

L1: Introduction



M/4-2-2019

* What is Artificial Intelligence ?

هذه الكثير من المصطلحات تندرج تحت الذكاء الاصطناعي منكم أن تكون :
Data science , machine learning , deep learning ,
expert systems etc .

في المعنى الذكاء الاصطناعي هو جعل الآلة أن يتصرف كما يتصرف الإنسان وليس كما
يتفكر الإنسان ، فالآلة لا تتعلم بنفسها .

The machine take the knowledge and decision by A set of
past practices .

The way of learning in the machine is done through knowledge
and past experiences .

- The simplest thing to imagine the machine learning or
AI is make (if-else statements) .

أداء ال machine أفضل من أداء الإنسان في كل عام ولكنه هذا لا يعني
أه ال machine تفكر أفضل من الإنسان .

- الآلة أفضل من الإنسان أو تتقلب عليه من خلال أمرين :

1. الآلة لا تنسى

2. الخبرات التي تحفل عليها الآلة من أمور صعبة تؤدي إلى توفير الأمور الباقية .

مثال :

فعل الآلة من خلال البصيرة ، التي تقدر سرعة واضحة في عدة جوانب والتي هي من ذلك

لأنه من الممكن في أي sample يكون مع البصيرة أسرع وأكثر على البصيرة من أي

من الممكن أنه يعتبرها من البصيرة لذلك فإنه عدة samples .

* If we increase the past experience or training of computer ,
the model will be more general and the error rate
will be lower .

- The intelligence is come from knowledge. (الذكاء يأتي من المعرفة)

* How we give the knowledge to the Computer? (كيف نمنح الحاسوب المعرفة؟)

- At first, we should give it the knowledge then, by the knowledge it should make the decisions.

Before answer the previous question we should know that the knowledge is problem dependent.

- There is no general knowledge, each problem have specific knowledge.

- Usually if we want to enter the data to the computer, we using file or keyboard ... etc.

But this is Data not knowledge.

* The difference between Data and knowledge is:

Data contain redundancy, noise and many things including the knowledge.

لذلك Data هي المعرفة Knowledge

- كل technique في AI لها مدخلتها في خزنة Knowledge، لذلك لا يوجد في AI طريقة عامة لمعالجتها خزنة ال Knowledge في problem. فهي تعتمد على المعرفة التي سيتم اختيارها من الخزانة.

- أهم شئ في AI هو تخزين Knowledge وما يتبعه من أدوات ومعرفة.

- Intelligence = knowledge + ability to perceive, feel, understand, process, communicate, judge, and learn.

معرفة + قدرة على اتخاذ القرارات

One of the most important applications on the AI is the Robotics.

* Main area of AI is ① Search : there's many algorithms on search.
② planning , ③ learning , ④ Robotics , ⑤ Natural language processing

* Intelligent System Should do :

- Should have the ability to automatically perform tasks.
- Should have the Flexibility in dealing with variability.

AI : acting rationally \Rightarrow rational agent.

Agent takes the status of environment from the sensors, but the sensors is not necessary be hardware.

وظيفة الاستشعار هي استقبال المدخلات من البيئة التي من خلالها يستدل على المحيط.
حل تطبيقية عن الإنترنت تتميز منطوق أمرية :

① IP ② port number of application.

ملاحظة ذلك عند إرسال البريد الإلكتروني إلى المرسل إليه مباشرة، بل إلى server
ثم نقل رسالة البريد الإلكتروني إلى عميل العميل أو العميل المرسل لها من خلال مستخدم (مستلم غير موثوق)
لذلك الآلية ليست بالضرورية أن يكون شيئاً مادياً.

العمل على نتيجة بديهية يعتمد على المدخل الذي من خلاله يتم الاستدلال والبرهان المنطقية المنجزة
من الآلة (Knowledge)

- Sensing
- decision \rightarrow which is done through sensing and experience.
- output

[sensors, agent program, actuators]

* Examples of Agents :-

- human agent : sensing → إحساس
decisions → تفكير (thinking)
actuators → أجزاء → (hands, fingers, ... etc.)
- Robot : sensors → camera, microphone, ... etc.
decisions → software agent.
actuators → motors → (wheels, lights, ... etc.)

[agents s/w \rightarrow sensing, actuators \rightarrow أجزاء \rightarrow AI \rightarrow تفكير]

- Environment is not fully observable
it's not static.

[environment, sensing, actuators]

L2: Problem Formulation and Solving by Search.

W/ 6-2-2019.

trees: there isn't any cycles in the trees, we can reach any node by one track. and the vice versa in the **Graphs**.

So the trees is special case of graphs.

We will apply the search on the graphs.

in search we are moving from state to another through a specific decision.

The question is what is the role of AI here ??

We compare between one algorithm and another through:

- 1] time complexity.
- 2] Space complexity.

The answer of question is:

AI adds to the decision - that searching algorithm do it - knowledge to reach to the best solution.

مبدأ كفاءة الخوارزميات تنبثق في أوقات قليلة من كفاءة الخوارزميات التي لها القدرة الزمنية التي
تسمح لها بالعثور على الحل الأمثل في وقت أقل من الخوارزميات التي ليس لديها Knowledge
الكافية في الذاكرة وقد استلزمنا وقتاً أطول للعثور على الحل.

AI will solve a real problems.

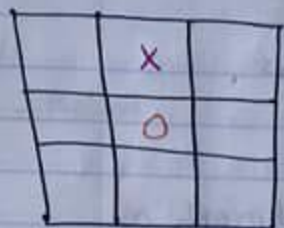
دور ال Search في AI هو البحث عن أفضل حل ممكن أو بين حل موجود

↳ (Optimization)

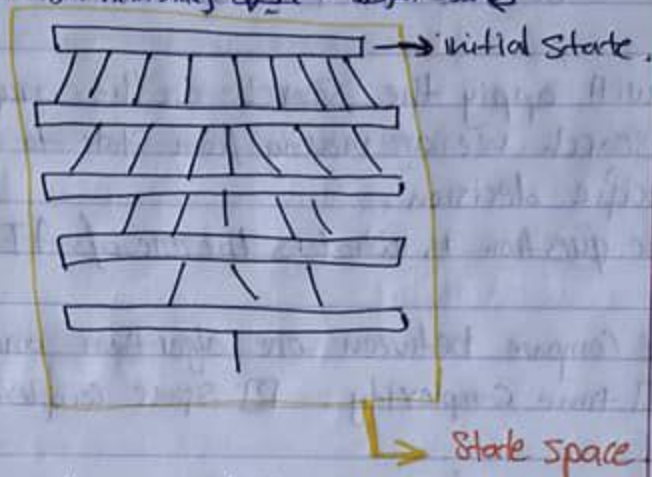
دور ال Search في AI هو البحث عن أفضل حل ممكن أو بين حل موجود

Example :

X/O game.



Initially you have 9 choices.
[generate for next state] هذه المرحلة تليها



- * We apply search algorithm on search space (graph)
- The different possibilities of states.

Next States : All possible states in the current status.

We don't need to generate all state space, we should know any next state.

* What should I know from the issue to solve it by search??

بما يجب ان نعرفه من المسألة لكي نحلها بالبحث

- 1 we should define the initial state (start point),
 - 2 we should know the goal state (stop point) termination state.
- for solving any search algorithm.

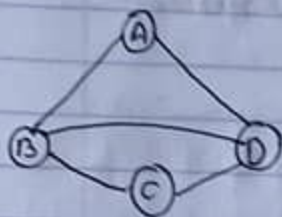
For AI Search techniques we also need to know about the Cost \rightarrow to reach to the optimal solution by comparing the cost.

3 Standards of differentiation. (معايير التمييز)

* Searching techniques depends on :
Knowledge base search or don't use knowledge on search.

* we implement the data structure of tree as linked list - as programming - and implement the data structure of graph as two dimensional array.

- the implementation of graph means collectivity between nodes.



The easiest way to implement the collections between nodes is two dimensional array.

* we will implement the nodes as objects and each object has pointers refer to another objects that collect with them.

الخوارزميات (next state) هي التي تقرر العرف بين Algorithm و Algorithm آخر.

it's not possible to define illegal next state.

في كل حالة يوجد شروط يجب انقضاءها ومن ثم نتنازل ان next state غير ممكنة مثال ذلك في لعبة الشطرنج لكل حرك يوجد حركات معينة تقبلها عن اللعبة.

- examples of algorithms that don't use knowledge on search: depth first search and breadth first search.

Some algorithms we don't need to know the path, but another algorithms, the path is the solution.

the algorithms that don't use knowledge \Rightarrow can't know if the solution is optimal or not.

But the algorithms that depends on knowledge \Rightarrow we can know the best solution.

Algorithm $\left\{ \begin{array}{l} \rightarrow \text{complete} \rightarrow \text{optimal / not optimal.} \\ \rightarrow \text{incomplete} \end{array} \right.$

operator / action : actions that move from current state to the next state.
Condition of operator and current state are in operator.

قبل أن يتم الأجراء التالي للحل يجب أن يكون الحل موجوداً
مثال ذلك : جسد المتغيرات المتصلة

- Some applications - problems
 - Route finding : الأكثر شيوعاً
 - VLSI layout : path isn't necessary.

There's many examples of problems on slides.

الهدف العام لل (Searching Algorithm) هو الانتقال من Initial state الى goal state
بناءً على معايير محددة يمكن أن تكون الـ Initial state أو goal state غير ذلك.

* **Uninformed Search** is not use knowledge but it's not Random Search.

L3: Uninformed Search / Breadth-First-Search.

M / 11-2-2019.

* Searching algorithm in general :

- Starts from initial state
- Test the state including initial state, if it's goal \Rightarrow state terminate.

* Search algorithm و آخره Search algorithm. الـ Search algorithm.

1. Search tree or generation

2. Visited nodes

3. goal or the path

- الجهد الأول هو عمل formalization الـ

Initial state

goal state

← operators (actions)

ليس جميع الـ Search algorithm. الـ Search algorithm.

* Searching Algorithm

→ informed Search : Using knowledge (heuristic)

→ Uninformed Search : Without using knowledge.

* We use knowledge: to reach to the solution faster or to get the optimal solution.

→ in informed Search

reduce the time to reach the goal.

- uninformed algorithms use the order of nodes to go to the next state.

\Rightarrow we guarantee that the algorithm will be complete.

complete \Rightarrow goal or the path

Uninformed Search : breadth first search, depth first search.

Standard of discrimination: معايير التمييز

- Time Complexity
- Space Complexity
- Completeness: in the worst case \rightarrow reaches to the all nodes.
- Optimality \rightarrow The best solution.

any optimal is complete but not any complete is optimal.

The cost will increase on time and space \leftarrow Complete \neq Optimality

Complete optimal \leftarrow جميع المقارنات بعين محترمة في الذاكرة لذلك Space, time \rightarrow يزيد

التوضيح

optimal solution \leftarrow يعني أقل Cost. Cost يمكن أن يكون time, space \rightarrow ... التي ولكن قد نشأ مسافة للوصول على (optimal solu) يكون عن هذا زيادة المسافة والوقت، ولكن لا يهمنا المسار الذي سنأخذ للوصول إلى optimal.

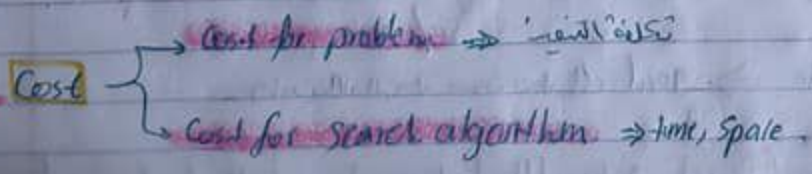
Searching Algorithm \leftarrow يتم عملها قبل التنفيذ (before running time).

فعل تنفيذ الخوارزمية من أجل أقل من الحل الذي أريد.

لذلك وظيفة (AI) البحث عن أفضل حل قبل البدء (Running).

هناك 2 مكونات التكلفة، قبل البدء يكون هناك تكلفة ولها هذه التكلفة (تكاليف التشغيل) problem cost تكلفتها عند تنفيذها تكلفتها التنفيذية.

عندما نتحدث عن optimal solution هو أقل Cost لل problem والذي لا search algorithm.



* Complete and optimal we can know them from the procedure of the search.

- Searching algorithms work on state space, as data structure is the graph.

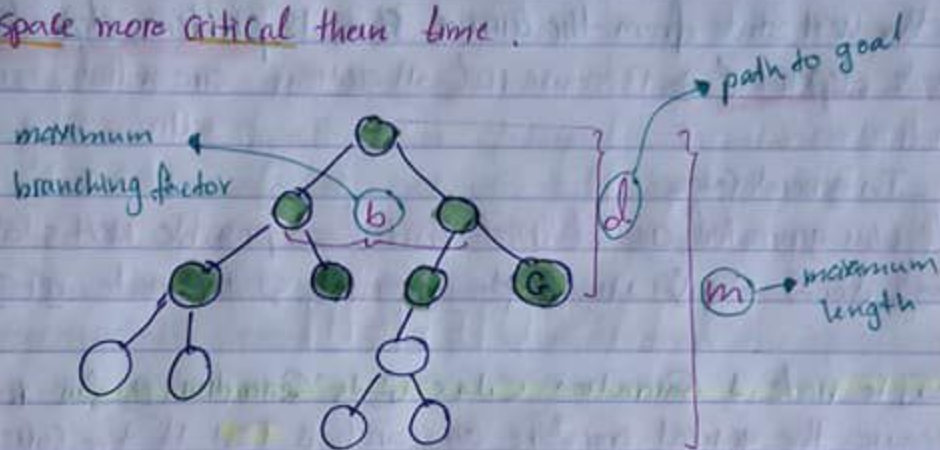
- The concept of any searching algorithm is "Search tree".

* How we calculate the time & space Complexity??

- depth \rightarrow number of levels \rightarrow affects on time.

- Branching factor (تعدد الفروع) \rightarrow affect on space.

- space more critical than time.



* Formalization \rightarrow moving to the mathematics or logic or Algorithm.
 نقل من الفهم إلى المنطق.

we should formalize it \Leftarrow Search by algorithmic process

goal state is حالة الهدف, initial state \rightarrow حالة البداية \Leftarrow Searching algorithm الخوارزمية
 (Set of operators) مجموعة العمليات

* River Example \rightarrow from slides.

initial state \rightarrow حالة البداية \rightarrow final state \rightarrow حالة الهدف

There is a Constraint \Rightarrow not all actions open in any state.

State space \rightarrow فضاء الحالات (List of actions) \rightarrow قائمة الإجراءات
 State space \rightarrow فضاء الحالات

Technically \Rightarrow The best implementation of state is "node".

in 'C' \Rightarrow "struct", in 'Java' \Rightarrow "object"

لللعبة داي في قبل اللب لعلبة (search) للعبة (Search Space) للعبة
وذلك بـ 3 و 3 في الـ (State space) للعبة و 3 في الـ (Search Space) للعبة
3 في الـ (Search Space) للعبة

What we can put in the node??

- Datatypes - Methods

[key], [array of objects]

What is the next state from the current state is problem dependent,
and state dependent \Rightarrow Because not all actions are valid.

* note \Rightarrow In general case.

- Any node was generated and does not visited \Rightarrow possible next state.

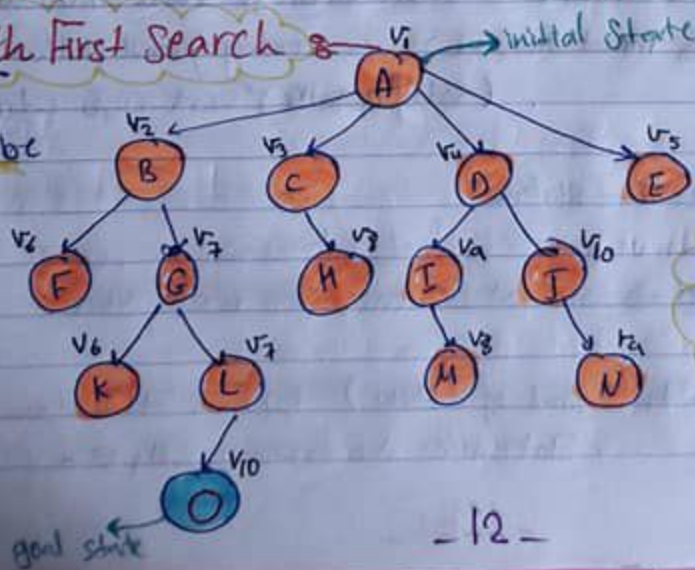
[next state \Rightarrow generated node \Rightarrow possible next state]

* note \Rightarrow In general case.

Because the general procedure on search is Test if the current state is Goal or not, if it isn't goal \Rightarrow then generate for next state. if algorithm goes to an next state and don't go to another Generation States, these states will not removed. (next state \Rightarrow possible next state)

Breadth First Search

guarantee to be Complete.



* not depends on knowledge but it depends on Order

next state \Rightarrow problem dependent & state dependent.

* Uninformed search looks to all nodes the same, but it doesn't choose the node randomly.

