

I Search

- Informed search: uses evaluation function F to choose which nodes to expand
- Heuristic search: informed search but F is an estimate which contains heuristic h
↳ (e.g.) Straight line distance



Computing $h(\text{node})$ is usually much easier.

Greedy best-first search: make a locally optimal choice based on a heuristic h

Problems:

- can get trapped in local minimum
- can get trapped in endless cycle
- can terminate with sub-optimal path

A* search :- uses $h(\text{node})$

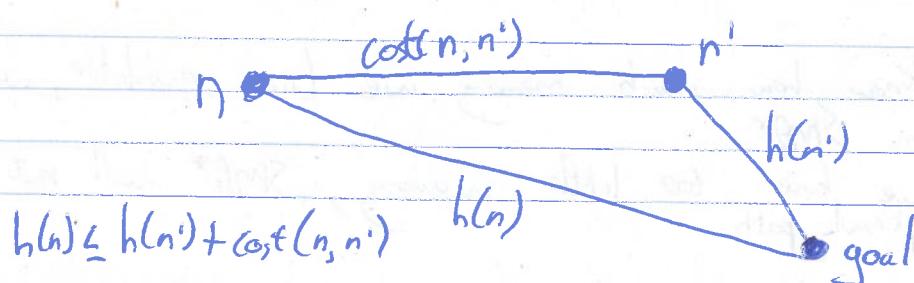
(- uses the cost to get to node n from the starting point)

$$f(n) = h(n) + g(n)$$

\uparrow
usually 1

Admissible : h never overestimates

Consistent : h satisfies the triangle inequality :



A consistent heuristic is always admissible but not vice versa

A^* search has a space complexity of $O(b^n)$

↓
We can get a linear complexity by using IDA*

Iterative deepening A^*

→ keeps at most 1 path in memory ~~at other cost~~

Iterative depth first searches to an increasing limit.

→ Normal Iterative deepening DFS: increase depth limit by 1
↓
until a solution is found

IDA*: depth limit is the path length f_{lim}
↳ Nodes with $f > f_{lim}$ will not be expanded.

Initial: $f_{lim} = f(\text{start-node}) = h(\text{start-node})$

Then: f_{lim} is increased to the lowest $f(\text{node}) > f_{lim}$

Tree-search A^* is optimal if h is admissible.

Graph-search A^* is optimal if h is consistent.

Simplified Memory-bounded A^* (SMFA*)

→ Used if we are very tight on memory or if we have more memory than IDA* uses.

↓

If we know how much memory we have available, we can use SMFA*

↳ if we have too little memory, SMFA* will not find an optimal path

- The current F value of each node is continuously updated
 - unexpanded/partially expanded: $f(n) = g(n) + h(n)$
 - fully expanded: $f(n) = \text{lowest } f\text{-value of any descendant}$
 - no-go nodes (or which path exceeds memory limit): $f(n) = \infty$

- Nodes also store the lowest F-value of any forgotten descendant
- If memory limit is reached, node with highest F-value is dropped

↓
of children
lowest F-value is forced at parent node

- When all actual $f(n) > \text{forgotten } f(n)$, the node is re-expanded.

Minimal spanning trees: useful for finding the cheapest way to connect N vertices given possible connections between them.

Can also be used for cheapest path problems
 ↳ Some MSTs are paths

MST finding algorithms:

- Prim → Pick the minimal edge from the connected subset
- Kruskal → Add lowest cost edges without making a loop

Greedy algorithms

II Local search algorithm

Local search : do not search the entire problem space

↳ Start with a random "possible solution" and improve it.

or chosen after some variable

Restart: changing the entire set of variables

Walk: changing a few variables at a time

→ NOT complete (not guaranteed to find the optimal solution)

→ Local search is very inefficient

→ We need something to make it more efficient

evaluation function

Iterative best improvement

→ Much like greedy best-first search

1. We find the neighbour states of the current state
(states with a one variable difference)

2. Select the neighbour with the best evaluation

3. Change the state to the neighbour state

→ Iterate

→ This method can get stuck in a local optimum

↳ No better neighbors and no adequate solution

Can be fixed by random restarts

Simulated annealing: iterative improvement with decreasing shakiness

- We shake a lot in the beginning and hope to get out of large valleys
- We shake less later on and hope to get out of all subvalleys

NOTE: only walks of decreasing randomness, no restart

1. Start with a random initialization
2. Pick a variable at random and change it
3. If this new assignment is better, change the current state
4. If this solution is worse, change the current state with a random probability

Probability for accepting worsening steps:

$$\text{accept} = e^{(h(A) - h(A'))/T}$$

- $h(A)$ = evaluation of the state
- $h(A')$ = evaluation of the new assignment
- T = temperature

If T is high, the probability of acceptance is also high

Maximum flow networks

Flow network: Edges with weights, nodes, a source and a sink node

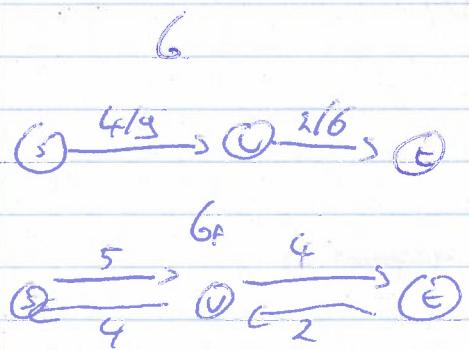
Directed graph

Max flow: maximum amount of outgoing and incoming flow from $s \rightarrow t$

Ford - Fulkerson method

1. Start with $f(u,v) = 0$ for all $u,v \in E$
2. Iteratively increase $|f|$ by some augmenting path in the residual network G_f .
3. Stop if G_f does not contain any augmenting paths anymore

Residual network:



Augmenting path: simple path in G_f from $s \rightarrow t$

Cut: partition of V into S and T such that
 $s \in S$ and $t \in T$

Net flow: flow across cut (S,T)

Capacity of $c(S,T)$ is the maximum flow $S \rightarrow T$

Minimum cut: cut with the lowest capacity of all cuts

III Complexity analysis

Complexity of an algorithm: Function relating the input length to the number of steps or storage locations

Time complexity: numbers of lines of code executed

Worst case: upper bound

Average case: average

Best case: lower bound

We are not interested in the exact number of instructions

We use the rate of growth

→ Big-Oh notation (asymptotic upper bound on running time)

Notation:

$$f(n) = O(g(n))$$

↳ f does not outgrow g modulo a constant factor or scalar

$$f(n) \leq c \cdot g(n) \text{ for all } n \geq n_0$$

Constant time: $O(1)$

Logarithmic time: $O(\log_2(n))$

Polynomial time: $O(n^c)$ for some $c > 0$

Exponential time: $O(c^n)$ for some $c > 1$

Decision problems: problems that have yes or no as an answer

Complexity class P: decision problems that can be decided in time polynomial in the input size.

Complexity class NP: decision problems that can be verified in time polynomial in the input size.

$$P \subseteq NP$$

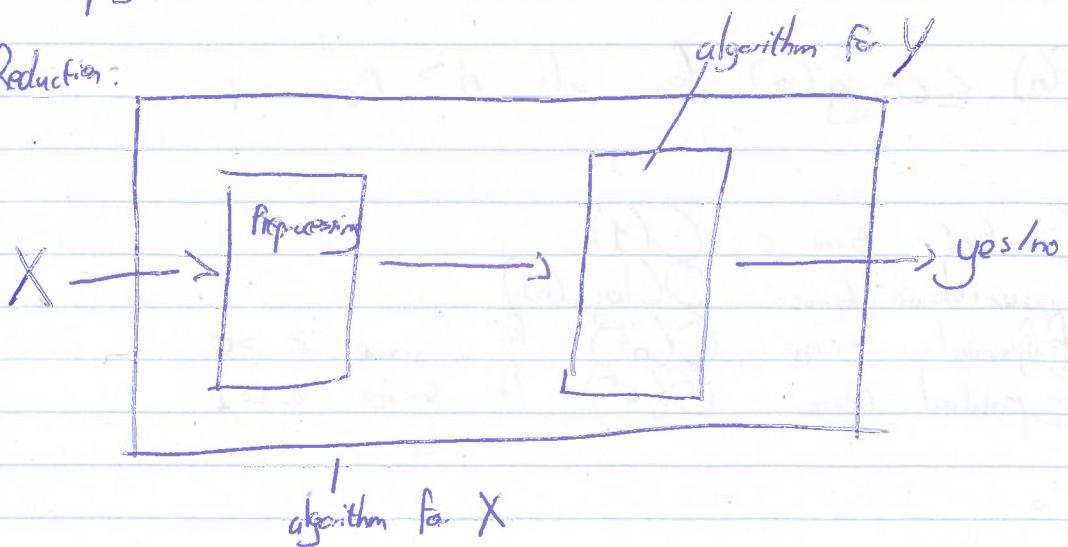
NP-hard: a problem that is at least as hard as any other problem in NP modulo a polynomial transformation

NP-complete: a problem is NP-complete if it is in NP and is NP-hard

Proving a problem is in NP-hard:

→ Transform problem X to a problem Y that is NP-hard in polynomial time

Reduction:



IV

Planning

CSP (Constraint satisfaction problem)

- Variables
- Domains (one for each variable)
- Constraints (predicate over the variables)
- Possible worlds

Ub. A sudoku puzzle:

- Cells that we can fill in (variables)
- Numbers we can put from 1-9 (domains)
- Restrictions on the values (constraints)
- Possible ways to assign numbers (possible worlds)

CSP solving:

1. list all possible worlds
2. For each world check the constraints

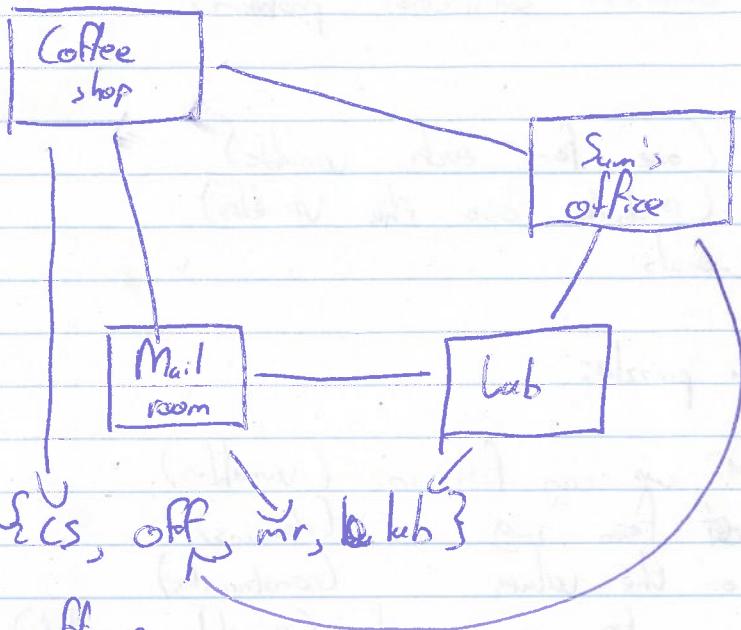
↳ Extremely time consuming

Planning: we want to find a sequence of actions that
is:
→ Starting from an empty state

- possible
- result in a state satisfying the goal

vb. Delivery robot example

Environment:



Features:

Location

Domain: {CS, OFF, MR, LAB}

(RHC) Rob has coffee

Domain: {true, false}

(SWE) Sam wants coffee

Domain: {true, false}

(MW) Mail is waiting

Domain: {true, false}

(RHM) Rob has mail

Domain: {true, false}

Actions: - Move clockwise (mc)

- Move anti-clockwise (muc)

- Pick up coffee (PUC) - must be at CS

- Deliver coffee (DeC) - must be at OFF and RHC

- Pick up mail (PUM) - must be in MR and MW

- Deliver mail (DelM) - must be in OFF and RHM

Solve a planning problem:

- see it as a pure search problem
- see it as a CSP problem

II Planning: intelligent decision making

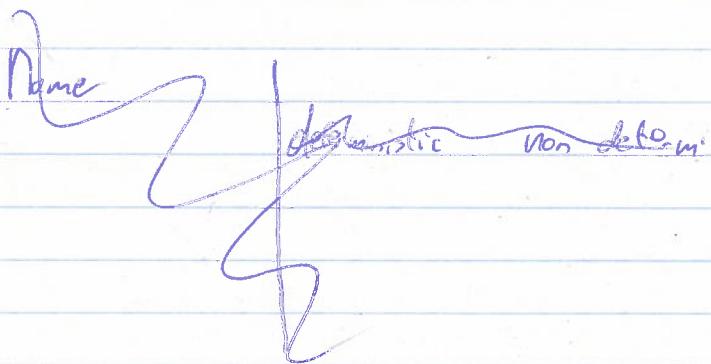
Solution to a planning problem: paths from an initial state to a goal state

↓
Solvable by Dijkstra but is infeasible for larger state spaces

) } g
constructing the transition graph

→ Planning algorithms try to avoid constructing whole graph

Classes of problems



1. Classical planning

Dynamics: deterministic

Observability: Full

Horizon: finite

2. Conditional planning with Full observability

Dynamics: nondeterministic

Observability: Full

Horizon: finite

3. Conditional planning with partial observability

Dynamics: nondeterministic

Observability: partial

Horizon: finite or infinite

4. Conformant planning

Dynamics: deterministic or nondeterministic

Observability: none

Horizon: finite

5. Markov decision processes (MDP)

Dynamics: probabilistic

Observability: full

Horizon: finite or infinite

6. Partially observable MDPs

Dynamics: probabilistic

Observability: ~~full~~ partial

Horizon: finite or infinite

Explicit state representation: explicitly enumerate each state

Ub. Delivery robot domain:

Place: {lab, cs, off, mw}

RHL: {rhc, rhc}

Sitx: {swc, swc}

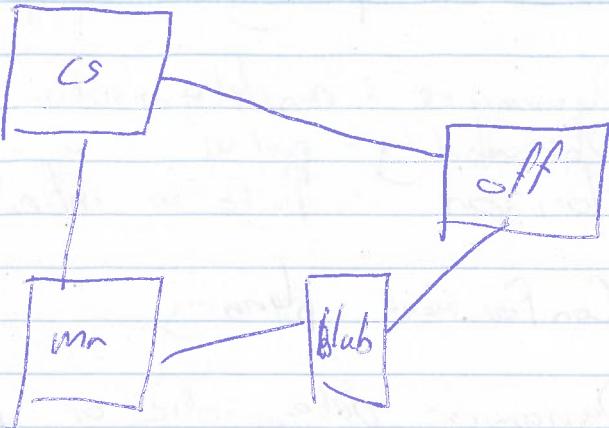
MW: {rmw, mw}

RHM: {rbm, rbm}

State is a five tuple of the above variables

Actions:-

- nc
- mnc
- pac
- dc
- pump
- dm



For each state, specify the actions possible and the resulting state from that action.

