothesis begins with a population, or collection of initial hypotheses.

members of the current population give rise to the next generation population by means of operations such as random 'mutations' and 'crossovers' which are patterned after processes in biological evaluations

→ at each step, the hypotheses in current population are evaluated related to a given measure of 'fitness' with the most fit hypotheses selected probabilistically as 'seeds' for producing the next generation.

Rather than ^{Search from} general to specific hypotheses or from simple to complex,

GA's generates successor hypotheses by repeatedly mutating and

recombining ^{parts} of the best currently known hypothesis

In fits, the best hypothesis is depend as the one that optimizes a predefined numerical measure for the problem at hand

Called hypothesis 'fitness'

Algorithm :-

GA (fitness, fitness-threshold, p, r, m)

fitness : A function that assigns an evaluation score, given a hypothesis.

fitness-threshold : A threshold specifying an termination criterian.

p : the number of hypotheses to be included in the population

r : The fraction of the population to be replaced by crossover at each step.

m : mutation rate

- Initialize population: $P \leftarrow$ Generate p hypotheses at random
- Evaluate : for each h in P , compute $\text{fitness}(h)$
- while $[\max_h \text{fitness}(h)] < \text{fitness-threshold}$ do:

Create a new generations, P_s :

1. Select: Probabilistically select $(1-r)p$ members of P to add P_s .

The probability $Pr(h_i)$ of selecting hypothesis h_i from P

is given by

$$Pr(h_i) = \frac{\text{fitness}(h_i)}{\sum_{j=1}^p \text{fitness}(h_j)}$$

2. Crossover: Probabilistically select $\frac{r}{2}p$ Pairs of hypotheses from P , according to $Pr(h_i)$ given above. for each Pair(h_1, h_2), Produce two offspring by applying The crossover Operator
Add all offspring to P_s .

3. Mutate :- choose m percent of the members of P_s with confirm probability for each invert one randomly selected bit in its representation.

4. Update : $P \leftarrow P_s$

5. Evaluate : for each h in P , compute $\text{Fitness}(h)$.

- Return the hypothesis from P that has the highest fitness

$P_s \leftarrow$ Successor Population.

Representing hypothesis :-

Hypotheses in GAs are often represented by bit strings, so that they can be easily manipulated by genetic operators such as mutation, and crossover.

→ Sets of 'if-then' rules can easily be represented in this way by choosing an encoding of rules that allocates specific substrings for each rule precondition & postcondition.

Eg :- Consider the attribute 'outlook' which can take on any of three values = sunny, overcast, Rain.

$$\text{outlook} = \{\text{sunny}, \text{overcast}, \text{Rain}\}$$

One way to represent a constraint on outlook is to use a bit string of length three where each bit ^{position} corresponds each bit position through one of its three possible values

$$\text{outlook} = \text{Overcast} \Rightarrow 010$$

$$\text{outlook} = \text{Overcast} \vee \text{Rain} \text{ (more general)} \Rightarrow 011$$

$$\text{outlook} = \text{Sunny} \vee \text{Overcast} \vee \text{Rain} \text{ (most general)} \Rightarrow 111$$

Consider a binary attribute 'wind'

wind = {strong, weak}

For example, a rule pre condition such as

(outlook = Overcast \wedge Rain) \wedge (wind = strong)

outlook wind
 011 10

Rule post condition (such as playTennis = Yes) can be represented in a similar fashion.

For example, the rule
IF wind = strong THEN play Tennis = Yes will be represented by the string.

outlook wind playTennis
 111 10 10

In some GAs hypotheses are represented by symbolic descriptions rather than bitstrings.