- **Breadth** choose the node from left to write \Rightarrow على الترتيب الهجائي
 \hookrightarrow moves level level. in graph

- Visit all nodes in the same level then goes to the next level.

- To Implement the Breadth as Data Structure we used the "Queue".

- Search algorithm on graph \Rightarrow 'Tree'

Using Queue as Data Structure of Breadth is guarantee to reach all levels, and this guarantee to be complete.

- if there is cycles is not guarantee to be complete, but in this algorithm avoid the cycles by using Queue. Because I can check any new input if it exists then it will not entered the queue to avoid the repetition, so I avoid the cycles.

* Space Complexity \Rightarrow b^{d+1} exponential

* Time Complexity \Rightarrow b^{d+1} exponential.

- Breadth First Search \Rightarrow 1 Complete

2 It's not guarantee to be optimal.

So we cannot decided if this solution optimal or not.

L4: Uniform-Cost / Depth Search.
 W/13-2-2019

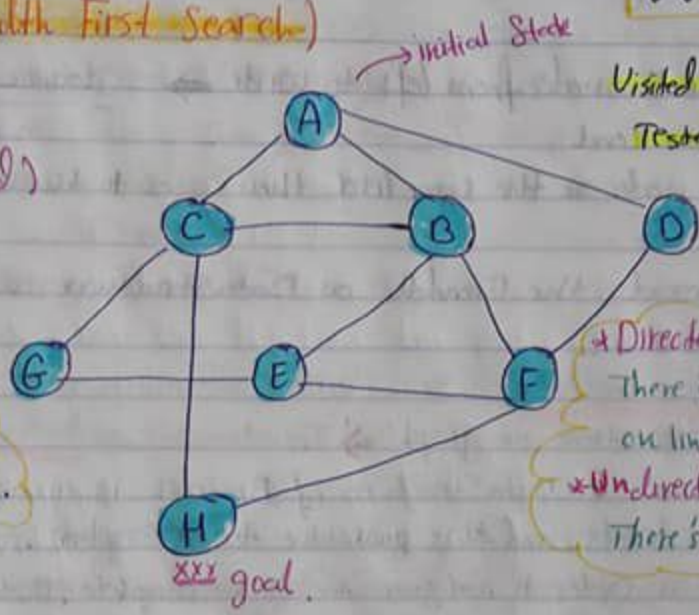
* Next State: Any links with node without parent.

* Example: (Breadth First Search)

Queue

Graph (Undirected)
 - State Space -

Can move from C to B or from B to C because the graph is Undirected.

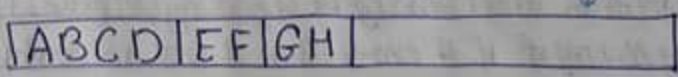


Visited nodes: Generated & Tested.

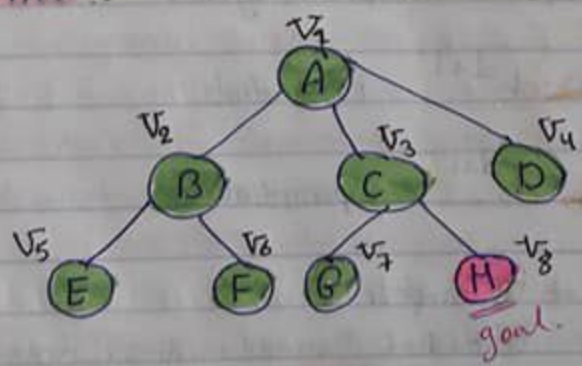
Any visited node is generated but not any generated is visited.

* Directed graph: There is a direction on links between nodes.

* Undirected graph: There's No directions.



Search Tree



* Visited Nodes:
 A, B, C, D, E, F, G, H.

* When the breadth search be optimal??
 - when the Cost is equal in the same level.

* Path:
 A, C, H.

solution doi njo qia Cost di njo njo *
 optimal. njo njo

2) Uniform-Cost-Search. Data Structure is "Queue".

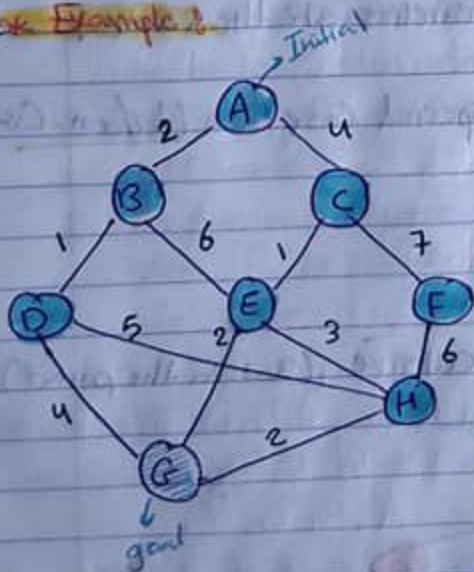
Cost → path cost (exact) ⇒ $\sum \text{cost of edges}$
 → expected cost → AI (heuristic): Remaining cost from next state to the goal.

* So we means in Uniform-Cost-Search the path cost.

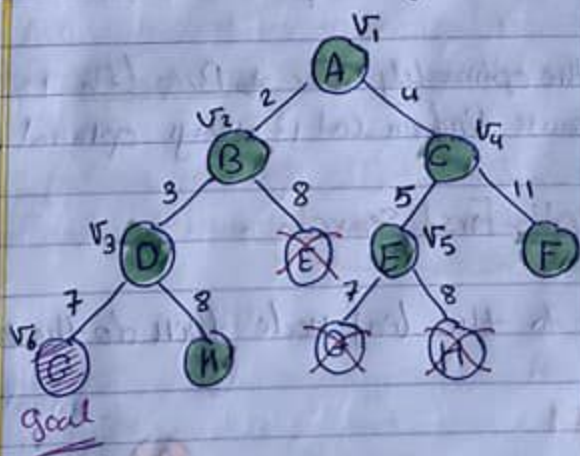
Uniform-Cost as dijkstra algorithm ⇒ returned the shortest path in the graph from initial state to the goal state.

* In Uniform-cost we calculate the path cost from initial state to the next state.

* Example 2



Search tree (order doesn't matter)



- * Visited Nodes: A, B, D, C, E, G
- * Paths: A, B, D, G

* إذا كان node زيارتها ما ينحجبها.
 * إذا تكررت node لا نأخذها مرة بل نقارنها مع القديمة إذا كان Cost الجديدة أقل نأخذ القديم ونحذف الجديد. في هذه الحالة يتغير path.
 * فنأخذ الثاني لأن Cost خيار القديم.

- * Uniform-Cost-Search 3 - Complete: prevents the cycles.
- Optimal

* Question 1

While the Uniform-Cost is optimal, why we need the AI?

The answer is for speed.

Because the number of visited nodes in AI less than Uniform-Cost, in addition to the optimality.

So, In AI the speed is larger than Uniform-Cost, because in AI we take into consideration the path cost and expected cost.

- Space Complexity: exponential. $[b^{c \cdot e}]$
- Time complexity: exponential. $[b^{c \cdot e}]$

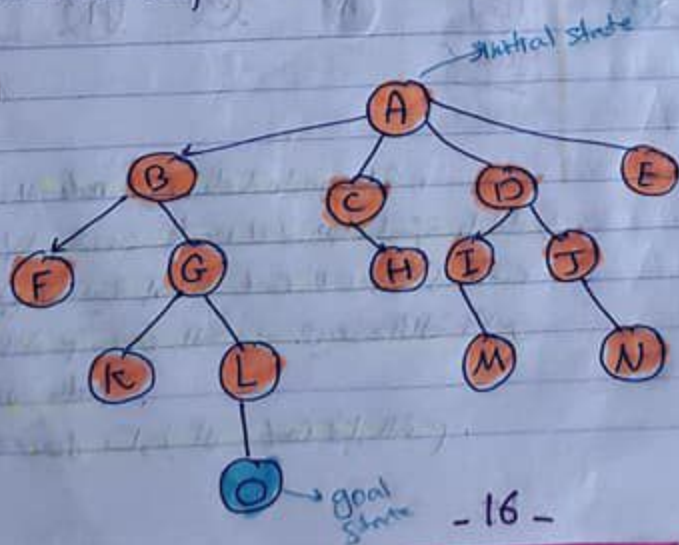
* Uniform-Cost is special case from Breadth First search.
When the Cost is equal then the two searches are the same.

* On the optimality side \Rightarrow Breadth is special case from Uniform-Cost.
Because Uniform-Cost is always optimal.

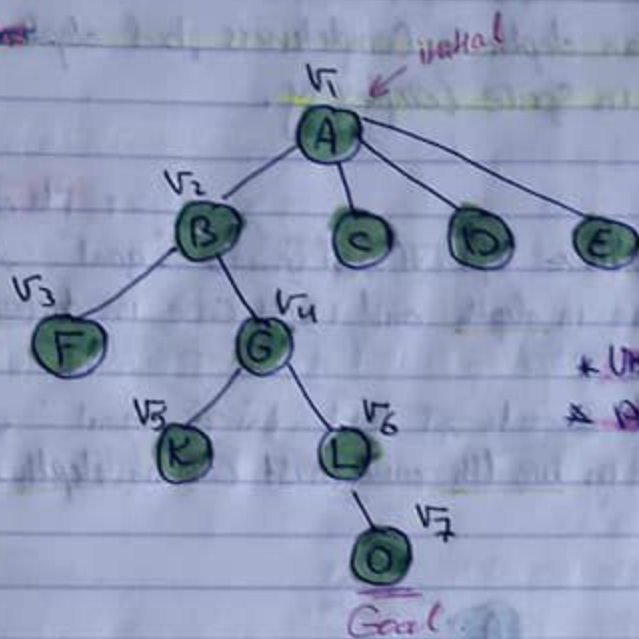
[3] * Depth-First Search

— Goes to the leaf node then do the recursive. (Back to the parent)

* Example 1



→ Search tree



* Visited nodes: A, B, F, G, K, L, O
 * Path: A, B, G, L, O

* The implementation of depth algorithm as data structure is "Stacks".

When we put the nodes on the stack we put nodes from write to the left, then when we visited nodes to the leaf node and it's not a goal we remove it from stack to go to the next node.

In previous example we put on the stack: E, D, C, B, G, F. When 'F' is not visited node we remove it to go to the 'G' node.

depth search is breadth search depth search. Space complexity is high. memory is high. nodes are removed from stack. memory is high.

depth search is not complete. cycle is not allowed. Complete search is better.

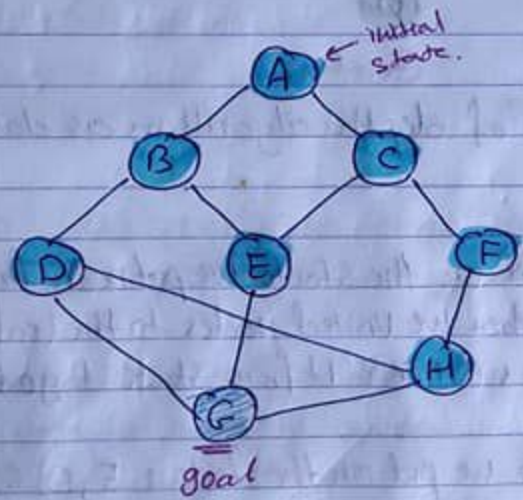
- * Space complexity: polynomial (b^m)
- * Time complexity: exponential (b^m)
- * Completeness: Not guarantee → not optimal.

* Breadth is better than depth in Completeness, but depth is better than breadth in space complexity.

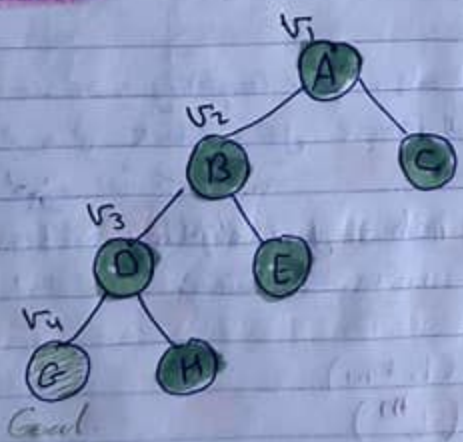
breadth vs depth \Leftarrow left \rightarrow right \rightarrow goal \sim 15/15
 This is the best case in depth and worst case in breadth.

depth vs breadth \Leftarrow node \rightarrow level \rightarrow goal \sim 15/15
 This is the best case in breadth and worst case in depth.

* Example:



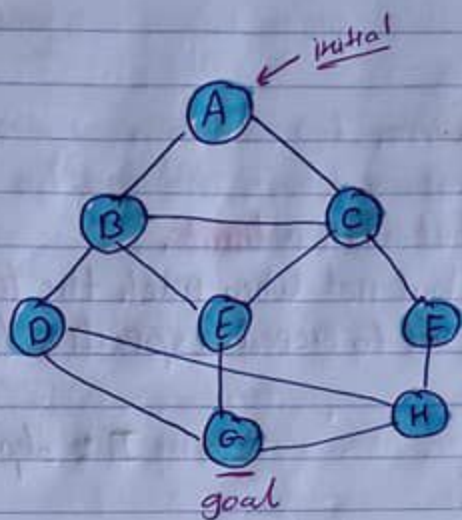
Search tree



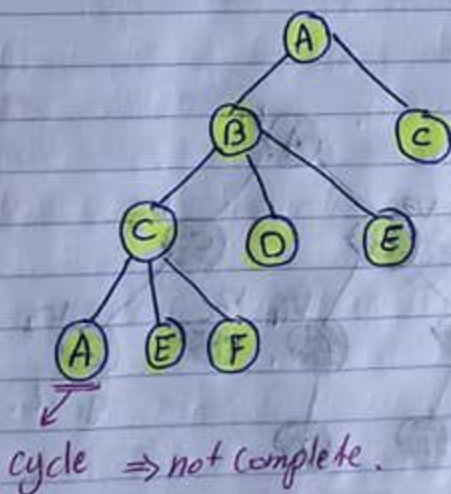
* Visited Nodes
 A, B, D, G.

* Path 1
 A, B, D, G.

* Example:



Search Tree.



L5: Uninformed (bfs + Informed (greedy + A*)
M/18/2/2019

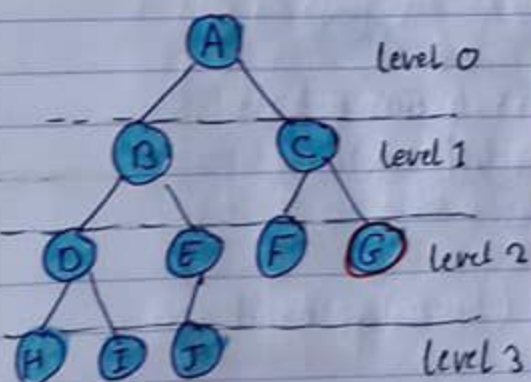
4) Depth-Limited Search

- Similar to depth-first, but with a limit.
- Depth-Limited do recursive but not when reach the leaf node, the recursive condition is not to exceed a specific limit.

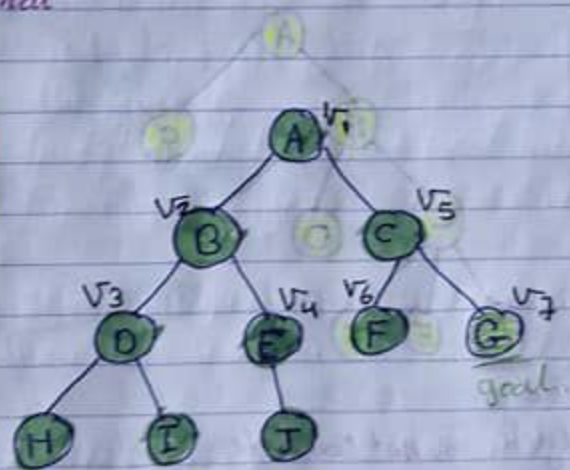
- Time Complexity: b^I
- Space Complexity: $b \cdot I$
- Completeness: Not complete
- Optimality: Not optimal

I : depth limit.

* Example:



limit = 2
goal = G



Search tree

- * Visited nodes: A, B, D, E, C, F, G
- * path: A, C, G

- Levels that larger than the limit, we don't visit them. So we conclude from the previous example that the depth-limited is working as depth- but with limit.

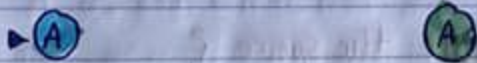
5 Iterative Deepening

- Applies Limited-depth with increasing depth limits (dynamic).
- Combines advantages of Breadth-first and depth-first methods.