T

State-space graph

Complexity:-
- finding a solution: polynomial
- finding the shortest solution: NP-hard

Feature based representation of actions: preconditions and effects

Causal rules: specify when a feature gets a new value

jh. $R_{loc}' = cs \leftarrow R_{loc} = off \wedge Act = mnc$

Frame rules: specify when a feature keeps its value

$R_{loc}' = cs \leftarrow R_{loc} = cs \wedge Act \neq nc \wedge Act \neq mnc$

Action centered representation: STRIPS

!

Stanford research institution problem solver

Representation:

- precondition
- effect

pac: - precondition: $[cs, rhc]$
- effect : $[rhc]$

④

If the effect is a boolean list:

effect := add-list + delete-list

→ features not mentioned stay unchanged

STRIPS representation:

- initial state
- goal state
- variables
- actions

Two types of goals:

- Achievement goal (true in the final state)
- Maintenance goal (true in all states)

Two types of planning:

- scalar fixing (any solution)
- optimal (optimal solution)

Forward planning : searches the state-space graph from the initial state

→ The relevant part of the graph can be created dynamically

Regression planning : searches the state-space graph from the goal states

VI

CSP Planning

CSP: constraint satisfaction problem

Action Features: variables over actions.

PUC: Boolean

DCL: Boolean

...

Move: domain ($\{m_0, m_1, m_2\}$)



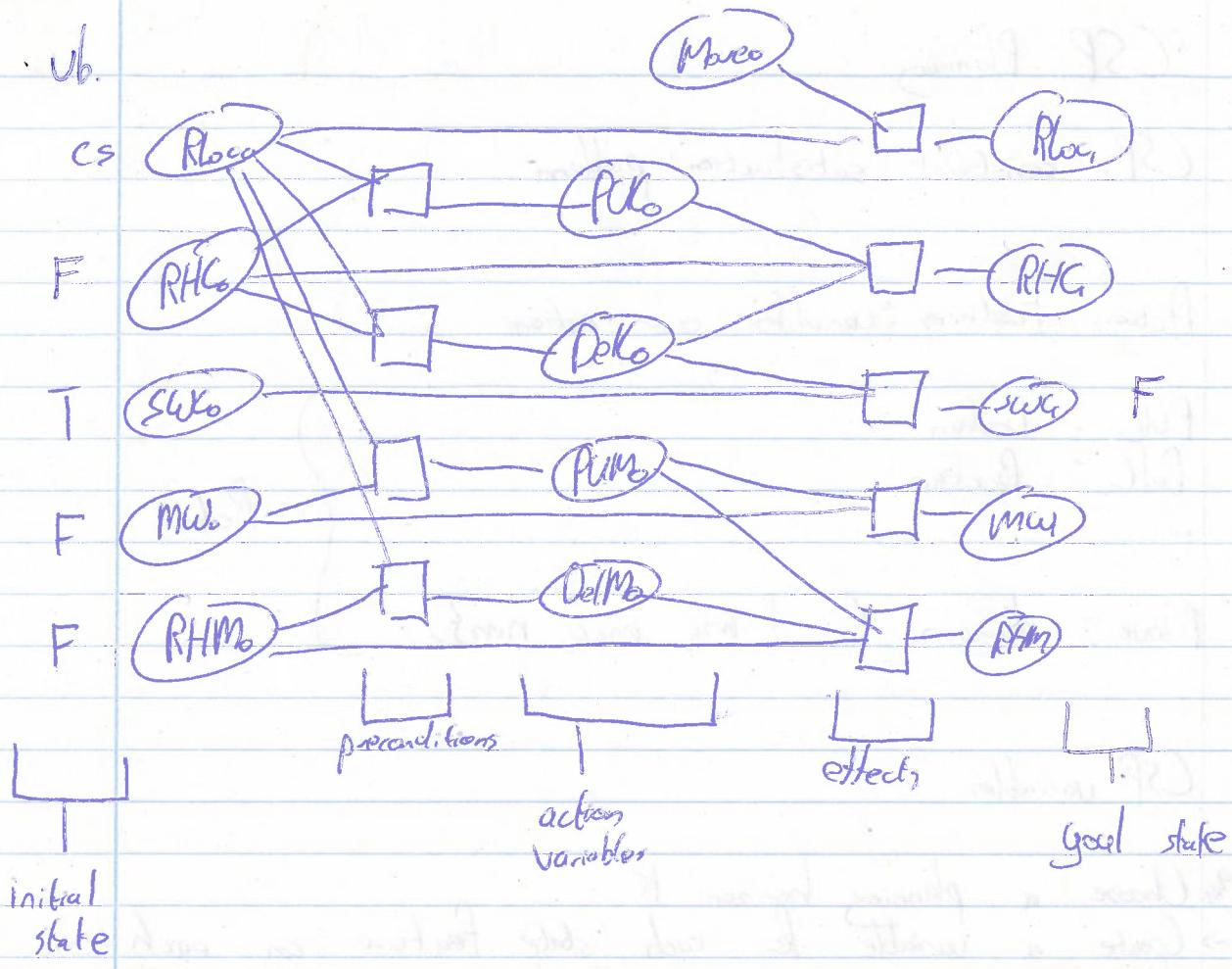
Robot

CSP variables

- Choose a planning horizon K
- Create a variable for each state feature at each time from 0 to K
- Create a variable for each action feature for each time in the range 0 to $K-1$

Constraint types:

- State constraint: constraint between variables in the same time step
- Precondition constraint: constraints between states at time t and actions at time t (what action is available at t ?)
- Effect constraint: between state variables at time t , action variables at time t and state variables at time $t+1$
- Action constraint: constraints between actions at time t (which actions can't co-occur)
- Initial state constraint: values of variables at time 0
- Goal constraint: values of variables at time $t+1$



Solving a CSP can be done by logically conducting steps.

Mutual relations: two actions that can't co-occur:

- One action negates the effects of the other
- One of the effects of one action is the negation of a precondition of the other
- One of the preconditions of one action is mutex with a precondition of the other

Two literals that can't co-occur:

- One of the literals is the negation of the other
- Each pair of actions that could achieve the two literals is mutex

Solving a CSP:

1. Run the algorithm for $K=1$
2. Run the algorithm for $K=2, 3, 4, \dots$ until a solution is found or it is clear that no solution exists

Partial order planning: executing multiple (non exclusive) actions in the same time step

$$\cup_b (Mov_e = mc, \{DelC = \text{True}, DelM = \text{True}\})$$

↳ partial order plan, the order of DelC and DelM does not matter but must happen after mc

POP as a search problem

→ States are (mostly unfinished) plans

Each plan has 4 components:

- a set of actions
- a set of ordering constraints $\text{act}_1 \xrightarrow{O} \text{act}_2$
- a set of causal links $\text{act}_1 \xrightarrow{P} \text{act}_2$
- a set of open preconditions

→ A plan is consistent if and only if there are no cycles in the ordering constraints and no conflicts with the causal links

→ A consistent plan with no open preconditions is a solution

→ A partial order plan is executed by repeatedly choosing any of the possible next actions

Solving POP

1. Assume propositional planning problem:
Initial plan: Start and Finish
 - the ordering constraint Start \rightarrow Finish
 - no causal links
 - all the preconditions in Finish are open

2. Successor function:
 - pick one open precondition p on an ~~act~~ action act,
 - generate a successor plan for every possible consistent way of choosing action acts that achieves p

3. Test the goal



- Generating a successor plan:
 - add causal link acts \rightarrow product
 - add the ordering constraint acts \rightarrow act.
 - resolve conflicts between new causal links and all existing actions
 - resolve conflicts between action acts

VII

Planning with costs as a search problem

- Nodes are planning states
- Edges correspond to operators

Optimal plan: plan with minimal cost

STRIPS

$$\text{Problem } P = \langle F, I, G, O, c \rangle$$

F: fluents (~~variables~~ Boolean variables)

I: initial state $I \subseteq F$

G: goal state $G \subseteq F$

O: operators $\langle \text{Pre}(o), \text{Add}(o), \text{Del}(o) \rangle$ each $\subseteq F$

c: costs $c(o)$

States are subsets of fluents that have the value true

Heuristic function (h): estimate of the distance from a state to the goal

- $h^*(s)$ = the cost of the lowest-cost path from s to a goal node

- h is admissible if for all $s \in S$: $h(s) \leq h^*(s)$
(never overestimate)

- h is consistent if for all $s \in S$ and edges (s, s') :
 $h(s) \leq c(s, s') + h(s')$

STRIPS:

- uninformed heuristic

Relaxation:

- a less constrained version of the problem \rightarrow heuristic

Classic STRIPS heuristic: uses the number of state variables that differ in the current state and a STRIPS goal.

- Ignores almost all problem structure
- Most successors have the same estimate

Relaxation: assumption about a problem that makes it easier to solve

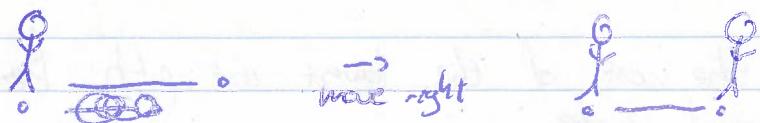
1b. Manhattan distance: spots between two indices

- ↳ does not consider barriers

Blind heuristic: $h_{\text{blind}}(s) = 0$ for all s

Relaxation P^+ : ignore bad effects of actions

- Once a variable has a certain value, that value can always be used to satisfy preconditions



P^+ is admissible

→ any solution to P is a solution to P^+ as well

In strips, we need to remove the delete list:

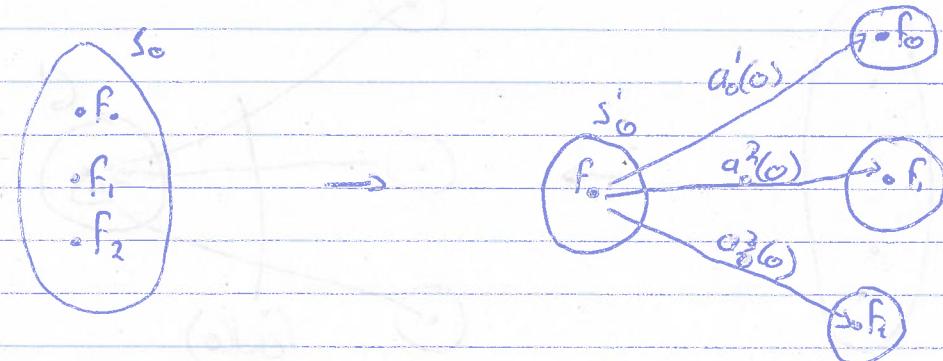
$$P^+ = \langle F, I, G, O^+, c \rangle$$

with $O^+ = \{\text{Pre}(c), \text{Add}(c), \phi\}$

Simplifying P^+

Ideally, $|S_0| = 1$ (number of fluents in a state equals 1)

To simplify P^+ :



A start state with size $|S_0|=n$ can be replaced with:

- a new start state containing a new fluent

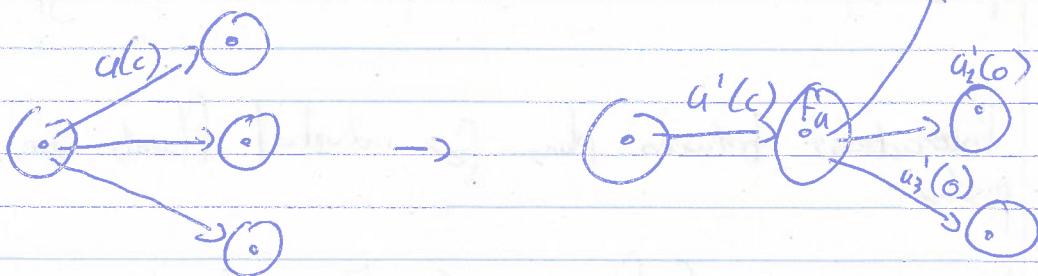
$$S'_0 = \{F\}$$

- a set of n zero-cost actions:

$$A'_0 = \{\langle \{F\}, \{f_i\}, \emptyset \rangle \mid f_i \in S_0\}$$

Ideally, $|\text{Add}(a)| = 1$

To simplify P^+ :