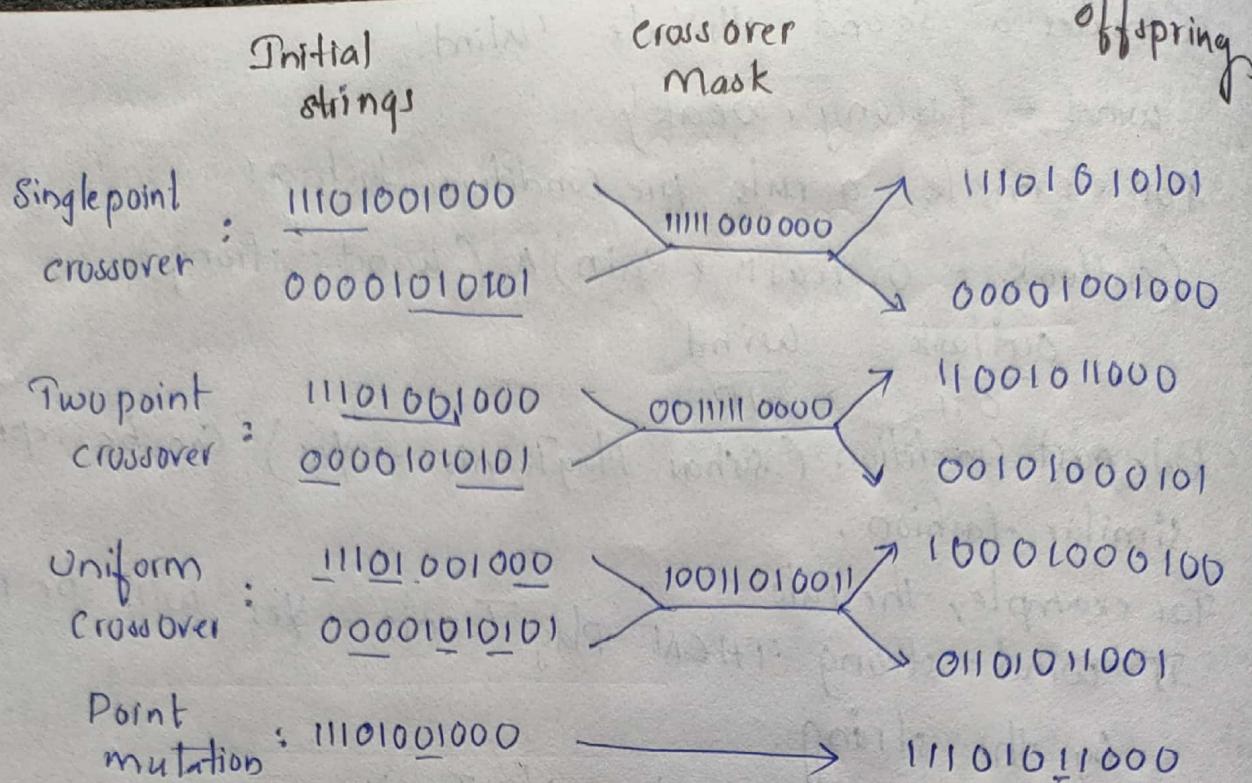
Genetic Operators

The generation of successors in a GA is determined by a set of operators that recombine and mutate selected members of the current population. The two most common operators are 'crossover' and 'mutation'. The 'crossover' operator produces two new offspring

from two parent strings by copying selected bits from each parent. The bit at position 'i' in each offspring is copied from the bit at position 'i' in one of the two parents.

The choice of which parent contributed the bit for position 'i' is determined by an additional string called the 'crossover mask'.

Ex:



In particular, the 'mutation' operator produces small random changes to the bitstring by choosing a single bit at random, then changing its value.

Mutation is often performed after crossover has been applied, as in our 'prototypical algorithm'

Genetic Programming

Genetic programming is a form of evolutionary computation in which individuals in the evolving population are computer programs rather than bit strings.

- Programs manipulated by GP are typically represented by trees corresponding to the parse tree of the program.
- Each func call is represented by a node in a tree and arguments to the function are given by its descending nodes.

Ex: Function: $\sin(x) + \sqrt{x^2 + y}$ is represented by the tree structure.

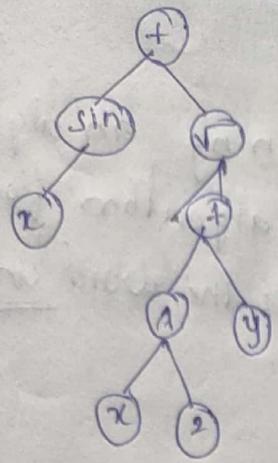


fig :- Program tree representation in EPL.

In the algo on each iteration, in the algo it produces a new generation of individuals using selection, crossover, and mutation.

- Crossover operations are performed by replacing a randomly chosen subtree of one parent program by a subtree from other parent program.