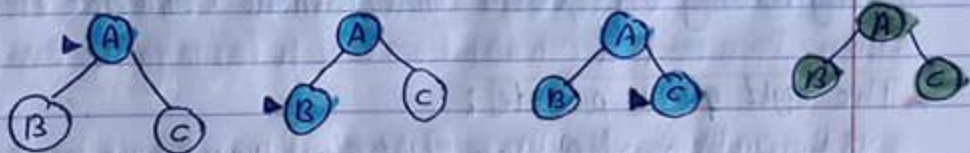
- Time Complexity: b^d
- Space Complexity: $b \times d$
- Completeness: Yes, it is complete
- Optimality: Yes, it is optimal

Example:

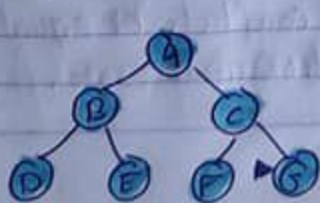
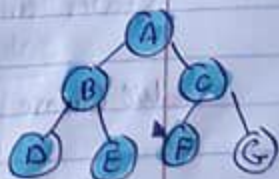
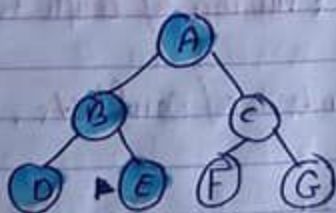
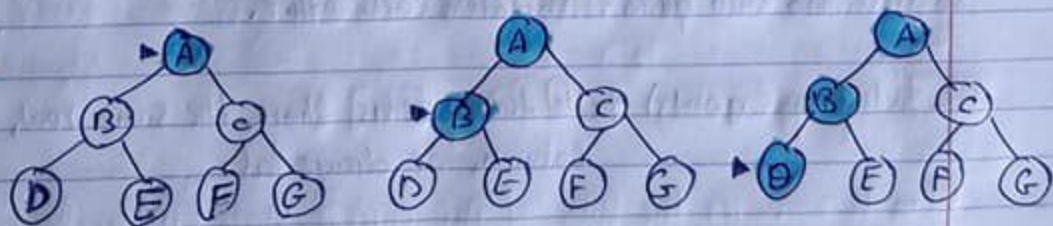
Limit = 0



Limit = 1



Limit = 2



The same way with limit 3.

6] Bi-directional Search.

- * Search Simultaneously from two directions.
forward from the initial and backward from the goal state.
- It is not important -

Informed Search.

- * Informed Search uses the knowledge to reduce the searching time.
- * Reduce the options.

+ In AI we care about Completeness and optimality.

How the knowledge reduce the space?

By removing the options that I know I will not need them.

The eight queens puzzle:

is the problem of placing eight chess queens on an 8x8 chess board so that no two queens threaten each other.

- Solution (goal): No two queens share the same row, column, or diagonal.

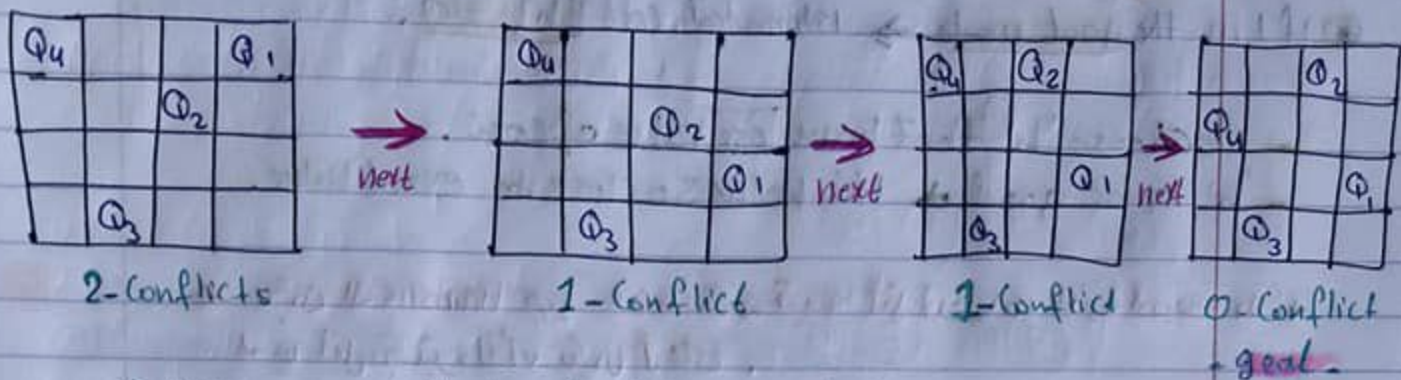
- The knowledge uses in this problem to compare between next states and choose the best choice is: number of conflicts.

• We choose next state that have less number of conflicts.

goal: number of conflicts equal zero.

next state: should have less number of conflicts or equal the number of conflicts in the current state.

Example: (4 queens)



- in 4x4 queen: it should not be any of the queens in the corners to solve the problem.

- I don't choose state worst than current state.

⊗ Informed Search may reach to the optimal solution or may not.

⊗ When we know more about the problem we solve it faster, because when we have enough knowledge we reduce the options of search.

⊗ The general name of knowledge in AI is - heuristic -

* heuristic: is function, with input: state, output: quantitative thing.

* Comparison → Subjective: two opposite answers.

→ Objective: one answer

is estimated cost

* heuristic: is function through it we make estimation for remaining cost.

* I have two types of cost

- path cost: exact.
- 'h': estimated cost from current state to the goal state.

- * if I in leaf node \rightarrow heuristic give me infinite output (∞). [$h = \infty$].
- * if I in the goal node \rightarrow estimated cost [h] = 0.

- we choose 'h' that have less value of cost. \leftarrow
- 'h' to be good \rightarrow It's value so close to exact value.

* يمكن لا heuristic جيد عندما يكون قيمة أقل زائد أي القيمة الأصغر ولكن ليس أقل بكثير ولا يكون أعلى.

- * If 'h' that I choose is closed to the exact that's mean the knowledge that I use is enough and not poor.

وقيمة 'h' قريبة من الـ exact value

التي نستخدمها في algorithm

- * Admissible : if h less or equal the exact cost (Real cost).
- * non admissible : if estimated cost of h large than exact cost.

- * if h admissible the solution is guaranteed to be optimal.
- * if h nonadmissible the solution will not be optimal.

Three approaches to defining [h] :-

- 1 - h measures the value of the current state (its "goodness"). This type of heuristic uses for optimization problems, since there is no reference solution.

* Optimization: لا يوجد الحل الأمثل الذي نستخدمه كمرجع

\rightarrow There is no optimal solution to be reference, so we used goodness for optimization problem.

- 2 - h measures the estimated cost of getting to the goal from the current state.

3) h measures the estimated cost of getting to the goal state from the current state and the cost of the existing path to it.
 path cost + estimated cost.

* ترتيب ال nodes حسب ال h value

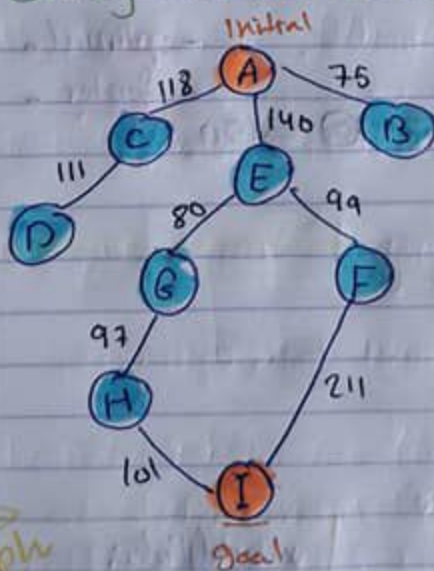
There is two algorithms in Informed Search :

1) greedy best-first search . 2) A^* search

* Now I know two things : path cost and estimated cost.
 we know that 'Uniform cost' use path cost.
 'greedy best first search' use estimated cost.
 ' A^* ' use path cost and estimated Cost.

Example :

1) * Greedy Best-First search . 2) * A^* . 3) * Uniform-Cost.



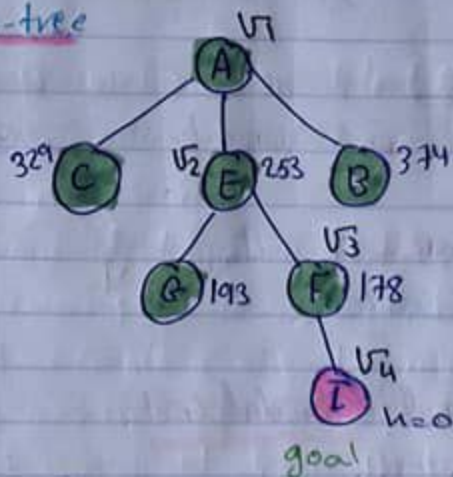
State	Heuristic (h(n))
A	366
B	374
C	329
D	244
E	253
F	178
G	193
H	98
I	0

graph

→ Solution

1 Greedy Search.

Search-tree

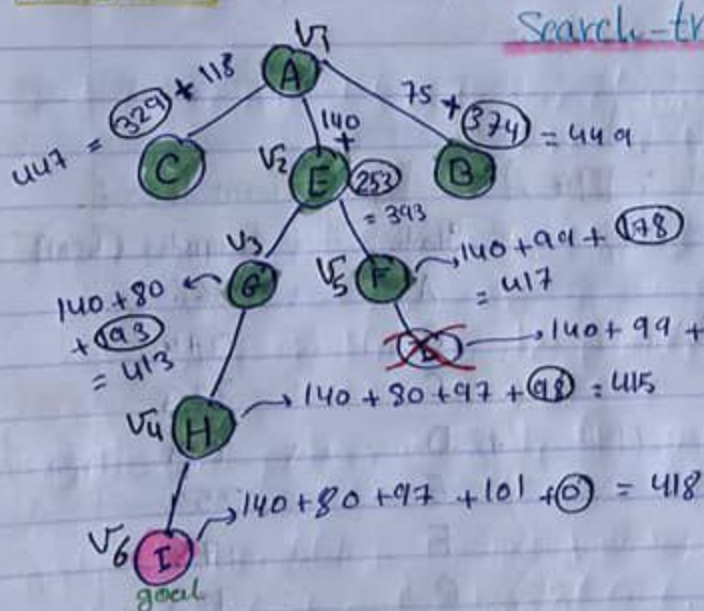


* Visited Nodes : A, E, F, I
* path : A, E, F, I.

- not optimal -

2 A* Search.

Search-tree



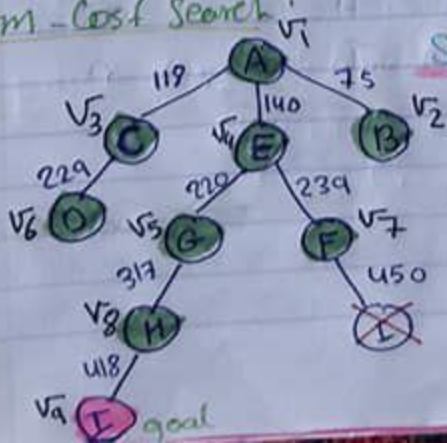
* Visited Nodes : A, E, G, H, F, I
* path : A, E, G, H, I.

- optimal -

- faster -

3 Uniform-Cost Search.

Search-tree



* Visited Nodes :
A, B, C, E, G, D, F, H, I.

* path :
A, E, G, H, I

- optimal -

L6: Examples on A* & greedy / IDA*

W/20-2-2019

إذا استخدمنا الخوارزميات الجشعة Algorithm الجشعة Solution الجشعة لا تضمن أن تكون الحل الأمثل Cost الجشعة لا تضمن أن تكون الحل الأمثل optimal الجشعة لا تضمن أن تكون الحل الأمثل optimal

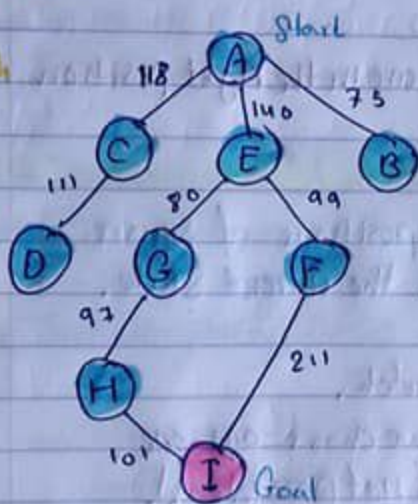
As we know 'Greedy-search' not guarantee to be optimal, but is it guaranteed to be Complete??

The answer is not guarantee to be complete.

Greedy Search

Example

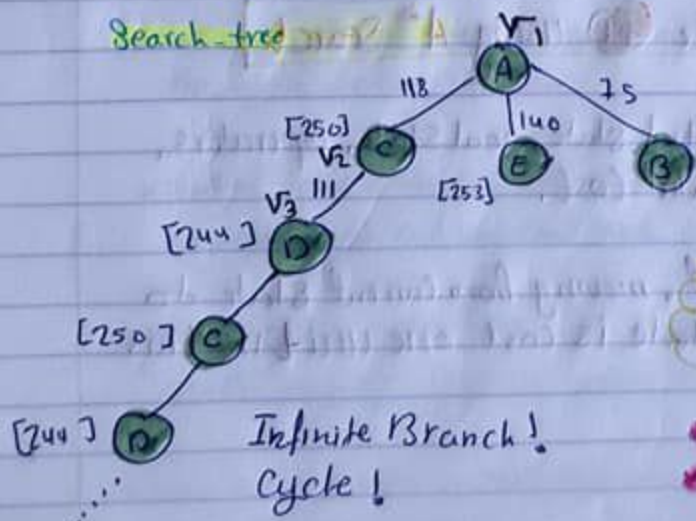
graph



State : Heuristic

A	: 366
B	: 374
C	: 250
D	: 244
E	: 253
F	: 178
G	: 193
H	: 98
I	: 0

Search tree



Data Structure is a priority queue.

[based on edge] add the queue then sort them on heuristic.

So, it's not guaranteed to be complete.

Time Complexity: b^m
Space Complexity: b^m [expanded]

• About **A* Search**: Data Structure is also "priority queue", but we sorted on [path cost + heuristic].

* There is an example on slide (26) [8-puzzle].

Some notes about the example.



on **greedy search**.

we should define initial state, goal state and operators.

operators: Up, down, left, Right.

but we can't apply 4 operators in each state, this is depends on the position of space.

and we should define the heuristic.

• The heuristic here is number of squares that are in the right position in the goal.

• to choose next state → number of wrong positions of squares should be less than the current state.

[Estimated] Cost: number of squares that are mismatch.

* if all next states have the same cost we choose any one.

• next state shouldn't be the parent, it's not possible !!

• The same problem on slide (49) using **A* Search**.

• we need in addition to initial state, goal state, operators, the path cost and estimated cost.

• path cost: displacement, moving from current state to the next state is cost one unit in this example.

if h_1, h_2 are admissible heuristic \Rightarrow

$h_3 = h_1 + h_2 \Rightarrow$ non admissible

$h_3 = h_1 - h_2 \Rightarrow$ non admissible [there is no cost in negative]

$h_3 = |h_1 - h_2| \Rightarrow$ admissible

$h_3 = \max(h_1, h_2) \Rightarrow$ admissible

$h_3 = \min(h_1, h_2) \Rightarrow$ admissible

admissible is better than nonadmissible, but if there's two admissible which one is better than another??

The answer is which is closer to the exact.

So, in previous example the best heuristic is $\max(h_1, h_2)$.

\rightarrow less than the exact and near to it,

** Improving A^* (Iterative deepening A^*):

* A^* complete, optimal, fast, but now we should improve the space complexity.

* Time complexity for all algorithm is exponential.

IDA* is optimal

To improve the space complexity we will combine between depth and A^* .

* We apply the iterative-deepening as a cut-off rather than the depth for the iteration.

* cut-off value is the smallest cost of any node.
we work depth but we don't reach the leaf, the limit is cut-off,

* we are going to the left, we do recursive (Back) if we reach cut-off (cost) larger than one I define.

* New cut-off value is level cost in generate nodes (not visited).
cut-off value = $\min(\text{cost of non visited nodes})$.