An action a with $|\text{Add}(a)| = n$ can be replaced with:

- a new fluent f_a representing that the action has been executed

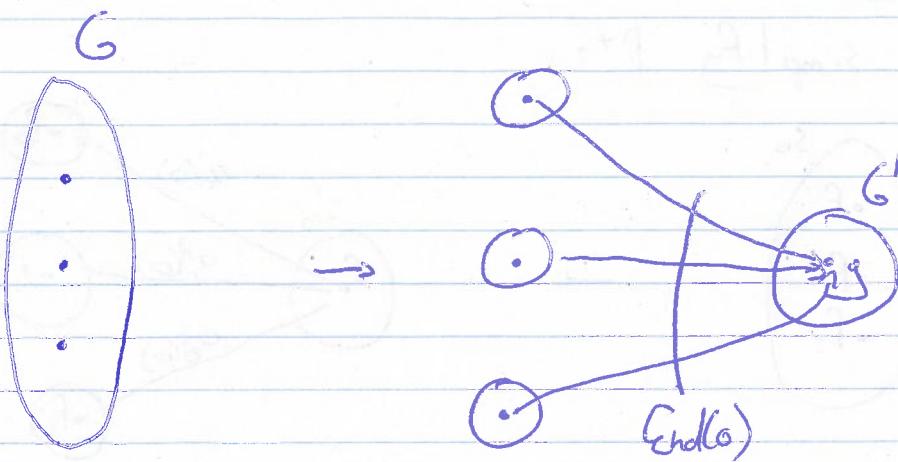
- an action $a' = \langle \text{Pre}(a), f_a, \emptyset \rangle$ with cost $c(a') = c(a)$

- a set of n zero-cost actions:

$$A' = \{\langle \{f_a\}, \{a'\}, \emptyset \rangle \mid a' \in \text{Add}(a)\}$$

Ideally $|G| = 1$

To simplify P^+ :



A goal with size $|G| = n$ can be replaced with:

- A new goal with a singly newly-introduced fluent $G' = \{g\}$
- A single action, $\text{End} = \langle G, g, \emptyset \rangle$

P^+ as a graph problem

Special cases of P^+ can be solved as graph problems:

- $|G| = |\text{Pre}(o)| = 1 \rightarrow$ Shortest path } tractable
- $|G| > 1, |\text{Pre}(o)| = 1 \rightarrow$ Directed Steiner tree } NP-hard
- $|\text{Pre}(o)| > 1 \rightarrow$ Optimal directed hyperpath } NP-hard

Interactions between plans for individual fluents are always positive:

For any set of fluents $G \subseteq F$,

$$\sum_{g \in G} h^+(g) \leq h^+(G)$$

Difficulty of P^+ : estimating cost of sets

Additive heuristic h_{add} : estimate the cost of a set as the sum of the costs of the individual fluents

Relies on the independence assumption for costs:

$$h_{add}(G) = \sum_{g \in G} h_{add}(g)$$

$$h_{add}(a; s) = \underbrace{cost(a)}_{\text{cost}} + h_{add}(\text{Pre}(a); s)$$

\rightarrow Cost to get to a from s is the cost to get to a plus the cost to get to all preconditions of a .

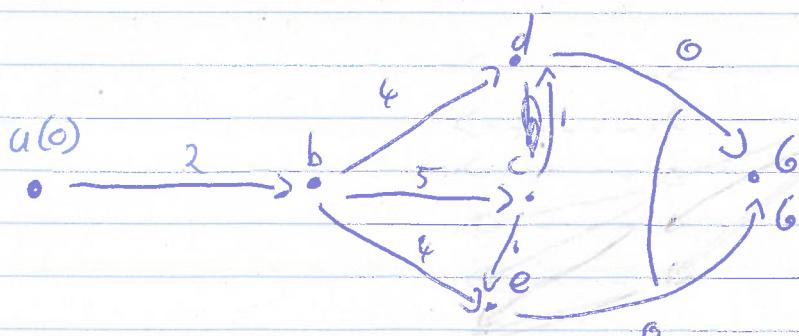
$$h_{add}(p; s) = \begin{cases} 0 & \text{if } p \in s \\ \min \text{ to get to } p \text{ from } s & \text{otherwise} \\ h_{add}(a_p(s)) \end{cases}$$

$$a_p(s) = \min_{a \in G(p)} h_{add}(a; s)$$

$$h_{add}(P; s) = \sum_{p \in P} h_{add}(p; s)$$

$$h_{add}(s) = h_{add}(G; s)$$

We can use Dijkstra to compute h_{add} :



a	b	c	d	e	f	g
0	2	7	6	6	12	

Problem with h_{add} :

- $h_{\text{add}}(G) = 12$

but a more optimal solution exists via c:

- $h^*(G) = 9$



1. h_{add} counts edges twice ($a \rightarrow b$)
2. Independence assumption:

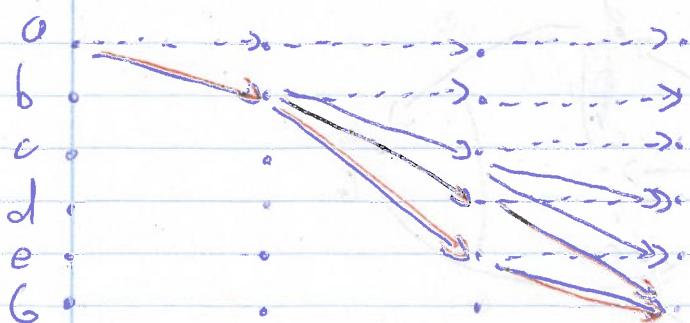
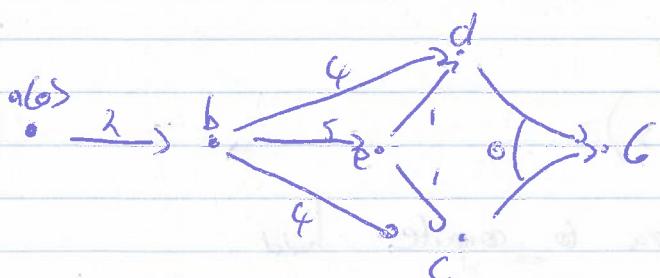
$$h_{\text{add}}(d, e) = h_{\text{add}}(d) + h_{\text{add}}(e)$$

↳ not optimal

Roller Heuristic h_{RF}

→ don't count edges twice

Relaxed planning graphs: graphically represent heuristic computation



Construct relaxed plan by starting from the goal, choose supports for each fluent at first layer from which it appears

$$h_{FF}(s) = \{a \rightarrow b, b \rightarrow d, b \rightarrow e, d, e \rightarrow G\} = 10$$

h_{FF} is NOT admissible

h_{max} heuristic : replaces \sum in hadd with max :

$$h_{max}(P; s) = \max_{p \in P} (h_{max}(p; s))$$

Estimates cost of set as cost of most expensive fluent in set

\rightarrow Admissible

heuristic	description	admissible?
hadd	sum of cost precondition plus sum of cost	no
h_{max}	max cost of precondition plus sum of cost	yes
h_{FF}	exploit positive interaction	no

VIII

Bayesian networks

→ Also known as probabilistic / Belief networks

Probability:

$P(A=a)$ or $P_{\sim}(A=a)$ denotes variable A takes value a

can also be non-binary

Meaning of probability:

- Frequentist: if A is repeated many times, the outcome a would occur $P(A=a)$ times (the amount of repetitions) (approximately)

- Bayesian!

Jugemental Subjective: The subjective degree of belief based on all available information, that one has in the event $A=a$

According to the information available to the subject

- Probabilities always sum up to one
- Prior and posterior probabilities are relative to an event

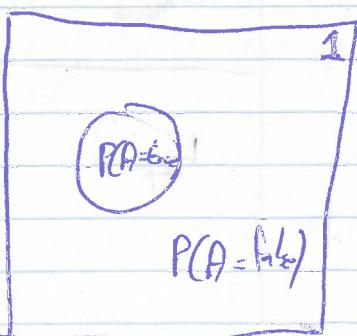
ub. X = "outcome of throwing a dice"

$$\Omega(X) \subset \{1, \dots, 6\}$$

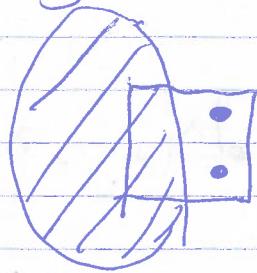
↓

of the probabilities

values of X



Prior probability: $P(X=4) = 1/6$



Posterior probability: $P(X=4) = 1/3$

Event: showing you a partially occluded die

Joint probabilities: probabilities of multiple events happening at the same time

v.b. $P(A=\text{true}, B=\text{true})$: proportion of outcomes where both A and B are true.

Conditional probabilities: probabilities of A conditioned on B

v.b. $P(A=\text{true} | B=\text{true})$: out of all the outcomes where B is true, how many also have A equal to true

Joint probability distribution: distribution of probabilities over multiple variables

For each combination of variables, assign a probability:

v.b.	A	B	$P(A, B)$
	true	true	0.1
	true	false	0.04
	false	true	0.3
	false	false	0.2

Product rule:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1})$$

→ For k boolean values, we need a table of size 3^k and have to specify $2^k - 1$ numbers

↓
We can solve this using independence

Variables A and B are independent if any of the following holds:

- $P(A|B) = P(A) \cdot P(B)$
- $P(A|B) = P(A)$
- $P(B|A) = P(B)$

Then we say $A \perp\!\!\!\perp B$

Conditional independence: two variables are conditionally independent if any of the following holds:

- $P(A, B | C) = P(A|C) \cdot P(B|C)$
- $P(A|B, C) = P(A|C)$
- $P(B|A, C) = P(B|C)$

Knowing C tells you everything about A 's information about B is not relevant anymore.

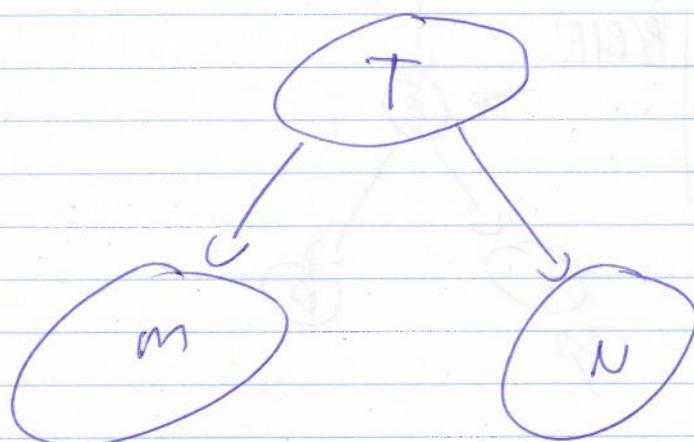
Then we say $A \perp\!\!\!\perp B | C$

v.b. Causal chain



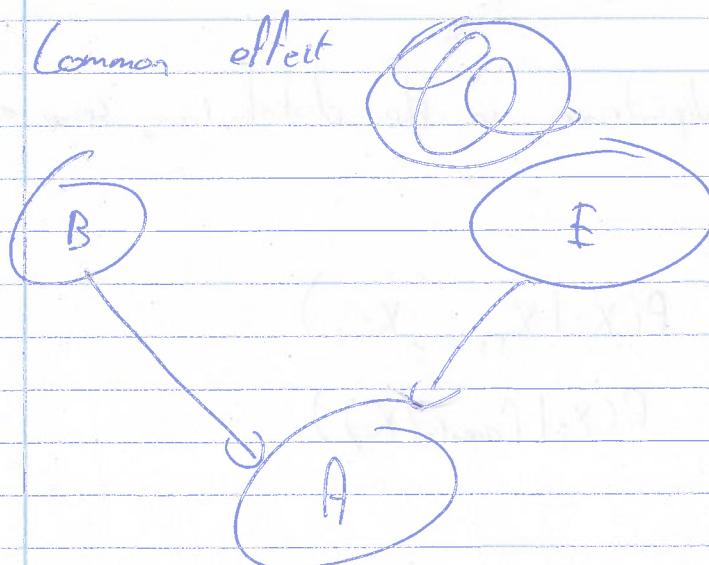
$A \nparallel B$ but $A \perp\!\!\!\perp B \mid W$