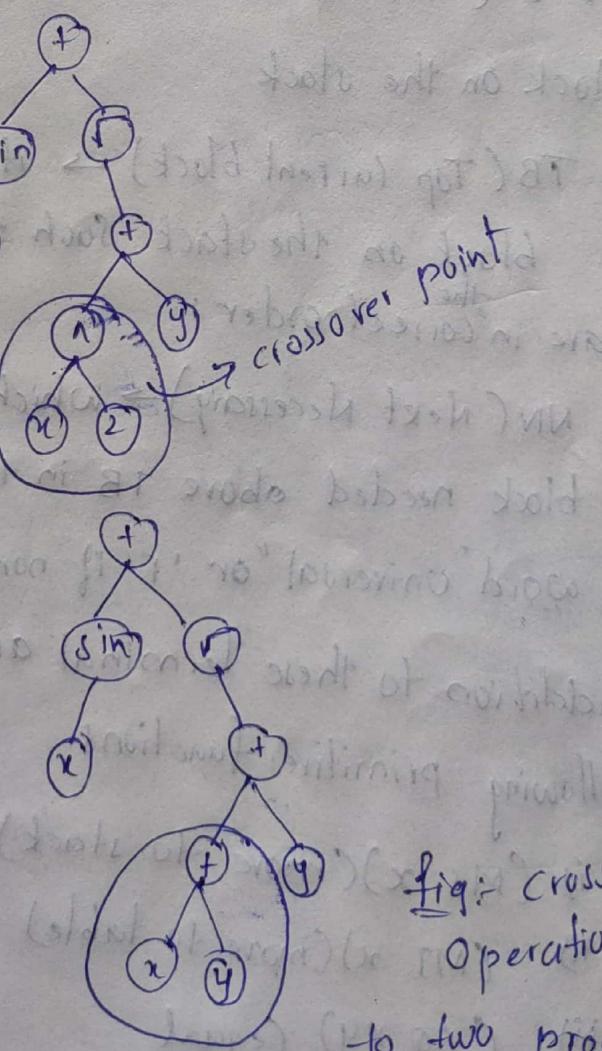
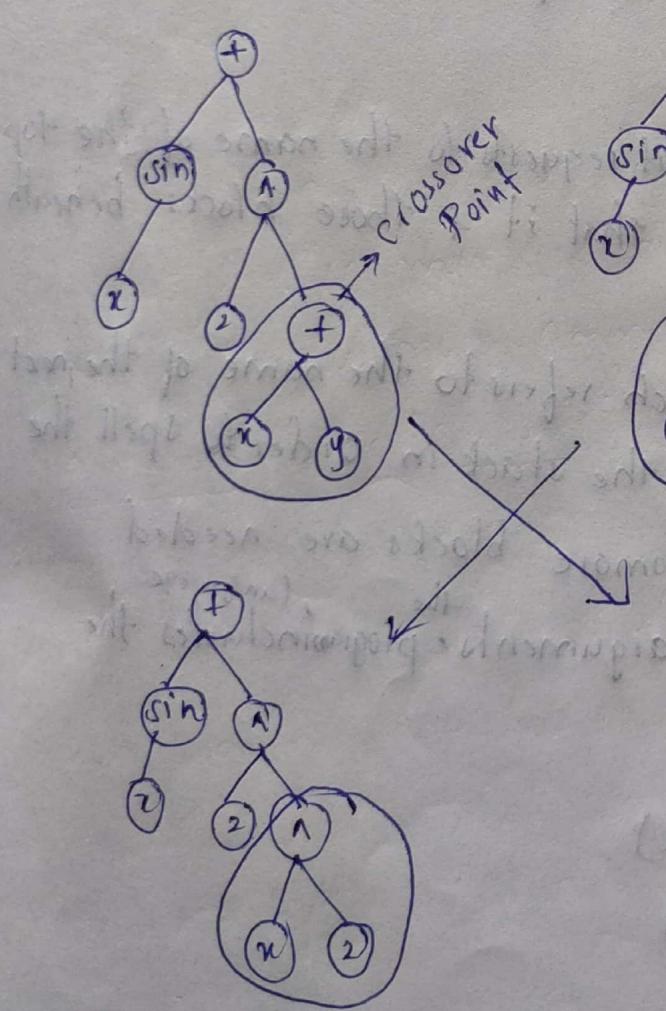
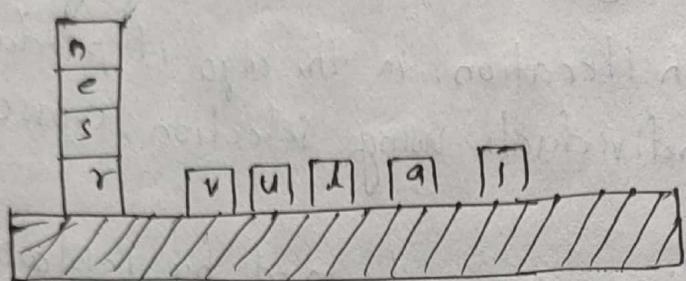


fig :- Crossover operation applied to two programs.

Illustrative Example:

One illustrative examples involves learning an algo for stacking blocks. The task is to develop a general algorithm for stacking the blocks into a single stack that spells the word "universal"



→ fig: A block stacking program.

The following are the three terminal arguments

- (i) CS (current stack) → which refers to the name of the top block on the stack
- (ii) TB (Top current block) → it refers to the name of the topmost block on the stack such that it & those blocks beneath it are in ^{the} correct order.
- (iii) NN (Next Necessary) → which refers to the name of the next block needed above TB in the stack in order to spell the word "Universal" or 'F' if no more blocks are needed.

In addition to these terminal arguments, the language program includes the following primitive functions.

- i) (MS α) (move to stack)
- ii) (MT α) (move to table)
- iii) ($\epsilon Q \alpha y$) (equal)
- iv) ($\text{NOT } x$)
- v) ($\text{DO } \alpha y$) (do until)

The population was initialized to a set of 300 random programs
After 10 generations the system discovered the following program,
which solved all training example.

(EQ(DU(MT(s))(NOT(s)))(DU(MS NN)(NOT NN)))