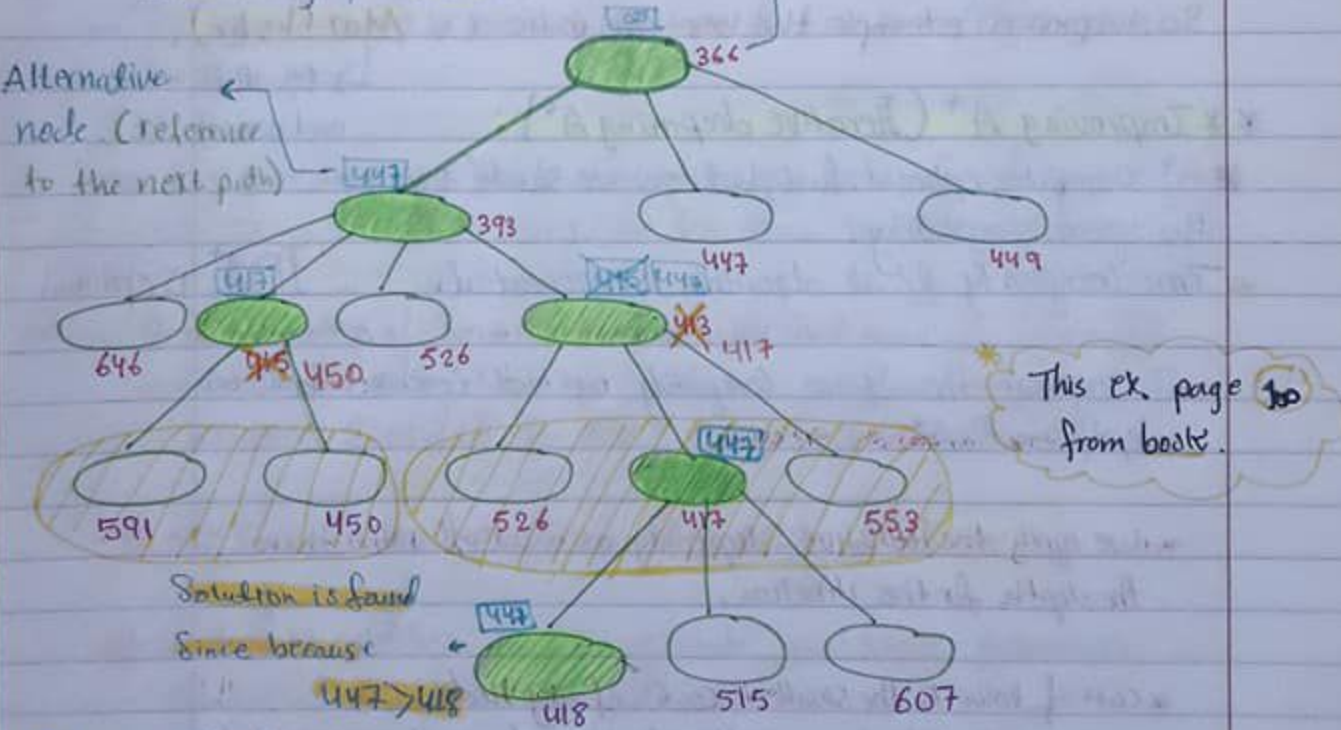
* Simple Recursive Best-First Search (RBFS).

- we will reduce the memory by:

1. If current f-values exceeds this alternative f-value then backtrack to alternative path.
2. Upon backtracking change f-value to best f-value of its children.

* RBFS, example.

- Routing problem -



* Some explanations:

At first, each node we visit it should have alternative path. How we choose the alternative path?? we choose it from the same level, if there is an alternative best than it we choose from the high levels.
 When we go to the alternative node?? if the current node worst than alternative we go to the alternative after update the current node and remove next generated nodes from memory.

The difference between A^* and RBFS, that A^* doesn't remove any nodes from the memory, but RBFS remove from memory when it moves to the alternative path.

RBFS is better about Space Complexity because it remove from memory. Specially when the branching factor is large. (compared with A^*). But IDA* is better than RBFS about Space Complexity, because IDA* works depth.

* Simplified Memory-Bounded A^* (SMA*)

Efficiency in memory \rightarrow SMA* is better than IDA* & RBFS. but it's more 'complicated' and expensive on time.

* How this algorithm works??

Bounded \rightarrow have specific number of states stay in the search space (search tree).

* At first, work A^* until reach to the limit, then we start remove from the tree, the worst nodes depends to the cost.

- in this way we keep two things, the first thing is the number of node don't exceed the bounded. the second thing is guarantee to reach to the optimal solution.

* when we reach the limit and remove the worst node from the tree we replace them by generated next states.

* *شبه A^* الذاكرة المقيدة SMA**

* SMA* is not guarantee to be complete if the number of nodes in the path greater than memory bound.

So, to reach to the solution the memory bound should be at least equal to the path.

* depth, breadth, Uniform Cost, Iterative deepening, greedy, A*
 IDA*, RBFS, SMA* \rightarrow Solution is the path

* Example 18

• Suppose you want to solve the following game, played with three digit numbers between 100 and 999. Let S and G be the start and goal numbers; a move consists in changing a digit in a number, by adding and subtracting 1. Hence, for example, a move changes 678 into 688 or changes 213 into 113. Moves, all of unitary cost one move one unit cost, are subject to the following constraints.

- 1 cannot be added to 9, nor subtracted from 0.
- Two consecutive moves cannot modify the digit in the same position. E.g, 313 \rightarrow 314 \rightarrow 315 is forbidden while 313 \rightarrow 314 \rightarrow 214 is allowed.
- Two consecutive moves cannot be both additions, even on different positions, E.g, 313 \rightarrow 314 \rightarrow 315 is forbidden while 313 \rightarrow 314 \rightarrow 304 214 is allowed.

1 Formalize the problem as a search problem (initial state, goal state, operators),

* initial state: S (100 - 999)

* goal state: G (100 - 999)

* operators: (add & sub) \rightarrow with constraints \rightarrow moves to the next state.

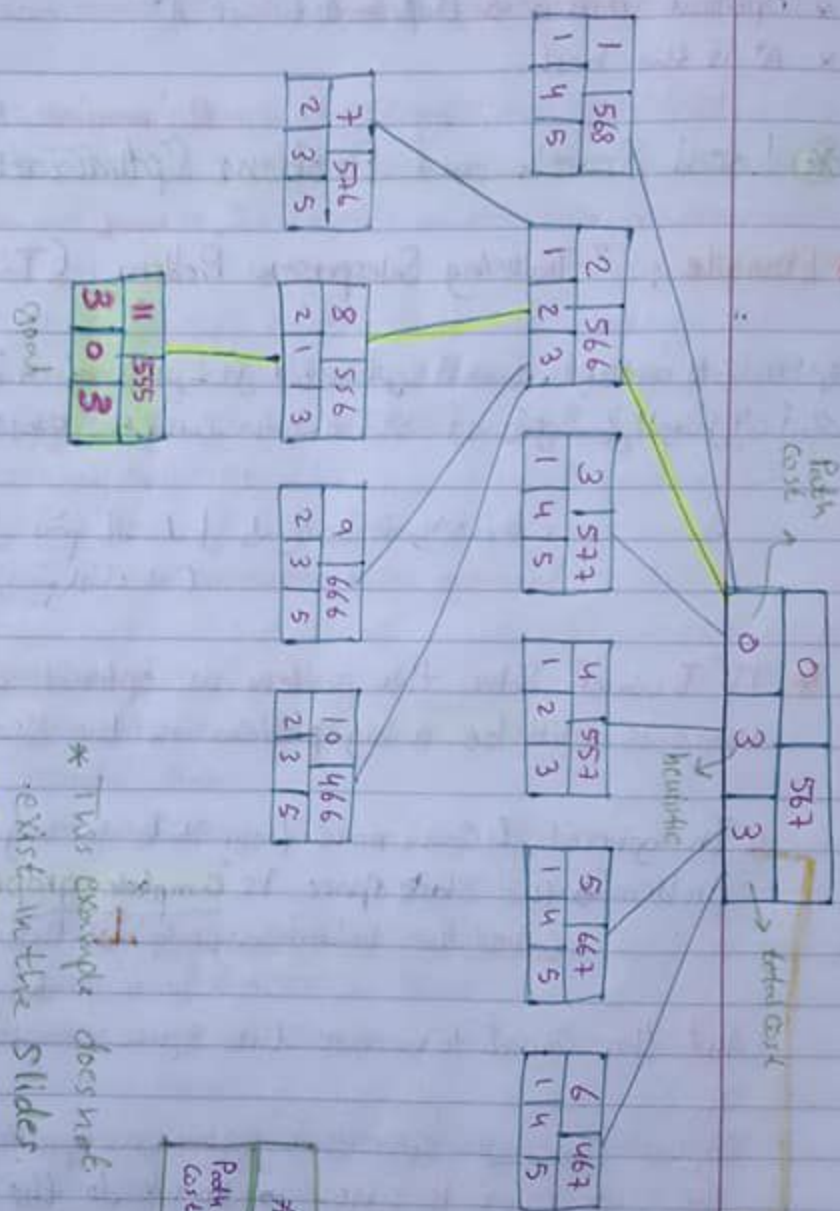
2 Give an admissible heuristic for this problem.

$$\begin{array}{ccc}
 S_1, S_2, S_3 & & G_1, G_2, G_3 \\
 S & \longrightarrow & G \\
 (100-999) & & (100-999)
 \end{array}$$

* $h = |S_1 - G_1| + |S_2 - G_2| + |S_3 - G_3|$

\hookrightarrow admissible heuristic function.

- 3 Draw the tree generated by A^* in the solution of the problem in the case $S = 567$, $G = 555$ up the reach of the goal node. Associate with each node of the tree: the state, f , g , and h , and a number indicating the order in which nodes are generated.



* This example does not exist in the slides.

* Solve it using depth, breadth & A^*

#	State	Path-Cost	heuristic	total cost
---	-------	-----------	-----------	------------

* If the solution is the path we can choose any one of previous algorithms. (9 algorithms)

* Optimal Solution \Rightarrow Uniform-cost or A^*

* A^* is the best.

* Local Search and Problem Optimization.

* Example: "Traveling Salesperson Problem" (TSP)

* اسم البرهان وطريقته توصيل الرسائل الى مكانها المخصص وبأقل التكاليف
الوقتية التي ابدت اسم البرهان فكل آفة نقل جميع الرسائل الى أماكنها.

* تسليم الرسائل يتم بالليل وليس نهاراً والهدف
(معالجة مشاكل التوزيع)

* If I want solve this problem as optimal solution using A^* , there is will be a big problem in the space. Why??

In general, I can move from state to any state, in TSP problem \rightarrow the state space is "complete graph" (fully connected graph) (direct links between any node) \rightarrow there is a path from any node to another one.

And this caused to increase the space complexity \rightarrow to have optimal solution

So we can not solve this problem as optimal solution because the state space is huge. \rightarrow So we do the optimization.

$$\# \text{ of Solutions} = \frac{(n-1)!}{2}$$

Informed Search

* Optimization: Improved existing solution but there is no optimality.

- any local search is complete, but space complexity is constant because there is no path in this problem.

LS: Local Search and problem optimization.

W/27-2-2019

- * Some problems, the path isn't important, we don't need it. In this case I can remove it from memory. So, the **Space Complexity will be order of constant**.
- * This type of Search is **Informed Search** (we still use the knowledge).
- * **Time complexity**, we can not know it. It depends on algorithm and the initial state.
- * **Completeness**, It has no meaning here \rightarrow because I work in these problem's optimization techniques. So, I have a solution and I want to improve it concerning the cost.
- * There is no optimality \rightarrow not guaranteed to be optimal.

* **Improving Search methods**

- make algorithms more efficient: **improve time & space**

- avoiding repeated states.
- utilizing memory efficiently.

- use additional knowledge about the problem.

- properties ("shape") of the search space.
 - more interesting areas are investigated first.
- pruning of irrelevant areas.
 - areas that are guaranteed not to contain a solution can be discarded.

(Impossible in this branch a solution exists)

* local search is **Informed Search** \rightarrow using **knowledge** (heuristic)

- ① goalness of state.
- ② estimated cost to reach the optimal solution
- ③ sum of path cost + heuristic

* Local Search use the goodness of state.

↳ It isn't heuristic.

- In optimization there is no definition for optimal solution.

- We look here to the quality of solution. (just stop)

It's possible to be error ratio.

* In 'TSP' problem → We will use the goodness: Total of all the distances that we crossed them until we reach to the our destination.

goodness here is knowledge but not heuristic → So there is no optimality.

'TSP' → return the total distances not the path

* Example 3, 'graph coloring'

- Starting with random coloring of nodes.

- Change color of one node to reduce the number of conflicts.

* We take care about → The country should be colored..

↳ The country should not have the same color with nearby country.

* another example VLSI *

* In genetic problems in general → use the result only, we don't care about experiences.

* In the previous algorithms *

- Formalization (initial, goal, operators).

- path cost / heuristic / path + heuristic.

- check the node if it's goal or not.

- if not → generate next states by using the operators.

* the common thing in these algorithms [a] → all of them save the searchspace in the memory to return the path.

* Local Search

• Key idea (surprisingly simple):

1. Select (random) initial state (generate an initial guess).
(any solution chosen randomly).

2. Make local modification to improve current state (evaluate current state and move to other states).

3. Repeat Step 2 until goal state found (or out of time).

Termination here is different \rightarrow have more than one condition.
(reach to the wanted requirement) / or there is no time

How to move from initial to the termination step??

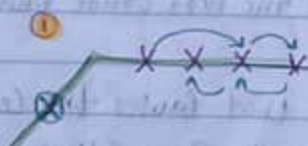
By set of action generate all possible next states, then choose the best next state. (less or equal current state).

then, after choose the next state, we will remove any previous states.

(It possible to enter in loop) \Rightarrow loop here means 'It doesn't improve itself'.

* If entered in loop \rightarrow for example: If all next state is in the same level, so it will go from state to another with same level, because it will not back to the worst state.

(Steady State)



error (local minimum)

improvement (profit) local max.

initial

(3)

(ridges)

local search used to
① improve time ② solve these problems.

local: because doesn't have global outlook

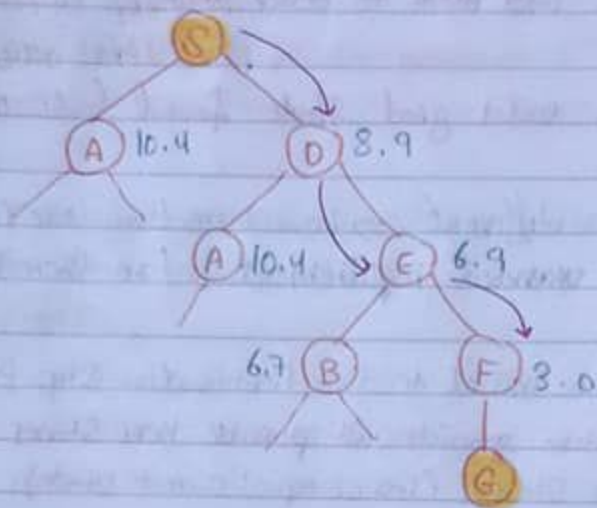
problems in local search:

- enter steady state
- local maximum
- ridges

II Hill Climbing Search

The general idea in this algorithm is local search.

Initial state \rightarrow termination test.



- 8-Queens example

- if I solve it using breadth algorithm \rightarrow in one state we have branching factor equal (8 * 56).

\rightarrow so this search space will reserve a huge size in memory.

- using local search \rightarrow will choose the next state that have less number of conflicts.

1 Change the position of queen that causes the conflict

2 when we change it, put it in the same column, because each queen will be in different column to avoid conflicts and reach to the solution. (we reduce branching factor from

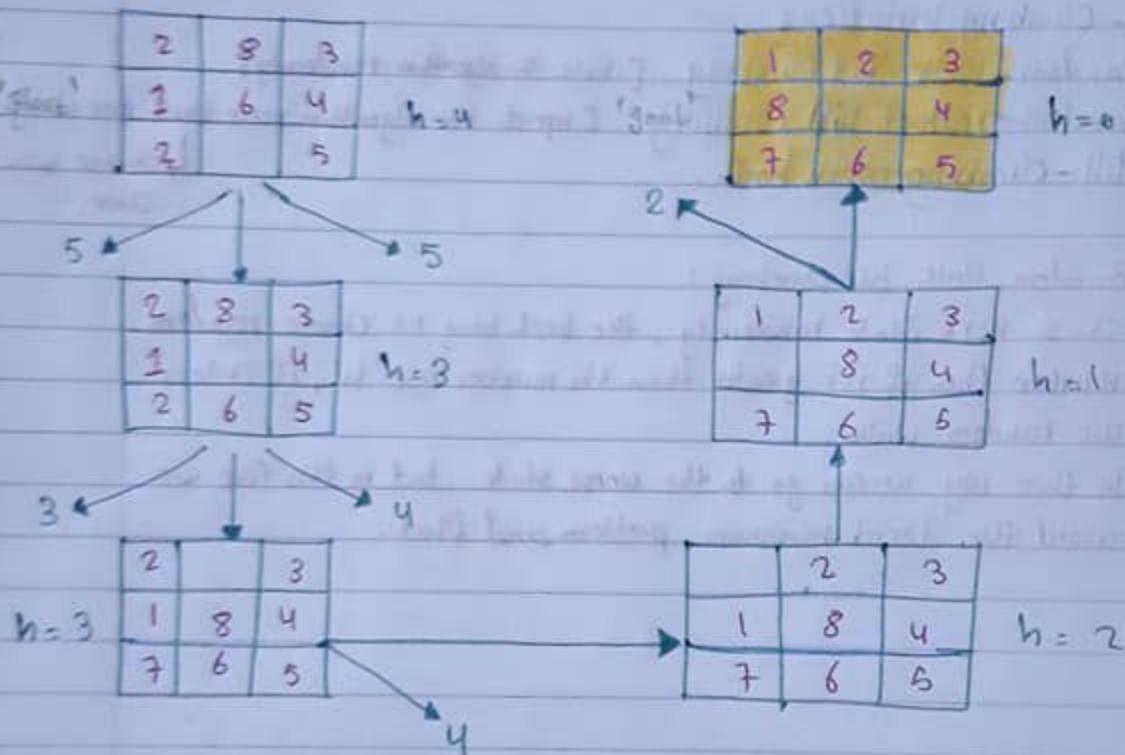
56 to 7). \rightarrow less memory.