Common cause



$M \nparallel N$ but $M \perp\!\!\!\perp N \mid T$

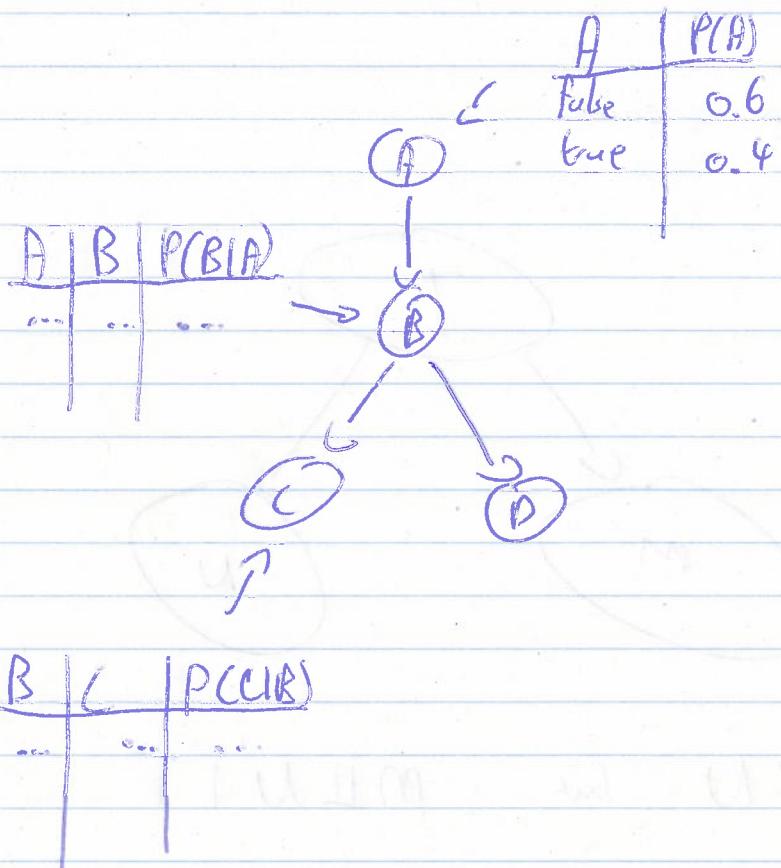
Common effect



$B \perp\!\!\!\perp E$ but $B \nparallel E \mid A$

Bayesian netw~~rk~~^t: a directed acyclic graph with nodes representing variables, each node has a probability table.

Each node ~~is~~ has a conditional probability table based on the parent nodes:



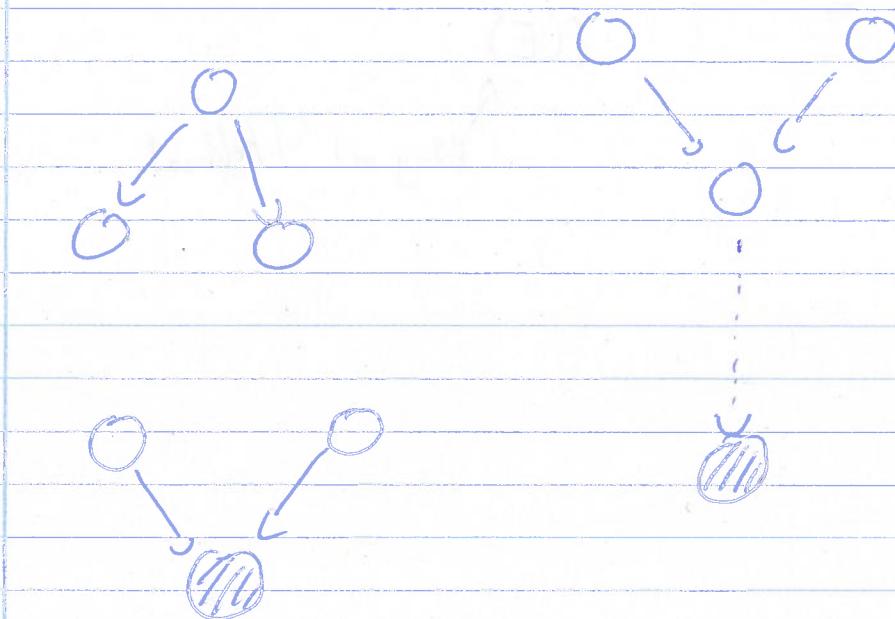
Because of the independences in the distribution, some arcs can be pruned:

$$\begin{aligned}
 P(X_1, \dots, X_n) &= \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1}) \\
 &= \prod_{i=1}^n P(X_i | \text{Parents}(X_i))
 \end{aligned}$$

D-separation (reachability): two nodes are D-separated, if all chains connecting them are inactive

Active triplets:

$$O \rightarrow O \rightarrow O$$



Inactive triplets:

$$O \rightarrow \text{XXX} \rightarrow O$$



Bayes rule

$$P(E|C) = P(E, C) / P(C)$$

$$P(C|E) = \frac{P(E|C) \cdot P(C)}{P(E)}$$

↑
Posterior

↓ Likelihood
Prior
Marginal likelihood

IX

Probability theory

- P is a probability measure over a set Ω
- P should obey three axioms:
 - $P(A) \geq 0$ for all events A
 - $P(\Omega) = 1$
 - $P(A \cup B) = P(A) + P(B)$ for disjoint events A and B
- Some consequences:
 - $P(A) = 1 - P(\Omega \setminus A)$
 - If $A \subseteq B$ then $P(A) \leq P(B)$
 - $P(A \cup B) = P(A) + P(B) - P(A \cap B) \leq P(A) + P(B)$

Probabilities over all disjoint subsets sum to 1:

$$\sum_{x \in X} P(X=x) = 1$$

Joint distribution: ~~set of~~ all possible combinations of events with their probability

	b_1	b_2	
a_1	0,05	0,54	0,59
a_2	0,08	0,04	0,12
a_3	0,17	0,12	0,29
		$\rightarrow P(A)$	
$\underbrace{0,30 \quad 0,70}$			
		\downarrow	
$P(B)$			

$$P(B) = P(a_1, b_1) + P(a_1, b_2) = 0,05 + 0,54 = 0,59$$

Conditional probability:

$$P(H|E) = \frac{P(H,E)}{P(E)}$$

Bayes rule:

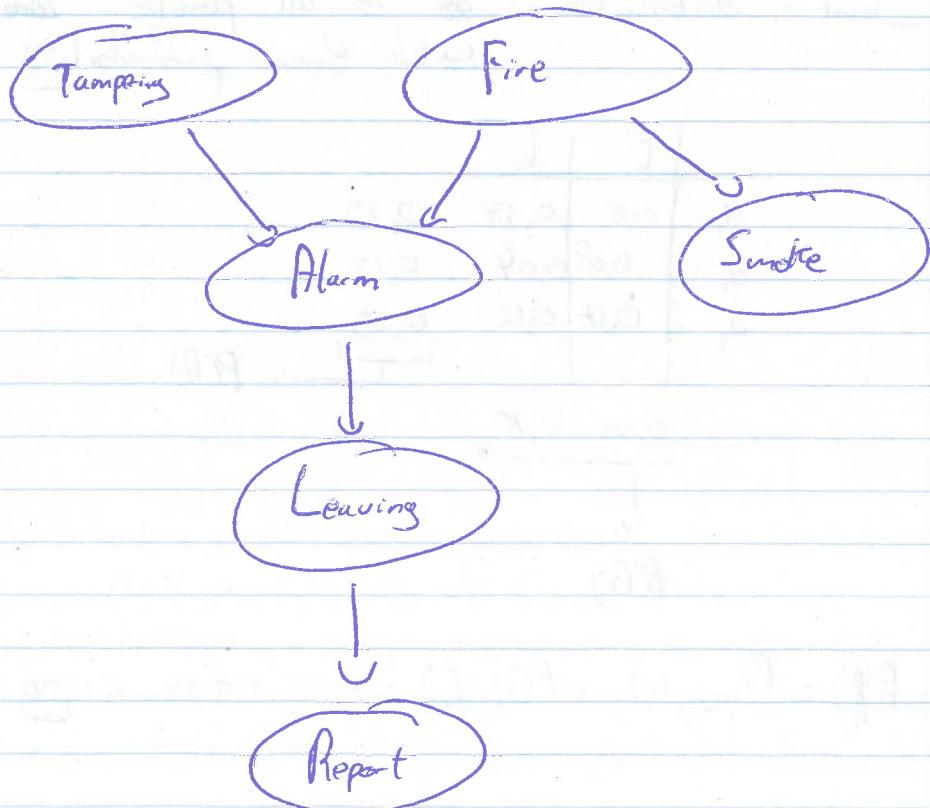
$$P(H|E) = \frac{P(H) \cdot P(E|H)}{P(E)}$$

Bayesian network:
- directed acyclic graph
- edges from parents to children

The joint probability distribution in a Bayesian network:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{Parents}_i(X_i))$$

b. Bayesian network



$$P(\text{tampering}) = 0,02$$

$$P(\text{fire}) = 0,01$$

$$P(\text{alarm} \mid \text{fire} \wedge \text{tampering}) = 0,5$$

$$P(\text{alarm} \mid \text{fire} \wedge \neg \text{tampering}) = 0,99$$

$$P(\text{alarm} \mid \neg \text{fire} \wedge \text{tampering}) = 0,85$$

$$P(\text{alarm} \mid \neg \text{fire} \wedge \neg \text{tampering}) = 0,0001$$

$$P(\text{smoke} \mid \text{fire}) = 0,9$$

$$P(\text{smoke} \mid \neg \text{fire}) = 0,01$$

$$P(\text{leaving} \mid \text{alarm}) = 0,98$$

$$P(\text{leaving} \mid \neg \text{alarm}) = 0,001$$

$$P(\text{report} \mid \text{leaving}) = 0,75$$

$$P(\text{report} \mid \neg \text{leaving}) = 0,01$$

$P(\text{Alarm})$ (composing $P(\text{Alarm})$):

→ We do not need leaving, report and smoke (branch nodes)

$$P(\text{Alarm}) = \sum_{\text{tampering, fire}} P(\text{Alarm} \mid \text{Tampering, Fire}) P(\text{Tampering}) P(\text{Fire}) =$$

$$P(\text{Alarm} \mid \text{tamp, fire}) P(\text{tamp}) P(\text{fire}) +$$

$$P(\text{Alarm} \mid \text{tamp, } \neg \text{fire}) P(\text{tamp}) P(\neg \text{fire}) +$$

$$P(\text{Alarm} \mid \neg \text{tamp, fire}) P(\neg \text{tamp}) P(\text{fire}) +$$

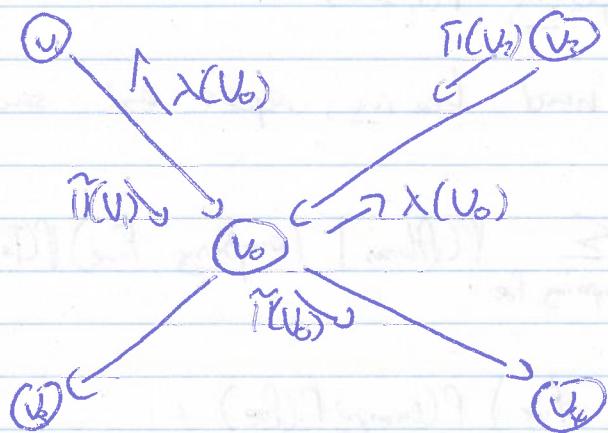
$$P(\text{Alarm} \mid \neg \text{tamp, } \neg \text{fire}) P(\neg \text{tamp}) P(\neg \text{fire}) = 0,0266$$

$$\text{And so: } P(\neg \text{Alarm}) = 0,9734$$

These computations can be done quickly:

- Look at nodes as objects that can pass around information to connected nodes. Keep computations local,
- Look at computations and find a smart way of reusing them, thus minimizing the amount of summations and multiplications

Message passing is focused on signals being sent between nodes



- $\tilde{\pi}(V)$ are messages sent to a child node
- $\lambda(V)$ are messages sent to a parent node

Every node knows its own local probability table

vb. Compute $P(V_0)$

$$P(V_0) = P(V_0|U_1, U_2) \cdot P(U_1) \cdot P(U_2)$$

V_0 knows $P(V_0|U_1, U_2)$ but does not know $P(U_1)$ and $P(U_2)$

V_1 sends $\tilde{\pi}(U_1) = P(U_1)$ to V_0
 V_2 sends $\tilde{\pi}(U_2) = P(U_2)$ to V_0

Message passing only works when a network is a tree, no cycles.

If a tree or cyclic:

- Loop-cut set conditioning \rightarrow condition on non-loop variables
- Clustering of variables

Variable elimination: solve a query by summing out variables in a logical order

\rightarrow Works on arbitrary networks (not only trees)

Factor: a factor $f(X_1, \dots, X_n)$ is a function from variables to a real value:

$$f(X_1, \dots, X_n) = r \in \mathbb{R}$$

\rightarrow A conditional/joint/marginal probability distribution can be represented by a factor.