- but when we change the position for queen we go to position with same conflicts or less than the previous, but we don't go to the more number of conflicts.

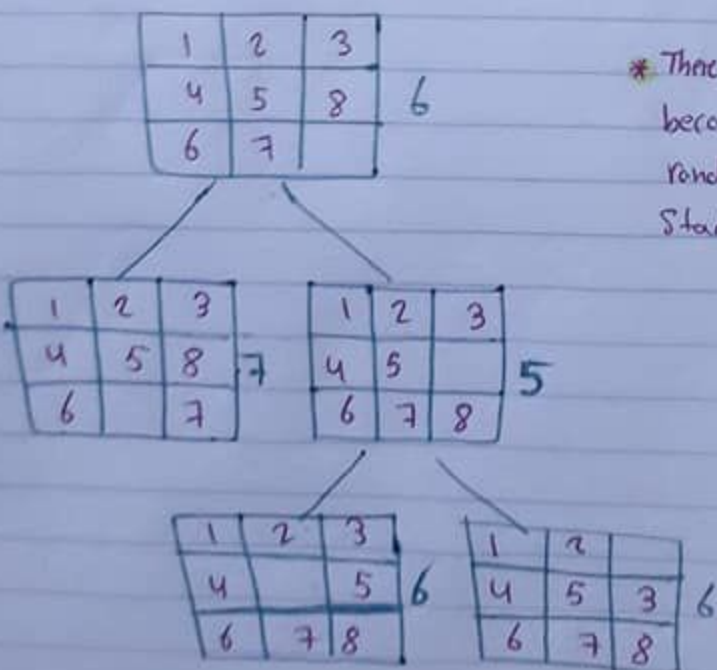
- If there is more than one queen causes the conflict, we change the one who causes more conflict.

• 8-Puzzle example

goalness : number of mismatch or conflict.



• Drawbacks of Hill Climbing.



* There is a problem in this example because the initial state chosen randomly. → Change the initial state. (Run randomly).
↳ Restart.

* Improve the hill-climbing search.

- Hill-Climbing Variations

- Random-walk hill climbing. [choose the state randomly]
- Random-restart hill climbing. [repeat the algorithm more than one time].
- Hill-climbing with both.

↳ change initial state.

- Random-walk hill climbing:

Choose next state randomly, the best way is choose random number then if it's greater than the number use "h", if it less use random walk.

- * In this way we can go to the worst state, but in this case we avoid the local minimum, problem and flat.

L9: Simulated Annealing & Population Based Algorithms (Local beam Search).

M/4-3-2019

2

* **Simulated Annealing**: It is hill-climbing search but it doesn't choose the next state as a greedy approach.

It has a function called 'Cost function', It choose the next state randomly, if it is better than the current state, the function jump to it, if it is worst, stay in the current state and then choose new random state.

[and the condition does not achieved]
↳ but if the condition is achieved \Rightarrow choose the state.

Before \Rightarrow next state have cost equal or less than current state.

But Simulated Annealing merge between greedy and random-walk.

The general idea is: if next state better than current state then go to it, else, go to it under a condition.

Specific Ratio \rightarrow greedy

Specific Ratio \rightarrow random-walk

* Condition:

$$R < e^{\frac{(-\Delta E)}{kT}}$$

* See Simulated Annealing Flow Charts - slide 37 - *

* If we replace the condition in Simulated annealing by threshold then, it becomes hill-climbing with random-walk.

So, the simulated annealing use cost function (condition) instead the threshold.

The expression (condition) $\left[e^{\frac{-\Delta E}{k_b T}} > R \right]$ have an important parameter and it 'T' \equiv temperature.

Why temperature ??

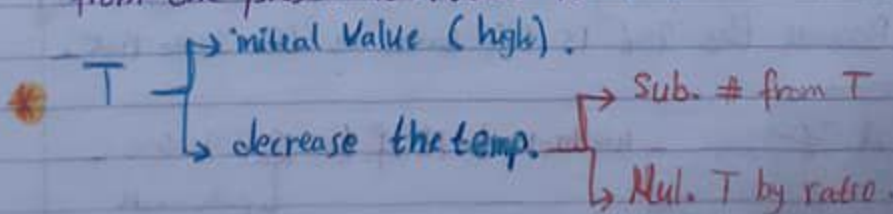
The idea came from 'cooling metals'.

The metals are usually to be solid, they are heated at high temperature such as glass. It becomes liquid at high temperatures, then cooled it gradually. At first, high leaps in cooling then the cooling slows down, so that the material becomes solid and its molecules bond together. (unimportant part :)).

* Initially, the value of T is large \Rightarrow the value in the above expression is small.

That is mean, if the next state is bad then it choose the state, but by the time, the probability of receiving a bad state is lower.

* If I want to apply the simulated annealing on a problem, the expression doesn't change. The only thing you specify from one problem to another is 'T'.



* Simulated annealing is local search and local search is informed search. So, before we determine the initial value of T and the scaling we should determine the heuristic.

* In hill climbing, the heuristic is enough to choose the next state, but in simulated annealing we need to edit the cost function.

T : initialization & scaling.

- simulated annealing solve the local minimum problem, because, it accepts bad states.

* Example: (TSP) - Traveling Salesman Problem -

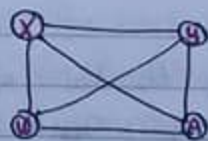
- we choose this problem to apply the local search, because we don't need the path and the search space is very large.

* Formalization:

- initial state: any route from $\frac{C_{n-1}}{2}$. (any random solution)

- How we implement the initial state??
implement it as a sequence.

→ I W Y A X ← initial state.



fully connected graph.

- Goal state: improve the initial solution to reach to the solution with lowest possible cost. [cost \equiv total distances]

- operators: interchange between nodes.

[but we don't change nodes that is next to each other]
- Because the cost is the same in the same link.

→ I W X A Y ←
fixed.

- we move to it if the cost is less → path length.

Tabu Search \Rightarrow (تجربی)

* Population Based Algorithms (Local Search Algorithms)

- Local Beam Search - Genetic Algorithms.

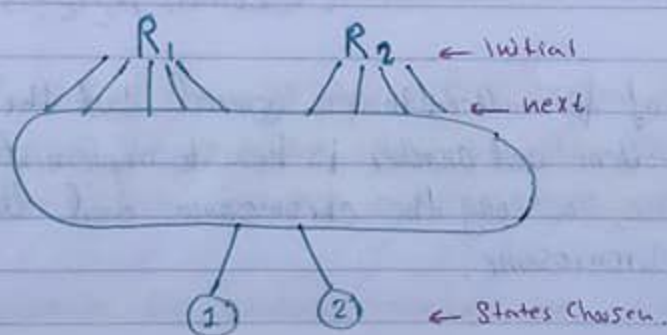
In these algorithm we will start with two initial states. but in hill-climbing and simulated annealing we start with one initial state.

Local Beam Search.

At first generate 'k' initial states.

'TSP' problem \rightarrow assume $k=2$, start from 2 routes.

then generate next state from both initial states, then choose best two from all next states.



Using 'k' states help us to make expand our search space.

It's look like see the fully search space.

In this way we reduce the local minimum problem.

- By visiting the states in many regions in search space -

L10: Genetic Algorithms

W/ 6-3-2019

2 Genetic Algorithms

- The idea came from the genetics in the biology.
- The genetics deal with more than one solution.
- The general idea is mating between two species with good qualities so that the children inherit good qualities.

we should know how to do the formalization on the problem as genetic.

* الجين يحمل الصفات، والبيرن موجود على كروموسوم، والآخر موجود على DNA، فنجد عمل تزاوج بين كروموسوم وكروموسوم آخر نحصل على جين أفضل من أبوين في حال أن الآباء يحملون صفات جيدة. والبيرن يرث من أبى الصفة.

* In genetics, the procedure of formalization is general, but the difference between one problem and another is how to implement the problem as chromosome, how long the chromosome and the form of the gene in the chromosome.

* مثال عملي: 'knapsack problem'

في السفر يجب التقيد بوزن محدود لايتحمله حماره. فهنا يجب الموازنة بين وزنة الشيء وقيمته. لذلك لتقدير الأختيار التي ستؤخذ تأخذ بعين الاعتبار الوزن والقيمة. فنحن نريد ان نأخذ ما لا يوجد 'optimal' لأنها تكون غير صالحة في آخر مصفاه الأختيار بعد الحلول لذلك نقدر عليها 'optimization'

ومثال ذلك أيضاً 'Summary' ليس لها حل optimal.

AI genetic algorithm, At first should determine the chromosome of problem (length of chromosome in terms number of genes, the implementation of genes) and the order, if we read the chromosome from left

to right or from right to left, does it make sense or not?

* في مسألة السفر ال order لا يهم ولكن في مسألة TSP - ال order مهم .

* The chromosome is the most important thing in formalization.

* مثال السفر ← ال chromosome بعدد الجينات فالجينات هنا هو عدد الأتجاهات ، وكل اتجاه إما أن يكون موجود أو لا يكون لذلك كل جين نمثله ب '0' أو '1' .
- طول ال chromosome هو عدد الجينات ونقل عدد الأتجاهات .

عدد الجينات 7 →

0	1	1	0	1	1	0
---	---	---	---	---	---	---

نرى ال مثال فثبت مكان الجين لأنه ال order ليس له معنى فنأخذ عن واحد
وعلم وعلوه ب '0' أو '1' .
إما مسألة TSP والترتيب مهم لذلك لا نثبت مكان الجين .

- ال chromosome يمثل 'Solutions' ، إذا وجد chromosome آخر فهذا يعني حل آخر .

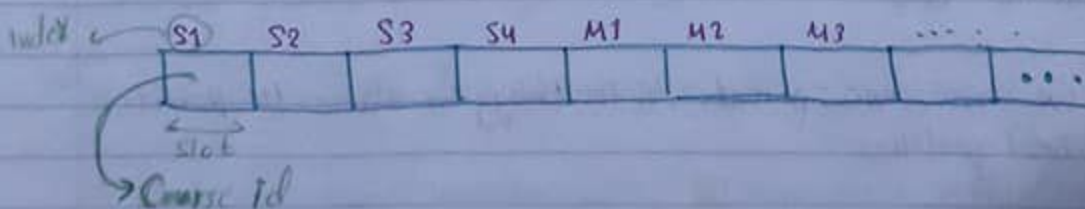
** TSP → # of genes is implement as number of cities (1) .
give each city a symbol like character (Unique code) → back to initial
the difference between one chromosome and another is the order.

** Optimize the Student Schedule ** (Example)

1 length of chromosome is (number of hours per day * # of days).
look

So the genes is number of hours work in all days .
each (slot) is hour with specific index.

2 We represent the gene by the Course ID or Course Name .
The constraint is prevent two courses be in the same time .



* Summarization problem *

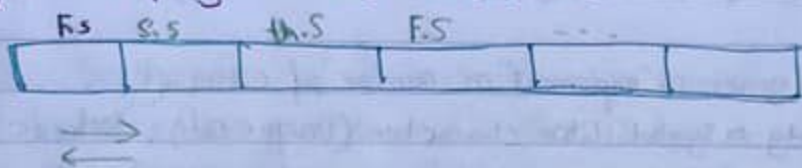
- 1 Length of chromosome is # of Sentences in the paragraph.
that's mean each gene is Sentences.

* Summary

- abstractive → using your own language.
- **Extractive** → chose Sentences from the paragraph.

- 2 Each gene is implemented by zero or one. That is mean either chose the sentence or not. Also, that's mean the position for each gene is fixed will not change.

The order is not important. [read from left to right or from right to left] in the chromosome.



* When the implementation is '0' or '1' ⇒ the order isn't important.

* How we chose the first two chromosomes (solutions) ??

- Randomly

* To generate a chromosome ^{with} type of implementation is binary, is use random function for each index in the chromosome.
if $\text{rand}() > 0.5$ put the value '1', if $\text{rand}() < 0.5$ put the value to '0'.

- At first we chose two parents, in the biology we chose the parents with the best qualities

1) * The first step in genetic is generation process by create set of chromosomes randomly. "Populations"

2) * Then choose the best two chromosomes to mating between them.
- choose these two parents depends on function called 'Fitness function'.

a. Fitness function: is evaluation function returns a value indicates to the best chromosome.

↳ Fitness function is like heuristic.

we didn't give a fitness function the heuristic name because the fitness function measures the quality of chromosome, not the cost to the goal.

There's no reference in genetic because the genetic is optimization techniques.

* The fitness function in the TSP problem is total distance.

* The fitness function in knapsack problem is the benefit (Value)

* The fitness function in Summmary problem is the features.

↳ what we meant about the features → for example:

- The sentence like the title is important,

- the sentence contains dates and numbers is important,

- the position of sentences: sentences in the conclusion are important, and the first sentence also it's important.

- the sentence content: (Keywords) also important.

- the sentence that start with 'For Example' is not important because it's a detailed sentence.

3) * The mating between the parents is called "Crossover"

Crossover means interchange. (تبادل)

- The crossover between two parent gives two children!

* The question is:

1) Which genes should interchanged?

2) How much?

There is no definitive answer.

- In 'Summary problem' we do the crossover randomly using random function that choose the index from the two chromosomes and interchange them.
- The crossover in knapsack problem depends on Value. We look up to reach to the solution with maximum value and not exceeding the required weight.
- The crossover in TSP problem is change the order of cities.

4) * we ask ourselves, the chromosomes that is produced from crossover is their fitness function achieves the requirement that is required? or repeats the previous steps again?

* Any optimization techniques, its termination is:

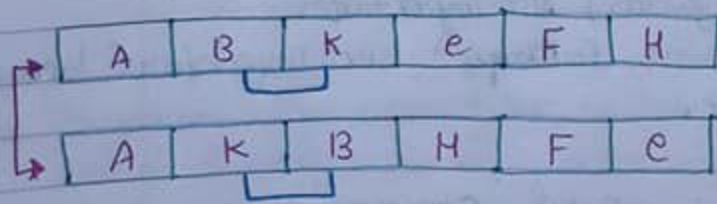
- time out
- reached to the requirement needed.
- There is no improvement

* in hill-climbing and simulated annealing we use the random-walk, in genetic we will use the 'mutation' *طرق* to exit the steady state.

- Mutation means: way using to exit the non-optimization mode.

- if the implementation of chromosome is binary, choose any gene in the chromosome and reflect its value \Rightarrow this make mutation, means exit the steady state.

* Example on crossover:



A K B e F H
A B K H F e

* we should avoid the repetition in the same chromosome.

* each problem has own rules should keep them.

L11: Constraint Satisfaction Problems (CSP)

M/11-3-2019

* Formalization :

given a set of variables, each variable takes at least one value.
Since, this assignment must satisfy all constraints.

* Example : Map-Coloring (Slide 3+4)

- Constraints : adjacent regions must have different colors.

- Given initially [Red, Green, Blue] \leftarrow Domains (Values).

- The regions (Countries) \leftarrow (Variables).

- The solution of this problem is give each variable a value with no violation of the constraints.

Representation :

* How to represent the CSP problem ??

As a graph, each node in the graph represents the variables, each link between node and another represents the constraints between variables.

* in the previous example, we can represent the cities as nodes and the constraints between cities as links between nodes, but the cities that no constraints between them, we don't put the link between them.

- CSP is Incremental Search Problem -

* initial state : Unassigned Variables.

* goal state : All variables have values assigned based on constraints.

* operators (successor function): the assignment.

* In general : The CSP solution is incremental approach.

that is mean, we choose a variable and make assignment then go to the next variable until reach the goal state.

* Depth is the best algorithm for the CSP. Why?

because the depth do 'Backtracking', but in these problems it doesn't reach the leaf. It make backtracking when there is a problem in the constraints.

We know that there's a set of variable and each variable have a set of values and when the variable is assigned we don't go to the next variable before make sure that the constraint is true or not. if true then we continue to the next variable but if not true, we do the backtracking. Also, the backtracking starts with the nearest variable.

* For example: if there is an problem in constraints in V_3 , we do backtracking in V_3 by choose another values. If the problem is still, we back to V_2 and so on.

So, we do the backtracking in two cases \Rightarrow

1. When constraints are violated.

2. reach to the variable doesn't have any values that I can choose.

* So, the algorithm that will use is 'Depth With Backtracking'.

* The generated search space is very very large.

$$[N! \times D^N] \leftarrow \text{Space Complexity}$$

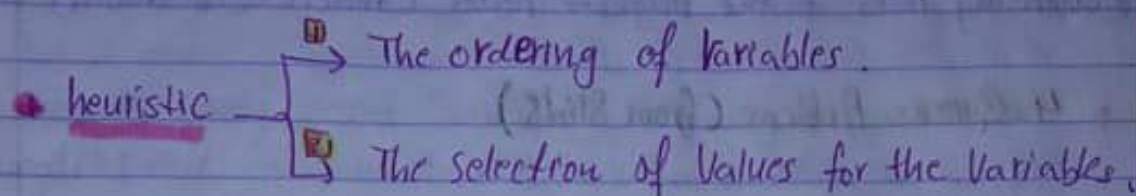
* There is no cycle problem in this algorithm, because it is incremental approach. So, it is impossible to assign a variable that assigned before.