However, a factor does not need to represent a particular distribution.

$$\begin{aligned} \text{Vb. } f_1(\text{Alarm}, \text{Tampering}) &\stackrel{\text{def}}{=} P(\text{Alarm} | \text{Tampering}) \\ \text{Vb. } f_2(\text{Alarm}) &\stackrel{\text{def}}{=} P(\text{Alarm} | \neg \text{temp}) \end{aligned}$$

Computations over factors:

$$f_1(A, B) \times f_2(B, C) = f_3(A, B, C) \text{ where}$$

$$f_3(a, b, c) = f_1(a, b) \cdot f_2(b, c)$$

u.b. $f_1(a_1, b_1) = 0,5$
 $f_2(b_1, c_1) = 0,3$

$$f_3(a_1, b_1, c_1) = 0,35$$

Factor marginalization: summing out a factor

$$\sum_B f_3(A, B, C) = f_4(A, C)$$

f_3

a_1	b_1	c_1	0,25
a_1	b_1	c_2	0,35
a_1	b_2	c_1	0,08
a_1	b_2	c_2	0,16
a_2	b_1	c_1	0,55
a_2	b_1	c_2	0,07
...

f_4

a_1	c_1	0,33
a_1	c_2	0,51
a_2	c_1	0,05
...

Factor reduction: reducing the amount of variables

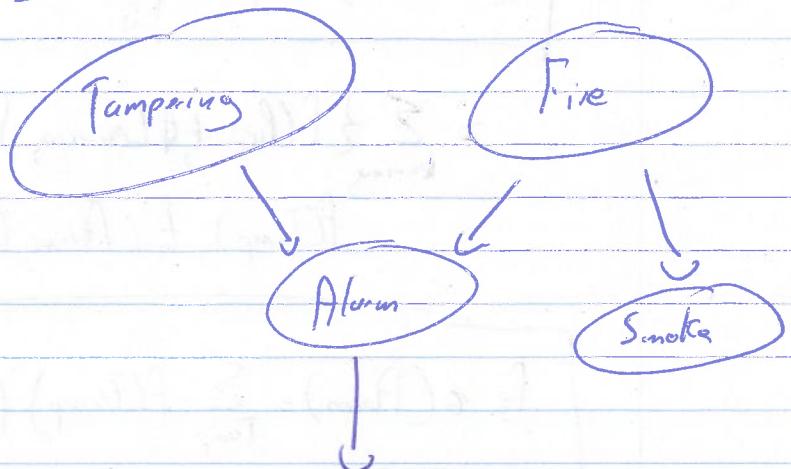
$$f_3(A, B, c_1) = f_5(A, B)$$

a_1	b_1	c_1	0,33
a_1	b_1	c_2	0,35
a_1	b_2	c_1	0,08
a_1	b_2	c_2	0,16
...

a_1	b_1	f_5	0,33
a_1	b_2		0,08
...

Variable elimination algorithm

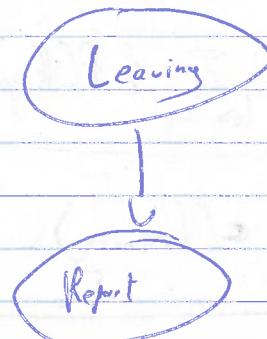
Qb. Compute $P(\text{Report})$



Elimination ordering:

$\text{Smoke} - \text{Fire} - \text{Tampering} -$

$\text{Alarm} - \text{Leaving}$



$$P(\text{Report}) = \sum_{\text{tamp, fire, alarm, smoke, leaving}} P(\text{Tampering}, \text{Fire}, \text{Alarm}, \text{Smoke}, \text{Leaving}, \text{Report})$$

$$= \sum_{\text{tamp, fire, alarm, smoke, leaving}} \{ P(\text{Tampering}) P(\text{Fire}) P(\text{Alarm}) P(\text{Smoke} | \text{Fire}) P(\text{Leaving} | \text{Alarm}) P(\text{Report} | \text{Leaving}) \}$$

Set summations
before puts
that need
them

$$= \sum_{\text{leaving}} \{ P(\text{Report} | \text{Leaving}) \cdot \sum_{\text{alarm}} \{ P(\text{Leaving} | \text{Alarm}) \} \cdot \sum_{\text{tamp}} \{ P(\text{Tamp}) \cdot \sum_{\text{fire}} \{ P(\text{Fire}) P(\text{Alarm} | \text{Fire}, \text{Tamp}) \sum_{\text{smoke}} \{ P(\text{Smoke} | \text{Fire}) \} \} \} \}$$

$$P(\text{Fire}) = \sum_{\text{smoke}} P(\text{Smoke} | \text{Fire}) = \sum_{\text{smoke}} \{ P(\text{Smoke} | \text{Fire}) \} = P(\text{Smoke} | \neg \text{Fire}) + P(\neg \text{Smoke} | \text{Fire})$$

$$= \sum_{\text{leaving}} \{ f(\text{Report} | \text{Leaving}) \sum_{\text{alarm}} \{ f(\text{Leaving} | \text{Alarm}) \} \sum_{\text{tamp}} \{ f(\text{Tamp}) \cdot \sum_{\text{fire}} \{ f(\text{Fire}) f(\text{Alarm} | \text{Fire}, \text{Tamp}) f(\text{Smoke} | \text{Fire}) \} \} \}$$

$$f_2(\text{Alarm}, \text{Tamp}) = \sum_{\text{Fire}} f(\text{Fire}) f(\text{Alarm}, \text{Tamp}, \text{Fire}) f_1(\text{Fire})$$

$$= \sum_{\text{Leaving}} \{ f(\text{Report}, \text{Leaving}) \sum_{\text{Alarm}} \{ f(\text{Leaving}, \text{Alarm}) \sum_{\text{Tamp}} \\ f(\text{Tamp}), f_3(\text{Alarm}, \text{Tamp}) \} \}$$

$$f_3 \otimes (\text{Alarm}) = \sum_{\text{Tamp}} f(\text{Tamp}) f_2(\text{Alarm}, \text{Tamp})$$

$$= \sum_{\text{Leaving}} \{ f(\text{Report}, \text{Leaving}) \otimes \sum_{\text{Alarm}} f(\text{Leaving}, \text{Alarm}) f_3(\text{Alarm}) \}$$

$$f_4 \otimes (\text{Leaving}) = \sum_{\text{Alarm}} f(\text{Leaving}, \text{Alarm}) f_3(\text{Alarm})$$

$$= \sum_{\text{Leaving}} \{ f(\text{Report}, \text{Leaving}) f_4(\text{Leaving}) \}$$

$$f_5(\text{Report}) = \sum_{\text{Leaving}} f(\text{Report}, \text{Leaving}) f_4(\text{Leaving})$$

$$= f_5(\text{Report})$$

\rightarrow Normalize if f_5 is not a probability distribution

- ① Write down query variables
- ② Write down observed variables
- ③ Write down:
 - the product \otimes to compute the query
 - the reduced formula based on the network structure
- ④ Identify factors and reduce observed variables
- ⑤ Fix an elimination ordering
- ⑥ For every variable Z in this ordering
 - Multiply factors containing Z
 - Sum out Z to obtain new factor f_2
 - Remove the multiplied factors from the list and add f_2

③ Normalize the result

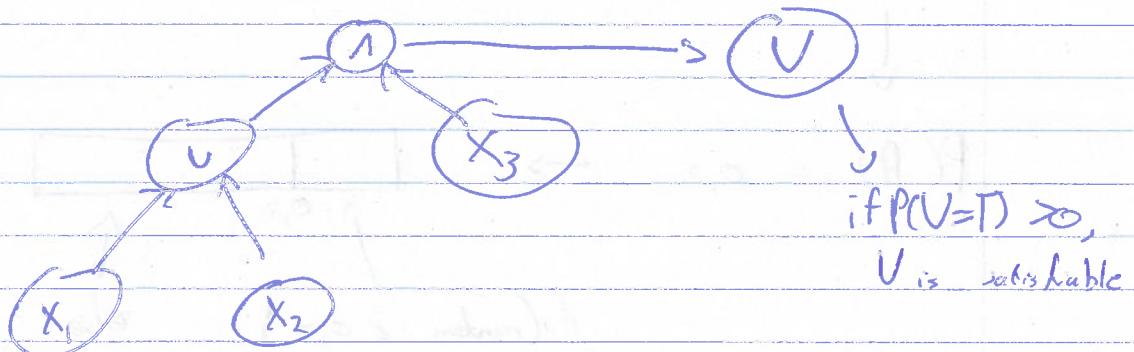
Finding good elimination orderings

- NP-Hard (optimal)
- Bayesian inference, even with an optimal ordering, is NP-Hard

Reduction from SAT:

For any formula ϕ we can construct a Bayesian network with designated variable V such that $P(V=\text{true}) > 0$ if ϕ is satisfiable.

ub. $(X_1 \vee X_2) \wedge X_3$:



X

Stochastic simulation: idea: get probabilities from samples

If we could sample from a variable's posterior probability distribution, we could estimate its posterior probability distribution.

Generating samples: a bayesian network is a-cyclic and directed

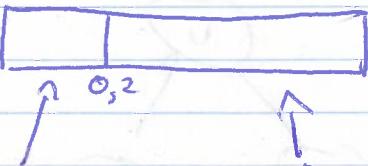


variables can be ordered such that node X is ordered after its parents

We have to generate a random numbers for each variable in the Bayesian network.



$$P(A) = 0,2 \rightarrow$$



if (random < 0,2)

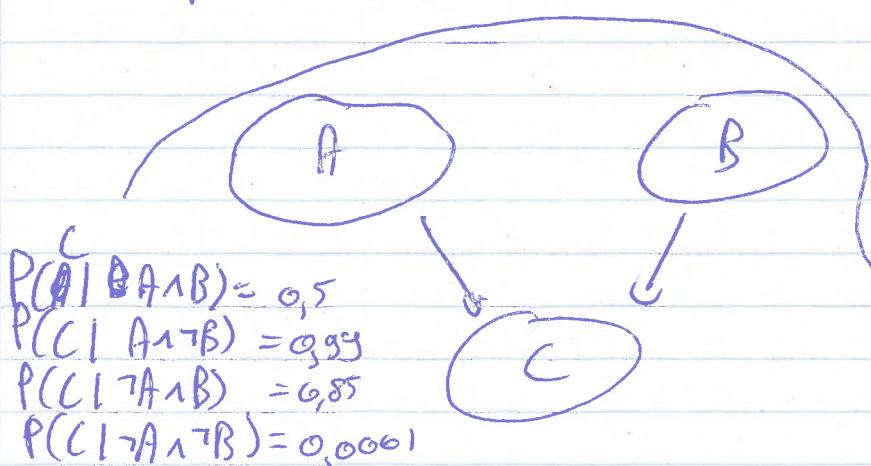
else



$P(A) = \text{True}$

$A = \text{False}$

with dependence:



1. First calculate A & B
2. Pick the probability from the options
3. Calculate C

Now we can repeat this a lot of times



Forward sampling

How many samples do we need?

$$P(|\hat{P}_{approx} - P_{real}| \geq \epsilon) \leq 2 \cdot e^{-2 \cdot \#samples \cdot \epsilon^2}$$

How many samples to ensure that in less than 0.8 fraction of all cases, the difference between estimated and actual distribution $\geq \epsilon$?



$$\text{Solve } 2 \cdot e^{-2 \cdot \#samples \cdot \epsilon^2} < 0.8 \text{ for } \#samples$$



$$\#samples \geq -\frac{\ln(0.8/2)}{2 \cdot \epsilon^2}$$



$$\text{if } \epsilon = 0.1, \delta = 0.05 \rightarrow \#samples > 184$$

Rejection sampling: reject samples that can't be used:

if, calculate $P(\text{sample} | \neg \text{smoke}, \text{true})$

Every sample that does not have "smoke or report will be rejected

Drawback: a lot of samples will be rejected, which is very inefficient

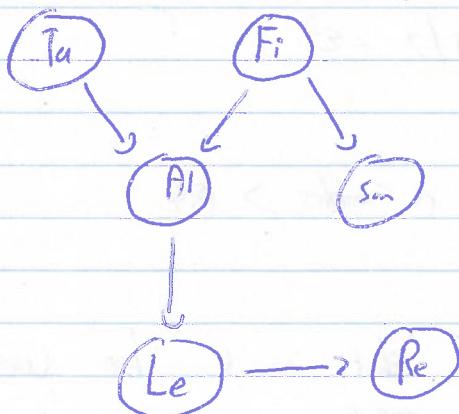
Importance sampling : don't reject sampler but weight sample according to how likely they are included

Probability of a proposition is weighted average of samples from a distribution

$$P(\alpha | \text{evidence}) \approx \frac{\sum_{\text{sample}} \text{weight}(\text{sample})}{\sum_{\text{sample}} \text{weight}(\text{sample})}$$

Likelihood weighting : special case of IS, sample only non-evidence variables and weight each example according to $P(\text{evidence} | \text{sample})$

ub. We want to know $P(\text{tall} | \text{sm, re})$



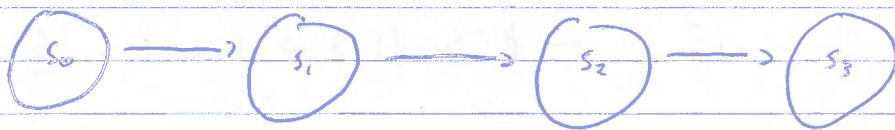
$$\begin{aligned} P(\text{sm} | \text{fi}) &= 0,9 \\ P(\text{sm} | \neg \text{fi}) &= 0,01 \\ P(\text{re} | \text{le}) &= 0,75 \\ P(\text{re} | \neg \text{le}) &= 0,01 \end{aligned}$$

	Ta	Fi	Al	Le	Weight
S ₁	true	false	true	false	$0,01 \times 0,01$
S ₂	false	true	false	false	$0,9 \times 0,01$
S ₃	false	true	true	true	$0,9 \times 0,75$
S ₄	true	true	false	true	$0,9 \times 0,75$

Drawback: many samples are needed

Approximate inference is in general NP-hard

Markov chains: special case of Bayesian network with chain structure



$$\text{Thus } P(S_{t+1} | S_0, \dots, S_t) = P(S_{t+1} | S_t)$$

Often S_t represents the state at time t

If we know S_t , we do not need to know S_0, \dots, S_{t-1} in order to compute S_{t+1} .