— There is no optimality in CSP —

- The order of variable to assigned is very important.
- The assignment of values is important.

We should improve the time complexity by reduce the backtracking.

- We should use the heuristic to reduce the backtracking, but how??
 - By choose the best order of variable assigned and the selection of values is also important by select the values that reduce the conflicts.



①

- Most - Constrained Variable (minimum remaining values)
 - Variable with the fewest possible values is selected to minimize the branching factor.

- Most - Constraining Variable
 - Variable with the largest number of constraints on other unassigned variables.
 - In the graph choose node with largest number of links.

②

- Least - Constraining Value
 - For a selected variable, choose the value that leaves more freedom for future choices.

• See examples on the slides about each type of heuristic

* You can choose one or more of heuristic, but how??
 At first we choose Most - Constraining Variable, then if two nodes are the same, we use minimum remaining values.
 You are free about choosing the heuristic, you can choose one or two or three or nothing.

- * We know that to reduce backtracking, we remove all values from the neighbor nodes that cause the conflicts with the assignment after assign the current node. and this step is called 'Forward Checking'.
- Forward \Rightarrow because I look to the nodes that unassigned.
- Forward checking looks to the neighbor nodes (there is a constraints) only.

Example : 4-Queens Problem (from Slides)

- * The Variables : Queens
- * The Values : Board.
- we assume that initially each queen should be in different column. So, each queen has 4-choices on the board.
- [we have 4-queen variables with 4-values for each one]
- graph will be fully connected graph.

- You can see the details from slides -

L12: Arc Consistency (Adversarial Search & Games)

W/13-3-2019

Heuristics and Forward checking See the neighbor nodes only. So, we can not know if there is will be a backtracking or not in the future nodes.

But there is a way called 'Arc Consistency' by using it we can expected if this assignment will cause a backtracking or not.

What is the arc consistency?

If there is two variables $[X, Y]$ and there is a link (Constraints) between them and each variable have list of values.

The Constraint is called 'Consistent' if there is a value of Y after X is assigned.

- Coloring Map example on slide -

* Arc Consistency reduces the backtracking. \Rightarrow preprocessing of assignment

1 At first we apply the forward checking (if node is assigned then its value remove from neighbor nodes).

2 After that, we apply the arc consistency on nodes that is lost some of their values. By checking the consistency on nodes that lost their values with their arcs. How we check the consistency??

We see if all values of the node is assigned one by one and the node that linked with it have values then these nodes are consistent.

arc consistency is stopped if there is no update on variables. update & node lost a value.

Heuristics + Forward checking + arc consistency \Rightarrow reduce time complexity.

- There is an example on the slides -

Adversarial Search & Games.

The Searching techniques in games are separated from the other Searching algorithms and this is due to games qualities.

- * for example:
- ① The search space of games is very large. So, we need algorithms that are more optimize in Complexity.
 - ② In games there are two players (agents) that are play in opposite of each other. but in the previous algorithm there are one agent searches for goal.
 - ③ In the environment, there is a chance and there is nothing fixed.

* When a State is excluded then we already excludes all next States from this state. (pruning)

- Two agents.
- Profit and loss.
- Efficiency in memory (pruning)

- Formalization :

- 1 Initial state depends on game.

- 2 Successor function : legal moves depends on game.

- 3 Utility function : produces a value for terminal state.

Ex: 'Chess' → outcome : Win / Loss / draw with values

$+1$, -1 , 0 → Draw
The winner is the first player The winner is the second player.

* There is no optimal solution in games. The goal is to Win. also some games to reach the goal you should get the highest points.

* Search Space as we learned before.

- we should think in
 - Current state
 - actions
 - Utility function (heuristic)

● MiniMax Algorithm.

- The first one choose the MAX and the second one choose the Min of first's MAX.

* The different between one player and another is the perspicacity (رؤية النظر)

L13: MiniMax & Alpha-Beta.

M/18-3-2019.

MiniMax algorithm is applied in two players like chess, checkers, tic-tac-toe, etc.

* In MiniMax algorithm there is some limit, we choose the Min & MAX from this limit then back to the beginning. Then, choose the path that you put in your mind. Sure, we assume that the two players are in the same thinking.

• Terminal nodes: values calculated from the utility function.

• Other nodes: values calculated via minimax algorithm.

- So, in MiniMax we can know the maximum benefit, before start the game.

At first, we put a plan and from this plan we know the sequence of actions. from this we get the path.

* The Search algorithm that use in MiniMax is 'Depth-First Search' \Rightarrow space: polynomial, time: exponential.

• We start with calculate the utility function.

• Then use the depth Search.

• We can not visit all nodes, because in some games there is a constraint on time and we know the search space of games is very large although we choose the depth search but the time is exponential.

⇒ We should make a pruning by delete some parts of Search Space to reduce the space & time.

■ If guarantee that there is no solution from some parts of search space we remove the branch.

In this way I reduce the time and space.

Q: But how I know that there is no solution from this branch ??

The answer is by Utility function.

✶ The Utility function is 'Alpha-Beta Pruning'

* Alpha: Max player. * Beta: Min player.

1. The pruning happens between Max and Min.

Max makes a pruning for Min or Min do the pruning for Max, but there is no pruning between Max & Max or between Min & Min.

Minimax & Alpha-Beta return the same solution.

2. The decision of pruning comes from high level.

3. The node can make the pruning if it has a value.

[3 bits]

4. Checks the answer by compare between Minimax and Alpha-Beta solutions. Their should be the same but the difference is time & space complexity.

— See examples on Minimax & Alpha Beta from slides —

L14: Introduction on machine learning

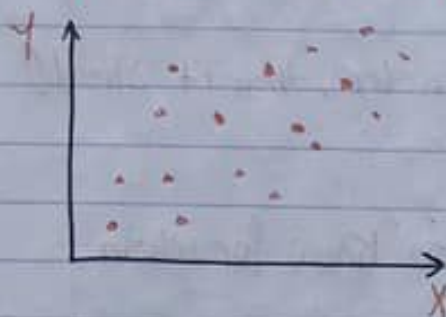
W/20-3-2019

- * Machine learning: Is a set of techniques that are classified to the three types: 1) Supervised 2) unsupervised 3) reinforcement

→ back to the Statistical → 'regression' or 'Interpolation' in numerical

* take the simplest case on regression → 'Linear regression'

→ input (X), output (Y) → $Y = aX + b$



→ these values of X & Y we get them from the experiment.

- * Why we do the regression if we have the value of X & Y??

→ to find the relation between the input & the output, but why I find the relation??

→ To avoid the repetition of the experiment.

So, we find the mathematical model between X & Y as a result from experiment to use it only for the next time and don't repeat the experiment.

- * Mathematical Models → Linear, exponential, logarithm, polynomial, ...

→ The first task in regression is determine the type of relation between X & Y (mathematical models) after get the values from experiment.

→ Find the exact relation between X & Y in Linear relation that means find a & b.

- * Learning is a more comprehensive concept of regression. by learning we find a model → mathematical, probabilistic, set of rules or Networks. → All these model aim to find the relation between X & Y by previous experiments.

→ The input of machine learning is the previous experiments.
I provide the data for the machine, then it learn the relation between the input and output.
After learn the relation, I can in the future use it in new things.

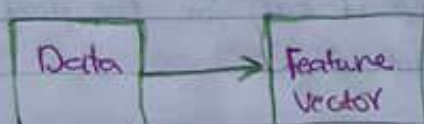
① * For example in email there is a spam and ham emails. the input of machine learning is the ^{text} messages. I can put with message an attachment or addresses (sender or receiver), images, links...
* So, input: message + meta data.

② * another example: give it a picture of a character, then it should know what is the character. ↳ input (image)

③ * another example: give it a ^{record} voice (wave), then know for whom. ↳ input

- In previous example the machine learning find the relation to make a prediction for your aim.

- The first step in machine learning is we need a data and these data are given from previous experiments.



machine learning is independent on input.

* The reason that the machine learning is independent on input because we convert the input to the formula that all machines learning understand it. and this formula is the feature vector.

→ Feature vector is an array of one dimension and this array is representation of all inputs in data.
each sample in input is converted to feature vector. and this vector contains many types of data maybe integer, string, ...

→ In the second example, If we want to write a program that its input is a picture of a character and the output is the type of the character. What should I think about? and how?

A B C D

→ Build it as set of rules, and these rules contains set of features for the character.

Some of features: Holes, Black to white ratio, number of straight lines, Angles,

to compare between one thing and another we based on a set of features.

— The basis of machine learning is the features, and these features is the input of machine learning not the row data because the machine learning doesn't understand the row data. So, we should convert the row data to set of features.

→ Row Data: pictures, text file, audio file, video file,

→ Process Data: is the data after modification.

— The first step to build the feature vector is think about the qualities of picture. [for each sample]

Because each sample should convert to features and the features are fixed on all samples.

Features (Columns)

Data Set

(rows) Samples	# Lines	Holes	B/W	# Angls	output
Value → 3	2	1	0.3	3	A
		2	0.35	4	B

* feature vector represents one sample so feature vector is 'row'

| ← X → | | ← Y → |

* Build two for loops → first loop for each sample (image). the second loop for features.

for
for

- Collecting the data.

- Convert the row data to features vectors } Your responsibility

↳ Features & Values.

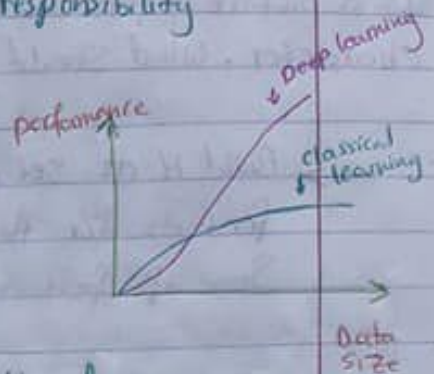
* Data Set is the input of machine learning.

table

↳ classical machine learning

* machine learning

↳ Deep learning.



- classical machine learning: you determine the features.

- Deep learning: input is row data not the feature vectors.
there is a tool build the feature vectors.

* Why we still work on classical machine learning while the deep learning is less effort.

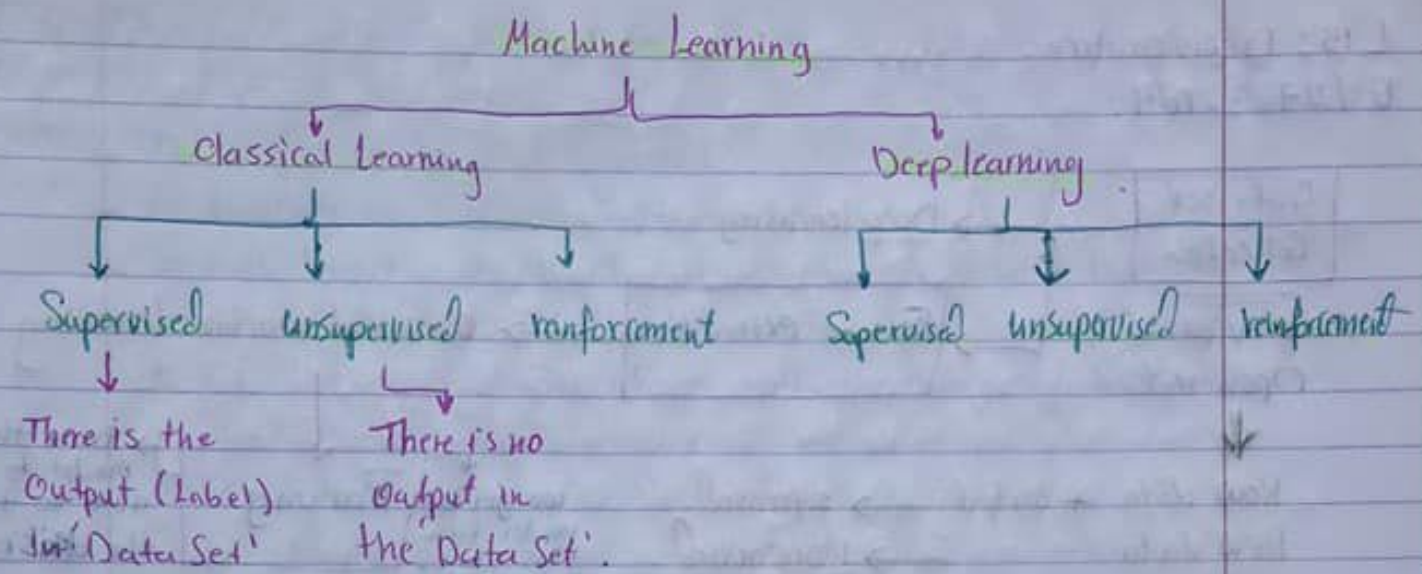
* Because you need more knowledge about deep learning because it's a lot. Also, as a technique it's harder than machine learning. Also, deep learning need a huge size of data, and you can't run the deep learning in your PC, because it's need a lot of computational. usually we run the deep learning on GPUs.

* On large data set → deep learning is good.
* On small data set → classical learning is better.

* In the course, we will talk about classical learning only.

- feature vector contain the values of features.

* The aim of machine learning is find the relation between input & output (your aim).



→ Why we sometimes put the output and sometimes we don't put it ??
 It depends to the size of data and the changes in data.

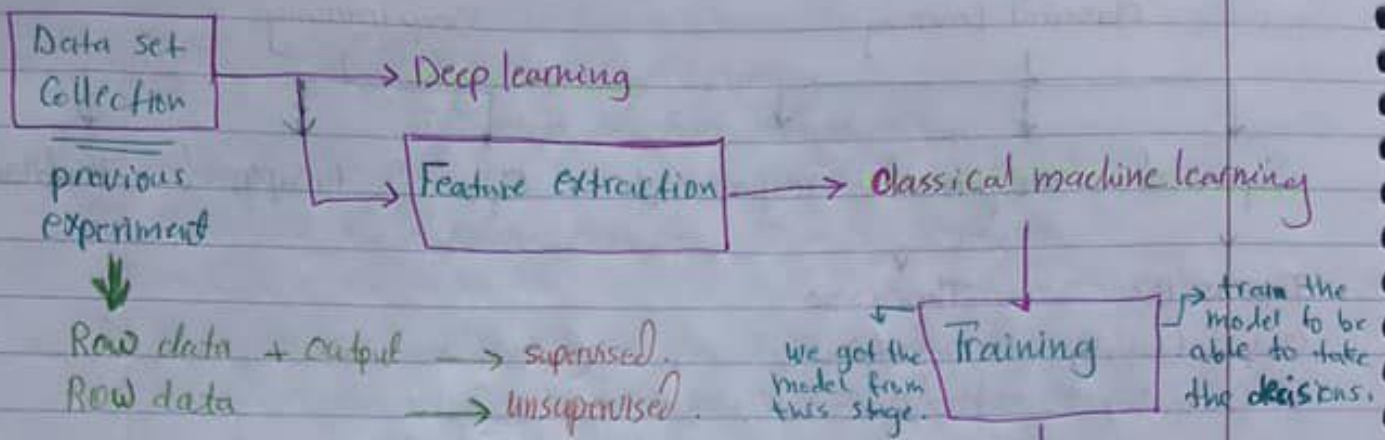
→ in unsupervised the relation doesn't between input & output because there is no column of output. So, we make the data as a groups or clusters since the elements in the same clusters are similar.

* reinforcement : There is two inputs (feature vector, feedback).

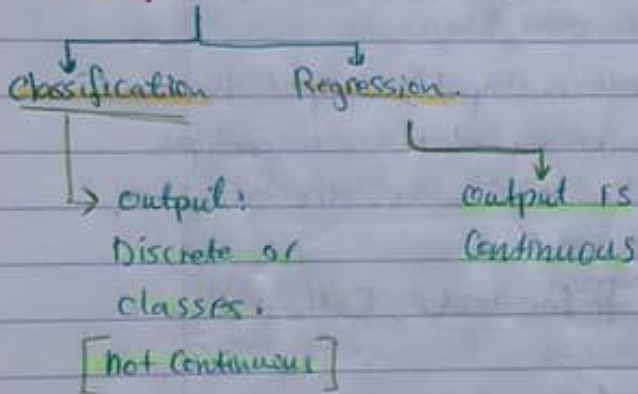
→ Feature vector is application dependent.

in classical machine learning, the feature extraction stage is very critical and it's person's responsibility.

L15: Decision tree.
W/27-3-2019



* Supervised :

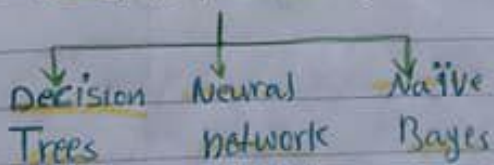


to test the model and know the quality of model.

there is a part of data doesn't enter in training stage.

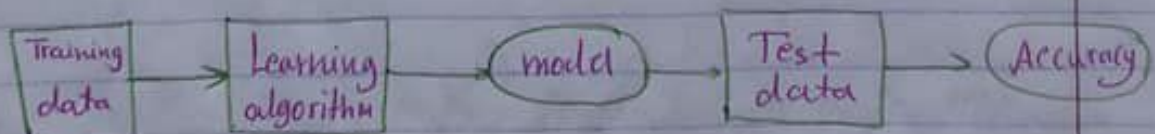
* Not all features have the same importance.

* Classification :



I divide the data into two groups. First part for training & the second part for testing.
70% 30%

* The Decision tree that I will get from learning is a set of rules [set of if-statements].



* **Decision Tree** → output is tree.

- The output is in leaf. (classes)
- It basically is 'if-else' statements.
- In each level we use the feature of each sample with values on tree.

* If the output doesn't change with the change of feature, then I don't need another feature.

↓
+ for example :

if own-house 'True' → output 'Yes' → don't need another feature.
if own-house 'False' → output '??' → see another feature like 'Has-Job'.

if Has-Job 'True' → output 'Yes' → don't need another feature.
if Has-Job 'False' → output 'No' → need another feature.

if (own-house):

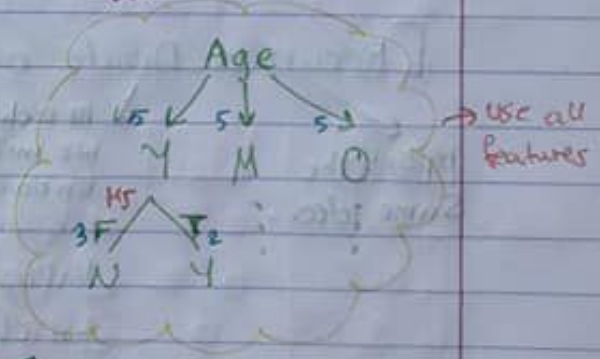
Yes

else if (Has-Job):

Yes

else ...

NO



own-house

F

T

HS

F

T

N

Y

use two features with high accuracy

→ general tree

and this what

we need → less number of features

with high accuracy.

L16: Decision tree Cont...
S11-4-2019.

→ Decision tree is subjective and waste time.

↳ if we build it manually.

* How to build the decision tree automatically?

The building of decision tree done Step by Step.

In each time we take the feature and from this feature, the tree is branching depending to the Value of feature.

[Recursive approach] → In each time repeat the same idea in each branch of feature.

→ initially I have all data Set, then start choosing the attributes. that's mean we will make a partitions on data set depends on Values of (Chosen attributes (features))

[Recursive Divide and Conquer Approach]

repeat the same idea.

↳ In each time we make partitioning on data set of sub data set based on Values of attributes that we choose.

→ We need general tree because we will use it in runing time on new samples that doesn't see it in training Stage.

So, we need the decision tree be 'Generalized'.

↳ less number of features on data set.
(less depth on tree)

We avoid the specification when we build the tree.

* We need automated method to build generalized tree.

→ To build generalized tree you should choose the most discriminative features.

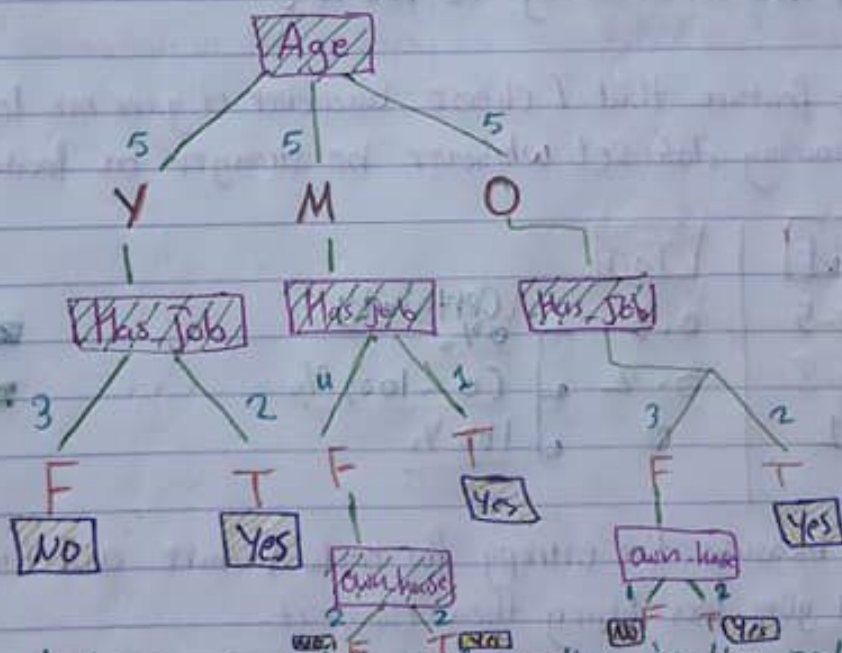
- * we choose one of the feature attribute to be a root of the tree. And the leaf represents the class.
- * the model that I get from decision tree is 'Hierarchy of if-else statement'

→ Instead you get the model manually, the decision tree as supervised learning method will get it automatically.

- * If we want to compare between the decision tree and other supervised learning method we look at speed. Since, decision tree is the fastest one in training phase.

→ Usually we compare between algorithms in three things:

- 1 Training.
- 2 performance.
- 3 The ability to return a generalized model.



- * The difference between one tree and another is the order of attributes and the most important attribute is the root.

→ How we choose the best attribute to get generalized tree?

The Feature that have more information to discriminate.

- Has Job - → The best feature

* Information Theory: measures the amount of information in the attributes and measure in bit Unit. [mathematical measure]

→ Information is subset of data.

↳ noise, redundancy, ...

* We measure the information theory by 'Entropy'. and the entropy measure the amount of blurriness in data (attributes).

→ example :

in a symmetrical coin → the probability of head = 50%, and the probability of tail is 50%.

But the Certainty of appearing exactly head = 0% & the Certainty of appearing exactly tail = 0%.

→ So, the uncertainty is 100%.

→ the feature that I choose, whenever it gave me less blurriness in remaining data set whenever be stronger as feature.

	Head	Tail	
probability:	0.5	0.5	Certainty
	0.8	0.2	(0-100)%
	1	0	100%

■ Certainty
■ probability

→ We measure the entropy for each feature and the feature that gives less blurriness then choose it.

* How we apply it in decision tree? How we know that this decision tree less uncertainty than another??

→ After partition process, if the remaining samples need less number of features until reach the leaf (class) then this tree is contain less amount of blurriness.

So, the first feature in decision tree contain a large amount of information.

~ The main challenge in building the decision tree is the order of attributes. (The operation hierarchy in building the decision tree) We choose these attributes from the strongest to the weakest to reach the generalized tree. [less # of features / less depth]. And to know the order of attributes we need 'information theory', and it's measure the amount of information in bit unit for attribute.

$$* \text{Entropy}(D) = - \sum_{j=1}^{|C|} \text{probability}(\text{class}_j) \cdot \log_2(\text{probability}(\text{class}_j))$$

→ Less entropy → less amount of blurry in feature. (Strong feature)
So, now I can order the features from the strongest to the weakest by calculating the entropy.

* How many times the entropy is calculated? and how many times the summation in the entropy is repeated (# of iterations)?

→ At first we calculate the entropy for all features (on all data set). Then, we ask how much is left of the uncertainty?
So, initially we have a specific uncertainty and through the time it will be less. Then, calculate the entropy for each feature.

~ The aim of the entropy is reduce the uncertainty.

* Uncertainty is the entropy of all data set. Then, when we choose a feature from data set that represent the root or sub root that it's mean how much - this root or sub root will reduce the uncertainty that I start from it.
Because, in each branch we get a new data set, and this new data set will choose for it a new sub root, and this sub root how much will reduce the uncertainty?

[It will be clear through an example]

Example on Slide (41):

1 Calculate the entropy for all data set (classes). [Data Set table on Slide (41)]

Class $\begin{cases} \rightarrow \text{Yes} \\ \rightarrow \text{No} \end{cases}$

$$p(\text{Yes}) = \frac{9}{15}, \quad p(\text{No}) = \frac{6}{15}$$

$$\rightarrow \text{Entropy (D)} = - \left[\frac{9}{15} \log_2 \left(\frac{9}{15} \right) + \frac{6}{15} \log_2 \left(\frac{6}{15} \right) \right] = \boxed{0.971} \quad \rightarrow \text{reference}$$

Then, we choose the best feature that reduce this entropy as much as possible.

2 Calculate the entropy for each feature to choose the root.

Age $\begin{cases} \rightarrow \text{Young} \rightarrow \text{new data set} \rightarrow p(\text{Young}) = \frac{5}{15} \\ \rightarrow \text{Middle} \rightarrow p(\text{Middle}) = \frac{5}{15} \\ \rightarrow \text{old} \rightarrow p(\text{old}) = \frac{5}{15} \end{cases}$

$$\star \text{Entrop}_{\text{Age}}(\text{D}) = p(\text{young}) \times \text{Entropy}(\text{young}) + p(\text{middle}) \times \text{Entropy}(\text{middle}) + p(\text{old}) \times \text{Entropy}(\text{old})$$

$$= \frac{5}{15} \cdot \left(- \left(p(\text{Yes}) \log_2(p(\text{Yes})) + p(\text{No}) \log_2(p(\text{No})) \right) \right) \quad \text{Young}$$

$$+ \frac{5}{15} \cdot \left(- \left(p(\text{Yes}) \log_2(p(\text{Yes})) + p(\text{No}) \log_2(p(\text{No})) \right) \right) \quad \text{Middle}$$

$$+ \frac{5}{15} \cdot \left(- \left(p(\text{Yes}) \log_2(p(\text{Yes})) + p(\text{No}) \log_2(p(\text{No})) \right) \right) \quad \text{old}$$

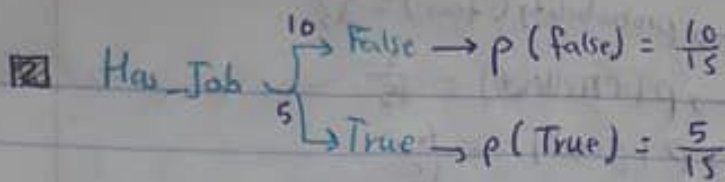
$$= \frac{5}{15} \cdot \left(- \left(\frac{2}{5} \log_2 \left(\frac{2}{5} \right) + \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right) \right)$$

$$+ \frac{5}{15} \cdot \left(- \left(\frac{3}{5} \log_2 \left(\frac{3}{5} \right) + \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right) \right)$$

$$+ \frac{5}{15} \cdot \left(- \left(\frac{4}{5} \log_2 \left(\frac{4}{5} \right) + \frac{1}{5} \log_2 \left(\frac{1}{5} \right) \right) \right)$$

$$= \frac{5}{15} (0.971) + \frac{5}{15} (0.971) + \frac{5}{15} (0.722)$$

$$= \boxed{0.888}$$



\odot Entropy (Has Job) = $p(\text{True}) \times \text{Entropy}(\text{True}) + p(\text{False}) \times \text{Entropy}(\text{False})$

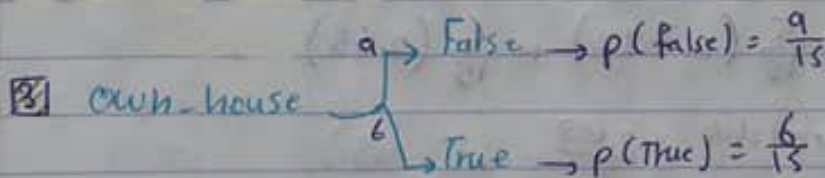
$$= \frac{5}{15} \left(- (p(\text{Yes}) \log_2(\text{Yes}) + p(\text{No}) \log_2(\text{No})) \right) + \frac{10}{15} \left(- (p(\text{Yes}) \log_2(\text{Yes}) + p(\text{No}) \log_2(\text{No})) \right)$$

$$= \frac{5}{15} \left(- \left(\frac{5}{5} \log_2(1) + \frac{0}{5} \log_2(0) \right) \right)$$

$$+ \frac{10}{15} \left(- \left(\frac{4}{10} \log_2(0.4) + \frac{6}{10} \log_2(0.6) \right) \right)$$

$$= \frac{5}{15} (0) + \frac{10}{15} (0.971)$$

$$= \boxed{0.647}$$



\odot Entropy (Own house) = $p(\text{True}) \times \text{Entropy}(\text{True}) + p(\text{False}) \times \text{Entropy}(\text{False})$

$$= \frac{6}{15} \left(- (p(\text{Yes}) \log_2(\text{Yes}) + p(\text{No}) \log_2(\text{No})) \right) + \frac{9}{15} \left(- (p(\text{Yes}) \log_2(\text{Yes}) + p(\text{No}) \log_2(\text{No})) \right)$$

$$= \frac{6}{15} \left(- \left(\frac{6}{6} \log_2(1) + \frac{0}{6} \log_2(0) \right) \right) + \frac{9}{15} \left(- \left(\frac{3}{9} \log_2\left(\frac{3}{9}\right) + \frac{6}{9} \log_2\left(\frac{6}{9}\right) \right) \right)$$

$$= \frac{6}{15} (0) + \frac{9}{15} (0.918) = \boxed{0.551}$$

4 → fair → probability (fair) = $\frac{4}{15}$
 5 → excellent → $p(\text{excellent}) = \frac{5}{15}$
 6 → good → $p(\text{good}) = \frac{6}{15}$

Entropy (D) = $p(\text{fair}) \times \text{Entropy}(\text{fair}) + p(\text{excellent}) \times \text{Entropy}(\text{excellent})$
 + $p(\text{good}) \times \text{Entropy}(\text{good})$

$= \frac{4}{15} \left(- (p(\text{yes}) \log_2(\text{yes}) + p(\text{no}) \log_2(\text{no})) \right)$
 + $\frac{5}{15} \left(- (p(\text{yes}) \log_2(\text{yes}) + p(\text{no}) \log_2(\text{no})) \right)$
 + $\frac{6}{15} \left(- (p(\text{yes}) \log_2(\text{yes}) + p(\text{no}) \log_2(\text{no})) \right)$

$= \frac{4}{15} \left(- \left(\frac{0}{4} \log_2(0) + \frac{4}{4} \log_2(1) \right) \right)$

+ $\frac{5}{15} \left(- \left(\frac{4}{5} \log_2\left(\frac{4}{5}\right) + \frac{1}{5} \log_2\left(\frac{1}{5}\right) \right) \right)$

+ $\frac{6}{15} \left(- \left(\frac{5}{6} \log_2\left(\frac{5}{6}\right) + \frac{1}{6} \log_2\left(\frac{1}{6}\right) \right) \right)$

$= \frac{4}{15} (0) + \frac{5}{15} (0.722) + \frac{6}{15} (0.65)$

$= \boxed{0.501} \sim \text{on slide} = \frac{0.608}{?}$

→ Reference = $\boxed{0.971}$

* $E_{\text{Age}}(D) = 0.888 \rightarrow \text{gain}(D, \text{Age}) = 0.971 - 0.888 = 0.083$

+ $E_{\text{Has_Job}}(D) = 0.647 \rightarrow \text{gain}(D, \text{Has_Job}) = 0.971 - 0.647 = 0.324$

⊕ $E_{\text{Own_House}}(D) = 0.551 \rightarrow \text{gain}(D, \text{Own_House}) = 0.971 - 0.551 = \boxed{0.42}$

+ $E_{\text{Credit_Rating}}(D) = \frac{0.501}{\frac{0.608}{\text{slide}}} \rightarrow \text{gain}(D, \text{Credit_Rating}) = 0.971 - \frac{0.501}{\frac{0.608}{\text{slide}}} = \frac{0.47}{0.608} = \boxed{0.363}$

⊛ From these results, the root is Own_House. Since, the Own_House has the highest information

I will consider this value

At this point, the decision tree looks like.



* Here, we observe that whenever the OWN_House is 'True', the class is always 'Yes', it's no coincidence by any chance, the simple tree resulted because of the highest information gain is given by the attribute OWN_House.

3] Next step, the new reference is $E(\text{OWN_House}(\text{false})) (D)$.
 - New Data Set.



Age	Has_Job	Credit_Rating	Class
Young	false	Fair	No
Young	false	excellent	No
Young	True	good	Yes
Young	false	fair	No
Middle	false	Fair	No
middle	false	good	No
old	true	good	Yes
old	true	excellent	Yes
old	false	fair	No

* Will repeat the same previous step in this new data set.