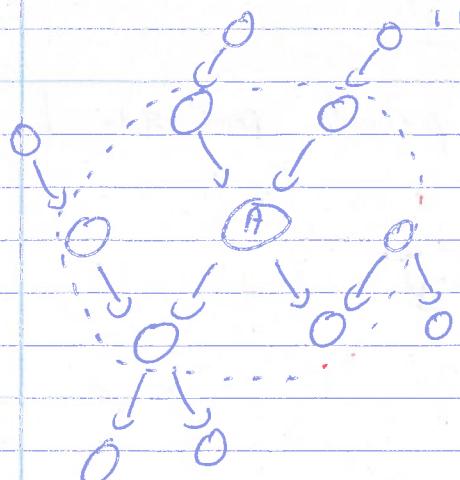
Conditionally independent

→ the future is independent of the past given the present



Markov property

Markov blanket: a set of nodes around a node shielding it from outside influences



- A is shielded from influences from outside by a joint value assignment to its parents, its children, and its co-parents

when these are all observed, all nodes to A from nodes outside the blanket are d-separated

Stationary Markov chain: a markov chain is stationary if:

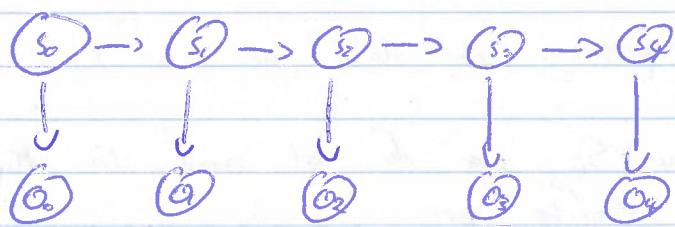
$$\text{for all } t \geq 0, t' \geq 0, P(S_{t+1} | S_t) = P(S_{t'+1} | S_{t'})$$

!

We specify only:

- $P(S_0)$ \rightarrow initial conditions
- $P(S_{t+1} | S_t)$ \rightarrow transition dynamics

Hidden Markov Model: special case of Bayesian network with the following structure



- $P(S_0)$ \rightarrow initial conditions
- $P(S_{t+1} | S_t)$ \rightarrow transition dynamics
- $P(O_t | S_t)$ \rightarrow sensor model

Queries to a HMM:

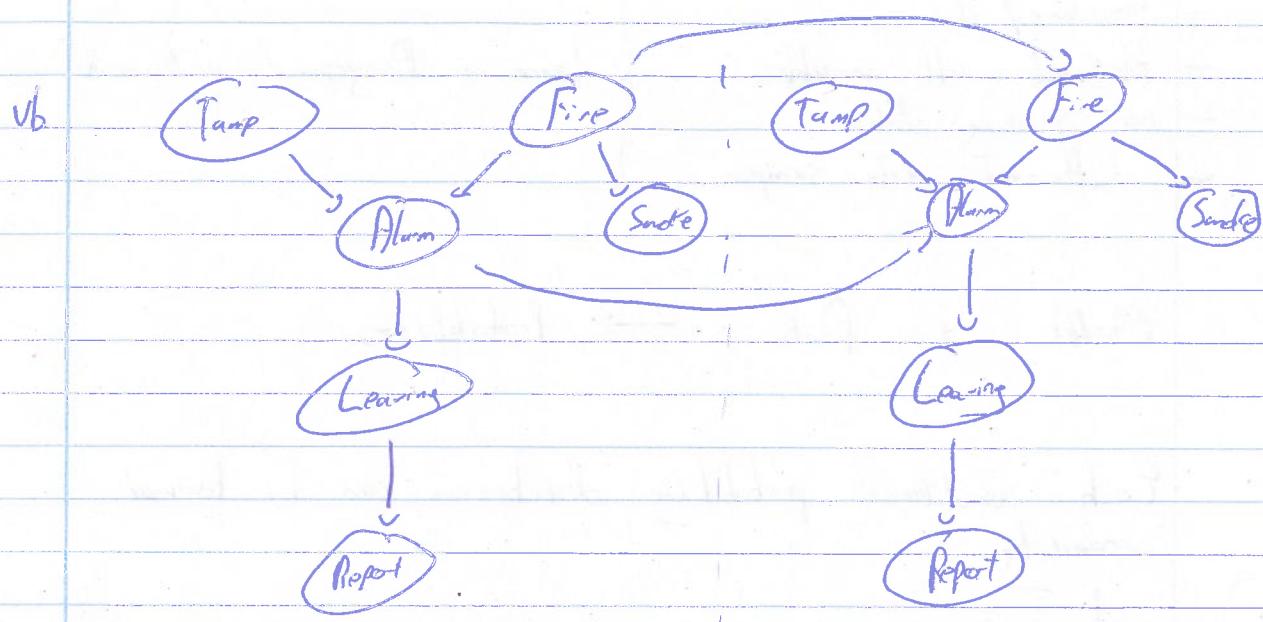
- What is the distribution of the current "true" state based on the observation history?

$P(S_t | O_0, \dots, O_t)$ (Filtering)

- What was the distribution of a particular past state based on the observation history

$P(S_t | O_0, \dots, O_T, O_t)$ (Smoothing)

Dynamic Bayesian Network: a HMM but does not need to satisfy Markov property



$$\mathbb{P}(t_0)$$

$$P(t_{n+1} | t_n)$$

Constructing Bayesian Networks: bayesian networks can be constructed or learned

Probability generation:

- Statistical domain knowledge
- Learning from data
- Knowledge from domain experts



Conditional probabilities are very hard to obtain

Learning Bayesian networks:

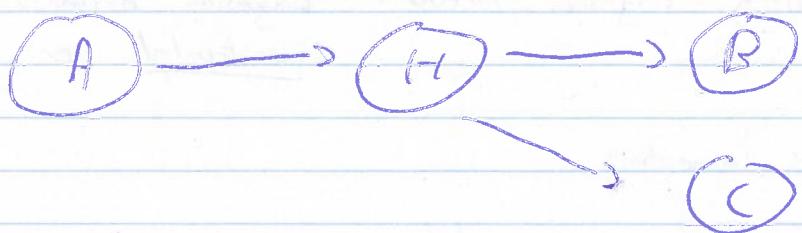
- structure / domain
 - observed all variables
 - no missing data
 - sufficient data samples
- } learn a Bayesian network

Model + Data \rightarrow Probabilities

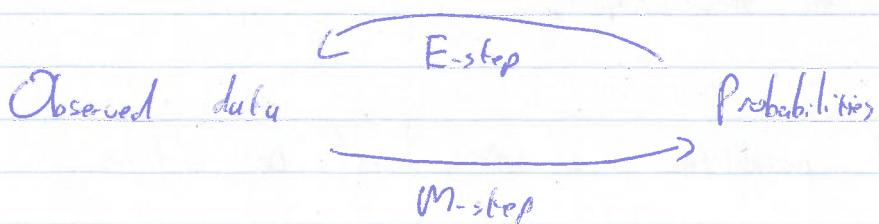
Each conditional probability distribution can be learned separately

! When there are many parents for a node, there can be little or no data for each probability estimate

Unobserved variables



What if we had only observed A, B & C?



Expectation - Maximization algorithm : used for computing probabilities of unobserved variables

1. Start with made-up probabilities (in this case to $P(H|A)$)

2. E-step : generate data points for unobserved variables, give the expected number of data points based on the given probability distribution

3. M-step : infer the probabilities from the generated data
4. Report.

Structure learning: learning the structure from data

→ Independence tests

XI Utility theory

Actions result in outcomes

if o_1 and o_2 are outcomes:

- $o_1 \succeq o_2$ means o_1 is at least as desirable as o_2
- $o_1 \sim o_2$ means $o_1 \succeq o_2$ and $o_2 \succeq o_1$
- $o_1 \succ o_2$ means $o_1 \succeq o_2$ and $o_2 \not\succeq o_1$

An agent may not know the outcome of its actions, but only have a probability distribution of the outcomes;

• Lottery: probability distribution over outcomes:

$$[p_1 : o_1, p_2 : o_2, \dots, p_K : o_K]$$

with p_1, \dots, p_K being probabilities

$$\Rightarrow \sum_i p_i = 1$$

and o_1, \dots, o_K being outcomes

Completeness: Agents have to act, so they must have preferences: H_0 : $o_1 \succeq o_2$ or $o_2 \succeq o_1$

For any given pair of courses, passing one is at least as desirable as passing the other.

Transitivity: preferences must be transitive

$$\text{if } o_1 \succeq o_2 \text{ and } o_2 \succ o_3 \text{ then } o_1 \succ o_3$$

Monotonicity: an agent prefers a large chance of getting a better outcome than a smaller chance.

If $\omega_1 \succ \omega_2$ and $p > q$ then $\left[p:\omega_1, 1-p:\omega_2 \right] \succ \left[q:\omega_1, 1-q:\omega_2 \right]$

Continuity: suppose $\omega_1 \succ \omega_2$ and $\omega_2 \succ \omega_3$, then there exist $p \in [0, 1]$ such that:

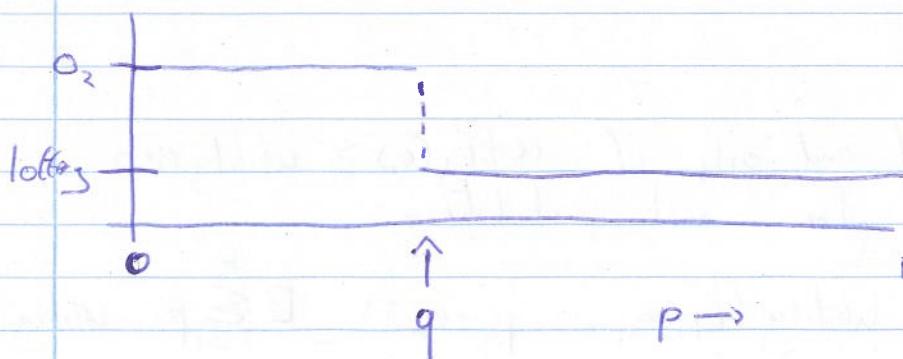
$$\omega_2 \sim [p:\omega_1, 1-p:\omega_3]$$

Suppose $\omega_1 \succ \omega_2$ and $\omega_2 \succ \omega_3$, consider the case where the agent would prefer:

$$-\omega_2 - \text{the lottery } [p:\omega_1, 1-p:\omega_3]$$

These for different values of $p \in [0, 1]$

Then:



Decomposability: an agent is indifferent between lotteries that have same probabilities and outcomes

$$[p_1: o_1, 1-p_1: [q: o_2, 1-q: o_3]] \sim [p: o_1, (1-p)q: o_2, (1-p)(1-q): o_3]$$

Substitutability: if $o_1 \sim o_2$ then the agent is indifferent between lotteries that only differ by o_1 and o_2

$$[p: o_1, 1-p: o_3] \sim [p: o_2, 1-p: o_3]$$

We need a measure of preference that can be combined with probabilities.

Preferences can be measured by a function:

utility: outcomes $\rightarrow [0, 1]$

such that:

- $o_1 \succ o_2$ if and only if utility(o_1) \geq utility(o_2)
- utilities are linear with probabilities:

$$\text{utility}([p_1: o_1, \dots, p_k: o_k]) = D \sum_{i=1}^k p_i \cdot \text{utility}(o_i)$$

Affairs paradox

What would you prefer?

A: \$1m

B: lottery: [0,10: \$2.5m, 0.89: \$1m, 0.01: \$0]

What would you prefer?

C: lottery: [0,11: \$1m, 0.89: \$0]

D: lottery: [0,10: \$2.5m, 0.9: \$0]

It is inconsistent with the axioms of preferences to have:

A \succ B, D \succ C

A, C : lottery [0,11: \$1m, 0.89: X]

B, D : lottery [0,10: \$2.5m, 0.01: \$0, 0.89: X]

Framing effect: by putting the question differently, people value certain choices more or less

What an agent should do depends on:

- The agents ability: possible options
- the agents belief: the way the world could be
- The agents preference

Decision variables: variables an agent gets to choose a value for

Decision theory

Expected value: the average value, weighting possible worlds by their probability

$f(\omega)$: value of f in world ω

$$E(F) = \sum_{\omega \in \Omega} P(\omega) \cdot f(\omega)$$

$$E(F|e) = \sum_{\omega \in e} P(\omega|e) \cdot f(\omega)$$

Single decisions: agent makes all decisions before acting

The agent can choose a value for each decision variable

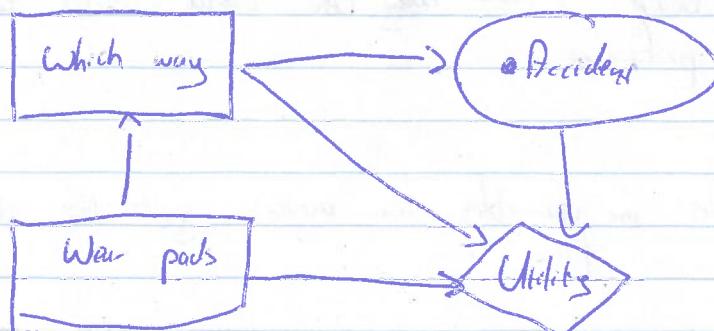
Expected utility of decision $D = d_i$: $E(u|D=d_i)$



Decision networks:

- Belief network
- Decision nodes: the agent chooses a value for these
- Utility nodes: parents are the variables on which the utility depends

vb.