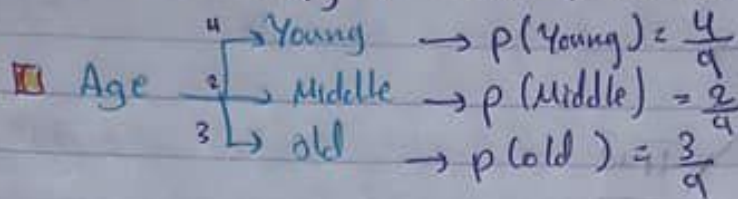
1] Calculate the entropy for all data set

$$\text{class} \begin{cases} \rightarrow \text{Yes} \\ \rightarrow \text{No} \end{cases} \rightarrow P(\text{Yes}) = \frac{3}{9}, \quad P(\text{No}) = \frac{6}{9}$$

$$\text{Entropy}(D) = - \left[\frac{3}{9} \log_2 \left(\frac{3}{9} \right) + \frac{6}{9} \log_2 \left(\frac{6}{9} \right) \right] = \boxed{0.918}$$

↓
New reference

② Calculate the entropy for each feature to choose the sub root.



$$* \text{Entropy (D)} = p(\text{Young}) \times \text{Entropy}(\text{Young}) + p(\text{Middle}) \text{Entropy}(\text{Middle}) + p(\text{Old}) \cdot \text{Entropy}(\text{Old})$$

$$= \frac{4}{9} \left(- (p(\text{Yes}) \cdot \log_2(\text{Yes}) + p(\text{No}) \log_2(\text{No})) \right)_{\text{Young}}$$

$$+ \frac{2}{9} \left(- (p(\text{Yes}) \log_2(\text{Yes}) + p(\text{No}) \log_2(\text{No})) \right)_{\text{Middle}}$$

$$+ \frac{3}{9} \left(- (p(\text{Yes}) \log_2(\text{Yes}) + p(\text{No}) \log_2(\text{No})) \right)_{\text{Old}}$$

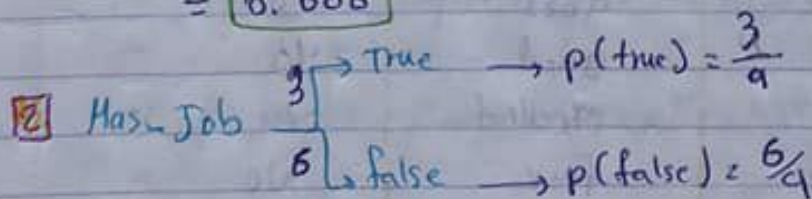
$$= \frac{4}{9} \left(- \left(\frac{1}{4} \log_2 \left(\frac{1}{4} \right) + \frac{3}{4} \log_2 \left(\frac{3}{4} \right) \right) \right)$$

$$+ \frac{2}{9} \left(- \left(\frac{0}{2} \log_2(0) + \frac{2}{2} \log_2(1) \right) \right)$$

$$+ \frac{3}{9} \left(- \left(\frac{2}{3} \log_2 \left(\frac{2}{3} \right) + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right) \right)$$

$$= \frac{4}{9} (0.811) + \frac{2}{9} (0) + \frac{3}{9} (0.918)$$

$$= \boxed{0.666}$$



$$* \text{Entropy (D)} = p(\text{true}) \times E(\text{true}) + p(\text{false}) \times E(\text{false})$$

$$\text{Has Job} = \frac{3}{9} \left(- \left(\frac{0}{3} \log_2(0) + \frac{3}{3} \log_2(1) \right) \right)$$

$$+ \frac{6}{9} \left(- \left(\frac{0}{6} \log_2(0) + \frac{6}{6} \log_2(1) \right) \right)$$

$$= \frac{3}{9} (0) + \frac{6}{9} (0) = \boxed{0}$$

31 Credit Rating

$$\begin{array}{l}
 4 \rightarrow \text{Fair} \rightarrow p(\text{fair}) = \frac{4}{9} \\
 2 \rightarrow \text{Excellent} \rightarrow p(\text{excellent}) = \frac{2}{9} \\
 3 \rightarrow \text{Good} \rightarrow p(\text{good}) = \frac{3}{9}
 \end{array}$$

Entropy (D) = $p(\text{fair}) \times E(\text{fair}) + p(\text{excellent}) \times E(\text{excellent}) + p(\text{good}) \times E(\text{good})$

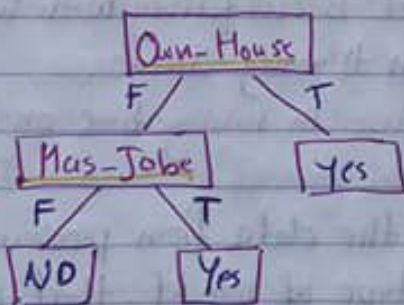
$$= -\frac{4}{9} \left(\frac{0}{4} \log_2 0 + \frac{4}{4} \log_2 1 \right) + -\frac{2}{9} \left(\frac{1}{2} \log_2 0.5 + \frac{1}{2} \log_2 0.5 \right) + -\frac{3}{9} \left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \left(\frac{1}{3} \right) \right)$$

$$= \frac{4}{9}(0) + \frac{2}{9}(1) + \frac{3}{9}(0.918) = 0.528$$

→ Reference = 0.918

- $g(D, \text{Age}) = 0.918 - 0.666 = 0.252$
- $g(D, \text{Has Job}) = 0.918 - 0 = 0.918$
- $g(D, \text{Credit Rating}) = 0.918 - 0.528 = 0.390$

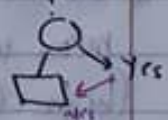
* From the result the subroot is 'Has Job'.



← The final decision tree that comes from training and this is a generalized tree with less number of features.

* We stop the algorithm when:

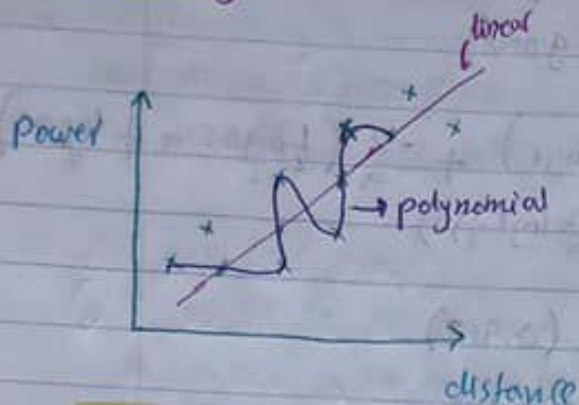
- 1 All Samples are classified.
- 2 There's no more attributes. → majority class →
- 3 There's no Samples → no problem → finish



L17: Overfitting & Neural Networks

M/8-4-2019

* Overfitting:



* polynomial with less error and more fitting.

(Regression)

Error: the difference between the expected curve and the actual points.

→ Overfitting → less error in regression or misclassification (less in classification) [in training data]

* now we need a good model in testing not training.

So we make model to be more fitting on testing.

I don't call about model is overfitting with high accuracy or performance in training stage.

Because, the model will use in running time on unseen samples.

→ Briefly, after collected the data from previous experiment, I should choose the type of model that will be overfitting in testing stage.

* In decision tree, it be overfitting when it be long. That is mean we use large number of attributes in building the tree (high depth) and this will reflect negatively on testing.

- In general, overfitting means that model give me high performance in training but bad in testing: one of reasons is use a lot of features or these features that use are not the optimized set.

→ This is optimization problem ~ what is the minimum depth of the tree or minimum number of features that are used in tree and give me a good accuracy on testing data?? ← This is what we are looking for.

* large number of features → good accuracy on training
↳ bad accuracy on testing

→ The reasons of overfitting

- noise → (تقلبات - تغير في القيمة)
- outliers → (قيمة متطرفة بعيدة (خاتة) قيمة متطرفة بعيدة)
- ↳ Data is not enough (It's poor in information)

* overfitting → you consider that the noise is something real in data then you do more fit to combine it to a class. and this is wrong, so we need generalized to avoid this problem.

- For more explanation: See example slide (14)

* overfitting: classification with a lot of details (تصنيف الجمل)
This good in training but bad in testing.

* How to do the generalization in tree?

□ prepruning
□ postpruning] to avoid overfitting.

* postpruning:

- Build the full tree. → (target accuracy in training)
- Then remove nodes from leaves. (lastest features in leaves)
- Try new tree in testing data → if the accuracy is the same, then remove more features. Stop → when there's a difference between testing & training accuracy.

* prepruning:

- if the feature reached to the information gain less than specific threshold then stop the building of tree. There is no need to build a full tree.
(faster)

→ postpruning is more generalized than prepruning and doesn't depend on application.

- So, postpruning is better than prepruning. (in general).

* Neural Networks

↳ its model is networks of neurons.

networks → interconnection between neurons.

* model → Network

↳ different from problem to another.

↳ # of layers, # of neurons in one layer, weights, # of inputs of neurons, # of outputs from neurons.

→ neuron in neural networks is the basic element.

the computations happen inside the neuron, link: to communicate between it and another neuron.

* The output from a neuron is input for another neurons.

- Networks in neural is not fully connected.

→ How do the customization of the network base on feature vector and the classes?

→ (# of input neurons) = (# of features in feature vector)
in input layer

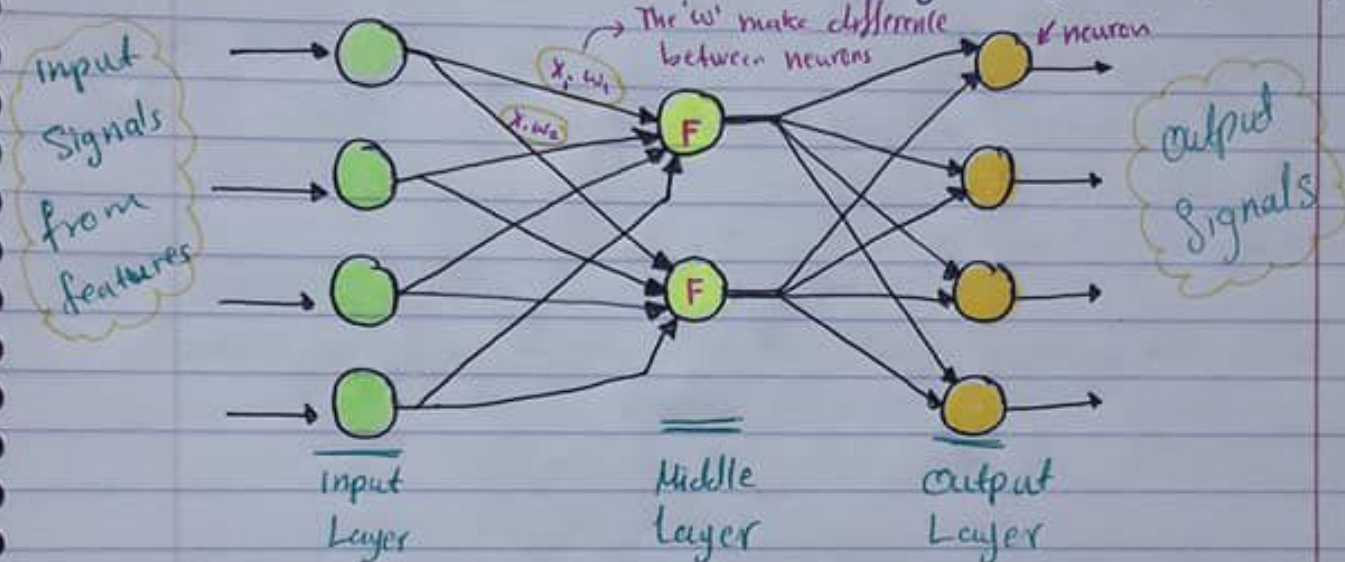
→ (# of output neurons) = (# of classes)
in output layer

of layers need optimization.
(middle layers)

* Neuron is function take the inputs from the output of previous layer or from features.

(# of inputs \rightarrow ~~one~~ feature or # of outputs in previous layer)
for one neuron

\rightarrow Fully Connected between neighbors Layers only. (layer X , layer $X+1$)



* Training Stage: Neural network finds the values of weights that through them I can make true classification.

\rightarrow But, initially the weights are chosen randomly.

* Function of neural network (neuron):

$$\rightarrow a = \sum_{i=1}^M x_i w_i - \phi$$

L18: Neural Networks Cont.....

W/10-4-2019

* Architecture of network is input parameter.

When we run the network, we set the number of layers and number of neurons in middle layer.

The network in automatical way, it build the input layer and output layer. Because, it comes from features vector & classes.

→ So, the architecture of network is not something trained. The training stage in neural network make a tuning on weights to do the true classification.

* How the neural network make the tuning or update on weight?
- initially start with random values of weights.

- neural network in deep learning → large number of layers.

- neural network in classical learning → small number of layers.

because in classical learning, the feature extraction is not part of network but in deep learning is network responsibility.

So, we need specific layers for feature extraction and specific layers in classification. That is the reason of needing to large number of layers in deep learning.

→ The tuning on weights based on → 1. What is the output of network
2. and the actual output.

* We make the tuning on weights after predict the class from network and from this prediction compared with actual output.

Then, make a tuning.

* Tuning → increase the weight
→ decrease the weight

* in classical learning usually number of layers are 1 or 2 or 3.

* The Perceptron :

- Linear Separable :

→ Need one straight line to separate between any two classes.

- Single layer is enough to build the classifier. [1]

- Non Linear Separable :

→ Can not separate between two classes in one straight line.

- need more than one layer neural networks to build the classifier.

of layers usually [2 or 3].

* perceptron : is single layer neural network.

- The training time of neural network is very large compared with decision tree.

* The decision tree is the fastest classifier.

* The neural network is the slowest classifier.

* perceptron - linear separable :

$$\left(\sum_{i=1}^n x_i w_i \right) - \phi$$

activation function

↳ because from this we decide the active class.

- one layer

- initialize the values of weights & threshold.

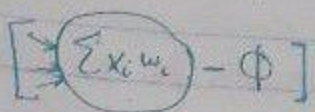
- choose type of activation function.

- The initial weights are randomly assigned, usually in range [0.5, 0.5] (recommendation)

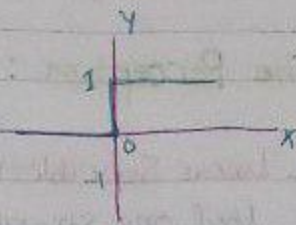
* Activation function of a neuron:

linear
separable

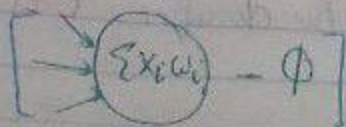
① Step function.



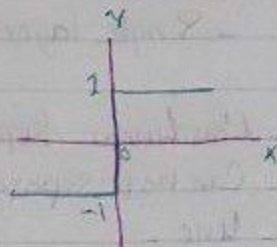
$$\begin{aligned} &\geq 0, 1 \\ &< 0, 0 \end{aligned}$$



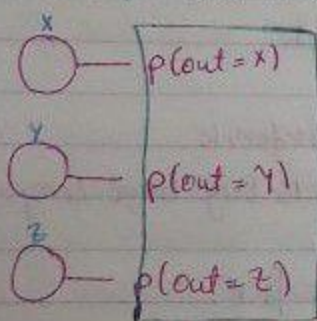
② Sign function.



$$\begin{aligned} &\geq 0, 1 \\ &< 0, -1 \end{aligned}$$



↑ For one class, if we have many classes (neurons out), then we need neuron for each class.



← Softmax layer

$$\sum \text{probability} = 1$$

→ choose the class with highest probability

* Error = Desired output on data set - predicted output (actual) from network

$$E = Y_d - Y_A$$

per one feature vector

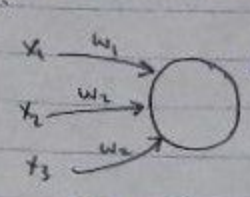
If error > 0 → increase the Y_A

If error < 0 → decrease the Y_A

→ neural network is incremental

→ each feature vector has an output, and from the predicted output, I find the error then update the weight. The new weight will use in next feature, and so on...

$$W_{i, \text{Next}} = W_{i, \text{Current}} + \text{Error} \Delta w_i$$



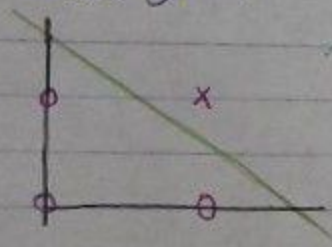
we don't change the weight in the same amount.

$$W_{i, \text{Next}} = W_{i, \text{Current}} + \text{learning rate} \cdot X_i \cdot \text{error}$$

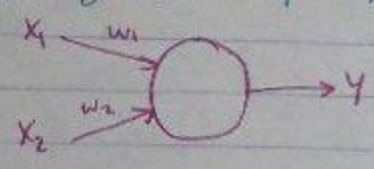
we update the weight scaled on X_i

we stop training when we reach the wanted requirement.

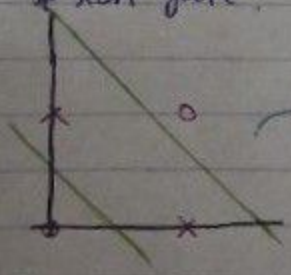
* AND gate.



→ Single layer, two inputs, output 0 or 1.



* XOR gate.



→ we cannot separate it using single layer.

* → AND gate example on slides:

→ Training Data is:

x_1	x_2	out ← Desired output (y_d)
0	0	0
0	1	0
1	0	0
1	1	1

" See complete example on slide :) "

Forward → get output

Backward → update the weights

↳ new weights → use in next sample.