• A single-stage decision network consists of:

- Directed acyclic graph with decision, random and utility nodes
- A domain for each decision variable and each random variable
- A unique utility node

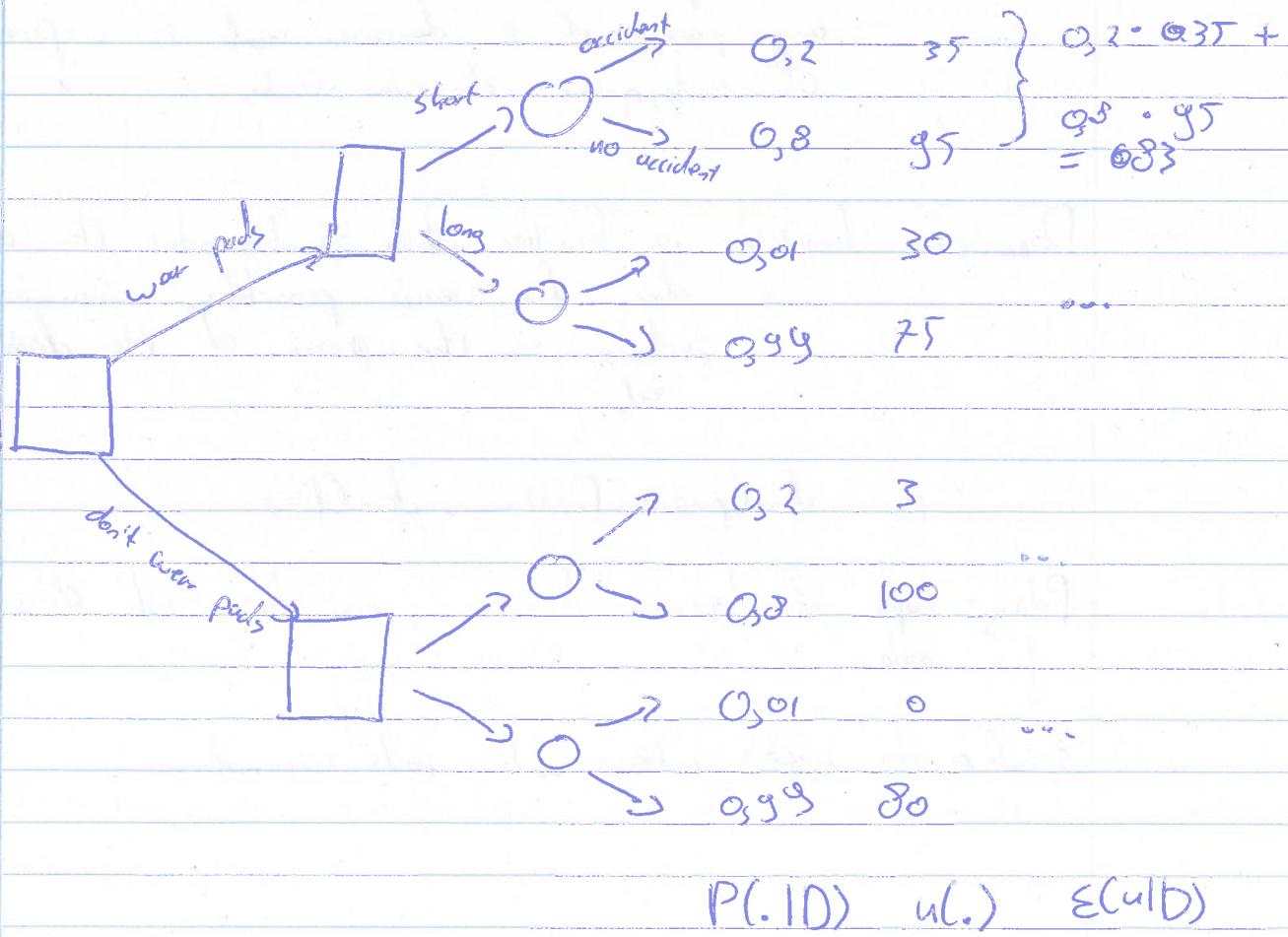
Factors:

- Utility function: factor on the parents of the utility node
- Conditional probability for each random variable given its parents

Finding an optimal decision:

- Create a factor for each conditional probability and for the utility
- Sum out all of the random variables
- This creates a factor on D that gives the expected utility for each value in the domain of D
- Choose D with the maximum value in the factor

Vb.



P(D) u(.) E(u|D)

Sequential decisions: actions can depend on what is observed

Sequential decision problem: sequence of decision variables

$$D_1, \dots, D_m$$

Each D_i has an information set of variables parents(D_i), whose value will be known at the time decision D_i is made.

No-forgetting: a no-forgetting decision network is a network where:

- the decisions are totally ordered. This is the order in which the actions will be taken
- all decision nodes that come before D_i are parents of decision node D_i .
- any parent of a decision node is a parent of subsequent decision nodes.

Decision function: a function which specifies what the agent does for each possible assignment of values to the parent of the decision node.

$$\delta_i : \text{dom}(\text{parents}(D_i)) \rightarrow \text{dom}(D_i)$$

Policy: tuple of decision functions, one for each decision node

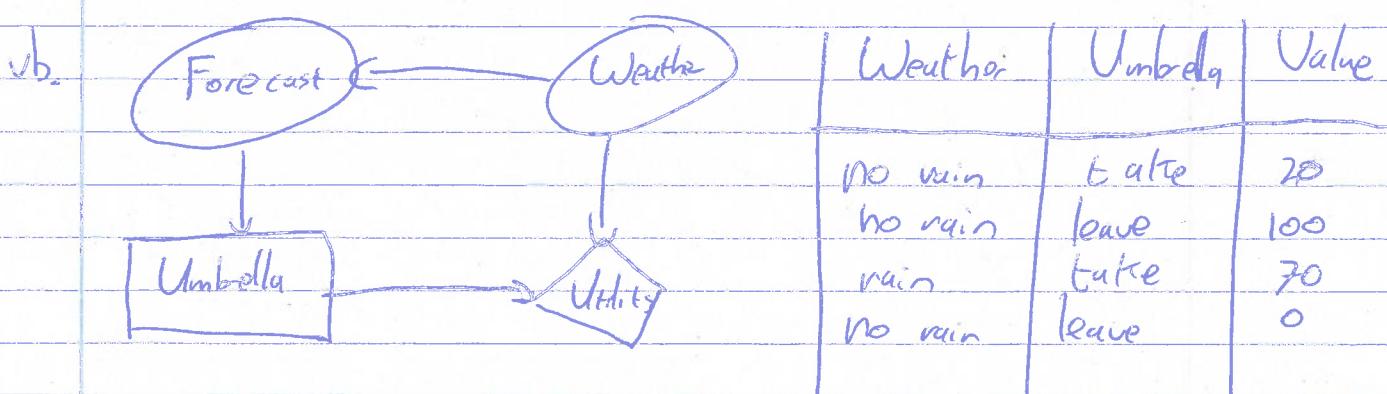
specifies an agent's actions under such circumstance

Optimal policy: policy with the highest utility

$$E(u|s) = \sum_{\omega \in S} u(\omega) \cdot P(\omega)$$

Finding an optimal policy

1. Create a factor for each conditional probability table and a factor for the utility
2. Sum out random variables that are not parents of a decision node
3. Let D be the last variable - D is only in a factor f with its parents
4. Eliminate D by maximizing:
 - an optimal decision function for D
 - a new factor $\max_D f$
5. Repeat 2 until there are no more decision nodes
6. Sum out the remaining random variables. Multiply the factors: this is the expected utility of an optimal policy



Weather	Forecast	value	Weather	value
no rain	sunny	0,7	no rain	0,7
no rain	cloudy	0,2	rain	0,3
no rain	rainy	0,1		
rain	sunny	0,15		
rain	cloudy	0,25		
rain	rainy	0,6		

forecast	umbrella	value
sunny	false	$0,7 \cdot 0,7 \cdot 20 + 0,3 \cdot 0,15 \cdot 0,7 = 12,95$
sunny	true	$0,7 \cdot 0,7 \cdot 100 + 0,3 \cdot 0,15 \cdot 0 = 89$
cloudy	false	$0,7 \cdot 0,2 \cdot 20 + 0,3 \cdot 0,75 \cdot 20 = 8,05$
cloudy	true	$0,7 \cdot 0,2 \cdot 100 + 0,3 \cdot 0,25 \cdot 0 = 14$
raining	false	$0,7 \cdot 0,1 \cdot 20 + 0,3 \cdot 0,6 \cdot 20 = 14$
raining	true	$0,7 \cdot 0,1 \cdot 100 + 0,3 \cdot 0,6 \cdot 0 = 7$

Now maximizing D:

Forecast	value
sunny	49
cloudy	14
raining	14

Max umbrella F

forecast	umbrella
sunny	leave
cloudy	leave
raining	false

ang max umbrella F

XII

Decision processes

Utility and time: how to trade off pleasures today with pleasures in the future?

Ongoing processes: checking utility at the end might not be useful



agent might not
get to the end

→ A sequence of rewards might be useful

Rewards:

- action costs
- prizes
- penalties

Finite horizon: fixed number of steps

Infinite horizon: never ending number of steps

Indefinite horizon: agent will eventually stop but does not know when

World state: information such that if the agent knew the world state no information about the past is relevant to the future

↳ Markovian assumption

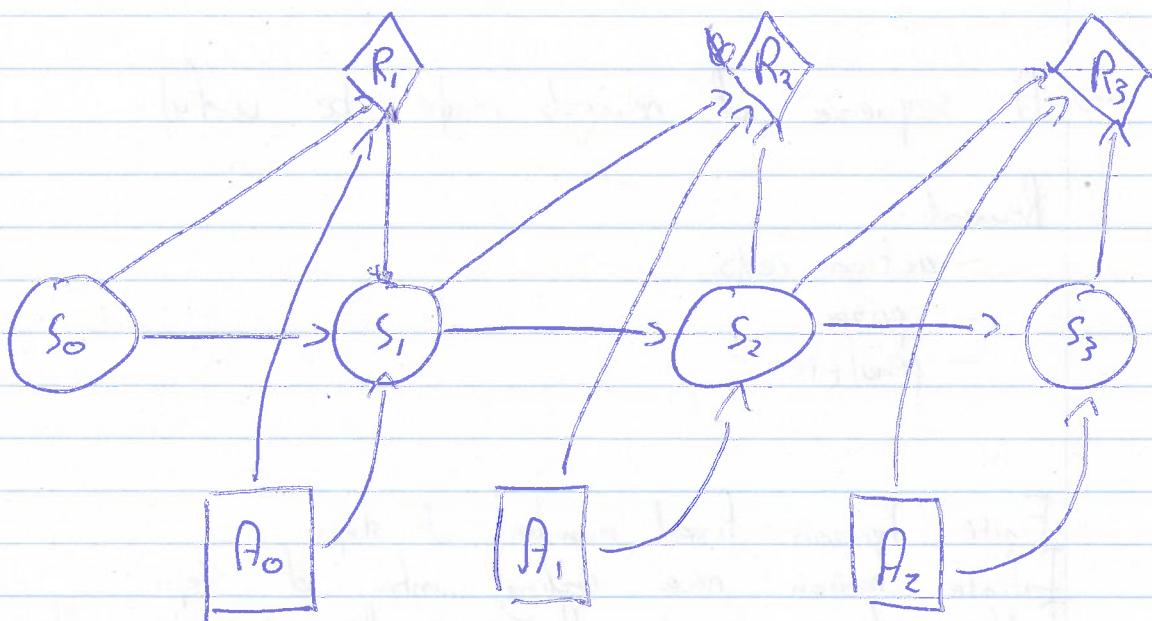
$$P(S_{t+1} | S_0, A_0, \dots, S_t, A_t) = P(S_{t+1} | S_t, A_t)$$

$P(S'_t | S_t, a)$: probability that the agent will be in state s' after performing action a in state s

Markov chain: belief network in which with random variables where each variable only directly depends on its predecessor in the sequence



Markov Decision Process: augments a Markov chain with actions and values



At each stage, the agent decides what to do, the reward and resulting state depend on both the previous state and the action performed

- MDP:
 - set of States S
 - set of Actions A
 - $P(S_{t+1} | S_t, A_t)$ specifies the dynamics
 - $R(S_t, A_t, S_{t+1})$ specifies the reward at time t
 - $R(s, a, s')$ is the expected reward received
 - $R(s, a) = \sum_s P(s' | s, a) R(s, a, s')$
 - γ is the discount factor

Planning horizon: how far ahead the planner looks to make a decision

Fully observable MDP: the agent gets to observe S_t when deciding on action A_t

Partially observable MDP: the agent has some noisy sensor of the state.

To decide what to do, agent compares different sequences of rewards



Agent receives sequence of rewards $r_1, r_2, r_3, \dots, r_n$



combining rewards

Total reward: $V = \sum_{i=1}^{\infty} r_i$

Average reward: $V = \lim_{n \rightarrow \infty} (r_1 + \dots + r_n)/n$

Discounted reward: $V = r_1 + \gamma r_2 + \gamma^2 r_3 + \gamma^3 r_4 + \dots$

with

where γ is the discount factor
 $0 \leq \gamma \leq 1$

Policy: specification of what the agent should do as a function of the state it is in.

Stationary policy: $\pi: S \rightarrow A$

An agent following π in state s will execute $\pi(s)$ next

Optimal policy: a policy with the maximum expected discounted reward

For a Fully observable MDP with stationary dynamics and rewards with infinite or indefinite horizon, there is always an optimal stationary policy

Ub. Policy one:

States: {Healthy, sick}

Actions: {party, relax}

Available policies: $| \text{Actions} |^{|\text{States}|} = 2^2 = 4$

1. $\pi[\text{Healthy}] = \text{party}, \pi[\text{sick}] = \text{party}$
2. $\pi[\text{Healthy}] = \text{relax}, \pi[\text{sick}] = \text{relax}$
3. $\pi[\text{Healthy}] = \text{party}, \pi[\text{sick}] = \text{relax}$
4. $\pi[\text{Healthy}] = \text{relax}, \pi[\text{sick}] = \text{party}$

- $Q^\pi(s, a)$ is the expected value of doing action a starting in state s , then following policy π
- $V^\pi(s)$ is the expected value of following policy π in state s

+ @
Bellman equations:

$$Q^\pi(s, a) = R(s, a) + \gamma \cdot \sum_{s'} P(s' | s, a) \cdot V^\pi(s')$$

$$\text{where } R(s, a) = \sum_{s'} P(s' | s, a) \cdot R(s, a, s')$$

$$\text{and } V^\pi(s) = Q^\pi(s, \pi(s))$$

- $Q^*(s, a)$: expected value of doing a in state s then following the optimal policy
- $V^*(s)$: expected value of following the optimal policy in state s
- Bellman optimality equations:

$$Q^*(s, a) = R(s, a) + \sum_{s'} P(s' | s, a) \gamma V^*(s')$$

$$V^*(s) = \max_a Q^*(s, a)$$

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

Value iteration: algorithm for finding the optimal policy

Inputs:

- set S of states

- set A of actions

- P is a state transition function specifying $P(s' | s, a)$

- R is a rewards function $R(s, a)$

Outputs:

- $\pi^*(s)$, the optimal policy
- $V^*(s)$, the optimal value function } approximately

Algorithm:

1. Assign $V_0^*(s)$ arbitrarily

2. For each k and s :

$$V_k^*(s) = \max_a Q_k^*(s, a)$$

$$Q_k^*(s, a) = R(s, a) + \gamma \cdot \sum_{s'} P(s' | s, a) \cdot V_{k-1}^*(s')$$

For each s :

$$\pi(s) = \arg \max_a Q_{k_{\max}}^{\pi}(s, a)$$

$$\text{with } Q_{k_{\max}}^{\pi}(s, a) = R(s, a) + \gamma \cdot \sum_{s'} P(s'|s, a) \cdot V_{k_{\max}}^{\pi}(s')$$

vb.

s	A	$P(\text{healthy} s, A)$	$P(\text{sick} s, A)$	R
healthy	party	0,7	0,3	10
healthy	relax	0,95	0,05	7
sick	party	0,1	0,9	2
sick	relax	0,5	0,5	0

b. Assign $V_0^{\pi}(s) = 0$

$$V_0^{\pi}(\text{healthy}) = 0$$

$$V_0^{\pi}(\text{sick}) = 0$$

2. For each k and s : $V_k^{\pi}(s) = \max_a Q_k^{\pi}(s, a)$

$$\text{where } Q_k^{\pi}(s, a) = R(s, a) + \gamma \cdot \sum_{s'} P(s'|s, a) \cdot V_{k+1}^{\pi}(s')$$

$$\begin{aligned}
 Q_1(\text{healthy}, \text{party}) &= R(\text{healthy}, \text{party}) + \gamma \cdot (P(\text{healthy}|\text{healthy, party}) \cdot V_0(\text{healthy}) \\
 &\quad + P(\text{sick}|\text{healthy, party}) \cdot V_0(\text{sick})) \\
 &= 10 + 0,8 \cdot (0,7 \cdot 0 + 0,3 \cdot 0) = 10
 \end{aligned}$$

$$Q_1(\text{healthy}, \text{relax}) = 0,7$$

$$V_1(\text{healthy}) = \max(Q_1(\text{healthy}, \text{party}), Q_1(\text{healthy}, \text{relax})) = \max(10, 0,7) = 10$$

$$\begin{aligned}
 \pi(\text{healthy}) &= \arg \max \{ Q_1(\text{healthy}, \text{party}), Q_1(\text{healthy}, \text{relax}) \} \\
 &= \arg \max \{ 10, 0,7 \} = \text{party}
 \end{aligned}$$

Policy iteration

1. Set π_0 arbitrarily, let $i=0$

2. Repeat:

- evaluate $Q^{\pi_i}(s, a)$

- let $\pi_{i+1}(s) = \underset{a}{\operatorname{argmax}} Q^{\pi_i}(s, a)$

- set $i = i + 1$

until $\pi_i(s) \approx \pi_{i-1}(s)$

Finding a solution to a set of $|S| \cdot |A|$ unknowns

Solve:

$$V^{\pi}(s) = R(s, \pi(s)) + \gamma \cdot \sum_{s'} P(s'|s, \pi(s)) \cdot V^{\pi}(s')$$

for each s :

$$V^*(s) = \max_a Q^{\pi}(s, a)$$

where $Q^{\pi}(s, a) = \underset{s'}{\operatorname{argmax}} Q^{\pi}(s, a)$

$$\text{where } Q^{\pi}(s, a) = R(s, a) + \gamma \cdot \sum_{s'} P(s'|s, a) \cdot V^{\pi}(s')$$

Vb.

1. Assign $\pi(s)$ arbitrarily:

$$\begin{aligned} \pi(\text{healthy}) &= \text{party} \\ \pi(\text{sick}) &= \text{relax} \end{aligned}$$

$$2. V^{\pi}(\text{healthy}) = R(\text{healthy}, \text{party}) + \gamma \cdot (P(\text{healthy} | \text{healthy, party}) \cdot V^{\pi}(\text{healthy}) + P(\text{sick} | \text{healthy, party}) \cdot V^{\pi}(\text{sick}))$$

$$V^{\pi}(\text{sick}) = R(\text{sick}, \text{relax}) + \gamma \cdot (P(\text{healthy} | \text{sick, relax}) \cdot V^{\pi}(\text{healthy}) + P(\text{sick} | \text{sick, relax}) \cdot V^{\pi}(\text{sick}))$$

$$\text{e.g. } 1 = 0 + 0.8 \cdot (0.7 \cdot V^{\pi}(\text{healthy}) + 0.3 \cdot V^{\pi}(\text{sick}))$$

drive them both

$$Q''(\text{healthy, party}) = R(\text{healthy, party}) + \gamma \cdot (P(\text{healthy}) \cdot Q''(\text{healthy, party}) \\ + P(\text{sick}) \cdot Q''(\text{sick}))$$

$$= 10 \cdot 0,8 \cdot (0,7 \cdot 35,71 + 0,3 \cdot 23,81) \\ = 35,71$$

$$Q''(\text{healthy, relax}) = 35,10$$

$$Q''(\text{sick, party}) = 22$$

$$Q''(\text{sick, relax}) = 23,81$$

]

$$\hat{a}_{\text{sick}} = \text{relax}$$

XIII

Learning types:

- Supervised learning: learning from a teacher
 - ↳ training data includes desired outputs
- Unsupervised learning:
 - ↳ training data does not include desired outputs
- Reinforcement learning: learning to act under evaluative feedback



- most like natural learning
- trial and error

→ Reward good behaviour; punish bad behaviour

Reinforcement learning agent

- The learning agent is given the possible states and set of actions it can carry out
- The agent observes the state of the environment and the reward received
- The goal of the agent is to maximize its reward

Experiences: each time the agent has to choose its action, requires it to:

- explore to gain knowledge
- exploit acquired knowledge

- The agent initially knows all states S and all actions A .
- The dynamics $P(s'|s, a)$ and $R(s, a)$ are not given to the agent



works like MDP

Reinforcement learning is hard

- A reward may be caused by a lot of actions before
- A long-term effect of an action depend on what the agent will do in the future
- The explore-exploit dilemma: at what time should the agent be greedy or inquisitive?

Main approaches to RL:

- Search through a space of policies
- Learn a model consisting of state transition function $P(s'|s, a)$ and reward function $R(s, a, s')$, solve this as an MDP
- Learn $Q^*(s, a)$, use this to guide action.

Suppose we have a sequence of values: v_1, v_2, \dots, v_k and we want a running average of the first k

$$A_k = \frac{v_1 + v_2 + \dots + v_k}{k}$$

TD Formula :

$$\alpha_k = \frac{1}{k}$$

$$\begin{aligned} A_k &= (1 - \alpha_k) A_{k-1} + \alpha_k v_k \\ &= A_{k-1} + \alpha_k [v_k - A_{k-1}] \end{aligned}$$

TD Error: specifies how different the new A_k is

Q-learning: an agent tries to learn the optimal policy from its history

$$\langle s_0, a_0, r_0, s_1, a_1, r_1, \dots \rangle$$

$$\text{Experience} = \langle s, a, r, s' \rangle$$

Q -learning uses temporal differences to estimate the value of $Q^*(s, a)$

↓
the agent maintains a table $Q(s, A)$ where $Q(s, a)$ represents its current estimate of $Q^*(s, a)$

↓
 (s, a, r, s')
an experience provides one data point for $Q(s, a)$

↓
shows that the agent received the future value of $r + \gamma V(s')$

$V(s') = \max_a Q(s', a) \rightarrow$ the actual current reward plus the discounted estimated future value
↳ data point is called a return

↓
We can also use the experience directly on $Q(s, a)$:

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Algorithm:

1. Initialize $Q(s, A)$ arbitrarily

2. Observe s

Repeat forever:

- select and carry out an action a

- observe reward r and state s'

- $Q(s, a) = Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a') - Q(s, a))$

- $s \leftarrow s'$

It does not yet explore knowledge

Exploration strategies:

- ϵ -greedy strategy: choose random action with probability ϵ
or best action with probability $1 - \epsilon$
- Softmax: choose action a with probability:

$$\frac{e^{Q(s,a)/T}}{\sum_a e^{Q(s,a)/T}}$$

where T is the temperature

- optimism in the face of uncertainty: initialize Q to values that encourage exploration
- upper confidence bounds: take into account average + variance

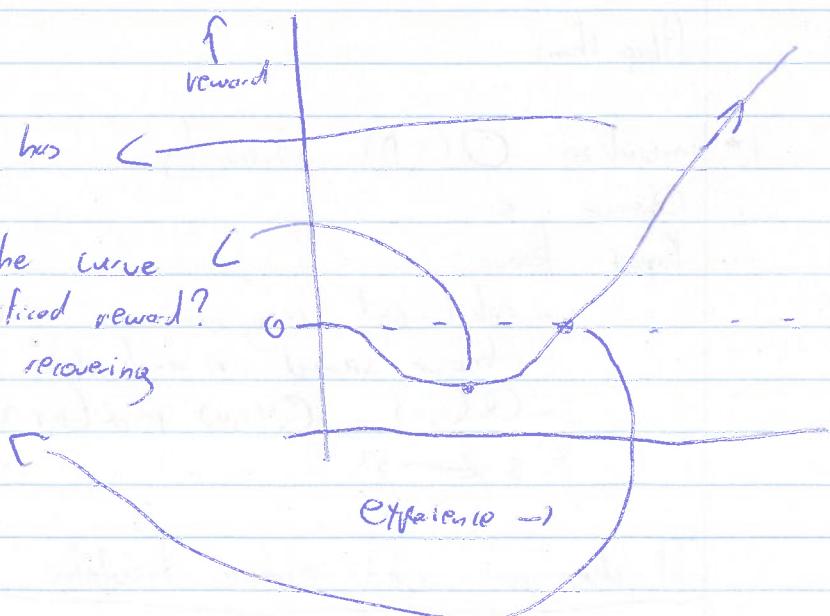
Problems with Q-learning:

- it does one backup between each experience
- it learns separately for each state

Evaluating RL algorithms: use the cumulative reward (plot)

Three statistics:

- slope after algorithm has stabilized
- variance minimum of the curve
 \hookrightarrow how much sacrificed reward?
- time it takes for recovering
 $(0 - \text{crossing})$



SARSA: on-policy learning



also uses a' in experience: $\langle s, a, r, s', a' \rangle$ to update $Q(s, a)$

↳ remember previous action a^*

